

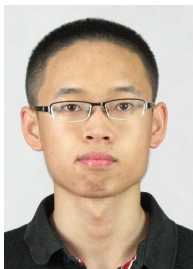
# Approximate Message Passing: A Non-asymptotic Framework And Beyond



Yuting Wei

Statistics & Data Science, Wharton  
University of Pennsylvania

Optimization Seminar @UPenn



Gen Li, UPenn → CUHK

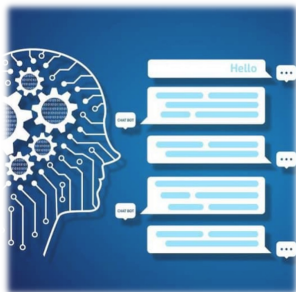
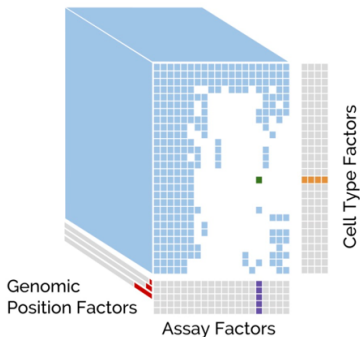


Wei Fan, UPenn

*"A non-asymptotic framework for approximate message passing in spiked models,"*  
Gen Li, Yuting Wei, *arxiv.2208.03313*

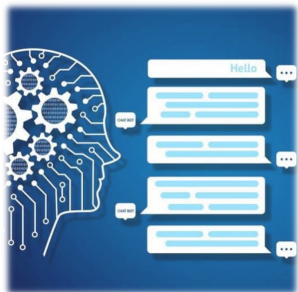
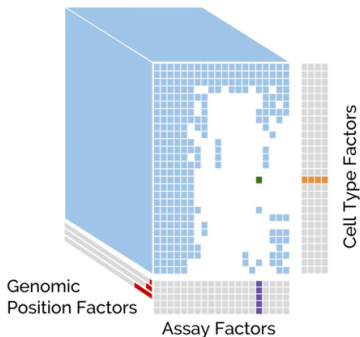
*"Approximate message passing from random initialization with applications to  $\mathbb{Z}_2$  synchronization,"* Gen Li, Wei Fan, Yuting Wei, *PNAS*, 2023

# High-dimensional statistical tasks



**Statistical tasks:** solution to convex/non-convex optimization problems  
*e.g. linear regression, generalized linear models, low-rank matrix estimation, phase retrieval, tensor decomposition...*

# High-dimensional statistical tasks

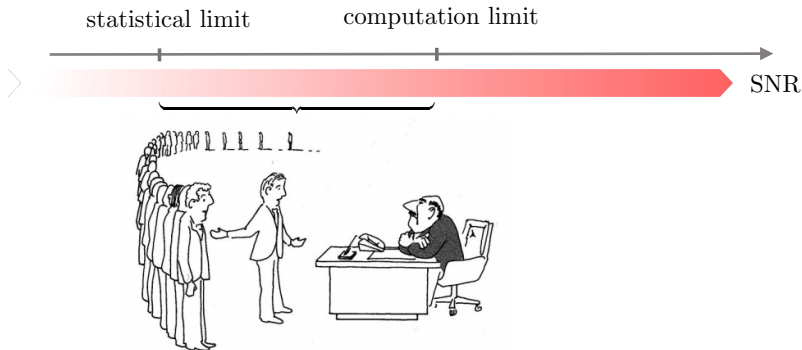


**Statistical tasks:** solution to convex/non-convex optimization problems  
*e.g. linear regression, generalized linear models, low-rank matrix estimation, phase retrieval, tensor decomposition...*

When problem sizes are large, **computation complexity** is an issue!

# Statistical accuracy vs. computation complexity

**statistical-to-computational gap** in problems with combinatorial nature (e.g. *community detection, planted cliques, sparse principal component analysis, structured matrix models, sparse tensor models...*)

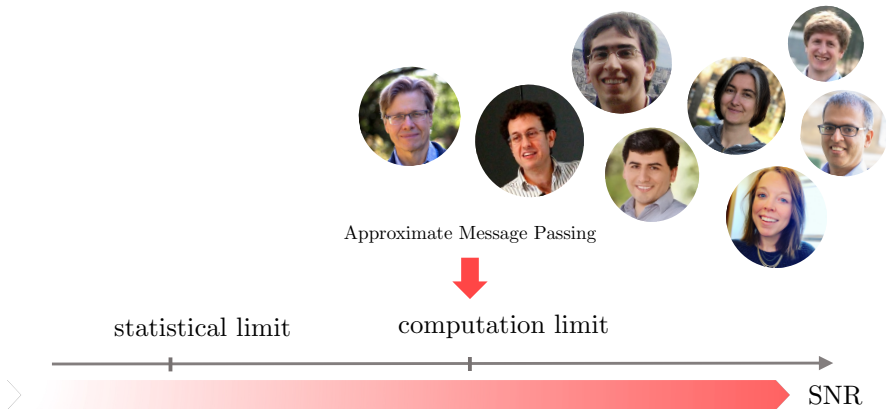


*"I can't find an efficient algorithm, but neither can all these people."*

— see survey [Bandeira, Perry, Wein'18](#)

# Statistical accuracy vs. computation complexity

**statistical-to-computational gap** in problems with combinatorial nature (e.g. community detection, planted cliques, sparse principal component analysis, structured matrix models, sparse tensor models...)



— see tutorial [Feng, Venkataramanan, Rush, Samworth' 22](#)

# A simple model: spiked Wigner model

$$M = \lambda \begin{matrix} \text{red} & \text{red} & \text{light red} & \text{red} & \text{light red} & \text{light red} & \text{red} \\ \text{light red} & & & & & & \\ \text{red} & & & & & & \\ \text{light red} & & & & & & \\ \text{light red} & & & & & & \\ \text{red} & & & & & & \end{matrix} v^{\star\top} + \begin{matrix} W \\ \text{blue} & \text{blue} & \text{blue} & \text{light blue} & \text{blue} & \text{light blue} & \text{light blue} \\ \text{light blue} & \text{blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} \\ \text{blue} & \text{light blue} & \text{blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} \\ \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} \\ \text{blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} \\ \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} & \text{light blue} \end{matrix}$$

$M = \lambda v^{\star} v^{\star\top} + W$



Johnstone (2001),

# A simple model: spiked Wigner model

$$M = \lambda \begin{matrix} \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \end{matrix} \begin{matrix} \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \\ \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \end{matrix} v^{\star\top} + \begin{matrix} \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \end{matrix} W$$

The diagram illustrates the matrix equation  $M = \lambda v^{\star} v^{\star\top} + W$ . The matrix  $M$  is shown as a 7x7 grid of red squares, with a vertical column of red squares labeled  $v^{\star}$  and a horizontal row of red squares labeled  $v^{\star\top}$ . The matrix  $W$  is shown as a 7x7 grid of blue squares. A plus sign is placed between the two matrices.

- $W_{ij} = W_{ji} \sim \mathcal{N}(0, \frac{1}{n})$  and  $W_{ii} \sim \mathcal{N}(0, \frac{2}{n})$







# Spiked Wigner model with structures

$$M = \lambda \begin{matrix} \text{[red gradient bar]} \\ v^* \end{matrix} + \begin{matrix} v^{*\top} \\ \text{[blue gradient bar]} \\ W \end{matrix}$$

**Applications:** spin-glass problems, community detection, image alignment, angular synchronization

# Spiked Wigner model with structures

$$M = \lambda \begin{matrix} \color{red}{\square} \\ \color{red}{\square} \\ \color{red}{\square} \\ \color{red}{\square} \\ \color{red}{\square} \\ \color{red}{\square} \\ \color{red}{\square} \end{matrix} \begin{matrix} \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} & \color{red}{\square} \end{matrix} + \begin{matrix} \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \\ \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} & \color{blue}{\square} \end{matrix}$$

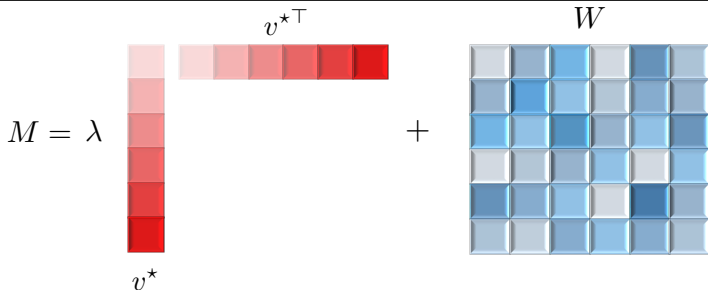
$v^*$   $v^{*\top}$   $W$

**Applications:** spin-glass problems, community detection, image alignment, angular synchronization

- $\mathbb{Z}_2$  synchronization:  $\sqrt{n}v_i^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}\{+1, -1\}$



# Spiked Wigner model with structures

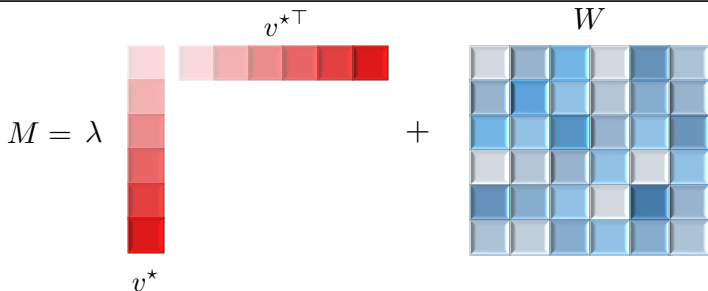


**Applications:** spin-glass problems, community detection, image alignment, angular synchronization

- $\mathbb{Z}_2$  synchronization:  $\sqrt{n}v_i^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}\{+1, -1\}$
- sparse Wigner model:  $\|v^*\|_0 = k$
- non-negative Wigner model:  $v_i^* \geq 0$

Singer (2011), Panchenko (2013), Deshpande, Abbe & Montanari (2016), Perry, Wein, Bandeira, Moitra (2018), Javanmard, Montanari & Ricci-Tersenghi (2016)...

# Spiked Wigner model with structures



**Applications:** spin-glass problems, community detection, image alignment, angular synchronization

- $\mathbb{Z}_2$  synchronization:  $\sqrt{n}v_i^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}\{+1, -1\}$
- sparse Wigner model:  $\|v^*\|_0 = k$
- non-negative Wigner model:  $v_i^* \geq 0$
- cone-constrained spiked models:  $v^* \in \mathcal{K}$  (e.g. [monotone](#), [convex](#))

Singer (2011), Panchenko (2013), Deshpande, Abbe & Montanari (2016), Perry, Wein, Bandeira, Moitra (2018), Javanmard, Montanari & Ricci-Tersenghi (2016)...

# An incomplete list of prior art

---

$\mathbb{Z}_2$  synchronization:

- Baik, Arous, P  ch  '05
- Panchenko'13
- Javanmard et al.'16
- Montanari & Sen'16
- Lelarge & Miolane'19
- Deshpande, Abbe, Montanari'17
- Celentano, Fan, Mei'21

general convex cones:

- Deshpande, Montanari, Richard'14
- Lesieur, Krzakala, Zdeborov  '17
- Bandeira, Kunisky, Wein'19

sparse PCA (Wigner / Wishart)

- Johnstone & Lu'09
- d'Aspremont et al.'04
- Amini & Wainwright'08
- Vu & Lei'12
- Berthet & Rigollet'13
- Ma'13
- Lesieur, Krzakala, Zdeborov  '15
- Deshpande & Montanari'14
- Wang, Berthet, Samworth'16
- Ding, Kunisky, Wein, Bandeira'19

positive Wigner models

- Montanari & Richard'16



# Idealistic estimators

---

Maximum likelihood estimator  $:= \arg \min_{\substack{v \in \mathcal{S}^{n-1} \\ v \text{ with structures}}} \|M - \lambda v v^\top\|_F^2$

Bayes optimal estimator  $:= \mathbb{E}[v v^\top \mid M]$

# AMP for spiked models

---

Maximum likelihood estimator  $:= \arg \min_{\substack{v \in \mathcal{S}^{n-1} \\ v \text{ with structures}}} \|M - \lambda v v^\top\|_F^2$

Bayes optimal estimator  $:= \mathbb{E}[v v^\top \mid M]$

— *in general, computationally infeasible...*

# AMP for spiked models

---

Approximate message passing (AMP) for spiked models:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ for } t \geq 1$$

where  $\langle x \rangle := \frac{1}{n} \sum_{i=1}^n x_i$ .

# AMP for spiked models

---

Approximate message passing (AMP) for spiked models:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ for } t \geq 1$$

where  $\langle x \rangle := \frac{1}{n} \sum_{i=1}^n x_i$ .

- Onsager correction term  $\langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1})$

# AMP for spiked models

Approximate message passing (AMP) for spiked models:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ for } t \geq 1$$

where  $\langle x \rangle := \frac{1}{n} \sum_{i=1}^n x_i$ .

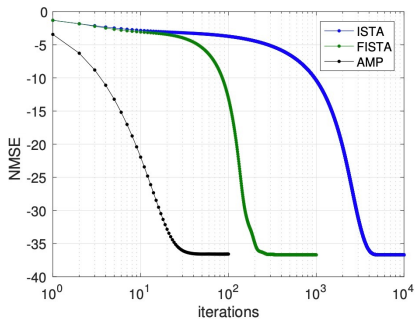
- Onsager correction term  $\langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1})$
- $\eta_t$ : denoising function selected *a priori* (tailored to structure of  $v^*$ )
  - ▶  $\mathbb{Z}_2$  **synchronization**:  $\eta_t(x) = \rho_t \tanh(x)$
  - ▶ **sparse estimation**:  $\eta_t(x) = \rho_t \cdot \text{sign}(x)(|x| - \tau_t)_+$
  - ▶ **general cone**:  $\eta_t(x) = \rho_t \cdot \text{Proj}_{\mathcal{K}}(x)$

# Some background of AMP

---

- AMP is a low-complexity, iterative algorithm

[Donoho, Maleki, Montanari (2009, 2010a, 2011b), Bayati & Montanari (2011)...]



AMP in computing LASSO

# Some background of AMP

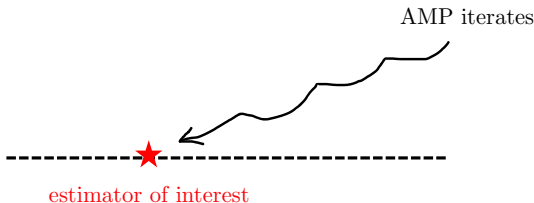
---

- AMP is a low-complexity, iterative algorithm  
[Donoho, Maleki, Montanari (2009, 2010a, 2011b), Bayati & Montanari (2011)...]
- Theoretically optimal vs. computationally feasible estimators  
[Reeves, Pfister (2019), Barbier et al. (2017), Lelarge & Miolane (2019), Montanari & Ramji (2019), Celentano & Montanari (2019)...]

# Some background of AMP

---

- AMP is a low-complexity, iterative algorithm  
[Donoho, Maleki, Montanari (2009, 2010a, 2011b), Bayati & Montanari (2011)...]
- Theoretically optimal vs. computationally feasible estimators  
[Reeves, Pfister (2019), Barbier et al. (2017), Lelarge & Miolane (2019), Montanari & Ramji (2019), Celentano & Montanari (2019)...]
- A useful tool to analyze other statistical procedures [Donoho, Maleki, Montanari (2009), Donoho & Montanari (2016), Sur, Chen, Candès. (2017), Bu et al. (2020), Fan & Wu (2021), Li & Wei (2021)...]

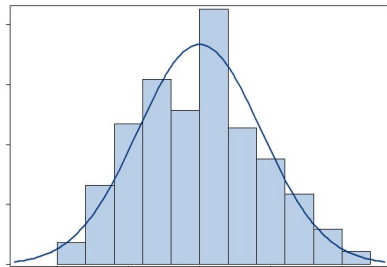




# Prior theory of AMP

---

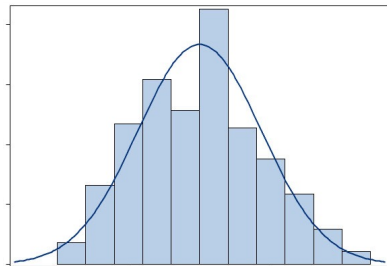
**Exact asymptotics:** for constant # iterations  $t$  (e.g.  $t = 20$ ), empirical distribution of the coordinates of AMP iterate  $x_t \in \mathbb{R}^n$  is approximately Gaussian ( $n \rightarrow \infty$ )



histogram of coordinates of  $x_t$

# Prior theory of AMP

**Exact asymptotics:** for constant # iterations  $t$  (e.g.  $t = 20$ ), empirical distribution of the coordinates of AMP iterate  $x_t \in \mathbb{R}^n$  is approximately Gaussian ( $n \rightarrow \infty$ )



histogram of coordinates of  $x_t$

Its variance is given by low-dimensional recursion:

state evolution:  $\tau_{t+1} = F(\tau_t)$

$\tau_t$  captures the variance at iteration  $t$

[Bayati & Montanari (2011), Javanmard & Montanari (2013), Schniter & Rangan (2014)]

## Prior results: exact asymptotics

---

### Theorem (Montanari & Venkataramanan'19)

Suppose the empirical distribution  $\{v_i^*\}_{i=1}^n \rightarrow \mu_V$  on  $\mathbb{R}$ , with  $\mathbb{E}[V^2] = 1$ . For constant # iterations  $t$  (*independent of  $n$* ), it satisfies,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (x_{t,i} - v_i^*)^2 = \mathbb{E} \left[ (\alpha_t V + \beta_t G - V)^2 \right], \quad \text{a.s.}$$

where  $V \sim \mu_V$  and  $G \sim \mathcal{N}(0, 1)$  are independent.

## Prior results: exact asymptotics

### Theorem (Montanari & Venkataramanan'19)

Suppose the empirical distribution  $\{v_i^*\}_{i=1}^n \rightarrow \mu_V$  on  $\mathbb{R}$ , with  $\mathbb{E}[V^2] = 1$ . For constant # iterations  $t$  (independent of  $n$ ), it satisfies,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (x_{t,i} - v_i^*)^2 = \mathbb{E}[(\alpha_t V + \beta_t G - V)^2], \quad \text{a.s.}$$

where  $V \sim \mu_V$  and  $G \sim \mathcal{N}(0, 1)$  are independent.

- State evolution (SE) via the recursion

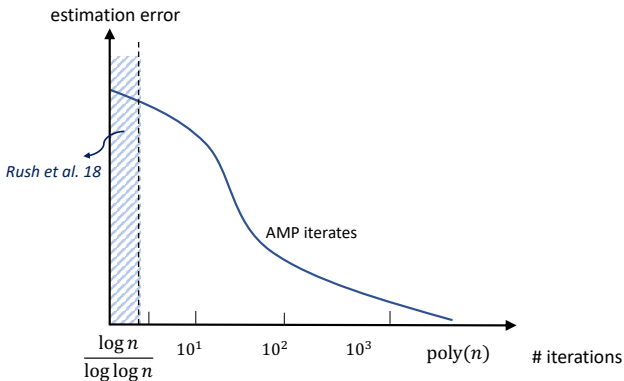
$$(\alpha_{t+1}, \beta_{t+1}) = F(\alpha_t, \beta_t) = \begin{cases} \alpha_{t+1} = \lambda \mathbb{E}[V \cdot \eta_t(\alpha_t V + \beta_t G)] \\ \beta_{t+1}^2 = \mathbb{E}[\eta_t^2(\alpha_t V + \beta_t G)] \end{cases}$$

*Non-asymptotic analyses are quite limited so far...*

- compared to other optimization methods
- compared to other analysis techniques



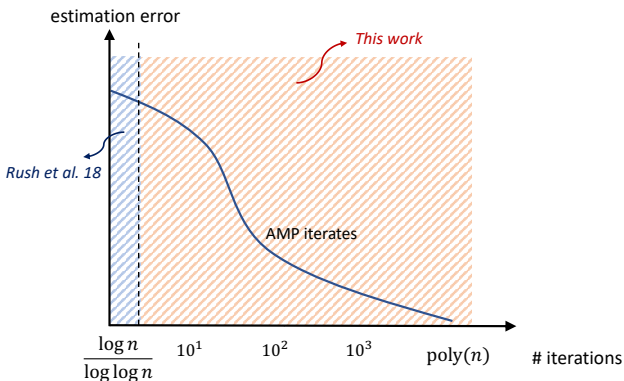
# Non-asymptotic analysis?



**Non-asymptotic result:** Rush & Venkataramanan (2018)

$\# \text{iterations} = o(\log n / \log \log n)$  (based on state-evolution analysis)

# Non-asymptotic analysis?



**Question:** Is it possible to develop non-asymptotic analysis of AMP beyond  $o(\log n / \log \log n)$  iterations?

*Our solution: a new decomposition for AMP iterates*



# This work: a new decomposition of AMP

## Theorem (Li & Wei'22)

Initialize AMP with  $x_1$  independent of  $W$ . For every  $1 \leq t \leq n$ , AMP yields the decomposition

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t, \quad (*)$$

for  $\phi_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{1}{n} I_n)$ .

# This work: a new decomposition of AMP

## Theorem (Li & Wei'22)

Initialize AMP with  $x_1$  independent of  $W$ . For every  $1 \leq t \leq n$ , AMP yields the decomposition

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t, \quad (*)$$

for  $\phi_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{1}{n} I_n)$ .

here  $(\alpha_{t+1}, \beta_t, \xi_t)$  obeys

$$\alpha_{t+1} = \lambda v^{*\top} \eta_t(x_t),$$

$$\beta_t^k = \langle \eta_t(x_t), z_k \rangle \quad \text{for an explicit-defined basis } \{z_k\}$$

$$\|\xi_t\|_2 = \left\langle \sum_{k=1}^{t-1} \mu^k \phi_k, \delta_t \right\rangle - \langle \delta'_t \rangle \sum_{k=1}^{t-1} \mu^k \beta_{t-1}^k + \Delta_t + O\left(\sqrt{\frac{t \log n}{n}} \|\beta_t\|_2\right) \quad \text{w.h.p.}$$

# This work: a new decomposition of AMP

## Theorem (Li & Wei'22)

Initialize AMP with  $x_1$  independent of  $W$ . For every  $1 \leq t \leq n$ , AMP yields the decomposition

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t, \quad (*)$$

for  $\phi_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{1}{n} I_n)$ .

- $x_t$  behaves like  $\alpha_t v^* + \sum_{k=1}^{t-1} \beta_{t-1}^k \phi_k$  if  $\|\xi_{t-1}\|_2$  is small

$$\text{Wasserstein}_1 \left( \mu \left( \frac{1}{\|\beta_{t-1}\|_2} \sum_{k=1}^{t-1} \beta_{t-1}^k \phi_k \right), \mathcal{N} \left( 0, \frac{1}{n} I_n \right) \right) \leq \sqrt{\frac{t \log n}{n}}.$$

# This work: a new decomposition of AMP

## Theorem (Li & Wei'22)

Initialize AMP with  $x_1$  independent of  $W$ . For every  $1 \leq t \leq n$ , AMP yields the decomposition

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t, \quad (*)$$

for  $\phi_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{1}{n} I_n)$ .

- $x_t$  behaves like  $\alpha_t v^* + \sum_{k=1}^{t-1} \beta_{t-1}^k \phi_k$  if  $\|\xi_{t-1}\|_2$  is small
- if  $\{\eta_t\}$  are nice (smooth & with finite jumps), we can track how  $\|\xi_t\|_2$  depends on  $\lambda, t, n$

# This work: a new decomposition of AMP

## Theorem (Li & Wei'22)

Initialize AMP with  $x_1$  independent of  $W$ . For every  $1 \leq t \leq n$ , AMP yields the decomposition

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t, \quad (\star)$$

for  $\phi_k \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \frac{1}{n} I_n)$ .

- $x_t$  behaves like  $\alpha_t v^* + \sum_{k=1}^{t-1} \beta_{t-1}^k \phi_k$  if  $\|\xi_{t-1}\|_2$  is small
- if  $\{\eta_t\}$  are nice (smooth & with finite jumps), we can track how  $\|\xi_t\|_2$  depends on  $\lambda, t, n$
- decomposition  $(\star)$  can be extended for spectral initialization

# Finite-sample error control

## Theorem (Li & Wei'22 (informal))

AMP iterates satisfy  $x_{t+1} = \alpha_{t+1}v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$  w.h.p. with

$$\alpha_{t+1} = \lambda v^{*\top} \int \eta_t \left( \alpha_t v^* + \frac{1}{\sqrt{n}} x \right) \varphi_n(dx) + \lambda \Delta_{\alpha,t}, \quad \|\beta_t\|_2 = 1,$$

where the residual terms obey

$$|\Delta_{\alpha,t}| \lesssim B_t + \rho \|\xi_{t-1}\|_2,$$

$$\|\xi_t\|_2 \leq \kappa_t \|\xi_{t-1}\|_2 + O\left( A_t + \rho \sqrt{\frac{t \log n}{n}} \|\xi_{t-1}\|_2 \right).$$

# Finite-sample error control

## Theorem (Li & Wei'22 (informal))

AMP iterates satisfy  $x_{t+1} = \alpha_{t+1}v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$  w.h.p. with

$$\alpha_{t+1} = \lambda v^{*\top} \int \eta_t \left( \alpha_t v^* + \frac{1}{\sqrt{n}} x \right) \varphi_n(dx) + \lambda \Delta_{\alpha,t}, \quad \|\beta_t\|_2 = 1,$$

where the residual terms obey

$$|\Delta_{\alpha,t}| \lesssim B_t + \rho \|\xi_{t-1}\|_2,$$

$$\|\xi_t\|_2 \leq \kappa_t \|\xi_{t-1}\|_2 + O\left( A_t + \rho \sqrt{\frac{t \log n}{n}} \|\xi_{t-1}\|_2 \right).$$

It suffices to control

- $\kappa_t < 1 - c$
- $A_t$  corresponds to an upper bound for quantity

$$\left| \sum_{k=1}^{t-1} \mu^k \underbrace{\left[ \langle \phi_k, \eta_t(v_t) \rangle - \langle \eta'_t(v_t) \rangle \beta_{t-1}^k \right]}_{Y_k} \right|, \quad \text{with } v_t := \alpha_t v^* + \sum_{k=1}^{t-1} \beta_{t-1}^k \phi_k$$

*Application in a concrete example:  $\mathbb{Z}_2$  synchronization*

— *for other examples refer to [Li and Wei, 2022](#)*



## Prior art: A hybrid procedure

---

- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\sqrt{n}v_i^* \sim \text{Unif}(\{\pm 1\})$
- Goal: recover  $v^*$  given  $M$

— AMP is approximately Gaussian in a fixed  $t$ , large  $n$  limit

## Connections of $Z_2$ and SBM

---

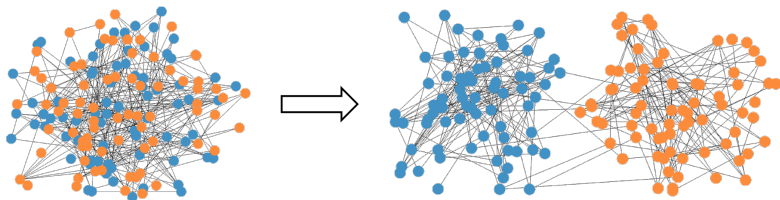
A symmetric two-group model:

- vertex set  $V = [n] = V_+ \cup V_-$  with  $\mathbb{P}(i \in V_+) = \mathbb{P}(i \in V_-) = 1/2$
- stochastic block model:  $(X, G) \sim \text{SBM}(n, p, q)$
- goal: characterize minimum mean square error/mutual information

# Connections of $Z_2$ and SBM

A symmetric two-group model:

- vertex set  $V = [n] = V_+ \cup V_-$  with  $\mathbb{P}(i \in V_+) = \mathbb{P}(i \in V_-) = 1/2$
- stochastic block model:  $(X, G) \sim \text{SBM}(n, p, q)$
- goal: characterize minimum mean square error/mutual information



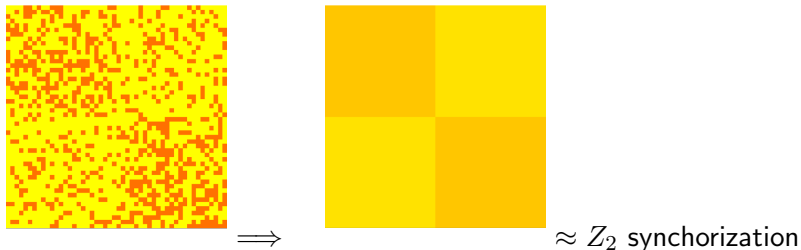
*"Asymptotic Mutual Information for the Two-Groups Stochastic Block Model,"* Deshpande, Abbe, Montanari 2017

*"Community Detection and Stochastic Block Models: Recent Developments,"* Abbe 2018

# Connections of $Z_2$ and SBM

A symmetric two-group model:

- vertex set  $V = [n] = V_+ \cup V_-$  with  $\mathbb{P}(i \in V_+) = \mathbb{P}(i \in V_-) = 1/2$
- stochastic block model:  $(X, G) \sim \text{SBM}(n, p, q)$
- goal: characterize minimum mean square error/mutual information



*"Asymptotic Mutual Information for the Two-Groups Stochastic Block Model,"* Deshpande, Abbe, Montanari 2017

*"Community Detection and Stochastic Block Models: Recent Developments,"* Abbe 2018

## Prior art: A hybrid procedure

---

- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\sqrt{n}v_i^* \sim \text{Unif}(\{\pm 1\})$
- Goal: recover  $v^*$  given  $M$

— AMP is approximately Gaussian in a fixed  $t$ , large  $n$  limit

## Prior art: A hybrid procedure

---

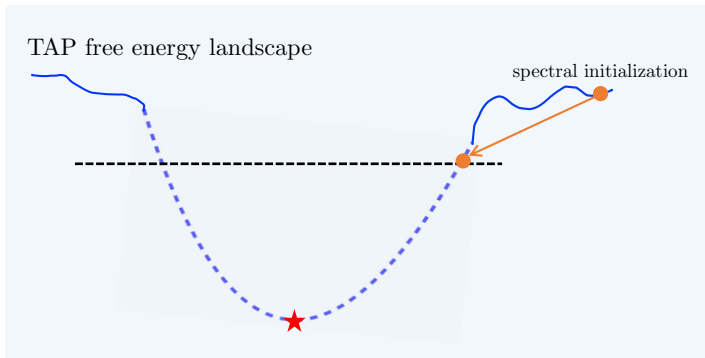
- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\sqrt{n}v_i^* \sim \text{Unif}(\{\pm 1\})$
- Goal: recover  $v^*$  given  $M$

A hybrid procedure proposed in [Celentano, Fan, Mei'21](#)

## Prior art: A hybrid procedure

- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\sqrt{n}v_i^* \sim \text{Unif}(\{\pm 1\})$
- Goal: recover  $v^*$  given  $M$

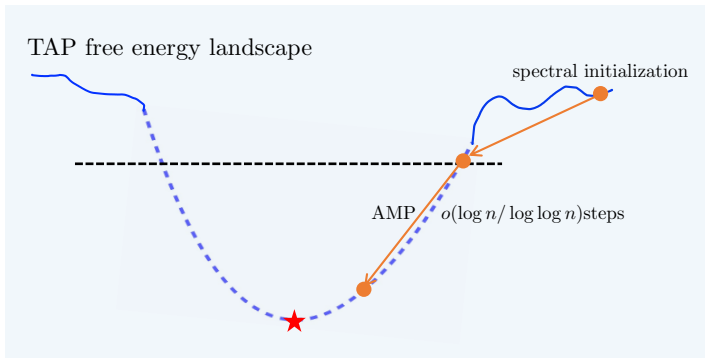
A hybrid procedure proposed in [Celentano, Fan, Mei'21](#)



# Prior art: A hybrid procedure

- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\sqrt{n}v_i^* \sim \text{Unif}(\{\pm 1\})$
- Goal: recover  $v^*$  given  $M$

A hybrid procedure proposed in [Celentano, Fan, Mei'21](#)

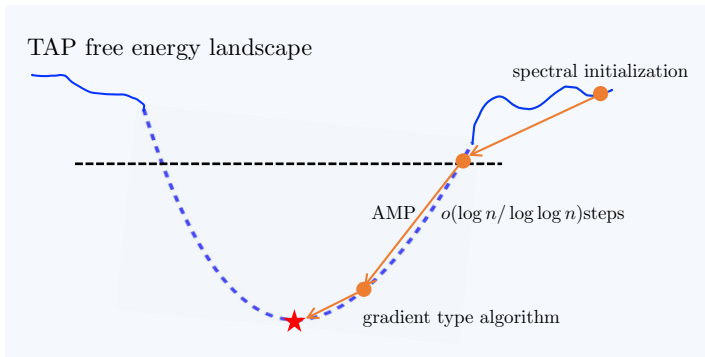




# Prior art: A hybrid procedure

- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\sqrt{n}v_i^* \sim \text{Unif}(\{\pm 1\})$
- Goal: recover  $v^*$  given  $M$

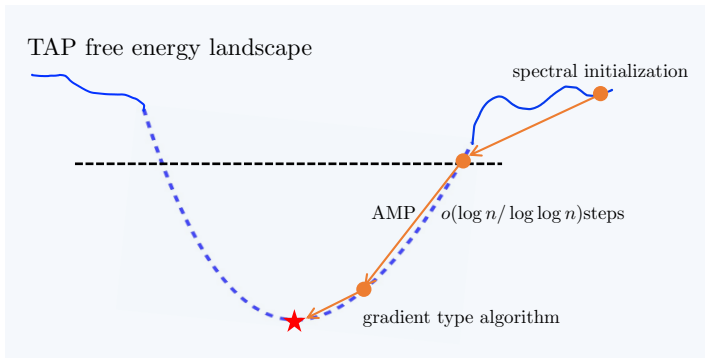
A hybrid procedure proposed in [Celentano, Fan, Mei'21](#)



# Prior art: A hybrid procedure

- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\sqrt{n}v_i^* \sim \text{Unif}(\{\pm 1\})$
- Goal: recover  $v^*$  given  $M$

A hybrid procedure proposed in [Celentano, Fan, Mei'21](#)



**Open question:** spectrally-initialized AMP is sufficient for  $\lambda > 1$ ?

## $\mathbb{Z}_2$ Synchronization: our results

### Theorem (Li & Wei'22)

Spectrally-initialized AMP satisfies

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t,$$

with

$$\alpha_{t+1} = \mathbb{E} \left[ \lambda v^{*\top} \eta_t \left( \alpha_t v^* + \frac{1}{\sqrt{n}} G \right) \right] + O \left( \sqrt{\frac{t \log n}{(\lambda - 1)^3 n}} \right),$$

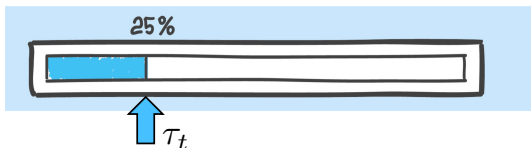
$$\|\beta_t\|_2 = 1, \quad \|\xi_t\|_2 \lesssim O \left( \sqrt{\frac{t \log n}{(\lambda - 1)^3 n}} + \sqrt{\frac{\log^7 n}{(\lambda - 1)^9 n}} \right)$$

w.h.p. provided that  $t \lesssim \frac{(\lambda - 1)^{10}}{\log^7 n} n$ .

- spectral initialization provides a warm-start with  $\alpha_1 \asymp \sqrt{\lambda^2 - 1}$

# Connection to state evolution

---

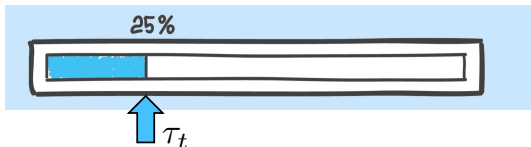


(asymptotic) state evolution [Deshpande, Abbe, Montanari \(2016\)](#):

$$\tau_{t+1} := \lambda^2 \int \tanh(\tau_t + \sqrt{\tau_t}x)\varphi(dx)$$

# Connection to state evolution

---



(asymptotic) state evolution [Deshpande, Abbe, Montanari \(2016\)](#):

$$\tau_{t+1} := \lambda^2 \int \tanh(\tau_t + \sqrt{\tau_t}x)\varphi(dx)$$

here

$$\alpha_t^2 - \tau_t = O\left(\sqrt{\frac{t \log n}{(\lambda - 1)^8 n}} + \sqrt{\frac{\log^7 n}{(\lambda - 1)^{14} n}}\right)$$

# Connection to state evolution

---



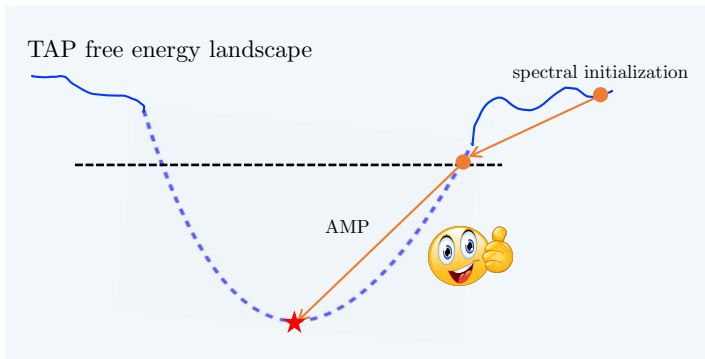
(asymptotic) state evolution [Deshpande, Abbe, Montanari \(2016\)](#):

$$\tau_{t+1} := \lambda^2 \int \tanh(\tau_t + \sqrt{\tau_t}x)\varphi(dx)$$

$$\alpha_t^2 - \tau^* = c(1 - (\lambda - 1))^t + O\left(\sqrt{\frac{t \log n}{(\lambda - 1)^8 n}} + \sqrt{\frac{\log^7 n}{(\lambda - 1)^{14} n}}\right)$$

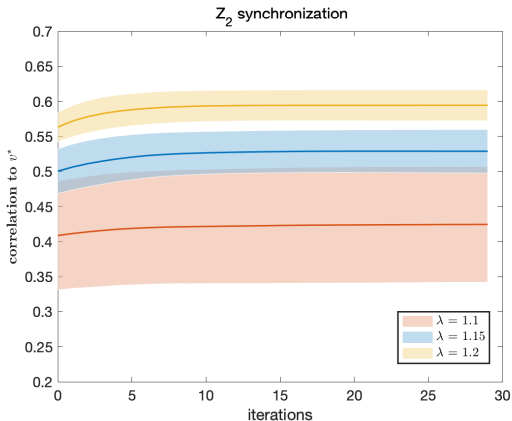
# Take-home message #1

---



- Answer the open question (Celentano, Fan & Mei (2021)) positively: spectrally-initialized AMP is enough!

## $\mathbb{Z}_2$ Synchronization: simulations

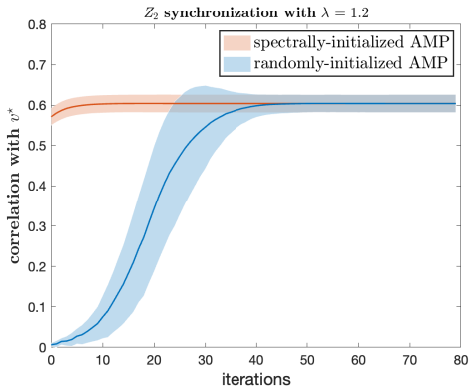


**Figure:** Convergence of spectrally-initialized AMP for different signal strengths with  $n = 10000$ . Repeat 40 times.



**Question:** *Is spectral initialization really necessary for AMP?*

# Simulation: AMP with random initialization



**Figure:** The correlation of  $\eta_t(x_t)$  and  $v^*$  vs. iteration count  $t$  for AMP with both random and spectral initialization. Here  $n = 10000$ . Repeat 20 times.

# AMP with random initialization

## Theorem (Li, Fan, Wei'23)

For  $t \leq \frac{cn(\lambda-1)^5}{\log^2 n}$ , *randomly-initialized* AMP satisfies, w.h.p,

- **(Decomposition)**  $x_{t+1} = \alpha_{t+1}v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$ , with

$$\alpha_{t+1} := \lambda v^{*\top} \eta_t(x_t),$$

$$\|\beta_t\|_2 = 1, \quad \|\xi_t\|_2 \lesssim \sqrt{\frac{t \log n}{n(\lambda-1)^2}} + \sqrt{\frac{\log^4 n}{n(\lambda-1)^3}};$$

$$- \tau_{t+1} := \lambda^2 \int \tanh(\tau_t + \sqrt{\tau_t}x) \varphi(dx)$$

# AMP with random initialization

## Theorem (Li, Fan, Wei'23)

For  $t \leq \frac{cn(\lambda-1)^5}{\log^2 n}$ , *randomly-initialized* AMP satisfies, w.h.p,

- **(Decomposition)**  $x_{t+1} = \alpha_{t+1}v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$ , with

$$\alpha_{t+1} := \lambda v^{*\top} \eta_t(x_t),$$

$$\|\beta_t\|_2 = 1, \quad \|\xi_t\|_2 \lesssim \sqrt{\frac{t \log n}{n(\lambda-1)^2}} + \sqrt{\frac{\log^4 n}{n(\lambda-1)^3}};$$

- **(Crossing time)**

$$\varsigma := \min\{t : |\alpha_t| \geq \frac{1}{2} \sqrt{\lambda^2 - 1}\} = O\left(\frac{\log n}{\lambda - 1}\right);$$

$$- \tau_{t+1} := \lambda^2 \int \tanh(\tau_t + \sqrt{\tau_t}x) \varphi(dx)$$

# AMP with random initialization

## Theorem (Li, Fan, Wei'23)

For  $t \leq \frac{cn(\lambda-1)^5}{\log^2 n}$ , *randomly-initialized* AMP satisfies, w.h.p,

- **(Decomposition)**  $x_{t+1} = \alpha_{t+1}v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$ , with

$$\alpha_{t+1} := \lambda v^{*\top} \eta_t(x_t),$$

$$\|\beta_t\|_2 = 1, \quad \|\xi_t\|_2 \lesssim \sqrt{\frac{t \log n}{n(\lambda-1)^2}} + \sqrt{\frac{\log^4 n}{n(\lambda-1)^3}};$$

- **(Crossing time)**

$$\varsigma := \min\{t : |\alpha_t| \geq \frac{1}{2} \sqrt{\lambda^2 - 1}\} = O\left(\frac{\log n}{\lambda - 1}\right);$$

- **(Non-asymptotic SE)** for any  $t \geq \varsigma$ ,

$$\alpha_t^2 = \left(1 + O\left(\sqrt{\frac{(t + \frac{\log^3 n}{\lambda-1}) \log n}{n(\lambda-1)^5}}\right)\right) \tau_{t+1}.$$

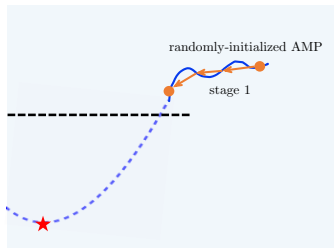
$$- \tau_{t+1} := \lambda^2 \int \tanh(\tau_t + \sqrt{\tau_t} x) \varphi(dx)$$

# Dynamics after random initialization

randomly-initialized AMP

- escape from random initialization

$$\alpha_{t+1} \approx \lambda \alpha_t + \lambda g_{t-1}$$



$$\alpha_t \approx n^{-1/4}$$

$O\left(\frac{\log n}{\lambda - 1}\right)$  #steps

# Dynamics after random initialization

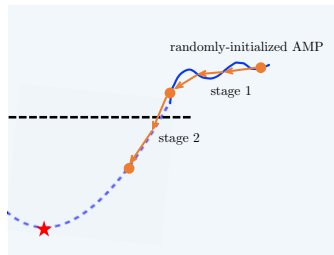
## randomly-initialized AMP

- escape from random initialization

$$\alpha_{t+1} \approx \lambda \alpha_t + \lambda g_{t-1}$$

- exponential growth

$$\alpha_{t+1} \geq (1 + c(\lambda - 1))^{1/2} \alpha_t$$



$$\alpha_t \approx n^{-1/4}$$

$$\alpha_t \approx \sqrt{\lambda^2 - 1}$$

$$O\left(\frac{\log n}{\lambda - 1}\right) \text{ \#steps}$$

$$O\left(\frac{\log n}{\lambda - 1}\right) \text{ \#steps}$$

# Dynamics after random initialization

randomly-initialized AMP

- escape from random initialization

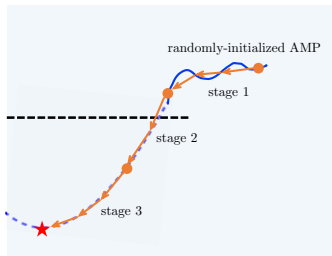
$$\alpha_{t+1} \approx \lambda \alpha_t + \lambda g_{t-1}$$

- exponential growth

$$\alpha_{t+1} \geq (1 + c(\lambda - 1))^{1/2} \alpha_t$$

- local refinement

$$|\alpha_t^2 - \tau^*| \lesssim (1 - (\lambda - 1))^{t-c} + \sqrt{\frac{t/n}{(\lambda - 1)^6}}$$



$$\alpha_t \approx n^{-1/4}$$

$$\alpha_t \approx \sqrt{\lambda^2 - 1}$$

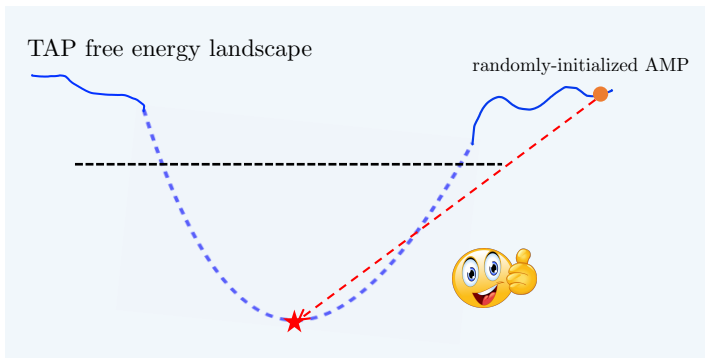
$$O\left(\frac{\log n}{\lambda - 1}\right) \text{ #steps}$$

$$O\left(\frac{\log n}{\lambda - 1}\right) \text{ #steps}$$

$$t \leq \frac{n(\lambda - 1)^5}{\log^2 n}$$

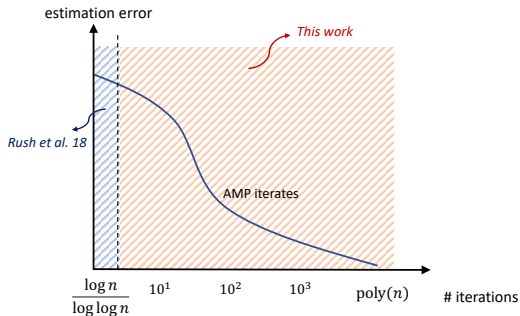


## Take-home message #2



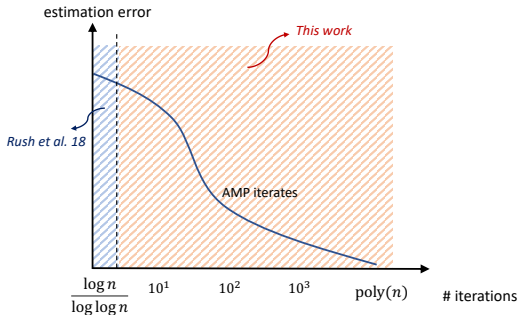
- It takes *randomly-initialized* AMP at most  $O\left(\frac{\log n}{\lambda-1}\right)$  iterations to get  $\tilde{O}\left(\sqrt{\frac{1}{n(\lambda-1)^6}}\right)$  close to the Bayes-optimal risk.

# Concluding remarks



- a new non-asymptotic framework of AMP that allows for # iterations  $O\left(\frac{n}{\text{poly}(\log n)}\right)$  given informative/spectral initialization

# Concluding remarks



- a new non-asymptotic framework of AMP that allows for # iterations  $O\left(\frac{n}{\text{poly}(\log n)}\right)$  given informative/spectral initialization
- analyze performance of randomly-initialized AMP for  $\mathbb{Z}_2$  synchronization

## Concluding remarks: future extensions

---

- other statistical settings
- apply these results for statistical inference
- connections to other polynomial-time algorithms
- universality results
- infinite number of iterations
- etc...

Thanks for your attention! Questions?

**Paper:**

“A non-asymptotic framework for approximate message passing in spiked models,” G. Li, Y. Wei, *arxiv.2208.03313*

“Approximate message passing from random initialization with applications to  $\mathbb{Z}_2$  synchronization,” G. Li, W. Fan, Y. Wei, *PNAS*, 2023

“Non-asymptotic analyses for approximate message passing with applications to sparse and robust regression,” G. Li, Y. Wei, *upcoming*

# Extensions to sparse/robust regression

---

- So far, AMP for spiked models  $M = v^*v^{*\top} + W$ :

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ for } t \geq 1$$

where  $\langle x \rangle := \frac{1}{n} \sum_{i=1}^n x_i$

# Extensions to sparse/robust regression

- So far, AMP for spiked models  $M = v^*v^{*\top} + W$ :

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ for } t \geq 1$$

$$\text{where } \langle x \rangle := \frac{1}{n} \sum_{i=1}^n x_i$$

- AMP for sparse/robust regression  $y = X\theta + \varepsilon$ :

$$s_t = XF_t(\beta_t) - \langle F'_t \rangle G_{t-1}(s_{t-1}),$$

$$\beta_{t+1} = X^\top G_t(s_t) - \langle G'_t \rangle F_t(\beta_t)$$

## Theorem (informal)

For  $t \lesssim n/\log^4 n$ , AMP iterates satisfy

$$\beta_{t+1} = \theta_{t+1} - \theta^* = \sum_{k=1}^t \alpha_t^k \psi_k + \zeta_t$$

where with probability at least  $1 - O(n^{-10})$ ,  $\|\zeta\|_2 \leq \left(\frac{t \log n}{n}\right)^{1/3}$ .

**A glimpse of our main proof idea...**

— *decomposition*:  $x_{t+1} = \alpha_{t+1}v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$



# Prior non-asymptotic guarantees

---

AMP for spiked models:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ for } t \geq 1$$

- **Challenges:** deal with statistical dependence between iterations

# Prior non-asymptotic guarantees

AMP for spiked models:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ for } t \geq 1$$

- **Challenges:** deal with statistical dependence between iterations
- [Rush & Venkataraman'16](#) #iterations =  $o(\log n / \log \log n)$   
— based on state-evolution analysis in [Bayati & Montanari'11](#)

$$\begin{aligned} & \mathbb{P}(\text{residual at time } t \geq \epsilon) \\ &= \mathbb{P}\left(\sum_{i=0}^{t-1} r_i^t \geq \epsilon\right) \leq \sum_{i=0}^{t-1} \mathbb{P}\left(r_i^t \leq \frac{\epsilon}{t}\right) \leq t C_{t-1} \exp\left(-\frac{c_{t-1}}{t^2} n \epsilon^2\right) \end{aligned}$$

statistical dependence      induction step

requires  $\frac{n}{(t!)^2} \rightarrow \infty \rightarrow t = o(\log n / \log \log n)$

## Main proof idea: a new decomposition

---

- $\{x_t\}$  is the sequence generated by AMP

## Main proof idea: a new decomposition

---

- $\{x_t\}$  is the sequence generated by AMP
- define an orthonormal basis  $\{z_t\}$  where

$$z_1 := \frac{\eta_1(x_1)}{\|\eta_1(x_1)\|_2} \in \mathbb{R}^n, \quad \text{and} \quad W_1 := W$$

## Main proof idea: a new decomposition

---

- $\{x_t\}$  is the sequence generated by AMP
- define an orthonormal basis  $\{z_t\}$  where

$$z_1 := \frac{\eta_1(x_1)}{\|\eta_1(x_1)\|_2} \in \mathbb{R}^n, \quad \text{and} \quad W_1 := W$$

- write  $U_{t-1} := [z_k]_{1 \leq k \leq t-1} \in \mathbb{R}^{n \times (t-1)}$  and denote

$$z_t := \frac{(I - U_{t-1}U_{t-1}^\top) \eta_t(x_t)}{\|(I - U_{t-1}U_{t-1}^\top) \eta_t(x_t)\|_2} \quad \text{Gram-Schmidt orthogonalization,}$$

$$W_t := (I - z_{t-1}z_{t-1}^\top) W_{t-1} (I - z_{t-1}z_{t-1}^\top)$$

# Main proof idea: a new decomposition

---

- $\{x_t\}$  is the sequence generated by AMP
- define an orthonormal basis  $\{z_t\}$  where

$$z_1 := \frac{\eta_1(x_1)}{\|\eta_1(x_1)\|_2} \in \mathbb{R}^n, \quad \text{and} \quad W_1 := W$$

- write  $U_{t-1} := [z_k]_{1 \leq k \leq t-1} \in \mathbb{R}^{n \times (t-1)}$  and denote

$$z_t := \frac{(I - U_{t-1}U_{t-1}^\top) \eta_t(x_t)}{\|(I - U_{t-1}U_{t-1}^\top) \eta_t(x_t)\|_2} \quad \text{Gram-Schmidt orthogonalization,}$$

$$W_t := (I - z_{t-1}z_{t-1}^\top) W_{t-1} (I - z_{t-1}z_{t-1}^\top)$$

- write  $\eta_t(x_t) = \sum_{k=1}^t \beta_t^k z_k$ , for  $\beta_t^k := \langle \eta_t(x_t), z_k \rangle$

# Main proof idea: a new decomposition

---

- AMP updates:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ where } M = \lambda v^* v^{*\top} + W$$

- Goal:  $x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$

# Main proof idea: a new decomposition

---

- AMP updates:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ where } M = \lambda v^* v^{*\top} + W$$

- Goal:  $x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$

$$M\eta_t(x_t)$$



# Main proof idea: a new decomposition

- AMP updates:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ where } M = \lambda v^* v^{*\top} + W$$

- Goal:  $x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$

$$M\eta_t(x_t)$$

$$= v^* \underbrace{\lambda v^{*\top} \eta_t(x_t)}_{\alpha_{t+1}} + \left\{ W_t + \underbrace{\sum_{k=1}^{t-1} \left[ W_k z_k z_k^\top + z_k z_k^\top W_k - z_k z_k^\top W_k z_k z_k^\top \right]}_{W_k - W_{k+1}} \right\} \cdot \underbrace{\sum_{k=1}^t \beta_t^k z_k}_{\eta_t(x_t)}$$

# Main proof idea: a new decomposition

- AMP updates:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ where } M = \lambda v^* v^{*\top} + W$$

- Goal:  $x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$

$$M\eta_t(x_t)$$

$$\begin{aligned} &= v^* \underbrace{\lambda v^{*\top} \eta_t(x_t)}_{\alpha_{t+1}} + \left\{ W_t + \underbrace{\sum_{k=1}^{t-1} \left[ W_k z_k z_k^\top + z_k z_k^\top W_k - z_k z_k^\top W_k z_k z_k^\top \right]}_{W_k - W_{k+1}} \right\} \cdot \underbrace{\sum_{k=1}^t \beta_t^k z_k}_{\eta_t(x_t)} \\ &= \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k W_k z_k + \dots \end{aligned}$$

# Main proof idea: a new decomposition

- AMP updates:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ where } M = \lambda v^* v^{*\top} + W$$

- Goal:  $x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$

$$M\eta_t(x_t)$$

$$\begin{aligned} &= v^* \underbrace{\lambda v^{*\top} \eta_t(x_t)}_{\alpha_{t+1}} + \left\{ W_t + \underbrace{\sum_{k=1}^{t-1} \left[ W_k z_k z_k^\top + z_k z_k^\top W_k - z_k z_k^\top W_k z_k z_k^\top \right]}_{W_k - W_{k+1}} \right\} \cdot \underbrace{\sum_{k=1}^t \beta_t^k z_k}_{\eta_t(x_t)} \\ &= \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k W_k z_k + \dots \\ &= \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \underbrace{(W_k z_k + \zeta_k)}_{\phi_k \sim \mathcal{N}(0, \frac{1}{n} \mathbf{I}_n)} + \dots \end{aligned}$$

# Main proof idea: a new decomposition

- AMP updates:

$$x_{t+1} = M\eta_t(x_t) - \langle \eta'_t(x_t) \rangle \cdot \eta_{t-1}(x_{t-1}), \text{ where } M = \lambda v^* v^{*\top} + W$$

- Goal:  $x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t$

$$M\eta_t(x_t)$$

$$= v^* \underbrace{\lambda v^{*\top} \eta_t(x_t)}_{\alpha_{t+1}} + \left\{ W_t + \underbrace{\sum_{k=1}^{t-1} \left[ W_k z_k z_k^\top + z_k z_k^\top W_k - z_k z_k^\top W_k z_k z_k^\top \right]}_{W_k - W_{k+1}} \right\} \cdot \underbrace{\sum_{k=1}^t \beta_t^k z_k}_{\eta_t(x_t)}$$

$$= \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k W_k z_k + \dots$$

$$= \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \underbrace{(W_k z_k + \zeta_k)}_{\phi_k \sim \mathcal{N}(0, \frac{1}{n} \mathbf{I}_n)} + \dots$$

$$\xi_t = \sum_{k=1}^{t-1} z_k \left[ \langle W_k z_k, \eta_t(x_t) \rangle - \langle \eta'_t(x_t) \rangle \beta_{t-1}^k - \beta_t^k z_k^\top W_k z_k \right] - \sum_{k=1}^t \beta_t^k \zeta_k$$

# sparse PCA in spiked models

- Setting:  $M = \lambda v^* v^{*\top} + W$  where  $\|v^*\|_0 = k$
- Goal: recover  $v^*$  given  $M$

$$\lambda \approx \sqrt{\frac{k \log n}{n}}$$

statistical limit

$$\lambda \approx \sqrt{\frac{k^2}{n}}$$

computation limit

reduction to planted cliques:  
Berthet & Rigollet (2013)

SNR



"I can't find an efficient algorithm, but neither can all these people."

Zou et al. (2006)  
Amini and Wainwright (2008)  
Ma (2013)  
Deshpande and Montanari (2014b)  
Hopkins et al. (2017)

# Sparse PCA: our results

## Theorem (Li & Wei'22)

Suppose  $0 < \lambda \lesssim 1$ . Given an informative initialization (with non-vanishing correlation with  $v^*$ ), AMP satisfies

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t,$$

with

$$\alpha_{t+1} = \mathbb{E} \left[ \lambda v^{*\top} \eta_t \left( \alpha_t v^* + \frac{1}{\sqrt{n}} G \right) \right] + \sqrt{\frac{k + t \log^3 n}{n}},$$

$$\|\beta_t\|_2 = 1, \quad \|\xi_t\|_2 \lesssim \sqrt{\frac{k + t \log^3 n}{n}} \quad \text{w.h.p.}$$

provided that  $\frac{t \log^3 n}{n \lambda^2} \lesssim 1$  and  $\frac{k \log n}{n \lambda^2} \lesssim 1$ .

# Sparse PCA: our results

## Theorem (Li & Wei'22)

Suppose  $0 < \lambda \lesssim 1$ . Given an informative initialization (with non-vanishing correlation with  $v^*$ ), AMP satisfies

$$x_{t+1} = \alpha_{t+1} v^* + \sum_{k=1}^t \beta_t^k \phi_k + \xi_t,$$

with

$$\alpha_{t+1} = \mathbb{E} \left[ \lambda v^{*\top} \eta_t \left( \alpha_t v^* + \frac{1}{\sqrt{n}} G \right) \right] + \sqrt{\frac{k + t \log^3 n}{n}},$$

$$\|\beta_t\|_2 = 1, \quad \|\xi_t\|_2 \lesssim \sqrt{\frac{k + t \log^3 n}{n}} \quad \text{w.h.p.}$$

provided that  $\frac{t \log^3 n}{n \lambda^2} \lesssim 1$  and  $\frac{k \log n}{n \lambda^2} \lesssim 1$ .

denoising functions:

$$\eta_t(x) = \gamma_t \text{sign}(x) (|x| - \tau_t)_+ \quad \text{where } \gamma_t^{-1} := \|(|x_t| - \tau_t)_+\|_2, \tau_t \asymp \sqrt{\frac{\log n}{n}}$$

## Several remarks

---

- recall the (asymptotic) state evolution:

$$\alpha_{t+1}^* := \frac{\lambda v^{*\top} \int \text{ST}_{\tau_t} \left( \alpha_t^* v^* + \frac{x}{\sqrt{n}} \right) \varphi_n(dx)}{\sqrt{\int \left\| \text{ST}_{\tau_t} \left( \alpha_t^* v^* + \frac{x}{\sqrt{n}} \right) \right\|_2^2 \varphi_n(dx)}}$$

then

$$|\alpha_{t+1} - \alpha_{t+1}^*| \lesssim \sqrt{\frac{k \log n + t \log^3 n}{n}}$$



## Several remarks

---

- recall the (asymptotic) state evolution:

$$\alpha_{t+1}^* := \frac{\lambda v^{*\top} \int \text{ST}_{\tau_t} \left( \alpha_t^* v^* + \frac{x}{\sqrt{n}} \right) \varphi_n(dx)}{\sqrt{\int \left\| \text{ST}_{\tau_t} \left( \alpha_t^* v^* + \frac{x}{\sqrt{n}} \right) \right\|_2^2 \varphi_n(dx)}}$$

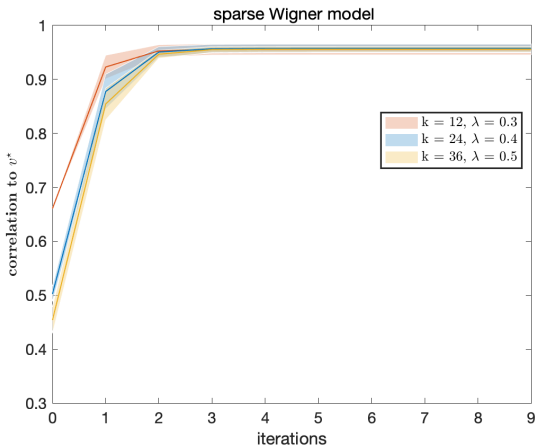
then

$$|\alpha_{t+1} - \alpha_{t+1}^*| \lesssim \sqrt{\frac{k \log n + t \log^3 n}{n}}$$

- two sufficient initialization schemes:

- ▶ AMP with **diagonal maximization**:  $\lambda \|v^*\|_\infty \gtrsim \sqrt{\frac{k \log n}{n}}$
- ▶ AMP with **sample-split initialization**:  $\lambda \gtrsim \sqrt{\frac{k^2}{n}}$  and  $\|v^*\|_\infty \lesssim \frac{\log n}{k}$

# Sparse PCA: simulations



**Figure:** Convergence of AMP with diagonal maximization for different signal strengths with  $n = 10000$ . Repeat 40 times.

## Auxiliary details

---

Define  $\zeta_k := \left(\frac{\sqrt{2}}{2} - 1\right) z_k z_k^\top W_k z_k + \sum_{i=1}^{k-1} g_i^k z_i$

$$W_k z_k + \zeta_k = \phi_k \stackrel{\text{i.i.d}}{\sim} \mathcal{N}\left(0, \frac{1}{n} I_n\right)$$

- conditioning on  $x_1, \{z_i\}_{i < k}$ ,  $W_k$  is a Wigner matrix in subspace  $U_{k-1}^\perp$
- $W_k z_k$  has zero variance along the directions of  $\{z_i\}_{i < k}$  and  $\frac{2}{n}$  variance along the direction of  $z_k$

# Conditioning technique

$$\begin{array}{ll} \text{AMP updates} & x_{t+1} = Wm_t - \gamma_t m_{t-1} \\ \text{where} & m^t = \eta_t(x_t), \quad \gamma_t = \langle \eta'_t(x_t) \rangle \end{array}$$

- $m_{-1} = 0, x_0 = 0$  and  $x_1 = W\eta_t(0)$
- $\sigma$ -algebra  $\mathcal{F}_t$  generated by  $\{x_0, x_1, \dots, x_t\}$ , conditioning on  $\mathcal{F}$  is equivalent to conditioning on event

$$\mathcal{E}_t := \left\{ x_1 + \gamma_0 m_{-1} = Wm_0, x_2 + \gamma_1 m_1 = Wm_1, \dots, x_t + \gamma_{t-1} m_{t-1} = Wm_{t-1} \right\}$$

- $W$  conditioning on linear observations

$$\begin{aligned} W|_{\mathcal{F}_t} &\stackrel{d}{=} \mathbb{E}[W|_{\mathcal{F}_t}] + P_t^\perp W^{\text{new}} P_t^\perp \\ W|_{\mathcal{F}_t} m^t &\stackrel{d}{=} \underbrace{W^{\text{new}} P_t^\perp m^t}_{\text{Gaussian term}} + \underbrace{W^{\text{new}} (I - P_t^\perp) m^t + \mathbb{E}[W|_{\mathcal{F}_t}] m^t}_{\text{non-Gaussian term}} \end{aligned}$$

Bolthausen (2006), Bayati & Montanari (2011), Rush & Venkataramanan (2016), Berthier, Montanari & Nguyen (2020)