


Accelerating Polarization via Alphabet Extension

Iwan Duursma   

University of Illinois Urbana-Champaign, IL, USA

Ryan Gabrys   

University of California San Diego, CA, USA

Venkatesan Guruswami   

University of California, Berkeley, CA, USA

Ting-Chun Lin  

University of California San Diego, CA, USA

Hon Hai (Foxconn) Research Institute, Taipei, Taiwan

Hsin-Po Wang   

University of California San Diego, CA, USA

Abstract

Polarization is an unprecedented coding technique in that it not only achieves channel capacity, but also does so at a faster speed of convergence than any other technique. This speed is measured by the “scaling exponent” and its importance is three-fold. Firstly, estimating the scaling exponent is challenging and demands a deeper understanding of the dynamics of communication channels. Secondly, scaling exponents serve as a benchmark for different variants of polar codes that helps us select the proper variant for real-life applications. Thirdly, the need to optimize for the scaling exponent sheds light on how to reinforce the design of polar code.

In this paper, we generalize the binary erasure channel (BEC), the simplest communication channel and the protagonist of many polar code studies, to the “tetrahedral erasure channel” (TEC). We then invoke Mori–Tanaka’s 2×2 matrix over \mathbb{F}_4 to construct polar codes over TEC. Our main contribution is showing that the dynamic of TECs converges to an almost–one-parameter family of channels, which then leads to an upper bound of 3.328 on the scaling exponent. This is the first non-binary matrix whose scaling exponent is upper-bounded. It also polarizes BEC faster than all known binary matrices up to 23×23 in size. Our result indicates that expanding the alphabet is a more effective and practical alternative to enlarging the matrix in order to achieve faster polarization.

2012 ACM Subject Classification Mathematics of computing → Coding theory; Theory of computation → Error-correcting codes

Keywords and phrases polar code, scaling exponent

Digital Object Identifier 10.4230/LIPIcs.APPROX/RANDOM.2022.17

Category RANDOM

Related Version *Full Version:* <https://arxiv.org/abs/2207.04522> [14]

Funding *Ryan Gabrys:* NSF CCF-2107346.

Venkatesan Guruswami: NSF CCF-2210823 and Simons Investigator Award.

Hsin-Po Wang: NSF CCF-1764104.

1 Introduction

A fundamental question at the center of the theory of communication is whether we can fully utilize a noisy channel to transmit information. In modern terminology, can error correcting codes achieve channel capacity? The answer is positive; in fact, multiple code constructions do so. Among them, polar code is a special one as it achieves capacity faster than any other known code.



© Iwan Duursma, Ryan Gabrys, Venkatesan Guruswami, Ting-Chun Lin, and Hsin-Po Wang; licensed under Creative Commons License CC-BY 4.0

Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2022).

Editors: Amit Chakrabarti and Chaitanya Swamy; Article No. 17; pp. 17:1–17:15



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

17:2 Accelerating Polarization via Alphabet Extension

Polar coding was invented by Arikan around 2008 [2]. During that time, Arikan was experimenting with channel combining and splitting. By treating two independent binary channels as a single quaternary channel (combining) and tasking ourselves with guessing certain linear combinations of the inputs (splitting), he synthesized two channels, denoted by W^{\sqcup} and W° , out of the original channel W . Arikan realized that, when combining and splitting is applied recursively, the channels undergo an intriguing dynamic that ultimately results in most synthetic channels being either almost noiseless or extremely noisy. This is *channel polarization*, the first ingredient underlying polar codes.

The second ingredient of polar codes, also given by Arikan in said seminal paper, is the relation between the dynamic of synthetic channels and the construction and performance of the code. Arikan's insight was that synthetic channels that become almost noiseless can be used to transmit information bits, and synthetic channels that become extremely noisy can be "frozen" to some fixed values. The rate at which we communicate meaningful bits is then the proportion of synthetic channels that are almost noiseless. So, whether we can achieve channel capacity becomes a problem of counting the numbers of good and bad synthetic channels.

It then became apparent, perhaps even appealing, that one can study the dynamic of synthetic channels by means of stochastic processes. Take the *binary erasure channel* (BEC) as an example. Let W be $\text{BEC}(\varepsilon)$, the BEC with erasure probability ε , where $0 < \varepsilon < 1$. The channels W^{\sqcup} and W° are $\text{BEC}(2\varepsilon - \varepsilon^2)$ and $\text{BEC}(\varepsilon^2)$, respectively. A process $\{H_n\}_n$ is thus defined by having $H_0 := \varepsilon$ and $H_{n+1} := 2H_n - H_n^2$ or H_n^2 with equal probability. It can be shown that if

$$\mathbb{P}\{H_n \leq f(n)\} = 1 - H_0 - g(n),$$

where $f, g > 0$ are functions in n , then there is a polar code of length 2^n , miscommunication probability $2^n f(n)$, and gap to capacity $g(n)$.

It was at this point that the study of polar codes branched. On one branch, called the *error exponent regime*, g is a constant and the asymptotics of f is examined. On the other branch, called the *scaling exponent regime*, f is a constant¹ and the asymptotics of g is examined. On the error exponent branch, it was shown that $f(n)$ is roughly $\exp(-e^{\beta n})$, where $\beta > 0$ is a constant depending on the matrix used in the code construction. The task of determining β for each matrix has been fully resolved; interested readers are referred to [3, 26, 22, 33].

On the scaling exponent branch, making progress is harder and slower. For BECs, [21, 25] managed to estimate that $g(n) \approx 2^{-n/3.527}$. For binary memoryless symmetric (BMS) channels, it was first shown that $g(n) < 2^{-n/\mu}$ for some constant $0 < \mu < \infty$ [20]. This makes polar codes the only known code family that converges to capacity at a polynomial rate in the block length. More realistic estimates of μ were given later: $3.553 < \mu$ [18], $3.579 < \mu < 6$ [23], $\mu < 5.702$ [17], $\mu < 4.714$ [31], and very recently $\mu < 4.63$ [49]. Now that we know the μ for polar code and the optimal value being $\mu \approx 2$ for random code [4, 24, 37], the discrepancy begs the question: Can one modify polar code to reach a smaller scaling exponent?

¹ Not always; sometimes $f \rightarrow 0$ but only exponentially fast in n . Note that $2^n f(n)$, the upper bound on the miscommunication probability, is allowed to exceed 1, so the corresponding code can be meaningless. Yet the asymptotics of g capture the behaviors of other meaningful codes.

The answer is positive: Arikan used the matrix $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ (this is called the *kernel* in literature) to combine and split channels. Instead, one can use a larger matrix, for instance

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix},$$

to combine and split channels. In [16, 50, 46, 45, 5, 28], binary matrices ranging from 3×3 to 64×64 are deployed and the scaling exponents over BECs are estimated. The best scaling exponent up to every matrix size is plotted in Figure 2. There are also meta-asymptotic results stating that $\mu \approx 2$ can be achieved using larger and larger matrices. This statement was proved over q -ary erasure channels [36], binary erasure channels [15], all BMS channels [19], and finally discrete memoryless channels [48].

As much as we want to lower polar code's scaling exponent, there is one caveat that renders large matrices impractical: the smallest matrix whose scaling exponent is strictly better than $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ is the 8×8 matrix above. Using this matrix takes twice more time to decode (estimate based on the method of [9]), whereas the benefit we gain is that μ slightly decreases from 3.627 to 3.577. As the matrix gets larger and deviates more from the tensor powers of $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$, the time complexity grows drastically. For this reason, it is unlikely that we will ever see polar code based on large matrices (unless it is for other concerns [6]).

Large matrix aside, many other techniques emerge with empirical evidence that they improve polar code – concatenation, cyclic redundancy check, and list decoder to name a few. But none of them sees a proof of improvements in the scaling exponent; in fact, quite the opposite was reported [30]. So we are back to the starting point where we want to improve polar codes' scaling exponent while minimizing the complexity penalty.

One approach that seems promising, albeit very little is known due to its innate technical difficulty, is to use a non-binary input alphabet. This line of research started from Şaşoğlu [42, 41, 13], wherein the goal was to find at least one way to polarize arbitrary finite alphabets regardless of the speed. In particular, the usual matrix $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ is known to polarize prime fields. Later, Sahebi–Pradhan [40] and Park–Barg [35] showed that $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ cannot polarize non-prime fields. Then, Mori–Tanaka [33] classified all matrices that can polarize finite fields (i.e., the alphabet size must be a prime power). One step forward, Nasser [34] classified all binary operators (i.e., bivariate functions) that can polarize arbitrary finite alphabets. In [7, 8], the authors showed that, for any polarizing matrix over prime fields, one has $\mu < \infty$. In [48], the authors showed that $\mu \approx 2$ is reachable over arbitrary finite alphabets.

Why is a non-binary input alphabet attractive? There are at least three reasons. First, modulation²: For quadrature amplitude modulation (QAM) and amplitude and phase-shift keying (APSK), a constellation point is more likely to be confused with constellation points nearer to it. A non-binary channel models this proximity relation more naturally than a series of correlated binary channels do [44, 11]. Second, two-stage polarization: If we weakly-polarize a binary channel with $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$, treat two binary channels as a quaternary channel, and strongly-polarize the quaternary channel with the 4×4 Reed–Solomon matrix, we can improve the asymptotics of $f(n)$ from $\exp(-2^{0.5n})$ to $\exp(-2^{0.5731n})$ [38] (see also [1, 12]). Third, and most importantly, scaling exponent: Several works have observed that non-binary matrices of the form $\begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}$ just polarize faster than $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ [51, 29, 43]. Could it be that the non-binary scaling exponents are smaller?

² Modulation means translating digital signals to analog signals. A digital signal will be mapped to a point on the complex plane, which represents a sine wave with certain amplitude and phase; such a point is called a constellation point, the union of all points a constellation diagram.

Consider [39]’s technique that uses $\begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}$ to polarize non-binary channels; their result has an implication that non-binary channels’ scaling exponent is at least as good as binary channels’. In this paper, we aim to answer the question of whether the former is strictly better than the latter. By defining a toy model that contains a pair of BECs as a special case and estimating the scaling exponent of $\begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}$, we provide a proof of concept result that an expansion in alphabet size does result in an improvement in scaling exponent. Recall that BECs form a one-parameter family and that this property makes its scaling behavior easy to analyze. This paper’s overall strategy is to show that the descendants of a quaternary channel converge to an almost-one-parameter family; we then analyze the scaling behavior of this family and conclude the following.

► **Theorem 1 (main theorem).** *Treating a pair of BECs as a quaternary channel, the 2×2 matrix $\begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}$ over \mathbb{F}_4 induces a scaling exponent less than 3.451. Here, $\omega^2 + \omega + 1 = 0$.*

This paper is organized as follows. Section 2 reviews the essence of polar code. Section 3 defines tetrahedral erasure channels (TECs), defines balanced TECs to be those that possess some symmetry, and defines edge-heavy TECs to be those that will be polarized faster. Section 4 defines serial combination and parallel combination that will be used to polarize TECs. Section 5 shows that unbalanced TECs tend to become very close to balanced TECs, so it suffices to consider the speed of polarization of the latter. Section 6 shows that edge-light TECs tend to become very close to edge-heavy TECs, so it suffices to consider the speed of polarization of the latter. Section 7 shows the speed of polarization of a generic TEC is faster than the classical BEC.

2 Polar Code

Readers who are familiar with polar code can safely skip this section. This section serves a simplified, high-level summary of classical polar code. More details are found in [47, Chapter 2]. We assume BEC throughout the section.

Let $X \in \mathbb{F}_2$ be a random variable following the uniform distribution. Let $Y \in \mathbb{F}_2 \cup \{?\}$ be a random variable with $\mathbb{P}\{Y = X\} = 1 - \varepsilon$ and $\mathbb{P}\{Y = ?\} = \varepsilon$. Here, $\varepsilon \in [0, 1]$ is called the *erasure probability*. The pair $(X|Y)$ is called a *binary erasure channel* (BEC) and denoted by $\text{BEC}(\varepsilon)$. The entropy $H(\text{BEC}(\varepsilon)) = H(X|Y) = \varepsilon$ is defined through Shannon’s mean.

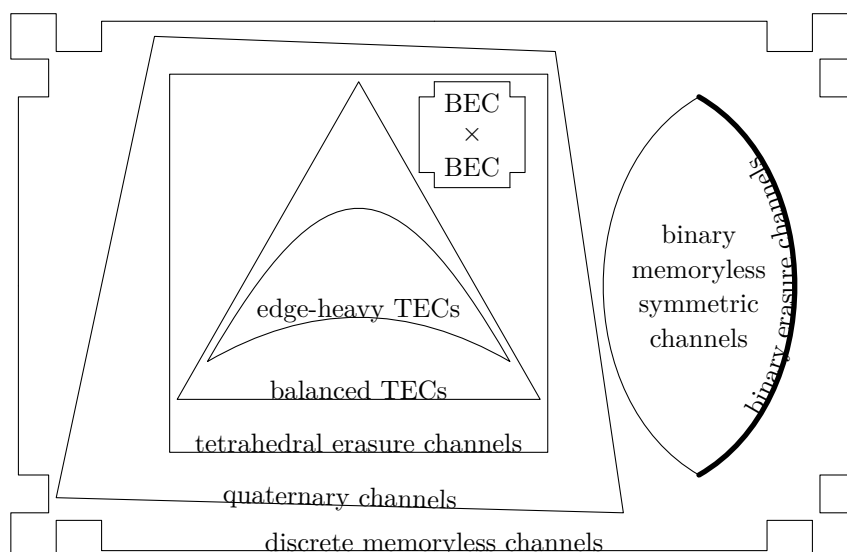
Let $(X_1|Y_1)$ and $(X_2|Y_2)$ be two iid copies of $\text{BEC}(\varepsilon)$. Define the serial combination $\text{BEC}(\varepsilon)^\top$ to be $(X_1 + X_2|Y_1, Y_2)$. That is, what do we know about $X_1 + X_2$ when given Y_1 and Y_2 ? One sees that it is information theoretically equivalent to $\text{BEC}(2\varepsilon - \varepsilon^2)$. Define the parallel combination $\text{BEC}(\varepsilon)^\circ$ to be $(X_1|Y_1, Y_2, X_1 + X_2)$. That is, what do we know about X_1 when given Y_1, Y_2 , and $X_1 + X_2$? One sees that it is information theoretically equivalent to $\text{BEC}(\varepsilon^2)$.

Serial and parallel combinations apply recursively. A polar code of block length 2^n consists of a subset of strings $\mathcal{I} \subseteq \{\top, \circ\}^n$. In this code, a synthetic channel

$$\left(\dots \left((\text{BEC}(\varepsilon)^{c_1})^{c_2} \right) \dots \right)^{c_n} \quad (1)$$

will be used to transmit useful information iff $(c_1, c_1, \dots, c_n) \in \mathcal{I}$. The code rate of this polar code is $|\mathcal{I}|/2^n$. The exact miscommunication probability of this polar code is hard to find, but has an upper bound of

$$\sum_{\mathcal{I}} H \left(\left(\dots \left((\text{BEC}(\varepsilon)^{c_1})^{c_2} \right) \dots \right)^{c_n} \right).$$



■ **Figure 1** The Euler diagram of channels featured in this paper. The cross is the set of pairs of BECs; it will converge to the set of balanced TECs (Section 5). The balanced TECs will then converge to edge-heavy TECs (Section 6). And then edge-heavy TECs polarize faster than BECs. Note that BECs are a one-parameter family of extreme BMS channels, hence the thick curve.

To define a good \mathcal{I} , choose a function $f(n)$ and collect all strings $(c_1, c_1, \dots, c_n) \in \{\Gamma, \circ\}^n$ such that $H(\text{formula (1)})$ is less than $f(n)$. The fact that the erasure probabilities undergo simple evolutions $\varepsilon \mapsto 2\varepsilon - \varepsilon^2$ and $\varepsilon \mapsto \varepsilon^2$ motivates the following stochastic process: define $\{H_n\}_n$ by initial value $H_0 := \varepsilon$ and evolution rule $H_{n+1} := 2H_n - H_n^2$ or H_n^2 with equal probability. Then the code rate $|\mathcal{I}|/2^n$ coincides with $\mathbb{P}\{H_n \leq f(n)\}$. The gap to capacity $g(n) := 1 - H_0 - |\mathcal{I}|/2^n = 1 - H_0 - \mathbb{P}\{H_n \leq f(n)\}$ is thus motivated.

In a way, the study of polar code over BEC is the study of the cdf of H_n , with emphasis put on the hard threshold at $1 - H_0$. Abusing the same logic, this paper is a study of a stochastic process $\{W_n\}_n$ that lives in $[0, 1]^5 \cap \{p + q + r + s + t = 1\}$, which happens to have peculiar implications in coding theory.

3 A New Channel Model

We are to define a type of quaternary channels in this section. This should be the smallest possible set of quaternary channels that meet the following: (a) it should model a pair of BECs as a special case; and (b) it should be closed under pre-processing the input using invertible linear transformations.

3.1 Tetrahedral erasure channel

Let the input alphabet be \mathbb{F}_2^2 ; and we assume the uniform input distribution throughout the paper. For any input $(x_1, x_2) \in \mathbb{F}_2^2$, the output will be in $(\mathbb{F}_2 \cup \{?\})^3$ and assume one of the following five erasure patterns:

- $(x_1, x_1 + x_2, x_2)$ with probability p ;
- $(x_1, ?, ?)$ with probability q ;
- $(?, x_1 + x_2, ?)$ with probability r ;
- $(?, ?, x_2)$ with probability s ;
- $(?, ?, ?)$ with probability t .

Here we call p, q, r, s, t the *subspace erasure probabilities* and they sum to 1. Such a channel is denoted by $\text{TEC}(p, q, r, s, t)$. For brevity, we say a TEC outputs (x_1, x_2) , outputs x_1 , outputs $x_1 + x_2$, outputs x_2 , and outputs nothing to represent the five erasure patterns.

A TEC can be related to a tetrahedron whose vertices are $(0, 0, 0)$, $(1, 1, 0)$, $(1, 0, 1)$, and $(0, 1, 1)$. Outputting (x_1, x_2) corresponds to the vertex $(x_1, x_1 + x_2, x_2)$. Outputting x_1 corresponds to the edge $(x_1, x_1, 0) - (x_1, 1 - x_1, 1)$. Outputting nothing corresponds to the tetrahedron per se. That is to say, a TEC takes a vertex as an input and outputs the same vertex with probability p , outputs an edge attached to that vertex with probability $q + r + s$, and output the entire tetrahedron with probability t .

There is another way to interpret a TEC. Consider \mathbb{F}_4 and let ω be a primitive element therein. A TEC takes $x := x_1\omega + x_2 \in \mathbb{F}_4$ as an input and outputs x , $\text{tr}(x)$, $\text{tr}(\omega x)$, $\text{tr}(x/\omega)$, or nothing, each with probability p, q, r, s , and t . Here, $\text{tr}: \mathbb{F}_4 \rightarrow \mathbb{F}_2$ is the field trace. It is the matrix trace if we use the matrices $\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$, $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, $\begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$, $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$ to represent $0, 1, \omega, 1 + \omega \in \mathbb{F}_4$.

TEC is not an ad hoc channel that we happen to know how to deal with. It relates to other channels that have been discussed in literature.

► **Proposition 2.** *The “ q -ary erasure channel with erasure probability ε ” [32, 36], when $q = 4$, is a TEC of the form $\text{TEC}(1 - \varepsilon, 0, 0, 0, \varepsilon)$.*

► **Proposition 3.** *When transmitting two bits x_1 and x_2 through $\text{BEC}(\delta)$ and $\text{BEC}(\varepsilon)$, respectively, the outputs can be simulated by $\text{TEC}((1 - \delta)(1 - \varepsilon), (1 - \delta)\varepsilon, 0, \delta(1 - \varepsilon), \delta\varepsilon)$.*

The proofs are trivial. The propositions imply that any scaling exponent estimate for TEC immediately generalizes to 4-ary erasure channels and BECs.

3.2 Channel functionals

The *conditional entropy* (hereafter *entropy*) of a TEC is defined by the following; it is meant to be compatible with Shannon’s definition:

$$H(\text{TEC}(p, q, r, s, t)) := \frac{q + r + s}{2} + t.$$

The *edge mass* of a TEC is defined by the following; it measures the “polarizability” of a TEC:

$$E(\text{TEC}(p, q, r, s, t)) := q + r + s.$$

The *Quetelet index* of a TEC W is defined by

$$Q(W) := \frac{E(W)}{H(W)(1 - H(W))}.$$

Clearly, $0 \leq E(W) \leq 2 \min(H(W), 1 - H(W))$ and $0 \leq Q(W) \leq 4$. We call a TEC W *edge-heavy* if $Q(W) \geq 2\sqrt{7} - 4$. Adolphe Quetelet invented the body mass index that determines if a person is overweight or underweight. Here, we use Quetelet index to determine if a TEC possesses too much edge mass (easy to polarize) or too little (hard to polarize).

A TEC is *balanced* if $q = r = s$. Put it another way, the edges of the tetrahedron weigh the same. It is not hard to see that H and E uniquely determine a balanced TEC by

$$p = 1 - H(W) - \frac{E(W)}{2}, \quad q = r = s = \frac{E(W)}{3}, \quad t = H(W) - \frac{E(W)}{2}.$$

The *moment of inertia* of a TEC is defined by

$$A(\text{TEC}(p, q, r, s, t)) := (q - r)^2 + (r - s)^2 + (s - q)^2.$$

A TEC is balanced iff its moment of inertia vanishes. See also the “symmetric over the product” condition in [10] and the “equidistance” condition in [42].

4 Channel Synthesis

TECs can be serially combined or parallelly combined as in the theory of density evolution [27].

4.1 Serial combination

Let $U := \text{TEC}(p, q, r, s, t)$ and $V := \text{TEC}(p', q', r', s', t')$ be two TECs. The *serial combination* of U and V is defined to be the task of guessing $(u_1 + v_1, u_2 + v_2)$ given the output of inputting (u_1, u_2) into U and the output of inputting (v_1, v_2) into V . Let us go over all 25 erasure patterns that are classified into five scenarios.

Scenario one – U outputs (u_1, u_2) and V outputs (v_1, v_2) : Now we know $(u_1 + v_1, u_2 + v_2)$ in its entirety. This scenario happens with probability pp' .

Scenario two – U outputs u_1 with or without u_2 , and V outputs v_1 with or without v_2 , but either u_2 or v_2 is missing: In this case, we can infer $u_1 + v_1$, but we cannot infer $u_2 + v_2$. So this case feels like $(x_1, x_2) := (u_1 + v_1, u_2 + v_2)$ underwent a TEC and only x_1 went through. The probability that only x_1 went through is $pq' + qq' + qp'$.

Scenario three – U outputs (u_1, u_2) or $u_1 + u_2$, and V outputs (v_1, v_2) or $v_1 + v_2$, but scenario one does not happen: For this case, we know neither $u_1 + v_1$ nor $u_2 + v_2$. But we can infer $(u_1 + v_1) + (u_2 + v_2)$. So this case feels like $(x_1, x_2) := (u_1 + v_1, u_2 + v_2)$ underwent a TEC and only $x_1 + x_2$ went through. The probability that only $x_1 + x_2$ went through is $pr' + rr' + rp'$.

Scenario four – U outputs u_2 with or without u_1 , and V outputs v_2 with or without v_1 , but either u_1 or v_1 is missing: In this case, we can infer $x_2 := u_2 + v_2$ but not $x_1 := u_1 + v_1$. So this case feels like x_1 is erased. This scenario happens with probability $ps' + ss' + sp'$.

Scenario five – U outputs one bit (u_1 or $u_1 + u_2$ or u_2) and V outputs one bit (v_1 or $v_1 + v_2$ or v_2) but the erasure patterns do not match; or at least one of U and V outputs nothing: We cannot infer $u_1 + v_1$ because either u_1 or v_1 is missing. We cannot infer $u_2 + v_2$ because either u_2 or v_2 is missing. We cannot infer $(u_1 + v_1) + (u_2 + v_2)$, either. So this case feels like both $x_1 := u_1 + v_1$ and $x_2 := u_2 + v_2$ are erased, so is $x_1 + x_2$. The probability that we learn nothing about (x_1, x_2) is $(q + r + s)(q' + r' + s') - qq' - rr' - ss' + t + t' - tt'$.

Note that these five scenarios correspond to the five erasure patterns in the definition of TEC. Denote by $U \boxtimes V$ the serial combination of U and V ; it is a TEC with subspace erasure probabilities

$$U \boxtimes V := \text{TEC}(pp', pq' + qq' + qp', pr' + rr' + rp', ps' + ss' + sp', 1 - \text{the four to the left}).$$

4.2 Parallel combination

The *parallel combination* of $U := \text{TEC}(p, q, r, s, t)$ and $V := \text{TEC}(p', q', r', s', t')$ is defined to be the task of guessing (u_1, u_2) given $(u_1 + v_1, u_2 + v_2)$ (the perfect output of $U \boxplus V$), the result of feeding (u_1, u_2) into U , and the result of feeding (v_1, v_2) into V .

Denote by $U \otimes V$ the parallel combination of U and V . One can go over its erasure scenarios like the previous subsection does. For instance, if U outputs u_1 and V outputs $v_1 + v_2$, then we can infer v_1 (using u_1 and $u_1 + v_1$), followed by v_2 (using v_1 and $v_1 + v_2$), and finally u_2 (using v_2 and $u_2 + v_2$); and hence we can completely recover u_1 and u_2 . Details omitted, it can be shown that $U \otimes V$ is a TEC with subspace erasure probabilities

$$U \otimes V := \text{TEC}(1 - \text{the four to the right}, tq' + qq' + qt', tr' + rr' + rt', ts' + ss' + st', tt').$$

Note that there is a duality between $\text{TEC}(p, q, r, s, t)$ and $\text{TEC}(t, s, r, q, p)$ that keeps E as is, maps H to $1 - H$, and swaps parallel and serial combinations. The duality grants us the convenience of proving half of a theorem and the other half follows by symmetry.

4.3 Mori–Tanaka’s twisting kernel

A 2×2 polarization *kernel* K over \mathbb{F}_4 is defined with a “twist” as follows: For a pair of inputs $u, v \in \mathbb{F}_4$, let K be the linear transformation that reads $(u, v) \mapsto (u + \omega v, v)$ or, equivalently,

$$\begin{bmatrix} u & v \end{bmatrix} \mapsto \begin{bmatrix} u & v \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}.$$

This kernel was studied by Mori–Tanaka [33] and is shown to be polarizing. If we treat \mathbb{F}_4 as \mathbb{F}_2^2 , then K reads $((u_1, u_2), (v_1, v_2)) \mapsto ((u_1 + v_1 + v_2, u_2 + v_1), (v_1, v_2))$ or, equivalently,

$$\begin{bmatrix} u_1 & u_2 & v_1 & v_2 \end{bmatrix} \mapsto \begin{bmatrix} u_1 & u_2 & v_1 & v_2 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix},$$

where $u_1, u_2, v_1, v_2 \in \mathbb{F}_2$. The kernel K combines two TECs U and V to synthesize $U \boxtimes (V\omega)$ and $U \otimes (V\omega)$, where $V\omega$ is the channel that multiplies the input by ω before feeding it into V . For brevity, $W \boxtimes (W\omega)$ and $W \otimes (W\omega)$ are denoted by W^\perp and W° , respectively.

Multiplying a TEC by ω behaves like a rotation of order 3 (after all, $\omega^3 = 1$ and it is rotating the tetrahedron). It maps $\text{TEC}(p, q, r, s, t)$ to $\text{TEC}(p, s, q, r, t)$. If W is balanced, rotation does not alter it: $W = W\omega$. If it is not balanced, then the rotation helps mis-match q, r, s so that a large probability is paired with a small probability. More precisely,

$$\text{TEC}(p, q, r, s, t)^\perp := \text{TEC}(p^2, ps + sq + qp, pq + qr + rp, pr + rs + sp, 1 - \text{the other four}),$$

$$\text{TEC}(p, q, r, s, t)^\circ := \text{TEC}(1 - \text{the other four}, ts + sq + qt, tq + qr + rt, tr + rs + st, t^2).$$

Twisting makes it easier to reduce q, r , and s and redistribute the mass to p and t .

4.4 Channel process

For a TEC W , we call W^Γ the *serial-child* of W and W° the *parallel-child* of W . Together, they are the *children* of W . The *descendants* of W are the children of W together with the descendants of the children of W . The *n th-generation* descendants of W are the $(n - 1)$ th-generation descendants of the children of W ; the 0th is W itself.

When W is understood from the context, let W_0 be W . For n a positive integer, let W_n be a random child of W_{n-1} with equal probability.

The common strategy used to estimate the scaling exponent concerns a concave function $\psi: [0, 1] \rightarrow \mathbb{R}$ such that $\psi(0) = \psi(1) = 0$ and is positive elsewhere. With ψ , one finds a $0 < \mu < \infty$ such that

$$\frac{\psi(H(W^\perp)) + \psi(H(W^\circ))}{2\psi(H(W))} \leq 2^{-1/\mu}$$

With this “eigenvalue,” a routine argument [47, Sections 5.8–5.10] will show that

$$\mathbb{P}\{H(W_n) < \exp(-e^{n^{1/3}})\} > 1 - H(W_0) - 2^{-n/\mu}.$$

5 Unbalanced TEC Becomes Balanced

In this section, we argue that TECs undergoing the polarization process tend to become more balanced than before. We do so by showing that the moments of inertia are decreasing.

► **Theorem 4** (uniform loss of inertia). $A(W^\square), A(W^\circ) \leq A(W)(1 - A(W)/3)$ for any TEC W .

A proof of the theorem is in Appendix A.1 of the full version [14]. Now the recurrence relation $A(W_{n+1}) \leq A(W_n)(1 - A(W_n)/3)$ is equivalent to $A(W_{n+1}) - A(W_n) \leq -A(W_n)^2/3$ and analogous to the ordinary differential equation $f'(n) \leq -f(n)^2/3$. Solving it, we get $f(n) = O(1/n)$; analogously, $A(W_n) = O(1/n)$.

► **Corollary 5** (ultimate loss of inertia). Fix a TEC W , then $A(W_n) = O(1/n)$ as $n \rightarrow \infty$.

Another way to look at it is the average decay of $A(W_n)$.

► **Proposition 6** (average loss of inertia). $A(W^\square) + A(W^\circ) \leq A(W)$ for any TEC W .

A proof of the proposition is in Appendix A.2 of the full version [14]. By the proposition, $A(W) \geq A(W^\perp) + A(W^\circ) \geq A(W^{\perp\perp}) + A(W^{\perp\circ}) + A(W^{\circ\perp}) + A(W^{\circ\circ}) \geq A(W^{\perp\perp\perp}) \dots$. Hence the expectation of $A(W_n)$ over all W_n is at most $A(W)/2^n$.

Corollary 5 and Proposition 6 imply that any unbalanced TEC will promptly become very similar to a balanced one.³ The speed of polarization of unbalanced TECs is thus dominated by that of balanced TECs. We now turn to the analysis of the polarization speed of balanced TECs.

6 Balanced TECs Hoard Edge Mass

In this section, we argue that the Quetelet index $Q(W_n) := E(W_n)/H(W_n)(1 - H(W_n))$ of a sufficiently deep descendant is about 1.6. Put another way, there is a “trap” that constrains the relation between $E(W_n)$ and $H(W_n)$.

Recall that a balanced TEC W is edge-heavy if $Q(W) \geq 2\sqrt{7} - 4$. Let $\alpha := 2\sqrt{7} - 4 \approx 1.3$.

► **Theorem 7** (trapping region). If W is balanced and edge-heavy, then its children are edge-heavy.

A proof of the theorem is in Appendix B.1 of the full version [14]. The theorem implies that all descendants of an edge-heavy are edge-heavy. For a TEC that is not edge-heavy, its descendants will still become “edge-heavier” by the following theorem.

► **Theorem 8** (attraction toward the trap). Fix any $\varepsilon > 0$; choose $\delta := 3\varepsilon/8$. Let W be any balanced TEC. We have that $Q(W) \leq \alpha - \varepsilon$ implies $Q(W^\square) \geq Q(W)(1 + H(W)\delta)$ and $Q(W^\circ) \geq Q(W)(1 + (1 - H(W))\delta)$.

A proof of the theorem is in Appendix B.2 of the full version [14]. It is clear that the factors $H(W)$ and $1 - H(W)$ before δ slow down the rate at which $Q(W_n)$ approaches $2\sqrt{7} - 4$, especially when $H(W)$ is close to 0 or 1, respectively. These factors cannot be optimized away. To see why, suppose that $H(W) = x \approx 1$ and $E(W) = y \approx 0$. Then

³ Note that Corollary 5 is a weak statement about every single descendant of W , while Proposition 6 implies a strong statement about $A(W_n)$ averaged over all n th-generation descendants. Only Corollary 5 will be used later.

17:10 Accelerating Polarization via Alphabet Extension

$H(W^\circ)$ is about $x^2 + O(y^2)$ and $E(W^\circ)$ is about $2xy + O(y^2)$. Hence $Q(W^\circ)$ is about $2xy/x^2(1-x^2) \approx y/x(1-x) = Q(W)$. That being the case, we would like to add that TECs whose Quetelet index can hardly be improved are already polarized, so we shall not worry about them. Besides, we can prove uniform attraction using Theorem 8.

► **Theorem 9** (uniform attraction). *Fix any $\varepsilon > 0$. For any balanced TEC W such that $Q(W) \leq \alpha - \varepsilon$, there exists an integer $m > 0$ such that $Q(W_n) \geq Q(W)(1 + \varepsilon/8)$ for all $n \geq m$.*

A proof of the theorem is in Appendix B.3 of the full version [14]. Uniform attraction means that every child is at least making some positive progress toward the trap. Small steps of the descendants accumulate to a giant leap of the family.

► **Corollary 10** (ultimate attraction). *For any $\varepsilon > 0$ and any balanced TEC W such that $Q(W) > 0$, there exists an integer $m > 0$ such that $Q(W_n) \geq \alpha - \varepsilon$ for all $n \geq m$.*

Proof. Apply the uniform attraction theorem repeatedly. Every application improves the Quetelet index by a factor of $1 + (\alpha - Q(W_n))/8$. So after a finite number of applications the Quetelet index can be made $\geq \alpha - \varepsilon$. ◀

To summarize this and the previous section, we have two trends: unbalanced TECs tend to become balanced; and “edge-light” TECs tend to become edge-heavy.

The following proposition is a bound on Quetelet index in the opposite direction.

► **Proposition 11** (attraction on the other side). *Let W be a balanced TEC with $Q(W) \leq 2$. Then $Q(W^\square) \leq 2$ and $Q(W^\circ) \leq 2$.*

Some comments on how to prove this proposition is in Appendix B.4 of the full version [14].

The following proposition gives a tighter trapping region than Theorem 8 and Proposition 11 do. A proof is omitted but similar to those of Theorem 8 and Proposition 11. For the optimal trapping region, see the discussion in Appendix D of the full version [14].

► **Proposition 12.** *Let $f(x) := x(1-x)(1.66 - 0.38x(1-x))$. Then $E(W) \leq f(H(W))$ implies $E(W^\square) \leq f(H(W^\square))$ and $E(W^\circ) \leq f(H(W^\circ))$. Let $g(x) := x(1-x)(2 - 2x(1-x)/3)$. Then $E(W) \geq g(H(W))$ implies $E(W^\square) \geq g(H(W^\square))$ and $E(W^\circ) \geq g(H(W^\circ))$.*

7 Edge-heavy TECs Polarize Faster

Let W be any balanced TEC with a fixed $H(W) = x$ and a variable $E(W) = y$. Then $H(W^\perp) = 2x - x^2 + y^2/12$ is increasing in y and $H(W^\circ) = x^2 - y^2/12$ is decreasing in y .

The monotonicity has two applications. Application one: If we know too little to lower bound $Q(W)$, we will upper bound $H(W^\circ)$ using x^2 . In this case, the speed of polarization is at least $\mu \approx 3.627$, the number induced by the standard polar code. Application two: If we know $Q(W) \geq \alpha$, we will upper bound $H(W^\circ)$ using $x^2 - (\alpha x(1-x))^2/12$. This time, $H(W^\circ)$ and $H(W^\top)$ are more separated so the speed of polarization is strictly better than $\mu \approx 3.627$. Any positive α , not necessarily $2\sqrt{7} - 4$, can improve the scaling. This is demonstrated by the following lemma that uses $9/7$ in place of α .

► **Lemma 13** (eigenfunction and eigenvalue). *Let $\psi(x) := (x(1-x))^{0.697}(5 - \sqrt{x(1-x)})$. For balanced TECs with $Q(W) \geq 9/7$,*

$$\frac{\psi(H(W^\square)) + \psi(H(W^\circ))}{2\psi(H(W))} < 0.818.$$

Comments on how to verify the lemma is in Appendix C of the full version [14].

► **Theorem 14** (main theorem). *Consider a pair of BECs treated as a TEC, or consider any TEC where $pqrst > 0$. The 2×2 matrix $\begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}$ over \mathbb{F}_4 induces a scaling exponent less than 3.451.*

Proof. Two iid copies of $\text{BEC}(\varepsilon)$ can be seen as $W := \text{TEC}((1-\varepsilon)^2, (1-\varepsilon)\varepsilon, 0, \varepsilon(1-\varepsilon), \varepsilon^2)$. If ε is 0 or 1, there is nothing to prove. Suppose $0 < \varepsilon < 1$, then both W^\perp and W° have five positive subspace erasure probabilities. (That is, their “ p, q, r, s, t ” are all positive). The descendants of a TEC with five positive subspace erasure probabilities satisfy the same property. In particular, all descendants have positive Quetelet index.

Let W be a TEC whose descendants all have positive Q . By Corollary 5, it takes W a finite number of generations to become very similar to a balanced TEC. That is, for any $\delta > 0$ there exists an $m > 0$ such that $A(W_m) < \delta$. Although W_n is never balanced, what we proved about balanced TECs still hold for “almost-balanced” TECs up to a diminishing error term. So we may proceed as if W_n is balanced for $n \geq m$.

By Corollary 10, it takes another finite number of generations to become “almost edge-heavy.” In particular, there exists an m' such that $Q(W_{m'}) \geq 9/7$ (note that $9/7 \approx 1.286$ and $2\sqrt{7} - 4 \approx 1.291$).

Before the m' th generation, the eigenvalues of the form

$$\frac{\psi(H(W_n^\Gamma)) + \psi(H(W_n^\circ))}{2\psi(H(W_n))}$$

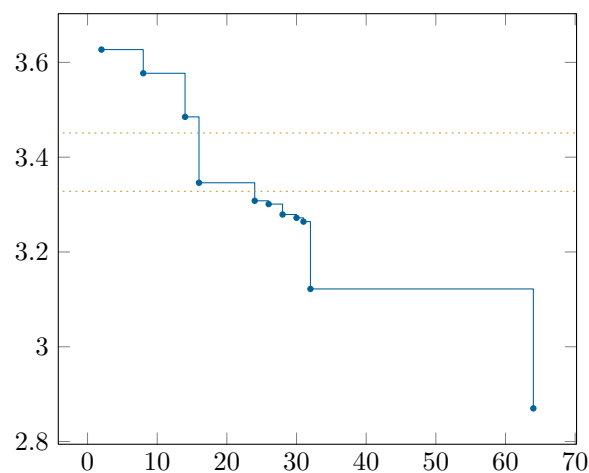
was less than 1. After the m' th generation, the eigenvalues of said form will be less than $0.818 < 2^{-1/3.451}$, by Lemma 13. As n goes to infinity, 0.818 dominates the overall scaling behavior. Hence W , and hence any BEC, enjoys scaling exponent less than 3.451. ◀

In the abstract, we claim that the scaling exponent of $\begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}$ over TECs (and hence BECs) is < 3.328 . This number will be derived in Appendix D of the full version [14] with more intense numerical calculations. In particular, there is a new trapping region that is bounded by two linear splines and is significantly smaller than the region bounded by $ax(1-x)$ for $a = 2\sqrt{7} - 4$ and 2; the attraction toward the new trap is witnessed by sampling TECs with low edge-mass. In Appendix E of the full version [14], we also examine the actual values of $H(W_n)$ and its asymptotic behavior aligns with the estimate 3.328.

8 Conclusions

In this paper, we argue that $\begin{bmatrix} 1 & 0 \\ \omega & 1 \end{bmatrix}$ polarizes BECs faster than $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ does. We first show that a pair of BECs will be transformed into balanced TECs. We then show that balanced TECs will be transformed into edge-heavy TECs. Finally, we show that edge-heavy TECs assume a better scaling exponent.

Our rigorous overestimate of the scaling exponent is 3.451; there is another overestimate of 3.328 with strong numerical evidence. Compared to Arıkan’s 2×2 matrix with $\mu \approx 3.627$, Fazeli–Vardy’s 8×8 matrix with $\mu \approx 3.577$ [16], Trofimiuk–Trifonov’s 16×16 matrix with $\mu \approx 3.346$ [46], and Yao–Fazeli–Vardy’s 32×32 matrix with $\mu \approx 3.122$ [50], our result suggests that one should consider expanding the alphabet size prior to enlarging the matrix size. More precisely, the rigorous estimate is analogous to a 15×15 binary matrix; the more accurate estimate is analogous to a 20×20 binary matrix (see Figure 2).



■ **Figure 2** Horizontal axis: matrix size; vertical axis: scaling exponent of the best known matrix [16, 50, 46, 45, 5]. A matrix size will be skipped if no known matrix outruns all smaller matrices. Underlying channel is BEC. Our estimates 3.451 and 3.328 are marked as dotted lines.

References

- 1 Fariba Abbasi, Hessam Mahdaviyar, and Emanuele Viterbo. Hybrid non-binary repeated polar codes. *IEEE Transactions on Wireless Communications*, pages 1–1, 2022. doi:10.1109/TWC.2022.3159807.
- 2 Erdal Arıkan. Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels. *IEEE Transactions on Information Theory*, 55(7):3051–3073, July 2009. doi:10.1109/TIT.2009.2021379.
- 3 Erdal Arıkan and Emre Telatar. On the rate of channel polarization. In *2009 IEEE International Symposium on Information Theory*, pages 1493–1495, June 2009. doi:10.1109/ISIT.2009.5205856.
- 4 D. Baron, M.A. Khojastepour, and R.G. Baraniuk. How quickly can we approach channel capacity? In *Conference Record of the Thirty-Eighth Asilomar Conference on Signals, Systems and Computers, 2004.*, volume 1, pages 1096–1100 Vol.1, November 2004. doi:10.1109/ACSSC.2004.1399310.
- 5 Manan Bhandari, Ishan Bansal, and V. Lalitha. On the polarizing behavior and scaling exponent of polar codes with product kernels. In *2020 National Conference on Communications (NCC)*, pages 1–6, February 2020. doi:10.1109/NCC48643.2020.9056096.
- 6 Valerio Bioglio, Frédéric Gabry, Ingmar Land, and Jean-Claude Belfiore. Multi-kernel polar codes: Concept and design principles. *IEEE Transactions on Communications*, 68(9):5350–5362, September 2020. doi:10.1109/TCOMM.2020.3006212.
- 7 Jarosław Błasiok, Venkatesan Guruswami, Preetum Nakkiran, Atri Rudra, and Madhu Sudan. General strong polarization. *J. ACM*, 69(2), March 2022. doi:10.1145/3491390.
- 8 Jarosław Błasiok, Venkatesan Guruswami, and Madhu Sudan. Polar Codes with Exponentially Small Error at Finite Block Length. In Eric Blais, Klaus Jansen, José D. P. Rolim, and David Steurer, editors, *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2018)*, volume 116 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 34:1–34:17, Dagstuhl, Germany, 2018. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. doi:10.4230/LIPIcs.APPROX-RANDOM.2018.34.
- 9 Sarit Buzaglo, Arman Fazeli, Paul H. Siegel, Veeresh Taranalli, and Alexander Vardy. Permuted successive cancellation decoding for polar codes. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 2618–2622, June 2017. doi:10.1109/ISIT.2017.8007003.

- 10 Eduardo Camps, Hiram H. López, Gretchen L. Matthews, and Eliseo Sarmiento. Polar decreasing monomial-cartesian codes. *IEEE Transactions on Information Theory*, 67(6):3664–3674, June 2021. doi:10.1109/TIT.2020.3047624.
- 11 Semih Cayci, Toshiaki Koike-Akino, and Ye Wang. Nonbinary polar coding for multilevel modulation. In *2019 Optical Fiber Communications Conference and Exhibition (OFC)*, pages 1–3, 2019.
- 12 Peiyao Chen, Baoming Bai, and Xiao Ma. Two-stage polarization-based nonbinary polar codes for 5g urlc, 2018. doi:10.48550/ARXIV.1801.08059.
- 13 Mao-Ching Chiu. Non-binary polar codes with channel symbol permutations. In *2014 International Symposium on Information Theory and its Applications*, pages 433–437, October 2014.
- 14 Iwan Duursma, Ryan Gabrys, Venkatesan Guruswami, Ting-Chun Lin, and Hsin-Po Wang. Accelerating polarization via alphabet extension, 2022. doi:10.48550/ARXIV.2207.04522.
- 15 Arman Fazeli, Hamed Hassani, Marco Mondelli, and Alexander Vardy. Binary linear codes with optimal scaling: Polar codes with large kernels. *IEEE Transactions on Information Theory*, 67(9):5693–5710, September 2021. doi:10.1109/TIT.2020.3038806.
- 16 Arman Fazeli and Alexander Vardy. On the scaling exponent of binary polarization kernels. In *2014 52nd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 797–804, September 2014. doi:10.1109/ALLERTON.2014.7028536.
- 17 Dina Goldin and David Burshtein. Improved bounds on the finite length scaling of polar codes. *IEEE Transactions on Information Theory*, 60(11):6966–6978, November 2014. doi:10.1109/TIT.2014.2359197.
- 18 Ali Goli, S. Hamed Hassani, and Rüdiger Urbanke. Universal bounds on the scaling behavior of polar codes. In *2012 IEEE International Symposium on Information Theory Proceedings*, pages 1957–1961, July 2012. doi:10.1109/ISIT.2012.6283641.
- 19 Venkatesan Guruswami, Andrii Riazanov, and Min Ye. Arıkan meets shannon: Polar codes with near-optimal convergence to channel capacity. *IEEE Transactions on Information Theory*, pages 1–1, 2022. doi:10.1109/TIT.2022.3146786.
- 20 Venkatesan Guruswami and Patrick Xia. Polar codes: Speed of polarization and polynomial gap to capacity. *IEEE Transactions on Information Theory*, 61(1):3–16, January 2015. doi:10.1109/TIT.2014.2371819.
- 21 S. Hamed Hassani, Kasra Alishahi, and Rudiger Urbanke. On the scaling of polar codes: Ii. the behavior of un-polarized channels. In *2010 IEEE International Symposium on Information Theory*, pages 879–883, June 2010. doi:10.1109/ISIT.2010.5513585.
- 22 S. Hamed Hassani, Ryuhei Mori, Toshiyuki Tanaka, and Rüdiger L. Urbanke. Rate-dependent analysis of the asymptotic behavior of channel polarization. *IEEE Transactions on Information Theory*, 59(4):2267–2276, April 2013. doi:10.1109/TIT.2012.2228295.
- 23 Seyed Hamed Hassani, Kasra Alishahi, and Rüdiger L. Urbanke. Finite-length scaling for polar codes. *IEEE Transactions on Information Theory*, 60(10):5875–5898, October 2014. doi:10.1109/TIT.2014.2341919.
- 24 Masahito Hayashi. Information spectrum approach to second-order coding rate in channel coding. *IEEE Transactions on Information Theory*, 55(11):4947–4966, November 2009. doi:10.1109/TIT.2009.2030478.
- 25 Satish Babu Korada, Andrea Montanari, Emre Telatar, and Rüdiger Urbanke. An empirical scaling law for polar codes. In *2010 IEEE International Symposium on Information Theory*, pages 884–888, June 2010. doi:10.1109/ISIT.2010.5513579.
- 26 Satish Babu Korada, Eren Şaşođlu, and Rüdiger Urbanke. Polar codes: Characterization of exponent, bounds, and constructions. *IEEE Transactions on Information Theory*, 56(12):6253–6264, December 2010. doi:10.1109/TIT.2010.2080990.
- 27 Ingmar Land and Johannes Huber. Information combining. *Foundations and Trends® in Communications and Information Theory*, 3(3):227–330, 2006. doi:10.1561/0100000013.

- 28 Liping Lin. On the construction of the kernel matrix by primitive bch codes for polar codes. *Communications and Network*, 14(1):23–35, 2021.
- 29 Guan-Chen Liu and Qi-Yue Yu. Non-binary polar coded system for the two-user multiple-access channel, 2021. doi:10.48550/ARXIV.2111.03839.
- 30 Marco Mondelli, S. Hamed Hassani, and Rüdiger L. Urbanke. Scaling exponent of list decoders with applications to polar codes. *IEEE Transactions on Information Theory*, 61(9):4838–4851, September 2015. doi:10.1109/TIT.2015.2453315.
- 31 Marco Mondelli, S. Hamed Hassani, and Rüdiger L. Urbanke. Unified scaling of polar codes: Error exponent, scaling exponent, moderate deviations, and error floors. *IEEE Transactions on Information Theory*, 62(12):6698–6712, December 2016. doi:10.1109/TIT.2016.2616117.
- 32 Ryuhei Mori and Toshiyuki Tanaka. Non-binary polar codes using reed-solomon codes and algebraic geometry codes. In *2010 IEEE Information Theory Workshop*, pages 1–5, August 2010. doi:10.1109/CIG.2010.5592755.
- 33 Ryuhei Mori and Toshiyuki Tanaka. Source and channel polarization over finite fields and reed-solomon matrices. *IEEE Transactions on Information Theory*, 60(5):2720–2736, May 2014. doi:10.1109/TIT.2014.2312181.
- 34 Rajai Nasser. An ergodic theory of binary operations—part i: Key properties. *IEEE Transactions on Information Theory*, 62(12):6931–6952, December 2016. doi:10.1109/TIT.2016.2616642.
- 35 Woomyoung Park and Alexander Barg. Polar codes for q-ary channels, $q = 2^r$. *IEEE Transactions on Information Theory*, 59(2):955–969, 2013. doi:10.1109/TIT.2012.2219035.
- 36 Henry D. Pfister and Rüdiger Urbanke. Near-optimal finite-length scaling for polar codes over large alphabets. In *2016 IEEE International Symposium on Information Theory (ISIT)*, pages 215–219, July 2016. doi:10.1109/ISIT.2016.7541292.
- 37 Yury Polyanskiy, H. Vincent Poor, and Sergio Verdú. Channel coding rate in the finite blocklength regime. *IEEE Transactions on Information Theory*, 56(5):2307–2359, May 2010. doi:10.1109/TIT.2010.2043769.
- 38 Noam Presman, Ofer Shapira, and Simon Litsyn. Mixed-kernels constructions of polar codes. *IEEE Journal on Selected Areas in Communications*, 34(2):239–253, February 2016. doi:10.1109/JSAC.2015.2504278.
- 39 Constantin Runge, Thomas Wiegart, Diego Lentner, and Tobias Prinz. Multilevel binary polar-coded modulation achieving the capacity of asymmetric channels, 2022. arXiv:2202.04010.
- 40 Aria G. Sahebi and S. Sandeep Pradhan. Multilevel polarization of polar codes over arbitrary discrete memoryless channels. In *2011 49th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 1718–1725, September 2011. doi:10.1109/Allerton.2011.6120375.
- 41 Eren Şaçoğlu. Polar codes for discrete alphabets. In *2012 IEEE International Symposium on Information Theory Proceedings*, pages 2137–2141, July 2012. doi:10.1109/ISIT.2012.6283740.
- 42 Eren Şaçoğlu, Emre Telatar, and Erdal Arıkan. Polarization for arbitrary discrete memoryless channels. In *2009 IEEE Information Theory Workshop*, pages 144–148, October 2009. doi:10.1109/ITW.2009.5351487.
- 43 Valentin Savin. Non-binary polar codes for spread-spectrum modulations. In *2021 11th International Symposium on Topics in Coding (ISTC)*, pages 1–5, August 2021. doi:10.1109/ISTC49272.2021.9594166.
- 44 Mathis Seidl, Andreas Schenk, Clemens Stierstorfer, and Johannes B. Huber. Polar-coded modulation. *IEEE Transactions on Communications*, 61(10):4108–4119, October 2013. doi:10.1109/TCOMM.2013.090513.130433.
- 45 Grigori Trofimuk. Shortened polarization kernels. In *2021 IEEE Globecom Workshops (GC Wkshps)*, pages 1–6, December 2021. doi:10.1109/GCWkshps52748.2021.9681982.

- 46 Grigorii Trofimiuk and Peter Trifonov. Window processing of binary polarization kernels. *IEEE Transactions on Communications*, 69(7):4294–4305, July 2021. doi:10.1109/TCOMM.2021.3072730.
- 47 Hsin-Po Wang. Complexity and second moment of the mathematical theory of communication, 2021. arXiv:2107.06420.
- 48 Hsin-Po Wang and Iwan M. Duursma. Polar codes' simplicity, random codes' durability. *IEEE Transactions on Information Theory*, 67(3):1478–1508, 2021.
- 49 Hsin-Po Wang, Ting-Chun Lin, Alexander Vardy, and Ryan Gabrys. Sub-4.7 scaling exponent of polar codes, 2022. doi:10.48550/ARXIV.2204.11683.
- 50 Hanwen Yao, Arman Fazeli, and Alexander Vardy. Explicit polar codes with small scaling exponent. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pages 1757–1761, 2019. doi:10.1109/ISIT.2019.8849741.
- 51 Peihong Yuan and Fabian Steiner. Construction and decoding algorithms for polar codes based on 2×2 non-binary kernels. In *2018 IEEE 10th International Symposium on Turbo Codes Iterative Information Processing (ISTC)*, pages 1–5, December 2018. doi:10.1109/ISTC.2018.8625284.