

General Strong Polarization

JAROSŁAW BŁASIOK, Department of Computer Science, Columbia University, USA VENKATESAN GURUSWAMI, Computer Science Department, Carnegie Mellon University, USA PREETUM NAKKIRAN, Halicioğlu Data Science Institute, University of California San Diego, USA ATRI RUDRA, Computer Science and Engineering Department, University at Buffalo, USA MADHU SUDAN, Harvard John A. Paulson School of Engineering and Applied Sciences, Harvard University, USA

Arıkan's exciting discovery of polar codes has provided an altogether new way to efficiently achieve Shannon capacity. Given a (constant-sized) invertible matrix M, a family of polar codes can be associated with this matrix and its ability to approach capacity follows from the *polarization* of an associated [0, 1]-bounded martingale, namely its convergence in the limit to either 0 or 1 with probability 1. Arıkan showed appropriate polarization of the martingale associated with the matrix $G_2 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$ to get capacity achieving codes. His analysis was later extended to all matrices M that satisfy an obvious necessary condition for polarization.

While Arıkan's theorem does not guarantee that the codes achieve capacity at small blocklengths (specifically in length, which is a polynomial in $1/\varepsilon$ where ε is the difference between the capacity of a channel and the rate of the code), it turns out that a "strong" analysis of the polarization of the underlying martingale would lead to such constructions. Indeed for the martingale associated with G_2 such a strong polarization was shown in two independent works (Guruswami and Xia (IEEE IT'15) and Hassani et al. (IEEE IT'14)), thereby resolving a major theoretical challenge associated with the efficient attainment of Shannon capacity.

In this work we extend the result above to cover martingales associated with all matrices that satisfy the necessary condition for (weak) polarization. In addition to being vastly more general, our proofs of strong

This article combines results presented in preliminary form at STOC 2018 [4] and RANDOM 2018 [5]. Jarosław Blasiok work was done when the author was a Ph.D. student at Harvard University, supported by ONR grant N00014-15-1-2388.

Venkatesan Guruswami portions of this work were done during visits by the author to the School of Physical and Mathematical Sciences, Nanyang Technological University, Singapore, and the Center for Mathematical Sciences and Applications, Harvard University. Research supported in part by NSF grants CCF-1422045, CCF-1563742 and CCF-1814603, a Packard Fellowship, and a Simons Investigator award.

Preetum Nakkiran work supported in part by a Simons Investigator Award, NSF Awards CCF 1565641 and CCF 1715187, the NSF Graduate Research Fellowship Grant No. DGE1144152, a Google PhD Fellowship, and the NSF/Simons Collaboration on the Theoretical Foundations of Deep Learning.

Atri Rudra research supported in part by NSF grant CCF-1717134.

Madhu Sudan work supported in part by a Simons Investigator Award and NSF Awards CCF 1565641 and CCF 1715187. Authors' addresses: J. Błasiok, Department of Computer Science, Columbia University, 500 West 120th Street, New York, NY, 10027; email: jb4451@columbia.edu; V. Guruswami, Computer Science Department, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213; email: venkat@cs.cmu.edu; P. Nakkiran, Halicioğlu Data Science Institute, University of California San Diego, 10100 Hopkins Dr, La Jolla, CA, 92093; email: preetum@ucsd.edu; A. Rudra, Department of Computer Science and Engineering, Davis 338, University at Buffalo, Buffalo, NY, 14260; email: atri@buffalo.edu; M. Sudan, Harvard John A. Paulson School of Engineering and Applied Sciences, Harvard University, 33 Oxford Street, Cambridge, MA, 02138; email: madhu@cs.harvard.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

0004-5411/2022/03-ART11 \$15.00

https://doi.org/10.1145/3491390

11:2 J. Błasiok et al.

polarization are (in our view) also much simpler and modular. Key to our proof is a notion of *local polarization* that only depends on the evolution of the martingale in a single time step. We show that local polarization always implies strong polarization. We then apply relatively simple reasoning about conditional entropies to prove local polarization in very general settings. Specifically, our result shows strong polarization over all prime fields and leads to efficient capacity-achieving source codes for compressing arbitrary i.i.d. sources, and capacity-achieving channel codes for arbitrary symmetric memoryless channels. We show how to use our analyses to achieve exponentially small error probabilities at lengths inverse polynomial in the gap to capacity. Indeed we show that we can essentially match any error probability while maintaining lengths that are only inverse polynomial in the gap to capacity.

CCS Concepts: • Theory of computation \rightarrow Error-correcting codes; • Mathematics of computing \rightarrow Coding theory; Stochastic processes;

Additional Key Words and Phrases: Polar codes, polarization, capacity-achieving codes

ACM Reference format:

Jarosław Błasiok, Venkatesan Guruswami, Preetum Nakkiran, Atri Rudra, and Madhu Sudan. 2022. General Strong Polarization. *J. ACM* 69, 2, Article 11 (March 2022), 67 pages. https://doi.org/10.1145/3491390

1 INTRODUCTION

Polar codes, proposed in Arıkan's remarkable work [2], gave a fresh information-theoretic approach to construct linear codes that achieve the Shannon capacity of symmetric channels, together with efficient encoding and decoding algorithms. About a decade after their discovery, there is now a vast and extensive body of work on polar coding spanning hundreds of papers. The underlying concept of polarizing transforms has emerged as a versatile tool to successfully attack a diverse collection of information-theoretic problems beyond the original channel and source coding applications, including wiretap channels [22], the Slepian-Wolf, Wyner-Ziv, and Gelfand-Pinsker problems [19], broadcast channels [11], multiple access channels [1, 8], and interference networks [31]. We recommend the survey by Şaşoğlu [7] for a nice treatment of the early work on polar codes. On the practical side, polar codes show impressive coding gains when a list decoding variant of the decoder is applied [29] and have been adopted for the enhanced mobile broadband control channels for the 5G NR (New Radio) interface.

Arıkan's original analysis was asymptotic and established that capacity can be achieved in the limit of large block lengths but did not quantify the speed of convergence to capacity. Effective finite-length convergence bounds were provided several years later in References [15–17], establishing that the polar coding approach leads to a family of codes of rate $C-\varepsilon$ for transmission over a channel of (Shannon) capacity C, where the block length of the code and the decoding time grow only polynomially in $1/\varepsilon$. In contrast, for all previous constructions of codes, the decoding algorithms required time exponential in $1/\varepsilon$. Getting a polynomial running time in $1/\varepsilon$ was one of the central theoretical challenges in the field of algorithmic coding theory, and polar codes were the first to overcome this challenge. Follow-up works have also investigated concrete bounds on the scaling exponent μ , i.e., the finite exponent μ for which the block length of the code can be bounded by $(1/\varepsilon)^{\mu}$ [12, 23], culminating in recent works that achieved $\mu \to 2$, which is the optimal value, first for the erasure channel [10, 25] and later for all channels [13, 30] using variants of polar codes.

The analyses of polar codes turn into questions about *polarizations* of certain *martingales* (which we refer to as Arıkan martingales in this work). The vast class of polar codes alluded to in the

previous paragraph all build on polarizing martingales, and the results of References [15–17] show that for one of the families of polar codes, the underlying martingale polarizes "extremely fast"—a notion we refer to as *strong polarization* and will define shortly.

The primary goal of this work is to understand the process of polarization of martingales and in particular to understand when a martingale polarizes strongly. In attempting to study this question, we come up with a local notion of polarization and show that this local notion is sufficient to imply strong polarization. Applying this improved understanding to the martingales arising in the study of polar codes, we show that a simple necessary condition for weak polarization of such martingales is actually sufficient for strong polarization. This allows us to extend the previous results on strong polarization, which only applied to a specific class of codes, to a broad class of codes and show essentially that all polarizing codes lead to polynomial convergence to capacity. We further show that this can be achieved while maintaining the same exponentially falling error probability achieved in the original asymptotic analyses that did not give any quantitative bounds on the convergence to capacity. Below we formally describe the notion of polarization of martingales and our results concerning them, along with their implications for quantitatively strong convergence to capacity of polar codes when applied to the associated Arıkan martingales. Figure 1 gives a detailed roadmap of this article with different columns indicating different categories of results and each column describing a hierarchy of results.

1.1 Polarization of [0,1]-martingales

Our interest is mainly in the (rate of) polarization of a specific family of martingales that we call the Arıkan martingales. We will define these objects later but first describe the notion of polarization for general [0, 1]-bounded martingales. The middle left (green) column in Figure 1 shows the various notions of polarization that we define in this section.

Recall that a sequence of random variables X_0, \ldots, X_t, \ldots is said to be a *martingale* if for every t and a_0, \ldots, a_t it is the case that $\mathbb{E}[X_{t+1}|X_0 = a_0, \ldots, X_t = a_t] = a_t$. We say that that a martingale is [0, 1]-bounded (or simply a [0, 1]-martingale) if $X_t \in [0, 1]$ for all $t \ge 0$.

Definition 1.1 (Weak Polarization). A [0, 1]-martingale sequence $X_0, X_1, \ldots, X_t, \ldots$ is defined to be weakly polarizing if $\lim_{t\to\infty} \{X_t\}$ exists with probability 1, and this limit is either 0 or 1.

Note that the limit of the martingale sequence $X_0, X_1, \dots, X_t, \dots$ is a Bernoulli random variable with expectation X_0 .¹

Thus, a polarizing martingale does not converge to a single value with probability 1 but rather converges to one of its extreme values. For the applications to constructions of polar codes, we need more explicit bounds on the rates of convergence leading to the notions of (regular) polarization and strong polarization defined below in Definitions 1.3 and 1.4, respectively.

Definition 1.2 $((\tau_{\ell}, \tau_h, \varepsilon)$ -Polarization). For functions $\tau_{\ell}, \tau_h, \varepsilon : \mathbb{Z}^+ \to \mathbb{R}^{\geq 0}$, a [0, 1]-martingale sequence $X_0, X_1, \ldots X_t, \ldots$ is defined to be $(\tau_{\ell}, \tau_h, \varepsilon)$ -polarizing if for all t we have

$$\Pr(X_t \in (\tau_\ell(t), 1 - \tau_h(t))) < \varepsilon(t).$$

Definition 1.3 (Regular Polarization). A [0, 1]-martingale sequence $X_0, X_1, \ldots, X_t, \ldots$ is defined to be regular polarizing if for all constant $\gamma > 0$ there exist $\varepsilon(t) = o(1)$, such that the martingale $\{X_t\}_{t\geq 0}$ is $(\gamma^t, \gamma^t, \varepsilon(t))$ -polarizing.

¹The claim on expectation follows, since by definition, $\mathbb{E}[X_{t+1}] = \mathbb{E}[X_t]$.

11:4 J. Błasiok et al.

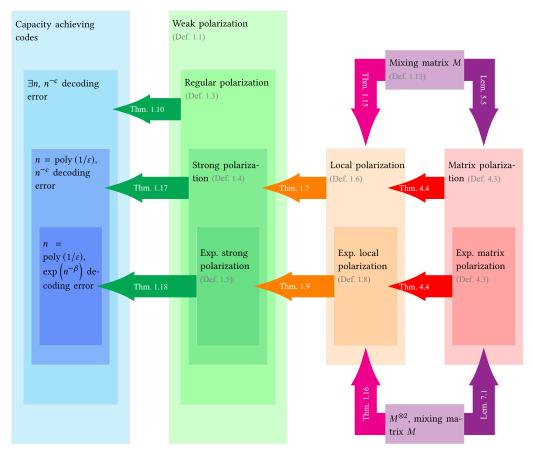


Fig. 1. Overview of our results (excluding those in Section 1.6). The blue boxes (on the extreme left) represent the various coding results (n is the code block length, c and $\beta < 1$ are absolute constants). The green boxes (middle left) are the various notations of polarizations that we study in the article. The orange boxes (middle right) are the two notions of local polarization and the red boxes (extreme right) are the two notions of matrix polarizations we use. Purple boxes (top and bottom on right) show the notions of mixing matrices that we use. All the arrows denote the various results we prove (except for Theorem 1.10, which is implicit in Arıkan [2]) in this article.

We refer to the above as being "sub-exponentially" close to the limit (since it holds for every $\gamma > 0$). While weak polarization by itself is an interesting phenomenon, regular polarization (of Arıkan martingales) leads to capacity-achieving codes (though without explicit bounds on the length of the code as a function of the gap to capacity) and thus regular polarization is well explored in the literature and tight necessary and sufficient conditions are known for regular polarization of Arıkan martingales [3, 20].

To get codes of block length polynomially small in the gap to capacity, an even stronger notion of polarization is needed, where we require that the sub-exponential closeness to the limit happens with *all but exponentially small probability*. We define this formally next.

Definition 1.4 (Strong Polarization). A [0, 1]-martingale sequence $X_0, X_1, \ldots, X_t, \ldots$ is defined to be strongly polarizing if for all $\gamma > 0$ there exist $0 < \eta < 1$ and $\beta < \infty$ such that the martingale $\{X_t\}_{t\geq 0}$ is $(\gamma^t, \gamma^t, \beta \cdot \eta^t)$ -polarizing.

Finally, to get codes where the decoding error probability is exponentially small in the block length, the codes need to polarize even more strongly. We abstract this notion as follows:

Definition 1.5 (Exponentially Strong Polarization). We say that X_t has Λ-exponentially strong polarization if for every $0 < \gamma < 1$ there exist constants $0 < \eta < 1$ and $\beta < \infty$ such that the martingale $\{X_t\}_{t\geq 0}$ is $(2^{-2^{\Lambda t}}, \gamma^t, \beta \eta^t)$ -polarizing.

Note that this definition is asymmetric with respect to the two boundaries and expects tighter polarization when $X_t \to 0$ than when $X_t \to 1$. The reasons for this un-aesthetic choice are the following: (1) For the strong decoding results, the tighter polarization when $X_t \to 0$ suffices. (2) Several of the martingales we consider do not achieve sufficiently tight polarization when $X_t \to 1$ (the Λ they achieve as $X_t \to 1$ is much smaller than what is needed in the decoding results). (3) The analysis of the best polarizations when $X_t \to 0$ is completely different than the analysis when $X_t \to 1$. Due to these reasons we work with this asymmetric definition of exponentially strong polarization.

In contrast to the rich literature on regular polarization, results on strong polarization and exponentially strong polarization are quite rare, reflecting a general lack of understanding of this phenomenon. Indeed, while (roughly) an Arıkan martingale can be associated with every invertible matrix over any finite field \mathbb{F}_q , the only concrete matrix for which exponentially strong polarization was known prior to this work was for $G_2 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$ [15–17].

Part of the reason behind the lack of understanding of strong polarization is that polarization is a "limiting phenomenon" in that one tries to understand $\lim_{t\to\infty} X_t$, whereas most stochastic processes, and the Arıkan martingales in particular, are defined by local evolution, i.e., one that relates X_{t+1} to X_t . The main contribution of this work is to give a local definitions of polarization (Definitions 1.6 and 1.8) and then showing that these definitions imply strong and exponentially strong polarization (Theorems 1.7 and 1.9). Later, we show that Arıkan martingales polarize locally whenever they satisfy a simple condition that is necessary even for weak polarization. And while the Arıkan martingale itself is not locally exponentially polarizing, we show that the "two-step" Arıkan martingale is exponentially locally polarizing under the same simple condition. (The "two step" version of a martingale X_0, X_1, X_2, \ldots , is just the martingales for which previously only regular polarization was known.

1.2 Results I: Local to Strong Global Polarization of Martingales

Before giving the definition of local polarization, we motivate our definition using some simple examples. Consider the martingale Z_0, Z_1, \ldots , where $Z_0 = 1/2$, and $Z_{t+1} = Z_t + Y_{t+1}2^{-(t+2)}$, where Y_1, \ldots, Y_t, \ldots are chosen uniformly and independently from $\{-1, +1\}$. Clearly, this sequence is not polarizing (the limit of Z_t is uniform in [0, 1]). One reason why this happens is that as time progresses, the martingale slows down and stops varying much. We would like to prevent this, but this is also inevitable if a martingale is polarizing and bounded. In particular, a polarizing martingale would be slowed at the boundaries (i.e., when X_t is close to 0 or close to 1) and cannot

 $^{^2}$ An exception is the work by Pfister and Urbanke [25], who showed that for the q-ary erasure channel for large-enough q, the martingale associated with a $q \times q$ Reed–Solomon based matrix proposed in Reference [24] polarizes strongly, and the resulting polar codes achieve scaling exponent tending to 2.

11:6 J. Błasiok et al.

vary much. The first condition in our definition of local polarization insists that this be the only reason a martingale slows down (we refer to this as *variance in the middle*).

Next we consider what happens when a martingale is close to the boundary. For this part, consider a martingale $Z_0 = 1/2$ and $Z_{t+1} = Z_t + \frac{1}{2}Y_{t+1} \min\{Z_t, 1 - Z_t\}$, where again Y_1, \ldots, Y_t, \ldots are chosen uniformly and independently from $\{-1, +1\}$. This martingale does polarize and even shows regular polarization, but it can also be easily seen that the probability that $Z_t < \frac{1}{2} \cdot 2^{-t}$ is zero (whereas we would like probability of being less than say 10^{-t} to go to 1). So this martingale definitely does not show strong polarization. This is so, since even in the best case the martingale is approaching the boundary at a fixed exponential rate and not a sub-exponential one. To overcome this obstacle we require that when the martingale is close to the boundary, with a fixed constant probability it should get much closer in a single step (a notion we refer to as *suction at the ends*).

The middle right (orange) column in Figure 1 shows the notions of local polarization we define in this section (the arrows from the orange column to the middle left (green) columns show the main theorems in this section).

The definition below makes the above requirements precise.

Definition 1.6 (Local Polarization). A [0, 1]-martingale sequence X_0, \ldots, X_j, \ldots , is locally polarizing if the following conditions hold:

- (1) **(Variance in the middle):** For every $\tau > 0$, there is a $\theta = \theta(\tau) > 0$ such that for all j, we have: If $X_j \in (\tau, 1 \tau)$, then $\mathbb{E}[(X_{j+1} X_j)^2 | X_j] \ge \theta$.
- (2) **(Suction at the ends):** There exists an $\alpha > 0$, such that for all $c < \infty$, there exists a $\tau = \tau(c) > 0$, such that:
 - (a) If $X_j \le \tau$, then $\Pr[X_{j+1} \le X_j/c|X_j] \ge \alpha$.
 - (b) Similarly, if $1 X_i \le \tau$, then $\Pr[(1 X_{i+1} \le (1 X_i)/c|X_i] \ge \alpha$.

We refer to condition (a) above as *Suction at the low end* and condition (b) as *Suction at the high end*.

When we wish to be more explicit, we refer to the sequence as $(\alpha, \tau(\cdot), \theta(\cdot))$ locally polarizing.

As such, it is not clear that this definition is of any use. For example, it (1) neither obviously implies strong polarization nor (2) is it obviously satisfiable by any interesting martingale. In this article, we address both these issues. First, we establish general theorems connecting local polarization to strong polarization, as described in Theorems 1.7 and 1.9 below. Then, we leverage this to prove quantitatively strong capacity-approaching properties of polar codes via the strong polarization of Arıkan martingales associated with polar codes (Section 1.3). By our local-to-strong conversion, this in turn follows from the local polarization of Arıkan martingales, which we establish in Theorems 1.15 and 1.16.

Theorem 1.7 (Local vs. Strong Polarization). If a [0,1]-martingale sequence X_0, \ldots, X_t, \ldots , is locally polarizing, then it is also strongly polarizing.

If the suction at the ends shows by the martingale is even stronger, then we can get even stronger polarization. The following definition captures the stronger suction property.

Definition 1.8 (Exponential Local Polarization). We say that X_t has (η, b) -exponential local polarization if it satisfies local polarization (Definition 1.6) and the following additional property:

(1) (Strong suction at the low end): There exists $\tau > 0$ such that if $X_j \le \tau$, then $\Pr[X_{j+1} \le X_j^b | X_j] \ge \eta$.

Note that the interesting range for the parameter b is b > 1 and that is the range on which most of our results will focus.

In the same way that local polarization implies strong global polarization of a martingale, this new stronger local condition implies a stronger global polarization behavior.

Theorem 1.9 (Local to Global Exponential Polarization). Let $\Lambda, b, \eta > 0$ be such that $\Lambda < \eta \log_2 b$. Then, if a [0,1]-bounded martingale X_0, X_1, X_2, \ldots satisfies (η, b) -exponential local polarization, then it also satisfies Λ -exponentially strong polarization.³

Theorems 1.7 and 1.9 are proved in Section 3. In the rest of this section, we turn to showing that the notions of local polarization are not vacuous. Indeed, in later sections we show that the Arıkan martingales polarize locally (under simple necessary conditions). First, we give some background on polar codes.

1.3 The Arikan Martingale and Capacity-achieving Polar Codes

The setting of polar codes considers an arbitrary symmetric memoryless channel and yields codes that aim to achieve the capacity of this channel. These notions are reviewed in Section 2.2.1. Given any q-ary memoryless channel $C_{Y|Z}$ and invertible matrix $M \in \mathbb{F}_q^{k \times k}$, the theory of polar codes implicitly defines a martingale, which we call the Arıkan martingale associated with $(M, C_{Y|Z})$ and studies its polarization. (An additional contribution of this work is that we give an explicit compact definition of this martingale, see Definition 4.1. Since we do not need this definition for the purposes of this section, we defer it to Section 4.) The consequences of regular polarization are described by the following remarkable theorem. (Below we use $M \otimes N$ to denote the tensor product of the matrix M and N. Further, we use $M^{\otimes t}$ to denote the tensor of a matrix M with itself t times.)

Theorem 1.10 (Asymptotic Convergence to Capacity; Implied by Arikan [2]). Let C be a q-ary symmetric memoryless channel and let $M \in \mathbb{F}_q^{k \times k}$ be an invertible matrix. If the Arikan martingale associated with (M,C) polarizes regularly, then given $\varepsilon > 0$ and $c < \infty$ there is a t_0 such that for every $t \ge t_0$ there is a code $C \subseteq \mathbb{F}_q^n$ for $n = k^t$ of dimension at least (Capacity $(C) - \varepsilon$) $\cdot n$ such that C is an affine code generated by the restriction of $(M^{-1})^{\otimes t}$ to a subset of its rows and an affine shift. Moreover, there is a polynomial time decoding algorithm for these codes that has failure probability bounded by n^{-c} .

To obtain codes with faster convergence to capacity, we will need stronger forms of polarization, and a more quantitative version of this theorem, with effective upper bounds on t_0 as a function of the gap ε to capacity. The following version relates parameters of polarization with the quality of the associated code.

Theorem 1.11 (Quantitative Convergence to Capacity [2, 16, 17]). Let C be a q-ary symmetric memoryless channel and let $M \in \mathbb{F}_q^{k \times k}$ be an invertible matrix. If the Arikan martingale associated with (M, C) satisfies $(\tau_\ell, \tau_h, \varepsilon)$ -polarization, then for every t, there is an affine code C, that is generated by the rows of $(M^{-1})^{\otimes t}$ and an affine shift, such that the rate of C is at least

Capacity(
$$C$$
) – $\varepsilon(t)$ – $\tau_h(t)$,

and C can be encoded and decoded⁵ in time $O(n \log n)$ where $n = k^t$ and failure probability of the decoder is at most $O(n \cdot \log q \cdot \tau_{\ell}(t))$.

 $^{^3}$ Note that to get $\eta \log_2 b > \Lambda > 0$ we need $\log_2 b > 0$ and so b > 1 .

⁴We remark that the encoding and decoding are not completely uniform as described above, since the subset of rows and the affine shift that are needed to specify the code are only guaranteed to exist. In the case of additive channels, where the shift can be assumed to be zero, the work of Tal and Vardy [28] (or Reference [16, Sec. V]) removes this non-uniformity by giving a polynomial time algorithm to find the subset.

⁵The running times count the number of floating point operations where real numbers are maintained with $O(\log n)$ bits

⁵The running times count the number of floating point operations where real numbers are maintained with $O(\log n)$ bits of precision.

11:8 J. Błasiok et al.

Remark 1.12. So, in particular, if $\tau_h(t)$, $\varepsilon(t) = O(\rho^t)$, then we get ε close to capacity at block lengths roughly $(1/\varepsilon)^{\log k/\log(1/\rho)}$, which is a polynomial in ε provided $\rho < 1$. Of course, for the code to be useful, we also need $\tau_\ell(t) \ll k^{-t}$. Both conditions are guaranteed by strong polarization. Λ -exponentially strong polarization guarantees decoding failure probability at most $O(n \cdot \log q \cdot \exp(-\Omega(n^{\Lambda/\log_2 k})))$.

This theorem is implicit in the works above, but for completeness we include a proof in Appendix A.2.2 and Appendix A.2.3.

For any binary input symmetric channel, Arıkan and Telatar [3] proved that the martingale associated with the matrix $G_2 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$, polarizes regularly (Arıkan's original paper [2] proved a weaker form of regular polarization with $\tau(t) < 2^{-5t/4}$, which also sufficed for decoding error going to 0). Subsequent work generalized this to other matrices with the work of Korada, Şaşoğlu, and Urbanke [20] giving a precise characterization of matrices M for which the Arıkan martingale polarizes (again over binary input channels). We will refer to such matrices as mixing, formally defined below for all finite fields.

Definition 1.13. (Mixing Matrix). A matrix $M \in \mathbb{F}_q^{k \times k}$ is said to be mixing if it is invertible and none of the permutations of the rows of M yields an upper triangular matrix, i.e., for every permutation $\pi : [k] \to [k]$ there exists $i, j \in [k]$ with $j < \pi(i)$ such that $M_{i,j} \neq 0$.

It is not too hard to show that the Arıkan martingale associated with non-mixing matrices do not polarize (even weakly). In contrast, Reference [20] shows that every mixing matrix over \mathbb{F}_2 polarizes regularly. Mori and Tanaka [24] show that the same result holds for all prime fields and give a slightly more complicated criterion that characterizes (regular) polarization for general fields. (These works show that the decoding failure probability of the resulting polar codes is at most $2^{-n^{\beta}}$ for some positive β determined by the structure of the mixing matrix—this follows from an even stronger decay in the first of the two parameters in the definition of polarization. However, they do *not* show strong polarization, which is what we achieve.)

As alluded to earlier, strong polarization is defined such that it yields codes with polynomial gap to capacity, via Theorem 1.11.

Theorem 1.14 (References [2, 16, 17]). Let C be a q-ary symmetric memoryless channel, and let $M \in \mathbb{F}_q^{k \times k}$ be an invertible matrix. Suppose that the Arikan martingale associated with (M, C) polarizes strongly.

Then, for every c there exists $t_0(x) = O_c(\log x)^7$ such that for every $\varepsilon > 0$ and every $t \ge t_0(1/\varepsilon)$ there is an affine code C that is generated by the rows of $(M^{-1})^{(\otimes t)}$ and an affine shift, with the property that the rate of C is at least Capacity(C) $-\varepsilon$, and C can be encoded and decoded in time $O(n \log n)$ where $n = k^t$ and failure probability of the decoder is at most n^{-c} .

If we assume that the Arikan martingale associated with (M,C) has exponentially strong polarization, then the failure probability of the decoder is at most $\exp(-n^{\beta})$ for some $\beta > 0.8$

The proof of this theorem, as a direct corollary from Theorem 1.11 is included in Appendix A.2.2 for completeness.

As alluded to earlier, the only Arıkan martingales that were known to polarize strongly were those where the underlying matrix was $G_2 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$. Specifically, Guruswami and Xia [16] and

⁶We use 1-indexing in this article.

⁷The notation $O_c(\cdot)$ hides a constant factor that only depends on c.

⁸Throughout this article, we use the notation $\exp(x)$ to denote a function of the form c^x for some constant c > 1. The exact value of c may be different in each usage but will always be bounded away from 1.

Hassani et al. [17] show strong polarization of the Arıkan martingale associated with this matrix over any binary input symmetric channel, and Guruswami and Velingker [15] extended to the case of q-ary input channels for prime q. By using the concept of local polarization, we are able to extend these results to all mixing matrices.

1.4 Results II: Local Polarization of Arıkan Martingales

The results in this subsection appear as the pink arrows (from top and bottom box on the right to the middle right (orange) boxes) in Figure 1.

In our second main result, we show that every mixing matrix gives rise to an Arıkan martingale that is locally polarizing:

THEOREM 1.15 (LOCAL POLARIZATION OF ARIKAN MARTINGALES). For every prime q, for every mixing matrix $M \in \mathbb{F}_q^{k \times k}$, and for every symmetric memoryless channel $C_{Y|Z}$ over \mathbb{F}_q , the associated Arikan martingale is locally polarizing.

Theorem 1.15 is proved in Section 5.5.

We also show that the "two-step martingale," or equivalently the martingale associated with $M^{\otimes 2}$ for mixing matrices M, is exponentially locally polarizing.

Theorem 1.16 (Exponential Local Polarization of Arikan Martingales). For every prime $q, \varepsilon > 0$, every mixing matrix $M \in \mathbb{F}_q^{k \times k}$, and for every symmetric memoryless channel $C_{Y|Z}$ over \mathbb{F}_q , the Arikan martingale sequence associated with $M^{\otimes 2}$ and $C_{Y|Z}$ is $(\frac{1}{k^2}, 2 - \varepsilon)$ -exponentially locally polarizing.

Theorem 1.16 is proved in Section 7.

1.5 Implications for Polar Codes with Polynomial Convergence to Capacity

Results in this section are the two bottom green arrows (from the middle left (green) boxes to the left most (blue) boxes) in Figure 1.

As a consequence of Theorems 1.7, 1.14, and 1.15, we have the following theorem.

Theorem 1.17 (Polynomially Fast Convergence to Capacity and Inverse Polynomial Error Probability). For every prime q, every mixing matrix $M \in \mathbb{F}_q^{k \times k}$, every symmetric memoryless channel C over \mathbb{F}_q , and every $c < \infty$, there is a polynomial p such that for every $\varepsilon > 0$, and every $n = k^t > p(1/\varepsilon)$, there is an affine code C that is generated by the rows of $(M^{-1})^{(\otimes t)}$ and an affine shift, with the property that the rate of C is at least Capacity $(C) - \varepsilon$, and C can be encoded and decoded in time $O(n \log n)$ and failure probability of the decoder is at most n^{-c} .

Again, as a consequence of Theorems 1.9, 1.11, and 1.16, we have the following theorem that achieves decoding failure probability that is $\exp(-n^{\beta})$ for some $\beta > 0$. We refer to such a function as *root-exponentially small*, and when $\beta \to 1$, we call it *near-exponentially small*.

Theorem 1.18 (Polynomial Convergence to Capacity & Root-exponentially Small Error Probability). For every prime q, every mixing matrix $M \in \mathbb{F}_q^{k \times k}$, and every symmetric memoryless channel C over \mathbb{F}_q , there is a polynomial p and $\beta > 0$ such that for every $\varepsilon > 0$ and every $n = k^t \ge p(1/\varepsilon)$, there is an affine code C that is generated by the rows of $(M^{-1})^{(\otimes t)}$ and an affine shift, with the property that the rate of C is at least Capacity $(C) - \varepsilon$, and C can be encoded and decoded in time $O(n \log n)$ and failure probability at most $\exp(-n^{\beta})$.

1.6 Additional Results Optimizing Decoding Error Probability

The above theorems shows that all polar codes associated with every mixing matrix achieves the Shannon capacity of a symmetric memoryless channel efficiently, thus vastly expanding on the

11:10 J. Błasiok et al.

class of polar codes known to satisfy this condition. By choosing the mixing matrix carefully, we can even achieve decoding error probability close to $2^{-\Omega(n)}$; specifically, we can get near-exponentially small decoding error probability, i.e., falling as $\exp(-n^{\beta})$ for any desired $\beta < 1$.

Theorem 1.19 (Near-exponentially Small Error Probability and Polynomial Convergence to Capacity). For every prime q, every symmetric memoryless channel C over \mathbb{F}_q , and every $\beta < 1$, there exists k, a mixing matrix $M \in \mathbb{F}_q^{k \times k}$, and a polynomial p such that for every $\varepsilon > 0$ and every $n = k^t \ge p(1/\varepsilon)$, there is an affine code C that is generated by the rows of $(M^{-1})^{(\otimes t)}$ and an affine shift, with the property that the rate of C is at least Capacity $(C) - \varepsilon$, and C can be encoded and decoded in time $O(n \log n)$ and failure probability at most $\exp(-n^{\beta})$.

Theorem 1.19 is proved in Section 8.2.

Finally, for a broad class of channels, we show that we achieve nearly the best possible error exponent for any given mixing matrix M, while achieving polynomial gap to capacity, using the proofs of this article.

Theorem 1.20 (Polynomial Convergence to Capacity at no Price in Decoding Error Probability). Suppose $M \in \mathbb{F}_q^{k \times k}$ and $\beta > 0$ satisfy the condition that for every q-ary symmetric channel C and for every $\varepsilon > 0$, for sufficiently large $n = k^s$, there is an affine code C of length n generated by the rows of $(M^{-1})^{(\otimes s)}$ of rate at least Capacity $(C) - \varepsilon$ such that C can be decoded with failure probability at most $\exp(-n^{\beta})$.

Then, for every $\beta' < \beta$ and every symmetric channel C' with inputs from \mathbb{F}_q , there is a polynomial p such that for every $\varepsilon > 0$ and every $n = k^t \ge p(1/\varepsilon)$ there is an affine code C that is generated by the rows of $(M^{-1})^{(\otimes t)}$ and an affine shift, with the property that the rate of C is at least Capacity $(C') - \varepsilon$, and C can be encoded and decoded in time $O(n \log n)$ and failure probability at most $\exp(-n^{\beta'})$.

Theorem 1.20 is proved in Section 8.3. It is worth emphasizing two desirable aspects about Theorem 1.20:

- (1) We only need to assume that polar codes based on M achieve capacity for the q-ary symmetric channel but get a conclusion for every symmetric channel (with \mathbb{F}_q inputs).
- (2) Further, we assume nothing about the speed of convergence to capacity for the q-ary symmetric channel and conclude polynomial convergence to capacity (positive scaling exponent) for arbitrary symmetric channels. We do assume root-exponential decoding error probability for the q-ary symmetric channel, but this has been established for all mixing matrices in the limit of $n \to \infty$ [20, 24]. Moreover, in this limit [20] gives a characterization of the best possible exponent β for any given matrix M. Theorem 1.20 asserts that essentially the same characterization applies with polynomial convergence to capacity.

1.7 Comparison with Previous Analyses of (Strong) Polarization

While most of the ingredients going into our eventual analysis of strong polarization are familiar in the literature on polar codes, our proofs end up being much simpler and modular. We describe some of the key steps in our proofs and contrast them with those in previous works.

Definition of Local Polarization. While we are not aware of a definition similar to local polarization being explicit in the literature before, such notions have been considered implicitly before. For instance, for the variation in the middle (where we require that $\mathbb{E}[(X_{t+1} - X_t)^2] \ge \theta$ if $X_t \in (\tau, 1-\tau)$) some of the previous analyses (e.g., in References [15, 16]) required θ be quadratic in

⁹A *q*-ary symmetric channel is one where the symbol is unaltered with probability $1 - \theta$, and flipped to a uniform value with probability θ , for a channel parameter $\theta \in [0, 1]$.

au. In contrast, our requirement on the variation is very weak and qualitative, allowing any function $\theta(\tau) > 0$. Similarly, our requirement in the *suction at the ends* case is relative mild and qualitative. In previous analyses the requirements were of the form "if $X_t \leq \tau$, then $X_{t+1} \leq X_t^2$ with positive probability." This high demand on the suction case prevented the analyses from relying only on the local behavior of the martingale X_0, \ldots, X_t, \ldots and instead had to look at other parameters associated with it that essentially depend on the entire sequence. (For the reader familiar with previous analyses, this is where the Bhattacharyya parameters enter the picture.) Our approach, in contrast, only requires arbitrarily large constant factor drop and thereby works entirely with the local properties of X_t .

Local Polarization Implies Strong Polarization. Our proof that local polarization implies strong polarization is short (about 3 pages) and comes in two parts. The first part uses a simple variance argument to shows that X_t is exponentially close (in t) to the limit except with probability exponentially small in t. The second part then amplifies X_t 's proximity to $\{0,1\}$ to subexponentially small values using the suction at the end guarantee of each local step, coupled with Doob's martingale inequality and standard concentration inequalities. Such a two-part breakdown of the analysis is not new; however, our technical implementation is more abstract, more general, and more compact all at the same time.

Local Polarization of Arıkan Martingales. We will elaborate further on the approach for this after defining the Arıkan martingales, but we can say a little bit already now: First, we essentially reduce the analysis of the polarization of Arıkan martingale associated with an arbitrary mixing matrix M to the analysis when $M = G_2$. This reduction loses in the parameters $(\alpha, \tau(\cdot), \theta(\cdot))$ specifying the level of local polarization, but since our strong polarization theorem works for any function, such loss in performance does not hurt the eventual result. Finally, local polarization for the case where the matrix is G_2 is of course standard, but even here our proofs (which we include for completeness) are simpler, since they follow from known entropic inequalities on sums of two independent random variables. We stress that even quantitatively weak forms of these inequalities meet our requirements of local polarization, and we do not need strong forms of such inequalities (like Mrs. Gerber's lemma for the binary case [7, 16] and an ad hoc one for the prime case [15]).

General vs. Prime Fields. One weaknesses in our analysis that, in contrast to the result of Mori and Tanaka [24], who characterize the set of matrices that lead to regular polarization over general fields, we only get a characterization (for strong polarization) over prime fields. We feel that this limitation is not inherent to our approach. The only (but crucial) place where the prime field plays a role is in the "variance in the middle" lemma (Lemma 5.3) for Arıkan's basic 2×2 kernel G_2 , which in fact does not polarize regularly over general fields due to the existence of subfields. There might be a way around this by reduction to a different 2×2 kernel that actually polarizes regularly.

Concrete Polynomial Upper Bounds on Block Length. A second weakness in our analysis is that, while we develop a general framework to prove strong polarization and polynomial convergence to capacity, the constants are not optimized and will lead to poor upper bounds on the exponent μ of the polynomial in the block length as a function of the gap to capacity. This quantity is called the scaling exponent, and our main goal in this work is to prove that for every mixing matrix M has a finite scaling exponent $\mu = \mu(M)$.

For the case of $M=G_2$ and binary alphabet (the original Arıkan setting), an upper bound of $\mu \leq 6$ was shown in Reference [17] and improved to 5.702 in Reference [12] and to 4.714 in Reference [23]. For the case of the **binary erasure channel (BEC)**, Reference [23] showed an upper bound of $\mu \leq 3.639$, which is close to the heuristic value of ≈ 3.627 reported in Reference [21]. This latter value is also argued as a *lower bound* on μ for the binary-erasure channel in Reference

11:12 J. Błasiok et al.

[17] (for the proof technique of bounding decoding error probability by the sum of Bhattacharyya parameters of the channels seen by the successive cancellation decoder). For kernels besides G_2 , we were unaware of any concrete (or even finite) upper bounds on μ besides our work (except for large random kernels discussed next).

Subsequent Work. Quantitative versions of Shannon's noisy coding theorem theorem show that one can achieve a scaling exponent of 2 for any discrete memoryless channel, and converse theorems show that this is optimal [27, 32]. For erasure channels over large alphabets, it was shown in Reference [25] that random $\ell \times \ell$ kernels for larger ℓ achieve a scaling exponent approaching 2. Such a result was then shown for the BEC in Reference [10].

While these results hinted at the potential of polar codes to achieve near-optimal scaling exponents, they only applied to erasure channels. Analyzing polar codes for more general channels, including the basic **binary symmetric channel (BSC)**, is significantly more complex. Variants of polar codes were shown to achieve a scaling exponent approaching 2 for all binary-input symmetric channels in Reference [13], together with polynomial time constructions and quasi-linear encoding/decoding complexity. A similar result was shown for all discrete memoryless channels over any finite alphabet in Reference [30], albeit the efficient construction of such codes remains to be worked out (but once constructed the codes admit efficient encoding/decoding). These results also use large random kernels. For concrete kernels, this work remains the only general approach to show strong polarization and finite scaling exponent.

1.8 Organization of the Rest of This Article

We first introduce some of the notation and probabilistic preliminaries used to define and analyze the Arıkan martingale in Section 2. We then prove Theorem 1.7 showing that local polarization implies strong polarization in Section 3. This is followed by the formal definition of the Arıkan martingale in Section 4. Section 5.3 then asserts conditions on the entropy of the sum of two independent variables and uses these to prove Theorem 1.15 asserting the local polarization of the Arıkan martingale. Section 6 proves these entropic conditions. Section 7 proves the exponential local polarization of the two-step Arıkan martingale (Theorem 1.16). In Section 8, we prove Theorems 1.19 and 1.20, which strengthen the error analysis for codes to nearly optimal. Finally, in Appendix A we show for completeness how the Arıkan martingale (and its convergence) can be used to construct capacity achieving codes.

2 PRELIMINARIES AND NOTATION

In this section, we introduce the notation needed to define the Arıkan martingale (which will be introduced in Section 4). We also include information-theoretic and probabilistic inequalities that will be necessary for the subsequent analysis.

2.1 Notation

The Arikan martingale is based on a recursive construction of a vector valued random variable. To cleanly describe this construction, it is useful to specify our notational conventions for vectors, tensors, and how to view the tensor products of matrices. These notations will be used extensively in the following sections.

¹⁰For erasure channels, all intermediate channels seen by the decoder of the recursive polar code construction are also erasure channels, with varying erasure probabilities. Even for the BSC, however, the intermediate channels become incredibly complex with huge alphabet sizes. So one must effectively argue about and find a construction that is able to handle a plethora of channels that do not admit analytically simple descriptions.

2.1.1 General Notation. For a prime power q, we use \mathbb{F}_q to denote the finite field with q elements and use \mathbb{F}_q^* to denote the non-zero elements in \mathbb{F}_q .

We will use $O(\cdot)$ for "Big-Oh" notation.

2.1.2 Probability Notation. Throughout this work, all random variables involved will be discrete. For a probability distribution D and random variable X, we write $X \sim D$ to mean that X is distributed according to D and independent of all other variables. Similarly, for a set S, we write $X \sim S$ to mean that X is independent and uniform over S. For a set S, let $\Delta(S)$ denote the set of probability distributions over S.

We occasionally abuse notation by treating distributions as random variables. That is, for $D \in \Delta(\mathbb{F}_q^k)$ and a matrix $M \in \mathbb{F}_q^{k \times k}$, we write DM to denote the distribution of the random variable $\{XM\}_{X \sim D}$. For a distribution D and an event E, we write D|E to denote the conditional distribution of D conditioned on E.

2.1.3 Tensor Notation. Here we introduce useful notation for dealing with scalars, vectors, tensors, and tensor-products. All scalars will be non-bold, for example, $X \in \mathbb{F}_q$. All our vectors will be row vectors (except when explicitly noted) and will be in bold. Any tensors of order ≥ 1 (including vectors) will be in bold, for example: $Y \in \mathbb{F}_q^k$. One exception to this is the matrix M used in the polarization transforms, which we do not show in bold.

Subscripts are used to index tensors, with indices starting from 1. For example, for Y as above, $Y_i \in \mathbb{F}_q$. Matrices and higher-order tensors are indexed with multiple subscripts: For $Z \in (\mathbb{F}_q^k)^{\otimes 3}$, we may write $Z_{1,2,1} \in \mathbb{F}_q$. We often index tensors by tuples (*multiindices*), which will be in bold: For $i = (1,2,1) \in [k]^3$, we write $Z_i = Z_{1,2,1}$. Let < be the lexicographic order on these indexing tuples.

When an index into a tensor is the concatenation of multiple tuples, we emphasize this by using brackets in the subscript. For example, for tensor Z as above, and i = (1, 2) and j = 1, we may write $Z_{[i,j]} = Z_{1,2,1}$.

For a given tensor Z, we can consider fixing some subset of its indices, yielding a *slice* of Z (a tensor of lower order). We denote this with brackets, using \cdot to denote unspecified indices. For example for tensor $Z \in (\mathbb{F}_q^k)^{\otimes 3}$ as above, we have $Z_{[1,2,\cdot]} \in \mathbb{F}_q^k$ and $Z_{[\cdot,1]} \in (\mathbb{F}_q^k)^{\otimes 2}$.

We somewhat abuse the indexing notation, using $Z_{< i}$ to mean the set of variables $\{Z_j : j < i\}$. Similarly, $Z_{[i, < j]} := \{Z_{[i,k]} : k < j\}$.

We occasionally unwrap tensors into vectors, using the correspondence between $(\mathbb{F}_q^k)^{\otimes t}$ and $\mathbb{F}_q^{k^t}$. Here, we unwrap according to the lexicographic order < on tuples.

Finally, for matrices specifically, $M_{i,j}$ specifies the entry in the ith row and jth column of matrix M. Throughout, all vectors will be row-vectors by default.

2.1.4 Tensor Product Recursion. The construction of polar codes and analysis of the Arıkan martingale rely crucially on the recursive structure of the tensor product. Here we review the definition of the tensor product and state its recursive structure.

For a linear transform $M: \mathbb{F}_q^k \to \mathbb{F}_q^k$, let $M^{\otimes t}: (\mathbb{F}_q^k)^{\otimes t} \to (\mathbb{F}_q^k)^{\otimes t}$ denote the *t*-fold tensor power of M. Explicitly (fixing basis for all the spaces involved), this operator acts on tensors $X \in (\mathbb{F}_q^k)^{\otimes t}$ as

$$[M^{\otimes t}(X)]_j = \sum_{i \in [k]^t} X_i M_{i_1, j_1} M_{i_2, j_2} \cdots M_{i_t, j_t}.$$

The tensor product has the following recursive structure: $M^{\otimes t} = (M^{\otimes t-1}) \otimes M$, which corresponds explicitly to

$$[M^{\otimes t}(X)]_{[a,j_t]} = \sum_{i_t \in [k]} M_{i_t,j_t} [M^{\otimes t-1}(X_{[\cdot,i_t]})]_a.$$
(1)

11:14 J. Błasiok et al.

In the above, if we define tensor

$$Y^{(i_t)} := M^{\otimes t-1}(X_{\lceil \cdot, i_t \rceil}),$$

then this becomes

$$[M^{\otimes t}(X)]_{[a,\cdot]} = M((Y_a^{(1)}, Y_a^{(2)}, \dots, Y_a^{(k)})),$$
(2)

where the vector $(Y_a^{(1)}, Y_a^{(2)}, \dots, Y_a^{(k)}) \in \mathbb{F}_q^k$. Finally, we use that $(M^{\otimes t})^{-1} = (M^{-1})^{\otimes t}$.

Information Theory Preliminaries

For the sake of completeness, we include the information-theoretic concepts and tools we use in this article.

For a discrete random variable X, let H(X) denote its binary entropy:

$$H(X) := \sum_{a \in Support(X)} p_X(a) \log \left(\frac{1}{p_X(a)}\right),$$

where $p_X(a) := \Pr(X = a)$ is the probability mass function of X. Throughout, $\log(\cdot)$ by default denotes $\log_2(\cdot)$.

For $p \in [0,1]$, we overload this notation, letting H(p) denote the entropy H(X) for $X \sim$ Bernoulli(p).

For arbitrary random variables X, Y, let H(X|Y) denote the conditional entropy:

$$H(X|Y) = \underset{V}{\mathbb{E}}[H(X|Y=y)].$$

For a *q*-ary random variable $X \in \mathbb{F}_q$, let $\overline{H}(X) \in [0,1]$ denote its (normalized) *q*-ary entropy:

$$\overline{H}(X) := \frac{H(X)}{\log(q)}.$$
(3)

Finally, the *mutual information* between jointly distributed random variables *X*, *Y* is

$$I(X; Y) := H(X) - H(X|Y) = H(Y) - H(Y|X).$$

We will use the following standard properties of entropy (see, for instance, Reference [6]):

(1) (Adding independent variables increases entropy): For any random variables X, Y, Zsuch that X, Y are conditionally independent given Z, we have

$$H(X+Y|Z) \ge H(X|Z). \tag{4}$$

(2) (Transforming Conditioning): For any random variables X, Y, any function f, and any bijection σ , we have

$$H(X|Y) = H(X + f(Y)|Y) = H(X + f(Y)|\sigma(Y)).$$
 (5)

- (3) **(Chain rule):** For arbitrary random variables X, Y: H(X, Y) = H(X) + H(Y|X).
- (4) (Conditioning does not increase entropy): For X, Y, Z arbitrary random variables, $H(X|Y,Z) \le H(X|Y).$
- (5) (Monotonicity): For $p \in [0, 1/2)$, the binary entropy H(p) is non-decreasing with p. And for $p \in (1/2, 1]$, the binary entropy H(p) is non-increasing with p.
- (6) (Deterministic postprocessing does not increase entropy): For arbitrary random variables X, Y and function f we have $H(X|Y) \ge H(f(X)|Y)$.

- (7) **(Conditioning on independent variables):** For random variables X, Y, Z where Z is independent from (X, Y), we have H(X|Y) = H(X|Y, Z).
- 2.2.1 Channels. Given a finite field \mathbb{F}_q , and output alphabet \mathcal{Y} , a q-ary channel $C_{Y|Z}$ is a probabilistic function from \mathbb{F}_q to \mathcal{Y} . Equivalently, it is given by q probability distributions $\{C_{Y|\alpha}\}_{\alpha\in\mathbb{F}_q}$ supported on \mathcal{Y} . We use notation C(Z) to denote the channel operating on inputs Z. A memoryless channel maps \mathbb{F}_q^n to \mathcal{Y}^n by acting independently (and identically) on each coordinate. A symmetric channel is a memoryless channel where for every $\alpha, \beta \in \mathbb{F}_q$ there is a bijection $\sigma: \mathcal{Y} \to \mathcal{Y}$ such that for every $y \in \mathcal{Y}$ it is the case that $C_{Y=y|\alpha} = C_{Y=\sigma(y)|\beta}$, and moreover for any pair $y_1, y_2 \in \mathcal{Y}$, we have $\sum_{x \in \mathbb{F}_q} C_{Y=y_1|x} = \sum_{x \in \mathbb{F}_q} C_{Y=y_2|x}$ (see, for example, Reference [6, Section 7.2]). As shown by Shannon every memoryless channel has a finite capacity, denoted Capacity $(C_{Y|Z})$. For symmetric channels, this is the mutual information I(Y;Z) between the input Z and output Y where Z is drawn uniformly from \mathbb{F}_q and Y is drawn from $C_{Y|Z}$ given Z.

2.3 Basic Probabilistic Inequalities

In this section, we collect a few useful probabilistic and information-theoretic inequalities, all of which are standard. The proofs are included for convenience.

We first show that a random variable with small-enough entropy will usually take its most-likely value and thus maximum likelihood recovery is successful with high probability. In fact, we show that even if the likelihoods are known only very approximately maximum likelihood decoding will still be quite successful.

Lemma 2.1. Let X be an arbitrary discrete random variable with range X. Then there exist $\hat{x} \in X$ such that

$$\Pr(X \neq \hat{x}) \leq H(X).$$

In particular, one can take $\hat{x} = \operatorname{argmax}_{\alpha} \{ \Pr(X = \alpha) \}.$

Furthermore, given \tilde{p}_{α} 's satisfying $|\tilde{p}_{\alpha} - \Pr(X = \alpha)| \le 1/4$ for every $\alpha \in X$, if we let $\tilde{x} = \operatorname{argmax}_{\alpha}\{\tilde{p}_{\alpha}\}$ then we have $\Pr(X \neq \tilde{x}) \le 3H(X)$.

PROOF. Let $\alpha := H(X)$ and let $p_i := \Pr_X (X = i)$. Let $\hat{x} = \operatorname{argmax}_i \{p_i\}$ be the value maximizing this probability. Let $p_{\hat{x}} = 1 - \gamma$. We wish to show that $\gamma \le \alpha$. If $\gamma \le 1/2$, then we have

$$\alpha = H(X) = \sum_{i} p_{i} \log \frac{1}{p_{i}}$$

$$\geq \sum_{i \neq \hat{x}} p_{i} \log \frac{1}{p_{i}} \qquad \text{(Since all summands are non-negative)}$$

$$\geq \sum_{i \neq \hat{x}} p_{i} \log \frac{1}{\sum_{j \neq \hat{x}} p_{j}} \qquad \text{(Since } p_{i} \leq \sum_{j \neq \hat{x}} p_{j}.)$$

$$= \left(\sum_{i \neq \hat{x}} p_{i}\right) \cdot \log \left(\frac{1}{\sum_{j \neq \hat{x}} p_{j}}\right)$$

$$= \gamma \cdot \log 1/\gamma$$

$$\geq \gamma \qquad \text{(Since } \gamma \leq 1/2 \text{ and so } \log 1/\gamma \geq 1)$$

11:16 J. Błasiok et al.

as desired. Now, if $\gamma > 1/2$, then we have a much simpler case, since now we have

$$\alpha = H(X) = \sum_{i} p_{i} \log \frac{1}{p_{i}}$$

$$\geq \sum_{i} p_{i} \log \frac{1}{p_{\hat{x}}} \qquad (Since \ p_{i} \leq p_{x})$$

$$= \log \frac{1}{p_{\hat{x}}} \qquad (Since \ \sum_{i} p_{i} = 1)$$

$$= \log \frac{1}{1 - \gamma}$$

$$\geq 1. \qquad (Since \ \gamma \geq 1/2)$$

But γ is always at most 1 so in this case also we have $\alpha \ge 1 \ge \gamma$ as desired.

For the furthermore part of the lemma statement, note that if $\gamma < 1/4$, then, by the condition $|\tilde{p}_{\alpha} - p_{\alpha}| \le 1/4$, we have $\tilde{p}_{\hat{x}} > 1/2$ while $\tilde{p}_{x'} < 1/2$ for every $x' \ne \hat{x}$. Thus in this case we have $\tilde{x} = \hat{x}$ and so by the first part above we have $\Pr(X \ne \hat{x}) = \Pr(X \ne \hat{x}) \le H(X)$. Now, if $\gamma > 1/4$ as in the second part above, then we have $H(X) \ge \log \frac{1}{1-\gamma} \ge .415 \ge 1/3$, and so we get $\Pr(X \ne \hat{x}) \le 1 \le 3H(X)$.

For the decoder, we will need a conditional version of Lemma 2.1, saying that if a variable X has low conditional entropy conditioned on Y, then X can be predicted well given the instantiation of variable Y.

Lemma 2.2. Let X, Y be arbitrary discrete random variables with range X, \mathcal{Y} respectively. Then there exists a function $\hat{X}: \mathcal{Y} \to X$ such that

$$\Pr_{X \mid Y} \left(X \neq \hat{X}(Y) \right) \le H(X|Y).$$

In particular, the following estimator satisfies this:

$$\hat{X}(y) := \underset{x}{\operatorname{argmax}} \left\{ \Pr \left(X = x | Y = y \right) \right\}.$$

Furthermore, given $\tilde{p}_{x,y}$'s satisfying $|\tilde{p}_{x,y} - \Pr(X = x|Y = y)| \le 1/4$ for every $x \in \mathcal{X}, y \in \mathcal{Y}$, if we let $\tilde{X}(y) = \operatorname{argmax}_{x} \{\tilde{p}_{x,y}\}$, then we have $\Pr(X \neq \tilde{X}(y)) \le 3H(X|Y)$.

PROOF. For every setting of Y = y, we can bound the error probability of this estimator using Lemma 2.1 applied to the conditional distribution X|Y = y:

$$\Pr_{X,Y} \left(X \neq \hat{X}(Y) \right) = \mathbb{E} \left[\Pr_{X|Y} \left(\hat{X}(Y) \neq X \right) \right] \\
\leq \mathbb{E} \left[H(X|Y = y) \right] \qquad \text{(Lemma 2.1)} \\
= H(X|Y) .$$

The furthermore part of the lemma statement follows similarly by using the furthermore part of Lemma 2.1.

We also use the well-known Fano's inequality, which works as a weak converse to the above lemma, asserting that if a random variable X is predictable given Y, then its conditional entropy is small.

LEMMA 2.3 (FANO'S INEQUALITY). For a pair of random variables $(X, Y) \in X \times \mathcal{Y}$, if there exists a function $\hat{X} : \mathcal{Y} \to X$ such that $\Pr(\hat{X}(Y) \neq X) \leq \delta$ with $\delta < \frac{1}{2}$, then $H(X|Y) \leq 2\delta(\log \delta^{-1} + \log |X|)$.

We will need an inverse to the usual Chebychev inequality. Recall that Chebychev shows that variables with small variance are concentrated close to their expectation:

$$\Pr(|Z - \mathbb{E}[Z]| \ge \lambda) \le \frac{\operatorname{Var}(Z)}{\lambda^2}.$$

The Paley–Zygmund inequality below can be used to invert it (somewhat): For a random variable W with comparable fourth and second central moment, by applying the lemma below to $Z = (W - \mathbb{E}[W])^2$ we can deduce that W has positive probability of deviating noticeably from the mean.

Lemma 2.4 (Paley-Zygmund). If $Z \ge 0$ is a random variable with finite variance, then

$$\Pr(Z > \lambda \mathbb{E}[Z]) \ge (1 - \lambda)^2 \frac{\mathbb{E}[Z]^2}{\mathbb{E}[Z^2]}.$$

Next, we define the notion of a sequence of random variables being adapted to another sequence of variables, which will be useful in our later proofs.

Definition 2.5. We say that a sequence $Y_1, Y_2 \dots$ of random variables is *adapted* to the sequence $X_1, X_2 \dots$ if and only if for every t, Y_t is completely determined given $X_1, \dots X_t$. We will use $\mathbb{E}[Z|X_{[1:t]}]$ as a shorthand $\mathbb{E}[Z|X_1, \dots X_t]$ and $\Pr(E|X_{[1:t]})$ as a shorthand for $\mathbb{E}[\mathbb{1}_E|X_1, \dots X_t]$. If the underlying sequence X is clear from context, then we will skip it and write just $\mathbb{E}[Z|\mathcal{F}_t]$.

LEMMA 2.6. Consider a sequence of non-negative random variables $Y_1, Y_2, \ldots, Y_t, \ldots$ adapted to the sequence X_1, X_2, \ldots If for every t we have $\Pr(Y_{t+1} > \lambda \mid X_{[1:t]}) \le \exp(-\lambda)$, then for every T > 0:

$$\Pr\left(\sum_{i < T} Y_i > CT\right) \le \exp(-\Omega(T))$$

for some universal constant C.

PROOF. First, observe that

$$\mathbb{E}[\exp(Y_{t+1}/2)|\mathcal{F}_t] = \int_0^\infty \Pr(\exp(Y_{t+1}/2) > \lambda | \mathcal{F}_t) \, d\lambda$$

$$\leq 1 + \int_1^\infty \exp(-2\log \lambda) \, d\lambda$$

$$= 1 + \int_1^\infty \lambda^{-2} \, d\lambda$$

$$\leq \exp(C_0)$$
(6)

for some constant C_0 . However, we have decomposition (where we apply Equation (6) in the first inequality):

$$\mathbb{E}\left[\exp\left(\sum_{i\leq T}\frac{Y_i}{2}\right)\right] = \mathbb{E}\left[\mathbb{E}\left[\exp\left(\sum_{i\leq T}\frac{Y_i}{2}\right)|\mathcal{F}_{T-1}\right]\right]$$

$$= \mathbb{E}\left[\exp\left(\sum_{i\leq T-1}Y_i/2\right)\mathbb{E}\left[\exp\left(Y_T/2\right)|\mathcal{F}_{T-1}\right]\right]$$

$$\leq \mathbb{E}\left[\exp\left(\sum_{i\leq T-1}Y_i/2\right)\right] \cdot \exp(C_0)$$

$$\leq \cdots$$

$$\leq \exp(C_0T).$$

11:18 J. Błasiok et al.

In the above, the second equality follows from the fact that the sequence Y_1, Y_2, \ldots is adapted to X_1, X_2, \ldots We can now apply Markov inequality to obtain the desired tail bound:

$$\Pr\left(\sum_{i \le T} Y_i > 4C_0 T\right) = \Pr\left(\exp\left(\frac{1}{2}\sum_{i \le T} Y_i\right) > \exp(2C_0 T)\right)$$

$$\leq \mathbb{E}\left[\exp\left(\frac{1}{2}\sum_{i \le T} Y_i\right)\right] \cdot \exp(-2C_0 T)$$

$$\leq \exp\left(-C_0 T\right).$$

The following bound for a moment generating function of a bounded random variable is standard and is commonly used in the proof of Bernstein inequality.

Lemma 2.7. For any random variable X such that |X| < 1 with probability 1, and every $0 < \lambda < \frac{1}{4}$, we have

$$\log \mathbb{E}[\exp(\lambda X)] \le \lambda \mathbb{E}[X] + C\lambda^2 \mathbb{E}[X^2],$$

where C is some universal constant.

PROOF. Since |X| < 1, we have $\mathbb{E} |X|^k \leq \mathbb{E} X^2$, and therefore

$$\mathbb{E} \exp(\lambda X) = \sum_{k} \frac{\lambda^{k}}{k!} \mathbb{E}[X^{k}]$$

$$\leq 1 + \lambda \mathbb{E}[X] + (\lambda^{2} + O(\lambda^{3})) \mathbb{E}[X^{2}].$$

Moreover, for some constant C, and every $|x| < \frac{1}{2}$, we have $\log(1+x) \le x + Cx^2$, and therefore

$$\log \mathbb{E}[\exp(\lambda X)] \le \lambda \mathbb{E}[X] + C\lambda^{2}(\mathbb{E}[X^{2}] + \mathbb{E}[X]^{2}) + O(\lambda^{3}) \mathbb{E}[X^{2}]$$
$$\le \lambda \mathbb{E}[X] + C'\lambda^{2} \mathbb{E}[X^{2}].$$

Lemma 2.8. Consider a sequence of random variables Y_1, Y_2, \ldots with $Y_i \in \{0, 1\}$, adapted to the sequence X_t . If $\Pr(Y_{t+1} = 1 | X_{[1:t]}) > \mu_{t+1}$ for some deterministic value μ_t , then for $\mu := \sum_{t \leq T} \mu_t$ and any $\varepsilon > 0$ we have

$$\Pr\left(\sum_{t \le T} Y_t < (1 - \varepsilon)\mu\right) \le \exp\left(-\Omega(\varepsilon^2 \mu)\right).$$

PROOF. Consider a random variable $M_{t+1} := \mathbb{E}[Y_{t+1}|X_{[1:t]}]$ (depending on $X_{[1:t]}$), we know that $M_t > \mu_t$ with probability 1, and let us take $Z_t := (1 - \varepsilon)M_t - Y_t$.

Standard calculation involving Markov inequality yields the following bound for any $\lambda > 0$:

$$\Pr\left(\sum_{t \leq T} Y_t < \sum_{t \leq T} (1 - \varepsilon)\mu_t\right) \leq \Pr\left(\sum_{t \leq T} Y_t < \sum_{t \leq T} (1 - \varepsilon)M_t\right)$$

$$= \Pr\left(\sum_{t \leq T} \lambda Z_t > 0\right)$$

$$= \Pr\left(\exp\left(\sum_{t \leq T} \lambda Z_t\right) > 1\right)$$

$$\leq \mathbb{E}\left[\exp\left(\sum_{t \leq T} \lambda Z_t\right)\right]. \tag{7}$$

To bound this latter quantity, we introduce conditioning on $X_{[1:T-1]}$ as follows:

$$\mathbb{E}\left[\exp\left(\sum_{t\leq T}\lambda Z_{t}\right)\right] = \mathbb{E}\left[\mathbb{E}\left[\exp\left(\sum_{t\leq T}\lambda Z_{t}\right)|X_{[1:T-1]}\right]\right]$$

$$= \mathbb{E}\left[\exp\left(\sum_{t\leq T-1}\lambda Z_{t}\right)\mathbb{E}\left[\exp(\lambda Z_{T})|X_{[1:T-1]}\right]\right],$$
(8)

where the second equality follows from the fact that Z_t is adapted to X_t .

By Lemma 2.7 for any $0 < \lambda < \frac{1}{4}$, we have

$$\mathbb{E}\left[\exp(\lambda Z_T)|X_{[1:T-1]}\right] \leq \exp(-\lambda \varepsilon M_T + C_1 \lambda^2 M_T)$$

for some constant C_1 . Now, if we chose $\lambda = \frac{1}{2C_1}\varepsilon$, then we get

$$\mathbb{E}\left[\exp\left(\lambda Z_{T}\right) | X_{[1:T-1]}\right] \leq \mathbb{E}\left[\exp\left(-C\varepsilon^{2} M_{T}\right)\right]$$

$$\leq \exp(-C\varepsilon^{2} \mu_{T}), \tag{9}$$

where $C = \frac{1}{8C_1}$, since $\mu_T \leq M_T$ deterministically.

Together with Equation (8), this yields

$$\mathbb{E}\left[\exp\left(\sum_{t\leq T}\lambda Z_{t}\right)\right] \leq \mathbb{E}\left[\exp\left(\sum_{t\leq T-1}\lambda Z_{t}\right)\right] \exp(-C\varepsilon^{2}\mu_{T})$$

$$\leq \cdots$$

$$\leq \mathbb{E}\left[\exp\left(\sum_{t\leq T}-C\varepsilon^{2}\mu_{t}\right)\right] = \exp(-\Omega(\varepsilon^{2}\mu)). \tag{10}$$

Finally, combining Equations (7) and (10) we have $\Pr(\sum_{t \leq T} Y_t < (1 - \varepsilon)\mu) \leq \exp(-\Omega(\varepsilon^2 \mu))$ as desired.

Finally, we will use the well-known Doob's martingale inequality:

Lemma 2.9 (Doob's Martingale Inequality [9, Theorem 5.4.2]). If a sequence X_0, X_1, \ldots is a martingale, then for every T we have

$$\Pr\left(\sup_{t\leq T}X_t>\lambda\right)\leq \frac{\mathbb{E}[|X_T|]}{\lambda}.$$

COROLLARY 2.10. If X_0, X_1, \ldots is a nonnegative martingale, then for every T we have

$$\Pr\left(\sup_{t\leq T}X_t>\lambda\right)\leq \frac{\mathbb{E}[X_0]}{\lambda}.$$

3 LOCAL TO GLOBAL POLARIZATION

In this section, we prove Theorems 1.7 and 1.9, which assert that every (exponentially) locally polarizing [0, 1]-martingale is also (exponentially) strongly polarizing. The proofs in this section depend on some basic probabilistic concepts and inequalities mentioned in Section 2.3.

The proof of both statements are implemented in two main steps. In the first step, common to both, we show that any locally polarizing martingale, is mildly polarizing, namely that it is $((1 - \frac{\nu}{2})^t, (1 - \frac{\nu}{2})^t, (1 - \frac{\nu}{4})^t)$ -polarizing for *some* constant ν depending only on the parameters α, τ, θ of local polarization. This means that, except with exponentially small probability, $\min\{X_{t/2}, 1 - X_{t/2}\}$ is exponentially small in t, which we can use to ensure that X_s for all $\frac{t}{2} \le s \le t$ stays in the range where the conditions of (*strong*) suction at the ends apply (again, except with exponentially small

11:20 J. Błasiok et al.

failure probability). In the second step, we show that if the martingale stays in the *suction at the ends* regime, then it will polarize strongly; i.e., if we have a [0,1]-martingale, such that in each step it has probability at least α to decrease by a factor of c, then we can deduce that at the end we have $\Pr(X_T > c^{-\alpha T/4}) \le \exp(-\Omega(\alpha T))$. A completely similar argument shows that when the martingale shows strong suction at the low end, we have $\Pr(X_T > \exp(-\Delta^T)) \le \exp(-\Omega(\alpha T))$, for some $\Delta > 1$, thus yielding exponentially strong polarization.

3.1 Mild Polarization

We start by showing that in the first t/2 steps we do get exponentially small polarization, with all but exponentially small failure probability. This is proved using a simple potential function $\min\{\sqrt{X_t}, \sqrt{1-X_t}\}$, which we show shrinks by a constant factor, 1-v for some v>0, in expectation at each step. Previous analyses in References [15, 16] tracked $\sqrt{X_t(1-X_t)}$ (or some tailormade algebraic functions [17, 23]) as potential functions and relied on quantitatively strong forms of variance in the middle to demonstrate that the potential diminishes by a constant factor in each step. While such analyses can lead to sharper bounds on the parameter v, which in turn translate to better *scaling exponents* in the polynomial convergence to capacity, e.g., see Reference [17, Thm. 18] or Reference [23, Thm. 1], these analyses are more complex, and less general.

LEMMA 3.1. If a [0, 1]-martingale sequence $X_0, \ldots X_t, \ldots$, is $(\alpha, \tau(\cdot), \theta(\cdot))$ -locally polarizing, then there exist v > 0, depending only on α, τ, θ , such that

$$\mathbb{E}\left[\min\left(\sqrt{X_t},\sqrt{1-X_t}\right)\right] \le (1-\nu)^t.$$

PROOF. Set $\tau_0 = \tau(4)$, $\theta_0 = \theta(\tau_0)$. We will show that $\mathbb{E}[\min(\sqrt{X_{t+1}}, \sqrt{1-X_{t+1}})|X_t] \le (1-\nu)\min(\sqrt{X_t}, \sqrt{1-X_t})$ for some $\nu > 0$ depending on τ_0 , θ_0 , and α . The statement of the lemma will follow by induction. The base case of t = 0 follows, since $X_0 \in [0, 1]$.

Let us condition on X_t and first consider the case $X_t \in (\tau_0, 1 - \tau_0)$. We know that

$$\mathbb{E}\left[\,\min\left(\sqrt{X_{t+1}},\sqrt{1-X_{t+1}}\right)\right] \leq \min\left(\,\mathbb{E}\left[\sqrt{X_{t+1}}\right],\mathbb{E}\left[\sqrt{1-X_{t+1}}\right]\right),$$

and we will show that $\mathbb{E}[\sqrt{X_{t+1}}] \leq (1-\nu)\sqrt{X_t}$. The proof of $\mathbb{E}[\sqrt{1-X_{t+1}}] \leq (1-\nu)\sqrt{1-X_t}$ is symmetric.

Indeed, let us take $R := \sqrt{\frac{X_{t+1}}{X_t}}$. Because $(X_t)_t$ is a martingale, we have $\mathbb{E}[R^2] = 1$, and by Jensen's inequality, we have that $\mathbb{E}[R] \le \sqrt{\mathbb{E}[R^2]} \le 1$, where all the expectations above are conditioned on X_t . Take δ such that $\mathbb{E}[R] = 1 - \delta$. We will show a lower bound on δ in terms of θ_0 , τ_0 , and α_0 .

We note that

$$Var(R) = \mathbb{E}[R^2] - (\mathbb{E}[R])^2 = 1 - (1 - \delta)^2 = 2\delta - \delta^2 \le 2\delta.$$
 (11)

The high-level idea of the proof is that we can show that local polarization criteria implies that T is relatively far from 1 with noticeable probability, but if $\mathbb{E}[R]$ were close to 1, then by Chebyshev inequality we would be able to deduce that R is far from its mean with much smaller probability. This implies that mean of R has to be bounded away from 1.

More concretely, observe first that by Chebyshev inequality, we have $\Pr(|R - \mathbb{E}[R]| > \lambda) < \frac{\operatorname{Var}(R)}{\lambda^2} \le \frac{2\delta}{\lambda^2}$, where the inequality follows from Equatino (11). Hence, for $C_0 = 4$, we have the following:

$$\Pr\left(|R-1| \ge \delta + C_0 \sqrt{\delta \theta_0^{-1} \tau_0^{-2}}\right) \le \frac{1}{8} \theta_0^2 \tau_0^4. \tag{12}$$

¹¹This is enough, since we pick c to be large enough (given γ) so that $c^{-\alpha T/4} \leq \gamma^T$, and we pick β and η such that $\beta \eta^T \geq \exp(-\Omega(\alpha T))$.

However, because of the Variation in the middle condition of local polarization, we have

$$\operatorname{Var}(R^2) = \frac{\mathbb{E}[X_{t+1}^2]}{X_t^2} - \frac{\mathbb{E}[X_{t+1}]^2}{X_t^2} = \frac{\mathbb{E}[X_{t+1}^2] - X_t^2}{X_t^2} \ge \frac{\theta_0}{X_t^2} \ge \theta_0,$$

where the second equality follows from the fact that $\mathbb{E}[X_{t+1}] = X_t$ and the last inequality follows, since $X_t \leq 1$. Moreover, $R < \frac{1}{\sqrt{\tau_0}}$, because $\sqrt{X_{t+1}} < 1$ and $\sqrt{X_t} > \sqrt{\tau_0}$.

Let us now consider $Z=(R^2-1)^2$. We have $\mathbb{E}[Z]=\operatorname{Var}(R^2)\geq \theta_0$, and, moreover, $\mathbb{E}[Z^2]<\tau_0^{-4}$ (because R is bounded and $\tau_0\leq 1$), hence by Lemma 2.4 (for $C_1=1/2$),

$$\Pr\left((1-R^2)^2 > C_1\theta_0\right) \ge \frac{1}{4}\theta_0^2\tau_0^4.$$

And also $1 - R^2 = -(1 - R)^2 + 2(1 - R) < 2(1 - R)$; hence, if $(1 - R^2)^2 > C_1\theta_0$, then $|1 - R| > \frac{\sqrt{C_1}}{2}\sqrt{\theta_0}$, which implies (for the choice of $C_2 = \sqrt{C_1}/2$):

$$\Pr(|R-1| > C_2 \sqrt{\theta_0}) \ge \frac{1}{4} \theta_0^2 \tau_0^4.$$
 (13)

By comparing Equations (12) and (13), we deduce that $C_2\sqrt{\theta_0} < \delta + C_0\sqrt{\delta}\theta_0^{-1}\tau_0^{-2}$, which in turn implies that $\delta \geq C_4\theta_0^3\tau_0^4$, (for $C_4 = C_2^2/(4C_0^2)$; note that with our choice of parameters, we have $C_0\sqrt{\delta}\theta_0^{-1}\tau_0^{-2} \geq \delta$), and by the definition of δ we have $\mathbb{E}[\sqrt{X_{t+1}}|X_t] \leq (1-\delta)\sqrt{X_t}$]. The same argument applies to show that $\mathbb{E}[\sqrt{1-X_{t+1}}|X_t] \leq (1-C_4\theta_0^3\tau_0^4)\sqrt{1-X_t}$.

Consider now the case when $X_t < \tau_0$. For T, δ as above (and again after conditioning on X_t), we have $\text{Var}(R) < 2\delta$ (note that the argument for this inequality from the previous case also holds here), and hence by Chebyshev inequality (for the choice of $C_5 = 2$),

$$\Pr\left(|R-1| \ge \delta + C_5\sqrt{\frac{\delta}{\alpha}}\right) \le \frac{\alpha}{2}.$$
 (14)

However, because of the *suction at the end* condition of local polarization, we know that with probability α , we have $R \leq \frac{1}{2}$, which means $|R-1| \geq \frac{1}{2}$ and by comparing this with Equation (14), we deduce that $\delta + C_5 \sqrt{\frac{\delta}{\alpha}} \geq \frac{1}{2}$, which in turn implies that $\delta \geq C_6 \alpha$ (for $C_6 = \frac{1}{16C_5^2}$; note that by our parameter choices we have $C_5 \sqrt{\frac{\delta}{\alpha}} \geq \delta$). Therefore, in the case $X_t < \tau_0$, we have $\mathbb{E}[\sqrt{X_{t+1}}|X_t] \leq (1 - C_6 \alpha) \sqrt{X_t} = (1 - C_6 \alpha) \min(\sqrt{X_t}, \sqrt{1 - X_t})$. The case $X_t > 1 - \tau_0$ is symmetric and is omitted. This implies the statement of the lemma with $\nu = \min(C_6 \alpha, C_4 \theta_0^3 \tau_0^2)$.

COROLLARY 3.2. If a [0, 1]-martingale sequence $X_0, \ldots X_t, \ldots$, is $(\alpha, \tau(\cdot), \theta(\cdot))$ -locally polarizing, then there exist v > 0, depending only on α, τ, θ , such that

$$\Pr\left(\min(X_{t/2}, 1 - X_{t/2}) > \lambda \left(1 - \frac{\nu}{2}\right)^t\right) \le \left(1 - \frac{\nu}{4}\right)^t \frac{1}{\sqrt{\lambda}}.$$

PROOF. By applying Markov Inequality to the bound from Lemma 3.1 (with t/2 instead of t), we get

$$\Pr\left(\min(X_{t/2}, 1 - X_{t/2}) > \lambda \left(1 - \frac{\nu}{2}\right)^{t}\right) = \Pr\left(\min\left(\sqrt{X_{t/2}}, \sqrt{1 - X_{t/2}}\right) > \sqrt{\lambda} \left(1 - \frac{\nu}{2}\right)^{t/2}\right)$$

$$\leq (1 - \nu)^{t/2} \left(1 - \frac{\nu}{2}\right)^{-t/2} \frac{1}{\sqrt{\lambda}}$$

$$\leq \left(1 - \frac{\nu}{4}\right)^{t} \frac{1}{\sqrt{\lambda}}.$$

11:22 J. Błasiok et al.

3.2 Strong Polarization

Next, we show that if a [0,1]-martingale indeed stays in the *suction at the ends* range for all steps $s \ge \frac{t}{2}$, i.e., in each step it has constant probability α of dropping by some large constant factor C, then at the end we may expect it to be $(C^{-\alpha t/8}, C^{-\alpha t/8}, \exp(-\Omega(\alpha t)))$ -polarizing.

LEMMA 3.3. There exists $c < \infty$, such that for all K, α with $K\alpha \ge c$ the following holds. Let X_t be a martingale satisfying $\Pr(X_{t+1} < e^{-K}X_t|X_t) \ge \alpha$, where $X_0 \in (0,1)$. Then $\Pr(X_T > \exp(-\alpha KT/4)) \le \exp(-\Omega(\alpha T))$.

PROOF. Consider $Y_{t+1} := \log \frac{X_{t+1}}{X_t}$, and note that sequence Y_t is adapted to sequence X_t in the sense of Definition 2.5. We have the following bounds on the upper tails of Y_{t+1} , conditioned on $X_{[1:t]}$, given by Markov inequality (and recalling that $\mathbb{E}[X_{t+1}|X_t] = X_t$):

$$\Pr(Y_{t+1} > \lambda \mid \mathcal{F}_t) = \Pr\left(\frac{X_{t+1}}{X_t} > \exp(\lambda) \mid X_{[1:t]}\right) = \Pr\left(X_{t+1} > \exp(\lambda)X_t \mid X_{[1:t]}\right) \le \exp(-\lambda).$$

Let us decompose $Y_{t+1} =: (Y_{t+1})_+ + (Y_{t+1})_-$, where $(Y_{t+1})_+ := \max(Y_{t+1}, 0)$. By Lemma 2.6 and the fact that $(Y_{t+1})_+ \ge Y_{t+1}$,

$$\Pr\left(\sum_{t\leq T} (Y_{t+1})_+ > CT\right) \leq \exp(-\Omega(T)).$$

However, let E_{t+1} be the indicator of $Y_{t+1} \leq -K$. It is again adapted to the sequence X_t , and we know that $\Pr(E_{t+1}|X_{[1:t]}) \geq \alpha$; hence, by Lemma 2.8 with probability at most $\exp(-\Omega(\alpha T))$ at most $\alpha T/2$ of those events holds. Note that $(Y_t)_- \leq 0$, which implies that if at least $\alpha T/2$ of the events E_t hold, then we have $\sum_{t \leq T} (Y_t)_- \leq -\alpha KT/2$. Thus, we have $\Pr(\sum_{t \leq T} (Y_t)_- > -\alpha KT/2) \leq \exp(-\Omega(\alpha T))$. Therefore, as long as $\alpha K/4 > C$ (which is true if we set c = 4C), we can conclude

$$\Pr\left(\sum_{t\leq T}Y_t > -\alpha KT/4\right) \leq \exp(-\Omega(T)) + \exp(-\Omega(\alpha T)) \leq \exp(-\Omega(\alpha T)).$$

The proof is complete by noting that $\sum_{t \le T} Y_t = \log(X_T/X_0)$ and recalling that $X_0 \le 1$.

We are now ready to show that local polarization leads to strong polarization:

PROOF OF THEOREM 1.7. For given γ , we take K to be large enough so that $\exp(-\alpha K/8) \le \gamma$ and moreover αK to be large enough to satisfy assumptions of Lemma 3.3. Let us also take $\tau_0 = \tau(e^K)$. We consider ν as in Corollary 3.2. We have

$$\Pr\left(\min(X_{t/2}, 1 - X_{t/2}) > \left(1 - \frac{\nu}{2}\right)^t \tau_0\right) \le \left(1 - \frac{\nu}{4}\right)^{-t} \frac{1}{\sqrt{\tau_0}}.$$

Now Doob's martingale inequality (Corollary 2.10) implies that, conditioned on $X_{t/2} < (1 - \frac{\nu}{4})^t \tau_0$, we have $\Pr(\sup_{i \in (t/2,t)} X_i > \tau_0) \le (1 - \frac{\nu}{4})^t$.

Finally, after conditioning on $X_i \leq \tau_0$, $\forall t/2 \leq i \leq t$, process X_i for $i \in (t/2,t)$ satisfies conditions of Lemma 3.3, because X_i always stays below τ_0 and as such *suction at the end* condition of local polarization corresponds exactly to the assumption in this lemma. Therefore, we can conclude that except with probability $\exp(-\Omega(\alpha t)) + (1 - \frac{\nu}{4})^{-t} \frac{1}{\sqrt{\tau_0}}$ (which is $\exp(-\Omega_{\alpha,\nu}(t))$), we have $X_t < \exp(-\alpha Kt/8) = \gamma^t$. The other case $(1 - X_{t/2} < (1 - \frac{\nu}{2})^t \tau_0)$ is symmetric, and in this case we get $1 - X_t < \exp(-\alpha Kt/8)$ except with probability $\exp(-\Omega_{\alpha,\nu}(t))$.

3.3 Exponentially Strong Polarization

In this section, we prove the analog of Theorem 1.7-Theorem 1.9. We first prove a helper lemma.

LEMMA 3.4. There exist $C < \infty$ such that for all $0 < \eta < 1, b \ge 1, 0 < \varepsilon < 1$ following holds. Let X_t be a martingale satisfying $\Pr(X_{t+1} < X_t^b | X_t) \ge \eta$, where $X_0 \in (0, 1)$. Then

$$\Pr(\log X_T > (\log X_0 + CT)b^{(1-\varepsilon)\eta T}) < \exp(-\Omega(\varepsilon^2 \eta T)).$$

PROOF. As in the proof of Lemma 3.3, let us consider random variables $Y_{t+1} := \log(X_{t+1}/X_t)$. This sequence of random variables is adapted to the sequence X_t in the sense of Definition 2.5. Let us decompose $Y_t = Y_t^+ + Y_t^-$, where $Y_t^+ = \max(Y_t, 0)$. Note that by Markov inequality

$$\Pr\left(Y_{t+1} > \lambda | X_{[1:t]}\right) = \Pr\left(X_{t+1} > X_t \exp(\lambda) | X_{[1:t]}\right) \le \exp(-\lambda) \frac{\mathbb{E}[X_{t+1} | X_{[1:t]}]}{X_t} = \exp(-\lambda).$$

By Lemma 2.6, we deduce that for some C, we have

$$\Pr\left(\sum_{i\leq T}Y_i^+>CT\right)\leq \exp\left(-\Omega(T)\right).$$

However, if we take Z_t to be the indicator variable for an event $X_t < X_{t-1}^b$, then note that the sequence z_t is adapted to the sequence X_t . By Lemma 2.8, we have

$$\Pr\left(\sum_{i < T} Z_i \le (1 - \varepsilon)\eta T\right) \le \exp\left(-\Omega(T\varepsilon^2\eta)\right).$$

If neither of these unlikely events hold, that is, we simultaneously have $\sum_{i \leq T} Y_i^+ \leq CT$ and $\sum_{i \leq T} Z_i > (1 - \varepsilon)\eta T$, then we can deduce that $\log X_T \leq (\log X_0 + CT)b^{(1-\varepsilon)\eta T}$, i.e., the largest possible value of X_T is obtained if all the initial Y_i were positive and added up to CT (at which point value of the martingale would satisfy $\log X_{T'} \leq \log X_0 + CT$), followed by $(1 - \varepsilon)\eta T$ steps indicated by variables Z_i ; for each of those steps, $\log X_{t+1} \leq b \log X_t$.

We are now ready to prove the analog of Lemma 3.3 for exponentially strong polarization.

LEMMA 3.5. For all $0 < \eta < 1, b \ge 1, 0 < \varepsilon < 1$ the following holds. Let X_t be a martingale with values in (0,1) satisfying $\Pr(X_{t+1} < X_t^b | X_t) \ge \eta$, where $X_0 < \exp(-\gamma T)$ for some $\gamma > 0$, then

$$\Pr(\log X_T \geq -b^{(1-\varepsilon)\eta T}) < \exp(-\Omega_{\varepsilon,\eta,\gamma}(T))).$$

PROOF. Consider sequence $t_0, t_1, \ldots t_m \in [T]$, where $t_0 = 0, t_m = T$, and $\frac{\gamma T}{C} \leq |t_i - t_{i-1}| \leq \frac{\gamma T}{2C}$, and therefore $m = O(C\gamma^{-1})$, where C is a constant appearing in the statement of Lemma 3.4. For each index $s \in [m]$ we consider a martingale $X_i^{(s)} := X_{t_s+i}$, and we will apply Lemma 3.4 to this martingale $X^{(s)}$, with $T = t_{s+1} - t_s$. We can union bound total failure probability by $m \exp(-\Omega(\gamma \varepsilon^2 \eta T))$, which is upper bounded by the claim bound of $\exp(-\Omega_{\varepsilon,\eta,\gamma}(T))$.

In case we succeed, we can deduce that for each i we have

$$\log X_{t_i} < (\log X_{t_{i-1}} + C(t_i - t_{i-1}))b^{(1-\varepsilon)\eta(t_i - t_{i-1})}. \tag{15}$$

We will show that by our choice of parameters, we can bound $C(t_i - t_{i-1}) \le -\frac{1}{2} \log X_{t_i}$. Let us first discuss how this is enough to complete the proof. Indeed, in such a case we have

$$\log X_{t_i} < \frac{1}{2} (\log X_{t_{i-1}}) b^{(1-\varepsilon)\eta(t_i - t_{i-1})}, \tag{16}$$

and by induction

$$\log X_{t_m} < \frac{1}{2^m} (\log X_0) b^{(1-\varepsilon)\eta t_m}.$$

11:24 J. Błasiok et al.

For fixed η , m, and T large enough (depending on η , m, ε), this yields $\log X_T < -b^{(1-2\varepsilon)\eta T}$, and the result follows up by changing ε by a factor of 2.

All we need to do is to show is that for every *i* we have

$$C(t_{i+1} - t_i) \le -\frac{1}{2} \log X_{t_i},\tag{17}$$

assuming that Equation (15) holds for every i. We will show this inductively, together with $\log X_{t_i} \leq -\gamma T$. Note that we assumed this inequality to be true for $X_{t_0} = X_0$. By our choice of parameters we have $C(t_{i+1} - t_i) \leq \frac{\gamma T}{2}$, and therefore for t_{i+1} the inequality (17) is satisfied.

We will now show that $\log X_{t_{i+1}} \leq \log X_{t_i} \leq -\gamma T$ to finish the proof by induction. We can apply Equation (16) to X_{t_i} to deduce that $\log X_{t_{i+1}} \leq \frac{1}{2}(\log X_{t_i})b^{\frac{1}{2}\frac{\gamma}{C}T}$ (which is true, since $b \geq 1$, $\eta \leq 1$, $\epsilon \geq 0$). This for large values of T (given parameters b, γ and C) yields $\log X_{t_{i+1}} < \log X_{t_i}$; indeed, this inequality will be true as soon as $b^{\frac{\gamma}{2C},T} > 2$, because both $\log X_{t_{i+1}}$ and $\log X_{t_i}$ are negative, which completes the proof.

We are now ready to prove local polarization to global polarization theorem for exponential polarization.

PROOF OF THEOREM 1.9. Consider exponentially locally polarizing martingale, and let us fix some $\varepsilon > 0$. By Corollary 3.2 with $t = 2\varepsilon T$ and $\lambda = 1$, we deduce that for some $\nu > 0$ we have

$$\Pr\left(\max\left(X_{\varepsilon T}, 1 - X_{\varepsilon T}\right) \ge \left(1 - \frac{\nu}{2}\right)^{2\varepsilon T}\right) < \exp\left(-\Omega_{\varepsilon, \nu}(T)\right).$$

We condition on $\max(X_{\varepsilon T}, 1 - X_{\varepsilon T}) < (1 - \frac{\nu}{2})^{2\varepsilon T}$. Now let K be a large-enough constant depending on α and γ , the target rate of polarization in the high end. Now let $\tau > 0$ be such that $\tau \leq \min\left(\tau\left(e^K\right), \tau_0\right)$, where τ_0 is given by the definition of suction at the low end and $\tau(\cdot)$ is from the suction at the high end. Note that this implies that (1) if $X_t < \tau$, then we have

$$\Pr\left(X_{t+1} < X_t^b | X_t\right) \ge \eta,\tag{18}$$

which holds, since $\tau \leq \tau_0$, and (2) if $1 - X_t < \tau$, then we have

$$\Pr((1 - X_{t+1}) < \exp(-K)(1 - X_t)|X_t) \ge \alpha,\tag{19}$$

which follows from the condition on suction at the high end. By Doob's martingale inequality (specifically Corollary 2.10), we deduce that $\Pr(\max_{t \in [\varepsilon T, T]} \max(X_t, 1 - X_t) > \tau) \leq \tau^{-1}(1 - \frac{\nu}{2})^{-2\varepsilon T} \leq \exp(-\Omega_{\tau, \nu, \varepsilon}(T))$. Let us now condition in turn on this event not happening.

We will consider first the case when $X_{\varepsilon T} < (1 - \frac{\nu}{2})^{2\varepsilon T}$, and let us put $\gamma_0 := -2\varepsilon \log(1 - \frac{\nu}{2})$ (note that $\gamma_0 > 0$), so that $X_{\varepsilon T} < \exp(-\gamma_0 T)$. We can now apply Lemma 3.5 to the martingale sequence starting with $X_{\varepsilon T}$. (Note that the assumptions of Lemma 3.5 are satisfied as long as X_t for $t \in [\varepsilon T, T]$ stays bounded by τ due to Equation (18).) Hence, we deduce that in this case, except with probability $\exp(-\Omega_{\gamma_0,\varepsilon,\eta}(T)) \le \exp(-\Omega_{\nu,\varepsilon,\eta}(T))$, we have

$$\log X_T < -b^{(1-\varepsilon)^2\eta T},$$

and therefore $X_T < 2^{-b^{(1-\varepsilon)^2\eta T}}$. Note that this implies that $\Lambda = \log_2(b^{(1-\varepsilon)^2\eta}) = (1-\varepsilon)^2\eta\log_2 b$ (hence for any $\Lambda < \eta\log_2 b$, we pick ε appropriately). We also pick β and η such that $\beta\eta^T \ge \exp(-\Omega_{\varepsilon,\eta,\mu,K}(T))$.

However, if $1 - X_t < \tau$ for all $\varepsilon T \le t \le T$, then the suction at the high end condition of local polarization applies (i.e., Equation (19) holds), and we can apply Lemma 3.3 (we pick K large enough so that $K\alpha > c$) to martingale $\tilde{X}_t \triangleq 1 - X_{\varepsilon T + t}$ to deduce that except with probability

 $\exp(-\Omega_{\alpha}(T))$, we have $1 - X_T < \exp(-\alpha K(1 - \varepsilon T)/4) < \gamma^T$ for suitable choice of K depending on γ and α . Finally, we pick β and η such that $\beta \eta^T \ge \exp(-\Omega_{\varepsilon,\eta,\mu,K,\alpha}(T))$.

4 ARIKAN MARTINGALE AND ITS LOCAL POLARIZATION

In this section, we formally describe the Arıkan martingale associated with an invertible matrix $M \in \mathbb{F}_a^{k \times k}$ and a channel $C_{Y|Z}$.

Before we proceed with the formal definition, to provide overview of the goals of this construction, we shall briefly point out its main features for the special case of Arıkan martingale $\{X_t\}_{t=0}^{\infty}$ associated with an additive channel C—where channel output Y = Z + U, with U being some random variable in \mathbb{F}_q not depending on Z.

- (1) For given t, marginal distribution X_t is distributed identically as $\overline{H}((UM^{\otimes t})_i|(UM^{\otimes t})_{< i})$ for uniformly random index i, where U is a vector of k^t i.i.d. random variables distributed as the error U.
- (2) Sequence X_t is a martingale—in particular, we provide coupling of the distributions above over different t in a non-trivial way.
- (3) Definition of the martingale X_t is "local" in some sense, which makes it manageable to analyze how X_t and X_{t+1} are related and eventually show local polarization.

In Appendix A.2.3, we elaborate on the connection of the Arıkan martingale with polar codes—specifically, the main link is a more general version of the first property for all symmetric channels and is proved as Lemma A.18.

Briefly, the Arıkan martingale measures at time t, the distribution of conditional entropy of a random variable A'_i , conditioned on the values of a vector of variables B' and on the values of A'_j for j smaller (according to \prec) than i for a random choice of the index i. Here A' is a vector of k^t random variables taking values in \mathbb{F}_q while $B' \in \mathcal{Y}^{k^t}$. The exact construction of the joint distribution of these $2k^t$ variables is the essence of the Arıkan construction of codes, and we describe it shortly. The hope with this construction is that eventually (for large values of t) the conditional entropies are either very close to 0 or very close to $\log q$ for most choices of i.

When t = 1, the process starts with k independent and identical pairs of variables $\{(A_i, B_i)\}_{i \in [k]}$, where $A_i \sim \mathbb{F}_q$ and $B_i \sim C_{Y|Z=A_i}$. (So each pair corresponds to an independent input/output pair from transmission of a uniformly random input over the channel $C_{Y|Z}$.) Let $A = (A_1, \ldots, A_k)$ and $B' = (B_1, \ldots, B_k)$, and note that the conditional entropies $H(A_i|A_{< i}, B')$ are all equal, and this entropy, divided by $\log_2 q$, will be our value of X_0 . However, if we now let $A' = A \cdot M$, then the conditional entropies $H(A'_i|A'_{< i}, B')$ are no longer equal (for most, and in particular for all mixing, matrices M). However, conservation of conditional entropy on application of an invertible transformation tells us that $\mathbb{E}_{i\sim [k]}[H(A'_i|A'_{< i}, B')/\log_2 q] = X_0$. Thus letting $X_1 = H(A'_i|A'_{< i}, B')/\log_2 q$ (for random i) gives us the martingale at time t = 1.

While this one step of multiplication by M differentiates among the k (previously identical) random variables, it does not yet polarize. The hope is that by iterating this process one can get polarization. But to get there we need to describe how to iterate this process. This iteration is conceptually simple (though notationally still complex) and illustrated in Figure 2. Roughly, the idea is that at the beginning of stage t, we have defined a joint distribution of k^t -dimensional vectors (A, B) along with a multi-index $i \in [k]^t$. We now sample k independent and identically distributed pairs of these random variables $\{(A^{(\ell)}, B^{(\ell)})\}_{\ell \in [k]}$ and view $(A^{(\ell)})_{\ell \in [k]}$ as a $k^t \times k$ matrix, which we multiply by M to get a new $k^t \times k$ matrix. Flattening this matrix into a k^{t+1} -dimensional

¹²In the context of polar coding, differentiation and polarization are good events, and hence our "hope."

11:26 J. Błasiok et al.

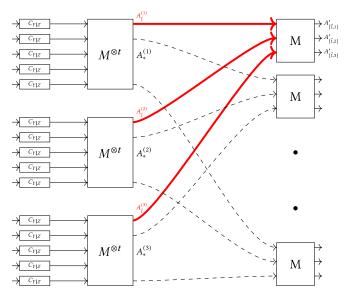


Fig. 2. Evolution of Arıkan martingale for 3×3 matrix M.

vector gives us a sample from the distribution of $A' \in \mathbb{F}_q^{k^{t+1}}$. B' is simply the concatenation of all the vectors $(B^{(\ell)})_{\ell \in [k]}$. And, finally, the new index $j \in [k]^{t+1}$ is simply obtained by extending $i \in [k]^t$ with a (t+1)th coordinate distributed uniformly at random in [k]. X_{t+1} is now defined to be $\overline{H}(A'_i|A'_{< i}, B')$, where $\overline{H}(\cdot)$ is the normalized q-ary entropy defined in Equation (3). The formal description is below.

Definition 4.1 (Arıkan Martingale). Given an invertible matrix $M \in \mathbb{F}_q^{k \times k}$ and a channel description $C_{Y|Z}$ for $Z \in \mathbb{F}_q, Y \in \mathcal{Y}$, the Arıkan-martingale $X_0, \ldots X_t, \ldots$ associated with it is defined as follows. For every $t \in \mathbb{N}$, let D_t be the distribution on pairs $\mathbb{F}_q^{k^t} \times \mathcal{Y}^{k^t}$ described inductively below: A sample (A, B) from D_0 supported on $\mathbb{F}_q \times \mathcal{Y}$ is obtained by sampling $A \sim \mathbb{F}_q$, and $B \sim C_{Y|Z=A}$. For $t \geq 0$, a sample $(A', B') \sim D_{t+1}$ supported on $\mathbb{F}_q^{k^{t+1}} \times \mathcal{Y}^{k^{t+1}}$ is obtained as follows:

- Draw k independent samples $(A^{(1)}, B^{(1)}), \ldots, (A^{(k)}, B^{(k)}) \sim D_t$. Let A' be given by $A'_{[i,\cdot]} = (A^{(1)}_i, \ldots, A^{(k)}_i) \cdot M$ for all $i \in [k]^t$ and $B' = (B^{(1)}, B^{(2)}, \ldots, B^{(k)})$.

Then, the sequence X_t is defined as follows: Sample $i_l \in [k]$ uniformly and independently for $l=1,2,\ldots,t$. Let $j=(i_1,\ldots,i_t)$, and let $X_t:=\overline{H}(A_j|A_{< j},B)$, where the entropies are with respect to the distribution $(A, B) \sim D_t$.¹³

Figure 2 illustrates the definition by highlighting the construction of the vector A' and in particular highlights the recursive nature of the construction.

It is easy (and indeed no different than in the case t=1) to show that $\mathbb{E}[X_{t+1}|X_t]=X_t$, and so the Arıkan martingale is indeed a martingale. This is shown below.

Proposition 4.2. For every matrix M and channel $C_{Y|Z}$, the Arikan martingale is a martingale and in particular a [0, 1]-martingale.

¹³We stress that the only randomness in the evolution of X_t is in the choice of i_1, \ldots, i_t, \ldots . The process of sampling Aand B is only used to define the distributions for which we consider the conditional entropies $H(A_j|A_{< j},B)$.

PROOF. The fact that $X_t \in [0,1]$ follows from the fact for $0 \le H(A_i|A_{< i},B) \le H(A_i) \le \log_2 q$ (the upper bound follows, since $A_{< i} \in \mathbb{F}_q$) and so $0 \le X_t = H(A_i|A_{< i},B)/\log_2 q \le 1$.

We turn to showing that $\mathbb{E}[X_{t+1}|X_t=a]=a$. To this end, consider a sequence of indices $i=(i_1,\ldots i_t)$, such that $\overline{H}(A_i\mid A_{< i},B)=a$. We wish to show that $\mathbb{E}_{i_{t+1}\sim [k]}[\overline{H}(A'_{[i,i_{t+1}]}\mid A'_{<[i,i_{t+1}]},B')]=a$.

Since the pairs $(A^{(s)}, B^{(s)})$ are independent samples from D_t , note that for any s, we have $\overline{H}(A_i^{(s)} | A_{< i}^{(s)}, B^{(s)}) = a$. Furthermore, because of the same independence, we have

$$\overline{H}\left(A_{i}^{(s)} \mid A_{\prec i}^{(s)}, B^{(s)}\right) = \overline{H}\left(A_{i}^{(s)} \mid \bigcup_{j \in [k]} A_{\prec i}^{(j)}, \bigcup_{j \in [k]} B^{(j)}\right)$$
and
$$\overline{H}\left(A_{i}^{(1)}, \dots, A_{i}^{(k)} \mid \bigcup_{j \in [k]} A_{\prec i}^{(j)}, \bigcup_{j \in [k]} B^{(j)}\right) = k \cdot a.$$

By the invertibility of M, we have

$$\overline{H}(A'_{[i,1]}, \dots, A'_{[i,k]} \mid \cup_{j \in [k]} A^{(j)}_{< i}, \cup_{j \in [k]} B^{(j)}) = \overline{H}(A^{(1)}_{i}, \dots, A^{(k)}_{i} \mid \cup_{j \in [k]} A^{(j)}_{< i}, \cup_{j \in [k]} B^{(j)}) = k \cdot a.$$

We can apply again invertibility of the matrix M to deduce that conditioning on $\bigcup_{j \in [k]} A_{< i}^{(j)}$ is the same as conditioning on $A'_{<[i,1]}$, i.e., for any multiindex i' < i variables $A_{i'}^{(1)}, \ldots A_{i'}^{(k)}$ and $A'_{[i',1]}, \ldots A'_{[i',k]}$ are related via invertible transform M. This yields

$$\overline{H}(A'_{[i,1]}, \dots A'_{[i,k]} \mid A'_{<[i,1]}, B') = \overline{H}(A'_{[i,1]}, \dots A'_{[i,k]} \mid \bigcup_{j \in [k]} A_{$$

Finally, by the Chain rule of entropy we have

$$\overline{H}(A'_{[i,1]}, \dots A'_{[i,k]} \mid A'_{<[i,1]}, B') = \sum_{i_{t+1}=1}^{k} \overline{H}(A'_{[i,i_{t+1}]} \mid A'_{[i,< i_{t+1}]}, A'_{<[i,1]}, B')$$

$$= \sum_{i_{t+1}=1}^{k} \overline{H}(A'_{[i,i_{t+1}]} \mid A'_{<[i,i_{t+1}]}, B')$$

Putting these together, we have $\mathbb{E}[X_{t+1}|X_t=a]=\mathbb{E}_{i_{t+1}}[\overline{H}(A'_{[i,i_{t+1}]}|A'_{\langle [i,i_{t+1}]},B')]=\frac{1}{k}\cdot ka=a.$

Finally, we remark that based on the construction it is not too hard to see that if M were an identity matrix, or more generally a non-mixing matrix, then X_t would deterministically equal X_0 . (There is no differentiation and thus no polarization.) The thrust of this article is to show that in all other cases we have strong polarization.

4.1 Matrix Polarization and the Arıkan Martingale

Note that the definition of the Arıkan martingale is itself complex, and in particular the distribution of X_t , the variable at the tth step, needs a description whose complexity grows with t. The essence of the polarization argument does not depend on this intricacy of the definition, most of which can be abstracted away. Indeed, we do so formally by considering a simpler (single step) randomization process associated with a matrix M. We define a matrix M to be *polarizing* if this single-step process satisfies properties similar to those of local polarization (see Definition 4.3). Then, in Theorem 4.4 we show that if a matrix M satisfies matrix polarization, then for every channel C the Arıkan martingale associated with M and C is locally polarizing. This is a notationally heavy but conceptually light proof, whose essence is to verify that certain variables are independent, and so conditioning on such variables does not change entropies. This will allow us focus on a simpler single step process in future sections to prove (exponential) local polarization.

11:28 J. Błasiok et al.

We start with the definition of matrix polarization.

Definition 4.3 (Matrix (Exponential) Polarization). We say that a matrix $M \in \mathbb{F}_q^{k \times k}$ satisfies matrix polarization if and only if for every pair of random variables (U, W), such that $U = (U_1, \dots, U_k) \in \mathbb{F}_q^k$, $W = (W_1, \dots, W_k)$ is supported on some finite set, and the pairs (U_i, W_i) are independently and identically distributed for $i \in [k]$, the vector $V := U \cdot M$ satisfies the following conditions:

(1) **(Variance in the middle):** There is some index $j \in [k]$ for which the following holds: For every $\tau > 0$, there exists $\varepsilon = \varepsilon(\tau) > 0$ such that if $\overline{H}(U_1|W) \in (\tau, 1 - \tau)$, then

$$\overline{H}((V)_j|V_{< j},W) \ge \overline{H}(U_1|W) + \varepsilon.$$

(2) **(Suction at the lower end):** There is some index $j \in [k]$ for which the following holds: For every $c < \infty$, there exists $\tau > 0$ such that if $\overline{H}(U_1|W) < \tau$, then

$$\overline{H}(V_j|V_{< j},W) \leq \frac{1}{c}\overline{H}(U_1|W).$$

(3) **(Suction at the high end):** Analogously to suction at the low end, there is some index $j \in [k]$ for which the following holds:

For every $c < \infty$, there exists $\tau > 0$ such that if $\overline{H}(U_1|W) > 1 - \tau$, then

$$1 - \overline{H}(V_j|V_{< j}, W) \le \frac{1}{c} (1 - \overline{H}(U_1|W)).$$

Additionally, we say that M satisfies (η, b) -exponential matrix polarization if we have the following property:

2'. (Strong Suction at the lower end): There exists $\tau > 0$ such that if $\overline{H}(U_1|W) < \tau$, then for at least η fraction of the indices $j \in [k]$ we have then

$$\overline{H}(V_j|V_{< j},W) \leq \overline{H}(U_1|W)^b.$$

Thus, the notion of matrix polarization is somewhat more general than polarization of the corresponding Arıkan martingale.

- (1) In the latter, the conditioning in the entropy prescribes some specific relations between U and W, where in the former W is arbitrary (subject to the condition that the pairs (U_j, W_j) are i.i.d.).
- (2) Furthermore, the definitions also make slight changes to the conditions of Variance in the middle and suction only requiring the existence of $j \in [k]$ having a certain property as opposed requiring that a random choice of j satisfy some condition.

The differences in (1) above allows for cleaner proofs, since the specific structure of W is not needed. The class of differences in (2) above changes some probabilities and/or variances by factors depending on k, but this difference is negligible. In the following section, we formally confirm that matrix polarization is a sufficient condition for martingale polarization.

4.2 Matrix Polarization Implies Local Polarization of Arıkan Martingale

In this section, we prove that matrix polarization implies local polarization of the corresponding Arıkan martingale.

Theorem 4.4. For every matrix $M \in \mathbb{F}_q^{k \times k}$ and every symmetric memoryless channel $C_{Y|Z}$, if M satisfies matrix polarization, then the Arikan martingale associated with M and C is satisfies local polarization. Furthermore, if M satisfies (η, b) -exponential matrix polarization, then the Arikan-martingale satisfies (η, b) -exponential local polarization.

We begin with a lemma that will be useful in the proof of Theorem 4.4:

LEMMA 4.5. Let $A^{(1)}, \ldots A^{(k)}$, and A' be defined as in Definition 4.1, and let V, W be arbitrary random variables. Then for any multiindex $i \in [k]^t$ and any $i_{t+1} \in [k]$ we have

$$\overline{H}\Big(V\mid A_{<[i,i_{t+1}]}',W\Big)=\overline{H}\Big(V\mid A_{[i,< i_{t+1}]}',A_{< i}^{(1)},A_{< i}^{(2)},\ldots A_{< i}^{(k)},W\Big).$$

Proof. Observe first that by definition of the order \prec we have that $A'_{\prec[i,i_{t+1}]} = (A'_{\prec[i,1]}, A'_{(i,\prec i_{t+1})})$, hence

$$\overline{H}\big(V\mid A_{<[i,i_{t+1}]}',W\big)=\overline{H}\big(V\mid A_{[i,< i_{t+1}]}',A_{<[i,1]}',W\big).$$

The definition of the sequence A' in terms of A (in Definition 4.1) reads

$$A'_{[j,\cdot]} = \left(A_j^{(1)}, \ldots, A_j^{(k)}\right) M.$$

Note that if random variables B, B' are related by invertible function B = f(B'), then $\overline{H}(A|B) = \overline{H}(A|B')$. By definition of mixing matrix, M is invertible, and hence variables $A'_{<[i,1]}$ and variables $A^{(1)}_{< i}, \ldots A^{(k)}_{< i}$ are indeed related by invertible (linear) transformation, which yields

$$\overline{H} \Big(V \mid A'_{[i, < i_{t+1}]}, A'_{<[i, 1]}, W \Big) = \overline{H} \Big(V \mid A'_{[i, < i_{t+1}]}, A^{(1)}_{< i}, A^{(2)}_{< i}, \dots A^{(k)}_{< i}, W \Big) \; . \qquad \qquad \Box$$

We now turn to the proof of Theorem 4.4.

PROOF OF THEOREM 4.4. Fix a matrix M, channel $C_{Y|Z}$, and time t. We start by recalling the definition of the variables X_t and X_{t+1} in the definition of the Arıkan martingale and also recall what local polarization entails for these variables.

Let $(A,B), (A^{(1)},B^{(1)}), \ldots (A^{(k)},B^{(k)}) \sim D_t$ denote independent random variables. Let (A',B') constructed from $(A^{(1)},B^{(1)}), \ldots (A^{(k)},B^{(k)})$ as in Definition 4.1, i.e., we have $A'_{[i',\cdot]} = (A^{(1)}_{i'},\ldots,A^{(k)}_{i'})\cdot M$ for every $i'\in [k]^t$ and $B'=(B^{(1)},B^{(2)},\ldots B^{(k)})$. Now let $i\triangleq (i_1,\ldots i_t)$ be sampled uniformly from $[k]^t$ and let $i_{t+1}\in [k]$ be chosen independently and uniformly.

Then by Definition 4.1 we have

$$X_t = \overline{H}(A_i \mid A_{< i}, B)$$

and
$$X_{t+1} = \overline{H}(A'_{[i,i_{t+1}]}|A'_{<[i,i_{t+1}]},B').$$

That is, for
$$U = (A_i^{(1)}, \dots, A_i^{(k)})$$
, we have $A'_{[i,\cdot]} = U \cdot M$, and $B' = (B^{(1)}, \dots, B^{(k)})$.

We will use the property of matrix polarization of M with $U = (U_1, \ldots, U_k)$, where $U_s = A_i^{(s)}$ and $W = (W_1, \ldots, W_k)$, where $W_s = (A_{>i}^{(s)}, B^{(s)})$ to deduce local polarization of Arıkan martingale. Note that the pairs $(U_1, W_1), \ldots, (U_k, W_k)$ are identically distributed and independent as required. We let $V = U \cdot M$. By the definition of Arıkan martingale we have $A'_{[i,\cdot]} = V = U \cdot M$. Thus, the matrix polarization of M implies bounds on the conditional entropy of $\overline{H}(V_j|V_{< j}, W)$, where $V_j = A'_{[i,j]}$, where X_{t+1} studies conditional entropy of $(V_{i_{t+1}}|A'_{<[i,i_{t+1}]}, B')$. In what follows, we verify that despite the difference the latter can be bounded as required for the condition of (exponential) local polarization of the Arıkan martingale. We tackle each of the conditions in order, but first we note that by Lemma 4.5 we have

$$\overline{H}(A'_{[i,j]} \mid A'_{<[i,j]}, B') = \overline{H}(A'_{[i,j]} \mid A'_{[i,
(20)$$

11:30 J. Błasiok et al.

Let $h \triangleq X_t = \overline{H}(A_i \mid A_{< i}, B)$. Note that for every $s \in [k]$ we have $\overline{H}(A_i^{(s)} \mid A_{< i}^{(s)}, B^{(s)}) = h$, because all the pairs $(A^{(s)}, B^{(s)})$ are distributed independently and identically to (A, B). Moreover, for every $j \in [k]$ we have $\overline{H}(U_j \mid W) = \overline{H}(A_i^{(s)} \mid A_{< i}^{(s)}, B^{(s)}) = h$, where the first equality follows from the fact that pairs $(A^{(s)}, B^{(s)})_{s \in [k]}$ are identically and independently distributed (so the index j does not matter).

We will start with the *Variance in the middle* condition of martingale local polarization (Definition 1.6). As a reminder, what we need to show is that if $h \in (\tau, 1 - \tau)$, then

$$\operatorname{Var}_{i_{t+1} \sim [k]} (\overline{H}(A'_{[i,i_{t+1}]} | A'_{<[i,i_{t+1}]}, B') - \overline{H}(A_i | A_{< i}, B)) > \theta(\tau).$$

Note that by the martingale property (Proposition 4.2) we have

$$\mathbb{E}_{i_{t+1} \sim [k]} [\overline{H}(A'_{[i,i_{t+1}]}|A'_{<[i,i_{t+1}]},B') - \overline{H}(A_i|A_{< i},B)] = 0,$$

and as such to obtain the lower bound on the variance it is enough to show that

$$\Pr_{i_{t+1} \sim [k]} \left(\overline{H}(A'_{[i,i_{t+1}]} | A'_{\sim [i,i_{t+1}]}, B') \ge h + \varepsilon(\tau) \right) \ge \frac{1}{k}. \tag{21}$$

This would allow us to deduce that the variance above is lower bounded by $\varepsilon(\tau)^2/k$. (Note that this lower bound is true for every h, and hence the actual variance needed in the statement of the *Variance in the middle* condition is also true.)

We now use the *Variance in the middle* condition of the matrix polarization (Definition 4.3) for M with variables (U, W). This condition asserts that for some index j we have entropy gain $\overline{H}((U \cdot M)_j | (U \cdot M)_{< j}, W) \ge h + \varepsilon(\tau)$. Combining this with Equation (20) proves inequality (21) and therefore shows *variation in the middle* for the Arıkan martingale.

Next, we turn to the proof of *suction at the upper end*. Here, we will show that for every c if $1-h < \tau(c)$, then with probability at least $\frac{1}{k}$ over the choice of i_{t+1} , we will have $1-\overline{H}(A'_{[i,i_{t+1}]}|A'_{<[i,i_{t+1}]},B') \leq \frac{1}{c}(1-h)$. The corresponding *suction at the upper end* condition of matrix polarization asserts the existence of an index j, such that $1-\overline{H}(V_j|V_{< j},W) \leq \frac{1}{c}(1-h)$. With probability at least $\frac{1}{k}$, we have $i_{t+1}=j$, and in this case we have

$$1 - \overline{H}(A'_{[i,i_{t+1}]} \mid A'_{<[i,j]}, B') = 1 - \overline{H}((U \cdot M)_j \mid (U \cdot M)_{< j}, W) \le \frac{1}{c}(1 - h),$$

where the first equality above is by Equation (20).

The proof of *suction at the lower end* is symmetric. We now turn to the proof of *strong suction at the low end* (Definition 1.8). Let M satisfy (η, b) -exponential matrix polarization. Recall that we wish to show that $\Pr_{i_{t+1} \sim [k]}(\overline{H}(A'_{[i,i_{t+1}]} \mid A'_{<[i,i_{t+1}]}, B') < h^b) \ge \eta$. Once again by Equation (20), we have, for every $j \in [k]$, $\overline{H}(A'_{[i,j]} \mid A'_{<[i,j]}, B') = \overline{H}(V_j \mid V_{<j}, W)$. This is exactly the property given by the *strong suction at the low end* property of matrix polarization (using $h = \overline{H}(U_1|W)$).

This concludes the proof.

Thus, to prove Theorems 1.15 and 1.16, we now need to prove that for every mixing matrix M, M satisfies matrix polarization and M^2 satisfies exponential polarization. We argue the former in Section 5 and the latter in Section 7.

5 PROOF OF MATRIX POLARIZATION

In this section, we prove that every mixing matrix satisfies matrix polarization, modulo some entropic inequalities whose proofs are deferred to Section 6. Combined with Theorem 4.4, this immediately yields Theorem 1.15, which asserts the local polarization of the Arıkan martingale.

Informally, this section can be viewed as reducing matrix polarization of a general (mixing) matrix to the matrix polarization of the matrix $G_2 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$. Formally, what we do is state three entropic inequalities (see Section 5.1) that arise naturally in the proof of the matrix polarization of G_2 . These inequalities relate the conditional entropy of a sum of two random variables to the entropy of each of those random variables. Indeed, these inequalities can be used to show immediately that G_2 satisfies matrix polarization, and we do so in Lemma 5.4. But the bulk of the work, and novelty, in this section is in Section 5.3 where we show (via carefully executed "Gaussian elimination") that these entropic inequalities suffice to show the matrix polarization of *every* mixing matrix.

5.1 Entropic Lemmas in the 2×2 Case

We state here the three entropic lemmas. The proofs of the first two are deferred to Section 6. The third lemma is well known and we provide a reference for its proof.

The first lemma arises from the analysis of the suction at the upper end (for $X_t > 1 - \tau$) condition of Definition 4.3.

LEMMA 5.1. For every finite field \mathbb{F}_q and every $\gamma > 0$, there exist τ , such that if (X_1, A_1) and (X_2, A_2) are independent random variables with $X_i \in \mathbb{F}_q$ such that $1 - \overline{H}(X_2 \mid A_2) \leq \tau$, then

$$1 - \overline{H}(X_1 + X_2 \mid A_1, A_2) \le \gamma (1 - \overline{H}(X_1 \mid A_1)).$$

The next lemma comes analogously from the *suction at the low end* (for $X_t < \tau$) condition of Definition 4.3.

LEMMA 5.2. For every finite field \mathbb{F}_q and every $\gamma > 0$, there exist τ such that the following holds. Let (X_1, A_1) and (X_2, A_2) be any pair of independent random variables with $X_i \in \mathbb{F}_q$, and such that A_1, A_2 are identically distributed, and moreover for every a we have $\overline{H}(X_1 \mid A_1 = a) = \overline{H}(X_2 \mid A_2 = a)$. If $\overline{H}(X_1 \mid A_1) = \overline{H}(X_2 \mid A_2) \leq \tau$, then we have

$$H(X_1 \mid X_1 + X_2, A_1, A_2) \le \gamma \overline{H}(X_1 \mid A_1).$$

Finally, we use the following lemma from Reference [7, Lemma 4.2], which corresponds to the *Variance in the middle* condition of Definition 4.3. This is the only place where we need the field size q to be prime.

Lemma 5.3 ([7, Lemma 4.2]). For every $\tau > 0$ and prime finite field \mathbb{F}_q , there exist $\varepsilon > 0$ such that if (X_1, A_1) and (X_2, A_2) are independent pairs of random variables (but not necessarily identically distributed), with $X_i \in \mathbb{F}_q$ for some prime q, then

$$\overline{H}(X_1 \mid A_1), \overline{H}(X_2 \mid A_2) \in (\tau, 1 - \tau)$$

implies

$$\overline{H}(X_1 + X_2 | A_1, A_2) \ge \max\{\overline{H}(X_1 \mid A_1), \overline{H}(X_2 \mid A_2)\} + \varepsilon.$$

5.2 Matrix Polarization of Arıkan's 2×2 Kernel

As an illustration of how the lemmas arise in the study of matrix polarization, we prove that G_2 satisfies matrix polarization. We remark that we do not need this lemma for the rest of this article—we present it purely as an example.

Lemma 5.4. Over every prime field \mathbb{F}_q , the matrix $G_2 = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$ satisfies matrix polarization.

PROOF. Note that we have

$$(V_1, V_2) = (U_1, U_2) \cdot \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix},$$

11:32 J. Błasiok et al.

i.e., $V_1 = U_1 + U_2$ and $V_2 = U_2$. By Lemma 5.3, we have that the choice j = 1 satisfies the *variance* in the middle condition of matrix polarization for G_2 . By Lemma 5.2, we have the choice j = 2 yields the suction at the low end condition (with $c = \frac{1}{\gamma}$) of matrix polarization for G_2 . Finally, by Lemma 5.1 we have that the choice j = 1 satisfies the Suction at the upper end condition (with $c = \frac{1}{\gamma}$) for G_2 .

5.3 Polarization of $k \times k$ Mixing Matrices

In this section, we prove the following.

LEMMA 5.5. For every prime field \mathbb{F}_q and every positive k every mixing matrix $M \in \mathbb{F}_q^{k \times k}$ satisfies matrix polarization.

We will apply Gaussian elimination on M to reduce to the entropic inequalities of the 2×2 case from Section 5.1. The high-level strategy for showing polarization of $k \times k$ mixing matrix M is as follows. Consider i.i.d. random variables $(U_1, W_1), \ldots (U_k, W_k)$, and linearly transformed variables $V = U \cdot M$, where $U = (U_1, \ldots, U_k)$.

In Section 5.4, we will show that

(1) There are some indices $j, \ell, s \in [k]$ and some $\alpha \in \mathbb{F}_q^*$ for which

$$\overline{H}(V_i|V_{\leq i},W) \geq \overline{H}(U_\ell + \alpha U_s|W).$$

(2) There are some indices $j, \ell, s \in [k]$ and some $\alpha \in \mathbb{F}_a^*$ for which

$$\overline{H}(V_j|V_{< j},W) \leq \overline{H}(U_\ell|U_\ell + \alpha U_s,W).$$

Those two, together with entropic inequalities stated in Section 5.1, are enough to show polarization of a given matrix.

Before we proceed with the formal proof of those two inequalities, we give an informal exposition of the main idea behind it. For the sake of this exposition, let us focus on the inequality $\overline{H}(V_i|V_{< j},W) \leq \overline{H}(U_\ell|U_\ell+\alpha U_s,W)$, and let us skip conditioning on W.

The main observation is that if $B_1, \ldots B_m, B_{m+1}$ are all linear combinations of variables $V_1, \ldots V_{j-1}$, then we have $\overline{H}(V_j|V_{< j}) = \overline{H}(V_j + B_{m+1}|V_{< j}, B_1, \ldots B_m) \leq \overline{H}(V_j + B_{m+1}|B_1, \ldots B_m)$. Here it is enough to instantiate this observation with m=1. Since variables V_i are themselves linear combinations of variables U (with coefficients given by matrix M), all we need to do is find an index j, and some indices ℓ , s, such that

$$U_{\ell} \in V_i + \operatorname{span}\{V_{\leq i}\} \quad \text{and} \quad U_{\ell} + \alpha U_s \in \operatorname{span}\{V_{\leq i}\}.$$
 (22)

In particular, we use $B_2 \in \text{span}\{V_{< j}\}$ for the first inclusion to set $U_\ell = V_j + B_2$ and use $B_1 \in \text{span}\{V_{< j}\}$ for the second inclusion to set $U_\ell + \alpha U_s = B_1$. Note that with the inequalities in the above paragraph would give us the desired inequality.

Turns out that if the matrix M is mixing, then this can be achieved by carefully applying Gaussian Elimination on this matrix, as we explain next.

5.4 Reduction to the 2×2 Case

This section will be devoted to proving following two lemmas.

LEMMA 5.6 (REDUCTION FOR SUCTION AT THE UPPER END AND VARIANCE). Let (U, W) be a joint distribution, where $U = (U_1, \ldots, U_k) \in \mathbb{F}_q^k$ (with U_i for $i \in [k]$ being independent conditioned on W), and let M be any mixing matrix. Then, there exist three indices $j, \ell, s \in [k]$, and $\alpha \in \mathbb{F}_q^*$, such that

$$\overline{H}((UM)_j \mid (UM)_{< j}, W) \ge \overline{H}(U_\ell + \alpha U_s \mid W).$$

LEMMA 5.7 (REDUCTION FOR SUCTION AT THE LOWER END). Let (U, W) be a joint distribution, where $U = (U_1, \ldots, U_k) \in \mathbb{F}_q^k$, and let M be any mixing matrix. Then, there exist three indices $j, \ell, s \in [k]$, and $\alpha \in \mathbb{F}_q^*$, such that

$$\overline{H}((UM)_i \mid (UM)_{< i}, W) \leq \overline{H}(U_\ell \mid U_\ell + \alpha U_s, W).$$

As discussed previously, to show Lemma 5.7 and Lemma 5.6, we will apply Gaussian Elimination to prove the following three lemmas about mixing matrices.

We start with a simple equivalent characterization of a mixing matrix:

Lemma 5.8. Invertible matrix M is mixing if and only if there exists j such that the support of the first j columns has size greater than j.

PROOF. We need to prove that if we let $S_j = \{i \in [k] | \text{ exists } j' \in [j] \text{ s.t. } M_{i,j'} \neq 0\}$, then there exists a j s.t. $|S_j| > j$. To see this, note that $|S_j|$ is invariant under permutation of the rows, and for upper triangular matrices $|S_j| \leq j$. So if M is not mixing, then for all j we have $|S_j| \leq j$. Conversely, if for every j we have $|S_j| \leq j$, then either we have $|S_j| < j$ for some j, and in which case M is not invertible, or $|S_j| = j$ for every j, in which case we can find a permutation $\pi : [k] \to [k]$ such that for every j $S_j = \{\pi(1), \ldots, \pi(j)\}$. Permuting the rows so that $\pi(j)$ is the jth row makes M upper triangular, and so once again we get M is not mixing.

We will now state the linear-algebraic properties of a mixing matrices that correspond directly to the entropic inequalities in Lemma 5.6 and Lemma 5.7. Specifically, it will not be too difficult to deduce Lemma 5.6 from Lemma 5.9 as we discussed before—the crux of the argument is that $\overline{H}((UM)_j \mid (UM)_{< j}, W) = \overline{H}((UM)_j + B \mid (UM)_{< j}, W)$, where B is some linear combination of variables $(UM)_1, \ldots (UM)_{j-1}$, and Lemma 5.9 describes how to find suitable B. Lemma 5.10 plays the same role in the proof of Lemma 5.7.

LEMMA 5.9. Let M be a mixing matrix, and let $a_1, \ldots a_k \in \mathbb{F}_q^k$ denote columns of M. Then there exists index j and a vector $\mathbf{v} \in a_j + \operatorname{span}\{a_1, \ldots a_{j-1}\}$, such that $|\operatorname{supp}(\mathbf{v})| \geq 2$ and $\operatorname{supp}(\mathbf{v}) \cap \operatorname{supp}(a_i) = \emptyset$ for i < j, where $\operatorname{supp}(\mathbf{v})$ is a set of non-zero coordinates of \mathbf{v} .

PROOF. Let $S_i = \bigcup_{t \le i} \operatorname{supp}(\boldsymbol{a}_t)$. By Lemma 5.8, this means that there exist a j such that $|S_j| > j$. Consider the smallest j satisfying $|S_j| > j$. By a straightforward inductive argument for any k < j, we have $\operatorname{span}\{\boldsymbol{a}_1, \dots \boldsymbol{a}_k\} = \operatorname{span}\{\boldsymbol{e}_\ell : \ell \in S_k\}$, where \boldsymbol{e}_i are the standard basis vectors. Now, we can decompose $\boldsymbol{a}_j = \boldsymbol{v} + \boldsymbol{w}$, where $\operatorname{supp}(\boldsymbol{w}) \subseteq S_{j-1}$ and $\operatorname{supp}(\boldsymbol{v}) \cap S_{j-1} = \emptyset$. Since $|S_{j-1}| = j-1$ and $|S_j| > j$, we have $|\operatorname{supp}(\boldsymbol{v})| \ge 2$, and by construction $\boldsymbol{w} \in \operatorname{span}\{\boldsymbol{e}_\ell : \ell \in S_{j-1}\} = \operatorname{span}\{\boldsymbol{a}_1, \dots \boldsymbol{a}_{j-1}\}$.

LEMMA 5.10. Let M be a mixing matrix, and let $a_1, \ldots a_k \in \mathbb{F}_q^k$ denote columns of M. Then, there exists three indices $j, \ell, s \in [k]$ and $\alpha_1, \alpha_2 \in \mathbb{F}_q^*$, such that $\alpha_1 \mathbf{e}_\ell \in \mathbf{a}_j + \operatorname{span}\{\mathbf{a}_1, \ldots \mathbf{a}_{j-1}\}$ and $\mathbf{e}_\ell + \alpha_2 \mathbf{e}_s \in \operatorname{span}\{\mathbf{a}_1, \ldots \mathbf{a}_{j-1}\}$, where $\mathbf{e}_i \in \mathbb{F}_q^k$ are the standard basis vectors.

In our proof of this lemma, we will use a process that is essentially the well-known Gaussian elimination applied to a matrix M. Specifically, the following proposition captures the properties of intermediate matrices in a Gaussian elimination process useful for our argument.

PROPOSITION 5.11. For any $k \times k$ invertible matrix M, there is a permutation $\pi : [k] \to [k]$ and a sequence of matrices $M^{(1)}, \dots M^{(k)}$ (we call matrix $M^{(j)}$ the jth step matrix) with the following properties.

For any j, if we use $a_1, \ldots a_k$ to denote columns of $M, a'_1, \ldots a'_k$ to denote columns of $M^{(j)}$, and $e_1, \ldots e_k$ to denote standard basis vectors, then we have

11:34 J. Błasiok et al.

- (1) For every $s \in [k]$, we have $a'_s \in a_s + \operatorname{span}\{a_1, \ldots a_j\}$.
- (2) For every $s \in [k]$ and every $\ell \leq j$, we have $\langle a'_s, e_{\pi(\ell)} \rangle = \begin{cases} 1 & \text{if } \ell = s \\ 0 & \text{otherwise} \end{cases}$.
- (3) $\operatorname{span}\{a_1, \ldots, a_j\} = \operatorname{span}\{a'_1, \ldots a'_j\}.$

For example, if j = 3, then $M^{(3)}$ up to some row permutation π must have the form:

$$M^{(3)} = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \star & \star & \star & \star & \dots \\ \star & \star & \star & \star & \dots \end{bmatrix}, \tag{23}$$

where each column of $M^{(3)}$ is a corresponding column of M shifted by some linear combination of the first three columns of M.

PROOF OF PROPOSITION 5.11. The proof proceeds by induction. For the base case we consider $M^{(0)} = M$, and it is easy to verify the properties for M.

Let $j \ge 1$. For the inductive hypothesis, we assume a matrix $M^{(j-1)}$ satisfying properties above and a one-to-one map $\pi: [j-1] \to [k]$. For the inductive step, we want to find $M^{(j)}$ and $\pi(j)$ as in the statement of this proposition. Note that at the end of the induction, when j = k the one-to-one map π is also onto, and hence $\pi: [k] \to [k]$ is a permutation as needed.

Let $a_1^{(j-1)}, \ldots a_k^{(j-1)}$ denote columns of $M^{(j-1)}$. Since M is invertible, we have $a_j \notin \text{span}\{a_1, \ldots a_{j-1}\}$. Using properties 1 and 3 for $M^{(j-1)}$, we conclude that $a_j^{(j-1)} \notin \text{span}\{a_1^{(j-1)}, \ldots a_{j-1}^{(j-1)}\}$. In particular, this implies that $a_j^{(j-1)} \neq 0$. Let $\pi(j)$ be such that $\langle a_j^{(j-1)}, e_{\pi(j)} \rangle \neq 0$. Note that $\pi(j) \neq \pi(s)$ for any s < j by property 2 of the matrix $M^{(j-1)}$, and therefore π is a one-to-one mapping.

For
$$s \neq j$$
, let us take $a'_s := a_s^{(j-1)} - \frac{\langle a_s^{(j-1)}, e_{\pi(j)} \rangle}{\langle a_j^{(j-1)}, e_{\pi(j)} \rangle} a_j^{(j-1)}$, and, finally, $a'_j := \frac{1}{\langle a_j^{(j-1)}, e_{\pi(j)} \rangle} a_j^{(j-1)}$.

Next, we verify that properties 1–3 hold for matrix $M^{(j)}$ given by the columns $a'_1, \ldots a'_k$ defined above.

Indeed, the first property holds, since for any s, we have $a'_s = a_s^{(j-1)} + \gamma a_j^{(j-1)}$, where γ is some scalar. By induction, we have $a_s^{(j-1)} \in a_s + \operatorname{span}\{a_1, \ldots a_{j-1}\}$, and $a_j^{(j-1)} \in \operatorname{span}\{a_1, \ldots, a_j\}$; therefore $a'_s \in a_s + \operatorname{span}\{a_1, \ldots, a_j\}$.

To show the second property, for any $s \neq j$ we have

$$\langle \boldsymbol{a}_s', \boldsymbol{e}_{\pi(\ell)} \rangle = \langle \boldsymbol{a}_s^{(j-1)}, \boldsymbol{e}_{\pi(\ell)} \rangle - \frac{\langle \boldsymbol{a}_s^{(j-1)}, \boldsymbol{e}_{\pi(j)} \rangle}{\langle \boldsymbol{a}_j^{(j-1)}, \boldsymbol{e}_{\pi(j)} \rangle} \langle \boldsymbol{a}_j^{(j-1)}, \boldsymbol{e}_{\pi(\ell)} \rangle.$$

If $\ell < j$, then by induction we have $\langle \boldsymbol{a}_j^{(j-1)}, \boldsymbol{e}_{\pi(\ell)} \rangle = 0$, and hence the second term vanishes, and we have $\langle \boldsymbol{a}_s', \boldsymbol{e}_{\pi(\ell)} \rangle = \langle \boldsymbol{a}_s^{(j-1)}, \boldsymbol{e}_{\pi(\ell)} \rangle$, which again by induction is 1 if $s = \ell$ and 0 otherwise. For $\ell = j$, we have $\langle \boldsymbol{a}_s', \boldsymbol{e}_{\pi(\ell)} \rangle = \langle \boldsymbol{a}_s^{(j-1)}, \boldsymbol{e}_{\pi(\ell)} \rangle - \langle \boldsymbol{a}_s^{(j-1)}, \boldsymbol{e}_{\pi(j)} \rangle = 0$. Further, when s = j, we have $\langle \boldsymbol{a}_j', \boldsymbol{e}_{\pi(\ell)} \rangle = \frac{\langle \boldsymbol{a}_j^{(j-1)}, \boldsymbol{e}_{\pi(\ell)} \rangle}{\langle \boldsymbol{a}_j^{(j-1)}, \boldsymbol{e}_{\pi(j)} \rangle}$ is 1 exactly when $\ell = j$ and 0 when $\ell < j$ (where the latter claim follows from property 2 for $M^{(j-1)}$.

Finally, for the third property, the inclusion $\operatorname{span}\{a_1',\ldots a_j'\}\subset \operatorname{span}\{a_1,\ldots a_j\}$ follows from the property 1, and the dim $\operatorname{span}\{a_1',\ldots a_j'\}=j$ by property 2, which implies $\operatorname{span}\{a_1',\ldots a_j'\}=\operatorname{span}\{a_1,\ldots a_j\}$.

PROOF OF LEMMA 5.10. By Lemma 5.8 a matrix M is mixing if for some index i the support of the first i columns has size strictly greater than i. Let j-1 be the largest index with this property, and note that $j \le k$ (since all the k columns trivially have support size of k).

Let $M^{(j-1)}$ be the (j-1)th step matrix of M defined as in Proposition 5.11.

By definition of j, the span of the first j columns of M must exactly equal span{ $e_{\pi(1)}, \ldots, e_{\pi(j)}$ }, since the total support of all those columns has size exactly j, the columns are linearly independent, and each of $\pi(1), \ldots \pi(j)$ is in this support by property 2 and 3 of matrix $M^{(j)}$.

Thus, all of the first j columns of $M^{(j-1)}$ can only be supported on coordinates $\{\pi(1), \ldots, \pi(j)\}$. Further, by the second property in Proposition 5.11 of the (j-1)th step matrix, the jth column of $M^{(j-1)}$ has zero on all coordinates $\pi(s)$ for s < j. Thus, it must be of form $\alpha_1 e_{\pi(j)}$ for some scalar $\alpha_1 \neq 0$ (since otherwise the jth column of $M^{(j-1)}$ would be $\mathbf{0}$, which would contradict the fact that M is invertible).

Finally, because total support of the first (j-1) columns is larger than j-1, there must exist some column $\ell < j$ of $M^{(j-1)}$ that is supported on the coordinate $\pi(j)$. This, along with the second property of $M^{(j-1)}$ in Proposition 5.11 implies that the ℓ th column of $M^{(j-1)}$ must be exactly ($\boldsymbol{e}_{\pi(\ell)} + \beta_2 \boldsymbol{e}_{\pi(j)}$) for some $\beta_2 \in \mathbb{F}_a^*$.

We can now conclude the statement of the lemma. We have shown that the jth column of $M^{(j-1)}$ is of form $\alpha_1 \boldsymbol{e}_{\pi(j)}$, and by the first property of $M^{(j-1)}$ in Proposition 5.11 it is contained in \boldsymbol{a}_j + span{ $\boldsymbol{a}_1,\ldots,\boldsymbol{a}_{j-1}$ }. This proves the first part of the statement of the lemma (by using $\ell \leftarrow \pi(j)$). The argument for the second part is as follows. We have shown that, on one hand, the ℓ th column of $M^{(j-1)}$ is of form $\boldsymbol{e}_{\pi(\ell)} + \beta_2 \boldsymbol{e}_{\pi(j)}$, on the other hand, it is contained in the span of first j-1 columns of the matrix M by the first property of $M^{(j-1)}$ in Proposition 5.11. In other words, we have $\beta_2^{-1} \cdot \boldsymbol{e}_{\pi(\ell)} + \boldsymbol{e}_{\pi(j)}$ is in the span of first j-1 columns of the matrix M. This shows the second part of the statement of the lemma (by picking $\alpha_2 \leftarrow \beta_2^{-1}$ and $s \leftarrow \pi(\ell)$ and recalling that in the first part we have already set $\ell \leftarrow \pi(j)$).

With Lemmas 5.9 and 5.10 in hand, we are well equipped to show Lemmas 5.6 and 5.7 accordingly.

PROOF OF LEMMA 5.6. Let $a_1, \ldots a_k$ be columns of matrix M. By Lemma 5.9, there is an index j, and a vector $\mathbf{v} = a_j + \mathbf{w}$, where $\mathbf{w} \in \text{span}\{a_1, \ldots a_{j-1}\}$ such that $\text{supp}(\mathbf{v}) \cap \text{supp}(a_i) = \emptyset$ for each i < j.

This implies

$$\overline{H}((UM)_{j} | (UM)_{< j}, W) = \overline{H}(\langle U, a_{j} \rangle | \langle U, a_{1} \rangle, \dots, \langle U, a_{j-1} \rangle, W)
= \overline{H}(\langle U, a_{j} \rangle + \langle U, w \rangle | \langle U, a_{1} \rangle, \dots, \langle U, a_{j-1} \rangle, W)
(Since $w \in \text{span}\{a_{1}, \dots, a_{j-1}\})
= \overline{H}(\langle U, v \rangle | W),$$$

where the last equality follows from the fact that $\langle U, v \rangle$ is independent from $\langle U, a_1 \rangle, \ldots \langle U, a_{j-1} \rangle$ conditioned on W (since v has disjoint support with all a_i for i < j).

11:36 J. Błasiok et al.

Now, since $|\operatorname{supp}(\boldsymbol{v})| > 2$, let us say that $\boldsymbol{v} = \alpha_{\ell} \boldsymbol{e}_{\ell} + \alpha_{2} \boldsymbol{e}_{s} + \boldsymbol{r}$, where $\operatorname{supp}(\boldsymbol{r}) \cap \{\ell, s\} = \emptyset$. We have

$$\begin{split} \overline{H}(\langle \boldsymbol{U}, \boldsymbol{v} \rangle \,|\, \boldsymbol{W}) &= \overline{H}(\alpha_1 \boldsymbol{U}_\ell + \alpha_2 \boldsymbol{U}_s + \langle \boldsymbol{U}, \boldsymbol{r} \rangle \,|\, \boldsymbol{W}) \\ &\geq \overline{H}(\alpha_1 \boldsymbol{U}_\ell + \alpha_2 \boldsymbol{U}_s + \langle \boldsymbol{U}, \boldsymbol{r} \rangle \,|\, \langle \boldsymbol{U}, \boldsymbol{r} \rangle, \boldsymbol{W}) \\ &\qquad \qquad \qquad \text{(Since conditioning does not increase entropy)} \\ &= \overline{H}(\boldsymbol{U}_\ell + \alpha_1^{-1} \alpha_2 \boldsymbol{U}_s \,|\, \langle \boldsymbol{U}, \boldsymbol{r} \rangle, \boldsymbol{W}) \\ &= \overline{H}(\boldsymbol{U}_\ell + \alpha_1^{-1} \alpha_2 \boldsymbol{U}_s \,|\, \langle \boldsymbol{U}, \boldsymbol{r} \rangle, \boldsymbol{W}) \\ &= \overline{H}(\boldsymbol{U}_\ell + \alpha_1^{-1} \alpha_2 \boldsymbol{U}_s \,|\, \boldsymbol{W}), \end{split}$$

where the last equality follows, since $\operatorname{supp}(\mathbf{r}) \cap \{\ell, s\} = \emptyset$. The proof is complete by setting $\alpha = \alpha_1^{-1}\alpha_2$.

PROOF OF LEMMA 5.7. Let $a_1, \ldots a_k$ be columns of matrix M. By Lemma 5.10, there are indices j, ℓ, s and $\alpha_1, \alpha_2 \in \mathbb{F}_q^*$, such that $\alpha_1 \cdot \boldsymbol{e}_{\ell} = \boldsymbol{a}_j + \boldsymbol{w}$, where $\boldsymbol{w} \in \text{span}\{\boldsymbol{a}_1, \ldots \boldsymbol{a}_{j-1}\}$, and $\boldsymbol{e}_{\ell} + \alpha_2 \boldsymbol{e}_s \in \text{span}\{\boldsymbol{a}_1, \ldots \boldsymbol{a}_{j-1}\}$.

This implies

$$\overline{H}((UM)_{j} \mid (UM)_{< j}, W) = \overline{H}(\langle U, a_{j} \rangle \mid \langle U, a_{1} \rangle, \dots, \langle U, a_{j-1} \rangle, W)$$

$$= \overline{H}(\langle U, \alpha_{1} e_{\ell} \rangle - \langle U, w \rangle \mid \langle U, a_{1} \rangle, \dots, \langle U, a_{j-1} \rangle, \langle U, e_{\ell} + \alpha_{2} e_{s} \rangle, W)$$

$$(Since e_{\ell} + \alpha_{2} e_{s} \text{ is in span}\{a_{1}, \dots a_{j-1}\})$$

$$= \overline{H}(\langle U, \alpha_{1} e_{\ell} \rangle \mid \langle U, a_{1} \rangle, \dots, \langle U, a_{j-1} \rangle, \langle U, e_{\ell} + \alpha_{2} e_{s} \rangle, W)$$

$$(Since w \text{ is in span}\{a_{1}, \dots a_{j-1}\})$$

$$\leq \overline{H}(\langle U, \alpha_{1} e_{\ell} \rangle \mid \langle U, e_{\ell} + \alpha_{2} e_{s} \rangle, W)$$

$$(Additional conditioning decreases entropy.)$$

$$= \overline{H}(\alpha_{1} U_{\ell} \mid U_{\ell} + \alpha_{2} U_{s}, W)$$

$$= \overline{H}(U_{\ell} \mid U_{\ell} + \alpha_{2} U_{s}, W),$$

where the last equality follows, since $\alpha_1 \in \mathbb{F}_q^*$ and hence the map $x \mapsto \alpha_1 \cdot x$ is a bijection. \square

We are now ready to prove that every mixing matrix is a polarizing matrix.

Proof of Lemma 5.5. The proof follows easily by combining Lemmas 5.1 to 5.3, 5.6, and 5.7 as we elaborate on below.

Let $M \in \mathbb{F}_q^{k \times k}$ be a mixing matrix, and let (U, W) be random variables such that $U = (U_1, \ldots, U_k) \in \mathbb{F}_q^k$, $W = (W_1, \ldots, W_k)$ is supported on some finite set, and the pairs (U_i, W_i) are independently and identically distributed for $i \in [k]$. Further, let the vector $V := U \cdot M$.

For the *Variance in the middle* condition, we need to show that there exists $j \in [k]$ such that for every $\tau > 0$ there exists $\varepsilon = \varepsilon(\tau) > 0$ such that if $\overline{H}(U_1|W) \in (\tau, 1 - \tau)$, then $\overline{H}((V)_j|V_{< j}, W) \ge \overline{H}(U_1|W) + \varepsilon$. By Lemma 5.6, we have that there exist $j, \ell, s \in [k]$ and $\alpha \in \mathbb{F}_q^*$ such that

$$\overline{H}((V)_{j}|V_{< j}, W) \ge \overline{H}(U_{\ell} + \alpha U_{s}|W) = \overline{H}(U_{\ell} + \alpha U_{s}|W_{\ell}, W_{s}), \tag{24}$$

where the equality above uses the fact that the (U_i, W_i) pairs are independent. By Lemma 5.3 applied with $X_1 = U_\ell$, $A_1 = W_\ell$, $X_2 = \alpha U_s$, and $A_2 = W_s$, we conclude that for every $\tau > 0$ there exists $\varepsilon > 0$ such that

$$\overline{H}(U_{\ell} + \alpha U_{s}|W_{\ell}, W_{s}) \geq \max\{\overline{H}(U_{\ell}|W_{\ell}), \overline{H}(\alpha U_{s}|W_{s})\} + \varepsilon = \overline{H}(U_{1}|W_{1}) + \varepsilon = \overline{H}(U_{1}|W) + \varepsilon,$$

Journal of the ACM, Vol. 69, No. 2, Article 11. Publication date: March 2022.

where the first equality uses the fact that the map $\alpha U_s \mapsto U_s$ is invertible and that the (U_i, W_i) pairs are identically distributed and the final equality uses the fact that these pairs are independent. The Variance in the middle condition follows by combining the two steps above.

For suction at the high end, we need to show there is some index $j \in [k]$ such that for every $c < \infty$, there exists $\tau > 0$ such that if $\overline{H}(U_1|W) > 1 - \tau$, then $1 - \overline{H}(V_j|V_{<j},W) \le \frac{1}{c}(1 - \overline{H}(U_1|W))$. Here again by Lemma 5.6, we have that there exist $j, \ell, s \in [k]$ and $\alpha \in \mathbb{F}_q^*$ such that Equation (24) holds. Now applying Lemma 5.1 to X_1, A_1, X_2, A_2 as in the previous paragraph and $\gamma = 1/c$, we get that there exists $\tau > 0$ such that the requirement for suction at the higher end holds.

Finally, for suction at the lower end we need to show there is some index $j \in [k]$, such that for every $c < \infty$ there exists $\tau > 0$ such that if $\overline{H}(U_1|W) < \tau$, then $\overline{H}(V_j|V_{< j},W) \le \frac{1}{c}\overline{H}(U_1|W)$. We first apply Lemma 5.7 to get that there exist $j, \ell, s \in [k]$ and $\alpha \in \mathbb{F}_q^*$ such that

$$\overline{H}((V)_{i}|V_{< i}, W) \leq \overline{H}(U_{\ell}|U_{\ell} + \alpha U_{s}, W) = \overline{H}(U_{\ell}|U_{\ell} + \alpha U_{s}, W_{\ell}, W_{s}).$$

Now applying Lemma 5.2 with X_1, A_1, X_2, A_2 , and γ as in the previous paragraph, we get that

$$\overline{H}(U_{\ell}|U_{\ell}+\alpha U_{s},W_{\ell},W_{s})\leq \frac{1}{c}\overline{H}(U_{\ell}|W_{\ell})=\frac{1}{c}\overline{H}(U_{1}|W_{1})=\frac{1}{c}\overline{H}(U_{1}|W).$$

This concludes our analysis of the reductions.

5.5 Proof of Theorem 1.15

For completeness and easy reference, we restate Theorem 1.15 below and include its proof.

Theorem 5.12. Local Polarization of Arikan Martingales For every prime q, for every mixing matrix $M \in \mathbb{F}_q^{k \times k}$, and for every symmetric memoryless channel $C_{Y|Z}$ over \mathbb{F}_q , the associated Arikan martingale is locally polarizing.

PROOF OF THEOREM 1.15. Let $M \in \mathbb{F}_q^{k \times k}$ be a mixing matrix. By Lemma 5.5, we have that M satisfies matrix polarization. Now, by Theorem 4.4, we have that for every symmetric memoryless channel $C_{Y|Z}$ over \mathbb{F}_q , the Arikan martingale associated with M and $C_{Y|Z}$ is locally polarizing. \square

6 PROOFS OF ENTROPIC LEMMAS

We now turn to the entropic lemmas stated and used in Section 5.1.

6.1 Suction at the Upper End

To establish Lemma 5.1, we will first show similar kind of statement for unconditional entropies. To this end, we first show that for random variables taking values in *small* set, having entropy close to maximal is essentially the same as being close to uniform with respect to L_2 distance. The L_2 distance of a probability distribution to uniform is controlled by the sum of squares of non-trivial Fourier coefficients of the distribution, and all the non-trivial Fourier coefficients are significantly reduced after adding two independent variables close to the uniform distribution.

Finally, a simple averaging argument is sufficient to lift this result to conditional entropies, establishing Lemma 5.1.

Lemma 6.1. If $X \in \mathbb{F}_q$ is a random variable with a distribution \mathcal{D}_X , then

$$d_2(\mathcal{D}_X, U)^2 \frac{1}{2\log q} \le 1 - \overline{H}(X) \le d_2(\mathcal{D}_X, U)^2 O(q^2),$$

where U is a uniform distribution over \mathbb{F}_q , and $d_p(\mathcal{D}_1, \mathcal{D}_2) := (\sum_{x \in \mathbb{F}_q} (\mathcal{D}_1(x) - \mathcal{D}_2(x))^p)^{1/p}$.

11:38 J. Błasiok et al.

PROOF. Pinskers inequality [26] yields $d_1(\mathcal{D}_X, U) \leq \sqrt{2 \log q} \cdot \sqrt{1 - \overline{H}(X)}$, and by standard relations between ℓ_p norms, we have $d_2(\mathcal{D}_X, U) \leq d_1(\mathcal{D}_X, U)$, which after rearranging yields the bound $d_2(\mathcal{D}_X, U)^2 \leq (2 \log q)(1 - \overline{H}(X))$, which in turn proves the claimed lower bound.

For the upper bound, given $i \in \mathbb{F}_q$ let us take δ_i such that $\mathcal{D}_X(i) \stackrel{\text{def}}{=} \Pr(X = i) = \frac{1+\delta_i}{q}$. Note that this implies (along with the fact that $\sum_{i \in \mathbb{F}_q} \mathcal{D}_X(i) = 1$):

$$\sum_{i \in \mathbb{F}_q} \delta_i = 0 \tag{25}$$

and

$$d_2(\mathcal{D}_X, U)^2 = \frac{1}{q^2} \sum_i \delta_i^2.$$
 (26)

Now

$$1 - \overline{H}(X) = 1 + \frac{1}{\log q} \sum_{i \in \mathbb{F}_q} \frac{(1 + \delta_i)}{q} \log \left(\frac{1 + \delta_i}{q} \right) = \frac{1}{\log q} \sum_{i \in \mathbb{F}_q} \frac{(1 + \delta_i)}{q} \log(1 + \delta_i),$$

where the second equality follows from Equation (25).

By Taylor expansion we have $\log(1 + \delta_i) = \delta_i + \mathcal{E}(\delta_i)$ with some error term $\mathcal{E}(\delta_i)$ such that $|\mathcal{E}(\delta_i)| \le 2\delta_i^2$ for $|\delta_i| < 1$. Therefore, in the case when all $\delta_i < 1$, we have (for some constant *C*):

$$1 - \overline{H}(X) = \frac{1}{q \log q} \sum_{i \in F_q} (1 + \delta_i) (\delta_i + \mathcal{E}(\delta_i))$$

$$\leq \frac{1}{q \log q} \sum_{i \in F_q} [\delta_i + \delta_i^2 + O(\delta_i^2)]$$

$$\leq \frac{1}{q \log q} \left[\sum_{i \in F_q} \delta_i + C \sum_{i \in F_q} \delta_i^2 \right]$$

$$\leq Cq \cdot d_2(\mathcal{D}_X, U)^2,$$

where the inequality follows from Equation (25) and Equation (26). If some $\delta_i \geq 1$, then the inequality is satisfied trivially: $d_2(\mathcal{D}_X, U) \geq \frac{1}{q}$, hence $1 - \overline{H}(X) \leq q^2 \cdot d_2(\mathcal{D}_X, U)^2$.

Lemma 6.2. If $X, Y \in \mathbb{F}_q$ are independent random variables, then $1 - \overline{H}(X + Y) \leq \text{poly}(q)(1 - \overline{H}(X))(1 - \overline{H}(Y))$.

PROOF. By Lemma 6.1, it is enough to show that $d_2(\mathcal{D}_{X+Y}, U)^2 \leq \operatorname{poly}(q) d_2(\mathcal{D}_X, U)^2 d_2(\mathcal{D}_Y, U)^2$. For a distribution \mathcal{D}_X , consider a Fourier transform of this distribution given by $\hat{\mathcal{D}}_X(k) = \mathbb{E}_{j \sim \mathcal{D}_X} \omega^{jk}$, where $\omega = \exp(-2\pi i/q)$. As usual, we have $\hat{\mathcal{D}}_{X+Y}(k) = \hat{\mathcal{D}}_X(k)\hat{\mathcal{D}}_Y(k)$.

Moreover, by Parseval's identity we will show that

$$d_2(\mathcal{D}_X, U)^2 = \frac{1}{q} \sum_{k \neq 0} \hat{\mathcal{D}}_X(k)^2.$$

Indeed, as in the proof of Lemma 6.1, define $\mathcal{D}_X(i) =: \frac{1+\delta_i}{q}$. Then, by Parseval's identity, we have

$$\frac{1}{q} \cdot \sum_{k \in \mathbb{F}_q} \hat{\mathcal{D}}_X(k)^2 = \sum_{i \in \mathbb{F}_q} \frac{(1 + \delta_i)^2}{q^2} = \frac{1}{q^2} \left(\sum_{i \in \mathbb{F}_q} 1 + 2 \sum_{i \in \mathbb{F}_q} \delta_i + \sum_{i \in \mathbb{F}_q} \delta_i^2 \right) \frac{1}{q} + d_2(\mathcal{D}_X, U)^2,$$

Journal of the ACM, Vol. 69, No. 2, Article 11. Publication date: March 2022.

which implies the claimed bound by noting that $\hat{\mathcal{D}}_X(0) = 1$. (In the above, the last equality follows from Equation (25) and Equation (26).)

This yields

$$d_{2}(\mathcal{D}_{X+Y}, U)^{2} = \frac{1}{q} \cdot \sum_{k \neq 0} \hat{\mathcal{D}}_{X}(k)^{2} \hat{\mathcal{D}}_{Y}(k)^{2}$$

$$\frac{1}{q} \cdot \leq \left(\sum_{k \neq 0} \hat{\mathcal{D}}_{X}(k)^{2}\right) \left(\sum_{k \neq 0} \hat{\mathcal{D}}_{Y}(k)^{2}\right) = q d_{2}(\mathcal{D}_{X}, U)^{2} d_{2}(\mathcal{D}_{Y}, U)^{2}.$$

LEMMA 6.3. Let $X_1, X_2 \in \mathbb{F}_q$ be a pair of random variables, and let A_1, A_2 be pair of discrete random variables, such that (X_1, A_1) and (X_2, A_2) are independent. Then

$$1 - \overline{H}(X_1 + X_2 | A_1, A_2) \le (1 - \overline{H}(X_1 | A_1))(1 - \overline{H}(X_2 | A_2)) \operatorname{poly}(q).$$

PROOF. We have

$$1 - \overline{H}(X_{1} + X_{2}|A_{1}, A_{2})$$

$$= \sum_{a_{1}, a_{2}} \Pr(A_{1} = a_{1}) \Pr(A_{2} = a_{2}) (1 - \overline{H}(X_{1} + X_{2}|A_{1} = a_{1}, A_{2} = a_{2}))$$

$$\leq \operatorname{poly}(q) \sum_{a_{1}, a_{2}} \Pr(A_{1} = a_{1}) \Pr(A_{1} = a_{1}) (1 - \overline{H}(X_{1}|A_{1} = a_{1}, A_{2} = a_{2})) (1 - \overline{H}(X_{2}|A_{1} = a_{1}, A_{2} = a_{2}))$$

$$= \operatorname{poly}(q) \sum_{a_{1}, a_{2}} \Pr(A_{1} = a_{1}) (1 - \overline{H}(X_{1}|A_{1} = a_{1})) \Pr(A_{2} = a_{2}) (1 - \overline{H}(X_{2}|A_{2} = a_{2}))$$

$$= \operatorname{poly}(q) \left(\sum_{a_{1}} \Pr(A_{1} = a_{1}) (1 - \overline{H}(X_{1}|A_{1} = a_{1})) \right) \left(\sum_{a_{2}} \Pr(A_{2} = a_{2}) (1 - \overline{H}(X_{2}|A_{2} = a_{2})) \right)$$

$$= \operatorname{poly}(q) (1 - \overline{H}(X_{1}|A_{1})) (1 - \overline{H}(X_{2}|A_{2})),$$

where the inequality follows from Lemma 6.2 and the second equality follows from independence of (X_1, A_1) and (X_2, A_2) .

We are now ready to prove Lemma 5.1.

PROOF OF LEMMA 5.1. Given γ , q, take $\tau = \gamma/P(q)$, where P(q) is the polynomial appearing in the statement of Lemma 6.3. By applying the conclusion of Lemma 6.3, we have

$$1 - \overline{H}(X_1 + X_2 | A_1, A_2) \leq (1 - \overline{H}(X_1 | A_1))(1 - \overline{H}(X_2 | A_2)P(q)$$

$$\leq (1 - \overline{H}(X_1 | A_1))\tau P(q)$$

$$= \gamma(1 - \overline{H}(X_1 | A_1)).$$

6.2 Suction at the Lower End

In this subsection, will show Lemma 5.2. To this end, we want to show that for pairs (X_1, A_1) and (X_2, A_2) with low conditional entropy $\overline{H}(X_1 \mid A_1) < \tau, \overline{H}(X_2 \mid A_2) < \tau$, the entropy of the sum is almost as big as sum of corresponding entropies, i.e., $\overline{H}(X_1 + X_2 \mid A_1, A_2) \ge (1 - \gamma)(\overline{H}(X_1 \mid A_1) + \overline{H}(X_2 \mid A_2))$ —and the statement of Lemma 5.2 will follow (as we show later) by application of chain rule. To this end, we first show the same type of statement for non-conditional entropies, i.e., if $\overline{H}(X_1) < \tau, \overline{H}(X_2) < \tau$, then $\overline{H}(X_1 + X_2) > (1 - \gamma)(\overline{H}(X_1) + \overline{H}(X_2))$; this fact can be deduced by reduction to the analogous fact for binary random variables, where it becomes just a simple computation. Then we proceed by lifting this statement to the corresponding statement about conditional entropies—this requires somewhat more effort than in Lemma 5.1.

11:40 J. Błasiok et al.

LEMMA 6.4. Let X, Y be independent random variables in \mathbb{F}_q . For any $\gamma < 1$, there exists $\alpha = \alpha(\gamma)$ such that if $\overline{H}(X) \leq \alpha$ and $\overline{H}(Y) \leq \alpha$, then

$$\overline{H}(X + Y) \ge (1 - \gamma)(\overline{H}(X) + \overline{H}(Y)).$$

First, we will show some preliminary useful lemmas.

Assumption 6.5. In the following, without loss of generality, let 0 be the most likely value for both random variables X, Y. (This shifting does not affect entropies).

LEMMA 6.6. Let X be a random variable over \mathbb{F}_q , such that 0 is the most-likely value of X. Then, for any q and any $\gamma < 1$, there exists $\alpha_2(q, \gamma) > 0$ such that

$$\overline{H}(X) \leq \alpha_2(q,\gamma) \implies \Pr[X \neq 0] \leq \gamma \overline{H}(X).$$

PROOF. Let $\beta := \Pr[X \neq 0]$, and $\alpha := \overline{H}(X)$. We have

$$\alpha \log q = H(X) \ge H(\overline{\delta}(X)) = H(\beta) \ge \beta \log(1/\beta).$$

In the above, the inequality follows from the fact that applying a deterministic function to a random variable can only decrease its entropy. Thus,

$$\Pr[X \neq 0] = \beta \le \frac{\alpha \log q}{\log(1/\beta)}$$
$$\le \frac{\alpha \log q}{\log(1/\alpha) - \log \log q},$$

where we used the fact that $\beta \le \alpha \log q$ from Lemma 2.1. Hence, as soon as $\log \frac{1}{\alpha} \ge \frac{\log q}{\gamma} + \log \log q$, the statement of the lemma holds.

LEMMA 6.7 (SUCTION-AT-LOWER-END IN THE BINARY CASE). Let U, V be independent binary random variables. There exists a function $\alpha_0(\gamma)$ such that, for any $0 < \gamma < 1$,

$$H(U), H(V) \le \alpha_0(\gamma) \implies H(U \oplus V) \ge (1 - \gamma)(H(U) + H(V))^{14}$$

PROOF. Let p_1 and p_2 be the biases of U,V, respectively, such that $U \sim \text{Bernoulli}(p_1)$ and $V \sim \text{Bernoulli}(p_2)$. Let $p_1 \circ p_2 = p_1(1-p_2) + (1-p_1)p_2$ be the bias of $U \oplus V$, that is $U \oplus V \sim \text{Bernoulli}(p_1 \circ p_2)$. We first describe some useful bounds on H(p). We have $H(p) \geq p \log 1/p$. For $p \leq 1/2$, we also have

$$-(1-p)\log(1-p) \le (1/\ln 2)(1-p)(p+p^2) \le (1/\ln 2)p \le 2p,$$

where the first inequality follows from the fact that $-\ln(1-x) \le x + x^2$ for $x \le \frac{1}{2}$, and so we have $H(p) \le p(2 + \log 1/p)$. Summarizing, we have

$$p\log(1/p) \le H(p) \le p\log(1/p) + 2p.$$

Suppose $H(p_1), H(p_2) \le \tau$. We now consider $H(p_1) + H(p_2) - H(p_1 \circ p_2)$. WLOG assume $p_1 \le p_2$. Note that this implies

$$p_1 \circ p_2 \leq p_1 + p_2 \leq 2p_2$$
.

¹⁴We note that we could have replaced \oplus by just + as those operations are over \mathbb{F}_2 but we chose to keep + for addition over reals in this lemma.

We have

$$\begin{split} &H(p_1) + H(p_2) - H(p_1 \circ p_2) \\ &\leq p_1(\log(1/p_1) + 2) + p_2(\log(1/p_2) + 2) - (p_1 \circ p_2) \log(1/(p_1 \circ p_2)) \\ &\leq p_1(\log(1/p_1) + 2) + p_2(\log(1/p_2) + 2) - (p_1 + p_2 - 2p_1p_2) \log(1/(2p_2)) \\ &= p_1 \log(2p_2/p_1) + p_2 \log(2p_2/p_2) + 2p_1p_2 \log(1/(2p_2)) + 2(p_1 + p_2) \\ &\leq p_1 \log(p_2/p_1) + 2p_1p_2 \log(1/(p_2)) + 6p_2 \\ &\leq 2p_1H(p_2) + 7p_2 \qquad \qquad \text{(Using } p_1 \log(p_2/p_1) \leq p_2) \\ &\leq 2p_1H(p_2) + 7H(p_2)/\log(1/p_2) \\ &\leq 9H(p_2)/\log(1/\tau). \end{split}$$

In the above, the last inequality follows from the assumption that $\tau \leq 1/8$ (which will be true in our case). Indeed, note that with this assumption $\tau \log(1/\tau) \leq 1$ (which along with the fact that $p_1 \leq \tau$ implies $p_1 \leq 1/\log(1/\tau)$) and $p_2 \leq \tau$ (since we have $p_2 \log(1/p_2) \leq \tau$). Thus, we have

$$H(U), H(V) \le \tau \implies H(U) + H(V) - H(U \oplus V) \le 9H(V)/\log(1/\tau).$$

This implies the desired statement, for $\alpha_0(\gamma) := 2^{-9/\gamma}$.

Let $\overline{\delta}: \mathbb{F}_q \to \{0,1\}$ be the complemented Kronecker delta function, $\overline{\delta}(x) := \mathbb{I}\{x \neq 0\}$. We show that for small-enough entropies, the entropy $H(\overline{\delta}(X))$ is comparable to H(X).

Lemma 6.8. There exists a function $\alpha_1(\gamma)$ such that for any given $0 < \gamma < 1$, and any arbitrary random variable $X \in \mathbb{F}_q$ such that 0 is the most likely value of X,

$$\overline{H}(X) \leq \alpha_1(\gamma) \implies \overline{H}(X) \geq \frac{1}{\log q} H(\overline{\delta}(X)) \geq (1-\gamma)\overline{H}(X).$$

PROOF. The first inequality $\overline{H}(X)\log q=H(X)\geq H(\overline{\delta}(X))$ always holds, by the fact that deterministic postprocessing does not increase entropy. Thus, we will now show the second bound: For small-enough entropies, $\frac{1}{\log q}H(\overline{\delta}(X))\geq (1-\gamma)\overline{H}(X)$. This is equivalent with showing that $H(\overline{\delta}(X))\geq (1-\gamma)H(X)$. Given γ , let $\alpha_1:=\alpha_2(q,\gamma)$ be the entropy guaranteed by Lemma 6.6, so that if $\overline{H}(X)\leq \alpha_1$, then $\Pr[\overline{\delta}(X)=1]=\Pr[X\neq 0]\leq \gamma\overline{H}(X)$. Now, for $H(X)\leq \alpha_1$, we have

$$H(X) = H(X, \overline{\delta}(X)) - H(\overline{\delta}(X)|X), \qquad (Chain rule)$$

$$= H(X, \overline{\delta}(X)), \qquad (as \overline{\delta}(X)|X \text{ is deterministic})$$

$$= H(\overline{\delta}(X)) + H(X|\overline{\delta}(X)), \qquad (Chain rule)$$

$$= H(\overline{\delta}(X)) + H(X|\overline{\delta}(X) = 1) \Pr[\overline{\delta}(X) = 1], \qquad (as H(X|\overline{\delta}(X) = 0) = 0, \text{ since } X|\overline{\delta}(X) = 0 \text{ is deterministic})$$

$$\leq H(\overline{\delta}(X)) + \log(q) \Pr[\overline{\delta}(X) = 1], \qquad (as X \in \mathbb{F}_q, \text{ so } H(X|\overline{\delta}(X) = 1) \leq H(X) \leq \log(q))$$

$$\leq H(\overline{\delta}(X)) + \log(q)\gamma \overline{H}(X), \qquad (by \text{ Lemma 6.6})$$

$$= H(\overline{\delta}(X)) + \gamma H(X).$$

Thus, if $H(X) \le \alpha_1$, then $(1 - \gamma)H(X) \le H(\overline{\delta}(X))$ as desired.

Now, by combining these, we can reduce suction-at-the-lower-end from \mathbb{F}_q to the binary case.

11:42 J. Błasiok et al.

PROOF OF LEMMA 6.4. Given y, we will set $\alpha \le 1/4$, to be determined later. Notice that we have

$$\overline{H}(X+Y) = \frac{1}{\log q}H(X+Y) \ge \frac{1}{\log q}H(\overline{\delta}(X+Y)). \tag{27}$$

We will proceed to show first that

$$H(\overline{\delta}(X+Y)) \ge H(\overline{\delta}(X) \oplus \overline{\delta}(Y)).$$
 (28)

This inequality is justified by comparing the distributions of $\overline{\delta}(X+Y)$ and $\overline{\delta}(X) \oplus \overline{\delta}(Y)$, both binary random variables, and noticing that

$$\Pr[\overline{\delta}(X+Y)=0] = \Pr[X+Y=0] \leq \Pr[\{X=0,Y=0\} \cup \{X\neq 0,Y\neq 0\}] = \Pr[\overline{\delta}(X)\oplus \overline{\delta}(Y)=0].$$

Moreover, let us observe that $\Pr[\overline{\delta}(X+Y)=0]=\Pr[X+Y=0]\geq 1/2$. Indeed,

$$\Pr[X + Y \neq 0] \leq H(X + Y) \leq H(X, Y) \leq H(X) + H(Y) \leq 2\alpha \leq 1/2.$$

In the above, the second inequality follows, since X+Y is a deterministic function of X,Y and the third inequality follows from the chain rule and the fact that conditioning can only decrease entropy. Therefore, by monotonicity of the binary entropy function H(p) for $1/2 \le p \le 1$, and since $\Pr[\overline{\delta}(X+Y)=0] \le \Pr[\overline{\delta}(X) \oplus \overline{\delta}(Y)=0]$ we have

$$H(\overline{\delta}(X+Y)) \ge H(\overline{\delta}(X) \oplus \overline{\delta}(Y)).$$

This justifies Equation (28).

Now we conclude by using the suction-lemma in the binary case, applied to $\overline{\delta}(X) \oplus \overline{\delta}(Y)$.

Let γ' be a small-enough constant, such that $(1 - \gamma')^2 \ge (1 - \gamma)$. Let $\alpha_0 := \alpha_0(\gamma')$ be the entropy bound provided by Lemma 6.7, and let $\alpha_1 := \alpha_1(\gamma')$ be the entropy bound provided by Lemma 6.8. Set $\alpha := \min\{\alpha_0, \alpha_1, 1/4\}$.

Then, for $\overline{H}(X)$, $\overline{H}(Y) \leq \alpha$, we have

$$\overline{H}(X+Y)\log q \geq H(\overline{\delta}(X+Y)), \qquad \text{(Equation (27))}$$

$$\geq H(\overline{\delta}(X) \oplus \overline{\delta}(Y)), \qquad \text{(Equation (28))}$$

$$\geq (1-\gamma')(H(\overline{\delta}(X)) + H(\overline{\delta}(Y))), \qquad \text{(Lemma 6.7 and } \overline{H}(\overline{\delta}(Z)) \leq \overline{H}(Z) \text{ for r.v. } Z)$$

$$\geq (1-\gamma')^2(\overline{H}(X) + \overline{H}(Y))\log q. \qquad \text{(Lemma 6.8)}$$

With our setting of γ' , this concludes the proof.

We will now see how Lemma 6.4 implies its strengthening for conditional entropies.

Lemma 6.9. Let (X_1, A_1) and (X_2, A_2) be independent random variables with $X_i \in \mathbb{F}_q$, and such that A_1, A_2 are identically distributed; and, moreover, for every a we have $\overline{H}(X_1|A_1=a)=\overline{H}(X_2|A_2=a)$. Then, for every $\gamma>0$, there exist τ such that if $\overline{H}(X_1|A_1)\leq \tau$, then

$$\overline{H}(X_1 + X_2 | A_1, A_2) \ge (1 - \gamma)(\overline{H}(X_1 | A_1) + \overline{H}(X_2 | A_2)). \tag{29}$$

PROOF. Let us take $\alpha:=\overline{H}(X_1|A_1)=\overline{H}(X_2|A_2)$. For given γ , we shall find τ such that if $\alpha<\tau$, then Equation (29) is satisfied. Let us now consider $G_A:=\{a:\overline{H}(X_1|A_1=a)<\alpha_1\}$, for $\alpha_1=\frac{\alpha}{\gamma}$. (In the remainder of the proof when we want to talk about a random variable from the identical distribution from which A_1 and A_2 are drawn, we will denote it by A.) By Markov inequality

$$\Pr(A \notin G_A) \leq \frac{\alpha}{\alpha_1} = \gamma.$$

Let us now fix τ , which appears in the statement of this lemma to be smaller than γ and, moreover, small enough, so that when $\alpha < \tau$ for every $a_1, a_2 \in G_A$ we can apply Lemma 6.4 to

distributions $(X_1|A_1 = a_1)$ and $(X_2|A_2 = a_2)$ to ensure that $H(X_1 + X_2|A_1 = a_1, A_2 = a_2) \ge (1 - \gamma)(H(X_1|A_1 = a_1) + H(X_2|A_2 = a_2))$.

Let us use shorthand $S(a_1, a_2) = \overline{H}(X_1 + X_2 | A_1 = a_1, A_2 = a_2) \Pr(A_1 = a_1, A_2 = a_2)$. We have

$$\overline{H}(X_1 + X_2 | A_1, A_2) = \sum_{\substack{a_1, a_2 \\ a_1 \in G_A \\ a_2 \in G_A}} S(a_1, a_2) + \sum_{\substack{a_1 \notin G_A \\ a_2 \in G_A}} S(a_1, a_2) + \sum_{\substack{a_1 \notin G_A \\ a_2 \notin G_A}} S(a_1, a_2) + \sum_{\substack{a_1 \in G_A \\ a_2 \notin G_A}} S(a_1, a_2).$$
(30)

If both a_1 and a_2 are in G_A , then by Lemma 6.4 we have

$$S(a_1, a_2) \ge (1 - \gamma)(\overline{H}(X_1 | A_1 = a_1) + H(X_2 | A_2 = a_2)) \Pr(A_1 = a_1, A_2 = a_2),$$

and therefore

$$\sum_{a_1 \in G_A, a_2 \in G_A} S(a_1, a_2) \ge 2(1 - \gamma) \Pr(A \in G_A) \sum_{a_1 \in G_A} H(X_1 | A_1 = a_1) \Pr(A_1 = a_1), \tag{31}$$

where in the above we have used the fact that A_1 and A_2 are identically distributed.

However, for $a_1 \notin G_A$, $a_2 \in G_A$ let us bound

$$S(a_1, a_2) = \overline{H}(X_1 + X_2 | A_1 = a_1, A_2 = a_2) \Pr(A_1 = a_1, A_2 = a_2)$$

$$\geq \overline{H}(X_1 + X_2 | A_1 = a_1, A_2 = a_2, X_2) \Pr(A_1 = a_1, A_2 = a_2)$$

$$= \overline{H}(X_1 | A_1 = a_1) \Pr(A_1 = a_1, A_2 = a_2),$$

where the inequality follows from the fact that additional conditioning decreases entropy and for the second equality we used the fact that, since X_1 and X_2 are independent, $\overline{H}(X_1+X_2|A_1=a_1,A_2=a_2,X_2)=\overline{H}(X_1|A_1=a_1,A_2=a_2,X_2)=\overline{H}(X_1|A_1=a_1,A_2=a_2)=\overline{H}(X_1|A_1=a_1)$. Summing this bound over all such pairs yields

$$\sum_{a_1 \notin G_A, a_2 \in G_A} S(a_1, a_2) \ge \Pr(A \in G_A) \sum_{a_1 \notin G_A} \overline{H}(X_1 | A_1 = a_1) \Pr(A_1 = a_1), \tag{32}$$

and symmetrically for the third summand, we get

$$\sum_{a_1 \in G_A, a_2 \notin G_A} S(a_1, a_2) \ge \Pr(A \in G_A) \sum_{a_2 \notin G_A} \overline{H}(X_2 | A_2 = a_2) \Pr(A_2 = a_2).$$
(33)

Plugging Equations (31)–(33) into Equation (30) (and using the fact that A_1 and A_2 are identically distributed), we find

$$\overline{H}(X_1 + X_2 | A_1, A_2) \ge 2(1 - \gamma) \Pr(A_1 \in G_A) \sum_{a_1} \overline{H}(X_1 | A_1 = a_1) \Pr(A_1 = a_1)$$

$$= 2(1 - \gamma) \Pr(A \in G_A) \overline{H}(X_1 | A_1).$$

We have $Pr(A \in G_A) \ge (1 - \gamma)$, which yields

$$\overline{H}(X_1 + X_2 | A_1, A_2) \ge 2(1 - \gamma)^2 \alpha \ge 2(1 - 2\gamma)\alpha$$

and the statement of the lemma follows, after rescaling y by half.

Finally, we are ready to prove Lemma 5.2.

11:44 J. Błasiok et al.

PROOF OF LEMMA 5.2. By chain rule, we have

$$\overline{H}(X_1 \mid X_1 + X_2, A_1, A_2) = \overline{H}(X_1, X_1 + X_2 \mid A_1, A_2) - \overline{H}(X_1 + X_2 \mid A_1, A_2)$$

$$= \overline{H}(X_1, X_2 \mid A_1, A_2) - \overline{H}(X_1 + X_2 \mid A_1, A_2)$$

$$= 2\overline{H}(X_1 \mid A_1) - \overline{H}(X_1 + X_2 \mid A_1, A_2),$$

where the last equality follows from the independence of (X_1, A_1) and (X_2, A_2) . Now we can apply Lemma 6.9 to get

$$\overline{H}(X_1 \mid X_1 + X_2, A_1, A_2) \le 2\overline{H}(X_1 \mid A_1) - (1 - \gamma)(2\overline{H}(X_1 \mid A_1) = 2\gamma\overline{H}(X_1 \mid A_1),$$

and the statement follows directly from Lemma 6.9 and rescaling γ by half.

7 EXPONENTIAL MATRIX POLARIZATION

The main result of this section shows the exponential matrix polarization of $M^{\otimes 2}$ for every mixing matrix.

LEMMA 7.1. For every prime p, every mixing matrix $M \in \mathbb{F}_q^{k \times k}$ and every $\varepsilon > 0$, the matrix $M^{\otimes 2}$ satisfies $(\frac{1}{k^2}, 2 - \varepsilon)$ -exponential matrix polarization.

Before turning to the proof, we first note that this immediately yields Theorem 1.16.

PROOF OF THEOREM 1.16. By Lemma 7.1, we have that for every prime q and mixing matrix $M \in \mathbb{F}_q^{k \times k}$, the matrix $M^{\otimes 2}$ satisfies $(\frac{1}{k^2}, 2 - \varepsilon)$ -exponential matrix polarization. By Theorem 4.4, we then have that for every symmetric memoryless channel $C_{Y|Z}$, the Arikan martingale associated with $M^{\otimes 2}$ and $C_{Y|Z}$ is $(\frac{1}{k^2}, 2 - \varepsilon)$ -exponentially locally polarizing.

The rest of the section is devoted to the proof of Lemma 7.1. We start with a simple proposition.

Proposition 7.2. For every field \mathbb{F}_q and every matrix $M \in \mathbb{F}_q^{k \times k}$, its tensor $M^{\otimes 2}$ is mixing if M is mixing.

PROOF. Let $S_j = \{i \in [k] | \exists j' \in [j] \text{ s.t. } M_{i,j'} \neq 0\}$, and then $\exists j \text{ s.t. } |S_j| > j$. By Lemma 5.8, there exists a j such that $|S_j| > j$. With this observation, the proposition follows easily. Given mixing M, let j be the index such that $|S_j| > j$. Recall that $M^{\otimes 2}$ is composed of k^2 submatrices of dimensions $k \times k$ each, with the i,jth submatrix being $M_{ij} \cdot M$. Let i be an index such that $M_{i1} \neq 0$. (Such an index must exist or else we have an all zero column that contradicts the invertibility of M.) Then the first k columns of $M^{\otimes 2}$ contain the $k \times k$ submatrix $M_{i1} \cdot M$, and in this submatrix itself we have that the support of the first j columns has size larger than j. We conclude the first j columns of $M^{\otimes 2}$ have support size larger than j and so by Lemma 5.8, $M^{\otimes 2}$ is mixing.

7.1 Exponential Polarization of a 2×2 Matrix

We will first prove that a single specific matrix, namely $\begin{pmatrix} 1 & 0 \\ \alpha & 1 \end{pmatrix}$, after taking second Kronecker power, satisfies exponential polarization. Recall that in Section 5.3 the local polarization of a mixing matrix was shown essentially by reducing to this case. We will follow a similar plan in this section.

LEMMA 7.3. Let q be a prime and let $H = \begin{pmatrix} 1 & 0 \\ \alpha & 1 \end{pmatrix}$ for $\alpha \in \mathbb{F}_q^*$. Then, for every $\varepsilon > 0$, the matrix $H^{\otimes 2}$ satisfies $(\frac{1}{4}, 2 - \varepsilon)$ exponential matrix polarization.

PROOF. Note that since H is mixing, by Proposition 7.2, we have that $H^{\otimes 2}$ is also mixing. And so, by Lemma 5.5, we have that $H^{\otimes 2}$ satisfies the conditions of matrix polarization (specifically, variance in the middle and suction at the upper and lower ends from Definition 4.3). It remains only to argue exponential matrix polarization, i.e., strong suction at the ends.

Given $\varepsilon > 0$, let $\tau > 0$ be such that for every $\delta < \tau$ we have $6(\log \frac{1}{3\delta^2} + \log q) \le \delta^{-\varepsilon}$. Note that the identity is satisfied for small-enough δ , since the LHS is $O(\log(\frac{1}{\delta}))$ while the RHS is $O((\frac{1}{\delta})^{\varepsilon})$. Now let $\delta < \tau$, and now consider arbitrary sequence of i.i.d. random variables $(U_1, W_1), \ldots, (U_4, W_4)$ with $H(U_i|W_i) = \delta$. We can explicitly write down matrix $H^{\otimes 2}$ as

$$H^{\otimes 2} = \left[\begin{array}{cccc} 1 & 0 & 0 & 0 \\ \alpha & 1 & 0 & 0 \\ \alpha & 0 & 1 & 0 \\ \alpha^2 & \alpha & \alpha & 1 \end{array} \right].$$

Matrix $H^{\otimes 2}$ has four rows. So, to achieve exponential polarization with $\eta = \frac{1}{4}$, we need to show that there is at least one index i satisfying the strong suction inequality (with parameter $b = 2 - \varepsilon$). We do so for i = 4. Let us consider vector $U = (U_1, \ldots, U_4)$ and similarly $W = (W_1, \ldots, W_4)$, and let $V = (V_1, \ldots, V_4) = U \cdot H^{\otimes 2}$. We want to bound

$$\overline{H}(V_4|V_{<4},W) = \overline{H}(U_4|U_1 + \alpha U_2 + \alpha U_3 + \alpha^2 U_4, U_2 + \alpha U_4, U_3 + \alpha U_4, W)$$

$$= \overline{H}(U_4|U_1 - \alpha^2 U_4, U_2 + \alpha U_4, U_3 + \alpha U_4, W),$$

where the equality follows, since $U_1 - \alpha^2 U_4 = U_1 + \alpha U_2 + \alpha U_3 + \alpha^2 U_4 - \alpha (U_2 + \alpha U_4) - \alpha (U_3 + \alpha U_4)$, and hence the map

$$\left(U_{1}+\alpha U_{2}+\alpha U_{3}+\alpha^{2} U_{4},U_{2}+\alpha U_{4},U_{3}+\alpha U_{4}\right)\mapsto\left(U_{1}-\alpha^{2} U_{4},U_{2}+\alpha U_{4},U_{3}+\alpha U_{4}\right)$$

is a bijection.

The main idea to bound the conditional entropy of U_4 above is that if any of U_i is "known" for $i \in \{1, 2, 3\}$, then given the variables being conditioned on, U_4 is also "known." Of course, none of the U_i 's are known, but each is predictable given W_i , and we use this predictability to bound the conditional entropy. Details follow.

Let Σ denote the domain of W_i 's. Using $\overline{H}(U_i|W_i)=\delta$, by Lemma 2.2, we have that there exists some function $f:\Sigma\to\mathbb{F}_q$, such that $\Pr(f(W_i)\neq U_i)\leq \delta$. Let $V_1':=-\alpha^2U_4+U_1$. We now give a predictor $g(V_1',V_2,V_3,W)$ for U_1 . Let $X_1=-\alpha^{-2}(V_1'-f(W_1)), X_2=\alpha^{-1}(V_2-f(W_2)),$ and $X_3=\alpha^{-1}(V_3-f(W_3)).$ Note that if for some i we have $f(W_i)=U_i$, then we have $X_i=U_4$. Using this we set g as follows: If two of X_1,X_2,X_3 have the same value, then we define $g(V_1',V_2,V_3,W)$ to be this value; otherwise, we set $g(V_1',V_2,V_3,W)$ arbitrarily.

By construction of g, we have that if there exist two choices of $i \in \{1, 2, 3\}$ satisfying $f(W_i) = U_i$, then $g(V_1', V_2, V_3, W) = U_4$. In turn, this implies $\Pr(g(V_1', V_2, V_3, W) \neq U_4) \leq 3\delta^2$, since by symmetry, we have

$$\Pr(g(V_1', V_2, V_3, W) \neq U_4) \leq 3\Pr(f(W_1) \neq U_1 \land f(W_2) \neq U_2) = 3\Pr(f(W_1) \neq U_1))^2 \leq 3\delta^2,$$

where the equality follows, since (U_i, W_i) are independent.

Converting the predictability of U_1 by $g(\cdots)$ into an entropy bound by Fano's inequality Lemma 2.3, we have $\overline{H}(U_4|U_1-\alpha^2U_4,U_2+\alpha U_4,U_3+\alpha U_4,W)\leq 6\delta^2(\log\frac{1}{3\delta^2}+\log q)$. By the choice of τ and $\delta<\tau$, we have $6(\log\frac{1}{3\delta^2}+\log q)\leq \delta^{-\varepsilon}$ and so

$$\overline{H}(V_4|V_{<4},W) \le \delta^{2-\varepsilon} = \left(\overline{H}(U_1|W_1)\right)^{2-\varepsilon},\tag{34}$$

as desired.

11:46 J. Błasiok et al.

7.2 Exponential Polarization of Any Mixing Matrix via Useful Containment

We will now proceed to show that exponential polarization of $M^{\otimes 2}$ for any mixing matrix M can be reduced to the lemma above. We first provide an intuitive explanation of the reasoning below.

To show that a matrix M' satisfies an exponential polarization (or just *suction at the lower end* condition of local polarization), one needs to show that for any i.i.d. variables U_i with entropy $H(U_i) = \delta$ and some index j, we can upper bound $\overline{H}((UM')_j|(UM')_{< j})$ (for the sake of the clarity of this exposition, we skip conditioning on W_i). If we write $V_i = (UM)_i$, then we wish to upper bound $\overline{H}(V_j|V_1,\ldots V_{j-1})$ (where all V_i are linear forms in $\{U_i\}_{i\in[k]}$). Now, for any B_1,\ldots,B_m that all can be expressed as linear combinations of $V_1,\ldots V_{j-1}$, we have

$$\overline{H}(V_j|V_1...V_{j-1}) = \overline{H}(V_j + B_m|V_1,...V_{j-1}, B_1,...B_{m-1}) \le \overline{H}(V_j + B_m|B_1,...,B_{m-1}).$$

In Section 5.4, we showed using Gaussian elimination that for any mixing matrix M, one can find j, ℓ, s , and linear forms W_1, W_2 s.t. $V_j + W_2 = U_\ell$ and $W_1 = \alpha U_\ell + U_s$, which implied $\overline{H}(V_j + W_2 | W_1) = \overline{H}(U_\ell | \alpha U_\ell + U_s)$. This can be thought of as showing that in some sense any mixing matrix M contains a matrix $H = \begin{pmatrix} 1 & 0 \\ \alpha & 1 \end{pmatrix}$ and reduces the problem of showing the local polarization of the former to understanding local polarization of the latter.

Here we introduce a technical notion of *useful containment* that is tailored to extend this reasoning in a way that has a convenient property expressed by Lemma 7.7, i.e., since matrix M contains H in this specific sense, the matrix $M^{\otimes 2}$ contains $H^{\otimes 2}$ and by the reasoning outlined in the previous paragraph, we can deduce exponential local polarization of $M^{\otimes 2}$ from this containment and the entropy upper bound proved in Lemma 7.3.

We wish to note here that the subsequent definition and lemmas are tailored to the specific statement we are proving. In particular, *useful containment* is not a transitive relation. More importantly, and unfortunately, it is not true that for any exponentially polarizing matrix R; if R is usefully contained in M, then M is exponentially polarizing. Lemma 7.8 asserts this property only for $R = H^{\otimes 2}$.

The following definition of containment relation for matrices will be used to implement the ideas outlined above.

Definition 7.4 (Matrix (Useful) Containment). For any finite field \mathbb{F}_q and integers $k \geq m \geq 1$, we say that a matrix $M \in \mathbb{F}_q^{k \times k}$ contains a matrix $R \in \mathbb{F}_q^{m \times m}$ if there exist some $T \in \mathbb{F}_q^{k \times m}$ and a permutation matrix $P \in \mathbb{F}_q^{k \times k}$, such that $PMT = \begin{bmatrix} R \\ 0 \end{bmatrix}$. We say that P and T witness the containment of R in M. If, moreover, the last non-zero row of T is scaling of the standard basis vector, i.e., $T_j = \alpha e_m$ for some $\alpha \in \mathbb{F}_q^*$, then we say that containment is R in M is useful, and we denote it by $R \sqsubseteq_u M$.

We emphasize that useful containment is *not* a partial order.

Comparing this definition to the exposition above, the permutation P is used to express the fact that we can freely permute labels of variables U_1, \ldots, U_k , whereas the matrix T encodes coefficients for linear forms $B_1, \ldots B_{m-1}, B_m + \alpha V_j$. Finally, the condition on the last non-zero row of T being of form αe_m is here to express the idea that V_j is not allowed to appear in any of the forms $B_1, \ldots B_m$.

The following fact about useful containment will be helpful.

PROPOSITION 7.5. If $R \sqsubset_u M$, then for every upper triangular matrix U with non-zero diagonal elements $U_{i,i}$, we also have $R \sqsubset_u MU^{-1}$.

PROOF. Consider matrix T and permutation P as in the definition of useful containment for $R \sqsubset_u M$. We can pick the very same permutation P and matrix T' = UT to witness $R \sqsubset_u MU^{-1}$.

All we have to show is that last non-zero row of T' is the (scaled) standard basis vector $\alpha \boldsymbol{e}_m$. Indeed, if j_0 is the last non-zero row of T, and $j > j_0$, then rows $(U)_j$ are supported exclusively on elements with indices larger than j_0 , and hence $(UT)_j = (U)_j T = 0$. However, $(UT)_{j_0} = \sum_i U_{j_0,i} T_i$. Since for $i < j_0$ the entry $U_{j_0,i} = 0$, and for $i > j_0$ we have $T_i = 0$, this implies $(UT)_{j_0} = U_{j_0,j_0} T_{j_0} = U_{j_0,j_0} \alpha \boldsymbol{e}_m$, where the last equality follows from the fact that T was useful—that is, $T_{j_0} = \alpha \boldsymbol{e}_m$ and $T_i = 0$ for $i > j_0$. Since both $U_{j_0,j_0} \neq 0$ and $\alpha \neq 0$, we have $U_{j_0,j_0} \alpha \neq 0$, as desired.

Lemma 5.7 can now be reinterpreted as the following lemma. We give a full new proof here, as we describe it now in the language of useful containment.

Lemma 7.6. Every mixing matrix
$$M \in \mathbb{F}_q^{k \times k}$$
 usefully contains matrix $H = \begin{pmatrix} 1 & 0 \\ \alpha & 1 \end{pmatrix}$ for some $\alpha \in \mathbb{F}_q^*$.

PROOF. For every matrix M, there is some permutation matrix P' and pair L, U, such that P'MLU, where L is lower triangular (such that its diagonal is all 1s) and U is upper triangular. ¹⁵ Matrix M being mixing is equivalent to the statement that L and U are invertible, and, moreover, L is not diagonal. (In particular, M is invertible if and only if L and U are and $M = (P')^{-1}LU$ is the permutation of an upper triangular matrix if and only if L is diagonal.) Thus, by Proposition 7.5, it suffices to show that every lower-triangular L, which is not diagonal, contains $H \sqsubset_u L$. Indeed, let s be the last column of L that contains more than a single non-zero entry, and let r to be the last row of non-zero entry in column $L_{\cdot,s}$. Note that column $L_{\cdot,r}$ has single non-zero entry $L_{r,r}=1$. We will show a matrix $T \in \mathbb{F}_q^{k \times 2}$ as in the definition of useful containment. Let us specify a second column of $T_{2} := e_{r}$; note that in this case $(LT)_{2} = e_{r}$. To specify the first column of T, we wish to find a linear combination of columns of L_1, \ldots, L_{r-1} , such that $\sum_{i \leq r-1} t_i L_{i, \cdot} = e_s + \alpha e_r$, where $\alpha = L_{r,s} \neq 0$. Then coefficients t_i can be used as the first column of matrix T, which would imply that $(LT)_{i,1} = e_s + \alpha e_r$. We can set those coefficients to $t_i = -L_{s,i}$ for $i \in [s+1,r-1]$, $t_s = 1$ and $t_i = 0$ for i < s; this setting is correct, because columns L_i . for $i \in [s+1, r-1]$ have only one non-zero entry $L_{i,i}$. As already observed, the first column of LT is $e_s + \alpha e_r$ while the second column is e_r . Thus, if P is any matrix corresponding to a permutation that maps $s \mapsto 1$ and $r \mapsto 2$, then the containment $H \sqsubset_{u} L$ is witnessed by pair P and T, as desired.

Lemma 7.7. If matrix
$$R \sqsubset_u M$$
 where $R \in \mathbb{F}_q^{s \times s}$ and $M \in \mathbb{F}_q^{k \times k}$, then $R^{\otimes 2} \sqsubset_u M^{\otimes 2}$.

PROOF. Consider matrix T and permutation P that witness the useful containment for $R \sqsubset_u M$. Note that by the mixed product property of tensors, $P^{\otimes 2}M^{\otimes 2}T^{\otimes 2}=(PMT)^{\otimes 2}$. As such, restriction of a matrix $P^{\otimes 2}M^{\otimes 2}T^{\otimes 2}$ to rows corresponding to $[k]\times[k]$ is exactly $R^{\otimes 2}$, and all remaining rows are zero. We can apply additional permutation matrix \tilde{P} so that those are exactly first k^2 rows of the matrix $\tilde{P}P^{\otimes 2}M^{\otimes 2}T^{\otimes 2}$ give matrix $R^{\otimes 2}$, and the remaining rows are zero. Finally, since the last non-zero row of T was a scaling of the standard basis vector, the same is true for $T^{\otimes 2}$.

Lemma 7.8. If matrix $M \in \mathbb{F}_q^{k \times k}$ usefully contains matrix $R = \begin{pmatrix} 1 & 0 \\ \alpha & 1 \end{pmatrix}^{\otimes 2}$, then matrix M satisfies the strong suction condition of $(\frac{1}{k}, 2 - \varepsilon)$ exponential polarization.

PROOF. By the definition of exponential matrix polarization, it suffices to prove that there exists an index $j \in [k]$ such that $\overline{H}((UM)_j|(UM)_{< j},W) \leq \overline{H}((UR)_4|(UR)_{< 4},W)$. Once we have this, the proof of Lemma 7.3 (specifically Equation (34)) asserts that the conditional entropy is bounded as required. So we turn to proving this.

 $[\]overline{}^{15}$ This, e.g., follows from Gaussian Elimination and the corresponding "LU decomposition" of any matrix. Also note that the the assumption on the diagonal elements of L holds without loss of generality.

11:48 J. Błasiok et al.

Take $P \in \mathbb{F}_q^{k \times k}$ and $T \in \mathbb{F}_q^{k \times 4}$ witness the containment $R \sqsubset_u M$. Let, moreover, j be the last non-zero row of T. We have

$$\begin{split} \overline{H}((UM)_{j}|(UM)_{< j}, W) &= \overline{H}((UM)_{j}T_{j,4} + (UM)_{< j}T_{< j,4}|(UM)_{< j}, W) \\ &= \overline{H}((UMT)_{4}|(UM)_{< j}, W) \\ &= \overline{H}((UMT)_{4}|(UM)_{< j}, (UM)_{< j}T_{< j, < 4}, W) \\ &\leq \overline{H}((UMT)_{4}|(UM)_{< j}T_{< j, < 4}, W). \end{split}$$

In the above, the first equality follows, since $T_{j,4} \neq 0$ (and hence the map $(UM)_j \mapsto (UM)_j T_{j,4}$ is a bijection) and the fact that $(UM)_{< j} T_{< j,4}$ is deterministic function of $(UM)_{< j}$. The second equality follows, since $M_{> j,\cdot} = \mathbf{0}$, the third one introduces conditioning on $(UM)_{< j} T_{< j,< 4}$, which is deterministic given $(UM)_{< j}$, and the inequality follows, because entropy is decreasing under additional conditioning. Observe now that $(UM)_{< j} T_{< j,< 4} = (UMT)_{< 4}$. Indeed, according to the definition of useful containment and because j is last non-zero row of T, we have $T_{j,< 4} = 0$ (jth row has only one non-zero entry $T_{j,4}$), as well as $T_{>j,< 4} = 0$. Therefore

$$\overline{H}((UM)_{j}|(UM)_{< j}, W) \leq \overline{H}((UMT)_{4}|(UMT)_{< 4}, W)
= \overline{H}((UP^{-1}R)_{4}|(UP^{-1}R)_{< 4}, W)
= \overline{H}((UR)_{4}|(UR)_{< 4}, W),$$

where the last inequality follows from the fact that variables U_i are i.i.d. hence for the permutation matrix P, UP^{-1} has the same distribution as U.

With the above ingredients in place we are ready to prove Lemma 7.1.

Proof of Lemma 7.1. Since M is mixing we have that $M^{\otimes 2}$ is also mixing (Proposition 7.2) and so by Lemma 5.5 we have that $M^{\otimes 2}$ satisfies the conditions of matrix polarization. So it suffices to prove $M^{\otimes 2}$ satisfies the conditions of $(\frac{1}{k^2}, 2 - \varepsilon)$ exponential matrix polarization.

By Lemma 7.6, we have that M usefully contains $H = \begin{pmatrix} 1 & 0 \\ \alpha & 1 \end{pmatrix}$. Then, by Lemma 7.7 we have that $M^{\otimes 2}$ usefully contains $H^{\otimes 2}$. Finally, by Lemma 7.8 applied to $M^{\otimes 2}$ (which is a $k^2 \times k^2$ matrix) we have that $M^{\otimes 2}$ satisfies $(1/k^2, 2-\varepsilon)$ exponential matrix polarization.

8 NEARLY OPTIMAL DECODING ERROR PROBABILITIES

Finally, we turn to the proofs of Theorems 1.19 and 1.20. Recall that the former yields codes achieving decoding error probability $\exp(-N^{\beta})$ for any $\beta < 1$ while doing so at block lengths polynomial in the gap to capacity. The latter result shows that the techniques in this article are essentially optimal (for a broad class of channels) by showing that any analysis that bounds the decoding error probability can be used as a black box to achieve a similar decoding error probability in our analysis framework while additionally guaranteeing convergence at polynomial lengths in the gap to capacity. We first present the former, though before doing so, we make a small digression to recollect some known definitions of linear codes that we will use in this section (for more details see, e.g., Reference [14, Chap. 2]).

8.1 Basics of Linear Error-correcting Codes

A linear q-ary error correcting code C of block length n_0 and dimension k_0 is a linear subspace of $\mathbb{F}_q^{n_0}$ of dimension k_0 . Equivalently, there exists a full rank $G \in \mathbb{F}_q^{k_0 \times n_0}$ such that $C = \{ \boldsymbol{v} \cdot G | \boldsymbol{v} \in \mathbb{F}_q^{k_0} \} - G$ is called the generator matrix of C. The kernel/null-space/dual of C, denoted by C^{\perp} or ker G, is given by $\{ \boldsymbol{w} | \langle \boldsymbol{w}, \boldsymbol{c} \rangle = 0 \text{ for all } \boldsymbol{c} \in C \}$. A generator matrix of C^{\perp} is called a parity-check matrix of

C. The distance of a code *C* is the minimum number of positions any two codewords in *C* differ in. For linear code *C*, its distance is exactly $\min_{c \in C \setminus \{0\}} \operatorname{wt}(c)$, where $\operatorname{wt}(x)$ is the number of non-zero elements in x.

8.2 Polar Codes with Decoding Failure Probability Approaching $2^{-N^{1-o(1)}}$

Theorem 1.19 is proved by giving a sufficient structural condition on matrices for very strong exponential polarization. The following lemma states this condition.

LEMMA 8.1. Let q be prime. If a mixing matrix $M \in \mathbb{F}_q^{k \times k}$ is decomposed as $M = [M_0|M_1]$, where $M_0 \in \mathbb{F}_q^{k \times (1-\eta)k}$ is such that $\ker M_0^T$ is a linear code of distance larger than 2b, then matrix M satisfies $(\eta, b - \varepsilon)$ -exponential matrix polarization for every $\varepsilon > 0$.

PROOF. By Lemma 5.5, we have that M satisfies the conditions of matrix polarization (specifically, variance in the middle and suction at the upper and lower ends from Definition 4.3). It remains only to argue exponential matrix polarization, i.e., strong suction at the lower end.

Let us again consider a sequence of i.i.d. pairs (U_i, W_i) for $i \in [k]$, such that $H(U_i|W_i) = \delta$. By Lemma 2.2, there is some $f: \Sigma \to \mathbb{F}_q$ such that $\Pr(f(W_i) \neq U_i) \leq \delta$ (for every $i \in [k]$). Let us define $\tilde{U}_i := U_i - f(W_i)$.

We will bound $\overline{H}((UM)_j|(UM)_{< j}, W)$, for all $j > (1 - \eta)k$. We have

$$\overline{H}((UM)_{i}|(UM)_{\leq i},W) \leq \overline{H}(U|(UM)_{\leq i},W) \leq \overline{H}(U|UM_{0},W) = H(\tilde{U}|\tilde{U}M_{0},W) \leq \overline{H}(\tilde{U}|\tilde{U}M_{0}),$$

where the first two inequalities follow from the fact that for random variables (X, Y, S, T) it is always the case that $\overline{H}(X|S,T) \leq \overline{H}(X,Y|S,T) \leq \overline{H}(X,Y|S)$ (the second inequality also uses the fact that $(UM)_{< j}$ is a sub-matrix of UM_0). The equality follows from the definition of \tilde{U}_i and the fact that $f(\cdot)$ is deterministic function. The final inequality follows from the fact that conditioning can only decrease the entropy.

Given $\tilde{U}M_0$ we can produce estimate $\hat{U} := \operatorname{argmin}_V \{ \operatorname{wt}(V) : VM_0 = \tilde{U}M_0 \}$, where $\operatorname{wt}(V) = |\{j : V_j \neq 0\}|$.

We note that if $\operatorname{wt}(\tilde{U}) \leq b$, then $\hat{U} = \tilde{U}$. Indeed, we have $\operatorname{wt}(\hat{U}) \leq \operatorname{wt}(\tilde{U})$, and therefore $\operatorname{wt}(\hat{U} - \tilde{U}) \leq 2\operatorname{wt}(\tilde{U}) \leq 2b$, but, however, $(\hat{U} - \tilde{U})M_0 = 0$, and by the assumption on distance of $\operatorname{ker} M_0^T$ we deduce that $\hat{U} - \tilde{U} = 0$. Therefore, $\operatorname{Pr}(\tilde{U} \neq \hat{U}) \leq \operatorname{Pr}(\operatorname{wt}(\tilde{U}) > b)$. All coordinates of \tilde{U} are independent, and each \tilde{U}_i is nonzero with probability at most δ , therefore

$$\Pr(\mathrm{wt}(\tilde{U}) > \beta_1) \leq \binom{k}{b} \delta^b.$$

Further, by Fano inequality (Lemma 2.3), we have

$$H(\tilde{U}|\tilde{U}M_0) \le 2C\delta^b(b\log\delta^{-1} + \log C + \log q),$$

where $C = \binom{k}{b}$. Again, for any ε , and small-enough δ (with respect to ε, b, k, q), we have $H(\tilde{U}|\tilde{U}M_0) \leq \delta^{b-\varepsilon}$.

This shows that for any $j > (1 - \eta)k$ (note that there are at least ηk such values of j) and small-enough δ we have

$$\overline{H}((UM)_j|(UM)_{< j},W) \le \delta^{b-\varepsilon},$$

which completes the proof of a exponential matrix polarization for M.

We are now almost ready to prove Theorem 1.19. We start with a corollary that uses standard results on existence of codes with good distance.

11:50 J. Błasiok et al.

COROLLARY 8.2. For every v > 0 and every prime field \mathbb{F}_q , there exist k, and matrix $M \in \mathbb{F}_q^{k \times k}$, such that matrix M satisfies $(1 - v, k^{1-v})$ exponential matrix polarization.

PROOF. Consider a parity check matrix M_0 of a BCH code with distance $2k^{1-\nu}$. We can achieve this with a matrix $M_0 \in \mathbb{F}_q^{k \times k_0}$, where $k_0 = O(k^{1-\nu} \log k)$ (see, e.g., Reference [14, Exercise 5.10]). Hence, as soon as $k > \Omega(2^{\nu^{-1} \log \nu^{-1}})$, we have $k_0 < \nu k$. Note that if $k_0 = \nu_0 k$, then by Lemma 8.1 we can hope for $(1 - \nu_0, k^{1-\nu_0} - \varepsilon)$ exponential matrix polarization. We can now complete M_0 to a mixing matrix to get overall $(1 - \nu, k^{1-\nu})$ exponential matrix polarization (since $\nu_0 < \nu$). To complete matrix M_0 to a mixing matrix, by Lemma 5.8 it is enough to complete it in arbitrary way to an invertible matrix, since already the first column of M_0 has support larger than 1.

Remark 8.3 (Exponential Polarization of Random Kernels). It is worth noting that by the same argument and standard results on the distance of random linear codes, a random matrix $M \in \mathbb{F}_q^{k \times k}$ with high probability satisfies a $(1 - \nu, k^{1-\nu})$ local polarization, with $\nu \to 0$ as $k \to \infty$. Thus polar codes arising from a large random matrix will usually have this property.

We now complete the proof of Theorem 1.19.

PROOF OF THEOREM 1.19. Given $\beta < 1$ and q, let $\nu = (1-\beta)/3$. Now let k and M be as given by Corollary 8.2. By Theorem 4.4, we have that for every channel $C_{Y|Z}$, M satisfies $(1-\nu,k^{1-\nu})$ -exponential local polarization. By Theorem 1.9, we have that the same martingale satisfies Λ -exponentially strong polarization for $\Lambda = (1-\nu)^2 \log_2 k \ge (1-2\nu) \log_2 k$. By Theorem 1.11 (in particular, Remark 1.12), we then get that the resulting codes have failure probability $O(N \cdot \log q \cdot \exp(-N^{1-2\nu})) \le \exp(-N^{1-3\nu}) = \exp(-N^{\beta})$, where the first inequality holds for sufficiently large N (as a function of ν).

8.3 Universality of Local Polarization

Suppose we know that polar codes associated with a matrix $M \in \mathbb{F}_q^{k \times k}$ achieve capacity with error probability $\exp(-n^\beta)$ in the limit of block lengths $n \to \infty$ (which may happen at lengths growing super polynomially in ε the gap to capacity). In this section, we prove a general result (previously stated as Theorem 1.20) that "lifts" (in a black box manner) such a statement to the claim that, for every $\beta' < \beta$, polar codes associated with M achieve polynomially fast convergence to capacity (i.e., the block length n can be as small as $\operatorname{poly}(1/\varepsilon)$ for rates within ε of capacity) and $\operatorname{exp}(-n^{\beta'})$ decoding error probability $\operatorname{simultaneously}$. Thus, convergence to capacity at finite block length comes with almost no price in the (exponent of) decoding failure probability.

Put differently, the result states that one can get polynomial convergence to capacity for free once one has a proof of convergence to capacity in the limit of $n \to \infty$ with root-exponential decoding error probability. Such proofs of convergence to capacity has been shown in Reference [20] for the binary alphabet and Reference [24] for general alphabets. Yet another way of viewing the results of this section are that every proof of convergence to capacity has a proof of local polarization embedded in it.

We get our result by proving a structural result that is roughly a converse to Lemma 8.1. Specifically, in Lemma 8.5 we show that if a matrix M leads to a polar code with exponentially small failure probability, then some high (but constant sized) tensor power $M^{\otimes t}$ of M contains the parity check matrix of a high distance code. In fact, more generally if a matrix in $\mathbb{F}_q^{k \times s}$ is the parity check matrix of a code that has a decoding algorithm that corrects errors from a q-symmetric channel with failure probability $\exp(-k^{\beta})$, then this code has high distance.

Combining Lemma 8.5 with Lemma 8.1, we get that every matrix that leads to a polar code with low error probability has a constant sized tensor that is a exponentially polarizing matrix. This immediately leads to a proof of Theorem 1.20.

To derive our results, we focus on a simple q-ary symmetric channel defined next.

Definition 8.4. For any finite field \mathbb{F}_q and $\gamma \in [0,1]$, we will denote by $B_q(\gamma)$ the distribution on \mathbb{F}_q such that for $Z \sim B_q(\gamma)$ we have $\Pr(Z=0) = 1 - \gamma$ and $\Pr(Z=k) = \frac{\gamma}{q-1}$ for any $k \neq 0$.

LEMMA 8.5. Consider a matrix $H \in \mathbb{F}_q^{k \times s}$ and arbitrary decoding algorithm $Dec : \mathbb{F}_q^s \to \mathbb{F}_q^k$, such that for independent random variables $U_1, \ldots U_i \sim B_q(\gamma)$ with $\gamma < \frac{1}{2}$, we have $Pr(Dec(UH) \neq U) < \exp(-k^{\beta})$. Then $\ker H$ is a code of distance at least $\frac{k^{\beta}}{\ln^{-1}(q/\gamma)}$.

PROOF. Consider maximum likelihood decoder $\mathrm{Dec}'(y) := \mathrm{argmax}_{x \in \mathbb{F}_q^k} \Pr(U = x | UH = y)$. By definition, we have $\Pr(\mathrm{Dec}'(UH) \neq U) < \Pr(\mathrm{Dec}(UH) \neq U) < \exp(-k^\beta)$.

Note that for *U* distributed according to $B_q(\gamma)$, we have $\mathrm{Dec}'(y) = \mathrm{argmin}_{x:xH=y} \mathrm{wt}(x)$, where $\mathrm{wt}(x)$ is number of non-zero elements of x.

Consider set $E = \{ \boldsymbol{x} \in \mathbb{F}_q^k \mid \text{, where there exists } \boldsymbol{h} \in \ker M, \operatorname{wt}(\boldsymbol{x} + \boldsymbol{h}) < \operatorname{wt}(\boldsymbol{x}) \}$, and observe that $\operatorname{Pr}(\operatorname{Dec}'(UH) \neq U) \geq \operatorname{Pr}(U \in E)$. We say that vector $\boldsymbol{u} \in \mathbb{F}_q^k$ is dominated by $\boldsymbol{v} \in \mathbb{F}_q^k$ (denoted by $\boldsymbol{u} \leq \boldsymbol{v}$) if and only if $\forall i \in \operatorname{supp}(\boldsymbol{u}), \ \boldsymbol{u}_i = \boldsymbol{v}_i$. We will argue that for any $\boldsymbol{w}_1 \in E$ and any $\boldsymbol{w}_2 \geq \boldsymbol{w}_1$, we have $\boldsymbol{w}_2 \in E$. Indeed, if $\boldsymbol{w}_1 \in E$, then there is some $\boldsymbol{h} \in \ker H$ such that $\operatorname{wt}(\boldsymbol{w}_1 + \boldsymbol{h}) < \operatorname{wt}(\boldsymbol{w}_1)$. We will show that $\operatorname{wt}(\boldsymbol{w}_2 + \boldsymbol{h}) < \operatorname{wt}(\boldsymbol{w}_2)$, which implies that $\boldsymbol{w}_2 \in E$. Given that $\boldsymbol{w}_1 \leq \boldsymbol{w}_2$, we can equivalently say that there is a vector \boldsymbol{d} with $\boldsymbol{w}_1 + \boldsymbol{d} = \boldsymbol{w}_2$ and $\operatorname{wt}(\boldsymbol{w}_2) = \operatorname{wt}(\boldsymbol{w}_1) + \operatorname{wt}(\boldsymbol{d})$. Hence,

$$\operatorname{wt}(w_2 + h) = \operatorname{wt}(w_1 + d + h) \le \operatorname{wt}(w_1 + h) + \operatorname{wt}(d) < \operatorname{wt}(w_1) + \operatorname{wt}(d) = \operatorname{wt}(w_2).$$

Consider now $w_0 \in \ker H$ to be minimum weight non-zero vector, and let us denote $A = \operatorname{wt}(w_0)$. We wish to show a lower bound for A. By definition of the set E, we have $w_0 \in E$, and by upward closure of E with respect to domination we have

$$\Pr(U \in E) \ge \Pr(w_0 \le U) = \left(\frac{\gamma}{q-1}\right)^A \ge \left(\frac{\gamma}{q}\right)^A.$$

However, we have

$$\Pr(U \in E) \le \Pr(\operatorname{Dec}'(UH) \ne U) \le \Pr(\operatorname{Dec}(UH) \ne U) \le \exp(-k^{\beta}).$$

By comparing these two inequalities we get

$$A \ge \frac{k^{\beta}}{\ln(q/\gamma)} \ .$$

PROOF OF THEOREM 1.20. Consider the channel that outputs X+Z on input X, where $Z \sim B_q(\gamma)$ for some $\gamma > 0$ (depending on β, β'). The hypothesis on M implies that for sufficiently large n the polar code of block length n corresponding to M will have failure probability at most $\exp(-n^\beta)$ on this channel. Using the well-known equivalence between correcting errors for this additive channel and linear compression schemes (see, e.g., Reference [14, Prop. 11.2.1]), we obtain that for all large-enough t there is some subset S of $(h_q(\gamma) + \varepsilon)k^t$ columns of $M^{\otimes t}$ that defines a linear compression scheme (for k^t i.i.d. copies of $B_q(\gamma)$), along with an accompanying decompression scheme with error probability (over the randomness of the source) at most $\exp(-k^{\beta t})$.

We now claim that for all $\beta' < \beta$, there exists $t_0 = t_0(\beta', \beta)$ such that the Arikan martingale associated with some column permuted version of $M^{\otimes t_0}$, is $\beta' t_0 \log_2 k$ -exponentially strongly polarizing.

The proof of this claim is in fact immediate, given the ingredients developed in previous sections. Apply the hypothesis about M in the theorem with the choice $\varepsilon = (\beta - \beta')/4$ and γ chosen small enough as a function β , β' so that $h_q(\gamma) \leq (\beta - \beta')/4$ and let t_0 be a larger than promised value of t in the statement and large enough so that $3 \ln(q/\gamma) < m^{(\beta-\beta')/2}$ with $m := k^{t_0}$. Take, moreover,

11:52 J. Błasiok et al.

 $\ell = (h_q(\gamma) + \varepsilon)m$ and $L = M^{\otimes \ell_0}$. Using Lemma 8.5 and the equivalence between linear coding for source and channel coding (mentioned above), we know there is submatrix $L' \in \mathbb{F}_q^{m \times \ell}$ of L such that $\ker((L')^T)$ defines a code of distance $\Delta \geq m^\beta / \ln(q/\gamma)$. Define $M_0 = [L' \mid \cdot] \in \mathbb{F}_q^{m \times m}$ to be any matrix obtained by permuting the columns of L such that the columns in L' occur first. By Lemma 8.1, the matrix M_0 is $(1 - \ell/m, \Lambda)$ -exponential matrix polarizing with $\Lambda = \Delta/2 - o(1) > \Delta/3$.

For our choice of γ , ε , we have $\ell/m \leq \frac{\beta-\beta'}{2}$ and for our choice of t_0 (and therefore m) we have $\Lambda \geq m^{(\beta+\beta')/2}$. Using Theorem 4.4 and Theorem 1.7, it follows that the Arikan martingale associated with M_0 exhibits $(\beta+\beta')/2 \times (1-\frac{\beta-\beta'}{2})\log_2 m$ -exponentially strong polarization. Since

$$\beta^{"} \stackrel{\mathrm{def}}{=} \frac{\beta + \beta'}{2} \cdot \left(1 - \frac{\beta - \beta'}{2}\right) = \beta' + \frac{\beta - \beta'}{2} \cdot \left(1 - \frac{\beta + \beta'}{2}\right) > \beta',$$

the claim follows (in the above we used the fact that $0 < \beta' < \beta < 1$).

Applying Theorem 1.11 (and Remark 1.12) to the matrix $M_0 = M^{\otimes t_0}$ we conclude that there is a polynomial p such that given the gap to capacity $\varepsilon > 0$, and for every s satisfying $N = k^{t_0 s} \ge p(\frac{1}{\varepsilon})$ there is an affine code generated by a subset of rows of $(M_0^{-1})^{\otimes s}$, which achieves ε -gap to capacity and has failure probability $\exp(-N^{\beta^n}) \cdot N \cdot \log q < \exp(-N^{\beta'})$ for large-enough N. But this resulting code is simply an affine code generated by a subset of the rows of $(M^{-1})^{\otimes t}$, for $t = st_0$, which concludes the proof.

APPENDIX

A CODES FROM POLARIZATION

In this section, we describe the construction of polar codes and analyze the failure probability of decoders by corresponding them to the Arıkan martingale. This proves Theorems 1.11 and 1.14.

Specifically, we first describe the polar encoder along with a fast $O(n \log n)$ -time implementation, where n is the blocklength. Then, in Appendix A.2 we define the (inefficient) successive-cancellation decoder and analyze its failure probability assuming a correspondence between polar coding and the Arıkan martingale. In Appendix A.2.2, we describe a fast $O(n \log n)$ -time decoder that is functionally equivalent to the successive-cancellation decoder. Finally, in Appendix A.2.3, we prove the required correspondence between polar coding and the Arıkan martingale.

Throughout this section, fix parameters $k \in \mathbb{N}$ as the dimension of the mixing matrix $M \in \mathbb{F}_q^{k \times k}$, \mathbb{F}_q as a finite field, and $n = k^t$ as the codeword length.

A.1 Polar Encoder

Given a set $S \subseteq [n]$ and a fixing $\alpha \in \mathbb{F}_q^{|S^c|}$, 16 we define the polar code of dimension |S| by giving the encoder mapping $\mathbb{F}_q^S \to \mathbb{F}_q^n$ as follows:

ALGORITHM 1: Polar Encoder

Constants: $M \in \mathbb{F}_q^{k \times k}, S \subseteq [n], \alpha \in \mathbb{F}_q^{S^c}$

Input: $U \in \mathbb{F}_q^S$ Output: $Z \in \mathbb{F}_q^n$

1: **procedure** Polar-Encoder($U; \alpha$)

- Extend U to $\overline{U} \in \mathbb{F}_q^n$ by letting $(\overline{U}_i)_{i \notin S} = \alpha$ for coordinates not in S
- 3: **Return** $Z = \overline{U} \cdot (M^{-1})^{\otimes t}$

¹⁶We use the notation $S^c = [n] \setminus S$.

The above gives a polynomial time algorithm for encoding. An $O_q(n \log n)$ algorithm can also be obtained by using the recursive structure imposed by the tensor powers.

Below, we switch to considering vectors in $\mathbb{F}_q^{k^t}$ as tensors in $(\mathbb{F}_q^k)^{\otimes t}$, indexed by multiindices $i \in [k]^t$. The following encoder takes as input the "extended" message \overline{U} , as defined above.

```
ALGORITHM 2: Fast Polar Encoder
```

```
Constants: M \in \mathbb{F}_q^{k \times k}
Input: \overline{U} \in (\mathbb{F}_q^k)^{\otimes t}
Output: Z = \overline{U} \cdot (M^{-1})^{\otimes t}
  1: procedure Fast-Polar-Encoder_t(U)
            If t = 0 then
  2:
                  Return U
  3:
            for all j \in [k] do
  4:
                  Z^{(j)} \leftarrow \text{Fast-Polar-Encoder}_{t-1}(\overline{U}_{[\cdot,j]})
  5:
            for all i \in [k]^{t-1} do
  6:
                  Z_{[i,\cdot]} \leftarrow (Z_i^{(1)}, Z_i^{(2)}, \dots, Z_i^{(k)}) \cdot M^{-1}
  7:
             Return Z
  8:
```

It is not too hard to verify that Algorithm 2 runs in $O_{k,q}(n \log n)$ time. Indeed, if T(n) is the runtime of the algorithm on inputs of size $n = k^t$, then each call results in k recursive calls to inputs of size $\frac{n}{k}$. Further, each recursive call solve $\frac{n}{k}$ systems of linear equations (each of which can be solved in $O_q(k^3)$ time). Thus we get the recurrence (using the fact that k is a constant) of $T(n) = k \cdot T(n/k) + O_{k,q}(n)$, which results in the desired $O_{k,q}(n \log n)$ runtime.

A.2 The Successive-Cancellation Decoder

Here we describe a successive-cancellation decoder. Note that this decoder is not efficient, but the fast decoder described later will nearly have the same error probability as this decoder.

For given channel outputs Y, let Z be the posterior distribution on channel inputs given outputs Y. Each $Z_i \in \Delta(\mathbb{F}_q)$ is the conditional distribution $Z_i|Y_i$ defined by the channel $C_{Y|Z}$ and the received output Y_i .

Now we define the decoder on the distribution vector Z and the fixing $\alpha \in (\mathbb{F}_q \cup \{\bot\})^n$ as follows. We implicitly represent the subset S^c of fixed positions by denoting $\alpha_i = \bot$ for those indices.

Remark A.1. We note that parts in brown are not needed for the algorithm itself and only used in the analysis. Further, unless explicitly stated otherwise, we will use SC-DECODER to just denote the \hat{U} part of the output (i.e., we will ignore P by default).

Note that several of the above steps, including computing the joint distribution of U and marginal distributions of U_i , are not computationally efficient though we will get efficient algorithms effectively approximating these distributions later. Even then, we will only get an algorithm that gets an estimate of the probabilities $\Pr_U(U_i = x)$ to within an additive error of 1/4 for every $x \in \mathbb{F}_q$. In what follows, we will use the following definition:

Definition A.2. We will term an algorithm that runs an SC-DECODER where the algorithm gets an estimate of the probabilities $\Pr_U(U_i = x)$ to within an additive error of 1/4 for every $x \in \mathbb{F}_q$ an Approximate-Successive-Cancellation Decoder.

11:54 J. Błasiok et al.

ALGORITHM 3: Successive-Cancellation Decoder

```
Constants: M \in \mathbb{F}_q^{k \times k}, n = k^s,
Input: Z \in \Delta(\mathbb{F}_q)^n, \alpha \in (\mathbb{F}_q \cup \{\bot\})^n
Output: \hat{U} \in \mathbb{F}_q^n, P \in (\Delta(\mathbb{F}_q) \cup \bot)^n
    1: procedure \overline{SC}-Decoder(Z; \alpha)
                Compute the distribution U \in \Delta(\mathbb{F}_q^n) defined by U \leftarrow ZM^{\otimes s}
   2:
                for all i \in [n] do
   3:
                       If \alpha_i = \bot then
   4:
                              For x \in \mathbb{F}_q, \hat{U}_i \leftarrow \operatorname{argmax}_{x \in \mathbb{F}_q} \{ \Pr_U (U_i = x) \} ; \underbrace{P_i(x)} \leftarrow \Pr_U (U_i = x) 
   5:
   6:
                              \hat{U}_i \leftarrow \alpha_i; P_i \leftarrow \bot
   7:
                       Update distribution U \leftarrow (U|U_i = \hat{U}_i)
   8:
                Return \hat{U}, P
   9:
```

A.2.1 Decoding Analysis. For this section, it will be useful to keep Remark A.1 in mind.

We will first reason about the "genie-aided" case, when the fixing $\alpha \in (\mathbb{F}_q \cup \{\bot\})^n$ of non-message bits is chosen uniformly at random, and revealed to both the encoder and decoder. Then, we will argue that it is sufficient to use a deterministic fixing $\alpha = \alpha_0$.

We now argue that over a uniform choice of message U_S , and a uniform fixing α of non-message bits, the probability of decoding failure is bounded as follows.

CLAIM A.3. For $S \subseteq [n]$, let $V \in (\mathbb{F}_q \cup \{\bot\})^n$ be given by $V_i \sim \mathbb{F}_q$ if $i \in S$ and \bot otherwise. Let $\alpha \in (\mathbb{F}_q \cup \{\bot\})^n$ be given by $\alpha_i \sim \mathbb{F}_q$ if $i \notin S$ and \bot otherwise. Let $Z := POLAR-ENCODER(V; \alpha)$ and Y sampled according to the channel $Y := C_{Y|Z}(Z)$. Let $U \in \mathbb{F}_q^n$ be given by $U_i = V_i$ if $i \in S$ and α_i if $i \notin S$. With this notation, we have

$$\Pr[SC\text{-Decoder}(Y; \boldsymbol{\alpha}) \neq \boldsymbol{U}] \leq \sum_{i \in S} H(\boldsymbol{U}_i \mid \boldsymbol{U}_{< i}, Y).$$

Furthermore, for every approximate-successive-cancellation decoder D we have

$$\Pr[D(Y; \boldsymbol{\alpha}) \neq \boldsymbol{U}] \leq 3 \sum_{i \in S} H(\boldsymbol{U}_i \mid \boldsymbol{U}_{< i}, Y).$$

PROOF. Note that *U* is uniform over \mathbb{F}_q^n . Now, we have

$$\begin{split} \Pr\left(\text{SC-Decoder}(Y; \boldsymbol{\alpha}) \neq \boldsymbol{U}\right) &= \Pr\left(\exists i \; \hat{\boldsymbol{U}}_i \neq \boldsymbol{U}_i\right) \\ &= \sum_{i \leq n} \Pr\left(\hat{\boldsymbol{U}}_i \neq \boldsymbol{U}_i \text{ and } \hat{\boldsymbol{U}}_{< i} = \boldsymbol{U}_{< i}\right) \\ &\leq \sum_{i \leq n} \Pr\left(\hat{\boldsymbol{U}}_i \neq \boldsymbol{U}_i \mid \hat{\boldsymbol{U}}_{< i} = \boldsymbol{U}_{< i}\right). \end{split}$$

Clearly, for $i \notin S$ we have $\Pr[\hat{U}_i \neq U_i] = 0$, since both are defined to be equal to α_i on those coordinates. It is enough to show that for $i \in S$ we have

$$\Pr(\hat{U}_i \neq U_i \mid U_{< i} = \hat{U}_{< i}) \leq H(U_i \mid U_{< i}, Y).$$

This follows directly from Lemma 2.2, as \hat{U}_i is defined exactly as a maximum likelihood estimator of U_i given channel outputs Y and conditioning on $U_{< i}$ (note that the conditioning is happening in Line 8).

The furthermore part of the claim follows from using the furthermore part of Lemma 2.2 in the final step above. \Box

CLAIM A.4. Let $n = k^t$, $U \sim \mathbb{F}_q^n$, $Z := U(M^{-1})^{\otimes t}$, $Y := C_{Y|Z}(Z)$, where $C_{Y|Z}$ is a symmetric channel. If Arikan Martingale associated with (M, C) satisfies $(\tau_\ell, \tau_h, \varepsilon)$ -polarization, then there exists a subset $S \subset [n]$ of size (Capacity $(C_{Y|Z}) - \varepsilon - \tau_h$)n, such that

$$\sum_{i \in S} H(U_i \mid U_{< i}, Y) \le \tau_{\ell} n \log q.$$

PROOF. Applying Lemma A.18, we can deduce that for uniformly random index $i \in [n]$, normalized entropies $\overline{H}(U_i|U_{\le i},Y)$ are distributed identically as X_t in the Arikan Martingale.

Now, for symmetric channels, the uniform distribution achieves capacity (see, e.g., Reference [6, Theorem 7.2.1]). In addition, since matrix $(M^{(-1)})^{\otimes t}$ is invertible, vector Z also has a uniform distribution. Thus, for uniform channel input Z,

$$n \cdot \text{Capacity}(C_{Y|Z}) = \overline{H}(Z) - \overline{H}(Z|Y) = n - \overline{H}(Z|Y).$$
 (35)

Let *S* be the set of all indices *i* such that $\overline{H}(U_i \mid U_{\leq i}, Y) < \tau_{\ell}$. By definition, we have

$$\sum_{i \in S} \overline{H}(U_i \mid U_{< i}, Y) \le \tau_{\ell} n,$$

as desired.

Now observe that polarization of martingale X_t and Lemma A.18 directly implies that we have at most εn indicies i satisfying $\overline{H}(U_i \mid U_{< i}) \in (\tau_\ell, 1 - \tau_h)$ (recall that in Lemma A.18 we pick one such index uniformly at random). Let S' be a set of indices for which $\overline{H}(U_i \mid U_{< i}, Y) > 1 - \tau_h$. We have

$$n(1 - \operatorname{Capacity}(C_{Y|Z})) = \overline{H}(U(M^{-1})^{\otimes t} \mid Y), \qquad (\text{Equation (35)})$$

$$= \overline{H}(U_1, \dots, U_n \mid Y), \qquad (\text{Since } (M^{-1})^{\otimes t} \text{ is full rank})$$

$$= \sum_{i \in [n]} \overline{H}(U_i \mid U_{< i}, Y), \qquad (\text{Chain rule})$$

$$\geq \sum_{i \in S'} \overline{H}(U_i \mid U_{< i}, Y),$$

$$\geq (1 - \tau_h) |S'| \geq |S'| - \tau_h n,$$

which implies that

$$|S'| \leq n(1 - \text{Capacity}(C_{Y|Z}) + \tau_h),$$

and, finally,

$$|S| \ge n - |S'| - \varepsilon n \ge n(\operatorname{Capacity}(C_{Y|Z}) - \varepsilon - \tau_h)$$
.

We can now combine the above to prove a version of Theorem 1.14 for the (inefficient) successive-cancellation decoder:

Theorem A.5. Let C be a q-ary symmetric memoryless channel, and let $M \in \mathbb{F}_q^{k \times k}$ be an invertible matrix. If the Arikan martingale associated with (M, C) satisfies $(\tau_\ell, \tau_h, \varepsilon)$ -polarization, then for every t, there is an affine code C that is generated by the rows of $(M^{-1})^{\otimes t}$ and an affine shift, such that the rate of C is at least Capacity $(C) - \varepsilon(t) - \tau_h(t)$, and C can be encoded in time $O(n \log n)$, where $n = k^t$. Furthermore, the successive-cancellation decoder succeeds with probability at least $1 - n \log(q) \tau_\ell$,

11:56 J. Błasiok et al.

and every approximate-successive-cancellation decoder succeeds with probability at least $1 - 3n \log(q) \tau_{\ell}$.

PROOF. Let $\overline{U} \sim \mathbb{F}_q^n$, $\overline{Z} := U(M^{-1})^{\otimes t}$, and $\overline{Y} := C_{Y|Z}(\overline{Z})$. By Claim A.4, there exist a set $S \subset [n]$ of size (Capacity($C_{Y|Z}$) $-\varepsilon - \tau_h$)n, such that

$$\sum_{i \in S} H(\overline{U}_i \mid \overline{U}_{< i}, \overline{Y}) \le \tau_{\ell} n \log q.$$

However, by Claim A.3, the failure probability of the successive-cancellation decoder is bounded by

$$\Pr_{U,\alpha,Y}[\text{SC-Decoder}(Y;\alpha)_S \neq U] \leq \sum_{i \in S} H(U_i \mid U_{< i}, Y), \tag{36}$$

where random variables U, Y, α are defined as in Claim A.3. Note that, in fact, the joint distributions of (U, Y, Z) and $(\overline{U}, \overline{Y}, \overline{Z})$ are the same, despite superficially more complicated way in which sampling from distribution (U, Y, Z) was defined. Therefore,

$$\sum_{i \in S} H(U_i \mid U_{< i}, Y), = \sum_{i \in S} H(\overline{U}_i \mid \overline{U}_{< i}, \overline{Y})$$

$$\leq \tau_{\ell} n \log q.$$

Note that this failure probability is an average over random choice of fixing α , but this implies there is some deterministic fixing $\alpha = \alpha_0$ with failure probability at least as good. Further, by linearity of the encoding (Algorithm 2) such a deterministic fixing yields an affine code. The rate of this code is $|S|/n \ge (\text{Capacity}(C_{Y|Z}) - \varepsilon - \tau_h)$ as desired.

If we replace the successive-cancellation decoder by an approximate successive cancellation decoder, then the theorem follows by using the furthermore part of Claim A.3 in Equation (36) above.

A.2.2 Fast Decoder. In this section, we will define the recursive FAST-DECODER algorithm. The observation that polar codes admit a recursive fast-decoder was made in the original work of Arıkan [2]. Our presentation is somewhat different in that it decodes general product distributions (and does not require the marginals to be identical).

Fast-Decoder will take on input descriptions of the posterior distributions on channel inputs $\{Z_i\}_{i\in[k]^s}$ for some s, where each individual $Z_i\in\Delta(\mathbb{F}_q)$ is a distribution over \mathbb{F}_q , as well as $\alpha\in(\mathbb{F}_q\cup\{\bot\})^{[k]^s}$, where $\alpha_i\in\mathbb{F}_q$ are the fixed values corresponding to non-message positions. The output of Fast-Decoder is a vector $\hat{Z}\in(\mathbb{F}_q^k)^{\otimes s}$ —the guess for the actual channel inputs. To recover the message, it is enough to apply $\hat{U}:=\hat{Z}M^{\otimes s}$ and restrict it to the positions where $\alpha_i=\bot$.

In Algorithm 4, for $W_i \in \Delta(\mathbb{F}_q^k)$, a description of joint probability distribution over \mathbb{F}_q^k , we will write $\pi_j(W_i) \in \Delta(\mathbb{F}_q)$ as a jth marginal of W_i for $j \in [k]$, i.e., projection on the jth coordinate. In addition, we will use $\pi_{\leq j}(W_i) \in \Delta(\mathbb{F}_q)^j$ to denote the projection of W to the first j marginal coordinates.

We make an remark analogous to Remark A.1 for FAST-DECODER:

Remark A.6. In the code above, the parts in brown are not needed for the running of the algorithm but included, since they help with the analysis. Further, unless explicitly stated otherwise, we will use SC-Decoder to just denote the \hat{Z} part of the output (i.e., we will ignore Q, \hat{U}^F by default).

ALGORITHM 4: Fast Decoder

```
Constants: M \in \mathbb{F}_q^{k \times k}
Input: Z = \{Z_i \in \Delta(\mathbb{F}_q)\}_{i \in [k]^s}, \ \alpha \in (\mathbb{F}_q \cup \{\bot\})^{[k]^s}
Output: \hat{Z} \in (\mathbb{F}_q^k)^{\otimes s}, Q \in (\Delta(\mathbb{F}_q^k) \cup \{\bot\})^{\otimes s}, \hat{U}^{\mathsf{F}} \in (\mathbb{F}_q^k)^{\otimes s}
   1: procedure Fast-Decoder<sub>s</sub>(Z; \alpha)
                If s = 0 then
                       If \alpha = \bot then
   3:
                              Return \hat{Z} = \operatorname{argmax}_{x \in \mathbb{F}_a} \Pr(Z = x), Q = Z, \hat{U}^{F} = \hat{Z}
   4:
   5:
                              Return \hat{Z} = \alpha, Q = \bot, \hat{U}^{F} = \alpha
   6:
               else
   7:
                       for all i \in [k]^{s-1} do
   8:
                              Compute joint distribution W_i \in \Delta(\mathbb{F}_a^k), given by W_i \leftarrow Z_{\lceil \cdot, i \rceil} M
   9:
                       for all j \in [k] do
 10:
                              Z^{\prime(j)} \leftarrow \{\pi_i(W_i)\}_{i \in [k]^{s-1}}
 11:
                              \hat{\pmb{V}}_{[j,\cdot]}, \underline{\pmb{Q}}_{[j,\cdot]}, \hat{\pmb{U}}_{[j,\cdot]}^{\mathrm{F}} \leftarrow \mathrm{Fast-Decoder}(\pmb{Z}'^{(j)}; \pmb{\alpha}_{[j,\cdot]}, s-1)
 12:
                              for all i \in [k]^{s-1} do
 13:
                                     Update distribution W_i \leftarrow (W_i | \pi_{< i}(W_i) = \hat{V}_{\lceil < i \mid i \rceil})
 14:
                       for all i \in [k]^{s-1} do
 15:
                              \hat{Z}_{\lceil \cdot, i \rceil} \leftarrow V_{\lceil \cdot, i \rceil} \cdot M^{-1}
 16:
                       Return \hat{Z}, O. \hat{U}^{F}
 17:
```

Analogously to Definition A.2, we define a similar approximate version of FAST-DECODER:

Definition A.7. We will term an algorithm that runs an Fast-Decoder where the algorithm gets an estimate of the probabilities $\Pr(Z = x)$ to within an additive error of 1/4 for every $x \in \mathbb{F}_q$ a precision-bounded Fast-Decoder.

The FAST-DECODER as described above runs in time $O(n \log n)$, where $n = k^s$ is block length if one assumes infinite precision arithmetic. Furthermore, even a bounded-precision model only requires $O(n \log n)$ operations in the "floating point RAM" model — the model where a non-negative real number $r \in [0,1]$ is represented with two $\ell = O(\log n)$ bit integers a,b as $a \cdot 2^b$ and two such numbers can be added, multiplied, or divided in a single step.

In bit more detail, the above representation is also known as the *Floating point number system* [18, Chapter 2]. Before we go into the details of the runtime analysis of FAST-DECODER, we quickly summarize the relevant properties of the floating point number system.

Floating point number system and floating point RAM model. We recall the definition of the floating point number system:

Definition A.8 ([18], Section 2.1). A floating point number system $F \subset \mathbb{R}$ is a subset of real numbers whose elements have the form

$$y = \pm a \cdot \beta^{e-\Delta},$$

where

- The integer $\beta \ge 2$ is the *base* or *radix*
- The natural number Δ is the *precision*

• The integer *e* is the *exponent* and has the range $e_{\min} \le e \le e_{\max}$ for integers $e_{\min} \le e_{\max}$

• The natural number *a* is the *significand*, and it is assumed that

$$\beta^{\Delta-1} \le a \le \beta^{\Delta} - 1.$$

The representation range of F is given by $[\beta^{e_{\min}-1}, \beta^{e_{\max}}(1-\beta^{-\Delta})]$.

Before we proceed, we note the simplications to the above definition that we use in our model:

Definition A.9. We use the floating point number system from Definition A.8 with the following simplifications/modifications:

- Set $\beta = 2$.
- $\Delta = \ell$.¹⁷
- $e_{\min} = -2^{\ell}$ and $e_{\max} = 2^{\ell}$.

For the rest of this discussion, we will assume the parameters that we have set in Definition A.9. Next, we recall some properties of the floating point number system that we will use as given in our runtime analysis of Fast-Decoder.

Before we present the results, we fix some more notation. For $x \in \mathbb{R}$ falling within the representation range of the floating point system, we will use $\mathrm{fl}(x)$ to denote the closest approximation of x in the floating point system. For any vector $\mathbf{y} \in \mathbb{R}^k$, we will overload notation and use $\mathrm{fl}(\mathbf{y})$ to denote the vector obtained by applying $\mathrm{fl}(\cdot)$ to each component of \mathbf{y} . This leads to the following definition, which defines a crucial quantity that will turn up in our approximation bounds.

Definition A.10. The unit roundoff is defined as

$$u=2^{-\Delta}$$
.

We first recall a bound on the approximation error that the rounding entails:

LEMMA A.11 (REFERENCE [18], THEOREM 2.2). Let $x \in \mathbb{R}$ be in the representation range of the floating point system. Then

fl
$$(x) = (1 + \delta) \cdot x$$
, where $|\delta| < u$.

We will also use $fl(\cdot)$ applied to a formula to denote a result of a floating-point evaluation of this formula. We will use the so-called *standard model* [18, Section 2.2]:

Definition A.12 (Standard Model). The standard model assumes the following precision bounds on binary operations. Given $x, y \in F$ and op $\{+, -, \times, \div\}$, we have

fl
$$(x \text{ op } y) = (x \text{ op } y) \cdot (1 + \delta)$$
 where $|\delta| \le u$,

as long as x op y is in the representation range.

In particular, even if x op y happens to have the exact representation in the floating point number system F, we do not require the result of this floating point operation to be exact.

For the rest of the section, we will assume the standard model in our floating point RAM model. Next, we present a technical lemma that will be useful for us:

LEMMA A.13 (SIMPLE GENERALIZATION OF LEMMA 3.1 IN REFERENCE [18]). Let $\delta_1, \ldots, \delta_n$ be such that $\sum_{i=1}^{n} |\delta_i| < 1$ and let $\rho_i \in \{-1, 1\}$ for all $i \in [n]$. Then we have

$$\prod_{i=1}^{n} (1+\delta_i)^{\rho_i} = 1+\theta,$$

¹⁷Since we are using ℓ bits to represent a.

where

$$|\theta| \le \frac{\sum_{i=1}^{n} |\delta_i|}{1 - \sum_{i=1}^{n} |\delta_i|}.$$

Finally, we present approximation error bounds for computing a bounded-degree rational function, which will be crucial in our runtime analysis of FAST-DECODER:

LEMMA A.14. Let $f(X_1, ..., X_N)$ and $g(X_1, ..., X_N)$ be multi-linear polynomials¹⁸ such that both satisfy the following properties:

- the degree is at most d
- there are at most m monomials
- all the coefficients are non-negative and have exact representation in the floating point number system.

Further, let $\mathbf{x}, \widetilde{\mathbf{x}} \in \mathbb{R}^N_{>0}$, be such that there exists an $\varepsilon > 0$ such that for every $i \in [N]$, we have

$$|\boldsymbol{x}_i - \widetilde{\boldsymbol{x}}_i| \leq \varepsilon \boldsymbol{x}_i,$$

and, moreover, let e_0 be such that all \tilde{x}_i and all coefficients of f, g lie in $[2^{-e_0}, 2^{e_0}]$. Then, assuming

$$4\left(d\cdot\varepsilon + \left(d + \log m\right)\cdot u\right) + 1 \leq \frac{1}{2},\tag{37}$$

$$e_1 := 2(d+1)(e_0+1) + 4\log m + 1 \le 2^{\ell},$$
 (38)

we have that

$$\left| \frac{f(x)}{g(x)} - \text{fl}\left(\frac{f(\widetilde{x})}{g(\widetilde{x})}\right) \right| \le 8 \cdot (d \cdot \varepsilon + (d + \log m + 1) \cdot u) \cdot \frac{f(x)}{g(x)},\tag{39}$$

and, moreover,

$$\left|\log \operatorname{fl}\left(\frac{f(\widetilde{\mathbf{x}})}{g(\widetilde{\mathbf{x}})}\right)\right| \le e_1,\tag{40}$$

where the $(\frac{f(\bar{x})}{g(\bar{x})})$ is computed by using pairwise operations (and paying for approximation error for each such operation as in the standard model).

PROOF. We will compute $(\frac{f(\widetilde{x})}{g(\widetilde{x})})$ by first computing each monomial in $f(\widetilde{x})$ and $g(\widetilde{x})$ and then summing the at-most m values in a depth $\log m$ tree fashion. Finally, we divide $f(\widetilde{x})$ by $g(\widetilde{x})$ to obtain our answer.

For notational convenience for each $i \in [N]$, define ε_i such that $\widetilde{\mathbf{x}}_i = (1 + \varepsilon_i) \cdot \mathbf{x}_i$. Note that we have $|\varepsilon_i| \leq \varepsilon$.

To see the error bound, consider an arbitrary monomial, which we assume WLOG to be $\prod_{i=1}^d X_i$. We compute $\prod_{i=1}^d \widetilde{\boldsymbol{x}}_i$ in the obvious way. It is easy to check that $(\prod_{i=1}^d \widetilde{\boldsymbol{x}}_i) = (\prod_{i=1}^d \widetilde{\boldsymbol{x}}_i) \cdot \prod_{i=1}^{d-1} (1+\delta_i)$, where $|\delta_i| \leq u$. Further, by definition of ε_i , we have

$$\operatorname{fl}\left(\prod_{i=1}^{d} \widetilde{\boldsymbol{x}}_{i}\right) = \left(\prod_{i=1}^{d} \boldsymbol{x}_{i}\right) \cdot \prod_{i=1}^{d} (1 + \delta_{i})(1 + \varepsilon_{i}),$$

where for notational simplicity define $\delta_d = 1$.

To apply the error bounds for the floating point operations, we need to argue that all the results of the multiplications in the computation above are within the representations range. Indeed, since $\tilde{x}_i \geq 2^{-e_0}$, all the intermediate results in the multiplication above are at least

 $^{^{18}}$ The result can be proven for general polynomials as well. However, since we only need the result for multi-linear polynomials and the notation for multi-linear polynomials is slightly cleaner, we stick with the multi-linear case.

11:60 J. Błasiok et al.

 $(\frac{2^{-e_0}}{1+u})^{d+1} \ge 2^{-(e_0+1)(d+1)}$. Similarly, for the upper bound: Since $\widetilde{x_i} \le 2^{e_0}$, all the intermediate results are at most $((1+u)2^{e_0})^{(d+1)} \le 2^{(e_0+1)(d+1)}$, which is assumed to be within the representation range (38).

Now let us consider the computation of $\mathrm{fl}\,(f(\widetilde{x}))$. Let \mathcal{M} be the collection of all subset of size at most d that correspond to the monomials in $f(X_1,\ldots,X_N)$. Then when computing $\mathrm{fl}\,(f(\widetilde{x}))$, for each $S\in\mathcal{M}$, we first compute $\mathrm{fl}\,(\prod_{i\in S}\widetilde{x}_i)$, which satisfies by the above discussion,

$$\hat{m}_S \stackrel{\text{def}}{=} \mathrm{fl}\left(\prod_{i \in S} \widetilde{\boldsymbol{x}}_i\right) = \left(\prod_{i \in S} \boldsymbol{x}_i\right) \cdot \prod_{i \in S} (1 + \delta_i)(1 + \varepsilon_i).$$

Now recall, we need to compute $\sum_{S \in \mathcal{M}} \hat{m}_S$. This in turn adds more error. In particular, if we use the algorithm that computes the sum in a recursive-pairwise manner, then we get that

$$\operatorname{fl}\left(f(\widetilde{\boldsymbol{x}})\right) = \sum_{S \in \mathcal{M}} \widetilde{m}_S,$$

where

$$\widetilde{m}_S = \hat{m}_S \prod_{j=1}^{\log |\mathcal{M}|} \left(1 + \delta_j^{(S)}\right),$$

where each $|\delta_i^{(S)}| \le u$. In other words, we have

$$\widetilde{m}_{S} = \left(\prod_{i \in S} \boldsymbol{x}_{i}\right) \cdot \left(\prod_{i \in S} (1 + \delta_{i})(1 + \varepsilon_{i})\right) \cdot \left(\prod_{j=1}^{\log |\mathcal{M}|} \left(1 + \delta_{j}^{(S)}\right)\right).$$

The above along with Lemma A.13, shows that for every $S \in \mathcal{M}$,

$$\left| \widetilde{m}_{S} - \left(\prod_{i \in S} x_{i} \right) \right| \leq \frac{|S| (\varepsilon + u) + \log m \cdot u}{1 - |S| (\varepsilon + u) + \log m \cdot u} \cdot \left(\prod_{i \in S} x_{i} \right)$$

$$\leq 2 \cdot (d(\varepsilon + u) + \log m \cdot u) \cdot \left(\prod_{i \in S} x_{i} \right),$$

here the first inequality follows from the facts that $|\mathcal{M}| \le m$, $|\delta_i| \le \varepsilon$, $|\delta_j^{(S)}| \le u$ and $|\varepsilon_i| \le \varepsilon$ and the second inequality follows from the fact that $|S| \le d$ and $d(\varepsilon + u) + \log m \cdot u \le \frac{1}{2}$ (which in turn follows from Lemma 37).

Now, using the fact that all cofficients in f(x) are non-negative and have an exact representation in the floating point number system, the above then implies that

$$\left|\operatorname{fl}\left(f(\widetilde{\boldsymbol{x}})\right) - f(\boldsymbol{x})\right| \leq 2 \cdot \left(d(\varepsilon + u) + \log m \cdot u\right) \cdot f(\boldsymbol{x}).$$

By a similar argument, we get

$$|\operatorname{fl}(q(\widetilde{x})) - q(x)| \le 2 \cdot (d(\varepsilon + u) + \log m \cdot u) \cdot q(x).$$

As earlier, to apply the error bounds on the result of each floating point addition in the calculation, we need to ensure that all results of all the intermediate computations are within the representation range. Since we are adding exactly represented non-negative values, the lower bound of the representation range is trivially smaller than any of those intermediate values. The largest intermediate value can appear at the end of the calculation and is upper bounded by $(1+u)^{\log m} m((1+u)2^{e_0})^{d+1}) \leq 2^{(d+1)(e_0+1)+2\log m}$, which is assumed to be in the representation range (38).

Then note that to compute the final answer, we divide $f(f(\widetilde{x}))$ by $f(g(\widetilde{x}))$, which along with Definition A.12, Lemma A.13, and Lemma 37, proves the claimed bound in Equation (39), as desired.

Journal of the ACM, Vol. 69, No. 2, Article 11. Publication date: March 2022.

Moreover, since $|\log \mathrm{fl}(g(\widetilde{x}))| \le (d+1)(e_0+1) + 2\log m$, and similarly for $|\log \mathrm{fl}(f(\widetilde{x}))|$, the quotient satisfy $|\log(\frac{f(\widetilde{x})}{g(\widetilde{x})})| \le 2(d+1)(e_0+1) + 4\log m + 1 = e_1$, proving Equation (40).

Finally, we state the definition of a floating point RAM:

Definition A.15 (Floating Point RAM Model). A floating point RAM works with numbers in the floating point system as in Definition A.9 with $\ell = O(\log n)$ for inputs of size n. Each arithmetic operation in the floating point number system is assumed to take unit time.

We note that in the above, each floating point number can be represented with constant many registers of $O(\log n)$ bits and that each of the basic floating operations translates to constant many operations over constant many registers of $O(\log n)$ bits. In other words, each such floating point operation can be done in O(1) time in the standard RAM model, and this justifies the assumption on floating point operations taking unit time in the above definition.

Runtime analysis of Fast-Decoder. We are now ready to do a runtime analysis of Fast-Decoder.

Lemma A.16. For $n=k^s$, Fast-Decoder runs in $O_{q,k}(n\log n)$ time assuming unit cost infinite precision arithmetic. Furthermore, it can be implemented in a bounded-precision floating point RAM model (of Definition A.15) to compute every intermediate real number to within an additive error of 1/4 in $O_{q,k}(n\log n)$ time, as long as the description of the channel $C_{Y|Z}$ is given in a floating point number system using $O(\log n)$ bits per conditional probability. In other words, bounded-precision Fast-Decoder can also be implemented in $O_{q,k}(n\log n)$ time in the floating point RAM.

PROOF. We first remark that we use a "truth-table" representation for each probability distribution, i.e., we store tables with q and q^k floating point numbers, respectively, to represent a distribution in $\Delta(\mathbb{F}_q)$ and $\Delta(\mathbb{F}_q^k)$, respectively. In other words, each Z_j for each $j \in [k]^s$ is a vector length q, and W_i for each $i \in [k]^{s-1}$ is a vector of length q^k .

Let us separate out the computing on real numbers and the rest. It is easy to see that for a recursive call with $n=k^s$, all the operations that do not involve floating point operations can be done in $O_{q,k}(n)$ time. We also note that Lines 9 and 14 are the only places where we have to perform floating point operations. Further, it can be checked that there are $O_{q,k}(n)$ such operation. Thus, the running time (in both infinite precision setting and floating point RAM model), T(n) of Fast-Decoder satisfies the recurrence $T(n) \leq kT(n/k) + O_{q,k}(n)$, which yields $T(n) = O_{q,k}(n \log n)$.

Finally, we prove the claim on the claimed precision in the floating point RAM model. We note that while our final desired precision is only an additive 1/4, intermediate precision needs to be high, since the precision goes down at each recursive call. More precisely, our goal is to use Lemma A.14 to bound this error. Before we can apply Lemma A.14, we verify that the preconditions of the lemma holds.

As mentioned, Lines 9 and 14 are the only places where we have to perform floating point operations are the only places to perform floating point operations. In particular, the input are the $N=q\cdot k^s$ probability values in Z (denote these N probability values by $\boldsymbol{p}=(p_1,\ldots,p_N)$). Line 9 computes for each of the q^k values in $\boldsymbol{W_i}$ a degree k multi-linear polynomial in k of the N variables (in fact, this polynomial is actually a monomial). Line 14 is where we update the q^k values of $\boldsymbol{W_i}$. In particular, each computed value is a rational function $\frac{f(p)}{g(p)}$, where $f(X_1,\ldots,X_N)$ is still a monomial in k variables and $g(X_1,\ldots,X_N)$ is a multilinear polynomial of degree k with at most q^k monomials each with a coefficient of 1. Note that f and g satisfy the pre-conditions of Lemma A.14.

11:62 J. Błasiok et al.

Now, consider a recursive call to FAST-DECODER with $s \leftarrow s - i$. We first note that we do not have access to \boldsymbol{p} but rather an approximation $\boldsymbol{\widetilde{p}}$ where each entry has an error bounded by $1 \pm \varepsilon_i$, where we define ε_i soon. Moreover, we will maintain the bound e_i on the magnitude of the exponents of the approximations at the ith level of the recursion, namely we shall ensure that on the ith level of recursion for each j we have $|\log \boldsymbol{\widetilde{p}}_i| \le e_i$; the e_i will be defined soon as well.

First, we note that by Lemma A.11, we have that $\varepsilon_0 \leq u$, and $e_0 \leq 2^{O(\log n)}$, since we assumed that the description of the channel is specified using $O(\log n)$ bits. Now applying Lemma A.14 with $d \leftarrow k, m \leftarrow q^k, \varepsilon \leftarrow \varepsilon_i, \mathbf{x} \leftarrow \mathbf{p}$, and $\widetilde{\mathbf{x}} \leftarrow \widetilde{\mathbf{p}}$ from Equation (39) (it can be verified that Lemma 37 will be satisfied with our parameter choice), we get

$$\varepsilon_{i+1} \le 8 (k \cdot \varepsilon_i + (k + k \log q + 1) \cdot u) \le 32 \cdot k \log q \cdot \varepsilon_i$$

where the inequality uses $k \ge 1$ and the fact that ε_i is increasing in i and hence $u \le \varepsilon_0 \le \varepsilon_i$. Thus, we have that

$$\varepsilon_s \leq (32 \cdot k \log q)^s \cdot u$$
.

Similarly, from Equation (40), we get

$$e_{i+1} \le 2(k+1)(e_i+1) + 4k\log q + 1 \le (13k\log q) \cdot e_i$$

and therefore $e_s \le (13k\log q)^s \cdot e_0$. Since $e_i \le e_s$ for each $i \le s$, to ensure condition Lemma 38 in all applications of Lemma A.14, it is enough to pick ℓ such that $(13k\log q)^s e_0 \le 2^{\ell}$, that is,

$$\ell \ge s \cdot (\log 13 + \log k + \log \log q) + \log e_0$$
.

However, note that at any stage the additive error for any probability value calculated by Fast-Decoder is upper bounded by ε_s . Thus, if we pick

$$\ell \ge s \cdot (\log k + \log \log q + 5) + 2$$
,

then we have $\varepsilon_s \leq \frac{1}{4}$ (since $u = 2^{-\ell}$). The proof is complete by noting that if we chose ℓ to be maximum of those two necessary lower bounds bounds, then we have $\ell = O_{k,q}(\log n)$ and hence we indeed are working with a floating point RAM model.

Correctness of FAST-DECODER. With the runtime analysis of FAST-DECODER out of the way, in the next lemma we show that FAST-DECODER is equivalent to the SC-DECODER on the same input. For this lemma, we assume that [n] is equated with $[k]^s$ and elements of $[k]^s$ are enumerated in lex order by SC-DECODER. Also it would be useful to keep Remark A.6 and Remark A.1 in mind.

Lemma A.17. Let Z be a product distribution (where each $Z_i \in \Delta(\mathbb{F}_q)$ is a distribution over \mathbb{F}_q), and let $\alpha \in (\mathbb{F}_q \cup \{\bot\})^{[k]^s}$. For $i \in [k]^s$, let P_i be the quantity defined on Line 5 of SC-Decoder for input $(Z; \alpha)$, and let Q_i be from the output of FAST-DECODER $(Z; \alpha, s)$. Then we have for every $i \in [k]^s$, $P_i = Q_i$ and

$$FAST-DECODER(Z; \alpha) \cdot M^{\otimes s} = SC-DECODER(Z; \alpha).$$

Furthermore, the output of the precision-bounded FAST-DECODER equals the output of an approximate-successive-cancellation decoder on $(Z; \alpha)$.

PROOF. We prove the lemma by induction on s. For s = 0, the lemma is immediate (from line 5 in SC-Decoder and line 4 in Fast-Decoder), so assume the lemma holds for s' < s.

Our proof will compare two sets of variables, \hat{U}^{F} from the definition of Fast-Decoder and \hat{U}^{SC} , which we define next. Given Z, α as in the statement of the lemma, let U be the joint distribution defined by

$$U := ZM^{\otimes s}$$
.

Now define \hat{U}^{SC} such that for all $i \in [k]^s$:

$$\hat{\boldsymbol{U}}_{i}^{SC} = \begin{cases} \operatorname{argmax}_{x \in \mathbb{F}_{q}} \operatorname{Pr} \left(\boldsymbol{U}_{i} = x | \boldsymbol{U}_{< i} = \hat{\boldsymbol{U}}_{< i}^{SC} \right) = \operatorname{argmax}_{x \in \mathbb{F}_{q}} \boldsymbol{X}_{i}(x) & \text{if } \boldsymbol{\alpha}_{i} = \bot \\ \boldsymbol{\alpha}_{i} & \text{if } \boldsymbol{\alpha}_{i} \in \mathbb{F}_{q} \end{cases}$$
(41)

We start by noting that $\hat{U}^{SC} = SC\text{-Decoder}(Z; \alpha)$ (this can be argued, e.g., by induction on i). If $\alpha_i \in \mathbb{F}_q$, then it is easy to check that $\hat{U}_i^F = \hat{U}_i^{SC}$, so for the rest of the proof we will assume this as given and the focus will be on indices i such that $\alpha_i = \bot$. Next, we note that the outputs \hat{Z} and \hat{U}^F of Fast-Decoder are related by the condition $\hat{Z} = \text{Fast-Polar-Encoder}(\hat{U}^F)$. (In particular, Lines 4, 12, and 16 correspond exactly to the code of Fast-Polar-Encoder.) Restated, this implies

$$\hat{Z} \cdot M^{\otimes s} = \hat{U}^{\mathrm{F}}.\tag{42}$$

Thus to prove the lemma, it suffices to prove that $\hat{\boldsymbol{U}}^{\mathrm{F}} = \hat{\boldsymbol{U}}^{\mathrm{SC}}$. To do so, we use the recursive structure of Fast-Decoder and prove that for every $j \in [k]$, $\hat{\boldsymbol{U}}^{\mathrm{F}}_{[j,\cdot]} = \hat{\boldsymbol{U}}^{\mathrm{SC}}_{[j,\cdot]}$. We do so by induction on j.

First, recall that $\hat{U}^{\mathrm{F}}_{[j,\,\cdot]}$ = Fast-Decoder $_{s-1}(Z'^{(j)};\pmb{lpha}_{[j,\,\cdot]})$ with

$$Z'^{(j)} = \{ (Z \cdot M)_{[j,\cdot]} | (Z \cdot M)_{[< j,\cdot]} = \hat{V}_{[< j,\cdot]} \},$$

where the equality follows from Lines 9 and 14. To compare with $\hat{U}^{SC}_{[j,\cdot]}$, we need a inductive structure on \hat{U}^{SC} , and we use a simple property that we describe informally first and then describe in formal notation. Informally, if the input stream to the successive cancellation decoder is split into three parts, the prefix A, the central part B, and the suffix C, then the decoding on the central part is independent of the suffix. Furthermore, the decoding of the central part is the output of the successive cancellation decoder on a modified input that incorporates the conditioning induced by the decoding of the prefix. Formally, the above can be expressed as the following: Let $A \in (\Delta(\mathbb{F}_q))^a$, $B \in \Delta(\mathbb{F}_q)^b$, and $C \in \Delta(\mathbb{F}_q)^c$ and $C \in \mathbb{F}_q \cup \{\bot\}^a$, $C \in \mathbb{F}_q \cup \{\bot\}^a$, $C \in \mathbb{F}_q \cup \{\bot\}^a$, and $C \in \mathbb{F}_q \cup \{\bot\}^a$. If $C \in \mathbb{F}_q \cup \{\bot\}^a$, and $C \in \mathbb{F}_q \cup \{\bot\}^a$, and $C \in \mathbb{F}_q \cup \{\bot\}^a$, and $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, and $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}^a$, where $C \in \mathbb{F}_q \cup \{\bot\}^a$ for $C \in \mathbb{F}_q \cup \{\bot\}$

By the (outer) inductive hypothesis (on s), it suffices to show that $Z^{r(j)} \cdot M^{\otimes s-1}$ is distributed identically 19 to $\tilde{U}^{(j)}$. We now simplify the former. We have

$$Z'^{(j)} \cdot M^{\otimes s-1} = \{ (Z \cdot M^{\otimes s})_{[j,\cdot]} | (Z \cdot M)_{[< j,\cdot]} = \hat{V}_{[< j,\cdot]} \} = \{ U_{[j,\cdot]} | (Z \cdot M)_{[< j,\cdot]} = \hat{V}_{[< j,\cdot]} \},$$
 where the first equality uses the fact that $(Z \cdot M)_{[j,\cdot]} \cdot M^{\otimes s-1} = (Z \cdot M^{\otimes s})_{[j,\cdot]}.$

Comparing with the definition of $\tilde{U}^{(j)} = \{U_{[j,\cdot]} | U_{[< j,\cdot]} = \hat{U}^{\text{SC}}_{[< j,\cdot]} \}$, it thus suffices to show that the conditioning events $(Z \cdot M)_{[< j,\cdot]} = \hat{V}_{[< j,\cdot]}$ and $U_{[< j,\cdot]} = \hat{U}^{\text{SC}}_{[< j,\cdot]}$ are identical. For every $\ell < j$, we have, by applying Equation (42) to the outputs of Fast-Decoder($Z'^{(j)}, \alpha_{[\ell,\cdot]}, s-1$) in Line 12, we have $\hat{V}_{[\ell,\cdot]} \cdot M^{\otimes s-1} = \hat{U}^{\text{F}}_{[\ell,\cdot]}$. Now using (inner) inductive hypothesis on $\ell < j$, we have $\hat{V}_{[\ell,\cdot]} \cdot M^{\otimes s-1} = \hat{U}^{\text{SC}}_{[\ell,\cdot]}$. We use this and the invertibility of $M^{\otimes s-1}$ to rephrase the event $(Z \cdot M)_{[< j,\cdot]} = \hat{V}_{[< j,\cdot]}$ as $(Z \cdot M)_{[< j,\cdot]} \cdot M^{\otimes s-1} = \hat{V}_{[< j,\cdot]} \cdot M^{\otimes s-1} = \hat{U}^{\text{SC}}_{[< j,\cdot]}$. Simplifying the left-hand side, we get

¹⁹Technically, we want $Z'^{(j)}$ to be identically distributed to $\tilde{U}^{(j)}$, but this condition is equivalent, since $M^{\otimes s-1}$ has full rank.

11:64 J. Błasiok et al.

 $(Z \cdot M)_{[< j, \cdot]} \cdot M^{\otimes s-1} = (Z \cdot M^{\otimes s})_{[< j, \cdot]} = U_{[< j, \cdot]}$. Thus we get that the two events are indeed identical and thus yield $\hat{U}_{[i, \cdot]}^{\mathbb{F}} = \hat{U}_{[i, \cdot]}^{SC}$.

The proof that $P_{[j,\cdot]} = Q_{[j,\cdot]}$ for every $j \in [k]$ is completely similar, and we omit the details. Furthermore, note that an equivalent view of Fast-Decoder is that it is an efficient algorithm to compute the Q_i 's, which it then uses to run SC-Decoder. Thus, if a bounded-precision Fast-Decoder computes every entry of P_i to within an additive error of 1/4, then the bounded precision Fast-Decoder implements an approximate-successive-cancellation decoder.

Proofs of Theorem 1.11 and Theorem 1.14. Now we can prove Theorem 1.11 (modulo Claim A.4, which we prove in the next sub-section).

PROOF OF THEOREM 1.11. In the model of infinite precision arithmetic, Theorem 1.11 follows from Theorem A.5 and the equivalence of SC-Decoder and Fast-Decoder from Lemma A.17 with the running time bound following from Lemma A.16.

In the bounded precision case, by Lemma A.17 we have that the bounded-precision Fast-Decoder implements an approximate-successive-cancellation decoder. Applying Theorem A.5 again in this setting, we have that the decoding error probability still remains $O(n\tau \log q)$, and the running time of $O(n \log n)$ from Lemma A.16 is now in the standard floating point RAM model. \square

Finally, Theorem 1.14 is essentially a corollary of Theorem 1.11 and the definition of (exponential) strong polarization.

Proof of Theorem 1.14. Fix some constant c, and take $\gamma < k^{-c-1} \log^{-1} q$, with $n = k^t$. Note that this implies that

$$\gamma^{t} = \frac{1}{(k^{t})^{c+1} \cdot \log^{t} q} = \frac{1}{(k^{t})^{c+1} \cdot \log^{t} q}.$$
 (43)

By the definition of strong polarization property, we know that for some constants β , η , martingale X_t is $(\gamma^t, \gamma^t, \beta \cdot \eta^t)$ -polarizing. Hence, by Theorem 1.11, the corresponding polar code has rate at least

Capacity(
$$C$$
) – $\beta \eta^t - \gamma^t$

for $t = \Theta_{\eta,\beta}(\log(1/\varepsilon))$, and we have $\beta \eta^t + \gamma^t \le \varepsilon$, where the inequality follow from Equation (43) and our choice of t.

The probability of decoding failure is at most

$$n\gamma^t\log q \le n(n)^{-c-1}\log^{-t+1}(q) \le n^{-c},$$

where the first inequality follows from Equation (43).

By the definition of strong polarization property, we know that for some constants β , η , Λ , martingale X_t is $(2^{-2^{\Lambda t}}, \gamma^t, \beta \cdot \eta^t)$ -polarizing. We use the same choice of t as in the strong polarizing case, and using the same argument as in that case we get that the polar code has the claimed rate. The probability of decoding error is at most

$$n\log q \cdot 2^{-2^{\Lambda t}} = n\log q \cdot 2^{-2^{\Lambda \frac{\log n}{\log k}}} = n\log q \cdot 2^{-n^{\frac{\Lambda}{\log k}}} \le 2^{-n^{\beta'}}$$

for some $\beta' = \Omega_{\Lambda,k,q}$ (1), as desired.

A.2.3 Arıkan Martingale and Polar Coding. Here we build a correspondence between the definition of the Arıkan Martingale and the process of polar coding, which was used in the proof of Claim A.4.

Lemma A.18. For a matrix $M \in \mathbb{F}_q^{k \times k}$ and symmetric channel $C_{Y|Z}$, let $\{X_t\}$ be the associated Arikan Martingale. For a given t, let $L = M^{\otimes t}$ be the polarization transform, and let $n = k^t$ be the block length. Let the channel inputs Z_i be i.i.d. uniform in \mathbb{F}_q and channel outputs $Y_i \sim C_{Y|Z}(Z_i)$.

Then, for a uniformly random index $i \in [n]$, the normalized entropy $\overline{H}((ZL)_i \mid (ZL)_{< i}, Y)$ is distributed identically as X_t .

PROOF. Throughout this proof, we will switch to considering vectors in $\mathbb{F}_q^{k^t}$ as tensors in $(\mathbb{F}_q^k)^{\otimes t}$, for convenience—this correspondence is induced by lexicographic ordering \prec on tuples $[k]^t$. Also, we will write H(Z) to mean the operator H acting on Z. More specifically, for a linear map defined by matrix H, we use H(Z) = ZH. In this notation, we wish to show that the distribution of X_t is identical to $\overline{H}((M^{\otimes t}(Z))_i \mid Y, (M^{\otimes t}(Z))_{\prec i})$ for a uniformly random multiindex $i \in [k]^t$.

We will show by induction that for all t, there is some permutation of coordinates $^{20} \sigma' : [k]^t \to [k]^t$ such that the joint distributions

$$\{(A', B')\}_{(A', B') \sim D_t} \equiv \{(M^{\otimes t}(Z), \sigma'(C(Z)))\}_{Z \sim (\mathbb{F}_a^k)^{\otimes t}}, \tag{44}$$

where $(A', B') \sim D_t$ are the distributions defined in the tth step of the Arıkan martingale, and $Z \sim (\mathbb{F}_q^k)^{\otimes t}$ is sampled with i.i.d. uniform coordinates. This is sufficient, because a permutation of the channel outputs does not affect the relevant entropies. That is,

$$\overline{H}(A'_i \mid A'_{\prec i}, B') = \overline{H}(A'_i \mid A'_{\prec i}, \sigma'(B')).$$

First, the base case t=0 follows by definition of the distribution D_0 in the Arıkan martingale (and the fact that $M(Z_1) \sim \mathbb{F}_q$).

For the inductive step, assume the claim holds for t-1. Let σ be the permutation guaranteed for t-1. For each $j \in [k]$, sample an independent uniform $Z^{(j)} \sim (\mathbb{F}_q^k)^{\otimes t-1}$ and define

$$(\mathbf{A}^{(j)}, \mathbf{B}^{(j)}) := (M^{\otimes t-1}(\mathbf{Z}^{(j)}), \ \sigma(C(\mathbf{Z}^{(j)}))).$$
 (45)

By the inductive hypothesis, $(\mathbf{A}^{(j)}, \mathbf{B}^{(j)}) \sim D_{t-1}$, for each $j \in [k]$.

As in the Arıkan martingale, define (A', B') deriving from $\{(A^{(j)}, B^{(j)})\}_{j \in [k]}$ as

$$A'_{[i,\cdot]} := M((A_i^{(1)}, \dots, A_i^{(k)}))$$
 and $B'_{[j,\cdot]} := B^{(j)}$. (46)

Note that B' can equivalently be written (unwrapped) as

$$B' := (B^{(1)}, B^{(2)}, \dots, B^{(k)}).$$

By definition of the Arıkan martingale, we have $(A', B') \sim D_t$.

Finally, define $Z \in (\mathbb{F}_q^k)^{\otimes t}$ by

$$Z_{[\cdot,j]} := Z^{(j)}. \tag{47}$$

To finish the proof, we will show that $(A', B') = (M^{\otimes t}(Z), \sigma'(C(Z)))$ for some permutation σ' . The main claim is the following.

CLAIM A.19. For every instantiation of the underlying randomness in Z, we have

$$A' = M^{\otimes t}(Z).$$

²⁰This is in fact just a reversal of the coordinates, i.e., $\sigma'((i_1, i_2, \dots i_t)) = (i_t, \dots, i_2, i_1)$.

11:66 J. Błasiok et al.

Proof of Claim A.19. Expanding the recursive definition of the tensor product, Equation (2), we have

$$[M^{\otimes t}(Z)]_{[i,\cdot]} = M((W_i^{(1)}, W_i^{(2)}, \dots W_i^{(k)})),$$

where

$$W^{(j)} := M^{\otimes t-1}(Z_{[\cdot,j]}) = M^{\otimes t-1}(Z^{(j)}) = A^{(j)}.$$

Here the last equality is by the inductive assumption. Thus,

$$[M^{\otimes t}(Z)]_{[i,\cdot]} = M((A_i^{(1)}, \dots, A_i^{(k)}))$$

$$= A'_{[i,\cdot]}.$$
 (By definition, given in Equation (46))

And so $M^{\otimes t}(Z) = A'$ as desired.

Continuing the proof of Lemma A.18, we now have

$$(A', B') = (A', (B^{(1)}, B^{(2)}, \dots, B^{(k)}))$$

$$(\sigma(C(Z^{(1)})), \sigma(C(Z^{(2)})), \dots, \sigma(C(Z^{(k)}))) \qquad \text{(Definition of sampling, Equation (45))}$$

$$= (A', \sigma'(C(Z))) \qquad (\star)$$

$$= (M^{\otimes t}(Z), \sigma'(C(Z))). \qquad \text{(Claim A.19)}$$

In the above, the equality in line (\star) follows by taking σ' to be the permutation that sorts $[k]^t$ in the order of *least significant symbol* first (based on our definition in Equation (47)) and then sorts each group (thought of as $[k]^{t-1}$ in the natural way) recursively according to σ . Unwinding this recursion, one can see that σ' is in fact the *symbol-reversal* permutation on $[k]^t$.

This establishes the equivalence of the distributions claimed in Equation (44) and completes the proof. \Box

REFERENCES

- [1] Emmanuel Abbe and Emre Telatar. 2012. Polar codes for the *m*-user multiple access channel. *IEEE Trans. Inf. Theory* 58, 8 (2012), 5437–5448.
- [2] Erdal Arıkan. 2009. Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels. IEEE Trans. Inf. Theory 55, 7 (July 2009), 3051–3073.
- [3] Erdal Arıkan and Emre Telatar. 2009. On the rate of channel polarization. In *Proceedings of 2009 IEEE International Symposium on Information Theory*. 1493–1495.
- [4] Jaroslaw Blasiok, Venkatesan Guruswami, Preetum Nakkiran, Atri Rudra, and Madhu Sudan. 2018. General strong polarization. In Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing (STOC'18), Ilias Diakonikolas, David Kempe, and Monika Henzinger (Eds.). ACM, 485–492. https://doi.org/10.1145/3188745.3188816
- [5] Jaroslaw Blasiok, Venkatesan Guruswami, and Madhu Sudan. 2018. Polar codes with exponentially small error at finite block length. In Proceedings of the Annual Conference on Approximation, Randomization, and Combinatorial Optimization and Algorithms and Techniques, (APPROX/RANDOM'18), Eric Blais, Klaus Jansen, José D. P. Rolim, and David Steurer (Eds.), LIPIcs, Vol. 116. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 34:1–34:17. https://doi.org/ 10.4230/LIPIcs.APPROX-RANDOM.2018.34
- [6] Thomas M. Cover and Joy A. Thomas. 2005. Elements of Information Theory (2nd ed.). John Wiley & Sons, Hoboken, NJ.
- [7] Eren Şaşoğlu. 2012. Polarization and polar codes. Found. Trends Commun. Inf. Theory 8, 4 (2012), 259–381. https://doi.org/10.1561/0100000041
- [8] Eren Şaşoğlu, Emre Telatar, and Edmund M. Yeh. 2013. Polar codes for the two-user multiple-access channel. IEEE Trans. Inf. Theory 59, 10 (2013), 6583–6592.
- [9] Rick Durrett. 2019. Probability: Theory and Examples (5th ed.). Cambridge University Press.
- [10] Arman Fazeli, Hamed Hassani, Marco Mondelli, and Alexander Vardy. 2021. Binary linear codes with optimal scaling: Polar codes with large kernels. IEEE Trans. Inf. Theory 67, 9 (2021), 5693–5710. https://doi.org/10.1109/TIT.2020. 3038806

- [11] Naveen Goela, Emmanuel Abbe, and Michael Gastpar. 2013. Polar codes for broadcast channels. In *Proceedings of the IEEE International Symposium on Information Theory*. 1127–1131.
- [12] Dina Goldin and David Burshtein. 2014. Improved bounds on the finite length scaling of polar codes. *IEEE Trans. Inf. Theory* 60, 11 (2014), 6966–6978.
- [13] Venkatesan Guruswami, Andrii Riazanov, and Min Ye. 2020. Arikan meets Shannon: Polar codes with near-optimal convergence to channel capacity. In *Proceedings of the 52nd ACM Symposium on Theory of Computing (STOC'20)*. 552–564. https://doi.org/10.1145/3357713.3384323
- [14] Venkatesan Guruswami, Atri Rudra, and Madhu Sudan. March 15, 2019. Essential Coding Theory. Retrieved June 13, 2021 from https://cse.buffalo.edu/faculty/atri/courses/coding-theory/book/index.html.
- [15] Venkatesan Guruswami and Ameya Velingker. 2015. An entropy sumset inequality and polynomially fast convergence to shannon capacity over all alphabets. In *Proceedings of the 30th Conference on Computational Complexity*. 42–57.
- [16] Venkatesan Guruswami and Patrick Xia. 2015. Polar codes: Speed of polarization and polynomial gap to capacity. *IEEE Trans. Inf. Theory* 61, 1 (2015), 3–16. Preliminary version in Proc. of FOCS 2013.
- [17] Seyed Hamed Hassani, Kasra Alishahi, and Rüdiger L. Urbanke. 2014. Finite-length scaling for polar codes. IEEE Trans. Inf. Theory 60, 10 (2014), 5875–5898. https://doi.org/10.1109/TIT.2014.2341919
- [18] Nicholas J. Higham. 2002. Accuracy and Stability of Numerical Algorithms (2nd ed.). SIAM. https://doi.org/10.1137/1. 9780898718027
- [19] Satish Babu Korada. 2010. Polar codes for Slepian-Wolf, Wyner-Ziv, and Gelfand-Pinsker. In *Proceedings of the IEEE Information Theory Workshop*. 1–5.
- [20] Satish Babu Korada, Eren Şaşoğlu, and Rüdiger L. Urbanke. 2010. Polar codes: Characterization of exponent, bounds, and constructions. *IEEE Trans. Inf. Theory* 56, 12 (2010), 6253–6264.
- [21] Satish Babu Korada, Andrea Montanari, Emre Telatar, and Rüdiger L. Urbanke. 2010. An empirical scaling law for polar codes. In Proceedings of the IEEE International Symposium on Information Theory. 884–888.
- [22] Hessam Mahdavifar and Alexander Vardy. 2011. Achieving the secrecy capacity of wiretap channels using polar codes. IEEE Trans. Inf. Theory 57, 10 (2011), 6428–6443.
- [23] Marco Mondelli, S. Hamed Hassani, and Rüdiger L. Urbanke. 2016. Unified scaling of polar codes: Error exponent, scaling exponent, moderate deviations, and error floors. IEEE Trans. Inf. Theory 62, 12 (2016), 6698–6712. https://doi.org/10.1109/TIT.2016.2616117
- [24] Ryuhei Mori and Toshiyuki Tanaka. 2014. Source and channel polarization over finite fields and reed-solomon matrices. *IEEE Trans. Inf. Theory* 60, 5 (2014), 2720–2736.
- [25] Henry D. Pfister and R\u00fcdiger L. Urbanke. 2019. Near-optimal finite-length scaling for polar codes over large alphabets. IEEE Trans. Inf. Theory 65, 9 (2019), 5643-5655. https://doi.org/10.1109/TIT.2019.2915595
- [26] M. S. Pinsker. 1964. Information and Information Stability of Random Variables and Processes. Holden-Day.
- [27] Volker Strassen. 1962. Asymptotische abschatzungen in shannon's informationstheories. In *Transactions of the 3rd Prague Conference on Information Theory*. 689–723.
- [28] Ido Tal and Alexander Vardy. 2013. How to construct polar codes. IEEE Trans. Inf. Theory 59, 10 (Oct 2013), 6562-6582.
- [29] Ido Tal and Alexander Vardy. 2015. List decoding of polar codes. IEEE Trans. Inf. Theory 61, 5 (2015), 2213-2226.
- [30] Hsin-Po Wang and Iwan M. Duursma. 2021. Polar codes' simplicity, random codes' durability. IEEE Trans. Inf. Theory 67, 3 (2021), 1478–1508. https://doi.org/10.1109/TIT.2020.3041570
- [31] Lele Wang and Eren Şaşoğlu. 2014. Polar coding for interference networks. In *Proceedings of the IEEE International Symposium on Information Theory*. 311–315. https://doi.org/10.1109/ISIT.2014.6874845
- [32] Jacob Wolfowitz. 1957. The coding of messages subject to chance errors. Illinois J. Math. 1 (1957), 591-606.

Received April 2019; revised June 2021; accepted October 2021