

# BP-VB-EP Based Static and Dynamic Sparse Bayesian Learning with Kronecker Structured Dictionaries

Christo Kurisummoottil Thomas, Dirk Slock

Communication Systems Department, EURECOM, Sophia Antipolis, France



ICASSP 2020, May 4-8,  
Special Session:

“AMP and other Approximate Bayesian Inference Techniques”

# Outline

- 1 Introduction
- 2 Combined BP-MF-EP Framework
- 3 Kronecker Structured Dictionary Learning using BP/VB
- 4 Numerical Results and Conclusion



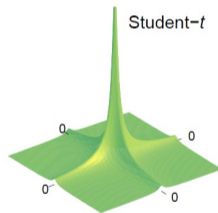
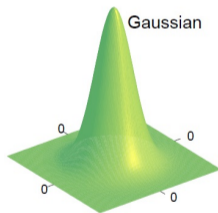
# Sparsification of the Innovation Sequence (Sparse Bayesian Learning-SBL)

- Bayesian Compressed Sensing: 2-layer hierarchical prior for  $\mathbf{x}$  as in [Tipping:JMLR01], inducing sparsity for  $\mathbf{x}$ .

$$p(x_{i,t}|\lambda_i) = \mathcal{N}(x_{i,t}; 0, \lambda_i^{-1}), \quad p(\lambda_i/a, b) = \Gamma^{-1}(a)b^a \lambda_i^{a-1} e^{-b\lambda_i}$$

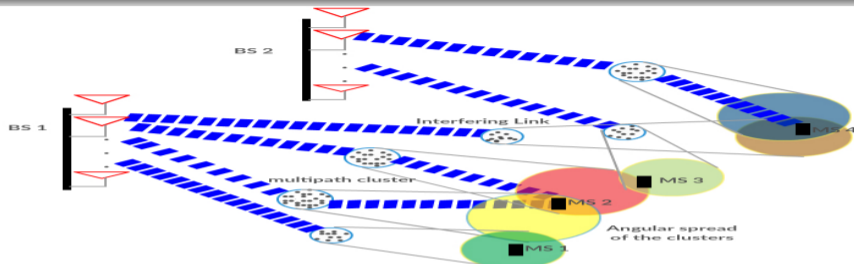
⇒ sparsifying Student-t marginal

$$p(x_{i,t}) = \frac{b^a \Gamma(a + \frac{1}{2})}{(2\pi)^{\frac{1}{2}} \Gamma(a)} (b + w_i^2/2)^{-(a+\frac{1}{2})}$$



- We apply the (Gamma) prior not to the precision of the state  $\mathbf{x}$  but it's innovation  $\mathbf{w}$ , allowing to sparsify at the same time the components of  $\mathbf{x}$  AND their variation in time (innovation).

# Tensor Representation (Channel Tracking in MaMIMO OFDM)

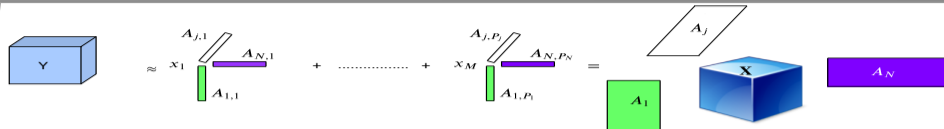


- Sampling across Doppler space and stacking all the subcarrier and sampled (in Doppler) elements as a vector

$$\text{vec}(\mathbf{H}(t)) = \sum_{i=1}^L x_{i,t} \mathbf{h}_t(\phi_i) \otimes \mathbf{h}_r(\theta_i) \otimes \mathbf{v}_f(\tau_i) \otimes \mathbf{v}_t(f_i) = \mathbf{A}(t)\mathbf{x}_t$$

- 4-D Tensor model, Delay, Doppler and Tx/Rx spatial dimensions.
- Array response itself: Kronecker structure in the case of polarization or in the case of 2D antenna arrays with separable structure [Sidiropoulos:icassp18].
- User mobility changes the scattering geometry and path coefficients.
- Tensor based KF proposed here **avoids the off-grid basis mismatch issues.**

# Kronecker Structured Tensor Models



- Tensor signals appear in many applications: massive multi-input multi-output (MIMO) radar, massive MIMO (MaMIMO) channel estimation, speech processing, image and video processing.
- Retaining the tensorial structure in the received signal can be beneficial compared to processing its unstructured matrix version.
- The signal model for the recovery of a time varying sparse signal under Kronecker structured (KS) dictionary matrix can be formulated as

$$\text{Observation: } \mathbf{y}_t = \underbrace{(\mathbf{A}_1(t) \otimes \mathbf{A}_2(t) \dots \otimes \mathbf{A}_N(t))}_{\mathbf{A}(t)} \mathbf{x}_t + \mathbf{v}_t, \quad \mathbf{y}_t = \text{vec}(\mathbf{Y}_t)$$

$$\text{State Update: } \mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \mathbf{w}_t,$$

$\mathbf{Y}_t \in \mathcal{C}^{l_1 \times l_2 \times \dots \times l_N}$  is the observations or data at time  $t$ ,  $\mathbf{A}_{j,i}(t) \in \mathcal{C}^{l_j}$ , the factor matrix  $\mathbf{A}_j(t) = [\mathbf{A}_{j,1}(t), \dots, \mathbf{A}_{j,p_j}(t)]$  and the overall unknown parameters are  $\llbracket \mathbf{A}_1(t), \dots, \mathbf{A}_N(t); \mathbf{x}_t \rrbracket$ ,  $\mathbf{x}_t$  is  $M (= \prod_{j=1}^N P_j)$ -dimensional sparse center tensor and  $\mathbf{w}_t$ ,  $\mathbf{v}_t$  are the state or measurement noise.

# Our Contribution Here Compared to State of the Art

- Extension of our previous work<sup>1</sup> [ThomasSlock:ICASSP2019] on Kronecker Structured (KS) dictionary learning (DL) - called as SAVED-KS DL (Space Alternating Variational Estimation for KS DL) to the case of time varying sparse state vector.
- Variational Bayes (VB) at the scalar level i.e mean field (MF) for  $\mathbf{x}_t$  and at the column level for the factor matrices in the measurement matrix is quite suboptimal since the posterior covariance computed does not converge to the true posterior [ThomasSlock:DSW2018]. One potential solution: **belief propagation (BP)**.
- Also, **better variational free energy (VFE) approximation**, e.g. Belief Propagation (BP) or Expectation Propagation (EP) instead of MF.
- Inspired by the framework of [RieglerFleury:TIT2013], we combine BP and MF approximations in such a way as to optimize the message passing (MP) framework.
- The main focus: **reduced complexity (e.g. VB leads to the low complexity) algorithms without sacrificing much on the performance (BP for performance improvement)**.
- Dynamic autoregressive SBL (DAR-SBL)**: a case of joint Kalman filtering (KF) with a linear time-invariant diagonal state-space model, and parameter estimation  $\implies$  **an instance of nonlinear filtering**.

---

<sup>1</sup>abscd

# Outline

- 1 Introduction
- 2 Combined BP-MF-EP Framework**
- 3 Kronecker Structured Dictionary Learning using BP/VB
- 4 Numerical Results and Conclusion



# Variational Free Energy (VFE) Framework

- Intractable joint posterior distribution of the parameters  $\theta = \{\mathbf{x}, \mathbf{A}, \mathbf{f}, \boldsymbol{\lambda}, \gamma\}$ .
- Actual posterior:  $p(\theta) = \frac{1}{Z} \underbrace{\prod_{a \in \mathcal{A}_{BP}} f_a(\theta_a) \prod_{b \in \mathcal{A}_{MF}} f_b(\theta_b)}_{\text{factor nodes}}$ , where  $\mathcal{A}_{BP}, \mathcal{A}_{MF}$  = set of variable nodes belonging to the BP/MF part with  $\mathcal{A}_{BP} \cap \mathcal{A}_{MF} = \emptyset$ .
- The whole  $\theta$  is partitioned into the set  $\theta_i$ , and we want to approximate the true posterior  $p(\theta)$  by an approximate posterior  $q(\theta) = \prod_i q_i(\theta_i)$ . The individual variable portions  $\theta_i$  are the variable nodes in a factor graph.
- $\mathcal{N}_{BP}(i), \mathcal{N}_{MF}(i)$  – the set of neighbouring factor nodes of variable node  $i$  which belong to the BP/MF part.
- $\mathcal{I}_{MF} = \bigcup_{a \in \mathcal{A}_{MF}} \mathcal{N}(a), \mathcal{I}_{BP} = \bigcup_{a \in \mathcal{A}_{BP}} \mathcal{N}(a)$ .  $\mathcal{N}(a)$  – the set of neighbouring variable nodes of any factor node  $a$ .
- The resulting Free Energy (Entropy – Average Energy) obtained by the combination of BP and MF are written as below (let  $q_i(\theta_i)$  represents the belief about  $\theta_i$  (the approximate posterior)),

$$F_{BP,MF} = \sum_{a \in \mathcal{A}_{BP}} \sum_{\theta_a} q_a(\theta_a) \ln \frac{q_a(\theta_a)}{f_a(\theta_a)} - \sum_{b \in \mathcal{A}_{MF}} \sum_{\mathbf{x}_b} \prod_{i \in \mathcal{N}(b)} q_i(\theta_i) \ln f_b(\theta_b) - \sum_{i \in \mathcal{I}} (|\mathcal{N}_{BP}(i)| - 1) \sum_{\theta_i} q_i(\theta_i) \ln q_i(\theta_i).$$

# Message Passing Expressions

- The beliefs have to satisfy the following **normalization and marginalization constraints**,

$$\sum_{\theta_i} q_i(\theta_i) = 1, \forall i \in \mathcal{I}_{MF} \setminus \mathcal{I}_{BP}, \quad \sum_{\theta_a} q_a(\theta_a) = 1, \forall a \in \mathcal{A}_{BP},$$

$$q_i(\theta_i) = \sum_{\theta_a \setminus \theta_i} q_a(\theta_a), \quad \forall a \in \mathcal{A}_{BP}, i \in \mathcal{N}(a).$$

- The **fixed point equations corresponding to the constrained optimization of the approximate VFE** can be written as follows [RieglerFleuryTIT2013],

$$q_i(\theta_i) = z_i \prod_{a \in \mathcal{N}_{BP}(i)} m_{a \rightarrow i}^{BP}(\theta_i) \prod_{a \in \mathcal{N}_{MF}(i)} m_{a \rightarrow i}^{MF}(\theta_i), \implies \text{Product of incoming beliefs}$$

$$n_{i \rightarrow a}(\theta_i) = \prod_{c \in \mathcal{N}_{BP}(i) \setminus a} m_{c \rightarrow i}(\theta_i) \prod_{d \in \mathcal{N}_{MF}(i)} m_{d \rightarrow i}(\theta_i), \implies \text{variable to factor nodes}$$

$$m_{a \rightarrow i}^{MF}(\theta_i) = \exp(\langle \ln f_a(\theta_a) \rangle_{j \in \mathcal{N}(a) \setminus i} n_{j \rightarrow a}(\theta_j)),$$

$$m_{a \rightarrow i}^{BP}(\theta_i) = \left( \int \prod_{j \in \mathcal{N}(a) \setminus i} n_{j \rightarrow a}(\theta_j) f_a(\theta_a) \prod_{j \neq i} d\theta_j \right), \implies \text{factor to variable nodes}$$

where  $\langle \rangle_q$  represents the expectation w.r.t distribution  $q$ .

# Gaussian BP-MF-EP KF

- Proposed Method: Alternating optimization between non linear KF for the sparse states (plus the hyperparameters) and BP for dictionary learning (DL).
- Diagonal AR(1) ( DAR(1) ) Prediction Stage:** Since there is no coupling between the scalars in the state update, it is enough to update the prediction stage using MF. However, the interaction between  $x_{l,t}$  and  $f_l$  requires Gaussian projection, using expectation propagation (EP). More details in [ThomasSlock:Asilomar2019].
- $y_n$  – factor node,  $x_l$  – variable node.  $(l, n)$  or  $(n, l)$  to represent the messages passed from  $l$  to  $n$  or viceversa. Gaussian messages from  $y_n$  to  $x_l$  parameterized by mean  $\hat{x}_{n,l}^{(t)}$  and variance  $\nu_{n,l}^{(t)}$ .
- Measurement Update (Filtering) Stage:** For the measurement update stage, the posterior for  $x_t$  is inferred using BP. In the measurement stage, the prior for  $x_{l,t}$  gets replaced by the belief from the prediction stage. We define  $d_{l,t} = \left( \sum_{n=1}^N \nu_{n,l}^{(t)-1} \right)^{-1}$ ,  $r_{l,t} = d_{l,t} \left( \sum_{n=1}^N \frac{\hat{x}_{n,l}^{(t)}}{\nu_{n,l}^{(t)}} + \frac{\hat{x}_{l,t|t-1}}{\sigma_{l,t|t-1}^2} \right)$ .

$$\sigma_{l,t|t}^{-2} = \lambda_{l,t} + d_{l,t}^{-1}, \quad \hat{x}_{l,t|t} = \frac{r_{l,t}}{1 + d_{l,t} \sigma_{l,t|t}^{-2}}.$$

# Lag-1 Smoothing Stage for Correlation Coefficient $\mathbf{f}$



$$\mathbf{y}_t = \mathbf{A}(t)\mathbf{F}\mathbf{x}_{t-1} + \tilde{\mathbf{v}}_t, \text{ where } \tilde{\mathbf{v}}_t = \mathbf{A}(t)\mathbf{w}_t + \mathbf{v}_t, \tilde{\mathbf{v}}_t \sim \mathcal{CN}(0, \tilde{\mathbf{R}}_t)$$

- We show in Lemma 1 of [ThomasSlock:Asilomar2019] that **KF is not enough to adapt the hyperparameters, instead we need at least a lag 1 smoothing** (i.e. the computation of  $\hat{\mathbf{x}}_{k,t-1|t}, \sigma_{k,t-1|t}^2$  through BP). For the smoothing stage, we use BP.
- Gaussian Posterior for  $\mathbf{x}_t$ :

$$\begin{aligned} \sigma_{k,t-1|t}^{-2,(i)} &= (\hat{f}_{k|t}^2 + \sigma_{f_{k|t}}^2)\mathbf{A}_k^H(t)\tilde{\mathbf{R}}_t^{-1}\mathbf{A}_k(t) + \sigma_{k,t-1|t-1}^{-2}, \\ \langle \mathbf{x}_{k,t-1|t}^{(i)} \rangle &= \sigma_{k,t-1|t}^{2,(i)}(\hat{f}_{k|t}^H\mathbf{A}_k^H(t)\tilde{\mathbf{R}}_t^{-1}(\mathbf{y}_t - \mathbf{A}_k(t)\mathbf{F}_{k|t}) \langle \mathbf{x}_{k,t-1|t}^{(i-1)} \rangle + \frac{\hat{\mathbf{x}}_{k,t-1|t-1}}{\sigma_{k,t-1|t-1}^2}). \end{aligned}$$

- Applying the MF rule, the resulting Gaussian distribution for  $\mathbf{f}$  has mean,  $\sigma_{f_i|t}^{-2}$  and variance,  $\hat{\mathbf{f}}_{i|t}$ , the detailed derivations for which are in [ThomasSlock:Asilomar2019].

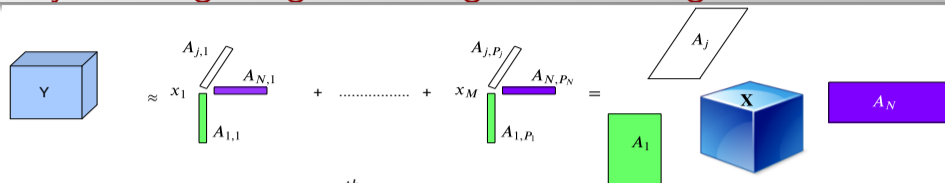
$$\begin{aligned} \sigma_{f_i|t}^{-2} &= (|\hat{\mathbf{x}}_{i,t-1|t}|^2 + \sigma_{i,t-1|t}^2)\mathbf{A}_i^T(t)\tilde{\mathbf{R}}_t^{-1}\mathbf{A}_i(t), \\ \hat{\mathbf{f}}_{i|t} &= \sigma_{f_i|t}^2\hat{\mathbf{x}}_{i,t-1|t}^H\mathbf{A}_i^H(t)\tilde{\mathbf{R}}_t^{-1}(\mathbf{y}_t - \mathbf{A}_i(t)\hat{\mathbf{F}}_{i|t}\hat{\mathbf{x}}_{i,t-1|t}). \end{aligned}$$

- $\tilde{\mathbf{R}}_t = \mathbf{A}(t)\mathbf{\Lambda}^{-1}\mathbf{A}(t)^H + \frac{1}{\gamma}\mathbf{I}$ .

# Outline

- 1 Introduction
- 2 Combined BP-MF-EP Framework
- 3 Kronecker Structured Dictionary Learning using BP/VB**
- 4 Numerical Results and Conclusion

# Dictionary Learning using Tensor Signal Processing



- Let  $Y_{i_1, \dots, i_N}$  represents the  $i_1 i_2 \dots i_N^{\text{th}}$  element of the tensor and  $\mathbf{y} = [y_{1,1, \dots, 1}, y_{1,1, \dots, 2, \dots}, y_{1, i_2, \dots, i_N}]^T$ , then it can be verified that [Sidiropoulos:TSP17],  
 $\mathbf{y}_t = (\mathbf{A}_1(t) \otimes \mathbf{A}_2(t) \dots \otimes \mathbf{A}_N(t)) \mathbf{x}_t + \mathbf{w}_t, \mathbf{w} \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$ ,  
 Matrix Unfolding:  $\mathbf{Y}^{(n)} = \mathbf{A}_n(t) \mathbf{X}^{(n)} (\mathbf{A}_N(t) \otimes \dots \otimes \mathbf{A}_{n+1}(t) \otimes \mathbf{A}_{n-1}(t) \dots \otimes \mathbf{A}_1(t))^T$ .
- $\mathbf{A}_j(t)$  is of dimension,  $l_j \times P_j$  and the resulting Tensor is  $\mathcal{C}^{l_1 \times \dots \times l_N}$ .
- Retaining the Tensor structure in the dictionary matrix leads to better estimates than using the matricized version for  $\mathbf{A}$  and learning it.
- Less free variables to be estimated in the Tensor structured case.
- Variational Bayesian Inference using the following approximate posterior

$$q(\mathbf{x}, \alpha, \gamma, \mathbf{A}) = q_\gamma(\gamma) \prod_{i=1}^M q_{x_i}(x_i) \prod_{i=1}^M q_{\alpha_i}(\alpha_i) \prod_{j=1}^N \prod_{i=1}^{P_j} q_{\mathbf{A}_{j,i}}(\mathbf{A}_{j,i}), \implies \text{SAVED-KS DL Or}$$

$$q(\mathbf{x}, \alpha, \gamma, \mathbf{A}) = q_\gamma(\gamma) \prod_{i=1}^M q_{x_i}(x_i) \prod_{i=1}^M q_{\alpha_i}(\alpha_i) \prod_{j=1}^N q_{\mathbf{A}_j}(\mathbf{A}_j) \implies \text{Joint VB for DL}$$

## Suboptimality of SAVED-KS DL and Joint VB

- From the expression for the error covariance in the estimation of the factor  $\mathbf{A}_{j,i}$  (SAVED-KS DL in [ThomasSlock:ICASSP19]) ( $\text{tr}\{(\bigotimes_{k=N, k \neq j}^1 \langle \mathbf{A}_k^T(t) \mathbf{A}_k^*(t) \rangle \langle \mathbf{X}^{(j)T} \mathbf{X}^{(j)} \rangle)\mathbf{I}\}$ ),  $\implies$  it does not take into account the estimation error in the other columns of  $\mathbf{A}_j(t)$ . The columns of  $\mathbf{A}_j(t)$  can be correlated, for e.g. if we consider two paths (say  $i, j$ ) with same DoA but with different delays, the delay responses  $\mathbf{v}_f(\tau_i(t))$  and  $\mathbf{v}_f(\tau_j(t))$  may be correlated.
- The joint VB estimates (mean and covariance) can be obtained as

$$\begin{aligned} \mathbf{M}_j^T &= \widehat{\mathbf{A}}_{1,j}^T(t) = \langle \gamma \rangle \Psi_j^{-1} \mathbf{B}_j^T, \\ \Psi_j &= (\langle \gamma \rangle \langle \mathbf{X}^{(j)} (\bigotimes_{k=N, k \neq j}^1 \langle \mathbf{A}_k^T(t) \mathbf{A}_k^*(t) \rangle) \mathbf{X}^{(j)T} \rangle), \end{aligned} \quad (1)$$

where  $\mathbf{V}_j = \langle \mathbf{X}^{(j)} \rangle \langle (\bigotimes_{k=N, k \neq j}^1 \mathbf{A}_k(t))^T \rangle$  and  $\mathbf{B}_j$  is defined as with the first row of  $(\mathbf{Y}^{(j)} \mathbf{V}_j^T)$  removed. However, the joint VB involves a matrix inversion and is not recommended for large system dimensions. Nevertheless, it is possible to estimate each columns of  $\mathbf{A}_j(t)$  by BP, since each column estimate can be expressed as the solution of a linear system of equation from (1),  $\widehat{\mathbf{A}}_{j,i}^T(t) = \Psi_j^{-1} \mathbf{b}_{j,i}$ .  $\mathbf{b}_{j,i}$  represents the  $i^{\text{th}}$  column of  $\mathbf{B}_j^T$ .

# Optimal Partitioning of the Measurement Stage and KS DL

## Lemma

For the measurement stage, an optimal partitioning is to apply BP for the sparse vector  $\mathbf{x}_t$  and VB (SAVED-KS) for the columns of the factor matrices  $\mathbf{A}_{j,i}(t)$  assuming the vectors  $\mathbf{A}_{j,i}(t)$  are independent and have zero mean. However, if the columns of  $\mathbf{A}_j(t)$  are correlated, then a joint VB, with the posteriors of the factor matrices assumed independent, should be done for an optimal performance.

- Proof: Follows from [KalyanSlock:EUSIPCO2019, Lemma 1], where the main message was that if the parameter partitioning in VB is such that the different parameter blocks are decoupled at the level of FIM, then VB is not suboptimal in terms of (mismatched) Cramer-Rao Bound (mCRB).

$$\mathbf{y}_t = \underbrace{\left(\sum_{r=1}^M x_{r,t} \mathbf{F}_r\right)}_{F(\mathbf{x}_t)} \underbrace{\left(\otimes_{j=1}^N \Phi_{j,t}\right)}_{f(\Phi_t)} + \mathbf{w}_t. \quad \mathbf{J}(\Phi_t) = [\mathbf{J}(\Phi_{1,t}) \dots \mathbf{J}(\Phi_{N,t})]$$

where,  $\mathbf{J}(\Phi_{j,t}) = \mathbf{F}(\mathbf{x}_t)(\Phi_{1,t} \otimes \dots \otimes \mathbf{I}_{l_j p_j} \otimes \dots \otimes \Phi_{N,t})$ ,

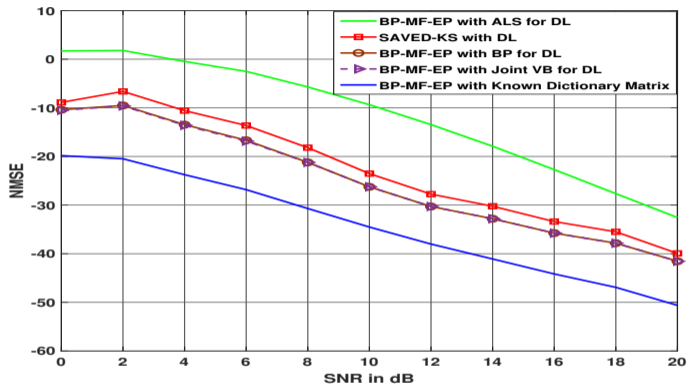
$$FIM = \begin{bmatrix} E(\gamma) \mathbf{J}(\Phi_t)^T \mathbf{J}(\Phi_t) & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & E(\gamma) \mathbf{J}(\mathbf{x}_t)^T \mathbf{J}(\mathbf{x}_t) + E(\Gamma^{-1}) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & aE(\Gamma^{-2}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & (N + c - 1)E(\gamma^{-2}) \end{bmatrix}$$



# Outline

- 1 Introduction
- 2 Combined BP-MF-EP Framework
- 3 Kronecker Structured Dictionary Learning using BP/VB
- 4 Numerical Results and Conclusion

# BP-MF-EP Outperforms SAVED-KS DL



Static SBL: NMSE as a function of  $N$ .

- ALS- Alternating Least Squares.
- Exponential power delay profile for  $x_t$ .
- 30 non zero elements in  $x_t$ , same support across all time.
- Dimensions: 3-D Tensor (4, 8, 8), with  $M = 200$ .

# Conclusions and Thank You!

- Further advancements from [ThomasSlock:ICASSP19]: VB with a too fine variable partitioning is quite suboptimal.
- Better approximation is message passing based methods such as belief propagation (BP) and expectation propagation (EP)
- BP or EP message passing can be implemented using low complexity methods such as AMP/GAMP/VAMP, which are proven to be Bayes optimal under certain conditions on  $\mathbf{A}$ .
- AMP - Approximate message passing. We also derived an Generalized Vector AMP (GVAMP-SBL) SBL version to take care of diagonal power delay profile.
- Further work to be done on learning a combination of structured and unstructured Kronecker factor matrices.