# Analysis of Saturated Belief Propagation Decoding of Low-Density Parity-Check Codes

Shrinivas Kudekar, Tom Richardson and Aravind Iyengar
Qualcomm, New Jersey, USA
Email: {skudekar,tomr,ariyengar}@qti.qualcomm.com

*Abstract*—We consider the effect of log-likelihood ratio saturation on belief propagation decoder low-density parity-check codes. Saturation is commonly done in practice and is known to have a significant effect on error floor performance. Our focus is on threshold analysis and stability of density evolution.

We analyze the decoder for standard low-density parity-check code ensembles and show that belief propagation decoding generally degrades gracefully with saturation. Stability of density evolution is, on the other hand, rather strongly effected by saturation and the asymptotic qualitative effect of saturation is similar to reduction by one of variable node degree.

We also show under what conditions the block threshold for the saturated belief propagation corresponds with the bit threshold.

## I. INTRODUCTION

Standard belief propagation (BP) decoding of binary low-density parity-check (LDPC) codes involves passing messages typically representing log-likelihood ratios (LLRs) which can take any value in $\overline{\mathbb{R}} \triangleq \mathbb{R} \cup \{\pm\infty\}$ [1]. The asymptotic analysis developed for BP decoding of LDPC codes inherently assumes that the messages have unbounded magnitude. In practice, however, decoders typically use uniformly quantized and bound LLRs. Density evolution can be applied directly to such decoders but analysis is often difficult and there are few general results. Hence, it is of interest to understand the effect of saturation of LLR magnitudes as a perturbation of full belief propagation. We call such a saturated decoder as a *saturating* belief propagation decoder (SatBP). Note that the decoder is strictly speaking not a BP decoder, but we adhere to the BP nomenclature as we view SatBP as a perturbation of BP.

In the design of capacity-achieving codes it is helpful to understand how practical decoder concessions, like saturation, affect performance. For this purpose, we will analyze the SatBP decoder in the asymptotic limit of the blocklength going to infinity. In particular, if LLRs are saturated at magnitude K then how much degradation from the BP threshold should be expected. Naturally, one expects that as $K \to +\infty$, that one can reliably transmit arbitrarily close to the BP threshold [1]. We will see that this is not entirely correct and that, in particular, saturation can undermine the stability of the perfect decoding fixed point if, for example, the fraction of degree two variable nodes in an irregular ensemble is non-zero. Our analysis shows that when the minimum variable node degree is at least three then there exists a large but finite saturation value K such that the SatBP decoder can achieve arbitrarily

small bit error rate whenever the full BP decoder can achieve arbitrarily small bit error rate. Furthermore, a more careful stability analysis shows that in fact one can achieve reliability in terms of the block error rate.

### A. Related Work

The papers [2]–[5] consider the effect of saturation on error floor performance. It is observed in these works that saturation can limit the ability of decoding to escape trapping set behavior, thereby worsening error floor performance. In [6], [7] some decoder variations are given that help reduce error floors. Here we see an explicit effort to ameliorate the effect of saturation. A related but distinct direction was taken in [8]. There the authors made modifications to discrete node update rules so as to reduce error floor failure events. They fine tune finite state message update rules to optimize performance on a particular graph structure. There have been other works that examine the effects of practical concessions. In [9] the authors consider the effect of quantization in LDPC coded flash memories. In [10] and[11] the effects of saturation and quantization are modeled as noise terms. Finally, in [12] an analysis is done to evaluate the effect on capacity on quantization of channel outputs. Although we take a different approach in this paper by focusing on asymptotic behavior, the fundamental conclusion is similar to the error floor results in [2]–[5]: saturation can dramatically effect the stability of the decoder.

The paper is organized as follows. In the next section we will briefly review the standard asymptotic analysis of the BP decoder using density evolution (DE). Then in sections III and IV we will introduce the SatBP decoder and perform perturbation analysis on the SatBP decoder using the Wasserstein metric [13]. In section V we will use stability analysis to examine block thresholds for SatBP. We will see that in many cases the block threshold will correspond with the bit threshold, but the conditions required are more stringent than in the non-saturated decoder case.

## II. BP DECODING, DENSITY EVOLUTION AND THE WASSERSTEIN DISTANCE

In this section we briefly review the BP decoder and the DE analysis [14] in the case of transmission over a general BMS channel using standard LDPC code ensemble. Most of the material presented here can be found in [1].

We assume transmission over a BMS channel. Let $X(= \pm 1)$ denote the input and let $Y$ be the output. Further, let $p(Y = y \mid X = x)$ denote the *transition probability* describing the channel. We generally characterize a BMS channel by its so-called $L$-distribution, c. More precisely, c is the distribution of

$$\ln \frac{p(Y \mid X = +1)}{p(Y \mid X = -1)}$$

conditioned that $X = +1$. Generally, we may assume that

$$Y = \ln \frac{p(Y \mid X = +1)}{p(Y \mid X = -1)}.$$

The symmetry of the channel is $p(Y = y \mid X = x) = p(Y = -y \mid X = -x)$ and the resulting densities c are symmetric, [1], which means $e^{-\frac{1}{2}x}\mathsf{c}(x)$ is an even function of $x$.

Given $Z$ distributed according to c, we write $\mathfrak{c}$ to denote the distribution of $\tanh(Z/2)$, and $|\mathfrak{c}|$ to denote the distribution of $|\tanh(Z/2)|$. We refer to these as the D and $|D|$ distributions respectively. We use $|\mathfrak{C}|$ to denote the corresponding cumulative $|D|$ distribution, see [1, Section 4.1.4]. Under symmetry, the distribution of $|Z|$ determines the distribution of $Z$.

For threshold analysis of LDPC ensembles we typically consider a parameterized *family* of channels. We write $\{\mathrm{BMS}(\sigma)\}$ to denote the family parameterized by the scalar $\sigma$. Often it will be more convenient to denote this family by $\{\mathsf{c}_\sigma\}$, i.e., to use the family of $L$-densities which characterize the channel family. One natural candidate for the parameter $\sigma$ is the entropy of the channel denoted by h. Thus, we also consider the characterization of the family given by $\mathrm{BMS}(\mathsf{h})$.

*A. Degradation, Symmetric Densities and Functionals of Densities*

Let $p_{Z \mid X}(z \mid x)$ denote the transition probability associated to a BMS channel $\mathsf{c}'$ and let $p_{Y \mid X}(y \mid x)$ denote the transition probability of another BMS channel c. We then say that $\mathsf{c}'$ is *degraded* with respect to c if there exists a channel $p_{Z \mid Y}(z \mid y)$ so that

$$p_{Z \mid X}(z \mid x) = \sum_y p_{Y \mid X}(y \mid x) p_{Z \mid Y}(z \mid y).$$

We will use the notation $\mathsf{c} \prec \mathsf{c}'$ to denote that $\mathsf{c}'$ is degraded with respect to c (as a mnemonic think of c as the erasure probability of a BEC and replace $\prec$ with $\leq$).

A useful characterization of degradation, see [15], [1, Theorem 4.74], is that $\mathsf{c} \prec \mathsf{c}'$ is equivalent to

$$\int_0^1 f(x)|\mathfrak{c}|(x)\,\mathrm{d}x \leq \int_0^1 f(x)|\mathfrak{c}'|(x)\,\mathrm{d}x \tag{1}$$

for all $f(x)$ that are non-increasing and concave on $[0, 1]$. In particular, this characterization implies that $F(\mathsf{a}) \leq F(\mathsf{b})$ for $\mathsf{a} \prec \mathsf{b}$ if $F(\cdot)$ is either the Battacharyya or the entropy functional. This is true since both are linear functionals of the distributions and their respective kernels in the $|D|$-domain are decreasing and concave, see [1]. An alternative characterization [1] of degradation in terms of the cumulative distribution functions $|\mathfrak{C}|(x)$ and $|\mathfrak{C}'|(x)$ is that for all $z \in [0, 1]$,

$$\int_z^1 |\mathfrak{C}|(x)\mathrm{d}x \leq \int_z^1 |\mathfrak{C}'|(x)\,\mathrm{d}x. \tag{2}$$

A BMS channel family $\{\mathrm{BMS}(\mathsf{h})\}_{\underline{\mathsf{h}}}^{\overline{\mathsf{h}}}$ is said to be *ordered* (by degradation) if $\mathsf{h}_1 \leq \mathsf{h}_2$ implies $\mathsf{c}_{\mathsf{h}_1} \prec \mathsf{c}_{\mathsf{h}_2}$. (The reverse order, $\mathsf{h}_1 \geq \mathsf{h}_2$, is also allowed but we generally stick to the stated convention.)

*Definition 1 (Symmetric Densities):* Let $A$ denote an $L$-distribution in $\overline{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$. Then $A$ is *symmetric* if it satisfies the following condition for every bounded, continuous function $f : \mathbb{R} \mapsto \mathbb{R}$,

$$\int f(x)\mathrm{d}A(x) = \int e^{-x} f(-x)\mathrm{d}A(x). \tag{3}$$

We say that an $L$-density a is *symmetric* if $\mathsf{a}(-y) = \mathsf{a}(y)e^{-y}$. We recall that all densities which stem from BMS channels are symmetric, see [1, Sections 4.1.4, 4.1.8 and 4.1.9]. ∎

Functionals of densities often used in analysis are the Battacharyya, the entropy, and the error probability functional. For a density a, these are denoted by $\mathfrak{B}(\mathsf{a})$, $\mathsf{H}(\mathsf{a})$, and $\mathfrak{E}(\mathsf{a})$, respectively and are defined by

$$\mathfrak{B}(\mathsf{a}) = \mathbb{E}(e^{-y/2}), \quad \mathsf{H}(\mathsf{a}) = \mathbb{E}(\log_2(1 + e^{-y}))$$
$$\mathfrak{E}(\mathsf{a}) = \mathbb{P}\{y < 0\} + \frac{1}{2}\mathbb{P}\{y = 0\}.$$

where $y$ is distributed according to a. Note that these definitions are valid even if a is not symmetric, although they lose some of their original meaning. We will apply these definitions to saturated densities that are not necessarily symmetric. It is not hard to see that $\mathfrak{E}(\mathsf{a}) \leq \mathfrak{B}(\mathsf{a})$ for any density a not necessarily symmetric. Hence in the paper the main functional of interest is the Battarcharyya parameter.

*B. BP Decoder, DE analysis and the Wasserstein metric*

The definition of the standard BP decoder can be found in [1]. The asymptotic performance of the BP decoder is given by the DE technique [1], [14]. Throughout the paper we will consider standard LDPC code ensembles as specified by their degree distributions [1]. The analysis can be applied to more sophisticated structures, but we restrict to this case for simplicity of presentation. Thus we let $\lambda(\cdot)$ and $\rho(\cdot)$ represent the variable node and check node degree profile respectively. The ensemble is then denoted by $(\lambda, \rho)$.

*Definition 2 (DE for BP Decoder cf. [1]):* For $\ell \geq 1$, the DE equation for a $(\lambda, \rho)$ ensemble is given by

$$\mathsf{x}_\ell = \mathsf{c} \circledast \lambda(\rho(\mathsf{x}_{\ell-1})).$$

Here, c is the $L$-density of the BMS channel over which transmission takes place and $\mathsf{x}_\ell$ is the density emitted by variable nodes in the $\ell$-th round of density evolution. Initially we have $\mathsf{x}_0 = \Delta_0$, the delta function at 0. The operators $\circledast$ and $\boxast$ correspond to the convolution of densities at variable and check nodes, respectively, see [1, Section 4.1.4]. The notation $\rho(\mathsf{x}_{\ell-1})$ represents the weighted check node convolution of the density $\mathsf{x}_{\ell-1}$. E.g., if $\rho(x) = x^{d_r-1}$, then $\rho(\mathsf{x}_{\ell-1}) = \mathsf{x}_{\ell-1}^{\boxast d_r-1}$. ∎

*Discussion*: For $(d_l, d_r)-$regular codes, the DE equation is given by $\mathsf{x}_\ell = \mathsf{c} \circledast (\mathsf{x}_{\ell-1}^{\boxast d_r-1})^{\circledast d_l-1}$. The DE analysis is simplified when we consider the class of symmetric message-passing

decoders. The definition of symmetric message-passing decoders can be found in [1]. Note that this definition of symmetry pertains to the actual messages in the decoder and not to the densities which appear in the DE analysis. We will see later that the saturated decoder is a symmetric message-passing decoder and hence its DE analysis is simplified by restricting to the use of the all zero (actually we use $+1$ for 'zero') codeword.

*Definition 3 (BP Threshold):* Consider an ordered and complete channel family $\{c_h\}$. Let $x_\ell(h)$ denote the distribution in the $\ell$-th round of DE when the channel is $c_h$. Then the *BP threshold* of the $(\lambda, \rho)$ ensemble is typically defined as

$$h^{\text{BP}}(\lambda, \rho, \{c_h\}) = \sup\{h : x_\ell(h) \overset{\ell \to \infty}{\to} \Delta_{+\infty}\}.$$

Here $\Delta_{+\infty}$ is the delta function at infinity representing the perfect decoding density. An equivalent definition is

$$h^{\text{BP}}(\lambda, \rho, \{c_h\}) = \sup\{h : \mathfrak{E}(x_\ell(h)) \overset{\ell \to \infty}{\to} 0\}.$$

The later form is more convenient for our purposes and it is the one we shall adopt. ∎

We will also say that for a given channel $c$, the BP decoder is *successful* if and only if $\mathfrak{E}(x_\ell(h)) \overset{\ell \to \infty}{\to} 0$ or $\mathfrak{B}(x_\ell(h)) \overset{\ell \to \infty}{\to} 0$. In other words, for any given $\epsilon > 0$, there exists $\ell$ such that $\mathfrak{B}(x_\ell(h)) < \epsilon$.

In the sequel we will use the Wasserstein metric to measure distance between distributions. We recall the definition of the Wasserstein metric below. For more properties of the Wasserstein metric see [16].

*Definition 4 (Wasserstein Metric – [17, Chapter 6]):* Let $|\mathfrak{a}|$ and $|\mathfrak{b}|$ denote two $|D|$-distributions. The Wasserstein metric, denoted by $d(|\mathfrak{a}|, |\mathfrak{b}|)$, is defined as

$$d(|\mathfrak{a}|, |\mathfrak{b}|) = \sup_{f(x) \in \text{Lip}(1)[0,1]} \left| \int_0^1 f(x)(|\mathfrak{a}|(x) - |\mathfrak{b}|(x)) \, dx \right|, \quad (4)$$

where $\text{Lip}(1)[0,1]$ denotes the class of Lipschitz continuous functions on $[0,1]$ with Lipschitz constant 1.

In [18] it is shown that the Wasserstein distance is equivalent to the $L_1$ norm of the difference between the $|D|$-distributions. ∎

## III. SATURATED BELIEF PROPAGATION DECODING

In this section we introduce the saturated BP decoder. More precisely, we consider decoding with BP update rules at the nodes but the outgoing messages are restricted to the domain $[-K, K]$ for some $K > 0$ by saturation.

### A. Saturated Decoder

*Definition 5 (Saturation):* We define the *saturation* operation at $\pm K$ for some $K \in \mathbb{R}^+$, denoted $\lfloor \cdot \rfloor_K$, by

$$\lfloor x \rfloor_K = \min(K, |x|) \cdot \text{sgn}(x), \quad (5)$$

where

$$\text{sgn}(x) = \begin{cases} -1, & x < 0 \\ 1, & x \geq 0 \end{cases}.$$

*Definition 6 (Saturated BP Decoder):* Consider the standard $(d_l, d_r)$-regular ensemble. The saturated BP decoder is defined by the following rules. Let $\phi^{(\ell)}(\mu_1, \ldots, \mu_{d_r-1})$ and $\psi^{(\ell)}(\mu_1, \ldots, \mu_{d_l-1})$ denote the outgoing message from the check node and the variable node side respectively. Abusing the notation above, $\mu_1, \ldots, \mu$. denotes the incoming messages on both the check node and the variable node side. Then,

$$\phi^{(\ell)}(\mu_1, \ldots, \mu_{d_r-1}) = \left\lfloor 2\tanh^{-1}\left( \prod_{i=1}^{d_r-1} \tanh(\mu_i/2) \right) \right\rfloor_K,$$

$$\psi^{(\ell)}(\mu_1, \ldots, \mu_{d_l-1}) = \left\lfloor \mu_0 + \sum_{i=1}^{d_l-1} \mu_i \right\rfloor_K,$$

where $\mu_0$ is the message coming from the channel. Also, we set $\phi^{(0)}(\mu_1, \ldots, \mu_{d_r-1}) = 0$.

*Lemma 7 (SatBP Decoder is symmetric):* The SatBP decoder given in Definition 6 is a symmetric message-passing decoder.

*Proof:* From Definition 4.83 in [1] it is not hard to see that variable-node symmetry is satisfied for $\ell = 0$. In general, variable node symmetry is the following condition (for $\ell \geq 1$) on the message update function

$$\psi^{(\ell)}(-\mu_0, -\mu_1, \ldots, -\mu_{d_l-1}) = -\psi^{(\ell)}(\mu_0, \mu_1, \ldots, \mu_{d_l-1}).$$

Since $\lfloor x \rfloor_K = -\lfloor -x \rfloor_K$ we see that variable node symmetry is preserved by saturation. Let $b_1 \in \{\pm 1\}, \ldots, b_{d_r-1} \in \{\pm 1\}$, then by Definition 4.83 in [1], for the check node symmetry we have

$$\phi^{(\ell)}(b_1\mu_1, \ldots, b_{d_r-1}\mu_{d_r-1})$$

$$= \min\left( 2\tanh^{-1}\left( \prod_{i=1}^{d_r-1} \tanh(|\mu_i|/2) \right), K \right) \text{sgn}\left( \prod_{i=1}^{d_r-1} b_i\mu_i \right)$$

$$= \min\left( 2\tanh^{-1}\left( \prod_{i=1}^{d_r-1} \tanh(|\mu_i|/2) \right), K \right) \text{sgn}\left( \prod_{i=1}^{d_r-1} \mu_i \right) \prod_{i=1}^{d_r-1} b_i$$

$$= \phi^{(\ell)}(\mu_1, \ldots, \mu_{d_r-1}) \left( \prod_{i=1}^{d_r-1} b_i \right).$$

and we see again that symmetry is preserved by saturation. ∎

*Discussion:* The symmetry of the message-passing decoder together with symmetry of the channel allows us to use the all-zero codeword assumption. This along with the concentration results (see Theorem 4.94 in [1]) allows to write down the density evolution of the SatBP decoder in the usual way. Note that if messages entering a check node are saturated in magnitude at K then outgoing messages are automatically saturated at K. This holds not just for BP but for many message passing algorithms such as the min-sum algorithm. Our analysis has two parts: bounding the effect of saturation over finitely many iterations and stability analysis. For the bounding analysis we focus on BP although the technique can be easily extended to other decoders. In the stability analysis we explicitly relax the assumptions to cover a variety of check node updates.

Given $X \sim \mathfrak{a}$, let $\lfloor \mathfrak{a} \rfloor_K$ denote the distribution of $\lfloor X \rfloor_K$. Note that the saturation operation can be viewed as a channel

taking $X$ to $\lfloor X \rfloor_K$. We have immediately

$$\mathsf{a} \prec \lfloor \mathsf{a} \rfloor_K .$$

In general $\lfloor \mathsf{a} \rfloor_K$ will not be symmetric even if $\mathsf{a}$ is symmetric since we will not typically have $\lfloor \mathsf{a} \rfloor_K(-K) = e^{-K}\lfloor \mathsf{a} \rfloor_K(K)$. If $\mathsf{a}$ is symmetric then we will have

$$\lfloor \mathsf{a} \rfloor_K(-K) \le e^{-K}\lfloor \mathsf{a} \rfloor_K(K). \tag{6}$$

Although using lemma 7 one can write down the DE recursion for the SatBP decoder, we know that in general the densities will not be symmetric. Two of the most useful properties of DE for BP are that it preserves both symmetry of densities and ordering by degradation. These properties are sacrificed by saturation, but can be recovered with a slight variation. There are two alternatives for this. One is to place the saturated probability mass at $\pm z$ instead at $\pm K$ where $z$ is chosen according to the actual LLR conditioned on magnitude K. The second alternative is to slightly degrade the density by moving some probability mass from K to $-K$. This can be interpreted operationally as flipping the sign of a message with magnitude K with some probability $\gamma$. The flipping rate $\gamma$ is chosen so that the resulting probability that the sign of the message is incorrect is $e^{-K}/(1 + e^{-K})$. In general $\gamma$ is upper bounded by this value and for large K this is a small perturbation. Of the two approaches the second is inferior in that it degrades the channel more than the first. On the other hand, the second approach preserves ordering by degradation while the first does not. We shall adopt the second approach.

Let us introduce the notation $D(p, z)$ to denote the density

$$D(p, z) = p\Delta_{-z} + (1 - p)\Delta_z .$$

Here $\Delta_z$ ($\Delta_{-z}$) is the delta function at $z$ ($-z$). We will sometimes denote $p\Delta_{-z}$ as $D_-(p, z)$ and $(1 - p)\Delta_z$ as $D_+(p, z)$. When $(p, z)$ is clear from context we may drop it from the notation. Using this notation we have for symmetric $\mathsf{a}$,

$$\lfloor \mathsf{a} \rfloor_K = \gamma D(q, z)(x) + \mathsf{a}(x)\mathbb{1}_{\{|x| < K\}} \tag{7}$$

where $\gamma = \mathbb{P}_\mathsf{a}\{|x| \ge K\}$ and $\gamma q = \mathbb{P}_\mathsf{a}\{x \le -K\}$.

*Lemma 8 (Symmetric SatBP):* Given a symmetric density $\mathsf{a}$ we define

$$\lfloor \mathsf{a} \rfloor_{K_{sym}} = \gamma D(p, z)(x) + \mathsf{a}(x)\mathbb{1}_{\{|x| < K\}}$$

where $p = e^{-K}/(1 + e^{-K})$ and $\gamma = \mathbb{P}_\mathsf{a}\{|x| \ge K\}$. Then,

(i) $\lfloor \mathsf{a} \rfloor_{K_{sym}}$ is a symmetric $L$-density.
(ii) $\lfloor \mathsf{a} \rfloor_K \prec \lfloor \mathsf{a} \rfloor_{K_{sym}}$.

*Proof:* Part (i) is immediate. To prove part (ii) we note that comparing with the non-symmetrized case we see that $p \ge q$. Thus, $\lfloor \mathsf{a} \rfloor_{K_{sym}}$ can be realized by taking messages with distribution $\lfloor \mathsf{a} \rfloor_K$ and flipping the sign of a message with magnitude K by a quantity $\lambda$ with $\lambda$ determined by

$$p = \frac{e^{-K}}{1 + e^{-K}} = \lambda(1 - q) + (1 - \lambda)q .$$

∎

As a consequence of Lemma 8, we will term the operation used to obtain $\lfloor \mathsf{a} \rfloor_{K_{sym}}$ from $\mathsf{a}$ as *symmetric-saturation*.

We summarize all the claims above in the following.

*Corollary 9 (Degradation Order):* For symmetric $\mathsf{a}$ we have

$$\mathsf{a} \prec \lfloor \mathsf{a} \rfloor_K \prec \lfloor \mathsf{a} \rfloor_{K_{sym}}.$$

It is fairly intuitive that as K becomes larger, the density $\lfloor \mathsf{a} \rfloor_{K_{sym}}$ should become close to the density $\mathsf{a}$. This is the content of the next lemma which uses the Wasserstein distance between distributions.

*Lemma 10:* Let $\mathsf{a}$ be a symmetric $L$-density. Then,

$$d(\mathsf{a}, \lfloor \mathsf{a} \rfloor_{K_{sym}}) \le 1 - \tanh(K/2),$$

where $d(\cdot, \cdot)$ is the Wasserstein distance defined previously.

*Proof:* For any $0 \le z < K$ we have $\mathbb{P}_\mathsf{a}\{x \le z\} = \mathbb{P}_{\lfloor \mathsf{a} \rfloor_K}\{x \le z\} = \mathbb{P}_{\lfloor \mathsf{a} \rfloor_{K_{sym}}}\{x \le z\}$ and for any $z \ge K$ we have $1 = \mathbb{P}_{\lfloor \mathsf{a} \rfloor_K}\{x \le z\} = \mathbb{P}_{\lfloor \mathsf{a} \rfloor_{K_{sym}}}\{x \le z\}$. Since $\tanh(x/2)$ is increasing and $\tanh(-x/2) = -\tanh(x/2)$ we have

$$|\lfloor \mathfrak{A} \rfloor_{K_{sym}}|(z) = \mathbb{1}_{\{z < \tanh(K/2)\}}|\mathfrak{A}|(z) + \mathbb{1}_{\{z \ge \tanh(K/2)\}} .$$

By [18], we have that the Wasserstein distance is equivalent to the $L_1$ norm of the difference between the $|D|$-distributions. Clearly, the distance is bounded by $1 - \tanh(K/2)$. ∎

Let $T(\cdot)$ denote a DE iteration for the full BP decoder, i.e.,

$$T(\mathsf{c}, \mathsf{x}) = \mathsf{c} \circledast \lambda(\rho(\mathsf{x})).$$

*Definition 11 (DE for Sym. and Non-Sym. Saturation):* Consider a BMS channel with $L$-density $\mathsf{c}$. Let $\Delta_0$ denote the perfectly noisy channel. Let $\mathsf{x}^{(0)} = \Delta_0$. Then the DE for symmetric SatBP decoder is defined as,

$$\mathsf{x}^{(\ell)} = \lfloor \mathsf{c} \circledast \lambda(\rho(\mathsf{x}^{(\ell-1)})) \rfloor_{K_{sym}}.$$

The DE for non-symmetric SatBP decoder is defined as,

$$\mathsf{x}^{(\ell)} = \lfloor \mathsf{c} \circledast \lambda(\rho(\mathsf{x}^{(\ell-1)})) \rfloor_K.$$

Finally, we use the notation $S_{K_{sym}}(\mathsf{c}, \mathsf{x}) = \lfloor T(\mathsf{c}, \mathsf{x}) \rfloor_{K_{sym}}$ and $S_K(\mathsf{c}, \mathsf{x}) = \lfloor T(\mathsf{c}, \mathsf{x}) \rfloor_K$. ∎

Now imagine that we run both the full DE and symmetric saturated DE starting with the density $\Delta_0$. In the next lemma we show that at every iteration the order of degradation between the full DE and symmetric saturated DE is preserved. We will use the notation $T^{(\ell)}(\mathsf{c}, \Delta_0)$ to denote the $\ell$ iteration of the full DE. More precisely, $T^{(\ell)}(\mathsf{c}, \Delta_0) = T(\mathsf{c}, T^{(\ell-1)}(\mathsf{c}, \Delta_0))$. As a shorthand, we will use $T^{(\ell)}(\mathsf{c}, \Delta_0) = T(T^{(\ell-1)}(\mathsf{c}, \Delta_0))$. We similarly define $S^{(\ell)}_{K_{sym}}(\mathsf{c}, \Delta_0)$ and $S^{(\ell)}_K(\mathsf{c}, \Delta_0)$.

*Lemma 12 (Degradation Order under DE):* For any $\ell \ge 0$ we have

$$T^{(\ell)}(\mathsf{c}, \Delta_0) \prec S^{(\ell)}_{K_{sym}}(\mathsf{c}, \Delta_0).$$

*Proof:* Let $\mathsf{x}^{(\ell)}$ denote the DE for usual BP decoder and $\mathsf{z}^{(\ell)}$ denote the DE for the symmetric saturation operation. Since $\mathsf{x}^{(0)} = \mathsf{z}^{(0)} = \Delta_0$, we have that $\mathsf{x}^{(1)} = \mathsf{c}$ and $\mathsf{z}^{(1)} = \lfloor \mathsf{c} \rfloor_{K_{sym}}$. From corollary 9 we get that $\mathsf{x}^{(1)} \prec \mathsf{z}^{(1)}$. Now, since DE preserves the order of degradation, we get

$$\mathsf{x}^{(2)} = T(\mathsf{c}, \mathsf{x}^{(1)}) \prec T(\mathsf{c}, \mathsf{z}^{(1)}) \overset{\text{Lem. 9}}{\prec} \lfloor T(\mathsf{c}, \mathsf{z}^{(1)}) \rfloor_{K_{sym}} = \mathsf{z}^{(2)}.$$

Continuing, for all $\ell$ we get $T^{(\ell)}(\mathsf{c}, \Delta_0) \prec S^{(\ell)}_{K_{sym}}(\mathsf{c}, \Delta_0)$. ∎

We now estimate the distance between the densities appearing in the DE of standard BP and the DE of the symmetric saturation operation. For this we again use the Wasserstein distance defined previously (for symmetric densities).

*Lemma 13 (Distance Between Symmetric SatBP and BP):* Consider $\ell$ iterations of the DE for the standard BP and the symmetric saturation operation. Then

$$d(T^{(\ell)}(\mathsf{c},\Delta_0), S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0)) \leq 2e^{-\mathrm{K}+\ell\cdot\ln(2(d_l-1)(d_r-1))}.$$

*Proof:* Let $T(\cdot)$ and $S_{\mathrm{K_{sym}}}(\cdot)$ be defined as in lemma 12 and consider the Wasserstein distance between them. We have,

$$d(T^{(\ell)}(\mathsf{c},\Delta_0), S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0))$$
$$= d(T(T^{(\ell-1)}(\mathsf{c},\Delta_0)), S_{\mathrm{K_{sym}}}(S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)))$$
$$\overset{\text{Trian. ineq.}}{\leq} d(T(T^{(\ell-1)}(\mathsf{c},\Delta_0)), T(S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)))$$
$$+ d(T(S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)), S_{\mathrm{K_{sym}}}(S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)))$$
$$\overset{\text{(viii), Lem. 13 in [18]}}{\leq} \alpha_\ell d(T^{(\ell-1)}(\mathsf{c},\Delta_0), S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0))$$
$$+ d(T(S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)), S_{\mathrm{K_{sym}}}(S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)))$$
$$\overset{(a)}{\leq} \alpha_\ell d(T^{(\ell-1)}(\mathsf{c},\Delta_0), S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)) + \left(1 - \tanh\left(\frac{\mathrm{K}}{2}\right)\right),$$

where

$$\alpha_\ell = 2(d_l-1)\sum_{j=1}^{d_r-1}(1-\mathfrak{B}^2(\mathsf{a}))^{\frac{d_r-1-j}{2}}(1-\mathfrak{B}^2(\mathsf{b}))^{\frac{j-1}{2}},$$
$$\leq 2(d_l-1)(d_r-1)$$

where $\mathsf{a} = T^{(\ell-1)}(\mathsf{c},\Delta_0)$ and $\mathsf{b} = S_{\mathrm{K_{sym}}}^{(\ell-1)}(\mathsf{c},\Delta_0)$ and $d_l$ and $d_r$ correspond to the average variable node and check node degrees. Also, the inequality $(a)$ is obtained by using lemma 10.

Continuing with the above inequality we get,

$$d(T^{(\ell)}(\mathsf{c},\Delta_0), S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0))$$
$$\leq (1-\tanh\left(\frac{\mathrm{K}}{2}\right))(1+\alpha_\ell+\alpha_\ell\alpha_{\ell-1}+\ldots+\alpha_\ell\alpha_{\ell-1}\cdots\alpha_2),$$

From the bound on $\alpha_\ell$ we obtain $(1+\alpha_\ell+\alpha_\ell\alpha_{\ell-1}+\ldots+\alpha_\ell\alpha_{\ell-1}\cdots\alpha_2) \leq (2(d_l-1)(d_r-1))^\ell$.

Combining with $1 - \tanh(\mathrm{K}/2) \leq 2e^{-\mathrm{K}}$ we get,

$$d(T^{(\ell)}(\mathsf{c},\Delta_0), S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0)) \leq 2e^{-\mathrm{K}+\ell\cdot\ln(2(d_l-1)(d_r-1))}.$$

∎

The above gives us a bound on the $\mathfrak{B}(S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0))$. Using (ix) Lemma 13 in [18] we get,

$$\mathfrak{B}(S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0)) \leq \mathfrak{B}(T^{(\ell)}(\mathsf{c},\Delta_0))$$
$$+ 2\sqrt{d(T^{(\ell)}(\mathsf{c},\Delta_0), S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0))}$$
$$\leq \mathfrak{B}(T^{(\ell)}(\mathsf{c},\Delta_0)) + 2\sqrt{2}e^{\frac{-\mathrm{K}+\ell\cdot\ln(2(d_l-1)(d_r-1))}{2}}. \quad (8)$$

*Discussion*: In the sequel, we will denote K by $\mathrm{K}^v$ to distinguish between the saturation levels appearing at variable nodes, check nodes and the channel. To summarize, we show that for $\mathrm{K}^v$ large enough, for every iteration the Battacharyya parameter of the symmetric saturated DE remains close to the

Battacharyya of the full DE. In the next section we will relate the symmetric saturated DE to the non-symmetric saturated DE to show that the Battacharyya parameter for the SatBP decoder can also be made small by choosing $\mathrm{K}^v$ large enough.

## IV. Convergence of Nonsymmetrized Saturated DE

The results of the previous section show that, when transmitting below the threshold of the full BP decoder and using sufficiently many iterations, the Battacharrya parameter of the densities in the symmetric SatBP decoder can be small by choosing $\mathrm{K}^v$ large enough. More precisely, consider transmission over a general BMS channel $\mathsf{c}$ such that we are transmitting below the BP threshold of the channel family. Let us assume transmission using $(\lambda,\rho)$ ensemble with average variable node and check node degree given by $d_l$ and $d_r$ respectively. Then, given an $\epsilon > 0$, there exists $\ell_0(\mathsf{c},\epsilon) \in \mathbb{N}$ such that for all $\ell \geq \ell_0$, $\mathfrak{B}(T^{(\ell)}(\mathsf{c},\Delta_0)) \leq \epsilon/2$. Then, by choosing $\mathrm{K}^v$ large enough, specifically $\mathrm{K}^v > \mathrm{K}_0 \triangleq l_0(\mathsf{c},\epsilon)\ln(2(d_l-1)(d_r-1)) + 2\ln\frac{4\sqrt{2}}{\epsilon}$, we have that $\mathfrak{B}(S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0)) \leq \epsilon$.

### A. Non-symmetrized SatBP Decoder

We now show that the Battacharrya parameter for the non-symmetric SatBP decoder can also be made small by choosing $\mathrm{K}^v$ large enough. We first consider a fixed computation tree and then average over the tree ensemble.

We begin with an operational description of symmetrization. Consider a fixed tree $\mathsf{T}$ of depth $\ell$. Let $Y$ denote the vector of received LLR values associated to the variable nodes under the all-zero codeword assumption. In addition, for each variable node we assume an independent random variable uniformly distributed on $[0,1]$. We denote the vector of these variables by $Z = \{Z_v\}$, where $v$ is the index for the variable nodes. Now, the node operations correspond to BP except that outgoing messages from the variable nodes are magnitude saturated at $\mathrm{K}^v$. The independent random variables are used for the flipping operation. The flipping probability for each node is determined by density evolution. If the outgoing message has magnitude $\mathrm{K}^v$ then its sign is flipped if $Z_v < \lambda_v$ where $\lambda_v$ is the appropriate flipping probability.

Let the received LLR magnitude of a variable node $v$ be $x$. The probability with which we flip the bit is such that the final error probability is equal to $\frac{e^{-\mathrm{K}^v}}{1+e^{-\mathrm{K}^v}}$. For received LLR magnitude of $x$, the probability that it is received correctly is $\frac{1}{1+e^{-x}}$. As a consequence we get,

$$\frac{e^{-\mathrm{K}^v}}{1+e^{-\mathrm{K}^v}} = \lambda_v\frac{1}{1+e^{-x}} + (1-\lambda_v)\frac{e^{-x}}{1+e^{-x}},$$

where $\lambda_v$ is the flipping probability of variable node $v$ and $x \geq \mathrm{K}^v$. Solving we get $\lambda_v = \frac{e^{-\mathrm{K}^v}}{1+e^{-\mathrm{K}^v}}\frac{1-e^{-x+\mathrm{K}^v}}{1-e^{-x}} \leq \frac{e^{-\mathrm{K}^v}}{1+e^{-\mathrm{K}^v}}$. Thus the probability that a variable node, with a received LLR magnitude greater than $\mathrm{K}^v$, is not flipped is at least $\frac{1}{1+e^{-\mathrm{K}^v}} \geq 1 - e^{-\mathrm{K}^v}$.

Let us denote the outgoing message at the variable node by $x$. From the above we see that the distribution of the outgoing message $x$ is $S_{\mathrm{K_{sym}}}^{(\ell)}(\mathsf{c},\Delta_0))$. Let us consider the conditional

distribution $p(x\,|\,Y,Z)$. We obtain $S_{K_{sym}}^{(\ell)}(c,\Delta_0))$ by averaging over $Y$, $Z$ and the code ensemble. Let $A_{K^v}$ denote the event that $Z_v \geq 1 - e^{-K^v}$ for each $v$. This is clearly independent of the received values. Assuming a fixed computation tree $T$ (i.e., we suppress dependence on $T$ in the notation) we have

$$p(x\,|\,Y) = p(x\,|\,Y,A_{K^v})p(A_{K^v}) + p(x\,|\,Y,\bar{A}_{K^v})(1 - p(A_{K^v})),$$

where $\bar{A}_{K^v}$ denotes the complement event and, by independence, we can averaging over $Y$ to obtain

$$p(x) = p(x\,|\,A_{K^v})p(A_{K^v}) + p(x\,|\,\bar{A}_{K^v})(1 - p(A_{K^v}))$$

hence

$$p(x\,|\,A_{K^v}) = \frac{p(x) - p(x\,|\,\bar{A}_{K^v})(1 - p(A_{K^v}))}{p(A_{K^v})}$$

Now $p(x\,|\,A_{K^v})$ is the distribution of the non-symmetric SatBP decoder. Intuitively one expects $p(x\,|\,\bar{A}_{K^v})$ to be inferior (higher probability of error, larger Battacharyya parameter) to $p(z\,|\,A_{K^v})$, but this appears difficult to prove. We have, however, $p(A_{K^v}) \geq (1 - e^{-K^v})^{|V(T)|} \geq 1 - e^{-K^v}|V(T)|$ where $|V(T)|$ is the number of variable nodes in the tree.

The above analysis is summarized in the following lemma.

*Lemma 14 (SatBP Decoder versus Symmetrized SatBP):* For any $0 < \epsilon < 1$ and $\ell \in \mathbb{N}$, there exists a $K^v$ large enough such that

$$\mathfrak{B}(S_K^{(\ell)}(c,\Delta_0)) \leq \frac{1}{1-\epsilon}\,\mathfrak{B}(S_{K_{sym}}^{(\ell)}(c,\Delta_0)).$$

*Proof:* From the above analysis we have that for a fixed tree $T$ of depth $\ell$,

$$p(x\,|\,A_{K^v}) = \frac{p(x) - p(x\,|\,\bar{A}_{K^v})(1 - p(A_{K^v}))}{p(A_{K^v})}$$
$$\leq \frac{p(x)}{p(A_{K^v})} \leq \frac{p(x)}{1 - e^{-K}|V(T)|}.$$

where $p(x\,|\,A_{K^v})$ is the distribution of the non-symmetric SatBP decoder. For any fixed number of iterations, the total maximum number of variable nodes in a computation tree is fixed. Hence we can take $K^v$ large enough so that $e^{-K^v}|V(T)| < \epsilon$ for all $T$. Note that the required $K^v$ grows linearly in the number of iterations. Averaging over the tree ensemble and multiplying by the kernel $e^{-x/2}$, we get the desired result. ∎

*Discussion:* Let us summarize. From the above analysis we have that for any $0 < \epsilon < 1/2$, there exists $K^v > 0$, large enough such that the Battacharyya parameter of the SatBP decoder is upper bounded by $\epsilon$. Note that the value of $K^v$ depends on the number of iterations of the full BP required to get its Battacharyya parameter to be at the most $\epsilon/2$. So given a channel $c$ such that the BP decoder is successful when transmitting over $c$, the number of such iterations required is fixed. Call it $\ell_0(c,\epsilon)$. Then, from the above analysis we have that for $K^v \geq K_0 \triangleq \ell_0(c,\epsilon)\ln(2(d_l-1)(d_r-1)) + 2\ln\frac{8\sqrt{2}}{\epsilon}$, $\mathfrak{B}(S_K^{(\ell)}(c,\Delta_0)) \leq \epsilon$. Note that we can make the Battacharyya as small as desired by increasing the number of iterations and consequently increasing $K^v$. But then the saturation value $K^v$ becomes infinite. Hence to make the Battacharyya arbitrarily

small we now need to show that once the Battacharyya parameter is made small enough, by choosing $K^v$ large but fixed, then the subsequent iterations of the SatBP decoder will drive the Battacharyya parameter down to zero. This is the content of the stability analysis done in the next section. We will see that in order to make the Battacharyya parameter arbitrarily small, it is sufficient to bring it close to the stability region. By choosing $\epsilon$ according to equation (13) and arguments following it, we can choose $K^v$ large enough so that we are guaranteed to be in the stability region. Furthermore, we have that $K_0$, defined above, now depends only on the channel $c$ and the degree distribution.

## V. Stability Analysis of the SatBP Decoder

An important part of the asymptotic analysis of LDPC codes involves the analysis of the convergence of DE to a zero error state. In this section we analyze the stability of the SatBP. We begin with some necessary conditions.

For stability of the zero error condition there must exist a positive invariant set of zero error distributions, i.e., a subset $S$ of distributions so that $\mathfrak{E}(s) = 0$ for all $s \in S$ and $S_K(c,s) \in S$. Existence of $S$ follows easily from the compactness of the space of densities and continuity of DE.

*Lemma 15:* Assume the channel $c$ has support at $-L$, $L > 0$. In an irregular ensemble with minimum variable degree $d_l$ the support of all densities in $S$ must lie in $[L/(d_l - 2), \infty]$.

*Proof:* It is obvious that $S = \emptyset$ in an irregular ensemble with $d_l = 1$, so we assume $d_l \geq 2$. We use $a^{(\ell)}$ and $b^{(\ell)}$ to denote the density of the message coming out of the variable nodes and check nodes respectively in the density evolution process. We claim that if $a^{(\ell)}$ has support on $(-\infty, z_\ell]$ with $z_\ell > 0$ then $a^{(\ell+1)}$ has support on $(-\infty, z_{\ell+1}]$ with $z_{\ell+1} = z_\ell - (L - (d_l - 2)z_\ell)$. To see the claim note that $b^{(\ell)}$ also has support on $(-\infty, z_\ell]$ and it follows that $a^{(\ell+1)}$ has support on $(-\infty, z_{\ell+1}]$ where $z_{\ell+1} = (d_l - 1)z_\ell - L = z_\ell - (L - z_\ell(d_l - 2))$.

Assume $a^{(0)} \in S$ has support on $(-\infty, z_0]$ where $z_0 < L/(d_l - 2)$ and define $\delta := L - (d_l - 2)z_0 > 0$. By the above claim it follows from an inductive argument that $a^{(\ell)} \in S$ has support on $(-\infty, z_\ell]$ where $z_\ell$ is a decreasing sequence satisfying $z_\ell \leq z_0 - \ell\delta$. For $\ell$ large enough the right hand side is negative, implying a non-zero error probability, and we obtain a contradiction with the definition of $S$. ∎

### A. Failure of Stability with Degree Two

From Lemma 15 we immediately have

*Lemma 16:* In an irregular ensemble with $\lambda_2 > 0$ no invariant set $S$ exists for any value of $K^v < \infty$ unless the channel is the BEC.

*Proof:* If $d_l = 2$ and the channel is not the BEC and hence has support on $(-\infty, 0)$, then Lemma 15 shows that there can be no positive invariant zero-error set of distributions with support on $[-K^v, K^v]$ for $K^v < \infty$. ∎

In the case of the BEC it can be seen that saturated DE matches unsaturated DE except that the mass at $+\infty$ in unsaturated DE is not placed at $+K^v$. Hence, stability is unaffected by saturation. If the channel has unbounded support

on $(-\infty, 0]$, then there is no possibility of stability under saturation no matter what the degree. A condition on the finite channel support is given in the section on stability with degree at least three.

### B. Near Stability

Even though stability with saturation cannot be achieved in irregular ensembles with degree two variable nodes, it is not surprising that for large $K^v$ the residual error rate can be made very small. For sufficiently large $K^v$ the residual error rate will have no practical consequence. In this section we quantify the residual error rate.

The stability analysis of standard irregular ensembles under BP decoding rests on the relations

$$\mathfrak{B}(\mathsf{c} \circledast \lambda(\mathsf{a})) = \mathfrak{B}(\mathsf{a})\lambda(\mathsf{a}) \qquad (9)$$

and

$$\mathfrak{B}(\rho(\mathsf{a})) \leq 1 - \rho(1 - \mathfrak{B}(\mathsf{a})). \qquad (10)$$

Equality (9) continues to hold without symmetry of $\mathsf{a}$ or $\mathsf{c}$. The inequality (10), however, does not hold without symmetry. In Appendix A we prove a more general form of the following.

*Lemma 17:* Let the incoming L-densities at a degree $d+1$ check node be $\mathsf{a}_1, ..., \mathsf{a}_d$ and let $\mathsf{b}$ be the outgoing density. Then

$$\mathfrak{B}(\mathsf{b}) \leq \sum_{i=1}^{d} \mathfrak{B}(\mathsf{a}_i).$$

*Discussion:* The above result holds for a wide range of check node update operations including BP and the min-sum decoder.

Throughout this section we will use $\mathsf{a}$ ($\mathsf{b}$) to denote the density coming out of a variable node (check node). We also use $\mathsf{a}^{(n)}$ and $\mathsf{b}^{(n)}$ to denote the densities coming out of the variable nodes and check nodes at the $n$th iteration of the saturated DE recursion. We prove the following result,

*Lemma 18:* Consider an irregular ensemble with minimum variable node degree $d_{\min} \geq 2$. Assume $\lambda_2 \rho'(1)\mathfrak{B}(\mathsf{c}) < 1$. Then, there exists a constant $x^*$, a constant $N$, and a constant $C(d_{\min})$ such that, for all $K^v$ large enough, if for some $n_0$ we have $\mathfrak{B}(\mathsf{a}^{(n_0)}) \leq x*$ then $\mathfrak{B}(\mathsf{a}^{(n)}) \leq C(d_{\min})e^{-K^v/2}$ for all $n \geq n_0 + N$. Moreover, if $d_{\min} > 2$ we can have $C(d_{\min}) = 3$.

*Proof:* To incorporate saturation into the analysis based on the Battacharyya parameter we have the inequality for any $K > 0$,

$$\mathfrak{B}(\lfloor \mathsf{a} \rfloor_K) \leq \mathfrak{B}(\mathsf{a}) + e^{-K/2}.$$

Indeed, we have

$$\mathfrak{B}(\lfloor \mathsf{a} \rfloor_K) = e^{K/2} \int_{-\infty}^{-K} \mathsf{a}(x)dx + \int_{-\infty}^{+\infty} \mathbb{1}_{\{|x|<K\}}\mathsf{a}(x)e^{-x/2}dx$$
$$+ e^{-K/2} \int_{K}^{\infty} \mathsf{a}(x)dx$$
$$\leq \int_{-\infty}^{-K} \mathsf{a}(x)e^{-x/2}dx + \int_{-\infty}^{+\infty} \mathbb{1}_{\{|x|<K\}}\mathsf{a}(x)e^{-x/2}dx$$
$$+ \int_{K}^{\infty} \mathsf{a}(x)e^{-x/2}dx + e^{-K/2}$$
$$= \mathfrak{B}(\mathsf{a}) + e^{-K/2}, \qquad (11)$$

where the last inequality follows since $e^{-K/2}\int_{K}^{\infty}\mathsf{a}(x)dx \leq e^{-K/2}\int_{-\infty}^{\infty}\mathsf{a}(x)dx = e^{-K/2}$. As a result of the saturation of messages, we see that the minimum value of the Battacharyya parameter is equal to $e^{-K/2}$ and we can therefore not hope to reach a smaller value.

*Minimum variable node degree equal to 2:* Let us assume $d_{\min} = 2$, i.e., $\lambda_2 > 0$. Let $\mathsf{a}^{(n_0)}$ be any $L$-density which need not be symmetric. Consider

$$g(x) := \lambda_2\,\mathfrak{B}(\mathsf{c})\rho'(1) + (1 - \lambda_2)\,\mathfrak{B}(\mathsf{c})(\rho'(1))^2 x.$$

Since $\lambda_2\,\mathfrak{B}(\mathsf{c})\rho'(1) < 1$, there exists an $x^* > 0$ such that $g(x^*) < 1$. Choose $x^*$ such that $g(x^*) < 1$ and $\rho'(1)x^* < 1$. Now assume $\mathfrak{B}(\mathsf{a}^{(n_0)}) \leq x^*$. Choose $K^v$ large enough such that $\frac{1}{1-g(x^*)}e^{-K^v/2} < x^*$.

Let us perform the saturated DE recursion once. We have,

$$\mathfrak{B}(\mathsf{a}^{(n_0+1)}) = \mathfrak{B}(\lfloor \mathsf{c} \circledast \lambda(\rho(\mathsf{a}^{(n_0)}))\rfloor_{K^v})$$
$$\overset{(11)}{\leq} \mathfrak{B}(\mathsf{c} \circledast \lambda(\rho(\mathsf{a}^{(n_0)}))) + e^{-K^v/2}$$
$$= \mathfrak{B}(\mathsf{c})\lambda\Big(\sum_i \rho_i\,\mathfrak{B}((\mathsf{a}^{(n_0)})^{\boxplus(i-1)})\Big) + e^{-K^v/2}$$
$$\overset{\text{Lemma 17}}{\leq} \mathfrak{B}(\mathsf{c})\lambda\Big(\mathfrak{B}(\mathsf{a}^{(n_0)})\sum_i(i-1)\rho_i\Big) + e^{-K^v/2}$$
$$\overset{\text{since } \rho'(1)\mathfrak{B}(\mathsf{a}^{(n_0)})<1}{\leq} \lambda_2\,\mathfrak{B}(\mathsf{c})\rho'(1)\mathfrak{B}(\mathsf{a}^{(n_0)})$$
$$+ (1 - \lambda_2)\,\mathfrak{B}(\mathsf{c})(\rho'(1)\mathfrak{B}(\mathsf{a}^{(n_0)}))^2 + e^{-K^v/2}$$
$$= g(\mathfrak{B}(\mathsf{a}^{(n_0)}))\mathfrak{B}(\mathsf{a}^{(n_0)}) + e^{-K^v/2}$$
$$\leq g(x^*)\mathfrak{B}(\mathsf{a}^{(n_0)}) + e^{-K^v/2} \qquad (12)$$
$$\leq g(x^*)x^* + e^{-K^v/2}$$
$$\leq x^*,$$

where the last inequality follows from the choice of $K^v$.

By induction, the above inequality gives $\mathfrak{B}(\mathsf{a}^{(n)}) \leq x^*$ for all $n \geq n_0$. Consider any $n = n_0 + k$. Also by induction on (12), we get

$$\mathfrak{B}(\mathsf{a}^{(n_0+k)}) \leq x^*(g(x^*))^k + e^{-K^v/2}\sum_{j=0}^{k-1}(g(x^*))^j$$
$$= x^*(g(x^*))^k + e^{-K^v/2}\frac{1 - (g(x^*))^k}{1 - g(x^*)}.$$

It follows that any $\epsilon > 0$ and all $k$ large enough we have

$$\mathfrak{B}(\mathsf{a}^{(n_0+k)}) \leq e^{-K^v/2}\frac{1 - \epsilon}{1 - g(x^*)}.$$

*Minimum variable node degree equal to 3:* Let us now assume that the minimum variable node degree is 3. Let us denote,

$$f(x) = \lambda_3\,\mathfrak{B}(\mathsf{c})\rho'(1)^2 x + (1 - \lambda_3)\,\mathfrak{B}(\mathsf{c})\rho'(1)^3 x^2. \qquad (13)$$

Choose $x^* > 0$ such that $f(x^*) \leq 1/2$ and $\rho'(1)x^* < 1$. Let $n_0$ be such that $\mathfrak{B}(\mathsf{a}^{(n_0)}) \leq x^*$. Choose $K^v$ large enough so that $2e^{-K^v/2} < x^*$. Following the previous analysis, we have for all $n \geq n_0$

$$\mathfrak{B}(\mathsf{a}^{(n+1)}) \leq \lambda_3\,\mathfrak{B}(\mathsf{c})(\rho'(1)\mathfrak{B}(\mathsf{a}^{(n)}))^2$$
$$+ (1 - \lambda_3)\,\mathfrak{B}(\mathsf{c})(\rho'(1)\mathfrak{B}(\mathsf{a}^{(n)}))^3 + e^{-K^v/2}$$

A little algebra then shows that there exists $N > n_0$ so that for all $n \geq N$ we have

$$\mathfrak{B}\left(\mathsf{a}^{(n)}\right) \leq 3e^{-\mathrm{K}^v/2} \tag{14}$$

$$\mathfrak{B}\left(\mathsf{b}^{(n)}\right) \leq 3\rho'(1)e^{-\mathrm{K}^v/2} \tag{15}$$

where $\mathsf{b}^{(n)}$ denotes the density coming out of the check nodes. Also, (15) follows from (14) and Lemma 17. ∎

The "near stability" analysis done above can clearly not show convergence to zero error although it can be used to show convergence to relatively small error rate. As we showed above, unlike the unsaturated case, zero error rate convergence cannot be achieved with the saturated decoder when degree two variable nodes are included. For degree three and higher, stability can be shown but a refined analysis is needed.

### C. Stability Analysis with Minimum Variable Node Degree Equal to Three

In this section we consider irregular ensembles where the minimum variable node degree is at least three. We generalize the standard stability analysis by separating out the saturated probability mass and tracking it through the variable node and check node updates. For simplicity we shall restrict to right regular ensembles. We show that convergence to zero error rate occurs and that convergence is exponential in iteration. In the unsaturated case this can be achieved with degree two variable nodes and with degree three and above doubly exponential convergence occurs. In subsequent sections we show that double exponential convergence can be attained in the saturated case for degree four and above although a modification is needed for degree four. For degree three doubly exponential convergence can be recovered but only with the dramatic and likely impractical step of erasing all received values near the end of the decoding.

We assume regular check nodes with degree $d_r$ and we let $\mathrm{K}^p$ denote the magnitude of an outgoing message when all incoming messages have magnitude $\mathrm{K}^v$. Although we focus on BP-like decoding our analysis applies to other algorithms such as min-sum, in which case we have $\mathrm{K}^p = \mathrm{K}^v$. In general, if $\mathrm{K}_1, ..., \mathrm{K}_{d_r-1}$ are incoming message magnitudes at a check node then we assume that the corresponding outgoing magnitude $\mathrm{K}_{\mathrm{out}}$ satisfies

$$-\ln \sum_{i=1}^{d_r-1} e^{-\mathrm{K}_i} \leq \mathrm{K}_{\mathrm{out}} \leq \min_i\{\mathrm{K}_i\} \tag{16}$$

Both conditions are satisfied by BP and min-sum. E.g., for BP we can write explicitly $\tanh(\mathrm{K}_i/2) = (1 - e^{-\mathrm{K}_i/2})/(1 + e^{-\mathrm{K}_i/2})$ and then some algebra[1] gives us (16). We note in passing that the left inequality implies $-\ln \sum_{i=1}^{d_r-1} e^{-\lambda \mathrm{K}_i} \leq \lambda \mathrm{K}_{\mathrm{out}}$ for all $\lambda \in [0, 1]$. We will make use of the case $\lambda = \frac{1}{2}$.

Messages entering a check node update $\mathsf{a}$ have the form

$$\mathsf{a} = \gamma D(p, \mathrm{K}^v) + \bar{\gamma}\mathsf{m}$$

[1]Indeed, it is not hard to see that $\frac{1-e^{-\mathrm{K}_{\mathrm{out}}}}{1+e^{-\mathrm{K}_{\mathrm{out}}}} = \frac{1-\sum_i e^{-\mathrm{K}_i}+A}{1+\sum_i e^{-\mathrm{K}_i}+B}$, where $A, B \geq 0$. Furthermore, one can show that $A(1 + \sum_i e^{-\mathrm{K}_i}) \geq B(1 - \sum_i e^{-\mathrm{K}_i})$, which implies that $\frac{1-e^{-\mathrm{K}_{\mathrm{out}}}}{1+e^{-\mathrm{K}_{\mathrm{out}}}} \geq \frac{1-\sum_i e^{-\mathrm{K}_i}}{1+\sum_i e^{-\mathrm{K}_i}}$ giving us the inequality.

where $\mathsf{m}$ is supported on $(-\mathrm{K}^v, \mathrm{K}^v)$ and has total mass 1 (if it has zero probability we have $\bar{\gamma} = 0$.)

Messages entering a variable node update $\mathsf{b}$ have the form

$$\mathsf{b} = \gamma D(p, \mathrm{K}^p) + \bar{\gamma}\mathsf{m}$$

where $\mathrm{K}^p \leq \mathrm{K}^v$ is the outgoing magnitude at a check when all incoming magnitudes equal $\mathrm{K}^v$ and $\mathsf{m}$ is supported on $(-\mathrm{K}^p, \mathrm{K}^p)$. From (16) we have $e^{-\mathrm{K}^p} \leq (d_r - 1)e^{-\mathrm{K}^v}$. We assume $\mathrm{K}^v > 2\ln(d_r - 1)$ large enough so that $2\mathrm{K}^p > \mathrm{K}^v$. In the subsequent analysis we also assume that the support of the channel $\mathsf{c}$ is restricted to $(-\mathrm{K}^c, \mathrm{K}^c)$ where we assume that $\mathrm{K}^c \leq 2\mathrm{K}^p - \mathrm{K}^v$.

The analysis tracks the quantities $\gamma p$ and $\bar{\gamma}\mathfrak{B}(\mathsf{m})$. For stability we aim to show that both quantities converge to 0. Note that this implies that $\gamma \to 1$. In the standard stability analysis of irregular ensembles and full BP, one tracks the Battacharyya parameter of the density through the DE iterations when the density is near $\Delta_\infty$. At the check node the Battacharyya parameter undergoes a constant factor gain with a factor of $\rho'(1)$. On the variable node side the parameter is raised to the power of the minimum variable node degree less one, and scaled the channel Battacharyya. Thus, one arrives at the stability condition $\lambda_2\rho'(1)\mathfrak{B}(\mathsf{c}) < 1$. If the minimum variable node degree is three then the update bound takes the form $\mathfrak{B}\left(\mathsf{a}^{(\ell+1)}\right) \leq C\,\mathfrak{B}\left(\mathsf{a}^{(\ell)}\right)^2$, for some positive constant $C$, and one obtains doubly exponential decay in $\mathfrak{B}\left(\mathsf{a}^{(\ell)}\right)$. For the saturated case we accomplish something similar, although the conditions are different. As a first step we show that we still have constant factor gain at check nodes.

*1) Check Node Analysis:* We assume a right regular ensemble with check degree $d + 1$. Let us represent the density entering the check node as $\gamma D(p, \mathrm{K}^v) + \bar{\gamma}\mathsf{m}$ where $\mathsf{m}$ is a density supported on $(-\mathrm{K}^v, \mathrm{K}^v)$. Then the density emerging out of the check node is given by $\gamma' D(p', \mathrm{K}^p) + \bar{\gamma}'\mathsf{m}' \triangleq (\gamma D(p, \mathrm{K}^v) + \bar{\gamma}\mathsf{m})^{\boxplus d}$, where $\mathrm{K}^p$ is the magnitude of the check output when all inputs are $\mathrm{K}^v$, which satisfies $\mathrm{K}^v - \ln d \leq \mathrm{K}^p \leq \mathrm{K}^v$, and support of $\mathsf{m}'$ is also $(-\mathrm{K}^p, \mathrm{K}^p)$. Let us now perform the computation explicitly. In this section we use $D$ to denote $D(p, \mathrm{K}^v)$. We have,

$$(\gamma D(p, \mathrm{K}^v) + \bar{\gamma}\mathsf{m})^{\boxplus d} = \sum_{k=0}^{d} \binom{d}{k} \gamma^k \bar{\gamma}^{d-k} D^{\boxplus k} \boxplus \mathsf{m}^{\boxplus d-k}$$

$$= \bar{\gamma}^d \mathsf{m}^{\boxplus d} + \sum_{k=1}^{d-1} \binom{d}{k} \gamma^k \bar{\gamma}^{d-k} D^{\boxplus k} \boxplus \mathsf{m}^{\boxplus d-k} + \gamma^d D^{\boxplus d}$$

where we have separated out two of the terms from the sum. Although we have indicated that density evolution for check node update is associative, which it is for min-sum and sum-product algorithms, we do not actually require the associative property and a density $D^{\boxplus k} \boxplus \mathsf{m}^{\boxplus d-k}$ can simply be understood as the outgoing one corresponding to $k$ incoming messages from density $D$ and $d - k$ messages from density $\mathsf{m}$.

By Lemma 24 we have for $1 \leq k \leq d - 1$,

$$\mathfrak{B}\left(D^{\boxplus k} \boxplus \mathsf{m}^{\boxplus d-k}\right) \leq (1 + k(e^{\frac{\mathrm{K}^v}{2}}\mathfrak{B}(D) - 1))(d - k)\mathfrak{B}(\mathsf{m})$$

$$\leq ke^{\frac{\mathrm{K}^v}{2}}\mathfrak{B}(D)(d - k)\mathfrak{B}(\mathsf{m}).$$

A little algebra shows that

$$\sum_{k=1}^{d-1}\binom{d}{k}\gamma^k\bar\gamma^{d-k}k(d-k) = \gamma\bar\gamma d(d-1)$$

and we now obtain

$$\mathfrak{B}\left(\sum_{k=1}^{d-1}\binom{d}{k}\gamma^k\bar\gamma^{d-k}\mathsf{D}^{\boxtimes k}\boxtimes\mathsf{m}^{\boxtimes d-k}\right)$$
$$\leq \gamma\bar\gamma d(d-1)e^{\frac{\mathrm{K}^v}{2}}\mathfrak{B}\left(\mathsf{D}\right)\mathfrak{B}\left(\mathsf{m}\right).$$

Lemma 24 also gives

$$\mathfrak{B}\left(\mathsf{m}^{\boxtimes d}\right)\leq d\,\mathfrak{B}\left(\mathsf{m}\right),$$

so we now have

$$\mathfrak{B}\left(\sum_{k=0}^{d-1}\binom{d}{k}\gamma^k\bar\gamma^{d-k}\mathsf{D}^{\boxtimes k}\boxtimes\mathsf{m}^{\boxtimes d-k}\right)$$
$$\leq d\Big((d-1)\gamma e^{\frac{\mathrm{K}^v}{2}}\mathfrak{B}\left(\mathsf{D}\right)+1\Big)\bar\gamma\,\mathfrak{B}\left(\mathsf{m}\right).$$

We have $\gamma' D(p',\mathrm{K}^p) = \gamma^d\mathsf{D}^{\boxtimes d}$ so $p' = \frac{1-(1-2p)^d}{2}\leq dp$ where we have used Lemma 23 to obtain the last inequality.

We summarize the results as follows.

*Lemma 19:* Let the incoming density to a degree $d+1$ check node be $\gamma D(p,\mathrm{K}^v) + \bar\gamma\mathsf{m}$. Then the outgoing density $\gamma' D(p',\mathrm{K}^p) + \bar\gamma'\mathsf{m}'$ satisfies the following

$$\begin{bmatrix}\bar\gamma'\,\mathfrak{B}(\mathsf{m}')\\ \gamma'p'\end{bmatrix}\leq d\begin{bmatrix}\xi & 0\\ 0 & 1\end{bmatrix}\begin{bmatrix}\bar\gamma\,\mathfrak{B}(\mathsf{m})\\ \gamma p\end{bmatrix}$$

where $\xi = \Big((d-1)\gamma e^{\frac{\mathrm{K}^v}{2}}\mathfrak{B}\left(D(p,\mathrm{K}^v)\right)+1\Big)$.
In the stability region we will have the bound $\xi\leq 3$ so we see that we have been able to obtain a linear growth bound for the check node density evolution update.

*2) Variable Node Analysis:* Consider a variable node of degree $d+1$ and incoming density

$$\mathsf{b} = \gamma D(p,\mathrm{K}^p) + \bar\gamma\mathsf{m}.$$

The outgoing density from the variable node has the form

$$\mathsf{a} = \gamma' D(p',\mathrm{K}^v) + \bar\gamma'\mathsf{m}'.$$

The density $\mathsf{a}$ is the saturation of

$$\sum_{k=0}^{d-2}\binom{d}{k}\gamma^k\bar\gamma^{d-k}\mathsf{c}\circledast\mathsf{D}^{\circledast k}\circledast\mathsf{m}^{\circledast(d-k)}$$
$$+ d\bar\gamma\gamma^{d-1}\mathsf{c}\circledast\mathsf{D}^{\circledast d-1}\circledast\mathsf{m} + \gamma^d\mathsf{c}\circledast\mathsf{D}^{\circledast d} \tag{17}$$

where in this section we use $\mathsf{D}$ to denote $D(p,\mathrm{K}^p)$. In particular $\gamma'p'$ is the total mass of this density on $(-\infty,-\mathrm{K}^v]$ and $\gamma'\mathsf{m}'$ is the restriction of this density to $(-\mathrm{K}^v,\mathrm{K}^v)$.

We see in the above decomposition that incoming messages either have magnitude $\mathrm{K}^p$, i.e. are drawn from $\mathsf{D}$, or they are drawn from $\mathsf{m}$ and therefore take values in $(-\mathrm{K}^p,\mathrm{K}^p)$. We can define a type for an outgoing message consisting of a triple of non-negative integers $(n_-,n_\mathsf{m},n_+)$ where $n_-+n_\mathsf{m}+n_+ = d$. Here $n_-$ represents the number of $-\mathrm{K}^p$ incoming messages, $n_+$ the number of $+\mathrm{K}^p$ incoming messages, and $n_\mathsf{m}$ the number of incoming message drawn from $\mathsf{m}$ that comprise the outgoing message. Our analysis will pay special attention

to the terms with $n_\mathsf{m} = 0$ and $n_\mathsf{m} = 1$ which is why we distinguished these terms.

A handy elementary result is the following.
*Lemma 20:* If $a,b\geq 0$ and $k\leq d$ then

$$\sum_{i=0}^{d-k}\binom{d}{i}a^{d-i}b^i\leq\binom{d}{k}a^k(a+b)^{d-k}$$

*Proof:* For $i\leq d-k$ we have,

$$\binom{d}{i}\leq\binom{d}{i}\binom{d-i}{k} = \binom{d}{k}\binom{d-k}{i}.$$

and the lemma follows from the binomial theorem. We remark that there is an alternate form since $\binom{d}{k} = \binom{d}{d-k}$. ∎

Let us consider the three parts of (17). The first part comprises messages types $(n_-,n_\mathsf{m},n_+)$ where $n_\mathsf{m}\geq 2$. The second part comprises messages types $(n_-,n_\mathsf{m},n_+)$ with $n_\mathsf{m} = 1$ and the third part comprises messages types $(n_-,n_\mathsf{m},n_+)$ with $n_\mathsf{m} = 0$. We will consider the contribution of each part to $\gamma'p'$ and to $\bar\gamma'\mathsf{m}'$.

Let us first consider $\gamma'p'$. We use the bound $\int_{-\infty}^{-\mathrm{K}}\mathsf{a}(x)dx\leq e^{-\frac{\mathrm{K}}{2}}\mathfrak{B}(\mathsf{a})$, which is valid for any density and any $\mathrm{K}\geq 0$, Lemma 20 and the multiplicative property of Battacharyya parameter at the variable node side to obtain

$$\int_{-\infty}^{-\mathrm{K}^v}\sum_{k=0}^{d-2}\binom{d}{k}\gamma^k\bar\gamma^{d-k}\mathsf{c}\circledast\mathsf{D}^{\circledast k}\circledast\mathsf{m}^{\circledast(d-k)}(x)dx$$
$$\leq e^{-\frac{\mathrm{K}^v}{2}}\frac{d(d-1)}{2}(\bar\gamma\,\mathfrak{B}(\mathsf{m}))^2\mathfrak{B}(\mathsf{c})\mathfrak{B}(\mathsf{b})^{d-2}. \tag{18}$$

Now we consider contributions from $n_\mathsf{m} = 1$. A message of type $(n_-,1,n_+)$ has value at most $(n_+-n_-)\mathrm{K}^p+(\mathrm{K}^p+\mathrm{K}^c)$ and at least $(n_+-n_-)\mathrm{K}^p-(\mathrm{K}^p+\mathrm{K}^c)$. Recall that $(-\mathrm{K}^c,\mathrm{K}^c)$ is the channel support. Hence if $n_+-n_- > 0$ then the message has value greater than $-\mathrm{K}^v$ and if $n_+-n_- < -1$ then the message has value less than $-\mathrm{K}^v$. If $n_+-n_- = 0$ then the message has value less than $-\mathrm{K}^v$ only if the contribution from $\mathsf{c}\circledast\mathsf{m}$ is less than $-\mathrm{K}^v$. If $n_+-n_- = -1$ then the message can have value less than $-\mathrm{K}^v$ only if the contribution from $\mathsf{c}\circledast\mathsf{m}$ is less than $0$. Hence, we obtain

$$\int_{-\infty}^{-\mathrm{K}^v}\mathsf{c}\circledast\mathsf{m}\circledast\mathsf{D}^{d-1}(x)dx\leq$$
$$\begin{cases}\sum_{j=0}^{\frac{d-4}{2}}\binom{d-1}{j}p^{d-1-j}\bar p^j\\ \quad+\binom{d-1}{\frac{d-2}{2}}p^{\frac{d}{2}}\bar p^{\frac{d-2}{2}}\mathfrak{E}(\mathsf{c}\circledast\mathsf{m}) & d\text{ even}\\[2mm]\sum_{j=0}^{\frac{d-3}{2}}\binom{d-1}{j}p^{d-1-j}\bar p^j\\ \quad+\binom{d-1}{\frac{d-1}{2}}p^{\frac{d-1}{2}}\bar p^{\frac{d-1}{2}}e^{-\frac{\mathrm{K}^v}{2}}\mathfrak{B}(\mathsf{c}\circledast\mathsf{m}) & d\text{ odd}\end{cases} \tag{19}$$

Note that for the case $d$ even, we use $\mathfrak{E}(\mathsf{c}\circledast\mathsf{m})$ to bound the contribution from $(\mathsf{c}\circledast\mathsf{m})(x)$ for $x\leq 0$. Now we consider contributions from $n_\mathsf{m} = 0$. A message of type $(n_-,0,n_+)$ has value at most $(n_+-n_-)\mathrm{K}^p+(\mathrm{K}^c)$ and at least $(n_+-n_-)\mathrm{K}^p-(\mathrm{K}^c)$. Hence if $n_+-n_-\geq 0$ then the message has value greater than $-\mathrm{K}^v$ and if $n_+-n_- < -1$ then the message has value less than $-\mathrm{K}^v$. If $n_+-n_- = -1$ then the message can have value less than $-\mathrm{K}^v$ only if the contribution from $\mathsf{c}$

is less than 0. Hence, we obtain

$$\int_{-\infty}^{-\mathrm{K}^v} \mathsf{c} \circledast \mathsf{D}^d(x)dx \leq$$

$$\begin{cases} \sum_{j=0}^{\frac{d-2}{2}} \binom{d}{j} p^{d-j}\bar{p}^j & d \text{ even} \\ \sum_{j=0}^{\frac{d-3}{2}} \binom{d}{j} p^{d-j}\bar{p}^j + \binom{d}{\frac{d-1}{2}} p^{\frac{d+1}{2}}\bar{p}^{\frac{d-1}{2}} \mathfrak{E}(\mathsf{c}) & d \text{ odd} \end{cases}$$
$$(20)$$

Using the bound $\mathfrak{E}(\mathsf{c} \circledast \mathsf{m}) \leq \mathfrak{B}(\mathsf{c} \circledast \mathsf{m})$ and Lemma 20 we obtain from (19)

$$\int_{-\infty}^{-\mathrm{K}^v} \mathsf{c} \circledast \mathsf{m} \circledast \mathsf{D}^{d-1}(x)dx \leq$$

$$\begin{cases} \binom{d-1}{\frac{d-2}{2}} p^{\frac{d}{2}} (p + \mathfrak{B}(\mathsf{c} \circledast \mathsf{m})) & d \text{ even} \\ \binom{d-1}{\frac{d-1}{2}} p^{\frac{d-1}{2}} (p + e^{-\frac{\mathrm{K}^v}{2}} \mathfrak{B}(\mathsf{c} \circledast \mathsf{m})) & d \text{ odd} \end{cases}$$

and using the bound $\mathfrak{E}(\mathsf{c}) \leq 1$ and Lemma 20 we obtain from (20)

$$\int_{-\infty}^{-\mathrm{K}^v} \mathsf{c} \circledast \mathsf{D}^d(x)dx \leq \binom{d}{\lfloor \frac{d-1}{2} \rfloor} p^{\lceil \frac{d+1}{2} \rceil}.$$

Combining the above into (17) we have

$$\gamma'p' \leq e^{-\frac{\mathrm{K}^v}{2}} \frac{d(d-1)}{2} (\bar{\gamma}\, \mathfrak{B}(\mathsf{m}))^2 \mathfrak{B}(\mathsf{c}) \mathfrak{B}(\mathsf{b})^{d-2}$$
$$+ d\binom{d-1}{\lfloor \frac{d-1}{2} \rfloor} (\gamma p)^{\lfloor \frac{d}{2} \rfloor} ((\gamma p) + \mathfrak{B}(\mathsf{c})(\bar{\gamma}\, \mathfrak{B}(\mathsf{m})))$$
$$+ \binom{d}{\lfloor \frac{d-1}{2} \rfloor} (\gamma p)^{\lceil \frac{d+1}{2} \rceil}$$

$$\leq e^{-\frac{\mathrm{K}^v}{2}} \frac{d(d-1)}{2} (\bar{\gamma}\, \mathfrak{B}(\mathsf{m}))^2 \mathfrak{B}(\mathsf{c}) \mathfrak{B}(\mathsf{b})^{d-2}$$
$$+ (d+1)(4\gamma p)^{\lfloor \frac{d}{2} \rfloor + 1} + d(4\gamma p)^{\lfloor \frac{d}{2} \rfloor} \mathfrak{B}(\mathsf{c})(\bar{\gamma}\, \mathfrak{B}(\mathsf{m}))$$
$$(21)$$

where we have used $\binom{d}{\lfloor \frac{d-1}{2} \rfloor} \leq 2^{d-1}$. We note that when $d$ is odd we can add another factor of $e^{-\frac{\mathrm{K}^v}{2}}$ to the last term.

Now we consider the contribution to $\bar{\gamma}'\mathsf{m}'$. Let us introduce the notation $\lfloor \mathsf{a} \rfloor_{\mathrm{K}}^\circ(x) = \mathsf{a}(x)\mathbb{1}_{\{|x| < \mathrm{K}\}}$. First we note that the contribution to $\mathfrak{B}(\mathsf{m}')$ from types with $n_\mathsf{m} \geq 2$ is upper bounded by

$$\mathfrak{B}\Big(\sum_{k=0}^{d-2} \binom{d}{k} \gamma^k \bar{\gamma}^{d-k} \mathsf{c} \circledast \mathsf{D}^{\circledast k} \circledast \mathsf{m}^{\circledast(d-k)}\Big) \leq$$
$$\frac{d(d-1)}{2} (\bar{\gamma}\, \mathfrak{B}(\mathsf{m}))^2 \mathfrak{B}(\mathsf{c}) \mathfrak{B}(\mathsf{b})^{d-2},$$

where we applied Lemma 20.

Let us introduce the notation $q = e^{\frac{\mathrm{K}^p}{2}} p$ and $\tilde{q} = e^{-\frac{\mathrm{K}^p}{2}} \bar{p}$. Note that for any density $\mathsf{a}$ we have $\mathfrak{B}(\mathsf{a} \circledast \Delta_{\mathrm{K}}) = e^{-\mathrm{K}} \mathfrak{B}(\mathsf{a})$.

Now we consider the contribution from types with $n_\mathsf{m} = 1$. A type $(n_-, 1, n_+)$ will have a non-zero contribution only if the interval centered on $(n_+ - n_-)\mathrm{K}^p$ of width $2(\mathrm{K}^\mathsf{c} + \mathrm{K}^p)$ intersects $(-\mathrm{K}^v, \mathrm{K}^v)$. Note that $\mathsf{m}' = \lfloor \mathsf{c} \circledast \mathsf{m} \circledast \mathsf{D}^{d-1} \rfloor_{\mathrm{K}^v}^\circ$. Since we assume $2\mathrm{K}^p \geq \mathrm{K}^\mathsf{c} + \mathrm{K}^v$ and $\mathrm{K}^p \leq \mathrm{K}^v$ we obtain

$$\mathfrak{B}(\lfloor \mathsf{c} \circledast \mathsf{m} \circledast \mathsf{D}^{d-1} \rfloor_{\mathrm{K}^v}^\circ) \leq$$

$$\mathfrak{B}(\mathsf{c})\, \mathfrak{B}(\mathsf{m}) \begin{cases} \sum_{j=\frac{d-2}{2}}^{\frac{d}{2}} \binom{d-1}{j} q^{d-1-j}\tilde{q}^j & d \text{ even} \\ \sum_{j=\frac{d-3}{2}}^{\frac{d+1}{2}} \binom{d-1}{j} q^{d-1-j}\tilde{q}^j & d \text{ odd} \end{cases}$$

Using the inequality $2\binom{d-1}{\frac{d-3}{2}} \geq \binom{d-1}{\frac{d-1}{2}}$ for odd $d$ we can write this as

$$\mathfrak{B}(\lfloor \mathsf{c} \circledast \mathsf{m} \circledast \mathsf{D}^{d-1} \rfloor_{\mathrm{K}^v}^\circ)$$

$$\leq \mathfrak{B}(\mathsf{c})\, \mathfrak{B}(\mathsf{m}) \begin{cases} \binom{d-1}{\frac{d}{2}} (q\tilde{q})^{\frac{d-2}{2}}(q + \tilde{q}) & d \text{ even} \\ \binom{d-1}{\frac{d-3}{2}} (q\tilde{q})^{\frac{d-3}{2}}(q + \tilde{q})^2 & d \text{ odd} \end{cases}$$

$$= \mathfrak{B}(\mathsf{c})\, \mathfrak{B}(\mathsf{m}) \begin{cases} \binom{d-1}{\frac{d}{2}} (p\bar{p})^{\frac{d-2}{2}} \mathfrak{B}(D(p, \mathrm{K}^p)) & d \text{ even} \\ \binom{d-1}{\frac{d-3}{2}} (p\bar{p})^{\frac{d-3}{2}} \mathfrak{B}(D(p, \mathrm{K}^p))^2 & d \text{ odd} \end{cases}$$

$$\leq \mathfrak{B}(\mathsf{c})\, \mathfrak{B}(\mathsf{m}) \begin{cases} (4p)^{\frac{d-2}{2}} \mathfrak{B}(D(p, \mathrm{K}^p)) & d \text{ even} \\ 2(4p)^{\frac{d-3}{2}} \mathfrak{B}(D(p, \mathrm{K}^p))^2 & d \text{ odd} \end{cases}$$

Finally we consider the contribution from types with $n_\mathsf{m} = 0$. A type $(n_-, 0, n_+)$ will have a non-zero contribution only if the interval centered on $(n_+ - n_-)\mathrm{K}^p$ of width $2\mathrm{K}^\mathsf{c}$ intersects $(-\mathrm{K}^v, \mathrm{K}^v)$. Hence we obtain

$$\mathfrak{B}(\lfloor \mathsf{c} \circledast \mathsf{D}^d \rfloor_{\mathrm{K}^v}^\circ)$$

$$\leq \mathfrak{B}(\mathsf{c}) \begin{cases} \binom{d}{\frac{d}{2}} q^{\frac{d}{2}}\tilde{q}^{\frac{d}{2}} & d \text{ even} \\ \sum_{j=\frac{d-1}{2}}^{\frac{d+1}{2}} \binom{d}{j} q^{d-j}\tilde{q}^j & d \text{ odd} \end{cases}$$

$$= \mathfrak{B}(\mathsf{c}) \begin{cases} \binom{d}{\frac{d}{2}} (p\bar{p})^{\frac{d}{2}} & d \text{ even} \\ \binom{d}{\frac{d-1}{2}} (p\bar{p})^{\frac{d-1}{2}} \mathfrak{B}(D(p, \mathrm{K}^p)) & d \text{ odd} \end{cases}$$

$$\leq \mathfrak{B}(\mathsf{c}) \begin{cases} (4p)^{\frac{d}{2}} & d \text{ even} \\ (4p)^{\frac{d-1}{2}} \mathfrak{B}(D(p, \mathrm{K}^p)) & d \text{ odd} \end{cases}$$

To get the final bound on $\gamma'\mathfrak{B}(\mathsf{m}')$ we need to multiply the above bounds by $d\bar{\gamma}\gamma^{d-1}$ when $n_\mathsf{m} = 1$ and by $\gamma^d$ when $n_\mathsf{m} = 0$. In the next section we will use $\mathfrak{B}(D(p, \mathrm{K}^p)) \leq \mathfrak{B}(\mathsf{b})$ to further bound the above expressions.

### D. Stability with Minimum Degree 3.

Let us assume that the minimum variable node degree, given by $d + 1$, is at least three and a right regular degree $d_r + 1$.

In view of (14) and (15) we may assume $\mathfrak{B}(\mathsf{a}^{(n)}) \leq 3e^{-\frac{\mathrm{K}^v}{2}}$ which implies $\mathfrak{B}(\mathsf{b}^{(n)}) \leq 3d_r e^{-\frac{\mathrm{K}^v}{2}}$, $\gamma^{(n)}p^{(n)}e^{\frac{\mathrm{K}^v}{2}} \leq 3e^{-\frac{\mathrm{K}^v}{2}}$ and $\mathfrak{B}(\mathsf{m}^{(n)}) \leq 3e^{-\frac{\mathrm{K}^v}{2}}$ for all $n \geq N$ for some $N \in \mathbb{N}$. Here we use the notation, $\mathsf{a}^{(n)} = \gamma^{(n)} D(p^{(n)}, \mathrm{K}^v) + \bar{\gamma}^{(n)}\mathsf{m}^{(n)}$. We assume $\mathrm{K}^v$ large enough so that for all $d$ we have

$$\frac{d(d-1)}{2} \mathfrak{B}(\mathsf{c}) \mathfrak{B}(\mathsf{b}^{(n)})^{d-2} \leq 1.$$

We put together everything done previously to bound the contributions to the density coming out of the variable nodes at the $(n + 1)$th iteration. To do this, we first use the check node analysis in Lemma 19 with incoming density given by $\mathsf{a}^{(n)}$. Then, using the variable node analysis of the previous section we obtain

$$\gamma^{(n+1)}p^{(n+1)} \leq e^{-\frac{\mathrm{K}^v}{2}} (d_r\xi\bar{\gamma}^{(n)} \mathfrak{B}(\mathsf{m}^{(n)}))^2$$
$$+ (d+1)(4d_r\gamma^{(n)}p^{(n)})^{\lfloor \frac{d}{2} \rfloor + 1}$$
$$+ d(4d_r\gamma^{(n)}p^{(n)})^{\lfloor \frac{d}{2} \rfloor} \mathfrak{B}(\mathsf{c})d_r\xi(\bar{\gamma}^{(n)} \mathfrak{B}(\mathsf{m}^{(n)})),$$
$$(22)$$

$$\bar{\gamma}^{(n+1)} \, \mathfrak{B}(\mathsf{m}^{(n+1)}) \leq (d_r \xi \bar{\gamma}^{(n)} \, \mathfrak{B}(\mathsf{m}^{(n)}))^2$$
$$+ 2d \, \mathfrak{B}(\mathsf{c}) d_r \xi (\bar{\gamma}^{(n)} \, \mathfrak{B}(\mathsf{m}^{(n)})) (d_r 4 \gamma^{(n)} p^{(n)})^{\lfloor \frac{d-2}{2} \rfloor} \mathfrak{B}(\mathsf{b}^{(n)})$$
$$+ \mathfrak{B}(\mathsf{c})(d_r 4 \gamma^{(n)} p^{(n)})^{\lfloor \frac{d}{2} \rfloor}$$

$$\tag{23}$$

To obtain the second inequality we use $\mathfrak{B}(\mathsf{b}^{(n)}) \leq 1$, where we assume $\mathrm{K}^v$ large enough so that $3 d_r e^{-\frac{\mathrm{K}^v}{2}} \leq 1$.

Now for any $\epsilon > 0$ we choose $\mathrm{K}^v$ large enough so that $(d_r 4 \gamma^{(n)} p^{(n)}) < 1$ and for all $d \geq 2$ we have

$$\epsilon \geq (d_r \xi)^2 \bar{\gamma}^{(n)} \, \mathfrak{B}(\mathsf{m}^{(n)}) + 2d \, \mathfrak{B}(\mathsf{c}) d_r \xi \, \mathfrak{B}(\mathsf{b}^{(n)}),$$
$$\epsilon \geq e^{-\frac{\mathrm{K}^v}{2}} \mathfrak{B}(\mathsf{c}) 4 d_r,$$
$$\epsilon \geq e^{-\frac{\mathrm{K}^v}{2}} 4 d_r (d+1) \Big( 1 + \mathfrak{B}(\mathsf{c}) d_r \xi (\bar{\gamma}^{(n)} \, \mathfrak{B}(\mathsf{m}^{(n)})) \Big),$$

which then yields

$$\begin{bmatrix} \bar{\gamma} \, \mathfrak{B}(\mathsf{m}) \\ e^{\frac{\mathrm{K}^v}{2}} \gamma p \end{bmatrix}^{(n+1)} \leq \epsilon \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \bar{\gamma} \, \mathfrak{B}(\mathsf{m}) \\ e^{\frac{\mathrm{K}^v}{2}} \gamma p \end{bmatrix}^{(n)}, \tag{24}$$

where $[\cdot]^{(n)}$ denotes the values at the $n$th iteration. We summarize our findings in the following.

*Theorem 21:* Consider an irregular ensemble with check regular degree $d_r$ and minimum variable node degree at least three. If a channel $\mathsf{c}$ is below the BP threshold then it is below the threshold for SatBP for $\mathrm{K}^v$ sufficiently large.

*Proof:* Assume the channel $\mathsf{c}$ is below the BP threshold. Let $x^*$ be the constant of Lemma 18. Under BP we have $\mathfrak{B}(T^{(\ell)}(\mathsf{c}, \Delta_0)) < x^*/2$ for some $\ell$ large enough. By Lemma 14 and Lemma 13 we have $\mathfrak{B}(S_{\mathrm{K}^v}^{(\ell)}(\mathsf{c}, \Delta_0)) \leq x^*$ for $\mathrm{K}^v$ large enough. By Lemma 18, and assuming $\mathrm{K}^v$ large enough, we have $\mathfrak{B}(S_{\mathrm{K}^v}^{(n)}(\mathsf{c}, \Delta_0)) \leq 3 e^{-\frac{\mathrm{K}^v}{2}}$ for all $n$ large enough. The stability analysis above then implies that $\lim_{n \to \infty} \mathfrak{E}(S_{\mathrm{K}^v}^{(n)}(\mathsf{c}, \Delta_0)) = 0$. ∎

## VI. BLOCK THRESHOLDS AND SPEED OF CONVERGENCE

Thresholds for iterative coding systems are usually *bit* thresholds. In some cases one can show that the iterative block error rate has the same threshold [19], [20]. For standard irregular ensembles it is sufficient that variable node degrees are at least three. The key observation for degree three and above is that below the bit threshold the bit error rate converges to zero doubly exponentially in iteration. One can maintain tree-like neighborhoods with blocklength growing exponentially in iteration and therefore the block error rate can be shown to converge to zero. In [19] it was shown that degree two variable nodes connected in an accumulate structure could be admitted while retaining the block threshold result provided an appropriate update schedule was adopted. The key idea there was that, by effectively updating a string of degree two updates in sequence for each iteration, one could achieve exponential decay in error probability with as large and exponent as required.

In this section we consider the impact of saturation on the block threshold. The stability analysis for ensembles with minimum variable node degree three shows exponential decay in iteration of bit error probability with arbitrarily large exponent. Consequently, we can show for a suitable ensemble that the block threshold coincides with the bit threshold. Nevertheless, saturation has a pronounced effect on stability and we observe this especially in the conditions required for doubly exponential convergence of the bit error probability. We show that doubly exponential convergence occurs for SatBP with minimum variable node degree five. With minimum variable node degree four doubly exponential convergence does not occur but can be recovered the addition of a single extra LLR magnitude and a two-tiered saturation. For minimum variable degree three doubly exponential convergence of the bit error rate can be recovered with a more radical modification of the decoding process (erase received values once the bit error rate is sufficiently small.)

Let us briefly review the standard block threshold arguments. For further details we refer to [19], [20]. Density evolution gives the bit error rate $P_b(\ell)$ as a function of iteration assuming tree-like neighborhoods up to iteration $\ell$. For block length $n$ the block error rate, assuming tree-like neighborhoods, is upper bounded by $n P_b(\ell)$. For the block error rate analysis we require that *all* computation trees are tree-like. This is accomplished through an expurgation or modification of the standard ensemble. The simplest approach, and the one we adopt, is to consider $n = n(\ell)$ large enough so that the fraction of variable nodes whose neighborhoods are not tree-like tends to zero as $\ell$ gets large. Then, we modify the code by declaring the associated bits as known and set to $0$. This lowers slightly the rate of the code and in effect modifies slightly the degree structure. The net effect is an improvement in bitwise performance. Asymptotically in large $\ell$ the modification is negligible so that full rate is recovered.

The basic calculation is as follows. Consider a computation tree associated to $\ell$ iterations. Let $\mathcal{M}_\ell$ denote the number of variable nodes in the computation tree. Let $n \gg \mathcal{M}_\ell$ denote the block length. It is not difficult to see that there exists a constant $\gamma$ independent of $\ell$ and $n$ such that the probability that the neighborhood is tree like is at least

$$(1 - \gamma \frac{\mathcal{M}_\ell}{n})^{\mathcal{M}_\ell} \geq (1 - \gamma \frac{\mathcal{M}_\ell^2}{n})$$

Now, we have a bound of the form $\mathcal{M}_\ell^2 \leq e^{M\ell}$ (where $M$ depends on the degree structure) and we choose $n = e^{N\ell}$ where $N > M$. Thus $N$ depends only on the degree structure of the code. It then follows that the fraction of variable nodes whose neighborhoods are not tree-like is tending to $0$ in $\ell$. To show that the block threshold equals the bit threshold it remains only to show that

$$\lim_{\ell \to \infty} e^{N\ell} P_b(\ell) = 0.$$

It is sufficient therefore to show that

$$\liminf_{\ell \to \infty} (-\ln P_b(\ell)) > N.$$

Let us consider $E(\ell) := \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \bar{\gamma} \, \mathfrak{B}(\mathsf{m}) \\ e^{\frac{\mathrm{K}^v}{2}} \gamma p \end{bmatrix}^{(\ell)}$. We clearly have

$$P_b(\ell) \leq E(\ell) = \bar{\gamma}^{(\ell)} \, \mathfrak{B}(\mathsf{m}^{(\ell)}) + e^{\frac{\mathrm{K}^v}{2}} \gamma^{(\ell)} p^{(\ell)}.$$

11

From the previous analysis we know that there exists an $\ell_0$ such that $E(\ell_0)$ is small. Recursing equation (24), we get $E(\ell + \ell_0) \leq (2\epsilon)^\ell E(\ell_0) = E(\ell_0) e^{-\ell \ln(1/(2\epsilon))}$. We can now make $\epsilon$ arbitrarily small by choosing $K^v$ large enough. Hence for sufficiently large $K^v$ we obtain

$$\liminf_{\ell \to \infty}(-\ln E(\ell)) > N$$

thus establishing the desired result.

### A. Variable nodes with Minimum Degree at least 5

In this section we show that SatBP does achieve doubly exponential convergence in $\ell$ of the error probability when the variable node degrees are at least five.

The rate of convergence depends largely on the variable node update. It is clear from (21) that, even with degree three, $\gamma' p'$ has quadratic dependence on $\gamma p$ and $\bar{\gamma} \mathfrak{B}(\mathsf{m})$. For doubly exponential convergence we can admit linear dependence of $\bar{\gamma} \mathfrak{B}(\mathsf{m})$ on $\gamma p$, but the dependence on $\bar{\gamma} \mathfrak{B}(\mathsf{m})$ must be of higher order. Let us make this more precise.

As before we assume $K^v$ and $N$ large enough so that for all $d$ and $n \geq N$ we have $\frac{d(d-1)}{2} \mathfrak{B}(\mathsf{c}) \mathfrak{B}(\mathsf{b}^{(n)})^{d-2} \leq 1$ and $(4 d_r \gamma^{(n)} p^{(n)}) < 1$. Then from (22) and (23), assuming $d \geq 4$, we get

$$e^{\frac{K^v}{2}} \gamma^{(n+1)} p^{n+1} \leq (d_r \xi \bar{\gamma}^{(n)} \mathfrak{B}(\mathsf{m}^{(n)}))^2$$
$$+ e^{-K^v}(4 d_r e^{\frac{K^v}{2}} \gamma^{(n)} p^{(n)})^3$$
$$+ e^{-\frac{K^v}{2}}(4 d_r e^{\frac{K^v}{2}} \gamma^{(n)} p^{(n)})^2 \mathfrak{B}(\mathsf{c}) d_r \xi (\bar{\gamma}^{(n)} \mathfrak{B}(\mathsf{m}^{(n)})),$$
$$(25)$$

$$\bar{\gamma}^{(n+1)} \mathfrak{B}(\mathsf{m})^{(n+1)} \leq (d_r \xi \bar{\gamma}^{(n)} \mathfrak{B}(\mathsf{m}^{(n)}))^2$$
$$+ 2d \mathfrak{B}(\mathsf{c}) d_r \xi (\bar{\gamma}^{(n)} \mathfrak{B}(\mathsf{m}^{(n)}) e^{-\frac{K^v}{2}} (d_r 4 e^{\frac{K^v}{2}} \gamma^{(n)} p^{(n)}) \mathfrak{B}(\mathsf{b}^{(n)})$$
$$+ e^{-K^v} \mathfrak{B}(\mathsf{c})(d_r 4 e^{\frac{K^v}{2}} \gamma^{(n)} p^{(n)})^2,$$
$$(26)$$

from which we easily obtain that for $K^v$ large enough we have

$$\bar{\gamma}^{(n+1)} \mathfrak{B}(\mathsf{m})^{(n+1)} + e^{\frac{K^v}{2}} \gamma^{(n+1)} p^{(n+1)} \leq$$
$$2(d_r \xi)^2 (\bar{\gamma}^{(n)} \mathfrak{B}(\mathsf{m})^{(n)} + e^{\frac{K^v}{2}} \gamma^{(n)} p^{(n)})^2,$$

which yields doubly exponential convergence in the iterations.

### B. Decoder Alteration for Degree Four

When $d = 3$ (degree four) the SatBP decoder does not yield doubly exponential stability convergence. The limiting effect arises in the variable node analysis from messages of type $(n_- = 0, n_{\mathsf{m}} = 1, n_+ = 2)$ which contribute a linear dependence of $\mathfrak{B}(\mathsf{m}')$ on $\mathfrak{B}(\mathsf{m})$. This occurs because $0 < 2K^p - (K^p + K^c) < K^v$. If the support of $\mathsf{m}$ were reduced to $[-\lambda K^v, \lambda K^v]$ where $2K^p - (\lambda K^v + K^c) > K^v$ then this term would be eliminated and doubly exponential convergence can be recovered.

Thus, for minimum degree four we consider a two step saturation at variable nodes where all messages with magnitude at least $K^v$ are saturated to $K^v$ and messages with magnitude between $\lambda K^v$ and $K^v$ are saturated to $\lambda K^v$. Hence, for this section we assume the inequality

$$2K^p - K^v \geq K^c + \lambda K^v.$$

We assume $\lambda \in (\frac{1}{2}, 1]$ and note that the above inequality then implies $K^c \leq (1 - \lambda) K^v$.

Note that an equivalent interpretation under scaling of the saturation levels is that we append an additional magnitude level to the SatBP decoder. Under this interpretation we identify $\lambda K^v$ with $K^v$ and $K^v$ with $\lambda^{-1} K^v$ where magnitudes above this level are saturated to $\lambda^{-1} K^v$. Under this interpretation the modification appears as an improvement on SatBP and, using this perspective, it is relatively easy to reproduce the results on the approximation of BP by the saturating decoder. Let us make this more precise. For notational purposes we will adhere to the original interpretation.

Let $\lfloor\!\lfloor \mathsf{a} \rfloor\!\rfloor_{\lambda, K}$ denote the double saturation of $\mathsf{a}$ and let $\lfloor\!\lfloor \mathsf{a} \rfloor\!\rfloor_{\lambda, K_{\text{sym}}}$ denote the symmetrized version. Let $S_{\lambda, K_{\text{sym}}}$ denote the corresponding one step density evolution update. We easily obtain the following generalization of Lemma 10

$$d(\mathsf{a}, \lfloor\!\lfloor \mathsf{a} \rfloor\!\rfloor_{\lambda, K_{\text{sym}}}) \leq d(\mathsf{a}, \lfloor \mathsf{a} \rfloor_{\lambda K_{\text{sym}}}) \leq 1 - \tanh(\lambda K/2),$$

where $\mathsf{a}$ is any symmetric $L$-density. It is not hard to see that we can also obtian the following generalization of Lemma 13,

$$d(T^{(\ell)}(\mathsf{c}, \Delta_0), S_{\lambda, K_{\text{sym}}}^{(\ell)}(\mathsf{c}, \Delta_0)) \leq 2 e^{-\lambda K + \ell \cdot \ln(2(d_l - 1)(d_r - 1))}.$$

The relationship between the symmetrized decoder and the non-symmetrized version as analyzed in in Lemma 14 remains essentially unchanged and we have that for any $0 < \epsilon < 1$ and $\ell \in \mathbb{N}$, there exists a $K^v$ large enough such that

$$\mathfrak{B}(S_{\lambda, K}^{(\ell)}(\mathsf{c}, \Delta_0)) \leq \frac{1}{1 - \epsilon} \mathfrak{B}(S_{\lambda, K_{\text{sym}}}^{(\ell)}(\mathsf{c}, \Delta_0)).$$

We can now focus our attention on the stability analysis. Let $\mathsf{a}$ be a density supported on $[-K^v, K^v]$. Then we have the two bounds,

$$\mathfrak{B}(\lfloor\!\lfloor \mathsf{a} \rfloor\!\rfloor_{\lambda, K^v}) \leq e^{\frac{K^v - \lambda K^v}{2}} \mathfrak{B}(\lfloor \mathsf{a} \rfloor_{K^v}), \qquad (27)$$

$$\mathfrak{B}(\lfloor\!\lfloor \mathsf{a} \rfloor\!\rfloor_{\lambda, K^v}) \leq \mathfrak{B}(\mathsf{a}) + e^{-\frac{\lambda K^v}{2}}. \qquad (28)$$

The first (multiplicative) inequality is new and will be used to establish doubly exponential convergence. Indeed, since $e^{\frac{K^v - \lambda K^v}{2}} \geq 1$, we have

$$\mathfrak{B}(\lfloor\!\lfloor \mathsf{a} \rfloor\!\rfloor_{\lambda, K^v}) \leq e^{\frac{K^v - \lambda K^v}{2}} e^{\frac{K^v}{2}} \int_{-\infty}^{-K^v} \mathsf{a}(x) dx$$
$$+ \int_{-K^v}^{\lambda K^v} e^{-\frac{x}{2}} \mathsf{a}(x) dx + e^{\frac{K^v - \lambda K^v}{2}} \int_{\lambda K^v}^{K^v} e^{-\frac{x}{2}} \mathsf{a}(x) dx$$
$$+ e^{\frac{K^v - \lambda K^v}{2}} e^{-\frac{K^v}{2}} \int_{K^v}^{\infty} \mathsf{a}(x) dx \leq e^{\frac{K^v - \lambda K^v}{2}} \mathfrak{B}(\lfloor \mathsf{a} \rfloor_{K^v}).$$

The second (additive) inequality allows us to reproduce the near stability analysis of Section V-B to obtain as in the derivation of 14 and 15 for the doubly saturated decoder the bounds

$$\mathfrak{B}(\mathsf{a}^{(n)}) \leq 3 e^{-\lambda K^v/2} \qquad (29)$$

$$\mathfrak{B}(\mathsf{b}^{(n)}) \leq 3\rho'(1) e^{-\lambda K^v/2}. \qquad (30)$$

which hold for $n \geq N$ (for some $N \in \mathbb{N}$) and $K^v$ large enough assuming the channel is below the BP threshold.

We assume that no additional saturation is performed at the check node so, in particular, Lemma 19 still applies. In

the variable node analysis we note that (21) still applies. The change in the analysis concerns the bound on $\bar{\gamma}'\mathsf{m}'$ in the variable node analysis. New considerations apply to the inner saturation of the density $\mathsf{m}'$. Further note that the incoming densities in to the variable nodes have support on $\pm\mathrm{K}^p\cup(-\lambda\mathrm{K}^v,\lambda\mathrm{K}^v)$. First we note the contribution from types with $n_\mathsf{m}\geq 2$. Let the notation $\lfloor\mathsf{a}\rceil^\circ_{\lambda,\mathrm{K}^v}$ denote the density on the support $[-\lambda\mathrm{K}^v,\lambda\mathrm{K}^v]$ which is equivalent, in this case, to the support on $(-\mathrm{K}^v,\mathrm{K}^v)$. Using analysis in the previous section and the inequality (27) we get,

$$\mathfrak{B}\left(\left\lfloor\sum_{k=0}^{d-2}\binom{d}{k}\gamma^k\bar{\gamma}^{d-k}\mathsf{c}\circledast\mathsf{D}^{\circledast k}\circledast\mathsf{m}^{\circledast(d-k)}\right\rceil^\circ_{\lambda,\mathrm{K}^v}\right)\leq$$
$$e^{\frac{\mathrm{K}^v-\lambda\mathrm{K}^v}{2}}\frac{d(d-1)}{2}(\bar{\gamma}\,\mathfrak{B}(\mathsf{m}))^2\,\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{b})^{d-2}.$$

Now we consider the contribution from types with $n_\mathsf{m}=1$. A type $(n_-,1,n_+)$ can have a non-zero contribution to $\mathsf{m}'$ only if the interval centered on $(n_+-n_-)\mathrm{K}^p$ of width $2(\mathrm{K}^\mathsf{c}+\lambda\mathrm{K}^v)$ intersects $(-\mathrm{K}^v,\mathrm{K}^v)$. Since we assume $2\mathrm{K}^p\geq\mathrm{K}^\mathsf{c}+\mathrm{K}^v+\lambda\mathrm{K}^v$ and $\mathrm{K}^p\leq\mathrm{K}^v$ we obtain

$$\mathfrak{B}(\lfloor\mathsf{c}\circledast\mathsf{m}\circledast\mathsf{D}^{d-1}\rceil^\circ_{\mathrm{K}^v})\leq$$
$$\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{m})\begin{cases}\sum_{j=\frac{d-2}{2}}^{\frac{d}{2}}\binom{d-1}{j}q^{d-1-j}\tilde{q}^j & d\text{ even}\\\binom{d-1}{\frac{d-1}{2}}q^{\frac{d-1}{2}}\tilde{q}^{\frac{d-1}{2}} & d\text{ odd},\end{cases}$$

where recall that $q=e^{\frac{\mathrm{K}^p}{2}}p$ and $\tilde{q}=e^{-\frac{\mathrm{K}^p}{2}}\bar{p}$. Again, combining the above with (27), we obtain

$$\mathfrak{B}\left(\lfloor\mathsf{c}\circledast\mathsf{m}\circledast\mathsf{D}^{d-1}\rceil^\circ_{\lambda,\mathrm{K}^v}\right)\leq$$
$$e^{\frac{\mathrm{K}^v-\lambda\mathrm{K}^v}{2}}\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{m})\begin{cases}\binom{d-1}{\frac{d}{2}}(p\bar{p})^{\frac{d-2}{2}}\,\mathfrak{B}(D(p,\mathrm{K}^p)) & d\text{ even}\\\binom{d-1}{\frac{d-1}{2}}(p\bar{p})^{\frac{d-1}{2}} & d\text{ odd}.\end{cases}$$

Finally we consider the contribution from types with $n_\mathsf{m}=0$. A type $(n_-,0,n_+)$ will have a non-zero contribution to $\mathsf{m}'$ only if the interval centered on $(n_+-n_-)\mathrm{K}^p$ of width $2\mathrm{K}^\mathsf{c}$ intersects $(-\mathrm{K}^v,\mathrm{K}^v)$. Hence we obtain

$$\mathfrak{B}(\lfloor\mathsf{c}\circledast\mathsf{D}^d\rceil^\circ_{\mathrm{K}^v})\leq\mathfrak{B}(\mathsf{c})\begin{cases}\binom{d}{\frac{d}{2}}q^{\frac{d}{2}}\tilde{q}^{\frac{d}{2}} & d\text{ even}\\\sum_{j=\frac{d-1}{2}}^{\frac{d+1}{2}}\binom{d}{j}q^{d-j}\tilde{q}^j & d\text{ odd}\end{cases}$$

which gives

$$\mathfrak{B}(\lfloor\mathsf{c}\circledast\mathsf{D}^d\rceil^\circ_{\lambda,\mathrm{K}^v})$$
$$\leq\mathfrak{B}(\mathsf{c})\begin{cases}\binom{d}{\frac{d}{2}}(p\bar{p})^{\frac{d}{2}} & d\text{ even}\\e^{\frac{\mathrm{K}^v-\lambda\mathrm{K}^v}{2}}\binom{d}{\frac{d-1}{2}}(p\bar{p})^{\frac{d-1}{2}}\,\mathfrak{B}(D(p,\mathrm{K}^p)) & d\text{ odd}\end{cases}$$

Since $\lambda>\frac{1}{2}$ we can assume for $d\geq 3$ and for $\mathrm{K}^v$ large enough that,

$$e^{\frac{\mathrm{K}^v-\lambda\mathrm{K}^v}{2}}\frac{d(d-1)}{2}\,\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{b})^{d-2}$$
$$\overset{(30)}{\leq}e^{-\frac{2\lambda-1}{2}\mathrm{K}^v}\frac{d(d-1)}{2}3\rho'(1)\,\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{b})^{d-3}$$
$$\leq 1.$$

Also,

$$\mathfrak{B}(\lfloor\mathsf{c}\circledast\mathsf{m}\circledast\mathsf{D}^{d-1}\rceil^\circ_{\lambda,\mathrm{K}^v})\leq$$

$$e^{\frac{\mathrm{K}^v-\lambda\mathrm{K}^v}{2}}\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{m})\begin{cases}\binom{d-1}{\frac{d}{2}}(p\bar{p})^{\frac{d-2}{2}}\,\mathfrak{B}(D(p,\mathrm{K}^p)) & d\text{ even}\\\binom{d-1}{\frac{d-1}{2}}(p\bar{p})^{\frac{d-1}{2}} & d\text{ odd}\end{cases}$$
$$\leq e^{\frac{\mathrm{K}^v-\lambda\mathrm{K}^v}{2}}\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{m})\begin{cases}(4p)^{\frac{d-2}{2}} & d\text{ even}\\(4p)^{\frac{d-1}{2}} & d\text{ odd}.\end{cases}$$

Finally,

$$\mathfrak{B}(\lfloor\mathsf{c}\circledast\mathsf{D}^d\rceil^\circ_{\lambda,\mathrm{K}^v})\leq\mathfrak{B}(\mathsf{c})\begin{cases}(4p\bar{p})^{\frac{d}{2}} & d\text{ even}\\e^{-\frac{2\lambda-1}{2}\mathrm{K}^v}3\rho'(1)(4p)^{\frac{d-1}{2}} & d\text{ odd}\end{cases}$$
$$\leq(4p)^{\lfloor\frac{d-1}{2}\rfloor}.$$

To get the final bound on $\gamma'\mathfrak{B}(\mathsf{m}')$ we need to multiply the above bounds by $d\bar{\gamma}\gamma^{d-1}$ when $n_\mathsf{m}=1$ and by $\gamma^d$ when $n_\mathsf{m}=0$. For $\mathrm{K}^v$ large enough we can make $4\gamma p\leq 1$. Thus we get,

$$\bar{\gamma}'\,\mathfrak{B}(\mathsf{m}')\leq\mathfrak{B}(\mathsf{c})\Big((\bar{\gamma}\,\mathfrak{B}(\mathsf{m}))^2+de^{\frac{\mathrm{K}^v-\lambda\mathrm{K}^v}{2}}\,\mathfrak{B}(\mathsf{c})(\bar{\gamma}\,\mathfrak{B}(\mathsf{m}))(4\gamma p)$$
$$+(4\gamma p)\Big)$$

Assuming $d\geq 3$ we also have from the previous analysis,

$$\gamma p'\leq e^{-\frac{\mathrm{K}^v}{2}}\frac{d(d-1)}{2}(\bar{\gamma}\,\mathfrak{B}(\mathsf{m}))^2\,\mathfrak{B}(\mathsf{c})\,\mathfrak{B}(\mathsf{b})^{d-2}$$
$$+(d+1)(4\gamma p)^{\lfloor\frac{d}{2}\rfloor+1}+d(4\gamma p)^{\lfloor\frac{d}{2}\rfloor}\,\mathfrak{B}(\mathsf{c})(\bar{\gamma}\,\mathfrak{B}(\mathsf{m})).$$

Thus we now obtain quadratic dependence and hence doubly exponential convergence even when minimum variable node degree is four.

### C. Decoder Alteration for Degree Three

In this section we will show that when the minimum variable node degree is 3, we can still have doubly exponential convergence of the bit error rate which implies an exponential (in blocklength) convergence of the block error rate with a decoder alteration. In this case, however, we require an iteration dependent alteration of the decoder. We alter the decoder only after the error rate is sufficiently small. Hence, for the analysis we assume operation in the near stability region. More precisely, we have $\mathfrak{B}(\mathsf{a})\leq 3e^{-\mathrm{K}^v/2}$, where $\mathsf{a}$ is the outgoing density at the variable nodes. Since $\mathsf{a}=\gamma D(p,\mathrm{K}^v)+\bar{\gamma}\mathsf{m}$, we further have $\bar{\gamma}\,\mathfrak{B}(\mathsf{m})\leq 3e^{-\mathrm{K}^v/2}$ and $\gamma p\leq 3e^{-\mathrm{K}^v/2}$.

We note that the previous technique of saturation at two levels does not yield the quadratic dependence we seek for the term $\mathfrak{B}(\mathsf{m}')$. Indeed, any incoming density having the type $(n_-=0,n_\mathsf{m}=1,n_+=1)$ will always contribute to the outgoing density of type $\mathsf{m}'$, implying linear dependence of $\mathfrak{B}(\mathsf{m}')$ on $\mathfrak{B}(\mathsf{m})$. To show doubly exponentially fast convergence of the bit error rate, we modify the decoder as follows. After the messages have become reasonably good, i.e., we are in the near stability region, we erase the channel information. The intuition is that at this point the extrinsic information is good enough for successful decoding. Then for every incoming message we make a hard-decision to either $+1$ or $-1$ based on the sign of its LLR value. The decoding algorithm then proceeds in a manner similar to the erasure decoder [1]. Let us explain this in more detail.

The decoder has now three messages $\{-1,0,+1\}$. At the variable node side, there is an erasure message on the outgoing

edge if and only if all the incoming messages are erasures or there is exactly one $+1$ and $-1$ message. The outgoing edge carries a $-1$ message if and only if all incoming messages are $-1$ or one message is an erasure and the other is $-1$. At the check node side, the outgoing message is an erasure if at least one incoming message is an erasure, else the outgoing message is the product of the incoming messages. We can now write the density evolution equation analysis for this decoder as follows. Let $x_\ell$ and $y_\ell$ represent the probability of the messages $0$ and $-1$, respectively, coming out of the variable node. Also, let $w_\ell$ and $z_\ell$ represent the probability of the messages $0$ and $-1$, coming out of the check node respectively. Since we are in the near stability region, it is not hard to see that $x_0 \leq \bar{\gamma}\,\mathfrak{B}(\mathsf{m}) \leq c e^{-K^v/2}$ and $y_0 \leq \mathfrak{B}(\mathsf{a}) \leq c e^{-K^v/2}$. Indeed, $y_0 = \int_{x<0} \mathsf{a}(x)dx \leq \int_{x \leq 0} \mathsf{a}(x)e^{-x/2}dx \leq \mathfrak{B}(\mathsf{a})$. From the decoder rules we immediately get,

$$x_\ell \overset{(a)}{\leq} w_\ell^2 + z_\ell,$$
$$y_\ell = z_\ell^2 + w_\ell z_\ell,$$
$$w_\ell = 1 - (1 - x_{\ell-1})^{d_r-1} \leq (d_r - 1)x_{\ell-1},$$
$$z_\ell \overset{(b)}{\leq} 1 - (1 - y_{\ell-1})^{d_r-1} \leq (d_r - 1)y_{\ell-1},$$

where $d_r$ is the check node degree. To obtain $(a)$ we simply upper bound the probability of message with value $+1$ by $1$. At the check node side, the outgoing message is $-1$ if there are odd number of incoming messages that are $-1$. This implies that at least one incoming message must be $-1$ and hence we obtain inequality $(b)$.

Combining the four inequalities above, it is not hard to see that $x_\ell + y_\ell \leq C(x_{\ell-2} + y_{\ell-2})^2$ for some positive constant $C$. This implies $x_\ell + y_\ell \leq (Ax_0)^{2^{n/2}}$, where $A$ is some positive constant and $n$ is the number of iterations of the erasure decoder. Hence we obtain the doubly exponential convergence.

## VII. THRESHOLD FOR THE SATBP DECODER AND CHANNELS WITH INFINITE SUPPORT

Consider a channel family, $\mathrm{BMS}(\mathsf{h})$, ordered by $\mathsf{h}$ and let $\mathsf{h}^{\mathrm{BP}}(\lambda, \rho)$ denote the BP threshold when transmitting over this channel family using a $(\lambda, \rho)$ ensemble. Also, a priori the channel has support on $(-\infty, \infty)$.

Let us describe the analysis of the SatBP decoder in this case. Consider transmission over a channel with $L$-density $\mathsf{c}$. From the previous analysis we have that the channel support must be finite for stability of the perfect decoding fixed point when we use the SatBP decoder. As a result, we saturate the channel $\mathsf{c}$ to a value $K^{\mathsf{c}} \leq 2K^p - K^v$ before we feed it to the SatBP decoder. The value $K^p$ is defined in section V-C. Thus we consider transmission over a channel $\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}}$.

For the purpose of analysis we also consider the corresponding symmetric channel, achieved via flipping as explained previously. Denote it by $\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}$. We have the following lemma.

*Lemma 22 (Stability Condition for Sym. Sat. Channels):* Consider transmission over a general BMS channel $\mathsf{c}$ using $(\lambda, \rho)$ ensemble. Let $\mathsf{c} \in \mathrm{BMS}(\mathsf{h})$ be such that it satisfies the following stability condition,

$$(\lambda'(0)\rho'(1))(\mathfrak{B}(\mathsf{c}) + 2e^{-K^{\mathsf{c}}/2}) < 1.$$

Then, the full BP decoder is successful when transmitting over the symmetric channel $\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}$. Furthermore, the loss in capacity is also bounded by $\frac{2}{\ln 2}e^{-K^{\mathsf{c}}/2}$.

*Proof:* We bound the Wasserstein distance between the DE with channel $\mathsf{c}$ and DE with channel $\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}$ as follows,

$$d(T^{(\ell)}(\mathsf{c}, \Delta_0), T^{(\ell)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0)) =$$
$$d(T(\mathsf{c}, T^{(\ell-1)}(\mathsf{c}, \Delta_0)), T(\mathsf{c}, T^{(\ell-1)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0)))+$$
$$d(T(\mathsf{c}, T^{(\ell-1)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0)), T(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, T^{(\ell-1)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0)))$$
$$\overset{(vi,viii),\ \mathrm{Lem.\ 13\ in\ [18]}}{\leq} \alpha_\ell d(T^{(\ell-1)}(\mathsf{c}, \Delta_0), T^{(\ell-1)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0))$$
$$+ 2d(\mathsf{c}, \lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}})$$
$$= \alpha_\ell d(T^{(\ell-1)}(\mathsf{c}, \Delta_0), T^{(\ell-1)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0)) + 2(1 - \tanh(\frac{K^{\mathsf{c}}}{2})),$$

where

$$\alpha_\ell = 2(d_l - 1)\sum_{j=1}^{d_r-1}(1 - \mathfrak{B}^2(\mathsf{a}))^{\frac{d_r-1-j}{2}}(1 - \mathfrak{B}^2(\mathsf{b}))^{\frac{j-1}{2}},$$

where $\mathsf{a} = T^{(\ell-1)}(\mathsf{c}, \Delta_0)$ and $\mathsf{b} = T^{(\ell-1)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0)$ and $d_l$ and $d_r$ correspond to the average variable node and check node degrees. Following the same steps as in the proof of lemma 13 we get

$$\mathfrak{B}(T^{(\ell)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0))$$
$$\leq \mathfrak{B}(T^{(\ell)}(\mathsf{c}, \Delta_0)) + 2\sqrt{2}e^{\frac{-K^{\mathsf{c}} + \ell \cdot \ln(2(d_l-1)(d_r-1))}{2}}.$$

Thus, for any $\xi > 0$, we can choose $K^{\mathsf{c}}$ large enough, such that $\mathfrak{B}(T^{(\ell)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0)) \leq \xi$ for all $\ell \geq \ell_0$. Here $\ell_0$ is such that $\mathfrak{B}(T^{(\ell_0)}(\mathsf{c}, \Delta_0)) \leq \xi/2$.

Let us denote $x_\ell = \mathfrak{B}(T^{(\ell)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0))$. Using extremes of information combining [1] we get $x_\ell \leq \mathfrak{B}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}})\lambda(1 - \rho(1 - x_{\ell-1}))$. Expanding around zero, we get $x_\ell \leq \mathfrak{B}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}})\lambda'(0)\rho'(1)x_{\ell-1} + O(x_{\ell-1}^2)$. Using the hypothesis of the lemma, lemma 10 and (ix), Lem. 13 in [18] we have, $\mathfrak{B}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}})\lambda'(0)\rho'(1) < 1$. Hence, there exists $\eta > 0$ such that $\mathfrak{B}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}})\lambda'(0)\rho'(1) + \eta < 1$. From above we know that there exists $\ell$ (and consequently $K^{\mathsf{c}}$ large enough) such that the second order term $O(x_{\ell-1}^2)$ is upper bounded by $\eta x_{\ell-1}$. Thus we get $x_\ell \leq (\mathfrak{B}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}})\lambda'(0)\rho'(1) + \eta)x_{\ell-1} < x_{\ell-1}$. Thus $x_\ell \to 0$ as $\ell \to \infty$ and we get the lemma.

The loss in capacity is bounded by using the Wasserstein distance. Thus $d(\mathsf{c}, \lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}) \leq 1 - \tanh(K^{\mathsf{c}}/2)$ implies $H(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}) \leq H(\mathsf{c}) + \frac{2}{\ln 2}e^{-K^{\mathsf{c}}/2}$. Above we have used $1 - \tanh(K^{\mathsf{c}}/2) \leq 2e^{-K^{\mathsf{c}}}$ and (ix), Lem. 13 in [18]. Thus, $1 - H(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}) \geq 1 - H(\mathsf{c}) - \frac{2}{\ln 2}e^{-K^{\mathsf{c}}/2}$. ∎

From the above lemma and the analysis in section IV we get[2] $\mathfrak{B}(T^{(\ell)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}}, \Delta_0)) \leq \frac{1}{1-\epsilon}\mathfrak{B}(T^{(\ell)}(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}, \Delta_0))$, for any $0 < \epsilon < 1$. Since $\mathsf{c} \prec \lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}} \prec \lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}}$, we have $H(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}}) \leq H(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}_{\mathrm{sym}}})$ which implies that $1 - H(\lfloor \mathsf{c} \rfloor_{K^{\mathsf{c}}}) \geq 1 - H(\mathsf{c}) - \frac{2}{\ln 2}e^{-K^{\mathsf{c}}/2}$.

---

[2]Recall that we associated a uniform random variable to each variable node which were used for the flipping operations for outgoing messages from the variable node side. For the present case, we can associate a random variable to each channel input which is used for the flipping operation for symmetrizing the saturated channel. These two operations are independent of each other. In section IV the event $A_{K^v}$ now corresponds to the event that there are no flips at both the variable node and channel input. This probability will be lower bounded by $1 - 2e^{-K^v}|V(\mathsf{T})|$.

Note that the stability analysis of section V does not rely on the symmetry of the channel. The symmetry allows us to that Battacharyya parameter of the channel is less than one, which is then used to show bounds. In the present case, since $\mathfrak{B}(\lfloor c \rfloor_{K^c}) \leq \mathfrak{B}(c) + e^{-K^c/2}$ we can proceed with the stability analysis as before and conclude that the SatBP decoder is successful when we first truncate the channel to a large but finite support. Furthermore, this truncation causes minimal loss in the maximum number of information bits that can be transmitted. Finally, we can also say that for any channel $c \prec c^{BP}$ such that $\mathfrak{B}(c) < \mathfrak{B}(c^{BP}) - 2e^{-K^c/2}$, the SatBP decoder is successful over the truncated channel. Thus, the loss in the BP threshold is also upper bounded by $Ce^{-K^c/2}$ for some constant $C$. Note that the threshold for the SatBP decoder is now defined with respect to the fixed point with Battacharyya parameter equal to $e^{-K^v/2}$.

## VIII. Conclusions and Outlook

In this paper we perform perturbation analysis of the standard LDPC code ensemble and BP decoder combination. Specifically, we show that saturating the messages arising in the BP decoding process affects the final success of the decoder. For general irregular LDPC code ensembles with minimum variable node degree three, we show that the saturation of the messages still allows for successful decoding as long as the saturation level $K^v$ is large enough. More precisely, whenever the channel is below the BP threshold, then there exists a saturation value $K^v$, which is large enough but finite, such that the SatBP decoder is also below its threshold. The stability of the SatBP decoder requires the support of the channel to be finite. In the case of channels with infinite support, we show that by saturating the channel first to a large enough value, we sacrifice little in terms of capacity. Then, on the saturated channel, the SatBP decoder is successful. Thus there is minimal sacrifice in the BP threshold of the LDPC code ensemble when we consider the SatBP decoder.

When the minimum variable node degree is two the saturated decoding system fails to have stability of perfect decoding. We show that the perfect decoding fixed point (the delta function at $K^v$) cannot be a stable fixed point of DE for the SatBP decoder unless the channel is the erasure channel. The key issue is that a density update at a degree two node variable nodes is convolution with the channel density. Repeated $k$ times, this involves to convolution of the channel density with itself $k$ times. In general this is equivalent to a channel density with support width $k$ times wider than the original channel. If the incoming density is saturated then for $k$ large enough a positive error probability is unavoidable. If the code structure (e.g. protograph designs) ensures that the number of successive degree two node updates in the density evolution is bounded, then the expansion $k$ is bounded and one can again recover stability with large enough saturation. Essentially, what is required is that each degree two variable node subgraph connected component (asymptotically a tree) have bounded size.

To give a more detailed indication of how this can work we consider the min-sum decoder and show that perfect decoding can be invariant even in the presence of degree two variable nodes. Let the maximum component size be denoted by $A$. For an edge $e$ connected to a degree two variable node let $2L_e + 1$ denote the maximum path length to the edge of the connected component. Note that $L_e + 1 \leq A$. To show invariance of a perfect decoding we assume $2(K^v - AK^c) - K^c \geq K^v$. Assume in some iteration that the following hold,

- The incoming message to a degree two variable node with edges $e_1, e_2$ on edge $e_i$ is at least $K^v - L_{e_i}K^c$.
- Incoming messages on a degree three or higher variable node are at least $K^v - AK^c$.

It is easy to check that this implies perfect decoding. Proceeding to the next iteration we obtain,

- The outgoing message on a degree two variable node on edge $e_2$ is at least $K^v - (L_{e_1} + 1)K^c$ (and vice-versa for $e_1$.)
- Outgoing messages on a degree three or higher variable node are at least $K^v$.

Now consider the subsequent incoming messages to the variable nodes. The minimum outgoing message from the previous iteration is at least $K^v - AK^c$ so incoming messages to a degree three or higher variable node are at least $K^v - AK^c$. Consider edge $e_1$ attached to a degree two variable node. The longest path, not traversing $e_1$, from its neighboring check node to a leaf check of the degree two connect component has edge length at most $2L_{e_1}$. Hence the minimum incoming message to the neighbor check node not from $e_1$ is $K^v - L_{e_1}K^c$. The minimum incoming message on edge $e_1$ to the degree two variable node is therefore at least $K^v - L_{e_1}K^c$. Thus, under the stated assumptions the above perfect decoding conditions are invariant.

*Future Directions:*

To complete the story of the analysis of the BP decoder under practical considerations, it would be nice to have the analysis of the quantized BP decoder. Thus, the messages are only allowed to take certain values on the real line. Every message is quantized to a bin and only the bin value is passed around. For the ease of analysis one can assume a uniformly quantized message space. It is not hard to see that such a quantized BP decoder is symmetric. Thus the standard DE analysis is applicable to the quantized BP decoder. A clear next step would be to see if the analysis performed for the SatBP decoder goes through for the quantized BP decoder. If yes, then it would be nice to see a unified perturbation analysis of saturated and quantized messages.

A nice side-effect of the analysis done above is that when there are degree three variable nodes present in the LDPC code, it is perhaps better to erase the channel information at those bits completely (after enough iterations are performed) to allow faster convergence to the correct codeword. This sheds some light on the practical design of BP decoders under saturation of messages. Could we glean similar lessons for practical decoder design when we consider the saturated and quantized BP decoder?

Another research direction would be to quantify the saturation and quantization levels in terms of gap to capacity.

Specifically, what should be the scaling of the saturation and quantization value when we backoff, say, $\delta$ from the BP capacity, $h^{BP}$. It seems intuitive that as we backoff more from $h^{BP}$ we should be able to attain the same error rate with smaller values of the saturation level and larger levels of quantization. In other words, as the gap to capacity increases, we should require lesser number of bits in the binary representation of the messages to get the desired error rate.

## APPENDIX A
### BATTACHARRYA PARAMETER INEQUALITY – LEMMA 17

We require the following inequality

*Lemma 23:* Let $p_1, ..., p_k$ each lie in $[0,1]$. Then

$$\frac{1 - \prod_{i=1}^{k}(1-2p_i)}{2} \leq \sum_{i=1}^{k} p_i$$

*Proof:* We have equality when $p_i = 0$ for each $i$. Differentiating the left hand side with respect to $p_j$ we obtain $\prod_{\{i \in [1:k] \setminus j\}}(1-2p_i)$ which has magnitude at most 1 and differentiating the left hand side with respect to $p_j$ we obtain 1. The inequality therefore follows by integration. ∎

The following generalizes Lemma 17.

*Lemma 24:* Let $D_1, D_2, ...D_k$ be L-densities of the form $D_i = D(p_i, K)$ and let $a_1, ..., a_{d-k}$ be L-densities. We do not assume that any of these densities are symmetric. Let b denote the density emerging from a check node update when the incoming densities are $D_1, ..., D_k, a_1, ..., a_{d-k}$, then

$$\mathfrak{B}(b) \leq \left(1 + \sum_{i=1}^{k}(e^{K/2}\mathfrak{B}(D_i) - 1)\right)\left(\sum_{i=j}^{d-k} \mathfrak{B}(a_j)\right).$$

(This holds even if $k = 0$ in which case we have only the second factor.) This generalizes a result from [21].

*Proof:* By averaging, we see that it is sufficient to prove the lemma for the case $a_i = D(q_i, z_i)$. With this assumption the outgoing message is of the form $b = D(s, r)$ where

$$s = \frac{1 - (\prod_{i=1}^{k}(1-2p_i))(\prod_{j=1}^{d-k}(1-2q_j))}{2},$$

and we have $r \leq \min\{K, q_1, ..., q_{d-k}\}$ and $e^{-r/2} \leq ke^{-K/2} + \sum_{j=1}^{d-k} e^{-q_i/2}$. We have $\mathfrak{B}(b) = se^{r/2} + (1-s)e^{-r/2}$.

Define

$$P = \frac{1 - \prod_{i=1}^{k}(1-2p_i)}{2}, \quad Q = \frac{1 - \prod_{j=1}^{d-k}(1-2q_j)}{2}$$

Then we have

$$1 - s = PQ + (1-P)(1-Q).$$

We claim the inequality

$$\mathfrak{B}(b) \leq (Pe^{K} + (1-P))(Qe^{r/2} + (1-Q)e^{-r/2}).$$

The claim follows from collecting terms and noting $e^{K}e^{r/2} \geq e^{-r/2}$, which is obvious, and $e^{K}e^{-r/2} \geq e^{r/2}$, which follows from $K^v \geq r$.

We now apply Lemma 23 to the left factor to obtain

$$Pe^{K} + (1-P) = 1 + P(e^{K} - 1)$$

$$\leq 1 + (\sum_{i=1}^{k} p_i)(e^{K} - 1)$$

$$= 1 + \sum_{i=1}^{k}(e^{K/2}(p_i e^{K/2} + (1-p_i)e^{-K/2}) - 1)$$

$$= 1 + \sum_{i=1}^{k}(e^{K/2}\mathfrak{B}(D_i) - 1).$$

Using $q_j \leq r$ and $\sum_{j=1}^{d-k} e^{-q_j/2} \geq e^{-r/2}$ and applying Lemma 23 to the right factor we obtain

$$Qe^{r/2} + (1-Q)e^{-r/2} = e^{-r/2} + Q(2\sinh(r/2))$$

$$\leq e^{-r/2} + (\sum_{j=1}^{d-k} q_j)(2\sinh(r/2))$$

$$\leq \sum_{j=1}^{d-k} e^{-q_j/2} + \sum_{j=1}^{d-k} q_j(2\sinh(q_j/2))$$

$$= \sum_{i=1}^{d-k} \mathfrak{B}(a_j).$$

∎

## REFERENCES

[1] T. Richardson and R. Urbanke, *Modern Coding Theory*. Cambridge University Press, 2008.

[2] X. Zhang and P. Siegel, "Will the real error floor please stand up?" in *Signal Processing and Communications (SPCOM), 2012 International Conference on*, 2012, pp. 1–5.

[3] B. Butler and P. Siegel, "Error floor approximation for ldpc codes in the awgn channel," in *Communication, Control, and Computing (Allerton), 2011 49th Annual Allerton Conference on*, 2011, pp. 204–211.

[4] C. Schlegel and S. Zhang, "On the dynamics of the error floor behavior in (regular) ldpc codes," *Information Theory, IEEE Transactions on*, vol. 56, no. 7, pp. 3248–3264, 2010.

[5] S. Zhang and C. Schlegel, "Controlling the error floor in ldpc decoding," *Communications, IEEE Transactions on*, vol. 61, no. 9, pp. 3566–3575, 2013.

[6] X. Zhang and P. Siegel, "Quantized min-sum decoders with low error floor for ldpc codes," in *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, 2012, pp. 2871–2875.

[7] ——, "Quantized iterative message passing decoders with low error floor for ldpc codes," pp. 1–14, 2013.

[8] B. Vasic, D. V. Nguyen, and S. K. Chilappagari, "Failures and error-floors of iterative decoders," *Channel Coding: Theory, Algorithms, and Applications, Academic Press Library in Mobile and Wireless, Communications, Elsevier, New York*, 2014.

[9] J. Wang, T. Courtade, H. Shankar, and R. Wesel, "Soft information for ldpc decoding in flash: Mutual-information optimized quantization," in *Global Telecommunications Conference (GLOBECOM 2011), 2011 IEEE*, 2011, pp. 1–6.

[10] N. Kanistras, I. Tsatsaragkos, I. Paraskevakos, A. Mahdi, and V. Paliouras, "Impact of llr saturation and quantization on ldpc min-sum decoders," in *Signal Processing Systems (SIPS), 2010 IEEE Workshop on*, 2010, pp. 410–415.

[11] N. Kanistras, I. Tsatsaragkos, and V. Paliouras, "Propagation of llr saturation and quantization error in ldpc min-sum iterative decoding," in *Signal Processing Systems (SiPS), 2012 IEEE Workshop on*, 2012, pp. 276–281.

[12] Y. Wu, L. Davis, and R. Calderbank, "On the capacity of the discrete-time channel with uniform output quantization," in *Information Theory, 2009. ISIT 2009. IEEE International Symposium on*, 2009, pp. 2194–2198.

[13] S. Kudekar, T. Richardson, and R. L. Urbanke, "Wave-like solutions of general one-dimensional spatially coupled systems," *CoRR*, vol. abs/1208.5273, 2012.

[14] T. Richardson and R. Urbanke, "The capacity of low-density parity check codes under message-passing decoding," *IEEE Trans. Inform. Theory*, vol. 47, no. 2, pp. 599–618, Feb. 2001.

[15] G. Hanoch and H. Levy, "The efficiency analysis of choices involving risk," *The Review of Economic Studies*, vol. 36, pp. 335–346, 1969.

[16] S. Kudekar, T. Richardson, and R. Urbanke, "Existence and Uniqueness of GEXIT curves via the Wasserstein Metric," in *Proc. of the IEEE Inform. Theory Workshop*, Paraty, Brazil, 2011.

[17] C. Villani, *Optimal transport, Old and New*. Springer, 2009, vol. 338.

[18] S. Kudekar, T. Richardson, and R. L. Urbanke, "Spatially coupled ensembles universally achieve capacity under belief propagation," *CoRR*, vol. abs/1201.2999, 2012.

[19] H. Jin and T. Richardson, "Block error iterative decoding capacity for ldpc codes," in *Information Theory, 2005. ISIT 2005. Proceedings. International Symposium on*, Sept 2005, pp. 52–56.

[20] M. Lentmaier, D. Truhachev, K. Zigangirov, and D. Costello, "An analysis of the block error probability performance of iterative decoding," *Information Theory, IEEE Transactions on*, vol. 51, no. 11, pp. 3834–3855, Nov 2005.

[21] K. Bhattad, V. Rathi, and R. Urbanke, "Degree optimization and stability condition for the min-sum decoderl," in *Proc. of the IEEE Inform. Theory Workshop*, 2007, conference, pp. 190–195.