

# Optimality of Approximate Message Passing Algorithms for Spiked Matrix Models with Rotationally Invariant Noise

Rishabh Dudeja\*      Songbin Liu†      Junjie Ma‡§

May 28, 2024

## Abstract

We study the problem of estimating a rank one signal matrix from an observed matrix generated by corrupting the signal with additive rotationally invariant noise. We develop a new class of approximate message-passing algorithms for this problem and provide a simple and concise characterization of their dynamics in the high-dimensional limit. At each iteration, these algorithms exploit prior knowledge about the noise structure by applying a non-linear matrix denoiser to the eigenvalues of the observed matrix and prior information regarding the signal structure by applying a non-linear iterate denoiser to the previous iterates generated by the algorithm. We exploit our result on the dynamics of these algorithms to derive the optimal choices for the matrix and iterate denoisers. We show that the resulting algorithm achieves the smallest possible asymptotic estimation error among a broad class of iterative algorithms under a fixed iteration budget.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Preliminaries</b>	<b>5</b>
<b>3</b>	<b>Main results</b>	<b>8</b>
3.1	Orthogonal Approximate Message Passing Algorithms . . . . .	8
3.2	The Optimal OAMP Algorithm . . . . .	9
3.3	An Optimality Result . . . . .	11
3.4	Information-Theoretic v.s. Computational Limits . . . . .	11
<b>4</b>	<b>Proof Ideas</b>	<b>13</b>
4.1	Heuristic Derivation of State Evolution (Theorem 1) . . . . .	13
4.2	Derivation of the Optimal OAMP Algorithm . . . . .	15
4.3	Proof of the Optimality Result (Theorem 2) . . . . .	17
<b>5</b>	<b>Numerical Experiments</b>	<b>20</b>
<b>A</b>	<b>Proofs for Preliminary Results</b>	<b>27</b>
A.1	Proofs for Random Matrix Theory Results . . . . .	27
A.2	Preliminaries on Gaussian Channels . . . . .	29
A.2.1	Proof of Lemma 2 . . . . .	31
A.2.2	Proof of Lemma 3 . . . . .	33
A.2.3	Proof of Lemma 4 . . . . .	35

---

\*Department of Statistics, University of Wisconsin–Madison. Email: [rdudeja@wisc.edu](mailto:rdudeja@wisc.edu)

†Academy of Mathematics and Systems Science, Chinese Academy of Sciences. Email: [liusongbin@lsec.cc.ac.cn](mailto:liusongbin@lsec.cc.ac.cn)

‡Academy of Mathematics and Systems Science, Chinese Academy of Sciences. Email: [majunjie@lsec.cc.ac.cn](mailto:majunjie@lsec.cc.ac.cn)

§The authors are alphabetically ordered.

<b>B</b>	<b>State Evolution for OAMP Algorithms (Theorem 1)</b>	<b>36</b>
B.1	Proof of Lemma 5	38
B.2	Proof of Lemma 6	41
B.3	Proof of Lemma 7	42
B.4	Proof of Lemma 8	43
<b>C</b>	<b>State Evolution for Optimal OAMP (Proposition 1)</b>	<b>44</b>
C.1	Proof of Lemma 9	47
<b>D</b>	<b>Replica Predictions for Quartic Potential (Proposition 2)</b>	<b>48</b>
D.1	Preliminaries of Trace Ensemble	48
D.2	Proof of Proposition 2	51
D.2.1	Auxiliary results	52
D.2.2	Proof of Proposition 2	55
<b>E</b>	<b>Omitted Proofs for the Optimality Result (Theorem 2)</b>	<b>57</b>
E.1	Proof of Proposition 3	57
E.1.1	Divergence Removal (Proof of Lemma 14)	59
E.1.2	Polynomial Approximation (Proof of Lemma 15)	61
E.2	Proof Proposition 4	64
E.2.1	Proof of Lemma 16	67
E.2.2	Proof of Lemma 17	70
E.2.3	Proof of Lemma 18	71
E.3	Proof of Proposition 5	72
E.3.1	Proof of Lemma 19	77
E.3.2	Proof of Lemma 20	80
<b>F</b>	<b>Some Miscellaneous Results</b>	<b>83</b>

## 1 Introduction

We consider the problem of estimating a symmetric rank-one matrix from a noisy  $N \times N$  observed matrix  $\mathbf{Y}$  generated from the *spiked matrix model* [35]:

$$\mathbf{Y} = \frac{\theta}{N} \mathbf{x}_* \mathbf{x}_*^\top + \mathbf{W}, \quad (1)$$

where  $\mathbf{x}_* \in \mathbb{R}^N$  is the  $N$ -dimensional unknown signal of interest,  $\theta \geq 0$  is the signal-to-noise ratio (SNR) parameter, and  $\mathbf{W}$  is a symmetric noise matrix. This model and its variants have been used to study a broad range of statistical inference problems, including sparse PCA [20, 36, 78], community detection [2, 21], and group synchronization [17, 30, 64, 73].

**Wigner Noise Model.** The most well-studied variant of the spiked model is the *Spiked Wigner Model* [6, 20, 27, 28, 54, 65], which assumes that the noise matrix has i.i.d. Gaussian entries. A rich line of work in high-dimensional statistics and random matrix theory has studied the problem from various perspectives.

- *Design and Analysis of Estimators.* A natural estimator for the signal is the PCA estimator or the leading eigenvector of the observed matrix. A line of work [4, 5, 12, 32, 39, 62, 63] initiated by Baik et al. [5] have obtained a sharp asymptotic characterization of the performance of this estimator in the high-dimensional limit and uncovered surprising properties of this estimator. One way to improve the performance of PCA is to exploit prior structural information about the signal, such as sparsity. Approximate Message-passing algorithms are a popular class of computationally efficient iterative algorithms designed to exploit such structural information. These algorithms were first discovered in the context of compressed sensing [9, 13, 23, 37] and have been adapted for the spiked Wigner model [38, 46, 47, 54, 60, 68]. A particularly attractive feature of these algorithms is that their performance in the high-dimensional limit is characterized by a simple deterministic recursion known as *state evolution*, which makes it possible to assess their information-theoretic optimality (or sub-optimality).

- *Information-theoretic Limits.* Under the assumption that the signal is drawn from a known prior distribution, the optimal estimator is the Bayes estimator  $\mathbb{E}[\mathbf{x}_*|\mathbf{Y}]$ . Using the powerful but non-rigorous cavity method from statistical physics, Lesieur et al. [44] derived a conjecture for the Bayes risk (or the asymptotic performance of the Bayes estimator) for the Spiked Wigner model. This conjecture has been proved rigorously in increasing generality in a series of influential works [6, 7, 10, 20, 27, 40, 43, 52]. This formula for asymptotic Bayes risk provides a fundamental information-theoretic limit on the performance of any estimator for the problem.
- *Computational Limits.* In general, computing the Bayes-optimal estimator is computationally intractable in high dimensions. However, a suitably designed AMP algorithm called the Bayes-optimal AMP algorithm can attain the Bayes risk for the problem for sufficiently high signal-to-noise ratios. For lower signal-to-noise ratios, the Bayes-optimal AMP algorithm fails to do so, and no computationally efficient algorithm is known to achieve the Bayes risk [7, 40, 45, 52]. This has led to a popular conjecture that the Bayes-optimal AMP algorithm is the optimal polynomial time algorithm for this problem [45, 52]. Celentano et al. [16], Montanari and Wu [56], and Montanari and Wein [55] (see also the earlier work of Schramm and Wein [71]) have provided evidence for this conjecture by showing the Bayes-optimal AMP algorithm achieves the lowest possible estimation error among a broad class of iterative algorithms and low-degree polynomial estimators.

Collectively, these works have enriched our understanding of the fundamental trade-offs between statistical optimality and computational efficiency in high-dimensional statistics.

**Rotationally Invariant Noise Model.** In this paper, we study a natural generalization of the i.i.d. Gaussian noise model known as the rotationally invariant noise model where one posits that the eigenvectors of the noise matrix are uniformly random orthogonal matrices independent of the eigenvalues. The rotationally invariant noise model is intended to model noise matrices with strong statistical dependence whose eigenvectors are generic. Several works (*see e.g.*, [22, 48, 59]) have observed that this can be a good model in some applications. Benaych-Georges and Nadakuditi [12] analyzed the performance of PCA (or spectral estimators) for this noise model. A line of work [29, 53, 58, 76] initiated by Oppor et al. [58] and Fan [29] has developed AMP algorithms for this problem, which can improve the performance of PCA by exploiting signal structure. However, our understanding of the fundamental information-theoretic and computational limits in this model is extremely limited. Recent work by Barbier et al. [8] studies the spiked matrix model with rotationally invariant noise under the assumption that the noise matrix is drawn from the *trace ensemble*. Under this model, the density of  $\mathbf{W}$  is given by:

$$p(\mathbf{W}) \propto \exp\left(-\frac{N}{2} \sum_{i=1}^N V(\lambda_i(\mathbf{W}))\right) \quad \text{where } \lambda_1(\mathbf{W}), \dots, \lambda_N(\mathbf{W}) \text{ denote the eigenvalues of } \mathbf{W}, \quad (2)$$

and the *potential function*  $V : \mathbb{R} \mapsto \mathbb{R}$  is a functional parameter for the noise model. This assumption ensures that the problem has a well-defined likelihood, enabling the study of information-theoretic limits. Moreover, an appropriate choice of  $V$  can capture a wide range of noise eigenvalue spectrums in different applications. Barbier et al. [8] make progress towards understanding the information-theoretic limits of the problem and provide important insights regarding the structure of optimal computationally efficient algorithms.

- *Information-theoretic Limits.* Barbier et al. [8] develop a general (although non-rigorous) recipe based on the replica method to derive conjectured formulas for Bayes risk under the assumption that the potential  $V$  is a polynomial function, providing explicit conjectured formulas for the asymptotic Bayes risk when  $V$  is a quartic (degree four) or a sextic (degree six) polynomial. As the degree of  $V$  grows, the derivation and the resulting formulas become increasingly complex, and a general conjectured formula for the Bayes risk is still unavailable.
- *Optimal Computationally Efficient Algorithms.* Surprisingly, Barbier et al. [8] demonstrate that a natural generalization of the Bayes-optimal AMP algorithm (this is the optimal iterative algorithm for i.i.d. Gaussian noise and is conjectured to be the optimal polynomial-time algorithm) to the rotationally invariant noise model is sub-optimal. The authors develop AMP algorithms that achieve improved performance by applying a non-linear matrix denoiser to the eigenvalues of the observed

matrix  $\mathbf{Y}$  and characterize their state evolution. Based on non-rigorous statistical physics techniques [57], the authors provide a procedure to derive good matrix denoisers for the problem. The authors propose matrix denoisers with explicit formulas when  $V$  is a quartic or sextic polynomial. However, the state evolution of the resulting AMP algorithm is quite complicated; hence, the resulting algorithm’s optimality (or sub-optimality) properties are not understood.

**Our Contributions.** Taking inspiration from the insights of Barbier et al. [8], we study the spiked matrix model with rotationally invariant noise from an algorithmic point of view. Our main contributions are:

- We develop a new class of approximate message-passing algorithms for this problem and provide a state evolution result that characterizes their dynamics in the high-dimensional limit (Theorem 1 in Section 3.1). At each iteration, these algorithms exploit the noise structure by applying a non-linear matrix denoiser to the eigenvalues of the observed matrix and the signal structure by applying a non-linear iterate denoiser to the previous iterates generated by the algorithm. These algorithms can be viewed as natural analogs of orthogonal [50] or vector approximate message-passing algorithms [69] developed for compressed sensing.
- A key feature of the AMP algorithms proposed in this work is that their state evolution is significantly simpler than that of existing AMP algorithms for this problem. Consequently, we are able to exploit our result to derive the optimal choices for the matrix and iterate denoisers (Section 3.2). Interestingly, we find that the matrix denoisers that optimize the performance of the AMP algorithm are closely related to the eigenvalue shrinkage estimators discovered by Bun et al. [15] in a separate line of work on denoising high-rank and unstructured signal matrices corrupted with rotationally invariant noise [15, 41, 42, 49, 67, 72].
- Building on the techniques developed by Celentano et al. [16] and Montanari and Wu [56] (in the context of i.i.d. Gaussian noise), we show that the AMP algorithm with the optimal choices for the matrix and iterate denoisers achieves the smallest possible asymptotic estimation error among a broad class of iterative algorithms under a fixed iteration budget (Theorem 2 in Section 3.3). This suggests that this algorithm might be the natural candidate for the optimal polynomial-time estimator for this problem.
- Finally, our results also suggest a general and concise conjecture for the Bayes risk for the problem (Section 3.4), which we hope can be proved rigorously in the future.

**Organization.** This paper is organized as follows. We start with some preliminary results in Section 2. The main results of this paper are presented in Section 3. Section 4 highlights the key ideas behind our results. Numerical experiments are presented in Section 5. The complete proofs of our main results are provided in the appendices.

**Notation.** We conclude this section by introducing the notations used in this paper.

*Some common sets:* The sets  $\mathbb{N}, \mathbb{R}, \mathbb{C}$  represent the set of positive integers, real numbers, and complex numbers, respectively. For  $N \in \mathbb{N}$ ,  $[N]$  is the set  $\{1, 2, 3, \dots, N\}$  and  $\mathbb{O}(N)$  denotes the set of  $N \times N$  orthogonal matrices.

*Linear Algebra:* For vectors  $u, v \in \mathbb{R}^k$ ,  $\|u\|$  is the  $\ell_2$  norm of  $u$ ,  $\langle u, v \rangle = \sum_{i=1}^k u_i v_i$  denotes the standard inner product on  $\mathbb{R}^k$ , and  $\text{diag}(u)$  represents the  $k \times k$  diagonal matrix constructed by placing the entries of  $u$  along the diagonal. For a matrix  $M \in \mathbb{R}^{k \times k}$ ,  $\text{Tr}[M]$ ,  $\|M\|_{\text{op}}$ ,  $\|M\|$  represent the trace, operator (spectral) norm, and Frobenius norm of  $M$  respectively. If  $M$  is symmetric, we denote the sorted eigenvalues of  $M$  by  $\lambda_1(M) \geq \dots \geq \lambda_k(M)$ . We use  $\mathbf{1}_k$  to denote the vector  $(1, 1, \dots, 1)$  in  $\mathbb{R}^k$ ,  $\mathbf{0}_k$  to denote the zero vector  $(0, 0, \dots, 0)$  in  $\mathbb{R}^k$ ,  $I_k$  denotes the  $k \times k$  identity matrix, and  $e_1, e_2, \dots, e_k$  to denote the standard basis vectors in  $\mathbb{R}^k$ . When the dimension is clear from the context, we will abbreviate  $\mathbf{1}_k, \mathbf{0}_k, I_k$  as  $\mathbf{1}, \mathbf{0}, I$ . We use the bold-face font for vectors and matrices whose dimensions diverge as  $N$  (the dimension of the signal vector) grows to  $\infty$ . For example, the signal  $\mathbf{x}_\star \in \mathbb{R}^N$ , and the noise matrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$  are bold-faced.

*Probability:* We use  $\mathbb{E}[\cdot]$ ,  $\text{Var}[\cdot]$ ,  $\text{Cov}[\cdot]$  to denote expectations, variances, and covariances of random variables. The Gaussian distribution on  $\mathbb{R}^k$  with mean vector  $\mu \in \mathbb{R}^k$  and covariance matrix  $\Sigma \in \mathbb{R}^{k \times k}$  is denoted by  $\mathcal{N}(\mu, \Sigma)$ . For a finite set  $A$ ,  $\text{Unif}(A)$  represents the uniform distribution on  $A$ . We will use  $\text{Unif}(\mathbb{O}(N))$  to denote the Haar measure on the orthogonal group  $\mathbb{O}(N)$ . For any  $x \in \mathbb{R}$ , the probability measure  $\delta_x$  on  $\mathbb{R}$

denotes the point mass at  $x$ . We use  $\xrightarrow{\mathbb{P}}$  and  $\xrightarrow{d}$  to denote convergence in probability and distribution, respectively. For a sequence of real-valued random variables  $(Y_N)_{N \in \mathbb{N}}$ , we say that  $\text{plim } Y_N = y$  if  $Y_N \xrightarrow{\mathbb{P}} y$ ,  $\text{plim sup } Y_N \leq y$  if for any  $\epsilon > 0$ ,  $\lim_{N \rightarrow \infty} \mathbb{P}(Y_N \geq y + \epsilon) = 0$ , and  $\text{plim inf } Y_N \geq y$  if for any  $\epsilon > 0$ ,  $\lim_{N \rightarrow \infty} \mathbb{P}(Y_N \leq y - \epsilon) = 0$ .

## 2 Preliminaries

We begin by introducing the assumptions under which we study the rank-1 matrix estimation problem in (1), along with some key concepts that play an important role in this paper.

**Convergence and Asymptotic Equivalence of High-Dimensional Vectors.** We rely on the following notion of convergence.

**Definition 1.** Let  $(\mathbf{v}_1, \dots, \mathbf{v}_\ell)$  be a collection of random vectors in  $\mathbb{R}^N$ . We say that the empirical distribution of the entries of the vectors  $(\mathbf{v}_1, \dots, \mathbf{v}_\ell)$  converges to random variables  $(\mathbf{V}_1, \dots, \mathbf{V}_\ell)$  as  $N \rightarrow \infty$  if, for any test function  $h : \mathbb{R}^\ell \mapsto \mathbb{R}$  which satisfies:

$$|h(\mathbf{v}) - h(\mathbf{v}')| \leq L \|\mathbf{v} - \mathbf{v}'\| \cdot (1 + \|\mathbf{v}\| + \|\mathbf{v}'\|) \quad \forall \mathbf{v}, \mathbf{v}' \in \mathbb{R}^\ell, \quad (3)$$

for some  $L < \infty$ , we have:

$$\frac{1}{N} \sum_{i=1}^N h(v_1[i], \dots, v_\ell[i]) \xrightarrow{\mathbb{P}} \mathbb{E}[h(\mathbf{V}_1, \dots, \mathbf{V}_\ell)] \text{ as } N \rightarrow \infty.$$

We denote convergence in this sense using the notation:  $(\mathbf{v}_1, \dots, \mathbf{v}_\ell) \xrightarrow{W_2} (\mathbf{V}_1, \dots, \mathbf{V}_\ell)$ .

The following definition introduces a related notion of asymptotic equivalence of high-dimensional vectors.

**Definition 2.** Two  $N$ -dimensional random vectors  $\mathbf{u}$  and  $\mathbf{v}$  are asymptotically equivalent if:

$$\frac{\|\mathbf{u} - \mathbf{v}\|^2}{N} \xrightarrow{\mathbb{P}} 0 \quad \text{as } N \rightarrow \infty.$$

We denote equivalence in this sense using the notation:  $\mathbf{u} \xrightarrow{N \rightarrow \infty} \mathbf{v}$ .

**Signal and Noise Models.** We posit the following assumptions on the signal  $\mathbf{x}_*$ , the side information  $\mathbf{a}$ , and the noise  $\mathbf{W}$ .

**Assumption 1** (Signal Model). The signal vector and side information satisfy  $(\mathbf{x}_*; \mathbf{a}) \xrightarrow{W_2} (X_*; A)$  for some limiting random variables  $(X_*, A)$  with joint distribution  $\pi$  which satisfies  $\mathbb{E}[X_*^2] = 1$  and  $\mathbb{E}[\|A\|^2] < \infty$ .

**Assumption 2** (Noise Model). We model the noise matrix  $\mathbf{W}$  as a random matrix with eigen-decomposition:

$$\mathbf{W} = \mathbf{U} \cdot \text{diag}(\lambda_1(\mathbf{W}), \dots, \lambda_N(\mathbf{W})) \cdot \mathbf{U}^\top,$$

where the matrix of eigenvectors  $\mathbf{U} \sim \text{Unif}(\mathbb{O}(N))$  is a Haar-distributed random orthogonal matrix and the eigenvalues  $\lambda_1(\mathbf{W}), \dots, \lambda_N(\mathbf{W})$  are deterministic. We will assume that the operator norm of  $\mathbf{W}$  is bounded by a  $N$ -independent constant  $C$  and the spectral measure  $\mu_N$  of  $\mathbf{W}$ :

$$\mu_N \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N \delta_{\lambda_i(\mathbf{W})}$$

converges weakly to a compactly supported distribution  $\mu$  on  $\mathbb{R}$ . We require the limiting spectral measure  $\mu$  to be absolutely continuous with respect to the Lebesgue measure. Furthermore, the density of  $\mu$ , which we denote using the same symbol  $\mu : \mathbb{R} \mapsto \mathbb{R}$  is assumed to be Holder continuous in the sense:

$$|\mu(\lambda) - \mu(\lambda')| \leq L \cdot |\lambda - \lambda'|^\alpha \quad \forall \lambda, \lambda' \in \mathbb{R}. \quad (4)$$

for some constants  $L < \infty$  and  $\alpha > 0$ .

**Stieltjes and Hilbert Transform.** For any probability measure (or more generally, finite measure)  $\chi$  on  $\mathbb{R}$ , the Stieltjes Transform  $\mathcal{S}_\chi : \mathbb{C} \setminus \mathbb{R} \mapsto \mathbb{C}$  of  $\chi$  is defined as:

$$\mathcal{S}_\chi(z) \stackrel{\text{def}}{=} \int_{\mathbb{R}} \frac{\chi(d\lambda)}{z - \lambda} \quad \forall z \in \mathbb{C} \setminus \mathbb{R}.$$

It is well-known that a probability measure is uniquely determined by its Stieltjes transform (see for e.g., [3, Theorem 2.4.3]). Suppose that  $\chi$  is compactly supported and absolutely continuous with respect to the Lebesgue measure with density  $\chi(\cdot)$ , which is Holder continuous (cf. (4)). Then, the Hilbert transform  $\mathcal{H}_\chi : \mathbb{R} \mapsto \mathbb{R}$  of  $\chi$  is defined as the Cauchy principal value of the following singular integral:

$$\mathcal{H}_\chi(z) \stackrel{\text{def}}{=} \lim_{\epsilon \rightarrow 0} \frac{1}{\pi} \int_{|z-\lambda| \geq \epsilon} \frac{\chi(\lambda)}{z - \lambda} d\lambda \quad \forall z \in \mathbb{R}. \quad (5)$$

Under the Holder continuity requirement (4) on  $\chi$ , the above limit exists and  $\mathcal{H}_\sigma : \mathbb{R} \mapsto \mathbb{R}$  is also Holder continuous in the sense of (4); see e.g., [61, Section 2.1] and [34, Theorem 14.11a].

**Spectral Measure in the Signal Direction.** Let  $\nu_N$  denote the spectral measure of the observed matrix  $\mathbf{Y}$  in the direction of the signal:

$$\nu_N \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N \langle \mathbf{u}_i(\mathbf{Y}), \mathbf{x}_\star \rangle^2 \cdot \delta_{\lambda_i(\mathbf{Y})}, \quad (6)$$

where  $\lambda_1(\mathbf{Y}), \dots, \lambda_N(\mathbf{Y})$  denote the eigenvalues of  $\mathbf{Y}$  and  $\mathbf{u}_1(\mathbf{Y}), \dots, \mathbf{u}_N(\mathbf{Y})$  denote the corresponding eigenvectors. The measure  $\nu_N$  plays an important role in our analysis. The following lemma collects some important results regarding the asymptotic behavior of  $\nu_N$  as  $N \rightarrow \infty$ . In order to state the result, we introduce the function  $\phi : \mathbb{R} \mapsto \mathbb{R}$ , which will play an important role throughout the paper:

$$\phi(\lambda) \stackrel{\text{def}}{=} (1 - \pi\theta\mathcal{H}_\mu(\lambda))^2 + \pi^2\theta^2\mu^2(\lambda) \quad \lambda \in \mathbb{R}. \quad (7)$$

In the above display  $\mathcal{H}_\mu$  denotes the Hilbert transform of  $\mu$ .

**Lemma 1.** *We have,*

1. *The measure  $\nu_N$  converges weakly in probability to a compactly supported probability measure  $\nu$  on  $\mathbb{R}$  with Stieltjes Transform:*

$$\mathcal{S}_\nu(z) = \frac{\mathcal{S}_\mu(z)}{1 - \theta\mathcal{S}_\mu(z)} \quad \forall z \in \mathbb{C} \setminus \mathbb{R},$$

*where  $\mathcal{S}_\mu$  denotes the Stieltjes transform of the limiting spectral measure  $\mu$ .*

*Let  $\nu = \nu_{\parallel} + \nu_{\perp}$  denote the Lebesgue decomposition of  $\nu$  into the absolutely continuous part  $\nu_{\parallel}$  and the singular part  $\nu_{\perp}$ . Then,*

2. *For Lebesgue-almost every  $\lambda$ ,  $\phi(\lambda) \neq 0$  and for  $\nu_{\perp}$ -almost every  $\lambda$ ,  $\phi(\lambda) = 0$ .*
3. *The density of the absolutely continuous part of  $\nu$  is given by  $\mu(\cdot)/\phi(\cdot)$  where  $\mu(\cdot)$  denotes the density of  $\mu$  and the function  $\phi(\cdot)$  is as defined in (7).*

**Gaussian Channels.** We will also use some basic notions regarding Gaussian channels, introduced in the definition below.

**Definition 3** (Gaussian Channel). A Gaussian channel is a collection of real-valued random variables  $(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$  where the signal  $\mathbf{X}_\star$  and the side information  $\mathbf{A}$  are drawn from the prior  $\pi$  from Assumption 1:  $(\mathbf{X}_\star, \mathbf{A}) \sim \pi$  and the observations  $\mathbf{X}_1, \dots, \mathbf{X}_t$  are given by:  $\mathbf{X}_i = \alpha_i \mathbf{X}_\star + \mathbf{Z}_i \quad \forall i \in [t]$ , where  $\alpha_1, \dots, \alpha_t$  are real numbers and  $(\mathbf{Z}_1, \dots, \mathbf{Z}_t) \sim \mathcal{N}(0, \Sigma)$  are zero mean and jointly Gaussian random variables, sampled independently from  $(\mathbf{X}_\star, \mathbf{A})$ . We introduce some important notions related to Gaussian channels.

*MMSE and MMSE Estimator.* The minimum mean squared error (MMSE) for estimating the signal  $X_*$  based on the observations and the side information  $(X_1, \dots, X_t; A)$ , denoted by  $\text{MMSE}(X_*|X_1, \dots, X_t; A)$  is defined as:

$$\text{MMSE}(X_*|X_1, \dots, X_t; A) \stackrel{\text{def}}{=} \min_{f \in L^2(X_1, \dots, X_t; A)} \mathbb{E}[\{X_* - f(X_1, \dots, X_t; A)\}^2], \quad (8)$$

where the minimum is over the set  $L^2(X_1, \dots, X_t; A)$ , which denotes the set of all measurable functions  $f : \mathbb{R}^{t+k} \mapsto \mathbb{R}$  which satisfy  $\mathbb{E}[f^2(X_1, \dots, X_t; A)] < \infty$ . The function  $f \in L^2(X_1, \dots, X_t; A)$  that minimizes the RHS in (8) is called the MMSE estimator for  $X_*$ .

*DMMSE and DMMSE Estimator.* The divergence-free minimum mean squared error (DMMSE) for estimating the signal  $X_*$  based on  $(X_1, \dots, X_t; A)$ , denoted by  $\text{DMMSE}(X_*|X_1, \dots, X_t; A)$  is defined as:

$$\text{DMMSE}(X_*|X_1, \dots, X_t; A) \stackrel{\text{def}}{=} \min_{f \in L^2(X_1, \dots, X_t; A)} \mathbb{E}[\{X_* - f(X_1, \dots, X_t; A)\}^2] \quad \text{subject to } \mathbb{E}[Z_i f(X_1, \dots, X_t; A)] = 0 \quad \forall i \in [t]. \quad (9)$$

The constraints  $\mathbb{E}[Z_i f(X_1, \dots, X_t; A)] \forall i \in [t]$  in (9) require the estimator to be uncorrelated with the noise and are called *divergence-free* constraints. The function  $f$  which minimizes the RHS in (9) is called the DMMSE estimator for  $X_*$ .

*Scalar Gaussian Channels.* A particularly important role is played by *scalar* Gaussian channels, which refers to a collection of real-valued random variables  $(X_*, X; A)$  generated as follows:

$$(X_*, A) \sim \pi, \quad X|X_*, A \sim \mathcal{N}(\sqrt{\omega} \cdot X_*, 1 - \omega)$$

where  $\omega \in [0, 1]$  is called the signal-to-noise ratio (SNR) of the channel. Notice that a scalar Gaussian channel is a Gaussian channel with a single observation  $X$  normalized to satisfy  $\mathbb{E}[X^2] = 1$ . For a scalar Gaussian channel  $(X_*, X; A)$ , we define several important functions which will play a key role in this paper.

*MMSE Function and MMSE Estimator for Scalar Gaussian Channels.* The function  $\text{mmse}_\pi : [0, 1] \mapsto [0, 1]$  represents MMSE of a scalar Gaussian channel as a function of the SNR  $\omega$ :

$$\text{mmse}_\pi(\omega) \stackrel{\text{def}}{=} \text{MMSE}(X_*|X, A).$$

The function  $\varphi(\cdot|\omega) : \mathbb{R} \times \mathbb{R}^k \mapsto \mathbb{R}$  denotes the MMSE estimator for the scalar Gaussian channel at SNR  $\omega$ :

$$\varphi(x; a|\omega) \stackrel{\text{def}}{=} \mathbb{E}[X_*|X = x, A = a] \quad \forall x \in \mathbb{R}, a \in \mathbb{R}^k. \quad (10)$$

*DMMSE Function and DMMSE Estimator for Scalar Gaussian Channels.* The function  $\text{dmmse}_\pi : [0, 1] \mapsto [0, 1]$  represents DMMSE of a scalar Gaussian channel as a function of the SNR  $\omega$ :

$$\text{dmmse}_\pi(\omega) \stackrel{\text{def}}{=} \text{DMMSE}(X_*|X, A).$$

The function  $\bar{\varphi}(\cdot|\omega) : \mathbb{R} \times \mathbb{R}^k \mapsto \mathbb{R}$  denotes the DMMSE estimator for the scalar Gaussian channel at SNR  $\omega$ . Ma and Ping [50] have shown that the DMMSE estimator of a scalar Gaussian channel at SNR  $\omega$  is given by (see Lemma 2 in Appendix A.2):

$$\bar{\varphi}(x; a|\omega) \stackrel{\text{def}}{=} \begin{cases} \left(1 - \frac{\sqrt{\omega}}{\sqrt{1-\omega}} \cdot \mathbb{E}[Z\varphi(X; A|\omega)]\right)^{-1} \cdot \left(\varphi(x; a|\omega) - \frac{\mathbb{E}[Z\varphi(X; A|\omega)]}{\sqrt{1-\omega}} \cdot x\right) & : \omega < 1 \\ x & : \omega = 1. \end{cases} \quad (11)$$

Finally, we will impose the following regularity condition in our analysis.

**Assumption 3.** For any  $\omega \in [0, 1]$ , the MMSE estimator  $\varphi(\cdot|\omega) : \mathbb{R} \times \mathbb{R}^k \mapsto \mathbb{R}$  for the scalar Gaussian channel  $(X_*, A) \sim \pi$ ,  $X|A, X_* \sim \mathcal{N}(\sqrt{\omega} \cdot X_*, 1 - \omega)$  is continuously differentiable and Lipschitz.

*Remark 1.* In the absence of any side information, a sufficient condition for Assumption 3 to hold is that the signal random variable  $X_*$  is compactly supported (see [54, Remark 2.3] and [31, Lemma 3.8]).

### 3 Main results

We now present the main results obtained in this paper.

#### 3.1 Orthogonal Approximate Message Passing Algorithms

We introduce a class of iterative methods for rank-1 matrix estimation problem, called *Orthogonal Approximate Message Passing* (OAMP) algorithms.

**Definition 4** (OAMP algorithms). An OAMP algorithm generates iterates  $\mathbf{x}_1, \mathbf{x}_2, \dots$  in  $\mathbb{R}^N$  according to the update rule:

$$\mathbf{x}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}, \quad (12a)$$

where  $\mathbf{Y} \in \mathbb{R}^{N \times N}$  is the observed noisy matrix and  $\mathbf{a} \in \mathbb{R}^{N \times k}$  denotes the side information available for estimating  $\mathbf{x}_*$ . The estimate of  $\mathbf{x}_*$  at iteration  $t$  is obtained by applying a post-processing function  $\psi_t$  to the iterates  $\mathbf{x}_1, \dots, \mathbf{x}_t$  and the side-information  $\mathbf{a}$ :

$$\hat{\mathbf{x}}_t = \psi_t(\mathbf{x}_1, \dots, \mathbf{x}_t; \mathbf{a}). \quad (12b)$$

In the above equations, for each  $t \in \mathbb{N}$ , the matrix denoiser  $\Psi_t : \mathbb{R} \mapsto \mathbb{R}$  acts on the matrix  $\mathbf{Y}$  using the usual convention: if  $\mathbf{Y} = \mathbf{O} \text{diag}(\lambda_1, \dots, \lambda_N) \mathbf{O}^\top$  is the eigen-decomposition of  $\mathbf{Y}$ , then,  $\Psi_t(\mathbf{Y}) = \mathbf{O} \text{diag}(\Psi_t(\lambda_1), \dots, \Psi_t(\lambda_N)) \mathbf{O}^\top$ . Likewise, the iterate denoisers  $f_t : \mathbb{R}^{t-1} \times \mathbb{R}^k \mapsto \mathbb{R}$  and the post-processing function  $\psi_t : \mathbb{R}^t \times \mathbb{R}^k \mapsto \mathbb{R}$  act entry-wise on the  $N$  components of its vector inputs.

*State Evolution Random Variables.* Each OAMP algorithm is associated with a collection of *state evolution random variables*  $(\mathbf{X}_*, (\mathbf{X}_t)_{t \in \mathbb{N}}; \mathbf{A})$ , which describes the joint asymptotic behavior of the signal  $\mathbf{x}_*$ , the iterates  $(\mathbf{x}_t)_{t \in \mathbb{N}}$ , and the side information  $\mathbf{a}$ . The distribution of these random variables is given by:

$$(\mathbf{X}_*, \mathbf{A}) \sim \pi, \quad \mathbf{X}_t = \beta_t \mathbf{X}_* + \mathbf{Z}_t \quad \forall t \in \mathbb{N}, \quad (13a)$$

where  $(\beta_t)_{t \in \mathbb{N}}$  is defined via the recursion:

$$\beta_t \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{X}_* f_t(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{A})] \cdot \mathbb{E}[\Psi_t(\Lambda_\nu)], \quad (13b)$$

and  $(\mathbf{Z}_t)_{t \in \mathbb{N}}$  are zero mean jointly Gaussian random variables, sampled independently of  $(\mathbf{X}_*; \mathbf{A})$ , whose covariance matrix is given by the recursion:

$$\mathbb{E}[\mathbf{Z}_s \mathbf{Z}_t] = \mathbb{E}[\mathbf{X}_* \mathbf{F}_s] \mathbb{E}[\mathbf{X}_* \mathbf{F}_t] \cdot \text{Cov}_{\Lambda_\nu \sim \nu}[\Psi_s(\Lambda_\nu), \Psi_t(\Lambda_\nu)] + (\mathbb{E}[\mathbf{F}_s \mathbf{F}_t] - \mathbb{E}[\mathbf{X}_* \mathbf{F}_s] \mathbb{E}[\mathbf{X}_* \mathbf{F}_t]) \cdot \text{Cov}_{\Lambda \sim \mu}[\Psi_s(\Lambda), \Psi_t(\Lambda)]. \quad (13c)$$

In the above display,  $\mathbf{F}_s \stackrel{\text{def}}{=} f_s(\mathbf{X}_1, \dots, \mathbf{X}_{s-1}; \mathbf{A})$  and  $\mathbf{F}_t \stackrel{\text{def}}{=} f_t(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{A})$ .

*Requirements on Matrix Denoisers.* The matrix denoisers  $(\Psi_t)_{t \in \mathbb{N}}$  used in the OAMP algorithm should be continuous functions which do not change with  $N$ , and are required to satisfy the *trace-free constraint*:

$$\text{plim}_{N \rightarrow \infty} \frac{\text{Tr}[\Psi_t(\mathbf{Y})]}{N} = \mathbb{E}[\Psi_t(\Lambda)] = 0, \quad \Lambda \sim \mu. \quad (14)$$

*Requirements on Iterate Denoisers and Post-processing Functions.* For each  $t \in \mathbb{N}$ , the iterate denoiser  $f_t : \mathbb{R}^{t-1} \times \mathbb{R}^k \mapsto \mathbb{R}$  and the post-processing function  $\psi_t : \mathbb{R}^t \times \mathbb{R}^k \mapsto \mathbb{R}$  used in the OAMP algorithm should be continuously differentiable and Lipschitz functions which do not change with  $N$ . In addition, the iterate denoisers  $(f_t)_{t \in \mathbb{N}}$  are required to satisfy the *divergence-free constraint*:

$$\mathbb{E}[\partial_s f_t(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{A})] = 0 \quad \forall s \in [t-1], t \in \mathbb{N}, \quad (15)$$

where  $\partial_s f_t$  denotes the partial derivative of  $f_t(x_1, \dots, x_s, \dots, x_{t-1}; a)$  with respect to  $x_s$ .



Our first main result is the following theorem, which provides a characterization of the dynamics of an OAMP algorithm in terms of the associated state evolution random variables.

**Theorem 1** (State evolution of OAMP). *Consider a general OAMP algorithm of the form (12) and let  $\{\mathbf{X}_\star, (\mathbf{X}_t)_{t \in \mathbb{N}}, \mathbf{A}\}$  be the associated state evolution random variables. Then for any  $t \in \mathbb{N}$ ,*

$$(\mathbf{x}_\star, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}). \quad (16)$$

We present a heuristic derivation of this result in Section 4.1, and the complete proof in Appendix B.

*Remark 2* (Connections to AMP Algorithms for Compressed Sensing). The OAMP algorithm introduced in Definition 4 can be viewed as a natural analog of the orthogonal AMP (OAMP) [50] or vector AMP (VAMP) algorithm [69, 74] developed for compressed sensing (or regularized linear regression). The key feature of these algorithms, which is shared by the algorithm in Definition 1, is the use of trace-free matrix denoisers (cf. (14)) and divergence-free iterate denoisers (cf. (15)). These features significantly simplify the algorithm's state evolution. However, unlike the OAMP algorithm introduced in this work, the compressed sensing algorithms compute a matrix-vector multiplication involving a rotationally invariant matrix at each iteration. In contrast, only the observed matrix  $\mathbf{Y}$  is available in the spiked matrix model, not the rotationally invariant noise matrix  $\mathbf{W}$ . Although  $\mathbf{Y}$  is a rank-1 perturbation of  $\mathbf{W}$ ,  $\Psi_t(\mathbf{Y})$ , the matrix used by the OAMP algorithm at iteration  $t$ , cannot be expressed as a simple perturbation of  $\mathbf{W}$  for a general matrix denoiser  $\Psi_t$ . This complicates the analysis of these algorithms. We refer the reader to Fan [29, Remark 3.3] for related discussions.

### 3.2 The Optimal OAMP Algorithm

Theorem 1 provides a characterization of the asymptotic mean squared error (MSE) of the estimator  $\hat{\mathbf{x}}_t$  computed by a general OAMP algorithm (12) in terms of the associated state evolution random variables:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{x}_\star - \hat{\mathbf{x}}_t\|^2}{N} = \mathbb{E}|\mathbf{X}_\star - \psi_t(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})|^2.$$

The limiting value of the MSE depends implicitly on the functions  $\{\Psi_t, f_t, \psi_t\}_{t \geq 1}$  used in the OAMP algorithm since these functions determine the joint distribution of  $(\mathbf{X}_\star, (\mathbf{X}_t)_{t \in \mathbb{N}}; \mathbf{A})$ . As our second contribution, we use the above characterization to derive the optimal choice for these functions which minimizes the MSE. We call the resulting algorithm the *optimal OAMP algorithm*, and introduce it below.

**The Optimal OAMP Algorithm.** We first introduce the optimal OAMP algorithm in some simple corner cases, and then consider the typical case.

*Corner Cases.* If  $\text{mmse}_\pi(0) = \mathbb{E}\text{Var}[\mathbf{X}_\star | \mathbf{A}] = 0$ , then it is possible to reconstruct the signal perfectly from the side information alone. Hence, the optimal OAMP algorithm outputs the estimator:

$$\hat{\mathbf{x}}_t = \varphi(\mathbf{a}|0) \quad \forall t \in \mathbb{N}, \quad (17a)$$

where  $\varphi(\cdot|0)$  is the MMSE estimator for the Gaussian channel  $(\mathbf{X}_\star; \mathbf{A})$  (which operates at SNR  $\omega = 0$ ). In this case, the optimal OAMP algorithm achieves zero asymptotic mean squared error:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{x}_\star - \hat{\mathbf{x}}_t\|^2}{N} \stackrel{\text{Thm. 1}}{=} \mathbb{E}|\mathbf{X}_\star - \varphi(\mathbf{A}|0)|^2 \stackrel{(10)}{=} \mathbb{E}|\mathbf{X}_\star - \mathbb{E}[\mathbf{X}_\star | \mathbf{A}]|^2 = \mathbb{E}\text{Var}[\mathbf{X}_\star | \mathbf{A}] = 0 \quad \forall t \in \mathbb{N}.$$

In the other extreme, if  $\text{mmse}_\pi(0) = \mathbb{E}\text{Var}[\mathbf{X}_\star | \mathbf{A}] = 1$  (recall  $\mathbb{E}[\mathbf{X}_\star^2] = 1$  from Assumption 1), the optimal OAMP algorithm returns the trivial estimator:

$$\hat{\mathbf{x}}_t = 0 \quad \forall t \in \mathbb{N}. \quad (17b)$$

In this case, the asymptotic MSE of the optimal OAMP algorithm is:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{x}_\star - \hat{\mathbf{x}}_t\|^2}{N} = \text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{x}_\star\|^2}{N} = \mathbb{E}\mathbf{X}_\star^2 = 1 \quad \text{by Assumption 1.}$$

*Typical Case.* In the typical situation when  $\text{mmse}_\pi(0) = \mathbb{E}\text{Var}[\mathbf{X}_*|\mathbf{A}] \in (0, 1)$ , the optimal OAMP algorithm generates a sequence of iterates  $\mathbf{x}_1, \mathbf{x}_2, \dots$  using the update rule:

$$\mathbf{x}_t = \frac{1}{\sqrt{\omega_t}} \left(1 + \frac{1}{\rho_t}\right) \cdot \Psi_*(\mathbf{Y}; \rho_t) \cdot \bar{\varphi}(\mathbf{x}_{t-1}; \mathbf{a}|\omega_{t-1}), \quad (17c)$$

The estimator returned by the optimal OAMP algorithm at iteration  $t$  is:

$$\hat{\mathbf{x}}_t \stackrel{\text{def}}{=} \varphi(\mathbf{x}_t; \mathbf{a}|\omega_t). \quad (17d)$$

In the above equations,

- The matrix denoiser  $\Psi_*$  used by the optimal OAMP algorithm is given by:

$$\Psi_*(\lambda; \rho) = 1 - \left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \frac{\phi(\lambda)}{\phi(\lambda) + \rho} \quad \forall \lambda \in \mathbb{R}, \rho \in (0, \infty), \quad (17e)$$

where  $\Lambda \sim \mu$  and the function  $\phi: \mathbb{R} \mapsto \mathbb{R}$  was introduced in (7).

- $\omega_t$  and  $\rho_t$  are computed using the recursion:

$$\rho_t = \frac{1}{\text{dmmse}_\pi(\omega_{t-1})} - 1, \quad \omega_t = 1 - \left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho_t} \right] \right)^{-1} \cdot \mathbb{E} \left[ \frac{1}{\phi(\Lambda) + \rho_t} \right] \quad (17f)$$

initialized with  $\omega_0 \stackrel{\text{def}}{=} 0$ .

- $\text{dmmse}_\pi(\omega)$  denotes the DMMSE for a scalar Gaussian channel with SNR  $\omega$ , and  $\varphi(\cdot|\omega), \bar{\varphi}(\cdot|\omega)$  denote the MMSE and DMMSE estimators for the Gaussian channel (Definition 3).

*Remark 3.* The additional scaling factor  $1/\sqrt{\omega_t} \cdot (1 + 1/\rho_t)$  in (17c) is introduced to normalize the iterates so that  $\|\mathbf{x}_t\|^2/N \xrightarrow{\mathbb{P}} 1$ .

The following result characterizes the asymptotic MSE of the optimal OAMP algorithm in this case. Its proof can be found in Appendix C.

**Proposition 1.** *Assume that  $\text{mmse}_\pi(0) = \mathbb{E}\text{Var}[\mathbf{X}_*|\mathbf{A}] \in (0, 1)$ . Let  $(\mathbf{X}_*, (\mathbf{X}_t)_{t \in \mathbb{N}}; \mathbf{A})$  denote the state evolution random variables associated with the optimal OAMP algorithm (17). Then,*

1. *For each  $t \in \mathbb{N}$ ,  $\omega_t \in [0, 1)$  and  $\rho_t \in (0, \infty)$ . Moreover, the state evolution random variables  $(\mathbf{X}_*, \mathbf{X}_t; \mathbf{A})$  form a scalar Gaussian channel with SNR  $\omega_t$  and, the asymptotic MSE of the estimator  $\hat{\mathbf{x}}_t$  returned by the optimal OAMP algorithm in (17) is given by:*

$$\text{plim}_{N \rightarrow \infty} \frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_*\|^2}{N} = \text{mmse}_\pi(\omega_t).$$

2. *The sequences  $(\omega_t)_{t \in \mathbb{N}}$  and  $(\rho_t)_{t \in \mathbb{N}}$  from (17f) are non-decreasing and converge to limit points  $\omega_* \in [0, 1)$  and  $\rho_* \in (0, \infty)$  as  $t \rightarrow \infty$ . The limit points  $(\omega_*, \rho_*)$  solve the following fixed point equation in  $(\omega, \rho)$ :*

$$\omega = 1 - \left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right], \quad \rho = \frac{1}{\text{dmmse}_\pi(\omega)} - 1, \quad \Lambda \sim \mu. \quad (18)$$

Consequently as  $t \rightarrow \infty$ , the asymptotic MSE of  $\hat{\mathbf{x}}_t$  converges to:

$$\lim_{t \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_*\|^2}{N} = \text{mmse}_\pi(\omega_*).$$

### 3.3 An Optimality Result

Our final main result shows that the optimal OAMP algorithm introduced in (17) attains the best possible asymptotic MSE within a broad class of iterative algorithms, which we define formally below.

**Definition 5** (Iterative Algorithms). An iterative algorithm generates iterates  $\mathbf{r}_1, \mathbf{r}_2, \dots$  in  $\mathbb{R}^N$  according to the update rule:

$$\mathbf{r}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{r}_1, \dots, \mathbf{r}_{t-1}; \mathbf{a}) + g_t(\mathbf{r}_1, \dots, \mathbf{r}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}. \quad (19a)$$

The estimate of  $\mathbf{x}_*$  at iteration  $t$  is obtained by applying a post-processing function  $\psi_t$  to the iterates  $\mathbf{r}_1, \dots, \mathbf{r}_t$  and the side-information  $\mathbf{a}$ :

$$\hat{\mathbf{r}}_t = \psi_t(\mathbf{x}_1, \dots, \mathbf{x}_t; \mathbf{a}). \quad (19b)$$

For each  $t \in \mathbb{N}$ , the matrix denoiser  $\Psi_t : \mathbb{R} \mapsto \mathbb{R}$  is required to be continuous and the functions  $f_t : \mathbb{R}^{t-1} \times \mathbb{R}^k \mapsto \mathbb{R}$ ,  $g_t : \mathbb{R}^{t-1} \times \mathbb{R}^k \mapsto \mathbb{R}$ , and  $\psi_t : \mathbb{R}^t \times \mathbb{R}^k \mapsto \mathbb{R}$  are required to be continuously differentiable and Lipschitz. Furthermore,  $(\Psi_t)_{t \in \mathbb{N}}$ ,  $(f_t)_{t \in \mathbb{N}}$ ,  $(g_t)_{t \in \mathbb{N}}$ , and  $(\psi_t)_{t \in \mathbb{N}}$  do not change with the dimension  $N$ .

The following theorem states our optimality result.

**Theorem 2.** Let  $\hat{\mathbf{r}}_t$  be the estimator returned by any iterative algorithm of the form (19) after  $t \in \mathbb{N}$  iterations. Let  $\hat{\mathbf{x}}_t$  be the estimator returned by the optimal OAMP algorithm in (17) after  $t$  iterations. Then,

$$\text{plim inf}_{N \rightarrow \infty} \frac{\|\hat{\mathbf{r}}_t - \mathbf{x}_*\|^2}{N} \geq \text{plim}_{N \rightarrow \infty} \frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_*\|^2}{N}.$$

### 3.4 Information-Theoretic v.s. Computational Limits

A natural question is whether the optimal OAMP algorithm introduced in (17) achieves the smallest asymptotic estimation error among *all estimators*, not just estimators computable using efficient iterative algorithms. To address this question, we begin by recalling the conjecture of Barbier et al. [8] regarding the fundamental information-theoretic limits of this problem.

**Replica conjecture for the Bayes risk.** Under the assumption that the entries of  $(\mathbf{x}_*, \mathbf{a})$  are drawn i.i.d. from a prior  $\pi$ , the information-theoretically optimal estimator is the Bayes estimator  $\mathbb{E}[\mathbf{x}_* | \mathbf{Y}, \mathbf{a}]$ . Barbier et al. [8] provide a conjecture for the asymptotic MSE of this estimator (also called the Bayes risk) assuming the noise matrix drawn from the *trace ensemble* [8, 61] with density:

$$p(\mathbf{W}) \propto \exp\left(-\frac{N}{2} \sum_{i=1}^N V(\lambda_i(\mathbf{W}))\right) \quad \text{where } \lambda_1(\mathbf{W}), \dots, \lambda_N(\mathbf{W}) \text{ denote the eigenvalues of } \mathbf{W}, \quad (20)$$

and the *potential function*  $V : \mathbb{R} \mapsto \mathbb{R}$  is a functional parameter for the noise model. The authors develop a recipe to derive conjectured formulas for the asymptotic Bayes risk based on the non-rigorous replica method for polynomial potentials  $V$ , providing explicit formulas for quartic and sextic polynomials. As the degree of  $V$  increases, the complexity of the derivation and resulting formulas also increases, with no general formula available. We state the replica conjecture for the asymptotic Bayes risk when the potential  $V : \mathbb{R} \mapsto \mathbb{R}$  is a quartic polynomial:

$$V(\lambda) = \frac{\gamma \lambda^2}{2} + \frac{\kappa \lambda^4}{4} \quad \forall \lambda \in \mathbb{R} \quad \text{where} \quad \kappa = \kappa(\gamma) = \frac{8 - 9\gamma + \sqrt{64 - 144\gamma + 108\gamma^2 - 27\gamma^3}}{27}, \quad (21)$$

and  $\gamma \in [0, 1]$  is a parameter for the noise model. The particular choice of  $\kappa$  in (21) ensures that the limiting spectral measure ( $\mu$ ) of  $\mathbf{W}$  drawn from the trace ensemble (20) has unit variance [8].

**Conjecture 1** (Replica Conjecture for Bayes risk, Barbier et al. 8). Suppose that the entries of  $(\mathbf{x}_*, \mathbf{a})$  are drawn i.i.d. from a prior  $\pi$  with  $\text{mmse}_\pi(0) \in (0, 1)$  and the noise matrix  $\mathbf{W}$  is drawn from the trace ensemble (20) with the quartic potential function  $V$  from (21). Then,

$$\text{plim}_{N \rightarrow \infty} \frac{\|\mathbb{E}[\mathbf{x}_* | \mathbf{Y}, \mathbf{a}] - \mathbf{x}_*\|^2}{N} = \text{mmse}_\pi(\omega_{\text{IT}}),$$

for some  $\omega_{\text{IT}} \in (0, 1)$  and  $\rho_{\text{IT}} \in (0, \infty)$  which solve the following system of fixed point equations<sup>1</sup> (in  $\omega, \rho$ ):

$$m = 1 - \text{mmse}_\pi(\omega), \quad (22a)$$

$$\mathbb{E}[\mathbf{H}] = 1 - m, \quad (22b)$$

$$\chi = \mathbb{E}[\Lambda \mathbf{Q} \mathbf{H}], \quad (22c)$$

$$\hat{m} \equiv \frac{\omega}{1 - \omega} = \kappa \theta^2 \left( \frac{m}{1 - m} \mathbb{E}[\Lambda^2 \mathbf{H}] + \frac{\chi}{1 - m} \mathbb{E}[\Lambda \mathbf{H}] + \mathbb{E}[\Lambda^2 \mathbf{Q} \mathbf{H}] \right) + \gamma \theta^2 m, \quad (22d)$$

where  $m \in \mathbb{R}$ ,  $\chi \in \mathbb{R}$  are two intermediate variables,  $\Lambda \sim \mu$ , and the random variables  $\mathbf{Q} = \mathbf{Q}(\Lambda, m, \chi)$ ,  $\mathbf{H} = \mathbf{H}(\Lambda, \rho)$  are defined by

$$\mathbf{Q} \stackrel{\text{def}}{=} \kappa \theta^2 m \Lambda^2 + \kappa \theta^2 \chi \Lambda - \frac{\kappa \theta^2}{1 - m} \mathbb{E} [m \Lambda^2 \mathbf{H} + \chi \Lambda \mathbf{H}] + \frac{m}{1 - m}, \quad (22e)$$

$$\mathbf{H} \stackrel{\text{def}}{=} (\rho + \theta^2 a^2 (\gamma + 2a^2 \kappa)^2 + 1 - \theta (\gamma \Lambda - \theta \kappa \Lambda^2 + \kappa \Lambda^3))^{-1}. \quad (22f)$$

In the above equations,  $(\gamma, \kappa)$  are the parameters for the quartic potential function (21),  $\theta$  is a parameter for the spiked matrix model, and  $a^2 \stackrel{\text{def}}{=} (\sqrt{\gamma^2 + 12\kappa} - \gamma)/(6\kappa)$ .

On the other hand, from Proposition 1 that the asymptotic MSE of the estimator  $\hat{\mathbf{x}}_t$  returned by the optimal OAMP algorithm in (17) satisfies:

$$\lim_{t \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_*\|^2}{N} = \text{mmse}_\pi(\omega_*),$$

where  $(\omega_*, \rho_*)$  denote the solution of the state evolution fixed point equations (in  $\omega \in (0, 1), \rho \in (0, \infty)$ ):

$$\omega = \mathcal{F}_1(\rho), \quad \rho = \mathcal{F}_2(\omega) \quad \text{with} \quad \mathcal{F}_1(\rho) \stackrel{\text{def}}{=} 1 - \frac{\mathbb{E}_{\Lambda \sim \mu} \left[ \frac{1}{\phi(\Lambda) + \rho} \right]}{\mathbb{E}_{\Lambda \sim \mu} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right]}, \quad \mathcal{F}_2(\omega) \stackrel{\text{def}}{=} \frac{1}{\text{dmmse}_\pi(\omega)} - 1, \quad (23)$$

found by the recursion  $\rho_t = \mathcal{F}_2(\omega_{t-1})$ ,  $\omega_t = \mathcal{F}_1(\rho_t)$ . The following proposition (proved in Appendix D.2) shows that for the quartic potential, the replica fixed point equations (22) and the state evolution fixed point equations (23) are equivalent.

**Proposition 2.** *Assume that the prior  $\pi$  satisfies  $\text{mmse}_\pi(0) \in (0, 1)$ . Any solution  $(\omega, \rho)$  to the state evolution fixed point equations (23) with  $\omega \in (0, 1), \rho \in (0, \infty)$  is also a solution to the replica fixed point equations (22), and vice-versa.*

We remark that an analogous result holds for the sextic (degree 6) potential. Since the replica equations for the sextic ensemble in Barbier et al. [8] are even more involved, we do not provide the details here. We anticipate that the fixed point equations in (23) are a unified and concise reformulation of the replica fixed point equations and characterize the asymptotic Bayes risk for general potential functions  $V$ . This reformulation of the replica conjecture for the asymptotic Bayes risk may be more amenable to rigorous proof than the significantly more complicated replica formulas derived using the approach of Barbier et al. [8].

**Information-Theoretic Optimality and Sub-Optimality of OAMP.** We can simplify the fixed point equations in (23) (which are equivalent to the replica fixed point equations (22)) by eliminating  $\rho$  to obtain a single equation  $\omega = \mathcal{F}_1(\mathcal{F}_2(\omega))$ . In some cases, this equation has a unique solution, implying the optimal OAMP algorithm matches the asymptotic MSE of the Bayes estimator and is information-theoretically optimal (see Figure 1(a)). In other cases, the equation may have multiple solutions, as illustrated in Figure 1(b), where it has two stable fixed points. Here, the optimal OAMP algorithm converges to the inferior fixed point  $\omega_*$ , while the Bayes risk corresponds to the superior fixed point  $\omega_{\text{IT}}$ . Since  $\omega_* < \omega_{\text{IT}}$ , the optimal OAMP or any iterative algorithm (in the form of (19a)) with a constant ( $N$ -independent) number of iterations fails to achieve the Bayes-optimal MSE. Such scenarios also occur in i.i.d. Gaussian noise models [16, 55, 56], and it is conjectured that no polynomial time algorithm can achieve the information-theoretically optimal MSE in these cases.

<sup>1</sup>When the fixed point equations in (22) have multiple solutions, the correct fixed point is the one that minimizes a certain free energy function calculated in [8].

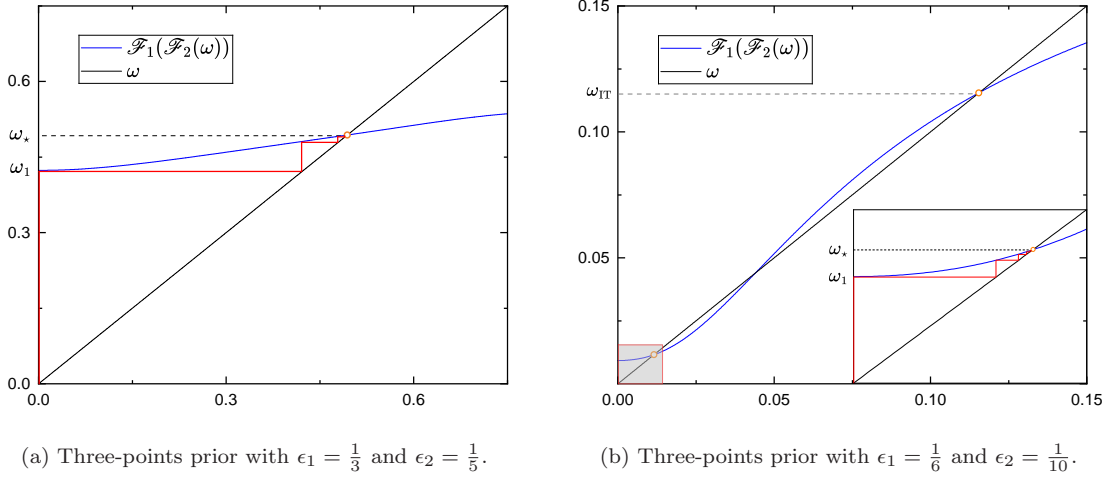


Figure 1: Plot of the fixed point equation  $\omega = \mathcal{F}_1(\mathcal{F}_2(\omega))$  for SNR  $\theta = 0.46$ , quartic noise model (21) with  $\gamma = 0$ , and a three-point signal prior  $\mathbf{X}_\star \sim \frac{\epsilon_1^2}{2}\delta_{\frac{1}{\epsilon_1}} + \frac{\epsilon_2^2}{2}\delta_{\frac{1}{\epsilon_2}} + (1 - \frac{\epsilon_1^2}{2} - \frac{\epsilon_2^2}{2})\delta_0$ .

## 4 Proof Ideas

This section presents some important intuitions and ideas used to obtain our main results.

### 4.1 Heuristic Derivation of State Evolution (Theorem 1)

We begin by an intuitive derivation of the state evolution result (Theorem 1) for OAMP algorithms, highlighting the key ideas involved in the proof. A rigorous proof of Theorem 1 is presented in Appendix B. Consider a general OAMP algorithm (Definition 4):

$$\mathbf{x}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}, \quad (24)$$

where  $(\Psi_t)_{t \in \mathbb{N}}$  are trace-free (cf. (14)) and  $(f_t)_{t \in \mathbb{N}}$  are divergence-free (cf. (15)). For ease of exposition, we present the key ideas assuming that the matrix denoisers  $(\Psi_t)_{t \in \mathbb{N}}$  are *polynomials* (the general case follows by a polynomial approximation argument). Let us decompose  $f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a})$  into a component along the direction of  $\mathbf{x}_\star$  and a component perpendicular to it:

$$f_t(\mathbf{x}_{<t}; \mathbf{a}) = \alpha_t \mathbf{x}_\star + \mathbf{f}_t^\perp \quad \text{where} \quad \alpha_t \stackrel{\text{def}}{=} \frac{\langle \mathbf{x}_\star, f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \rangle}{\|\mathbf{x}_\star\|^2}, \quad \mathbf{f}_t^\perp \stackrel{\text{def}}{=} f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) - \alpha_t \mathbf{x}_\star.$$

Using this decomposition, we can write the new iterate  $\mathbf{x}_t$  in (24) as

$$\mathbf{x}_t = \alpha_t \cdot \Psi_t(\mathbf{Y}) \mathbf{x}_\star + \Psi_t(\mathbf{Y}) \mathbf{f}_t^\perp. \quad (25)$$

As pointed out in Remark 2, the main difficulty in analyzing the algorithm above is that the matrix  $\Psi_t(\mathbf{Y})$  is not rotationally invariant. To address this, we express the update equation (25) in terms of the rotationally invariant noise matrix  $\mathbf{W}$  by expanding the polynomial  $\Psi_t(\mathbf{Y})$  in terms of  $\mathbf{W}$  using the fact that  $\mathbf{Y} = \frac{\theta}{N} \mathbf{x}_\star \mathbf{x}_\star^\top + \mathbf{W}$ . Although this expansion can initially seem complicated, the key insight which makes it tractable is to show that the following approximations hold:

$$\Psi_t(\mathbf{Y}) \mathbf{x}_\star \stackrel{N \rightarrow \infty}{\simeq} \tilde{\Psi}_t(\mathbf{W}) \mathbf{x}_\star, \quad \Psi_t(\mathbf{Y}) \mathbf{f}_t^\perp \stackrel{N \rightarrow \infty}{\simeq} \Psi_t(\mathbf{W}) \mathbf{f}_t^\perp, \quad (26)$$

where  $\tilde{\Psi}_t : \mathbb{R} \mapsto \mathbb{R}$  is polynomial obtained by appropriately transforming  $\Psi_t$ ; we refer the reader to Appendix B for details of the above approximations. The transformed matrix denoiser  $\tilde{\Psi}_t$  is determined by  $\Psi_t$  via a

complicated recursion and is not necessarily trace-free. Substituting (26) into (25), we get

$$\begin{aligned} \mathbf{x}_t &\approx \alpha_t \cdot \tilde{\Psi}_t(\mathbf{W})\mathbf{x}_* + \Psi_t(\mathbf{W})\mathbf{f}_t^\perp = \alpha_t \cdot \frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \cdot \mathbf{x}_* + \alpha_t \cdot \left( \tilde{\Psi}_t(\mathbf{W}) - \frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \cdot \mathbf{I}_N \right) \mathbf{x}_* + \Psi_t(\mathbf{W})\mathbf{f}_t^\perp \\ &\stackrel{\text{def}}{=} \underbrace{\alpha_t \cdot \frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \cdot \mathbf{x}_*}_{\text{Signal Component}} + \underbrace{\left( \tilde{\Psi}_t(\mathbf{W}) - \frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \cdot \mathbf{I}_N \right) \mathbf{x}_* + \Psi_t(\mathbf{W})\mathbf{f}_t^\perp}_{\text{Eff. Noise}}, \end{aligned} \quad (27a)$$

where the effective noise  $\mathbf{z}_t$  is defined as:

$$\begin{aligned} \mathbf{z}_t &\stackrel{\text{def}}{=} \alpha_t \cdot \left( \tilde{\Psi}_t(\mathbf{W}) - \frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \cdot \mathbf{I}_N \right) \mathbf{x}_* + \Psi_t(\mathbf{W})\mathbf{f}_t^\perp \\ &= \alpha_t \cdot \hat{\Psi}_t(\mathbf{W})\mathbf{x}_* + \Psi_t(\mathbf{W})\mathbf{f}_t^\perp \quad \text{where} \quad \hat{\Psi}_t(\mathbf{W}) \stackrel{\text{def}}{=} \tilde{\Psi}_t(\mathbf{W}) - \frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \cdot \mathbf{I}_N. \end{aligned} \quad (27b)$$

To derive the limiting distribution of  $\mathbf{x}_t$ , we analyze the signal and noise components separately.

**Signal Component.** Notice that the signal component in (27a) involves the factor  $\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]/N$ . Computing the limiting value of this factor is challenging because the transformed matrix denoiser  $\tilde{\Psi}_t$  does not have a convenient formula and is determined by  $\Psi_t$  via a complicated recursion (see Appendix B). The key idea which yields a concise formula for this factor is to observe:

$$\frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \stackrel{\text{(a)}}{\approx} \frac{\mathbf{x}_*^\top \tilde{\Psi}_t(\mathbf{W})\mathbf{x}_*}{N} \stackrel{\text{(b)}}{\approx} \frac{\mathbf{x}_*^\top \Psi_t(\mathbf{Y})\mathbf{x}_*}{N} \stackrel{\text{(c)}}{=} \int_{\mathbb{R}} \Psi_t(\lambda) \nu_N(d\lambda) \stackrel{\text{(d)}}{\approx} \mathbb{E}_{\Lambda_\nu \sim \nu}[\Psi_t(\Lambda_\nu)], \quad (28)$$

where the approximation in (a) follows from standard concentration results for quadratic forms of rotationally invariant matrices (see Fact 2 in Appendix F), (b) follows from the defining property of  $\tilde{\Psi}_t(\mathbf{W})$  in (26), (c) follows by recalling the definition of the spectral measure in the signal direction ( $\nu_N$ ) from (6), and (d) follows from the weak convergence of the compactly supported measure  $\nu_N$  to  $\nu$  (Lemma 1, item (1)).

**Noise Component.** The update rule (27b) for the effective noise  $\mathbf{z}_t$  is written in a form for which existing state evolution results apply [25, 26, 29, 50, 69, 74, 75]. Indeed,  $\hat{\Psi}_t(\mathbf{W})$  and  $\Psi_t(\mathbf{W})$  are functions of the rotationally-invariant matrix  $\mathbf{W}$  and have zero trace; and  $\mathbf{f}_t^\perp$  is a divergence-free function of  $\{\mathbf{x}_i\}_{i < t}$ . By appealing to existing results on the dynamics of AMP algorithms driven by rotationally invariant matrices, we show that the effective noise component  $\mathbf{z}_t$  converges to a centered Gaussian random variable  $\mathbf{Z}_t$ . Moreover, the covariance of  $\mathbf{Z}_t$  and  $\mathbf{Z}_s$  can be heuristically calculated based on a convenient property of OAMP algorithms. Specifically, for any  $G_1(\mathbf{W}), G_2(\mathbf{W}) \in \{\hat{\Psi}_t(\mathbf{W}), \Psi_t(\mathbf{W}), \Psi_s(\mathbf{W})\}$  and  $\mathbf{v}_1, \mathbf{v}_2 \in \{\mathbf{x}_*, \mathbf{f}_s^\perp, \mathbf{f}_t^\perp\}$ , we have

$$\frac{\mathbf{v}_1^\top G_1(\mathbf{W})G_2(\mathbf{W})\mathbf{v}_2}{N} \approx \frac{\text{Tr}[G_1(\mathbf{W})G_2(\mathbf{W})]}{N} \cdot \frac{\langle \mathbf{v}_1, \mathbf{v}_2 \rangle}{N}. \quad (29)$$

The intuition for the above property is that the random vectors  $\mathbf{v}_1, \mathbf{v}_2$  behave as if they are independent of the noise matrix  $\mathbf{W}$ , and the divergence-free and trace-free requirements imposed on OAMP algorithms are crucial for the validity of the above approximation. Using the above property, we immediately have that

$$\mathbb{E}[\mathbf{Z}_s \mathbf{Z}_t] \approx \frac{\langle \mathbf{z}_s, \mathbf{z}_t \rangle}{N} = \frac{\left\langle \alpha_t \hat{\Psi}_t(\mathbf{W})\mathbf{x}_* + \Psi_t(\mathbf{W})\mathbf{f}_t^\perp, \alpha_s \hat{\Psi}_s(\mathbf{W})\mathbf{x}_* + \Psi_s(\mathbf{W})\mathbf{f}_s^\perp \right\rangle}{N} \quad (30a)$$

$$\approx \alpha_s \alpha_t \cdot \frac{\text{Tr}[\hat{\Psi}_t(\mathbf{W}), \hat{\Psi}_s(\mathbf{W})]}{N} \cdot \frac{\|\mathbf{x}_*\|^2}{N} + \frac{\text{Tr}[\Psi_t(\mathbf{W})\Psi_s(\mathbf{W})]}{N} \cdot \frac{\langle \mathbf{f}_t^\perp, \mathbf{f}_s^\perp \rangle}{N}, \quad (30b)$$

where the cross terms vanish due to the orthogonality of  $\mathbf{x}_*$  and  $\mathbf{f}_t^\perp$  (and  $\mathbf{f}_s^\perp$ ). The limiting value of  $\text{Tr}[\hat{\Psi}_t(\mathbf{W}), \hat{\Psi}_s(\mathbf{W})]/N$  can be expressed in terms of  $\nu$  (the spectral measure in the signal direction) using the argument from (28). Replacing the various normalized inner products by their limiting values eventually leads to the claimed state evolution (13).

## 4.2 Derivation of the Optimal OAMP Algorithm

Next, we present an intuitive derivation of the optimal OAMP algorithm introduced in (17). While this derivation does not show that this algorithm attains the smallest estimation error among all iterative algorithms (as claimed in Theorem 2), it provides a simple and natural approach to derive the matrix denoisers and the iterate denoisers used by the algorithm. For simplicity, we consider the class of simple memory-free OAMP algorithms in which the iterate denoiser  $f_t$  only depends on the most recent iterate (cf. (12)):

$$\mathbf{x}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{x}_{t-1}; \mathbf{a}), \quad \forall t \in \mathbb{N}, \quad (31)$$

We will use a natural greedy heuristic to derive reasonable choices for the matrix and iterate denoisers. Specifically, we will derive a sensible choice of  $\Psi_t$  and  $f_t$ , assuming we have already specified the iterate and matrix denoisers for the first  $t - 1$  iterations. For any candidate  $\Psi_t$  and  $f_t$ , the distribution of the state evolution random variable  $\mathbf{X}_t$  associated with the  $\mathbf{x}_t$  is given by:

$$(\mathbf{X}_*, \mathbf{A}) \sim \pi, \quad \mathbf{X}_t = \beta_t \mathbf{X}_* + \mathbf{Z}_t, \quad (32a)$$

where  $\mathbf{Z}_t$  is a centered Gaussian random variable independent of  $(\mathbf{X}_*, \mathbf{A})$ , and

$$\beta_t = \mathbb{E}[\mathbf{X}_* f_t(\mathbf{X}_{t-1}; \mathbf{A})] \cdot \mathbb{E}[\Psi_t(\Lambda_\nu)], \quad (32b)$$

$$\mathbb{E}[\mathbf{Z}_t^2] = (\mathbb{E}[\mathbf{X}_* f_t(\mathbf{X}_{t-1}; \mathbf{A})])^2 \cdot \text{Var}[\Psi_t(\Lambda_\nu)] + \left( \mathbb{E}[f_t(\mathbf{X}_{t-1}; \mathbf{A})^2] - (\mathbb{E}[\mathbf{X}_* f_t(\mathbf{X}_{t-1}; \mathbf{A})])^2 \right) \cdot \mathbb{E}[\Psi_t^2(\Lambda)]. \quad (32c)$$

In the above definitions, the random variables  $\Lambda \sim \mu$ ,  $\Lambda_\nu \sim \nu$  are independent of all other random variables. A natural design principle is to choose  $f_t, \Psi_t$  to maximize the SNR  $\omega_t$  of the Gaussian channel  $(\mathbf{X}_*, \mathbf{X}_t; \mathbf{A})$ , which is given by the squared cosine similarity between  $\mathbf{X}_t$  and  $\mathbf{X}_*$ :

$$\omega_t \stackrel{\text{def}}{=} \frac{(\mathbb{E}[\mathbf{X}_* \mathbf{X}_t])^2}{\mathbb{E}[\mathbf{X}_*^2] \cdot \mathbb{E}[\mathbf{X}_t^2]} = \frac{\beta_t^2}{\beta_t^2 + \mathbb{E}[\mathbf{Z}_t^2]} = \frac{(\mathbb{E}[\Psi_t(\Lambda_\nu)])^2}{\mathbb{E}[\Psi_t^2(\Lambda_\nu)] + \frac{\delta_t}{1-\delta_t} \cdot \mathbb{E}[\Psi_t^2(\Lambda)]} \quad \text{where} \quad \delta_t \stackrel{\text{def}}{=} 1 - \frac{(\mathbb{E}[\mathbf{X}_* f_t(\mathbf{X}_{t-1}; \mathbf{A})])^2}{\mathbb{E}[\mathbf{F}_t^2]}. \quad (33)$$

We consider the two maximization problems over  $\Psi_t$  and  $f_t$  one-by-one.

**Optimal choice of  $\Psi_t$ .** Choosing  $\Psi_t$  to maximize  $\omega_t$  is equivalent to solving the following optimization problem:

$$\max_{\Psi} \frac{(\mathbb{E}[\Psi(\Lambda_\nu)])^2}{\mathbb{E}[\Psi^2(\Lambda_\nu)] + \rho^{-1} \cdot \mathbb{E}[\Psi^2(\Lambda)]} \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0. \quad (34)$$

with the parameter  $\rho \in (0, \infty)$  fixed at  $\rho = \delta_t^{-1} - 1$  (notice  $\delta_t$  does not depend on  $\Psi_t$ ). The optimization problem in (34) can be written in the more convenient but equivalent form:

$$\min_{\Psi} 1 - \frac{(\mathbb{E}[\Psi(\Lambda_\nu)])^2}{\mathbb{E}[\Psi^2(\Lambda_\nu)] + \rho^{-1} \cdot \mathbb{E}[\Psi^2(\Lambda)]} \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0 \quad (35a)$$

$$\stackrel{(a)}{=} \min_{\Psi} \min_{c \in \mathbb{R}} \mathbb{E} \left[ (1 - c\Psi(\Lambda_\nu))^2 \right] + \rho^{-1} \cdot \mathbb{E}[c^2 \Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[c\Psi(\Lambda)] = 0 \quad (35b)$$

$$\stackrel{(b)}{=} \min_{\Psi} \mathbb{E} \left[ (1 - \Psi(\Lambda_\nu))^2 \right] + \rho^{-1} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0, \quad (35c)$$

where step (a) can be verified by solving the inner quadratic minimization problem over  $c$  explicitly, and step (b) is due to a change of variable  $\Psi \mapsto c\Psi$ ; note that  $c\Psi$  also satisfies the zero-trace constraint. Note that any minimizer of (35c) is also a maximizer of (34) since the objective function of (34) is invariant to a rescaling of  $\Psi$ . Let  $\nu = \nu_\perp + \nu_\parallel$  be the Lebesgue decomposition of  $\nu$  into the singular part  $\nu_\perp$  and the absolutely continuous part  $\nu_\parallel$ . Let  $S$  denote the set:

$$S \stackrel{\text{def}}{=} \{\lambda \in \mathbb{R} : \phi(\lambda) \neq 0\} \quad \text{where} \quad \phi(\lambda) \stackrel{\text{def}}{=} (1 - \pi\theta\mathcal{H}_\mu(\lambda))^2 + \pi^2\theta^2\mu^2(\lambda). \quad (36)$$

From Lemma 1 (items (2) and (3)), we know that  $\nu_\perp(S) = 0$  and the density of  $\nu_\parallel$  with respect to the Lebesgue measure is given by  $\mu(\cdot)/\phi(\cdot)$  where  $\mu(\cdot)$  is the density of  $\mu$  (recall Assumption 2). Hence, the Lagrangian of the minimization problem (35c) is given by:

$$\begin{aligned} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi^2(\Lambda)] - 2\gamma\mathbb{E}[\Psi(\Lambda)] &= \int_{\mathbb{R}} |\Psi(\lambda) - 1|^2 \nu(d\lambda) + \int_{\mathbb{R}} \left( \frac{\Psi^2(\lambda)}{\rho} - 2\gamma\Psi(\lambda) \right) \mu(\lambda) d\lambda \\ &= \int_{S^c} (\Psi(\lambda) - 1)^2 \nu_\perp(d\lambda) + \int_S \left( (\Psi(\lambda) - 1)^2 \cdot \frac{1}{\phi(\lambda)} + \frac{\Psi^2(\lambda)}{\rho} - 2\gamma\Psi(\lambda) \right) \mu(d\lambda), \end{aligned}$$

where  $\gamma \in \mathbb{R}$  is the Lagrange multiplier. The Lagrangian is minimized by minimizing the integrand pointwise, leading to the following formula for the minimizer  $\Psi_\gamma$ :

$$\Psi_\gamma(\lambda) \stackrel{\text{def}}{=} \frac{\gamma + \frac{1}{\phi(\lambda)}}{\frac{1}{\phi(\lambda)} + \frac{1}{\rho}} = 1 + \frac{(\rho\gamma - 1) \cdot \phi(\lambda)}{\rho + \phi(\lambda)} \quad \forall \lambda \in S \text{ and } \Psi_\gamma(\lambda) = 1 \quad \forall \lambda \in S^c.$$

The solution of the constrained minimization problem (35c) is obtained by choosing the Lagrange multiplier  $\gamma$  so that the constraint  $\mathbb{E}_{\Lambda \sim \mu}[\Psi_\gamma(\Lambda)] = 0$  is fulfilled. Hence, the optimizer of (35c) and (34) is:

$$\Psi(\lambda) = 1 - \left( \mathbb{E}_{\Lambda \sim \mu} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \frac{\phi(\lambda)}{\phi(\lambda) + \rho} \stackrel{\text{def}}{=} \Psi_\star(\lambda; \rho) \quad \forall \lambda \in \mathbb{R}. \quad (37)$$

This is precisely the matrix denoiser used by the optimal OAMP algorithm in (17).

**Optimal choice of  $f_t$ .** Notice from (33) the iterate denoiser  $f_t$  which maximizes the SNR  $\omega_t$  can be derived by minimizing  $\delta_t$  with respect to  $f_t$ . Hence, we consider the following minimization problem:

$$\min_f 1 - \frac{(\mathbb{E}[\mathbf{X}_\star f(\mathbf{X}_{t-1}; \mathbf{A})])^2}{\mathbb{E}[f^2(\mathbf{X}_{t-1}; \mathbf{A})]} \quad \text{subject to} \quad \mathbb{E}[f'(\mathbf{X}_{t-1}; \mathbf{A})] = 0. \quad (38)$$

In the above definition,  $\mathbf{X}_{t-1} = \beta_{t-1}\mathbf{X}_\star + \sigma_{t-1}\mathbf{Z}$  is the state evolution random variable for  $\mathbf{x}_{t-1}$ ,  $\mathbf{Z} \sim \mathcal{N}(0, 1)$  is independent of  $(\mathbf{X}_\star, \mathbf{A}) \sim \pi$ , and  $\beta_{t-1} \in \mathbb{R}$  and  $\sigma_{t-1} > 0$  are fixed parameters (determined by the matrix and iterate denoisers used in the first  $t-1$  iterations). Repeating the arguments used in (35), we obtain:

$$\min_f 1 - \frac{(\mathbb{E}[\mathbf{X}_\star f(\mathbf{X}_{t-1}; \mathbf{A})])^2}{\mathbb{E}[f^2(\mathbf{X}_{t-1}; \mathbf{A})]}, \quad \text{subject to} \quad \mathbb{E}[f'(\mathbf{X}_{t-1}; \mathbf{A})] = 0 \quad (39a)$$

$$= \min_f \min_{c \in \mathbb{R}} \mathbb{E} \left[ (\mathbf{X}_\star - cf(\mathbf{X}_{t-1}; \mathbf{A}))^2 \right], \quad \text{subject to} \quad \mathbb{E}[cf'(\mathbf{X}_{t-1}; \mathbf{A})] = 0 \quad (39b)$$

$$= \min_f \mathbb{E} \left[ (\mathbf{X}_\star - f(\mathbf{X}_{t-1}; \mathbf{A}))^2 \right], \quad \text{subject to} \quad \mathbb{E}[f'(\mathbf{X}_{t-1}; \mathbf{A})] = 0 \quad (39c)$$

Recalling the definition of the DMMSE estimator from Definition 3 (see also Lemma 2 in Appendix A.2), the minimizer of (39c) is:

$$f_t(x; a) = \bar{\varphi} \left( \frac{x}{\sqrt{\beta_{t-1}^2 + \sigma_{t-1}^2}}; a \mid \omega_{t-1} \right) \quad \text{where} \quad \omega_{t-1} \stackrel{\text{def}}{=} \frac{\beta_{t-1}^2}{\beta_{t-1}^2 + \sigma_{t-1}^2}. \quad (40)$$

**Optimal OAMP algorithm.** Using the matrix and iterate denoisers derived above, we obtain the following memory-free OAMP algorithm (cf. (31)):

$$\mathbf{x}_t = c_t \cdot \Psi_\star(\mathbf{Y}; \rho_t) \cdot \bar{\varphi}_t(\mathbf{x}_{t-1}; \mathbf{a} \mid \omega_{t-1}), \quad (41)$$

where defined  $c_t \stackrel{\text{def}}{=} 1/\sqrt{\beta_t^2 + \sigma_t^2}$  and made a slight change of variable by absorbing the normalization factor (acting on the  $x$  input of  $\bar{\varphi}$ ) into the definition of  $\mathbf{x}_t$ . This scaling parameter  $c_t$  ensures state evolution random variable  $\mathbf{X}_t$  corresponding to  $\mathbf{x}_t$  satisfies  $\mathbb{E}[\mathbf{X}_t^2] = 1$ . This is precisely the optimal OAMP algorithm introduced in (17). The recursions for  $\omega_t, \rho_t$ , as well as the alternative formula  $c_t = 1/\sqrt{\omega_t} \cdot (1 + 1/\rho_t)$  can be derived by specializing the general state evolution equations in (32) to the optimal OAMP algorithm; see the proof of Proposition 1 in Appendix C for details.



### 4.3 Proof of the Optimality Result (Theorem 2)

We now present the proof of Theorem 2. We will introduce the key ideas involved in the form of some intermediate results, whose proofs are deferred to Appendix E.

**Implementing Iterative Algorithms using OAMP Algorithms.** Our proof builds on the techniques introduced in the work of Celentano et al. [16] and Montanari and Wu [56]. These works show that when the noise matrix  $\mathbf{W}$  has i.i.d. Gaussian entries, any estimator computed using  $t \in \mathbb{N}$  iterations of an iterative algorithm can also be computed using  $t$  iterations of a suitably designed AMP algorithm. Consequently, we can restrict the search space for the optimal  $t$ -iteration iterative algorithm to the space of  $t$ -iteration AMP algorithms, whose dynamics are significantly easier to analyze. In light of these results, it is natural to hope that any estimator computed using  $t$  iterations of a general iterative algorithm (as defined in Definition 5) can be computed using  $t$  iterations of a suitably designed OAMP algorithm (as defined in Definition 4). Unfortunately, the argument used by Celentano et al. [16] and Montanari and Wu [56] breaks down in our case because we require that the iterate denoisers used in the OAMP algorithm satisfy the divergence-free requirement in (15). To illustrate the difficulty caused by imposing the divergence-free requirement (15) and motivate our proof strategy, we consider the first two iterations of a simple iterative algorithm:

$$\mathbf{r}_1 = \Psi_1(\mathbf{Y}) \cdot f_1(\mathbf{a}), \quad \mathbf{r}_2 = \Psi_2(\mathbf{Y})f_2(\mathbf{r}_1; \mathbf{a}), \quad (42)$$

which we refer to as the *target algorithm*, and try to design an OAMP algorithm, which can reconstruct the iterates of the target algorithm. For simplicity, we assume that the matrix denoisers used by the target algorithm respect the trace-free requirement  $\mathbb{E}_{\Lambda \sim \mu}[\Psi_1(\Lambda)] = \mathbb{E}_{\Lambda \sim \mu}[\Psi_2(\Lambda)] = 0$  (requirement (14) in Definition 4). Since the divergence-free constraint (15) is inactive for the first iteration, we can construct a valid OAMP algorithm which simply copies the first iteration of the target algorithm  $\mathbf{x}_1 = \Psi_1(\mathbf{Y})f_1(\mathbf{a})$ . Let  $(\mathbf{X}_*, \mathbf{X}_1; \mathbf{A})$  be the state evolution random variables corresponding to the first iteration of this OAMP algorithm. Next, we consider the second iteration of the target algorithm (42). Notice that the denoiser  $f_2$  might not satisfy the divergence-free requirement  $\mathbb{E}[\partial_1 f_2(\mathbf{X}_1; \mathbf{A})] = 0$  imposed in the definition of OAMP algorithms (Definition 4). To ensure this requirement is fulfilled, we are forced to design the second iteration of the OAMP algorithm as:

$$\mathbf{x}_2 = \Psi_2(\mathbf{Y}) \cdot \bar{f}_2(\mathbf{x}_1; \mathbf{a}) \quad \text{where} \quad \bar{f}_2(x; a) \stackrel{\text{def}}{=} f_2(x; a) - \mathbb{E}[\partial_1 f_2(\mathbf{X}_1; \mathbf{A})] \cdot x \quad \forall x \in \mathbb{R}, a \in \mathbb{R}^k.$$

By construction  $\bar{f}_2$  satisfies the divergence-free constraint  $\mathbb{E}[\partial_1 \bar{f}_2(\mathbf{X}_1; \mathbf{A})] = 0$ . We can express the second iterate of the target algorithm (42) as:

$$\mathbf{r}_2 = \mathbf{x}_2 + \mathbb{E}[\partial_1 f_1(\mathbf{X}_1; \mathbf{A})] \cdot \Psi_2(\mathbf{Y}) \cdot \mathbf{x}_1 \stackrel{\text{(a)}}{=} \mathbf{x}_2 + \mathbb{E}[\partial_1 f_1(\mathbf{X}_1; \mathbf{A})] \cdot \underbrace{\Psi_2(\mathbf{Y})\Psi_1(\mathbf{Y}) \cdot f_1(\mathbf{a})}_{\text{uncomputed}}, \quad (43)$$

where in step (a) we recalled that  $\mathbf{x}_1 = \Psi_1(\mathbf{Y}) \cdot f_1(\mathbf{a})$ . In particular, we are unable to reconstruct the second iterate of the target algorithm in (42) using the first two iterations of the OAMP algorithm we have designed since the matrix-vector product  $\Psi_2(\mathbf{Y})\Psi_1(\mathbf{Y}) \cdot f_1(\mathbf{a})$  in (43) has not been computed by the OAMP algorithm in the sense that it cannot be expressed in terms of the OAMP iterates  $\mathbf{x}_1, \mathbf{x}_2$ .

**Lifted OAMP Algorithms.** To address the issue highlighted above, we introduce a more powerful class of algorithms called *lifted OAMP algorithms*, which are parameterized by a degree parameter  $D \in \mathbb{N}$ , which modulates the computational power of these algorithms. A degree- $D$  OAMP algorithm can compute  $D$  matrix-vector products per iteration. For instance, in the first iteration, a degree- $D$  lifted OAMP algorithm computes the matrix-vector products:

$$(\mathbf{Y} - \mathbb{E}[\Lambda] \cdot \mathbf{I}_N) \cdot f_1(\mathbf{a}), (\mathbf{Y}^2 - \mathbb{E}[\Lambda^2] \cdot \mathbf{I}_N) \cdot f_1(\mathbf{a}), \dots, (\mathbf{Y}^D - \mathbb{E}[\Lambda^D] \cdot \mathbf{I}_N) \cdot f_1(\mathbf{a}). \quad (44)$$

Since continuous functions can be approximated arbitrarily well by polynomial functions on compact sets, for any continuous function  $\Phi: \mathbb{R} \mapsto \mathbb{R}$ , we can approximate the matrix-vector product  $\Phi(\mathbf{Y}) \cdot f_1(\mathbf{a})$  using a linear combination of the matrix-vector products in (44). In particular, the missing vector-matrix product in (43) needed to reconstruct the second iteration of the target algorithm (42) can also be approximated using the matrix-vector products (44) computed by the lifted OAMP algorithm. The following definition formally introduces the class of lifted OAMP algorithms.

**Definition 6** (Degree- $D$  Lifted OAMP Algorithm). For any  $D \in \mathbb{N}$  (independent of  $N$ ), a degree- $D$  lifted OAMP algorithm maintains a sequence of iterates  $(\mathbf{w}_{ti})_{t \in \mathbb{N}, i \in [D]}$  indexed by a time index  $t \in \mathbb{N}$  and a degree index  $i \in [D]$ . At step  $t$ , the lifted OAMP computes the iterates  $\mathbf{w}_{t,1}, \dots, \mathbf{w}_{t,D}$  as follows:

$$\mathbf{w}_{t,i} = (\mathbf{Y}^i - \mathbb{E}[\Lambda^i] \cdot \mathbf{I}_N) \cdot f_t(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t-1,\bullet}; \mathbf{a}) \quad \forall i \in [D], t \in \mathbb{N}. \quad (45)$$

The estimate of  $\mathbf{x}_*$  at iteration  $t$  is obtained by applying a post-processing function  $\psi_t$  to the iterates  $\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t,\bullet}$  generated so far and the side-information  $\mathbf{a}$ :

$$\widehat{\mathbf{w}}_t = \psi_t(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t-1,\bullet}; \mathbf{a}), \quad (46)$$

In the above equations, the notation  $\mathbf{w}_{t,\bullet}$  is a shorthand for the collection of vectors  $\mathbf{w}_{t,1}, \mathbf{w}_{t,2}, \dots, \mathbf{w}_{t,D}$ . Furthermore for each  $t \in \mathbb{N}$ , the iterate denoiser  $f_t : \mathbb{R}^{(t-1)D} \times \mathbb{R}^k \mapsto \mathbb{R}$  and the postprocessing function  $\psi_t : \mathbb{R}^{tD} \times \mathbb{R}^k \mapsto \mathbb{R}$  are continuously differentiable, Lipschitz, act entry-wise on their vector inputs, and do not change with  $N$ . Each lifted OAMP algorithm is associated with a collection of *state evolution random variables*  $(\mathbf{X}_*, (\mathbf{W}_{t,i})_{t \in \mathbb{N}, i \in [D]}; \mathbf{A})$ , whose distribution is given by:

$$(\mathbf{X}_*, \mathbf{A}) \sim \pi, \quad \mathbf{W}_{t,i} = \alpha_t q_i \cdot \mathbf{X}_* + \mathbf{Z}_{ti} \quad \forall t \in \mathbb{N}, i \in [D], \quad (47a)$$

where  $\{\mathbf{Z}_{ti}\}_{t \in \mathbb{N}, i \in \mathbb{N}}$  are centered Gaussian random variables, sampled independently of  $(\mathbf{X}_*, \mathbf{A})$  with covariance:

$$\mathbb{E}[\mathbf{Z}_{si} \mathbf{Z}_{tj}] = \alpha_s \alpha_t \cdot \text{Cov}_{\Lambda_\nu \sim \nu}[\Lambda_\nu^i, \Lambda_\nu^j] + (\Sigma_{st} - \alpha_s \alpha_t) \cdot \text{Cov}_{\Lambda \sim \mu}[\Lambda^i, \Lambda^j]. \quad (47b)$$

In the above equations,  $(\alpha_t)_{t \in \mathbb{N}}$  and  $(\Sigma_{st})_{s \in \mathbb{N}, t \in \mathbb{N}}$  are given by the recursion:

$$\alpha_t = \mathbb{E}[\mathbf{X}_* f_t(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t-1,\bullet}; \mathbf{A})], \quad \Sigma_{st} = \mathbb{E}[f_s(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{s-1,\bullet}; \mathbf{A}) f_t(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t-1,\bullet}; \mathbf{A})]. \quad (47c)$$

Finally, we require the iterate denoisers  $(f_t)_{t \in \mathbb{N}}$  to satisfy the *divergence-free* constraint:

$$\mathbb{E}[\partial_{s,i} f_t(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t-1,\bullet}; \mathbf{A})] = 0 \quad \forall s \in [t-1], i \in [D], t \in \mathbb{N}, \quad (48)$$

where  $\partial_{s,i} f_t(w_{1,\bullet}, \dots, w_{s,\bullet}, \dots, w_{t-1,\bullet}; \mathbf{a})$  denotes the partial derivative of  $f_t(w_{1,\bullet}, \dots, w_{s,\bullet}, \dots, w_{t-1,\bullet}; \mathbf{a})$  with respect to  $w_{s,i}$ .

A degree- $D$  lifted OAMP algorithm which runs for  $t$  iterations can be viewed as an instance of an OAMP algorithm (Definition 4) which runs for  $tD$  iterations. Hence, the following result as an immediate corollary of Theorem 1.

**Corollary 1** (State evolution of lifted OAMP algorithms). *Consider a general degree- $D$  lifted OAMP algorithm of the form (45) and let  $(\mathbf{X}_*, (\mathbf{W}_{t,i})_{t \in \mathbb{N}, i \in [D]}; \mathbf{A})$  be the associated state evolution random variables. Then for any  $t \in \mathbb{N}$ ,*

$$(\mathbf{x}_*, \mathbf{w}_{1,\bullet}, \mathbf{w}_{2,\bullet}, \dots, \mathbf{w}_{t,\bullet}; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A}).$$

The following proposition (proved in Appendix E.1), shows that any estimator that is computed by an iterative algorithm (Definition 5) can be approximated arbitrarily well by a suitably designed degree- $D$  lifted OAMP algorithm in the limit  $D \rightarrow \infty$ .

**Proposition 3.** *Let  $\widehat{\mathbf{r}}_t$  be the estimator returned by any iterative algorithm (Definition 5) after  $t \in \mathbb{N}$  iterations. Then, for each  $D \in \mathbb{N}$ , there is a degree- $D$  lifted OAMP algorithm which returns an estimator  $\widetilde{\mathbf{w}}_t^{(D)}$  after  $t$  iterations, which satisfies:*

$$\lim_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{r}}_t - \widetilde{\mathbf{w}}_t^{(D)}\|^2}{N} = 0.$$

**The Optimal Lifted OAMP Algorithm.** Since Proposition 3 guarantees that any estimator computed using an iterative algorithm can be approximated by a suitably designed lifted OAMP algorithm, we focus on characterizing the optimal lifted OAMP algorithms. The optimal degree- $D$  lifted OAMP algorithm takes the form:

$$\mathbf{w}_{t,i} = (\mathbf{Y}^i - \mathbb{E}[\Lambda^i] \cdot \mathbf{I}_N) \cdot f_t^*(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t-1,\bullet}; \mathbf{a}) \quad \forall i \in [D], t \in \mathbb{N}, \quad (49a)$$

and returns the following estimator after  $t \in \mathbb{N}$  iterations:

$$\widehat{\mathbf{w}}_t^{(D)} = h_t^*(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t-1,\bullet}; \mathbf{a}). \quad (49b)$$

The description of the iterate denoisers  $(f_t^*)_{t \in \mathbb{N}}$  and the postprocessing functions  $(h_t^*)_{t \in \mathbb{N}}$  used by the optimal lifted OAMP algorithm is recursive. Suppose that the functions  $f_1^*, \dots, f_t^*$  have been specified. Let  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  denote the state evolution random variables corresponding to the first  $t$  iterations of the resulting lifted OAMP algorithm. Then:

1. The post-processing function  $h_t^*$  used at iteration  $t$  is the MMSE estimator (recall Definition 3) for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$ .
2.  $f_{t+1}^*$ , the iterate denoiser used by the lifted OAMP algorithm in the next iteration is the DMMSE estimator for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$ .

The following proposition (proved in Appendix E.2), shows that the estimator (49) achieves the lowest mean squared error among all estimators that can be computed using  $t$  iterations of a degree- $D$  lifted OAMP algorithm.

**Proposition 4.** Let  $\widehat{\mathbf{w}}_t^{(D)}$  be the estimator described in (49). Let  $\widetilde{\mathbf{w}}_t^{(D)}$  be any other estimator that can be computed using  $t$  iterations of some degree- $D$  lifted OAMP algorithm. Then, we have:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\|^2}{N} \leq \text{plim}_{N \rightarrow \infty} \frac{\|\widetilde{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\|^2}{N}.$$

*Proof Sketch.* Consider a general lifted OAMP algorithm:

$$\mathbf{w}_{t,i} = (\mathbf{Y}^i - \mathbb{E}[\Lambda^i] \cdot \mathbf{I}_N) \cdot f_t(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t-1,\bullet}; \mathbf{a}) \quad \forall i \in [D], t \in \mathbb{N},$$

which returns the estimator  $\widehat{\mathbf{w}}_t^{(D)} = h_t(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t,\bullet}; \mathbf{a})$  after  $t$  iterations. By Corollary 1, the asymptotic mean squared error for this estimator is  $\mathbb{E}[\|\mathbf{X}_* - h_t(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})\|^2]$ , where  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  denote the state evolution random variables associated with the lifted OAMP algorithm. Recalling the definition of the MMSE denoiser for a Gaussian channel from Definition 3, we find that the optimal choice for the postprocessing function is  $h_t = h_t^*$ , where  $h_t^*$  is the MMSE estimator for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$ . With this optimal choice of the post-processing function, the limiting MSE of the resulting estimator is  $\text{MMSE}(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$ . This depends implicitly on the iterate denoisers  $f_{1:t}$  used in the lifted OAMP algorithm since  $f_{1:t}$  determine the joint distribution of the state evolution random variables. To make this dependence explicit, we define  $\mathcal{M}_t(f_1, \dots, f_t) \stackrel{\text{def}}{=} \text{MMSE}(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  for any  $t \in \mathbb{N}$ . Deriving the optimal choice of the iterate denoisers is equivalent to minimizing the functional  $\mathcal{M}_t(f_1, \dots, f_t)$  with respect to  $f_{1:t}$ . The crux of the argument is to show that the iterate denoisers  $(f_t^*)_{t \in \mathbb{N}}$  used by the optimal lifted OAMP algorithm (49) correspond to a *greedy approach* to minimizing  $\mathcal{M}_t$ :

$$f_1^* \in \arg \min_{f_1} \mathcal{M}_1(f_1), \quad f_2^* \in \arg \min_{f_2} \mathcal{M}_2(f_1^*, f_2), \quad \dots \quad f_t^* \in \arg \min_{f_t} \mathcal{M}_t(f_1^*, \dots, f_{t-1}^*, f_t) \quad \forall t \in \mathbb{N}. \quad (50)$$

Moreover, we show that the functional  $\mathcal{M}_t$  has the property that for any choice of iterate denoisers  $f_{1:t-1}$ :

$$\min_{f_t} \mathcal{M}_t(f_1, \dots, f_t) = \mathcal{A}(\mathcal{M}_{t-1}(f_1, \dots, f_{t-1})), \quad (51)$$

for some explicit *non-decreasing* function  $\mathcal{A} : [0, 1] \mapsto [0, 1]$ . This property guarantees that the greedy approach (50) finds the global minimizer of  $\mathcal{M}_t$ . Indeed, if we assume as our induction hypothesis that  $\mathcal{M}_{t-1}(f_{1:t-1}^*) \leq \mathcal{M}_{t-1}(f_{1:t-1})$  for any choice of  $f_{1:t-1}$ , we can conclude that:

$$\mathcal{M}_t(f_{1:t}) \geq \min_{f_t} \mathcal{M}_t(f_{1:t}) \stackrel{(51)}{=} \mathcal{A}(\mathcal{M}_{t-1}(f_{1:t-1})) \stackrel{(a)}{\geq} \mathcal{A}(\mathcal{M}_{t-1}(f_{1:t-1}^*)) \stackrel{(51)}{=} \min_{f_t} \mathcal{M}_t(f_{1:t-1}^*, f_t) \stackrel{(50)}{=} \mathcal{M}_t(f_{1:t}^*),$$

where the inequality (a) follows from the induction hypothesis and the monotonicity of  $\mathcal{A}$ .  $\square$

**Simplifying the optimal lifted OAMP algorithm.** Observe that the optimal degree- $D$  lifted OAMP algorithm introduced in (49) is not an iterative algorithm in the sense of Definition 5. Indeed, an iterative algorithm is allowed a single matrix-vector multiplication per iteration, whereas the optimal lifted OAMP computes  $D$  matrix-vector products per iteration. Fortunately, the optimal lifted OAMP algorithm (49) can be simplified to eliminate these extra matrix-vector multiplications. This simplification leads to the following result, which shows that the optimal OAMP algorithm introduced in (17) can match the infinite degree ( $D \rightarrow \infty$ ) limit of the performance of the optimal lifted OAMP algorithm in (49).

**Proposition 5.** *Let  $\hat{\mathbf{x}}_t$  be the estimator returned by the optimal OAMP algorithm in (17) after  $t$  iterations. Let  $\hat{\mathbf{w}}_t^{(D)}$  be the estimator returned by optimal degree- $D$  lifted OAMP algorithm in (49) after  $t$  iterations. Then, for any  $t \in \mathbb{N}$ :*

$$\text{plim}_{N \rightarrow \infty} \frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_*\|^2}{N} = \lim_{D \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \frac{\|\hat{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\|^2}{N}.$$

*Proof Sketch.* We show that the estimator  $\hat{\mathbf{w}}_t^{(D)}$  returned by the optimal lifted OAMP algorithm (49) after  $t$  iterations depends only on a linear combination  $\mathbf{x}_t^{(D)} = \sum_{i=1}^D v_i \cdot \mathbf{w}_{t,i}$  of the iterates  $\mathbf{w}_{t,\bullet}$  generated at the last step of the algorithm. Thanks to this property, the  $\mathbf{x}_t^{(D)}$  (and hence, the estimator  $\hat{\mathbf{w}}_t$ ) can be computed using a single (instead of  $D$ ) matrix-vector multiplication(s):

$$\mathbf{x}_t^{(D)} = \sum_{i=1}^D v_i \cdot \mathbf{w}_{t,i}^{(D)} \stackrel{(49)}{=} \underbrace{\left( \sum_{i=1}^D v_i \cdot (\mathbf{Y}^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i] \cdot \mathbf{I}_N) \right)}_{\stackrel{\text{def}}{=} \Psi_t^{(D)}(\mathbf{Y})} \cdot f_t^*(\mathbf{w}_{<t,\bullet}; \mathbf{a}) = \Psi_t^{(D)}(\mathbf{Y}) \cdot f_t^*(\mathbf{w}_{<t,\bullet}; \mathbf{a}).$$

We argue that as  $D \rightarrow \infty$  the matrix denoiser  $\Psi_t^{(D)}$  that appears in the above equation converges to the optimal matrix denoiser used by the optimal OAMP algorithm (17) in an  $L^2$  sense, and hence the claimed result follows.  $\square$

Theorem 2 is an immediate consequence of the three intermediate results introduced above.

*Proof of Theorem 2.* Let  $\hat{\mathbf{x}}_t$  be the estimator returned by the optimal OAMP algorithm from (17) and let  $\hat{\mathbf{r}}_t$  be any estimator that can be computed using  $t$  iterations of some iterative algorithm. Proposition 3 guarantees the existence of a sequence of estimators  $(\tilde{\mathbf{w}}_t^{(D)})_{D \in \mathbb{N}}$  that approximate  $\hat{\mathbf{r}}_t$  and can be computed using  $t$  iterations of a degree- $D$  lifted OAMP algorithm. Finally, let  $\hat{\mathbf{w}}_t^{(D)}$  be the estimator returned by the optimal degree- $D$  lifted OAMP algorithm from (49). By the reverse triangle inequality, for any  $D \in \mathbb{N}$ ,  $\|\hat{\mathbf{r}}_t - \mathbf{x}_*\| \geq \|\tilde{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\| - \|\hat{\mathbf{w}}_t^{(D)} - \hat{\mathbf{r}}_t\|$ . We let  $N \rightarrow \infty$  and then  $D \rightarrow \infty$  to obtain:

$$\begin{aligned} \text{plim}_{N \rightarrow \infty} \inf \frac{\|\hat{\mathbf{r}}_t - \mathbf{x}_*\|}{\sqrt{N}} &\geq \liminf_{D \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \inf \frac{\|\tilde{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\|}{\sqrt{N}} - \limsup_{D \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \sup \frac{\|\hat{\mathbf{w}}_t^{(D)} - \hat{\mathbf{r}}_t\|}{\sqrt{N}} \\ &\stackrel{\text{Prop. 3}}{=} \liminf_{D \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \inf \frac{\|\hat{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\|}{\sqrt{N}} \stackrel{\text{Prop. 4}}{\geq} \liminf_{D \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \inf \frac{\|\hat{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\|}{\sqrt{N}} \stackrel{\text{Prop. 5}}{=} \text{plim}_{N \rightarrow \infty} \inf \frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_*\|}{\sqrt{N}}. \end{aligned}$$

Comparing the first and final expressions in the above chain of inequalities gives us the claim of Theorem 2.  $\square$

## 5 Numerical Experiments

We conclude the paper with some numerical experiments to demonstrate the performance of the optimal OAMP (17) algorithm. The codes to reproduce the results in this section are publicly available [1].

**Synthetic Noise Matrices.** Following [50, 69], we generate the noise matrix  $\mathbf{W}$  in a structured manner, making it feasible to efficiently simulate the dynamics of iterative algorithms on very high-dimensional matrices ( $N = 2 \times 10^6$ ). Specifically, we use the noise matrix  $\mathbf{W} = \mathbf{O} \text{diag}(\lambda_1, \dots, \lambda_N) \mathbf{O}^\top$  with  $\mathbf{O} = \mathbf{S}_1 \mathbf{F} \mathbf{S}_2 \mathbf{F}^\top \mathbf{S}_3$ , where  $\mathbf{S}_1, \mathbf{S}_2, \mathbf{S}_3$  are independent random signed diagonal matrices, and  $\mathbf{F}$  is the discrete cosine transform matrix. The eigenvalues  $\{\lambda_i\}_{i \in [N]}$  are sampled independently from the spectral  $\mu$  corresponding to the quartic noise model (21) with parameter  $\gamma = 0$ . Known universality results [24, 26, 75] guarantee that using this structured noise matrix (instead of random rotationally invariant noise drawn from the quartic model) does not change the asymptotic dynamics of AMP algorithms.

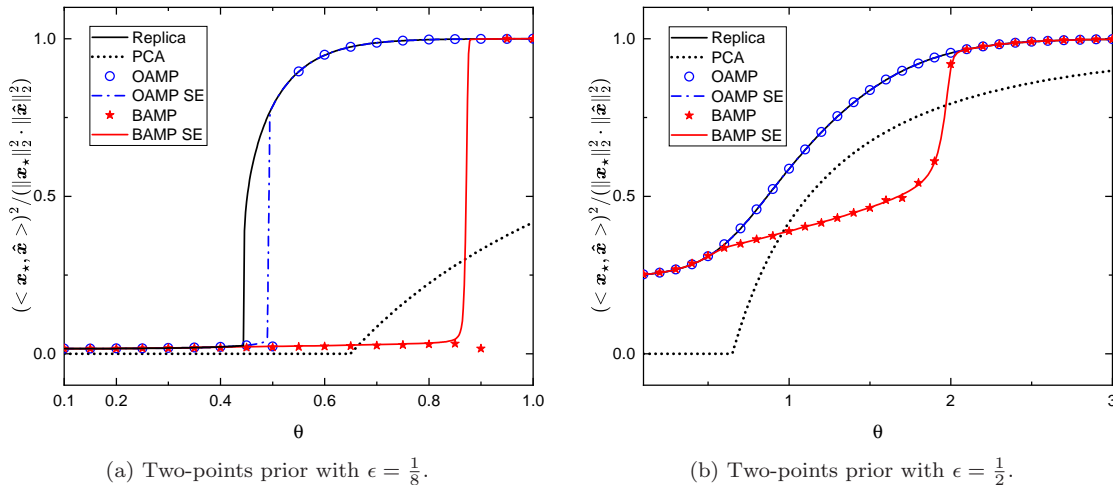


Figure 2: Performance (as measured by the normalized overlap with signal) of the PCA estimator, the Bayes-optimal estimator, the optimal OAMP algorithm (15 iterations), and the BAMP algorithm of Barbier et al. [8] (15 iterations) as a function of the SNR  $\theta$ .

Fig. 2 compares the performance of optimal OAMP with existing methods, including vanilla PCA and the BAMP algorithm introduced by Barbier et al. [8], along with conjectured replica formula (Conjecture 1) for the Bayes-optimal estimator (marked “Replica”). We sampled the signal  $\mathbf{x}_*$  from an i.i.d. two-point prior  $\epsilon^2 \delta_{1/\epsilon} + (1 - \epsilon^2) \delta_0$  (no side information is provided). Fig. 2 shows that the performances of the optimal OAMP algorithm and the BAMP algorithm match well with their corresponding state evolution (SE) predictions. The optimal OAMP algorithm outperforms vanilla PCA and the BAMP algorithm in both cases. We also see a statistical-computational gap for this problem for small  $\epsilon$ . Namely, there is a performance gap between the conjectured Bayes-optimal performance and the best achievable performance by iterative algorithms (attained by the optimal OAMP algorithm). In contrast, no statistical-computational gap exists for the two-point prior with large enough  $\epsilon$ . A similar phenomenon occurs when the noise matrix is i.i.d. Gaussian [54].

The performance of BAMP depends on the quality of initialization. Fig. 3 compares the dynamics of the optimal OAMP algorithm and BAMP when initialized at  $\mathbf{x}_0 = \sqrt{\omega_0} \mathbf{x}_* + \sqrt{1 - \omega_0} \mathbf{z}$  with  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_N)$ . We use the same noise model as before and sample the signal from an i.i.d. two-point prior with  $\epsilon = 0.5$  and  $\theta = 1.8$ . For a good initialization (Fig. 3(a)), the optimal OAMP algorithm and BAMP achieve the same MSE on convergence. However, with uninformative initialization (Fig. 3(b)), the performance of BAMP degrades after a few iterations, and the optimal OAMP algorithm achieves a superior MSE on convergence. Notice that the performance of BAMP is well predicted by its asymptotic state evolution. Hence, the observed behavior of BAMP is not caused by finite- $N$  effects or numerical stability issues and seems to be an inherent limitation of the algorithm.

**Realistic Noise Matrices.** The theoretical prediction of the asymptotic performance of OAMP requires the noise matrix to be rotationally-invariant. The noise matrix used in the experiments for Fig. 2 and Fig. 3 is structured but still synthetic. We now consider a more realistic model where the noise matrix is obtained from real datasets. Taking inspiration from [77], we generate the noise matrix from the covariance matrix of

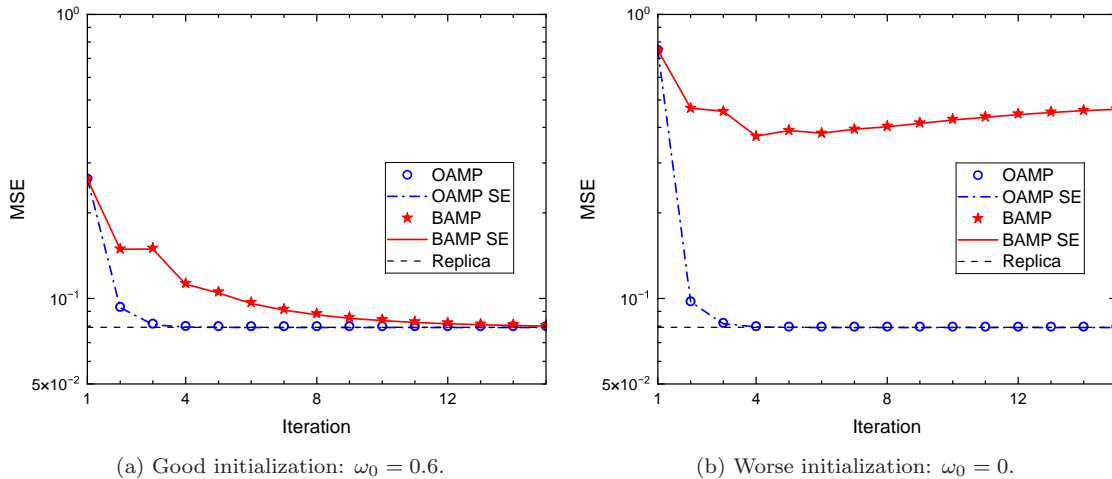


Figure 3: MSE performance of the optimal OAMP algorithm and the BAMP algorithm of Barbier et al. [8] as a function of the iteration number compared with the replica conjecture for the Bayes-optimal performance.

real datasets in bioinformatics, namely, the 1000 Genomes Project (1000G) [18] and International HapMap Project (Hapmap3) [19]. Both datasets undergo the preprocessing steps outlined in [77], producing a dataset containing 100,000 common single nucleotide polymorphisms (SNPs) for 2504 individuals for 1000G and a dataset comprising 142,186 SNPs for 1397 individuals for Hapmap3. In our experiments, we randomly select 3000 SNPs and compute the corresponding covariance matrix, resulting in a  $2054 \times 2054$  symmetric matrix for 1000G and a  $1397 \times 1397$  symmetric matrix for Hapmap3. We then obtain the noise matrix by further removing the “ground truth” signal components. Following [77], we estimate the “ground truth” signal components by forming covariance matrices using all SNPs from the original datasets (instead of the sub-sampled ones). The noise matrix is obtained by the following pre-processing steps: (1) we remove the “ground truth” signal components; (2) we further extract the principal components that correspond to outlying eigenvalues; (3) We center and scale the data properly.

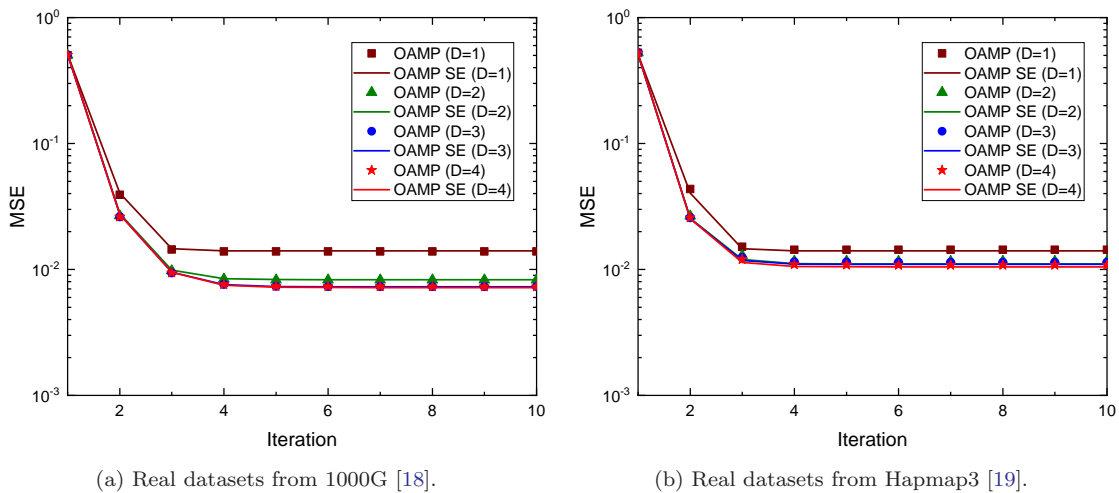


Figure 4: Performance of a data-driven implementation of the optimal degree- $D$  lifted OAMP algorithm from (49) on noise matrices derived from real datasets and signals randomly drawn from  $\mathbf{X}_* \sim 0.1\mathcal{N}(0, 10) + 0.9\delta_0$ .

Running the optimal OAMP algorithm requires estimating the spectral measure  $\mu$  and the associated optimal matrix denoiser (17e) from the observed data. While feasible, we found it more convenient to design a data-driven implementation of the optimal degree- $D$  lifted OAMP algorithm introduced in (49). Fig. 4 displays the MSE performance of this algorithm using the above-described realistic noise matrices and a randomly generated signal. The performance of the lifted OAMP improves as the degree  $D$  increases, saturating around  $D \approx 4$ , indicating that this algorithm can match the performance of the optimal OAMP algorithm in a data-driven way at the cost of 3 additional matrix multiplications per iteration. Fig. 4 also displays the asymptotic MSE predicted by the state evolution result for the corresponding rotationally invariant noise model with a matching spectrum. The performance of the lifted OAMP algorithm closely aligns with the state evolution prediction, suggesting an underlying universality phenomenon. We have observed that the pre-processing steps used are critical for universality. Understanding the precise conditions for universality (or its absence) is an interesting future direction.

## References

- [1] Accompanying code for optimality of message passing algorithms for spiked matrix models with rotationally invariant noise. <https://github.com/songIce/OAMP>.
- [2] Emmanuel Abbe. Community detection and stochastic block models: recent developments. *Journal of Machine Learning Research*, 18(177):1–86, 2018.
- [3] Greg W Anderson, Alice Guionnet, and Ofer Zeitouni. *An introduction to random matrices*. Cambridge University Press, 2010.
- [4] Jinho Baik and Jack W Silverstein. Eigenvalues of large sample covariance matrices of spiked population models. *Journal of Multivariate Analysis*, 97(6):1382–1408, 2006.
- [5] Jinho Baik, Gérard Ben Arous, and Sandrine Péché. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *Annals of Probability*, 33(5):1643–1697, 2005.
- [6] Jean Barbier and Nicolas Macris. The adaptive interpolation method for proving replica formulas. applications to the Curie–Weiss and Wigner spike models. *Journal of Physics A: Mathematical and Theoretical*, 52(29):294002, 2019.
- [7] Jean Barbier, Mohamad Dia, Nicolas Macris, Florent Krzakala, Thibault Lesieur, and Lenka Zdeborová. Mutual information for symmetric rank-one matrix estimation: A proof of the replica formula. *Advances in Neural Information Processing Systems*, 29, 2016.
- [8] Jean Barbier, Francesco Camilli, Marco Mondelli, and Manuel Sáenz. Fundamental limits in structured principal component analysis and how to reach them. *Proceedings of the National Academy of Sciences*, 120(30):e2302028120, 2023.
- [9] Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2):764–785, 2011.
- [10] Joshua K Behne and Galen Reeves. Fundamental limits for rank-one matrix estimation with groupwise heteroskedasticity. In *International Conference on Artificial Intelligence and Statistics*, pages 8650–8672. PMLR, 2022.
- [11] Serban Teodor Belinschi. The Lebesgue decomposition of the free additive convolution of two probability distributions. *Probability Theory and Related Fields*, 142(1):125–150, 2008.
- [12] Florent Benaych-Georges and Raj Rao Nadakuditi. The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, 2011.
- [13] Erwin Bolthausen. An iterative construction of solutions of the TAP equations for the Sherrington–Kirkpatrick model. *Communications in Mathematical Physics*, 325(1):333–366, 2014.

- [14] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013.
- [15] Joël Bun, Romain Allez, Jean-Philippe Bouchaud, and Marc Potters. Rotational invariant estimator for general noisy matrices. *IEEE Transactions on Information Theory*, 62(12):7475–7490, 2016.
- [16] Michael Celentano, Andrea Montanari, and Yuchen Wu. The estimation error of general first order methods. In *Conference on Learning Theory*, pages 1078–1141. PMLR, 2020.
- [17] Michael Celentano, Zhou Fan, and Song Mei. Local convexity of the TAP free energy and AMP convergence for  $\mathbb{Z}_2$ -synchronization. *The Annals of Statistics*, 51(2):519–546, 2023.
- [18] 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*, 526(7571):68, 2015.
- [19] International HapMap 3 Consortium. Integrating common and rare genetic variation in diverse human populations. *Nature*, 467(7311):52, 2010.
- [20] Yash Deshpande and Andrea Montanari. Information-theoretically optimal sparse PCA. In *2014 IEEE International Symposium on Information Theory*, pages 2197–2201. IEEE, 2014.
- [21] Yash Deshpande, Emmanuel Abbe, and Andrea Montanari. Asymptotic mutual information for the balanced binary stochastic block model. *Information and Inference: A Journal of the IMA*, 6(2):125–170, 2017.
- [22] David Donoho and Jared Tanner. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 367(1906):4273–4293, 2009.
- [23] David L Donoho, Arian Maleki, and Andrea Montanari. Message-passing algorithms for compressed sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919, 2009.
- [24] Rishabh Dudeja and Milad Bakhshizadeh. Universality of linearized message passing for phase retrieval with structured sensing matrices. *IEEE Transactions on Information Theory*, 68(11):7545–7574, 2022.
- [25] Rishabh Dudeja, Subhabrata Sen, and Yue M Lu. Spectral universality of regularized linear regression with nearly deterministic sensing matrices. *arXiv preprint arXiv:2208.02753*, 2022.
- [26] Rishabh Dudeja, Yue M. Lu, and Subhabrata Sen. Universality of approximate message passing with semirandom matrices. *The Annals of Probability*, 51(5):1616–1683, 2023.
- [27] Ahmed El Alaoui and Florent Krzakala. Estimation in the spiked Wigner model: a short proof of the replica formula. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pages 1874–1878. IEEE, 2018.
- [28] Ahmed El Alaoui, Florent Krzakala, and Michael Jordan. Fundamental limits of detection in the spiked Wigner model. *The Annals of statistics*, 48(2):863–885, 2020.
- [29] Zhou Fan. Approximate message passing algorithms for rotationally invariant matrices. *The Annals of Statistics*, 50(1):197–224, 2022.
- [30] Zhou Fan, Song Mei, and Andrea Montanari. TAP free energy, spin glasses and variational inference. *The Annals of Probability*, 49(1), 2021.
- [31] Oliver Y Feng, Ramji Venkataramanan, Cynthia Rush, and Richard J Samworth. A unifying tutorial on approximate message passing. *Foundations and Trends® in Machine Learning*, 15(4):335–536, 2022.
- [32] Delphine Féral and Sandrine Péché. The largest eigenvalue of rank one deformation of large Wigner matrices. *Communications in Mathematical Physics*, 272:185–228, 2007.



- [33] Dongning Guo, Yihong Wu, Shlomo S Shitz, and Sergio Verdú. Estimation in Gaussian noise: Properties of the minimum mean-square error. *IEEE Transactions on Information Theory*, 57(4):2371–2385, 2011.
- [34] Peter Henrici. *Applied and computational complex analysis, Volume 3: Discrete Fourier analysis, Cauchy integrals, construction of conformal maps, univalent functions*, volume 41. John Wiley & Sons, 1993.
- [35] Iain M Johnstone. On the distribution of the largest eigenvalue in principal components analysis. *The Annals of statistics*, 29(2):295–327, 2001.
- [36] Iain M Johnstone and Arthur Yu Lu. On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association*, 104(486):682–693, 2009.
- [37] Yoshiyuki Kabashima. A cdma multiuser detection algorithm on the basis of belief propagation. *Journal of Physics A: Mathematical and General*, 36(43):11111, 2003.
- [38] Yoshiyuki Kabashima, Florent Krzakala, Marc Mézard, Ayaka Sakata, and Lenka Zdeborová. Phase transitions and sample complexity in Bayes-optimal matrix factorization. *IEEE Transactions on information theory*, 62(7):4228–4265, 2016.
- [39] Antti Knowles and Jun Yin. The isotropic semicircle law and deformation of Wigner matrices. *Communications on Pure and Applied Mathematics*, 66(11):1663–1749, 2013.
- [40] Florent Krzakala, Jiaming Xu, and Lenka Zdeborová. Mutual information in rank-one matrix estimation. In *2016 IEEE Information Theory Workshop (ITW)*, pages 71–75. IEEE, 2016.
- [41] Olivier Ledoit and Sandrine Péché. Eigenvectors of some large sample covariance matrix ensembles. *Probability Theory and Related Fields*, 151(1):233–264, 2011.
- [42] Olivier Ledoit and Michael Wolf. Nonlinear shrinkage estimation of large-dimensional covariance matrices. *The Annals of Statistics*, 40(2):1024–1060, 2012.
- [43] Marc Lelarge and Léo Miolane. Fundamental limits of symmetric low-rank matrix estimation. In *Conference on Learning Theory*, pages 1297–1301. PMLR, 2017.
- [44] Thibault Lesieur, Florent Krzakala, and Lenka Zdeborová. Mmse of probabilistic low-rank matrix estimation: Universality with respect to the output channel. In *2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 680–687. IEEE, 2015.
- [45] Thibault Lesieur, Florent Krzakala, and Lenka Zdeborová. Phase transitions in sparse PCA. In *2015 IEEE International Symposium on Information Theory (ISIT)*, pages 1635–1639. IEEE, 2015.
- [46] Gen Li and Yuting Wei. A non-asymptotic framework for approximate message passing in spiked models. *arXiv preprint arXiv:2208.03313*, 2022.
- [47] Gen Li, Wei Fan, and Yuting Wei. Approximate message passing from random initialization with applications to z2 synchronization. *Proceedings of the National Academy of Sciences*, 120(31):e2302930120, 2023.
- [48] Yufan Li and Pragya Sur. Spectrum-aware adjustment: A new debiasing framework with applications to principal components regression. *arXiv preprint arXiv:2309.07810*, 2023.
- [49] Panagiotis Lolas and Lexing Ying. Shrinkage estimation of functions of large noisy symmetric matrices. *arXiv preprint arXiv:2106.05183*, 2021.
- [50] Junjie Ma and Li Ping. Orthogonal AMP. *IEEE Access*, 5:2020–2033, 2017.
- [51] James A Mingo and Roland Speicher. *Free probability and random matrices*, volume 35. Springer, 2017.
- [52] Léo Miolane. Fundamental limits of low-rank matrix estimation: the non-symmetric case. *arXiv preprint arXiv:1702.00473*, 2017.

- [53] Marco Mondelli and Ramji Venkataramanan. PCA initialization for approximate message passing in rotationally invariant models. *Advances in Neural Information Processing Systems*, 34:29616–29629, 2021.
- [54] Andrea Montanari and Ramji Venkataramanan. Estimation of low-rank matrices via approximate message passing. *The Annals of Statistics*, 49(1), 2021.
- [55] Andrea Montanari and Alexander S Wein. Equivalence of approximate message passing and low-degree polynomials in rank-one matrix estimation. *arXiv preprint arXiv:2212.06996*, 2022.
- [56] Andrea Montanari and Yuchen Wu. Statistically optimal first order algorithms: A proof via orthogonalization. *arXiv preprint arXiv:2201.05101*, 2022.
- [57] Manfred Opper and Ole Winther. Adaptive and self-averaging Thouless-Anderson-Palmer mean-field theory for probabilistic modeling. *Physical Review E*, 64(5):056131, 2001.
- [58] Manfred Opper, Burak Cakmak, and Ole Winther. A theory of solving TAP equations for Ising models with general invariant random matrices. *Journal of Physics A: Mathematical and Theoretical*, 49(11):114002, 2016.
- [59] Samet Oymak and Babak Hassibi. A case for orthogonal measurements in linear inverse problems. In *2014 IEEE International Symposium on Information Theory*, pages 3175–3179. IEEE, 2014.
- [60] Jason T Parker, Philip Schniter, and Volkan Cevher. Bilinear generalized approximate message passing—part i: Derivation. *IEEE Transactions on Signal Processing*, 62(22):5839–5853, 2014.
- [61] Leonid Andreevich Pastur and Mariya Shcherbina. *Eigenvalue distribution of large random matrices*. American Mathematical Soc., 2011.
- [62] Debashis Paul. Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statistica Sinica*, pages 1617–1642, 2007.
- [63] Sandrine Péché. The largest eigenvalue of small rank perturbations of Hermitian random matrices. *Probability Theory and Related Fields*, 134:127–173, 2006.
- [64] Amelia Perry, Alexander S Wein, Afonso S Bandeira, and Ankur Moitra. Message-passing algorithms for synchronization problems over compact groups. *Communications on Pure and Applied Mathematics*, 71(11):2275–2322, 2018.
- [65] Amelia Perry, Alexander S Wein, Afonso S Bandeira, and Ankur Moitra. Optimality and sub-optimality of PCA i: Spiked random matrix models. *The Annals of Statistics*, 46(5):2416–2451, 2018.
- [66] Marc Potters and Jean-Philippe Bouchaud. *A First Course in Random Matrix Theory: For Physicists, Engineers and Data Scientists*. Cambridge University Press, 2020.
- [67] Farzad Pourkamali and Nicolas Macris. Rectangular rotational invariant estimator for general additive noise matrices. In *2023 IEEE International Symposium on Information Theory (ISIT)*, pages 2081–2086. IEEE, 2023.
- [68] Sundeep Rangan and Alyson K Fletcher. Iterative estimation of constrained rank-one matrices in noise. In *2012 IEEE International Symposium on Information Theory Proceedings*, pages 1246–1250. IEEE, 2012.
- [69] Sundeep Rangan, Philip Schniter, and Alyson K Fletcher. Vector approximate message passing. *IEEE Transactions on Information Theory*, 65(10):6664–6684, 2019.
- [70] Konrad Schmüdgen. *The moment problem*, volume 9. Springer, 2017.
- [71] Tselil Schramm and Alexander S Wein. Computational barriers to estimation from low-degree polynomials. *The Annals of Statistics*, 50(3):1833–1858, 2022.

- [72] Guilhem Semerjian. Matrix denoising: Bayes-optimal estimators via low-degree polynomials. *arXiv preprint arXiv:2402.16719*, 2024.
- [73] Amit Singer. Angular synchronization by eigenvectors and semidefinite programming. *Applied and computational harmonic analysis*, 30(1):20–36, 2011.
- [74] Keigo Takeuchi. Rigorous dynamics of expectation-propagation-based signal recovery from unitarily invariant measurements. *IEEE Transactions on Information Theory*, 66(1):368–386, 2019.
- [75] Tianhao Wang, Xinyi Zhong, and Zhou Fan. Universality of approximate message passing algorithms and tensor networks. *arXiv preprint arXiv:2206.13037*, 2022.
- [76] Xinyi Zhong, Tianhao Wang, and Zhou Fan. Approximate message passing for orthogonally invariant ensembles: Multivariate non-linearities and spectral initialization. *arXiv preprint arXiv:2110.02318*, 2021.
- [77] Xinyi Zhong, Chang Su, and Zhou Fan. Empirical Bayes PCA in high dimensions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(3):853–878, 2022.
- [78] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of computational and graphical statistics*, 15(2):265–286, 2006.

## A Proofs for Preliminary Results

### A.1 Proofs for Random Matrix Theory Results

This section is devoted to the proof of Lemma 1.

*Proof of Lemma 1.* We consider each of the claims made in the lemma one by one.

**Weak convergence of  $\nu_N$  (Claim (1)).** Let  $\tilde{\nu}_N$  be a slightly modified version of  $\nu_N$ :

$$\tilde{\nu}_N = \frac{1}{N} \sum_{i=1}^N (\tilde{\mathbf{x}}_*^\top \mathbf{u}_i)^2 \cdot \delta_{\lambda_i(\mathbf{Y})}$$

where  $\tilde{\mathbf{x}}_* \stackrel{\text{def}}{=} \frac{\sqrt{N}\mathbf{x}_*}{\|\mathbf{x}_*\|}$ . This modification makes  $\tilde{\nu}_N$  a probability measure, i.e.,  $\tilde{\nu}_N(\mathbb{R}) = 1$ . Notice that for any bounded, continuous function  $f : \mathbb{R} \mapsto \mathbb{R}$ , we have:

$$\int_{\mathbb{R}} f(\lambda) \tilde{\nu}_N(\lambda) = \frac{N}{\|\mathbf{x}_*\|^2} \cdot \int_{\mathbb{R}} f(\lambda) \nu_N(\lambda).$$

By Assumption 1,  $\|\mathbf{x}_*\|^2/N \xrightarrow{\mathbb{P}} 1$ . Hence, it suffices to show that  $\tilde{\nu}_N$  converges weakly to a probability measure  $\nu$  in probability. To do so, it suffices to verify that the Stieltjes transform  $\mathcal{S}_{\tilde{\nu}_N}(z)$  of  $\tilde{\nu}_N$  converges pointwise to the Stieltjes transform of a probability measure  $\nu$  for any  $z \in \mathbb{C} \setminus \mathbb{R}$  (see for e.g., [3, Theorem 2.4.4]). To this end, we note that:

$$\mathcal{S}_{\tilde{\nu}_N}(z) = \frac{1}{\|\mathbf{x}_*\|^2} \sum_{i=1}^N \frac{\langle \mathbf{u}_i(\mathbf{Y}), \mathbf{x}_* \rangle^2}{z - \lambda_i(\mathbf{Y})} = \frac{1}{\|\mathbf{x}_*\|^2} \cdot \mathbf{x}_*^\top (z\mathbf{I}_N - \mathbf{Y})^{-1} \mathbf{x}_*,$$

where  $\lambda_{1:N}(\mathbf{Y})$  and  $\mathbf{u}_{1:N}(\mathbf{Y})$  denote the eigenvalues and eigenvectors of  $\mathbf{Y}$ . Hence,

$$\mathcal{S}_{\tilde{\nu}_N}(z) = \frac{1}{\|\mathbf{x}_*\|^2} \mathbf{x}_*^\top (z\mathbf{I}_N - \mathbf{Y})^{-1} \mathbf{x}_* = \frac{1}{\|\mathbf{x}_*\|^2} \mathbf{x}_*^\top \left( z\mathbf{I}_N - \frac{\theta \mathbf{x}_* \mathbf{x}_*^\top}{N} - \mathbf{W} \right)^{-1} \mathbf{x}_* \stackrel{(a)}{=} \frac{\frac{1}{\|\mathbf{x}_*\|^2} \mathbf{x}_*^\top (z\mathbf{I}_N - \mathbf{W})^{-1} \mathbf{x}_*}{1 - \frac{\theta}{N} \mathbf{x}_*^\top (z\mathbf{I}_N - \mathbf{W})^{-1} \mathbf{x}_*}, \quad (52)$$

where step (a) uses the Sherman-Morrison formula for the rank-one update to the matrix inverse. Using standard concentration results regarding quadratic forms of rotationally invariant matrices (stated as Fact 2 in Appendix F for convenience), we find that:

$$\frac{\mathbf{x}_*^\top (z\mathbf{I}_N - \mathbf{W})^{-1} \mathbf{x}_*}{\|\mathbf{x}_*\|^2} \xrightarrow{\mathbb{P}} \int_{\mathbb{R}} \frac{\mu(d\lambda)}{z - \lambda} = \mathcal{S}_\mu(z) \quad \forall z \in \mathbb{C} \setminus \mathbb{R}. \quad (53)$$

By Assumption 1,  $\|\mathbf{x}_*\|^2/N \xrightarrow{\mathbb{P}} 1$ . Hence,

$$1 - \frac{\theta}{N} \mathbf{x}_*^\top (z\mathbf{I}_N - \mathbf{W})^{-1} \mathbf{x}_* = 1 - \theta \cdot \frac{\|\mathbf{x}_*\|^2}{N} \cdot \frac{\mathbf{x}_*^\top (z\mathbf{I}_N - \mathbf{W})^{-1} \mathbf{x}_*}{\|\mathbf{x}_*\|^2} \xrightarrow{\mathbb{P}} 1 - \theta \mathcal{S}_\mu(z) \stackrel{(a)}{\neq} 0 \quad \forall z \in \mathbb{C} \setminus \mathbb{R}, \quad (54)$$

where the non-equality marked (a) follows by observing that  $\Im(\mathcal{S}_\mu(z)) \neq 0$  whenever  $\text{Im}(z) \neq 0$ , and so  $1 - \theta \mathcal{S}_\mu(z) \neq 0$ . Plugging (53) and (54) into (52), we obtain:

$$\mathcal{S}_{\tilde{\nu}_N}(z) \xrightarrow{\mathbb{P}} \frac{\mathcal{S}_\mu(z)}{1 - \theta \mathcal{S}_\mu(z)} \quad \forall z \in \mathbb{C} \setminus \mathbb{R}. \quad (55)$$

To prove the weak convergence of  $\tilde{\nu}_N$  to a *probability measure*  $\nu$ , we need to show additionally that the RHS of the above equation is the Stieltjes transform of some probability measure. By [51, Theorem 10, Chapter 3], it is sufficient to verify that:

$$\limsup_{y \rightarrow \infty} y \left| \frac{\mathcal{S}_\mu(iy)}{1 - \theta \mathcal{S}_\mu(iy)} \right| = 1. \quad (56)$$

Since  $\mu$  is a probability measure, we have  $y|\mathcal{S}_\mu(iy)| \rightarrow 1$  and  $\mathcal{S}_\mu(iy) \rightarrow 0$  as  $y \rightarrow \infty$  (see [51, Lemma 3, Chapter 3]). Hence, (56) holds. Consequently, the RHS of (55) is a Stieltjes transform of a probability measure  $\nu$  and  $\tilde{\nu}_N$  and  $\nu_N$  converge weakly to  $\nu$  in probability.

**Lebesgue Decomposition of  $\nu$  (Claims (2) and (3)).** Belinschi [11, Lemma 2.17] has shown that for any probability measure  $\chi$  on  $\mathbb{R}$  with Lebesgue decomposition  $\chi = \chi_{\parallel} + \chi_{\perp}$ , the Stieltjes transform of  $\nu$  satisfies:

$$\lim_{\epsilon \downarrow 0} \mathcal{S}_\chi(\lambda + i\epsilon) = -\infty \quad \text{for } \chi_{\perp}\text{-almost every } \lambda \in \mathbb{R}, \quad (57)$$

$$\lim_{\epsilon \downarrow 0} \Im[\mathcal{S}_\chi(x + i\epsilon)] = -\pi \frac{d\chi_{\parallel}}{d\lambda}(\lambda) \quad \text{for Lebesgue-almost every } \lambda \in \mathbb{R}. \quad (58)$$

In the display above,  $\frac{d\chi_{\parallel}}{d\lambda}$  denotes the density of  $\chi_{\parallel}$  with respect to the Lebesgue measure and the limit on the LHS of (58) is guaranteed to exist and is finite for Lebesgue-almost every  $\lambda$ . In light of this result, we study the limit:

$$\lim_{\epsilon \downarrow 0} \Im[\mathcal{S}_\nu(\lambda + i\epsilon)]$$

for an arbitrary  $\lambda \in \mathbb{R}$ . For any  $\epsilon > 0$ , we have:

$$\Im[\mathcal{S}_\nu(\lambda + i\epsilon)] = \Im \left[ \frac{\mathcal{S}_\mu(\lambda + i\epsilon)}{1 - \theta \mathcal{S}_\mu(\lambda + i\epsilon)} \right] = \frac{\Im[\mathcal{S}_\mu(\lambda + i\epsilon)]}{(1 - \theta \Re[\mathcal{S}_\mu(\lambda + i\epsilon)])^2 + \theta^2 \Im^2[\mathcal{S}_\mu(\lambda + i\epsilon)]}.$$

Since  $\mu$  is absolutely continuous with respect to the Lebesgue measure (Assumption 2), (58) implies:

$$\Im[\mathcal{S}_\mu(\lambda + i\epsilon)] \rightarrow -\pi \mu(\lambda) \quad \text{as } \epsilon \downarrow 0.$$

where  $\mu(\cdot)$  denotes the density of  $\mu$ . Moreover, by the Sokhotsky-Plemelj formula (see for instance, [61, Section 2.1]),

$$\Re[\mathcal{S}_\mu(\lambda + i\epsilon)] \rightarrow \pi \mathcal{H}_\mu(\lambda) \quad \text{as } \epsilon \downarrow 0.$$

Recall that:

$$\phi(\lambda) \stackrel{\text{def}}{=} (1 - \theta\pi\mathcal{H}_\mu(\lambda))^2 + \theta^2\pi^2 \cdot \mu(\lambda).$$

Hence,

$$\begin{aligned} \lim_{\epsilon \downarrow 0} \Im[\mathcal{H}_\nu(\lambda + i\epsilon)] &= -\infty \quad \text{iff } \phi(\lambda) = 0, \\ \lim_{\epsilon \downarrow 0} \Im[\mathcal{H}_\nu(\lambda + i\epsilon)] &= -\frac{\pi\mu(\lambda)}{\phi(\lambda)} \in (-\infty, 0] \quad \text{iff } \phi(\lambda) \neq 0. \end{aligned}$$

Combining the above results with (57) and (58), we conclude that:

$$\begin{aligned} \phi(\lambda) &\neq 0 \quad \text{for Lebesgue-almost every } \lambda \in \mathbb{R}, \\ \phi(\lambda) &= 0 \quad \text{for } \nu_\perp\text{-almost every } \lambda \in \mathbb{R}, \\ \frac{d\nu_\parallel}{d\lambda}(\lambda) &= \frac{\mu(\lambda)}{\phi(\lambda)} \quad \text{for Lebesgue-almost every } \lambda \in \mathbb{R}, \end{aligned}$$

as claimed. This concludes the proof of the lemma.  $\square$

## A.2 Preliminaries on Gaussian Channels

We state and prove some important properties of Gaussian channels (Definition 3) in this appendix.

**The DMMSE Estimator for Scalar Gaussian Channels.** The following lemma, due to Ma and Ping [50] provides a formula for the DMMSE estimator for a scalar Gaussian channel with SNR  $\omega$ . We provide a proof of this result in Section A.2.1 for completeness.

**Lemma 2** (Ma and Ping [50]). *Let  $(\mathbf{X}_*, \mathbf{X}; \mathbf{A})$  be a scalar Gaussian channel with SNR  $\omega \in [0, 1]$ . Then,*

1. *The function  $\bar{\varphi}(\cdot|\omega) : \mathbb{R}^{1+k} \mapsto \mathbb{R}$ :*

$$\bar{\varphi}(x; a|\omega) \stackrel{\text{def}}{=} \begin{cases} \left(1 - \frac{\sqrt{\omega}}{\sqrt{1-\omega}} \cdot \mathbb{E}[\mathbf{Z}\varphi(\mathbf{X}; \mathbf{A}|\omega)]\right)^{-1} \cdot \left(\varphi(x; a|\omega) - \frac{\mathbb{E}[\mathbf{Z}\varphi(\mathbf{X}; \mathbf{A}|\omega)]}{\sqrt{1-\omega}} \cdot x\right) & : \omega < 1 \\ x & : \omega = 1 \end{cases} \quad \forall x \in \mathbb{R}, a \in \mathbb{R}^k. \quad (59)$$

*is the DMMSE estimator for  $(\mathbf{X}_*, \mathbf{X}; \mathbf{A})$ .*

2. *The DMMSE estimator satisfies the identities:*

$$\mathbb{E}[\mathbf{X}_* \bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega) \bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)^2] = 1 - \text{dmmse}_\pi(\omega).$$

3. *When  $\omega < 1$ , the  $\text{mmse}_\pi$  and  $\text{dmmse}_\pi$  functions are related by the identity:*

$$\frac{1}{\text{dmmse}_\pi(\omega)} = \frac{1}{\text{mmse}_\pi(\omega)} - \frac{\omega}{1-\omega}.$$

**MMSE and DMMSE Estimators for General Gaussian Channels.** Next, we provide formulas for the MMSE and DMMSE estimators for general (multivariate) Gaussian channels. Consider a Gaussian channel  $(\mathbf{X}_*, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$ :

$$(\mathbf{X}_*, \mathbf{A}) \sim \pi, \quad (Z_1, \dots, Z_t) \sim \mathcal{N}(0, \Sigma), \quad \mathbf{X}_i = \alpha_i \cdot \mathbf{X}_* + Z_i \quad \forall i \in [t],$$

where the Gaussian noise  $(Z_1, \dots, Z_t)$  is sampled independently of the signal and side information  $(\mathbf{X}_*; \mathbf{A})$ . Observe that for any vector  $v \in \mathbb{R}^t$  which satisfies:

$$\langle v, \alpha \rangle^2 + v^\top \Sigma v = 1, \quad (60)$$

$(\mathbf{X}_\star, \langle v, \mathbf{X}_{\leq t} \rangle; \mathbf{A})$  forms a scalar Gaussian channel with SNR  $\omega = \langle v, \alpha \rangle^2$ . Among all vectors  $v \in \mathbb{R}^t$  which satisfy (60), the maximum SNR for the corresponding scalar channel is achieved by the vector:

$$v_{\text{opt}}(\mathbf{X}_\star | \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \stackrel{\text{def}}{=} \begin{cases} (\alpha^\top \Sigma^\dagger \alpha + (\alpha^\top \Sigma^\dagger \alpha)^2)^{-\frac{1}{2}} \cdot \Sigma^\dagger \alpha & \text{if } \alpha \in \text{Range}(\Sigma), \\ \langle P_\perp[\alpha], \alpha \rangle^{-1} \cdot P_\perp[\alpha] & \text{if } \alpha \notin \text{Range}(\Sigma), \end{cases} \quad (61a)$$

where  $\Sigma^\dagger$  denotes the pseudo-inverse of  $\Sigma$ ,  $\text{Range}(\Sigma)$  is the range space of  $\Sigma$ , and  $P_\perp[\cdot]$  denotes the projector on the orthogonal complement of  $\text{Range}(\Sigma)$ . We will refer to  $v_{\text{opt}}(\mathbf{X}_\star | \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$  as the *optimal linear combination* for the Gaussian channel  $(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$ . The maximum SNR is given by:

$$\omega_{\text{eff}}(\mathbf{X}_\star | \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \stackrel{\text{def}}{=} \langle \alpha, v_{\text{opt}}(\mathbf{X}_\star | \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \rangle^2 = \begin{cases} (1 + \alpha^\top \Sigma^\dagger \alpha)^{-1} \cdot \alpha^\top \Sigma^\dagger \alpha & \text{if } \alpha \in \text{Range}(\Sigma), \\ 1 & \text{if } \alpha \notin \text{Range}(\Sigma), \end{cases} \quad (61b)$$

which we call the *effective SNR* of the Gaussian channel. The following lemma shows that the MMSE and DMMSE of the Gaussian channel  $(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$  is same as the MMSE and DMMSE of the scalar Gaussian channel  $(\mathbf{X}_\star, \langle \mathbf{X}_{\leq t}, v_{\text{opt}} \rangle; \mathbf{A})$ , which operates at the SNR  $\omega_{\text{eff}}$ .

**Lemma 3.** *Let  $(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$  be a Gaussian channel. Let  $v_{\text{opt}}$  and  $\omega_{\text{eff}}$  denote the optimal linear combination and effective SNR for the Gaussian channel, as defined in (61). Then:*

1.  $\text{MMSE}(\mathbf{X}_\star | \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) = \text{mmse}_\pi(\omega_{\text{eff}})$  and the function:

$$f_\star(x; a) \stackrel{\text{def}}{=} \varphi(\langle x, v_{\text{opt}} \rangle; a | \omega_{\text{eff}}) \quad \forall x \in \mathbb{R}^t, \quad a \in \mathbb{R}^k,$$

is the MMSE estimator:

$$f_\star \in \arg \min_{f \in L^2(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})\}^2].$$

2.  $\text{DMMSE}(\mathbf{X}_\star | \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) = \text{dmmse}_\pi(\omega_{\text{eff}})$  and the function:

$$\bar{f}_\star(x; a) \stackrel{\text{def}}{=} \bar{\varphi}(\langle x, v_{\text{opt}} \rangle; a | \omega_{\text{eff}}) \quad \forall x \in \mathbb{R}^t, \quad a \in \mathbb{R}^k,$$

is the DMMSE estimator:

$$\bar{f}_\star \in \arg \min_{f \in L^2(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})\}^2] \quad \text{subject to } \mathbb{E}[Z_i f(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})] = 0 \quad \forall i \in [t].$$

*Proof.* The proof is deferred to Section A.2.2. □

**Monotonicity of MMSE and DMMSE.** Lastly, we will rely on the following natural monotonicity property of the  $\text{mmse}_\pi$  function.

*Fact 1* (Guo et al. 33, Proposition 9). The function  $\text{mmse}_\pi : [0, 1] \mapsto [0, 1]$  is non-increasing. If  $\text{mmse}_\pi(0) = \mathbb{E}[\text{Var}[\mathbf{X}_\star | \mathbf{A}]] \neq 0$ ,  $\text{mmse}_\pi : [0, 1] \mapsto [0, 1]$  is strictly decreasing.

The following lemma, whose proof is provided in Section A.2.3, shows that the  $\text{dmmse}_\pi$  function is also non-decreasing.

**Lemma 4.** *The function  $\text{dmmse}_\pi : [0, 1] \mapsto [0, 1]$  is non-increasing.*

The remainder of this section is devoted to the proofs of Lemma 2, Lemma 3, and Lemma 4 introduced above.

### A.2.1 Proof of Lemma 2

*Proof of Lemma 2.* This result is due to Ma and Ping [50], we provide a proof here for completeness. We consider a scalar Gaussian channel at SNR  $\omega$ :

$$(\mathbf{X}_*; \mathbf{A}) \sim \pi, \quad \mathbf{Z} \sim \mathcal{N}(0, 1), \quad \mathbf{X} = \sqrt{\omega} \cdot \mathbf{X}_* + \sqrt{1 - \omega} \cdot \mathbf{Z}, \quad (62)$$

where the Gaussian noise  $\mathbf{Z}$  is sampled independently of the signal and side information  $(\mathbf{X}_*; \mathbf{A})$ . Recall from Definition 3 that:

$$\text{dmmse}_\pi(\omega) \stackrel{\text{def}}{=} \text{DMMSE}(\mathbf{X}_* | \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \stackrel{\text{def}}{=} \min_{f \in L^2(\mathbf{X}; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_* - f(\mathbf{X}; \mathbf{A})\}^2] \quad \text{subject to } \mathbb{E}[\mathbf{Z}f(\mathbf{X}; \mathbf{A})] = 0. \quad (63)$$

Moreover, the function that achieves the lowest MSE in (63) is the DMMSE estimator for the channel. Let  $f$  be any function in  $L^2(\mathbf{X}; \mathbf{A})$  which satisfies the divergence-free requirement  $\mathbb{E}[\mathbf{Z}f(\mathbf{X}; \mathbf{A})] = 0$ . We can expand the mean squared error (MSE) of  $f$  as follows:

$$\mathbb{E}[(\mathbf{X}_* - f(\mathbf{X}; \mathbf{A}))^2] = 1 + \mathbb{E}[f^2(\mathbf{X}; \mathbf{A})] - 2\mathbb{E}[\mathbf{X}_* f(\mathbf{X}; \mathbf{A})] = 1 + \mathbb{E}[f^2(\mathbf{X}; \mathbf{A})] - 2\mathbb{E}[\varphi(\mathbf{X}; \mathbf{A})f(\mathbf{X}; \mathbf{A})], \quad (64)$$

where we used the tower property in the last step. For convenience, we define:

$$\beta \stackrel{\text{def}}{=} \frac{\mathbb{E}[\mathbf{Z}\varphi(\mathbf{X}; \mathbf{A}|\omega)]}{\sqrt{1 - \omega}}, \quad \widehat{\varphi}(x; a|\omega) \stackrel{\text{def}}{=} \varphi(x; a|\omega) - \beta x. \quad (65)$$

Observe that  $\widehat{\varphi}$  satisfies the divergence-free requirement:

$$\mathbb{E}[\mathbf{Z}\widehat{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = 0. \quad (66)$$

We consider the cross-term in (64):

$$\begin{aligned} \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})] &= \mathbb{E}[\widehat{\varphi}(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})] + \beta\mathbb{E}[\mathbf{X}f(\mathbf{X}; \mathbf{A})] \\ &\stackrel{\text{(a)}}{=} \mathbb{E}[\widehat{\varphi}(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})] + \beta\sqrt{\omega} \cdot \mathbb{E}[\mathbf{X}_* f(\mathbf{X}; \mathbf{A})] \\ &\stackrel{\text{(b)}}{=} \mathbb{E}[\widehat{\varphi}(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})] + \beta\sqrt{\omega} \cdot \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})]. \end{aligned} \quad (67)$$

In the above display step (a) follows from the fact that  $\mathbb{E}[\mathbf{Z}f(\mathbf{X}; \mathbf{A})] = 0$  and step (b) uses the tower property. We claim that:

$$\omega < 1 \implies \beta\sqrt{\omega} \neq 1. \quad (68)$$

Assuming the above claim, we can split our analysis into two cases.

**Case 1:**  $\omega < 1$ . We first prove each of the claims made in the lemma assuming  $\omega < 1$ .

**Proof of Claim 1.** Notice that (68) allows us to rearrange (67) to obtain:

$$\mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})] = \frac{\mathbb{E}[\widehat{\varphi}(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})]}{1 - \beta\sqrt{\omega}} = \mathbb{E}[\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)f(\mathbf{X}; \mathbf{A})], \quad (69)$$

where we recalled from the statement of the lemma that  $\bar{\varphi}$  was defined as:

$$\begin{aligned} \bar{\varphi}(x; a|\omega) &\stackrel{\text{def}}{=} \left(1 - \frac{\sqrt{\omega}}{\sqrt{1 - \omega}} \cdot \mathbb{E}[\mathbf{Z}\varphi(\mathbf{X}; \mathbf{A}|\omega)]\right)^{-1} \cdot \left(\varphi(x; a|\omega) - \frac{\mathbb{E}[\mathbf{Z}\varphi(\mathbf{X}; \mathbf{A}|\omega)]}{\sqrt{1 - \omega}} \cdot x\right) \\ &\stackrel{\text{(65)}}{=} \frac{\varphi(x; a|\omega) - \beta x}{1 - \beta\sqrt{\omega}} \stackrel{\text{(65)}}{=} \frac{\widehat{\varphi}(x; a|\omega)}{1 - \beta\sqrt{\omega}} \end{aligned} \quad (70)$$

Substituting the formula (69) in (64) and completing the square yields the following MSE decomposition:

$$\mathbb{E}[(\mathbf{X}_* - f(\mathbf{X}; \mathbf{A}))^2] = 1 - \mathbb{E}[\bar{\varphi}(\mathbf{X}; \mathbf{A})^2] + \mathbb{E}[\{f(\mathbf{X}; \mathbf{A}) - \bar{\varphi}(\mathbf{X}; \mathbf{A})\}^2] \quad (71)$$

Since  $\bar{\varphi}$  satisfies the divergence-free requirement  $\mathbb{E}[\mathbf{Z}\bar{\varphi}(\mathbf{X}; \mathbf{A})] = 0$  (cf. (66)), the above decomposition shows that  $\bar{\varphi}$  is the DMMSE estimator.

**Proof of Claim 2.** Setting  $f = \bar{\varphi}(\cdot; \cdot|\omega)$  in (69) and (71) yields:

$$\mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\bar{\varphi}^2(\mathbf{X}; \mathbf{A}|\omega)] = 1 - \mathbb{E}[(\mathbf{X}_* - \bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega))^2] = 1 - \text{dmmse}_\pi(\omega).$$

By the Tower property,  $\mathbb{E}[\mathbf{X}_*\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)]$ . Hence,

$$\mathbb{E}[\mathbf{X}_*\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\bar{\varphi}^2(\mathbf{X}; \mathbf{A}|\omega)] = 1 - \text{dmmse}_\pi(\omega), \quad (72)$$

which is the second claim made in the lemma.

**Proof of Claim 3.** Notice from (70) that the MMSE and DMMSE estimators are related by:

$$\varphi(x; a|\omega) = (1 - \beta\sqrt{\omega}) \cdot \bar{\varphi}(x; a|\omega) + \beta x. \quad (73)$$

Hence, we can relate  $\text{mmse}_\pi(\omega)$  and  $\text{dmmse}_\pi(\omega)$  as follows:

$$\begin{aligned} \text{mmse}_\pi(\omega) &= \mathbb{E}[(\mathbf{X}_* - \varphi(\mathbf{X}; \mathbf{A}|\omega))^2] = \mathbb{E}[\mathbf{X}_*^2] - 2\mathbb{E}[\mathbf{X}_*\varphi(\mathbf{X}; \mathbf{A}|\omega)] + \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)^2] \\ &\stackrel{(a)}{=} 1 - \mathbb{E}[\mathbf{X}_*\varphi(\mathbf{X}; \mathbf{A}|\omega)] \end{aligned} \quad (74)$$

$$\begin{aligned} &\stackrel{(73)}{=} 1 - (1 - \beta\sqrt{\omega}) \cdot \mathbb{E}[\mathbf{X}_*\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] - \beta\mathbb{E}[\mathbf{X}_*\mathbf{X}] \\ &= (1 - \beta\sqrt{\omega}) \cdot (1 - \mathbb{E}[\mathbf{X}_*\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)]) \\ &\stackrel{(72)}{=} (1 - \beta\sqrt{\omega}) \cdot \text{dmmse}_\pi(\omega). \end{aligned} \quad (75)$$

In the above display, step (a) relies on the assumption  $\mathbb{E}[\mathbf{X}_*^2] = 1$  (cf. Assumption 1) and the tower property  $\mathbb{E}[\mathbf{X}_*\varphi(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)^2]$ . Moreover, notice that:

$$\begin{aligned} \beta &\stackrel{\text{def}}{=} \frac{\mathbb{E}[\mathbf{Z}\varphi(\mathbf{X}; \mathbf{A}|\omega)]}{\sqrt{1 - \omega}} \stackrel{(62)}{=} \frac{\mathbb{E}[\mathbf{X}\varphi(\mathbf{X}; \mathbf{A}|\omega)] - \sqrt{\omega}\mathbb{E}[\mathbf{X}_*\varphi(\mathbf{X}; \mathbf{A}|\omega)]}{1 - \omega} \\ &\stackrel{(a)}{=} \frac{\mathbb{E}[\mathbf{X}\mathbf{X}_*] - \sqrt{\omega}\mathbb{E}[\mathbf{X}_*\varphi(\mathbf{X}; \mathbf{A}|\omega)]}{1 - \omega} = \frac{\sqrt{\omega}(1 - \mathbb{E}[\mathbf{X}_*\varphi(\mathbf{X}; \mathbf{A}|\omega)])}{1 - \omega} \stackrel{(74)}{=} \frac{\sqrt{\omega} \cdot \text{mmse}_\pi(\omega)}{1 - \omega}. \end{aligned}$$

In the above display, step (a) again uses the tower property  $\mathbb{E}[\mathbf{X}_*\varphi(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)^2]$ . Finally, we substitute the formula for  $\beta$  obtained in the above display in (75) to obtain:

$$\frac{1}{\text{dmmse}_\pi(\omega)} = \frac{1}{\text{mmse}_\pi(\omega)} - \frac{\omega}{1 - \omega}$$

after rearrangement. This is exactly the third claim made in the lemma.

**Case 2:**  $\omega = 1$ . When  $\omega = 1$ ,  $\mathbf{X} = \mathbf{X}_*$ . The estimator  $f(\mathbf{X}; \mathbf{A}) = \mathbf{X}$  has zero MSE and satisfies the divergence-free requirement  $\mathbb{E}[\mathbf{Z}\mathbf{X}] = \mathbb{E}[\mathbf{Z}]\mathbb{E}[\mathbf{X}_*] = 0$ . Hence,  $\bar{\varphi}(x; a) \stackrel{\text{def}}{=} x$  is the DMMSE estimator. Moreover, the identity (which is the second claim made in the lemma):

$$\mathbb{E}[\mathbf{X}_*\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\varphi(\mathbf{X}; \mathbf{A}|\omega)\bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega)] = \mathbb{E}[\bar{\varphi}^2(\mathbf{X}; \mathbf{A}|\omega)] = 1 - \text{dmmse}_\pi(\omega)$$

holds trivially in this case since  $\mathbf{X}_* = \bar{\varphi}(\mathbf{X}; \mathbf{A}|\omega) = \mathbf{X}$ ,  $\text{dmmse}_\pi(1) = 0$ , and  $\mathbb{E}[\mathbf{X}_*^2] = 1$ .

To finish the proof of the lemma, we need to prove the claim in (68).

**Proof of (68).** We prove the contrapositive of (68):  $\beta\sqrt{\omega} = 1 \implies \omega = 1$ . When  $\beta\sqrt{\omega} = 1$ , (67) shows that for any  $f \in L^2(\mathbf{X}; \mathbf{A})$  which satisfies the divergence-free requirement  $\mathbb{E}[\mathbf{Z}f(\mathbf{X}; \mathbf{A})] = 0$ , we have  $\mathbb{E}[\bar{\varphi}(\mathbf{X}; \mathbf{A})f(\mathbf{X}; \mathbf{A})] = 0$ . Taking  $f = \bar{\varphi}$  (cf. (66)), we conclude that  $\mathbb{E}[\bar{\varphi}^2] = 0$ . Recalling the definition of  $\bar{\varphi}$ , we obtain the following formula for the MMSE estimator for the scalar Gaussian channel:

$$\varphi(\mathbf{X}; \mathbf{A}) = \beta\mathbf{X} = \frac{\mathbf{X}}{\sqrt{\omega}},$$

which gives us the following formula for the MMSE:

$$\text{MMSE}(\mathbf{X}_*|\mathbf{X}; \mathbf{A}) = \frac{1 - \omega}{\omega}.$$



Since  $\text{MMSE}(\mathbf{X}_\star|\mathbf{X};\mathbf{A})$  is bounded by the MSE of the linear estimator  $\sqrt{\omega} \cdot \mathbf{X}$  (this is the linear estimator with the lowest MSE),

$$\text{MMSE}(\mathbf{X}_\star|\mathbf{X};\mathbf{A}) = \frac{1-\omega}{\omega} \leq \mathbb{E}[(\sqrt{\omega}\mathbf{X} - \mathbf{X}_\star)^2] = 1 - \omega.$$

Since  $\omega \in [0, 1]$ , by rearranging the above inequality we conclude that  $\omega = 1$ .  $\square$

### A.2.2 Proof of Lemma 3

*Proof of Lemma 3.* It will be convenient to define the random vector  $\mathbf{X}_{\leq t} \stackrel{\text{def}}{=} (\mathbf{X}_1, \dots, \mathbf{X}_t)^\top$  and the random variable  $\mathbf{S} \stackrel{\text{def}}{=} \langle \mathbf{X}_{\leq t}, v_{\text{opt}} \rangle$ . Observe that the conditional joint distribution of  $(\mathbf{X}_{\leq t}, \mathbf{S})$  given  $(\mathbf{X}_\star; \mathbf{A})$  is:

$$\begin{bmatrix} \mathbf{X}_{\leq t} \\ \mathbf{S} \end{bmatrix} \mid (\mathbf{X}_\star; \mathbf{A}) \sim \mathcal{N} \left( \mathbf{X}_\star \cdot \begin{bmatrix} \alpha \\ \sqrt{\omega_{\text{eff}}} \end{bmatrix}, \begin{bmatrix} \Sigma & v_{\text{opt}}^\top \Sigma \\ \Sigma v_{\text{opt}} & 1 - \omega_{\text{eff}} \end{bmatrix} \right).$$

Using Gaussian conditioning, we find that the conditional distribution of  $\mathbf{X}_{\leq t}$  given  $(\mathbf{S}, \mathbf{X}_\star; \mathbf{A})$  is given by:

$$\mathbf{X}_{\leq t} \mid (\mathbf{S}, \mathbf{X}_\star; \mathbf{A}) \sim \mathcal{N} \left( \frac{\mathbf{S}\alpha}{\sqrt{\omega_{\text{eff}}}}, \Sigma - \frac{1 - \omega_{\text{eff}}}{\omega_{\text{eff}}} \cdot \alpha\alpha^\top \right).$$

Since the RHS does not have any dependence on  $(\mathbf{X}_\star; \mathbf{A})$ , this implies that  $\mathbf{X}_{\leq t}$  and  $(\mathbf{X}_\star; \mathbf{A})$  are conditionally independent given  $\mathbf{S}$ :

$$\mathbf{X}_{\leq t} \perp (\mathbf{X}_\star; \mathbf{A}) \mid \mathbf{S}. \quad (76)$$

**Analysis of MMSE.** We first prove the MMSE formula and identify the MMSE estimator. To do so, we begin by showing the lower bound:

$$\text{MMSE}(\mathbf{X}_\star|\mathbf{X}_{\leq t}; \mathbf{A}) = \min_{f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] \geq \text{mmse}_\pi(\omega_{\text{eff}}).$$

To do so, we consider any  $f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})$ :

$$\begin{aligned} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] &\stackrel{(a)}{\geq} \mathbb{E}[\{\mathbf{X}_\star - \mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A})|\mathbf{S}, \mathbf{X}_\star, \mathbf{A}]\}^2] \\ &\stackrel{(76)}{=} \mathbb{E}[\{\mathbf{X}_\star - \mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A})|\mathbf{S}, \mathbf{A}]\}^2] \\ &\stackrel{(b)}{\geq} \min_{g \in L^2(\mathbf{S}; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - g(\mathbf{S}; \mathbf{A})\}^2] \\ &= \text{MMSE}(\mathbf{X}_\star|\mathbf{S}; \mathbf{A}) \\ &\stackrel{(c)}{=} \text{mmse}_\pi(\omega_{\text{eff}}), \end{aligned}$$

where step (a) used Jensen's inequality, step (b) follows by observing  $\mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A})|\mathbf{S}, \mathbf{A}] \in L^2(\mathbf{S}; \mathbf{A})$ , and step (c) is justified by noticing that  $(\mathbf{X}_\star, \mathbf{S}; \mathbf{A})$  forms a scalar Gaussian channel with SNR  $\omega_{\text{eff}}$ . Since  $f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})$  was arbitrary, we can minimize the LHS of the above display with respect to  $f$ , which yields:

$$\text{MMSE}(\mathbf{X}_\star|\mathbf{X}_{\leq t}; \mathbf{A}) \geq \text{mmse}_\pi(\omega_{\text{eff}}). \quad (77)$$

To obtain the upper bound, we observe that:

$$\text{mmse}_\pi(\omega_{\text{eff}}) \stackrel{(a)}{=} \mathbb{E}[\{\mathbf{X}_\star - \varphi(\mathbf{S}; \mathbf{A}|\omega_{\text{eff}})\}^2] \stackrel{(b)}{\geq} \min_{f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] \stackrel{\text{def}}{=} \text{MMSE}(\mathbf{X}_\star|\mathbf{X}_{\leq t}; \mathbf{A}), \quad (78)$$

where the equality in (a) follows by recalling the formula for the MMSE estimator in a scalar Gaussian channel and the inequality in (b) follows by observing that:

$$\varphi(\mathbf{S}; \mathbf{A}|\omega_{\text{eff}}) = \varphi(\langle \mathbf{X}_{\leq t}, v_{\text{opt}} \rangle; \mathbf{A}|\omega_{\text{eff}}) \stackrel{\text{def}}{=} f_\star(\mathbf{X}_{\leq t}; \mathbf{A}) \in L^2(\mathbf{X}_{\leq t}; \mathbf{A}).$$

Combining the conclusions of (77) and (78) yields:

$$\text{MMSE}(\mathbf{X}_\star|\mathbf{X}_{\leq t}; \mathbf{A}) = \mathbb{E}[\{\mathbf{X}_\star - f_\star(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] = \text{mmse}_\pi(\omega_{\text{eff}}),$$

which gives us the claimed formula for the MMSE and the MMSE estimator.

**Analysis of DMMSE.** We begin by observing that since  $(\mathbf{X}_\star, \mathbf{S} \stackrel{\text{def}}{=} \langle v_{\text{opt}}, \mathbf{X}_{\leq t} \rangle; \mathbf{A})$  form a scalar Gaussian channel with SNR  $\omega_{\text{eff}}$ , we can write  $\mathbf{S}$  as:

$$\mathbf{S} \stackrel{\text{def}}{=} \langle v_{\text{opt}}, \mathbf{X}_{\leq t} \rangle = \sqrt{\omega_{\text{eff}}} \mathbf{X}_\star + \sqrt{1 - \omega_{\text{eff}}} \mathbf{W}, \quad \mathbf{W} \stackrel{\text{def}}{=} \frac{\langle v_{\text{opt}}, \mathbf{Z}_{\leq t} \rangle}{\sqrt{1 - \omega_{\text{eff}}}} \sim \mathcal{N}(0, 1). \quad (79)$$

As before, we begin by showing the lower bound:

$$\begin{aligned} \text{DMMSE}(\mathbf{X}_\star | \mathbf{X}_{\leq t}; \mathbf{A}) &\stackrel{\text{def}}{=} \left( \min_{f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] \text{ subject to } \mathbb{E}[\mathbf{Z}_i f(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})] = 0 \forall i \in [t] \right) \\ &\geq \text{dmmse}_\pi(\omega_{\text{eff}}). \end{aligned}$$

For any  $f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})$  which satisfies the divergence-free constraints  $\mathbb{E}[\mathbf{Z}_i f(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})] = 0 \forall i \in [t]$ , we have:

$$\begin{aligned} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] &\stackrel{(a)}{\geq} \mathbb{E}[\{\mathbf{X}_\star - \mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{X}_\star, \mathbf{A}]\}^2] \\ &\stackrel{(76)}{=} \mathbb{E}[\{\mathbf{X}_\star - \mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{A}]\}^2] \\ &\stackrel{(b)}{\geq} \min_{g \in L^2(\mathbf{S}; \mathbf{A})} (\mathbb{E}[\{\mathbf{X}_\star - g(\mathbf{S}; \mathbf{A})\}^2] \text{ subject to } \mathbb{E}[\mathbf{W}g(\mathbf{S}; \mathbf{A})] = 0) \\ &= \text{DMMSE}(\mathbf{X}_\star | \mathbf{S}; \mathbf{A}) \\ &= \text{dmmse}_\pi(\omega_{\text{eff}}). \end{aligned}$$

In the above display:

1. Step (a) uses Jensen's inequality.
2. In step (b) we observed that  $\mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{A}] \in L^2(\mathbf{S}; \mathbf{A})$ . Moreover, the estimator  $\mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{A}]$  satisfies the divergence-free constraint  $\mathbb{E}[\mathbf{W} \mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{A}]] = 0$ . This can be verified as follows:

$$\mathbb{E}[\mathbf{W} \mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{A}]] \stackrel{(76)}{=} \mathbb{E}[\mathbf{W} \mathbb{E}[f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{X}_\star, \mathbf{A}]] \stackrel{(i)}{=} \mathbb{E}[\mathbb{E}[\mathbf{W} f(\mathbf{X}_{\leq t}; \mathbf{A}) | \mathbf{S}, \mathbf{X}_\star, \mathbf{A}]] = \mathbb{E}[\mathbf{W} f(\mathbf{X}_{\leq t}; \mathbf{A})] \stackrel{(ii)}{=} 0,$$

where we used the fact that  $\mathbf{W}$  is measurable with respect to  $\mathbf{X}_\star, \mathbf{S}$  (recall (79)) in step (i). To obtain the equality marked (ii), we expressed  $\mathbf{W}$  as a linear combination of  $\mathbf{Z}_{\leq t}$  (recall (79)) and used the fact that  $f$  satisfies the divergence-free constraints.

Since  $f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})$  was an arbitrary function that satisfies the divergence-free constraints, we can minimize the LHS of the above display with respect to  $f$  and obtain:

$$\text{DMMSE}(\mathbf{X}_\star | \mathbf{X}_{\leq t}; \mathbf{A}) \geq \text{dmmse}_\pi(\omega_{\text{eff}}). \quad (80)$$

To obtain the upper bound, we observe that:

$$\begin{aligned} \text{dmmse}_\pi(\omega_{\text{eff}}) &\stackrel{(a)}{=} \mathbb{E}[\{\mathbf{X}_\star - \bar{\varphi}(\mathbf{S}; \mathbf{A} | \omega_{\text{eff}})\}^2] \\ &\stackrel{(b)}{\geq} \left( \min_{f \in L^2(\mathbf{X}_{\leq t}; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - f(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] \text{ subject to } \mathbb{E}[\mathbf{Z}_i f(\mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})] = 0 \forall i \in [t] \right) \\ &\stackrel{\text{def}}{=} \text{DMMSE}(\mathbf{X}_\star | \mathbf{X}_{\leq t}; \mathbf{A}). \end{aligned} \quad (81)$$

In the above display:

1. The equality in (a) follows by recalling the formula for the DMMSE estimator in a scalar Gaussian channel (Lemma 2).
2. To obtain the inequality in (b), we observed that:

$$\bar{\varphi}(\mathbf{S}; \mathbf{A} | \omega_{\text{eff}}) = \bar{\varphi}(\langle \mathbf{X}_{\leq t}, v_{\text{opt}} \rangle; \mathbf{A} | \omega_{\text{eff}}) \stackrel{\text{def}}{=} \bar{f}_\star(\mathbf{X}_{\leq t}; \mathbf{A}) \in L^2(\mathbf{X}_{\leq t}; \mathbf{A}).$$

Moreover,  $\bar{f}_*(\mathbf{X}_{\leq t}; \mathbf{A}) \stackrel{\text{def}}{=} \bar{\varphi}(\langle \mathbf{X}_{\leq t}, v_{\text{opt}} \rangle; \mathbf{A} | \omega_{\text{eff}})$  satisfies the divergence-free requirements:

$$\mathbb{E}[\mathbf{Z}_i \bar{\varphi}(\langle \mathbf{X}_{\leq t}, v_{\text{opt}} \rangle; \mathbf{A} | \omega_{\text{eff}})] = 0 \quad \forall i \in [t].$$

This can be verified as follows:

$$\begin{aligned} \mathbb{E}[\mathbf{Z}_i \bar{\varphi}(\langle \mathbf{X}_{\leq t}, v_{\text{opt}} \rangle; \mathbf{A} | \omega_{\text{eff}})] &\stackrel{(79)}{=} \mathbb{E}[\mathbf{Z}_i \bar{\varphi}(\sqrt{\omega_{\text{eff}}} \mathbf{X}_* + \sqrt{1 - \omega_{\text{eff}}} \mathbf{W}; \mathbf{A} | \omega_{\text{eff}})] \\ &= \mathbb{E}[\bar{\varphi}(\sqrt{\omega_{\text{eff}}} \mathbf{X}_* + \sqrt{1 - \omega_{\text{eff}}} \mathbf{W}; \mathbf{A} | \omega_{\text{eff}}) \cdot \mathbb{E}[\mathbf{Z}_i | \mathbf{X}_*, \mathbf{A}, \mathbf{W}]] \\ &\stackrel{(i)}{=} c \mathbb{E}[\bar{\varphi}(\sqrt{\omega_{\text{eff}}} \mathbf{X}_* + \sqrt{1 - \omega_{\text{eff}}} \mathbf{W}; \mathbf{A} | \omega_{\text{eff}}) \cdot \mathbf{W}] \\ &\stackrel{(ii)}{=} 0, \end{aligned}$$

where the equality marked (i) follows by observing that  $(\mathbf{Z}_i, \mathbf{W})$  are centered Gaussian random variables independent of  $(\mathbf{X}_*, \mathbf{A})$  and hence  $\mathbb{E}[\mathbf{Z}_i | \mathbf{X}_*, \mathbf{A}, \mathbf{W}] = \mathbb{E}[\mathbf{Z}_i | \mathbf{W}] = c\mathbf{W}$  for some constant  $c$  (the exact formula for  $c$  is not important for the argument). The equality (ii) follows because  $\bar{\varphi}(\mathbf{S}; \mathbf{A})$  satisfies the divergence-free requirement  $\mathbb{E}[\mathbf{W} \bar{\varphi}(\mathbf{S}; \mathbf{A})] = 0$ .

Combining the conclusions of (80) and (81) yields:

$$\text{DMMSE}(\mathbf{X}_* | \mathbf{X}_{\leq t}; \mathbf{A}) = \mathbb{E}[\{\mathbf{X}_* - \bar{f}_*(\mathbf{X}_{\leq t}; \mathbf{A})\}^2] = \text{dmmse}_\pi(\omega_{\text{eff}}),$$

which gives us the claimed formula for the DMMSE and the DMMSE estimator. This concludes the proof of this lemma.  $\square$

### A.2.3 Proof of Lemma 4

*Proof of Lemma 4.* It is sufficient to show that for any  $\rho \in [0, 1]$  and any  $\omega \in [0, 1]$ ,  $\text{dmmse}_\pi(\omega) \leq \text{dmmse}_\pi(\rho\omega)$ . Consider a scalar Gaussian channel  $(\mathbf{X}_*, \mathbf{X}_0; \mathbf{A})$  at SNR  $\omega$ :

$$(\mathbf{X}_*; \mathbf{A}) \sim \pi, \quad \mathbf{X}_0 = \sqrt{\omega} \cdot \mathbf{X}_* + \sqrt{1 - \omega} \cdot \mathbf{Z}_0, \quad \mathbf{Z}_0 \sim \mathcal{N}(0, 1),$$

where  $(\mathbf{X}_*; \mathbf{A})$  and  $\mathbf{Z}_0$  are sampled independently. Let  $\mathbf{W}$  be a  $\mathcal{N}(0, 1)$  random variable, sampled independently of  $(\mathbf{X}_*; \mathbf{A})$  and  $\mathbf{Z}_0$ . We introduce the random variable  $\mathbf{X}_1 = \sqrt{\rho} \cdot \mathbf{X}_* + \sqrt{1 - \rho} \cdot \mathbf{W}$ . This construction ensures that  $(\mathbf{X}_*, \mathbf{X}_1; \mathbf{A})$  is a scalar Gaussian channel with SNR  $\rho\omega$ :

$$\mathbf{X}_1 = \sqrt{\rho\omega} \cdot \mathbf{X}_* + \sqrt{1 - \rho\omega} \cdot \mathbf{Z}_1, \quad \mathbf{Z}_1 \stackrel{\text{def}}{=} \frac{\sqrt{\rho(1 - \omega)} \cdot \mathbf{Z}_0 + \sqrt{1 - \rho} \cdot \mathbf{W}}{\sqrt{1 - \rho\omega}} \sim \mathcal{N}(0, 1). \quad (82)$$

To obtain the claim of the lemma, we observe that:

$$\begin{aligned} \text{dmmse}_\pi(\omega) &\stackrel{(a)}{=} \text{DMMSE}(\mathbf{X}_* | \mathbf{X}_0, \mathbf{X}_1; \mathbf{A}) \\ &\stackrel{\text{def}}{=} \left( \min_{f \in L^2(\mathbf{X}_0, \mathbf{X}_1; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_* - f(\mathbf{X}_0, \mathbf{X}_1; \mathbf{A})\}^2] \text{ subject to } \mathbb{E}[\mathbf{Z}_i f(\mathbf{X}_0, \mathbf{X}_1; \mathbf{A})] = 0 \forall i \in \{0, 1\} \right) \\ &\stackrel{(b)}{\leq} \left( \min_{f \in L^2(\mathbf{X}_1; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_* - f(\mathbf{X}_1; \mathbf{A})\}^2] \text{ subject to } \mathbb{E}[\mathbf{Z}_1 f(\mathbf{X}_1; \mathbf{A})] = 0 \right) \\ &\stackrel{\text{def}}{=} \text{DMMSE}(\mathbf{X}_* | \mathbf{X}_1; \mathbf{A}) \\ &= \text{dmmse}_\pi(\rho\omega). \end{aligned}$$

In the above display, the equality in step (a) follows from Lemma 3 (the effective SNR of the Gaussian channel  $(\mathbf{X}_*, \mathbf{X}_0, \mathbf{X}_1; \mathbf{A})$  is  $\omega$ ) and the inequality in step (b) follows by observing that the feasible set in the definition of  $\text{DMMSE}(\mathbf{X}_* | \mathbf{X}_1; \mathbf{A})$  is a subset of the feasible set in the definition of  $\text{DMMSE}(\mathbf{X}_* | \mathbf{X}_0, \mathbf{X}_1; \mathbf{A})$ . Indeed if  $f \in L^2(\mathbf{X}_1; \mathbf{A})$  satisfies  $\mathbb{E}[\mathbf{Z}_1 f(\mathbf{X}_1; \mathbf{A})] = 0$ , then  $f \in L^2(\mathbf{X}_0, \mathbf{X}_1; \mathbf{A})$ , and:

$$\mathbb{E}[\mathbf{Z}_0 f(\mathbf{X}_1; \mathbf{A})] = \mathbb{E}[f(\mathbf{X}_1; \mathbf{A}) \mathbb{E}[\mathbf{Z}_0 | \mathbf{X}_*, \mathbf{Z}_1, \mathbf{A}]] \stackrel{(i)}{=} \mathbb{E}[f(\mathbf{X}_1; \mathbf{A}) \mathbb{E}[\mathbf{Z}_0 | \mathbf{Z}_1]] \stackrel{(82)}{=} \frac{\sqrt{1 - \rho\omega}}{\sqrt{\rho(1 - \omega)}} \cdot \mathbb{E}[f(\mathbf{X}_1; \mathbf{A}) \mathbf{Z}_1] = 0,$$

where the equality marked (i) follows by observing that  $(\mathbf{Z}_0, \mathbf{Z}_1)$  are centered Gaussian random variables independent of  $(\mathbf{X}_*; \mathbf{A})$  (cf. (82)). This completes the proof of this lemma.  $\square$

## B State Evolution for OAMP Algorithms (Theorem 1)

In this appendix, we present the proof of Theorem 1. We begin by introducing the key ideas involved in the proof in the form of some intermediate results.

**Polynomial Approximation.** Consider a general OAMP algorithm:

$$\mathbf{x}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}. \quad (83)$$

By a polynomial approximation argument, we may assume that the matrix denoisers  $\{\Psi_t\}_{t \in \mathbb{N}}$  are polynomials. This argument is summarized in the following lemma, whose proof is deferred to Appendix B.1.

**Lemma 5.** *It is sufficient to prove Theorem 1 under the additional assumption that for each  $t \in \mathbb{N}$ , the matrix denoiser  $\Psi_t : \mathbb{R} \mapsto \mathbb{R}$  is a polynomial.*

**Orthogonal Decomposition.** Let  $(\mathbf{X}_\star, (\mathbf{X}_t)_{t \in \mathbb{N}}; \mathbf{A})$  denote the state evolution random variables associated with the OAMP algorithm in (83). A key idea involved in the proof of Theorem 1 is to consider the following decomposition of the functions  $(f_t)_{t \in \mathbb{N}}$ :

$$f_t(x_1, \dots, x_{t-1}; a) = \alpha_t x_\star + f_t^\perp(x_1, \dots, x_{t-1}; x_\star, a), \quad (84)$$

where:

$$\alpha_t \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{X}_\star f_t(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{A})], \quad f_t^\perp(x_1, \dots, x_{t-1}; x_\star, a) \stackrel{\text{def}}{=} f_t(x_1, \dots, x_{t-1}; a) - \alpha_t x_\star.$$

Notice that by construction,  $f_t^\perp$  is orthogonal to the signal state evolution random variable  $\mathbf{X}_\star$ :

$$\mathbb{E}[\mathbf{X}_\star f_t^\perp(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{X}_\star, \mathbf{A})] = 0. \quad (85)$$

Using this decomposition, we have:

$$\mathbf{x}_t = \alpha_t \cdot \Psi_t(\mathbf{Y}) \cdot \mathbf{x}_\star + \Psi_t(\mathbf{Y}) \cdot f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{x}_\star, \mathbf{a}). \quad (86)$$

Recall that the observed matrix  $\mathbf{Y}$  consists of a signal part and a noise part:  $\mathbf{Y} = \frac{\theta}{N} \mathbf{x}_\star \mathbf{x}_\star^\top + \mathbf{W}$ . The decomposition above allows us to disentangle the action of the signal component and the noise component using the following lemma, whose proof is deferred to Appendix B.2

**Lemma 6.** *Let  $\Psi : \mathbb{R} \mapsto \mathbb{R}$  be a polynomial function.*

1. *There exists a polynomial function  $\tilde{\Psi} : \mathbb{R} \mapsto \mathbb{R}$  associated with  $\Psi$  such that  $\Psi(\mathbf{Y}) \cdot \mathbf{x}_\star \stackrel{N \rightarrow \infty}{\simeq} \tilde{\Psi}(\mathbf{W}) \cdot \mathbf{x}_\star$ , where  $\stackrel{N \rightarrow \infty}{\simeq}$  denotes asymptotic equivalence of random vectors (Definition 2).*
2. *Let  $\mathbf{v}$  be a  $N$ -dimensional random vector with the property that:*

$$\frac{\langle \mathbf{W}^i \mathbf{v}, \mathbf{x}_\star \rangle}{N} \xrightarrow{\mathbb{P}} 0 \quad \forall i \in \mathbb{N}. \quad (87)$$

*Then,  $\Psi(\mathbf{Y}) \cdot \mathbf{v} \stackrel{N \rightarrow \infty}{\simeq} \Psi(\mathbf{W}) \cdot \mathbf{v}$ .*

We will use the lemma above to show that  $f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_t; \mathbf{x}_\star, \mathbf{a})$  does not interact with the signal component of  $\Psi_t(\mathbf{Y})$ , thanks to the orthogonality property in (85), which will guarantee that  $\mathbf{v} = f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_t; \mathbf{x}_\star, \mathbf{a})$  satisfies the orthogonality condition in (87). Formally, we will show that:

$$\Psi_t(\mathbf{Y}) \cdot f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{x}_\star, \mathbf{a}) \stackrel{N \rightarrow \infty}{\simeq} \Psi_t(\mathbf{W}) \cdot f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{x}_\star, \mathbf{a}), \quad (88a)$$

On the other hand,  $\mathbf{x}_\star$  does interact with the signal component in  $\Psi_t(\mathbf{Y})$  and this interaction is captured by the *transformed polynomial*  $\tilde{\Psi}_t$  associated with  $\Psi_t$ , whose existence is guaranteed by the lemma above:

$$\Psi_t(\mathbf{Y}) \cdot \mathbf{x}_\star \stackrel{N \rightarrow \infty}{\simeq} \tilde{\Psi}_t(\mathbf{W}) \cdot \mathbf{x}_\star. \quad (88b)$$

**Auxiliary OAMP Algorithm.** Using the approximations from (88) in the update equation of the OAMP algorithm given in (86), leads us to introduce an auxiliary OAMP algorithm which generates iterates  $\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots$  via the update rule:

$$\tilde{\mathbf{x}}_t = \alpha_t \cdot \tilde{\Psi}_t(\mathbf{W}) \cdot \mathbf{x}_\star + \Psi_t(\mathbf{W}) \cdot f_t^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{t-1}; \mathbf{x}_\star, \mathbf{a}) \quad \forall t \in \mathbb{N}. \quad (89)$$

This auxiliary OAMP algorithm will serve as an easy-to-analyze approximation to the original OAMP algorithm in (83). Since the auxiliary OAMP algorithm uses the rotationally invariant noise matrix  $\mathbf{W}$  in its iterations (rather than the non-rotationally invariant matrix  $\mathbf{Y}$  used by the original OAMP algorithm), its dynamics can be easily analyzed using known results on the state evolution of AMP algorithms for rotationally invariant matrices and the associated universality class [25, 26, 29, 50, 69, 74, 75]. Applying these results requires us to compute:

$$\text{plim}_{N \rightarrow \infty} \frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N}, \quad \text{plim}_{N \rightarrow \infty} \frac{\text{Tr}[\tilde{\Psi}_s(\mathbf{W})\tilde{\Psi}_t(\mathbf{W})]}{N} \quad \forall s, t \in \mathbb{N}.$$

At first glance, computing the above limits appears challenging as the transformed polynomials  $(\tilde{\Psi}_t)_{t \in \mathbb{N}}$  constructed in Lemma 6 have a complicated recursive characterization. Fortunately, the above limits have a simple formula in terms of  $\nu$ , the limiting spectral measure of  $\mathbf{Y}$  in the direction of the signal (recall Lemma 1).

**Lemma 7.** *For any  $s, t \in \mathbb{N}$ , we have:*

$$\frac{\text{Tr}[\tilde{\Psi}_t(\mathbf{W})]}{N} \xrightarrow{\mathbb{P}} \mathbb{E}[\tilde{\Psi}_t(\Lambda)] = \mathbb{E}[\Psi_t(\Lambda_\nu)], \quad \frac{\text{Tr}[\tilde{\Psi}_s(\mathbf{W})\tilde{\Psi}_t(\mathbf{W})]}{N} \xrightarrow{\mathbb{P}} \mathbb{E}[\tilde{\Psi}_s(\Lambda)\tilde{\Psi}_t(\Lambda)] = \mathbb{E}[\Psi_s(\Lambda_\nu)\Psi_t(\Lambda_\nu)],$$

where  $\Lambda \sim \mu$  and  $\Lambda_\nu \sim \nu$ .

The proof of this lemma is presented in Appendix B.3. An immediate consequence of this lemma and an existing result on the dynamics of AMP algorithms driven by rotationally invariant matrices [25, Theorem 2] is the following characterization of the dynamics of the auxiliary OAMP algorithm.

**Lemma 8.** 1. *For any  $t \in \mathbb{N}$ , the iterates generated by the auxiliary OAMP algorithm in (89) satisfy*

$$(\mathbf{x}_\star, \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_t; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}),$$

where  $(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$  are the state evolution random variables associated with the original OAMP algorithm in (83).

2. *Moreover, for any  $t \in \mathbb{N}$  and any  $i \in \mathbb{N}$ ,*

$$\frac{\langle \mathbf{W}^i f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{x}_\star, \mathbf{a}), \mathbf{x}_\star \rangle}{N} \xrightarrow{\mathbb{P}} 0.$$

The proof of the above lemma is deferred to Appendix B.4. We have now introduced all the key ideas involved in the proof of Theorem 1 and are in a position to present its proof.

*Proof of Theorem 1.* Our goal is to show that for any  $t \in \mathbb{N}$ ,

$$(\mathbf{x}_\star, \mathbf{x}_1, \dots, \mathbf{x}_t; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}).$$

In light of Lemma 8, it suffices to show that the auxiliary OAMP algorithm in (89) approximates the given OAMP algorithm in the sense:

$$\mathbf{x}_t \xrightarrow{N \rightarrow \infty} \tilde{\mathbf{x}}_t \quad \forall t \in \mathbb{N}. \quad (90)$$

We show the above claim by induction. As our induction hypothesis, we assume that  $\mathbf{x}_s \stackrel{N \rightarrow \infty}{\simeq} \tilde{\mathbf{x}}_s$  for all  $s < t$ , and verify the claim at iteration  $t$ :

$$\begin{aligned} \mathbf{x}_t &\stackrel{(86)}{=} \alpha_t \cdot \Psi_t(\mathbf{Y}) \cdot \mathbf{x}_* + \Psi_t(\mathbf{Y}) \cdot f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{x}_*, \mathbf{a}) \\ &\stackrel{(a)}{\stackrel{N \rightarrow \infty}{\simeq}} \alpha_t \cdot \Psi_t(\mathbf{Y}) \cdot \mathbf{x}_* + \Psi_t(\mathbf{Y}) \cdot f_t^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{t-1}; \mathbf{x}_*, \mathbf{a}) \\ &\stackrel{(b)}{\stackrel{N \rightarrow \infty}{\simeq}} \alpha_t \cdot \tilde{\Psi}_t(\mathbf{W}) \cdot \mathbf{x}_* + \Psi_t(\mathbf{W}) \cdot f_t^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{t-1}; \mathbf{x}_*, \mathbf{a}) \\ &\stackrel{(89)}{=} \tilde{\mathbf{x}}_t \end{aligned}$$

In the above display, step (a) follows from the induction hypothesis and the Lipschitz continuity of  $f_t^\perp$  (which is implied by the Lipschitz continuity of  $f_t$ ; see Definition 4 and (84)). In step (b), we appealed to Lemma 6. Indeed, Lemma 8 (claim (2)) guarantees that the orthogonality requirement (87) imposed in Lemma 6 is met. This proves the claim (90) by induction and concludes the proof of Theorem 1.  $\square$

## B.1 Proof of Lemma 5

*Proof of Lemma 5.* To prove Lemma 5, we assume that the claim of Theorem 1 holds under the additional assumption that the matrix denoisers used in the OAMP algorithm are polynomial functions. We will show that Theorem 1 continues to hold even without this additional assumption. To this end, we consider a general OAMP algorithm from Definition 4:

$$\mathbf{x}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}, \quad (91)$$

where the matrix denoisers  $(\Psi_t)_{t \in \mathbb{N}}$  are possibly non-polynomial functions. Let  $(\mathbf{X}_*, \{\mathbf{X}_t\}_{t \in \mathbb{N}}; \mathbf{A})$  denote the state evolution random variables associated with the OAMP algorithm above. Our goal is to show that:

$$(\mathbf{x}_*, \mathbf{x}_1, \dots, \mathbf{x}_t; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_*, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \quad \forall t \in \mathbb{N}.$$

**Polynomial Approximation.** We begin by approximating the matrix denoisers  $\Psi_t : \mathbb{R} \mapsto \mathbb{R}$  by polynomials. Recall from Assumption 2, that there is an  $N$ -independent finite constant  $C$  such that  $\|\mathbf{W}\|_{\text{op}} \leq C$ . Furthermore, thanks to Assumption 1,

$$\|\mathbf{Y}\|_{\text{op}} = \left\| \frac{\theta}{n} \mathbf{x}_* \mathbf{x}_*^\top + \mathbf{W} \right\|_{\text{op}} \leq C + \frac{\theta \cdot \|\mathbf{x}_*\|^2}{N} \xrightarrow{\mathbb{P}} C + \theta.$$

Hence, the event:

$$\mathcal{E}_N \stackrel{\text{def}}{=} \{\|\mathbf{Y}\|_{\text{op}} \leq K\} \quad \text{with} \quad K \stackrel{\text{def}}{=} C + \theta + 1$$

occurs with probability tending to 1. We restrict ourselves to this good event in the remainder of the proof. For each  $t \in \mathbb{N}$ , the Weierstrass Approximation Theorem guarantees the existence of a sequence of approximating polynomial functions  $(\Psi_t^{(D)})_{D \in \mathbb{N}}$  of increasing degree, where  $\Psi_t^{(D)}$  is a degree  $D$  polynomial, with the approximation guarantee:

$$\lim_{D \rightarrow \infty} \sup_{|\lambda| \leq K} \left| \Psi_t^{(D)}(\lambda) - \Psi_t(\lambda) \right| = 0, \quad (92)$$

Consider the OAMP algorithm which uses the degree- $D$  approximations  $(\Psi_t^{(D)})_{t \in \mathbb{N}}$  as the matrix denoisers:

$$\mathbf{x}_t^{(D)} = \left( \Psi_t^{(D)}(\mathbf{Y}) - \mathbb{E}[\Psi_t^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( f_t(\mathbf{x}_{<t}^{(D)}; \mathbf{a}) - \sum_{s=1}^{t-1} \mathbb{E}[\partial_s f_t(\mathbf{X}_{<t}^{(D)}; \mathbf{A})] \cdot \mathbf{x}_s^{(D)} \right) \quad \forall t \in \mathbb{N}, \quad (93)$$

where  $(\mathbf{X}_*, (\mathbf{X}_t^{(D)})_{t \in \mathbb{N}}; \mathbf{A})$  are the state evolution random variables associated with the OAMP algorithm above. Since we assume that Theorem 1 holds for OAMP algorithms with polynomial matrix denoisers, we have that:

$$(\mathbf{x}_*, \mathbf{x}_1^{(D)}, \dots, \mathbf{x}_t^{(D)}; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_*, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_t^{(D)}; \mathbf{A}) \quad \forall t \in \mathbb{N}, D \in \mathbb{N}.$$

We make the following claims regarding the polynomial approximation to the original OAMP algorithm and its state evolution random variables.

**Claim 1** (Convergence of State Evolution Random Variables). For any  $t \in \mathbb{N}$ ,

$$(\mathbf{X}_\star, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_t^{(D)}; \mathbf{A}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \quad \text{as } D \rightarrow \infty,$$

where  $\xrightarrow{W_2}$  denotes Wasserstein-2 convergence for random variables. Formally, this means that for any test function  $h : \mathbb{R}^{t+1} \times \mathbb{R}^k$ , independent of  $N$ , which satisfies the smoothness condition:

$$|h(x; a) - h(x'; a')| \leq L \cdot (\|x - x'\| + \|a - a'\|) \cdot (1 + \|x\| + \|x'\| + \|a\| + \|a'\|) \quad \forall x, x' \in \mathbb{R}^{t+1}, a, a' \in \mathbb{R}^k, \quad (94)$$

for some  $L < \infty$ , we have:

$$\lim_{D \rightarrow \infty} \mathbb{E}[h(\mathbf{X}_1^{(D)}, \dots, \mathbf{X}_t^{(D)}; \mathbf{A})] = \mathbb{E}[h(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})].$$

**Claim 2** (Convergence of Iterates). For any  $t \in \mathbb{N}$ ,

$$\limsup_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \frac{\|\mathbf{x}_t - \mathbf{x}_t^{(D)}\|^2}{N} = 0.$$

We first prove Lemma 5 assuming the above results. To do so, we consider a test function  $h : \mathbb{R}^t \times \mathbb{R}^k \mapsto \mathbb{R}$ , which satisfies the smoothness hypothesis in (94) and show that:

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N h(x_\star[i], x_1[i], \dots, x[t]; a[i]) = \mathbb{E}h(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}).$$

Notice that for any  $D \in \mathbb{N}$ , we have the decomposition:

$$\left| \frac{1}{N} \sum_{i=1}^N h(x_\star[i], x_1[i], \dots, x[t]; a[i]) - \mathbb{E}h(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \right| \leq (i) + (ii) + (iii), \quad (95)$$

where,

$$\begin{aligned} (i) &\stackrel{\text{def}}{=} \left| \frac{1}{N} \sum_{i=1}^N h(x_\star[i], x_1[i], \dots, x[t]; a[i]) - \frac{1}{N} \sum_{i=1}^N h(x_\star[i], x_1^{(D)}[i], \dots, x_t^{(D)}[i]; a[i]) \right| \\ (ii) &\stackrel{\text{def}}{=} \left| \frac{1}{N} \sum_{i=1}^N h(x_\star[i], x_1^{(D)}[i], \dots, x_t^{(D)}[i]; a[i]) - \mathbb{E}[h(\mathbf{X}_\star, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_t^{(D)}; \mathbf{A})] \right| \\ (iii) &\stackrel{\text{def}}{=} \left| \mathbb{E}[h(\mathbf{X}_\star, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_t^{(D)}; \mathbf{A})] - \mathbb{E}[h(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})] \right|. \end{aligned}$$

We consider each of the terms above. From Claim 2 and (94), we conclude that:

$$\lim_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} (i) = 0.$$

Since we assume that Theorem 1 holds for OAMP algorithms with polynomial matrix denoisers,

$$\text{plim}_{N \rightarrow \infty} (ii) = 0.$$

Finally, Claim 1 implies that:

$$\lim_{D \rightarrow \infty} (iii) = 0.$$

We let  $N \rightarrow \infty$  and  $D \rightarrow \infty$  in (95) and use the results above to conclude that:

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N h(x_\star[i], x_1[i], \dots, x[t]; a[i]) = \mathbb{E}h(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}),$$

as desired. To finish the proof of the lemma, we need to prove Claim 1 and Claim 2.

**Convergence of State Evolution Random Variables (Proof of Claim 1).** Our goal is to show that:

$$(\mathbf{X}_\star, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_t^{(D)}; \mathbf{A}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}) \quad \text{as } D \rightarrow \infty \quad \forall t \in \mathbb{N}, \quad (96)$$

We show this by induction. As our induction hypothesis, we assume that (96) holds for some  $t \in \mathbb{N}$  and show that:

$$(\mathbf{X}_\star, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_{t+1}^{(D)}; \mathbf{A}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_{t+1}; \mathbf{A}) \quad \text{as } D \rightarrow \infty.$$

To this end, we consider a test function  $h : \mathbb{R}^{t+2} \times \mathbb{R}^k \mapsto \mathbb{R}$  which satisfies the smoothness requirement in (94) and verify that:

$$\lim_{D \rightarrow \infty} \mathbb{E}[h(\mathbf{X}_\star, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_{t+1}^{(D)}; \mathbf{A})] = \mathbb{E}[h(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_{t+1}; \mathbf{A})]. \quad (97)$$

Recall from Definition 4 that the joint distribution of the state evolution random variables is given by:

$$(\mathbf{X}_\star; \mathbf{A}) \sim \pi, \quad \mathbf{X}_i = \beta_i \mathbf{X}_\star + \mathbf{Z}_i, \quad \mathbf{X}_i^{(D)} = \beta_i^{(D)} \cdot \mathbf{X}_\star + \mathbf{Z}_i^{(D)} \quad \forall i \in [t+1],$$

where  $(\mathbf{Z}_1, \dots, \mathbf{Z}_{t+1})$  and  $(\mathbf{Z}_1^{(D)}, \dots, \mathbf{Z}_{t+1}^{(D)})$  are Gaussian random vectors, sampled independently of  $(\mathbf{X}_\star; \mathbf{A})$ :

$$(\mathbf{Z}_1, \dots, \mathbf{Z}_{t+1}) \sim \mathcal{N}(0, \Sigma_{t+1}), \quad (\mathbf{Z}_1^{(D)}, \dots, \mathbf{Z}_{t+1}^{(D)}) \sim \mathcal{N}(0, \Sigma_{t+1}^{(D)}).$$

In the above equations, coefficients  $(\beta_i)_{i \in [t+1]}$  and  $(\beta_i^{(D)})_{i \in [t+1]}$  are given by:

$$\beta_i \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i] \cdot \mathbb{E}[\Psi_i(\Lambda_\nu)], \quad \beta_i^{(D)} \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i^{(D)}] \cdot \left( \mathbb{E}[\Psi_i^{(D)}(\Lambda_\nu)] - \mathbb{E}[\Psi_i^{(D)}(\Lambda)] \right) \quad \forall i \in [t+1],$$

and the entries of the covariance matrices  $\Sigma_{t+1}, \Sigma_{t+1}^{(D)}$  are given by:

$$\begin{aligned} (\Sigma_{t+1})_{i,j} &= \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i] \cdot \mathbb{E}[\mathbf{X}_\star \mathbf{F}_j] \cdot \text{Cov}[\Psi_i(\Lambda_\nu), \Psi_j(\Lambda_\nu)] \\ &\quad + (\mathbb{E}[\mathbf{F}_i \mathbf{F}_j] - \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i] \mathbb{E}[\mathbf{X}_\star \mathbf{F}_j]) \cdot \text{Cov}[\Psi_i(\Lambda), \Psi_j(\Lambda)] \quad \forall i, j \in [t+1], \\ \left( \Sigma_{t+1}^{(D)} \right)_{i,j} &= \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i^{(D)}] \cdot \mathbb{E}[\mathbf{X}_\star \mathbf{F}_j^{(D)}] \cdot \text{Cov}[\Psi_i^{(D)}(\Lambda_\nu), \Psi_j^{(D)}(\Lambda_\nu)] \\ &\quad + (\mathbb{E}[\mathbf{F}_i^{(D)} \mathbf{F}_j^{(D)}] - \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i] \mathbb{E}[\mathbf{X}_\star \mathbf{F}_j]) \cdot \text{Cov}[\Psi_i^{(D)}(\Lambda), \Psi_j^{(D)}(\Lambda)] \quad \forall i, j \in [t+1], \end{aligned}$$

where  $\Lambda \sim \mu$  and  $\Lambda_\nu \sim \nu$  and:

$$\mathbf{F}_i \stackrel{\text{def}}{=} f_i(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}; \mathbf{A}), \quad \mathbf{F}_i^{(D)} \stackrel{\text{def}}{=} f_i(\mathbf{X}_1^{(D)}, \dots, \mathbf{X}_{i-1}^{(D)}; \mathbf{A}) - \sum_{j=1}^{i-1} \mathbb{E}[\partial_j f_i(\mathbf{X}_1^{(D)}, \dots, \mathbf{X}_{i-1}^{(D)}; \mathbf{A})] \cdot \mathbf{X}_j^{(D)}.$$

By a weak convergence and uniform integrability argument, (97) follows if we can show that:

$$\lim_{D \rightarrow \infty} \beta_i^{(D)} = \beta_i, \quad \lim_{D \rightarrow \infty} \left( \Sigma_{t+1}^{(D)} \right)_{i,j} = (\Sigma_{t+1})_{i,j} \quad \forall i, j \in [t+1]. \quad (98)$$

Indeed, from (92) we deduce that:

$$\begin{aligned} \lim_{D \rightarrow \infty} \mathbb{E}[\Psi_i^{(D)}(\Lambda)] &= \mathbb{E}[\Psi_i(\Lambda)], \quad \lim_{D \rightarrow \infty} \mathbb{E}[\Psi_i^{(D)}(\Lambda_\nu)] = \mathbb{E}[\Psi_i(\Lambda_\nu)], \\ \lim_{D \rightarrow \infty} \text{Cov}[\Psi_i^{(D)}(\Lambda), \Psi_j^{(D)}(\Lambda)] &= \text{Cov}[\Psi_i(\Lambda), \Psi_j(\Lambda)], \quad \lim_{D \rightarrow \infty} \text{Cov}[\Psi_i^{(D)}(\Lambda_\nu), \Psi_j^{(D)}(\Lambda_\nu)] = \text{Cov}[\Psi_i(\Lambda_\nu), \Psi_j(\Lambda_\nu)]. \end{aligned}$$

Furthermore, using the induction hypothesis  $(\mathbf{X}_\star, \mathbf{X}_1^{(D)}, \dots, \mathbf{X}_t^{(D)}; \mathbf{A}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$  and the fact that the functions  $f_1, \dots, f_t$  are continuously differentiable and Lipschitz, we conclude that:

$$\lim_{D \rightarrow \infty} \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i^{(D)}] = \mathbb{E}[\mathbf{X}_\star \mathbf{F}_i], \quad \lim_{D \rightarrow \infty} \mathbb{E}[\mathbf{F}_i^{(D)} \mathbf{F}_j^{(D)}] = \mathbb{E}[\mathbf{F}_i \mathbf{F}_j].$$

Plugging in these limits into the formulae for  $(\beta_i^{(D)})_{i \in [t+1]}$  and  $\Sigma_{t+1}, \Sigma_{t+1}^{(D)}$ , immediately yields the desired conclusion (98) and hence proves (97). This completes the proof of the claim (96) by induction.



**Convergence of Iterates (Proof of Claim 2).** We need to show that the iterates generated by the algorithm in (93) approximate the iterates generated by the original OAMP algorithm (91) in the sense that,

$$\limsup_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \frac{\|\mathbf{x}_t - \mathbf{x}_t^{(D)}\|^2}{N} = 0 \quad \forall t \in \mathbb{N}.$$

We will shorthand asymptotic approximation statements like the above using the notation:

$$\mathbf{x}_t^{(D)} \stackrel{N, D \rightarrow \infty}{\simeq} \mathbf{x}_t \quad \forall t \in \mathbb{N}. \quad (99)$$

We will show the above claim using induction. As our induction hypothesis, we will assume that the claim holds for all the iterates generated before step  $t$  of the OAMP algorithm:

$$\mathbf{x}_s^{(D)} \stackrel{N, D \rightarrow \infty}{\simeq} \mathbf{x}_s \quad \forall s < t. \quad (100)$$

To complete the inductive proof, we need to show that the claim applies to the iterate generated at step  $t$ :

$$\mathbf{x}_t^{(D)} \stackrel{N, D \rightarrow \infty}{\simeq} \mathbf{x}_t.$$

Recalling the update equation for  $\mathbf{x}_t$  from (91), we have:

$$\begin{aligned} \mathbf{x}_t &= \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \\ &\stackrel{(a)}{\stackrel{N, D \rightarrow \infty}{\simeq}} \left( \Psi_t^{(D)}(\mathbf{Y}) - \mathbb{E}[\Psi_t^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \\ &\stackrel{(b)}{\stackrel{N, D \rightarrow \infty}{\simeq}} \left( \Psi_t^{(D)}(\mathbf{Y}) - \mathbb{E}[\Psi_t^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) - \sum_{s=1}^{t-1} \mathbb{E}[\partial_s f_t(\mathbf{X}_1^{(D)}, \dots, \mathbf{X}_{t-1}^{(D)}; \mathbf{A})] \cdot \mathbf{x}_s \right) \\ &\stackrel{(c)}{\stackrel{N, D \rightarrow \infty}{\simeq}} \left( \Psi_t^{(D)}(\mathbf{Y}) - \mathbb{E}[\Psi_t^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( f_t(\mathbf{x}_1^{(D)}, \dots, \mathbf{x}_{t-1}^{(D)}; \mathbf{a}) - \sum_{s=1}^{t-1} \mathbb{E}[\partial_s f_t(\mathbf{X}_1^{(D)}, \dots, \mathbf{X}_{t-1}^{(D)}; \mathbf{A})] \cdot \mathbf{x}_s^{(D)} \right) \\ &\stackrel{(93)}{=} \mathbf{x}_t^{(D)}. \end{aligned}$$

In the above display:

1. The approximation in step (a) relies on the observations:

$$\begin{aligned} \limsup_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \|\Psi_t(\mathbf{Y}) - \Psi_t^{(D)}(\mathbf{Y})\|_{\text{op}} &\leq \lim_{D \rightarrow \infty} \sup_{|\lambda| \leq K} \left| \Psi_t^{(D)}(\lambda) - \Psi_t(\lambda) \right| \stackrel{(92)}{=} 0, \\ \lim_{D \rightarrow \infty} \mathbb{E}[\Psi_t^{(D)}(\Lambda)] &= \mathbb{E}[\Psi_t(\Lambda)] \stackrel{\text{Def. 4}}{=} 0. \end{aligned}$$

2. The approximation in step (b) uses Claim 1 and the fact that  $f_t$  is Lipschitz to conclude that:

$$\lim_{D \rightarrow \infty} \mathbb{E}[\partial_s f_t(\mathbf{X}_1^{(D)}, \dots, \mathbf{X}_{t-1}^{(D)}; \mathbf{A})] = \mathbb{E}[\partial_s f_t(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{A})] = 0,$$

where the last equality follows by recalling that the iterate denoisers satisfy the divergence-free requirement  $\mathbb{E}[\partial_s f_t(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{A})] = 0$  (recall Definition 4).

3. The approximation in step (c) relies on the induction hypothesis (100) and the fact that the functions  $f_t$  are Lipschitz.

Hence, (99) holds by induction. This concludes the proof of Lemma 5.  $\square$

## B.2 Proof of Lemma 6

*Proof of Lemma 6.* . Since any polynomial can be written as a linear combination of monomials, it suffices to prove the claims of the lemma when  $\Psi(\lambda) = \lambda^d \forall \lambda \in \mathbb{R}$  for some  $d \in \mathbb{N}$ .

**Proof of Claim (1).** We need to show that there exists a sequence of polynomials  $(Q_d)_{d \in \mathbb{N}} : \mathbb{R} \mapsto \mathbb{R}$  such that:

$$\mathbf{Y}^d \cdot \mathbf{x}_* \stackrel{N \rightarrow \infty}{\simeq} Q_d(\mathbf{W}) \cdot \mathbf{x}_* \quad \forall d \in \mathbb{N}. \quad (101)$$

We will show that (101) holds by induction on  $d$ . We assume that (101) holds for some  $d \in \mathbb{N}$  and show that the claim also holds for  $d + 1$ . Recalling that  $\mathbf{Y} = \theta \cdot \frac{\mathbf{x}_* \mathbf{x}_*^\top}{N} + \mathbf{W}$ , we have:

$$\begin{aligned} \mathbf{Y}^{d+1} \cdot \mathbf{x}_* &= \left( \theta \cdot \frac{\mathbf{x}_* \mathbf{x}_*^\top}{N} + \mathbf{W} \right) \cdot \mathbf{Y}^d \cdot \mathbf{x}_* \\ &\stackrel{(a)}{\stackrel{N \rightarrow \infty}{\simeq}} \left( \theta \cdot \frac{\mathbf{x}_* \mathbf{x}_*^\top}{N} + \mathbf{W} \right) \cdot Q_d(\mathbf{W}) \cdot \mathbf{x}_* \\ &= \theta \cdot \frac{\mathbf{x}_*^\top Q_d(\mathbf{W}) \mathbf{x}_*}{N} + \mathbf{W} Q_d(\mathbf{W}) \cdot \mathbf{x}_* \\ &\stackrel{(b)}{\stackrel{N \rightarrow \infty}{\simeq}} \theta \cdot \mathbb{E}[Q_d(\Lambda)] + \mathbf{W} Q_d(\mathbf{W}) \cdot \mathbf{x}_*. \end{aligned}$$

In the above display, step (a) follows from the induction hypothesis  $\mathbf{Y}^d \cdot \mathbf{x}_* \stackrel{N \rightarrow \infty}{\simeq} Q_d(\mathbf{W}) \cdot \mathbf{x}_*$  and step (b) uses a standard result on the concentration of quadratic forms of rotationally invariant matrices (see Fact 2 in Appendix F). We define the polynomial  $Q_{d+1} : \mathbb{R} \mapsto \mathbb{R}$  as:

$$Q_{d+1}(\lambda) = \theta \cdot \mathbb{E}[Q_d(\Lambda)] + \lambda Q_d(\lambda) \quad \forall \lambda \in \mathbb{R}.$$

Hence,  $\mathbf{Y}^{d+1} \cdot \mathbf{x}_* \stackrel{N \rightarrow \infty}{\simeq} Q_{d+1}(\mathbf{W}) \cdot \mathbf{x}_*$ , as desired. This proves the first claim made in the lemma by induction.

**Proof of Claim (2).** The proof of the second claim follows from an analogous induction argument. Specifically, we show by induction that:

$$\mathbf{Y}^d \cdot \mathbf{v} \stackrel{N \rightarrow \infty}{\simeq} \mathbf{W}^d \cdot \mathbf{v} \quad \forall d \in \mathbb{N}.$$

As before, we assume the claim holds for some  $d \in \mathbb{N}$ , and we verify it for  $d + 1$ :

$$\begin{aligned} \mathbf{Y}^{d+1} \cdot \mathbf{v} &= \left( \theta \cdot \frac{\mathbf{x}_* \mathbf{x}_*^\top}{N} + \mathbf{W} \right) \cdot \mathbf{Y}^d \cdot \mathbf{v} \\ &\stackrel{(a)}{\stackrel{N \rightarrow \infty}{\simeq}} \left( \theta \cdot \frac{\mathbf{x}_* \mathbf{x}_*^\top}{N} + \mathbf{W} \right) \cdot \mathbf{W}^d \cdot \mathbf{v} \\ &= \theta \cdot \frac{\mathbf{x}_*^\top \mathbf{W}^d \mathbf{v}}{N} + \mathbf{W}^{d+1} \cdot \mathbf{v} \\ &\stackrel{(b)}{\stackrel{N \rightarrow \infty}{\simeq}} \mathbf{W}^{d+1} \cdot \mathbf{v}, \end{aligned}$$

where step (a) follows from the induction hypothesis and step (b) uses the assumption  $\mathbf{x}_*^\top \mathbf{W}^d \mathbf{v} / N \xrightarrow{\mathbb{P}} 0$ . This concludes the proof of this lemma.  $\square$

### B.3 Proof of Lemma 7

*Proof of Lemma 7.* We prove the second claim, the proof of the first claim is analogous. Since the spectral measure of  $\mathbf{W}$  converges to  $\mu$  (Assumption 2),

$$\frac{\text{Tr}[\tilde{\Psi}_s(\mathbf{W}) \tilde{\Psi}_t(\mathbf{W})]}{N} \xrightarrow{\mathbb{P}} \mathbb{E}[\tilde{\Psi}_s(\Lambda) \tilde{\Psi}_t(\Lambda)], \quad \frac{\text{Tr}[\tilde{\Psi}_s(\mathbf{W}) \tilde{\Psi}_t(\mathbf{W})]}{N} \xrightarrow{\mathbb{P}} \mathbb{E}[\tilde{\Psi}_s(\Lambda) \tilde{\Psi}_t(\Lambda)], \quad \Lambda \sim \mu.$$

To obtain a more convenient formula for the limit which does not involve the transformed polynomials  $(\tilde{\Psi}_t)_{t \in \mathbb{N}}$ , we observe that:

$$\mathbb{E}[\tilde{\Psi}_s(\Lambda)\tilde{\Psi}_t(\Lambda)] \stackrel{(a)}{=} \text{plim}_{N \rightarrow \infty} \frac{\mathbf{x}_*^\top \tilde{\Psi}_s(\mathbf{W})\tilde{\Psi}_t(\mathbf{W})\mathbf{x}_*}{N} \stackrel{\text{Lem. 6}}{=} \text{plim}_{N \rightarrow \infty} \frac{\mathbf{x}_*^\top \Psi_s(\mathbf{Y})\Psi_t(\mathbf{Y})\mathbf{x}_*}{N}.$$

where the equality in step (a) follows from standard concentration results for quadratic forms of rotationally invariant matrices (see Fact 2 in Appendix F). Notice that the RHS of the above display can be written as an expectation with respect to  $\nu_N$ , the spectral measure of  $\mathbf{Y}$  in the direction of  $\mathbf{x}_*$ . Indeed, if  $\lambda_{1:N}(\mathbf{Y})$  and  $\mathbf{u}_{1:N}(\mathbf{Y})$  denote the eigenvalues and corresponding eigenvectors of  $\mathbf{Y}$ :

$$\frac{\mathbf{x}_*^\top \Psi_s(\mathbf{Y})\Psi_t(\mathbf{Y})\mathbf{x}_*}{N} = \frac{1}{N} \sum_{i=1}^N \langle \mathbf{x}_*, \mathbf{u}_i(\mathbf{Y}) \rangle^2 \cdot \Psi_s(\lambda_i(\mathbf{Y}))\Psi_t(\lambda_i(\mathbf{Y})) = \int_{\mathbb{R}} \Psi_s(\lambda)\Psi_t(\lambda) \nu_N(d\lambda).$$

By Lemma 1 (item (1)), we know that  $\nu_N$  converges weakly to  $\nu$ . Hence,

$$\text{plim}_{N \rightarrow \infty} \frac{\text{Tr}[\tilde{\Psi}_s(\mathbf{W})\tilde{\Psi}_t(\mathbf{W})]}{N} = \text{plim}_{N \rightarrow \infty} \frac{\mathbf{x}_*^\top \tilde{\Psi}_s(\mathbf{W})\tilde{\Psi}_t(\mathbf{W})\mathbf{x}_*}{N} = \mathbb{E}[\Psi_s(\Lambda_\nu)\Psi_t(\Lambda_\nu)], \quad \Lambda_\nu \sim \nu,$$

as claimed.  $\square$

## B.4 Proof of Lemma 8

*Proof.* Recall that we started with a general OAMP algorithm (cf. Definition 4) with polynomial matrix denoising functions:

$$\mathbf{x}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}, \quad (102)$$

with associated state evolution random variables

$$(\mathbf{X}_*, (\mathbf{X}_t)_{t \in \mathbb{N}}; \mathbf{A}),$$

and constructed an auxiliary OAMP algorithm:

$$\tilde{\mathbf{x}}_t = \alpha_t \cdot \tilde{\Psi}_t(\mathbf{W}) \cdot \mathbf{x}_* + \Psi_t(\mathbf{W}) \cdot f_t^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{t-1}; \mathbf{x}_*, \mathbf{a}) \quad \forall t \in \mathbb{N}, \quad (103)$$

which is intended as an easy-to-analyze approximation to the original OAMP algorithm in (102). In the above display:

$$\alpha_t \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{X}_*^\top f_t(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{A})], \quad f_t^\perp(x_1, \dots, x_{t-1}; \mathbf{x}_*, \mathbf{a}) \stackrel{\text{def}}{=} f_t(x_1, \dots, x_{t-1}; \mathbf{a}) - \alpha_t \mathbf{x}_*, \quad (104)$$

and  $(\tilde{\Psi}_t)_{t \in \mathbb{N}}$  are the transformed polynomials associated with matrix denoisers  $(\Psi_t)_{t \in \mathbb{N}}$ , which were constructed in Lemma 6. Our goal is to show that for any  $t \in \mathbb{N}$  and any  $i \in \mathbb{N}$ ,

$$(\mathbf{x}_*, \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_t; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_*, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}), \quad \frac{\langle \mathbf{W}^i f_t^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{t-1}; \mathbf{x}_*, \mathbf{a}), \mathbf{x}_* \rangle}{N} \xrightarrow{\mathbb{P}} 0.$$

We will show that the auxiliary OAMP algorithm can be re-written as an algorithm whose dynamics have characterized in prior work [25, Theorem 2]. To this end, we fix a  $i \in \mathbb{N}$  and a  $t \in \mathbb{N}$ . Introduce the random variables:

$$\Lambda \sim \mu, \quad \Lambda_\nu \sim \nu.$$

For any  $s \leq t$ , we define:

$$\mathbf{v}_s \stackrel{\text{def}}{=} \left( \tilde{\Psi}_s(\mathbf{W}) - \mathbb{E}[\tilde{\Psi}_s(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \mathbf{x}_*, \quad \mathbf{w}_s \stackrel{\text{def}}{=} \Psi_s(\mathbf{W}) \cdot f_s^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{s-1}; \mathbf{x}_*, \mathbf{a}), \quad (105a)$$

$$\mathbf{w}_{t+1} \stackrel{\text{def}}{=} (\mathbf{W}^i - \mathbb{E}[\Lambda^i] \cdot \mathbf{I}_N) \cdot f_t^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{t-1}; \mathbf{x}_*, \mathbf{a}). \quad (105b)$$

Comparing (103) and (105), we conclude that:

$$\tilde{\mathbf{x}}_s = \beta_s \cdot \mathbf{x}_\star + \alpha_s \cdot \mathbf{v}_s + \mathbf{w}_s \quad \forall s \leq t, \quad (106)$$

where we defined:

$$\beta_s \stackrel{\text{def}}{=} \alpha_s \cdot \mathbb{E}[\tilde{\Psi}_s(\Lambda)] \stackrel{\text{Lem. 7}}{=} \alpha_s \cdot \mathbb{E}[\Psi_s(\Lambda_\nu)] \quad \forall s \leq t.$$

Hence, (105) can be re-expressed as:

$$\begin{aligned} \mathbf{v}_s &= \left( \tilde{\Psi}_s(\mathbf{W}) - \mathbb{E}[\tilde{\Psi}_s(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \mathbf{x}_\star, \\ \mathbf{w}_s &\stackrel{(106)}{=} \tilde{\Psi}_s(\mathbf{W}) \cdot f_s^\perp(\beta_1 \mathbf{x}_\star + \alpha_1 \mathbf{v}_1 + \mathbf{w}_1, \dots, \beta_{s-1} \mathbf{x}_\star + \alpha_{s-1} \mathbf{v}_{s-1} + \mathbf{w}_{s-1}; \mathbf{x}_\star, \mathbf{a}). \end{aligned}$$

The above algorithm is an instance of a Vector Approximate Message Passing (VAMP) algorithm, as defined in [25, Section 4.1.2]. Moreover, [25, Theorem 2] shows that:

$$(\mathbf{x}_\star, \mathbf{v}_1, \dots, \mathbf{v}_t, \mathbf{w}_1, \dots, \mathbf{w}_{t+1}; \mathbf{a}) \xrightarrow{W_2} (\mathbf{X}_\star, \mathbf{V}_1, \dots, \mathbf{V}_t, \mathbf{W}_1, \dots, \mathbf{W}_{t+1}; \mathbf{A}), \quad (107)$$

where the collection of random variables:

$$(\mathbf{X}_\star; \mathbf{A}), \quad (\mathbf{V}_1, \dots, \mathbf{V}_t), \quad (\mathbf{W}_1, \dots, \mathbf{W}_{t+1})$$

are mutually independent. The random variables  $(\mathbf{V}_1, \dots, \mathbf{V}_t)$  have zero mean and are jointly Gaussian, and the same is true for the random variables  $(\mathbf{W}_1, \dots, \mathbf{W}_{t+1})$ . The covariance matrices of  $(\mathbf{V}_1, \dots, \mathbf{V}_t)$  and  $(\mathbf{W}_1, \dots, \mathbf{W}_{t+1})$  are given by the recursion:

$$\mathbb{E}[\mathbf{V}_s \mathbf{V}_\tau] = \text{Cov}[\tilde{\Psi}_s(\Lambda), \tilde{\Psi}_\tau(\Lambda)] \stackrel{\text{Lem. 7}}{=} \text{Cov}[\Psi_s(\Lambda_\nu), \Psi_\tau(\Lambda_\nu)] \quad \forall s, \tau \leq t, \quad (108a)$$

$$\mathbb{E}[\mathbf{W}_s \mathbf{W}_\tau] = \text{Cov}[\Psi_s(\Lambda), \Psi_\tau(\Lambda)] \cdot (\mathbb{E}[\mathbf{F}_s \mathbf{F}_\tau] - \alpha_s \alpha_\tau) \quad \forall s, \tau \leq t. \quad (108b)$$

In the above display, for any  $s \leq t$ ,  $\mathbf{F}_s \stackrel{\text{def}}{=}} f_s(\beta_1 \mathbf{X}_\star + \alpha_1 \mathbf{V}_1 + \mathbf{W}_1, \dots, \beta_{s-1} \mathbf{X}_\star + \alpha_{s-1} \mathbf{V}_{s-1} + \mathbf{W}_{s-1}; \mathbf{A})$ . Consequently,

$$\begin{aligned} (\mathbf{x}_\star, \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_t; \mathbf{a}) &\stackrel{(106)}{=} (\mathbf{x}_\star, \beta_1 \cdot \mathbf{x}_\star + \alpha_1 \cdot \mathbf{v}_1 + \mathbf{w}_1, \dots, \beta_t \cdot \mathbf{x}_\star + \alpha_t \cdot \mathbf{v}_t + \mathbf{w}_t; \mathbf{a}) \\ &\xrightarrow{W_2} (\mathbf{X}_\star, \beta_1 \cdot \mathbf{X}_\star + \alpha_1 \cdot \mathbf{V}_1 + \mathbf{W}_1, \dots, \beta_t \cdot \mathbf{X}_\star + \alpha_t \cdot \mathbf{V}_t + \mathbf{W}_t; \mathbf{A}) \quad [\text{Using (107)}] \\ &\stackrel{d}{=} (\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A}), \end{aligned} \quad (109)$$

where the distributional equality in the last step follows by comparing (108) and the description of the joint distribution of the state evolution random variables  $(\mathbf{X}_\star, \mathbf{X}_1, \dots, \mathbf{X}_t; \mathbf{A})$  corresponding to the original OAMP algorithm (102) provided in Definition 4. This proves the first claim made in the lemma. To prove the second claim, we observe that:

$$\begin{aligned} \frac{\langle \mathbf{W}^i f_t^\perp(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}; \mathbf{x}_\star, \mathbf{a}), \mathbf{x}_\star \rangle}{N} &\stackrel{(105)}{=} \frac{\langle \mathbf{w}_{t+1}, \mathbf{x}_\star \rangle}{N} + \mathbb{E}[\Lambda^i] \cdot \frac{\langle f_t^\perp(\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{t-1}; \mathbf{x}_\star, \mathbf{a}), \mathbf{x}_\star \rangle}{N} \\ &\xrightarrow{\mathbb{P}} \mathbb{E}[\mathbf{W}_{t+1} \mathbf{X}_\star] + \mathbb{E}[\Lambda^i] \cdot \mathbb{E}[\mathbf{X}_\star f_t^\perp(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{X}_\star, \mathbf{A})] \\ &= 0. \end{aligned}$$

In the above display, the convergence in the second step follows from (107) and (109). The equality in the last step follows by recalling that  $\mathbf{W}_{t+1}$  is a mean-zero Gaussian random variable independent of  $\mathbf{X}_\star$  and the definition of  $f_t^\perp$  in (104) ensures  $\mathbb{E}[\mathbf{X}_\star f_t^\perp(\mathbf{X}_1, \dots, \mathbf{X}_{t-1}; \mathbf{X}_\star, \mathbf{A})] = 0$ .  $\square$

## C State Evolution for Optimal OAMP (Proposition 1)

Under the assumption  $\mathbb{E}\text{Var}[\mathbf{X}_\star | \mathbf{A}] \in (0, 1)$ , for the readers convenience, we recall that the optimal OAMP algorithm is given by (cf. (17)):

$$\mathbf{x}_t = \frac{1}{\sqrt{\omega_t}} \left( 1 + \frac{1}{\rho_t} \right) \cdot \Psi_\star(\mathbf{Y}; \rho_t) \cdot \bar{\varphi}(\mathbf{x}_{t-1}; \mathbf{a} | \omega_{t-1}). \quad (110)$$

The estimator returned by the optimal OAMP algorithm at iteration  $t$  is:

$$\widehat{\mathbf{x}}_t \stackrel{\text{def}}{=} \varphi(\mathbf{x}_t; \mathbf{a}|\omega_t). \quad (111)$$

In the above equations, the matrix denoiser  $\Psi_*$  used by the optimal OAMP algorithm is given by:

$$\Psi_*(\lambda; \rho) = 1 - \left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \frac{\phi(\lambda)}{\phi(\lambda) + \rho} \quad \forall \lambda \in \mathbb{R}, \rho \in (0, \infty), \quad (112)$$

where  $\Lambda \sim \mu$  and the function  $\phi : \mathbb{R} \mapsto \mathbb{R}$  was introduced in (7). The parameters  $\omega_t$  and  $\rho_t$  are computed using the following recursion, initialized with  $\omega_0 \stackrel{\text{def}}{=} 0$ :

$$\rho_t = \mathcal{F}_2(\omega_{t-1}), \quad \omega_t = \mathcal{F}_1(\rho_t) \quad \text{where} \quad \mathcal{F}_2(\omega) \stackrel{\text{def}}{=} \frac{1}{\text{dmmse}_\pi(\omega)} - 1, \quad \mathcal{F}_1(\rho) \stackrel{\text{def}}{=} 1 - \frac{\mathbb{E} \left[ \frac{1}{\phi(\Lambda) + \rho} \right]}{\mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right]}. \quad (113)$$

Let  $(\mathbf{X}_*, (\mathbf{X}_t)_{t \in \mathbb{N}}; \mathbf{A})$  denote the state evolution random variables associated with the above algorithm. We will rely on the following intermediate lemma.

**Lemma 9.** *We have,*

1. *The function  $\mathcal{F}_1$  is a continuous non-decreasing function which maps  $(0, \infty)$  to  $[0, 1)$  and satisfies  $\lim_{\rho \rightarrow \infty} \mathcal{F}_1(\rho) < 1$ .*
2. *The function  $\mathcal{F}_2$  is a continuous non-decreasing function which maps  $[0, 1)$  to  $(0, \infty)$  and satisfies  $\lim_{\omega \rightarrow 1} \mathcal{F}_2(\omega) = \infty$ .*

We defer the proof of this lemma to the end of this appendix (Appendix C.1), and present the proof of Proposition 1.

*Proof of Proposition 1.* We prove each of the two claims below.

**Proof of Claim (1).** We will show by induction that:

$$\omega_t \in [0, 1), \quad \rho_t \in (0, \infty), \quad (\mathbf{X}_*, \mathbf{X}_t; \mathbf{A}) \text{ form a scalar Gaussian channel with SNR } \omega_t. \quad (114)$$

Assuming this claim, we can compute the asymptotic MSE of the estimator  $\widehat{\mathbf{x}}_t = \varphi(\mathbf{x}_t; \mathbf{a}|\omega_t)$  returned by the OAMP algorithm:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{x}}_t - \mathbf{x}_*\|^2}{N} \stackrel{\text{Thm. 1}}{=} \mathbb{E} |\varphi(\mathbf{X}_t; \mathbf{A}|\omega_t) - \mathbf{X}_*|^2 \stackrel{\text{(a)}}{=} \text{mmse}_\pi(\omega_t).$$

In the above display, step (a) follows by observing that  $\varphi(\mathbf{X}_t; \mathbf{A}|\omega_t)$  is the MMSE estimator for the scalar Gaussian channel  $(\mathbf{X}_*, \mathbf{X}_t; \mathbf{A})$  (which has SNR  $\omega_t$ ). We now focus on proving (114) by induction. We assume that the claim (114) holds at step  $t$  as the induction hypothesis and show that it also holds at step  $t+1$ . Indeed, since  $\omega_t \in [0, 1)$ , Lemma 9 guarantees that  $\rho_{t+1} = \mathcal{F}_2(\omega_t) \in (0, \infty)$  and  $\omega_{t+1} = \mathcal{F}_1(\rho_{t+1}) \in [0, 1)$ . Next, we verify that  $(\mathbf{X}_*, \mathbf{X}_{t+1}; \mathbf{A})$  form a scalar Gaussian channel with SNR  $\omega_{t+1}$ . From Definition 4, we know that:

$$\mathbf{X}_{t+1} | (\mathbf{X}_*; \mathbf{A}) \sim \mathcal{N}(\beta_{t+1} \mathbf{X}_*, \sigma_{t+1}^2),$$

where:

$$\beta_{t+1} = \frac{\mathbb{E}[\mathbf{X}_* \cdot \bar{\varphi}(\mathbf{X}_t; \mathbf{A}|\omega_t)] \cdot \mathbb{E}[\Psi_*(\Lambda_\nu; \rho_{t+1})]}{\sqrt{\omega_{t+1}}(1 - \text{dmmse}_\pi(\omega_{t+1}))},$$

$$\sigma_{t+1}^2 = \frac{\{\mathbb{E}[\mathbf{X}_* \cdot \bar{\varphi}(\mathbf{X}_t; \mathbf{A}|\omega_t)]\}^2 \cdot \text{Var}[\Psi_*(\Lambda_\nu; \rho_{t+1})]}{\omega_{t+1}(1 - \text{dmmse}_\pi(\omega_{t+1}))^2} + \frac{\{\mathbb{E}[\bar{\varphi}^2(\mathbf{X}_t; \mathbf{A}|\omega_t)] - (\mathbb{E}[\mathbf{X}_* \cdot \bar{\varphi}(\mathbf{X}_t; \mathbf{A}|\omega_t)])^2\} \cdot \mathbb{E}[\Psi_*^2(\Lambda; \rho_{t+1})]}{\omega_{t+1}(1 - \text{dmmse}_\pi(\omega_{t+1}))^2}$$

In the above equations  $\Lambda \sim \mu$  and  $\Lambda_\nu \sim \nu$ . We begin by simplifying  $\beta_{t+1}$ . By the induction hypothesis,  $(\mathbf{X}_*, \mathbf{X}_t; \mathbf{A})$  forms a Gaussian channel with SNR  $\omega_t$ . A convenient property of the DMMSE estimator for a scalar Gaussian channel is the following identity (see Lemma 2 in Appendix A.2):

$$\mathbb{E}[\mathbf{X}_* \cdot \bar{\varphi}(\mathbf{X}_t; \mathbf{A}|\omega_t)] = \mathbb{E}[\bar{\varphi}^2(\mathbf{X}_t; \mathbf{A}|\omega_t)] = 1 - \text{dmmse}_\pi(\omega_t). \quad (115)$$

To compute  $\mathbb{E}[\Psi_*(\Lambda_\nu; \rho_{t+1})]$ , we consider the Lebesgue decomposition of  $\nu$  into the absolutely continuous part  $\nu_{\parallel}$  and the singular part  $\nu_{\perp}$ :

$$\begin{aligned} \mathbb{E}[\Psi_*(\Lambda_\nu; \rho_{t+1})] &\stackrel{(112)}{=} 1 - \left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \mathbb{E} \left[ \frac{\phi(\Lambda_\nu)}{\phi(\Lambda_\nu) + \rho_{t+1}} \right] \\ &\stackrel{(a)}{=} 1 - \left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \int_{\mathbb{R}} \frac{\phi(\lambda)}{\phi(\lambda) + \rho_{t+1}} \nu_{\parallel}(\mathrm{d}\lambda) \\ &\stackrel{(b)}{=} 1 - \frac{\mathbb{E} \left[ \frac{1}{\phi(\Lambda) + \rho_{t+1}} \right]}{\mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho_{t+1}} \right]} \stackrel{(113)}{=} \omega_{t+1} \end{aligned} \quad (116)$$

In the above display, step (a) follows by recalling that  $\phi(\lambda) = 0$  for  $\nu_{\perp}$ -a.e.  $\lambda$  (Lemma 1, item (2)) and step (b) follows from recalling that the density of  $\nu_{\parallel}$  is given by  $\mu(\cdot)/\phi(\cdot)$  (Lemma 1, item (3)). Plugging in the identities (115) and (116) into the expression for  $\beta_{t+1}$ , we obtain:

$$\beta_{t+1} = \frac{(1 - \text{dmmse}_\pi(\omega_t)) \cdot \omega_{t+1}}{\sqrt{\omega_{t+1}}(1 - \text{dmmse}_\pi(\omega_{t+1}))} = \sqrt{\omega_{t+1}}. \quad (117)$$

Similarly we can also simplify  $\sigma_{t+1}^2$  by exploiting the identities (115) and (116):

$$\sigma_{t+1}^2 \stackrel{(115)}{=} \frac{\text{Var}[\Psi_*(\Lambda_\nu; \rho_{t+1})] + \frac{\text{dmmse}_\pi(\omega_{t+1})}{1 - \text{dmmse}_\pi(\omega_{t+1})} \cdot \mathbb{E}[\Psi_*^2(\Lambda; \rho_{t+1})]}{\omega_{t+1}} \quad (118a)$$

$$\stackrel{(113)}{=} \frac{\mathbb{E}[\Psi_*^2(\Lambda_\nu; \rho_{t+1})] - (\mathbb{E}[\Psi_*(\Lambda_\nu; \rho_{t+1})])^2 + \frac{1}{\rho_{t+1}} \cdot \mathbb{E}[\Psi_*^2(\Lambda; \rho_{t+1})]}{\omega_{t+1}} \quad (118b)$$

$$\stackrel{(116)}{=} \frac{\mathbb{E}[\Psi_*^2(\Lambda_\nu; \rho_{t+1})] + \frac{1}{\rho_{t+1}} \cdot \mathbb{E}[\Psi_*^2(\Lambda; \rho_{t+1})] - \omega_{t+1}^2}{\omega_{t+1}}. \quad (118c)$$

After a straightforward (but tedious) calculation analogous to the one used to derive (116), we obtain the identity:

$$\mathbb{E}[\Psi_*^2(\Lambda_\nu; \rho_{t+1})] + \frac{1}{\rho_{t+1}} \cdot \mathbb{E}[\Psi_*^2(\Lambda; \rho_{t+1})] = \omega_{t+1}, \quad (118d)$$

which implies that  $\sigma_{t+1}^2 = 1 - \omega_{t+1}$ . Hence,

$$\mathbf{X}_{t+1} | (\mathbf{X}_*, \mathbf{A}) \sim \mathcal{N}(\sqrt{\omega_{t+1}} \cdot \mathbf{X}_*, 1 - \omega_{t+1}), \quad (119)$$

which verifies that  $(\mathbf{X}_*, \mathbf{X}_{t+1}; \mathbf{A})$  forms a Gaussian channel with SNR  $\omega_{t+1}$ .

**Proof of Claim (2).** Recall that the sequences  $\{\omega_t\}_{t \in \mathbb{N}}$  and  $\{\rho_t\}_{t \in \mathbb{N}}$  are generated via the recursion:

$$\rho_t = \mathcal{F}_2(\omega_{t-1}) \stackrel{\text{def}}{=} \frac{1}{\text{dmmse}_\pi(\omega_{t-1})} - 1, \quad (120a)$$

$$\omega_t = \mathcal{F}_1(\rho_t) \stackrel{\text{def}}{=} 1 - \left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho_t} \right] \right)^{-1} \cdot \mathbb{E} \left[ \frac{1}{\phi(\Lambda) + \rho_t} \right], \quad (120b)$$

where  $\omega_0 = 0$ . We first prove the monotonicity of  $\{\omega_t\}_{t \in \mathbb{N}}$  using an induction argument. Note that  $\omega_t = \mathcal{F}_1 \circ \mathcal{F}_2(\omega_{t-1})$ . First, note that  $\omega_1 \geq \omega_0 = 0$  (recall we have already shown  $\omega_t \geq 0$  for all  $t \in \mathbb{N}$ ). Since the composite function  $\mathcal{F}_1 \circ \mathcal{F}_2$  is non-decreasing (Lemma 9), it follows that  $\mathcal{F}_1 \circ \mathcal{F}_2(\omega_1) \geq \mathcal{F}_1 \circ \mathcal{F}_2(\omega_0)$ ,

namely,  $\omega_2 \geq \omega_1$ . Repeating this argument proves that  $\{\omega_t\}_{t \in \mathbb{N}}$  is non-decreasing. Moreover,  $\{\omega_t\}_{t \in \mathbb{N}}$  is bounded sequence taking values in  $[0, 1]$  (proved in the first claim). Hence,  $\{\omega_t\}_{t \in \mathbb{N}}$  converges to a limit point  $\omega_* \in [0, 1]$  which satisfies the equation  $\omega_* = \mathcal{F}_1 \circ \mathcal{F}_2(\omega_*)$ . We eliminate the possibility that  $\omega_* = 1$  by observing that Lemma 9 guarantees that:

$$\lim_{\omega \rightarrow 1} \mathcal{F}_1 \circ \mathcal{F}_2(\omega) = \lim_{\rho \rightarrow \infty} \mathcal{F}_1(\rho) < 1.$$

Since the sequence  $\{\rho_t\}_{t \in \mathbb{N}}$  is obtained by applying a continuous and non-decreasing function  $\mathcal{F}_2$  to a non-decreasing and convergent sequence  $\{\omega_t\}_{t \in \mathbb{N}}$ , we conclude that  $\{\rho_t\}_{t \in \mathbb{N}}$  is non-decreasing and converges to a limit  $\rho_* \stackrel{\text{def}}{=} \mathcal{F}_2(\omega_*) \in (0, \infty)$ . Observe that the limit points satisfy the equations:  $\omega_* = \mathcal{F}_1 \circ \mathcal{F}_2(\omega_*)$ ,  $\rho_* = \mathcal{F}_2(\omega_*)$ , which can be expressed as  $\omega_* = \mathcal{F}_1(\rho_*)$ ,  $\rho_* = \mathcal{F}_2(\omega_*)$  as claimed. This concludes the proof of the proposition.  $\square$

## C.1 Proof of Lemma 9

*Proof of Lemma 9.* We prove the claims made about  $\mathcal{F}_1$  and  $\mathcal{F}_2$  below.

**Properties of  $\mathcal{F}_1$ .** Recall that:

$$\mathcal{F}_1(\rho) \stackrel{\text{def}}{=} 1 - \frac{\mathbb{E} \left[ \frac{1}{\phi(\Lambda) + \rho} \right]}{\mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right]} \quad \forall \rho \in (0, \infty).$$

The continuity of  $\mathcal{F}_1(\rho)$  follows from the Dominated Convergence Theorem. To verify the monotonicity of  $\mathcal{F}_1$ , we calculate its derivative:

$$\mathcal{F}'_1(\rho) = \frac{\mathbb{E} \left[ \left( \frac{1}{\phi(\Lambda) + \rho} \right)^2 \right] \cdot \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] - \mathbb{E} \left[ \frac{1}{\phi(\Lambda) + \rho} \right] \cdot \mathbb{E} \left[ \frac{\phi(\Lambda)}{(\phi(\Lambda) + \rho)^2} \right]}{\left( \mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^2} \geq 0,$$

where the inequality follows from Chebyshev's association inequality [14, Theorem 2.14], which states that for any two random variables  $(X, Z)$  such that  $Z \geq 0$  and any two non-decreasing functions  $f, g$  we have:

$$\mathbb{E}[Z] \mathbb{E}[f(X)g(X)Z] \geq \mathbb{E}[f(X)Z] \cdot \mathbb{E}[g(X)Z].$$

Indeed, applying this inequality with  $X = \phi(\Lambda)$ ,  $f(X) = X$ ,  $g(X) = X + \rho$  and  $Z = 1/(\phi(\Lambda) + \rho)^2$  shows that the numerator of  $\mathcal{F}'_1(\rho)$  is non-negative. Finally, we check that  $\mathcal{F}_1 : (0, \infty) \mapsto [0, 1]$ . Notice from (C.1) that for any  $\rho \in (0, \infty)$ ,  $\mathcal{F}_1(\rho) < 1$ . On the other hand, since  $\mathcal{F}_1(\rho) \geq 0$  since, thanks to the monotonicity of  $\mathcal{F}_1$  we have:

$$\mathcal{F}_1(\rho) \geq \mathcal{F}_1(0) = 1 - \mathbb{E} \left[ \frac{1}{\phi(\Lambda)} \right] \stackrel{(a)}{=} 1 - \nu_{\parallel}(\mathbb{R}) \geq 0.$$

In the above display, the step (a) follows by recalling from Lemma 1 (item (3)) that the density of  $\nu_{\parallel}$ , the absolutely continuous part of  $\nu$  is given by  $\mu(\cdot)/\phi(\cdot)$  and by observing that  $\nu_{\parallel}(\mathbb{R}) \leq \nu(\mathbb{R}) = 1$ . Last we compute:

$$\lim_{\rho \rightarrow \infty} \mathcal{F}_1(\rho) = 1 - \lim_{\rho \rightarrow \infty} \frac{\mathbb{E} \left[ \frac{1}{\phi(\Lambda) + \rho} \right]}{\mathbb{E} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right]} = 1 - \lim_{\rho \rightarrow \infty} \frac{\mathbb{E} \left[ \frac{1}{\rho^{-1} \cdot \phi(\Lambda) + 1} \right]}{\mathbb{E} \left[ \frac{\phi(\Lambda)}{\rho^{-1} \cdot \phi(\Lambda) + 1} \right]} = 1 - \frac{1}{\mathbb{E}[\phi(\Lambda)]} < 1,$$

as claimed.

**Properties of  $\mathcal{F}_2$ .** Recall that:

$$\mathcal{F}_2(\omega) \stackrel{\text{def}}{=} \frac{1}{\text{dmmse}_\pi(\omega)} - 1.$$

The continuity and monotonicity of the DMMSE function  $\text{dmmse}_\pi(\cdot)$  (see From Lemma 4 in Appendix A.2) implies that  $\mathcal{F}_2$  continuous and non-decreasing. Next, we verify that  $\mathcal{F}_2 : [0, 1) \mapsto (0, \infty)$ . Notice that for any  $\omega \in [0, 1)$ ,  $\mathcal{F}_2(\omega) \geq \mathcal{F}_2(0) = (\text{dmmse}_\pi(0))^{-1} - 1 = (\text{mmse}_\pi(0))^{-1} - 1 > 0$  (recall that  $\text{mmse}_\pi(0) \in (0, 1)$  was one of the hypothesis assumed in Proposition 1). Moreover, for any  $\omega < 1$ ,  $\mathcal{F}_2(\omega) < \infty$  since:

$$\text{dmmse}_\pi(\omega) \geq \text{mmse}_\pi(\omega) \stackrel{\text{(a)}}{>} \text{mmse}_\pi(1) = 0,$$

where (a) follows from the fact that  $\text{mmse}_\pi(\cdot)$  is strictly decreasing if  $\text{mmse}_\pi(0) \in (0, 1)$  (see Fact 1 in Appendix A.2). Hence,  $\mathcal{F}_2 : [0, 1) \mapsto (0, \infty)$ . This concludes the proof of the lemma.  $\square$

## D Replica Predictions for Quartic Potential (Proposition 2)

We begin with some preliminary results about the trace ensemble. The Proof of Proposition 2 will be presented in Appendix D.2.

### D.1 Preliminaries of Trace Ensemble

When the potential function  $V : \mathbb{R} \mapsto \mathbb{R}$  is a polynomial, the function  $\phi(\lambda)$  in (7) corresponding to the trace ensemble is a polynomial with degree one less than  $V$ , for  $\lambda \in \text{Supp}(\mu)$ . Lemma 10 below summarizes some useful properties of  $\phi$ .

**Lemma 10** (The function  $\phi(\lambda)$  for trace ensemble). *Consider the trace ensemble defined in (20). Assume that  $V : \mathbb{R} \mapsto \mathbb{R}$  is a degree  $k$  polynomial and independent of  $N$ , where  $k \geq 2$ . Let  $\mu$  be the limiting spectral measure corresponding to the trace ensemble with  $V$ . Let  $\phi : \mathbb{R} \mapsto \mathbb{R}$  be defined as in (7). Then, the following holds:*

1. The spectral measure  $\mu$  is absolutely continuous with density function:

$$\mu(\lambda) = \frac{1}{2\pi} \sqrt{(4Q(\lambda) - (V'(\lambda))^2)_+}, \quad \forall \lambda \in \mathbb{R}, \quad (121)$$

where  $(x)_+ = \max\{x, 0\}$  and  $Q(\lambda)$  satisfies

$$Q(\lambda) = \int \frac{V'(\lambda) - V'(t)}{\lambda - t} \mu(t) dt. \quad (122)$$

The support of  $\mu$  is a finite union of finite intervals, and the density function  $\mu(\lambda)$  is Holder continuous.

2. Define  $\phi_{\text{poly}} : \mathbb{R} \mapsto \mathbb{R}$  as

$$\phi_{\text{poly}}(\lambda) \stackrel{\text{def}}{=} 1 - \theta \cdot V'(\lambda) + \theta^2 \cdot Q(\lambda), \quad \forall \lambda \in \mathbb{R}. \quad (123)$$

Then,  $\phi_{\text{poly}}$  is a degree  $k - 1$  polynomial and

$$\phi(\lambda) = \phi_{\text{poly}}(\lambda), \quad \forall \lambda \in \text{Supp}(\mu), \quad (124)$$

where  $\phi(\lambda)$  is defined in (7).

*Proof.* The expression of the density function of  $\mu$  can be found in, e.g., [61, Theorem 11.2.1]. Since  $V(x)$  is a degree- $k$  polynomial,  $V'(\lambda)$  is a degree  $k - 1$  polynomial and  $Q(\lambda)$  defined in (122) is a degree  $k - 2$  polynomial. Hence, the support  $\text{Supp}(\mu) = \{\lambda | 4Q(\lambda) - (V'(\lambda))^2 \geq 0\}$  must be a union of finite intervals [61, Eq. (11.2.14)]. Furthermore,  $\mu(\lambda)$  is a composition of Holder continuous functions and is therefore Holder



continuous. To prove the second claim, we note that the Hilbert transform of  $\mu(\lambda)$  is related to the potential function via [61, Theorem 11.2.1]:

$$\pi \cdot \mathcal{H}_\mu(\lambda) = \frac{1}{2} V'(\lambda), \quad \forall \lambda \in \text{Supp}(\mu). \quad (125)$$

Substituting (121) and (125) into (7), we get

$$\phi(\lambda) \stackrel{\text{def}}{=} [1 - \theta \pi \mathcal{H}_\mu(\lambda)]^2 + \pi^2 \theta^2 \mu^2(\lambda) = 1 - \theta \cdot V'(\lambda) + \theta^2 \cdot Q(\lambda), \quad \forall \lambda \in \text{Supp}(\mu). \quad (126)$$

To check that  $\phi_{\text{poly}}$  is a degree  $k-1$  polynomial, we note that  $V'(\lambda)$  is a degree  $k-1$  polynomial and  $Q(\lambda)$  in (121) is a degree  $k-2$  polynomial.  $\square$

Lemma 10 shows that  $\phi(\lambda)$  agrees with a polynomial  $\phi_{\text{poly}}(\lambda)$  for  $\lambda \in \text{Supp}(\mu)$ , but they are generally different outside of the support. The following proposition shows that we can replace  $\phi(\cdot)$  by the polynomial  $\phi_{\text{poly}}(\cdot)$  in the optimal OAMP algorithm with no performance loss asymptotically.

**Proposition 6** (A variant of optimal OAMP). *Assume that  $\mathbf{W}$  is drawn from the trace ensemble with a polynomial potential. Consider a variant of the optimal OAMP algorithm in (17c) where  $\phi(\cdot)$  in the matrix denoiser is replaced by  $\phi_{\text{poly}}(\cdot)$ . Then, Proposition 1 still holds for this OAMP algorithm.*

*Proof.* First of all, the variant of the optimal OAMP algorithm satisfies the required trace-free condition:  $\mathbb{E}[\Psi_{\text{poly}}(\Lambda)] = 0$ , where

$$\Psi_{\text{poly}}(\lambda) \stackrel{\text{def}}{=} 1 - \left( \mathbb{E} \left[ \frac{\phi_{\text{poly}}(\Lambda)}{\phi_{\text{poly}}(\Lambda) + \rho} \right] \right)^{-1} \cdot \frac{\phi_{\text{poly}}(\lambda)}{\phi_{\text{poly}}(\lambda) + \rho}, \quad \forall \lambda \in \mathbb{R}, \rho \in (0, \infty).$$

Therefore, the state evolution framework applies to this variant of optimal OAMP. To prove that the state evolution equation for this new variant of OAMP is identical to the original one, it is not difficult to show that it suffices to check the following identity:

$$\frac{(\mathbb{E}[\Psi_{\text{poly}}(\Lambda_\nu)])^2}{\mathbb{E}[\Psi_{\text{poly}}^2(\Lambda_\nu)] + \rho^{-1} \cdot \mathbb{E}[\Psi_{\text{poly}}^2(\Lambda)]} = \frac{(\mathbb{E}[\Psi(\Lambda_\nu)])^2}{\mathbb{E}[\Psi^2(\Lambda_\nu)] + \rho^{-1} \cdot \mathbb{E}[\Psi^2(\Lambda)]}.$$

To this end, we shall verify  $\mathbb{E}[\Psi_{\text{poly}}(\Lambda_\nu)] = \mathbb{E}[\Psi(\Lambda_\nu)]$ ,  $\mathbb{E}[\Psi_{\text{poly}}^2(\Lambda_\nu)] = \mathbb{E}[\Psi^2(\Lambda_\nu)]$ , and  $\mathbb{E}[\Psi_{\text{poly}}(\Lambda)] = \mathbb{E}[\Psi(\Lambda)]$ . From item 3 of Lemma 10,  $\phi(\lambda) = \phi_{\text{poly}}(\lambda)$  for  $\lambda \in \text{Supp}(\mu)$ . Hence,  $\mathbb{E}[\Psi_{\text{poly}}(\Lambda)] = \mathbb{E}[\Psi(\Lambda)]$ . It remains to prove  $\mathbb{E}[\Psi_{\text{poly}}(\Lambda_\nu)] = \mathbb{E}[\Psi(\Lambda_\nu)]$ . (The proof for  $\mathbb{E}[\Psi_{\text{poly}}^2(\Lambda_\nu)] = \mathbb{E}[\Psi^2(\Lambda_\nu)]$  is similar.) From Lemma 1,

$$\begin{aligned} \mathbb{E}[\Psi_{\text{poly}}(\Lambda_\nu)] &= \int_{\mathbb{R}} \Psi_{\text{poly}}(\lambda) \nu_\perp(d\lambda) + \int_{\mathbb{R}} \Psi_{\text{poly}}(\lambda) \nu_\parallel(d\lambda) \\ &= \int_{\mathbb{R}} \Psi_{\text{poly}}(\lambda) \nu_\perp(d\lambda) + \int_{\mathbb{R}} \Psi(\lambda) \nu_\parallel(d\lambda), \end{aligned}$$

where the second equality is from the fact that  $\nu_\parallel$  has the same support as  $\mu$ , and  $\phi(\lambda) = \phi_{\text{poly}}(\lambda)$  for  $\lambda \in \text{Supp}(\mu)$ . It remains to prove

$$\int_{\mathbb{R}} \Psi_{\text{poly}}(\lambda) \nu_\perp(d\lambda) = \int_{\mathbb{R}} \Psi(\lambda) \nu_\perp(d\lambda).$$

Lemma 1 shows that the measure  $\nu_\perp$  concentrates on the set  $S \stackrel{\text{def}}{=} \{x \in \mathbb{R} : \phi(x) = 0\}$ . It is therefore sufficient to prove  $\phi_{\text{poly}}(\lambda) = \phi(\lambda)$  for  $\lambda \in S$ . From (7), we have

$$\pi \mathcal{H}_\mu(\lambda) = \frac{1}{\theta} \quad \text{and} \quad \mu(\lambda) = 0, \quad \forall \lambda \in S.$$

On the other hand, the Stieltjes transform of  $\mu$  for the trace ensemble with a polynomial potential  $V$  satisfies [61, Theorem 11.2.1]

$$\mathcal{S}_\mu^2(z) - V'(z) \mathcal{S}_\mu(z) + Q(z) = 0, \quad \forall z \in \mathbb{C} \setminus \mathbb{R},$$

where the function  $Q$  is defined in (122). By the Sokhotsky-Plemelj formula [61, Section 2.1],

$$\lim_{\epsilon \downarrow 0} \mathcal{S}_\mu(\lambda + i\epsilon) = \pi \mathcal{H}_\mu(\lambda) - i \cdot \pi \mu(\lambda), \quad \forall \lambda \in \mathbb{R}.$$

Combining the above three equations leads to the following equation:

$$\frac{1}{\theta^2} - V'(\lambda) \cdot \frac{1}{\theta} + Q(\lambda) = 0, \quad \forall \lambda \in S.$$

Therefore,

$$\phi_{\text{poly}}(\lambda) \stackrel{\text{def}}{=} 1 - \theta \cdot V'(\lambda) + \theta^2 \cdot Q(\lambda) = 0, \quad \forall \lambda \in S.$$

Hence,  $\phi_{\text{poly}}(\lambda) = \phi(\lambda)$  for  $\lambda \in S$ , which completes the proof.  $\square$

We next define a function closely related to  $\phi_{\text{poly}}$ , which recovers the matrix processing function used in the Bayes-optimal AMP (BAMP) algorithm proposed in Barbier et al. [8].

**Definition 7** (Definition of  $J(\lambda)$ ). Let  $k \geq 2$  and  $V : \mathbb{R} \mapsto \mathbb{R}$  be a degree  $k$  polynomial. Let the degree  $k-1$  polynomial  $\phi_{\text{poly}}$  in (123) be represented as

$$\phi_{\text{poly}}(\lambda) = C_\phi + \sum_{i=1}^{k-1} C_{\phi,i} \lambda^i, \quad \forall \lambda \in \mathbb{R}. \quad (127a)$$

We define the function  $J : \mathbb{R} \mapsto \mathbb{R}$  as

$$J(\lambda) \stackrel{\text{def}}{=} - \sum_{i=1}^{k-1} C_{\phi,i} \lambda^i, \quad \forall \lambda \in \mathbb{R}. \quad (127b)$$

**Examples of  $\phi_{\text{poly}}(\lambda)$  and  $J(\lambda)$ .** We now provide explicit formulas for the functions  $\phi_{\text{poly}}(\lambda)$  and  $J(\lambda)$  for three polynomial potential functions considered in Barbier et al. [8].

1. *Quadratic potential:*

$$V(\lambda) = \frac{\lambda^2}{2}. \quad (128a)$$

This corresponds to the Gaussian orthogonal ensemble (GOE). The density function and its Hilbert transform are  $\mu(\lambda) = \frac{1}{2\pi} \sqrt{4 - \lambda^2}$  and  $\mathcal{H}_\mu(\lambda) = V'(\lambda)/2\pi$ ,  $\forall \lambda \in [-2, 2]$ . From (123) and (127b), we have

$$\phi_{\text{poly}}(\lambda) = 1 + \theta^2 - \theta \cdot \lambda, \quad (128b)$$

$$J(\lambda) = \theta \cdot \lambda. \quad (128c)$$

2. *Quartic potential:*

$$V(\lambda) = \frac{\gamma \lambda^2}{2} + \frac{\kappa \lambda^4}{4}, \quad (129a)$$

where  $\gamma \in [0, 1]$  is a parameter and  $\kappa = \kappa(\gamma) = \left(8 - 9\gamma + \sqrt{64 - 144\gamma + 108\gamma^2 - 27\gamma^3}\right) / 27 \in [0, \frac{16}{27}]$ . This choice of  $\kappa$  ensures that the limiting eigenvalue distribution  $\mu$  has unit variance. Its density function is given by  $\mu(\lambda) = (\gamma + 2a^2\kappa + \kappa\lambda^2) \sqrt{4a^2 - \lambda^2} / (2\pi)$ ,  $\forall \lambda \in [-2a, 2a]$ , where  $a^2 := (\sqrt{\gamma^2 + 12\kappa - \gamma}) / (6\kappa)$ . The functions  $\phi_{\text{poly}}$  and  $J$  read [66]

$$\phi_{\text{poly}}(\lambda) = \theta^2 a^2 (\gamma + 2a^2 \kappa)^2 + 1 - \theta (\gamma \lambda - \theta \kappa \lambda^2 + \kappa \lambda^3), \quad (129b)$$

$$J(\lambda) = \theta (\gamma \lambda - \theta \kappa \lambda^2 + \kappa \lambda^3). \quad (129c)$$

3. *Purely sextic potential:*

$$V(\lambda) = \frac{\xi\lambda^6}{6}, \quad (130a)$$

where  $\xi = \frac{27}{80}$ . Again, this value of  $\xi$  ensures that the spectral measure  $\mu$  has unit variance. Its density function is given by  $\mu(\lambda) = (6b^4\xi + 2b^2\xi\lambda^2 + \xi\lambda^4)\sqrt{4b^2 - \lambda^2}/(2\pi)$ ,  $\forall \lambda \in [-2b, 2b]$ , where  $b^2 := 2/3$ . The functions  $\phi_{\text{poly}}$  and  $J$  are given by [8]

$$\phi_{\text{poly}}(\lambda) = \frac{27}{50}\theta^2 + 1 - \theta(-\theta\xi\lambda^2 - \theta\xi\lambda^4 + \xi\lambda^5), \quad (130b)$$

$$J(\lambda) = \theta(-\theta\xi\lambda^2 - \theta\xi\lambda^4 + \xi\lambda^5). \quad (130c)$$

Fig. 5 plots the functions  $\phi(\cdot)$  and  $\phi_{\text{poly}}(\cdot)$  corresponding to the quartic potential. In this setup, the measure  $\mu$  is supported on the interval  $[-2a, 2a]$  with  $2a \approx 1.7$ , and it can be shown that  $\nu$  has a point mass (corresponding to the outlying eigenvalue of  $\mathbf{Y}$ ) at  $\lambda_o \approx 2.3$ . We see from Fig. 5 that  $\phi(\cdot)$  and  $\phi_{\text{poly}}(\cdot)$  coincide inside the support of  $\mu$  but differ outside. Moreover, the curves of  $\phi(\cdot)$  and  $\phi_{\text{poly}}(\cdot)$  intersect precisely at the location of the outlying eigenvalue:  $\phi(\lambda_o) = \phi_{\text{poly}}(\lambda_o) = 0$ ; see the proof of Proposition 6.

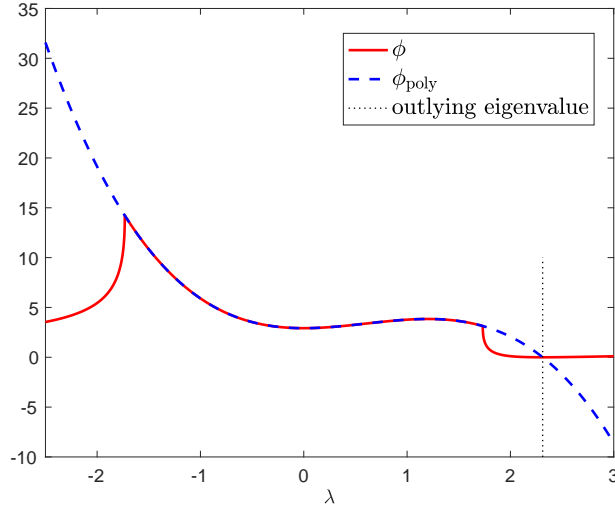


Figure 5: Plot of  $\phi$  and  $\phi_{\text{poly}}$  for the quartic ensemble with  $\gamma = 0$ .  $\theta = 1.8$ .  $\text{Supp}(\mu) = [-2a, 2a]$  with  $2a \approx 1.7$ . The singular part of the measure  $\nu$  contains a single point mass at  $\lambda_o \approx 2.3$ .

## D.2 Proof of Proposition 2

We will show that any solution  $(\rho, \omega)$  to the state evolution (SE) fixed point equation (17f) is a solution to the replica fixed point equation of (22). The proof for the other direction, namely, each solution to the replica fixed point equation is also a solution to the SE equations is similar and thus omitted.

We will introduce a series of manipulations of the SE equations (17f) that eventually lead to the replica equations (22), for the *quartic potential* case. Our proof consists of the following steps:

1. In Lemma 11, we rewrite (17f) as (131). Note that unlike the original SE equation (17f) which only involves the measure  $\mu$ , the new formulation (131) involves both  $\mu$  and  $\nu$ . The introduction of the new formulation (131) is somewhat ad hoc, and is only for the purpose of deriving the replica equations as presented in [8].
2. We recognize that the term in (131a) that involves  $\nu$  is a resolvent operator for  $J(\Lambda_\nu)$ . It can be shown that the expectation of this term is the limit of  $\frac{1}{N}\mathbf{x}_*^\top(\rho\mathbf{I}_N + C_\phi\mathbf{I}_N - J(\mathbf{Y}))^{-1}\mathbf{x}_*$  as  $N \rightarrow \infty$ . In the quartic potential case, the function  $J(\cdot)$  is a degree three polynomial; see (129c). It is possible

to invoke the matrix inversion lemma to express the term into a resolvent in terms of  $J(\mathbf{W})$  plus a low-rank perturbation term, and subsequently calculate the asymptotics of each term involved. This step is summarized in Lemma 13.

3. Using Lemma 13 and some algebraic manipulations, we finally arrive at the replica equations we aim to prove in Proposition 2. The derivations are postponed to Section D.2.2.

### D.2.1 Auxiliary results

To compare with the replica calculations in [8], we rewrite the fixed point equation (18) into an equivalent form in the following lemma.

**Lemma 11** (Alternative fixed point equations for SE). *Assume that  $\mathbb{E}\text{Var}[\mathbf{X}_\star|\mathbf{A}] \in (0, 1)$ . Let  $(\omega, \rho)$  be defined as in (18). Then,  $(\omega, \rho)$  satisfies the following equations:*

$$\frac{\omega}{1-\omega} = \left(1 - \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right) \cdot \left( \left( \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right)^{-1} - \left( \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda_\nu)} \right] \right)^{-1} \right), \quad (131a)$$

$$\mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] = \text{mmse}_\pi(\omega), \quad (131b)$$

where  $\Lambda \sim \mu$  and  $\Lambda_\nu \sim \nu$ .

*Proof.* We prove the two equations separately.

**Proof of (131a).** We rewrite the first equation of (18) as

$$\frac{\omega}{1-\omega} = \left( \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right)^{-1} - (1 + \rho) \quad (132a)$$

$$= \left(1 - \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right) \cdot \left( \left( \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right)^{-1} - \rho \cdot \left(1 - \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right)^{-1} \right). \quad (132b)$$

Comparing (131a) and (132a), we see that it suffices to prove the following

$$\mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda_\nu)} \right] = \frac{1}{\rho} \cdot \left(1 - \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right), \quad \forall \rho > 0. \quad (133)$$

To prove (133), we note that

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda_\nu)} \right] &= \int_{\mathbb{R}} \frac{1}{\rho + \phi(\lambda)} \nu(d\lambda) \\ &= \int_{\mathbb{R}} \frac{1}{\rho + \phi(\lambda)} \nu_\perp(d\lambda) + \int_{\mathbb{R}} \frac{1}{\rho + \phi(\lambda)} \nu_\parallel(d\lambda) \\ &\stackrel{(a)}{=} \frac{1}{\rho} \cdot \nu_\perp(\mathbb{R}) + \int_{\mathbb{R}} \frac{1}{\rho + \phi(\lambda)} \nu_\parallel(d\lambda) \\ &= \frac{1}{\rho} \cdot (1 - \nu_\parallel(\mathbb{R})) + \int_{\mathbb{R}} \frac{1}{\rho + \phi(\lambda)} \nu_\parallel(d\lambda) \\ &\stackrel{(b)}{=} \frac{1}{\rho} \cdot \left(1 - \int_{\mathbb{R}} \frac{1}{\phi(\lambda)} \mu(d\lambda)\right) + \int_{\mathbb{R}} \frac{1}{\rho + \phi(\lambda)} \cdot \frac{1}{\phi(\lambda)} \mu(d\lambda) \\ &= \frac{1}{\rho} \cdot \left(1 - \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right), \quad \Lambda \sim \mu, \end{aligned}$$

where step (a) and (b) are due to items 2 and 3 of Lemma 1, respectively.

**Proof of (131b).** From the second equation of (18), we have

$$\rho = \frac{1}{\text{dmmse}_\pi(\omega)} - 1 \stackrel{(a)}{=} \frac{1}{\text{mmse}_\pi(\omega)} - \frac{1}{1-\omega},$$

where step (a) follows from Lemma 2. Hence,

$$\left(\rho + \frac{1}{1-\omega}\right)^{-1} = \text{mmse}_\pi(\omega). \quad (134)$$

Now recall the first equation of (18):

$$\omega = 1 - \left(\mathbb{E}\left[\frac{\phi(\Lambda)}{\phi(\Lambda) + \rho}\right]\right)^{-1} \cdot \mathbb{E}\left[\frac{1}{\phi(\Lambda) + \rho}\right].$$

By simple algebra, we have

$$\left(\rho + \frac{1}{1-\omega}\right)^{-1} = \mathbb{E}\left[\frac{1}{\phi(\Lambda) + \rho}\right]. \quad (135)$$

Combining (134) and (135) leads to (131b).  $\square$

Note that (131a) involves the term  $\mathbb{E}\left[\frac{1}{\rho + \phi(\Lambda_\nu)}\right]$ , which depends on the function  $\phi(\cdot)$ . To facilitate further analysis, we rewrite this term using the function  $\phi_{\text{poly}}(\lambda) \stackrel{\text{def}}{=} C_\phi(\lambda) + J(\lambda)$  defined in Appendix D.1; see Lemma 10 and Definition 7. Lemma 12 below summarizes this result. Its proof is similar to the proof of Proposition 6 and omitted.

**Lemma 12.** *Let  $\mu$  be the limiting spectral measure associated with the quartic potential. Let  $C_\phi$  and  $J(\cdot)$  be defined as in Definition (7). Then, the following holds for any  $\rho > 0$ :*

$$\mathbb{E}\left[\frac{1}{\rho + \phi(\Lambda_\nu)}\right] = \mathbb{E}\left[\frac{1}{\rho + C_\phi - J(\Lambda_\nu)}\right], \quad \Lambda_\nu \sim \nu. \quad (136)$$

The following lemma reformulate the term  $\mathbb{E}\left[\frac{1}{\rho + C_\phi - J(\Lambda_\nu)}\right]$ , which involves the measure  $\Lambda_\nu \sim \nu$ , into a (complicated) expression that only depends on  $\Lambda \sim \mu$ .

**Lemma 13.** *Consider the quartic potential for which the function  $J(\lambda)$  is given in (129c). Then, the following identity holds for any  $\rho > 0$ :*

$$\mathbb{E}\left[\frac{1}{\rho + C_\phi - J(\Lambda_\nu)}\right] = -\frac{1}{(\kappa\theta^2)^2 d_4 + \gamma\theta^2 + \frac{3\kappa\theta^2 d_2 + (\kappa\theta^2)^2(2d_1 d_3 - 3d_2^2) + (\kappa\theta^2)^3(d_2^3 + d_0 d_3^2 - 2d_1 d_2 d_3) - 1}{d_0 + \kappa\theta^2 d_1^2 - \kappa\theta^2 d_0 d_2}}, \quad (137)$$

where  $\Lambda_\nu \sim \nu$ , and  $\kappa$  and  $\gamma$  are parameters for the quartic potential function and  $\theta$  is the parameter for the observation model, and

$$d_i \stackrel{\text{def}}{=} \mathbb{E}\left[\frac{\Lambda^i}{\rho + \phi(\Lambda)}\right], \quad \forall i \in \mathbb{N} \cup \{0\}. \quad (138)$$

*Proof.* Recall that the probability measure  $\nu$  is the limit of the empirical measure  $\nu_N \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N \langle \mathbf{u}_i(\mathbf{Y}), \mathbf{x}_* \rangle^2 \delta_{\lambda_i(\mathbf{Y})}$ , where  $\{\mathbf{u}_i(\mathbf{Y})\}_{i \in [N]}$  and  $\{\lambda_i(\mathbf{Y})\}_{i \in [N]}$  are the eigenvectors and eigenvalues of  $\mathbf{Y}$ . By the convergence of  $\nu_N$  (Lemma 1), we have

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{x}_*^\top (c_* \mathbf{I}_N - J(\mathbf{Y}))^{-1} \mathbf{x}_* = \mathbb{E}\left[\frac{1}{c_* - J(\Lambda_\nu)}\right]. \quad (139)$$

where  $c_* \stackrel{\text{def}}{=} \rho + C_\phi$ . Recall that  $J(\cdot)$  is a degree three polynomial for quartic potential:

$$J(\lambda) = \theta(\gamma\lambda - \theta\kappa\lambda^2 + \kappa\lambda^3). \quad (140)$$

Based on the above definition of  $J(\cdot)$ , we can decompose  $J(\mathbf{Y})$  as  $J(\mathbf{W})$  plus a low-rank perturbation term:

$$\begin{aligned}
J(\mathbf{Y}) &= J\left(\frac{\theta \mathbf{x}_* \mathbf{x}_*^\top}{N} + \mathbf{W}\right) \\
&= \gamma \theta \mathbf{W} - \kappa \theta^2 \mathbf{W}^2 + \kappa \theta \mathbf{W}^3 + \frac{\gamma \theta^2}{N} \mathbf{x}_* \mathbf{x}_*^\top + \frac{1}{N} \kappa \theta^3 \mathbf{x}_*^\top \mathbf{W} \mathbf{x}_* \cdot \frac{\mathbf{x}_* \mathbf{x}_*^\top}{N} + \frac{1}{N} \kappa \theta^2 (\mathbf{W}^2 \mathbf{x}_* \mathbf{x}_*^\top + \mathbf{x} \mathbf{x}^\top \mathbf{W}^2 + \mathbf{W} \mathbf{x}_* \mathbf{x}_*^\top \mathbf{W}) \\
&= J(\mathbf{W}) + \frac{1}{N} \mathbf{x}_* \mathbf{x}_*^\top (\gamma \theta^2 \mathbf{I}_N + \kappa \theta^2 \mathbf{W}) + \frac{1}{N} \kappa \theta^2 \mathbf{W} \mathbf{x}_* \mathbf{x}_*^\top \mathbf{W} + \frac{1}{N} \kappa \theta^2 \mathbf{W}^2 \mathbf{x}_* \mathbf{x}_*^\top + \frac{1}{N} \kappa \theta^3 \mathbf{x}_*^\top \mathbf{W} \mathbf{x}_* \cdot \frac{\mathbf{x}_* \mathbf{x}_*^\top}{N} \\
&= J(\mathbf{W}) + \frac{1}{N} \mathbf{P} \mathbf{Q}^\top,
\end{aligned}$$

where

$$\mathbf{P} \stackrel{\text{def}}{=} [\mathbf{x}_*, \mathbf{W} \mathbf{x}_*, \mathbf{W}^2 \mathbf{x}_*] \in \mathbb{R}^{N \times 3}, \quad (141)$$

$$\mathbf{Q} \stackrel{\text{def}}{=} \left[ \left( \left( \gamma \theta^2 + \kappa \theta^3 \frac{\mathbf{x}_*^\top \mathbf{W} \mathbf{x}_*}{N} \right) \mathbf{I}_N + \kappa \theta^2 \mathbf{W} \right) \mathbf{x}_*, \kappa \theta^2 \mathbf{W} \mathbf{x}_*, \kappa \theta^2 \mathbf{x}_* \right] \in \mathbb{R}^{N \times 3}. \quad (142)$$

Hence,

$$\begin{aligned}
(c_* \mathbf{I} - J(\mathbf{Y}))^{-1} &= \left( c_* \mathbf{I} - J(\mathbf{W}) - \frac{1}{N} \mathbf{P} \mathbf{Q}^\top \right)^{-1} \\
&\stackrel{(a)}{=} \left( \mathbf{A} - \frac{1}{N} \mathbf{P} \mathbf{Q}^\top \right)^{-1} \\
&\stackrel{(b)}{=} \mathbf{A}^{-1} + \frac{1}{N} \mathbf{A}^{-1} \mathbf{P} \left( \mathbf{I}_3 - \frac{\mathbf{Q}^\top \mathbf{A}^{-1} \mathbf{P}}{N} \right)^{-1} \mathbf{Q}^\top \mathbf{A}^{-1},
\end{aligned}$$

where the matrix  $\mathbf{A}$  in step (a) denotes

$$\mathbf{A} \stackrel{\text{def}}{=} c_* \mathbf{I} - J(\mathbf{W}), \quad (143)$$

and step (b) is from the Morrison–Woodbury formula. Using the above identity, we obtain the following representation of the LHS of (139):

$$\frac{1}{N} \mathbf{x}^\top (c_* \mathbf{I} - J(\mathbf{Y}))^{-1} \mathbf{x} = \frac{\mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{x}_*}{N} + \frac{1}{N^2} \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{P} \left( \mathbf{I}_3 - \frac{\mathbf{Q}^\top \mathbf{A}^{-1} \mathbf{P}}{N} \right)^{-1} \mathbf{Q}^\top \mathbf{A}^{-1} \mathbf{x}_*. \quad (144)$$

Therefore, to calculate the limit of  $\frac{1}{N} \mathbf{x}^\top (c_* \mathbf{I} - J(\mathbf{Y}))^{-1}$  as  $N \rightarrow \infty$ , it suffices to derive the limit of the terms  $\frac{1}{N} \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{x}_*$ ,  $\frac{1}{N} \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{P} \in \mathbb{R}^{1 \times 3}$ ,  $\frac{1}{N} \mathbf{Q}^\top \mathbf{A}^{-1} \mathbf{P} \in \mathbb{R}^{3 \times 3}$  and  $\frac{1}{N} \mathbf{Q}^\top \mathbf{A}^{-1} \mathbf{x}_* \in \mathbb{R}^{3 \times 1}$ , which essentially involves computing the limit of  $\frac{1}{N} \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{W}^i \mathbf{x}_*$ , for  $i = 1, 2$ ; cf. (141) and (142). For this purpose, we can apply standard results in random matrix theory (e.g., [12, Proposition 9.3]) to obtain

$$\begin{aligned}
\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{W}^i \mathbf{x}_* &= \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{x}_*^\top (c_* \mathbf{I} - J(\mathbf{W}))^{-1} \mathbf{W}^i \mathbf{x}_*, \quad \forall i \in \mathbb{N} \cup \{0\} \\
&= \mathbb{E} \left[ \frac{\Lambda^i}{c_* - J(\Lambda)} \right] \\
&\stackrel{(a)}{=} \mathbb{E} \left[ \frac{\Lambda^i}{\rho + C_\phi - J(\Lambda)} \right] \\
&\stackrel{\text{def}}{=} d_i,
\end{aligned} \quad (145)$$

where  $\Lambda \sim \mu$  denotes the limiting eigenvalue distribution of  $\mathbf{W}$ , and step (a) is from the definition of  $c_* \stackrel{\text{def}}{=} \rho + C_\phi$ . Similarly,

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{x}_*^\top \mathbf{W} \mathbf{x}_* = \mathbb{E}[\Lambda] = 0. \quad (146)$$

Using (145) and (146), and recalling the definitions of the matrices  $\mathbf{P}$  and  $\mathbf{Q}$  in (141) and (142), we arrive at

$$\begin{aligned} \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{P} &= [d_0, d_1, d_2], \\ \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{Q} &= [\gamma\theta^2 d_0 + \kappa\theta^2 d_2, \kappa\theta^2 d_1, \kappa\theta^2 d_2], \\ \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{Q}^\top \mathbf{A}^{-1} \mathbf{P} &= \begin{bmatrix} \kappa\theta^2 d_2 + \gamma\theta^2 d_0 & \kappa\theta^2 d_3 + \gamma\theta^2 d_2 & \kappa\gamma^2 d_4 + \gamma\theta^2 d_2 \\ \kappa\theta^2 d_1 & \kappa\theta^2 d_2 & \kappa\theta^2 d_3 \\ \kappa\theta^2 d_0 & \kappa\theta^2 d_1 & \kappa\theta^2 d_2 \end{bmatrix}. \end{aligned} \quad (147)$$

Combining (144), (145) and (147), and after lengthy but straightforward calculations, we obtain

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{x}_*^\top (c_* \mathbf{I} - J(\mathbf{Y}))^{-1} \mathbf{x}_* = - \frac{1}{(\kappa\theta^2)^2 d_4 + \gamma\theta^2 + \frac{3\kappa\theta^2 d_2 + (\kappa\theta^2)^2 (2d_1 d_3 - 3d_2^2) + (\kappa\theta^2)^3 (d_3^2 + d_0 d_3^2 - 2d_1 d_2 d_3) - 1}{d_0 + \kappa\theta^2 d_1^2 - \kappa\theta^2 d_0 d_2}}. \quad (148)$$

Combining (148) and (139) yields the desired identity in (137).  $\square$

## D.2.2 Proof of Proposition 2

In what follows, we prove that any solution  $(\rho, \omega)$  to the state evolution (SE) fixed point equation (17f), with  $\omega \in (0, 1), \rho \in (0, \infty)$ , is a solution to the replica fixed point equation of (22). This is essentially achieved by combining Lemma 11 and Lemma 13. More specifically, the replica fixed point equations in (22), which we aim to prove, can be obtained by replacing the term  $\mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda_\nu)} \right]$  in (131a) by the alternative representation of  $\mathbb{E} \left[ \frac{1}{\rho + C_\phi - J(\Lambda_\nu)} \right]$  in (137), which is valid because of Lemma 12, together with additional algebraic manipulations.

Before we present the detailed proof, we first introduce a few definitions. Let

$$\chi \equiv \chi(\rho) \stackrel{\text{def}}{=} \frac{(1 - d_0)(d_1 + \kappa\theta^2 d_0 d_3 - \kappa\theta^2 d_1 d_2)}{d_0 - \kappa\theta^2 d_0 d_2 + \kappa\theta^2 d_1^2}, \quad (149)$$

where

$$\begin{aligned} d_i &\stackrel{\text{def}}{=} \mathbb{E}[\mathbf{H}\Lambda^i] \quad \forall i \in \mathbb{N} \cup \{0\}, \\ \mathbf{H} &\stackrel{\text{def}}{=} \mathbf{H}(\Lambda, \rho) = \frac{1}{\rho + \phi_{\text{poly}}(\Lambda)} = \frac{1}{\rho + C_\phi - J(\Lambda)} \\ &\stackrel{(129b)}{=} \frac{1}{\rho + \theta^2 a^2 (\gamma + 2a^2 \kappa)^2 + 1 - \theta (\gamma \Lambda - \theta \kappa \Lambda^2 + \kappa \Lambda^3)}. \end{aligned}$$

The rationale for the above definition of  $\chi$  will be clear from subsequent derivations. The notations  $\chi(\rho)$  and  $\mathbf{H}(\Lambda, \rho)$  above highlight the fact that  $\chi$  and  $\mathbf{H}$  both depend on the variable  $\rho$ , which is defined to be a solution to the SE fixed point equation. We further define

$$m \equiv m(\omega) \stackrel{\text{def}}{=} 1 - \text{mmse}_\pi(\omega), \quad (150)$$

$$\mathbf{Q} \equiv \mathbf{Q}(\Lambda, \mathbf{H}, \rho, \omega) \stackrel{\text{def}}{=} \kappa\theta^2 m \Lambda^2 + \kappa\theta^2 \chi \Lambda - \frac{\kappa\theta^2}{1 - m} \mathbb{E} [m \Lambda^2 \mathbf{H} + \chi \Lambda \mathbf{H}] + \frac{m}{1 - m}. \quad (151)$$

Again, the notation  $m(\omega)$  and  $\mathbf{Q}(\Lambda, \mathbf{H}, \rho, \omega)$  indicate that  $m$  depends on  $\omega$ , and  $\mathbf{Q}$  depends on the random variables  $\Lambda$  and  $\mathbf{H}$  (which is a function of  $\rho$ ) and the variables  $(\rho, \omega)$ .

To finish the proof, it remains to verify that the variables  $\chi, m, \mathbf{H}$  and  $\mathbf{Q}$  as defined above satisfy the replica fixed point equations (22).

**Proof of (22b).** By the definition of  $\mathbf{H}$ :

$$\begin{aligned}\mathbb{E}[\mathbf{H}] &\stackrel{(22f)}{=} \mathbb{E}\left[\frac{1}{\rho + C_\phi - J(\Lambda)}\right] \\ &\stackrel{(124)}{=} \mathbb{E}\left[\frac{1}{\rho + \phi(\Lambda)}\right] \\ &\stackrel{(131b)}{=} \text{mmse}_\pi(\omega) \\ &\stackrel{(150)}{=} 1 - m.\end{aligned}$$

This proves (22b).

**Proof of (22c).** We first verify equation (22b):  $\chi = \mathbb{E}[\Lambda\mathbf{Q}\mathbf{H}]$ . To this end, we will substitute the definitions of  $\mathbf{Q}$  and  $\chi$  into  $\mathbb{E}[\Lambda\mathbf{Q}\mathbf{H}]$  and then show that it is equal to  $\chi$ :

$$\begin{aligned}\mathbb{E}[\Lambda\mathbf{Q}\mathbf{H}] &\stackrel{(a)}{=} \mathbb{E}\left[\Lambda\left(\kappa\theta^2 m\Lambda^2 + \kappa\theta^2 \chi\Lambda - \frac{\kappa\theta^2}{1-m}\mathbb{E}[m\Lambda^2\mathbf{H} + \chi\Lambda\mathbf{H}] + \frac{m}{1-m}\right)\mathbf{H}\right] \\ &= \kappa\theta^2\left(m\mathbb{E}[\Lambda^3\mathbf{H}] + \chi\mathbb{E}[\Lambda^2\mathbf{H}] - \frac{1}{1-m}\mathbb{E}[\Lambda\mathbf{H}]\left(m\mathbb{E}[\Lambda^2\mathbf{H}] + \chi\mathbb{E}[\Lambda\mathbf{H}]\right)\right) + \frac{m}{1-m}\mathbb{E}[\Lambda\mathbf{H}] \\ &\stackrel{(b)}{=} \kappa\theta^2\left(md_3 + \chi d_2 - \frac{1}{1-m}d_1(md_2 + \chi d_1)\right) + \frac{m}{1-m}d_1 \\ &\stackrel{(c)}{=} \kappa\theta^2\left((1-d_0)d_3 + \chi d_2 - \frac{1}{d_0}(1-d_0)d_1d_2 + \chi\frac{1}{d_0}d_1^2\right) + \frac{1-d_0}{d_0}d_1 \\ &\stackrel{(d)}{=} \frac{(1-d_0)(d_1 + \kappa\theta^2 d_0 d_3 - \kappa\theta^2 d_1 d_2)}{d_0 - \kappa\theta^2 d_0 d_2 + \kappa\theta^2 d_1^2} \\ &\stackrel{(e)}{=} \chi,\end{aligned}\tag{152}$$

where step (a) is from the definition of  $\mathbf{Q}$  in (151); step (b) is from the definition of the shorthand  $d_i$ :  $d_i = \mathbb{E}[\Lambda^i\mathbf{H}]$ ; step (c) is due to the identity  $d_0 = 1 - m$ , which further follows from definition  $d_0 = \mathbb{E}[\mathbf{H}]$  and the equation  $\mathbb{E}[\mathbf{H}] = 1 - m$  (see (22b)); step (d) follows from the definition of  $\chi$  in (149) and straightforward calculations; step (e) is from the definition of  $\chi$  in (149). This completes the proof of (22c).

**Proof of (22d).** The term  $\mathbb{E}[\Lambda^2\mathbf{Q}\mathbf{H}]$  appears on the RHS of (22d). We first rewrite it as an explicit function of  $(d_0, d_1, d_2, d_3)$  and  $\chi$ :

$$\begin{aligned}\mathbb{E}[\Lambda^2\mathbf{Q}\mathbf{H}] &= \mathbb{E}\left[\Lambda^2\left(\kappa\theta^2 m\Lambda^2 + \kappa\theta^2 \chi\Lambda - \frac{\kappa\theta^2}{1-m}\mathbb{E}[m\Lambda^2\mathbf{H} + \chi\Lambda\mathbf{H}] + \frac{m}{1-m}\right)\mathbf{H}\right] \\ &\stackrel{(a)}{=} \kappa\theta^2 md_4 + \kappa\theta^2 d_3 \chi - \frac{\kappa\theta^2}{1-m}(md_2^2 + d_1 d_2 \chi) + \frac{m}{1-m}d_2 \\ &\stackrel{(b)}{=} (1-d_0)\left(\kappa\theta^2 d_4 + \frac{\kappa\theta^2 d_3 \chi}{1-d_0} - \frac{\kappa\theta^2}{d_0}\left(d_2^2 + \frac{\chi}{1-d_0}d_1 d_2\right) + \frac{1}{d_0}d_2\right),\end{aligned}\tag{153}$$

where step (a) is from the definition of  $\mathbf{Q}$  in (151); step (b) uses the identity  $d_0 = 1 - m$ , similar to step (c) of (152). Using (153), we can continue to rewrite the RHS of (22d) as a function of  $(d_0, d_1, d_2, d_3)$  and  $\chi$ :

$$\begin{aligned}\text{RHS of (22d)} &= \kappa\theta^2\left(\frac{m}{1-m}\mathbb{E}[\Lambda^2\mathbf{H}] + \frac{\chi}{1-m}\mathbb{E}[\Lambda\mathbf{H}] + \mathbb{E}[\Lambda^2\mathbf{Q}\mathbf{H}]\right) + \gamma\theta^2 m \\ &= \kappa\theta^2\left(\frac{1}{d_0}(1-d_0)d_2 + \frac{1}{d_0}d_1\chi + \mathbb{E}[\Lambda^2\mathbf{Q}\mathbf{H}]\right) + (1-d_0)\gamma\theta^2 \\ &= (1-d_0)\left((\kappa\theta^2)^2 d_4 + \gamma\theta^2 + \frac{\kappa\theta^2}{d_0(1-d_0)}(d_1 + \kappa\theta^2(d_0 d_3 - d_1 d_2))\chi + \frac{\kappa\theta^2}{d_0}(2d_2 - \kappa\theta^2 d_2^2)\right).\end{aligned}\tag{154}$$



Note that  $\chi$  is itself a function of  $(d_0, d_1, d_2, d_3)$ ; see (149). Substituting (149) into (154) eliminates the variable  $\chi$ . After some algebra we finally obtain

$$\begin{aligned}
\text{RHS of (22d)} &\stackrel{(a)}{=} (1 - d_0) \left( (\kappa\theta^2)^2 d_4 + \gamma\theta^2 + \frac{3\kappa\theta^2 d_2 + (\kappa\theta^2)^2 (2d_1 d_3 - 3d_2^2) + (\kappa\theta^2)^3 (d_2^3 + d_0 d_3^2 - 2d_1 d_2 d_3) - 1}{d_0 + \kappa\theta^2 d_1^2 - \kappa\theta^2 d_0 d_2} + \frac{1}{d_0} \right) \\
&\stackrel{(b)}{=} (1 - d_0) \left( - \left( \mathbb{E} \left[ \frac{1}{\rho + C_\phi - J(\Lambda_\nu)} \right] \right)^{-1} + \frac{1}{d_0} \right) \\
&\stackrel{(c)}{=} \left( 1 - \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right) \left\{ - \left( \mathbb{E} \left[ \frac{1}{\rho + C_\phi - J(\Lambda_\nu)} \right] \right)^{-1} + \left( \mathbb{E} \left[ \frac{1}{\rho + \phi(\Lambda)} \right] \right)^{-1} \right\} \\
&\stackrel{(d)}{=} \frac{\omega}{1 - \omega},
\end{aligned} \tag{155}$$

where step (a) is a consequence of (149) and (154); step (b) is from Lemma 13; step (c) is from the definition of  $d_0$  (see (138)); and step (d) is due to (131) (Lemma 11). This proves (22d).

To conclude, we have verified that the stationary point of the state evolution of OAMP  $(\rho, \omega)$  satisfies equations (22a)-(22d). The proof for the other direction, namely, each fixed point of (22) such that  $\rho > 0$  and  $\omega \in (0, 1)$  also satisfies the SE equations (17f) is similar and thus omitted.

## E Omitted Proofs for the Optimality Result (Theorem 2)

### E.1 Proof of Proposition 3

We consider a given iterative algorithm (Definition 5) of the form:

$$\mathbf{r}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{r}_1, \dots, \mathbf{r}_{t-1}; \mathbf{a}) + g_t(\mathbf{r}_1, \dots, \mathbf{r}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}, \tag{156}$$

which returns the following estimator after  $t$  iterations:

$$\hat{\mathbf{r}}_t = \psi_t(\mathbf{r}_1, \dots, \mathbf{r}_t; \mathbf{a}), \tag{157}$$

Our goal will be to construct a lifted OAMP algorithm (Definition 6), which can approximate the estimator above. Throughout the proof, we will refer to the given iterative algorithm (156) as the *target algorithm*. Our construction of the desired lifted OAMP algorithm proceeds in two steps.

**Step 1 (Divergence Removal).** We first implement the target algorithm (156) using an *intermediate OAMP algorithm* which takes the form:

$$\mathbf{v}_{t,S} = \bar{\Psi}_S(\mathbf{Y}) \cdot \bar{F}_t(\mathbf{v}_{<t,\bullet}; \mathbf{a}) \quad \forall t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty. \tag{158}$$

In the above display,

1. The iterates generated by the algorithm are indexed by two indices, a time index  $t$ , which takes values in  $\mathbb{N}$ , and a set index  $S$ , which can be any finite and non-empty subset of  $\mathbb{N}$ .
2. For any finite subset  $S = \{s_1, s_2, \dots, s_k\}$  with sorted elements  $s_1 < s_2 < \dots < s_k$ , we define the matrix denoiser:

$$\Psi_S(\lambda) \stackrel{\text{def}}{=} \Psi_{s_k}(\lambda) \cdot \Psi_{s_{k-1}}(\lambda) \cdot \dots \cdot \Psi_{s_1}(\lambda), \quad \bar{\Psi}_S(\lambda) \stackrel{\text{def}}{=} \Psi_S(\lambda) - \mathbb{E}_{\Lambda \sim \mu} [\Psi_S(\Lambda)] \quad \forall \lambda \in \mathbb{R}.$$

3. The iterates are generated in the following sequence:

- (a) In step (1), we generate the iterate  $\mathbf{v}_{1,\{1\}}$ .
- (b) In step (2), we generate the new iterates  $\mathbf{v}_{1,\{1,2\}}$  and  $\mathbf{v}_{2,\{1\}}, \mathbf{v}_{2,\{1,2\}}$ .

(c) At step  $(t)$ , we generate the new iterates:

$$\{\mathbf{v}_{1,S} : S \subset [t], t \in S\}, \{\mathbf{v}_{2,S} : S \subset [t], t \in S\}, \dots, \{\mathbf{v}_{t-1,S} : S \subset [t], t \in S\} \text{ and } \{\mathbf{v}_{t,S} : S \subset [t]\}.$$

We use the notation  $\mathbf{v}_{\leq t, \bullet}$  to denote the collection of iterates that have been generated at the end of  $t$  steps:

$$\mathbf{v}_{\leq t, \bullet} = \{\mathbf{v}_{s,S} : s \leq t, S \subset [t]\}.$$

The notation  $\mathbf{v}_{< t, \bullet}$  denotes the collection of iterates that have been generated before the  $t$ th step, which is the same as the collection  $\mathbf{v}_{\leq (t-1), \bullet}$ .

4. The functions  $(F_t)_{t \in \mathbb{N}}$  are continuously differentiable, Lipschitz and satisfy the divergence-free requirement:

$$\mathbb{E}[\partial_{s,S} \bar{F}_t(\mathbf{V}_{< t, \bullet}; \mathbf{A})] = 0 \quad \forall S \subset [t-1], s < t,$$

where  $(\mathbf{X}_*, \{\mathbf{V}_{t,S} : t \in \mathbb{N}, S \subset \mathbb{N}\}; \mathbf{A})$  are the state evolution random variables associated with the intermediate OAMP algorithm in (158).

5. Furthermore, the iterates generated by the intermediate OAMP algorithm in (158) can approximate the iterates of the target algorithm in (156). Specifically, for each  $t \in \mathbb{N}$ , there is a postprocessing function  $H_t : \mathbb{R}^{t \cdot 2^t + k} \mapsto \mathbb{R}$ , which is continuously differentiable and Lipschitz, which satisfies:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{r}_t - H_t(\mathbf{v}_{\leq t, \bullet}, \mathbf{a})\|^2}{N} = 0. \quad (159)$$

The construction of the intermediate OAMP algorithm with the properties stated above is presented in the proof of the following Lemma 14, whose proof is deferred to Section E.1.1.

**Lemma 14.** *For any target algorithm of the form (156), there is an intermediate OAMP algorithm of the form (158), whose iterates can approximate the iterates generated by the target algorithm in the sense of (159).*

**Step 2: Polynomial Approximation.** Next, by approximating the functions  $(\Psi_S)_{S \subset \mathbb{N}}$  by polynomials, we construct a lifted OAMP algorithm, whose iterates can approximate the iterates generated by the intermediate OAMP algorithm (158). This construction is presented in the following Lemma 5, whose proof is deferred to Section E.1.2.

**Lemma 15.** *For each  $D \in \mathbb{N}$ , there is a degree- $D$  lifted AMP algorithm (Definition 6):*

$$\mathbf{w}_{t,i}^{(D)} = (\mathbf{Y}^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i] \cdot \mathbf{I}_N) \cdot \bar{F}_t^{(D)}(\mathbf{w}_{1,\bullet}^{(D)}, \dots, \mathbf{w}_{t-1,\bullet}^{(D)}; \mathbf{a}) \quad t \in \mathbb{N}, i \in [D], \quad (160)$$

which can approximate the iterates produced by the intermediate OAMP algorithm (158) in the following sense: for any  $D \in \mathbb{N}$ ,  $t \in \mathbb{N}$ , and any finite subset  $S \subset \mathbb{N}$ , there is a homogeneous linear postprocessing function  $h_{t,S}^{(D)} : \mathbb{R}^D \mapsto \mathbb{R}$  which has the property that the vector  $h_{t,S}^{(D)}(\mathbf{w}_{t,\bullet}^{(D)})$  approximates the intermediate OAMP iterate  $\mathbf{v}_{t,S}$  in the sense:

$$\lim_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \frac{\|h_{t,S}^{(D)}(\mathbf{w}_{t,\bullet}^{(D)}) - \mathbf{v}_{t,S}\|^2}{N} = 0 \quad (161)$$

With these preliminary results, we can now present the proof of Proposition 3.

*Proof of Proposition 3.* Throughout the proof, we shorthand an approximation result like (159) using the notation:

$$\mathbf{r}_t \stackrel{N \rightarrow \infty}{\simeq} H_t(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \quad \forall t \in \mathbb{N}. \quad (162)$$

Likewise, we shorthand an approximation result like (161) using:

$$\mathbf{v}_{t,S} \stackrel{N, D \rightarrow \infty}{\simeq} \mathbf{v}_{t,S}^{(D)} \quad \forall t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty. \quad (163)$$

Lemma 14 and Lemma 15 together show that the lifted OAMP algorithm on (160) can approximate the iterates of the target algorithm in (156):

$$\mathbf{r}_t \stackrel{(a)}{\stackrel{N}{\simeq}} H_t(\mathbf{v}_{1,\bullet}, \dots, \mathbf{v}_{t,\bullet}; \mathbf{a}) \stackrel{(b)}{\stackrel{N, D \rightarrow \infty}{\simeq}} H_t(\{h_{s,R}^{(D)}(\mathbf{w}_{s,\bullet}^{(D)}) : s \leq t, R \subset [t]\}; \mathbf{a}) \quad \forall t \in \mathbb{N}. \quad (164)$$

In the above display, the approximation in step (a) follows from (159) and the approximation in step (b) follows by replacing each iterate  $\mathbf{v}_{s,R}$ ,  $s \in [t], R \subset [t]$  of the intermediate OAMP algorithm by its approximation in (161), and exploiting the fact that  $H_t$  is a Lipschitz function. For convenience, we define the composite postprocessing functions  $\{H_t^{(D)} : t \in \mathbb{N}, D \in \mathbb{N}\}$ :

$$H_t^{(D)}(\mathbf{w}_{\leq t,\bullet}; \mathbf{a}) \stackrel{\text{def}}{=} H_t(\{h_{s,R}^{(D)}(\mathbf{w}_{s,\bullet}) : s \leq t, R \subset [t]\}; \mathbf{a}) \quad \forall \mathbf{w}_{\leq t,\bullet} \in \mathbb{R}^{tD}, \mathbf{a} \in \mathbb{R}^k.$$

With this definition, (164) can be expressed as:

$$\mathbf{r}_t \stackrel{N, D \rightarrow \infty}{\simeq} H_t^{(D)}(\mathbf{w}_{\leq t,\bullet}^{(D)}; \mathbf{a}) \quad \forall t \in \mathbb{N}.$$

Plugging in the above approximation in (157), we have the following approximation for the estimator  $\hat{\mathbf{r}}_t$  returned by the target algorithm:

$$\hat{\mathbf{r}}_t \stackrel{N, D \rightarrow \infty}{\simeq} \tilde{\mathbf{w}}_t^{(D)} \stackrel{\text{def}}{=} \psi_t(H_1^{(D)}(\mathbf{w}_{1,\bullet}^{(D)}; \mathbf{a}), H_2^{(D)}(\mathbf{w}_{\leq 2,\bullet}^{(D)}; \mathbf{a}), \dots, H_t^{(D)}(\mathbf{w}_{\leq t,\bullet}^{(D)}; \mathbf{a}); \mathbf{a}). \quad (165)$$

Hence, we have constructed an estimator  $\tilde{\mathbf{w}}_t^{(D)}$  which can be computed by running a degree- $D$  lifted OAMP algorithm for  $t$  steps, and approximates the estimator  $\hat{\mathbf{r}}_t$  returned by the target iterative algorithm after  $t$  steps. This proves the claim of Proposition 3.  $\square$

We now present the proofs of Lemma 14 and Lemma 15.

### E.1.1 Divergence Removal (Proof of Lemma 14)

*Proof of Lemma 14.* Consider a target iterative algorithm of the form:

$$\mathbf{r}_t = \Psi_t(\mathbf{Y}) \cdot f_t(\mathbf{r}_1, \dots, \mathbf{r}_{t-1}; \mathbf{a}) + g_t(\mathbf{r}_1, \dots, \mathbf{r}_{t-1}; \mathbf{a}) \quad \forall t \in \mathbb{N}. \quad (166)$$

Our goal is to construct an intermediate OAMP algorithm of the form:

$$\mathbf{v}_{t,S} = \bar{\Psi}_S(\mathbf{Y}) \cdot \bar{F}_t(\mathbf{v}_{\leq t,\bullet}; \mathbf{a}) \quad \forall t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty. \quad (167)$$

For each  $t \in \mathbb{N}$ , we wish to design a postprocessing function  $H_t : \mathbb{R}^{t \cdot 2^t + k} \mapsto \mathbb{R}$  which allows us to approximate the iterates of the target algorithm with the iterates of the intermediate OAMP algorithm in the sense:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{r}_t - H_t(\mathbf{v}_{\leq t,\bullet}, \mathbf{a})\|^2}{N} \rightarrow 0 \quad \forall t \in \mathbb{N}. \quad (168)$$

We will shorthand approximation statements of the form (159) using the notation:

$$\mathbf{r}_t \stackrel{N}{\simeq} H_t(\mathbf{v}_{\leq t,\bullet}; \mathbf{a}) \quad \forall t \in \mathbb{N}. \quad (169)$$

The construction proceeds by induction. Suppose that we have specified the first  $t$  functions  $(F_s)_{s \leq t}$  and the post-processing functions  $(H_s)_{s \leq t}$  for  $t$  iterations. We now extend our construction to the  $(t+1)$ th iteration. Recall the update rule for  $\mathbf{r}_{t+1}$ :

$$\begin{aligned} \mathbf{r}_{t+1} &= \Psi_{t+1}(\mathbf{Y}) \cdot f_{t+1}(\mathbf{r}_{\leq t}; \mathbf{a}) + g_{t+1}(\mathbf{r}_{\leq t}; \mathbf{a}) \\ &\stackrel{N}{\simeq} \Psi_{t+1}(\mathbf{Y}) \cdot f_{t+1}(H_1(\mathbf{v}_{\leq 1,\bullet}; \mathbf{a}), \dots, H_t(\mathbf{v}_{\leq t,\bullet}; \mathbf{a})) + g_{t+1}(H_1(\mathbf{v}_{\leq 1,\bullet}; \mathbf{a}), \dots, H_t(\mathbf{v}_{\leq t,\bullet}; \mathbf{a})), \end{aligned}$$

where the approximation in the last step follows from the induction hypothesis  $\mathbf{r}_s \stackrel{N \rightarrow \infty}{\simeq} H_s(\mathbf{v}_{\leq s, \bullet}; \mathbf{a})$  for any  $s \leq t$ , along with the fact that  $\|\Psi_t(\mathbf{Y})\|_{\text{op}}$  is uniformly bounded as  $N \rightarrow \infty$  and the functions  $f_{t+1}, g_{t+1}$  are Lipschitz. We introduce the composite functions:

$$\begin{aligned} F_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) &\stackrel{\text{def}}{=} f_{t+1}(H_1(\mathbf{v}_{\leq 1, \bullet}; \mathbf{a}), H_2(\mathbf{v}_{\leq 2, \bullet}; \mathbf{a}), \dots, H_t(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}); \mathbf{a}) \quad \mathbf{v}_{\leq t, \bullet} \in \mathbb{R}^{t \cdot 2^t}, \mathbf{a} \in \mathbb{R}^k, \\ G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) &\stackrel{\text{def}}{=} g_{t+1}(H_1(\mathbf{v}_{\leq 1, \bullet}; \mathbf{a}), H_2(\mathbf{v}_{\leq 2, \bullet}; \mathbf{a}), \dots, H_t(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}); \mathbf{a}) \quad \mathbf{v}_{\leq t, \bullet} \in \mathbb{R}^{t \cdot 2^t}, \mathbf{a} \in \mathbb{R}^k. \end{aligned}$$

With these definitions we have that,

$$\mathbf{r}_{t+1} \stackrel{N \rightarrow \infty}{\simeq} \Psi_{t+1}(\mathbf{Y}) \cdot F_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}). \quad (170)$$

The above expression suggests that we should use the function  $F_{t+1}$  in the  $t+1$  step of the intermediate OAMP algorithm. However, since  $F_{t+1}$  is not divergence-free (cf. Definition 4), we define:

$$\bar{F}_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \stackrel{\text{def}}{=} F_t(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) - \sum_{s=1}^t \sum_{R \subset [t]} b_{s,R} \cdot \mathbf{v}_{s,R} \quad \forall \mathbf{v}_{\leq t, \bullet} \in \mathbb{R}^{t \cdot 2^t}, \mathbf{a} \in \mathbb{R}^k.$$

In the above display, we introduced the scalars:

$$b_{s,R} \stackrel{\text{def}}{=} \mathbb{E}[\partial_{s,R} F_{t+1}(\mathbf{V}_{\leq t, \bullet}; \mathbf{A})] \quad \forall s \leq t, R \subset [t],$$

where  $(\mathbf{X}_*, \{\mathbf{V}_{s,R} : s \leq t, R \subset [t]\}; \mathbf{A})$  denote the state evolution random variables corresponding to the iterates of the intermediate OAMP algorithm that have already been constructed (by the induction hypothesis). Notice that by construction,  $\bar{F}_{t+1}$  satisfies the divergence-free requirement:

$$\mathbb{E}[\partial_{s,S} \bar{F}_{t+1}(\mathbf{V}_{\leq t, \bullet}; \mathbf{A})] = 0 \quad \forall s \leq t, S \subset [t].$$

We can now extend our construction of the intermediate OAMP algorithm to the  $(t+1)$ th step:

$$\mathbf{v}_{t+1,S} = \bar{\Psi}_S(\mathbf{Y}) \cdot \bar{F}_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \quad \forall S \subset \mathbb{N} \quad 1 \leq |S| < \infty. \quad (171)$$

Next, we define the postprocessing function  $H_{t+1}$ . For convenience, we introduce the vectors  $(\mathbf{F}_\tau)_{\tau \in [t+1]}$ :

$$\bar{\mathbf{F}}_\tau \stackrel{\text{def}}{=} \bar{F}_\tau(\mathbf{v}_{< \tau, \bullet}; \mathbf{a}) \quad \forall \tau \in [t+1]. \quad (172)$$

With these definitions, we can express the approximation in (170) as:

$$\begin{aligned} \mathbf{r}_{t+1} &\stackrel{N \rightarrow \infty}{\simeq} \Psi_{t+1}(\mathbf{Y}) \cdot F_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \\ &\stackrel{N \rightarrow \infty}{\simeq} \Psi_{t+1}(\mathbf{Y}) \cdot \mathbf{F}_{t+1} + \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot \Psi_{t+1}(\mathbf{Y}) \cdot \mathbf{v}_{s,S} + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}). \end{aligned} \quad (173)$$

For any  $S \subset [t+1]$  we define:

$$m_S \stackrel{\text{def}}{=} \mathbb{E}_{\Lambda \sim \mu}[\Psi_S(\Lambda)].$$

With this definition, we can express the approximation in (173) as:

$$\begin{aligned}
\mathbf{r}_{t+1} &\stackrel{N \rightarrow \infty}{\simeq} \Psi_{t+1}(\mathbf{Y}) \cdot \mathbf{F}_{t+1} + \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot \Psi_{t+1}(\mathbf{Y}) \cdot \mathbf{v}_{s,S} + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \\
&= \Psi_{t+1}(\mathbf{Y}) \cdot \mathbf{F}_{t+1} + \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot \Psi_{t+1}(\mathbf{Y}) \cdot (\Psi_S(\mathbf{Y}) - m_S \mathbf{I}_N) \cdot \mathbf{F}_s + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \\
&= \Psi_{t+1}(\mathbf{Y}) \cdot \mathbf{F}_{t+1} + \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot (\Psi_{\{t+1\} \cup S}(\mathbf{Y}) - m_S \Psi_{t+1}(\mathbf{Y})) \cdot \mathbf{F}_s + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \\
&\stackrel{(167), (171)}{=} \mathbf{v}_{t+1, \{t+1\}} + \left( \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot (\mathbf{v}_{s, S \cup \{t+1\}} - m_S \cdot \mathbf{v}_{s, \{t+1\}}) \right) + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \\
&\quad + m_{\{t+1\}} \cdot \mathbf{F}_{t+1} + \left( \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot (m_{S \cup \{t+1\}} - m_S \cdot m_{\{t+1\}}) \cdot \mathbf{F}_s \right). \tag{174}
\end{aligned}$$

Recalling (172),

$$\begin{aligned}
\mathbf{r}_{t+1} &\stackrel{N \rightarrow \infty}{\simeq} \mathbf{v}_{t+1, \{t+1\}} + \left( \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot (\mathbf{v}_{s, S \cup \{t+1\}} - m_S \cdot \mathbf{v}_{s, \{t+1\}}) \right) + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) \\
&\quad + m_{\{t+1\}} \cdot \bar{\mathbf{F}}_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) + \left( \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot (m_{S \cup \{t+1\}} - m_S \cdot m_{\{t+1\}}) \cdot \bar{\mathbf{F}}_s(\mathbf{v}_{\leq s, \bullet}; \mathbf{a}) \right).
\end{aligned}$$

In light of the above expression, we define the postprocessing function for the  $t+1$ th iteration as:

$$\begin{aligned}
H_{t+1}(\mathbf{v}_{\leq t+1, \bullet}; \mathbf{a}) &\stackrel{\text{def}}{=} \mathbf{v}_{t+1, \{t+1\}} + \left( \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot (\mathbf{v}_{s, S \cup \{t+1\}} - m_S \cdot \mathbf{v}_{s, \{t+1\}}) \right) \\
&\quad + G_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) + m_{\{t+1\}} \cdot \check{\mathbf{F}}_{t+1}(\mathbf{v}_{\leq t, \bullet}; \mathbf{a}) + \left( \sum_{s=1}^t \sum_{S \subset [t]} b_{s,S} \cdot (m_{S \cup \{t+1\}} - m_S \cdot m_{\{t+1\}}) \cdot \check{\mathbf{F}}_s(\mathbf{v}_{\leq s, \bullet}; \mathbf{a}) \right).
\end{aligned}$$

Hence, we can approximate the  $t+1$ th iterate of the target algorithm using the iterates of the intermediate OAMP algorithm:

$$\mathbf{r}_{t+1} \stackrel{N \rightarrow \infty}{\simeq} H_{t+1}(\mathbf{v}_{\leq t+1, \bullet}; \mathbf{a}).$$

This completes the construction of the intermediate OAMP algorithm with all the desired properties by induction.  $\square$

### E.1.2 Polynomial Approximation (Proof of Lemma 15)

*Proof of Lemma 15.* Recall the intermediate OAMP algorithm takes the form:

$$\mathbf{v}_{t,S} = \bar{\Psi}_S(\mathbf{Y}) \cdot \bar{\mathbf{F}}_t(\mathbf{v}_{< t, \bullet}; \mathbf{a}) \quad \forall t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty. \tag{175}$$

To implement this algorithm using a degree- $D$  lifted OAMP algorithm, we begin by repeating the arguments used in the proof of Lemma 5 to construct an OAMP algorithm of the form:

$$\begin{aligned}
\mathbf{v}_{t,S}^{(D)} &= \left( \Psi_S^{(D)}(\mathbf{Y}) - \mathbb{E}_{\Lambda \sim \mu} [\Psi_S^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( \bar{\mathbf{F}}_t(\mathbf{v}_{< t, \bullet}^{(D)}; \mathbf{a}) - \sum_{s=1}^{t-1} \sum_{R \subset [t-1]} \mathbb{E}[\partial_{s,R} \bar{\mathbf{F}}_t(\mathbf{V}_{< t, \bullet}^{(D)}; \mathbf{A})] \cdot \mathbf{v}_{s,R}^{(D)} \right) \tag{176} \\
&\quad \forall t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty,
\end{aligned}$$

which approximates the iterates generated by the intermediate OAMP algorithm in the sense that,

$$\limsup_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \frac{\|\mathbf{v}_{t,S}^{(D)} - \mathbf{v}_{t,S}\|^2}{N} = 0 \quad \forall t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty. \quad (177)$$

In the above display:

- The matrix denoiser  $\Psi_S^{(D)} : \mathbb{R} \mapsto \mathbb{R}$  is a degree- $D$  polynomial which approximates the matrix denoiser  $\Psi_S : \mathbb{R} \mapsto \mathbb{R}$  and is constructed using the Weierstrass approximation theorem.
- $(\mathbf{X}_*, \{\mathbf{V}_{t,S}^{(D)} : t \in \mathbb{N}, S \subset \mathbb{N}\}; \mathbf{A})$  denote the state evolution random variables associated with (176).

We claim that there is a degree- $D$  lifted AMP algorithm (Definition 6):

$$\mathbf{w}_{t,i}^{(D)} = (\mathbf{Y}^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i] \cdot \mathbf{I}_N) \cdot \bar{F}_t^{(D)}(\mathbf{w}_{<t,\bullet}^{(D)}; \mathbf{a}) \quad \forall t \in \mathbb{N}, i \in [D], \quad (178)$$

along with homogeneous linear postprocessing maps

$$h_{t,S}^{(D)} : \mathbb{R}^D \mapsto \mathbb{R} \quad t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty,$$

which can exactly recover the iterates generated by the degree- $D$  approximation to the intermediate OAMP algorithm (176):

$$\mathbf{v}_{t,S}^{(D)} = h_{t,S}^{(D)}(\mathbf{w}_{t,\bullet}^{(D)}) \quad \forall t \in \mathbb{N}, S \subset \mathbb{N}, 1 \leq |S| < \infty. \quad (179)$$

Observe that (177) and (179) immediately imply Lemma 15. Indeed, for any  $t \in \mathbb{N}$  and any finite subset  $S \subset \mathbb{N}$ :

$$\limsup_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \frac{\|h_{t,S}^{(D)}(\mathbf{w}_{t,\bullet}^{(D)}) - \mathbf{v}_{t,S}\|^2}{N} \stackrel{(179)}{=} \limsup_{D \rightarrow \infty} \text{plim sup}_{N \rightarrow \infty} \frac{\|\mathbf{v}_{t,S}^{(D)} - \mathbf{v}_{t,S}\|^2}{N} \stackrel{(177)}{=} 0,$$

as desired. Hence, the remainder of the proof consists of constructing the lifted OAMP algorithm and proving (179).

**Construction of the lifted OAMP Algorithm.** We will construct the desired lifted OAMP algorithm (178) by induction. As our induction hypothesis, we will assume that we have specified the functions  $\bar{F}_1^{(D)}, \dots, \bar{F}_{t-1}^{(D)}$ , as well as the homogenous linear postprocessing maps:

$$h_{s,S}^{(D)} : \mathbb{R}^D \mapsto \mathbb{R} \quad s < t, S \subset [t-1], S \geq 1,$$

so that the iterates resulting lifted OAMP algorithm  $\{\mathbf{w}_{s,i}^{(D)} : s < t, i \in [D]\}$  can reconstruct the following iterates of the degree- $D$  approximation to the intermediate OAMP algorithm:

$$\mathbf{v}_{s,S}^{(D)} = h_{s,S}^{(D)}(\mathbf{w}_{s,\bullet}^{(D)}) \quad \forall s < t, S \subset [t-1], |S| > 1. \quad (180)$$

We now extend our construction to step  $t$ . To do so, we need to construct the function  $\bar{F}_t^{(D)}$  used in the  $t$ th iteration of the lifted OAMP algorithm, as well as the linear postprocessing maps:

$$\begin{aligned} h_{t,S}^{(D)} &: \mathbb{R}^D \mapsto \mathbb{R} \quad S \subset [t], S \geq 1, \\ h_{s,S}^{(D)} &: \mathbb{R}^D \mapsto \mathbb{R} \quad S \subset [t], S \ni t, \end{aligned}$$

which ensure that:

$$\mathbf{v}_{t,S}^{(D)} = h_{t,S}^{(D)}(\mathbf{w}_{t,\bullet}^{(D)}) \quad \forall S \subset [t], |S| > 1, \quad (181)$$

$$\mathbf{v}_{s,S}^{(D)} = h_{s,S}^{(D)}(\mathbf{w}_{s,\bullet}^{(D)}) \quad \forall s < t, S \subset [t], S \ni t. \quad (182)$$

We will present the proof for (181), the proof for (182) is analogous. We begin by recalling the update rule for  $\mathbf{v}_{t,S}^{(D)}$ :

$$\begin{aligned} \mathbf{v}_{t,S}^{(D)} &= \left( \Psi_S^{(D)}(\mathbf{Y}) - \mathbb{E}_{\Lambda \sim \mu}[\Psi_S^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( \bar{F}_t(\mathbf{v}_{<t,\bullet}^{(D)}; \mathbf{a}) - \sum_{s=1}^{t-1} \sum_{R \subset [t-1]} \mathbb{E}[\partial_{s,R} \bar{F}_t(\mathbf{V}_{<t,\bullet}^{(D)}; \mathbf{A})] \cdot \mathbf{v}_{s,R}^{(D)} \right) \\ &\stackrel{(a)}{=} \left( \Psi_S^{(D)}(\mathbf{Y}) - \mathbb{E}_{\Lambda \sim \mu}[\Psi_S^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( \bar{F}_t(\{h_{s,R}^{(D)}(\mathbf{w}_{s,\bullet}^{(D)})\}; \mathbf{a}) - \sum_{s=1}^{t-1} \sum_{R \subset [t-1]} \mathbb{E}[\partial_{s,R} \bar{F}_t(\mathbf{V}_{<t,\bullet}^{(D)}; \mathbf{A})] \cdot h_{s,R}^{(D)}(\mathbf{w}_{s,\bullet}^{(D)}) \right), \end{aligned} \quad (183)$$

where step (a) follows from induction hypothesis in (180). In light of the above expression, we define  $F_{t,D}$  as the composite function:

$$F_t^{(D)}(\mathbf{w}_{<t,\bullet}; \mathbf{a}) \stackrel{\text{def}}{=} \bar{F}_t(\{h_{s,R}^{(D)}(\mathbf{w}_{s,\bullet}) : s < t, R \subset [t-1], |R| > 1\}; \mathbf{a}) \quad \forall \mathbf{w}_{<t,\bullet} \in \mathbb{R}^{(t-1)D}, \mathbf{a} \in \mathbb{R}^k, \quad (184)$$

and the divergence-free iterate denoiser:

$$\bar{F}_t^{(D)}(\mathbf{w}_{<t,\bullet}; \mathbf{a}) = F_t^{(D)}(\mathbf{w}_{<t,\bullet}; \mathbf{a}) - \sum_{s=1}^{t-1} \sum_{i=1}^D \mathbb{E}[\partial_{s,i} F_t^{(D)}(\mathbf{W}_{<t,\bullet}^{(D)}; \mathbf{A})] \cdot \mathbf{w}_{s,i} \quad \forall \mathbf{w}_{<t,\bullet} \in \mathbb{R}^{(t-1)D}, \mathbf{a} \in \mathbb{R}^k, \quad (185)$$

where  $(\mathbf{X}_s, \{\mathbf{W}_{s,i} : s < t, i \in [D]\}; \mathbf{A})$  are the state evolution random variables associated with the lifted OAMP algorithm (178) that have been constructed so far (as a part of the induction hypothesis). We now extend our construction of the desired lifted OAMP algorithm to step  $t$ . At step  $t$ , the lifted OAMP algorithm generates the iterates:

$$\mathbf{w}_{t,i}^{(D)} = (\mathbf{Y}^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i] \cdot \mathbf{I}_N) \cdot \bar{F}_t^{(D)}(\mathbf{w}_{<t,\bullet}^{(D)}; \mathbf{a}) \quad i \in [D]. \quad (186)$$

Hence (183) can be expressed as:

$$\begin{aligned} \mathbf{v}_{t,S}^{(D)} &\stackrel{(184)}{=} \left( \Psi_S^{(D)}(\mathbf{Y}) - \mathbb{E}_{\Lambda \sim \mu}[\Psi_S^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( F_t^{(D)}(\mathbf{w}_{<t,\bullet}^{(D)}; \mathbf{a}) - \sum_{s=1}^{t-1} \sum_{R \subset [t-1]} \mathbb{E}[\partial_{s,R} \bar{F}_t(\mathbf{V}_{<t,\bullet}^{(D)}; \mathbf{A})] \cdot h_{s,R}^{(D)}(\mathbf{w}_{s,\bullet}^{(D)}) \right) \\ &\stackrel{(b)}{=} \left( \Psi_S^{(D)}(\mathbf{Y}) - \mathbb{E}_{\Lambda \sim \mu}[\Psi_S^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \left( F_t^{(D)}(\mathbf{w}_{<t,\bullet}^{(D)}; \mathbf{a}) - \sum_{s=1}^{t-1} \sum_{i=1}^D \mathbb{E}[\partial_{s,i} F_t^{(D)}(\mathbf{W}_{<t,\bullet}^{(D)}; \mathbf{A})] \cdot \mathbf{w}_{s,i}^{(D)} \right) \\ &\stackrel{(185)}{=} \left( \Psi_S^{(D)}(\mathbf{Y}) - \mathbb{E}_{\Lambda \sim \mu}[\Psi_S^{(D)}(\Lambda)] \cdot \mathbf{I}_N \right) \cdot \bar{F}_t^{(D)}(\mathbf{w}_{<t,\bullet}^{(D)}; \mathbf{a}) \end{aligned} \quad (187)$$

where the step marked (b) follows by relating the derivatives of the function  $F_t^{(D)}$  and the function  $\bar{F}_t$  using the chain rule and by exploiting the linearity of postprocessing maps  $\{h_{s,R}\}$ . Recall that  $\Psi_S^{(D)} : \mathbb{R} \mapsto \mathbb{R}$  is a degree- $D$  polynomial, and so, it can be expressed as a linear combination of the monomials  $\{\lambda \mapsto \lambda^i : i \in [D]\}$ :

$$\Psi_S^{(D)}(\lambda) - \mathbb{E}_{\Lambda \sim \mu}[\Psi_S^{(D)}(\Lambda)] = \sum_{i=1}^D \chi_S^{(D)}[i] \cdot (\lambda^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i]) \quad i \in [D].$$

Using this representation, the formula for  $\mathbf{v}_{t,S}^{(D)}$  obtained in (187) can be expressed as a linear combination of the newly constructed lifted OAMP iterates in (186):

$$\mathbf{v}_{t,S}^{(D)} = \sum_{i=1}^D \chi_S^{(D)}[i] \cdot \mathbf{w}_{t,i}^{(D)}.$$

Hence, we define the postprocessing function  $h_{t,S}^{(D)}$  as:

$$h_{t,S}^{(D)}(\mathbf{w}_{t,\bullet}) = \sum_{i=1}^D \chi_S^{(D)}[i] \cdot \mathbf{w}_{t,i} \quad \forall \mathbf{w}_{t,\bullet} \in \mathbb{R}^D.$$

which is linear and ensures that  $\mathbf{v}_{t,S}^{(D)} = h_{t,S}^{(D)}(\mathbf{w}_{t,\bullet}^{(D)})$ , as desired. This completes the construction of the desired lifted OAMP algorithm by induction and concludes the proof of this lemma.  $\square$

## E.2 Proof Proposition 4

Consider any degree- $D$  lifted OAMP algorithm (Definition 6) of the form:

$$\mathbf{w}_{t,i} = (\mathbf{Y}^i - \mathbb{E}_{\Lambda \sim \nu}[\Lambda^i] \cdot \mathbf{I}_N) \cdot f_t(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t-1,\bullet}; \mathbf{a}) \quad i \in [D], t \in \mathbb{N}, \quad (188)$$

which returns the following estimator after  $t$  iterations:

$$\widehat{\mathbf{w}}_t = h_t(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t,\bullet}; \mathbf{a}).$$

Let  $(\mathbf{X}_\star, W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$  denote the state evolution random variables associated with the first  $t$  iterations of the lifted OAMP algorithm, which form a Gaussian channel in the sense of Definition 3. We observe that, by Corollary 1:

$$\begin{aligned} \text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{x}_\star - \widehat{\mathbf{w}}_t\|^2}{N} &= \mathbb{E}[\{\mathbf{X}_\star - h(W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})\}^2] \\ &\stackrel{(a)}{\geq} \min_{h \in L^2(W_{\leq t,\bullet}; \mathbf{A})} \mathbb{E}[\{\mathbf{X}_\star - h(W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})\}^2] \\ &\stackrel{\text{def}}{=} \text{MMSE}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A}). \end{aligned} \quad (189)$$

In the above display, the inequality in step (a) follows by observing that since the function  $h$  is a Lipschitz function (as required by Definition 6), and hence,  $h \in L^2(W_{\leq t,\bullet}; \mathbf{A})$ . Furthermore, if Assumption 3 holds, the lower bound above is tight since we can choose  $h$  as the MMSE estimator for the Gaussian channel  $(\mathbf{X}_\star, W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$ .

In light of (189), for a  $t$  iteration lifted OAMP algorithm which uses functions  $f_{1:t}$  in its iterations, we define:

$$\mathcal{M}_t(f_1, \dots, f_t) \stackrel{\text{def}}{=} \text{MMSE}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A}),$$

where  $(\mathbf{X}_\star, W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$  are the state evolution random variables associated with the lifted OAMP algorithm. Notice that the RHS of the above display depends on the functions  $f_{1:t}$  implicitly, since the distribution of the state evolution random variables is determined by  $f_{1:t}$  (recall Definition 6).

**A convenient formula for  $\mathcal{M}_t(f_1, \dots, f_t)$ .** To determine the optimal lifted OAMP algorithm, we need to optimize the functional  $\mathcal{M}_t$  with respect to  $f_{1:t}$ . To do so, we first derive a convenient formula for  $\mathcal{M}_t(f_1, \dots, f_t) = \text{MMSE}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$ . Lemma 3 in Appendix A.2 shows that the MMSE (and other relevant properties) of a multivariate Gaussian channel, such as  $(\mathbf{X}_\star, W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$  can be expressed in terms of the MMSE of an *effective scalar Gaussian channel*:

$$(\mathbf{X}_\star, \mathbf{S} \stackrel{\text{def}}{=} \langle v_{\text{opt}}, W_{\leq t,\bullet} \rangle; \mathbf{A}), \quad (190)$$

whose observation random variable  $\mathbf{S}$  is a linear combination of the observation random variables of the original multivariate Gaussian channel. The weights of linear combination are given by  $v_{\text{opt}}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$  (which we will sometimes shorthand as  $v_{\text{opt}}$ , when the multivariate Gaussian channel in question is clear from the context). The effective scalar Gaussian channel operates at SNR  $\omega_{\text{eff}}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$ , and the MMSE of the original multivariate channel and the effective scalar channel are connected via the formula:

$$\text{MMSE}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A}) = \text{MMSE}(\mathbf{X}_\star | \mathbf{S}; \mathbf{A}) = \text{mmse}_\pi(\omega_{\text{eff}}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})),$$

where for  $\omega \in [0, 1]$ ,  $\text{mmse}_\pi(\omega)$  is the MMSE of a scalar Gaussian channel at SNR  $\omega$  (recall Definition 3). The following lemma provides a formula for the effective SNR  $\omega_{\text{eff}}(\mathbf{X}_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$ . To state this formula, recall from Definition 6 that the joint distribution of the state evolution random variables  $(\mathbf{X}_\star, W_{\leq t,\bullet}; \mathbf{A})$  is given by:

$$(\mathbf{X}_\star; \mathbf{A}) \sim \pi, \quad (W_{1,\bullet}, \dots, W_{t,\bullet}) | (\mathbf{X}_\star; \mathbf{A}) \sim \mathcal{N}(\mathbf{X}_\star \cdot \alpha \otimes q, \alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma). \quad (191a)$$



In the above display, the entries of  $q \in \mathbb{R}^D$ ,  $Q \in \mathbb{R}^{D \times D}$ , and  $\Gamma \in \mathbb{R}^{D \times D}$  are given by:

$$q_i = \mathbb{E}[\Lambda_\nu^i] - \mathbb{E}[\Lambda^i], \quad Q_{ij} = \mathbb{E}[(\Lambda_\nu^i - \mathbb{E}[\Lambda^i]) \cdot (\Lambda_\nu^j - \mathbb{E}[\Lambda^j])], \quad \Gamma_{ij} = \text{Cov}[\Lambda^i, \Lambda^j] \quad \text{where } \Lambda \sim \mu, \Lambda_\nu \sim \nu, \quad (191b)$$

the entries of  $\alpha \in \mathbb{R}^t$  and  $\Sigma \in \mathbb{R}^{t \times t}$  are given by:

$$\alpha_s \stackrel{\text{def}}{=} \mathbb{E}[X_\star \cdot f_s(W_{<s, \bullet}; \mathbf{A})], \quad \Sigma_{s\tau} \stackrel{\text{def}}{=} \mathbb{E}[f_s(W_{<s, \bullet}; \mathbf{A}) \cdot f_\tau(W_{<\tau, \bullet}; \mathbf{A})] \quad \forall s, \tau \in [t], \quad (191c)$$

and  $\otimes$  denotes the Kronecker (or tensor) product for matrices.

**Lemma 16.** *Consider a lifted OAMP algorithm which uses functions  $f_1, f_2, \dots, f_t$  with  $f_i : \mathbb{R}^{(i-1) \cdot D + k} \mapsto \mathbb{R}$  in the first  $t$  iterations. Let  $(X_\star, W_{1, \bullet}, \dots, W_{t, \bullet}; \mathbf{A})$  be the state evolution random variables associated with the algorithm. Then,*

$$\omega_{\text{eff}}(X_\star | W_{1, \bullet}, \dots, W_{t, \bullet}; \mathbf{A}) = \begin{cases} q^\top \left[ Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right]^\dagger q & : \kappa_t \in [0, 1) \\ 0 & : \kappa_t = 1. \end{cases}$$

Moreover, if  $\kappa_t < 1$ , then,

$$v_{\text{opt}}(X_\star | W_{1, \bullet}, \dots, W_{t, \bullet}; \mathbf{A}) = \frac{1}{\sqrt{\omega_{\text{eff}}(X_\star | W_{\leq t, \bullet}; \mathbf{A})}} \cdot (\Sigma^\dagger \alpha) \otimes [(1 - \kappa_t)Q + \kappa_t \Gamma]^\dagger q,$$

where the scalar  $\kappa_t \in [0, 1]$  is defined as<sup>2</sup>:

$$\kappa_t \stackrel{\text{def}}{=} \min \left\{ \mathbb{E} [\{X_\star - g(W_{<t, \bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_t) \right\}.$$

*Proof.* The proof of this lemma is presented in Section E.2.1. □

Combining the effective SNR formula from the lemma above with (190), we obtain the following expression for  $\mathcal{M}(f_1, \dots, f_t)$ :

$$\mathcal{M}(f_1, \dots, f_t) \stackrel{\text{def}}{=} \text{MMSE}(X_\star | W_{1, \bullet}, \dots, W_{t, \bullet}; \mathbf{A}) = \text{mmse}_\pi(\omega_{\text{eff}}(X_\star | W_{1, \bullet}, \dots, W_{t, \bullet}; \mathbf{A})) = m(\kappa_t), \quad (192a)$$

where we introduced the function  $m : [0, 1] \mapsto [0, 1]$  defined as follows:

$$m(\kappa) \stackrel{\text{def}}{=} \begin{cases} \text{mmse}_\pi \left( q^\top \left[ Q + \frac{\kappa}{1 - \kappa} \cdot \Gamma \right]^\dagger q \right) & : \kappa \in [0, 1) \\ \text{mmse}_\pi(0) & : \kappa = 1, \end{cases} \quad (192b)$$

where  $q, Q, \Gamma$  are as defined in (47).

**Monotonicity of  $m(\cdot)$**  The following monotonicity property of the function  $m$  introduced in (192) will play an important role in the proof of Proposition 4.

**Lemma 17.** *The function  $m : [0, 1] \mapsto [0, 1]$  is non-decreasing.*

*Proof.* See Section E.2.2. □

<sup>2</sup>For a collection of functions  $f_1, \dots, f_t$ ,  $\text{Span}(f_1, \dots, f_t)$  denotes the set of all functions that can be written as a linear combination of  $f_1, \dots, f_t$ .

**A Greedy Approach Minimizing  $\mathcal{M}_t$ .** To derive the optimal OAMP algorithm, our goal is to solve the variational problem:

$$\min_{f_1} \min_{f_2} \cdots \min_{f_t} \mathcal{M}_t(f_1, \dots, f_t), \quad (193)$$

Recall that Proposition 4 claims that the iterate denoisers  $(f_t^*)_{t \in \mathbb{N}}$  that optimize (193) are defined recursively: once the optimal denoisers  $f_1^*, \dots, f_t^*$  have been specified, the optimal denoiser for the  $t + 1$  iteration is the DMMSE estimator for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$ , where  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  denote the state evolution random variables corresponding to the first  $t$  iterations. We will show that these choices for the iterate denoisers correspond to a *greedy approach* to solving the optimization problem in (193):

$$f_1^* \in \arg \min_{f_1} \mathcal{M}_1(f_1), \quad f_2^* \in \arg \min_{f_2} \mathcal{M}_2(f_1^*, f_2), \quad \dots, \quad f_t^* \in \arg \min_{f_t} \mathcal{M}_t(f_1^*, f_2^*, \dots, f_{t-1}^*, f_t).$$

Proving this requires us to characterize the solution of the simpler (compared to (193)) variational problem:

$$\min_{f_t} \mathcal{M}_t(f_1, \dots, f_t),$$

which is done in the following lemma.

**Lemma 18.** *Let:*

1.  $f_1, \dots, f_{t-1}$  with  $f_i : \mathbb{R}^{(i-1) \times D+k} \mapsto \mathbb{R}$  be a fixed collection of iterate denoisers which satisfy the requirements imposed in the definition of lifted OAMP algorithms (Definition 6) ;
2.  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t-1,\bullet}; \mathbf{A})$  be the state evolution random variables associated with the first  $t-1$  iterations of a lifted OAMP algorithm which uses the functions  $f_1, \dots, f_{t-1}$  in the first  $t-1$  iterations;
3.  $f_t^\sharp : \mathbb{R}^{(t-1) \times D+k} \mapsto \mathbb{R}$  denote the DMMSE estimator for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t-1,\bullet}; \mathbf{A})$ .

Then, for any iterate denoiser  $f_t : \mathbb{R}^{(t-1) \times D+k} \mapsto \mathbb{R}$  which satisfies the requirements imposed in Definition 6, we have:

$$\mathcal{M}_t(f_1, \dots, f_t) \geq \mathcal{M}_t(f_1, \dots, f_t^\sharp) = m(\kappa_t^\sharp), \quad (194)$$

where  $\kappa_t^\sharp \stackrel{\text{def}}{=} \text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\mathcal{M}_{t-1}(f_{1:t-1}))$ .

*Proof.* The proof of this lemma is presented in Section E.2.3. □

Using the above intermediate lemmas, we can now present the proof of Proposition 4.

*Proof of Proposition 4.* To prove the claim of Proposition 4, we need to show that for any  $t \in \mathbb{N}$  and for any choice of iterate denoisers  $f_{1:t}$  which satisfy the requirements imposed in Definition 6,

$$\mathcal{M}_t(f_1, \dots, f_t) \geq \mathcal{M}_t(f_1^*, \dots, f_t^*), \quad (195)$$

where  $f_{1:t}^*$  denote the functions used in the iterations of the optimal lifted OAMP algorithm in (49), which are specified recursively as follows: after the functions  $f_1^*, \dots, f_{t-1}^*$  have been specified,  $f_t^*$  was defined as the DMMSE estimator for the Gaussian channel corresponding to the first  $t-1$  iterations of the resulting lifted OAMP algorithm. By the Lemma 18,  $f_{1:t}^*$  are precisely the functions derived from the greedy heuristic:

$$f_1^* \in \arg \min_{f_1} \mathcal{M}_1(f_1), \quad f_2^* \in \arg \min_{f_2} \mathcal{M}_2(f_1^*, f_2), \quad \dots, \quad f_t^* \in \arg \min_{f_t} \mathcal{M}_t(f_1^*, f_2^*, \dots, f_{t-1}^*, f_t),$$

where the minimization is understood to be over iterate denoisers which satisfy the requirements imposed in Definition 6. We will obtain the desired conclusion in (195) by induction. As our induction hypothesis, we assume that for any choice of the iterate denoisers  $f_1, \dots, f_{t-1}$ , we have:

$$\mathcal{M}_{t-1}(f_1, \dots, f_{t-1}) \geq \mathcal{M}_{t-1}(f_1^*, \dots, f_{t-1}^*). \quad (196)$$

In order to obtain the desired conclusion (195), we observe that:

$$\begin{aligned} \mathcal{M}_t(f_1, \dots, f_t) &\stackrel{(a)}{\geq} m(\kappa_t^\sharp), \text{ where } \kappa_t^\sharp \stackrel{\text{def}}{=} \text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\mathcal{M}_{t-1}(f_{1:t-1})) \\ &\stackrel{(b)}{\geq} m(\kappa_t^*), \text{ where } \kappa_t^* \stackrel{\text{def}}{=} \text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\mathcal{M}_{t-1}(f_{1:t-1}^*)) \\ &\stackrel{(c)}{=} \mathcal{M}_t(f_1^*, \dots, f_t^*). \end{aligned}$$

In the above display:

- Step (a) used the lower bound for  $\mathcal{M}_t(f_1, \dots, f_t)$  provided in Lemma 18.
- In step (b), we observed by the induction hypothesis that  $\mathcal{M}_{t-1}(f_{1:t-1}^*) \leq \mathcal{M}_{t-1}(f_{1:t-1})$ . Since  $\text{mmse}_\pi$  and  $\text{dmmse}_\pi$  are non-increasing functions (Lemma 1 and Lemma 4), we conclude that  $\kappa_t^* \leq \kappa_t^\sharp$ . Finally, since  $m$  is non-increasing (Lemma 17), we concluded that  $m(\kappa_t^\sharp) \geq m(\kappa_t^*)$ .
- In step (c), we again appealed to Lemma 18:

$$\mathcal{M}_t(f_1^*, \dots, f_t^*) = \min_{f_t} \mathcal{M}_t(f_1^*, \dots, f_t) = m(\kappa_t^*), \text{ where } \kappa_t^* \stackrel{\text{def}}{=} \text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\mathcal{M}_{t-1}(f_{1:t-1}^*)).$$

This concludes the proof of Proposition 4.  $\square$

### E.2.1 Proof of Lemma 16

*Proof of Lemma 16.* Recall from Definition 6 that the joint distribution of the state evolution random variables is given by:

$$(\mathbf{X}_*; \mathbf{A}) \sim \pi, \quad (\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}) | (\mathbf{X}_*; \mathbf{A}) \sim \mathcal{N}(\mathbf{X}_* \cdot \alpha \otimes q, \alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma). \quad (197)$$

In the above display, the entries of  $q \in \mathbb{R}^D$ ,  $Q \in \mathbb{R}^{D \times D}$ , and  $\Gamma \in \mathbb{R}^{D \times D}$  are given by:

$$q_i = \mathbb{E}[\Lambda_\nu^i] - \mathbb{E}[\Lambda^i], \quad Q_{ij} = \mathbb{E}[(\Lambda_\nu^i - \mathbb{E}[\Lambda^i]) \cdot (\Lambda_\nu^j - \mathbb{E}[\Lambda^j])], \quad \Gamma_{ij} = \text{Cov}[\Lambda^i, \Lambda^j] \quad \text{where } \Lambda \sim \mu, \Lambda_\nu \sim \nu, \quad (198)$$

and the entries of  $\alpha \in \mathbb{R}^t$  and  $\Sigma \in \mathbb{R}^{t \times t}$  are given by:

$$\alpha_s \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{X}_* \cdot f_s(\mathbf{W}_{<s,\bullet}; \mathbf{A})], \quad \Sigma_{s\tau} \stackrel{\text{def}}{=} \mathbb{E}[f_s(\mathbf{W}_{<s,\bullet}; \mathbf{A}) \cdot f_\tau(\mathbf{W}_{<\tau,\bullet}; \mathbf{A})] \quad \forall s, \tau \in [t]. \quad (199)$$

The proof of the claim is obtained by applying the effective SNR formula for general multivariate Gaussian channels stated as Lemma 3 in Appendix A.2 to the Gaussian channel  $(\mathbf{X}_*; \mathbf{W}_{\leq t, \bullet}; \mathbf{A})$ . According to this formula:

$$\omega_{\text{eff}}(\mathbf{X}_* | \mathbf{W}_{\leq t, \bullet}; \mathbf{A}) = \begin{cases} \frac{\vartheta}{1+\vartheta} & \text{if } \alpha \otimes q \in \text{Range}(\alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma), \\ 1 & \text{if } \alpha \otimes q \notin \text{Range}(\alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma), \end{cases} \quad (200a)$$

where:

$$\vartheta \stackrel{\text{def}}{=} \alpha \otimes q^\top (\alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma)^\dagger \alpha \otimes q. \quad (200b)$$

The claim of the lemma is obtained by simplifying the above formula.

**Orthogonalization.** To obtain the simplified formula stated in the lemma, it will be convenient to construct an orthonormal basis for  $\text{Span}(f_{1:t})$ . Let  $r = \text{Rank}(\Sigma)$  and  $\hat{f}_1, \hat{f}_2, \dots, \hat{f}_r : \mathbb{R}^{(t-1)D} \mapsto \mathbb{R}$  denote the basis functions constructed by orthogonalizing the functions  $f_{1:t}$  with respect to the state evolution random variables  $(\mathbf{W}_{<t,\bullet}; \mathbf{A})$ , via the Gram-Schmidt procedure. These basis functions are orthonormal in the sense that:

$$\mathbb{E}[\hat{f}_s(\mathbf{W}_{<t,\bullet}; \mathbf{A}) \hat{f}_\tau(\mathbf{W}_{<t,\bullet}; \mathbf{A})] = \begin{cases} 1 & : s = \tau \\ 0 & : \text{otherwise.} \end{cases}$$

We can express the functions  $f_{1:t}$  as a linear combination of the basis functions:

$$f_\tau = \sum_{s=1}^r L_{\tau s} \widehat{f}_s \quad \forall \tau \in [t], \quad L_{\tau s} \stackrel{\text{def}}{=} \mathbb{E}[f_\tau(\mathbf{W}_{<t,\bullet}; \mathbf{A}) \widehat{f}_s(\mathbf{W}_{<t,\bullet}; \mathbf{A})] \quad \forall s \in [r], \tau \in [t]. \quad (201)$$

Let  $L$  denote the  $t \times r$  matrix whose entries are given by the coefficients  $(L_{\tau s})_{\tau \in [t], s \in [r]}$  defined above. Using the decomposition in (201), along with the definitions of  $\alpha, \Sigma$  in (199), we obtain the following alternative expressions for  $\alpha$  and  $\Sigma$ :

$$\Sigma = LL^\top, \quad \alpha = L\widehat{\alpha}, \quad \widehat{\alpha} \stackrel{\text{def}}{=} (\mathbb{E}[\mathbf{X}_* \widehat{f}_1(\mathbf{W}_{<t,\bullet}; \mathbf{A})], \mathbb{E}[\mathbf{X}_* \widehat{f}_2(\mathbf{W}_{<t,\bullet}; \mathbf{A})] \cdots, \mathbb{E}[\mathbf{X}_* \widehat{f}_t(\mathbf{W}_{<t,\bullet}; \mathbf{A})])^\top. \quad (202)$$

Similarly, by exploiting the fact that  $\widehat{f}_{1:t}$  form an orthonormal basis for  $\text{Span}(f_{1:t})$ , we have the following alternate formula for  $\kappa_t$ :

$$\begin{aligned} \kappa_t &\stackrel{\text{def}}{=} \min \left\{ \mathbb{E} [\{\mathbf{X}_* - g(\mathbf{W}_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_t) \right\} \\ &= 1 - \|\widehat{\alpha}\|^2 \in [0, 1]. \end{aligned} \quad (203)$$

These alternate expressions for  $\alpha, \Sigma, \kappa_t$  will be useful in our subsequent calculations.

**Useful Linear Algebraic Results.** The proof relies on the following intermediate linear algebraic claims.

$$\text{For any } \lambda \geq 0, \quad q \in \text{Range}(Q + \lambda\Gamma), \quad (204)$$

$$\text{If } \kappa_t < 1, \quad q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \neq 1 \Leftrightarrow \alpha \otimes q \in \text{Range}(\alpha\alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha\alpha^\top) \otimes \Gamma). \quad (205)$$

We defer the justification of these claims to the end of the proof.

**Corner Cases.** We begin by verifying the claim of the lemma for two corner cases. Consider the situation where  $\kappa_t = 1$ . In this case, (202) and (203) imply that  $\alpha = 0$ . As a result, using the formula for the effective SNR for a general multivariate Gaussian channel given in Lemma 3 (Appendix A.2), the effective SNR of the Gaussian channel in (197) is  $\omega_{\text{eff}} = 0$ , as claimed. Next, we consider the case when:

$$\kappa_t < 1, \quad q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q = 1.$$

From (205), we conclude that:

$$\alpha \otimes q \notin \text{Range}(\alpha\alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha\alpha^\top) \otimes \Gamma).$$

Hence, from (200), we conclude that the effective SNR of the Gaussian channel in (197) is  $\omega_{\text{eff}} = 1$ , as claimed.

**Derivation of the claimed formulae.** The remaining task is to prove the claim of the lemma when:

$$\kappa_t < 1, \quad q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \neq 1,$$

and we work under the above assumption for the remainder of the proof. From (205), we conclude that:

$$\alpha \otimes q \in \text{Range}(\alpha\alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha\alpha^\top) \otimes \Gamma).$$

From (200), we have that:

$$\omega_{\text{eff}}(\mathbf{X}_* | \mathbf{W}_{\leq t, \bullet}; \mathbf{A}) = \frac{\vartheta}{1 + \vartheta}, \quad (206)$$

$$v_{\text{opt}}(\mathbf{X}_* | \mathbf{W}_{\leq t, \bullet}; \mathbf{A}) = \frac{1}{\sqrt{\vartheta + \vartheta^2}} \cdot (\alpha\alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha\alpha^\top) \otimes \Gamma)^\dagger \cdot \alpha \otimes q \quad (207)$$

where:

$$\vartheta \stackrel{\text{def}}{=} \alpha \otimes q^\top (\alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma)^\dagger \alpha \otimes q. \quad (208)$$

To compute the pseudoinverse  $(\alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma)^\dagger$ , we define  $\bar{\alpha}$  as the unit vector along  $\hat{\alpha}$  and  $P$  as the projector orthogonal to  $\hat{\alpha}$ :

$$\bar{\alpha} \stackrel{\text{def}}{=} \frac{\hat{\alpha}}{\|\hat{\alpha}\|}, \quad P \stackrel{\text{def}}{=} I - \bar{\alpha} \bar{\alpha}^\top. \quad (209)$$

We have,

$$\begin{aligned} & (\alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma)^\dagger \stackrel{(202)}{=} (L^\dagger \otimes I)^\top \cdot [\hat{\alpha} \hat{\alpha}^\top \otimes (Q - qq^\top - \Gamma) + I \otimes \Gamma]^\dagger \cdot (L^\dagger \otimes I) \\ & \stackrel{(203),(209)}{=} (L^\dagger \otimes I)^\top \cdot [\bar{\alpha} \bar{\alpha}^\top \otimes \{(1 - \kappa_t) \cdot (Q - qq^\top) + \kappa_t \cdot \Gamma\} + P \otimes \Gamma]^\dagger \cdot (L^\dagger \otimes I) \\ & = (L^\dagger \otimes I)^\top \cdot [\bar{\alpha} \bar{\alpha}^\top \otimes \{(1 - \kappa_t) \cdot (Q - qq^\top) + \kappa_t \cdot \Gamma\}^\dagger + P \otimes \Gamma]^\dagger \cdot (L^\dagger \otimes I). \end{aligned}$$

The above formula for the pseudo-inverse yields:

$$\begin{aligned} & (\alpha \alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha \alpha^\top) \otimes \Gamma)^\dagger \cdot \alpha \otimes q = (L^{\dagger\top} \hat{\alpha}) \otimes [\{(1 - \kappa_t) \cdot (Q - qq^\top) + \kappa_t \cdot \Gamma\}^\dagger \cdot q] \\ & \stackrel{(a)}{=} (1 - \kappa_t)^{-1} \cdot \left\{ 1 - q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \right\}^{-1} \cdot (L^{\dagger\top} \hat{\alpha}) \otimes \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \\ & \stackrel{(202)}{=} (1 - \kappa_t)^{-1} \cdot \left\{ 1 - q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \right\}^{-1} \cdot (\Sigma^\dagger \alpha) \otimes \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q, \end{aligned} \quad (210)$$

where we used the Sherman-Morrison formula in (a). Combining the computation above with the definition of  $\vartheta$  in (208) yields the following formula for  $\vartheta$ :

$$\begin{aligned} \vartheta & \stackrel{(203)}{=} (\alpha^\top \Sigma^\dagger \alpha) \cdot (1 - \kappa_t)^{-1} \cdot \left\{ 1 - q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \right\}^{-1} \cdot q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \\ & \stackrel{(a)}{=} \left\{ 1 - q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \right\}^{-1} \cdot q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q \end{aligned}$$

where the equality marked (a) follows by observing that  $\alpha^\top \Sigma^\dagger \alpha = \|\hat{\alpha}\|^2 = 1 - \kappa_t$ . Plugging the expression of the previous display in (206) yields:

$$\begin{aligned} \omega_{\text{eff}}(\mathbf{X}_* | \mathbf{W}_{\leq t, \bullet}; \mathbf{A}) & = q^\top \left( Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma \right)^\dagger q, \\ v_{\text{opt}}(\mathbf{X}_* | \mathbf{W}_{\leq t, \bullet}; \mathbf{A}) & = \frac{1}{\sqrt{\omega_{\text{eff}}(\mathbf{X}_* | \mathbf{W}_{\leq t, \bullet}; \mathbf{A})}} \cdot (\Sigma^\dagger \alpha) \otimes ((1 - \kappa_t) \cdot Q + \kappa_t \cdot \Gamma)^\dagger q \end{aligned}$$

as claimed. To finish the proof, we need to prove the linear algebraic claims in (204) and (205).

**Proof of claim (204).** Fix any  $\lambda \geq 0$ . To show that  $q \in \text{Range}(Q + \lambda\Gamma)$ , it is sufficient to verify that  $q$  is orthogonal to  $\text{Null}(Q + \lambda\Gamma)$ . Towards this goal, we consider any  $v \in \text{Null}(Q + \lambda\Gamma)$  and aim to show that  $\langle v, q \rangle = 0$ . Since  $v$  lies in the null space of  $Q + \lambda\Gamma$ ,

$$v^\top (Q + \lambda\Gamma) v = 0 \implies v^\top Q v = 0,$$

where the implication is obtained by observing that  $Q \succeq 0$  and  $\Gamma \succeq 0$  (cf. (198)). Recalling the definition of  $q$  from (198), we obtain:

$$|\langle q, v \rangle|^2 = \left| \mathbb{E} \left[ \sum_{i=1}^D v_i \cdot (\Lambda_\nu^i - \mathbb{E}[\Lambda^i]) \right] \right|^2 \stackrel{(a)}{\leq} \mathbb{E} \left[ \left| \sum_{i=1}^D v_i \cdot (\Lambda_\nu^i - \mathbb{E}[\Lambda^i]) \right|^2 \right] \stackrel{(198)}{=} v^\top Q v = 0.$$

In the above display,  $\Lambda \sim \mu, \Lambda_\nu \sim \nu$  and the inequality in (a) follows from the Cauchy–Schwarz inequality. Hence,  $\langle q, v \rangle = 0$  for any  $v \in \text{Null}(Q + \lambda\Gamma)$ , which implies the claim in (204).

**Proof of claim (205).** For convenience, we introduce the definitions:

$$S \stackrel{\text{def}}{=} \alpha\alpha^\top \otimes (Q - qq^\top) + (\Sigma - \alpha\alpha^\top) \otimes \Gamma, \quad T \stackrel{\text{def}}{=} Q + \frac{\kappa_t}{1 - \kappa_t} \cdot \Gamma.$$

Assuming  $\kappa_t < 1$ , our goal is to show the equivalence:

$$q^\top T^\dagger q \neq 1 \Leftrightarrow \alpha \otimes q \in \text{Range}(S).$$

Consider the vector:

$$v \stackrel{\text{def}}{=} (L^\dagger \hat{\alpha}) \otimes (T^\dagger q).$$

A short computation reveals that:

$$Sv = (1 - \kappa_t) \cdot (1 - q^\top T^\dagger q) \cdot \alpha \otimes q, \quad \langle q, v \rangle = (1 - \kappa_t) \cdot q^\top T^\dagger q.$$

If  $q^\top T^\dagger q \neq 1$ , we have found a vector  $u = (1 - \kappa_t)^{-1} \cdot (1 - q^\top T^\dagger q)^{-1} \cdot v$  such that  $Su = \alpha \otimes q$ , which certifies that  $\alpha \otimes q \in \text{Range}(S)$ . On the other hand, if  $q^\top T^\dagger q = 1$ , then  $Sv = 0$  and we have found a vector  $v \in \text{Null}(S)$  such that  $\langle v, q \rangle = (1 - \kappa_t) \neq 0$ , which certifies that  $\alpha \otimes q \notin \text{Range}(S)$ . This concludes the proof of this lemma.  $\square$

## E.2.2 Proof of Lemma 17

*Proof of Lemma 17.* Recall from (192) that:

$$m(\kappa) \stackrel{\text{def}}{=} \begin{cases} \text{mmse}_\pi \left( q^\top \left[ Q + \frac{\kappa}{1 - \kappa} \cdot \Gamma \right]^\dagger q \right) & : \kappa \in [0, 1) \\ \text{mmse}_\pi(0) & : \kappa = 1. \end{cases} \quad (211)$$

The expressions for the vector  $q \in \mathbb{R}^D$  and the *positive semi-definite* matrices  $Q, \Gamma \in \mathbb{R}^{D \times D}$  appear in (47), but will not be needed in the proof. To show that  $m$  is a non-decreasing function, it suffices to verify that for any  $0 \leq \kappa \leq \kappa' < 1$ , we have  $m(\kappa) \leq m(\kappa') \leq m(1)$ . Indeed we have,

$$\begin{aligned} 0 \leq \kappa \leq \kappa' < 1 &\implies 0 \leq \frac{\kappa}{1 - \kappa} \leq \frac{\kappa'}{1 - \kappa'} \xrightarrow{\text{(a)}} 0 \preceq Q + \frac{\kappa}{1 - \kappa} \cdot \Gamma \preceq Q + \frac{\kappa'}{1 - \kappa'} \cdot \Gamma & (212) \\ &\xrightarrow{\text{(b)}} q^\top \left( Q + \frac{\kappa}{1 - \kappa} \cdot \Gamma \right)^\dagger q \geq q^\top \left( Q + \frac{\kappa'}{1 - \kappa'} \cdot \Gamma \right)^\dagger q \geq 0 \\ &\xrightarrow{\text{(c)}} \text{mmse}_\pi \left( q^\top \left[ Q + \frac{\kappa}{1 - \kappa} \cdot \Gamma \right]^\dagger q \right) \leq \text{mmse}_\pi \left( q^\top \left[ Q + \frac{\kappa'}{1 - \kappa'} \cdot \Gamma \right]^\dagger q \right) \leq \text{mmse}_\pi(0), \\ &\xrightarrow{\text{(d)}} m(\kappa) \leq m(\kappa') \leq m(1). \end{aligned}$$

In the above display, step (a) follows by observing that  $Q, \Gamma$  are positive semidefinite ( $\preceq$  denotes the standard Loewner partial ordering of symmetric matrices defined by the positive semidefinite cone), step (c) relies on the monotonicity of  $\text{mmse}_\pi(\cdot)$  (Fact 1). Finally, step (b) follows from the Fenchel conjugate formula for convex quadratic forms (recall from (204) that  $q \in \text{Range}(Q + \lambda\Gamma)$  for any  $\lambda \geq 0$ ):

$$\begin{aligned} q^\top \left( Q + \frac{\kappa}{1 - \kappa} \cdot \Gamma \right)^\dagger q &= \sup_{\xi \in \mathbb{R}^D} 2 \langle q, \xi \rangle - \xi^\top \left( Q + \frac{\kappa}{1 - \kappa} \cdot \Gamma \right) \xi \\ &\stackrel{(212)}{\geq} \sup_{\xi \in \mathbb{R}^D} 2 \langle q, \xi \rangle - \xi^\top \left( Q + \frac{\kappa'}{1 - \kappa'} \cdot \Gamma \right) \xi = q^\top \left( Q + \frac{\kappa'}{1 - \kappa'} \cdot \Gamma \right)^\dagger q. \end{aligned}$$

This verifies that  $m$  is non-decreasing.  $\square$

### E.2.3 Proof of Lemma 18

*Proof of Lemma 18.* Given a collection of iterate denoisers  $f_1, \dots, f_t$ , which satisfy the requirements of Definition 6, our goal is to show that:

$$\mathcal{M}_t(f_1, f_2, \dots, f_{t-1}, f_t) \geq m(\kappa_t^\sharp) \text{ and,} \quad (213)$$

$$\mathcal{M}_t(f_1, f_2, \dots, f_{t-1}, f_t^\sharp) = m(\kappa_t^\sharp), \quad (214)$$

where  $\kappa_t^\sharp \stackrel{\text{def}}{=} \text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\mathcal{M}_{t-1}(f_{1:t-1}))$  and  $f_t^\sharp$  is the DMMSE estimator for the Gaussian channel  $(\mathbf{X}_\star, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t-1,\bullet}; \mathbf{A})$  generated by the state evolution random variables associated with the lifted OAMP algorithm which uses the denoisers  $f_1, \dots, f_{t-1}$  in the first  $t-1$  iterations. From Lemma 16, we know that:

$$\mathcal{M}_t(f_1, f_2, \dots, f_t) = m(\kappa_t) \quad (215a)$$

where the scalar  $\kappa_t \in [0, 1]$  was defined as:

$$\kappa_t \stackrel{\text{def}}{=} \min \left\{ \mathbb{E} [\{\mathbf{X}_\star - g(\mathbf{W}_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_t) \right\}. \quad (215b)$$

We consider the two claims of the lemma (213) and (214) one by one.

**Proof of Claim (213).** In light of the monotonicity of  $m$  (Lemma 17), it suffices to show that

$$\kappa_t \stackrel{\text{def}}{=} \min \left\{ \mathbb{E} [\{\mathbf{X}_\star - g(\mathbf{W}_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_t) \right\} \geq \kappa_t^\sharp \stackrel{\text{def}}{=} \text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\mathcal{M}_{t-1}(f_{1:t-1})). \quad (216)$$

We can lower bound  $\kappa_t$  as follows:

$$\begin{aligned} \kappa_t &\stackrel{\text{def}}{=} \min \left\{ \mathbb{E} [\{\mathbf{X}_\star - g(\mathbf{W}_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_t) \right\} \\ &\stackrel{\text{(a)}}{\geq} \min \left\{ \mathbb{E} [\{\mathbf{X}_\star - g(\mathbf{W}_{<t,\bullet}; \mathbf{A})\}^2] : g \in L^2(\mathbf{W}_{<t,\bullet}; \mathbf{A}), \mathbb{E}[\mathbf{Z}_{s,i}g(\mathbf{W}_{<t,\bullet}; \mathbf{A})] = 0 \forall s \in [t-1], i \in [D] \right\} \\ &\stackrel{\text{def}}{=} \text{DMMSE}(\mathbf{X}_\star | \mathbf{W}_{<t,\bullet}; \mathbf{A}). \end{aligned} \quad (217)$$

In the above display, the random variables  $\mathbf{Z}_{<t,\bullet}$  represent the Gaussian noise random variables for the Gaussian channel  $(\mathbf{X}_\star, \mathbf{W}_{<t,\bullet}; \mathbf{A})$ . To obtain the inequality in step (a), we used the fact that any  $g \in \text{Span}(f_1, \dots, f_t)$  satisfies the divergence-free constraints:

$$\mathbb{E}[\mathbf{Z}_{s,i} \cdot g(\mathbf{W}_{<t,\bullet}; \mathbf{A})] = 0 \forall s \in [t-1], i \in [D].$$

Indeed, by Gaussian integration by parts, for any  $i < t$ ,  $s < t$  and  $j \in [D]$ ,

$$\mathbb{E}[\mathbf{Z}_{s,j} \cdot f_i(\mathbf{W}_{<i,\bullet}; \mathbf{A})]$$

can be expressed as a linear combination of the expected derivatives

$$\{\mathbb{E}[\partial_{\tau,\ell} f_i(\mathbf{W}_{<i,\bullet}; \mathbf{A})] : \tau < i, \ell \in [D]\},$$

which are all identically zero thanks to the divergence-free requirement imposed on the iterate denoisers  $f_1, \dots, f_{t-1}$  in Definition 6. We now cast the lower bound on  $\kappa_t$  obtained in (217) in the desired form (216). Using the DMMSE formula for general multivariate Gaussian channels (stated as Lemma 3 in Appendix Appendix A.2) we obtain:

$$\kappa_t \geq \text{DMMSE}(\mathbf{X}_\star | \mathbf{W}_{<t,\bullet}; \mathbf{A}) = \text{dmmse}_\pi(\omega_{\text{eff}}),$$

where  $\omega_{\text{eff}}$  is the effective SNR of the Gaussian channel  $(\mathbf{X}_\star, \mathbf{W}_{<t,\bullet}; \mathbf{A})$ . On the other hand, the MMSE formula for general multivariate Gaussian channels (stated as Lemma 3 in Appendix Appendix A.2) yields:

$$\mathcal{M}_{t-1}(f_1, \dots, f_{t-1}) \stackrel{\text{def}}{=} \text{MMSE}(\mathbf{X}_\star | \mathbf{W}_{<t,\bullet}; \mathbf{A}) = \text{mmse}_\pi(\omega_{\text{eff}}).$$

Since  $\text{mmse}_\pi(\cdot)$  is a strictly decreasing function (Fact 1), it has a well-defined inverse<sup>3</sup>  $\text{mmse}_\pi^{-1}$ , which can be used to relate the previous two displays:

$$\kappa_t \geq \text{dmmse}_\pi(\omega_{\text{eff}}) = \text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\mathcal{M}_{t-1}(f_1, \dots, f_{t-1})) \stackrel{\text{def}}{=} \kappa_t^\sharp, \quad (218)$$

which proves the claim made in (216) and concludes the proof of the first part of the lemma.

**Proof of Claim (214).** In light of the formula in (215), it suffices to show that:

$$\min \left\{ \mathbb{E} [\{X_\star - g(W_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_{t-1}, f_t^\sharp) \right\} = \kappa_t^\sharp.$$

From (218), we already know that:

$$\min \left\{ \mathbb{E} [\{X_\star - g(W_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_{t-1}, f_t^\sharp) \right\} \geq \kappa_t^\sharp$$

Hence, we only need to show the upper bound:

$$\min \left\{ \mathbb{E} [\{X_\star - g(W_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_{t-1}, f_t^\sharp) \right\} \leq \kappa_t^\sharp. \quad (219)$$

We upper bound the LHS in (219) by taking  $g = f_t^\sharp$  and use the DMMSE formula for a general multivariate Gaussian channel (see Lemma 3 in Appendix A.2) to obtain:

$$\begin{aligned} \min \left\{ \mathbb{E} [\{X_\star - g(W_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1, \dots, f_{t-1}, f_t^\sharp) \right\} &\leq \mathbb{E}[(X_\star - f_t^\sharp(W_{<t,\bullet}; \mathbf{A}))^2] \\ &\stackrel{\text{(a)}}{=} \text{DMMSE}(X_\star | W_{<t,\bullet}; \mathbf{A}) \\ &\stackrel{\text{Lem. 3}}{=} \text{dmmse}_\pi(\omega_{\text{eff}}) \stackrel{\text{(218)}}{=} \kappa_t^\sharp, \end{aligned}$$

where step (a) follows by recalling that  $f_t^\sharp$  was the DMMSE estimator for  $(X_\star, W_{1,\bullet}, \dots, W_{t-1,\bullet}; \mathbf{A})$ . This concludes the proof of the lemma.  $\square$

### E.3 Proof of Proposition 5

We begin by reminding the reader that the optimal degree- $D$  lifted OAMP algorithm was given by:

$$\mathbf{w}_{t,i} = (\mathbf{Y}^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i] \cdot \mathbf{I}_N) \cdot f_t^\star(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t-1,\bullet}; \mathbf{a}) \quad \forall i \in [D], t \in \mathbb{N}, \quad (220)$$

At the end of  $t$  iterations, the optimal lifted OAMP algorithm estimates  $\mathbf{x}_\star$  by:

$$\widehat{\mathbf{w}}_t^{(D)} = h_t^\star(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t,\bullet}; \mathbf{a}).$$

**Description of the iterate denoisers and post-processing functions.** The iterate denoisers functions  $f_1^\star, f_2^\star, \dots$  and the post-processing functions  $h_1^\star, h_2^\star, \dots$  were defined recursively: once  $f_1^\star, \dots, f_t^\star$  have been specified, let  $(X_\star, W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$  denote the state evolution random variables corresponding to the first  $t$  iterations of the resulting lifted OAMP algorithm. Then, the iterate denoiser for step  $t+1$  is the DMMSE estimator for the Gaussian channel  $(X_\star, W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A})$ . Lemma 3 in Appendix A.2 provides a general formula for the DMMSE estimator of a multivariate Gaussian channel. An application of this formula yields:

$$f_{t+1}^\star(w; a) \stackrel{\text{def}}{=} \bar{\varphi} \left( \left\langle v_t^{(D)}, w \right\rangle; a \mid \omega_t^{(D)} \right) \quad \forall w \in \mathbb{R}^{tD}, a \in \mathbb{R}^k \quad (221a)$$

<sup>3</sup>Strictly speaking, Fact 1 guarantees that  $\text{mmse}_\pi : [0, 1] \mapsto [0, 1]$  is strictly increasing only when  $\text{mmse}_\pi(0) = \mathbb{E}[\text{Var}[X_\star | \mathbf{A}]] > 0$ . In the corner case when  $\text{mmse}_\pi(0) = \mathbb{E}[\text{Var}[X_\star | \mathbf{A}]] = 0$ , notice that  $\text{dmmse}_\pi(0) = \mathbb{E}[\text{Var}[X_\star | \mathbf{A}]] = 0$ . Since  $\text{mmse}_\pi$  and  $\text{dmmse}_\pi$  are non-negative, non-increasing functions (cf. Fact 1 and Lemma 4), this means that  $\text{mmse}_\pi(\omega) = \text{dmmse}_\pi(\omega) = 0$  for any  $\omega \in [0, 1]$ . In this situation, we can define  $\text{mmse}_\pi^{-1}$  arbitrarily. Irrespective of our convention,  $\text{dmmse}_\pi \circ \text{mmse}_\pi^{-1}(\omega) = 0$ , since  $\text{dmmse}_\pi$  is a constant function which takes the value 0.



where:

$$v_t^{(D)} \stackrel{\text{def}}{=} v_{\text{opt}}(\mathbf{X}_* | \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A}), \quad \omega_t^{(D)} \stackrel{\text{def}}{=} \omega_{\text{eff}}(\mathbf{X}_* | \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A}) \quad (221b)$$

denote the optimal linear combination and effective SNR of the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  (as defined in Lemma 3). Moreover, the estimator at the end of  $t$  iterations is obtained by applying the MMSE estimator for the channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  to the iterates:

$$\widehat{\mathbf{w}}_t^{(D)} = h_t^*(\mathbf{w}_{\leq t, \bullet}; a), \quad (221c)$$

$$h_t^*(w; a) \stackrel{\text{def}}{=} \varphi\left(\left\langle v_t^{(D)}, w \right\rangle; a \mid \omega_t^{(D)}\right) \quad \forall w \in \mathbb{R}^{tD}, a \in \mathbb{R}^k \quad (221d)$$

**Dynamics of the Optimal Lifted OAMP Algorithm.** By Corollary 1, the joint distribution of  $(\mathbf{X}_*, \mathbf{W}_{\leq t, \bullet}; \mathbf{A})$  is given by:

$$(\mathbf{X}_*; \mathbf{A}) \sim \pi, \quad (222a)$$

$$(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}) | (\mathbf{X}_*; \mathbf{A}) \sim \mathcal{N}\left(\mathbf{X}_* \cdot \alpha_t \otimes q^{(D)}, \alpha_t \alpha_t^\top \otimes (Q^{(D)} - q^{(D)} q^{(D)\top}) + (\Sigma_t - \alpha_t \alpha_t^\top) \otimes \Gamma^{(D)}\right). \quad (222b)$$

In the above display, the entries of  $q^{(D)} \in \mathbb{R}^D$ ,  $Q^{(D)} \in \mathbb{R}^{D \times D}$ , and  $\Gamma^{(D)} \in \mathbb{R}^{D \times D}$  are given by:

$$q_i^{(D)} = \mathbb{E}[\Lambda_\nu^i] - \mathbb{E}[\Lambda^i], \quad Q_{ij}^{(D)} = \mathbb{E}\left[(\Lambda_\nu^i - \mathbb{E}[\Lambda^i]) \cdot (\Lambda_\nu^j - \mathbb{E}[\Lambda^j])\right], \quad \Gamma_{ij}^{(D)} = \text{Cov}[\Lambda^i, \Lambda^j] \quad \text{where } \Lambda \sim \mu, \Lambda_\nu \sim \nu, \quad (222c)$$

the entries of  $\alpha_t \in \mathbb{R}^t$  and  $\Sigma_t \in \mathbb{R}^{t \times t}$  are given by:

$$(\alpha_t)_s \stackrel{\text{def}}{=} \mathbb{E}[\mathbf{X}_* \cdot f_s^*(\mathbf{W}_{< s, \bullet}; \mathbf{A})], \quad (\Sigma_t)_{s\tau} \stackrel{\text{def}}{=} \mathbb{E}[f_s^*(\mathbf{W}_{< s, \bullet}; \mathbf{A}) \cdot f_\tau^*(\mathbf{W}_{< \tau, \bullet}; \mathbf{A})] \quad \forall s, \tau \in [t], \quad (222d)$$

and  $\otimes$  denotes the Kronecker (or tensor) product for matrices. Finally, the asymptotic MSE of the estimator  $\widehat{\mathbf{w}}_t^{(D)}$  returned after  $t$  iterations is given by:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_*\|^2}{N} = \mathbb{E}\left[(\mathbf{X}_* - h_t^*(\mathbf{W}_{\leq t, \bullet}; \mathbf{A}))^2\right] = \text{MMSE}(\mathbf{X}_* | \mathbf{W}_{\leq t, \bullet}; \mathbf{A}) \stackrel{\text{Lem. 3}}{=} \text{mmse}_\pi(\omega_t^{(D)}). \quad (222e)$$

The following lemma provides a recursive formula for the effective SNR  $\omega_t^{(D)} \stackrel{\text{def}}{=} \omega_{\text{eff}}(\mathbf{X}_* | \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  and the optimal linear combination  $v_t^{(D)} \stackrel{\text{def}}{=} v_{\text{opt}}(\mathbf{X}_* | \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$  at step  $t$ .

**Lemma 19.** Let  $\omega_0^{(D)} \stackrel{\text{def}}{=} 0$  and  $\mathbf{d}_0^{(D)} \stackrel{\text{def}}{=} \text{dmmse}_\pi(0) = \text{dmmse}_\pi(0) = \mathbb{E}[\text{Var}[\mathbf{X}_* | \mathbf{A}]]$ .

1. If  $\mathbf{d}_0^{(D)} < 1$ , then  $\omega_t^{(D)}$  and  $v_t^{(D)}$  admit the following recursive characterization:

$$\begin{aligned} \mathbf{d}_t^{(D)} &= \text{dmmse}_\pi(\omega_{t-1}^{(D)}), \\ \omega_t^{(D)} &= q^{(D)\top} \left[ Q^{(D)} + \frac{\mathbf{d}_t^{(D)}}{1 - \mathbf{d}_t^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)}, \\ v_t^{(D)} &= (\omega_t^{(D)})^{-\frac{1}{2}} \cdot (1 - \mathbf{d}_t^{(D)})^{-1} \cdot e_t \otimes \left[ Q^{(D)} + \frac{\mathbf{d}_t^{(D)}}{1 - \mathbf{d}_t^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)}, \end{aligned}$$

where  $e_t = (0, 0, \dots, 1)^\top \in \mathbb{R}^t$  denotes the last standard basis vector. Moreover, for each  $t \in \mathbb{N}$ ,

$$\omega_t^{(D)} > 0, \quad \mathbf{d}_t^{(D)} < 1 \quad \forall t \in \mathbb{N}.$$

2. If  $\mathbf{d}_0^{(D)} = 1$ , then  $\omega_t^{(D)} = 0$  for all  $t \in \mathbb{N}$ .

*Proof.* See Appendix E.3.1. □

**Eliminating Corner Cases.** Next, we observe that the proof of Proposition 5 is quite simple in the corner cases  $\text{mmse}_\pi(0) = \mathbb{E}[\text{Var}[X_\star|A]] = 1$  and  $\text{mmse}_\pi(0) = \mathbb{E}[\text{Var}[X_\star|A]] = 0$ .

**Case 1:**  $\text{mmse}_\pi(0) = 1$ . In this case, combining (222e) and Lemma 19 (Claim 2), we conclude that for any  $t \in \mathbb{N}$ ,

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_\star\|^2}{N} = \text{mmse}_\pi(\omega_t^{(D)}) = \text{mmse}_\pi(0) = 1.$$

On the other hand, the optimal OAMP algorithm (17) returns the estimator  $\widehat{\mathbf{x}}_t = \mathbf{0}$  for any  $t \in \mathbb{N}$ . Hence,

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{x}}_t - \mathbf{x}_\star\|^2}{N} = \text{plim}_{N \rightarrow \infty} \frac{\|\mathbf{x}_\star\|^2}{N} = \mathbb{E}[X_\star^2] = 1 = \text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_\star\|^2}{N},$$

as claimed.

**Case 2:**  $\text{mmse}_\pi(0) = 0$ . In this situation, we observe that for any  $t \in \mathbb{N}$ ,

$$0 \leq \text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_\star\|^2}{N} = \text{mmse}_\pi(\omega_t^{(D)}) \stackrel{\text{Fact 1}}{\leq} \text{mmse}_\pi(0) = 0 \implies \text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_\star\|^2}{N} = 0.$$

On the other hand, the simplified OAMP algorithm returns the estimator  $\widehat{\mathbf{x}}_t = \varphi(\mathbf{a}|0)$  for any  $t \in \mathbb{N}$ . Hence,

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{x}}_t - \mathbf{x}_\star\|^2}{N} = \mathbb{E}[(X_\star - \varphi(A|0))^2] = \text{mmse}_\pi(0) = 0 = \text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_\star\|^2}{N},$$

as claimed.

Hence, in the remainder of the proof, we will assume that:

$$\text{mmse}_\pi(0) = \text{dmmse}_\pi(0) \in (0, 1). \quad (223)$$

**Simplifying the Optimal Lifted OAMP Algorithm.** Recall the estimator computed at the end of  $t$  iterations of the optimal lifted OAMP algorithm is:

$$\widehat{\mathbf{w}}_t = h_t^\star(\mathbf{w}_{1,\bullet}, \dots, \mathbf{w}_{t,\bullet}),$$

where  $h_t^\star$  is the MMSE estimator for the channel  $(X_\star, W_{1,\bullet}, \dots, W_{t,\bullet})$ :

$$h_t^\star(w; a) \stackrel{\text{def}}{=} \varphi\left(\left\langle v_t^{(D)}, w \right\rangle; a | \omega_t^{(D)}\right) \quad \forall w \in \mathbb{R}^{tD}, a \in \mathbb{R}^k.$$

To compute this estimator, we only need to track the following vector  $\mathbf{x}_t^{(D)}$ , which is a linear combination of the lifted OAMP iterates  $\mathbf{w}_{\leq t, \bullet}$  with weights given by  $v_t^{(D)}$ :

$$\mathbf{x}_t^{(D)} \stackrel{\text{def}}{=} \sum_{s=1}^t \sum_{i=1}^D (v_t^{(D)})_{s,i} \cdot \mathbf{w}_{s,i}. \quad (224)$$

Indeed, the estimator returned by the optimal lifted OAMP algorithm can be computed using  $\mathbf{x}_t^{(D)}$ :

$$\widehat{\mathbf{w}}_t^{(D)} = \varphi(\mathbf{x}_t^{(D)}; \mathbf{a} | \omega_t^{(D)}).$$

In fact,  $\mathbf{x}_t^{(D)}$  can be tracked using a simpler iterative algorithm. To see this, we observe using the formula for  $v_t^{(D)}$  from Lemma 19 (Claim 1) that we have:

$$\mathbf{x}_t^{(D)} \stackrel{\text{def}}{=} \sum_{i=1}^D (v_t^{(D)})_{s,i} \cdot \mathbf{w}_{s,i} = \frac{1}{(1 - \mathbf{d}_t^{(D)}) \sqrt{\omega_t^{(D)}}} \sum_{i=1}^D \left\{ \left( Q^{(D)} + \frac{\mathbf{d}_t^{(D)}}{1 - \mathbf{d}_t^{(D)}} \cdot \Gamma^{(D)} \right)^\dagger q^{(D)} \right\}_i \cdot \mathbf{w}_{t,i} \quad (225a)$$

$$\stackrel{(220)}{=} \frac{1}{(1 - \mathbf{d}_t^{(D)}) \sqrt{\omega_t^{(D)}}} \left[ \sum_{i=1}^D \left\{ \left( Q^{(D)} + \frac{\mathbf{d}_t^{(D)}}{1 - \mathbf{d}_t^{(D)}} \cdot \Gamma^{(D)} \right)^\dagger q^{(D)} \right\}_i \cdot \{ \mathbf{Y}^i - \mathbb{E}_{\Lambda \sim \mu} [\Lambda^i] \cdot \mathbf{I}_N \} \right] \cdot f_t^*(\mathbf{w}_{<t;\bullet}; \mathbf{a}). \quad (225b)$$

In light of the above formula, we introduce the degree- $D$  polynomial matrix denoising function:

$$\Psi_\star^{(D)}(\lambda; \rho) \stackrel{\text{def}}{=} \sum_{i=1}^D \left\{ \left( Q^{(D)} + \frac{1}{\rho} \cdot \Gamma^{(D)} \right)^\dagger q^{(D)} \right\}_i \cdot (\lambda^i - \mathbb{E}_{\Lambda \sim \mu} [\Lambda^i]) \quad \forall \lambda \in \mathbb{R}, \rho \in (0, \infty], \quad (226)$$

as well as the sequence  $(\rho_t^{(D)})_{t \in \mathbb{N}}$ :

$$\rho_t^{(D)} \stackrel{\text{def}}{=} \frac{1}{\mathbf{d}_t^{(D)}} - 1, \quad \mathbf{d}_t^{(D)} = \frac{1}{\text{dmmse}_\pi(\omega_t^{(D)})} - 1 \quad \forall t \in \mathbb{N}. \quad (227)$$

This leads to the following formula for  $\mathbf{x}_t^{(D)}$ :

$$\begin{aligned} \mathbf{x}_t^{(D)} &= \frac{1}{\sqrt{\omega_t^{(D)}}} \left( \frac{1}{\rho_t^{(D)}} + 1 \right) \cdot \Psi_\star^{(D)}(\mathbf{Y}; \rho_t^{(D)}) \cdot f_t^*(\mathbf{w}_{1;\bullet}, \dots, \mathbf{w}_{t-1;\bullet}; \mathbf{a}) \\ &\stackrel{(221),(224)}{=} \frac{1}{\sqrt{\omega_t^{(D)}}} \left( \frac{1}{\rho_t^{(D)}} + 1 \right) \cdot \Psi_\star^{(D)}(\mathbf{Y}; \rho_t^{(D)}) \cdot \bar{\varphi}(\mathbf{x}_{t-1}^{(D)}; \mathbf{a} | \omega_{t-1}^{(D)}) \end{aligned}$$

In summary, we have obtained an OAMP algorithm (in the sense of Definition 4):

$$\mathbf{x}_t^{(D)} = \frac{1}{\sqrt{\omega_t^{(D)}}} \left( \frac{1}{\rho_t^{(D)}} + 1 \right) \cdot \Psi_\star^{(D)}(\mathbf{Y}; \rho_t^{(D)}) \cdot \bar{\varphi}(\mathbf{x}_{t-1}^{(D)}; \mathbf{a} | \omega_{t-1}^{(D)}) \quad \forall t \in \mathbb{N}, \quad (228)$$

which can reconstruct the estimator  $\hat{\mathbf{w}}_t^{(D)}$  computed by the optimal degree- $D$  lifted OAMP algorithm in (220):

$$\hat{\mathbf{w}}_t^{(D)} = \varphi(\mathbf{x}_t^{(D)}; \mathbf{a} | \omega_t^{(D)}) \quad \forall t \in \mathbb{N}.$$

**Large-degree limit.** Notice the similarity between the simplified version of the optimal degree- $D$  lifted OAMP algorithm in (228) and the optimal OAMP algorithm introduced in (17). Our goal will be to show that the performance of the simplified OAMP algorithm in (228) converges to the performance of the optimal OAMP algorithm in (17) as  $D \rightarrow \infty$ . To do so, we will rely on the following lemma which shows that the matrix denoiser  $\Psi_\star^{(D)}$  used in the simplified OAMP algorithm (228) can be viewed as the minimizer of a variational problem over degree- $D$  polynomials. By studying the same variational problem over  $L^2(\mu + \nu)$ , we obtain a candidate for the optimal matrix denoiser in the large degree limit, which turns out to be precisely the matrix denoiser (17e) used by the optimal OAMP algorithm in (17).

**Lemma 20.** Let  $\mathcal{L}_D$  and  $\mathcal{L}_\infty$  denote the following functions on the domain  $(0, \infty)$ :

$$\mathcal{L}_D(\rho) \stackrel{\text{def}}{=} \inf_{\Psi \in \mathcal{P}_D} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0, \quad (229)$$

$$\mathcal{L}_\infty(\rho) \stackrel{\text{def}}{=} \inf_{\Psi \in L^2(\mu + \nu)} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0, \quad (230)$$

where  $\Lambda \sim \mu, \Lambda_\nu \sim \nu$ ,  $\mathcal{P}_D$  denotes the set of all polynomial functions on  $\mathbb{R}$  with degree at most  $D$ , and  $L^2(\mu + \nu)$  denotes the set of all real valued functions on  $\mathbb{R}$  which are square integrable with respect to  $\mu + \nu$ . Then, we have:

1. For any  $\rho \in (0, \infty)$ , the function  $\Psi_\star^{(D)}(\cdot; \rho)$  defined in (226) is the minimizer of the variational problem in (229) and,

$$\mathcal{L}_D(\rho) = 1 - q_D^\top \cdot \left( Q^{(D)} + \frac{1}{\rho} \cdot \Gamma^{(D)} \right)^\dagger \cdot q^{(D)} = 1 - \mathbb{E}[\Psi_\star^{(D)}(\Lambda_\nu; \rho)].$$

2. For any  $\rho \in (0, \infty)$ , the function:

$$\Psi_\star(\lambda; \rho) = 1 - \left( \mathbb{E}_{\Lambda \sim \mu} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \frac{\phi(\lambda)}{\phi(\lambda) + \rho} \quad \forall \lambda \in \mathbb{R}, \rho \in (0, \infty),$$

is a minimizer of the variational problem in (230). In the above display, for any  $\lambda \in \mathbb{R}$ ,  $\phi(\lambda) \stackrel{\text{def}}{=} (1 - \pi\theta\mathcal{H}_\mu(\lambda))^2 + \pi^2\theta^2\mu^2(\lambda)$ , where  $\mathcal{H}_\mu$  is the Hilbert transform of  $\mu$  and  $\theta$  is the SNR of the spiked matrix model (1).

3. For any  $D \in \mathbb{N} \cup \{\infty\}$ , the function  $\mathcal{L}_D$  maps the interval  $(0, \infty)$  to  $(0, 1]$ .
4. Let  $(\rho_D)_{D \in \mathbb{N}}$  be any non-decreasing sequence in  $(0, \infty)$  which converges to  $\rho \in (0, \infty)$  as  $D \rightarrow \infty$ . Then,  $\mathcal{L}_D(\rho_D) \downarrow \mathcal{L}_\infty(\rho)$  as  $D \rightarrow \infty$ . Moreover,  $\Psi_\star^{(D)}(\cdot; \rho_D)$  converges to  $\Psi_\star(\cdot; \rho)$  in  $L^2(\mu + \nu)$  as  $D \rightarrow \infty$ .
5. For any  $\rho \in (0, \infty)$ , the function  $\Psi_\star(\cdot; \rho)$  satisfies  $\mathbb{E}[\Psi_\star(\Lambda_\nu; \rho)] = 1 - \mathcal{L}_\infty(\rho)$ .

*Proof.* The proof of this lemma is provided in Appendix E.3.2. □

**Proof of Proposition 5.** We now have all the ingredients to complete the proof of Proposition 5. Recall that the optimal estimator that can be computed using  $t$  iterations of a degree- $D$  lifted OAMP algorithm is:

$$\widehat{\mathbf{w}}_t^{(D)} = \varphi(\mathbf{x}_t^{(D)}; \mathbf{a}|\omega_t^{(D)}),$$

where the iterates  $\mathbf{x}_t$  are computed using the update rule:

$$\mathbf{x}_t^{(D)} = \frac{1}{\sqrt{\omega_t^{(D)}}} \left( \frac{1}{\rho_t^{(D)}} + 1 \right) \cdot \Psi_\star^{(D)}(\mathbf{Y}; \rho_t^{(D)}) \cdot \bar{\varphi}(\mathbf{x}_{t-1}^{(D)}; \mathbf{a}|\omega_{t-1}^{(D)}) \quad \forall t \in \mathbb{N}. \quad (231)$$

Furthermore, the limiting mean-squared error of the estimator  $\widehat{\mathbf{w}}_t^{(D)}$  is given by:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_\star\|^2}{N} = \text{mmse}_\pi(\omega_t^{(D)}). \quad (232)$$

In the above equations, the parameters  $\{\omega_t^{(D)}\}_{t \in \mathbb{N}}, \{\mathbf{d}_t^{(D)}\}_{t \in \mathbb{N}}, \{\rho_t^{(D)}\}_{t \in \mathbb{N}}$  are updated according to the following recursion from Lemma 19 (also recall the definition of  $\rho_t^{(D)}$  from (227)) initialized with  $\omega_0^{(D)} = 0$ :

$$\mathbf{d}_t^{(D)} = \text{dmmse}_\pi(\omega_{t-1}^{(D)}), \quad \rho_t^{(D)} = \frac{1}{\mathbf{d}_t^{(D)}} - 1, \quad \omega_t^{(D)} = (q^{(D)})^\top \left[ Q^{(D)} + \frac{1}{\rho_t^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)} \stackrel{\text{(a)}}{=} 1 - \mathcal{L}_D(\rho_t^{(D)}), \quad (233)$$

where the equality marked (a) follows from Lemma 20 (Claim (2)). To prove Proposition 5, we need to show that in the large degree limit ( $D \rightarrow \infty$ ), the performance of the estimator  $\widehat{\mathbf{w}}_t^{(D)}$  is identical to the performance of the estimator  $\widehat{\mathbf{x}}_t$ :

$$\widehat{\mathbf{x}}_t \stackrel{\text{def}}{=} \varphi(\mathbf{x}_t; \mathbf{a}|\omega_t) \quad \forall t \in \mathbb{N},$$

returned by the optimal OAMP algorithm from (17):

$$\mathbf{x}_t = \frac{1}{\sqrt{\omega_t}} \left( 1 + \frac{1}{\rho_t} \right) \cdot \Psi_\star(\mathbf{Y}; \rho_t) \cdot \bar{\varphi}(\mathbf{x}_{t-1}; \mathbf{a} | \omega_{t-1}) \quad \forall t \in \mathbb{N}. \quad (234)$$

In the above equation, the parameters  $\{\omega_t\}_{t \in \mathbb{N}}, \{\mathbf{d}_t\}_{t \in \mathbb{N}}, \{\rho_t\}_{t \in \mathbb{N}}$  are updated according to the following recursion initialized with  $\omega_0 = 0$ :

$$\mathbf{d}_t = \text{dmmse}_\pi(\omega_{t-1}), \quad \rho_t = \frac{1}{\mathbf{d}_t} - 1, \quad \omega_t = 1 - \mathcal{L}_\infty(\rho_t). \quad (235)$$

Let  $(\mathbf{X}_\star, (\mathbf{X}_t)_{t \in \mathbb{N}}; \mathbf{A})$  denote the state evolution random variables corresponding to the above algorithm. We will show by induction that for each  $t \in \mathbb{N}$ ,

1. for any  $D \in \mathbb{N}$ , the sequences the sequences  $\{\omega_t^{(D)}\}_{t \in \mathbb{N}}, \{\mathbf{d}_t^{(D)}\}_{t \in \mathbb{N}}, \{\rho_t^{(D)}\}_{t \in \mathbb{N}}$  generated by (233) satisfy:

$$\omega_t^{(D)} \in [0, 1), \quad \mathbf{d}_t^{(D)} \in (0, 1), \quad \rho_t^{(D)} \in (0, \infty) \quad \forall t \in \mathbb{N}. \quad (236)$$

In particular, the update equation for  $\rho_t^{(D)}$  in (233) is well-defined (avoids division by zero).

2.  $\omega_t^{(D)} \uparrow \omega_t$ ,  $\mathbf{d}_t^{(D)} \downarrow \mathbf{d}_t$ , and  $\rho_t^{(D)} \uparrow \rho_t$  as  $D \rightarrow \infty$ .

These claims immediately imply Proposition 5:

$$\text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{x}}_t - \mathbf{x}_\star\|^2}{N} \stackrel{\text{Prop. 1}}{=} \text{mmse}_\pi(\omega_t) \stackrel{(a)}{=} \lim_{D \rightarrow \infty} \text{mmse}_\pi(\omega_t^{(D)}) \stackrel{(232)}{=} \lim_{D \rightarrow \infty} \text{plim}_{N \rightarrow \infty} \frac{\|\widehat{\mathbf{w}}_t^{(D)} - \mathbf{x}_\star\|^2}{N},$$

as claimed. In the above display step (a) follows from the claim  $\omega_t^{(D)} \uparrow \omega_t$  as  $D \rightarrow \infty$ . We now prove the three claims made above. As the induction hypothesis, we assume the claims hold for some  $t \in \mathbb{N}$  and verify that they continue to hold at step  $t + 1$ .

**Proof of Claim (1).** We focus on proving (236) holds at iteration  $t + 1$ , assuming it holds at iteration  $t$ . The induction hypothesis  $\omega_t^{(D)} \geq 0$  and the monotonicity of  $\text{dmmse}_\pi(\cdot)$  (see Lemma 4 in Appendix A.2) imply that  $\mathbf{d}_{t+1}^{(D)} = \text{dmmse}_\pi(\omega_t^{(D)}) \leq \text{dmmse}_\pi(0) = \text{mmse}_\pi(0) < 1$  (recall (223)). On the other hand, since  $\omega_t^{(D)} < 1$ , by the strict monotonicity of  $\text{mmse}_\pi$  (see Fact 1 in Appendix A.2), we conclude that  $\mathbf{d}_{t+1}^{(D)} = \text{dmmse}_\pi(\omega_t^{(D)}) \geq \text{mmse}_\pi(\omega_t^{(D)}) > \text{mmse}_\pi(1) = 0$ . Hence,  $\mathbf{d}_{t+1}^{(D)} \in (0, 1)$ . Since  $\rho_{t+1}^{(D)} = (\mathbf{d}_{t+1}^{(D)})^{-1} - 1$ , we conclude that  $\rho_{t+1}^{(D)} \in (0, \infty)$ . Finally, since  $\omega_{t+1} = 1 - \mathcal{L}_D(\rho_{t+1}^{(D)})$  and Lemma 20 guarantees that  $\mathcal{L}_D : (0, \infty) \mapsto (0, 1]$ , we conclude that  $\omega_{t+1}^{(D)} \in [0, 1)$ , as claimed in (236). This proves the first claim.

**Proof of Claim (2).** We recall from (233) that  $\mathbf{d}_{t+1}^{(D)} = \text{dmmse}_\pi(\omega_t^{(D)})$ . By the induction hypothesis  $\omega_t \uparrow \omega_t \in (0, 1)$ . Since  $\text{dmmse}_\pi(\cdot)$  is a non-increasing continuous function (see Lemma 4 in Appendix A.2), we conclude that  $\mathbf{d}_{t+1} \downarrow \text{dmmse}_\pi(\omega_t) \stackrel{(235)}{=} \mathbf{d}_{t+1} \in (0, 1)$  and  $\rho_{t+1}^{(D)} \stackrel{(233)}{=} (\mathbf{d}_{t+1}^{(D)})^{-1} - 1 \uparrow (\mathbf{d}_{t+1})^{-1} - 1 \stackrel{(235)}{=} \rho_t$ . Finally, recalling the update rule for  $\omega_{t+1}^{(D)}$  from (233) and using Lemma 20 (Claim 4), we conclude that:  $\omega_{t+1}^{(D)} = 1 - \mathcal{L}_D(\rho_{t+1}^{(D)}) \uparrow 1 - \mathcal{L}_\infty(\rho_{t+1}) \stackrel{(235)}{=} \omega_{t+1}$ . This proves the second claim. This concludes the proof of Proposition 5.

### E.3.1 Proof of Lemma 19

The proof of Lemma 19 relies on the following intermediate result.

**Claim 3.** For any  $s, t \in \mathbb{N}$  with  $s \leq t$ :

$$\mathbb{E}[f_s^\star(W_{1,\bullet}, \dots, W_{s-1,\bullet}; \mathbf{A}) \cdot f_t^\star(W_{1,\bullet}, \dots, W_{t-1,\bullet}; \mathbf{A})] = \mathbb{E}[X_\star \cdot f_s^\star(W_{1,\bullet}, \dots, W_{s-1,\bullet}; \mathbf{A})].$$

We will first prove Lemma 19 assuming this claim. The proof of Claim 3 is provided at the end of this section.

*Proof of Lemma 19.* Recall that for each  $t \in \mathbb{N}$ ,

$$v_t^{(D)} \stackrel{\text{def}}{=} v_{\text{opt}}(X_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A}), \quad \omega_t^{(D)} \stackrel{\text{def}}{=} \omega_{\text{eff}}(X_\star | W_{1,\bullet}, \dots, W_{t,\bullet}; \mathbf{A}),$$

where  $(X_\star, (W_{t,\bullet})_{t \in \mathbb{N}}; \mathbf{A})$  are the state evolution random variables associated with the optimal lifted OAMP algorithm. In addition, we will find it helpful to introduce the quantity:

$$d_t^{(D)} \stackrel{\text{def}}{=} \min \left\{ \mathbb{E} [\{X_\star - g(W_{<t,\bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1^\star, \dots, f_t^\star) \right\}.$$

Since  $\text{Span}(f_1^\star) \subset \text{Span}(f_1^\star, f_2^\star) \subset \dots \subset \text{Span}(f_1^\star, \dots, f_t^\star) \subset \dots$ ,  $(d_t^{(D)})_{t \in \mathbb{N}}$  is a non-increasing sequence.

**Case 1:**  $d_0^{(D)} < 1$ . Notice that since  $(d_t^{(D)})_{t \in \mathbb{N}}$  is a non-increasing sequence and  $d_0^{(D)} < 1$ , we conclude that  $d_t^{(D)} < 1$  for all  $t \in \mathbb{N}$ , which is one of the claims pertaining to this case made in the lemma. To prove the other claims, we will show by induction on  $t$  that:

$$d_t^{(D)} = \text{dmmse}_\pi(\omega_{t-1}^{(D)}), \tag{237a}$$

$$\omega_t^{(D)} = q_D^\top \left[ Q^{(D)} + \frac{d_t^{(D)}}{1 - d_t^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)}, \tag{237b}$$

$$\omega_t^{(D)} > 0, \tag{237c}$$

$$v_t^{(D)} = (\omega_t^{(D)})^{-\frac{1}{2}} \cdot (1 - d_t^{(D)})^{-1} \cdot e_1 \otimes \left[ Q^{(D)} + \frac{d_t^{(D)}}{1 - d_t^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)}. \tag{237d}$$

As our induction hypothesis we assume that (237) holds, and show that (237) also holds for  $t + 1$ . Indeed, since  $f_{t+1}^\star$  is the DMMSE estimator for the Gaussian channel  $(X_\star, W_{\leq t, \bullet}; \mathbf{A})$ , by Lemma 18,

$$d_{t+1}^{(D)} \stackrel{\text{def}}{=} \min \left\{ \mathbb{E} [\{X_\star - g(W_{\leq t, \bullet}; \mathbf{A})\}^2] : g \in \text{Span}(f_1^\star, \dots, f_{t+1}^\star) \right\} = \text{dmmse}_\pi(\omega_t^{(D)}).$$

By appealing to Lemma 16, we conclude that:

$$\omega_{t+1}^{(D)} = q_D^\top \left[ Q^{(D)} + \frac{d_{t+1}^{(D)}}{1 - d_{t+1}^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)} > 0$$

$$v_{t+1}^{(D)} = (\omega_{t+1}^{(D)})^{-\frac{1}{2}} \cdot (1 - d_{t+1}^{(D)})^{-1} \cdot (\Sigma_{t+1}^\dagger \alpha_{t+1}) \otimes \left[ Q^{(D)} + \frac{d_{t+1}^{(D)}}{1 - d_{t+1}^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)}.$$

From Claim 3, we know that:

$$(\alpha_{t+1})_s = (\Sigma_{t+1})_{s,t+1} \quad \forall s \in [t+1] \implies \alpha_{t+1} = \Sigma_{t+1} e_{t+1} \implies \Sigma_{t+1}^\dagger \alpha_{t+1} = e_{t+1}.$$

Hence,

$$v_{t+1}^{(D)} = (\omega_{t+1}^{(D)})^{-\frac{1}{2}} \cdot (1 - d_{t+1}^{(D)})^{-1} \cdot e_{t+1} \otimes \left[ Q^{(D)} + \frac{d_{t+1}^{(D)}}{1 - d_{t+1}^{(D)}} \cdot \Gamma^{(D)} \right]^\dagger q^{(D)},$$

as desired. This proves the claim of the lemma when  $d_0^{(D)} < 1$ .

**Case 2:**  $d_0^{(D)} = 1$ . In this situation, we will show by induction that for any  $t \in \mathbb{N}$ ,

$$d_t^{(D)} = 1, \quad \omega_t^{(D)} = 0.$$

We assume that the above claim holds at step  $t$  as the induction hypothesis. As in the previous case, by Lemma 18, we have:

$$d_{t+1}^{(D)} = \text{dmmse}_\pi(\omega_t^{(D)}) = \text{dmmse}_\pi(0) = d_0^{(D)} = 1.$$

Appealing to Lemma 16, we conclude that  $\omega_{t+1}^{(D)} = 0$ , as desired. This proves the claim of the lemma.  $\square$

We now present the proof of Claim 3.

*Proof of Claim 3.* By the Tower property, we have:

$$\mathbb{E}[\mathbf{X}_* \cdot f_s^*(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{s-1,\bullet}; \mathbf{A})] = \mathbb{E}[\mathbb{E}[\mathbf{X}_* | \mathbf{W}_{<t,\bullet}, \mathbf{A}] f_s^*(\mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{s-1,\bullet}; \mathbf{A})]. \quad (238)$$

Since  $\mathbb{E}[\mathbf{X}_* | \mathbf{W}_{<t,\bullet}, \mathbf{A}]$  is the MMSE estimator for the channel  $(\mathbf{X}_*, \mathbf{W}_{<t,\bullet}; \mathbf{A})$ . Using the formula for the MMSE estimator for a general multivariate Gaussian channel (see Lemma 3 in Appendix A.2), we obtain:

$$\mathbb{E}[\mathbf{X}_* | \mathbf{W}_{<t,\bullet}, \mathbf{A}] = \varphi(\langle \mathbf{W}_{<t,\bullet}, v_{\text{opt}} \rangle; \mathbf{A} | \omega_{\text{eff}}), \quad (239)$$

where we use the shorthand notations  $\omega_{\text{eff}}$  and  $v_{\text{opt}}$  denote the effective SNR and optimal linear combination<sup>4</sup> for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t-1,\bullet}; \mathbf{A})$ . Similarly, Lemma 3 gives the following formula for  $f_t^*$ , the DMMSE estimator for this channel:

$$f_t^*(\mathbf{W}_{<t,\bullet}; \mathbf{A}) = \bar{\varphi}(\langle \mathbf{W}_{<t,\bullet}, v_{\text{opt}} \rangle; \mathbf{A} | \omega_{\text{eff}}). \quad (240)$$

Recall the formula for  $\bar{\varphi}$  (the DMMSE estimator in a scalar Gaussian channel) provided in Definition 3 (cf. (11)):

$$\bar{\varphi}(x; a | \omega) \stackrel{\text{def}}{=} \begin{cases} (1 - \sqrt{\omega} \cdot \beta(\omega))^{-1} \cdot (\varphi(x; a | \omega) - \beta(\omega) \cdot x) & : \omega < 1 \\ \varphi(x; a | \omega) & : \omega = 1 \end{cases}, \quad (241)$$

where:

$$\beta(\omega) \stackrel{\text{def}}{=} \frac{1}{\sqrt{1-\omega}} \cdot \mathbb{E}[\mathbf{Z} \varphi(\sqrt{\omega} \mathbf{X}_* + \sqrt{1-\omega} \mathbf{Z}; \mathbf{A} | \omega)] \quad \text{where } (\mathbf{X}_*; \mathbf{A}) \sim \pi, \mathbf{Z} | \mathbf{X}_*; \mathbf{A} \sim \mathcal{N}(0, 1).$$

Combining (239), (240), and (241) yields the following formula relating the MMSE and DMMSE estimators for the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$ :

$$\mathbb{E}[\mathbf{X}_* | \mathbf{W}_{<t,\bullet}, \mathbf{A}] = \begin{cases} (1 - \sqrt{\omega_{\text{eff}}} \cdot \beta(\omega_{\text{eff}})) \cdot f_t^*(\mathbf{W}_{<t,\bullet}; \mathbf{A}) + \beta(\omega_{\text{eff}}) \cdot \langle v_{\text{opt}}, \mathbf{W}_{<t,\bullet} \rangle & : \omega_{\text{eff}} < 1 \\ f_t^*(\mathbf{W}_{<t,\bullet}; \mathbf{A}) & : \omega_{\text{eff}} = 1. \end{cases}$$

We substitute the above expression in (238). If  $\omega_{\text{eff}} = 1$ , the claim of the lemma is immediate. If  $\omega_{\text{eff}} < 1$ , we have:

$$\begin{aligned} \mathbb{E}[\mathbf{X}_* \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})] &= (1 - \sqrt{\omega_{\text{eff}}} \cdot \beta(\omega_{\text{eff}})) \cdot \mathbb{E}[f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A}) \cdot f_t^*(\mathbf{W}_{<t,\bullet}; \mathbf{A})] + \beta(\omega_{\text{eff}}) \cdot \mathbb{E}[\langle v_{\text{opt}}, \mathbf{W}_{<t,\bullet} \rangle \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})]. \end{aligned}$$

Let  $\mathbf{Z}_{1,\bullet}, \dots, \mathbf{Z}_{t-1,\bullet}$  denote the Gaussian noise random variables in the Gaussian channel  $(\mathbf{X}_*, \mathbf{W}_{1,\bullet}, \dots, \mathbf{W}_{t,\bullet}; \mathbf{A})$ . Since  $(\mathbf{X}_*, \langle v_{\text{opt}}, \mathbf{W}_{<t,\bullet} \rangle; \mathbf{A})$  forms a scalar Gaussian channel with SNR  $\omega_{\text{eff}}$ :

$$\langle v_{\text{opt}}, \mathbf{W}_{<t,\bullet} \rangle = \sqrt{\omega_{\text{eff}}} \cdot \mathbf{X}_* + \sqrt{1-\omega_{\text{eff}}} \cdot \mathbf{Z}_{\text{eff}}, \quad \mathbf{Z}_{\text{eff}} \stackrel{\text{def}}{=} \frac{\langle v_{\text{opt}}, \mathbf{Z}_{<t,\bullet} \rangle}{\sqrt{1-\omega_{\text{eff}}}}.$$

Hence,

$$\begin{aligned} \mathbb{E}[\mathbf{X}_* \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})] &= (1 - \sqrt{\omega_{\text{eff}}} \cdot \beta(\omega_{\text{eff}})) \cdot \mathbb{E}[f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A}) \cdot f_t^*(\mathbf{W}_{<t,\bullet}; \mathbf{A})] \\ &\quad + \beta(\omega_{\text{eff}}) \cdot \mathbb{E}[(\sqrt{\omega_{\text{eff}}} \cdot \mathbf{X}_* + \sqrt{1-\omega_{\text{eff}}} \cdot \mathbf{Z}_{\text{eff}}) \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})] \\ &= (1 - \sqrt{\omega_{\text{eff}}} \cdot \beta(\omega_{\text{eff}})) \cdot \mathbb{E}[f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A}) \cdot f_t^*(\mathbf{W}_{<t,\bullet}; \mathbf{A})] + \sqrt{\omega_{\text{eff}}} \cdot \beta(\omega_{\text{eff}}) \cdot \mathbb{E}[\mathbf{X}_* \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})]. \quad (242) \end{aligned}$$

The final equality in the previous display follows by observing that:

$$\mathbb{E}[\mathbf{Z}_{\text{eff}} \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})] = \mathbb{E}[\mathbb{E}[\mathbf{Z}_{\text{eff}} | \mathbf{X}_*, \mathbf{Z}_{<s,\bullet}, \mathbf{A}] \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})] = \mathbb{E}[\mathbb{E}[\mathbf{Z}_{\text{eff}} | \mathbf{Z}_{<s,\bullet}] \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})] = 0.$$

This is because  $\mathbb{E}[\mathbb{E}[\mathbf{Z}_{\text{eff}} | \mathbf{Z}_{<s,\bullet}]]$  is a linear combination of  $\mathbf{Z}_{<s,\bullet}$  (since  $(\mathbf{Z}_{<s,\bullet}, \mathbf{Z}_{\text{eff}})$  are jointly Gaussian) and the DMMSE estimator  $f_s^*$  is uncorrelated with the channel noise  $\mathbf{Z}_{<s,\bullet}$ . Rearranging (242) gives us:

$$\mathbb{E}[\mathbf{X}_* \cdot f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A})] = \mathbb{E}[f_s^*(\mathbf{W}_{<s,\bullet}; \mathbf{A}) \cdot f_t^*(\mathbf{W}_{<t,\bullet}; \mathbf{A})],$$

which is precisely the claim set out to prove.  $\square$

<sup>4</sup>We suppress the dependence of these quantities on the iteration number  $t$  and the degree  $D$  of the lifted OAMP for notational convenience since it does not play a role in the proof.

### E.3.2 Proof of Lemma 20

*Proof of Lemma 20.* We consider each of the claims made in the lemma.

**Proof of Claim (1).** To solve the variational problem:

$$\mathcal{L}_D(\rho) \stackrel{\text{def}}{=} \inf_{\Psi \in \mathcal{P}_D} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0,$$

we parameterize any polynomial  $f$  of degree at most  $D$  which satisfies  $\mathbb{E}[\Psi(\Lambda)] = 0$  as:

$$\Psi(\lambda) = \sum_{i=0}^D v_i \cdot (\lambda^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i]) \quad \forall \lambda \in \mathbb{R}.$$

Hence,

$$\mathcal{L}_D(\rho) = \min_{v \in \mathbb{R}^D} \mathbb{E} \left[ \sum_{i=1}^D v_i \cdot (\Lambda_\nu^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i]) - 1 \right]^2 + \frac{1}{\rho} \cdot \mathbb{E} \left[ \sum_{i=1}^D v_i \cdot (\Lambda^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i]) \right]^2. \quad (243)$$

Expanding the square norms and recalling that the entries of  $q^{(D)} \in \mathbb{R}^D$ ,  $Q^{(D)} \in \mathbb{R}^{D \times D}$ , and  $\Gamma^{(D)} \in \mathbb{R}^{D \times D}$  are given by:

$$q_i^{(D)} = \mathbb{E}[\Lambda_\nu^i] - \mathbb{E}[\Lambda^i], \quad Q_{ij}^{(D)} = \mathbb{E}[(\Lambda_\nu^i - \mathbb{E}[\Lambda^i]) \cdot (\Lambda_\nu^j - \mathbb{E}[\Lambda^j])], \quad \Gamma_{ij}^{(D)} = \text{Cov}[\Lambda^i, \Lambda^j] \quad \text{where } \Lambda \sim \mu, \Lambda_\nu \sim \nu, \quad (244)$$

we obtain the following formula for  $\mathcal{L}_D(\rho)$ :

$$\mathcal{L}_D(\rho) = \min_{v \in \mathbb{R}^D} 1 + v^\top Q^{(D)} v - 2 \langle q^{(D)}, v \rangle + \frac{v^\top \Gamma^{(D)} v}{\rho}.$$

The optimizer of the above quadratic form is:

$$v = \left( Q^{(D)} + \frac{1}{\rho} \cdot \Gamma^{(D)} \right)^\dagger \cdot q^{(D)}.$$

Hence, a minimizer of (243) is:

$$\Psi(\lambda) = \sum_{i=1}^D \left\{ \left( Q^{(D)} + \frac{1}{\rho} \cdot \Gamma^{(D)} \right)^\dagger \cdot q^{(D)} \right\}_i \cdot (\lambda^i - \mathbb{E}_{\Lambda \sim \mu}[\Lambda^i]) \stackrel{(226)}{=} \Psi_\star^{(D)}(\lambda; \rho) \quad \forall \lambda \in \mathbb{R},$$

and:

$$\mathcal{L}_D(\rho) = 1 - q^{(D)\top} \left( Q^{(D)} + \frac{1}{\rho} \cdot \Gamma_D \right)^\dagger q^{(D)} = 1 - \langle q^{(D)}, v \rangle = 1 - \mathbb{E}[\Psi_\star^{(D)}(\Lambda_\nu; \rho)] \quad \text{where } \Lambda_\nu \sim \nu,$$

as claimed.

**Proof of Claim (2).** Our goal is to solve the optimization problem:

$$\mathcal{L}_\infty(\rho) \stackrel{\text{def}}{=} \inf_{\Psi \in L^2(\mu + \nu)} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0. \quad (245)$$

This is exactly the problem considered in (35c), where we showed that the minimizer of (245) is:

$$\Psi(\lambda) = 1 - \left( \mathbb{E}_{\Lambda \sim \mu} \left[ \frac{\phi(\Lambda)}{\phi(\Lambda) + \rho} \right] \right)^{-1} \cdot \frac{\phi(\lambda)}{\phi(\lambda) + \rho} \stackrel{\text{def}}{=} \Psi_\star(\lambda; \rho) \quad \forall \lambda \in \mathbb{R}, \quad (246)$$

as claimed.



**Proof of Claim (3).** We need to show that for any  $D \in \mathbb{N} \cup \{\infty\}$ , the function  $\mathcal{L}_D$  maps the interval  $(0, \infty)$  to  $(0, 1]$ . We will show this claim for  $\mathcal{L}_\infty$  and the exact same argument works for  $\mathcal{L}_D$  for any  $D \in \mathbb{N}$ . The definition of  $\mathcal{L}_\infty(\rho)$  (cf. (245)) implies that  $\mathcal{L}_\infty(\rho) \geq 0$ . Taking  $\Psi = 0$  in (245) shows that  $\mathcal{L}_\infty(\rho) \leq 1$ . To prove Claim (3), we need to verify that if  $\rho \in (0, \infty)$  then,  $\mathcal{L}_\infty(\rho) > 0$ . We prove this by contradiction. If  $\mathcal{L}_\infty(\rho) = 0$ , then the function  $\Psi_\star(\cdot; \rho)$ , which minimizes the objective in (245) satisfies  $\mathbb{E}[|\Psi_\star(\Lambda_\nu; \rho) - 1|^2] + \rho^{-1} \cdot \mathbb{E}[|\Psi_\star(\Lambda; \rho)|^2] = 0$ . This means that:

$$\nu(\{\lambda \in \mathbb{R} : \Psi_\star(\lambda; \rho) = 1\}) = 1, \quad \mu(\{\lambda \in \mathbb{R} : \Psi_\star(\lambda; \rho) = 0\}) = 1.$$

In particular  $\mu, \nu$  are mutually singular. However, this contradicts Lemma 1 which tells us that  $\nu_{\parallel}$ , the absolutely continuous part of  $\nu$  with respect to the Lebesgue measure has density  $\frac{d\nu_{\parallel}}{d\lambda}(\lambda) = \frac{\mu(\lambda)}{\phi(\lambda)}$ , and hence  $\nu_{\parallel} \ll \mu$ .

**Proof of Claim (4).** Consider a non-decreasing sequence  $(\rho_D)_{D \in \mathbb{N}}$  which converges to  $\rho \in (0, \infty)$  as  $D \rightarrow \infty$ . We begin by observing that  $\mathcal{L}_D(\rho_D)$  is a non-increasing sequence:

$$\begin{aligned} \mathcal{L}_{D+1}(\rho_{D+1}) &\stackrel{(243)}{=} \inf_{\Psi \in \mathcal{P}_{D+1}} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho_{D+1}} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0 \\ &\stackrel{(a)}{\leq} \inf_{\Psi \in \mathcal{P}_{D+1}} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho_D} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0 \\ &\stackrel{(b)}{\leq} \inf_{\Psi \in \mathcal{P}_D} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho_D} \cdot \mathbb{E}[\Psi^2(\Lambda)] \quad \text{subject to} \quad \mathbb{E}[\Psi(\Lambda)] = 0 \stackrel{(243)}{=} \mathcal{L}_D(\rho_D). \end{aligned}$$

In the above display, step (a) follows by observing that  $\rho_D \leq \rho_{D+1}$ . Step (b) relies on the observation that  $\mathcal{P}_D \subset \mathcal{P}_{D+1}$  and the infimum over larger sets is smaller. Hence  $(\mathcal{L}_D(\rho_D))_{D \in \mathbb{N}}$  is a non-increasing sequence. To show that  $\mathcal{L}_D(\rho_D) \downarrow \mathcal{L}_\infty(\rho)$ , we begin by observing that by repeating the argument used in the previous display, we can also obtain:

$$\mathcal{L}_D(\rho_D) \geq \mathcal{L}_\infty(\rho).$$

Hence, we have obtained a lower bound for  $\mathcal{L}_D(\rho_D)$  in terms of  $\mathcal{L}_\infty(\rho)$ . To show that in fact  $\mathcal{L}_D(\rho_D) \downarrow \mathcal{L}_\infty(\rho)$  we will also need an upper bound. Towards this goal, we consider a sequence of functions  $\{\widehat{\Psi}_D\}_{D \in \mathbb{N}}$  indexed by  $D \in \mathbb{N}$  such that  $\widehat{\Psi}_D : \mathbb{R} \mapsto \mathbb{R}$  is a degree- $D$  polynomial and:

$$\lim_{D \rightarrow \infty} \|\widehat{\Psi}_D - \Psi_\star(\cdot; \rho)\|_{\mu+\nu}^2 = 0 \quad \Leftrightarrow \quad \lim_{D \rightarrow \infty} \left( \mathbb{E}_{\Lambda \sim \mu} |\widehat{\Psi}_D(\Lambda) - \Psi_\star(\Lambda; \rho)|^2 + \mathbb{E}_{\Lambda_\nu \sim \nu} |\widehat{\Psi}_D(\Lambda_\nu) - \Psi_\star(\Lambda_\nu; \rho)|^2 \right) = 0. \quad (247)$$

The existence of these polynomial approximations follows by the fact that polynomials form a dense subset of  $L^2(\mu + \nu)$  since  $\mu + \nu$  is a compactly supported measure [70, Corollary 14.24, Definition 14.1]. We can use the function  $\lambda \mapsto \widehat{\Psi}_D(\lambda) - \mathbb{E}[\widehat{\Psi}_D(\Lambda)]$  as a candidate minimizer to upper bound  $\mathcal{L}_D(\rho_D)$ :

$$\mathcal{L}_\infty(\rho) \leq \mathcal{L}_D(\rho_D) \leq \mathbb{E} \left[ \left| \widehat{\Psi}_D(\Lambda_\nu) - 1 - \mathbb{E}[\widehat{\Psi}_D(\Lambda)] \right|^2 \right] + \frac{1}{\rho_D} \cdot \mathbb{E} \left[ \left| \widehat{\Psi}_D(\Lambda) - \mathbb{E}[\widehat{\Psi}_D(\Lambda)] \right|^2 \right]$$

We let  $D \rightarrow \infty$  in the above display and exploit (247) and  $\rho_D \rightarrow \rho \in (0, \infty)$  to conclude that:

$$\lim_{D \rightarrow \infty} \mathcal{L}_D(\rho_D) = \mathbb{E}[|\Psi_\star(\Lambda_\nu; \rho) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi_\star^2(\Lambda; \rho)] = \mathcal{L}_\infty(\rho),$$

as claimed. Finally, we show that:

$$\lim_{D \rightarrow \infty} \mathbb{E}_{\Lambda \sim \mu} |\Psi_\star^{(D)}(\Lambda; \rho_D) - \Psi_\star(\Lambda; \rho)|^2 = 0, \quad \lim_{D \rightarrow \infty} \mathbb{E}_{\Lambda_\nu \sim \nu} |\Psi_\star^{(D)}(\Lambda_\nu) - \Psi_\star(\Lambda_\nu; \rho_D)|^2 = 0$$

To do so, we will exploit the strong convexity of the objective function:

$$\mathcal{O}(\Psi) \stackrel{\text{def}}{=} \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi^2(\Lambda)]. \quad (248)$$

Notice that  $\mathcal{O}(\Psi) - \mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2]$  is convex and hence for any  $\lambda \in (0, 1)$  and any  $\Psi, \widehat{\Psi} \in L^2(\mu + \nu)$ ,

$$\begin{aligned} & \mathcal{O}(\lambda\Psi + (1-\lambda)\widehat{\Psi}) - \mathbb{E}[|\lambda\Psi(\Lambda_\nu) + (1-\lambda)\widehat{\Psi}(\Lambda_\nu) - 1|^2] \\ & \leq \lambda\mathcal{O}(\Psi) + (1-\lambda)\mathcal{O}(\widehat{\Psi}) - \lambda\mathbb{E}[|\Psi(\Lambda_\nu) - 1|^2] - (1-\lambda)\mathbb{E}[|\widehat{\Psi}(\Lambda_\nu) - 1|^2]. \end{aligned}$$

Rearranging yields:

$$\mathbb{E}[|\Psi(\Lambda_\nu) - \widehat{\Psi}(\Lambda_\nu)|^2] \leq \frac{\lambda\mathcal{O}(\Psi) + (1-\lambda)\mathcal{O}(\widehat{\Psi}) - \mathcal{O}(\lambda\Psi + (1-\lambda)\widehat{\Psi})}{\lambda(1-\lambda)}.$$

We instantiate the above inequality with  $\lambda = 1/2$ ,  $\Psi = \Psi_\star(\cdot; \rho)$  and  $\widehat{\Psi} = \Psi_\star^{(D)}(\cdot; \rho_D)$ . For notational convenience, we define the ‘‘midpoint’’ between these two matrix denoisers as  $\Phi^{(D)} \stackrel{\text{def}}{=} (\Psi_\star(\cdot; \rho) + \Psi_\star^{(D)}(\cdot; \rho_D))/2$ :

$$\mathbb{E}[|\Psi(\Lambda_\nu) - \widehat{\Psi}(\Lambda_\nu)|^2] \leq 2 \cdot \left[ \mathcal{O}(\Psi_\star(\cdot; \rho)) + \mathcal{O}(\Psi_\star^{(D)}(\cdot; \rho_D)) - 2\mathcal{O}(\Phi^{(D)}) \right] \quad (249)$$

Recalling that  $\Psi_\star(\cdot; \rho)$  is the minimizer of the variational problem that defines  $\mathcal{L}_\infty(\rho)$  (claim 2 of the lemma):

$$\mathcal{O}(\Psi_\star(\cdot; \rho)) \stackrel{(248)}{=} \mathbb{E}[|\Psi_\star(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi_\star^2(\Lambda)] = \mathcal{L}_\infty(\rho). \quad (250)$$

Moreover, since  $\Phi^{(D)}$  is a feasible point for the optimization problem that defines  $\mathcal{L}_\infty(\rho)$  (recall (245)), we conclude that:

$$\mathcal{O}(\Phi^{(D)}) \stackrel{(248)}{=} \mathbb{E}[|\Phi^{(D)}(\Lambda_\nu) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Phi^{(D)}(\Lambda)^2] \stackrel{(245)}{\geq} \mathcal{L}_\infty(\rho). \quad (251)$$

Recall that  $\rho_D \uparrow \rho$  and that  $\Psi_\star^{(D)}(\cdot; \rho_D)$  is the minimizer of the variational problem that defines  $\mathcal{L}_D(\rho_D)$ . Hence,

$$\begin{aligned} & \mathcal{O}(\Psi_\star^{(D)}(\cdot; \rho_D)) \stackrel{(248)}{=} \mathbb{E}[|\Psi_\star^{(D)}(\Lambda_\nu; \rho_D) - 1|^2] + \frac{1}{\rho} \cdot \mathbb{E}[\Psi_\star^{(D)}(\Lambda; \rho_D)^2] \\ & \leq \mathbb{E}[|\Psi_\star^{(D)}(\Lambda_\nu; \rho_D) - 1|^2] + \frac{1}{\rho_D} \cdot \mathbb{E}[\Psi_\star^{(D)}(\Lambda; \rho_D)^2] = \mathcal{L}_D(\rho_D). \end{aligned} \quad (252)$$

Plugging in (250), (251), and (252) into (249), we obtain:

$$\mathbb{E}[|\Psi(\Lambda_\nu) - \widehat{\Psi}(\Lambda_\nu)|^2] \leq 2 \cdot (\mathcal{L}_D(\rho_D) - \mathcal{L}_\infty(\rho)).$$

We have already shown that  $\mathcal{L}_D(\rho_D) \rightarrow \mathcal{L}_\infty(\rho)$  as  $D \rightarrow \infty$ . Hence,

$$\lim_{D \rightarrow \infty} \mathbb{E}[|\Psi(\Lambda_\nu) - \widehat{\Psi}(\Lambda_\nu)|^2] = 0. \quad (253)$$

An analogous argument yields  $\lim_{D \rightarrow \infty} \mathbb{E}[|\Psi(\Lambda) - \widehat{\Psi}(\Lambda)|^2] = 0$ , which completes the proof of Claim (4) in the statement of the lemma.

**Proof of Claim (5).** Recall that from the first claim of the lemma:

$$1 - \mathbb{E}_{\Lambda \sim \nu}[\Psi_\star^{(D)}(\Lambda; \rho)] = \mathcal{L}_D(\rho).$$

Taking  $D \rightarrow \infty$  and using the fact that  $\mathcal{L}_D(\rho) \rightarrow \mathcal{L}_\infty(\rho)$  and (253) yields the desired conclusion. This concludes the proof of the lemma.  $\square$

## F Some Miscellaneous Results

*Fact 2* (Benaych-Georges and Nadakuditi 12, Proposition 9.3). Let  $\mathbf{W}$  be a  $N \times N$  random noise matrix which satisfies Assumption 2 and let  $\mathbf{u}$  be a  $N$ -dimensional random vector which is independent of  $\mathbf{W}$  and satisfies:  $\|\mathbf{u}\|^2/N \xrightarrow{\mathbb{P}} r$  as  $N \rightarrow \infty$ . Then, for any continuous function  $f : \mathbb{R} \mapsto \mathbb{R}$ ,

$$\frac{\mathbf{u}^\top f(\mathbf{W})\mathbf{u}}{N} \xrightarrow{\mathbb{P}} r \cdot \mathbb{E}[f(\Lambda)] \quad \Lambda \sim \mu,$$

where  $\mu$  is the limiting spectral distribution of  $\mathbf{W}$ .

*Proof.* Recall from Assumption 2 that the eigen-decomposition of  $\mathbf{W}$  is given by:

$$\mathbf{W} = \mathbf{U} \cdot \text{diag}(\lambda_1(\mathbf{W}), \dots, \lambda_N(\mathbf{W})) \cdot \mathbf{U}^\top,$$

where the matrix of eigenvectors  $\mathbf{U} \sim \text{Unif}(\mathbb{O}(N))$  is a Haar-distributed random orthogonal matrix independent of the eigenvalues  $\lambda_1(\mathbf{W}), \dots, \lambda_N(\mathbf{W})$ . Defining  $\mathbf{v} \stackrel{\text{def}}{=} \mathbf{U}^\top \mathbf{u} / \|\mathbf{u}\| \sim \text{Unif}(\mathbb{S}^{N-1})$ , we find that:

$$\frac{\mathbf{u}^\top f(\mathbf{W})\mathbf{u}}{N} = \frac{\|\mathbf{u}\|^2}{N} \cdot \sum_{i=1}^N v_i^2 \cdot f(\lambda_i(\mathbf{W})).$$

Since  $\mathbf{v} \sim \text{Unif}(\mathbb{S}^{N-1})$ , by [12, Proposition 9.3] we have that:

$$\sum_{i=1}^N v_i^2 \cdot f(\lambda_i(\mathbf{W})) \xrightarrow{\mathbb{P}} \mathbb{E}[f(\Lambda)], \quad \Lambda \sim \mu.$$

The claim follows by combining the above conclusion with the hypothesis  $\|\mathbf{u}\|^2/N \xrightarrow{\mathbb{P}} r$  using Slutsky's theorem.  $\square$