# Soft-input, soft-output joint detection and GRAND

Hadi Sarieddeen
*Research Laboratory of Electronics*
*Massachusetts Institute of Technology*
Cambridge, MA 02139, USA
hadisari@mit.edu

Muriel Médard
*Research Laboratory of Electronics*
*Massachusetts Institute of Technology*
Cambridge, MA 02139, USA
medard@mit.edu

Ken. R. Duffy
*Hamilton Institute*
*Maynooth University*
Ireland
ken.duffy@mu.ie

*Abstract*—Guessing random additive noise decoding (GRAND) is a maximum likelihood (ML) decoding method that identifies the noise effects corrupting code-words of arbitrary codebooks. In a joint detection and decoding framework, this work demonstrates how GRAND can leverage crude soft information in received symbols and channel state information to generate, through guesswork, soft bit reliability outputs in log-likelihood ratios (LLRs). The LLRs are generated via successive computations of Euclidean-distance metrics corresponding to candidate noise-recovered words. Noting that the entropy of noise is much smaller than that of information bits, a small number of noise effect guesses generally suffices to hit a code-word, which allows generating LLRs for critical bits; LLR saturation is applied to the remaining bits. In an iterative (turbo) mode, the generated LLRs at a given soft-input, soft-output GRAND iteration serve as enhanced a priori information that adapts noise-sequence guess ordering in a subsequent iteration. Simulations demonstrate that a few turbo-GRAND iterations match the performance of ML-detection-based soft-GRAND in both AWGN and Rayleigh fading channels at a complexity cost that, on average, grows linearly (instead of exponentially) with the number of symbols.

*Index Terms*—GRAND, soft-GRAND, turbo-GRAND

## I. INTRODUCTION

The upcoming sixth generation (6G) of wireless communications promises to support a plethora of data-demanding and delay-sensitive applications [1], requiring both ultra-broadband high-frequency connectivity [2] and ultra-reliable low-latency communication (URLLC) [3]. Such variety in requirements compels a paradigm shift from structured, code-specific channel-code decoding to universal and practical decoding that is efficient for different code rates and lengths. Early universal near-maximum-likelihood (ML) decoders for linear codes adopted a list-decoding principle [4], [5]. Recently, guessing random additive noise decoding (GRAND) [6], [7] is celebrated as a novel and practical universal decoder suited for block-code constructions of moderate redundancy.

GRAND is a capacity-achieving channel-code decoder that has demonstrated ML decoding performance on arbitrary (even unstructured) code-books. Instead of directly identifying the transmitted code-word, GRAND aims at identifying the noise effect that corrupts the code-word; it successively reverses the noise effects from the received signal to recover candidate transmitted words. By leveraging information on channel and noise models, the candidate noise sequences are ordered and queried in decreasing likelihood, guaranteeing the first recovered code-word to be the ML decoding solution, even for channels with memory in the absence of interleaving. The guesswork literature [8]–[10] establishes the computational feasibility of GRAND for all moderate redundancy codes, where the Shannon entropy rate of noise is typically less than that of information symbols [7]. Furthermore, GRAND's computational efficiency and modularity have resulted in highly efficient circuit design, as demonstrated in a recent 65 nm [11] synthesis and a 40 nm [12] CMOS implementation.

Incorporating soft-detection symbol reliability information into decoding decisions enhances decoding accuracy [13], [14]. At one end, GRAND can leverage as soft information a one-bit mask designating whether a channel use is reliable or not [15]; specifying reliable bits via a channel-fading-induced mask is also proposed in [16]. At the other end, the complete information in continuous channel outputs serves as soft symbol-reliability information in the soft-GRAND (SGRAND) scheme [17]. A compromise between one-bit soft-GRAND and SGRAND is ordered reliability bits GRAND (ORBGRAND) [18], which matches the decoding accuracy of SGRAND through code-book-independent quantization of soft information in a hardware-friendly algorithm.

With fading channels, generating high-resolution soft information (as opposed to masks [16]) through exhaustive ML detection is computationally demanding, and low-complexity soft-output detectors only generate sub-optimal LLRs. In the absence of soft information (or with low-quality information), iterative decoding schemes can intrinsically generate soft-decoding reliability information to be fed as soft-input decoding information in subsequent iterations [19]. Such information can be updated over both the detection and decoding iterations [20], introducing more degrees of freedom in versatile, adaptive communication systems. However, how to generate soft-reliability outputs through GRAND remains unclear. Consequently, iterative soft-input, soft-output (SISO) GRAND (turbo-GRAND) has not yet been investigated.

In this work, we propose a variation of GRAND that does not avail of input soft-detection bit reliability information but leverages complex received symbols (soft-information in a crude, unprocessed form), channel state information (CSI), and demapped bits (linear detector outputs) to generate soft-

decoding bit-reliability information in log-likelihood ratios (LLRs). We compute the LLRs by populating Euclidean distance metrics corresponding to a list of candidate words–not necessarily code-words–resulting from noise guesswork. Because the number of guesses before hitting a code-word is limited, we cannot compute LLRs for all bits; we propose LLR thresholding. We propose a turbo-GRAND scheme in which the generated extrinsic LLRs are fed as a priori LLRs in subsequent SISO decoding iterations. With access to complex received vectors and continuous CSI, the Euclidean distance computations are common to both detection and decoding; turbo-GRAND thus realizes joint detection and decoding.

The paper is organized as follows. The problem formulation is first presented in Sec. II. Then, the proposed turbo-GRAND scheme is detailed in Sec. III. Performance and complexity results are reported in Sec. IV; conclusions are drawn in Sec. V. Regarding notation, bold upper case, bold lower case, and lower case letters correspond to matrices, vectors, and scalars, respectively. Scalar and vector $L_2$ norms are denoted by $|\cdot|$ and $\|\cdot\|$, respectively. $\mathsf{E}[\cdot]$, $(\cdot)^T$, and $(\cdot)^*$, stand for the expected value, transpose, and conjugate transpose, respectively. $\mathbf{I}_M$ is an identity matrix of size $M$, $\mathbf{0}_N$ is a vector of zeros of size $N$, $\mathbb{F}_u$ denotes a Galois field with $u$ elements, $\Pr(\cdot)$ is the probability function, and $\odot$ is the Hadamard product.

## II. System Model and Problem Formulation

### A. System Model

We consider a communication system of equivalent baseband input-output relation, $\mathbf{y} = \mathbf{Hx} + \mathbf{n}$, where $\mathbf{y} \in \mathbb{C}^{M\times 1}$ is the received symbol vector, $\mathbf{H} \in \mathbb{C}^{M\times M}$ is the channel matrix, $\mathbf{x} = [x_1 \cdots x_i \cdots x_M]^T \in \mathbb{C}^{M\times 1}$ is the transmitted symbol vector, and $\mathbf{n} = [n_1 \cdots n_i \cdots n_M]^T \in \mathbb{C}^{M\times 1}$ is the additive–possibly colored–noise vector $\left(\mathsf{E}[n_i n_i^*] = \sigma_i^2\right)$. Note that $\mathbf{H}$ can be an identity matrix in an additive white Gaussian noise (AWGN) system, a diagonal matrix under point-to-point fading, or a complete matrix under spatial diversity/multiplexing schemes. Furthermore, we assume the information symbol, $x_i$, to belong to a normalized complex constellation, $\mathcal{X}_i$ $\left(\mathsf{E}[x_i^* x_i] = 1\right)$. Consequently, $\mathbf{x} \in \bar{\mathcal{X}} \subset \mathbb{C}^{M\times 1}$, where $\bar{\mathcal{X}}$ is the finite $M$-dimensional lattice of all possible modulated symbol vectors. The bit-representation of $x_i$ is $\mathbf{c}_i = [c_{i,1} \cdots c_{i,j} \cdots c_{i,q_i}]^T \in \mathbb{F}_2^{q_i}$, where $q_i = \lceil \log_2(|\mathcal{X}_i|) \rceil$. The bit-representation of $\mathbf{x}$ is thus $\mathbf{c} = [\mathbf{c}_1 \cdots \mathbf{c}_i \cdots \mathbf{c}_M]^T \in \mathbb{F}_2^N$, where $N = \sum_{i=1}^M q_i$. We assume $\mathbf{c}$ to be a code-word encoded with an error correcting code $\alpha : \mathbb{F}_2^K \rightarrow \mathbb{F}_2^N$, of code-rate $R = K/N$. A code-book $C \triangleq \{\mathbf{c} : \mathbf{c} = \alpha(\mathbf{b}), \mathbf{b} \in \mathbb{F}_2^K\}$ includes all possible code-words, where $\mathbf{b}$ is the uncoded bit vector. We denote by $\mathbf{v} \in \mathbb{R}^{+N}$ a vector of $\sigma_{i,j}^2$ values, the second-order noise statistics per bit. At the receiver side, assuming perfect CSI, a hard detector, $\beta : \mathbb{C}^M \rightarrow \bar{\mathcal{X}}$, equalizes the channel and recovers a symbol vector, $\hat{\mathbf{x}}$, from $\mathbf{y}$; a demapper recovers a word, $\hat{\mathbf{c}}$, from $\hat{\mathbf{x}}$.

### B. Problem Formulation

The ML decoder computes the conditional probability of the demapped word, $\hat{\mathbf{c}}$, for each of the $2^K$ code-words, $\mathbf{c}$, in code-book, $C$. The $\mathbf{c}$ with the highest conditional likelihood of trans-

---

**Algorithm 1** Hard GRAND
**Input:** Demapped bits $\hat{\mathbf{c}}$; ordered noise-generating function $\Pi$; abandonment threshold $B$
**Output:** Decoded $\bar{\mathbf{c}}^{\mathrm{GRAND}}$
1: $k \leftarrow 0$
2: **while** $k < 2^N$ **do**
3: $\quad k \leftarrow k + 1$
4: $\quad \mathbf{w} \leftarrow \Pi(k)$ $\qquad\qquad \triangleright$ $k$th likely noise sequence
5: $\quad$ **if** $\hat{\mathbf{c}} \ominus \mathbf{w} \in C$ or $k = B$ **then**
6: $\qquad \bar{\mathbf{c}}^{\mathrm{GRAND}} \leftarrow \hat{\mathbf{c}} \ominus \mathbf{w}$
7: $\qquad$ **return** $\bar{\mathbf{c}}^{\mathrm{GRAND}}$
8: $\quad$ **end if**
9: **end while**

---

mission given what was received is the ML solution, $\bar{\mathbf{c}}^{\mathrm{ML}} = \arg\max\{\Pr(\hat{\mathbf{c}} \mid \mathbf{c}) : \mathbf{c} \in C\}$. Instead of searching code-words, GRAND searches putative, not necessarily memoryless, noise effect sequences that corrupt $\mathbf{c}$, $\mathbf{w} = [w_1 \cdots w_{i,j} \cdots w_N]^T \in \mathbb{F}_2^N$, with non-increasing probability. We express the channel's action at the bit level through function $\oplus$, where $\bar{\mathbf{c}} = \mathbf{c} \oplus \mathbf{w}$. We can write $\Pr(\bar{\mathbf{c}}|\mathbf{c}) = \Pr(\bar{\mathbf{c}} = \mathbf{c} \oplus \mathbf{w})$, and it follows that

$$\bar{\mathbf{c}}^{\mathrm{GRAND}} = \arg\max\{\Pr(\mathbf{w} = \bar{\mathbf{c}} \ominus \mathbf{c}) : \mathbf{c} \in C\}. \quad (1)$$

The receiver creates a list of noise effect sequences of decreasing order of likelihood through a noise generating function $\Pi : \{1, \cdots, 2^N\} \rightarrow \mathbf{w} \in \mathbb{F}_2^N$; the sequences are queried until the first code-word hit, $\mathbf{w} = \hat{\mathbf{c}} \ominus \mathbf{c}$ (block-code syndrome computations). GRAND is thus a ML decoder that executes Algorithm 1 and returns $\bar{\mathbf{c}}^{\mathrm{GRAND}}$; information bits are retrieved as $\bar{\mathbf{b}}^{\mathrm{GRAND}} = \alpha^{-1}(\bar{\mathbf{c}}^{\mathrm{GRAND}})$. Because the entropy of noise is small in most communication systems, GRAND is low-complexity. GRAND's efficiency is further guaranteed by abandoning guessing after a computational cut-off [7], $B$.

Soft GRAND accepts, in addition to $\hat{\mathbf{c}}$, bit-reliability information in a vector $\mathbf{\Lambda} = [\lambda_{1,1} \cdots \lambda_{i,j} \cdots \lambda_{M,q}]^T \in \mathbb{R}^N$ (assuming $q_i = q, \forall i$). We can generate a weight metric per putative noise sequence by multiplying the noise sequences by $|\mathbf{\Lambda}|$; noise sequences with smaller weights are more likely to occur. However, $\mathbf{\Lambda}$ is not always available (or available but with bad quality), as soft-output detectors are computationally demanding. We aim at generating soft-reliability outputs within GRAND, proposing a low-complexity noise-centric soft-output decoding algorithm, a function $\bar{\alpha} : \{\mathbb{F}_2^K, \mathbb{C}^M\} \rightarrow \mathbb{R}^N$, that accepts $\hat{\mathbf{c}}$ and $\mathbf{y}$ and generates output LLRs, $\bar{\mathbf{\Lambda}} = [\bar{\lambda}_{1,1} \cdots \bar{\lambda}_{i,j} \cdots \bar{\lambda}_{M,q}]^T \in \mathbb{R}^N$, with the knowledge of $\mathbf{H}$. By further integrating knowledge of noise statistics in LLR computations (scaling by noise variance), our proposal can account for noise bursts (Markovian channel noise, for example), foregoing interleaves, and whitening filters [7], [21]. The extrinsic LLRs, $\bar{\mathbf{\Lambda}}$, can then be fed as intrinsic LLRs, $\mathbf{\Lambda}$, in a subsequent soft GRAND; an iteration that is repeated in the proposed turbo-GRAND.

## III. Proposed SISO Turbo-GRAND

We propose a SISO variation of GRAND that leverages, in addition to the resources of Algorithm 1, the received symbols,

Fig. 1. A block diagram of turbo-GRAND.

CSI, and optional noise statistics, to generate extrinsic bit-reliability LLRs through joint detection and decoding. The LLRs are fed as input soft-decoding information to a subsequent SISO GRAND iteration (up to $T$ iterations), resulting in turbo-GRAND (Fig. 1). Turbo-GRAND aims to approach the performance of a soft GRAND with soft-input information from an exhaustive soft-output ML detector. Therefore, turbo-GRAND is helpful in the absence of soft-input information or the presence of sub-optimal soft information.

The LLR of bit $j$ of symbol $i$ is defined as

$$\lambda_{i,j} = \log\left(\Pr\left(c_{i,j}=1, \mathbf{y}, \mathbf{H}\right) / \Pr\left(c_{i,j}=0, \mathbf{y}, \mathbf{H}\right)\right). \quad (2)$$

Near-optimal ML-detection log-max LLRs [22] are computed by exhaustively searching the lattice $\bar{\mathcal{X}}$, computing $|\mathcal{X}_1| \times |\mathcal{X}_i| \times \cdots \times |\mathcal{X}_M|$ Euclidean distance metrics to solve for [23]

$$\lambda_{i,j}^{\mathrm{ML}} \approx \frac{1}{\sigma^2}\left(\min_{\mathbf{x}\in\bar{\mathcal{X}}^{i,j,1}} \|\mathbf{y}-\mathbf{Hx}\|^2 - \min_{\mathbf{x}\in\bar{\mathcal{X}}^{i,j,0}} \|\mathbf{y}-\mathbf{Hx}\|^2\right), \quad (3)$$

where $\bar{\mathcal{X}}^{i,j,1} \triangleq \{\mathbf{x} \in \bar{\mathcal{X}} : c_{i,j}=1\}$ and $\bar{\mathcal{X}}^{i,j,0} \triangleq \{\mathbf{x} \in \bar{\mathcal{X}} : c_{i,j}=0\}$ are subsets of symbol vectors in $\bar{\mathcal{X}}$, having in the corresponding $j$th bit of the $i$th symbol a value of 1 and 0, respectively. We have further assumed the case of white noise, $\sigma_{i,j} = \sigma, \forall i, j$. For colored noise, $\|\mathbf{y}-\mathbf{Hx}\|^2$ is replaced by $(\mathbf{y}-\mathbf{Hx})^* \Gamma^{-1}(\mathbf{y}-\mathbf{Hx})$, where $\Gamma = \mathrm{diag}(\mathbf{v})$. The hard ML detection output is $\hat{\mathbf{x}}^{\mathrm{ML}} = \arg\min_{\mathbf{x}\in\bar{\mathcal{X}}} \|\mathbf{y}-\mathbf{Hx}\|^2$. Furthermore, soft information can be extracted per symbol in linear detectors [24], which are near-optimal in point-to-point systems at a high signal-to-noise ratio (SNR) but sub-optimal under symbol interference/correlation. In particular, a zero-forcing (ZF) detector equalizes the channel by multiplying by its pseudo-inverse, $\hat{\mathbf{y}}^{\mathrm{ZF}} = (\mathbf{H}^*\mathbf{H})^{-1}\mathbf{H}^*\mathbf{y}$. The per-symbol ZF LLRs in $\mathbf{\Lambda}^{\mathrm{ZF}}$ are

$$\lambda_{i,j}^{\mathrm{ZF}} = \frac{1}{\sigma_i^{\mathrm{ZF}2}}\left(\min_{x_i\in\mathcal{X}_i^{j,1}} \left|\hat{y}_i^{\mathrm{ZF}}-x_i\right|^2 - \min_{x_i\in\mathcal{X}_i^{j,0}} \left|\hat{y}_i^{\mathrm{ZF}}-x_i\right|^2\right), \quad (4)$$

where $\mathcal{X}_i^{j,1} \triangleq \{x_i \in \mathcal{X}_i : c_{i,j}=1\}$ and $\mathcal{X}_i^{j,0} \triangleq \{x_i \in \mathcal{X}_i : c_{i,j}=0\}$ are subsets of symbols in the one-dimensional constellation $\mathcal{X}_i$, having a $j$th bit of 1 and 0, respectively, and $\sigma_i^{\mathrm{ZF}2} = \sigma_i^2(h_i^* h_i)^{-1}$ is a scaled noise variance. The ZF hard-outputs are computed per symbol as $\hat{x}_i^{\mathrm{ZF}} = \lfloor \hat{y}_i^{\mathrm{ZF}} - x_i \rfloor_{\mathcal{X}_i}$, where $\lfloor \eta \rceil_{\mathcal{X}_i} \triangleq \arg\min_{x\in\mathcal{X}_i} |\eta - x|$ is the slicing operator over $\mathcal{X}_i$. The demapped version of $\hat{\mathbf{x}}^{\mathrm{ZF}}$, $\hat{\mathbf{c}}^{\mathrm{ZF}}$, is the initial vector over which noise effects are guessed in turbo-GRAND.

Contrary to conventional ML and list decoders that query all or multiple code-words, the basic implementation of GRAND achieves ML decoding performance by querying noise sequences and recovering a single code-word. This scarcity in code-word hits across turbo-GRAND iterations is mitigated

---

**Algorithm 2** Joint detection and decoding - turbo-GRAND

**Input:** Demapped bits $\hat{\mathbf{c}} = \hat{\mathbf{c}}^{\mathrm{ZF}}$; received symbols $\mathbf{y}$; channel matrix $\mathbf{H}$; noise matrix $\mathbf{W}$; noise statistics $\mathbf{v}$; noise generating/sorting function $\Pi/\acute{\Pi}$; input LLRs $\mathbf{\Lambda}$ ($\mathbf{\Lambda} = \mathbf{\Lambda}^{\mathrm{ZF}}$ or $\mathbf{\Lambda} = \mathbf{0}_N$); abandonment threshold $B$

**Output:** Output LLRs $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}$; demapped $\hat{\mathbf{c}}^{\mathrm{GRAND}}$; decoded $\bar{\mathbf{c}}^{\mathrm{GRAND}}$

1:  $d^{\mathrm{ML}} \leftarrow \infty$; $\bar{d}^{\mathrm{ML}} \leftarrow \infty$; $\mathbf{d}^{\mathrm{cML}} \leftarrow N \odot (1/\mathbf{v})$
2:  $\hat{\mathbf{c}}^{\mathrm{GRAND}} \leftarrow \hat{\mathbf{c}}$; $\bar{\mathbf{c}}^{\mathrm{GRAND}} \leftarrow \hat{\mathbf{c}}$; $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}} \leftarrow \mathbf{\Lambda}$
3:  $t \leftarrow 0$;
4:  **while** $t < T$ **do**
5:      $\mathbf{s} \leftarrow \acute{\Pi}\left(\mathbf{W} \times \left|\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}\right|\right)$
6:      $k \leftarrow 0$
7:      **while** $k < B$ **do**
8:         $k \leftarrow k + 1$; $\mathbf{w} \leftarrow \Pi(\mathbf{s}(k))$     ▷ $k$th likely noise
9:         $\bar{\mathbf{c}} \leftarrow \hat{\mathbf{c}}^{\mathrm{GRAND}} \ominus \mathbf{w}$; $\bar{\mathbf{x}} \leftarrow \mathrm{mod}(\bar{\mathbf{c}})$
10:       $d \leftarrow (\mathbf{y}-\mathbf{H}\bar{\mathbf{x}})^* \Gamma^{-1}(\mathbf{y}-\mathbf{H}\bar{\mathbf{x}})$    ▷ $\Gamma = \mathrm{diag}(\mathbf{v})$
11:       $n \leftarrow 0$
12:       **while** $n < N$ **do**
13:          $n \leftarrow n + 1$
14:          **if** $\bar{c}_n \neq \hat{c}_n^{\mathrm{GRAND}}$ **then**   ▷ complementary bits
15:            **if** $d < d^{\mathrm{ML}}$ **then**
16:              $d_n^{\mathrm{cML}} \leftarrow d^{\mathrm{ML}}$   ▷ update $\mathbf{d}^{\mathrm{cML}}$ values
17:            **else if** $d < d_n^{\mathrm{cML}}$ **then**
18:              $d_n^{\mathrm{cML}} \leftarrow d$
19:            **end if**
20:          **end if**
21:       **end while**
22:       **if** $d < d^{\mathrm{ML}}$ **then**          ▷ detection
23:         $\hat{\mathbf{c}}^{\mathrm{GRAND}} \leftarrow \bar{\mathbf{c}}$; $d^{\mathrm{ML}} \leftarrow d$
24:       **end if**
25:       **if** $\bar{\mathbf{c}} \in C$ **or** $k = B$ **then**   ▷ code-word hit
26:         **if** $d < \bar{d}^{\mathrm{ML}}$ **then**       ▷ decoding
27:            $\bar{\mathbf{c}}^{\mathrm{GRAND}} \leftarrow \bar{\mathbf{c}}$; $\bar{d}^{\mathrm{ML}} \leftarrow d$
28:         **end if**
29:         $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}} \leftarrow \left(2\hat{\mathbf{c}}^{\mathrm{GRAND}} - 1\right) \odot \left(d^{\mathrm{ML}} - \mathbf{d}^{\mathrm{cML}}\right)$
30:         $t \leftarrow t + 1$
31:         **break**   ▷ go to next turbo iteration - line 4
32:       **end if**
33:      **end while**
34:  **end while**
35:  **return** $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}$, $\hat{\mathbf{c}}^{\mathrm{GRAND}}$, and $\bar{\mathbf{c}}^{\mathrm{GRAND}}$

---

in joint detection and decoding. We propose generating soft-output LLRs by populating a number of Euclidean distance computations equal to the number of noise guesses, as opposed to the exponential number of distance computations in (3). In particular, through guesswork, we aim at extracting an updated hard-detected vector, $\hat{\mathbf{c}}^{\mathrm{GRAND}}$, and a reliability metric for each bit in $\hat{\mathbf{c}}^{\mathrm{GRAND}}$, accumulated in the soft-decoding LLR vector, $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}$, all while recovering the decoded output, $\bar{\mathbf{c}}^{\mathrm{GRAND}}$.

Algorithm 2 illustrates turbo-GRAND in more detail. For detection, we keep track of an ML-detection distance metric, $d^{\mathrm{ML}}$; for decoding, we keep track of an ML-decoding distance metric, $\bar{d}^{\mathrm{ML}}$. After each turbo-GRAND iteration, the candidate vector corresponding to $d^{\mathrm{ML}}$ is the updated detected vector,

$\hat{\mathbf{c}}^{\mathrm{GRAND}}$, and the candidate vector corresponding to $\bar{d}^{\mathrm{ML}}$ is the updated decoded vector, $\bar{\mathbf{c}}^{\mathrm{GRAND}}$. Furthermore, for LLR computations, we accumulate a vector of counter-ML-detection distance metrics, $\mathbf{d}^{\mathrm{cML}} = [d_{1,1}^{\mathrm{cML}} \cdots d_{i,j}^{\mathrm{cML}} \cdots d_{M,q}^{\mathrm{cML}}]^T \in \mathbb{R}^N$, which tracks vectors closest to $\hat{\mathbf{c}}^{\mathrm{GRAND}}$, but with bit-flips at corresponding indices. Searching the entire lattice $\bar{\mathcal{X}}$ results in $d^{\mathrm{ML}} = \min_{\mathbf{x} \in \bar{\mathcal{X}}} \|\mathbf{y} - \mathbf{Hx}\|^2$; we can re-express (3) as

$$\lambda_{i,j}^{\mathrm{ML}} \approx \begin{cases} \frac{1}{\sigma^2} d^{\mathrm{ML}} - \frac{1}{\sigma^2} \min_{\mathbf{x} \in \bar{\mathcal{X}}^{i,j,0}} \|\mathbf{y} - \mathbf{Hx}\|^2 & \text{if } \hat{c}_{i,j}^{\mathrm{ML}} = 1 \\ \frac{1}{\sigma^2} \min_{\mathbf{x} \in \bar{\mathcal{X}}^{i,j,1}} \|\mathbf{y} - \mathbf{Hx}\|^2 - \frac{1}{\sigma^2} d^{\mathrm{ML}} & \text{if } \hat{c}_{i,j}^{\mathrm{ML}} = 0, \end{cases}$$
(5)

where $d_{i,j}^{\mathrm{cML}} = \min_{\mathbf{x} \in \bar{\mathcal{X}}^{i,j,0}} \|\mathbf{y} - \mathbf{Hx}\|^2$ if $\hat{c}_{i,j}^{\mathrm{ML}} = 1$ and $d_{i,j}^{\mathrm{cML}} = \min_{\mathbf{x} \in \bar{\mathcal{X}}^{i,j,1}} \|\mathbf{y} - \mathbf{Hx}\|^2$ if $\hat{c}_{i,j}^{\mathrm{ML}} = 0$ (note that colored noise scaling is embedded into $\mathbf{d}^{\mathrm{ML}}$ and $\mathbf{d}^{\mathrm{cML}}$ in Alg. 2). However, turbo-GRAND only searches a limited number of $\bar{\mathbf{x}}$ vectors that are extracted from the recovered $\bar{\mathbf{c}} = \hat{\mathbf{c}} \ominus \mathbf{w}$ words via guesswork, and that can be accumulated in a set $\mathcal{S}$ ($|\mathcal{S}| < T \times B \ll |\bar{\mathcal{X}}|$).

We first initialize $d^{\mathrm{ML}}$ to $\infty$ and $\mathbf{d}^{\mathrm{cML}}$ to saturation noise-scaled thresholds. Then, $d^{\mathrm{ML}}$ and $\mathbf{d}^{\mathrm{cML}}$ are updated iteratively, upon every new noise guess (up to $B$ guesses per iteration $t$ in GRAND with abandonment). For every guessed word, $\bar{\mathbf{c}}$, we re-generate a modulated vector, $\bar{\mathbf{x}} = \mathrm{mod}\,(\bar{\mathbf{c}})$, and add it to $\mathcal{S}$. Hence, for turbo-GRAND, $d^{\mathrm{ML}} = \min_{\bar{\mathbf{x}} \in \mathcal{S}} \|\mathbf{y} - \mathbf{H\bar{x}}\|^2$, and

$$\bar{\lambda}_{i,j}^{\mathrm{GRAND}} \approx \begin{cases} \frac{1}{\sigma^2} d^{\mathrm{ML}} - \frac{1}{\sigma^2} \min_{\bar{\mathbf{x}} \in \mathcal{S}^{i,j,0}} \|\mathbf{y} - \mathbf{H\bar{x}}\|^2 & \text{if } \hat{c}_{i,j}^{\mathrm{GRAND}} = 1 \\ \frac{1}{\sigma^2} \min_{\bar{\mathbf{x}} \in \mathcal{S}^{i,j,1}} \|\mathbf{y} - \mathbf{H\bar{x}}\|^2 - \frac{1}{\sigma^2} d^{\mathrm{ML}} & \text{if } \hat{c}_{i,j}^{\mathrm{GRAND}} = 0, \end{cases}$$
(6)

where $d_{i,j}^{\mathrm{cML}} = \min_{\bar{\mathbf{x}} \in \mathcal{S}^{i,j,0}} \|\mathbf{y} - \mathbf{H\bar{x}}\|^2$ if $\hat{c}_{i,j}^{\mathrm{GRAND}} = 1$ and $d_{i,j}^{\mathrm{cML}} = \min_{\bar{\mathbf{x}} \in \mathcal{S}^{i,j,1}} \|\mathbf{y} - \mathbf{H\bar{x}}\|^2$ if $\hat{c}_{i,j}^{\mathrm{GRAND}} = 0$. For the $i, j$ indices where $d_{i,j}^{\mathrm{cML}}$ cannot be computed over $\mathcal{S}$, owing to the absence of a corresponding $\bar{\mathbf{x}}$, the initial saturated value of $N/\sigma_i^2$ is retained. Note that we use a separate $\bar{d}^{\mathrm{ML}}$ for the decoded $\bar{\mathbf{c}}^{\mathrm{GRAND}}$ because not every $\bar{\mathbf{x}}$ in $\mathcal{S}$ corresponds to a code-word. Each turbo-GRAND iteration $t$ ends with a single code-word hit $\bar{\mathbf{c}}^{\mathrm{GRAND}}(t)$; the final decoded vector after $T$ iterations is

$$\bar{\mathbf{c}}^{\mathrm{GRAND}} = \underset{\bar{\mathbf{c}}^{\mathrm{GRAND}}(t);\ t \in \{1, \cdots, T\}}{\arg\min} \left\| \mathbf{y} - \mathbf{H}\,\mathrm{mod}\left(\bar{\mathbf{c}}^{\mathrm{GRAND}}(t)\right) \right\|^2.$$
(7)

Note that Algorithm 2 does not explicitly construct $\mathcal{S}$ and store it in memory for post-processing but instead updates the ML and counter-ML distance metrics on the fly.

The core of turbo-GRAND is a soft-input decoding mechanism which rank-orders candidate noise sequences according to $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}$. Let $\mathbf{W} \in \mathbb{F}_2^{2^N \times N}$ be a matrix containing in its rows all possible noise sequences, and let $\Pi : \mathbb{R}^{2^N} \to \{1, \cdots, 2^N\}^{2^N}$ be a sorting function (increasing order). Then, $\mathbf{s} = \Pi\,(\mathbf{W} \times |\mathbf{\Lambda}|)$ is a vector of sorted noise-sequence indices, and $\mathbf{w} = \Pi(\mathbf{s}(k))$ retrieves the $k$th likely noise sequence. However, populating noise sequences in a single matrix is not hardware-friendly nor computationally efficient. Alternatively, SGRAND [17] recursively constructs a max-heap for each combination of reliabilities in $\mathbf{\Lambda}$ to dynamically generate $\mathbf{w}$ vectors with increasing likelihoods. Also, ORBGRAND [18] builds a bit permutation map based on the decreasing rank

order of bit reliability to generate a pre-determined series of putative noise queries. Most probably, we have $\mathbf{s}(i) \leq \mathbf{s}(j)$ when $\Pr(\mathbf{w} = \mathbf{w}_i) \geq \Pr(\mathbf{w} = \mathbf{w}_j)$. In SGRAND [17], the latter is an "if and only if" condition. Thus, with either SGRAND or ORBGRAND, turbo-GRAND always queries the all-zeros noise sequence first and is biased to give higher priority to noise sequences of smaller Hamming weight, but not necessarily so, depending on the reliability information in $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}$. In the absence of soft information or further channel knowledge, noise query follows increasing Hamming weights for a probability of bit flip less than $1/2$. The latter is captured by an SGRAND core of turbo-GRAND, which reduces to hard-GRAND in the absence of soft information. Hence, SGRAND can be used in the first turbo-GRAND iteration when $\mathbf{\Lambda} = \mathbf{0}_N$. ORBGRAND requires some soft input information to match the performance of SGRAND, so it can be adopted when $\mathbf{\Lambda} = \mathbf{\Lambda}^{\mathrm{ZF}}$ in the first iteration. However, ORBGRAND entails more noise guesses on average, so it has the potential to generate richer LLRs.

Turbo-GRAND can be modified to support iterative disjoint detection and decoding, where $\bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}$ is fed all the way back to a separate detector. Towards this end, the detector should itself be SISO. For instance, the modified distance metric of a turbo MAP detector can be [25]

$$\varphi(\mathbf{x}) = -\frac{\|\mathbf{y} - \mathbf{Hx}\|^2}{\sigma^2} + \sum_{i=1}^{N} \sum_{j=1}^{q} (2\mathbf{c_x}(i, j) - 1)\, \bar{\mathbf{\Lambda}}^{\mathrm{GRAND}}(i, j),$$
(8)

where $\mathbf{c_x}$ is the bit representation of $\mathbf{x}$. The a posteriori LLRs can be calculated as

$$\lambda_{i,j}^{\mathrm{MAP}} = \left( \max_{\mathbf{x} \in \mathcal{X}^{i,j,1}} \varphi(\mathbf{x}) - \max_{\mathbf{x} \in \mathcal{X}^{i,j,0}} \varphi(\mathbf{x}) \right).$$
(9)

Disjoint detection and GRAND offers good modularity, where on every detection iteration, the noise can be filtered towards new signal subspaces [26], yielding more efficient GRAND.

## IV. PERFORMANCE AND COMPLEXITY EVALUATION

Following the system model of Sec. II, the block-error rate (BLER) performance of turbo-GRAND in decoding Bose–Chaudhuri–Hocquenghem (BCH) [ 127,113] codes is compared to reference GRAND/ORBGRAND schemes (OR-BGRAND exhibited superior performance and complexity tradeoffs compared to other code-specific decoders [18]). Assuming AWGN and normalized transmit power (SNR of $1/\sigma^2$), with the absence of input soft information, Fig. 2a illustrates that SGRAND-based turbo-GRAND generates soft information that matches ML-detection soft information. Turbo-GRAND can outperform ML-detection-based ORBGRAND because: 1) turbo iterations can introduce a list-decoding gain and 2) distance computations are over an entire code-word of many channel uses; feasible ML detection search routines only span a few channel uses. A similar relative performance is noted under Rayleigh fading in Fig. 2b and Fig. 2c, using binary phase-shift keying (BPSK) and 16-quadrature amplitude modulation (16-QAM) (Gray mapping),

(a) AWGN - SGRAND turbo-GRAND with nil soft input - BPSK.

(b) Rayleigh - SGRAND turbo-GRAND with nil soft input - BPSK.

(c) Rayleigh - SGRAND turbo-GRAND with nil soft input - 16QAM.

(d) Rayleigh - ORBGRAND turbo-GRAND (ZF soft input) - 90% CSI - BPSK.

Fig. 2. BLER performance evaluation of the proposed turbo-GRAND schemes with BCH [127,113] codes.

respectively. The gains in soft-GRAND schemes are larger under fading (more than 6 dB SNR gains at a BLER of $10^{-2}$), highlighting the importance of turbo-GRAND in the absence of soft information. We further note that the turbo-GRAND gains are captured in two iterations, beyond which diminishing returns are expected, as both SGRAND and ORBGRAND converge faster on every new iteration, $t$, guessing over an enhanced initial vector, $\hat{\mathbf{c}}^{\mathrm{GRAND}}(t)$. The achievable gains of turbo-GRAND can be improved by tuning a fixed guess budget, $B$. Alternatively, some sort of reactive taboo search [27] can be adopted, in which every iteration starts with a pseudo-random initial vector upon which noise is guessed.

Several communication system scenarios further highlight the turbo-GRAND gains, especially under symbol interference in spatial/path diversity schemes, where ML soft information is significantly better than ZF soft information. In such scenarios, starting from ZF soft information, ORBGRAND-based turbo-

GRAND can bridge the gap to the much more complex ML-detection-based soft-GRAND (this paper only covers uncorrelated point-to-point channels). In another scenario, under imperfect CSI, joint detection and decoding in turbo-GRAND outperforms conventional soft decoding. We assume 10% CSI error in Fig. 2d, where $\mathbf{H}_{\mathrm{err}} = 0.9\mathbf{H}+0.1\tilde{\mathbf{H}}$, and $\tilde{\mathbf{H}}$ has the same distribution as $\mathbf{H}$ but is independently and randomly generated. Starting from ZF soft information as input LLRs, ORBGRAND-based turbo-GRAND outperforms both ML-soft- and ZF-soft-ORBGRAND. Note, however, that our ML-detection implementation in these simulations only undergoes an exhaustive search over subsets of symbols in $\mathbf{x}$ of size four.

We next analyze the complexity in terms of floating-point operations in complex multiplication (CMT) and complex addition (CAD). We compare the additional processing in turbo-GRAND over soft-GRAND to the complexity of soft-output ZF and ML detection. The search complexities are dominated

TABLE I
SOFT-OUTPUT DATA DETECTION COMPLEXITY COMPARISON

| Detector | Complexity |
|---|---|
| ML | $\|\mathcal{X}\|^M \left((M^2+M)\text{CMT} + (M^2+M)\text{CAD}\right)$ |
| ZF | $\|\mathcal{X}\|\left(M\text{CMT} + M\text{CAD}\right)$ |
| turbo-GRAND | $T \times B \left((M^2+M)\text{CMT} + (M^2+M)\text{CAD}\right)$ |

by Euclidian distance computations and the complex matrix multiplications they entail. The complexity is exponential (in the symbol-vector length, $M$) with ML detection ((3) and (5)), polynomial with turbo-GRAND (6), and linear with ZF (4).

Table I illustrates the approximate worst-case complexity of generating soft information in one channel use. The average complexity of turbo-GRAND is much less because the guess budget, $B$, is not exploited on every iteration; with more iteration and higher SNR, a few or a single guess can recover a code-word. Turbo-GRAND is thus much less complex than conventional iterative list-based detection and decoding. Even with larger guess budgets, the recovered words on different iterations often overlap, and redundant computations can be saved. Furthermore, because noise-sequence guessing typically follows increasing Hamming weights, the vector Euclidean distance computations can reduce to simple symbol-based scalar distance computations/updates, resulting in an average linear turbo-GRAND complexity of $T \times B\left(M\text{CMT} + M\text{CAD}\right)$. Hence, for large modulation orders, a hardware-optimized turbo-GRAND can even prove to be less complex than soft-output ZF detection followed by soft-GRAND; much less complex than ML detection. Further simplifications can be made if the channel remains static over multiple uses.

## V. CONCLUSIONS

We proposed a mechanism for using GRAND to extract soft information in a joint detection and decoding framework. By leveraging access to complex received symbols, hard demapped bits, CSI, and possibly noise statistics, we generate LLRs by populating a shortlist of Euclidean distance computations. Such LLRs can be used in subsequent soft-decoding GRAND iterations, giving rise to turbo-GRAND. Compared to hard GRAND, a few iterations of turbo-GRAND introduce an excess of 6 dB SNR gains at a BLER of $10^{-2}$ (much higher gains at lower BLERs), under practical communication system scenarios of Rayleigh fading channels. Furthermore, turbo-GRAND can match and even outperform exhaustive ML-detection-based soft-GRAND at a much-reduced average linear complexity. This work can extend into a generic joint detection and decoding framework for future investigations.

## REFERENCES

[1] N. Rajatheva, I. Atzeni, E. Bjornson, A. Bourdoux, S. Buzzi, J.-B. Dore, S. Erkucuk, M. Fuentes, K. Guan *et al.*, "White paper on broadband connectivity in 6G," *arXiv preprint arXiv:2004.14247*, 2020.

[2] H. Sarieddeen, M.-S. Alouini, and T. Y. Al-Naffouri, "An overview of signal processing techniques for terahertz communications," *Proceedings of the IEEE*, vol. 109, no. 10, pp. 1628–1665, 2021.

[3] G. Durisi, T. Koch, and P. Popovski, "Toward massive, ultrareliable, and low-latency wireless communication with short packets," *Proceedings of the IEEE*, vol. 104, no. 9, pp. 1711–1726, 2016.

[4] D. Gazelle and J. Snyders, "Reliability-based code-search algorithms for maximum-likelihood decoding of block codes," *IEEE Trans. Inf. Theory*, vol. 43, no. 1, pp. 239–249, 1997.

[5] A. Valembois and M. Fossorier, "Box and match techniques applied to soft-decision decoding," *IEEE Trans. Inf. Theory*, vol. 50, no. 5, pp. 796–810, 2004.

[6] K. R. Duffy, J. Li, and M. Médard, "Guessing noise, not code-words," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2018, pp. 671–675.

[7] ——, "Capacity-achieving guessing random additive noise decoding," *IEEE Trans. Inf. Theory*, vol. 65, no. 7, pp. 4023–4040, 2019.

[8] M. M. Christiansen and K. R. Duffy, "Guesswork, large deviations, and shannon entropy," *IEEE Trans. Inf. Theory*, vol. 59, no. 2, pp. 796–802, 2013.

[9] A. Beirami, R. Calderbank, M. M. Christiansen, K. R. Duffy, and M. Médard, "A characterization of guesswork on swiftly tilting curves," *IEEE Trans. Inf. Theory*, vol. 65, no. 5, pp. 2850–2871, 2019.

[10] E. Arikan, "An inequality on guessing and its application to sequential decoding," *IEEE Trans. Inf. Theory*, vol. 42, no. 1, pp. 99–105, 1996.

[11] S. M. Abbas, T. Tonnellier, F. Ercan, and W. J. Gross, "High-throughput VLSI architecture for GRAND," in *Proc. IEEE Workshop on Signal Process. Syst.*, 2020, pp. 1–6.

[12] A. Riaz, V. Bansal, A. Solomon, W. An, Q. Liu, K. Galligan, K. R. Duffy, M. Medard, and R. T. Yazicigil, "Multi-code multi-rate universal maximum likelihood decoder using GRAND," in *IEEE 47th European Solid State Circuits Conf. (ESSCIRC)*, 2021, pp. 239–246.

[13] T. Kaneko, T. Nishijima, and S. Hirasawa, "An improvement of soft-decision maximum-likelihood decoding algorithm using hard-decision bounded-distance decoding," *IEEE Trans. Inf. Theory*, vol. 43, no. 4, pp. 1314–1319, 1997.

[14] M. Fossorier and S. Lin, "Soft-decision decoding of linear block codes based on ordered statistics," *IEEE Trans. Inf. Theory*, vol. 41, no. 5, pp. 1379–1396, 1995.

[15] K. R. Duffy and M. Médard, "Guessing random additive noise decoding with soft detection symbol reliability information - SGRAND," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2019, pp. 480–484.

[16] S. M. Abbas, M. Jalaleddine, and W. J. Gross, "GRAND for Rayleigh fading channels," *arXiv preprint arXiv:2205.00030*, 2022.

[17] A. Solomon, K. R. Duffy, and M. Médard, "Soft maximum likelihood decoding using GRAND," in *Proc. IEEE Int. Conf. Commun. (ICC)*, 2020, pp. 1–6.

[18] K. R. Duffy, "Ordered reliability bits guessing random additive noise decoding," in *Proc. IEEE Int. Conf. Acoustics, Speech, and Signal Process. (ICASSP)*, 2021, pp. 8268–8272.

[19] C. Berrou and A. Glavieux, "Near optimum error correcting coding and decoding: turbo-codes," *IEEE Trans. Commun.*, vol. 44, no. 10, pp. 1261–1271, 1996.

[20] A. Tomasoni, M. Siti, M. Ferrari, and S. Bellini, "Low complexity, quasi-optimal MIMO detectors for iterative receivers," *IEEE Trans. Commun.*, vol. 9, no. 10, pp. 3166–3177, 2010.

[21] W. An, M. Médard, and K. R. Duffy, "Keep the bursts and ditch the interleavers," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2020, pp. 1–6.

[22] M. Ivanov, C. Häger, F. Brännström, A. G. i Amat, A. Alvarado, and E. Agrell, "On the information loss of the max-log approximation in BICM systems," *IEEE Trans. Inf. Theory*, vol. 62, no. 6, pp. 3011–3025, 2016.

[23] H. Sarieddeen, M. M. Mansour, and A. Chehab, "Large MIMO detection schemes based on channel puncturing: Performance and complexity analysis," *IEEE Trans. Commun.*, no. 99, pp. 1–1, Dec. 2017.

[24] C. Studer, S. Fateh, and D. Seethaler, "ASIC implementation of soft-input soft-output MIMO detection using MMSE parallel interference cancellation," *IEEE J. Solid-State Circuits*, vol. 46, no. 7, pp. 1754–1765, 2011.

[25] A. Tomasoni, M. Siti, M. Ferrari, and S. Bellini, "Hardware oriented, quasi-optimal detectors for iterative and non-iterative MIMO receivers," *EURASIP Journal on Wireless Communications and Networking*, vol. 2012, no. 1, p. 62, 2012.

[26] H. Sarieddeen, M. M. Mansour, and A. Chehab, "Channel-punctured large MIMO detection," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2018, pp. 2147–2151.

[27] T. Datta, N. Srinidhi, A. Chockalingam, and B. S. Rajan, "Random-restart reactive tabu search algorithm for detection in large-MIMO systems," *IEEE Commun. Lett.*, vol. 14, no. 12, pp. 1107–1109, Dec. 2010.