# Bilinear Recovery using Adaptive Vector-AMP

**Subrata Sarkar (Ohio State)**
**Philip Schniter (Ohio State)**

(supported by NSF 1716388)

THE OHIO STATE UNIVERSITY

Asilomar 2019

**Alyson K. Fletcher (UCLA)**
**Sundeep Rangan (NYU)**

## Bilinear Recovery Problem

- Observations:

$$Y = \sum_{i=1}^{Q} b_i A_i X + W$$

where,

$$X : \text{unknown random matrix in } \mathbb{R}^{N \times L}$$
$$A_1, \ldots, A_Q : \text{known matrices in } \mathbb{R}^{M \times N}$$
$$b_1, \ldots, b_Q : \text{unknown deterministic parameters}$$
$$W : \text{white Gaussian noise.}$$

- Prior:

$$X_{nl} \overset{\text{i.i.d}}{\sim} p_X(\cdot; \boldsymbol{\theta}_X) \quad \text{deterministic unknown parameters } \boldsymbol{\theta}_X.$$
$$W_{ml} \overset{\text{i.i.d}}{\sim} \mathcal{N}(0, \sigma_w^2) \quad \text{unknown variance } \sigma_w^2.$$

**Goal**: jointly infer $X$ and estimate $\boldsymbol{\theta} \triangleq \{b, \boldsymbol{\theta}_X, \sigma_w^2\}$
**Approach**: combine variational inference with ML estimation.
**Applications**: Self-calibration, CS+matrix uncertainty, dictionary learning, ...

## Variational Inference

- For now, let's suppose that $\boldsymbol{\theta}$ is known.

- We would like to compute the posterior density

$$p(X|Y) = \frac{p(X; \boldsymbol{\theta}) p(Y|X; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} \quad \text{for } Z(\boldsymbol{\theta}) \triangleq \int p(X; \boldsymbol{\theta}) p(Y|X; \boldsymbol{\theta}) \, dX,$$

but the high-dimensional integral in $Z(\boldsymbol{\theta})$ is difficult to compute.

- We can avoid computing $Z(\boldsymbol{\theta})$ through variational optimization:

$$p(X|Y) = \arg\min_b D\big(b(X) \big\| p(X|Y)\big) \text{ where } D(\cdot\|\cdot) \text{ is KL divergence}$$
$$= \arg\min_b \underbrace{D\big(b(X) \big\| p(X; \boldsymbol{\theta})\big) + D\big(b(X) \big\| p(Y|X; \boldsymbol{\theta})\big) + H\big(b(X)\big)}_{\text{Gibbs free energy}}$$
$$= \arg\min_{b_1, b_2, q} \underbrace{D\big(b_1(X) \big\| p(X; \boldsymbol{\theta})\big) + D\big(b_2(X) \big\| p(Y|X; \boldsymbol{\theta})\big) + H\big(q(X)\big)}_{\triangleq J(b_1, b_2, q; \boldsymbol{\theta})}$$

such that $b_1 = b_2 = q$,

but the density constraint keeps the problem difficult.

- Expectation consistent approximation (EC) [1] relaxes the density constraint to moment-matching constraints:

$$p(X|Y) \approx \arg\min_{b_1, b_2, q} J(b_1, b_2, q; \boldsymbol{\theta})$$

$$\text{such that } \forall l \begin{cases} \mathbb{E}\{x_l|b_1\} = \mathbb{E}\{x_l|b_2\} = \mathbb{E}\{x_l|q\} \\ \text{tr}[\text{Cov}\{x_l|b_1\}] = \text{tr}[\text{Cov}\{x_l|b_2\}] = \text{tr}[\text{Cov}\{x_l|q\}]. \end{cases}$$

- The stationary points of EC are the densities

$$b_1(X) \propto \prod_{l=1}^{L} p(x_l; \boldsymbol{\theta}) \mathcal{N}(x_l; r_{1,l}, I/\gamma_{1,l})$$
$$b_2(X) \propto \prod_{l=1}^{L} p(y_l|x_l; \boldsymbol{\theta}) \mathcal{N}(x_l; r_{2,l}, I/\gamma_{2,l}) \text{ s.t.}$$
$$q(X) = \prod_{l=1}^{L} \mathcal{N}(x_l; \widehat{x}_l, I/\eta_l)$$

$$\begin{cases} \mathbb{E}\{x_l|b_1\} = \mathbb{E}\{x_l|b_2\} = \widehat{x}_l \\ \text{tr}[\text{Cov}\{x_l|b_1\}] \\ = \text{tr}[\text{Cov}\{x_l|b_2\}] = N/\eta_l. \end{cases}$$

## Vector AMP (VAMP)

- There exist several algorithms (e.g., EC, ADATAP [2], S-AMP [3]) whose fixed points coincide with the EC stationary points, but often they don't converge.

- An exception is Vector AMP [4], which can be derived using a form of approximate message passing on the vector-valued factor graph

$$p(X_1; \boldsymbol{\theta}) \blacksquare \!-\! \bigcirc \overset{X_1}{\underset{\delta(X_1 - X_2)}{\phantom{X}}} \overset{X_2}{\bigcirc} \!-\! \blacksquare \, p(Y|X_2; \boldsymbol{\theta})$$

In particular, VAMP is provably convergent under either
1) strictly log-concave prior $p(X; \boldsymbol{\theta})$ and arbitrary $A$ (after damping),
2) iid prior $p(X; \boldsymbol{\theta})$ and large, right-rotationally invariant $A$.

## VAMP algorithm

Initialize $\{R_1, \gamma_1\}$ and define the estimation functions

$$g_1(r_{1,l}, \gamma_{1,l}) \triangleq \mathbb{E}\{x_l|b_1; r_{1,l}, \gamma_{1,l}\} \text{ or any Lipschitz function}$$
$$g_2(r_{2,l}, \gamma_{2,l}) \triangleq \mathbb{E}\{x_l|b_2; r_{2,l}, \gamma_{2,l}\}$$

For $t = 1, 2, 3, \ldots$

$$\widehat{x}_{1,l} \leftarrow g_1(r_{1,l}, \gamma_{1,l}), \forall l \quad \text{denoising}$$
$$\eta_{1,l} \leftarrow \gamma_{1,l} N / \text{tr}[\partial g_1(r_{1,l}; \gamma_{1,l})/\partial r_{1,l}], \forall l$$
$$r_{2,l} \leftarrow (\eta_{1,l}\widehat{x}_{1,l} - \gamma_{1,l} r_{1,l})/(\eta_{1,l} - \gamma_{1,l}), \forall l \quad \text{pseudo-measurement}$$
$$\gamma_{2,l} \leftarrow \eta_{1,l} - \gamma_{1,l}, \forall l$$

$$\widehat{x}_{2,l} \leftarrow g_2(r_{2,l}, \gamma_{2,l}), \forall l \quad \text{LMMSE estimation}$$
$$\eta_{2,l} \leftarrow \gamma_{2,l} N / \text{tr}[\partial g_2(r_{2,l}; \gamma_{2,l})/\partial r_{2,l}], \forall l$$
$$r_{1,l} \leftarrow (\eta_{2,l}\widehat{x}_{2,l} - \gamma_{2,l} r_{2,l})/(\eta_{2,l} - \gamma_{2,l}), \forall l \quad \text{pseudo-prior}$$
$$\gamma_{1,l} \leftarrow \eta_{2,l} - \gamma_{2,l}, \forall l$$

## Expectation maximization (EM)

- We now return to the case where $\boldsymbol{\theta} = \{b, \boldsymbol{\theta}_X, \sigma_w^2\}$ is unknown.

- The maximum-likelihood (ML) estimate is

$$\widehat{\boldsymbol{\theta}} = \arg\max_{\boldsymbol{\theta}} p(Y; \boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta}} \{-\ln p(Y; \boldsymbol{\theta})\}.$$

- EM algorithm iteratively minimizes a tight upper bound on $-\ln p(Y; \boldsymbol{\theta})$:

$$\widehat{\boldsymbol{\theta}}^{t+1} = \arg\min_{\boldsymbol{\theta}} \mathbb{E}\big\{ -\ln p(X, Y; \boldsymbol{\theta}) | Y; \widehat{\boldsymbol{\theta}}^t \big\}$$
$$= \arg\min_{\boldsymbol{\theta}} \Big\{ -\ln p(Y; \boldsymbol{\theta}) + \underbrace{D\big(b^t(X) \big\| p(X|Y; \boldsymbol{\theta})\big)}_{\geq 0} \Big\}$$
$$\text{with } b^t(X) = p(X|Y; \widehat{\boldsymbol{\theta}}^t)$$

- The upper bound can also be rewritten in terms of Gibbs free energy

$$Q(\boldsymbol{\theta}, b^t) \triangleq -\ln p(Y; \boldsymbol{\theta}) + D(b^t(X) \| p(X|Y; \boldsymbol{\theta}))$$
$$= D\big(b^t(X) \big\| p(X; \boldsymbol{\theta})\big) + D\big(b^t(X) \big\| p(Y|X; \boldsymbol{\theta})\big) + H\big(b^t(X)\big)$$
$$= J(b^t, b^t, b^t; \boldsymbol{\theta})$$

which yields a variational interpretation of EM [5].

## Variance Auto-Tuning

- In VAMP, the precisions $\{\gamma_{1,l}, \gamma_{2,l}\}_{l=1}^{L}$ are imperfect when $\widehat{\boldsymbol{\theta}}$ is imperfect.

- So we estimate these precisions jointly with $\boldsymbol{\theta}$. E.g., for parameters $\boldsymbol{\theta}_X$:

$$(\boldsymbol{\gamma}_1, \widehat{\boldsymbol{\theta}}_X) \leftarrow \arg\max_{\boldsymbol{\gamma}_1, \boldsymbol{\theta}_X} p(R_1; \boldsymbol{\gamma}_1, \boldsymbol{\theta}_X) \text{ under } r_{1,l} = x_l + \mathcal{N}(0, I/\gamma_{1,l}), \; x_l \sim p(\cdot; \boldsymbol{\theta}_X)$$

- In practice, inner iterations of EM are used to solve the above "variance auto-tuning" problem.

- Under identifiability conditions, this leads to asymptotically consistent $\widehat{\boldsymbol{\theta}}_X$ [6].

## The proposed Bilinear Adaptive (BAd)-VAMP Algorithm

- Recall that VAMP iteratively computes a posterior approximation $b^t(X)$ by minimizing $J(b_1, b_2, q; \boldsymbol{\theta})$ (under moment constraints) with known $\boldsymbol{\theta}$.

- Likewise, EM iteratively estimates $\boldsymbol{\theta}$ by minimizing $J(b^t, b^t, b^t; \boldsymbol{\theta})$, assuming the posterior approximation $b^t(X) = p(X|Y; \boldsymbol{\theta}^t)$ is available.

- In BAd-VAMP, we combine VAMP, EM, and variance auto-tuning.

- In the denoising ($i = 1$) and LMMSE ($i = 2$) steps of VAMP, we infer $X$ and jointly estimate $(\boldsymbol{\gamma}, \boldsymbol{\theta})$ by running several inner iterations of

$$\forall l : \widehat{x}_{i,l} \leftarrow g_i(r_{i,l}, \gamma_{i,l}; \widehat{\boldsymbol{\theta}}_i), \quad \eta_{i,l} \leftarrow \gamma_{i,l} N / \text{tr}[\partial g_i(r_{i,l}, \gamma_{i,l}; \widehat{\boldsymbol{\theta}}_i)/\partial r_{i,l}] \quad \text{(denoising)}$$

$$\forall l : \gamma_{i,l} \leftarrow \Big( \frac{1}{N}\|\widehat{x}_{i,l} - r_{i,l}\|^2 + 1/\eta_{1,l} \Big)^{-1} \quad \text{(auto-tuning)}$$

$$q_i(X) \propto \prod_{l=1}^{L} f_i(x_{i,l}; \boldsymbol{\theta}_i) \mathcal{N}(x_l; r_{i,l}, I/\gamma_{i,l}) \quad \text{(belief update)}$$

$$\widehat{\boldsymbol{\theta}}_i \leftarrow \arg\max_{\boldsymbol{\theta}_i} \sum_{l=1}^{L} \mathbb{E}\left[\ln f_i(x_{i,l}, \boldsymbol{\theta}_i)|q_i\right] \quad \text{(EM update)}$$

with $i \in \{1, 2\}$, $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_X, \boldsymbol{\theta}_2 = \{b, \sigma_w^2\}$ and

$$f_1(x, \boldsymbol{\theta}_1) = p(x; \boldsymbol{\theta}_X), \quad f_2(x, \boldsymbol{\theta}_2) = \mathcal{N}\left(y; \sum_{j=1}^{Q} b_j A_j x, \sigma_w^2 I\right).$$
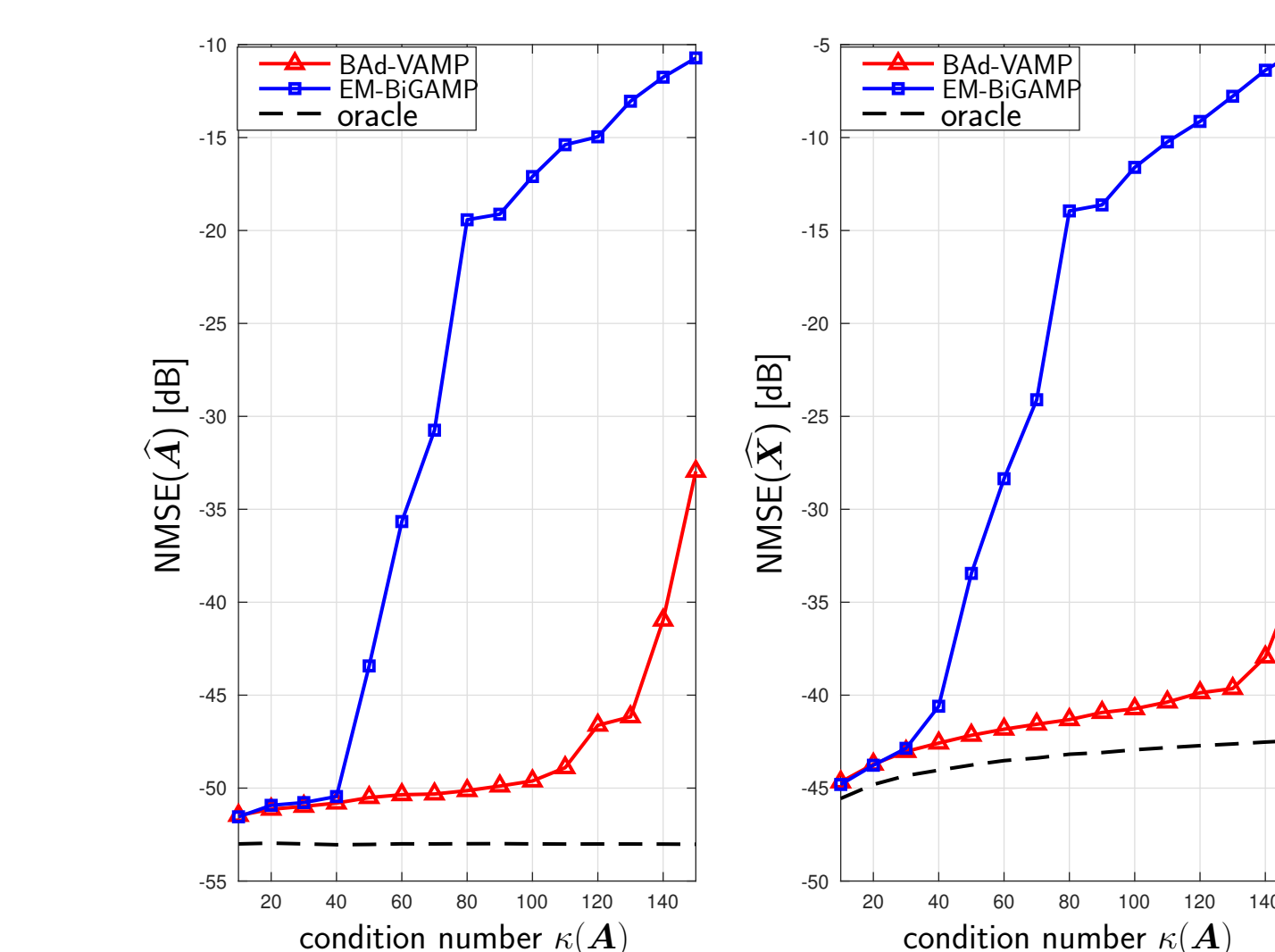
## Numerical Experiments

**CS with Matrix Uncertainty**: Recover $N = 256$-length sparse $x$ and $b_i \overset{\text{i.i.d}}{\sim} \mathcal{N}(0, 1)$ from $M$-length $y = (A_0 + \sum_{i=1}^{10} b_i A_i)x + w$, where $[A_i]_{m,n} \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$,



$x$ is sampled from the Bernoulli-Gaussian distribution, $\epsilon = 0.04$,

$$x_i \sim (1 - \epsilon)\delta(x) + \epsilon \mathcal{N}(x; 0, 1)$$

Plots show median NMSE on signal $\widehat{x}$ and parameter $\widehat{b}$ estimates versus sampling ratio $M/N$.

BAd-VAMP performs much better than WSS-TLS [7] and on par with EM-PBiGAMP [8] and VAMP-Lift [9].

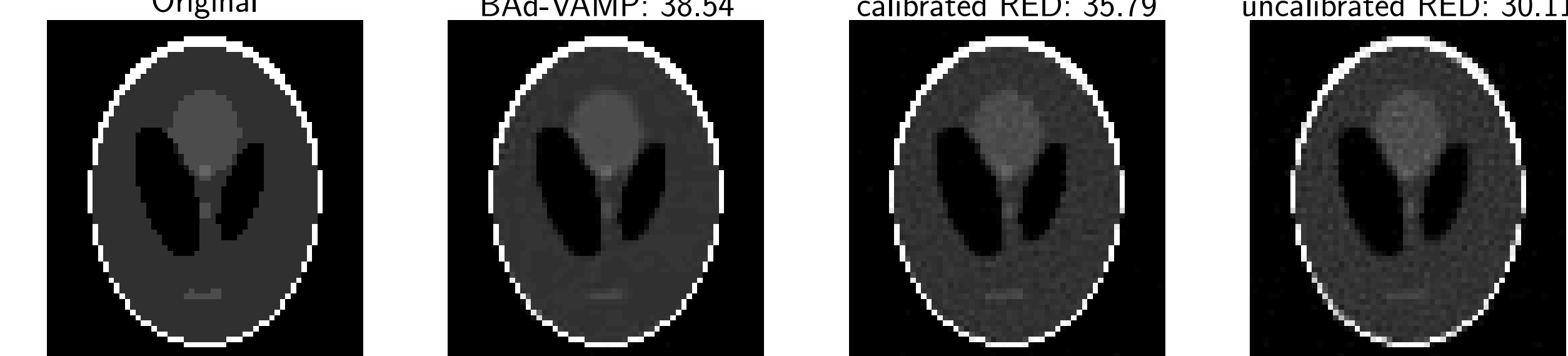**Dictionary Learning**: Given $Y$, the goal is to estimate $A$ & $X$ s.t. $Y \approx AX$.



Median NMSE versus condition number $\kappa(A)$ for $A \in \mathbb{R}^{64 \times 64}$, i.i.d. Bernoulli-Gaussian $X$ with $\epsilon = 0.2$, training length $L = 1331$, and SNR $= 40$ dB.

BAd-VAMP is much more robust to $\kappa(A)$ than EM-BiGAMP [10].

**Self-Calibration in Tomography**: Reconstruct image $x$ and calibration parameters $b_i \sim \mathcal{N}(1, \sigma_b^2)$ from measurements $y = \big[b_1\Psi_{\omega_1}^\top \; \ldots \; b_{25}\Psi_{\omega_{25}}^\top\big]^\top x + w$, where $\Psi_\omega$ is the Radon transform for angle $\omega$ and $\sigma_b = 0.06$. BM3D was used for $g_1(\cdot)$.



Original | BAd-VAMP: 38.54 | calibrated RED: 35.79 | uncalibrated RED: 30.11

PSNR (dB) of $64 \times 64$ Shepp-Logan phantom from 25 equally spaced tomographic projections.

## References

[1] M. Opper and O. Winther, "Expectation consistent approximate inference," *J. Mach. Learning Res.*, 2005. .

[2] M. Opper and O. Winther, "Adaptive and self-averaging Thouless-Anderson-Palmer mean-field theory for probabilistic modeling," *Phy. Rev. E*, 2001. .

[3] B. Çakmak, O. Winther, and B.H. Fleury, "S-AMP: Approximate message passing for general matrix ensembles," *ISIT* 2014.

[4] S. Rangan, P. Schniter, and A. K. Fletcher, "Vector approximate message passing," *arXiv:1610.03082.*

[5] R. Neal and G. Hinton, "A view of the EM algorithm that justifies incremental, sparse, and other variants," in *Learning in Graphical Models*, 1998.

[6] A. K. Fletcher, M. Sahraee-Ardakan, S. Rangan, and P. Schniter, "Rigorous dynamics and consistent estimation in arbitrarily conditioned linear systems," *Proc. NeurIPS*, 2017.

[7] H. Zhu, G. Leus, and G. B. Giannakis, "Sparsity-cognizant total least-squares for perturbed compressive sampling," *IEEE TSP*, 2011. .

[8] J. T. Parker, P. Schniter, "Parametric bilinear generalized approximate message passing", *IEEE JSTSP*, 2016.

[9] A. K. Fletcher, P. Pandit, S. Rangan, S. Sarkar, and P. Schniter, "Plug-in estimation in high-dimensional linear inverse problems: A rigorous analysis," *Proc. NeurIPS*, 2018.

[10] J. T. Parker, P. Schniter and V. Cevher, "Bilinear generalized approximate message passing", *IEEE TSP*, 2014.

[11] Y. Romano, M. Elad, and P. Milanfar, "The little engine that could: Regularization by denoising (RED)," *SIAM J. Imaging Science*, 2017.