

Low-Latency Communication with Computational Complexity Constraints

Hasan Basri Celebi, Antonios Pitarokoilis, Mikael Skoglund
School of Electrical Engineering and Computer Science
KTH Royal Institute of Technology, Stockholm, Sweden

Abstract—Low-latency communication is one of the most important application scenarios in next-generation wireless networks. Often in communication-theoretic studies latency is defined as the time required for the transmission of a packet over a channel. However, with very stringent latency requirements and complexity constrained receivers, the time required for the decoding of the packet cannot be ignored and must be included in the total latency analysis through accurate modeling. In this paper, we first present a way to calculate decoding time using *per bit* complexity metric and introduce an empirical model that accurately describes the trade-off between the decoding complexity versus the performance of state-of-the-art codes. By considering various communication parameters, we show that including the decoding time in latency analyses has a significant effect on the optimum selection of parameters.

Index Terms—Low-complexity receivers, low-latency, IoT, channel coding, ordered statistics decoder

I. INTRODUCTION

Ultra-reliable, low-latency communications (URLLC) have recently attracted significant research interest, due to emerging delay-critical applications, such as machine-to-machine communication, remote medical surgery, factory automation, and automated traffic control [1]. The performance of latency-constrained communication systems has mostly been evaluated in terms of outage probability, i.e., the probability that the instantaneous mutual information falls below a desired rate. However, it has been shown that, for short packets such an approximation can provide inaccurate estimates since it becomes valid when the length of the transmitted packet grows very large. Accurate bounds on the maximal coding rate that are non-asymptotic with respect to the length of the transmitted packet were given in [2].

Both the outage probability and the non-asymptotic bounds assume that decoding happens instantaneously, i.e., the time required for the decoding of a packet is negligible [3]. This assumption can be justified for small transmission rates, unlimited computational power at the decoder side or loose latency requirements [4]. However, for low-latency communication with complexity-constrained receivers, such as low-budget IoT devices, due to the slower processor capabilities, the decoding time must be accurately modeled and included in the total latency analysis.

To the best of the authors' knowledge, there is no generally accepted model of the computational complexity of a typical channel decoder. It appears, however, the number of

operations per information bit is often selected as a metric for the computational complexity [5]. A brief summary on decoding complexities of practical codes is presented in [6]. An overview on the recent developments for short block-length codes is made in [7] where the authors also discuss the decoding complexities of several state-of-the-art decoders. It is remarked that codes which approach the theoretical limits require higher computational complexities. Recently, in [8] the authors study the computational complexity, defined as the total number of binary operations, of some practical decoders and compare their block-error-rate (BLER) performance for short block-lengths. In [9] the *per information bit* computational complexity, i.e., number of binary operations per information bit, is studied as a function of BLER. It is further shown that an excess power with respect to the normal approximation in [2] must be spent to achieve a fixed allowed BLER at a fixed code-rate, when a particular code is chosen.

In this paper we propose a comprehensive model that relates various parameters of low-latency communication systems with computational complexity constraints, such as decoding complexity (in number of binary operations per information bit), BLER, signal-to-noise ratio (SNR), code rate, and codeword block-length, in an accurate and tractable way. In particular, we first consider the complexity of an ordered statistics (OS) decoder and derive a bound on the complexity that is mathematically tractable. A consistent way to calculate the decoding time of a complexity constrained receiver is presented. Then, a model which reveals the trade-off between the complexity versus performance, in terms of BLER, of OS decoders is proposed. Using the model, the minimum amount of power penalty that is required to meet the reliability constraint is derived. Finally, based on the proposed model we study some interesting communication scenarios that reveal the effect of decoding complexity in latency constrained communication.

II. SYSTEM MODEL

Communication over a binary-input, additive white Gaussian noise (BI-AWGN) channel is considered. Let x_i be the input symbol at the i -th channel use of duration T_s seconds, selected from the set $\{-1, +1\}$. Then, the received sample $y_i \in \mathbb{R}$ is given by

$$y_i = x_i + z_i, \quad (1)$$

where $z_i \sim \mathcal{N}(0, \sigma^2)$ and the signal-to-noise ratio is $\rho = \sigma^{-2}$. Transmission occurs in codewords of n symbols and the

This work was funded in part by the Swedish foundation for strategic research.

transmission duration is $d_T = nT_s$ seconds. Each codeword is the output of the channel encoder at an input of k information bits. Hence, the information rate of the code is $r = \frac{k}{n} \in (0, 1]$.

In this work, the normal approximation to bounds from finite blocklength information theory is used as the benchmark for the maximum information rate over the BI-AWGN [2]. The normal approximation for the channel in (1) for a codeword of length n , with a codeword error probability not exceeding $\epsilon \in (0, 1)$ is given by

$$R(n, \epsilon, \rho) = C(\rho) - \sqrt{\frac{V(\rho)}{n}} Q^{-1}(\epsilon) \log_2 e + O\left(\frac{1}{n}\right), \quad (2)$$

where $C(\rho)$ denotes the capacity, $V(\rho)$ the dispersion and $Q^{-1}(\cdot)$ is the inverse to the Q -function [10].

III. COMPUTATIONAL COMPLEXITY OF DECODERS

Given a codeword of n symbols, the total communication latency is given by

$$d_t = nT_s + d_D, \quad (3)$$

where d_D is the time required for the decoding. The goal of this section is to provide a mathematical expression for d_D , that can summarize in a simple and intuitive way various parameters of decoding algorithms that influence the decoding time.

A. Maximum Likelihood Decoder

The maximum likelihood decoder, which minimizes the probability of codeword error for equiprobable codewords, compares the vector of observations with every codeword and decides for the one minimizing the Euclidean distance. Since there are $M = 2^k$ codewords, the number of operations (not necessarily binary) per information bit is M/k , which implies that the complexity of the optimal decoder is exponential in k . This can be considered as an upper bound on the computational complexity of any practical code.

A simple expression that can exactly characterize the computational complexity of every code is unlikely to be found. In this work, we propose a model for the computational complexity that is based on using Bose, Chaudhuri and Hocquenghem (BCH) codes with OS decoders. Before justifying this choice a brief description of the OS decoders is in order.

B. Ordered Statistics (OS) Decoding

OS decoding is a near-ML soft-decision universal decoding algorithm for any (n, k, d_{\min}) linear block code, where d_{\min} is the minimum Hamming distance of the code. For a given observation vector \mathbf{y} , the log-likelihood ratios (LLRs) are computed and sorted in descending order. The generator matrix is permuted in the corresponding way and transformed to systematic form via Gauss–Jordan elimination. A detailed exposition of the OS decoding can be found in [11]. Here, we will focus on the complexity.

The key parameter of OS decoders is the order, $s \in \mathbb{R}^+$, that is intimately related with the computational complexity of the algorithm. Let E_s be the set of error patterns with Hamming weight up to s . These are the error patterns that are checked by

the OS decoding algorithm. When s is small only a few error patterns are checked and the error rate is high. As s increases, the computational complexity increases and the performance approaches the ML decoder. The cardinality of E_s is $|E_s| = \sum_{i=0}^s \binom{k}{i}$. In fact, this number is the total number of codeword comparisons and can be very high even for low s . For the AWGN channel, the recommended s for near-ML performance is [11]

$$s_r = \min \left\{ \left\lceil \frac{d_{\min}}{4} - 1 \right\rceil, k \right\}, \quad (4)$$

where $\lceil \cdot \rceil$ is the ceiling function.

C. Computational Complexity

The computational complexity of BCH codes with OS decoding is a reasonable choice for the modeling of the computational complexity of more general codes for various reasons. In [9] it was shown that the extended BCH (eBCH) codes with OS decoders come very close to the normal approximation (2) for short block-length codes. The OS decoding algorithm allows for a simple parameterization of the decoding complexity via a single parameter, i.e., the order s . In [9] it was also shown that as s increases, the computational complexity rapidly increases and the performance approaches the ML. On the other hand, for small s the computational complexity is reduced and the performance gracefully degrades. Further, the BCH codes are reasonably flexible in terms of the choice of the coding rate.

Focusing on the computation-intensive operations performed by the OS decoding algorithm can provide an immediate estimate of the total number of binary operations per information bit, i.e.,

$$c = \frac{k^2}{8} + \frac{n}{2} \sum_{i=0}^s \binom{k}{i}, \quad (5)$$

where the first term is due to the Gauss-Jordan elimination of the permuted generator matrix and the second term is the sum of the vector multiplications and comparisons for each error pattern [9]. Note that, for $s \leq 2$ the computational complexity is dominated by the first term, otherwise the second term dominates the complexity. Thus, the order of the complexity can be shown as

$$c = \begin{cases} \mathcal{O}(k^2), & \text{if } s \leq 2, \\ \mathcal{O}(k^s), & \text{if } s > 2. \end{cases} \quad (6)$$

It is worth to note that (5) is hard to interpret and increasingly difficult for mathematical tractability due to the sum of binomial coefficients. Therefore, we derive the following upper bound on c for further analyses.

$$c \leq \frac{k^2}{8} + \frac{n}{2} 2^{kh(\frac{s}{k})} \quad (7)$$

where $h(z) = -z \log_2(z) - (1-z) \log_2(1-z)$ is the binary entropy function. See Appendix for the proof. Note that (7) gets tighter with higher values of s .

D. Decoding Duration

Let T_b be the time required for a binary operation on the hardware platform that the decoder operates. Then the total latency for the transmission of a codeword of blocklength n is given by

$$d_t = nT_s + d_D = nT_s + kcT_b. \quad (8)$$

The decoding time is influenced by the particular hardware platform. For simplicity and generality we assume a linear relation between d_D and T_b . Accuracy of d_D can be improved further by evaluating the hardware technology. But since this is beyond the scope of this paper, we confine to (8) for further analysis.

Suppose that a latency constraint on d_t is applied such as

$$d_t = nT_s + kcT_b \leq d_m, \quad (9)$$

where d_m is the maximum latency deadline for d_t . Such a latency constraint restricts s as follows

$$s_m = \arg \max_{\{s|s \in \mathbb{Q}^+, nT_s + kcT_b \leq d_m\}} c, \quad (10)$$

where s_m denotes the maximum allowed order. Hence, using (7) and (9), the following inequality follows

$$h\left(\frac{s}{k}\right) \leq \frac{1}{k} \log_2 \gamma, \quad (11)$$

where $\gamma = \frac{8d_m - 8nT_s - k^3T_b}{4nkT_b}$. Note that (11) is not an upper bound on s , but meeting (11) guarantees the latency deadline. By using a tight approximation for binary entropy function, given as $h(q) \approx (4q(1-q))^{3/4}$, s_m can be approximated as

$$s_m \approx \frac{k}{2} \left(1 - \sqrt{1 - \left(\frac{\log_2 \gamma}{k} \right)^{4/3}} \right). \quad (12)$$

IV. COMPUTATIONAL COMPLEXITY VS POWER PENALTY

The selection of an order s for a particular code of fixed n , k and ρ can be used to control the total latency d_t of the communication, albeit at the expense of reduced reliability. One way to satisfy a desired target of reliability of ϵ , i.e., codeword error probability, a power penalty has to be paid. Hence, an interesting, yet complex, relation between total latency, computational complexity of the decoder and power spent for transmission arises.

In general, a direct proportion between complexity and BLER performance is expected for a decoder. Thus, a decoder will perform better when its complexity increases, and vice versa. Empirical results shown in [9] reveal that the complexity of the OS decoder exponentially increases as its performance approaches to the normal approximation. Therefore, in this section, we aim to model the trade-off between computational complexity and performance of OS decoders in the finite block-length regime with a tractable expression. Consequently, we first analyze the performance of OS decoders over BI-AWGN channel with different orders at different coding rates for $n = \{64, 128\}$ where the information bits are encoded with eBCH encoder. Results for $n = 128$ are plotted in Fig. 1

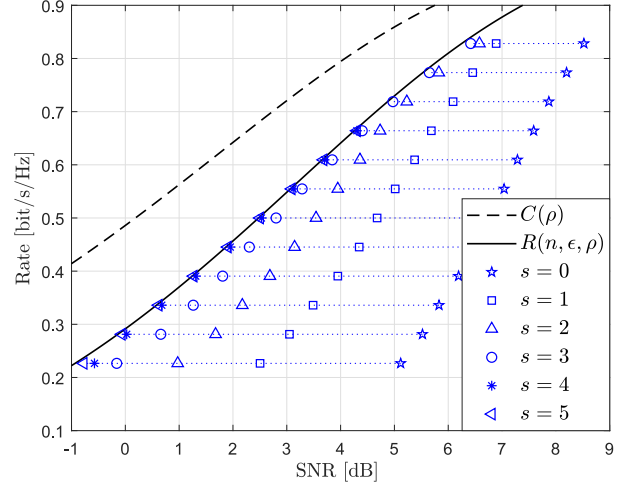


Fig. 1: Power requirements of OS decoders with different orders at different rates for $\epsilon = 10^{-3}$ when $n = 128$.

where the rate as a function of the SNR in dB is plotted. The dashed line is the ergodic capacity in the asymptotic regime, the solid line is the normal approximation (2) for $n = 128$ and $\epsilon = 10^{-3}$. Each blue, horizontal line joins the operating points of eBCH codes with OS decoding, BLER $\epsilon = 10^{-3}$ and a fixed rate. Starting from right to left the order of the decoder increases from $s = 0$ to $s = 5$ and each marker shows the required SNR of a decoder at a fixed rate. Note that these SNR values are computed by starting the BLER analyses from low SNR and detecting the required amount by gradually increasing it until BLER reaches ϵ . Fig. 1 illustrates some significant results. First of all, it shows that performance of OS decoders achieve $R(n, \epsilon, \rho)$ at any rate if s is sufficiently high. It is also clear that as the computational complexity of the decoder increases with increasing s , the power penalty required for the desired codeword error probability decreases. Conversely, a computational constraint due to stringent total latency constraints leads to a corresponding power penalty. Similar results are obtained from the analyses when $n = 64$. However, these results are not shown due to space-limitations.

In Fig. 2 the number of binary operations per bit, c , is plotted as a function of the power penalty $\Delta\rho$ for two codes with blocklength $n = \{64, 128\}$ and $k = \{36, 64\}$, respectively. The individual points correspond to simulation results with decoders of order $s = \{0, 1, \dots, 5\}$. It can be observed that in both cases, the relation between computational complexity and power penalty is closely approximated by a law of the type

$$\log_2 c = \frac{1}{a(\Delta\rho)^\gamma + b}, \quad (13)$$

for appropriate choices of the positive constants α , γ and b . This law is plotted as the dashed lines. For the case of $n = 64$ the parameters were found to be $a = 0.05$, $b = 0.03$, $\gamma = 0.4$ and for $n = 128$ the parameters are given by $a = 0.03$, $b = 0.03$, $\gamma = 0.6$.

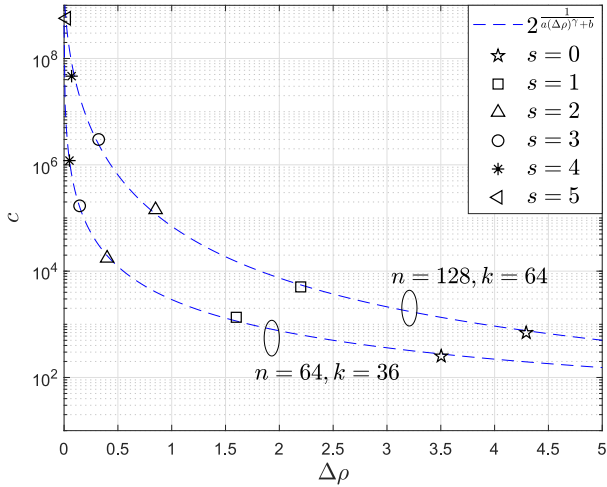


Fig. 2: Comparison of the proposed model with actual results.

The expression in (13) summarizes in a simple way an intuitive trade-off between computational complexity and power penalty for a fixed reliability constraint. Even though it was derived based on the OS decoder, numerical results in previous works [9, Fig. 6] show that other families of codes, such as tail-biting convolutional codes and polar codes, follow a similar law when it comes to the relation between computational complexity and power penalty. Hence, it can be advocated that (13) is a useful proxy for the study of URLLC systems with computational complexity constraints.

V. NUMERICAL EXAMPLES

In this section we present three interesting communication scenarios with latency, reliability and computational complexity constraints. These scenarios reveal that when the decoding time is included in the latency modeling of URLLC communications, the choice of optimal operating point becomes a non-trivial task and gives rise to an abundance of interesting problem formulations.

A. Maximal Information Rate

In Fig. 3 the information rate is plotted as a function of the SNR in dB. The dashed line corresponds to the ergodic capacity and the solid line to the normal approximation (2) for $n = 128$ and $\epsilon = 10^{-3}$. The remaining three plots correspond to maximal information rate when a total latency constraint of $d_m = \{10, 1, 0.25\}$ ms is imposed. It is assumed that the symbol interval is $T_s = 1 \mu\text{s}$ and the time required for a binary operation is $T_b = 1 \text{ ns}$. In particular, for each rate and blocklength, n , the maximum allowable decoding time is calculated using (8). This in turn yields the required power penalty $\Delta\rho$ via (13). Finally, the point on the plot is determined by shifting the point on the normal approximation by $\Delta\rho$ to the right.

B. Maximization of k

We consider the case that a message is subject to a total latency constraint of d_m , with codeword error probability ϵ

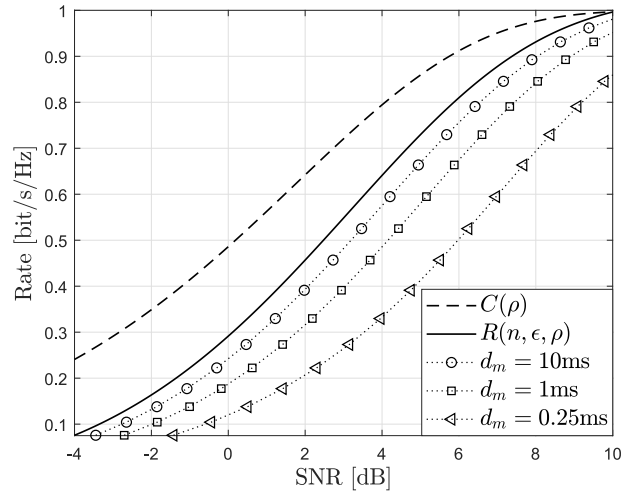


Fig. 3: New achievability bounds under latency and complexity constraints for $n = 128$, $\epsilon = 10^{-3}$, $T_s = 1 \mu\text{s}$, and $T_b = 1 \text{ ns}$.

and there is also a total power budget of P_m at the transmitter. We intend to maximize the number of information bits k that can be transmitted within the total latency and power budget.

When there is unlimited computational power, a codeword can be decoded instantaneously and therefore all the total latency budget can be used for the transmission of the message, i.e., $n = d_m/T_s$ symbols can be transmitted at a rate that is determined by (2), $R(d_m/T_s, \epsilon, P_m)$, which yields

$$k_m = \left\lfloor \frac{d_m R(d_m/T_s, \epsilon, P_m)}{T_s} \right\rfloor. \quad (14)$$

However, for a computational complexity constrained receiver an interesting trade-off arises. If n is selected small, the available duration for decoding can be sufficient so that a high rate code can be used. As n increases, the available duration for decoding shrinks and a code with decreasing coderate must be selected so that the total latency constraint is satisfied.

In Fig. 4 numerical results that correspond to the investigated scenario are plotted for $d_m = 1 \text{ ms}$, $P_m = 5 \text{ dB}$, and $\epsilon = 10^{-3}$. Three different choices for execution times for a binary operation are shown that correspond to $T_b = \{10, 1, 0.1, 0\} \text{ ns}$, where $T_b = 0 \text{ ns}$ stands for infinite computation power. The previously introduced trade-off is clear here and the maximums appear at $n = \{121, 217, 362, 1000\}$, respectively. Corresponding k_m values are $k_m = \{48, 91, 159, 803\}$.

Note that ratios of k_m values found for complexity constrained receivers to the k_m of infinite computation power receiver are 0.06, 0.11, 0.2, respectively. Thus, one can conclude that if complexity constraints and decoding duration are taken into account, one can transmit even less than 20% of the theoretical limits, depending on the receiver capabilities.

C. Minimization of d_t

Another interesting trade-off arises when we intend to transmit a fixed number of information symbols, k , subject to a codeword error probability constraint, ϵ and a maximum

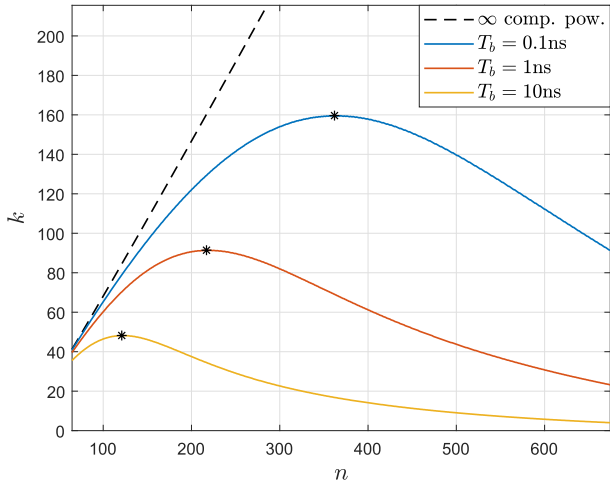


Fig. 4: Maximum k for several complexity constrained receivers where $d_m = 1$ ms, $P_m = 5$ dB, and $\epsilon = 10^{-3}$.

power constraint, P_m . In Fig. 5 the total latency is plotted as a function of the codeword length, n . It can be seen that for small n the coderate of the selected codebook must be very high. Hence, either the transmission is not possible when the required coderate exceeds (2) or the required decoder must operate very close to the normal approximation, which yields a very high required computational complexity. This translates to very high total latency. As n increases, the required rate is decreasing, hence it is more likely that it can be supported by the power budget or a rate sufficiently far from the normal approximation can be selected. In this case, a decoder with low complexity can be selected and the total latency is dominated by the codeword transmission latency. In Fig. 5 the optimal codeword lengths are $n_{\text{opt}} = \{226, 149, 78\}$ for power constraints $P_m = \{3, 5, 10\}$ dB, respectively. Infinite P_m implies that the symbols are transmitted error free and $n_{\text{opt}} = k$ since from (13), $d_D = kT_b \approx 0$ s and hence $d_t = nT_s$ and linearly increases in n .

VI. CONCLUSION

In this study, we focus on latency caused by signal transmission and decoding. We investigate their effect on the transmission parameters by modeling the behavior of OS decoders with an accurate and tractable mathematical expression. Results show that decoding time has a considerable effect on the bounds of the short block-length codes if there is a latency constraint on the system. It is shown that if complexity constraints and decoding duration are considered in the maximization of the total number of information bits to be sent under latency and error rate constraints, less than 20% can be achievable compared to the theoretical limits.

APPENDIX

Let $k \geq 1$ and $\frac{s}{k} \leq \frac{1}{2}$. It is true that

$$1 = \left(\frac{s}{k} + \left(1 - \frac{s}{k}\right)\right)^k \geq \sum_{i=0}^s \binom{k}{i} \left(\frac{s}{k}\right)^i \left(1 - \frac{s}{k}\right)^{k-i}. \quad (15)$$

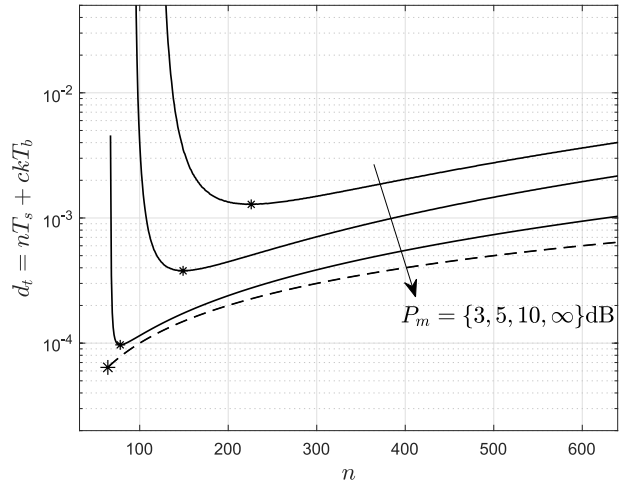


Fig. 5: Minimum d_t with respect to n for several P_m where $k = 64$, $\epsilon = 10^{-3}$, and $T_b = 10^{-9}$ s.

Define $A_i = \left(\frac{s}{k}\right)^i \left(1 - \frac{s}{k}\right)^{k-i}$ for $i \in [0, s]$. Then

$$\begin{aligned} \log_2 A_i &= i \log_2 \left(\frac{s}{nk}\right) + (k-i) \log_2 \left(1 - \frac{s}{k}\right) \geq \\ &= s \log_2 \left(\frac{s}{k}\right) + (k-s) \log_2 \left(1 - \frac{s}{k}\right) = -kh \left(\frac{s}{k}\right). \end{aligned} \quad (16)$$

Using (15) and (16) yields, after some algebra, $2^{kh \left(\frac{s}{k}\right)} \geq \sum_{i=0}^s \binom{k}{i}$.

REFERENCES

- [1] H. Ji, S. Park, J. Yeo, Y. Kim, J. Lee, and B. Shim, "Ultra-reliable and low-latency communications in 5g downlink: Physical layer aspects," *IEEE Wireless Communications*, vol. 25, no. 3, Jun. 2018.
- [2] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2307–2359, May 2010.
- [3] P. Grover, A. Goldsmith, and A. Sahai, "Fundamental limits on the power consumption of encoding and decoding," in *2012 IEEE International Symposium on Information Theory Proceedings*, July 2012.
- [4] C. Rächinger, J. B. Huber, and R. R. Müller, "Comparison of convolutional and block codes for low structural delay," *IEEE Transactions on Communications*, vol. 63, no. 12, Dec 2015.
- [5] A. Khandekar and R. J. McEliece, "On the complexity of reliable communication on the erasure channel," in *Proceedings. 2001 IEEE International Symposium on Information Theory*, June 2001.
- [6] K. Niu, K. Chen, J. Lin, and Q. T. Zhang, "Polar codes: Primary concepts and practical decoding algorithms," *IEEE Communications Magazine*, vol. 52, no. 7, pp. 192–203, July 2014.
- [7] G. Liva, L. Gaudio, and T. Ninnacs, "Code design for short blocks: A survey," in *Proc. EuCNC*, Athens, Greece, Jun 2016.
- [8] M. Sybis, K. Wesolowski, K. Jayasinghe, V. Venkatasubramanian, and V. Vukadinovic, "Channel coding for ultra-reliable low-latency communication in 5g systems," in *2016 IEEE 84th Vehicular Technology Conference (VTC-Fall)*, Sept 2016, pp. 1–5.
- [9] M. Shirvanimoghaddam, M. S. Mohammadi, R. Abbas, A. Minja, C. Yue, B. Matuz, G. Han, Z. Lin, W. Liu, Y. Li, S. Johnson, and B. Vucetic, "Short block-length codes for ultra-reliable low latency communications," *IEEE Communications Magazine*, Feb. 2019.
- [10] T. Erseghe, "Coding in the finite-blocklength regime: Bounds based on laplace integrals and their asymptotic approximations," *IEEE Transactions on Information Theory*, vol. 62, no. 12, Dec 2016.
- [11] M. P. C. Fossorier and S. Lin, "Soft-decision decoding of linear block codes based on ordered statistics," *IEEE Transactions on Information Theory*, vol. 41, no. 5, pp. 1379–1396, Sep. 1995.