

Deliverable 5.2: Large System Analysis Techniques and Outcomes

Zilu Zhao, Christian Forsch, Laura Cottatellucci, Dirk Slock

January 30, 2026

Abstract

Generalized Approximate Message Passing (GAMP) algorithms are highly effective for signal recovery and can be derived as asymptotic approximations of Expectation Propagation (EP). EP constructs approximate posteriors by iteratively combining extrinsic information with prior factors, whereas low-complexity algorithms such as GAMP derive extrinsics directly from posterior beliefs. In the Gaussian case, we show that extrinsics are closely related to Component-Wise Conditionally Unbiased MMSE (CWCU-MMSE) estimation, while posterior beliefs yield standard MMSE (linear MMSE) estimators. We rederive the revisited GVAMP algorithm as an asymptotic alternating minimization of the Kullback–Leibler divergence and analyze extrinsics via asymptotic perturbations linking posterior beliefs and extrinsic messages.

We further study AMP from the perspective of the Bethe Free Energy (BFE) of generalized linear models. By applying large-system-limit (LSL) approximations directly to belief propagation, we clarify the relationship between posterior distributions and extrinsic messages and rederive fundamental deterministic variance results. This interpretation explains the structure of augmented Lagrangian formulations used in convergent AMP variants and facilitates extensions to more complex models.

Finally, we apply investigate the LSL in bilinear systems by studying semi-blind channel estimation in cell-free massive MIMO uplink communications. We propose a simplified, decentralized EP-based approach that combines orthogonal pilots, central limit theory, and scale-aware updates, significantly reducing computational complexity while mitigating pilot contamination.

Contents

1	Introduction	3
1.1	Prior Work	3
1.2	Main Contributions	3
2	Bethe Free Energy of the Generalized Linear Model	4
2.1	Bethe Free Energy (BFE)	4
2.2	BFE of the GLM for BP	4
3	reGVAMP	5
3.1	reGVAMP from (Minka) EP	6
4	Relation to CWCU MMSE Estimator	8
5	GAMP from LSL Belief Propagation	10
5.1	Output Node	11
5.2	Input Node	12
6	LSL BFE and EP	13
7	Iterative Solution leading to GAMP	14
7.1	Update of λ_{μ_z}	14
7.2	Update of λ_{μ_w}	15
7.3	The update of λ_τ and τ_p	15
8	Iterative Solution leading to AMBGAMP	15
8.1	Update of λ_{μ_z}	15
8.2	Update of λ_{μ_w}	16
8.3	Update of \mathbf{u}	16
9	Bilinear System Model	16
9.1	Orthogonal Pilot sequences	17
9.2	Expectation Propagation on Semi-Blind structure	17
10	Bilinear Message Passing Derivations	17
10.1	Message from $\Psi_{2,lt}$ to x_{kt}	18
10.2	Message from $\Psi_{2,lt}$ to \mathbf{h}_{lk}	19
10.3	Message form $\Psi_{3,lq}$ to \mathbf{h}_{lk}	19
11	Asymptotic Behaviors in Bilinear Large Systems	20
12	Simplification of the Messages in Bilinear EP	22
13	Decentralized Method for Bilinear EP	22
14	Bilinear EP Simulation Results	23
15	Concluding Remarks	24
16	References	26

1 Introduction

Sparse signal recovery is a fundamental problem in signal processing with a wide range of applications. Many of these problems can be framed as the task of estimating a latent vector \mathbf{x} based on a correlated observation vector \mathbf{y} [1]. In the Bayesian framework, the complexity of Canonical Methods such as MMSE and MAP experiences exponential growth as the dimension of the problem grows.

By exploiting the structure of the models, graphical model based methods prove to be effective. Belief Propagation (BP) transforms the global inference problem into a local inference problem as outlined by [2]. Loopy Belief Propagation (LBP) extends BP by directly employing BP on a factorization scheme for $p(\mathbf{x}|\mathbf{y})$ that may involve loops [3]. In comparison to BP, LBP can be considered as an approximation method.

A limitation of (L)BP is that the (iterative) updating scheme leads to pdfs that correspond to the product of a large number of messages, leading to high complexity. To address this issue, Expectation Propagation (EP) was introduced [4]. EP has been shown to share a similar updating scheme as (L)BP, but for computational efficiency, the messages in (L)BP are projected into a suitable member of the family of exponential distributions [4].

Variational Bayes (VB) [1] provides a method seemingly parallel to (L)BP. VB aims to approximate the true posterior $p(\mathbf{x}|\mathbf{y})$ as a simpler distribution $\hat{q}(\mathbf{x})$. It introduces variational free energy $D[\hat{q}(\mathbf{x})||p(\mathbf{x}|\mathbf{y})]$ which is the Kullback-Leibler (KL) divergence. The approximate distribution is obtained by minimizing the KL divergence.

1.1 Prior Work

In both [1] and [5], the authors unify EP and BP within the framework of minimizing variational free energy. They demonstrate the close relationship between the fixed points of various message-passing algorithms and the stationary points of Bethe Free Energy (BFE).

EP can serve as an inference method in the generalized linear model (GLM). However, the computational cost corresponds to propagating $2MN$ messages as in Fig. 2 when the data matrix \mathbf{A} is of size $M \times N$. Generalized Approximate Message Passing (GAMP) [6] builds upon EP, but through the application of large system approximations (LSA), it effectively reduces the number of messages to $M+N$ extrinsics and (marginal) posteriors, providing a more computationally efficient approach.

In [7], the authors investigated the fixed points of the Generalized AMP (GAMP) algorithm for GLMs. They discovered that GAMP shares the same fixed point as the stationary points of the Large System Limit Bethe Free Energy (LSL BFE). In [8] we then proposed AMBGAMP which is guaranteed to converge. Building upon the works of [1], [9], [8], [10], [11], and [12].

The Component-Wise Conditionally Unbiased (CWCU) Minimum Mean Squared Error (MMSE) estimator is introduced in [13] and rederived in [14] for both joint Gaussian models and linear models. This concept was also used in [15], where the authors call it individual bias compensation. The connection between CWCU MMSE estimation and extrinsic information is explored in [16] specifically for linear Gaussian models.

1.2 Main Contributions

We rederive the reGVAMP algorithm that we introduced in [16], [11], from the point of view of alternating minimization of a LSL version of a desirable KLD. The asymptotics here involve only the CLT for extrinsics. We then derive the GAMP algorithm by directly introducing LSL simplifications in the LBP algorithm. This leads us to relate extrinsic messages to posterior pdfs by first order Taylor series expansion based perturbations. We also apply LSL approximations to the variances of the various Gaussians involved, which in fact leads to a rederivation of a fundamental LSA theorem describing the deterministic limit of LMMSE posterior variances.

We then investigate the LSL in Bilinear systems. We present a simplified, decentralized, EP-based method designed to address the Semi-Blind estimation problem in communication systems.

By utilizing orthogonal pilots, we are able to decouple the channels for different users into mutually exclusive groups, which reduces computational complexity. To further decrease computational demands, we integrate Expectation Propagation (EP) with Central Limit Theory (CLT), treating the interference as noise. Drawing inspiration from [17], we introduce further simplifications through scale analysis. Additionally, to lessen the load on the central processing unit (CPU), we explore a decentralized scheme.

2 Bethe Free Energy of the Generalized Linear Model

2.1 Bethe Free Energy (BFE)

Consider a pdf factorization

$$p(\mathbf{x}, \mathbf{y}) \propto \prod_{\alpha} f_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}), \quad (1)$$

where \mathbf{x}_{α} is a subvector of \mathbf{x} . In case of a tree-structured factor graph, an alternative equivalent form is [2]

$$p(\mathbf{x}|\mathbf{y}) = \frac{\prod_{\alpha} p(\mathbf{x}_{\alpha})}{\prod_i p(x_i)^{M_i-1}}, \quad (2)$$

where M_i is the number of subvectors \mathbf{x}_{α} that contain x_i . In (2), the $p(\mathbf{x}_{\alpha})$ and $p(x_i)$ are the exact factor (subvector) resp. variable marginals.

The concept of variational free energy suggests that to infer the marginals from a tree structured $p(\mathbf{x}, \mathbf{y})$ given in (1), we can use as trial distribution

$$q_{\mathbf{x}}(\mathbf{x}) = \frac{\prod_{\alpha} q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha})}{\prod_i q_{x_i}(x_i)^{M_i-1}}. \quad (3)$$

The true marginals can be obtained by [1]

$$\begin{aligned} \min_{q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}), q_{x_i}(x_i)} F &= D[q(\mathbf{x}) \| \prod_{\alpha} f_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha})]; \\ \text{s.t. } \forall \alpha, \forall i \in \mathcal{I}_{\alpha}, q_{x_i}(x_i) &= \int q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}) d\mathbf{x}_{\bar{i}}, \end{aligned} \quad (4)$$

where we define the shorthand notation (for arbitrary nonnegative functions q, p) $D(q\|p) = \int q(x) \ln \frac{q(x)}{p(x)} dx$ (which is the Kullback-Leibler Divergence (KLD) in case of normalized q, p) and $\mathbf{x}_{\bar{i}}$ denotes all \mathbf{x} except x_i . The free energy can be expanded as

$$F = \sum_{\alpha} D[q_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha}) \| f_{\mathbf{x}_{\alpha}}(\mathbf{x}_{\alpha})] + \sum_i (M_i - 1) H[q_{x_i}(x_i)], \quad (5)$$

where $H(\cdot)$ denotes entropy in nats. Note that this representation only holds for a tree structured distribution. For general graphs that contain loops, (2) no longer holds. Thus, in cases with loops, (5) is only an approximation of the variational free energy. The expression (5) is instead called Bethe free energy.

2.2 BFE of the GLM for BP

We consider a GLM with

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i), \quad \mathbf{z} = \mathbf{A}\mathbf{x}, \quad p(\mathbf{y}|\mathbf{z}) = \prod_{j=1}^M p(y_j|z_j), \quad (6)$$

where the ratio N/M is a constant for large system considerations. We interpret the linear mixing as a conditional probability

$$p(\mathbf{z}|\mathbf{x}) = \delta(\mathbf{z} - \mathbf{A}\mathbf{x}). \quad (7)$$

From this general linear model, a joint (loopy) factorization scheme comes up naturally:

$$p(\mathbf{x}, \mathbf{z}|\mathbf{y}) \propto p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(\mathbf{y}|\mathbf{z}) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) p(\mathbf{x}). \quad (8)$$

According to the definition of BFE (5), the associated BFE based on the joint factorization scheme (8) is calculated [1] as

$$F = D[q_{\mathbf{x}}(\mathbf{x})\|p(\mathbf{x})] + D[q_{\mathbf{z}}(\mathbf{z})\|p(\mathbf{y}|\mathbf{z})] + \sum_i H[q_{x_i}(x_i)] \\ + D[b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z})\|\delta(\mathbf{z} - \mathbf{A}\mathbf{x})] + \sum_j H[q_{z_j}(z_j)], \quad (9)$$

where $q_{\mathbf{x}}$, $q_{\mathbf{z}}$, $b_{\mathbf{x},\mathbf{z}}$, q_{x_i} and q_{z_j} are only approximate posteriors because of the loops in the factor graph. Since we need to minimize the BFE given by (9), the distribution function $b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z})$ must be of the form

$$b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z}) = b_{\mathbf{x}}(\mathbf{x})\delta(\mathbf{z} - \mathbf{A}\mathbf{x}), \quad (10)$$

to avoid an infinite value of the KLD, leading to $D[b_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z})\|\delta(\mathbf{z} - \mathbf{A}\mathbf{x})] = -H[b_{\mathbf{x}}]$. For BP, the BFE (9) needs to be minimized w.r.t. marginal consistency constraints $q_{\mathbf{x}}(x_i) = b_{\mathbf{x}}(x_i) = q_{x_i}(x_i)$, $q_{\mathbf{z}}(z_j) = q_{z_j}(z_j)$. Given the independent priors for \mathbf{x} , \mathbf{z} , minimization of the BFE leads to $q_{\mathbf{x}}(\mathbf{x}) = \prod_i q_{x_i}(x_i)$, $q_{\mathbf{z}}(\mathbf{z}) = \prod_j q_{z_j}(z_j)$. Furthermore, the maximization of $H[b_{\mathbf{x}}]$ under marginal constraints leads to $b_{\mathbf{x}}(\mathbf{x}) = \prod_i b_{x_i}(x_i)$. Together with the marginal constraints, this leads to the cancellation of the entropy terms in \mathbf{x} in the BFE, which becomes $F =$

$$\sum_i D[q_{x_i}(x_i)\|p(x_i)] + \sum_j D[q_{z_j}(z_j)\|p(y_j|z_j)] + \sum_j H[q_{z_j}(z_j)] \quad (11)$$

which needs to be minimized under the constraint $\mathbf{z} = \mathbf{A}\mathbf{x}$.

3 reGVAMP

reGVAMP (revisited Generalized Vector AMP) is motivated by only a *single asymptotic approximation*: the asymptotic Gaussianity of extrinsics. The extrinsic pdf of a variable x_i is the conditional pdf $p(\mathbf{y}|x_i)$, in which x_i is treated as a deterministic variable (no prior information), but the other variables $\mathbf{x}_{\bar{i}}$ remain random and their prior pdf is exploited to eliminate them from the joint pdf. The randomness of \mathbf{x} and \mathbf{A} will quickly lead to Gaussianity of $p(\mathbf{y}|x_i)$ by the CLT (think of asymptotic Gaussianity of Maximum Likelihood estimates).

reVAMP introduces both Gaussian and non-Gaussian marginal posteriors from Gaussian extrinsics and the true prior. This involves also the introduction of Gaussian approximations for the priors. Which in turn also leads to a multivariate Gaussian posterior approximation, which exhibits the posterior correlations between the variables. reGVAMP postulates a factored posterior approximation of the form

$$q_{\mathbf{x},\mathbf{z}}(\mathbf{x}, \mathbf{z}) = \prod_i q_{x_i|\mathbf{y}}(x_i) \prod_j q_{z_j|\mathbf{y}}(z_j) \\ = \prod_i q_{x_i}(x_i)m_{x_i}(x_i) \prod_j q_{z_j}(z_j)m_{z_j}(z_j), \quad (12)$$

where q_{x_i} and q_{z_j} are the Gaussian approximations for the priors while m_{x_i} and m_{z_j} are the Gaussian extrinsics for x_i and z_j .

A byproduct are non-Gaussian posterior marginals, e.g. of the form $m_i(x_i)p(x_i)$ where $p(x_i)$ is the true prior for x_i . Note that involving the true priors is something that could also be considered in Variational Bayes (VB) [1]. reVAMP attempts to optimize the better $\text{KLD}(p, q)$ whereas VB optimizes $\text{KLD}(q, p)$.

So, reGVAMP performs alternating minimization of the following KLD

$$\arg \min_{q_{\mathbf{x},\mathbf{z}|\mathbf{y}}} \text{KLD}(p(\mathbf{x}, \mathbf{z}|\mathbf{y})\|q_{\mathbf{x},\mathbf{z}|\mathbf{y}}(\mathbf{x}, \mathbf{z})), \quad (13)$$

with the approximate posterior as in (12). The KLD becomes

$$\begin{aligned}
& \text{KLD}(p(\mathbf{x}, \mathbf{z}|\mathbf{y})\|q_{\mathbf{x}|\mathbf{y}}(\mathbf{x})) + \text{KLD}[p(\mathbf{x}, \mathbf{z}|\mathbf{y})\|q_{\mathbf{z}|\mathbf{y}}(\mathbf{z})] + c^t \\
&= \sum_i \text{KLD}(p(\mathbf{x}, \mathbf{z}|\mathbf{y})\|q_{x_i|\mathbf{y}}(x_i)) \\
&+ \sum_j \text{KLD}(p(\mathbf{x}, \mathbf{z}|\mathbf{y})\|q_{z_j|\mathbf{y}}(z_j)) + c^t \\
&= \sum_i \text{KLD}(p(x_i|\mathbf{y})\|q_{x_i|\mathbf{y}}(x_i)) \\
&+ \sum_j \text{KLD}(p(z_j|\mathbf{y})\|q_{z_j|\mathbf{y}}(z_j)) + c^t
\end{aligned} \tag{14}$$

where c^t denotes some constant. In the last equality, we marginalized out the irrelevant variables. The marginalized posteriors $p(x_i|\mathbf{y})$ and $p(z_j|\mathbf{y})$ are

$$p(x_i|\mathbf{y}) \propto \underbrace{p_{x_i}(x_i)}_{\text{prior}} \underbrace{\int p(\mathbf{y}|\mathbf{z})p(\mathbf{z}|\mathbf{x}) \prod_{k \neq i} p_{x_k}(x_k) d\mathbf{z} d\mathbf{x}_{\bar{i}}}_{\text{extrinsic } p(\mathbf{y}|x_i)} \tag{15}$$

$$\begin{aligned}
p(z_j|\mathbf{y}) &\propto p(\mathbf{y}, z_j) = \int p(\mathbf{y}, \mathbf{z}) d\mathbf{z}_{\bar{j}} \\
&= \underbrace{p_{y_j|z_j}(z_j)}_{\text{prior}} \underbrace{\int \prod_{k \neq j} p_{y_k|z_k}(z_k) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) p(\mathbf{x}) d\mathbf{x} d\mathbf{z}_{\bar{j}}}_{\text{extrinsic } p(\mathbf{y}_{\bar{j}}, z_j)}.
\end{aligned} \tag{16}$$

In order to see which probability the extrinsic for z corresponds to, consider the short hand notation

$$p(\mathbf{z}) = \int \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = p(\mathbf{z}_{\bar{j}}|z_j) p(z_j) \tag{17}$$

which depends only on the prior for \mathbf{x} . Therefore, in (16),

$$\begin{aligned}
& \int \prod_{k \neq j} p_{y_k|z_k}(z_k) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\
&= p_{\mathbf{y}_{\bar{j}}|\mathbf{z}_{\bar{j}}}(\mathbf{z}_{\bar{j}}) p(\mathbf{z}_{\bar{j}}|z_j) p(z_j) = p(\mathbf{y}_{\bar{j}}, \mathbf{z}_{\bar{j}}, z_j),
\end{aligned} \tag{18}$$

Thus, we have

$$\begin{aligned}
p(x_i|\mathbf{y}) &\simeq p_{x_i}(x_i) m_{x_i}(x_i), \\
p(z_j|\mathbf{y}) &\simeq p_{y_j|z_j}(z_j) m_{z_j}(z_j).
\end{aligned} \tag{19}$$

Due to the CLT, the extrinsics can be approximated as Gaussian when system dimensions increase. The marginal KLDs become

$$\begin{aligned}
& \arg \min_{q_{x_i|\mathbf{y}}} \text{KLD}(p(x_i|\mathbf{y})\|q_{x_i|\mathbf{y}}(x_i)) \\
&\simeq \arg \min_{q_{x_i}} \text{KLD}(p_{x_i}(x_i) m_{x_i}(x_i)\|q_{x_i}(x_i) m_{x_i}(x_i)),
\end{aligned} \tag{20}$$

$$\begin{aligned}
& \arg \min_{q_{z_j|\mathbf{y}}} \text{KLD}(p(z_j|\mathbf{y})\|q_{z_j|\mathbf{y}}(z_j)) \\
&\simeq \arg \min_{q_{z_j}} \text{KLD}(p_{y_j|z_j}(z_j) m_{z_j}(z_j)\|q_{z_j}(z_j) m_{z_j}(z_j)).
\end{aligned} \tag{21}$$

3.1 reGVAMP from (Minka) EP

We can arrive at the same point (20),(21) by Minka-style EP. Approximate p by q at factor level, with

$$\begin{aligned}
p(\mathbf{x}, \mathbf{z}|\mathbf{y}) &= 1/Z_p \prod_i p_{x_i}(x_i) \prod_j p_{y_j|z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}), \\
q(\mathbf{x}, \mathbf{z}) &= 1/Z_q \prod_i q_{x_i}(x_i) \prod_j q_{z_j}(z_j) m(\mathbf{x}, \mathbf{z}).
\end{aligned} \tag{22}$$

What is $m(\mathbf{x}, \mathbf{z})$? The tilted pdf $\tilde{p}_\delta = 1/Z_q \prod_i q_{x_i}(x_i) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x})$ is already Gaussian, hence is unchanged after Gaussian projection. So we can take $m(\mathbf{x}, \mathbf{z}) = \delta(\mathbf{z} - \mathbf{A}\mathbf{x})$ and we get $q(\mathbf{x}, \mathbf{z}) = 1/Z_q \prod_i q_{x_i}(x_i) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x})$. For the optimization of a factor $q_{x_i}(x_i)$, fit a Gaussian to the *tilted/target pdf*

$$\tilde{p}_{x_i}(\mathbf{x}, \mathbf{z}) = 1/Z_{\tilde{p}_{x_i}} p_{x_i}(x_i) \prod_{k \neq i} q_{x_k}(x_k) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) \quad (23)$$

We get:

$$\begin{aligned} & \operatorname{argmin}_{f(x_i)} \operatorname{KLD}(\tilde{p}_{x_i} \| q) \\ &= \operatorname{arg} \min_{q_{x_i}(x_i)} \operatorname{KLD}(p(x_k) \prod_{k \neq i} q_{x_k}(x_k) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) \| \\ & \quad \prod_k q_{x_k}(x_k) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x})) = \\ & \operatorname{arg} \min_{q_{x_i}(x_i)} \int p(x_i) \prod_{k \neq i} q_{x_k}(x_k) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) \ln\left(\frac{p(x_i)}{q_{x_i}(x_i)}\right) d\mathbf{x} d\mathbf{z} \\ &= \operatorname{arg} \min_{q_{x_i}(x_i)} \int p(x_i) \ln\left(\frac{p(x_i)}{q_{x_i}(x_i)}\right) \\ & \quad \left[\int \prod_{k \neq i} q_{x_k}(x_k) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) d\mathbf{x}_k d\mathbf{z} \right] dx_i \\ & \quad = m_{x_i}(x_i) \text{ Gaussian extrinsic} \\ &= \operatorname{arg} \min_{q_{x_i}(x_i)} \int p_{x_i}(x_i) m_{x_i}(x_i) \ln \frac{p_{x_i}(x_i) m_{x_i}(x_i)}{q_{x_i}(x_i) m_{x_i}(x_i)} dx_i \\ &= \operatorname{arg} \min_{q(x_i)/m_{x_i}(x_i)} \int p_{x_i}(x_i) m_{x_i}(x_i) \ln \frac{p_{x_i}(x_i) m_{x_i}(x_i)}{q(x_i)} dx_i \end{aligned} \quad (24)$$

Since $m_{x_i}(x_i)$ is Gaussian, it will suffice to fit a Gaussian in x_i , say $q(x_i)$, via

$$\begin{aligned} \operatorname{KLD}(p(x_i|\mathbf{y}) \| q(x_i)) &= \operatorname{KLD}(p_{x_i}(x_i) p(\mathbf{y}|x_i) / Z_i \| q(x_i)) \\ &\approx \operatorname{KLD}(p_{x_i}(x_i) m_{x_i}(x_i) / Z_i \| q(x_i)) \\ &= \operatorname{KLD}(p_{x_i}(x_i) m_{x_i}(x_i) / Z_i \| q_{x_i}(x_i) m_{x_i}(x_i) / Z'_i). \end{aligned} \quad (25)$$

The reVAMP algorithm [16] approximates the posterior to Gaussian with the approximated Gaussian extrinsic:

$$p(x_i|\mathbf{y}) \approx \frac{p_{x_i}(x_i) m_{x_i}(x_i)}{Z_{x_i}(\mathbf{y})} \approx \mathcal{N}(x_i; \hat{x}_i, \tau_{x_i}) = q(x_i). \quad (26)$$

where $m_{x_i}(x_i) = \mathcal{N}(x_i; r_i, \tau_{r_i})$. The approximate Gaussian posterior $q(x_i)$ is obtained by moment matching with the better posterior approximation $p_{x_i}(x_i) m_{x_i}(x_i) / Z_{x_i}$.

We interpret the quotient of the approximated posterior and the approximate extrinsic as the approximated Gaussian prior.

$$\begin{aligned} p_{x_i}(x_i) &\approx q_{x_i}(x_i) = \mathcal{N}(x_i; m_{x_i}, \sigma_{x_i}^2) \propto \frac{\mathcal{N}(x_i; \hat{x}_i, \tau_{x_i})}{\mathcal{N}(x_i; r_i, \tau_{r_i})}, \\ 1/\sigma_{x_i}^2 &= 1/\tau_{x_i} - 1/\tau_{r_i}, \quad m_{x_i} = \sigma_{x_i}^2 (\hat{x}_i/\tau_{x_i} - r_i/\tau_{r_i}). \end{aligned} \quad (27)$$

This Gaussian approximation $q_{x_i}(x_i)$ does not correspond to direct moment matching of the true prior $p_{x_i}(x_i)$. So, reGVAMP admits two points of view:

- (1) minimize $\operatorname{KLD}(p \| q)$ with $q = \prod_i q(x_i) \prod_j q(z_j)$
- (2) do Minka EP with $q = \prod_i q_{x_i}(x_i) \prod_j q_{z_j}(z_j) \delta(\mathbf{z} - \mathbf{A}\mathbf{x})$

Both points of view lead to the same results!

The sense of the Gaussian prior approximations $q_{x_i}(x_i)$, $q_{z_j}(z_j)$ is that they are the equivalent Gaussian priors that, in the presence of the Gaussian extrinsics, produce the exact (nonlinear) MMSE estimate and variance that the original non-Gaussian prior would do! Direct Gaussian approximation of the priors is very suboptimal because that would only produce the correct LMMSE

estimate and variance!!!

In the case the true priors are Gaussian, the two are the same of course.

Apart from the *improved marginal posteriors* $m_{x_i}(x_i)p_{x_i}(x_i)/Z'_i$ (and similar for the z_j), together with $\delta(\mathbf{z} - \mathbf{A}\mathbf{x})$, the Gaussian prior approximations $q_{x_i}(x_i)$, $q_{z_j}(z_j)$ in reGVAMP lead to an equivalent overall Gaussian linear model. This can be used for Large System Analysis (random \mathbf{A} model) for the resulting posterior variances (MSEs), as obtained by GAMP. reVAMP does *alternating minimization of KLD(p, q)* which becomes iterative because an extrinsic $m_{x_i}(x_i)$ depends on the approximate Gaussian priors $\prod_{j \neq i} q_{x_j}(x_j)$, $\prod_j q_{z_j}(z_j)$. Since alternating minimization of a convex cost function converges, reVAMP can be expected to converge.

The Gaussian extrinsics approximations $p(x_i|\mathbf{y}) \approx m_i(x_i)$ are *asymptotically tight*. The Gaussian approximations that are not tight and that constitute the variational approximations are approximating marginal posteriors by Gaussian $q(x_i)$ or what follows from that, approximating priors $p_{x_i}(x_i)$ by Gaussian $q_{x_i}(x_i)$. Or the overall multivariate Gaussian posterior approximation is not tight also, but at least *captures full second-order moments*.

Hence in one point of view, re(G)VAMP minimized the desirable KLD(p, q), which becomes feasible thanks to asymptotic Gaussianity of the extrinsics like $p(\mathbf{y}|x_i) \approx m_{x_i}(x_i)$. However, re(G)VAMP can also be derived using EP, using a different formal posterior approximation.

4 Relation to CWCU MMSE Estimator

The algorithm proposed by [16] can be interpreted as an iterative method of finding consistent extrinsic and posterior messages for the case of a AWGN $p(\mathbf{y}|\mathbf{z})$. [16] also shows the close relation between CWCU LMMSE estimation [13] and the extrinsic. In the following, we will interpret the extrinsic as CWCU LMMSE estimation based on the Gauss-Markov theorem.

Based on the discussion of the previous section, when deriving the extrinsic for \mathbf{z} and \mathbf{x} , we find the system to be equivalent to a Gaussian linear model. Therefore, we can use the approximate prior and approximate likelihood as if they are the true prior and likelihood when deriving the extrinsics without large system approximations [14].

Consider jointly Gaussian \mathbf{y} and x (scalar)

$$\begin{bmatrix} \mathbf{y} \\ x \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_y \\ m_x \end{bmatrix}, \begin{bmatrix} \mathbf{C}_{yy} & \mathbf{C}_{yx} \\ \mathbf{C}_{xy} & C_{xx} \end{bmatrix} \right) \quad (28)$$

Then the extrinsic $p(\mathbf{y}|x)$ is Gaussian and based on Gaussi-Markov theorem

$$\begin{aligned} -2 \ln p(\mathbf{y}|x) &= c + (\mathbf{y} - \mathbf{m}_{y|x})^T \mathbf{C}_{y|x}^{-1} (\mathbf{y} - \mathbf{m}_{y|x}), \text{ with} \\ \mathbf{m}_{y|x} &= \mathbf{m}_y + \mathbf{C}_{yx} C_{xx}^{-1} (x - m_x), \\ \mathbf{C}_{y|x} &= \mathbf{C}_{yy} - \mathbf{C}_{yx} \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \end{aligned} \quad (29)$$

Interpreting (29) as a pdf in x (which Fisher called fiducial statistics), we can rewrite this quadratic exponent as

$$\begin{aligned} -2 \ln p(\mathbf{y}|x) &= c(\mathbf{y}) + (x - \hat{x}_{CL})^2 / \mathbf{C}_{\tilde{x}_{CL}\tilde{x}_{CL}}, \\ \hat{x}_{CL} &= m_x + d \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} (\mathbf{y} - \mathbf{m}_y) = d \hat{x}_L + (1 - d) m_x \\ \mathbf{C}_{\tilde{x}_{CL}\tilde{x}_{CL}} &= d \mathbf{C}_{\tilde{x}_L\tilde{x}_L}, \\ \text{with} & \\ \hat{x}_L &= m_x + \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} (\mathbf{y} - \mathbf{m}_y), \quad \mathbf{C}_{\tilde{x}_L\tilde{x}_L} = C_{xx} - \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} \mathbf{C}_{yx} \\ d &= \frac{C_{xx}}{\mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} \mathbf{C}_{yx}} \geq 1, \end{aligned} \quad (30)$$

where \hat{x}_{CL} , $\mathbf{C}_{\tilde{x}_{CL}\tilde{x}_{CL}}$ are the CWCU LMMSE estimate and error variance, and \hat{x}_L , $\mathbf{C}_{\tilde{x}_L\tilde{x}_L}$ are the LMMSE (and hence MMSE since Gaussian) estimate and error variance.

Now we will investigate the vector case. Define the operation $Diag(\mathbf{C}) = \text{diag}[\text{diag}(\mathbf{C})]$, which returns a diagonal matrix from the vector $\text{diag}(\mathbf{C})$, composed of the diagonal elements of square matrix \mathbf{C} .

Interpreting the previous x as a component x_i of a vector \mathbf{x} , we can write

$$\begin{aligned}\hat{\mathbf{x}}_{CL} &= \mathbf{m}_x + \mathbf{D} \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} (\mathbf{y} - \mathbf{m}_y) = \mathbf{D} \hat{\mathbf{x}}_L + (\mathbf{I} - \mathbf{D}) \mathbf{m}_x \\ \mathbf{C}_{\tilde{\mathbf{x}}_{CL} \tilde{\mathbf{x}}_{CL}} &= \mathbf{C}_{\tilde{\mathbf{x}}_L \tilde{\mathbf{x}}_L} + (\mathbf{D} - \mathbf{I}) \mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L} (\mathbf{D} - \mathbf{I}) \\ &\text{with} \\ \mathbf{D} &= \text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}) [\text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L})]^{-1}, \mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L} = \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} \mathbf{C}_{yx}\end{aligned}\quad (31)$$

where the expression for $\mathbf{C}_{\tilde{\mathbf{x}}_{CL} \tilde{\mathbf{x}}_{CL}}$ follows from

$$\tilde{\mathbf{x}}_{CL} = \mathbf{x} - \hat{\mathbf{x}}_{CL} = \tilde{\mathbf{x}}_L - (\mathbf{D} - \mathbf{I}) \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} (\mathbf{y} - \mathbf{m}_y), \quad (32)$$

and the two terms in this difference are decorrelated by the orthogonality property of LMMSE estimation.

Next, we'll show: $\mathbf{D} = \text{diag}(\boldsymbol{\tau}_{CL} / \boldsymbol{\tau}_L)$, where $\boldsymbol{\tau}_L = \text{diag}(\mathbf{C}_{\tilde{\mathbf{x}}_L \tilde{\mathbf{x}}_L})$ and $\boldsymbol{\tau}_{CL} = \text{diag}(\mathbf{C}_{\tilde{\mathbf{x}}_{CL} \tilde{\mathbf{x}}_{CL}})$, and " ./ " denotes element-wise division.

$$\begin{aligned}\mathbf{C}_{\tilde{\mathbf{x}}_{CL} \tilde{\mathbf{x}}_{CL}} &= \mathbf{C}_{\tilde{\mathbf{x}}_L \tilde{\mathbf{x}}_L} + (\mathbf{D} - \mathbf{I}) \mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L} (\mathbf{D} - \mathbf{I}) \\ &= \mathbf{C}_{\mathbf{x}\mathbf{x}} - \mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L} \mathbf{D} - \mathbf{D} \mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L} + \mathbf{D} \mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L} \mathbf{D}\end{aligned}\quad (33)$$

Calculate the diagonal elements

$$\begin{aligned}\text{diag}(\boldsymbol{\tau}_{CL}) &= \text{Diag}(\mathbf{C}_{\tilde{\mathbf{x}}_{CL} \tilde{\mathbf{x}}_{CL}}) = \text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}) \\ &+ \mathbf{D} \text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L}) \mathbf{D} - \text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L}) \mathbf{D} - \mathbf{D} \text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L}) \\ &= \text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}) [\text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L})]^{-1} \text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}) - \text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}),\end{aligned}\quad (34)$$

where we use $\mathbf{D} = \text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}) [\text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L})]^{-1}$ in (31).

Now we want to show $\mathbf{D} \text{diag}(\boldsymbol{\tau}_L) = \text{diag}(\boldsymbol{\tau}_{CL})$:

$$\begin{aligned}\mathbf{D} \text{diag}(\boldsymbol{\tau}_L) &= \mathbf{D} \text{Diag}(\mathbf{C}_{\tilde{\mathbf{x}}_L \tilde{\mathbf{x}}_L}) \\ &= \text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}) [\text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L})]^{-1} \\ &\cdot [\text{Diag}(\mathbf{C}_{\mathbf{x}\mathbf{x}}) - \text{Diag}(\mathbf{C}_{\hat{\mathbf{x}}_L \hat{\mathbf{x}}_L})] = \text{diag}(\boldsymbol{\tau}_{CL})\end{aligned}\quad (35)$$

The extrinsic for \mathbf{x} without large system approximations can be interpreted as CWCU MMSE estimation from the Gaussian model

$$\begin{bmatrix} \mathbf{m}_z \\ \mathbf{x} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{A} \mathbf{m}_x \\ \mathbf{m}_x \end{bmatrix}, \begin{bmatrix} \mathbf{A} \mathbf{D}_{\sigma_x^2} \mathbf{A}^T + \mathbf{D}_{\sigma_z^2} & \mathbf{A} \mathbf{D}_{\sigma_x^2} \\ \mathbf{D}_{\sigma_x^2} \mathbf{A}^T & \mathbf{D}_{\sigma_x^2} \end{bmatrix} \right). \quad (36)$$

The underlying equivalent Gaussian linear model is

$$\mathbf{m}_z = \mathbf{A} \mathbf{x} + \mathbf{v}_z \quad (37)$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_x, \mathbf{D}_{\sigma_x^2})$ and $\mathbf{v}_z \sim \mathcal{N}(\mathbf{0}, \mathbf{D}_{\sigma_z^2})$.

Likewise, we can interpret the extrinsic for \mathbf{z} as CWCU MMSE estimation from

$$\begin{bmatrix} \mathbf{A} \mathbf{m}_x \\ \mathbf{z} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_z \\ \mathbf{m}_z \end{bmatrix}, \begin{bmatrix} \mathbf{D}_{\sigma_z^2} + \mathbf{A} \mathbf{D}_{\sigma_x^2} \mathbf{A}^T & \mathbf{D}_{\sigma_z^2} \\ \mathbf{D}_{\sigma_x^2} \mathbf{A}^T & \mathbf{D}_{\sigma_x^2} \end{bmatrix} \right). \quad (38)$$

The underlying equivalent Gaussian linear model is

$$\mathbf{A} \mathbf{m}_x = \mathbf{z} + \mathbf{v}_z \quad (39)$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{m}_z, \mathbf{D}_{\sigma_z^2})$ and $\mathbf{v}_z \sim \mathcal{N}(\mathbf{0}, \mathbf{A} \mathbf{D}_{\sigma_x^2} \mathbf{A}^T)$.

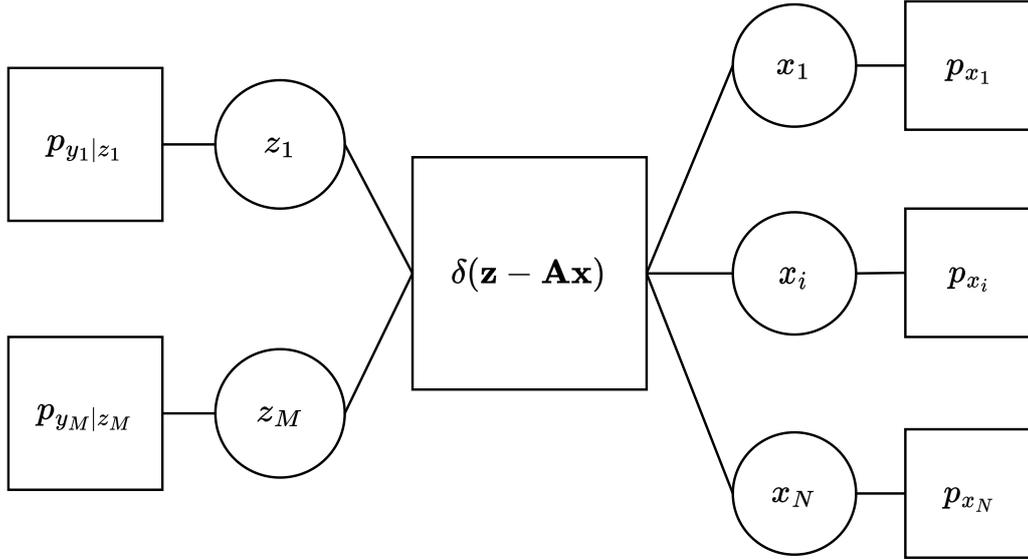


Figure 1: Factor Graph for the GLM used by reGVAMP. Circles: variable nodes, squares: factor nodes.

5 GAMP from LSL Belief Propagation

In reGVAMP, extrinsics in the GLM are built from the equivalent Gaussian linear model, which introduces equivalent Gaussian priors from Gaussian posterior approximations and Gaussian extrinsics.

GAMP exploits LSL simplifications of reGVAMP for a random \mathbf{A} with i.i.d. signs which leads to

- (i) Gaussianity of extrinsics (also in reGVAMP), and
- (ii) independence of marginals (extra w.r.t. reGVAMP).

(ii) leads to the large system simplifications of the variances, avoiding covariance matrix inverses. But also posterior and extrinsic estimates $\hat{\mathbf{x}}$, $\hat{\mathbf{z}}$ and \mathbf{r} , \mathbf{p} that are constructed by combining decoupled pieces of information. These estimates are non-linear MMSE and CWCU MMSE estimates in general. Extrinsics are not obtained as linear perturbations of corresponding MMSE estimates because those are not necessarily close to each other. Rather the interplay between \mathbf{x} and \mathbf{z} is exploited with perturbations due to the small effect of a single term in \mathbf{A} in the LSL. In both reGVAMP and GAMP, we have:

Gaussian extrinsics: $\mathcal{N}(\mathbf{x}; \mathbf{r}, \boldsymbol{\tau}_r)$, $\mathcal{N}(\mathbf{z}; \mathbf{p}, \boldsymbol{\tau}_p)$, and

Posterior marginals proportional to: $p_{\mathbf{x}}(\mathbf{x})\mathcal{N}(\mathbf{x}; \mathbf{r}, \boldsymbol{\tau}_r)$, $p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z})\mathcal{N}(\mathbf{z}; \mathbf{p}, \boldsymbol{\tau}_p)$ with Gaussian approximations $\mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \boldsymbol{\tau}_x)$, $\mathcal{N}(\mathbf{z}; \hat{\mathbf{z}}, \boldsymbol{\tau}_z)$.

reGVAMP considers the joint pdf factorization into $M + N + 1$ factors

$$p(\mathbf{x}, \mathbf{z}, \mathbf{y}) = \delta(\mathbf{z} - \mathbf{A}\mathbf{x}) \prod_{i=1}^N p_{x_i}(x_i) \prod_{k=1}^M p_{y_k|z_k}(y_k|z_k) \quad (40)$$

where $\delta(\mathbf{z} - \mathbf{A}\mathbf{x}) = \prod_{k=1}^M \delta(z_k - \mathbf{a}_k^T \mathbf{x})$, $\mathbf{A}^T = [\mathbf{a}_1 \cdots \mathbf{a}_M]$. The factor graph in Fig. 1 is without cycles. The factor graph considered determines the associated Belief or Expectation Propagation algorithms for minimizing the Bethe Free Energy [5]. GVAMP on the other hand considers the following joint pdf factorization into $2M + N$ factors

$$p(\mathbf{x}, \mathbf{z}, \mathbf{y}) = \prod_{i=1}^N p_{x_i}(x_i) \prod_{k=1}^M p_{y_k|z_k}(y_k|z_k) \delta(z_k - \mathbf{a}_k^T \mathbf{x}) \quad (41)$$

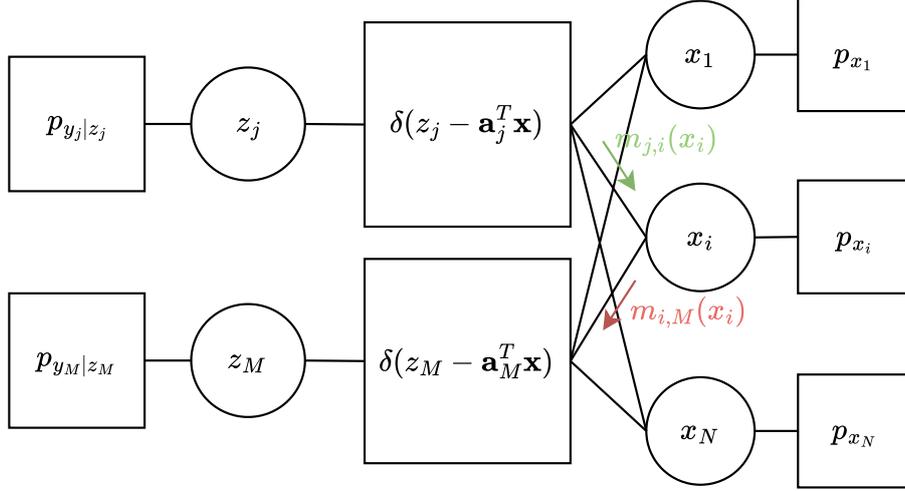


Figure 2: Factor Graph for the GLM used by GVAMP.

which leads to the factor graph in Fig. 2 which does contain cycles.

Message passing in the GLM scalar level factor graph of Fig. 2 alternates between the following message updates:

$$\begin{aligned}
 m_{k,n}(x_n) &\sim \int p(y_k|z_k) \delta(z_k - \mathbf{A}_{k,:} \mathbf{x}) \prod_{m \neq n} m_{m,k}(x_m) dz_k d\mathbf{x}_{\bar{n}} \\
 m_{n,k}(x_n) &\sim p_{x_n}(x_n) \prod_{i \neq k} m_{i,n}(x_n)
 \end{aligned} \tag{42}$$

where \sim denotes equality up to a normalization factor. This results in:

marginal posteriors: $m_n(x_n) \sim p_{x_n}(x_n) \prod_i m_{i,n}(x_n)$,
 extrinsic $z_k : \sim \int \delta(z_k - \mathbf{A}_{k,:} \mathbf{x}) \prod_n m_{n,k}(x_n) d\mathbf{x}$,
 extrinsic $x_n : \sim \prod_i m_{i,n}(x_n)$.

Like reGVAMP, GAMP uses Gaussian approximations for extrinsics. This requires Gaussian models for the messages. GAMP applies Gaussian approximations in 2 steps: (middle expression = prior \times Gaussian extrinsic)

$$\begin{aligned}
 m_{k,n}(x_n) &\rightarrow \int p(y_k|z_k) \delta(z_k - \mathbf{A}_{k,:} \mathbf{x}) \prod_{m \neq n} q_{m,k}(x_m) dz_k d\mathbf{x}_{\bar{n}} \\
 &\rightarrow q_{k,n}(x_n) = \mathcal{N}(x_n; \hat{x}_{k,n}, \hat{\tau}_{k,n}^x)
 \end{aligned} \tag{43}$$

$$\begin{aligned}
 m_{n,k}(x_n) &\rightarrow p_{x_n}(x_n) \prod_{i \neq k} q_{i,n}(x_n) \\
 &\rightarrow q_{n,k}(x_n) = \mathcal{N}(x_n; \hat{x}_{n,k}, \hat{\tau}_{n,k}^x)
 \end{aligned} \tag{44}$$

5.1 Output Node

We get for the incomplete extrinsic for z_k :

$$\int \delta(z_k - \mathbf{A}_{k,:} \mathbf{x}) \prod_{m \neq n} q_{m,k}(x_m) d\mathbf{x}_{\bar{n}} \sim \mathcal{N}(z_k; p_{k,n} + \mathbf{A}_{k,n} x_n, \hat{\tau}_{k,n}^p) \tag{45}$$

with

$$\begin{aligned} p_{k,n} &= \mathbf{A}_{k,\bar{n}} \widehat{\mathbf{x}}_{\bar{n},k} \\ \tau_{k,n}^p &= \mathbf{S}_{k,\bar{n}} \tau_{\bar{n},k}^x \approx \mathbf{S}_{k,\bar{n}} \tau_{\bar{n}}^x. \end{aligned} \quad (46)$$

Define $p_k = \mathbf{A}_{k,:} \widehat{\mathbf{x}}_{:,k} \Rightarrow p_{k,n} = p_k - \mathbf{A}_{k,n} \widehat{\mathbf{x}}_{n,k}$.

And $\tau_{k,n}^p = \tau_k^p - \mathbf{S}_{k,n} \tau_{n,k}^x$ where $\tau_k^p = \mathbf{S}_{k,:} \tau_x$.

Neglecting terms of order $\mathbf{S}_{k,n}$, we get

$$\mathcal{N}(z_k; p_{k,n} + \mathbf{A}_{k,n} x_n, \tau_{k,n}^p) \approx \mathcal{N}(z_k; p_k + \mathbf{A}_{k,n} \tilde{x}_n, \tau_k^p), \quad (47)$$

with $\tilde{x}_n = x_n - \widehat{x}_n$.

Then $m_{k,n}(x_n) \approx Z_z(p_k + \mathbf{A}_{k,n} \tilde{x}_n, y_k, \tau_k^p)$ with

$$\begin{aligned} Z_z(p, y, \tau_p) &= \int p_{y|z}(y|z) e^{-\frac{1}{2\tau_p}(z-p)^2} dz \\ \frac{\partial \ln Z_z}{\partial p} &= \frac{Z'_z}{Z_z} = s = \frac{\widehat{z}-p}{\tau_p}, \widehat{z} = \frac{1}{Z_z} \int z p_{y|z}(y|z) e^{-\frac{1}{2\tau_p}(z-p)^2} dz \\ \frac{\partial^2 \ln Z_z}{\partial p^2} &= -\tau_s = -\tau_s = \frac{Z''_z}{Z_z} - \left(\frac{Z'_z}{Z_z}\right)^2 = -(1 - \tau_z/\tau_p)/\tau_p \end{aligned}$$

Then up to second order in $\mathbf{A}_{k,n} \tilde{x}_n$ (Laplacian approximation in MAP case, Gaussian moment matching in MMSE case), a single measurement extrinsic for x_n becomes: $\ln m_{k,n}(x_n)$

$$\begin{aligned} &\approx \ln Z_z(p_k, y_k, \tau_k^p) + \frac{\partial \ln Z_z}{\partial p} \mathbf{A}_{k,n} \tilde{x}_n + \frac{\partial^2 \ln Z_z}{2 \partial p^2} \mathbf{A}_{k,n}^2 \tilde{x}_n^2 \\ &= c^t + [\mathbf{s}_k \mathbf{A}_{k,n} + \mathbf{A}_{k,n}^2 \tau_k^s \widehat{x}_n] x_n - \frac{1}{2} \tau_k^s \mathbf{A}_{k,n}^2 x_n^2. \end{aligned}$$

Now

$$\begin{aligned} \ln m_{n,k}(x_n) &= c^t + \ln p_{x_n}(x_n) + \sum_{i \neq k} \ln m_{i,n}(x_n) \\ &= c^t + \ln p_{x_n}(x_n) - \frac{1}{2\tau_{n,k}^r} (x_n - r_{n,k})^2 \end{aligned} \quad (48)$$

$$\text{with } \frac{1}{\tau_{n,k}^r} = \mathbf{S}_{k,n}^T \tau_k^s \quad (\approx \mathbf{S}_{:,n}^T \tau_s = \frac{1}{\tau_n^r})$$

$$\text{and } r_{n,k} = \tau_{n,k}^r (\mathbf{s}_k^T \mathbf{A}_{k,n} + \mathbf{S}_{k,n}^T \tau_k^s \widehat{x}_n) = \widehat{x}_n + \tau_{n,k}^r \mathbf{s}_k^T \mathbf{A}_{k,n}.$$

5.2 Input Node

We now get for the approximate posterior

$$m_n(x_n) = \frac{1}{Z_x(r_n, \tau_n^r)} p_{x_n}(x_n) e^{-\frac{1}{\tau_n^r} (\frac{x_n^2}{2} - x_n r_n)}, \quad (49)$$

with

$$\begin{aligned} Z_x(r, \tau_r) &= \int p_x(x) e^{-\frac{1}{\tau_r} (\frac{x^2}{2} - x r)} dx \\ \tau_r \frac{\partial \ln Z_x}{\partial r} &= \int x m(x) dx = \mathbb{E}(x|r, \tau_r) = \widehat{x} = \widehat{x}(r, \tau_r) \\ \tau_r^2 \frac{\partial^2 \ln Z_x}{\partial r^2} &= \tau_r \frac{\partial \widehat{x}}{\partial r} = \tau_x \end{aligned}$$

Now, with

$$r_n = \widehat{x}_n + \tau_n^r \mathbf{s}^T \mathbf{A}_{:,n}, \quad (50)$$

we can write

$$r_{n,k} \approx \widehat{x}_n + \tau_n^r \mathbf{s}_k^T \mathbf{A}_{k,n} = r_n - \tau_n^r s_k \mathbf{A}_{k,n}. \quad (51)$$

We get similarly for the mean $\widehat{x}_{n,k}$ of $m_{n,k}(x_n)$:

$$\begin{aligned}\widehat{x}_{n,k} &= \widehat{x}_n(r_{n,k}, \tau_n^r) = \widehat{x}_n(r_n - \tau_n^r s_k \mathbf{A}_{k,n}, \tau_n^r) \\ &\approx \widehat{x}_n(r_n, \tau_n^r) - \frac{\partial}{\partial r_n} \widehat{x}_n(r_n, \tau_n^r) \tau_n^r s_k \mathbf{A}_{k,n} = \widehat{x}_n - \tau_n^x s_k \mathbf{A}_{k,n}\end{aligned}$$

Plugging this in, we get

$$p_k = \mathbf{A}_{k,:} \widehat{\mathbf{x}}_{:,k} = \mathbf{A}_{k,:} \widehat{\mathbf{x}} - \mathbf{S}_{k,:} \boldsymbol{\tau}_x s_k = \mathbf{A}_{k,:} \widehat{\mathbf{x}} - \tau_k^p s_k$$

which completes the message passing. We may note that the variance derivations in the LSL of BP are equivalent to the large random matrix analysis of the MSE of LMMSE in the equivalent Gaussian linear model.

6 LSL BFE and EP

After the LSL simplifications [12], the BFE from (11) with marginal pdf consistency constraints can be seen to become equivalent to the following LSL-BFE [8], [10] :

$$\begin{aligned}\min_{q_{\mathbf{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p, \mathbf{u}} D[q_{\mathbf{x}} \| p_{\mathbf{x}}] + D[q_{\mathbf{z}} \| p_{\mathbf{y}|\mathbf{z}}] + \frac{1}{2} \sum_k \left[\frac{\text{var}(z_k | q_{\mathbf{z}})}{\tau_{p_k}} + \ln(\tau_{p_k}) \right] \\ \text{s.t. } E[\mathbf{z} | q_{\mathbf{z}}] = \mathbf{A} \mathbf{u} \\ E[\mathbf{x} | q_{\mathbf{x}}] = \mathbf{u} \\ \boldsymbol{\tau}_p = \mathbf{S} \text{var}(\mathbf{x} | q_{\mathbf{x}}).\end{aligned}\tag{52}$$

We will exploit some useful relations

$$\begin{aligned}\forall \boldsymbol{\tau}, \mathbf{c}^T \text{var}(\mathbf{x} | q_{\mathbf{x}}) &= \int \|\mathbf{x} - \mathbf{u}\|_{\mathbf{1}, \boldsymbol{\tau}}^2 q_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \\ \sum_k \frac{\text{var}(z_k | q_{\mathbf{z}})}{\tau_{p_k}} &= \int \|\mathbf{z} - \mathbf{A} \mathbf{u}\|_{\boldsymbol{\tau}_p}^2 q_{\mathbf{z}}(\mathbf{z}) d\mathbf{z}.\end{aligned}\tag{53}$$

The Lagrangian of (52) becomes

$$\begin{aligned}L &= D[q_{\mathbf{x}} \| p_{\mathbf{x}}] + D[q_{\mathbf{z}} \| p_{\mathbf{y}|\mathbf{z}}] + \frac{1}{2} \sum_k \left[\frac{\text{var}(z_k | q_{\mathbf{z}})}{\tau_{p_k}} + \ln(\tau_{p_k}) \right] \\ &+ \boldsymbol{\lambda}_{\mu_{\mathbf{z}}}^T \left(\mathbf{A} \mathbf{u} - \int \mathbf{z} q_{\mathbf{z}}(\mathbf{z}) d\mathbf{z} \right) + \boldsymbol{\lambda}_{\mu_{\mathbf{x}}}^T \left(\mathbf{u} - \int \mathbf{x} q_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \right) \\ &- \frac{1}{2} \boldsymbol{\lambda}_{\boldsymbol{\tau}}^T (\boldsymbol{\tau}_p - \mathbf{S} \text{var}(\mathbf{x} | q_{\mathbf{x}}))\end{aligned}\tag{54}$$

The derivatives w.r.t. $q_{\mathbf{x}}, q_{\mathbf{z}}, \boldsymbol{\tau}_p, \mathbf{u}$ become

$$\begin{aligned}\frac{\partial L}{\partial q_{\mathbf{x}}} &= \ln(q_{\mathbf{x}}) - \ln(p_{\mathbf{x}}) - \boldsymbol{\lambda}_{\mu_{\mathbf{x}}}^T \mathbf{x} + \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_{\mathbf{1}, (\mathbf{S}^T \boldsymbol{\lambda}_{\boldsymbol{\tau}})}^2 \\ \frac{\partial L}{\partial q_{\mathbf{z}}} &= \ln(q_{\mathbf{z}}) - \ln(p_{\mathbf{y}|\mathbf{z}}) - \boldsymbol{\lambda}_{\mu_{\mathbf{z}}}^T \mathbf{z} + \frac{1}{2} \|\mathbf{z} - \mathbf{A} \mathbf{u}\|_{\boldsymbol{\tau}_p}^2 \\ \frac{\partial L}{\partial \tau_{p_k}} &\propto -\frac{\text{var}(z_k | q_{\mathbf{z}})}{\tau_{p_k}^2} + \frac{1}{\tau_{p_k}} - \lambda_{\tau_k} \\ \frac{\partial L}{\partial \mathbf{u}} &= -\mathbf{A}^T \mathbf{D}_{\boldsymbol{\tau}_p}^{-1} (\widehat{\mathbf{z}} - \mathbf{A} \mathbf{u}) + \mathbf{A}^T \boldsymbol{\lambda}_{\mu_{\mathbf{z}}} + \boldsymbol{\lambda}_{\mu_{\mathbf{x}}} \\ &- \mathbf{D}_{\mathbf{S}^T \boldsymbol{\lambda}_{\boldsymbol{\tau}}} (\widehat{\mathbf{x}} - \mathbf{u}),\end{aligned}\tag{55}$$

where $\hat{\mathbf{z}} = \mathbb{E}(\mathbf{z}|q_{\mathbf{z}})$ and $\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x}|q_{\mathbf{x}})$. Zeroing derivatives:

$$q_{\mathbf{x}}(\mathbf{x}) \propto p_{\mathbf{x}}(\mathbf{x}) e^{-\frac{1}{2}\|\mathbf{x}-\mathbf{u}\|_{1./(S^T\lambda_{\tau})}^2} e^{\lambda_{\mu_{\mathbf{x}}}^T \mathbf{x}} \quad (56)$$

$$q_{\mathbf{z}}(\mathbf{z}) \propto p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2}\|\mathbf{z}-\mathbf{A}\mathbf{u}\|_{\tau_p}^2} e^{\lambda_{\mu_{\mathbf{z}}}^T \mathbf{z}} \quad (57)$$

$$\lambda_{\tau_k} = \frac{1}{\tau_{p_k}} - \frac{\tau_{z_k}}{\tau_{p_k}^2} \quad (58)$$

$$\begin{aligned} \left[\mathbf{A}^T \mathbf{D}_{\tau_p}^{-1} \mathbf{A} + \mathbf{D}_{S^T \lambda_{\tau}} \right] \mathbf{u} &= \mathbf{A}^T \mathbf{D}_{\tau_p}^{-1} \hat{\mathbf{z}} \\ &+ \mathbf{D}_{S^T \lambda_{\tau}} \hat{\mathbf{x}} - \mathbf{A}^T \lambda_{\mu_{\mathbf{z}}} - \lambda_{\mu_{\mathbf{x}}} \end{aligned} \quad (59)$$

where $\tau_{z_k} = \mathbb{E}[(z_k - \hat{z}_k)^2 | q_{\mathbf{z}}]$. By satisfying the two mean constraints in (52), the equation (59) becomes

$$\mathbf{A}^T \lambda_{\mu_{\mathbf{z}}} = -\lambda_{\mu_{\mathbf{x}}} \quad (60)$$

A solution can be obtained by solving the system of 7 equations containing (56), (57), (58), (60) and the three constraint equations in (52).

7 Iterative Solution leading to GAMP

We ignore pdf normalization for simplicity. Furthermore, we use red symbols to indicate parameters to be updated.

7.1 Update of $\lambda_{\mu_{\mathbf{z}}}$

Consider (57) and the two mean constraints in (52)

$$\mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2}\|\mathbf{z}-\hat{\mathbf{z}}\|_{\tau_p}^2} e^{\lambda_{\mu_{\mathbf{z}}}^{(\text{new})T} \mathbf{z}} \right] \quad (61)$$

$$= \mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2}\|\mathbf{z}-\mathbf{A}\hat{\mathbf{x}}\|_{\tau_p}^2} e^{\lambda_{\mu_{\mathbf{z}}}^T \mathbf{z}} \right] = \hat{\mathbf{z}} \quad (62)$$

We first use (62) to obtain $\hat{\mathbf{z}}$. Then we use this newly obtained $\hat{\mathbf{z}}$ to update $\lambda_{\mu_{\mathbf{z}}}^{(\text{new})}$, since we need to keep the exponential factor identical in order not to change the mean, i.e.,

$$e^{-\frac{1}{2}\|\mathbf{z}-\hat{\mathbf{z}}\|_{\tau_p}^2} e^{\lambda_{\mu_{\mathbf{z}}}^{(\text{new})T} \mathbf{z}} = e^{-\frac{1}{2}\|\mathbf{z}-\mathbf{A}\hat{\mathbf{x}}\|_{\tau_p}^2} e^{\lambda_{\mu_{\mathbf{z}}}^T \mathbf{z}}. \quad (63)$$

If we want to bridge the GAMP from [8] and BFE, we can denote

$$\mathbf{p} = \mathbf{A}\hat{\mathbf{x}} + \mathbf{D}(\tau_p)\lambda_{\mu_{\mathbf{z}}}. \quad (64)$$

With definition (64), the expression (62) can be written as

$$\hat{\mathbf{z}} = \mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2}\|\mathbf{z}-\mathbf{p}\|_{\tau_p}^2} \right]. \quad (65)$$

Therefore, the updating for $\lambda_{\mu_{\mathbf{z}}}^{(\text{new})}$ according to (63) becomes

$$\lambda_{\mu_{\mathbf{z}}}^{(\text{new})} = \mathbf{D}(\tau_p^{-1})(\mathbf{p} - \hat{\mathbf{z}}). \quad (66)$$

It is now clear that we can relate to the LSL BP GAMP above (or [8]) if we define

$$\mathbf{s} = -\lambda_{\mu_{\mathbf{z}}}. \quad (67)$$

For further use, we also state the computation of $\boldsymbol{\tau}_{\hat{\mathbf{z}}}$ explicitly:

$$\boldsymbol{\tau}_{\hat{\mathbf{z}}} = \mathbb{E} \left[(\mathbf{z} - \hat{\mathbf{z}})^2 |p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{p}\|_{\boldsymbol{\tau}_p}^2} \right] \quad (68)$$

$$= \mathbb{E} \left[(\mathbf{z} - \hat{\mathbf{z}})^2 |p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|_{\boldsymbol{\tau}_p}^2} e^{\boldsymbol{\lambda}_{\mu_{\mathbf{z}}}^{(\text{new})T} \mathbf{z}} \right] \quad (69)$$

where \mathbf{z}^2 denotes element-wise square of vector \mathbf{z} . (68) and (69) result in the same solution.

7.2 Update of $\boldsymbol{\lambda}_{\mu_{\mathbf{x}}}$

According to (60), we can update $\boldsymbol{\lambda}_{\mu_{\mathbf{x}}}$ by

$$\boldsymbol{\lambda}_{\mu_{\mathbf{x}}} = -\mathbf{A}^T \boldsymbol{\lambda}_{\mu_{\mathbf{z}}} . \quad (70)$$

To show the relation between this paper and [8], we define

$$\begin{aligned} \boldsymbol{\tau}_{\mathbf{r}} &= \mathbf{1} / (\mathbf{S}^T \boldsymbol{\lambda}_{\boldsymbol{\tau}}) \\ \mathbf{r} &= \hat{\mathbf{x}}^{\text{old}} + \mathbf{D}_{\boldsymbol{\tau}_{\mathbf{r}}} \mathbf{A}^T \mathbf{s} \end{aligned} \quad (71)$$

Then the updated posterior mean and variance becomes

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbb{E} \left[\mathbf{x} | p_{\mathbf{x}}(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\boldsymbol{\tau}_{\mathbf{r}}}^2} \right] \\ \boldsymbol{\tau}_{\hat{\mathbf{x}}} &= \mathbb{E} \left[(\mathbf{x} - \hat{\mathbf{x}})^2 | p_{\mathbf{x}}(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\boldsymbol{\tau}_{\mathbf{r}}}^2} \right], \end{aligned} \quad (72)$$

where we also used the mean constraint for \mathbf{x} in (52).

7.3 The update of $\boldsymbol{\lambda}_{\boldsymbol{\tau}}$ and $\boldsymbol{\tau}_p$

The updates of these two variables are quite straightforward. They are already explicitly given by (58) and the variance constraint in (52). To show the relation with GAMP in [8] explicitly, we can define $\boldsymbol{\tau}_{\mathbf{s}}$ and then get

$$\boldsymbol{\tau}_{\mathbf{s}} = \boldsymbol{\lambda}_{\boldsymbol{\tau}} \text{ from which } \boldsymbol{\tau}_p = \mathbf{S} \boldsymbol{\tau}_{\mathbf{s}}, \boldsymbol{\tau}_{s_k} = \frac{1}{\boldsymbol{\tau}_{p_k}} - \frac{\boldsymbol{\tau}_{z_k}}{\boldsymbol{\tau}_{p_k}^2}. \quad (73)$$

8 Iterative Solution leading to AMBGAMP

GAMP does not use the extra variable \mathbf{u} in (52) (hence uses $\mathbf{u} = \hat{\mathbf{x}}$) and as result is an algorithm that does not correspond to alternating optimization of a BFE, with the resulting convergence issues. For AMBGAMP, we keep variable \mathbf{u} , and use (56)-(60) along with the three constraints in (52) as a system of 8 equations to be solved.

8.1 Update of $\boldsymbol{\lambda}_{\mu_{\mathbf{z}}}$

Consider (57) and the mean constraint in (52)

$$\mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|_{\boldsymbol{\tau}_p}^2} e^{\boldsymbol{\lambda}_{\mu_{\mathbf{z}}}^{(\text{new})T} \mathbf{z}} \right] \quad (74)$$

$$= \mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{A}\mathbf{u}\|_{\boldsymbol{\tau}_p}^2} e^{\boldsymbol{\lambda}_{\mu_{\mathbf{z}}}^T \mathbf{z}} \right] = \hat{\mathbf{z}} \quad (75)$$

To make the connection with AMBGAMP in [8], we define

$$\mathbf{p} = \mathbf{A}\mathbf{u} + \boldsymbol{\tau}_p \cdot \boldsymbol{\lambda}_{\mu_{\mathbf{z}}} \quad (76)$$

Similar to the previous section, we have the update

$$\boldsymbol{\lambda}_{\mu_{\mathbf{z}}}^{(\text{new})} = (\mathbf{p} - \widehat{\mathbf{z}}) / \boldsymbol{\tau}_{\mathbf{p}}. \quad (77)$$

Substitute (76) into (77), and we have

$$\boldsymbol{\lambda}_{\mu_{\mathbf{z}}}^{(\text{new})} = \boldsymbol{\lambda}_{\mu_{\mathbf{z}}} + (\mathbf{A}\mathbf{u} - \widehat{\mathbf{z}}) / \boldsymbol{\tau}_{\mathbf{p}}. \quad (78)$$

If we define

$$\mathbf{s} = -\boldsymbol{\lambda}_{\mu_{\mathbf{z}}}, \quad (79)$$

it then follows

$$\mathbf{s}^{(\text{new})} = \mathbf{s} + (\widehat{\mathbf{z}} - \mathbf{A}\mathbf{u}) / \boldsymbol{\tau}_{\mathbf{p}}. \quad (80)$$

For the convenience of the further discussion, we write the update for the posterior mean and variance of \mathbf{z}

$$\widehat{\mathbf{z}} = \mathbb{E} \left[\mathbf{z} | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{p}\|_{\boldsymbol{\tau}_{\mathbf{p}}}^2} \right] \quad (81)$$

$$\boldsymbol{\tau}_{\widehat{\mathbf{z}}} = \mathbb{E} \left[(\mathbf{z} - \widehat{\mathbf{z}})^2 | p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}|\mathbf{z}) e^{-\frac{1}{2} \|\mathbf{z} - \mathbf{p}\|_{\boldsymbol{\tau}_{\mathbf{p}}}^2} \right]. \quad (82)$$

8.2 Update of $\boldsymbol{\lambda}_{\mu_{\mathbf{x}}}$

We can use (60) and (79) to obtain the update

$$\boldsymbol{\lambda}_{\mu_{\mathbf{x}}} = -\mathbf{A}^T \boldsymbol{\lambda}_{\mu_{\mathbf{z}}} = \mathbf{A}^T \mathbf{s} \quad (83)$$

If we define (and note that $\boldsymbol{\lambda}_{\boldsymbol{\tau}} = \boldsymbol{\tau}_{\mathbf{s}}$)

$$\boldsymbol{\tau}_{\mathbf{r}} = \mathbf{1} / (\mathbf{S}^T \boldsymbol{\lambda}_{\boldsymbol{\tau}}), \quad \mathbf{r} = \mathbf{u} + \boldsymbol{\tau}_{\mathbf{r}} \cdot \boldsymbol{\lambda}_{\mu_{\mathbf{x}}} = \mathbf{u} + \boldsymbol{\tau}_{\mathbf{r}} \cdot (\mathbf{A}^T \mathbf{s}) \boldsymbol{\lambda}_{\mu_{\mathbf{x}}}, \quad (84)$$

we have the explicit update for $\widehat{\mathbf{x}}$ and $\boldsymbol{\tau}_{\widehat{\mathbf{x}}}$:

$$\begin{aligned} \widehat{\mathbf{x}} &= \mathbb{E} \left[\mathbf{x} | p_{\mathbf{x}}(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\boldsymbol{\tau}_{\mathbf{r}}}^2} \right] \\ \boldsymbol{\tau}_{\widehat{\mathbf{x}}} &= \mathbb{E} \left[(\mathbf{x} - \widehat{\mathbf{x}})^2 | p_{\mathbf{x}}(\mathbf{x}) e^{-\frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_{\boldsymbol{\tau}_{\mathbf{r}}}^2} \right]. \end{aligned} \quad (85)$$

8.3 Update of \mathbf{u}

By combining (59) and (60), we get the solution

$$\mathbf{u} = \left[\mathbf{A}^T \mathbf{D}_{\boldsymbol{\tau}_{\mathbf{p}}}^{-1} \mathbf{A} + \mathbf{D}_{\boldsymbol{\tau}_{\mathbf{r}}}^{-1} \right]^{-1} (\mathbf{A}^T \mathbf{D}_{\boldsymbol{\tau}_{\mathbf{p}}}^{-1} \widehat{\mathbf{z}} + \mathbf{D}_{\boldsymbol{\tau}_{\mathbf{r}}}^{-1} \widehat{\mathbf{x}}). \quad (86)$$

AMBGAMP actually updates \mathbf{u} by applying SGD with stepsize by linesearch to the quadratic cost function that (86) is the solution of. The update of $\boldsymbol{\lambda}_{\boldsymbol{\tau}}$ and $\boldsymbol{\tau}_{\mathbf{p}}$ are identical to the updates in GAMP in (73).

9 Bilinear System Model

In the bilinear Cell Free systems, we consider a semi-blind signal model containing L APs. At the l -th AP,

$$[\mathbf{Y}_{p,l} \quad \mathbf{Y}_l] = \mathbf{H}_l [\mathbf{X}_p \quad \mathbf{X}] + [\mathbf{V}_{p,l} \quad \mathbf{V}_l]. \quad (87)$$

The received signals are composed of pilot part $\mathbf{Y}_{p,l} \in \mathbb{C}^{N \times P}$ and data part $\mathbf{Y}_l \in \mathbb{C}^{N \times T}$. The channels between different users are considered independent Gaussian i.e. $\text{vec}(\mathbf{H}_l) \sim \mathcal{CN}(\mathbf{0}, \boldsymbol{\Xi}_l)$ where $\boldsymbol{\Xi}_l \in \mathbb{C}^{NK \times NK}$ is a block diagonal matrix of K blocks $\boldsymbol{\Xi}_{\mu_{ik}} \in \mathbb{C}^{N \times N}$. The transmitted symbols can be decomposed as pilot symbols $\mathbf{X}_p \in \mathcal{S}^{K \times P}$ and data symbols $\mathbf{X} \in \mathcal{S}^{K \times T}$, where \mathcal{S} is the constellation set. We assume that the elements x_{kt} in \mathbf{X} follow the categorical distribution $p(x_{kt})$. The signal power is denoted as σ_x^2 . The noise is considered as i.i.d. Gaussian distribution, and thus, $\text{vec}([\mathbf{V}_{p,l} \quad \mathbf{V}_l]) \sim \mathcal{CN}(0, \sigma_v^2 \mathbf{I})$.

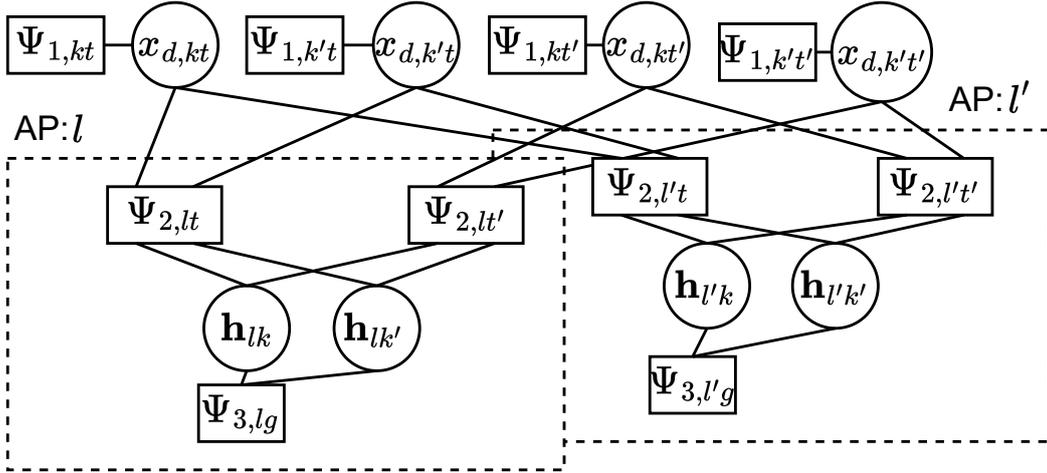


Figure 3: Partial factor graph

9.1 Orthogonal Pilot sequences

If orthogonal pilot sequences are used, we can first preprocess the pilot observation by right multiplying it with $\mathbf{x}_{p,g}^*$ which is the conjugated g -th pilot sequence. This results in an equivalent observation $\mathbf{y}_{p,lg}$

$$\mathbf{y}_{p,lg} = \mathbf{Y}_{p,l} \mathbf{x}_{p,g}^* = \sum_{k \in G_g} P \sigma_x^2 \mathbf{h}_{lk} + \mathbf{v}_{p,lg} \quad (88)$$

where $\mathbf{v}_{p,lg} = \mathbf{V}_{p,l} \mathbf{x}_{p,g}^* \sim \mathcal{N}(\mathbf{v}_{p,lg} | \mathbf{0}, P \sigma_x^2 \sigma_v^2 \mathbf{I})$, G_g denote the set of users using the g -th pilot sequence. We observe that every \mathbf{h}_{lk} occurs only in one group G_g , and the cross-correlation $\mathbb{E}[\mathbf{v}_{p,lg} \mathbf{v}_{p,l'g'}^H]$ is an all-zero matrix for all $g \neq g'$. Therefore, the observations $\mathbf{y}_{p,lg}$ and $\mathbf{y}_{p,l'g'}$ are independent. With orthogonal pilots, the factorization scheme is derived as

$$\begin{aligned} & p(\{\mathbf{y}_{p,lg}\}, \{\mathbf{Y}_l\}, \{\mathbf{H}_l\}, \mathbf{X}, \{\mathbf{V}_l\}) \\ &= \prod_{k,t} p(x_{kt}) \prod_l \prod_{t_1=1}^T p(\mathbf{y}_{lt_1} | \mathbf{H}_l, \mathbf{x}_{:t_1}) \prod_g p(\mathbf{y}_{p,lg}, \mathbf{H}_{lg}) \end{aligned} \quad (89)$$

where \mathbf{H}_{lg} is a matrix collecting all $k \in G_g$, \mathbf{h}_{lk} as its column vectors, and $\mathbf{x}_{:t_1}$ denotes the t_1 -th column of \mathbf{X} . We will base our EP (BP) algorithm based on this factorization scheme.

9.2 Expectation Propagation on Semi-Blind structure

For simplicity, we denote the factors in the factorization scheme (89) as

$$\Psi_{1,kt} = p(x_{kt}); \quad \Psi_{2,lt} = p(\mathbf{y}_{lt} | \mathbf{H}_l, \mathbf{x}_{:t}); \quad \Psi_{3,lg} = p(\mathbf{y}_{p,lg}, \mathbf{H}_{lg}).$$

The factor graph for (89) is illustrated in Fig. 3.

10 Bilinear Message Passing Derivations

This paper uses EP to estimate channel coefficients \mathbf{h}_{lk} and BP to estimate the data symbols x_{kt} . Furthermore, we specify the projection family of EP in this paper to be Gaussian distributions with diagonal covariance matrices. Now, we will examine each factor and derive its outbound message.

The message from $\Psi_{1,kt}$ to x_{kt} can be computed directly since no projection is needed, i.e., $\mu_{\Psi_{1,kt};x_{kt}}(x_{kt}) = p(x_{kt})$.

10.1 Message from $\Psi_{2,lt}$ to x_{kt}

Following EP rules, the extrinsic at node $\Psi_{2,lt}$ is updated by

$$\begin{aligned}\mu_{x_{kt};\Psi_{2,lt}}(x_{kt}) &\propto p(x_{kt}) \prod_{l' \neq l} \mu_{\Psi_{2,l't};x_{kt}}(x_{kt}) \\ \mu_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk}) &\propto \mu_{\Psi_{3,lg};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) \prod_{t' \neq t} \mu_{\Psi_{2,lt'};\mathbf{h}_{lk}}(\mathbf{h}_{lk}),\end{aligned}\quad (90)$$

where the extrinsic of \mathbf{h}_{lk} can be computed as a Gaussian $\mu_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk}) = \mathcal{CN}(\mathbf{h}_{lk} | \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}, \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}})$ with

$$\begin{aligned}\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} &= \left(\mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}}^{-1} + \sum_{t' \neq t} \mathbf{C}_{\Psi_{2,lt'};\mathbf{h}_{lk}}^{-1} \right)^{-1} \\ \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} &= \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} \left(\mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}}^{-1} \mathbf{m}_{\Psi_{3,lg};\mathbf{h}_{lk}} \right. \\ &\quad \left. + \sum_{t' \neq t} \mathbf{C}_{\Psi_{2,lt'};\mathbf{h}_{lk}}^{-1} \mathbf{m}_{\Psi_{2,lt'};\mathbf{h}_{lk}} \right)\end{aligned}$$

According to the EP rule, the message from $\Psi_{2,lt}$ to x_{kt} is

$$\mu_{\Psi_{2,lt};x_{kt}}(x_{kt}) \propto \frac{\text{proj}[b_{\Psi_{2,lt};x_{kt}}(x_{kt})]}{\mu_{x_{kt};\Psi_{2,lt}}(x_{kt})}, \quad (91)$$

where the belief (approximated posterior) at $\Psi_{2,lt}$ is defined as $b_{\Psi_{2,lt};x_{kt}}(x_{kt})$ with

$$\begin{aligned}b_{\Psi_{2,lt};x_{kt}}(x_{kt}) &\propto \mu_{x_{kt};\Psi_{2,lt}}(x_{kt}) \sum_{\mathbf{x}_{\bar{k}t}} \int p(\mathbf{y}_{lt} | x_{kt} \mathbf{h}_{lk} + \sum_{i \neq k} x_{it} \mathbf{h}_{li}) \\ &\quad \cdot \mu_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk}) \prod_{i \neq k} \mu_{\mathbf{h}_{li};\Psi_{2,lt}}(\mathbf{h}_{li}) \mu_{x_{it};\Psi_{2,lt}}(x_{it}) d\mathbf{H}_l.\end{aligned}\quad (92)$$

We use the notation $\mathbf{x}_{\bar{k}t}$ to denote all the elements in $\mathbf{x}_{:t}$ except the k -th element. The integral (and summation) in (92) can be considered as a marginalization operation. Furthermore, we can view the extrinsic messages as hypothetical priors. Due to CLT, we approximate $\sum_{i \neq k} x_{it} \mathbf{h}_{li}$ to a Gaussian where $x_{it} \sim \mu_{x_{it};\Psi_{2,lt}}(x_{it})$, $\mathbf{h}_{li} \sim \mu_{\mathbf{h}_{li};\Psi_{2,lt}}(\mathbf{h}_{li})$ [18]. Therefore, (92) becomes

$$\begin{aligned}b_{\Psi_{2,lt};x_{kt}}(x_{kt}) &\propto \mu_{x_{kt};\Psi_{2,lt}}(x_{kt}) \\ &\quad \cdot \int \int p(\mathbf{y}_{lt} | x_{kt} \mathbf{h}_{lk} + \mathbf{z}_{lkt}) \mu_{\mathbf{z}_{lkt}}(\mathbf{z}_{lkt}) \mu_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk}) d\mathbf{z}_{lkt} d\mathbf{h}_{lk},\end{aligned}\quad (93)$$

where $\mu_{\mathbf{z}_{lkt}}(\mathbf{z}_{lkt}) = \mathcal{CN}(\mathbf{z}_{lkt} | \mathbf{m}_{\mathbf{z}_{lkt}}, \mathbf{C}_{\mathbf{z}_{lkt}})$ with

$$\begin{aligned}\mathbf{m}_{\mathbf{z}_{lkt}} &= \sum_{i \neq k} m_{x_{it};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{li};\Psi_{2,lt}} \\ \mathbf{C}_{\mathbf{z}_{lkt}} &= \sum_{i \neq k} r_{x_{it};\Psi_{2,lt}} \mathbf{C}_{\mathbf{h}_{li};\Psi_{2,lt}} + \tau_{x_{it};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{it};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{it};\Psi_{2,lt}}^H\end{aligned}\quad (94)$$

where $m_{x_{it};\Psi_{2,lt}}$, $\tau_{x_{it};\Psi_{2,lt}}$ and $r_{x_{it};\Psi_{2,lt}}$ are the mean, variance and second-order moment of the normalized message $\mu_{x_{it};\Psi_{2,lt}}$. By applying the Gaussian reproduction lemma [18] and the fact that Gaussian distribution integrates to one, the belief (93) becomes

$$\begin{aligned}b_{\Psi_{2,lt};x_{kt}}(x_{kt}) &\propto \mathcal{CN}(\mathbf{0} | \mathbf{y}_{lt} - \mathbf{m}_{\mathbf{z}_{lkt}} - x_{kt} \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}, \\ &\quad \mathbf{C}_v + \mathbf{C}_{\mathbf{z}_{lkt}} + |x_{kt}|^2 \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}) \cdot \mu_{x_{kt};\Psi_{2,lt}}(x_{kt}).\end{aligned}\quad (95)$$

Therefore, by BP rules, the outbound message is

$$\begin{aligned} \mu_{\Psi_{2,lt};x_{kt}}(x_{kt}) &\propto \mathcal{CN}(\mathbf{0}|\mathbf{y}_{lt}-\mathbf{m}_{\mathbf{z}_{lkt}}-x_{kt}\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}, \\ &\mathbf{C}_v+\mathbf{C}_{\mathbf{z}_{lkt}}+|x_{kt}|^2\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}) \end{aligned} \quad (96)$$

10.2 Message from $\Psi_{2,lt}$ to \mathbf{h}_{lk}

Based on EP rules, the message to \mathbf{h}_{lk} is

$$\mu_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) \propto \frac{\text{proj}[b_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk})]}{\mu_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk})}, \quad (97)$$

where the belief is defined as

$$\begin{aligned} b_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) &\propto \sum_{\mathbf{x}:t} \int p(\mathbf{y}_{lt}|\sum_i x_{it}\mathbf{h}_{li}) \\ &\cdot \prod_i \mu_{\mathbf{h}_{li};\Psi_{2,lt}}(\mathbf{h}_{li}) \mu_{x_{it};\Psi_{2,lt}}(x_{it}) d\mathbf{h}_{l\bar{k}} \end{aligned} \quad (98)$$

We use $\mathbf{h}_{l\bar{k}}$ to denote all the column vectors in \mathbf{H}_l except the k -th column. By using the same approach from (92) to (95), and separating the terms that contains only x_{kt} [18] [19], the belief (98) becomes

$$\begin{aligned} &b_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) \\ &= \mathbb{E}_{b_{\Psi_{2,lt};x_{kt}}} \{ \mathcal{CN}[\mathbf{h}_{lk}|\mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x_{kt}), \mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x_{kt})] \} \end{aligned} \quad (99)$$

where $\mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(\cdot)$ and $\mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(\cdot)$ are defined as

$$\begin{aligned} \mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) &= [|x|^2(\mathbf{C}_v+\mathbf{C}_{\mathbf{z}_{lkt}})^{-1}+\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{-1}]^{-1} \\ \mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) &= \mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) \left[\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{-1} \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} \right. \\ &\left. + |x|^2(\mathbf{C}_v+\mathbf{C}_{\mathbf{z}_{lkt}})^{-1} \frac{\mathbf{y}_{lt}-\mathbf{m}_{\mathbf{z}_{lkt}}}{x} \right], \end{aligned} \quad (100)$$

where $\mathbf{C}_v = \sigma_v^2 \mathbf{I}$. The mean $\mathbf{m}_{\hat{\mathbf{h}}_{lk}^2}$ and covariance $\mathbf{C}_{\hat{\mathbf{h}}_{lk}^2}$ of the belief distribution (99) are

$$\begin{aligned} \mathbf{m}_{\hat{\mathbf{h}}_{lk}^2} &= \mathbb{E}_{b_{\Psi_{2,lt};x_{kt}}} [\mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x_{kt})] \\ \mathbf{C}'_{\hat{\mathbf{h}}_{lk}^2} &= \mathbb{E}_{b_{\Psi_{2,lt};x_{kt}}} [\mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x_{kt}) \\ &+ \mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x_{kt}) \mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x_{kt})^H] - \mathbf{m}_{\hat{\mathbf{h}}_{lk}^2} \mathbf{m}_{\hat{\mathbf{h}}_{lk}^2}^H. \end{aligned} \quad (101)$$

We project the belief at $\Psi_{2,lt}$ to a Gaussian with diagonal covariance matrix $\text{proj}[b_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk})] = \mathcal{CN}(\mathbf{h}_{lk}|\mathbf{m}_{\hat{\mathbf{h}}_{lk}^2}, \mathbf{C}_{\hat{\mathbf{h}}_{lk}^2})$, where $\mathbf{C}_{\hat{\mathbf{h}}_{lk}^2}$ is a digonal matrix with the same diagonal elements as $\mathbf{C}'_{\hat{\mathbf{h}}_{lk}^2}$. Finally, the message from $\Psi_{2,lt}$ to \mathbf{h}_{lk} is

$$\begin{aligned} \mu_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) &= \mathcal{CN}(\mathbf{h}_{lk}|\mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}}, \mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}) \\ &\propto \frac{\mathcal{CN}(\mathbf{h}_{lk}|\mathbf{m}_{\hat{\mathbf{h}}_{lk}^2}, \mathbf{C}_{\hat{\mathbf{h}}_{lk}^2})}{\mathcal{CN}(\mathbf{h}_{lk}|\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}, \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}})}. \end{aligned} \quad (102)$$

10.3 Message form $\Psi_{3,lg}$ to \mathbf{h}_{lk}

We assume $k \in G_g$. The extrinsic at $\Psi_{3,lg}$ is updated by

$$\mu_{\mathbf{h}_{lk};\Psi_{3,lg}}(\mathbf{h}_{lk}) \propto \prod_t \mu_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk}). \quad (103)$$

We denote this extrinsic message as $\mu_{\mathbf{h}_{lk};\Psi_{3,lg}}(\mathbf{h}_{lk}) = \mathcal{CN}(\mathbf{h}_{lk}|\mathbf{m}_{\mathbf{h}_{lk};\Psi_{3,lg}}, \mathbf{C}_{\mathbf{h}_{lk};\Psi_{3,lg}})$ with

$$\begin{aligned}\mathbf{C}_{\mathbf{h}_{lk};\Psi_{3,lg}} &= \left(\sum_t \mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}^{-1} \right)^{-1} \\ \mathbf{m}_{\mathbf{h}_{lk};\Psi_{3,lg}} &= \mathbf{C}_{\mathbf{h}_{lk};\Psi_{3,lg}} \left(\sum_t \mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}^{-1} \mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}} \right).\end{aligned}$$

The belief of \mathbf{h}_{lg} at the $\Psi_{3,lg}$ is

$$b_{\Psi_{3,lg}}(\mathbf{h}_{lg}) \propto p(\mathbf{y}_{p,lg}, \mathbf{h}_{lg}) \prod_{k \in G_g} p(\mathbf{h}_{lk}) \mu_{\mathbf{h}_{lk};\Psi_{3,lg}}(\mathbf{h}_{lk}). \quad (104)$$

All the factors appearing in (104) are Gaussian pdfs with diagonal covariance matrices. Therefore, the projection of $b_{\Psi_{3,lg}}(\mathbf{h}_{lg})$ results to itself. For simplicity, we define a hypothetical prior $q_{\mathbf{h}_{lk}|\mathbf{Y}_d}$ for \mathbf{h}_{lk} in (104) as

$$\begin{aligned}q_{\mathbf{h}_{lk}|\mathbf{Y}_d}(\mathbf{h}_{lk}) &= \mathcal{N}(\mathbf{h}_{lk} | \mathbf{m}_{\mathbf{h}_{lk}|\mathbf{Y}_d}, \mathbf{C}_{\mathbf{h}_{lk}|\mathbf{Y}_d}) \\ &\propto p(\mathbf{h}_{lk}) \mu_{\mathbf{h}_{lk};\Psi_{3,lg}}(\mathbf{h}_{lk}),\end{aligned} \quad (105)$$

where

$$\begin{aligned}\mathbf{C}_{\mathbf{h}_{lk}|\mathbf{Y}_d} &= (\Xi_{\mathbf{h}_{lk}}^{-1} + \mathbf{C}_{\mathbf{h}_{lk};\Psi_{3,lg}}^{-1})^{-1} \\ \mathbf{m}_{\mathbf{h}_{lk}|\mathbf{Y}_d} &= \mathbf{C}_{\mathbf{h}_{lk}|\mathbf{Y}_d} \mathbf{C}_{\mathbf{h}_{lk};\Psi_{3,lg}}^{-1} \mathbf{m}_{\mathbf{h}_{lk};\Psi_{3,lg}}\end{aligned} \quad (106)$$

The message from factor node $\Psi_{3,lg}$ to \mathbf{h}_{lk} can be derived as

$$\begin{aligned}\mu_{\Psi_{3,lg};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) &\propto \frac{\int b_{\Psi_{3,lg}}(\mathbf{h}_{lg}) d\mathbf{h}_{lg}}{\mu_{\mathbf{h}_{lk};\Psi_{3,lg}}(\mathbf{h}_{lk})} \\ &\propto p(\mathbf{h}_{lk}) \frac{\int p(\mathbf{y}_{p,lg}, \mathbf{h}_{lg}) \prod_{k \in G_g} q_{\mathbf{h}_{lk}|\mathbf{Y}_d}(\mathbf{h}_{lk}) d\mathbf{h}_{lg}}{q_{\mathbf{h}_{lk}|\mathbf{Y}_d}(\mathbf{h}_{lk})}\end{aligned} \quad (107)$$

The fraction operation in the second line of (107) can be interpreted as component-wise conditionally-unbiased LMMSE estimation [20]. Therefore, the message from $\Psi_{3,lg}$ to \mathbf{h}_{lk} is

$$\mu_{\Psi_{3,lg};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) = \mathcal{CN}(\mathbf{h}_{lk} | \mathbf{m}_{\Psi_{3,lg};\mathbf{h}_{lk}}, \mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}}), \quad (108)$$

where

$$\begin{aligned}\mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}} &= \left[\Xi_{\mathbf{h}_{lk}}^{-1} + \left(\frac{\sigma_v^2}{\sigma_x^2 P} \mathbf{I} + \sum_{k' \in G_g / \{k\}} \mathbf{C}_{\mathbf{h}_{lk'}|\mathbf{Y}_d} \right)^{-1} \right]^{-1} \\ \mathbf{m}_{\Psi_{3,lg};\mathbf{h}_{lk}} &= \Xi_{\mathbf{h}_{lk}} \left(\frac{\sigma_v^2}{\sigma_x^2 P} \mathbf{I} + \sum_{k' \in G_g / \{k\}} \mathbf{C}_{\mathbf{h}_{lk'}|\mathbf{Y}_d} + \Xi_{\mathbf{h}_{lk}} \right)^{-1} \\ &\quad \cdot \left(\frac{1}{\sigma_x^2 P} \mathbf{y}_{p,lg} - \sum_{k' \in G_g / \{k\}} \mathbf{m}_{\mathbf{h}_{lk'}|\mathbf{Y}_d} \right).\end{aligned} \quad (109)$$

11 Asymptotic Behaviors in Bilinear Large Systems

For scalable systems, $L \rightarrow \infty$ while $K = c_1 L$, $P = c_2 L$, $T = c_3 L$, where c_1, c_2, c_3 are some positive constants, we assume the channel coefficients $\forall l, n, k, \mathbb{E}[|h_{lnk}|^2] = O(\frac{1}{L})$ data constellation symbols $\forall s \in \mathcal{S}, s = O(1), 1/s = O(1)$. Furthermore, we assume the noise power scales as $\sigma_v^2 = O(1)$, $\sigma_v^{-2} = O(1)$. For simplicity, we define big-O-notations with matrix parameters to represent the element-wise asymptotic behavior, i.e., for matrices \mathbf{A}, \mathbf{B} of the same size, we have $\mathbf{A} = O(\mathbf{B}) \Leftrightarrow \forall i, j, [\mathbf{A}]_{ij} = O([\mathbf{B}]_{ij})$.

Assumption 11.1. We assume that $\mathbb{E}(\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}) = \mathbf{0}$, $\mathbb{E}[m_{x_{it};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{li};\Psi_{2,lt}}] = \mathbf{0}$, and $\forall i \neq j, \mathbb{E}([\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}^H]_{ij}) = 0$.

Property 11.2. For invertible matrices \mathbf{A}, \mathbf{B} , we have $(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} = \mathbf{A}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{B} = \mathbf{B}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{A}$.

Lemma 11.3. With proper initialization and Assumption 11.1, in each iteration, the updates satisfy $\mathbf{m}_{\mathbf{z}_{lkt}} = O(\mathbf{1})$, $\mathbf{C}_{\mathbf{z}_{lkt}} = O(\mathbf{I}) + O(\frac{1 \cdot 1^H}{\sqrt{L}})$, $\mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) = O(\frac{1}{\sqrt{L}})$, $\mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) = O(\frac{\mathbf{I}}{L}) + O(\frac{1 \cdot 1^H}{L^2})$, $\mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}} = O(\frac{1}{\sqrt{L}})$, $\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}} = O(\mathbf{I})$, $\mathbf{m}_{\Psi_{3,lg};\mathbf{h}_{lk}} = O(\frac{1}{\sqrt{L}})$, $\mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}} = O(\frac{\mathbf{I}}{L})$, $\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} = O(\frac{1}{\sqrt{L}})$, $\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} = O(\frac{\mathbf{I}}{L})$. Furthermore, $\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}^{-1} = O(\mathbf{I})$.

Proof. We prove this lemma by mathematical induction. Due to a proper initialization, we can assume the messages $\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}$, $\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}$, $\mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}}$, $\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}$, $\mathbf{m}_{\Psi_{3,lg};\mathbf{h}_{lk}}$, $\mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}}$ are initialized with the above-mentioned scales.

Then, we assume the lemma holds for the previous iterations and investigate the updates in the next iteration.

We first look at the update of $\mathbf{m}_{\mathbf{z}_{lkt}}$, $\mathbf{C}_{\mathbf{z}_{lkt}}$, which are updated according to (94). Similar to [21], we assume that $\forall i, m_{x_{it};\Psi_{2,lt}}$ are weakly independent of $\mathbf{m}_{\mathbf{h}_{li};\Psi_{2,lt}}$. Since the elements in the constellation set scale with $O(1)$, we know $m_{x_{it};\Psi_{2,lt}} = O(1)$. According to induction assumptions, the extrinsic mean $\mathbf{m}_{\mathbf{h}_{li};\Psi_{2,lt}} = O(\frac{1}{\sqrt{L}})$. Due to Assumption 11.1, we use the results from [21, Lemma 1] to obtain $\mathbf{m}_{\mathbf{z}_{lkt}} = \sum_{i \neq k} m_{x_{it};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{li};\Psi_{2,lt}} = O(\mathbf{1})$ and the covariance matrix $\mathbf{C}_{\mathbf{z}_{lkt}} = O(\mathbf{I}) + O(\frac{1 \cdot 1^H}{\sqrt{L}})$. For simplicity, we denote the diagonal terms of $\mathbf{C}_{\mathbf{z}_{lkt}}$ in (94) as $\mathbf{D}_{\mathbf{z}_{lkt}} = \sum_{i \neq k} \tau_{x_{it};\Psi_{2,lt}} \mathbf{C}_{\mathbf{h}_{li};\Psi_{2,lt}}$, and denote $\mathbf{B}_{\mathbf{z}_{lkt}} = \sum_{i \neq k} \tau_{x_{it};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{li};\Psi_{2,lt}} \mathbf{m}_{\mathbf{h}_{li};\Psi_{2,lt}}^H = O(\mathbf{I}) + O(\frac{1 \cdot 1^H}{\sqrt{L}})$. Thus, with these notations, $\mathbf{C}_{\mathbf{z}_{lkt}} = \mathbf{D}_{\mathbf{z}_{lkt}} + \mathbf{B}_{\mathbf{z}_{lkt}}$.

Now we investigate the update of $\mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x)$, $\mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x)$ in (100). By matrix inversion lemma,

$$\begin{aligned} \mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) &= \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} \\ &\quad - \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} \mathbf{Q}_{lkt}^{-\frac{H}{2}} \mathbf{W}_{lkt} (\mathbf{\Lambda}_{lkt} + \mathbf{I})^{-1} \mathbf{W}_{lkt}^H \mathbf{Q}_{lkt}^{-\frac{1}{2}} \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}, \end{aligned} \quad (110)$$

where we define the positive semi-definite diagonal matrix $\mathbf{Q}_{lkt} = \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} + \frac{1}{|x|^2} (\mathbf{C}_v + \mathbf{D}_{\mathbf{z}_{lkt}}) = O(\mathbf{I})$ and $\mathbf{Q}_{lkt}^{\frac{1}{2}} \mathbf{Q}_{lkt}^{\frac{H}{2}} = \mathbf{Q}_{lkt}$. By eigendecomposition, we define $\frac{1}{|x|^2} \mathbf{Q}_{lkt}^{-\frac{1}{2}} \mathbf{B}_{\mathbf{z}_{lkt}} \mathbf{Q}_{lkt}^{-\frac{H}{2}} = \mathbf{W}_{lkt} \mathbf{\Lambda}_{lkt} \mathbf{W}_{lkt}^H$, $\mathbf{I} = \mathbf{W}_{lkt} \mathbf{W}_{lkt}^H$. Therefore, $\mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) = O(\frac{\mathbf{I}}{L}) + O(\frac{1 \cdot 1^H}{L^2})$. We find the update of $\mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}}$ in (100) is dominated by the first term $\mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{-1} \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}$. By neglecting higher order infinitesimal terms, we have

$$\begin{aligned} \mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}} &\simeq \mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x) \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{-1} \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} = [\mathbf{I} + \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{-\frac{1}{2}} \\ &\quad \cdot (\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{\frac{H}{2}} \mathbf{C}_{\text{in},lkt}^{-1} \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{\frac{1}{2}} + \mathbf{I})^{-1} \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{-\frac{H}{2}} \mathbf{C}_{\text{in},lkt}^{-1}] \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} \\ &\simeq \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} = O(\frac{1}{\sqrt{L}}). \end{aligned}$$

To study the messages $\mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}}$, $\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}$, we first investigate the approximated (projected) belief of \mathbf{h}_{lk} at $\Psi_{2,lt}$. From the previous discussion, $\mathbf{C}'_{\hat{\mathbf{h}}_{lk}} \simeq \mathbb{E}_{b_{\Psi_{2,lt};x_{kt}}} [\mathbf{C}_{\hat{\mathbf{h}}_{lk}|x_{kt}}(x_{kt})]$, and thus,

$$\mathbf{C}'_{\hat{\mathbf{h}}_{lk}} \simeq \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} - \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} \cdot \mathbf{F} \cdot \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}, \quad (111)$$

where $\mathbf{F} = \sum_{x \in \mathcal{S}} b_{\Psi_{2,lt};x_{kt}}(x) \left[\frac{1}{|x|^2} \mathbf{B}_{\mathbf{z}_{lkt}} + \mathbf{Q}_{lkt} \right]^{-1}$. Thanks to the projection in EP, we are only interested in the diagonal elements of $\mathbf{C}'_{\hat{\mathbf{h}}_{lk}}$. With the approximation $\mathbf{m}_{\hat{\mathbf{h}}_{lk}|x_{kt}} \simeq \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}$, the n -th diagonal term reads

$$\begin{aligned} [\mathbf{C}'_{\hat{\mathbf{h}}_{lk}}]_{nn} &= [\mathbf{C}_{\hat{\mathbf{h}}_{lk}}]_{nn} \simeq \tau_{\mathbf{h}_{lnk};\Psi_{2,lt}} - \tau_{\mathbf{h}_{lnk};\Psi_{2,lt}}^2 [\mathbf{F}]_{nn} \\ &= [([\mathbf{F}]_{nn}^{-1} - \tau_{\mathbf{h}_{lnk};\Psi_{2,lt}})^{-1} + \tau_{\mathbf{h}_{lnk};\Psi_{2,lt}}^{-1}]^{-1} \end{aligned} \quad (112)$$

From the same analysis in (110), we know \mathbf{Q}_{lkt} is asymptotic upper and lower bounded. Since $\mathbf{B}_{\mathbf{z}_{lkt}}$ is positive semi-definite, we have $[\mathbf{F}]_{nn} = O(1)$ and $[\mathbf{F}]_{nn}^{-1} = O(1)$. Substitute (112) into (102), and we obtain $[\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}]_{nn} \simeq [\mathbf{F}]_{nn}^{-1} - \tau_{\mathbf{h}_{lk};\Psi_{2,lt}}$. Thus, $\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}} = O(\mathbf{I})$ and $\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}^{-1} = O(\mathbf{I})$. Since $\mathbf{m}_{\widehat{\mathbf{h}}_{lk}|x_{kt}} \simeq \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}$, it is straightforward to see $\mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}} = O(\frac{1}{\sqrt{L}})$.

The message covariance matrices $\mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}}$ and $\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}$ are both diagonal matrices. Due to the scalable-system assumption $T = O(L)$ and the elements in $\mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}}$ being asymptotically upper and lower bounded, one can show $\mathbf{C}_{\Psi_{3,lg};\mathbf{h}_{lk}} = O(\frac{1}{L})$ is upper bounded, and $\mathbf{m}_{\Psi_{3,lg};\mathbf{h}_{lk}} = O(\frac{1}{\sqrt{L}})$ according to (106)-(109). We can then show $\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} = O(\frac{1}{L})$ and $\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} = O(\frac{1}{\sqrt{L}})$ according to (90). \square

12 Simplification of the Messages in Bilinear EP

We define beliefs at the variable nodes as

$$\begin{aligned} b_{x_{kt}}(x_{kt}) &\propto p(x_{kt}) \prod_l \mu_{\Psi_{2,lt};x_{kt}}(x_{kt}) \\ b_{\mathbf{h}_{lk}}(\mathbf{h}_{lk}) &\propto \mu_{\Psi_{3,lg};\mathbf{h}_{lk}}(\mathbf{h}_{lk}) \prod_t \mu_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk}). \end{aligned} \quad (113)$$

Compared to the extrinsic messages in (90), the beliefs in (113) only differ by one factor. Since Loopy BP is used for estimating x_{kt} , it has been shown in [17] that we can assume $b_{x_{kt}}(x_{kt}) \simeq \mu_{x_{kt};\Psi_{2,lt}}(x_{kt})$.

This work estimates the channel coefficients \mathbf{h}_{lk} using EP. Therefore, a separate analysis from [17] is needed. We investigate the mean and covariance matrix difference between $b_{\mathbf{h}_{lk}}(\mathbf{h}_{lk})$ and $\mu_{\mathbf{h}_{lk};\Psi_{2,lt}}$ based on the Lemma 2.

Substitute (90) into (113), and we obtain the belief at \mathbf{h}_{lk} as

$$b_{\mathbf{h}_{lk}}(\mathbf{h}_{lk}) \propto \mu_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk}) \mu_{\Psi_{2,lt};\mathbf{h}_{lk}}(\mathbf{h}_{lk}). \quad (114)$$

Denote $\mathbf{m}_{\widehat{\mathbf{h}}_{lk}}$ and $\mathbf{C}_{\widehat{\mathbf{h}}_{lk}}$ as the mean and covariance matrix of $b_{\mathbf{h}_{lk}}(\mathbf{h}_{lk})$. From Lemma 2, we have

$$\begin{aligned} \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} - \mathbf{C}_{\widehat{\mathbf{h}}_{lk}} &= \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^2 (\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} + \mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}})^{-1} \\ \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} - \mathbf{m}_{\widehat{\mathbf{h}}_{lk}} &= \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} (\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} + \mathbf{C}_{\Psi_{2,lt};\mathbf{h}_{lk}})^{-1} \\ &\cdot (\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} - \mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}}). \end{aligned} \quad (115)$$

It has been shown in the proof of Lemma 2 that the difference $\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} - \mathbf{m}_{\Psi_{2,lt};\mathbf{h}_{lk}}$ is a higher order infinitesimal relative to $\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}$. Thus, the quotients $(\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}} - \mathbf{C}_{\widehat{\mathbf{h}}_{lk}}) \mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}^{-1}$ and $(\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}} - \mathbf{m}_{\widehat{\mathbf{h}}_{lk}}) / \mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}$ tend to zero as the system grows larger. Therefore, the difference in (115) are higher order infinitesimals relative to $\mathbf{C}_{\mathbf{h}_{lk};\Psi_{2,lt}}$ and $\mathbf{m}_{\mathbf{h}_{lk};\Psi_{2,lt}}$, respectively. Therefore, we have $b_{\mathbf{h}_{lk}}(\mathbf{h}_{lk}) \simeq \mu_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk})$. Based on the above discussion, we propose to replace the extrinsic (90) at $\Psi_{2,lt}$ by (116) and (117) to reduce complexity further,

$$\mu'_{\mathbf{h}_{lk};\Psi_{2,lt}}(\mathbf{h}_{lk}) = b_{\mathbf{h}_{lk}}(\mathbf{h}_{lk}); \quad (116)$$

$$\mu'_{x_{kt};\Psi_{2,lt}}(x_{kt}) = b_{x_{kt}}(x_{kt}). \quad (117)$$

13 Decentralized Method for Bilinear EP

To obtain the belief of x_{kt} , we need to combine the message from all the AP. We consider the case where all the L AP are connected, and the AP network has a tree structure. A decentralized message-passing method can be used based on the consensus propagation framework [22]. Define the normalized message from AP l to AP l' :

$$\nu_{l \rightarrow l'}(x_{kt}) \propto \mu_{\Psi_{2,lt};x_{kt}}(x_{kt}) \prod_{\bar{l} \in N(l)/\{l'\}} \nu_{\bar{l} \rightarrow l}(x_{kt}),$$

Algorithm 1: One Iteration of Decentralized EP

Require: $\Xi_{\mathbf{h}_{lk}}, \mathbf{y}_{p,lg}, \mathbf{y}_{lt}, p(x_{kt}), \sigma_x^2, \sigma_v^2, G_g$

- 1: Initialize $\mu_{\Psi_{3,lg}; \mathbf{h}_{lk}}, \mu_{\Psi_{2,lt}; x_{kt}}, \mu_{\Psi_{2,lt}; \mathbf{h}_{lk}}, \nu_{l \rightarrow l'}(x_{kt})$
- 2: At all the APs, $\forall k, t$, update $b_{x_{kt}}$ according to (118)
- 3: **for** $l=1:L$ **do**
- 4: $\forall k$, update $\mu_{\mathbf{h}_{lk}; \Psi_{3,lg}}$ based on (103)
- 5: $\forall k, t$, update $\mu'_{\mathbf{h}_{lk}; \Psi_{2,lt}}$ based on (113) and (117)
- 6: $\forall k, t$, update $\mu'_{x_{kt}; \Psi_{2,lt}}$ based on (116)
- 7: $\forall k$, update $\mu_{\Psi_{3,lg}; \mathbf{h}_{lk}}$ based on (108)-(109)
- 8: $\forall k, t$, update $\mu_{\Psi_{2,lt}; x_{kt}}$ based on (94)-(96)
- 9: $\forall k, t$, update $\mu_{\Psi_{2,lt}; \mathbf{h}_{lk}}$ based on (100)-(102)
- 10: $\forall l' \in N(l), k, t$, update $\nu_{l \rightarrow l'}(x_{kt})$ based on (119)
- 11: **end for**

where $N(l)$ denotes the set of connected neighbors of AP l . At convergence, the belief in (113) can be obtained by any AP l as

$$b'_{x_{kt}}(x_{kt}) \propto p(x_{kt}) \mu_{\Psi_{2,lt}; x_{kt}}(x_{kt}) \prod_{l' \in N(l)} \nu_{l' \rightarrow l}(x_{kt}). \quad (118)$$

Therefore, for a decentralized algorithm, we can replace the update of belief $b_{x_{kt}}$ in (113) by $b'_{x_{kt}}$ in (118). After updating the message $\mu_{\Psi_{2,lt}; x_{kt}}$, we update the shared message by

$$\nu_{l \rightarrow l'}^{new}(x_{kt}) \propto \mu_{\Psi_{2,lt}; x_{kt}}^{new}(x_{kt}) \prod_{\bar{l} \in N(l)/\{l'\}} \nu_{\bar{l} \rightarrow l}^{old}(x_{kt}), \quad (119)$$

where we use *new* and *old* to distinguish the message of different iterations. One possible ordering method is suggested in Algorithm 1.

14 Bilinear EP Simulation Results

Our study simulates an environment within a 400×400 square meter area, equipped with 16 APs and 8 User Terminals (UTs). Each AP features $N = 2$ antennas and is positioned at coordinates $(\frac{400}{3}i, \frac{400}{3}j)$, $i, j \in \{0, 1, 2, 3\}$. The UTs are uniformly distributed throughout the area. We denote the distance between each UT k and AP l as d_{lk} . Channel covariances for each user k at AP l are modeled using $N \times N$ diagonal matrices, represented as $\sigma_{h_{lk}}^2 \mathbf{I}$, where $10 \log_{10}(\sigma_{h_{lk}}^2) = -30 - 36.7 \log_{10}(d_{lk})$.

All the neighboring APs within $\frac{400}{3}$ meters are connected and can exchange information of the estimated data symbols. Furthermore, as illustrated in Algorithm 1, a synchronized message-exchanging scheme is used.

The length of the orthogonal pilot sequences is set to $P = 6$ to introduce pilot contamination.

We employ a 4QAM constellation of length $T = 10$ for signal transmission and assume a noise power of -96 dBm. The signal-to-noise ratio (SNR) is adjusted by varying the transmitted power. We base our results on 100 different realizations, which are illustrated in Figure 4. The normalized mean squared error (NMSE) of the channel estimates is defined as $\text{NMSE} = \frac{\text{tr}[(\hat{\mathbf{H}} - \mathbf{H})^2]}{\text{tr}[\mathbf{H}^2]}$, where $\hat{\mathbf{H}}$ are synthesized from the mean of $b_{\mathbf{h}_{lk}}(\mathbf{h}_{lk})$ defined in (113) and the operation $|\cdot|^2$ is defined as $|\mathbf{H}|^2 = \mathbf{H}^H \mathbf{H}$.

In the VL-EP scenario, we generate data symbols drawn from i.i.d. Gaussian distribution and apply the VL-EP algorithm [23] for channel estimation. In the Genie-Aided scenario, we implement the proposed algorithm as if the data symbols are known. In the MMSE Genie-Aided scenario, all the APs estimate the channel coefficients using the MMSE estimator with known channel coefficients.

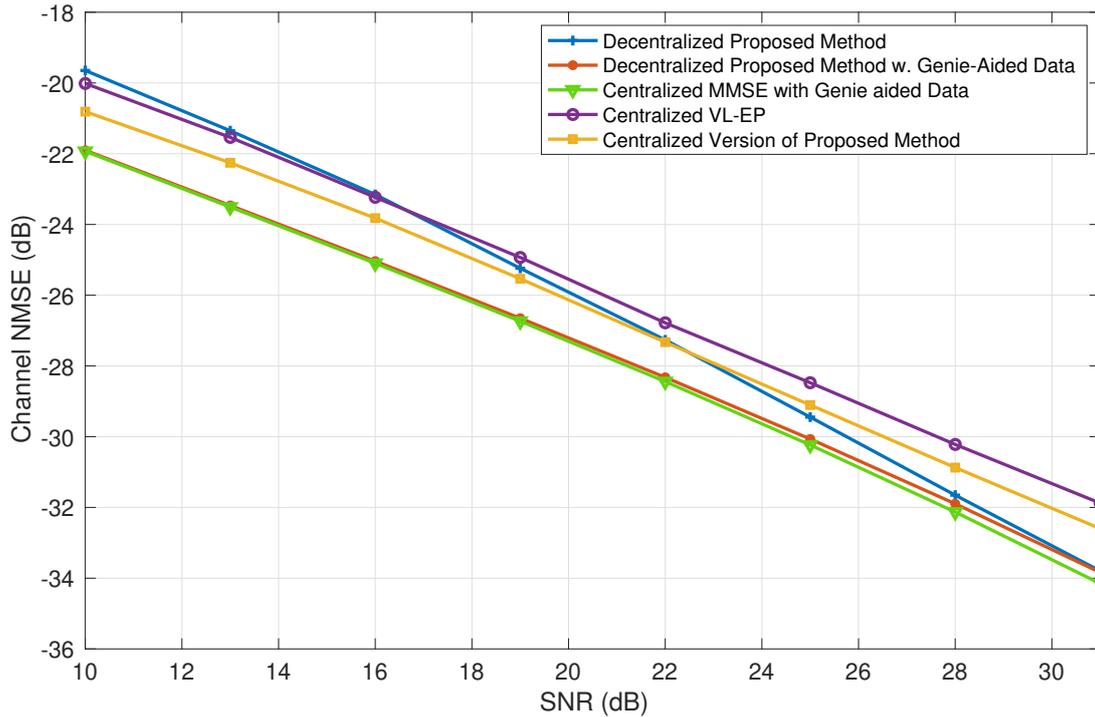


Figure 4: NMSE vs SNR

15 Concluding Remarks

In this paper, we studied the BFE of GLMs using a joint factorization scheme. This factorization allows us to extract approximate priors and likelihood. By looking at the stationary point in LSL we replace the non-separable constraints with separable ones. This leads to the LSL BFE. This paper also interprets extrinsics for both input and output nodes as CWCU LMMSE estimation operations.

We rederived the reGVAMP algorithm from the point of view of alternating minimization of a LSL version of a desirable KLD. The asymptotics here involve only the CLT for extrinsics. We then derive the GAMP algorithm by directly introducing LSL simplifications in the LBP algorithm. This leads us to relate extrinsic messages to posterior pdfs by first order Taylor series expansion based perturbations. We also apply LSL approximations to the variances of the various Gaussians involved, which in fact leads to a rederivation of a fundamental LSA theorem describing the deterministic limit of LMMSE posterior variances.

We have also shown that it is possible to derive the convergent AMBGAMP algorithm by analyzing the KKT conditions for optimizing the LSL BFE. And this while avoiding the quadratic augmentation terms of the Method of Moments, which require a very particular choice in their weights, and circumventing the ADMM-style update of a Lagrange multiplier. This is thanks to the introduction of the auxiliary variable \mathbf{u} in the mean consistency constraints. On the other hand, another solution to the LSL BFE, which eliminates \mathbf{u} via $\mathbf{u} = \hat{\mathbf{x}}$, leads to GAMP and corresponds to the original LSL BP based derivation, optimizing BFE with LSL approximations. Hence we have reconciled these seemingly different approaches.

The variance predictions in (AMB)GAMP are based on a sign i.i.d. model for \mathbf{A} , which leads to decorrelation and Gaussianity after multiplication of a vector with \mathbf{A} or \mathbf{A}^T , similar to spreading and despreading in CDMA. Another somewhat popular model for \mathbf{A} is the Right Rotationally Invariant class, in which (only) the right singular vectors of \mathbf{A} are modeled as random, and in particular as Haar distributed. This is the motivation for Vector AMP (VAMP) [24]. To keep complexity low however, VAMP has to restrict diagonal covariances to multiples of identity, which

e.g. is not useful for Sparse Bayesian Learning [25]. GAMP-style low complexity algorithms can be derived also, but they require some correction terms in the variance predictions, stemming from the Haar distribution [26], [27].

To investigate the bilinear asymptotics, this paper also introduces a simplified, decentralized EP-based algorithm for bilinear joint estimation. To simplify the factorization scheme, we leverage orthogonal pilots and the CLT. Through asymptotic analysis, we further refine the message update scheme within the algorithm. Although originally developed for an acyclic network of APs, our simulation results confirm the algorithm's effectiveness even when the APs are interconnected in a cyclic network.

16 References

- [1] D. Zhang, X. Song, W. Wang, G. Fettweis, and X. Gao, "Unifying Message Passing Algorithms Under the Framework of Constrained Bethe Free Energy Minimization," *IEEE Transactions on Wireless Communications*, vol. 20, no. 7, 2021.
- [2] M. J. Wainwright, M. I. Jordan, *et al.*, "Graphical Models, Exponential Families, and Variational Inference," *Foundations and Trends® in Machine Learning*, vol. 1, no. 1–2, 2008.
- [3] K. Murphy, Y. Weiss, and M. I. Jordan, "Loopy Belief Propagation for Approximate Inference: An Empirical Study," *arXiv preprint arXiv:1301.6725*, 2013.
- [4] T. Minka *et al.*, "Divergence Measures and Message Passing," tech. rep., Citeseer, 2005.
- [5] T. Heskes, M. Opper, W. Wiegner, O. Winther, and O. Zoeter, "Approximate Inference Techniques with Expectation Constraints," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2005, no. 11, 2005.
- [6] Q. Zou and H. Yang, "A Concise Tutorial on Approximate Message Passing," *arXiv preprint arXiv:2201.07487*, 2022.
- [7] S. Rangan, P. Schniter, E. Riegler, A. K. Fletcher, and V. Cevher, "Fixed Points of Generalized Approximate Message Passing with Arbitrary Matrices," *IEEE Transactions on Information Theory*, vol. 62, no. 12, 2016.
- [8] C. Kurisumoothil Thomas, Z. Zhao, and D. Slock, "Towards Convergent Approximate Message Passing by Alternating Constrained Minimization of Bethe Free Energy," in *IEEE Information Theory Workshop (ITW)*, (Saint Malo, France), 2023.
- [9] S. Rangan, A. Fletcher, P. Schniter, and U. Kamilov, "Inference for Generalized Linear Models via Alternating Directions and Bethe Free Energy Minimization," *IEEE Trans. Info. Theory*, Jan. 2017.
- [10] Z. Zhao and D. Slock, "Bethe Free Energy and Extrinsic in Approximate Message Passing," in *IEEE Asilomar Conf. Signals, Systems and Computers*, 2023.
- [11] Z. Zhao, F. Xiao, and D. Slock, "Vector approximate message passing for not so large N.I.I.D. generalized I/O linear models," in *IEEE Int'l Conf. Acoustics, Speech and Sig. Proc. (ICASSP)*, (Seoul), 2024.
- [12] Z. Zhao, F. Xiao, and D. Slock, "Extrinsic and Linearized Component-Wise Conditionally Unbiased MMSE Estimation as in GAMP," in *IEEE Asilomar Conf. Signals, Systems and Computers*, 2024.
- [13] M. Triki and D. T. Slock, "Component-Wise Conditionally Unbiased Bayesian Parameter Estimation: General Concept and Applications to Kalman Filtering and LMMSE Channel Estimation," in *Conference Record of the Thirty-Ninth Asilomar Conference on Signals, Systems and Computers, 2005.*, IEEE.
- [14] M. Huemer and O. Lang, "CWCU LMMSE Estimation: Prerequisites and Properties," *arXiv preprint arXiv:1412.1567*, 2014.
- [15] C. Sippel and R. F. Fischer, "Variants of VAMP for Signal Recovery in Wireless Sensor Networks," in *2022 IEEE International Mediterranean Conference on Communications and Networking (MeditCom)*, 2022.
- [16] Z. Zhao, F. Xiao, and D. Slock, "Approximate Message Passing for Not So Large niid Generalized Linear Models," in *Proc. Int'l Workshop Signal Processing Advances in Wireless Comm's (SPAWC)*, Sept. 2023.

-
- [17] J. T. Parker, P. Schniter, and V. Cevher, “Bilinear Generalized Approximate Message Passing—Part I: Derivation,” *IEEE Transactions on Signal Processing*, 2014.
- [18] Q. Zou, H. Zhang, C.-K. Wen, S. Jin, and R. Yu, “Concise Derivation for Generalized Approximate Message Passing Using Expectation Propagation,” *IEEE Signal Processing Letters*, 2018.
- [19] A. Karataev, C. Forsch, and L. Cottatellucci, “Bilinear Expectation Propagation for Distributed Semi-Blind Joint Channel Estimation and Data Detection in Cell-Free Massive MIMO,” *IEEE Open Journal of Signal Processing*, 2024.
- [20] Z. Zhao, F. Xiao, and D. Slock, “Approximate Message Passing for Not So Large NIID Generalized Linear Models,” in *Int’l Workshop on Signal Processing Advances in Wireless Comm’s (SPAWC)*, 2023.
- [21] P. Schniter, “A Simple Derivation of AMP and its State Evolution via First-Order Cancellation,” *IEEE Transactions on Signal Processing*, 2020.
- [22] C. C. Moallemi and B. Van Roy, “Consensus Propagation,” *IEEE Transactions on Information Theory*, 2006.
- [23] R. Gholami, L. Cottatellucci, and D. Slock, “Message Passing for a Bayesian Semi-Blind Approach to Cell-Free Massive MIMO,” in *Asilomar Conference on Signals, Systems, and Computers*, IEEE, 2021.
- [24] S. Rangan, P. Schniter, and A. K. Fletcher, “Vector Approximate Message Passing,” *IEEE Trans. On Info. Theo.*, Oct. 2019.
- [25] C. K. Thomas and D. Slock, “SAVE - Space alternating variational estimation for sparse Bayesian learning,” in *IEEE Data Science Workshop*, June 2018.
- [26] Z. Zhao and D. Slock, “Variance Predictions in VAMP/UAMP with Right Rotationally Invariant Measurement Matrices for niid Generalized Linear Models,” in *European Sig. Proc. Conf. (EUSIPCO)*, (Helsinki, Finland), 2023.
- [27] Z. Zhao and D. Slock, “Improved Variance Predictions in Approximate Message Passing,” in *IEEE Int’l Workshop Machine Learning and Sig. Proc. (MLSP)*, (Rome, Italy), 2023.