# Efficient evaluation of the error probability for pilot-assisted URLLC with Massive MIMO

A. Oguz Kislal, Alejandro Lancho, *Member, IEEE*, Giuseppe Durisi, *Senior Member, IEEE*, and Erik G. Ström, *Fellow, IEEE*

*Abstract*—We propose a numerically efficient method for evaluating the random-coding union bound with parameter $s$ on the error probability achievable in the finite-blocklength regime by a pilot-assisted transmission scheme employing Gaussian codebooks and operating over a memoryless block-fading channel. Our method relies on the saddlepoint approximation, which, differently from previous results reported for similar scenarios, is performed with respect to the number of fading blocks (a.k.a. diversity branches) spanned by each codeword, instead of the number of channel uses per block. This different approach avoids a costly numerical averaging of the error probability over the realizations of the fading process and of its pilot-based estimate at the receiver and results in a significant reduction of the number of channel realizations required to estimate the error probability accurately. Our numerical experiments for both single-antenna communication links and massive multiple-input multiple-output (MIMO) networks show that, when two or more diversity branches are available, the error probability can be estimated accurately with the saddlepoint approximation with respect to the number of fading blocks using a numerical method that requires about two orders of magnitude fewer Monte-Carlo samples than with the saddlepoint approximation with respect to the number of channel uses per block.

## I. Introduction

Next-generation wireless communication systems will support mission-critical links operating under stringent reliability and latency constraints. Denoted as ultra-reliable low-latency communications (URLLC), this type of links will enable applications such as vehicle-to-everything communication, factory automation [2], autonomous driving [3], and haptic communications [4].

One crucial characteristic of the URLLC traffic is that it often involves small information payloads combined with short packets, i.e., packets consisting of a small number of coded symbols. To understand why short packets are needed, it is worth recalling that the length of a data packet depends on the product of the available bandwidth and the signal duration. In URLLC, the signal duration is limited because

of the latency constraint of the targeted applications (e.g., control of automated factories, critical internet-of-things services). The bandwidth is often also limited, because of the need to orthogonalize the transmission of different users to avoid multiuser interference, which has a negative impact on the packet error probability. As pointed out in, e.g., [5], the classic asymptotic performance metrics used to design communication systems, i.e., the ergodic and the outage rates, are unsuitable in the short-packet regime. Thus, a much more precise characterization of the tradeoff between, transmission rate and error probability is required.

Finite-blocklength information theory, a field whose relevance to URLLC has become apparent after the seminal works in [6], [7], provides a precise characterization of such tradeoffs, in terms of nonasymptotic upper (achievability) and lower (converse) bounds on the smallest error probability compatible with a given SNR, transmission rate and packet size.

To satisfy the reliability requirements over fading channels in URLLC, under the above-mentioned diversity limitations in both time and frequency, it becomes crucial to leverage on the spatial diversity offered by multiple antennas. A promising approach is to use massive multiple-input multiple-output (MIMO)—a wireless cellular network architecture in which a base station (BS) with a large number of active antennas serves multiple users on the same time-frequency resources [8]. The benefits of massive MIMO are well understood [9], and this technology has been incorporated into the 5G standard.

Focusing on communication over memoryless block-fading channels, we present in this paper a numerically efficient method to evaluate information-theoretic upper bounds on the finite-blocklength error probability achievable in practically relevant scenarios, including massive MIMO deployments. Methods such as the one presented in this paper are necessary since evaluating most of the available information-theoretic error-probability bounds and approximations that are accurate for scenarios of interest for URLLC [10]–[12] is—as we shall see—extremely time consuming. This prevents the use of such expressions within URLLC optimization routines such as resource-allocation and scheduling algorithms.

*State of the art:* Throughout the paper, we will focus on the upper bound on the error probability obtained by using the random-coding union bound with parameter $s$ (RCUs) proposed in [13]. As discussed in, e.g., [10], this bound is particularly suited for transmission over fading channels because it provides achievability results that hold both for the optimal noncoherent maximum-likelihood (ML) decoder,

as well as for more practically relevant transmission schemes that rely on pilot-assisted transmission (PAT). For example, PAT schemes include the case in which the acquired channel estimate at the receiver is treated as perfect via the use of a mismatched scaled nearest-neighbor (SNN) decoder [14].

The RCUs bound involves the computation of a certain tail probability, which is not known in closed form and needs to be evaluated numerically. If performed naively, this step is time consuming because of the low error probabilities of interest in URLLC. A common approach to circumvent this issue encompasses the following two steps. One starts by noting that, given the realization of the fading channel and of its estimate at the receiver, the random variable whose tail probability is of interest can be written as a sum of independent random variables. Then, one uses the central-limit theorem to approximate this tail probability by a Gaussian tail probability. The resulting approximation, which is typically referred to as *normal approximation* (see, e.g., [7, Sec.IV]), is, however, not accurate for the error probabilities of interest in the URLLC regime [12], [13]. Furthermore, this approach still requires one to perform a Monte-Carlo averaging over the channel realizations and their estimate at the receiver, which is time consuming.

As shown in, e.g., [13], [15], a much more accurate approximation can be obtained by using the so-called *saddlepoint method* [16]. Consider a memoryless block-fading channel where each packet is assumed to span $n_b$ fading blocks. Assume that, within each fading block, $n_s$ coded symbols are transmitted. For this scenario, the saddlepoint method can be applied in two different ways: we can fix $n_b$ and perform a saddlepoint expansion with respect to (w.r.t.) $n_s$, i.e., perform an expansion that is accurate when the number of symbols per block is large. Alternatively, we can fix $n_s$ and perform a saddlepoint expansion w.r.t $n_b$, i.e., perform an expansion that is accurate when the number of fading blocks (a.k.a. diversity branches) is large.

For PAT transmission and SNN decoding, the first approximation has been recently studied in [10], [12] for the special case $n_b = 1$. As discussed in [12], this approach yields an approximation on the conditional error probability given the channel and its estimate, which needs then to be averaged w.r.t. the channel realizations and their estimates at the receiver. A different approach to evaluate this conditional error probability is described in [17].

The second approximation was studied in [11], but only for the case of optimal ML decoder. This approximation pertains the unconditional tail probability, and, hence, does not require an additional averaging over the realizations of the channel and its estimate. As we shall see, this makes the numerical evaluation of this approximation computationally efficient.

*Contributions:* Focusing on independent and identically distributed (i.i.d.) Gaussian codebooks, we present in this paper two saddlepoint approximations on the RCUs for the practically relevant case of PAT and SNN decoding: the one w.r.t. $n_s$ generalizes the one reported in [12] to arbitrary $n_b$ values; the one w.r.t. to $n_b$ generalizes the one reported in [11] to PAT and SNN decoding. Considering the URLLC regime, we then provide a detailed analysis of the accuracy

and the computational complexity of both approximations in: i) a single-input single-output (SISO) setup; ii) the uplink of a two-user single-cell massive MIMO network; iii) the uplink of a multi-user multi-cell massive MIMO network. This progression allows us to understand the impact of the number of BS and users in the network on both the accuracy and the numerical complexity of the considered approximations. Our analysis shows that, despite being developed under the assumption of large $n_b$, the saddlepoint w.r.t. $n_b$ is accurate for $n_b$ values as small as 2 in both SISO and multi-user MIMO scenarios. Furthermore, it entails a much smaller computational complexity than the saddlepoint w.r.t. $n_s$. Specifically, whenever $n_b \geq 2$, the number of samples required to evaluate the saddlepoint w.r.t. $n_b$ via Monte-Carlo simulation is typically around 2 orders of magnitude smaller than the number of samples required to evaluate the saddlepoint approximation w.r.t. $n_s$, once the averaging over the channel and its estimate is accounted for. We also show that, for the scenarios considered in the paper, the normal approximation is typically not accurate.

*Notation:* We denote random vectors and random scalars by upper-case boldface letters such as $\boldsymbol{X}$ and upper-case standard letters, such as $X$, respectively. Their realizations are indicated by lower-case letters of the same font. We use upper-case letters of two special fonts to denote deterministic matrices (e.g., $\mathsf{Y}$) and random matrices (e.g., $\mathbb{Y}$). To avoid ambiguities, we use another font, such as $\mathrm{R}$ for rate, to denote constants that are typically capitalized in the literature. The identity matrix of size $a \times a$ is written as $\mathsf{I}_a$. The circularly-symmetric Gaussian distribution is denoted by $\mathcal{CN}(0, \sigma^2)$, where $\sigma^2$ denotes the variance. The superscripts $(\cdot)^T$, $(\cdot)^H$, and $(\cdot)^*$ denote transposition, Hermitian transposition, and complex conjugation, respectively. We write $\log(\cdot)$ to denote the natural logarithm, $\|\cdot\|$ stands for the $\ell^2$-norm, $\mathbb{P}[\cdot]$ for the probability of an event, $\mathbb{E}[\cdot]$ for the expectation operator, $\mathbb{Var}[\cdot]$ for the variance of a random variable, and $Q(\cdot)$ for the Gaussian $Q$-function. Finally, for two functions $f(n)$ and $g(n)$, the notation $f(n) = o(g(n))$ means that $\lim_{n \to \infty} f(n)/g(n) = 0$ and the notation $f(n) = \mathcal{O}(g(n))$ means that $\lim \sup_{n \to \infty} |f(n)/g(n)| < \infty$.

*Organization of the paper:* In Section II, we present a finite-blocklength upper bound on the error probability for the SISO Rayleigh block-fading channel. We then introduce different methods to evaluate this bound and discuss their accuracy and computational complexity in the URLLC regime, with the help of numerical examples. The extension of our framework to multicell, multiuser massive MIMO networks is presented in Section III. Concluding remarks are provided in Section IV.

## II. A NON-ASYMPTOTIC UPPER BOUND ON THE ERROR PROBABILITY

### A. *The SISO System Model*

We start by considering a SISO memoryless block-fading channel. Specifically, the channel is assumed to stay constant over the transmission of a block of $n_c$ channel uses and to change independently across blocks. Each transmitted packet

spans $n_b$ such fading blocks. Hence, each packet consists of $n_b n_c$ complex-valued symbols. We assume that the first $n_p$ symbols within each block are pilots known to the receiver. These pilots are used to estimate the fading channel within the corresponding block. The input-output relation corresponding to the pilot transmission phase within block $\ell = 1, \ldots, n_b$ is modeled as

$$\boldsymbol{Y}_\ell^{(p)} = H_\ell \mathbf{x}_\ell^{(p)} + \boldsymbol{W}_\ell^{(p)}. \tag{1}$$

Here, $\mathbf{x}_\ell^p$ denotes the deterministic $n_p$-dimensional vector of pilot symbols, which we assume to satisfy the power constraint $\|\mathbf{x}_\ell^{(p)}\|^2 = \rho n_p$, where $\rho$ denotes the average transmit power per symbol. Furthermore, $H_\ell$ denotes the scalar random fading complex channel gain, and $\boldsymbol{W}_\ell^{(p)}$ denotes the $n_p$-dimensional additive noise vector, which may depend on the fading process.[1] We assume that the entries of $\boldsymbol{W}_\ell^{(p)}$ are conditionally independent and $\mathcal{CN}(0, \sigma_\ell^2)$-distributed given the realization of the fading process.

The received vector $\boldsymbol{Y}_\ell^{(p)}$ and the pilot sequence $\mathbf{x}_\ell^{(p)}$ are used by the receiver to obtain an estimate $\hat{H}_\ell$ of the channel $H_\ell$. Note that we have not specified the fading distribution or the algorithm used by the receiver to estimate the fading channel. Indeed, the error probability bounds we shall present in this section hold for arbitrary fading distributions and arbitrary channel-estimation algorithms.

Within each block, the pilot-transmission phase is followed by a data-transmission phase involving $n_s = n_c - n_p$ symbols per block, i.e., a total of $n_b n_s$ symbols. The input-output relation for the $\ell$th block in the data phase is given by

$$\boldsymbol{Y}_\ell = H_\ell \mathbf{x}_\ell + \boldsymbol{W}_\ell. \tag{2}$$

We assume that the $n_b n_s$-dimensional vector $[\mathbf{x}_1^T, \ldots, \mathbf{x}_{n_b}^T]^T$ is selected from a codebook $\mathcal{C}$ of size $\lceil \exp(n_b n_c \mathrm{R}) \rceil$, where $\mathrm{R}$ denotes the transmission rate in nats per channel use.[2]

To perform decoding, the receiver seeks the codeword in the codebook that is closest to the received signal, once each part of the codeword corresponding to a different fading block is scaled by the available channel estimate. Mathematically, given the received vector $[\mathbf{y}_1^T, \ldots, \mathbf{y}_{n_b}^T]^T$ and the channel estimates $\{\hat{h}_1, \ldots, \hat{h}_{n_b}\}$, the decoded codeword $\hat{\mathbf{x}} = [\hat{\mathbf{x}}_1^T, \ldots, \hat{\mathbf{x}}_{n_b}^T]^T$ is determined as follows:

$$\hat{\mathbf{x}} = \arg\min_{\bar{\mathbf{x}} = [\bar{\mathbf{x}}_1^T, \ldots, \bar{\mathbf{x}}_{n_b}^T]^T \in \mathcal{C}} \sum_{\ell=1}^{n_b} \|\mathbf{y}_\ell - \hat{h}_\ell \bar{\mathbf{x}}_\ell\|^2. \tag{3}$$

This decoder, which is known as mismatched SNN decoder, coincides with the ML decoder only when the receiver has perfect channel-state information, i.e., $\hat{h}_\ell = h_\ell$ for $\ell = 1, \ldots, n_b$. The attractive feature of this decoder is that information-theoretic bounds on its error probability can be approached in practice using good channel codes for the nonfading AWGN channel [10]. In contrast, approaching information-theoretic error-probability bounds for the optimal ML decoder considered in [11] with low-complexity coding schemes is still an

open problem (note, however, the recent progress reported in [18]).

### B. The RCUs Finite-Blocklength Bound

Like most of the achievablity results in information theory, the RCUs bound [13] we shall focus on in this paper is obtained by means of a random-coding argument. Specifically, one evaluates the average error probability averaged over a randomly constructed ensemble of codebooks. In this paper, we consider the i.i.d. Gaussian ensemble, in which each symbol of each codeword is drawn independently from a $\mathcal{CN}(0, \rho)$ distribution, where $\rho$ models the average transmit power per symbol in the data phase (same power as in the pilot phase). Although suboptimal at finite blocklength [19] in the nonfading SISO case, the i.i.d. Gaussian ensemble is often used in the literature because it leads tractable expressions when applied to PAT, SNN decoding, and multiuser MIMO scenarios.

Specialized to our setup, the RCUs bound results in the following upper bound $\epsilon_{ub}$ on the *packet error probability* $\epsilon$:

$$\epsilon \le \epsilon_{ub} = \mathbb{P}\left[ \frac{\log U}{n_c n_b} + \frac{1}{n_c n_b} \sum_{\ell=1}^{n_b} \imath_s\left( \boldsymbol{X}_\ell; \boldsymbol{Y}_\ell, \hat{H}_\ell \right) \le \mathrm{R} \right]. \tag{4}$$

Here, $U$ is a random variable that is uniformly distributed on $[0, 1]$ and independent of all other quantities,

$$\imath_s\left( \boldsymbol{X}_\ell; \boldsymbol{Y}_\ell, \hat{H}_\ell \right) = \sum_{k=1}^{n_s} \imath_s\left( X_{k,\ell}; Y_{k,\ell}, \hat{H}_\ell \right) \tag{5}$$

where $X_{k,\ell}$ and $Y_{k,\ell}$ are the $k$th element of $\boldsymbol{X}_\ell$ and $\boldsymbol{Y}_\ell$ respectively, and $\imath_s(x; y, \hat{h})$ is the so-called *generalized information density*, which, for the case of i.i.d. $\mathcal{CN}(0, \rho)$ codebooks and SNN decoding, is given by [12, App. A]

$$\imath_s(x; y, \hat{h}) = -s \left| y - \hat{h}x \right|^2 + \frac{s |y|^2}{1 + s\rho \left| \hat{h} \right|^2} + \log\left( 1 + s\rho \left| \hat{h} \right|^2 \right). \tag{6}$$

Finally, $s > 0$ is an optimization parameter that can be used to tighten the bound.

In general, no closed-form expression is available for (4). Hence, this probability needs to be evaluated with numerical methods. A naïve implementation of this step results in time-consuming simulations, given the low target error probabilities of interest in URLLC. We next discuss two approaches to compute (4) efficiently: one is based on asymptotic expansions of (4) applied w.r.t. the number of data symbols per block $n_s$, and the other is based on asymptotic expansions of (4) applied w.r.t. the number of blocks $n_b$. For each approach, we will present an expansion based on the central-limit theorem, which will result in the so-called normal approximation, and an expansion based on the saddlepoint method.

### C. Asymptotic Expansion w.r.t. $n_s$

The idea behind this approach, which for the case $n_b = 1$ has been explored in [12], is to analyze first a conditional

---

[1] Allowing for such dependency will turn out crucial to extend the SISO analysis to the multiuser MIMO case.

[2] With an abuse of notation, we will use $\mathrm{R}$ to denote also the rate measured in bits per channel use when presenting numerical experiments in Sections II-E and III-C

version of the probability in (4), in which the channel and its estimate within each block are given. Specifically, one focuses on

$$\epsilon_{\text{ub}}(\mathbf{h}, \hat{\mathbf{h}}) = \mathbb{P}\Bigg[\frac{\log U}{n_{\text{c}} n_{\text{b}}} + \frac{1}{n_{\text{c}} n_{\text{b}}} \sum_{\ell=1}^{n_{\text{b}}} \sum_{k=1}^{n_{\text{s}}} \imath_s\Big(X_{k,\ell}; Y_{k,\ell}, \hat{h}_\ell\Big)$$
$$\leq R \Big| \boldsymbol{H} = \mathbf{h}, \hat{\boldsymbol{H}} = \hat{\mathbf{h}}\Bigg] \quad (7)$$

where $\boldsymbol{H} = [H_1, \ldots, H_{n_{\text{b}}}]^T$ and $\hat{\boldsymbol{H}} = [\hat{H}_1, \ldots, \hat{H}_{n_{\text{b}}}]^T$. [3] Then one seeks an asymptotic approximation to this conditional probability that is easy to evaluate numerically because it depends on quantities that can be evaluated in closed form. Finally, one performs the averaging

$$\epsilon_{\text{ub}} = \mathbb{E}_{\boldsymbol{H}, \hat{\boldsymbol{H}}}\Big[\epsilon_{\text{ub}}(\boldsymbol{H}, \hat{\boldsymbol{H}})\Big] \quad (8)$$

over the channel and its estimate numerically.

*1) Normal Approximation:* One way to numerically approximate (7) is to perform a normal approximation w.r.t. $n_{\text{s}}$ based on the Berry-Esseen central-limit theorem [20, Ch. XVI.5]. Specifically, note that, given $\boldsymbol{H}$ and $\hat{\boldsymbol{H}}$, the $n_{\text{b}} n_{\text{s}}$ samples of $\{\imath_s(X_{k,\ell}; Y_{k,\ell}, \hat{h}_\ell)\}$ in (7) are conditionally independent and identically distributed within each fading block. By applying the Berry-Esseen central-limit theorem [20, Ch. XVI.5] to the tail probability in (7), we obtain

$$\epsilon_{\text{ub}}(\mathbf{h}, \hat{\mathbf{h}}) = Q\left(\frac{n_s \sum_{\ell=1}^{n_{\text{b}}} I_s(h_\ell, \hat{h}_\ell) - n_{\text{c}} n_{\text{b}} R}{\sqrt{n_s \sum_{\ell=1}^{n_{\text{b}}} V_s\Big(h_\ell, \hat{h}_\ell\Big)}}\right) + \mathcal{O}\left(\frac{1}{\sqrt{n_s}}\right). \quad (9)$$

Here,

$$I_s(h_\ell, \hat{h}_\ell) = \mathbb{E}\Big[\imath_s\Big(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell\Big)\Big] \quad (10)$$
$$= \Big(1 + s\rho\Big|\hat{h}_\ell\Big|^2\Big) + (\beta_\ell - \alpha_\ell) \quad (11)$$

where

$$\alpha_\ell = s\Big(\rho\Big|h_\ell - \hat{h}_\ell\Big|^2 + \sigma^2\Big) \quad (12)$$
$$\beta_\ell = \frac{s}{1 + s\rho\Big|\hat{h}\Big|^2}\Big(\rho|h_\ell|^2 + \sigma^2\Big) \quad (13)$$

is the so-called *generalized mutual information.* Furthermore, the variance of the information density is

$$V_s\Big(h_\ell, \hat{h}_\ell\Big) = \mathbb{V}\text{ar}\Big[\imath_s\Big(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell\Big)\Big] \quad (14)$$
$$= (\beta_\ell - \alpha_\ell)^2 + 2\alpha_\ell \beta_\ell(1 - \nu_\ell) \quad (15)$$

where

$$\nu_\ell = \frac{s^2\Big|\rho|h_\ell|^2 + \sigma^2 - h_\ell^* \hat{h}_\ell \rho\Big|^2}{\alpha_\ell \beta_\ell\Big(1 + s\rho\Big|\hat{h}_\ell\Big|^2\Big)}. \quad (16)$$

[3]Note that $Y_{k,\ell}$ depends on $H_\ell$ via (2).

Note that in (11) and (15) we set $k = 1$ without loss of generality since, given $\boldsymbol{H}$ and $\hat{\boldsymbol{H}}$ and for a fixed $\ell$, the random variables $\{\imath_s(X_{k,\ell}; Y_{k,\ell}, \hat{h}_\ell)\}$ are conditionally i.i.d. over $k$.

Strictly speaking, for (9) to hold, we need to verify that the third central moment of $\imath_s\Big(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell\Big)$ exists for every $\ell \in \{1, \ldots, n_{\text{b}}\}$. We will show in Section II-C2 that this is indeed the case. The normal approximation w.r.t. $n_{\text{s}}$ is finally obtained by ignoring the $\mathcal{O}(\cdot)$ term in (9) and by averaging the resulting approximation over $\boldsymbol{H}$ and $\hat{\boldsymbol{H}}$. This is typically done via Monte-Carlo simulations.

*2) Saddlepoint Approximation:* Since it is based on a central-limit theorem, the normal approximation is typically accurate only in the regime in which the target rate $R$ is close to the mean of the information density [11]. However, this regime may be of limited interest in URLLC, since it may correspond to packet error probability values above the URLLC target (see, e.g., [12, Fig. 1]).

A more refined approximation can be obtained by using the so-called saddlepoint method. It results in an error probability expansion given in terms of a leading factor that decays exponentially with $n_{\text{s}}$ and captures the behavior of the error probability in the large-deviation regime, and a subexponential factor, which is obtained by applying a refined normal approximation, and which makes the resulting approximation accurate in the short-packet regime. This method allows one to obtain an approximation that is accurate for a large range of target error probabilities and rates, including the ones relevant in URLLC scenarios.

We now state this approximation. Consider again the conditional probability given in (7). We fix again $k = 1$, without loss of generality, and let $\kappa(\zeta)$ be the cumulant generating function (CGF) of the random variable $-\sum_{\ell=1}^{n_{\text{b}}} \imath_s(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell)$:

$$\kappa(\zeta) = \log \mathbb{E}\Big[e^{-\zeta \sum_{\ell=1}^{n_{\text{b}}} \imath_s\big(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell\big)}\Big] \quad (17)$$
$$= \sum_{\ell=1}^{n_{\text{b}}} \log \mathbb{E}\Big[e^{-\zeta \imath_s\big(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell\big)}\Big]. \quad (18)$$

Note that $\kappa(\zeta)$ depends on $\mathbf{h}$ and $\hat{\mathbf{h}}$, but we choose not to make this dependence explicit, to keep the notation compact. Each term in (18) admits a closed-form expression. Specifically, let

$$g(\zeta, h_\ell, \hat{h}_\ell) = \mathbb{E}\Big[e^{-\zeta \imath_s\big(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell\big)}\Big] \quad (19)$$

be the moment generating function (MGF) of the random variable $-\imath_s\Big(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell\Big)$. Then [12, Eq. (56)]

$$g(\zeta, h_\ell, \hat{h}_\ell) = \frac{\Big(1 + (\beta_\ell - \alpha_\ell)\zeta - \alpha_\ell \beta_\ell(1 - \nu_\ell)\zeta^2\Big)^{-1}}{\Big(1 + s\rho\Big|\hat{h}_\ell\Big|^2\Big)^\zeta} \quad (20)$$

where $\alpha_\ell$, $\beta_\ell$, and $\nu_\ell$ where defined in (12), (13), and (16).

Note that, by substituting (20) into (18), one can obtain a closed-form expression not only for $\kappa(\zeta)$, but also for its first and second derivatives, which we shall denote as $\kappa'(\zeta)$ and $\kappa''(\zeta)$, and we shall need shortly. Specifically, let $\kappa_\ell(\zeta) =$

$\log g(\zeta, h_\ell, \hat{h}_\ell)$, so that $\kappa(\zeta) = \sum_{\ell=1}^{n_b} \kappa_\ell(\zeta)$. We have that (see [12, Eqs. (16)–(18)])

$$
\begin{aligned}
\kappa_\ell(\zeta) = & -\zeta \log(1 + s\rho |\hat{h}_\ell|^2) \\
& - \log(1 + (\beta_\ell - \alpha_\ell)\zeta - \alpha_\ell\beta_\ell(1 - \nu_\ell)\zeta^2)
\end{aligned}
\tag{21}
$$

$$
\begin{aligned}
\kappa'_\ell(\zeta) = & -\log(1 + s\rho |\hat{h}_\ell|^2) \\
& - \frac{(\beta_\ell - \alpha_\ell) - 2\alpha\beta_\ell(1 - \nu_\ell)\zeta}{1 + (\beta_\ell - \alpha_\ell)\zeta - \alpha_\ell\beta_\ell(1 - \nu_\ell)\zeta^2}
\end{aligned}
\tag{22}
$$

$$
\begin{aligned}
\kappa''_\ell(\zeta) = & \left[\frac{(\beta_\ell - \alpha_\ell) - 2\alpha_\ell\beta_\ell(1 - \nu_\ell)\zeta}{1 + (\beta_\ell - \alpha_\ell)\zeta - \alpha_\ell\beta_\ell(1 - \nu_\ell)\zeta^2}\right]^2 \\
& + \frac{2\alpha_\ell\beta_\ell(1 - \nu_\ell)}{1 + (\beta_\ell - \alpha_\ell)\zeta - \alpha_\ell\beta_\ell(1 - \nu_\ell)\zeta^2}.
\end{aligned}
\tag{23}
$$

A saddlepoint expansion can be established provided that the third derivative of the MGF of $-\imath_s(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell)$ exists in a neighborhood of zero for every $\ell \in \{1, \ldots, n_b\}$. Specifically, for every $\ell \in \{1, \ldots, n_b\}$, we require that there exist two values $\underline{\zeta}_\ell < 0 < \overline{\zeta}_\ell$ such that

$$
\sup_{\underline{\zeta}_\ell \le \zeta \le \overline{\zeta}_\ell} \left|\frac{d^3}{d\zeta^3} g(\zeta, h_\ell, \hat{h}_\ell)\right| < \infty.
\tag{24}
$$

As shown in [12, Appendix B], this condition holds in our setup with

$$
\underline{\zeta}_\ell = -\frac{\sqrt{(\beta_\ell - \alpha_\ell)^2 + 4\alpha_\ell\beta_\ell(1 - \nu_\ell)} + \alpha_\ell - \beta_\ell}{2\alpha_\ell\beta_\ell(1 - \nu_\ell)}
\tag{25}
$$

$$
\overline{\zeta}_\ell = -\frac{\sqrt{(\beta_\ell - \alpha_\ell)^2 + 4\alpha_\ell\beta_\ell(1 - \nu_\ell)} - \alpha_\ell + \beta_\ell}{2\alpha_\ell\beta_\ell(1 - \nu_\ell)}.
\tag{26}
$$

This implies in particular that the third moment of $-\imath_s(X_{1,\ell}; Y_{1,\ell}, \hat{h}_\ell)$, which can be obtained by evaluating the third derivative in (24) at $\zeta = 0$, exists—a condition we required to establish the normal approximation in Section II-C1.

By taking $\underline{\zeta} = \max\{\underline{\zeta}_1, \ldots, \underline{\zeta}_{n_b}\}$ and $\overline{\zeta} = \min\{\overline{\zeta}_1, \ldots, \overline{\zeta}_{n_b}\}$, we ensure that (24) holds simultaneously for every $\ell \in \{1, \ldots, n_b\}$. The saddlepoint expansion w.r.t. $n_s$ is stated in the following theorem.

*Theorem 1:* Assume that there exists a $\zeta \in [\underline{\zeta}, \overline{\zeta}]$ satisfying $R = -\kappa'(\zeta)n_s/(n_c n_b)$. If $\zeta \in [0, 1]$, then

$$
\begin{aligned}
\epsilon_{ub}(\mathbf{h}, \hat{\mathbf{h}}) = & \\
e^{n_s(\kappa(\zeta) - \zeta\kappa'(\zeta))} & \left[\Psi_{n_s, \zeta}(\zeta) + \Psi_{n_s, \zeta}(1 - \zeta) + o\left(\frac{1}{\sqrt{n_s}}\right)\right]
\end{aligned}
\tag{27}
$$

where

$$
\Psi_{b, \zeta}(u) = e^{b\frac{u^2}{2}\kappa''(\zeta)} Q\left(u\sqrt{b\kappa''(\zeta)}\right).
\tag{28}
$$

If $\zeta > 1$, then

$$
\begin{aligned}
\epsilon_{ub}(\mathbf{h}, \hat{\mathbf{h}}) = & \\
e^{n_s[\kappa(1) - \kappa'(\zeta)]} & \left[\tilde{\Psi}_{n_s}(1, 1) + \tilde{\Psi}_{n_s}(0, -1) + \mathcal{O}\left(\frac{1}{\sqrt{n_s}}\right)\right]
\end{aligned}
\tag{29}
$$

where

$$
\tilde{\Psi}_b(a_1, a_2) = e^{ba_1\left[-\kappa'(1) - R + \frac{\kappa''(1)}{2}\right]}
$$
$$
\times Q\left(a_1\sqrt{b\kappa''(1)} - a_2\frac{b(\kappa'(1) + R)}{\sqrt{b\kappa''(1)}}\right).
\tag{30}
$$

If $\zeta < 0$, then

$$
\begin{aligned}
\epsilon_{ub}(\mathbf{h}, \hat{\mathbf{h}}) = & 1 - \\
e^{n_s[\kappa(\zeta) - \zeta\kappa'(\zeta)]} & \left[\Psi_{n_s, \zeta}(-\zeta) - \Psi_{n_s, \zeta}(1 - \zeta) + o\left(\frac{1}{\sqrt{n_s}}\right)\right].
\end{aligned}
\tag{31}
$$

*Proof:* Although a direct proof of this theorem is not available in the literature, the desired expansions can be established following steps similar to the ones reported in [21, App. E] for the case of abstract channels and generic mismatch decoding rules and in [11, App. I] for the case of memoryless block-fading channels and ML decoding rule. ∎

We obtain the desired saddlepoint approximation of $\epsilon_{ub}(\mathbf{h}, \hat{\mathbf{h}})$ in (7) by omitting the $o(\cdot)$ and the $\mathcal{O}(\cdot)$ terms in (27), (29), and (31).

### D. Asymptotic Expansion w.r.t. $n_b$

We next present a different approach, which avoids the conditioning w.r.t. $\mathbf{H}$ and $\hat{\mathbf{H}}$ and the associated, often time-consuming, Monte-Carlo step. The idea is to exploit directly that the random variables $\{\imath_s(\mathbf{X}_\ell; \mathbf{Y}_\ell, \hat{H}_\ell)\}$ in (4) are i.i.d. across the block index $\ell$, and perform an asymptotic expansion of the tail probability in (4) w.r.t. the number of blocks $n_b$.

*1) Normal Approximation:* Proceeding similar to Section II-C1, we can obtain an asymptotic expansion—this time directly of $\epsilon_{ub}$ in (4)—by applying the Berry-Esseen central-limit theorem. Specifically,

$$
\begin{aligned}
\epsilon_{ub} = Q & \left(\frac{\sqrt{n_b}\left(n_s \mathbb{E}\left[I_s(H_1, \hat{H}_1)\right] - n_c R\right)}{\sqrt{n_s \mathbb{E}\left[V_s(H_1, \hat{H}_1)\right] + n_s^2 \mathbb{Var}\left[I_s(H_1, \hat{H}_1)\right]}}\right) \\
& + \mathcal{O}\left(\frac{1}{\sqrt{n_b}}\right).
\end{aligned}
\tag{32}
$$

Here, we have fixed $\ell = 1$ without loss of generality. We then obtain a normal approximation of $\epsilon_{ub}$ w.r.t. $n_b$ by neglecting the $\mathcal{O}(\cdot)$ term in (32). Note that, differently from the normal approximation provided in (9), the one provided in (32) applies directly to $\epsilon_{ub}$ and not to the conditional probability $\epsilon_{ub}(\mathbf{h}, \hat{\mathbf{h}})$. Hence, no Monte-Carlo averaging step is required at the end. On the negative side, although $I_s(\cdot, \cdot)$ and $V_s(\cdot, \cdot)$ are available in closed form (see (11) and (15)), the terms $\mathbb{E}\left[I_s(H_1, \hat{H}_1)\right]$, $\mathbb{E}\left[V_s(H_1, \hat{H}_1)\right]$, and $\mathbb{Var}\left[I_s(H_1, \hat{H}_1)\right]$ need to be evaluated numerically using, e.g., Monte-Carlo methods. Similarly, the existence of the third central moment of $\imath_s(\mathbf{X}_1; \mathbf{Y}_1, \hat{H}_1)$, which is required for (32) to hold, needs to be ensured with numerical methods.

*2) Saddlepoint Approximation:* We now proceed as in Section II-C2 and obtain a saddlepoint approximation of $\epsilon_{\mathrm{ub}}$ w.r.t. $n_{\mathrm{b}}$. As pointed out in Section II-C2, to establish a saddlepoint asymptotic expansions we need that the third derivative of the MGF of the random variables at hand exists in a neighborhood of zero. Specifically, to establish a saddlepoint approximation of $\epsilon_{\mathrm{ub}}$ w.r.t. $n_{\mathrm{b}}$, we shall require that, for some $\underline{\zeta} < 0 < \overline{\zeta}$,

$$\sup_{\underline{\zeta} < \zeta < \overline{\zeta}} \left| \frac{\mathrm{d}^3}{\mathrm{d}\zeta^3} \mathbb{E}\left[ e^{-\zeta \imath_s(\boldsymbol{X}_1; \boldsymbol{Y}_1, \hat{H}_1)} \right] \right| < \infty. \tag{33}$$

Unfortunately, differently from (24), the moment-generating function $\mathbb{E}\left[ e^{-\zeta \imath_s(\boldsymbol{X}_1; \boldsymbol{Y}_1, \hat{H}_1)} \right]$ is not known in closed form. Hence, no closed-form expressions for $\underline{\zeta}$ and $\overline{\zeta}$ are available and these quantities need to be estimated with numerical methods.

To state the saddlepoint approximation, we shall need the CGF $\gamma(\zeta)$ of the random variable $-\imath_s(\boldsymbol{X}_1; \boldsymbol{Y}_1, \hat{H}_1)$

$$\gamma(\zeta) = \log \mathbb{E}\left[ e^{-\zeta \imath_s(\boldsymbol{X}_1; \boldsymbol{Y}_1, \hat{H}_1)} \right] \tag{34}$$

and its first and second derivatives $\gamma'(\zeta)$ and $\gamma''(\zeta)$. Note that, given $H_1$ and $\hat{H}_1$, the random variable $\imath_s(\boldsymbol{X}_1; \boldsymbol{Y}_1, \hat{H}_1)$ can be decomposed into the sum of $n_{\mathrm{s}}$ conditionally i.i.d. terms (see (5)). Hence,

$$\gamma(\zeta) = \log \mathbb{E}_{H_1, \hat{H}_1}\left[ \prod_{k=1}^{n_{\mathrm{s}}} \mathbb{E}\left[ e^{-\zeta \imath_s(X_{k,1}; Y_{k,1}, \hat{H}_1)} \,\Big|\, H_1, \hat{H}_1 \right] \right]$$

$$= \log \mathbb{E}_{H_1, \hat{H}_1}\left[ g(\zeta, H_1, \hat{H}_1)^{n_{\mathrm{s}}} \right] \tag{35}$$

where the function $g$ was defined in (19). Let $p(\zeta)$ be the MGF of $-\imath_s(\boldsymbol{X}_1; \boldsymbol{Y}_1, \hat{H}_1)$. Then $p(\zeta)$ and its first and second derivatives with respect to $\zeta$ are given as

$$p(\zeta) = \mathbb{E}_{H_1, \hat{H}_1}\left[ g(\zeta, H_1, \hat{H}_1)^{n_{\mathrm{s}}} \right] \tag{36}$$

$$p'(\zeta) = n_{\mathrm{s}} \mathbb{E}_{H_1, \hat{H}_1}\left[ g(\zeta, H_1, \hat{H}_1)^{n_{\mathrm{s}}-1} g'(\zeta, H_1, \hat{H}_1) \right] \tag{37}$$

$$p''(\zeta) = n_{\mathrm{s}} \mathbb{E}_{H_1, \hat{H}_1}\Big[ (n_{\mathrm{s}}-1) g(\zeta, H_1, \hat{H}_1)^{n_{\mathrm{s}}-2} g'(\zeta, H_1, \hat{H}_1)^2$$
$$+ g(\zeta, H_1, \hat{H}_1)^{n_{\mathrm{s}}-1} g''(\zeta, H_1, \hat{H}_1) \Big] \tag{38}$$

where $g'$, $g''$, $p'$ and $p''$ denote the first and second derivatives with respect to $\zeta$ of the functions $g$ and $p$, respectively. We can then write $\gamma(\zeta)$ and its first and second derivatives as

$$\gamma(\zeta) = \log p(\zeta) \tag{39}$$

$$\gamma'(\zeta) = \frac{p'(\zeta)}{p(\zeta)} \tag{40}$$

$$\gamma''(\zeta) = \frac{p''(\zeta) p(\zeta) - p'(\zeta)^2}{p(\zeta)^2}. \tag{41}$$

We are now ready to state the saddlepoint expansion w.r.t. $n_{\mathrm{b}}$.

*Theorem 2:* Assume that there exists a $\zeta \in [\underline{\zeta}, \overline{\zeta}]$ satisfying $\mathrm{R} = -\gamma'(\zeta)/n_{\mathrm{c}}$. If $\zeta \in [0, 1]$ then

$$\epsilon_{\mathrm{ub}} = e^{n_{\mathrm{b}}[\gamma(\zeta) - \zeta\gamma'(\zeta)]}$$
$$\times \left[ \Phi_{n_{\mathrm{b}},\zeta}(\zeta) + \Phi_{n_{\mathrm{b}},\zeta}(1-\zeta) + o\left(\frac{1}{\sqrt{n_{\mathrm{b}}}}\right) \right] \tag{42}$$

where

$$\Phi_{b,\zeta}(u) = e^{b \frac{u^2}{2} \gamma''(\zeta)} Q\left( u\sqrt{b\gamma''(\zeta)} \right). \tag{43}$$

If $\zeta > 1$, then

$$\epsilon_{\mathrm{ub}} = e^{n_{\mathrm{b}}[\gamma(1) - \gamma'(\zeta)]}$$
$$\times \left[ \tilde{\Phi}_{n_{\mathrm{b}}}(1,1) + \tilde{\Phi}_{n_{\mathrm{b}}}(0,-1) + \mathcal{O}\left(\frac{1}{\sqrt{n_{\mathrm{b}}}}\right) \right] \tag{44}$$

where

$$\tilde{\Phi}_b(a_1, a_2) = e^{b a_1 \left[ -\gamma'(1) - \mathrm{R} + \frac{\gamma''(1)}{2} \right]}$$
$$\times Q\left( a_1 \sqrt{b\gamma''(1)} - a_2 \frac{b(\gamma'(1) + \mathrm{R})}{\sqrt{b\gamma''(1)}} \right). \tag{45}$$

Finally, if $\zeta < 0$,

$$\epsilon_{\mathrm{ub}} = 1 - e^{n_{\mathrm{b}}[\gamma(\zeta) - \zeta\gamma'(\zeta)]}$$
$$\times \left[ \Phi_{n_{\mathrm{b}},\zeta}(-\zeta) - \Phi_{n_{\mathrm{b}},\zeta}(1-\zeta) + o\left(\frac{1}{\sqrt{n_{\mathrm{b}}}}\right) \right]. \tag{46}$$

*Proof:* The proof follows by combining the steps in the proofs of [11, App. I] and [21, App. E]. Note that the saddlepoint approximation in [11] was developed for Rayleigh SISO block-fading channels but under the assumption of ML decoding. Furthermore, this expansion was provided only for the case $\zeta \in [0, 1]$. The saddlepoint derived in [21] holds only for channels whose input and output belong to finite-cardinality alphabets, but applies to arbitrary mismatched decoding rules and arbitrary values of $\zeta$. Our result is obtained by carefully combining the proof techniques used in these two papers. ∎

We obtain the desired saddlepoint approximation of $\epsilon_{\mathrm{ub}}$ in (4) w.r.t. $n_{\mathrm{b}}$ by neglecting the $\mathcal{O}(\cdot)$ and the $o(\cdot)$ terms in (42), (44), and (46). Note that, differently from the asymptotic expansion provided in Theorem 1, the one provided in Theorem 2 applies directly to $\epsilon_{\mathrm{ub}}$ and not on the conditional probability $\epsilon_{\mathrm{ub}}(\mathbf{h}, \hat{\mathbf{h}})$. Hence, no Monte-Carlo averaging step is required at the end. On the negative side, the function $\gamma(\cdot)$ and its first and second derivatives are not available in closed form and one needs to resort to numerical methods, such as Monte-Carlo averaging, to evaluate them. Specifically, one needs to evaluate numerically the expectation over the channel $H_1$ and its estimate $\hat{H}_1$ appearing in the definition of $p(\zeta)$ and of its first and second derivatives in (36), (37), and (38).

### E. Numerical Experiments

We next evaluate numerically the two saddlepoint approximations and the two normal approximations just introduced, with the goal of shedding light on the following three questions:

1) In typical scenarios, the number $n_{\mathrm{s}}$ of symbols per block is much larger than the number $n_{\mathrm{b}}$ of blocks spanned by a codeword. How large should $n_{\mathrm{b}}$ be for the approximations w.r.t. to $n_{\mathrm{b}}$ to be accurate?
2) Is the normal approximation (either w.r.t. $n_{\mathrm{s}}$ or w.r.t. $n_{\mathrm{b}}$) sufficiently accurate in the URLLC regime, or should one use instead the saddlepoint approximations?
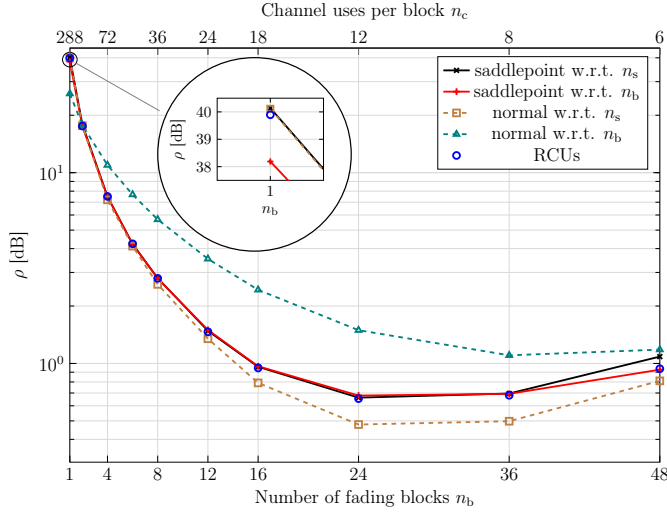
Fig. 1: Upper bound on the required transmit power $\rho$ to achieve $\epsilon = 10^{-5}$. Here, $n_b n_c = 288$, R = 0.104 bit per channel use; $n_p$ and $s$ are optimized.

3) All the approximations presented in Section II require numerical methods such as Monte-Carlo averaging for the evaluation of terms that are not available in closed form. Which method has lower complexity for a given targeted accuracy?

To answer these questions, we perform the following numerical experiments. We start by considering a Rayleigh-fading scenario where the $\{H_\ell\}_{\ell=1}^{n_b}$ are generated independently from a $\mathcal{CN}(0,1)$ distribution. Furthermore, we assume ML estimation of the channel at the receiver. Specifically, we set

$$\hat{H}_\ell = \frac{1}{\rho n_p}\left(\mathbf{x}_\ell^{(p)}\right)^H \mathbf{Y}_\ell^{(p)}. \qquad (47)$$

We assume that the noise variance $\sigma_\ell^2$ is equal to 1 for $\ell \in \{1, \ldots, n_b\}$, and consider a blocklength $n_b n_c$ of 288 channel uses. The results reported in this section are obtained after an optimization over the parameter $s > 0$ in (4) and over the number of pilots $n_p$ within each block of $n_c$ channel uses.

*Accuracy:* In Fig. 1, we report the transmit-power value $\rho$ needed to achieve an error probability $\epsilon = 10^{-5}$ for R = 0.104 bit per channel use, as a function of the number of fading blocks $n_b$ spanned by each codeword. Note that since the blocklength is fixed, $n_c$ decreases when $n_b$ is increased. This implies that fewer symbols are available in each block for pilot and data transmissions. The value of $\rho$ is estimated by means of the RCUs bound in (4), evaluated via a Monte-Carlo simulation involving the generation of $2 \times 10^{10}$ real Gaussian random variables, as well as via its normal and saddlepoint approximations w.r.t. $n_s$ and w.r.t. $n_b$, in which all expectations that need to be evaluated numerically are computed via Monte-Carlo averaging $4 \times 10^8$ real Gaussian random variables.

We see from the figure that the saddlepoint approximation w.r.t. $n_s$ is accurate for $n_b$ as high as 36, which corresponds to $n_c = 8$ symbols per block, optimally split into $n_p = 3$ pilots and $n_s = 5$ data symbols. This implies that $n_s = 5$ is sufficient for the saddlepoint approximation w.r.t. $n_s$ to be accurate for this setup. If $n_b$ is increased further, and, hence, $n_c$ is reduced,

this approximation loses accuracy. The normal approximation w.r.t. $n_s$ is accurate only for $n_b \leq 8$.

Moving to the approximations w.r.t. to $n_b$, we note that the normal approximation does not provide accurate results even when $n_b = 48$. The saddlepoint approximation slightly underestimates the required transmit power $\rho$ for $n_b = 1$, but, perhaps surprisingly, returns accurate results already for $n_b$ as small as 2.

Note finally that for the scenario considered in the figure, the required $\rho$ is large for $n_b = 1$ and decreases rapidly until $n_b = 24$, after which it increases again. This behavior can be explained as follows. Increasing $n_b$ for a fixed product $n_b n_c$ yields an increase of the number of diversity branches, which is beneficial, but also of the total number of pilot symbols $n_p n_b$ which is detrimental because it increases the effective rate of the channel code one can use to protect the information bits. The first effect dominates for $n_b \leq 24$, whereas the second effect dominates when $n_b > 24$.

*Complexity:* To address the third question, we assume again that all expectations that need to be evaluated numerically in the normal and saddlepoint approximations are computed via Monte-Carlo averaging. Since the channel and its estimate are jointly Gaussian random variables, as a proxy for numerical complexity, we count the minimum number of real Gaussian random variables that need to be generated to guarantee that the normalized mean squared difference between the error-probability value returned by the considered approximation and the actual error probability bound in (4) is less than a given threshold. Specifically, we compute each approximation $\mathsf{N}_{\text{sim}}$ times for the $\rho$ value achieving $\epsilon_{\text{ub}} = 10^{-5}$ in (4), and let $\epsilon_{\text{app}}^{(i)}(\mathsf{N})$ be the error-probability estimate obtained in the $i$th trial, when evaluating the considered approximation for the case in which the Monte-Carlo averaging is performed using $\mathsf{N}$ real Gaussian random variables. The normalized mean-squared difference is evaluated as follows:

$$e(\mathsf{N}) = \frac{1}{\mathsf{N}_{\text{sim}}} \sum_{i=1}^{\mathsf{N}_{\text{sim}}} \left(\frac{\epsilon_{\text{ub}} - \epsilon_{\text{app}}^{(i)}(\mathsf{N})}{\epsilon_{\text{ub}}}\right)^2. \qquad (48)$$

Clearly the smaller $e(\mathsf{N})$, the higher the accuracy.

In Fig. 2, we report the smallest value of $\mathsf{N}$ necessary to guarantee that $e(\mathsf{N})) \leq 0.5\%$ when $\mathsf{N}_{\text{sim}} = 100$, as a function of $n_b$. In the figure, we assumed that $n_b n_c = 288$, R = 0.104 bit per channel use, and that a target error probability of $10^{-5}$ needs to be guaranteed for all values of $n_b$. The transmit power is set according to the RCUs curve in Fig. 1. As shown in Fig. 2, the saddlepoint approximations w.r.t. $n_s$ requires around $2 \times 10^7$ real Gaussian samples. This is not surprising, since we want to evaluate accurately, via a Monte-Carlo procedure, an error probability of $10^{-5}$. Although the number of random variables that need to be generated increases with $n_b$, this increase translates in a larger value of $\mathsf{N}$ only for small values of $n_b$.

On the contrary, the saddlepoint w.r.t. $n_b$ requires only around $10^5$ samples whenever $n_b \geq 4$. This is around two orders of magnitude fewer samples than the saddlepoint w.r.t. $n_s$. This suggests that the complexity of the Monte-Carlo
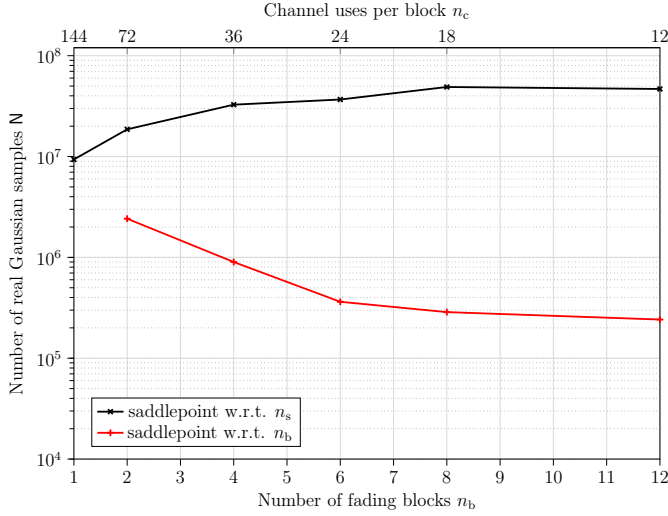
Fig. 2: Required number of real Gaussian samples to guarantee that $e(\mathrm{N}) \leq 0.5\%$. Here $n_\mathrm{b}n_\mathrm{c} = 288$, $\mathrm{R} = 0.104$ bit per channel use, $\mathrm{N}_\mathrm{sim} = 100$, and $\epsilon = 10^{-5}$.

procedure required to evaluate numerically the expectations in (36), (37), and (38), to the level of accuracy considered in this experiment, is much smaller than the complexity of the Monte-Carlo procedure required to evaluate numerically the expectation over $\boldsymbol{H}$ and $\hat{\boldsymbol{H}}$ in (8).

We do not report the complexity of the normal approximations since they do not achieve the targeted $e(\mathrm{N})$ because of their limited accuracy.

## III. MASSIVE MIMO NETWORK

In this section, we consider a multiuser massive MIMO cellular network with L cells, each served by a BS with M antennas. We assume there are K single-antenna users in each cell and focus on uplink transmission. As in Section II-A, we consider transmission over memoryless block-fading channels and use $n_\mathrm{c}$ and $n_\mathrm{b}$ to denote the number of symbols per block and the number of blocks spanned by each transmitted packet, respectively. We denote by $\boldsymbol{H}_{\ell,i,k}^j \in \mathbb{C}^\mathrm{M}$ the channel gain vector within the $\ell$th fading block between user $k$ in cell $i$ and the BS in cell $j$. We consider a spatially correlated Rayleigh fading model where $\boldsymbol{H}_{\ell,i,k}^j \sim \mathcal{CN}(\mathbf{0}_\mathrm{M}, \mathsf{R}_{i,k}^j)$. The normalized trace $\beta_{i,k}^j = \mathrm{tr}(\mathsf{R}_{i,k}^j)/\mathrm{M}$ determines the average channel gain between user $k$ in cell $i$ and the BS in cell $j$, while the eigenstructure of $\mathsf{R}_{i,k}^j$ describes its spatial channel correlation [9, Sec. 2.2].

### A. Uplink pilot transmission

The $n_\mathrm{p}$-dimensional pilot sequence transmitted by user $k$ in cell $j$ during fading block $\ell$ is denoted by the vector $\mathbf{x}_{\ell,j,k}^{(p)} \in \mathbb{C}^{n_\mathrm{P}}$. We assume that this vector satisfies $\|\mathbf{x}_{\ell,j,k}^{(p)}\|^2 = n_\mathrm{p}\rho$. Furthermore, we assume that the LK users employ mutually orthogonal pilot sequences during each fading block. In particular, we set $n_\mathrm{p} = \mathrm{LK}$. During the pilot-transmission

phase, the received signal $\mathbb{Y}_{\ell,j}^{(p)} \in \mathbb{C}^{\mathrm{M} \times n_\mathrm{p}}$ at the BS serving cell $j$ for fading block $\ell$ is given by

$$
\mathbb{Y}_{\ell,j}^{(p)} = \sum_{k=1}^\mathrm{K} \boldsymbol{H}_{\ell,j,k}^j \left(\mathbf{x}_{\ell,j,k}^{(p)}\right)^T
$$
$$
+ \sum_{i=1,i\neq j}^\mathrm{L} \sum_{k=1}^\mathrm{K} \boldsymbol{H}_{\ell,i,k}^j \left(\mathbf{x}_{\ell,i,k}^{(p)}\right)^T + \mathbb{W}_{\ell,j}^{(p)} \quad (49)
$$

where $\mathbb{W}_{\ell,j}^{(p)} \in \mathbb{C}^{\mathrm{M} \times n_\mathrm{p}}$ is the additive noise with i.i.d. $\mathcal{CN}(0, \sigma^2)$ entries.

We assume that the BS knows $\mathsf{R}_{i,k}^j$, and that it can compute the MMSE channel estimates [9, Sec. 3.2]

$$
\hat{\boldsymbol{H}}_{\ell,i,k}^j = \mathsf{R}_{i,k}^j \mathsf{Q}_{\ell,i,k}^j \left(\mathbb{Y}_{\ell,j}^{(p)} \left(\mathbf{x}_{\ell,i,k}^{(p)}\right)^*\right) \quad (50)
$$

where

$$
\mathsf{Q}_{\ell,i,k}^j = \left(\sum_{i'=1}^\mathrm{L} \sum_{k'=1}^\mathrm{K} \mathsf{R}_{i',k'}^j \left(\mathbf{x}_{\ell,i',k'}^{(p)}\right)^H \mathbf{x}_{\ell,i,k}^{(p)} + \sigma^2 \mathsf{I}_\mathrm{M}\right)^{-1}.
$$
$$ (51) $$

### B. Uplink data transmission

To decode the signal transmitted from user $k$ in cell $j$ over the $\ell$th fading block, which we denote by $\mathbf{x}_{\ell,j,k} \in \mathbb{C}^{n_\mathrm{s}}$, where $n_\mathrm{s} = n_\mathrm{c} - n_\mathrm{p}$, the BS serving cell $j$ uses the combining vector $\boldsymbol{V}_{\ell,j,k} \in \mathbb{C}^\mathrm{M}$ to compute the received vector $\boldsymbol{Y}_{\ell,j,k} \in \mathbb{C}^{n_\mathrm{s}}$ as follows:

$$
\boldsymbol{Y}_{\ell,j,k} = \left(\boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,j,k}^j\right) \mathbf{x}_{\ell,i,k} + \sum_{k'=1,k'\neq k}^\mathrm{K} \left(\boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,j,k'}^j\right) \mathbf{x}_{\ell,j,k'}
$$
$$
+ \sum_{i=1,i\neq j}^\mathrm{L} \sum_{k'=1}^\mathrm{K} \left(\boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,i,k'}^j\right) \mathbf{x}_{\ell,i,k'} + \boldsymbol{V}_{\ell,j,k}^H \mathbb{W}_{\ell,j}. \quad (52)
$$

Here, $\mathbb{W}_{\ell,j} \in \mathbb{C}^{\mathrm{M} \times n_\mathrm{s}}$ is the additive Gaussian noise on the $\ell$th fading block at the BS serving cell $j$ with i.i.d. $\mathcal{CN}(0, \sigma^2)$ entries.

We assume that the BS uses multicell-MMSE combiners, i.e.,

$$
\boldsymbol{V}_{\ell,j,k} = \left(\sum_{i=1}^\mathrm{L} \sum_{k'=1}^\mathrm{K} \hat{\boldsymbol{H}}_{\ell,i,k'}^j \left(\hat{\boldsymbol{H}}_{\ell,i,k'}^j\right)^H + \mathsf{Z}_{\ell,j}\right)^{-1} \hat{\boldsymbol{H}}_{\ell,j,k}^j
$$
$$ (53) $$

where

$$
\mathsf{Z}_{\ell,j} = \sum_{i=1}^\mathrm{L} \sum_{k=1}^\mathrm{K} \rho n_\mathrm{p} \mathsf{R}_{i,k}^j \mathsf{Q}_{\ell,i,k}^j \mathsf{R}_{i,k}^j + \frac{\sigma^2}{\rho} \mathsf{I}_\mathrm{M}. \quad (54)
$$

Note that (52) has the same form as (2). Indeed, set $\boldsymbol{Y}_\ell = \boldsymbol{Y}_{\ell,j,k}$, $\mathbf{x}_\ell = \mathbf{x}_{\ell,j,k}$, $H_\ell = \boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,j,k}^j$, $\hat{H}_\ell = \boldsymbol{V}_{\ell,j,k}^H \hat{\boldsymbol{H}}_{\ell,j,k}^j$, and $\boldsymbol{W}_\ell = \sum_{k'=1,k'\neq k}^\mathrm{K} \boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,j,k'}^j \mathbf{x}_{\ell,j,k'} + \sum_{i=1,i\neq j}^\mathrm{L} \sum_{k'=1}^\mathrm{K} \boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,i,k'}^j \mathbf{x}_{\ell,i,k'} + \boldsymbol{V}_{\ell,j,k}^H \mathbb{W}_{\ell,j}$. Note also that, given $\{\boldsymbol{H}_{\ell,i,k'}^j, \hat{\boldsymbol{H}}_{\ell,i,k'}^j\}$, the entries of the newly defined

vector $\boldsymbol{W}_\ell$ are conditionally i.i.d. and follow a $\mathcal{CN}(0, \sigma_\ell^2)$ distribution, with

$$\sigma_\ell^2 = \sigma^2 \|\boldsymbol{V}_{\ell,j,k}\|^2 + \rho \sum_{k'=1, k' \neq k}^{\mathrm{K}} \left| \boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,j,k'}^j \right|^2$$
$$+ \rho \sum_{i=1, i \neq j}^{\mathrm{L}} \sum_{k'=1}^{\mathrm{K}} \left| \boldsymbol{V}_{\ell,j,k}^H \boldsymbol{H}_{\ell,i,k'}^j \right|^2. \tag{55}$$

Hence, we can evaluate the uplink per-user error probability by using the information-theoretic bound in (4) and its normal and saddlepoint approximations discussed in Section II.

### C. Numerical Experiments

*1) Setup:* Our simulation setup consists of L square cells, each of size $75\,\mathrm{m} \times 75\,\mathrm{m}$, containing K users each. The BSs, which are equipped with a uniform linear array with M antenna elements separated by half a wavelength, are placed in the center of each cell. The antennas and the users are located in the same horizontal plane. Thus, the azimuth angle is sufficient to determine the directivity. We assume that the scatterers are uniformly distributed in the angular interval $[\varphi_{i,k} - \Delta, \varphi_{i,k} + \Delta]$, where $\varphi_{i,k}$ is the nominal angle of arrival of user $k$ in cell $i$ and $\Delta$ is the angular spread, which we set to $\Delta = 25°$. The $(m_1, m_2)$th entry of the matrix $\mathsf{R}_{i,k}^j$ is then given by [9, Sec. 2.6]

$$\left[ \mathsf{R}_{i,k}^j \right]_{m_1, m_2} = \frac{\beta_{i,k}^j}{2\Delta} \int_{-\Delta}^{\Delta} e^{\mathrm{j}\pi(m_1 - m_2)\sin(\varphi_{i,k} + \bar\varphi)} \mathrm{d}\bar\varphi. \tag{56}$$

Here, $\beta_{i,k}^j$ denotes the large-scale fading coefficient measured in dB

$$\beta_{i,k}^j = -35.3 - 37.6 \log_{10}\left( \frac{d_{i,k}^j}{1\,\mathrm{m}} \right) \tag{57}$$

with $d_{i,k}^j$ being the distance between the BS in the cell $j$ and the user $k$ in cell $i$. The communication takes place over a $20\,\mathrm{MHz}$ bandwidth with a total receiver noise power of $\sigma^2 = -94\,\mathrm{dBm}$ consisting of thermal noise and a noise figure of $7\,\mathrm{dB}$ in the receiver hardware.

In the next two subsections, we extend the accuracy and complexity analysis of the error-probability approximations performed for a SISO link in Section II-E to the massive MIMO uplink. Since repeating the study carried out for SISO is unfeasible in a multi-cell multi-user MIMO setting, because of complexity constraints, we first focus in Section III-C2 on a single-cell massive MIMO network with two users. We will then provide in Section III-C3 an extension of this analysis to the multi-cell multiuser massive MIMO network for the special case in which the number of blocks $n_{\mathrm{b}}$ is equal to 3.

*2) Accuracy and Complexity Analysis for the Two-User Case:* We consider the uplink of a single-cell massive MIMO network in which the BS serves two users (L = 1 and K = 2). The distance between the two users and the BS is $d_{1,1}^1 = d_{1,2}^1 = 36.4\,\mathrm{m}$. The nominal angle of user 1 w.r.t. the BS is $30°$, and the nominal angle of user 2 w.r.t. the BS is $40°$. We also assume that orthogonal pilot sequences are assigned to each user, that $n_{\mathrm{p}} = 2$, and that MMSE spatial combining based on MMSE channel estimation is used at the
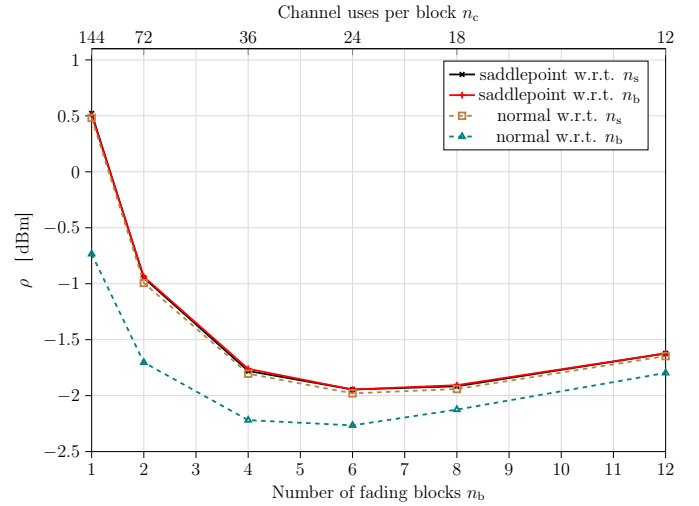


Fig. 3: Two-user, single-cell massive MIMO scenario: required transmit power $\rho$ to achieve $\epsilon = 10^{-5}$. Here, $n_{\mathrm{b}} n_{\mathrm{c}} = 144$, R = 2 bit per channel use, and $n_{\mathrm{p}} = 2$.
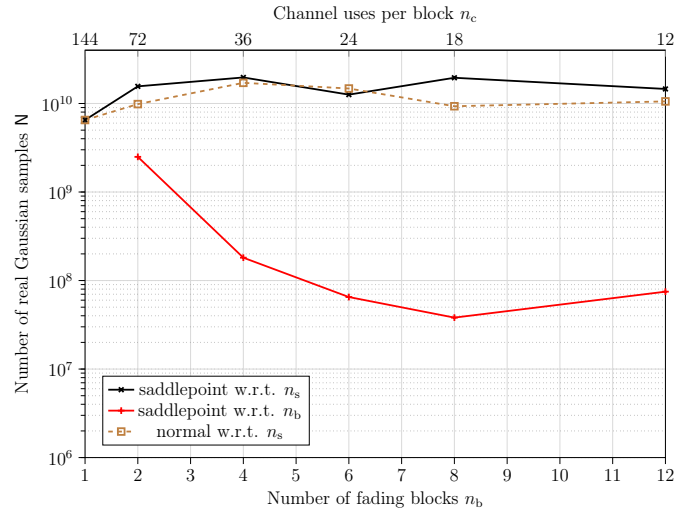


Fig. 4: Two-user, single-cell massive MIMO scenario: required number of real Gaussian samples to guarantee that $e(\mathrm{N}) \leq 0.5\%$. Here, $n_{\mathrm{b}} n_{\mathrm{c}} = 144$, R = 2 bit per channel use, $n_{\mathrm{p}} = 2$, $\mathrm{N}_{\mathrm{sim}} = 100$, and $\epsilon = 10^{-5}$.

BS. Finally, we set $n_{\mathrm{b}} n_{\mathrm{c}} = 144$ and R = 2 bit per channel use, which corresponds to 288 bits per packet.

In Fig. 3 we report the smallest $\rho$ value needed to achieve an error probability of $10^{-5}$, as a function of $n_{\mathrm{b}}$. All curves in the figures are obtained by performing a Monte-Carlo simulation involving the generation of $8 \times 10^{10}$ real Gaussian random variables. We observe that both saddlepoint approximations as well as the normal approximation over $n_{\mathrm{s}}$ agree and are therefore assumed to be accurate for the $n_{\mathrm{b}}$ values considered in the figure, including $n_{\mathrm{b}} = 1$. On the contrary, the normal approximation w.r.t. $n_{\mathrm{b}}$ is not does not appear to be accurate. Note that, because of the large spatial diversity available in this setup, increasing $n_{\mathrm{b}}$ from 2 to 6 has only a limited benefit in terms of $\rho$ and increasing $n_{\mathrm{b}}$ beyond 6 is actually deleterious, because of the reduction in the number of channel uses per block available for data transmission.

Focusing on both saddlepoint approximations and on the

normal approximation w.r.t. $n_s$, we illustrate in Fig. 4, the minimum number of real Gaussian samples that need to be generated in the Monte-Carlo step required in all approximations, to guarantee that $e(N) < 0.5\%$ for a target error probability of $10^{-5}$. Note that, unlike the SISO case, evaluating the transmit power required to achieve $\epsilon = 10^{-5}$ with the RCUs bound (4) is not feasible due to its computational complexity. Thus, in our complexity analysis, the transmit power is set to the arithmetic average of the saddlepoint approximation curves in Fig. 3.

We see from Fig. 4 that the number of real Gaussian samples required by all approximations is more than two orders of magnitude larger than in the SISO case (cf. Fig 2). This is expected since the channel within each fading block is now characterized by 200 (dependent) complex Gaussian random variables, instead of the single complex Gaussian random variable needed in the SISO case. We also see that the saddlepoint approximation w.r.t. $n_b$ requires between 1 and 2 orders of magnitude fewer samples than both normal approximation and saddlepoint approximation w.r.t. $n_s$. This observation is also in agreement with the results presented in Section II-E for the SISO case.

*3) Multi-Cell Multi-User Setup:* We finally consider a massive MIMO network consisting of L = 4 cells and K = 10 users and consider a wrap-around topology (for details, see [9, Sec. 4.1.3]). We assume for simplicity that the users within each cell are regularly spaced on a circle around the BS of radius $d_{j,k}^j = 33.75\,\mathrm{m}$. We consider a scenario in which $n_b = 3$ and $n_c = 70$, assume that all users transmit orthogonal pilots over each block and set $n_p = 40$. Finally, we set R = 2 bit per channel use, and consider a target error probability of $10^{-5}$. This target error probability is achieved by $\rho = -1.63\,\mathrm{dBm}$ for the saddlepoint approximation w.r.t. $n_b$, $\rho = -1.61\,\mathrm{dBm}$ for the saddlepoint approximation w.r.t. $n_s$, $\rho = -1.82\,\mathrm{dBm}$ the for the normal approximation w.r.t. $n_b$ and $\rho = -1.70\,\mathrm{dBm}$ for the normal approximation w.r.t. $n_s$. These values of $\rho$ are estimated using a Monte-Carlo procedure involving $10^{12}$ real Gaussian random variables. Similar to the 2-user massive MIMO case, the RCUs bound cannot be evaluated due to its computational complexity. The reported results suggest that, in agreement with the results obtained for the SISO and for the two-user massive MIMO cases, both saddlepoint approximations are accurate, whereas both normal approximations are not accurate.

To assess the complexity of the two saddlepoint approximation, we take the arithmetic average of the transmit power evaluated with the two saddlepoint approximations as reference transmit power. In Fig. 5, we depict the normalized mean-square difference $e(N)$, defined in (48), as a function of the number of real Gaussian samples N used in the Monte-Carlo step required for both saddlepoint approximations. We can observe that the saddlepoint approximation w.r.t. $n_b$ requires approximately 30 times fewer samples than the saddlepoint approximation w.r.t. $n_s$ to achieve $e(N) \leq 0.5\%$. This is again in accordance with the results reported in Fig. 2 for the SISO case, and Fig. 4 for the single-cell, two-user massive MIMO case.
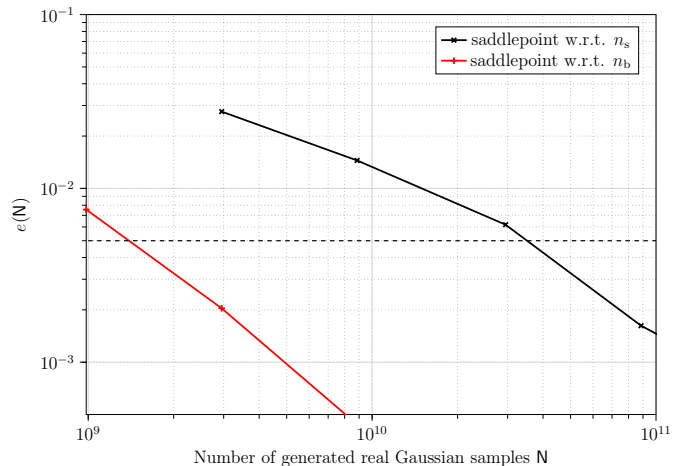


Fig. 5: Multi-cell, multi-user massive MIMO scenario. Here, $n_s = 30$, $n_p = 40$, $n_b = 3$, K = 10, L = 4, $N_{sim} = 100$, R = 2 bit per channel use, and $\epsilon = 10^{-5}$.

## IV. CONCLUSION

We presented numerically efficient methods to evaluate an upper bound on the error probability achievable over SISO and massive MIMO memoryless block-fading channels when pilot-assisted transmission, scaled nearest-neighbor decoding, and i.i.d. Gaussian codebooks are used. Our methods include both normal and saddlepoint approximations w.r.t. to the number of data symbols per block $n_s$, as well as novel normal and saddlepoint approximations w.r.t. the number of fading blocks $n_b$ spanned by each codeword. All approximations involve the numerical evaluations of expectations that are not known in closed form and can be evaluated using Monte-Carlo methods. Our numerical experiments reveal that the saddlepoint approximation w.r.t. to $n_b$ yield accurate estimates of the error probability in URLLC scenarios of practical relevance. Furthermore, it involves a numerical complexity (measured in terms of total number of Monte-Carlo samples required to achieve a given accuracy) roughly two orders of magnitude lower than the complexity of the saddlepoint approximations in $n_s$. This holds for a variety of scenarios ranging from SISO to multicell, multiuser massive MIMO. Hence, this approximation should be preferred when evaluating error probabilities within URLLC optimization routines such as resource-allocation and scheduling algorithms. The normal approximations w.r.t. $n_b$ and $n_s$ are not viable alternatives as they often provide inaccurate results for the scenarios considered in this paper, at no advantage in terms of complexity compared to the saddlepoint approximation.

## REFERENCES

[1] A. O. Kislal, A. Lancho, G. Durisi, and E. Ström, "Efficient evaluation of the error probability for pilot-assisted finite-blocklength transmission," in *Proc. Asilomar Conf. Signals, Syst., Comput*, Pacific Grove, CA, U.S.A., Nov. 2022.

[2] H. Ren, C. Pan, Y. Deng, M. Elkashlan, and A. Nallanathan, "Joint power and blocklength optimization for URLLC in a factory automation scenario," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 1786–1801, Mar. 2020.

[3] X. Song and M. Yuan, "Performance analysis of one-way highway vehicular networks with dynamic multiplexing of eMBB and URLLC traffics," *IEEE Access*, vol. 7, pp. 118 020–118 029, 2019.

[4] D. Van Den Berg, R. Glans, D. De Koning, F. A. Kuipers, J. Lugtenburg, K. Polachan, P. T. Venkata, C. Singh, B. Turkovic, and B. Van Wijk, "Challenges in haptic communications over the tactile internet," *IEEE Access*, vol. 5, pp. 23 502–23 518, 2017.

[5] G. Durisi, T. Koch, and P. Popovski, "Towards massive, ultra-reliable, and low-latency wireless communication with short packets," *Proc. IEEE*, vol. 104, no. 9, pp. 1711–1726, Sep. 2016.

[6] M. Hayashi, "Information spectrum approach to second-order coding rate in channel coding," *IEEE Trans. Inf. Theory*, vol. 55, no. 11, pp. 4947–4966, Nov. 2009.

[7] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2307–2359, May 2010.

[8] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, Nov. 2010.

[9] E. Björnson, J. Hoydis, and L. Sanguinetti, "Massive MIMO networks: Spectral, energy, and hardware efficiency," *Foundations and Trends® in Signal Processing*, vol. 11, no. 3-4, pp. 154–655, 2018.

[10] J. Östman, G. Durisi, E. G. Ström, M. C. Coşkun, and G. Liva, "Short packets over block-memoryless fading channels: Pilot-assisted or noncoherent transmission?" *IEEE Trans. Commun.*, vol. 67, no. 2, pp. 1521–1536, Feb. 2019.

[11] A. Lancho, J. Östman, G. Durisi, T. Koch, and G. Vazquez-Vilar, "Saddlepoint approximations for short-packet wireless communications," *IEEE Trans. Wireless Commun.*, vol. 19, no. 7, pp. 4831–4846, Jul. 2020.

[12] J. Östman, A. Lancho, G. Durisi, and L. Sanguinetti, "URLLC with massive MIMO: Analysis and design at finite blocklength," *IEEE Trans. Wireless Commun.*, vol. 20, no. 10, pp. 6387–6401, Oct. 2021.

[13] A. Martinez and A. Guillén i Fàbregas, "Saddlepoint approximation of random-coding bounds," in *Proc. Inf. Theory Appl. Workshop*, San Diego, CA, USA, Feb. 2011, pp. 257–262.

[14] A. Lapidoth and S. Shamai (Shitz), "Fading channels: how perfect need "perfect side information" be?" *IEEE Trans. Inf. Theory*, vol. 48, no. 5, pp. 1118–1134, May 2002.

[15] J. Font-Segura, G. Vazquez-Vilar, A. Martinez, A. Guillén i Fàbregas, and A. Lancho, "Saddlepoint approximations of lower and upper bounds to the error probability in channel coding," in *Proc. Conf. Inf. Sci. Sys. (CISS)*, Princeton, NJ, Mar. 2018.

[16] J. L. Jensen, *Saddlepoint approximations*. Oxford, U.K.: Oxford Univ. Press, 1995.

[17] G. Taricco, "A simple method to calculate random-coding union bounds for ultra-reliable low-latency communications," *IEEE Wireless Commun. Lett.*, Feb. 2022.

[18] P. Yuan, M. C. Coşkun, and G. Kramer, "Polar-coded non-coherent communication," *IEEE Commun. Lett.*, vol. 25, no. 6, pp. 1786–1790, Jun. 2021.

[19] J. Scarlett, V. Y. F. Tan, and G. Durisi, "The dispersion of nearest-neighbor decoding for additive non-Gaussian channels," *IEEE Trans. Inf. Theory*, vol. 63, no. 1, pp. 81–92, Jan. 2017.

[20] W. Feller, *An introduction to probability theory and its applications*. Wiley, 1971, vol. 2.

[21] J. Scarlett, A. Martinez, and A. Guillén i Fàbregas, "Mismatched decoding: Error exponents, second-order rates and saddlepoint approximations," *IEEE Trans. Inf. Theory*, vol. 60, no. 5, pp. 2647–2666, May 2014.