# TOWARDS CONVERGENT APPROXIMATE MESSAGE PASSING BY ALTERNATING CONSTRAINED MINIMIZATION OF BETHE FREE ENERGY

*Christo Kurisummoottil Thomas[1] and Zilu Zhao and Dirk Slock[2]*

[1]Wireless@VT, Bradley Dept of ECE, Virginia Tech, Arlington, VA, USA,
[2]Communication Systems Department, EURECOM, France
Emails: christokt@vt.edu, zilu.zhao@eurecom.fr, Dirk.Slock@eurecom.fr

## ABSTRACT

Generalized Approximate Message Passing (GAMP) allows for Bayesian inference in linear models with non-identically independently distributed (n.i.i.d.) priors and n.i.i.d. measurements of the linear mixture outputs. It represents an efficient technique for approximate inference, which becomes accurate when both rows and columns of the measurement matrix can be treated as sets of independent vectors and both dimensions become large. It has been shown that the fixed points of GAMP correspond to the extrema of a large system limit of the Bethe Free Energy (LSL-BFE), which represents a meaningful approximation optimization criterion regardless of whether the measurement matrix exhibits the independence properties. However, the convergence of (G)AMP can be problematic for certain measurement matrices. In this paper, we revisit the GAMP algorithm by applying a simplified version of the Alternating Direction Method of Multipliers (ADMM) to minimizing the LSL-BFE. We show convergence of the mean and variance subsystems in AMBGAMP and in the Gaussian case, convergence of mean and LSL variance to the Minimum Mean Squared Error (MMSE) quantities.

## 1. INTRODUCTION

In the Gaussian noise case, a sparse signal vector $\boldsymbol{x}$ can be recovered using the signal model: $\boldsymbol{y} = \mathbf{A}\,\boldsymbol{x} + \boldsymbol{v}$, where $\boldsymbol{y}$ is the observed data, $\mathbf{A}$ is the known measurement or sensing matrix of dimension $M \times N$, typically with $M < N$. In the sparse model case, $\boldsymbol{x}$ contains only $K$ non-zero entries, with $K < M < N$. Sparse Bayesian Learning (SBL) is a Bayesian inference algorithm proposed by [1] and [2]. SBL is based on a hierarchical prior on the sparse coefficients $\boldsymbol{x}$, inducing sparsity by choosing priors for the hyperparameters that make a portion of the estimate $\widehat{\boldsymbol{x}}$ zero. The Linear Minimum Mean Squared Error (LMMSE) estimation step in SBL at each iteration involves matrix inversion, which makes it computationally complex [3].

The Approximate Message Passing (AMP) algorithm has been introduced to reduce the complexity of Belief Propagation, from $2MN$ to $M + N$ messages. Generalized AMP (GAMP) allows non-Gaussian priors and measurement processes. But convergence of (G)AMP can be problematic for some matrices $\mathbf{A}$. Existing converging AMP versions introduced so far: 1) adding the Alternating Direction Method of Multipliers (ADMM) [4] leading to a higher complexity ADMM-GAMP, 2) exploiting part of the singular value decomposition (SVD) of the measurement matrix in Vector AMP (VAMP) [5], [6] or esp. Unitarily Transformed UT-AMP [7] (but which do not allow to handle n.i.i.d. priors conveniently), 3) introducing damping [8], but with typically difficult to determine damping requirements.

### 1.1. Contributions of this paper

• We propose a convergent version of GAMP, AMBGAMP, which applies alternating minimization to an augmented Lagrangian of a large system limit of the Bethe free Energy (BFE). AMBGAMP can be interpreted as applying a simplified ADMM to the BFE, with a constrained Lagrange multiplier parameterization for the mean constraint, and a quadratic optimization subproblem being solved by a gradient update with line search. The ADMM is complemented with a fixed

point iteration for the variance constraint.
• We show that AMBGAMP converges to the LMMSE estimate in the Gaussian case.
• We provide a convergence analysis of the variance subsystem.
• We show that in the Gaussian case the variances converge to the optimal MSE values in the large system limit.
• We provide a convergence analysis of the mean subsystem.

## 2. SYSTEM MODEL

The data model considered in GAMP is essentially a linear mixing model represented by

$$\mathbf{z} = \mathbf{A}\,\boldsymbol{x}\,,\ p_{\boldsymbol{x}}(\boldsymbol{x})\,,\ p_{\boldsymbol{y}|\mathbf{z}}(\boldsymbol{y}|\mathbf{z}) \tag{1}$$

with (possibly non) identically independently distributed (n.i.i.d.) prior $p_{\boldsymbol{x}}(\boldsymbol{x}) = \prod_{i=1}^{N} p_{x_i}(x_i)$ and n.i.i.d. measurements $p_{\boldsymbol{y}|\mathbf{z}}(\boldsymbol{y}|\mathbf{z}) = \prod_{k=1}^{M} p_{y_k|z_k}(y_k|z_k)$. In the case of Gaussian measurement noise, we have $\boldsymbol{y} = \mathbf{z} + \boldsymbol{v}$ with independent zero-mean n.i.i.d. Gaussian noise $\boldsymbol{v}$ with variance vector $\boldsymbol{\sigma}_v^2 = [\sigma_{v1}^2, \cdots, \sigma_{vM}^2]^T$. We shall also consider the case of a zero mean Gaussian prior for $\boldsymbol{x}$ with variances denoted as $\sigma_{xi}^2$. We represent the vector $\boldsymbol{\sigma}_x^2 = [\sigma_{x1}^2, \cdots \sigma_{xN}^2]^T$. In Bayesian estimation, we are interested in the posterior, which is given by

$$p_{\boldsymbol{x},\mathbf{z}|\boldsymbol{y}}(\boldsymbol{x}, \mathbf{z}|\boldsymbol{y}) = \frac{e^{-\sum_{i=1}^{N} f_{x_i}(x_i) - \sum_{k=1}^{M} f_{z_k}(z_k)}}{Z(\boldsymbol{y})}\ \mathbb{1}_{\{\mathbf{z} = \mathbf{A}\boldsymbol{x}\}}, \tag{2}$$

with the negative log-likelihoods defined as $f_{x_i}(x_i) = -\ln p_{x_i}(x_i)$, $f_{z_k}(z_k) = -\ln p_{y_k|z_k}(y_k|z_k)$, where the equality in case of $f_{z_k}(z_k)$ is up to constants that may depend on $\boldsymbol{y}$ (and which are absorbed in the normalization constant $Z(\boldsymbol{y})$). The problem in Bayesian estimation is the computation of this constant $Z(\boldsymbol{y})$ and of the posterior means and variances. Belief propagation is a message passing technique that allows to compute the posterior marginals. However, due to loops in the factor graph, loopy belief propagation may have convergence issues and is furthermore still relatively complex. GAMP is an approximate belief propagation technique which is motivated by asymptotic considerations in which the rows and columns of the measurement matrix $\mathbf{A}$ are considered as random and independent, in which case GAMP can actually produce the correct posterior marginals. In any case, GAMP computes a separable approximate posterior of the form

$$q_{\boldsymbol{x},\mathbf{z}}(\boldsymbol{x}, \mathbf{z}) = q_{\boldsymbol{x}}(\boldsymbol{x})\, q_{\mathbf{z}}(\mathbf{z}) = \prod_{i=1}^{N} q_{x_i}(x_i)\ \prod_{k=1}^{M} q_{z_k}(z_k), \tag{3}$$

in which the dependence on $\boldsymbol{y}$ has been omitted. The GAMP algorithm [9], [8] appears in the table for Algorithm 1. We only consider here Sum-Product GAMP (for MMSE estimation, as opposed to Max-Sum GAMP for MAP estimation).

## 3. PROPOSED AMBGAMP

The abbreviation AMB stands for ACM-LSL-BFE, which stands for Alternating Constrained Minimization of the LSL of the BFE. AMBGAMP employs most of the same updates as GAMP, but GAMP does not apply a strict alternating minimization (block coordinate descent) principle, particularly in the presence of constraints. Previous

**Algorithm 1** GAMP

---
**Require:** $\boldsymbol{y}$, $\mathbf{A}$, $\mathbf{S} = \mathbf{A}.\mathbf{A}$, $f_{\boldsymbol{x}}(\boldsymbol{x})$, $f_{\mathbf{z}}(\mathbf{z})$
1: Initialize: $t = 0$, $\widehat{\boldsymbol{x}}^t$, $\boldsymbol{\tau}_x^t$, $\mathbf{s}^{-1} = \mathbf{0}$
2: **repeat**
3:     [Output node update]
4:     $\boldsymbol{\tau}_p^t = \mathbf{S}\,\boldsymbol{\tau}_x^t$
5:     $\boldsymbol{p}^t = \mathbf{A}\,\widehat{\boldsymbol{x}}^t - \mathbf{s}^{t-1}.\boldsymbol{\tau}_p^t$
6:     $\widehat{\mathbf{z}}^t = \mathbb{E}(\mathbf{z}|\boldsymbol{p}^t, \boldsymbol{\tau}_p^t)$
7:     $\boldsymbol{\tau}_z^t = \text{var}(\mathbf{z}|\boldsymbol{p}^t, \boldsymbol{\tau}_p^t)$
8:     $\mathbf{s}^t = (\widehat{\mathbf{z}}^t - \boldsymbol{p}^t)./\boldsymbol{\tau}_p^t$
9:     $\boldsymbol{\tau}_s^t = (\mathbf{1} - \boldsymbol{\tau}_z^t./\boldsymbol{\tau}_p^t)./\boldsymbol{\tau}_p^t$
10:     [Input node update]
11:     $\boldsymbol{\tau}_r^t = \mathbf{1}./(\mathbf{S}^T\boldsymbol{\tau}_s^t)$
12:     $\boldsymbol{r}^t = \widehat{\boldsymbol{x}}^t + \boldsymbol{\tau}_r^t.\mathbf{A}^T\mathbf{s}^t$
13:     $\widehat{\boldsymbol{x}}^{t+1} = \mathbb{E}(\boldsymbol{x}|\mathbf{r}^t, \boldsymbol{\tau}_r^t)$
14:     $\boldsymbol{\tau}_x^{t+1} = \text{var}(\boldsymbol{x}|\mathbf{r}^t, \boldsymbol{\tau}_r^t)$
15: **until** Convergence

---

work [10] has demonstrated that any fixed point of the GAMP algorithm is a critical point of the following constrained minimization of a LSL of the BFE (see also [8] and references therein):

$$\min_{q_{\boldsymbol{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p} J_{LSL-BFE}(q_{\boldsymbol{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p)$$
$$s.t. \ \mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) = \mathbf{A}\,\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}}) \tag{4}$$
$$\boldsymbol{\tau}_p = \mathbf{S}\,\text{var}(\boldsymbol{x}|q_{\boldsymbol{x}}),$$

where the LSL BFE is given by

$$J_{LBFE}(q_{\boldsymbol{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p) = D(q_{\boldsymbol{x}}||e^{-f_{\boldsymbol{x}}}) + D(q_{\mathbf{z}}||e^{-f_{\mathbf{z}}})$$
$$+ H_G(q_{\mathbf{z}}, \boldsymbol{\tau}_p), \ \text{with} \ H_G(q_{\mathbf{z}}, \boldsymbol{\tau}_p) = \frac{1}{2}\sum_{k=1}^{M}\left[\frac{\text{var}(z_k|q_{z_k})}{\tau_{p_k}} + \ln(2\pi\tau_{p_k})\right] \tag{5}$$

and where $D(q||p) = \mathbb{E}_q(\ln(\frac{q}{p}))$ is the KLD and $H_G(q_{\mathbf{z}}, \boldsymbol{\tau}_p)$ is a sum of a KLD and an entropy of Gaussians with identical means but different variances. The LSL BFE optimization problem (5) can be reformulated with the following augmented Lagrangian

$$\min_{q_{\boldsymbol{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p, \boldsymbol{u}} \ \max_{\mathbf{s}, \boldsymbol{\tau}_s} L(q_{\boldsymbol{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p, \boldsymbol{u}, \mathbf{s}, \boldsymbol{\tau}_s) \ \text{with}$$
$$L = D(q_{\boldsymbol{x}}||e^{-f_{\boldsymbol{x}}}) + D(q_{\mathbf{z}}||e^{-f_{\mathbf{z}}}) + H_G(q_{\mathbf{z}}, \boldsymbol{\tau}_p)$$
$$+ \mathbf{s}^T(\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) - \mathbf{A}\,\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}})) - \frac{1}{2}\boldsymbol{\tau}_s^T(\boldsymbol{\tau}_p - \mathbf{S}\,\text{var}(\boldsymbol{x}|q_{\boldsymbol{x}})) \tag{6}$$
$$+ \frac{1}{2}\|\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}}) - \boldsymbol{u}\|_{\boldsymbol{\tau}_r}^2 + \frac{1}{2}\|\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) - \mathbf{A}\,\boldsymbol{u}\|_{\boldsymbol{\tau}_p}^2,$$

where $\mathbf{s}, \boldsymbol{\tau}_s$ are Lagrange multipliers, and $\boldsymbol{\tau}_r = \mathbf{1}./(\mathbf{S}^T\boldsymbol{\tau}_s)$ is just a short-hand notation for a quantity that depends on $\boldsymbol{\tau}_s$. We also use the notations: $\|\boldsymbol{u}\|_{\boldsymbol{\tau}}^2 = \sum_i u_i^2/\tau_i$, element-wise multiplication as in $\mathbf{s}.\boldsymbol{\tau}$ and element-wise division as in $\mathbf{1}./\boldsymbol{\tau}$, and $\mathbf{1}$ is a vector of ones. In [11], [12], a careful updating schedule was considered with partial optimization steps on subsets of primal and dual variables. However, that approach is not guaranteed to converge in general. In [13] we continued to consider an alternating optimization approach in which the schedule is less critical and some of the optimizations are reduced to gradient updates. The resulting algorithm can be considered an extended and generalized version of the ADMM algorithm (extended: there are more than two primal variable groups, generalized: the quadratic augmentation term does not exactly correspond to the linear (mean) constraint). However, there is an alternative point of view, based on [4], where a double mean constraint was introduced leading to the ADMM-GAMP augmented Lagrangian

$$\min_{q_{\boldsymbol{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p, \boldsymbol{u}} \ \max_{\mathbf{q}, \mathbf{s}, \boldsymbol{\tau}_s} L_A(q_{\boldsymbol{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p, \boldsymbol{u}, \mathbf{q}, \mathbf{s}, \boldsymbol{\tau}_s) \ \text{with}$$
$$L_A = D(q_{\boldsymbol{x}}||e^{-f_{\boldsymbol{x}}}) - \frac{1}{2}\boldsymbol{\tau}_s^T(\boldsymbol{\tau}_p - \mathbf{S}\,\text{var}(\boldsymbol{x}|q_{\boldsymbol{x}})) + D(q_{\mathbf{z}}||e^{-f_{\mathbf{z}}})$$
$$+ H_G(q_{\mathbf{z}}, \boldsymbol{\tau}_p) + \mathbf{q}^T(\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}}) - \boldsymbol{u})) + \mathbf{s}^T(\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) - \mathbf{A}\,\boldsymbol{u})) \tag{7}$$
$$+ \frac{1}{2}\|\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}}) - \boldsymbol{u}\|_{\boldsymbol{\tau}_r}^2 + \frac{1}{2}\|\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) - \mathbf{A}\,\boldsymbol{u}\|_{\boldsymbol{\tau}_p}^2,$$

For ADMM, the first two terms are the cost function for $q_{\boldsymbol{x}}$, the next two terms constitute the cost function for $q_{\mathbf{z}}$. The two groups of primal variables are $\{q_{\boldsymbol{x}}, q_{\mathbf{z}}\}$ and $\boldsymbol{u}$ (and the optimization of $L_A$ is decoupled between $q_{\boldsymbol{x}}, q_{\mathbf{z}}$). The two linear constraints together constitute a single extended set of linear constraints with extended Lagrange

multiplier $[\mathbf{q}^T\mathbf{s}^T]^T$. The appropriately weighted quadratic augmentation terms correspond exactly to the set of linear constraints. The optimization in [4] is organized with the usual ADMM algorithm alternating between minimizations over the two groups of primal variables, followed by the ADMM specific Lagrange multiplier update. The optimization over the remaining variable $\tau_p, \tau_s$ is then performed in an outer loop. We show here (by the variance subsystem convergence analysis) that this organization in two levels is not necessary. Furthermore, there is a redundancy between the linear and quadratic constraint terms in (7). Indeed, if we impose the constrained Lagrange multiplier structure $\mathbf{q}^T = -\mathbf{s}^T\mathbf{A}$, then we obtain the proposed $L$ in (6). This is constrained enough since the Lagrange multiplier $\mathbf{s}$ will lead to $\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) = \mathbf{A}\,\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}})$, in which case the quadratic augmentation terms are minimized by $\boldsymbol{u} = \mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}})$ and disappear. However, constraining $\mathbf{q}^T = -\mathbf{s}^T\mathbf{A}$ leads to a deviation from the strict ADMM structure and requires separate convergence analysis, which we provide here.

At iteration $t$ we propose the following updating sequence

$$\{\boldsymbol{u}^t\} = \arg\min_{\boldsymbol{u}} L(q_{\boldsymbol{x}}^{t-1}, q_{\mathbf{z}}^{t-1}, \boldsymbol{\tau}_p^{t-1}, \boldsymbol{u}, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1}) \tag{8}$$

$$\{q_{\boldsymbol{x}}^t\} = \arg\min_{q_{\boldsymbol{x}}} L(q_{\boldsymbol{x}}, q_{\mathbf{z}}^{t-1}, \boldsymbol{\tau}_p^{t-1}, \boldsymbol{u}^t, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1}) \tag{9}$$

$$\{q_{\mathbf{z}}^t\} = \arg\min_{q_{\mathbf{z}}} L(q_{\boldsymbol{x}}^t, q_{\mathbf{z}}, \boldsymbol{\tau}_p^t, \boldsymbol{u}^t, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1}) \tag{10}$$

$$\{\mathbf{s}^t\} = \arg\max_{\mathbf{s}} L(q_{\boldsymbol{x}}^t, q_{\mathbf{z}}^t, \boldsymbol{\tau}_p^t, \boldsymbol{u}^t, \mathbf{s}, \boldsymbol{\tau}_s^{t-1}) \tag{11}$$

$$\{\boldsymbol{\tau}_p^t, \boldsymbol{\tau}_s^t\} = \arg\min_{\boldsymbol{\tau}_p} \max_{\boldsymbol{\tau}_s} L(q_{\boldsymbol{x}}^t, q_{\mathbf{z}}^t, \boldsymbol{\tau}_p, \boldsymbol{u}^t, \mathbf{s}^t, \boldsymbol{\tau}_s) \tag{12}$$

The result appears in Algorithm 2.

### 3.1. Update of $\boldsymbol{u}$

To update $\boldsymbol{u}$, we use a gradient descent method with line search to optimize the step-size. From (6), (8), we get

$$L(q_{\boldsymbol{x}}^{t-1}, q_{\mathbf{z}}^{t-1}, \boldsymbol{\tau}_p^{t-1}, \boldsymbol{u}, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1})$$
$$= \frac{1}{2}\|\widehat{\boldsymbol{x}}^{t-1} - \boldsymbol{u}\|_{\boldsymbol{\tau}_r^{t-1}}^2 + \frac{1}{2}\|\widehat{\mathbf{z}}^{t-1} - \mathbf{A}\,\boldsymbol{u}\|_{\boldsymbol{\tau}_p^{t-1}}^2 + const. \tag{13}$$

where $const.$ denotes constants w.r.t. $\boldsymbol{u}$. The minimizing update can be obtained as
$$\boldsymbol{u}^t = \boldsymbol{u}^{t-1} - \eta^t\,\mathbf{g}^t \tag{14}$$
with gradient $\mathbf{g}^t = \mathbf{g}^t(\boldsymbol{u}^{t-1})$ where

$$\mathbf{g}^t(\boldsymbol{u}) = \nabla_{\boldsymbol{u}} L(q_{\boldsymbol{x}}^{t-1}, q_{\mathbf{z}}^{t-1}, \boldsymbol{\tau}_p^{t-1}, \boldsymbol{u}, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1})$$
$$= -\mathbf{A}^T((\widehat{\mathbf{z}}^{t-1} - \mathbf{A}\boldsymbol{u})./\boldsymbol{\tau}_p^{t-1}) - (\widehat{\boldsymbol{x}}^{t-1} - \boldsymbol{u})./\boldsymbol{\tau}_r^{t-1} \tag{15}$$
$$= \mathbf{g}^t(\mathbf{0}) + \mathcal{H}^t\,\boldsymbol{u}, \quad \mathcal{H}^t = \mathbf{D}(\mathbf{1}./\boldsymbol{\tau}_r^{t-1}) + \mathbf{A}^T\mathbf{D}(\mathbf{1}./\boldsymbol{\tau}_p^{t-1})\mathbf{A}$$

where $\mathbf{D}(\boldsymbol{\tau})$ denotes a diagonal matrix with diagonal elements $\boldsymbol{\tau}$. The step-size $\eta^t$ gets optimized for maximum descent :

$$\frac{\partial L(q_{\boldsymbol{x}}^{t-1}, q_{\mathbf{z}}^{t-1}, \boldsymbol{\tau}_p^{t-1}, \boldsymbol{u}^t, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1})}{\partial \eta^t} = 0$$
$$\Rightarrow \eta^t = \|\mathbf{g}^t\|^2/\mathbf{g}^{t\,T}\mathcal{H}^t\mathbf{g}^t. \tag{16}$$

### 3.2. Update of $q_{\boldsymbol{x}}$

For the update of $q_{\boldsymbol{x}}$, consider the relevant terms in the augmented Lagrangian (and remember that $\mathbf{1}./\boldsymbol{\tau}_r^{t-1} = \mathbf{S}^T\boldsymbol{\tau}_s^{t-1}$)

$$L(q_{\boldsymbol{x}}, q_{\mathbf{z}}^{t-1}, \boldsymbol{\tau}_p^{t-1}, \boldsymbol{u}^t, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1})$$
$$= D(q_{\boldsymbol{x}}||e^{-f_{\boldsymbol{x}}}) - \mathbf{s}^{t-1\,T}\mathbf{A}\,\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}})$$
$$+ \frac{1}{2}\boldsymbol{\tau}_s^{t-1\,T}\mathbf{S}\,\text{var}(\boldsymbol{x}|q_{\boldsymbol{x}}) + \frac{1}{2}\|\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}}) - \boldsymbol{u}^t\|_{\boldsymbol{\tau}_r^{t-1}}^2 + const.$$
$$= D(q_{\boldsymbol{x}}||e^{-f_{\boldsymbol{x}}}) + \frac{1}{2}(\mathbf{1}./\boldsymbol{\tau}_r^{t-1})^T\,\mathbb{E}(\boldsymbol{x}.\boldsymbol{x}|q_{\boldsymbol{x}}) \tag{17}$$
$$- \mathbf{s}^{t-1\,T}\mathbf{A}\,\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}}) - (\boldsymbol{u}^t./\boldsymbol{\tau}_r^{t-1}))^T\,\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}}) + const.$$
$$= D(q_{\boldsymbol{x}}||e^{-f_{\boldsymbol{x}}}) + \frac{1}{2}(\mathbf{1}./\boldsymbol{\tau}_r^{t-1})^T\,\mathbb{E}(\boldsymbol{x}.\boldsymbol{x}|q_{\boldsymbol{x}})$$
$$- (\boldsymbol{u}^t + \boldsymbol{\tau}_r^{t-1}.\mathbf{A}^T\mathbf{s}^{t-1})^T(\mathbb{E}(\boldsymbol{x}|q_{\boldsymbol{x}})./\boldsymbol{\tau}_r^{t-1}) + const.$$
$$= D(q_{\boldsymbol{x}}||e^{-f_{\boldsymbol{x}}}) + \frac{1}{2}\,\mathbb{E}(\|\boldsymbol{x} - \mathbf{r}^t\|_{\boldsymbol{\tau}_r^t}^2|q_{\boldsymbol{x}}) + const.$$

where $const.$ denotes constants w.r.t. $\boldsymbol{x}$, and $\mathbf{r}^t = \boldsymbol{u}^t + \boldsymbol{\tau}_r^{t-1}.\mathbf{A}^T \mathbf{s}^{t-1}$. The Lagrangian in (17) is separable. We get per component

$$\min_{q_{x_i}} D(q_{x_i}||g_{x_i}^t/Z_{x_i}^t) \Rightarrow q_{x_i}^t = g_{x_i}^t/Z_{x_i}^t, \; Z_{x_i}^t = \int g_{x_i}^t(x_i)\, dx_i,$$
$$-\ln g_{x_i}^t(x_i) = f_{x_i}(x_i) + \frac{1}{2\tau_{r_i}^t}[(x_k - r_i^t)^2 - r_i^{t\,2}].$$
(18)

The partition function $Z_{x_i}^t$ acts as cumulant generating function:

$$\tau_{r_i}^t \frac{\partial \ln Z_{x_i}^t}{\partial r_i^t} = \mathbb{E}(x_i|q_{x_i}^t) = \mathbb{E}(x_i|r_i^t, \tau_{r_i}^t) = \widehat{x}_i^t$$
$$(\tau_{r_i}^t)^2 \frac{\partial^2 \ln Z_{x_i}^t}{\partial r_i^{t\,2}} = \mathrm{var}(x_i|r_i^t, \tau_{r_i}^t) = \tau_{x_i}^t.$$
(19)

In the Gaussian prior case, we get a Gaussian posterior $q_{\boldsymbol{x}}^t$ with

$$\mathbf{1}./\boldsymbol{\tau}_x^t = \mathbf{1}./\boldsymbol{\tau}_r^{t-1} + \mathbf{1}./\boldsymbol{\sigma}_x^2, \; \widehat{\boldsymbol{x}}^t = \boldsymbol{\tau}_x^t.(\mathbf{r}^t./\boldsymbol{\tau}_r^{t-1}).$$
(20)

### 3.3. Update of $\{q_{\mathbf{z}}\}$
The relevant terms in the augmented Lagrangian are

$$L(q_{\boldsymbol{x}}^t, q_{\mathbf{z}}, \boldsymbol{\tau}_p^{t-1}, \boldsymbol{u}^t, \mathbf{s}^{t-1}, \boldsymbol{\tau}_s^{t-1})$$
$$= D(q_{\mathbf{z}}||e^{-f_{\mathbf{z}}}) + \frac{1}{2}\mathrm{var}(\mathbf{z}|q_{\mathbf{z}})./\boldsymbol{\tau}_p^{t-1}$$
$$+\mathbf{s}^{t-1\,T}\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) + \frac{1}{2}\|\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}) - \mathbf{A}\,\boldsymbol{u}^t\|_{\boldsymbol{\tau}_p^{t-1}}^2 + const.$$
$$= D(q_{\mathbf{z}}||e^{-f_{\mathbf{z}}}) + \frac{1}{2}\mathbb{E}(\mathbf{z}^T\mathbf{z}|q_{\mathbf{z}})./\boldsymbol{\tau}_p^{t-1}$$
$$-(\mathbb{E}(\mathbf{z}|q_{\mathbf{z}}))^T((\mathbf{A}\,\boldsymbol{u}^t)./\boldsymbol{\tau}_p^{t-1} - \mathbf{s}^{t-1}) + const.$$
$$= D(q_{\mathbf{z}}||e^{-f_{\mathbf{z}}}) + \frac{1}{2}\mathbb{E}(\|\mathbf{z} - \boldsymbol{p}^t\|_{\boldsymbol{\tau}_p^{t-1}}^2|q_{\mathbf{z}}) + const.$$
(21)

where $const.$ denotes constants w.r.t. $\mathbf{z}$ and $\boldsymbol{p}^t = \mathbf{A}\,\boldsymbol{u}^t - \mathbf{s}^{t-1}.\boldsymbol{\tau}_p^{t-1}$. The Lagrangian in (21) is again separable. We get per component

$$\min_{q_{z_k}} D(q_{z_k}||g_{z_k}^t/Z_{z_k}^t) \Rightarrow q_{z_k}^t = g_{z_k}^t/Z_{z_k}^t$$
$$Z_{z_k}^t = \int g_{z_k}^t(z_k)\, dz_k, \; -\ln g_{z_k}^t(z_k) =$$
$$f_{z_k}(z_k) + \frac{1}{2\tau_{p_k}^{t-1}}[(z_k - p_k^t)^2 - (p_k^t)^2].$$
(22)

The partition function $Z_{z_k}^t$ acts again as cumulant generating function:

$$-\frac{\partial \ln Z_{z_k}^t}{\partial s_k^{t-1}} = \mathbb{E}(z_k|q_{z_k}^t) = \mathbb{E}(z_k|p_k^t, \tau_{p_k}^{t-1}) = \widehat{z}_k^t$$
$$\frac{\partial^2 \ln Z_{z_k}^t}{\partial s_k^{t-1\,2}} = \mathrm{var}(z_k|p_k^t, \tau_{p_k}^{t-1}) = \tau_{z_k}^t$$
$$-\frac{\partial^3 \ln Z_{z_k}^t}{\partial s_k^{t-1\,3}} = \mathbb{E}((z_k - \mathbb{E}\,z_k)^3|q_{z_k}^t).$$
(23)

The case of Gaussian noise leads again to a Gaussian posterior $q_{\mathbf{z}}$ with

$$\mathbf{1}./\boldsymbol{\tau}_z^t = \mathbf{1}./\boldsymbol{\tau}_p^{t-1} + \mathbf{1}./\boldsymbol{\sigma}_v^2, \; \widehat{\mathbf{z}}^t = \boldsymbol{\tau}_z^t.(\boldsymbol{y}./\boldsymbol{\sigma}_v^2 + \boldsymbol{p}^t./\boldsymbol{\tau}_p^{t-1}).$$
(24)

### 3.4. Update of $\{\mathbf{s}\}$ (ADMM style)
Although the quadratic augmentation terms in the Lagrangian do not correspond exactly to a weighted quadratic version of the linear mean constraint, due to the introduction of the auxiliary variable $\boldsymbol{u}$ which streamlines the derivation of the updates of $q_{\boldsymbol{x}}$ and $q_{\mathbf{z}}$, nevertheless an ADMM style update of the mean constraint Lagrange multiplier $\mathbf{s}$ is possible. Indeed, the terms in (21) that contains $\mathbf{s}$ or $\widehat{\mathbf{z}}$ are

$$\widehat{\mathbf{z}}^T((\frac{1}{2}\widehat{\mathbf{z}} - \boldsymbol{p}^t)./\boldsymbol{\tau}_p^{t-1}) = \widehat{\mathbf{z}}^T(\mathbf{s}^{t-1} + (\frac{1}{2}\widehat{\mathbf{z}} - \mathbf{A}\boldsymbol{u}^t)./\boldsymbol{\tau}_p^{t-1}) \quad (25)$$

Taking the gradient w.r.t. $\widehat{\mathbf{z}}$ (as part of the optimization over $q_{\mathbf{z}}$) leads to the RHS of

$$\mathbf{s}^t = \mathbf{s}^{t-1} + (\widehat{\mathbf{z}}^t - \mathbf{A}\boldsymbol{u}^t)./\boldsymbol{\tau}_p^{t-1}.$$
(26)

Hence, if we use this update for $\mathbf{s}$, then (25) reduces to $\widehat{\mathbf{z}}^T\mathbf{s}^t$, as if the quadratic augmentation terms have disappeared! This is the main characteristic of the Lagrange multiplier update in ADMM, which corresponds to a gradient ascent with a particular choice of (diagonal) step-size.

---

**Algorithm 2** AMBGAMP
**Require:** $\boldsymbol{y}, \mathbf{A}, \mathbf{S} = \mathbf{A}.\mathbf{A}, f_{\boldsymbol{x}}(\boldsymbol{x}), f_{\mathbf{z}}(\mathbf{z})$
1: Initialize: $t = 0, \boldsymbol{u}^0 = \mathbf{0}, \widehat{\boldsymbol{x}}^0 = \mathbf{0}, \widehat{\mathbf{z}}^0 = \mathbf{0}, \mathbf{s}^0 = \mathbf{0}, \boldsymbol{\tau}_r^0 = \mathbf{1}, \boldsymbol{\tau}_p^0 = \mathbf{1}$
2: **repeat** (t=1,2,...)
3: $\quad \boldsymbol{u}^t = \boldsymbol{u}^{t-1} - \eta^t\,\mathbf{g}^t$, with $\mathbf{g}^t, \eta^t$ from (15), (16)
4: $\quad$ [Input node update]
5: $\quad \mathbf{r}^t = \boldsymbol{u}^t + \boldsymbol{\tau}_r^{t-1}.(\mathbf{A}^T\mathbf{s}^{t-1})$
6: $\quad \widehat{\boldsymbol{x}}^t = \mathbb{E}(\boldsymbol{x}|\mathbf{r}^t, \boldsymbol{\tau}_r^{t-1}),$    Gaussian $p_{\boldsymbol{x}}: \mathbf{1}./\boldsymbol{\tau}_x^t = \mathbf{1}./\boldsymbol{\tau}_r^{t-1} + \mathbf{1}./\boldsymbol{\sigma}_x^2$
7: $\quad \boldsymbol{\tau}_x^t = \mathrm{var}(\boldsymbol{x}|\mathbf{r}^t, \boldsymbol{\tau}_r^{t-1}),$    Gaussian $p_{\boldsymbol{x}}: \widehat{\boldsymbol{x}}^t = \boldsymbol{\tau}_x^t.(\mathbf{r}^t./\boldsymbol{\tau}_r^{t-1})$
8: $\quad \boldsymbol{\tau}_p^t = \mathbf{S}\,\boldsymbol{\tau}_x^t$
9: $\quad$ [Output node update]
10: $\quad \boldsymbol{p}^t = \mathbf{A}\,\boldsymbol{u}^t - \mathbf{s}^{t-1}.\boldsymbol{\tau}_p^t$
11: $\quad \widehat{\mathbf{z}}^t = \mathbb{E}(\mathbf{z}|\boldsymbol{p}^t, \boldsymbol{\tau}_p^t),$    Gaussian $p_{\boldsymbol{y}|\mathbf{z}}: \mathbf{1}./\boldsymbol{\tau}_z^t = \mathbf{1}./\boldsymbol{\tau}_p^t + \mathbf{1}./\boldsymbol{\sigma}_v^2$
12: $\quad \boldsymbol{\tau}_z^t = \mathrm{var}(\mathbf{z}|\boldsymbol{p}^t, \boldsymbol{\tau}_p^t),$    Gaussian $p_{\boldsymbol{y}|\mathbf{z}}: \widehat{\mathbf{z}}^t = \boldsymbol{\tau}_z^t.(\boldsymbol{y}./\boldsymbol{\sigma}_v^2 + \boldsymbol{p}^t./\boldsymbol{\tau}_p^t)$
13: $\quad \mathbf{s}^t = \mathbf{s}^{t-1} + (\widehat{\mathbf{z}}^t - \mathbf{A}\boldsymbol{u}^t)./\boldsymbol{\tau}_p^t$
14: $\quad \boldsymbol{\tau}_s^t = (\mathbf{1} - \boldsymbol{\tau}_z^t./\boldsymbol{\tau}_p^t)./\boldsymbol{\tau}_p^t,$    Gaussian $p_{\boldsymbol{y}|\mathbf{z}}: \boldsymbol{\tau}_s^t = \mathbf{1}./(\boldsymbol{\sigma}_v^2 + \boldsymbol{\tau}_p^t)$
15: $\quad \boldsymbol{\tau}_r^t = \mathbf{1}./(\mathbf{S}^T\boldsymbol{\tau}_s^t)$
16: **until** Convergence

---

### 3.5. Update of $\{\boldsymbol{\tau}_p, \boldsymbol{\tau}_s\}$
In [11], [12], the carefully chosen updating schedule made the quadratic augmentation terms inactive when updating $\{\boldsymbol{\tau}_p, \boldsymbol{\tau}_s\}$. Here these terms only become inactive at convergence. Nevertheless, these terms only play an active role for the means and not for the variances. Hence, we shall ignore them here. Thus, the terms of interest in (6) for (12) are

$$L(q_{\boldsymbol{x}}^t, q_{\mathbf{z}}^t, \boldsymbol{\tau}_p, \boldsymbol{u}^t, \mathbf{s}^t, \boldsymbol{\tau}_s)$$
$$= H_G(q_{\mathbf{z}}^t, \boldsymbol{\tau}_p) - \frac{1}{2}\boldsymbol{\tau}_s^T(\boldsymbol{\tau}_p - \mathbf{S}\,\boldsymbol{\tau}_x^t) + const. = const. +$$
$$\frac{1}{2}\sum_{k=1}^M \left[\frac{\tau_{z_k}^t}{\tau_{p_k}} + \ln(2\pi\,\tau_{p_k})\right] - \frac{1}{2}\sum_{k=1}^M \tau_{s_k}(\tau_{p_k} - \mathbf{S}_{k,:}\,\boldsymbol{\tau}_x^t)$$
(27)

where $const.$ denotes constants w.r.t. $\{\boldsymbol{\tau}_p, \boldsymbol{\tau}_s\}$. Deriving w.r.t. $\{\boldsymbol{\tau}_p, \boldsymbol{\tau}_s\}$ yields the feasibility conditions

$$\frac{\partial L}{\partial \tau_{s_k}} = 0 \Rightarrow \tau_{p_k}^t = \mathbf{S}_{k,:}\,\boldsymbol{\tau}_x^t$$
(28)

$$\frac{\partial L}{\partial \tau_{p_k}} = 0 \Rightarrow \tau_{s_k}^t = \frac{1}{\tau_{p_k}^t}\left(1 - \frac{\tau_{z_k}^t}{\tau_{p_k}^t}\right).$$
(29)

which we run as a fixed-point sub-algorithm. The position of these updates in the updating schedule is less important. Nevertheless we shall update $\{\boldsymbol{\tau}_p, \boldsymbol{\tau}_s\}$ as soon as the quantities on which they depend have been updated.

## 4. CONVERGENCE TO LMMSE
In the case of Gaussian $p_{\boldsymbol{x}}, p_{\boldsymbol{y}|\mathbf{z}}$, the cost function is quadratic in $\boldsymbol{x}$ etc., and we check convergence to the LMMSE estimate. At convergence we have

$$\mathbf{s}^t = \mathbf{s}^{t-1} \Rightarrow \widehat{\mathbf{z}} = \mathbf{A}\,\boldsymbol{u}$$
$$\Rightarrow \boldsymbol{u} = \widehat{\boldsymbol{x}} \text{ from } \mathbf{g}(\boldsymbol{u}) = \mathbf{0} = \mathbf{g}(\mathbf{0}) + \mathcal{H}\,\boldsymbol{u}$$
$$\widehat{\mathbf{z}} = \boldsymbol{\tau}_z.(\boldsymbol{y}./\boldsymbol{\sigma}_v^2 + (\mathbf{A}\,\boldsymbol{u})./\boldsymbol{\tau}_p - \mathbf{s})$$
$$\Rightarrow \mathbf{s} = \boldsymbol{y}./\boldsymbol{\sigma}_v^2 + (\mathbf{A}\,\widehat{\boldsymbol{x}})./\boldsymbol{\tau}_p - (\mathbf{A}\,\widehat{\boldsymbol{x}})./\boldsymbol{\tau}_z = (\boldsymbol{y} - \mathbf{A}\,\widehat{\boldsymbol{x}})./\boldsymbol{\sigma}_v^2 \quad (30)$$
$$\widehat{\boldsymbol{x}} = (\boldsymbol{\tau}_x./\boldsymbol{\tau}_r).(\boldsymbol{u} + \boldsymbol{\tau}_r.(\mathbf{A}^T\mathbf{s}))$$
$$\Rightarrow (1 - \boldsymbol{\tau}_x./\boldsymbol{\tau}_r)\widehat{\boldsymbol{x}} = \boldsymbol{\tau}_x.\widehat{\boldsymbol{x}}./\boldsymbol{\sigma}_x^2 = \boldsymbol{\tau}_x.(\mathbf{A}^T\mathbf{s}) \text{ or}$$
$$\widehat{\boldsymbol{x}} = [\mathbf{A}^T\mathbf{D}^{-1}(\boldsymbol{\sigma}_v^2)\mathbf{A} + \mathbf{D}^{-1}(\boldsymbol{\sigma}_x^2)]^{-1}\mathbf{A}^T\mathbf{D}^{-1}(\boldsymbol{\sigma}_v^2)\,\boldsymbol{y}.$$

Note that at convergence $\widehat{\boldsymbol{x}}$ does not depend on the various variance estimates that the algorithm produces. One can also get the following convergence values

$$\mathbf{s} = \boldsymbol{R}_{yy}^{-1}\boldsymbol{y}, \; \widehat{\boldsymbol{x}} = \mathbf{D}(\boldsymbol{\sigma}_x^2)\mathbf{A}^T\mathbf{s}, \; \mathbf{r} = \mathbf{D}(\boldsymbol{\sigma}_x^2 + \boldsymbol{\tau}_r)\mathbf{A}^T\mathbf{s} \quad (31)$$

where $\mathbf{r}$ corresponds to the componentwise conditionally unbiased MMSE estimate of $\boldsymbol{x}$ [14], [15] if $\boldsymbol{\tau}_r$ converges to its correct value.

Below we shall analyze the convergence of the proposed AMBGAMP algorithm. Note that the updates of $q_x$, $q_z$ in (18), (22) imply that these approximate posteriors inherit the higher order cumulants of their respective priors (cf. Edgeworth expansions around a Gaussian). Only the means and variances are affected by the iterative algorithmn. In the Gaussian case, the mean subsystem depends on the variances, but the variance subsystem runs independently. Hence their convergence can be analyzed separately. In the non-Gaussian case, their coupling may need to be reconsidered though.

## 5. CONVERGENCE OF THE VARIANCE SUBSYSTEM

In the Gaussian priors case, the updates of the variances can be checked to result in the following variance subsystem

$$1./\boldsymbol{\tau}_x^t = 1./\boldsymbol{\sigma}_x^2 + \mathbf{S}^T \boldsymbol{\tau}_s^{t-1},$$
$$1./\boldsymbol{\tau}_s^t = \boldsymbol{\sigma}_v^2 + \mathbf{S}\, \boldsymbol{\tau}_x^t. \tag{32}$$

To analyze convergence, we investigate the contractiveness of the mappings via their Jacobians

$$\frac{\partial \boldsymbol{\tau}_x^t}{\partial \boldsymbol{\tau}_s^{t-1\,T}} = -\mathbf{D}_x^{t,\,2}\,\mathbf{S}^T, \; \frac{\partial \boldsymbol{\tau}_s^t}{\partial \boldsymbol{\tau}_x^{t\,T}} = -\mathbf{D}_s^{t,\,2}\,\mathbf{S} \tag{33}$$

where we introduced the notation $\mathbf{D}_x^t = \mathbf{D}(\boldsymbol{\tau}_x^t)$ etc., $\mathbf{D}_x^{t,\,2} = (\mathbf{D}_x^t)^2$ etc. By the chain rule, we get

$$\frac{\partial \boldsymbol{\tau}_x^t}{\partial \boldsymbol{\tau}_x^{t-1\,T}} = \frac{\partial \boldsymbol{\tau}_x^t}{\partial \boldsymbol{\tau}_s^{t-1\,T}} \frac{\partial \boldsymbol{\tau}_s^{t-1}}{\partial \boldsymbol{\tau}_x^{t-1\,T}}$$
$$\frac{\partial \boldsymbol{\tau}_x^t}{\partial \boldsymbol{\tau}_x^{1\,T}} = \mathbf{D}_x^{t,\,2}\,\mathbf{S}^T\mathbf{D}_s^{t-1,\,2}\,\mathbf{S}\,\mathbf{D}_x^{t-1,\,2}\,\mathbf{S}^T\mathbf{D}_s^{t-2,\,2}\,\mathbf{S}\ldots. \tag{34}$$

Note that the cascade of Jacobians involves a cascade of the following matrices

$$\mathbf{B}^t = \mathbf{B}_x^t\,\mathbf{B}_s^{t-1} = (\mathbf{D}_x^t\,\mathbf{S}^T\mathbf{D}_s^{t-1})\,(\mathbf{D}_s^{t-1}\,\mathbf{S}\,\mathbf{D}_x^{t-1}) \tag{35}$$

We can investigate the contractivity of $\mathbf{B}^t$ using any norm since all valid norms are commensurate. A judicious choice here is the infinity norm

$$\|\mathbf{B}^t\|_\infty = \max_{\boldsymbol{x}\neq\boldsymbol{0}} \frac{\|\mathbf{B}^t\,\boldsymbol{x}\|_\infty}{\|\boldsymbol{x}\|_\infty} = \max_{\|\boldsymbol{x}\|_\infty=1} \|\mathbf{B}^t\boldsymbol{x}\|_\infty = \|\mathbf{B}^t\mathbf{1}\|_\infty \tag{36}$$

where the last identity follows from the non-negativity of the elements of $\mathbf{B}$ (and $\mathbf{B}_x$ or $\mathbf{B}_s$) which implies that $\mathbf{B}^t\boldsymbol{x} \preceq \mathbf{B}^t\mathbf{1}$ for any $\boldsymbol{x}$ with $\|\boldsymbol{x}\|_\infty = 1$ (where the relation $\boldsymbol{x} \preceq \boldsymbol{y}$ indicates that $\boldsymbol{x}$ is element-wise not larger than $\boldsymbol{y}$). Now we have

$$\mathbf{B}^t\mathbf{1} = \mathbf{B}_x^t\,\mathbf{B}_s^{t-1}\mathbf{1} \preceq \|\mathbf{B}_s^{t-1}\|_\infty\,\mathbf{B}_x^t\,\mathbf{1} \preceq \|\mathbf{B}_x^t\|_\infty\,\|\mathbf{B}_s^{t-1}\|_\infty\,\mathbf{1}$$
$$\Rightarrow\; \|\mathbf{B}^t\|_\infty \le \|\mathbf{B}_x^t\|_\infty\,\|\mathbf{B}_s^{t-1}\|_\infty < 1 \tag{37}$$

where we assume that at least one of $\|\mathbf{B}_x^t\|_\infty \le 1$, $\|\mathbf{B}_s^{t-1}\|_\infty \le 1$ is strictly smaller than one. So, (37) implies converges of the variance subsystem. The statements in (37) hold in the general non-Gaussian case. In the Gaussian case, we get from (35), (32)

$$\|\mathbf{B}_s^{t-1}\|_\infty = \|\mathbf{B}_s^{t-1}\mathbf{1}\|_\infty = \|\mathbf{D}_s^{t-1}\,\mathbf{S}\,\mathbf{D}_x^{t-1}\,\mathbf{1}\|_\infty$$
$$= \|\mathbf{D}_s^{t-1}\,\mathbf{S}\,\boldsymbol{\tau}_x^{t-1}\|_\infty = \max_k \frac{\mathbf{S}_{k,:}\,\boldsymbol{\tau}_x^{t-1}}{\sigma_{v\,k}^2 + \mathbf{S}_{k,:}\,\boldsymbol{\tau}_x^{t-1}} \le 1,$$
$$\|\mathbf{B}_x^t\|_\infty = \|\mathbf{B}_x^t\mathbf{1}\|_\infty = \|\mathbf{D}_x^t\,\mathbf{S}^T\mathbf{D}_s^{t-1}\,\mathbf{1}\|_\infty$$
$$= \|\mathbf{D}_x^t\,\mathbf{S}^T\boldsymbol{\tau}_s^{t-1}\|_\infty = \max_i \frac{\mathbf{S}_{:,i}^T\,\boldsymbol{\tau}_s^{t-1}}{1/\sigma_{x\,i}^2 + \mathbf{S}_{:,i}^T\boldsymbol{\tau}_s^{t-1}} \le 1. \tag{38}$$

## 6. LARGE SYSTEM ANALYSIS WITH N.I.I.D. A

To show that at convergence, the variance subsystem $\boldsymbol{\tau}_x$ converges in the large system limit to the optimal MSE in the Gaussian case. We use the following result from [16], [17] :

**Theorem 1.** *Let* $\boldsymbol{Q}_N, \mathbf{D}_N \in \mathbb{R}^{N\times N}$ *be deterministic symmetric matrices and* $\mathbf{Y}_N = \mathbf{X}_N\mathbf{D}\mathbf{X}_N^H = \sum_{i=1}^M d_i\boldsymbol{x}_i\boldsymbol{x}_i^H$*, with diagonal* $\mathbf{D}$ *and* $\mathbf{X}_N$ *containing* $M$ *independent columns* $\boldsymbol{x}_i$ *with covariance matrix* $\boldsymbol{\Theta}_i$*. Also, assume that* $\boldsymbol{Q}_N$*,* $\mathbf{D}_N$*,* $\boldsymbol{\Theta}_i$ *have uniformly bounded spectral norms. Then, as* $M, N \to \infty$ *at constant ratio*

$$\frac{1}{N}tr\left[\boldsymbol{Q}_N(\mathbf{Y}_N + \mathbf{D}_N)^{-1}\right] - \frac{1}{N}tr[\boldsymbol{Q}_N\mathbf{T}_N] \xrightarrow{a.s.} 0, \; with \tag{39}$$

$$\mathbf{T}_N = \left(\sum_{i=1}^M \frac{d_i\boldsymbol{\Theta}_\mathbf{i}}{1+e_i} + \mathbf{D}_N\right)^{-1}, \; where\; the\; e_i\; satisfy \tag{40}$$

$$e_k = tr\left[d_k\boldsymbol{\Theta}_k\left(\sum_{i=1}^M \frac{d_i\boldsymbol{\Theta}_i}{1+e_i} + \mathbf{D}_N\right)^{-1}\right], \; k = 1,\ldots,M. \tag{41}$$

The convergence in (39) is the convergence of a scalar to its mean, by LLN. Note the presence of the weights $1+e_i$ in the denominator of the sum in $\mathbf{T}_N$ in (40), which reflect that the expected value of a matrix inverse is not the inverse of its expected value. Note that the $\mathrm{tr}\,[\boldsymbol{\Theta}_i]$ should be of order 1, which means that the sum in $\mathbf{T}_N$ is implicitly normalized. The $e_i$ satisfy the implicit equations (41), and can be obtained as the fixed points of the RHS interpreted as a mapping (with global convergence).

We assume the columns of $\mathbf{A}^T = [\boldsymbol{a}_1 \ldots \boldsymbol{a}_M]$ to be zero mean and independent with diagonal covariance matrix $\mathbb{E}\left(\boldsymbol{a}_i\boldsymbol{a}_i^T\right) = \boldsymbol{\Theta}_i$. The optimal MSE in the Gaussian case is given by (with e.g. $\mathbf{D}_{\sigma_x^2} = \mathbf{D}(\boldsymbol{\sigma}_x^2)$)

$$\mathrm{MSE} = \frac{1}{N}\mathrm{tr}\left\{\left[\mathbf{A}^T\mathbf{D}_{\sigma_v^2}^{-1}\mathbf{A} + \mathbf{D}_{\sigma_x^2}^{-1}\right]^{-1}\right\}$$
$$\xrightarrow{a.s.} \frac{1}{N}\mathrm{tr}\left\{\left[\sum_{i=1}^M \frac{\boldsymbol{\Theta}_i}{\sigma_{v,i}^2(1+e_i)} + \mathbf{D}_{\sigma_x^2}^{-1}\right]^{-1}\right\}, \; with \tag{42}$$

$$e_k = \mathrm{tr}\left\{\frac{1}{\sigma_{v,k}^2}\boldsymbol{\Theta}_k\left[\sum_{i=1}^M \frac{\boldsymbol{\Theta}_i}{\sigma_{v,i}^2(1+e_i)} + \mathbf{D}_{\sigma_x^2}^{-1}\right]^{-1}\right\}. \tag{43}$$

On the other hand the GAMP variance subsystem converges to (32), without iteration indices. With large $\mathbf{A}$, $\mathbf{S}\boldsymbol{\tau}_x$ and $\mathbf{S}^T\boldsymbol{\tau}_s$ converge to their expected values

$$\mathbb{E}\left[\mathbf{S}\boldsymbol{\tau}_x\right]_i = \mathbb{E}\left[\mathbf{A}\mathbf{D}_x\mathbf{A}^T\right]_{ii} = \mathrm{tr}\{\boldsymbol{\Theta}_i\mathbf{D}_x\}; \tag{44}$$

$$\mathbb{E}\,\mathbf{D}\left(\mathbf{S}^T\boldsymbol{\tau}_s\right) = \mathbb{E}\,\mathrm{diag}\left(\mathbf{A}^T\mathbf{D}_s\mathbf{A}\right) = \sum_{k=1}^M \tau_{s,k}\boldsymbol{\Theta}_k. \tag{45}$$

Therefore, the empirical mean of the posterior variance $\boldsymbol{\tau}_x$ becomes

$$\frac{1}{N}\mathrm{tr}\{\mathbf{D}_x\} = \frac{1}{N}\mathrm{tr}\left\{\left[\mathbf{D}_{\sigma_x^2}^{-1} + \sum_{k=1}^M \tau_{s,k}\boldsymbol{\Theta}_k\right]^{-1}\right\}. \tag{46}$$

From (32), (44), it follows that

$$\tau_{s,k} = \frac{1}{\sigma_{v,k} + \mathrm{tr}\{\boldsymbol{\Theta}_k\mathbf{D}_x\}}. \tag{47}$$

Define $e_k' = \frac{\mathrm{tr}\{\boldsymbol{\Theta}_k\mathbf{D}_x\}}{\sigma_{v,k}^2}$ and substituting (47) into (46), we obtain

$$\frac{1}{N}\mathrm{tr}\{\mathbf{D}(\boldsymbol{\tau}_x)\} = \frac{1}{N}\mathrm{tr}\left\{\left[\mathbf{D}_{\sigma_x^2}^{-1} + \sum_{i=1}^M \frac{\boldsymbol{\Theta}_i}{\sigma_{v,i}(1+e_i')}\right]^{-1}\right\};$$
$$e_k' = \mathrm{tr}\left\{\frac{\boldsymbol{\Theta}_k}{\sigma_{v,k}^2}\left[\mathbf{D}_{\sigma_x^2}^{-1} + \sum_{i=1}^M \frac{\boldsymbol{\Theta}_i}{\sigma_{v,i}(1+e_i')}\right]^{-1}\right\}. \tag{48}$$

This shows that the system of equations defined by (48) is the same as the system defined by (42), (43). Therefore, the empirical mean of $\boldsymbol{\tau}_x$ converges to the optimal MSE in the large system limit. Note that above we have applied the Theorem with $\boldsymbol{Q}_N = \mathbf{I}_N$ but the same results hold for any $\boldsymbol{Q}_N$, corresponding to deterministic limits for variably weighted MSEs $\mathrm{tr}\{\boldsymbol{Q}_N\mathbf{D}_x\}/\mathrm{tr}\{\boldsymbol{Q}_N\}$.

## 7. CONVERGENCE OF THE MEAN SUBSYSTEM

We consider here the convergence proof for the case in which the update of $\boldsymbol{u}$ minimizes its quadratic cost function, i.e. $\mathbf{g}^t(\boldsymbol{u}^t) = \mathbf{g}^t(\mathbf{0}) + \mathcal{H}^t \boldsymbol{u}^t = \mathbf{0} \Rightarrow \boldsymbol{u}^t = -(\mathcal{H}^t)^{-1}\mathbf{g}^t(\mathbf{0})$. We shall investigate the convergence of the mean subsystem once the variance subsystem has converged. Similar to the convergence proof in [4], we will derive the Jacobian of the updating function and prove the convergence by showing that the Jacobian is contractive. We define the short hand notations

$$\boldsymbol{\tau} = \begin{bmatrix} \boldsymbol{\tau}_r^T & \boldsymbol{\tau}_p^T \end{bmatrix}^T, \mathbf{D} = \mathbf{D}(1./\boldsymbol{\tau}), \mathbf{B} = \mathbf{D}^{\frac{1}{2}}\begin{bmatrix} \mathbf{I} & \mathbf{A}^T \end{bmatrix}^T,$$
$$\mathbf{C} = \mathbf{D}^{-\frac{1}{2}}\begin{bmatrix} \mathbf{A} & -\mathbf{I} \end{bmatrix}^T, \mathbf{H} = \begin{bmatrix} \mathbf{0}_{M \times N} & \mathbf{I} \end{bmatrix}\mathbf{D}^{\frac{1}{2}}, \quad (49)$$
$$\boldsymbol{w}^t = \mathbf{D}^{\frac{1}{2}}\begin{bmatrix} \widehat{\boldsymbol{x}}^{t\,T} & \widehat{\boldsymbol{z}}^{t\,T} \end{bmatrix}^T, \boldsymbol{P} = \mathbf{B}\left(\mathbf{B}^T\mathbf{B}\right)^{-1}\mathbf{B}^T.$$

Furthermore, we define the update function for $\begin{bmatrix} \widehat{\boldsymbol{x}}^{t\,T} & \widehat{\boldsymbol{z}}^{t\,T} \end{bmatrix}^T$ as

$$\begin{bmatrix} \widehat{\boldsymbol{x}}^t \\ \widehat{\boldsymbol{z}}^t \end{bmatrix} = \mathbf{g}\left(\begin{bmatrix} \mathbf{r}^t \\ \mathbf{p}^t \end{bmatrix}\right) = \begin{bmatrix} \mathbf{g}_x(\mathbf{r}^t) \\ \mathbf{g}_z(\mathbf{p}^t) \end{bmatrix} = \begin{bmatrix} \mathbb{E}(\boldsymbol{x}|\mathbf{r}^t, \boldsymbol{\tau}_r) \\ \mathbb{E}(\mathbf{z}|\mathbf{p}^t, \boldsymbol{\tau}_p) \end{bmatrix}. \quad (50)$$

In the following we will show that the system $\boldsymbol{\theta}^t = \begin{bmatrix} \boldsymbol{w}^{t,T} & \mathbf{s}^{t,T} \end{bmatrix}^T$ is converging. With notations defined above, we can rewrite $\boldsymbol{u}^t = -(\mathcal{H}^t)^{-1}\mathbf{g}^t(\mathbf{0})$ as

$$\boldsymbol{u}^t = \left(\mathbf{B}^T\mathbf{B}\right)^{-1}\mathbf{B}^T\boldsymbol{w}^{t-1}. \quad (51)$$

The vector $\boldsymbol{w}^t$ is updated via

$$\boldsymbol{w}^t = \tilde{\mathbf{g}}\left(\boldsymbol{P}\boldsymbol{w}^{t-1} + \mathbf{C}\mathbf{s}^{t-1}\right) = \tilde{\mathbf{g}}\left(\begin{bmatrix} \boldsymbol{P} & \mathbf{C} \end{bmatrix}\boldsymbol{\theta}^{t-1}\right), \quad (52)$$

where $\tilde{\mathbf{g}}(\boldsymbol{v}) = \mathbf{D}^{\frac{1}{2}}\mathbf{g}\left(\mathbf{D}^{-\frac{1}{2}}\boldsymbol{v}\right)$. The update of $\mathbf{s}^t$ can be written as

$$\begin{aligned} \mathbf{s}^t &= \mathbf{s}^{t-1} + \mathbf{H}\left(\boldsymbol{w}^t - \mathbf{B}\boldsymbol{u}^t\right) \\ &= \mathbf{s}^{t-1} + \mathbf{H}\left[\tilde{\mathbf{g}}\left(\begin{bmatrix} \boldsymbol{P} & \mathbf{C} \end{bmatrix}\boldsymbol{\theta}^{t-1}\right) - \boldsymbol{P}\boldsymbol{w}^{t-1}\right]. \end{aligned} \quad (53)$$

Combining (52) and (53), we obtain the update equation for $\boldsymbol{\theta}^t$,

$$\boldsymbol{\theta}^t = \mathbf{h}(\boldsymbol{\theta}^{t-1}) = \begin{bmatrix} \mathbf{I} \\ \mathbf{H} \end{bmatrix}\tilde{\mathbf{g}}\left(\begin{bmatrix} \boldsymbol{P}\mathbf{C} \end{bmatrix}\boldsymbol{\theta}^{t-1}\right) + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ -\mathbf{H}\boldsymbol{P} & \mathbf{I} \end{bmatrix}\boldsymbol{\theta}^{t-1}, \quad (54)$$

where $\mathbf{h}(\boldsymbol{\theta}^{t-1})$ denotes the update function. We get for the Jacobian $\tilde{\mathbf{g}}'^t = \tilde{\mathbf{g}}'\left(\begin{bmatrix} \boldsymbol{P}\mathbf{C} \end{bmatrix}\boldsymbol{\theta}^{t-1}\right) = \mathbf{D}^{\frac{1}{2}}\mathbf{g}'^t\mathbf{D}^{-\frac{1}{2}} = \mathbf{g}'^t$ which is diagonal since $\mathbf{g}(.)$ is an elementwise function. As mentioned in [4], $\mathbf{g}'^t$ is a positive semi-definite diagonal matrix with all elements smaller than 1. Furthermore, in the Gaussian case $\mathbf{g}'^t$ is a constant matrix. The Jacobian of $\mathbf{h}(\boldsymbol{\theta}^{t-1})$ is given by

$$\begin{aligned} \mathbf{h}'^t &= \begin{bmatrix} \mathbf{I} \\ \mathbf{H} \end{bmatrix}\mathbf{g}'^t\begin{bmatrix} \boldsymbol{P} & \mathbf{C} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ -\mathbf{H}\boldsymbol{P} & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{g}'^t\boldsymbol{P} & \mathbf{g}'^t\mathbf{C} \\ \mathbf{H}(\mathbf{g}'^t - \mathbf{I})\boldsymbol{P} & \mathbf{H}\mathbf{g}'^t\mathbf{C} + \mathbf{I} \end{bmatrix}. \end{aligned} \quad (55)$$

We calculate the terms $\mathbf{H}\mathbf{g}'^t\mathbf{C}$ and $\mathbf{H}(\mathbf{g}'^t - \mathbf{I})$, which leads to

$$\mathbf{H}\mathbf{g}'^t\mathbf{C} = -\mathbf{g}_z'^t, \quad \mathbf{H}(\mathbf{g}'^t - \mathbf{I}) = (\mathbf{g}_z'^t - \mathbf{I})\mathbf{H}. \quad (56)$$

Hence, the Jacobian $\mathbf{h}'^t$ becomes

$$\mathbf{h}'^t = \begin{bmatrix} \mathbf{g}'^t\boldsymbol{P} & \mathbf{g}'^t\mathbf{C} \\ (\mathbf{g}_p'^t - \mathbf{I})\mathbf{H}\boldsymbol{P} & \mathbf{I} - \mathbf{g}_z'^t \end{bmatrix} = \begin{bmatrix} \mathbf{g}'^t & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \mathbf{g}_z'^t \end{bmatrix}\begin{bmatrix} \boldsymbol{P} & \mathbf{C} \\ -\mathbf{H}\boldsymbol{P} & \mathbf{I} \end{bmatrix}. \quad (57)$$

Since all the elements of $\mathbf{g}'^t$ range from 0 to 1,

$$\mathbf{0} \preceq \begin{bmatrix} \mathbf{g}'^t & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \mathbf{g}_z'^t \end{bmatrix} \preceq \max\{\|\mathbf{g}'^t\|_\infty, 1 - \|\mathbf{g}_z'^t\|_\infty\}\mathbf{I} \prec \mathbf{I}. \quad (58)$$

We calculate the eigenvalues of the second matrix at the right of (57) via

$$\det\begin{bmatrix} \lambda\mathbf{I} - \boldsymbol{P} & -\mathbf{C} \\ \mathbf{H}\boldsymbol{P} & \lambda\mathbf{I} - \mathbf{I} \end{bmatrix} = 0. \quad (59)$$

For the determinant of block matrices, we have

$$\begin{aligned} &\det\begin{bmatrix} \lambda\mathbf{I} - \boldsymbol{P} & -\mathbf{C} \\ \mathbf{H}\boldsymbol{P} & \lambda\mathbf{I} - \mathbf{I} \end{bmatrix} \\ &= \det(\lambda\mathbf{I} - \boldsymbol{P})\det[\lambda\mathbf{I} - \mathbf{I} + \mathbf{H}\boldsymbol{P}(\lambda\mathbf{I} - \boldsymbol{P})^{-1}\mathbf{C}]. \end{aligned} \quad (60)$$
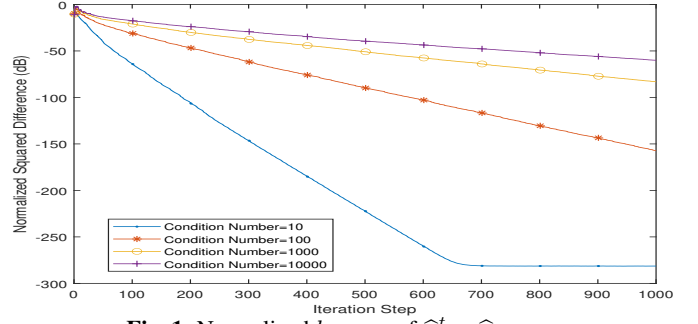


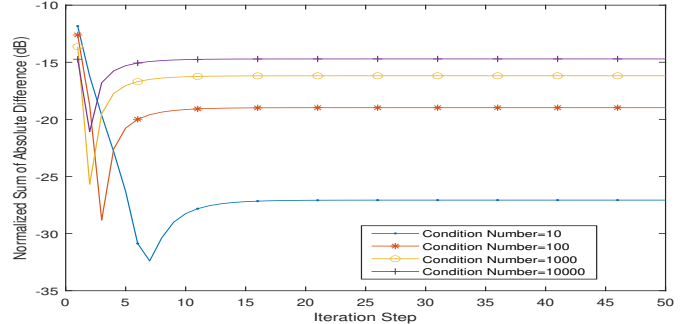**Fig. 1**. Normalized $l_2$ norm of $\widehat{\boldsymbol{x}}^t - \widehat{\boldsymbol{x}}_{\mathrm{MMSE}}$



**Fig. 2**. Normalized $l_1$ norm of $\boldsymbol{\tau}_x^t - \boldsymbol{\tau}_{\mathrm{MMSE}}$

By the matrix inverse lemma and the definition of $\boldsymbol{P}$, we get

$$(\lambda\mathbf{I} - \boldsymbol{P})^{-1} = \frac{\mathbf{I}}{\lambda} - \frac{\boldsymbol{P}}{\lambda - \lambda^2}. \quad (61)$$

Note that from the definition in (49), $\boldsymbol{P}\mathbf{C} = \mathbf{B}(\mathbf{B}^T\mathbf{B})^{-1}\mathbf{B}^T\mathbf{C} = \mathbf{0}$. Hence (60) becomes

$$\det\begin{bmatrix} \lambda\mathbf{I} - \boldsymbol{P} & -\mathbf{C} \\ \mathbf{H}\boldsymbol{P} & \lambda\mathbf{I} - \mathbf{I} \end{bmatrix} = \det(\lambda\mathbf{I} - \boldsymbol{P})\det(\lambda\mathbf{I} - \mathbf{I}). \quad (62)$$

Set this determinant to 0 and solve for $\lambda$, we find that the eigenvalues are 0 and 1. Therefore, the update operation $\mathbf{h}(\boldsymbol{\theta})$ is contractive.

## 8. SIMULATION RESULTS

For the simulations, we set the SNR to 20dB, and the system dimensions to $M \times N = 512 \times 1024$. We consider a Gaussian setting with white noise and $\boldsymbol{x}$ is drawn from an n.i.i.d. Gaussian distribution with zero mean and variance profile $\sigma_{x_i}^2 = 0.991^{i-1}$, $i = 1, \ldots, N$. For $\mathbf{A}$, we follow the setup in [18]. Namely first $\mathbf{A}$ gets generated as i.i.d. zero mean Gaussian, its SVD gets computed and the singular values $\{s_1 \geq \cdots \geq s_M\}$ are changed to a geometric series with a specific condition number $\frac{s_1}{s_M}$. We compare the results to the LMMSE estimator. Fig. 1 illustrates the difference between the $\widehat{\boldsymbol{x}}^t$ and the LMMSE $\widehat{\boldsymbol{x}}_{\mathrm{MMSE}}$, whereas Fig. 2 compares the difference between $\boldsymbol{\tau}_x^t$ and $\boldsymbol{\tau}_{\mathrm{MMSE}}$. The "normalization" mentioned in the captions refers to division by $N$. These simulations show that the AMBGAMP algorithm continues to work in unrealistically severe scenarios (in which AMP diverges).

## 9. CONCLUDING REMARKS

We propose a convergent version of GAMP, AMBGAMP, which applies alternating minimization to an augmented Lagrangian of a large system limit of the Bethe free Energy (BFE). AMBGAMP can be interpreted as applying a simplified ADMM to the BFE, with a constrained Lagrange multiplier parametrerization for the mean constraint, and a quadratic optimization subproblem being solved by a gradient update with line search. The ADMM is complemented with a fixed point iteration for the variance constraint.

## 10. REFERENCES

[1] M.E. Tipping, "Sparse Bayesian Learning and the Relevance Vector Machine," *J. Mach. Learn. Res.*, 2001.

[2] D. P. Wipf and B. D. Rao, "Sparse Bayesian Learning for Basis Selection ," *IEEE Trans. on Sig. Proc.*, vol. 52, no. 8, Aug. 2004.

[3] C. K. Thomas and D. Slock, "Low Complexity Static and Dynamic Sparse Bayesian Learning Combining BP, VB and EP Message Passing," in *Asilomar Conf. on Sig., Sys., and Comp.*, CA, USA, 2019.

[4] S. Rangan, A. Fletcher, P. Schniter, and U.S. Kamilov, "Inference for Generalized Linear Models via Alternating Directions and Bethe Free Energy Minimization," *IEEE Trans. Info. Theory*, Jan. 2017.

[5] S. Rangan, P. Schniter, and A. K. Fletcher, "Vector Approximate Message Passing," *IEEE Trans. On Info. Theo.*, Oct. 2019.

[6] P. Schniter, S. Rangan, and A. K. Fletcher, "Vector approximate Message Passing for the Generalized Linear Model," in *In IEEE Asilomar Conf. Signals, Systems and Computers*, 2016.

[7] Q. Guo and J. Xi, "Approximate message passing with unitary transformation," *arxiv.org/abs/1504.04799*, 2015.

[8] S. Rangan, P. Schniter, A. Fletcher, and V. Cevher, "Fixed Points of Generalized Approximate Message Passing with Arbitrary Matrices," *IEEE Trans. Info. Theory*, Dec. 2016.

[9] S. Rangan, "Generalized Approximate Message Passing for Estimation with Random Linear Mixing," in *Proc. IEEE Int. Symp. Inf. Theory*, Saint Petersburg, Russia, 2011, extended version: arxiv1010.5141.

[10] F. Krzakala, A. Manoel, E. W. Tramel, and L. Zdeborova, "Variational Free Energies for Compressed Sensing," in *IEEE Intl. Sympo. Info. Theo.*, Honolulu, HI, USA, Jun. 2014.

[11] D. Slock", "Convergent Approximate Message Passing by Alternating Constrained Minimization of Bethe Free Energy," in *IEEE 7th Forum on Research and Technologies for Society and Industry Innov. (RTSI)*, Paris, France, Aug. 2022.

[12] D. Slock", "Convergent Approximate Message Passing," in *IEEE Int'l Mediterr. Conf. Comm's and Netw'ing (MeditCom)*, Athens, Greece, Sep. 2022.

[13] C.K. Thomas and D. Slock, "Alternating Constrained Minimization based Approximate Message Passing," in *IEEE Int'l Conf. Acoustics, Speech and Sig. Proc. (ICASSP)*, 2023.

[14] M. Triki and D. Slock, "Component-Wise Conditionally Unbiased Bayesian Parameter Estimation: General Concept and Applications to Kalman Filtering and LMMSE Channel Estimation," in *IEEE Asilomar Conf. on Signals, Systems, and Computers*, Pacific Grove, USA, 2005.

[15] M. Triki and D. Slock, "Investigation of Some Bias and MSE Issues in Block-Component-Wise Conditionally Unbiased LMMSE," in *IEEE Asilomar Conf. on Signals, Systems, and Computers*, Pacific Grove, USA, 2006.

[16] S. Wagner, R. Couillet, M. Debbah, and D.T.M. Slock, "Large System Analysis of Linear Precoding in Correlated MISO Broadcast Channels Under Limited Feedback," *IEEE Trans. Info. Theory*, July 2012.

[17] Dirk TM Slock, "Nonlinear MMSE using Linear MMSE Bricks and Application to Compressed Sensing and Adaptive Kalman Filtering," IEEE ICASSP Expert to Non Expert (ETON) Primer, 2020.

[18] P. Schniter, S. Rangan, and A.K. Fletcher, "Vector Approximate Message Passing for the Generalized Linear Model," in *IEEE Asilomar Conf. on Signals, Systems and Computers*, 2016.