



HAL
open science

Marvels and pitfalls of the Langevin algorithm in noisy high-dimensional inference

Stefano Sarao Mannelli, Giulio Biroli, Chiara Cammarota, Florent Krzakala, Pierfrancesco Urbani, Lenka Zdeborova

► **To cite this version:**

Stefano Sarao Mannelli, Giulio Biroli, Chiara Cammarota, Florent Krzakala, Pierfrancesco Urbani, et al.. Marvels and pitfalls of the Langevin algorithm in noisy high-dimensional inference. *Physical Review X*, 2020, 10, pp.011057. cea-02009687

HAL Id: cea-02009687

<https://cea.hal.science/cea-02009687v1>

Submitted on 6 Feb 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Marvels and Pitfalls of the Langevin Algorithm in Noisy High-dimensional Inference

Stefano Sarao Mannelli^a, Giulio Biroli^b, Chiara Cammarota^c, Florent Krzakala^b,
Pierfrancesco Urbani^a, and Lenka Zdeborová^a

^aInstitut de physique théorique, Université Paris Saclay, CNRS, CEA, 91191 Gif-sur-Yvette, France

^bLaboratoire de Physique Statistique, CNRS & Université Pierre & Marie Curie & Ecole Normale Supérieure & PSL Université, 75005 Paris, France

^cDepartment of Mathematics, King's College London, Strand London WC2R 2LS, UK

Abstract

Gradient-descent-based algorithms and their stochastic versions have widespread applications in machine learning and statistical inference. In this work we perform an analytic study of the performances of one of them, the Langevin algorithm, in the context of noisy high-dimensional inference. We employ the Langevin algorithm to sample the posterior probability measure for the spiked matrix-tensor model. The typical behaviour of this algorithm is described by a system of integro-differential equations that we call the Langevin state evolution, whose solution is compared with the one of the state evolution of approximate message passing (AMP). Our results show that, remarkably, the algorithmic threshold of the Langevin algorithm is sub-optimal with respect to the one given by AMP. We conjecture this phenomenon to be due to the residual glassiness present in that region of parameters. Finally we show how a landscape-annealing protocol, that uses the Langevin algorithm but violate the Bayes-optimality condition, can approach the performance of AMP.

Contents

1	Motivation	2
2	The spiked matrix-tensor model	3
3	Bayes-optimal estimation and message-passing	4
4	Langevin Algorithm and its Analysis	5
5	Behavior of the Langevin algorithm	6
6	Glassy nature of the Langevin-hard phase	8
7	Better performance by annealing the landscape	8
8	Perspectives	9
	Appendix	15
A	Definition of the spiked matrix-tensor model	15

B	Approximate Message Passing, state evolution and phase diagrams	16
B.1	Approximate Message Passing and Bethe free entropy	16
B.2	Averaged free entropy and its proof	18
B.3	State evolution of AMP and its analysis	24
B.4	Phase diagrams of spiked matrix-tensor model	26
C	Langevin Algorithm and its state evolution	28
C.1	Dynamical Mean-Field Equations	29
C.2	Integro-differential equations	31
D	Numerical solution of the LSE equations	32
D.1	Fixed time-grid $(2 + p)$ -spin	32
D.2	Dynamical time-grid $(2 + p)$ -spin	33
D.3	Numerical checks on the dynamical algorithm	39
D.4	Extrapolation procedure	41
D.5	Annealing protocol	41
E	Glassy nature of the Langevin hard phase: the replica approach	43
E.1	Computation of the complexity through the replica method	43
E.2	Breakdown of the fluctuation-dissipation theorem in the Langevin hard phase	48

1 Motivation

Algorithms based on noisy variants of gradients descent [1, 2] stand at the roots of many modern applications of data science, and are being used in a wide range of high-dimensional non-convex optimization problems. The widespread use of stochastic gradient descent in deep learning [3] is certainly one of the most prominent examples. For such algorithms, the existing theoretical analysis mostly concentrate on convex functions, convex relaxations or on regimes where spurious local minima become irrelevant. For problems with complicated landscapes where, instead, useful convex relaxations are not known and spurious local minima cannot be ruled out, the theoretical understanding of the behaviour of gradient-descent-based algorithm remains poor and represents a major avenue of research.

The goal of this paper is to contribute to such an understanding in the context of statistical learning, and to transfer ideas and techniques developed for glassy dynamics [4] to the analysis of non-convex relaxations in high-dimensional inference. In statistical learning, the minimization of a cost function is not the goal per se, but rather a way to uncover an unknown structure in the data. One common way to model and analyze this situation is to generate data with a hidden structure, and to see if the structure can be recovered. This is easily set up as a *teacher-student* scenario [5, 6]: *First* a teacher generates latent variables and uses them as input of a prescribed model to generate a synthetic dataset. *Then*, the student observes the dataset and tries to infer the values of the latent variables. The analysis of this setting has been carried out rigorously in a wide range of teacher-student models for high-dimensional inference and learning tasks as diverse as planted clique [7], generalized linear models such as compressed sensing or phase retrieval [8], factorization of matrices and tensors [9, 10] or simple models of neural networks [11]. In these works, the information theoretically optimal performances—the one obtained by an ideal Bayes-optimal estimator, not limited in time and memory—have been computed.

The main question is, of course, how *practical algorithms*—operating in polynomial time with respect to the problem size—compare to these ideal performances. The last decade brought remarkable progress into our understanding of the performances achievable computationally. In particular, many algorithms based on message passing [12, 6], spectral methods [13], and semidefinite programs (SDP) [14] were analyzed. Depending on the signal-to-noise ratio, these algorithms were shown to be very efficient in many of those task. Interestingly, all these algorithm fail to reach good performance in the same region of the parameter space, and this striking

observation has led to the identification of a well-defined *hard phase*. This is a regime of parameters in which the underlying statistical problem can be information-theoretically solved, but no efficient algorithms are known, rendering the problem essentially unsolvable for large instances. This stream of ideas is currently gaining momentum and impacting research in statistics, probability, and computer science.

The performance of the noisy-gradient descent algorithms — that are certainly the one currently most used in practice— remains an entirely open question. Do they allow to reach the same performances as message passing and SDPs? Can they enter the hard phase, do they stop to be efficient at the same moment as the other approaches, or are they worse? The ambition of the present paper is to address these questions by analyzing the performance of the Langevin algorithm in the high-dimensional limit of a spiked matrix-tensor model, defined in detail in the next section.

We argue that this spiked matrix-tensor problem is both generic and relevant. Similar models have played a fundamental role in statistics and random matrix theory [15, 16]. Tensor factorization is also an important topic in mathematics and machine learning [17, 18, 19, 20, 21], and is a widely used in data analysis [22]. At variance with the pure spiked tensor case [17], this mixed matrix-tensor model has the advantage that algorithmic threshold appears at the same scale as the information-theoretic one, similarly to what is observed in simple models of neural networks [8, 11].

In the present paper, we analyse the behavior of a noisy gradient descent algorithm, also called *Langevin algorithm*. We explicitly compare its performance to the one of the Bayes optimal estimator and to the best known efficient algorithm so-far – the approximate message passing algorithm [12, 6]. In particular, contrary to what has been anticipated in [23, 24], but as surmised in [25], we observe that the performance of the Langevin algorithm is hampered by the many spurious metastable states still present in the AMP-easy phase. We show that the Langevin algorithm can approach AMP performances if one uses a landscape-annealing protocol in which instead of reducing the temperature, as done in thermal annealing, one instead anneals the strength of the contribution of the tensor. A number of intriguing conclusions can be drawn by these results and are likely to be valid beyond the considered model.

Finally, the possibility to describe analytically the behavior of the Langevin algorithm in this model is enabled by the existence of the Crisanti-Horner-Sommers-Cugliandolo-Kurchan (CHSCK) equations in spin glass theory, describing the behavior of the Langevin dynamics in the so-called spherical p -spin model [26, 27], where the method can be rigorously justified [28]. These equations were a key development in the field of statistical physics of disordered systems that lead to detailed understanding and predictions about the slow dynamics of glasses [4]. In this paper, we bring these powerful methods and ideas into the realm of statistical learning.

2 The spiked matrix-tensor model

We now detail the spiked matrix-tensor problem: a *teacher* generates a N -dimensional vector x^* by choosing each of its components independently from a normal Gaussian distribution of zero mean and unit variance. In the large N limit this is equivalent to have a flat distribution over the N -dimensional hypersphere \mathcal{S}_{N-1} defined by $|x^*|^2 = N$. The teacher then also generates a symmetric matrix Y_{ij} and a symmetric order- p tensor T_{i_1, \dots, i_p} as

$$\begin{aligned} Y_{ij} &= \frac{1}{\sqrt{N}} x_i^* x_j^* + \xi_{ij} \quad \forall i < j, \\ T_{i_1 \dots i_p} &= \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_{i_1}^* \dots x_{i_p}^* + \xi_{i_1 \dots i_p} \quad \forall i_1 < \dots < i_p, \end{aligned} \tag{1}$$

where ξ_{ij} and ξ_{i_1, \dots, i_p} are iid Gaussian components of a symmetric random matrix and tensor of zero mean and variance Δ_2 and Δ_p , respectively; ξ_{ij} and ξ_{i_1, \dots, i_p} correspond to noises corrupting the signal of the teacher. In the limit $\Delta_2 \rightarrow 0$, and $\Delta_p \rightarrow 0$, the above model reduces to the canonical spiked Wigner model [29], and spiked tensor model [17], respectively. The goal of the student is to infer the vector x^* from the knowledge of the matrix Y , of the tensor T , of the values Δ_2 and Δ_p , and the knowledge of the spherical prior. The scaling with

N as specified in Eq. (1) is chosen in such a way that the information-theoretically best achievable error varies between perfectly reconstructed spike x^* and random guess from the flat measure on \mathcal{S}_{N-1} . Here, and in the rest of the paper we denote $x \in \mathcal{S}_{N-1}$ the N -dimensional vector, and x_i with $i = 1, \dots, N$ its components.

Analogous matrix-tensor models, where next to a order- p tensor one observes a matrix created from the same spike are studied e.g. in [22] in the context of topic modeling, or in [17]. From the optimization-theory point of view, this model is highly non-trivial being high-dimensional and non-convex. For the purpose of the present paper this model is chosen because it involves three key ingredients: (a) It is in the class of models for which the Langevin algorithm can be analyzed exactly in large N limit. (b) The different phase transitions, both algorithmic and information theoretic, discussed hereafter, all happen at $\Delta_2 = \mathcal{O}(1)$, $\Delta_p = \mathcal{O}(1)$. (c) The AMP algorithm is in this model conjectured to be optimal among polynomial algorithms. It is this second and third ingredient that are not present in the pure spiked tensor model [17], making it unsuitable for our present study. We note that the Langevin algorithm was recently analyzed for the pure spiked tensor model in [20] in a regime where the noise variance is very small $\Delta \sim N^{-p/2}$, but we also note that in that model algorithms such as tensor unfolding and semidefinite programming work better, roughly up to $\Delta \sim N^{-p/4}$ [17, 18],

3 Bayes-optimal estimation and message-passing

In this section we present the performance of the Bayes-optimal estimator and of the approximate message passing algorithm. This theory is based on a straightforward adaptation of analogous results known for the pure spiked matrix model [29, 30, 9] and for the pure spiked tensor model [17, 10].

The Bayes-optimal estimator \hat{x} is defined as the one that among all estimators minimizes the mean-squared error (MSE) with the spike x^* . Starting from the posterior probability distribution

$$P(x|Y, T) = \frac{1}{Z(Y, T)} \left[\prod_{i=1}^N e^{-x_i^2/2} \right] \prod_{i < j} e^{-\frac{1}{2\Delta_2} \left(Y_{ij} - \frac{x_i x_j}{\sqrt{N}} \right)^2} \prod_{i_1 < \dots < i_p} e^{-\frac{1}{2\Delta_p} \left(T_{i_1 \dots i_p} - \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_{i_1} \dots x_{i_p} \right)^2}, \quad (2)$$

the Bayes-optimal estimator reads

$$\hat{x}_i = \mathbb{E}_{P(x|Y, T)}(x_i). \quad (3)$$

To simplify notation, and to make contact with the energy landscape and the statistical physics notations, it is convenient to introduce the energy cost function, or Hamiltonian, as

$$\mathcal{H}(x) = \mathcal{H}_2 + \mathcal{H}_p = -\frac{1}{\Delta_2 \sqrt{N}} \sum_{i < j} Y_{ij} x_i x_j - \frac{\sqrt{(p-1)!}}{\Delta_p N^{(p-1)/2}} \sum_{i_1 < \dots < i_p} T_{i_1 \dots i_p} x_{i_1} \dots x_{i_p} \quad (4)$$

so that keeping in mind that for $N \rightarrow \infty$ the spherical constraint is satisfied $|x|^2 = N$, the posterior is written as $P(x|Y, T) = \exp[-\mathcal{H}(x)] / \tilde{Z}(Y, T)$, where \tilde{Z} is the normalizing partition function.

With the use of the replica theory and its recent proofs from [9, 31, 10] one can establish rigorously that the mean squared error achieved by the Bayes-optimal estimator is given as $\text{MMSE} = 1 - m^*$ where $m^* \in \mathbb{R}$ is the global maximizer of the so-called free entropy of the problem

$$\Phi_{\text{RS}}(m) = \frac{1}{2} \log(1 - m) + \frac{m}{2} + \frac{1 + m^2}{4\Delta_2} + \frac{1 + m^p}{2p\Delta_p}. \quad (5)$$

This expression is derived, and proven, in the appendix Sec. B.1.

We now turn to the approximate message-passing (AMP) [17, 10], that is the best known so far for this problem. AMP is an iterative algorithm inspired from the work of Thouless-Anderson and Palmer in statistical physics [32]. We explicit its form in the Sec. B.1. Most remarkably performance of AMP can be evaluated by tracking its evolution with the iteration time and it is given in terms of the (possibly local) maximum of the above free entropy that is reached as a fixed point of the following iterative process

$$m^{t+1} = 1 - \frac{1}{1 + m^t / \Delta_2 + (m^t)^{p-1} / \Delta_p} \quad (6)$$

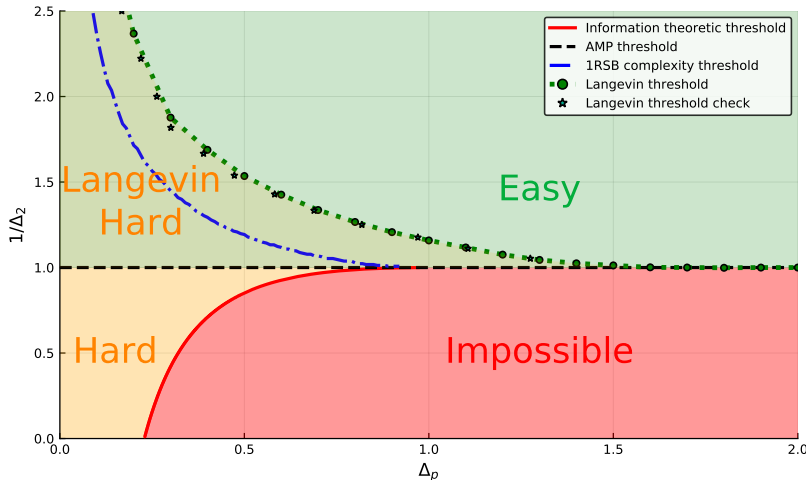


Figure 1: Phase diagram of the spiked 2 + 3-spin model (matrix plus order 3 tensor are observed). In the easy (green) region AMP achieves the optimal error smaller than random pick from the prior. In the impossible region (red) the optimal error is as bad as random pick from the prior, and AMP achieves it as well. In the hard region (orange) the optimal error is low, but AMP does not find an estimator better than random pick from the prior. In the case of Langevin algorithm the performance is strictly worse than that for AMP in the sense that the hard region increases up to line depicted in green dots. The blue dashed-dotted line delimits the region of existence of stable 1RSB states.

with initial condition $m^{t=0} = \epsilon$ with $0 < \epsilon \ll 1$. Eq. (6) is called the *State Evolution* of AMP and its validity is proven for closely related models in [33]. We denote the corresponding fixed point m_{AMP} and the corresponding estimation error $\text{MSE}_{\text{AMP}} = 1 - m_{\text{AMP}}$.

The phase diagram presented in Fig. 1 summarizes this theory for the spiked 2 + 3-spin model. It is deduced by investigating the local maxima of the scalar function (5). Notably we observe that the phase diagram in terms of Δ_2 and Δ_p splits into three phases

- **Easy** in green for $\Delta_2 < 1$ and any Δ_p : The fixed point of the state evolution (6) is the global maximizer of the free entropy (5), and $m^* = m_{\text{AMP}} > 0$.
- **Hard** in orange for $\Delta_2 > 1$ and low $\Delta_p < \Delta_p^{\text{IT}}(\Delta_2)$: The fixed point of the state evolution (6) is not the global maximizer of the free entropy (5), and $m^* > m_{\text{AMP}} = 0$.
- **Impossible** in red for $\Delta_2 > 1$ and high $\Delta_p > \Delta_p^{\text{IT}}(\Delta_2)$: The fixed point of the state evolution (6) is the global maximizer of the free entropy (5), and $m^* = m_{\text{AMP}} = 0$.

For the 2 + p -spin model with $p > 3$ the phase diagram is slightly richer and is presented in the appendix.

4 Langevin Algorithm and its Analysis

We now turn to the core of the paper and the analysis of the Langevin algorithm. In statistics, the most commonly used way to compute the Bayes-optimal estimator (3) is to attempt to sample the posterior distribution (2) and use several independent samples to compute the expectation in (3). In order to do that one needs to set up a stochastic dynamics on x that has a stationary measure at long times given by the posterior measure (2). The Langevin algorithm is one of the possibilities (others include notably Monte Carlo Markov chain). The common bottleneck is that the time needed to achieve stationarity can be in general exponential in the system size. In which case the algorithm is practically useless. However, this is not always the case and there are regions in parameter space where one can expect that the relaxation to the posterior measure happens on finite

timescales. Therefore it is crucial to understand where this happens and what are the associated relaxation timescales.

The Langevin algorithm on the hypersphere with Hamiltonian given by Eq. (4) reads

$$\dot{x}_i(t) = -\mu(t)x_i(t) - \frac{\partial \mathcal{H}}{\partial x_i} + \eta_i(t), \quad (7)$$

where $\eta_i(t)$ is a zero mean noise term, with $\langle \eta_i(t)\eta_j(t') \rangle = 2\delta_{ij}\delta(t-t')$ where the average $\langle \cdot \rangle$ is with respect to the realizations of the noise. The Lagrange multiplier $\mu(t)$ is chosen in such a way that the dynamics remains on the hypersphere. In the large N -limit one finds $\mu(t) = 1 - 2\mathcal{H}_2(t) - p\mathcal{H}_p(t)$ where the $\mathcal{H}_2(t)$ is the 1st term from (4) evaluated at $x(t)$, and $\mathcal{H}_p(t)$ is the value of the 2nd term from (4).

The presented spiked matrix-tensor model falls into the particular class of spherical $2+p$ -spin glasses [34, 35] for which the performance of the Langevin algorithm can be tracked exactly in the large- N limit via a set of integro-partial differential equations [26, 27], beforehand dubbed CHSCK. We call this generalised version of the CHSCK equations *Langevin State Evolution* (LSE) equations in analogy with the state evolution of AMP.

In order to write the LSE equations, we defined three dynamical correlation functions

$$C_N(t, t') \equiv \frac{1}{N} \sum_{i=1}^N x_i(t)x_i(t'), \quad (8)$$

$$\bar{C}_N(t) \equiv \frac{1}{N} \sum_{i=1}^N x_i(t)x_i^*, \quad (9)$$

$$R_N(t, t') \equiv \frac{1}{N} \sum_{i=1}^N \partial x_i(t)/\partial h_i(t')|_{h_i=0}, \quad (10)$$

where h_i is a pointwise external field applied at time t' to the Hamiltonian as $\mathcal{H} + \sum_i h_i x_i$. We note that the correlation functions defined above depend on the realization of the thermal history (i.e. of the noise $\eta(t)$) and on the disorder (here the matrix Y and tensor T). However, in the large- N limit they all concentrate around their averages. We thus define $C(t, t') = \lim_{N \rightarrow \infty} \mathbb{E}_{Y, T} \langle C_N(t, t') \rangle_\eta$ and analogously for $\bar{C}(t)$ and $R(t, t')$. Standard field theoretical methods [36] or dynamical cavity method arguments [37] can then be used to obtain a closed set of integro-differential equations for the averaged dynamical correlation functions, describing the average *global* evolution of the system under the Langevin algorithm. The resulting LSE equations are (see the appendix Sec. C for a complete derivation)

$$\begin{aligned} \frac{\partial}{\partial t} C(t, t') &= 2R(t', t) - \mu(t)C(t, t') + Q'(\bar{C}(t))\bar{C}(t') + \int_0^t dt'' R(t, t'')Q''(C(t, t''))C(t', t'') + \\ &\quad \int_0^{t'} dt'' R(t', t'')Q'(C(t, t'')), \\ \frac{\partial}{\partial t} R(t, t') &= \delta(t-t') - \mu(t)R(t, t') + \int_{t'}^t dt'' R(t, t'')Q''(C(t, t''))R(t'', t'), \\ \frac{\partial}{\partial t} \bar{C}(t) &= -\mu(t)\bar{C}(t) + Q'(\bar{C}(t)) + \int_0^t dt'' R(t, t'')\bar{C}(t'')Q(C(t, t'')), \end{aligned} \quad (11)$$

where we have defined $Q(x) = x^2/(2\Delta_2) + x^p/(p\Delta_p)$. The Lagrange multiplier, $\mu(t)$, is fixed by the spherical constraint, through the condition $C(t, t) = 1 \forall t$. Furthermore causality implies that $R(t, t') = 0$ if $t < t'$. Finally the Ito convention on the stochastic equation (7) gives $\forall t \lim_{t' \rightarrow t^-} R(t, t') = 1$.

5 Behavior of the Langevin algorithm

In order to assess the performances of the Langevin algorithm and compare it with AMP, we notice that the correlation function $\bar{C}(t)$ is directly related to accuracy of the algorithm. We solve the differential equations (11) numerically along the lines of [38, 39], for a detailed procedure see the appendix Sec. D, codes available at [40]. In Fig. 2 we plot the correlation with the spike $\bar{C}(t)$ as a function of the running time t for $p = 3$, fixed

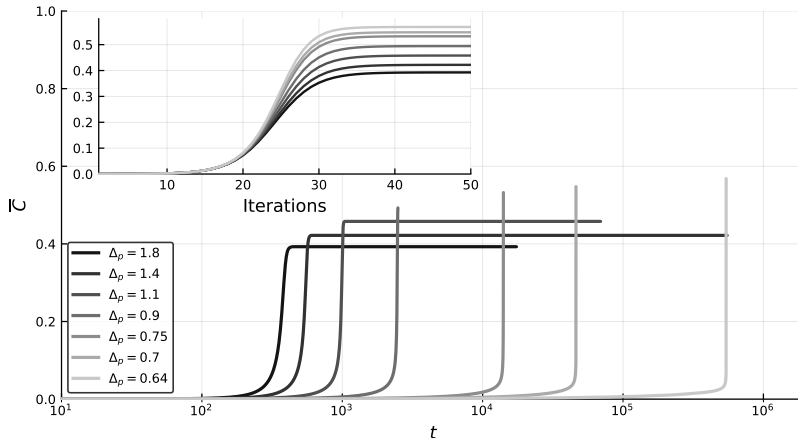


Figure 2: Evolution of the correlation with the signal $\overline{C}(t)$ in the Langevin algorithm at fixed noise on the matrix ($\Delta_2 = 0.7$) and different noises on the tensor (Δ_p). As we approach the transition, estimated to $\Delta_p^* \simeq 0.58$, the time required to jump to the solution diverges. Inset: the behavior of $\overline{C}(t)$ as a function of the iteration time for the AMP algorithm for the same values of Δ_p .

$\Delta_2 = 0.7$ and several values of Δ_p . In the inset of the plot we compare it to the same quantity obtained from the state evolution of the AMP algorithm. For the Langevin algorithm in Fig. 2 we see a pattern that is striking. One would expect that as the noise Δ_p decreases the inference problem is getting easier, the correlation with the signal is larger and is reached sooner in the iteration. This is, after all, exactly what we observe for the AMP algorithm in the inset of Fig. 2. Also for the Langevin algorithm the plateau reached for large times t becomes higher (better accuracy) as the noise Δ_p is reduced. Furthermore the height of the plateau coincides with that reached by AMP, thus testifying the algorithm reached equilibrium. However, contrary to AMP, the relaxation time for the Langevin algorithm increases dramatically when diminishing Δ_p (notice the log scale on x-axes of Fig. 2, as compared to the linear scale of the inset). We will define τ as the time it takes for the correlation to reach a value $\overline{C}_{\text{plateau}}/2$. We then plot the value of this equilibration time in the insets of Fig. 3 as a function of the noise Δ_p or Δ_2 having fixed Δ_2 (upper panel) or Δ_p (lower panel) respectively. The data are clearly consistent with a divergence of τ at a certain finite value of Δ_p^* and Δ_2^* . We fit the data by a power law fit $\tau(\Delta) = \left| \frac{1}{\Delta} - \frac{1}{\Delta^*} \right|^{-\gamma}$ and obtain in the particular case of fixed $\Delta_2 = 0.7$ a fit with $\gamma = 1.75$ and $\Delta_p^* = 0.58$, whereas for fixed $\Delta_p = 1.0$ we obtain $\gamma = 1.42$ and $\Delta_2^* = 0.86$. However, we are not able to strictly prove that the divergence of the relaxation time truly occurs, but at least our results imply that for $\Delta_p < \Delta_p^*$ and $\Delta_2 > \Delta_2^*$ the Langevin algorithm (7) is not a practical solver for the spiked mixed matrix-tensor problem. We will call the region $\Delta_p < \Delta_p^*$ and $\Delta_2 < 1$ where the AMP algorithm works optimally without problems yet Langevin algorithm does not, the *Langevin-hard region*. Both Δ_p^* and Δ_2^* are then plotted in Fig. 1 with green points (circles and stars) and consistently delimit the Langevin-hard region that extends considerably into the region where the AMP algorithm works optimally in a small number of iterations. Our main conclusion is thus that the Langevin algorithm designed to sample the posterior measure works efficiently in a considerably smaller region of parameters than the AMP as quantified in Fig. 1.

Fig. 4 presents another way to depict the observed data, the correlation $\overline{C}(t)$ reached after time t is plotted as a function of the tensor noise variance Δ_p . The results of AMP are depicted with dotted lines and, as one would expect, decrease monotonically as the noise Δ_p increases. The equilibrium value (black dashed) is reached within few dozens of iterations. On the contrary, the correlation reached by the Langevin algorithm after time t is non-monotonic and close to zero for small values of noise Δ_p signalling again a diverging relaxation time when Δ_p is decreased.

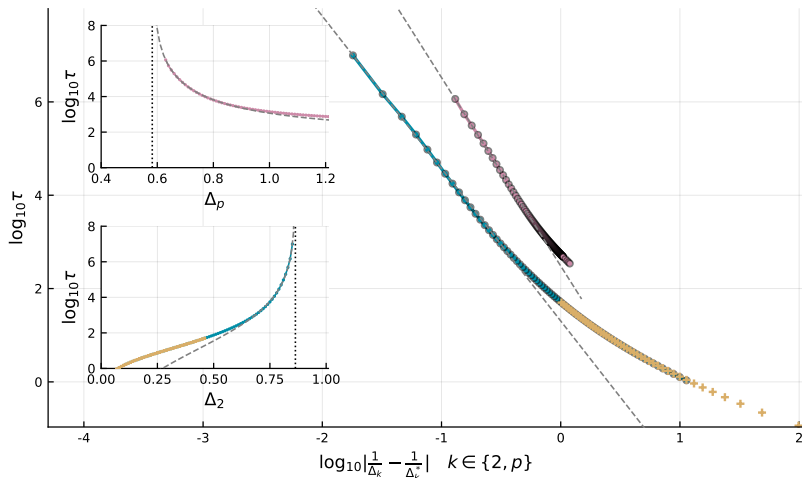


Figure 3: Fit using a power law of the relaxation times of the Langevin algorithm, the point of divergence marks the limit of the Langevin easy region. These fits have been performed both, for fixed $\Delta_2 = 0.7$ (blue circles and yellow crosses) and for fixed $\Delta_p = 1$ (pink circles). The circles are obtained with numerical solution of LSE that uses the dynamical grid while crosses are obtained using a fixed-grid (details in Sec. D.1).

6 Glassy nature of the Langevin-hard phase

We notice that for $\Delta_p \rightarrow 0$ and any finite Δ_2 , and after a suitable rescaling of time and temperature, the LSE equations coincide with the ones of the pure spiked-tensor ($p = 3$) at very low temperature and slightly perturbed by additional terms. The glassy nature of the landscape of such model [21, 20] qualitatively justifies why the Langevin algorithm remains trapped in one of the spurious minima, hence leading to a hard-Langevin phase for $\Delta_p \rightarrow 0$ and any finite Δ_2 (in particular $\Delta_2 < 1$).

In order to obtain a quantitative estimate of the extent of the glassy region, we repeated the analysis of [25] in the present model (details in the appendix Sec. E). In particular, we compute the entropy of metastable states (finite temperature extensions of spurious minima), *a.k.a.* the *complexity*, using the one-step replica symmetry breaking (1RSB) assumption for the structure of these states. We then analyze the stability of the 1RSB solution with respect to further levels of RSB and draw the region of parameters Δ_2, Δ_p for which states of positive complexity exist and are stable. This happens to be below the blue dashed-dotted line depicted in Fig. 1.

We observe that indeed in the regime where the 1RSB solution indicates the existence of an exponential number of spurious metastable states, the Langevin algorithm does not work. Moreover, the curve delimiting existence of the 1RSB stable states has the same trend as the boundary of the Langevin-hard regime. Yet the Langevin-hard phase extends to larger values of Δ_p , see Fig. 1. This quantitative discrepancy can be due to either the fact that we did not take into account effects of full-step replica symmetry breaking [41] or rather because the relation between the static replica calculation and the behavior of the Langevin algorithm is more complex than anticipated in the literature [41]. We let this point for investigation in future work.

7 Better performance by annealing the landscape

A generic strategy to avoid metastable states is by simulated annealing [42] where the thermal noise is given by $\langle \eta_i(t) \eta_j(t') \rangle = 2\delta_{ij} T(t) \delta(t - t')$ and the temperature $T(t)$ is time dependent and slowly decreased during the dynamics starting from very large values towards the target temperature, $T = 1$ in our case. We tried this algorithm and found that it does not succeed to recover the signal.

As discussed above, it is the tensor part of the Hamiltonian that induces glassiness. Therefore, we consider a protocol in which contribution of the tensor is weak at initial stages of the dynamics and increases gradually into

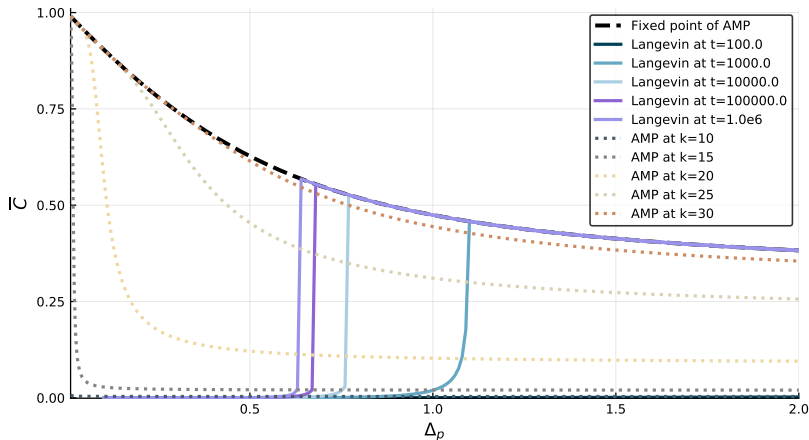


Figure 4: Correlation with the signal of AMP and Langevin at the k th iteration (at time t) for fixed $\Delta_2 = 0.7$.

its Bayes-optimal strength. Specifically, we focus on a time-dependent Hamiltonian, $\mathcal{H}_2(Y) + \frac{1}{T_p(t)}\mathcal{H}_p(T)$, and change dynamically $T_p(t)$ according to $T_p(t) = 1 + \frac{C}{\Delta_p} e^{-t/\tau_{\text{ann}}}$, C being a constant. The target value of $T_p(t \rightarrow \infty) = 1$ is the Bayes-optimal one and τ_{ann} is the characteristic timescale of the annealing procedure. In Fig. 5 we show that if τ_{ann} is sufficiently large, the Langevin hard phase is destroyed and this procedure approaches the performances as AMP. More specifically, before the $\mathcal{H}_p(T)$ component becomes relevant, i.e. for times that are $\tau_{\text{ann}} \lesssim \log \frac{C}{\Delta_p}$, the modified Langevin algorithm tends to the equilibrium solution of the pure spiked matrix model. The optimal AMP reconstruction is subsequently allowed when the tensor component gains its full weight, which occurs later for slower annealing. Conversely if τ_{ann} is too small, the Langevin algorithm does not have the time for escaping the glassy region of the tensor component before its own contribution to the dynamics becomes relevant and irremediably hampers the final reconstruction of the signal.

It is interesting to underline that the above finding is somewhat paradoxical from the point of view of Bayesian inference. In the present setting we know perfectly the model that generated the data and all its parameters, yet we see that for the Langevin algorithm it is computationally advantageous to mismatch the parameter T_p in order to reach faster convergence to equilibrium. This is particularly striking given that for AMP it has been proven in [7] that mismatching the parameters can never improve the performance.

8 Perspectives

In this work we have investigated the performances of the Langevin algorithm considered as a tool to sample the posterior measure in the spiked matrix-tensor model. We have shown that the Langevin algorithm fails to find the signal in part of the AMP-easy region. Our analysis is based on the Langevin State Evolution equations that describe the evolution of the algorithm in the large size limit.

In this work we managed to find the landscape-annealing protocol under which the Langevin algorithm is able to match the performance of AMP by relying on knowledge about generative model. It would be an interesting direction for future work to investigate whether the performance of the Langevin algorithm can be improved with some model-blind manner.

While we studied here the spiked matrix-tensor model, we expect that our findings are universal because they are due to the glassiness of the hard phase and therefore should apply to any local sampling dynamics, e.g. to Monte Carlo Markov chains. An interesting extension of this work would investigate algorithms closer to stochastic gradient descent and models closer to current neural network architectures.

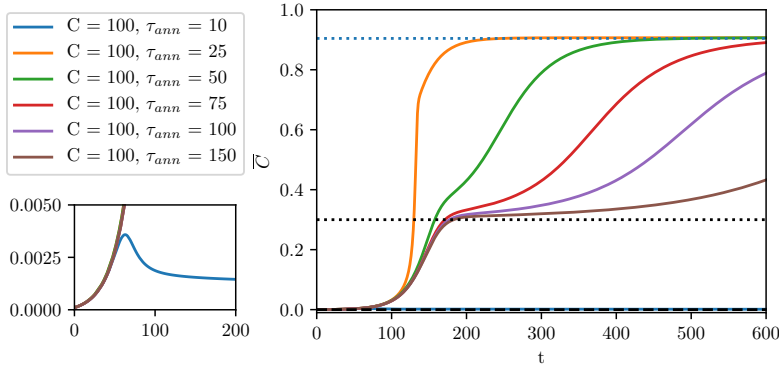


Figure 5: Behaviour of the correlation with the signal in Langevin algorithm where the tensor-related temperature T_p is annealed as $T_p(t) = 1 + \frac{C}{\Delta_p} e^{-t/\tau_{\text{ann}}}$ with $C = 100$ in the Langevin-hard regime with $\Delta_2 = 0.70$, $\Delta_p = 0.10$. We show the behavior for several rates τ (solid lines) and compare to the quenched Langevin algorithm (dashed line close to zero), and to the value reached by AMP for the full model (upper dotted line) and for the pure matrix model (lower dotted line). We see that unless the annealing is very fast, it reaches the AMP value. When the annealing is very slow it takes time to reach the AMP value.

Acknowledgments

We thank G. Folena and G. Ben Arous for precious discussions. We thank K. Miyazaki for sharing his code for the numerical integration of CHSCK equations. We acknowledge funding from the ERC under the European Union’s Horizon 2020 Research and Innovation Programme Grant Agreement 714608-SMiLe; from the French National Research Agency (ANR) grant PAIL; from ”Investissements d’Avenir” LabEx PALM (ANR-10-LABX-0039-PALM) (SaMURai and StatPhysDisSys); and from the Simons Foundation (#454935, Giulio Biroli).

References

- [1] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT’2010*, pages 177–186. Springer, 2010.
- [2] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 681–688, 2011.
- [3] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [4] Jean-Philippe Bouchaud, Leticia F Cugliandolo, Jorge Kurchan, and Marc Mézard. Out of equilibrium dynamics in spin-glasses and other glassy systems. *Spin glasses and random fields*, pages 161–223, 1998.
- [5] H. S. Seung, H. Sompolinsky, and N. Tishby. Statistical mechanics of learning from examples. *Phys. Rev. A*, 45:6056–6091, Apr 1992.
- [6] Lenka Zdeborová and Florent Krzakala. Statistical physics of inference: thresholds and algorithms. *Advances in Physics*, 65(5):453–552, 2016.
- [7] Yash Deshpande and Andrea Montanari. Finding hidden cliques of size $\sqrt{(N/e)}$ in nearly linear time. *Foundations of Computational Mathematics*, 15(4):1069–1128, 2015.
- [8] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Phase transitions, optimal errors and optimality of message-passing in generalized linear models. in *COLT’18, arXiv preprint arXiv:1708.03395*, 2017.

- [9] 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. In *Advances in Neural Information Processing Systems 29*, page 424–432. 2016.
- [10] Thibault Lesieur, Léo Miolane, Marc Lelarge, Florent Krzakala, and Lenka Zdeborová. Statistical and computational phase transitions in spiked tensor estimation. In *Information Theory (ISIT), 2017 IEEE International Symposium on*, pages 511–515. IEEE, 2017.
- [11] Benjamin Aubin, Antoine Maillard, Jean Barbier, Florent Krzakala, Nicolas Macris, and Lenka Zdeborová. The committee machine: Computational to statistical gaps in learning a two-layers neural network. In *Advances in Neural Information Processing Systems*, 2018.
- [12] 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, Nov 2009.
- [13] F. Krzakala, C. Moore, E. Mossel, J. Neeman, A. Sly, L. Zdeborová, and P. Zhang. Spectral redemption in clustering sparse networks. *Proceedings of the National Academy of Science*, 110:20935–20940, December 2013.
- [14] Samuel B Hopkins and David Steurer. Bayesian estimation from few samples: community detection and related problems. *arXiv preprint arXiv:1710.00264*, 2017.
- [15] Jinho Baik, Gérard Ben Arous, Sandrine Péché, et al. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *The Annals of Probability*, 33(5):1643–1697, 2005.
- [16] 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.
- [17] Emile Richard and Andrea Montanari. A statistical model for tensor PCA. In *Advances in Neural Information Processing Systems*, pages 2897–2905, 2014.
- [18] Samuel B Hopkins, Jonathan Shi, and David Steurer. Tensor principal component analysis via sum-of-square proofs. In *Conference on Learning Theory*, pages 956–1006, 2015.
- [19] Rong Ge and Tengyu Ma. On the optimization landscape of tensor decompositions. In *Advances in Neural Information Processing Systems*, pages 3653–3663, 2017.
- [20] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Algorithmic thresholds for tensor PCA. *arXiv preprint arXiv:1808.00921*, 2018.
- [21] Valentina Ros, Gerard Ben Arous, Giulio Biroli, and Chiara Cammarota. Complex energy landscapes in spiked-tensor and simple glassy models: ruggedness, arrangements of local minima and phase transitions. *arXiv preprint arXiv:1804.02686*, to appear in *PRX*, 2018.
- [22] Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *The Journal of Machine Learning Research*, 15(1):2773–2832, 2014.
- [23] Florent Krzakala and Lenka Zdeborová. Hiding quiet solutions in random constraint satisfaction problems. *Physical review letters*, 102(23):238701, 2009.
- [24] Aurelien Decelle, Florent Krzakala, Cristopher Moore, and Lenka Zdeborová. Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications. *Physical Review E*, 84(6):066106, 2011.

- [25] Fabrizio Antenucci, Silvio Franz, Pierfrancesco Urbani, and Lenka Zdeborová. On the glassy nature of the hard phase in inference problems. *to appear in Phys. Rev. X*, *arXiv preprint arXiv:1805.05857*, 2018.
- [26] A Crisanti, H Horner, and H-J Sommers. The spherical p -spin interaction spin-glass model. *Zeitschrift für Physik B Condensed Matter*, 92(2):257–271, 1993.
- [27] Leticia F Cugliandolo and Jorge Kurchan. Analytical solution of the off-equilibrium dynamics of a long-range spin-glass model. *Physical Review Letters*, 71(1):173, 1993.
- [28] Gerard Ben Arous, Amir Dembo, and Alice Guionnet. Cugliandolo-Kurchan equations for dynamics of spin-glasses. *Probability theory and related fields*, 136(4):619–660, 2006.
- [29] Yash Deshpande and Andrea Montanari. Information-theoretically optimal sparse PCA. In *Information Theory (ISIT), 2014 IEEE International Symposium on*, pages 2197–2201. IEEE, 2014.
- [30] Thibault Lesieur, Florent Krzakala, and Lenka Zdeborová. Constrained low-rank matrix estimation: Phase transitions, approximate message passing and applications. *Journal of Statistical Mechanics: Theory and Experiment*, 2017(7):073403, 2017.
- [31] Marc Lelarge and Léo Miolane. Fundamental limits of symmetric low-rank matrix estimation. *Probability Theory and Related Fields*, pages 1–71, 2016.
- [32] David J Thouless, Philip W Anderson, and Robert G Palmer. Solution of solvable model of a spin glass'. *Philosophical Magazine*, 35(3):593–601, 1977.
- [33] Adel Javanmard and Andrea Montanari. State evolution for general approximate message passing algorithms, with applications to spatial coupling. *Information and Inference: A Journal of the IMA*, 2(2):115–144, 2013.
- [34] Andrea Crisanti and Luca Leuzzi. Spherical $2+ p$ spin-glass model: An exactly solvable model for glass to spin-glass transition. *Physical review letters*, 93(21):217203, 2004.
- [35] Andrea Crisanti and Luca Leuzzi. Spherical $2+ p$ spin-glass model: An analytically solvable model with a glass-to-glass transition. *Physical Review B*, 73(1):014412, 2006.
- [36] Paul Cecil Martin, ED Siggia, and HA Rose. Statistical dynamics of classical systems. *Physical Review A*, 8(1):423, 1973.
- [37] Marc Mézard, Giorgio Parisi, and Miguel-Angel Virasoro. *Spin glass theory and beyond*. World Scientific Publishing, 1987.
- [38] Bongsoo Kim and Arnulf Latz. The dynamics of the spherical p -spin model: From microscopic to asymptotic. *EPL (Europhysics Letters)*, 53(5):660, 2001.
- [39] Ludovic Berthier, Giulio Biroli, J-P Bouchaud, Walter Kob, Kunimasa Miyazaki, and DR Reichman. Spontaneous and induced dynamic fluctuations in glass formers. I. General results and dependence on ensemble and dynamics. *The Journal of chemical physics*, 126(18):184503, 2007.
- [40] Stefano Sarao Mannelli, Giulio Biroli, Chiara Cammarota, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Langevin state evolution integrators, 2018. Available at: https://github.com/sphinxteam/spiked_matrix_tensor.
- [41] Tommaso Rizzo. Replica-symmetry-breaking transitions and off-equilibrium dynamics. *Physical Review E*, 88(3):032135, 2013.

- [42] Scott Kirkpatrick, C Daniel Gelatt, and Mario P Vecchi. Optimization by simulated annealing. *science*, 220(4598):671–680, 1983.
- [43] Andrea Crisanti and Luca Leuzzi. Exactly solvable spin–glass models with ferromagnetic couplings: The spherical multi-p-spin model in a self-induced field. *Nuclear Physics B*, 870(1):176–204, 2013.
- [44] Marc Mézard and Andrea Montanari. *Information, physics, and computation*. Oxford University Press, 2009.
- [45] Stéphane Boucheron, Gábor Lugosi, and Olivier Bousquet. Concentration inequalities. In *Advanced Lectures on Machine Learning*, pages 208–240. Springer, 2004.
- [46] Satish Babu Korada and Nicolas Macris. Exact solution of the gauge symmetric p-spin glass model on a complete graph. *Journal of Statistical Physics*, 136(2):205–230, 2009.
- [47] Florent Krzakala, Jiaming Xu, and Lenka Zdeborová. Mutual information in rank-one matrix estimation. In *2016 IEEE Information Theory Workshop (ITW)*, pages 71–75, Sept 2016.
- [48] Michael Aizenman, Robert Sims, and Shannon L Starr. Extended variational principle for the sherrington-kirkpatrick spin-glass model. *Physical Review B*, 68(21):214403, 2003.
- [49] Jean Barbier and Nicolas Macris. The adaptive interpolation method: A simple scheme to prove replica formulas in bayesian inference. *to appear in Probability Theory and Related Fields*, *arXiv preprint arXiv:1705.02780*, 2017.
- [50] 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, June 2018.
- [51] Jean-Christophe Mourrat. Hamilton-Jacobi equations for mean-field disordered systems. *arXiv preprint arXiv:1811.01432*, 2018.
- [52] Dongning Guo, S. Shamai, and S. Verdú. Mutual information and minimum mean-square error in gaussian channels. *IEEE Transactions on Information Theory*, 51(4):1261–1282, April 2005.
- [53] Hans-Otto Georgii. *Gibbs measures and phase transitions*, volume 9. Walter de Gruyter, 2011.
- [54] Nicolas Macris. Griffith-Kelly-Sherman correlation inequalities: A useful tool in the theory of error correcting codes. *IEEE Transactions on Information Theory*, 53(2):664–683, Feb 2007.
- [55] Andrea Montanari. Estimating random variables from random sparse observations. *European Transactions on Telecommunications*, 19(4):385–403, 2008.
- [56] Amin Coja-Oghlan, Florent Krzakala, Will Perkins, and Lenka Zdeborova. Information-theoretic thresholds from the cavity method. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 146–157, 2017.
- [57] Hidetoshi Nishimori. *Statistical physics of spin glasses and information processing: an introduction*, volume 111. Clarendon Press, 2001.
- [58] Francesco Guerra and Fabio Lucio Toninelli. The thermodynamic limit in mean field spin glass models. *Communications in Mathematical Physics*, 230(1):71–79, 2002.
- [59] Federico Ricci-Tersenghi, Guilhem Semerjian, and Lenka Zdeborová. Typology of phase transitions in bayesian inference problems. preprint:[arXiv:1806.11013](https://arxiv.org/abs/1806.11013).

- [60] Leticia F Cugliandolo. Course 7: Dynamics of glassy systems. In *Slow Relaxations and nonequilibrium dynamics in condensed matter*, pages 367–521. Springer, 2003.
- [61] Tommaso Castellani and Andrea Cavagna. Spin glass theory for pedestrians. *Journal of Statistical Mechanics: Theory and Experiment*, 2005:P05012, 2005.
- [62] Elisabeth Agoritsas, Giulio Biroli, Pierfrancesco Urbani, and Francesco Zamponi. Out-of-equilibrium dynamical mean-field equations for the perceptron model. *Journal of Physics A: Mathematical and Theoretical*, 51(8):085002, 2018.
- [63] Ludovic Berthier, Giulio Biroli, Jean-Philippe Bouchaud, Walter Kob, Kunimasa Miyazaki, and David R Reichman. Spontaneous and induced dynamic correlations in glass formers. II. Model calculations and comparison to numerical simulations. *The Journal of chemical physics*, 126(18):184504, 2007.
- [64] Rémi Monasson. Structural glass transition and the entropy of the metastable states. *Physical review letters*, 75(15):2847, 1995.
- [65] Francesco Zamponi. Mean field theory of spin glasses. *arXiv preprint [arXiv:1008.4844](https://arxiv.org/abs/1008.4844)*, 2010.
- [66] Andrea Crisanti and H-J Sommers. The spherical p-spin interaction spin glass model: the statics. *Zeitschrift für Physik B Condensed Matter*, 87(3):341–354, 1992.
- [67] LF Cugliandolo and J Kurchan. Weak ergodicity breaking in mean-field spin-glass models. *Philosophical Magazine B*, 71(4):501–514, 1995.
- [68] LF Cugliandolo and J Kurchan. On the out-of-equilibrium relaxation of the Sherrington-Kirkpatrick model. *Journal of Physics A: Mathematical and General*, 27(17):5749, 1994.
- [69] Leticia F Cugliandolo and Jorge Kurchan. Aging and effective temperatures in the low temperature mode-coupling equations. *Progress of Theoretical Physics Supplement*, 126:407–414, 1997.

Appendix

A Definition of the spiked matrix-tensor model

We consider a teacher-student setting in which the teacher constructs a matrix and a tensor from a randomly sampled signal and the student is asked to recover the signal from the observation of the matrix and tensor provided by the teacher [6].

The signal, x^* is an N -dimensional vector whose entries are real i.i.d. random variables sampled from the normal distribution (i.e. the prior is $P_X \sim \mathcal{N}(0, 1)$). The teacher generates from the signal a symmetric matrix and a symmetric tensor of order p . Those two objects are then transmitted through two noisy channels with variances Δ_2 and Δ_p , so that at the end one has two noisy observations given by

$$Y_{ij} = \frac{x_i^* x_j^*}{\sqrt{N}} + \xi_{ij}, \quad (12)$$

$$T_{i_1, \dots, i_p} = \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_{i_1}^* \dots x_{i_p}^* + \xi_{i_1, \dots, i_p}, \quad (13)$$

where, for $i < j$ and $i_1 < \dots < i_p$, ξ_{ij} and ξ_{i_1, \dots, i_p} are i.i.d. random variables distributed according to $\xi_{ij} \sim \mathcal{N}(0, \Delta_2)$ and $\xi_{i_1, \dots, i_p} \sim \mathcal{N}(0, \Delta_p)$. The ξ_{ij} and ξ_{i_1, \dots, i_p} are symmetric random matrix and tensor, respectively. Given Y_{ij} and T_{i_1, \dots, i_p} the inference task is to reconstruct the signal x^* .

In order to solve this problem we consider the Bayesian approach. This starts from the assumption that both the matrix and tensor have been produced from a process of the same kind of the one described by Eq. (12-13). Furthermore we assume to know the statistical properties of the channel, namely the two variances Δ_2 and Δ_p , and the prior on x . Given this, the posterior probability distribution over the signal is obtained through the Bayes formula

$$P(X|Y, T) = \frac{P(Y, T|X)P(X)}{P(Y, T)}, \quad (14)$$

where

$$\begin{aligned} P(Y, T|X) &= \prod_{i < j} P_Y \left(Y_{ij} \left| \frac{x_i x_j}{\sqrt{N}} \right. \right) \prod_{i_1 < \dots < i_p} P_T \left(T_{i_1 \dots i_p} \left| \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_{i_1} \dots x_{i_p} \right. \right) = \\ &\propto \prod_{i < j} e^{-\frac{1}{2\Delta_2} \left(Y_{ij} - \frac{x_i x_j}{\sqrt{N}} \right)^2} \prod_{i_1 < \dots < i_p} e^{-\frac{1}{2\Delta_p} \left(T_{i_1 \dots i_p} - \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_{i_1} \dots x_{i_p} \right)^2}. \end{aligned} \quad (15)$$

Therefore we have

$$P(X|Y, T) = \frac{1}{Z(Y, T)} \prod_i e^{-\frac{1}{2} x_i^2} \prod_{i < j} e^{-\frac{1}{2\Delta_2} \left(Y_{ij} - \frac{x_i x_j}{\sqrt{N}} \right)^2} \prod_{i_1 < \dots < i_p} e^{-\frac{1}{2\Delta_p} \left(T_{i_1 \dots i_p} - \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_{i_1} \dots x_{i_p} \right)^2}, \quad (16)$$

where $Z(Y, T)$ is a normalization constant.

Plugging Eqs. (12-13) into Eq. (16) allows to rewrite the posterior measure in the form of a *Boltzmann distribution* of the mixed $2 + p$ -spin Hamiltonian [34, 35, 43]

$$\begin{aligned} \mathcal{H} &= -\frac{1}{\Delta_2 \sqrt{N}} \sum_{i < j} \xi_{ij} x_i x_j - \frac{\sqrt{(p-1)!}}{\Delta_p N^{\frac{p-1}{2}}} \sum_{i_1 < \dots < i_p} \xi_{i_1 \dots i_p} x_{i_1} \dots x_{i_p} - \frac{N}{2\Delta_2} \left(\frac{1}{N} \sum_i x_i x_i^* \right)^2 + \\ &- \frac{N}{p\Delta_p} \left(\frac{1}{N} \sum_i x_i x_i^* \right)^p - \frac{1}{2} \sum_{i=1}^N x_i^2 + \text{const.} \end{aligned} \quad (17)$$

so that

$$P(X|Y, T) = \frac{1}{\tilde{Z}(Y, T)} e^{-\mathcal{H}}. \quad (18)$$

In the following we will refer to $\tilde{Z}(Y, T)$ as the *partition function*. We note here that in the large N limit, using a Gaussian prior on the variables x_i is equivalent to consider a flat measure over the N -dimensional hypersphere $\sum_{i=1}^N x_i^2 = N$. This choice will be used when we will describe the Langevin algorithm and in this case the last term in the Hamiltonian will become an irrelevant constant.

B Approximate Message Passing, state evolution and phase diagrams

Approximate Message Passing (AMP) is a powerful iterative algorithms to compute the local *magnetizations* $\langle x_i \rangle$ given the observed matrix and tensor. It is rooted in the cavity method of statistical physics of disordered systems [32, 37] and it has been recently developed in the context of statistical inference [12], where in the Bayes optimal case it has been conjectured to be optimal among all local iterative algorithms. Among the properties that make AMP extremely useful is the fact that its performances can be analyzed in the thermodynamic limit. Indeed in such limit, its dynamical evolution is described by the so called State Evolution (SE) equations [12]. In this section we derive the AMP equations and their SE description for the spiked matrix-tensor model and solve them to obtain the phase diagram of the model as a function of the variances Δ_2 and Δ_p of the two noisy channels.

B.1 Approximate Message Passing and Bethe free entropy

AMP can be obtained as a relaxed Gaussian closure of the Belief Propagation (BP) algorithm. The derivation that we present follows the same lines of [10, 30]. The posterior probability can be represented as a factor graph where all the variables are represented by circles and are linked to squares representing the interactions [44].

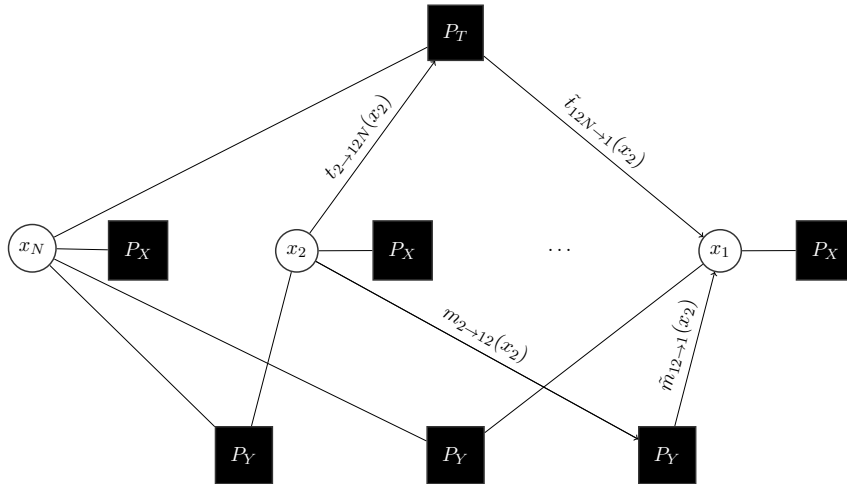


Figure 6: The factor graph representation of the posterior measure of the matrix-tensor factorization model. The variable nodes represented with white circles are the components of the signal while black squares are factor nodes that denote interactions between the variable nodes that appear in the interaction terms of the Boltzmann distribution in Eqs. (17-18). There are three types of factor nodes: P_X is the prior that depends on a single variable, P_Y that is the probability of observing a matrix element Y_{ij} given the values of the variables x_i and x_j , and finally P_T that is the probability of observing a tensor element T_{i_1, \dots, i_p} . The posterior, apart from the normalization factor, is simply given by the product of all the factor nodes.

This representation is very convenient to write down the BP equations. In the BP algorithm we iteratively update until convergence a set of variables, which are beliefs of the (cavity) magnetization of the nodes. The

intuitive underlying reasoning behind how BP works is the following. Given the current state of the variable nodes, take a factor node and exclude one node among its neighbors. The remaining neighbors through the factor node express a belief on the state of the excluded node. This belief is mathematically described by a probability distribution called *message*, $\tilde{m}_{ij \rightarrow i}^t(x_i)$ and $\tilde{t}_{ii_2 \dots i_p \rightarrow i}^t(x_i)$ depending on which factor node is selected. At the same time, another belief on the state of the excluded node is given by the rest of the network but the factor node previously taken into account, $m_{i \rightarrow ij}(x_i)$ and $t_{i \rightarrow ii_2 \dots i_p}(x_i)$ respectively. All these *messages* travel in the factor graph carrying partial information on the real magnetization of the single nodes, and they are iterated until convergence¹. The iterative scheme is described by the following equations

$$\tilde{m}_{ij \rightarrow i}^t(x_i) \propto \int dx_j m_{j \rightarrow ij}^t(x_j) P_Y \left(Y_{ij} \left| \frac{x_i x_j}{\sqrt{N}} \right. \right), \quad (19)$$

$$m_{i \rightarrow ij}^{t+1}(x_i) \propto P_X(x_i) \prod_{l \neq j} \tilde{m}_{il \rightarrow i}^t(x_i) \prod_{i_2 < \dots < i_p} \tilde{t}_{ii_2 \dots i_p \rightarrow i}^t(x_i), \quad (20)$$

$$\tilde{t}_{ii_2 \dots i_p \rightarrow i}^t(x_i) \propto \int \prod_{l=2 \dots p} \left(dx_l t_{il \rightarrow ii_2 \dots i_p}^t(x_l) \right) P_T \left(T_{ii_2 \dots i_p} \left| \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_i x_{i_2} \dots x_{i_p} \right. \right), \quad (21)$$

$$t_{i \rightarrow ii_2 \dots i_p}^{t+1}(x_i) \propto P_X(x_i) \prod_l \tilde{m}_{il \rightarrow i}^t(x_i) \prod_{k_2 < \dots < k_p \neq i_2 \dots i_p} \tilde{t}_{ik_2 \dots k_p \rightarrow i}^t(x_i) \quad (22)$$

and we have omitted the normalization constants that guarantee that the messages are probability distributions. When the messages have converged to a fixed point, the estimation of the local magnetizations can be obtained through the computation of the real marginal probability distribution of the variables given by

$$\mu_i(x_i) = \int \left[\prod_{j(\neq i)} dx_j \right] P(X|Y, T) = P_X(x_i) \prod_l \tilde{m}_{il \rightarrow i}^t(x_i) \prod_{i_2 < \dots < i_p} \tilde{t}_{ii_2 \dots i_p \rightarrow i}^t(x_i). \quad (23)$$

We note that the computational cost to produce an iteration of BP scales as $O(N^p)$. Furthermore Eqs. (19-22) are iterative equations for continuous functions and therefore are extremely hard to solve when dealing with continuous variables. The advantage of AMP is to reduce drastically the computational complexity of the algorithm by closing the equations on a Gaussian ansatz for the messages. This is justified in the present context since the factor graph is fully connected and therefore each iteration step of the algorithm involves sums of a large number of independent random variables that give rise to Gaussian distributions. Gaussian random variables are characterized by their mean and covariance that are readily obtained for $N \gg 1$ expanding the factor nodes for small $\omega_{ij} = x_i x_j / \sqrt{N}$ and $\omega_{i_1 \dots i_p} = \sqrt{(p-1)!} x_1 \dots x_p / N^{\frac{p-1}{2}}$.

Once the BP equations are relaxed on Gaussian messages, the final step to obtain the AMP algorithm is the so-called TAPyfication procedure [30, 32], which exploits the fact that the procedure of removing one node or one factor produces only a weak perturbation to the real marginals and therefore can be described in terms of the real marginals of the variable nodes themselves. By applying this scheme we obtain the AMP equations, which are described by a set of auxiliary variables $A^{(k)}$ and $B_i^{(k)}$ and by the mean $\langle x_i \rangle$ and variance $\sigma_i = \langle x_i^2 \rangle$ of the marginals of variable nodes. The AMP iterative equations are

$$B_i^{(2),t} = \frac{1}{\Delta_2 \sqrt{N}} \sum_k Y_{ki} \hat{x}_k^t - \frac{1}{\Delta_2} \left(\frac{1}{N} \sum_k \sigma_k^t \right) \hat{x}_i^{t-1}; \quad (24)$$

$$A^{(2),t} = \frac{1}{\Delta_2 N} \sum_k \left(\hat{x}_k^t \right)^2; \quad (25)$$

$$B_i^{(p),t} = \frac{\sqrt{(p-1)!}}{\Delta_p N^{(p-1)/2}} \sum_{k_2 \dots k_p} T_{ik_2 \dots k_p} \left(\hat{x}_{k_2}^t \dots \hat{x}_{k_p}^t \right) - \frac{p-1}{\Delta_p} \left[\left(\frac{1}{N} \sum_k \sigma_k^t \right) \left[\frac{1}{N} \sum_k \left(\hat{x}_k^t \right)^2 \right]^{p-2} \right] \hat{x}_i^{t-1}; \quad (26)$$

¹Note that the pointwise convergence of the algorithm depends on the situations.

$$A^{(p),t} = \frac{1}{\Delta_p} \left[\frac{1}{N} \sum_k (\hat{x}_k^t)^2 \right]^{p-1}; \quad (27)$$

$$\hat{x}_i^{t+1} = f(A^{(2)} + A^{(p)}, B_i^{(2)} + B_i^{(p)}); \quad (28)$$

$$\sigma_i^{t+1} = \frac{\partial}{\partial B} f(A, B) \Big|_{A=A^{(2)}+A^{(p)}, B=B_i^{(2)}+B_i^{(p)}}, \quad (29)$$

$$f(A, B) \equiv \int dx \frac{1}{\mathcal{Z}(A, B)} x P_X(x) e^{Bx - \frac{1}{2}Ax^2} = \frac{B}{1+A}. \quad (30)$$

It can be shown that these equations can be obtained as saddle point equations from the so called Bethe free entropy defined as $\Phi_{\text{Bethe}} = \log Z^{\text{Bethe}}(Y, T)/N$ where Z^{Bethe} is the Bethe approximation to the partition function which is defined as the normalization of the posterior measure. The expression of the Bethe free entropy per variable can be computed in a standard way (see [44]) and it is given by

$$\Phi_{\text{Bethe}} = \frac{1}{N} \left(\sum_i \log Z_i + \sum_{i \leq j} \log Z_{ij} + \sum_{i_1 \leq \dots \leq i_p} \log Z_{i_1 \dots i_p} - \sum_{i(j)} \log Z_{i,j} - \sum_{i(i_2 \dots i_p)} \log Z_{i(i_2 \dots i_p)} \right), \quad (31)$$

where

$$\begin{aligned} Z_i &= \int dx_i P_X(x_i) \prod_j \tilde{m}_{ij \rightarrow i}(x_i) \prod_{(i_2 \dots i_p)} \tilde{t}_{ii_2 \dots i_p \rightarrow i}(x_i), \\ Z_{ij} &= \int \prod_{j(\neq i)} [dx_j m_{j \rightarrow ij}(x_j)] \prod_{i < j} e^{-\frac{1}{2\Delta_2} (Y_{ij} - \frac{x_i x_j}{\sqrt{N}})^2}, \\ Z_{i_1 \dots i_p} &= \int \prod_{l=1}^p [dx_{i_l} t_{i_l \rightarrow i_1 \dots i_p}(x_{i_l})] \prod_{i_1 < \dots < i_p} e^{-\frac{1}{2\Delta_p} \left(T_{i_1 \dots i_p} - \frac{\sqrt{(p-1)!}}{N^{(p-1)/2}} x_{i_1} \dots x_{i_p} \right)^2}, \\ Z_{i(j)} &= \int dx_i m_{i \rightarrow ij}(x) \tilde{m}_{ij \rightarrow i}(x_i), \\ Z_{i(i_2 \dots i_p)} &= \int dx_i t_{i \rightarrow ii_2 \dots i_p}(x) \tilde{t}_{ii_2 \dots i_p \rightarrow i}(x_i) \end{aligned}$$

are a set of normalization factors. Using the Gaussian approximation for the messages and employing the same TAPyfication procedure used to get the AMP equations we obtain the Bethe free entropy density as

$$\begin{aligned} \Phi_{\text{Bethe}} &= \frac{1}{N} \sum_i \log \mathcal{Z}(A^{(p)} + A^{(2)}, B_i^{(p)} + B_i^{(2)}) + \frac{p-1}{p} \frac{1}{N} \sum_i \left[-B_i^{(p)} \hat{x}_i + A_i^{(p)} \frac{\hat{x}_i^2 + \sigma_i}{2} \right] + \\ &+ \frac{p-1}{2p\Delta_p} \left(\frac{\sum_i \hat{x}_i^2}{N} \right)^{p-1} \left(\frac{\sum_i \sigma_i}{N} \right) + \frac{1}{2N} \sum_i \left[-B_i^{(2)} \hat{x}_i + A_i^{(2)} \frac{\hat{x}_i^2 + \sigma_i}{2} \right] + \frac{1}{4\Delta_2} \left(\frac{\sum_i \hat{x}_i^2}{N} \right) \left(\frac{\sum_i \sigma_i}{N} \right), \end{aligned} \quad (32)$$

where we used the variables defined in eqs. (24-27) for sake of compactness and $\mathcal{Z}(A, B)$ is defined as

$$\mathcal{Z}(A, B) = \int dx P_X(x) e^{Bx - \frac{Ax^2}{2}} = \frac{1}{\sqrt{A+1}} e^{\frac{B^2}{2(A+1)}}. \quad (33)$$

B.2 Averaged free entropy and its proof

Eq. (32) represents the Bethe free entropy for a single realization of the factor nodes in the large size limit. Here we wish to discuss the actual, exact, value of this free entropy, that is:

$$f_N(Y, T) = \frac{\log Z(Y, T)}{N}$$

This is a random variable, since it depends a priori on the planted signal and the noise in the tensor and matrices. However one expects that, since free entropy is an intensive quantity, we expect from the statistical physics intuition that it should be self averaging and concentrate around its mean value in the large N limit [37]. In fact, this is easily proven. First, since the spherical model has a rotational symmetry, one may assume the planted assignment could be any vector on the hyper-sphere, and we might as well suppose it is the uniform one $x_i^* = 1 \forall i$: the true source of fluctuation comes from the noise Y and T . These can be controlled by noticing that the free entropy is a Lipschitz function of the Gaussian random variable Y and T . Indeed:

$$\partial_{Y_{ij}} f_N(Y, T) = \frac{1}{\Delta_2 N \sqrt{N}} \langle x_i x_j \rangle$$

so that the free energy f_N is Lipschitz with respect to Y with constant

$$L = \frac{1}{\Delta_2 N \sqrt{N}} \sqrt{\sum_{i < j} \langle x_i x_j \rangle^2} \leq \frac{1}{\Delta_2 N \sqrt{N}} \sqrt{\frac{1}{2} \sum_{i, j} \langle x_i x_j \rangle^2} = \frac{1}{\Delta_2 N \sqrt{N}} \sqrt{\frac{1}{2} \sum_{i, j} \langle x_i \tilde{x}_i x_j \tilde{x}_j \rangle}$$

where \tilde{x} represent a copy (or replica) of the system. In this case

$$L \leq \frac{1}{\Delta_2 N \sqrt{N}} \sqrt{N^2 \left\langle \left(\frac{\sum_i x_i \tilde{x}_i}{N} \right)^2 \right\rangle} = \frac{\sqrt{\langle q^2 \rangle}}{\Delta_2 \sqrt{N}}$$

where q is the overlap between the two replica x and \tilde{x} , that is bounded by one on the sphere, so $L \leq \frac{1}{\Delta_2 \sqrt{N}}$. Therefore, by Gaussian concentration of Lipschitz functions (the Tsirelson-Ibragimov-Sudakov inequality [45]), we have for some constant K :

$$\Pr [|f_n - \mathbb{E}_Y f_n| \geq t] \leq 2e^{-Nt^2/K} \quad (34)$$

and it particular any fluctuation larger than $O(1/\sqrt{N})$ is (exponentially) rare. A similar computation shows that f_N also concentrates with respect to the tensor T . This shows that in the large size limit, we can consider the averaged free entropy:

$$\mathcal{F}_N \equiv \frac{1}{N} \mathbb{E} [\log Z_N]$$

With our (non-rigorous) statistical physics tools, this can be obtained by averaging Eq. (32) over the disorder, see for instance [30], and this yields an expression for the free energy called the replica symmetric (RS) formula:

$$\Phi_{\text{RS}} = \lim_{N \rightarrow \infty} \mathbb{E}_{Y, T} \frac{\log Z(Y, T)}{N}. \quad (35)$$

We now state precisely the form of Φ_{RS} and prove the validity of Eq. (35). The RS free entropy for any prior distribution P_X reads as

$$\begin{aligned} \Phi_{\text{RS}} &\equiv \max_m \tilde{\Phi}_{\text{RS}}(m) \quad \text{where} \\ \tilde{\Phi}_{\text{RS}}(m) &= \mathbb{E}_{W, x^*} \left[\log \left[\mathcal{Z} \left(\frac{m}{\Delta_2} + \frac{m^{p-1}}{\Delta_p}, \left(\frac{m}{\Delta_2} + \frac{m^{p-1}}{\Delta_p} \right) x^* + \sqrt{\frac{m}{\Delta_2} + \frac{m^{p-1}}{\Delta_p}} W \right) \right] \right] - \frac{1}{4\Delta_2} m^2 - \frac{p-1}{2p\Delta_p} m^p, \end{aligned} \quad (36)$$

where W is a Gaussian random variable of zero mean and unit variance and x^* is a random variables taken from the prior P_X . We remind that the function $\mathcal{Z}(A, B)$ is defined via Eq. (33).

For Gaussian prior P_X , which is the one of interest here, we obtain

$$\tilde{\Phi}_{\text{RS}}(m) = -\frac{1}{2} \log \left(\frac{m}{\Delta_2} + \frac{m^{p-1}}{\Delta_p} + 1 \right) + \frac{1}{2} \left(\frac{m}{\Delta_2} + \frac{m^{p-1}}{\Delta_p} \right) - \frac{1}{4\Delta_2} m^2 - \frac{p-1}{2p\Delta_p} m^p. \quad (37)$$

The expression given in the main text is slightly different but can be obtained as follow. First notice that the extremization condition for $\tilde{\Phi}_{\text{RS}}(m)$ reads

$$m = 1 - \frac{1}{1 + \frac{m}{\Delta_2} + \frac{m^{p-1}}{\Delta_p}} \quad (38)$$

and by plugging this expression in Eq. (37) we recover the more compact expression $\Phi_{\text{RS}}(m)$ showed in the main text:

$$\Phi_{\text{RS}}(m) = \frac{1}{2} \log(1 - m) + \frac{m}{2} + \frac{m^2}{4\Delta_2} + \frac{m^p}{2p\Delta_p}. \quad (39)$$

The two expressions $\Phi_{\text{RS}}(m)$ and $\tilde{\Phi}_{\text{RS}}(m)$ are thus equal for each value of m that satisfy Eq. (38). The parameter m can be interpreted as the average correlation between the true and the estimated signal

$$m = \frac{1}{N} \sum_{i=1}^N x_i^* \hat{x}_i. \quad (40)$$

The average minimal mean squared error (MMSE) can be obtained from the maximizer m of the average Bethe free entropy as

$$\text{MMSE} \equiv \frac{1}{N} \sum_{i=1}^N \overline{(x_i^* - \hat{x}_i)^2} = 1 - m^*, \quad \text{where } m^* = \text{argmax}_m \tilde{\Phi}_{\text{RS}}(m). \quad (41)$$

where the overbar stands for the average over the signal x^* and the noise of the two Gaussian channels.

The validity of Eq. (36) can be proven rigorously for every prior having a bounded second moment. The proof we shall present is a straightforward generalization of the one presented in [10] for the pure tensor case, and in [31] for the matrix case, and it is based on two main ingredients. The first one is the Guerra interpolation method applied on the Nishimori line [46, 47, 31], in which we construct an interpolating Hamiltonian that depends on a parameter $t \in [0; 1]$ that is used to move from the original Hamiltonian of Eq. (17), to the one corresponding to a scalar denoising problem whose free entropy is given by the first term in Eq. (36). The second ingredient is the Aizenman-Sims-Starr method [48] which is the mathematical version of the cavity method (note that other techniques could also be employed to obtain the same results, see [9, 49, 50, 51]). The theorem we want to prove is:

Theorem 1 (Replica-Symmetric formula for the free energy). *Let P_X be a probability distribution over \mathbb{R} , with finite second moment Σ_X . Then, for all $\Delta_2 > 0$ and $\Delta_p > 0$*

$$\mathcal{F}_N \equiv \frac{1}{N} \mathbb{E} [\log Z_N] \xrightarrow{N \rightarrow \infty} \sup_{m \geq 0} \tilde{\Phi}_{\text{RS}}(m) \equiv \Phi_{\text{RS}}(\Delta_2, \Delta_p). \quad (42)$$

For almost every $\Delta_2 > 0$ and $\Delta_p > 0$, $\tilde{\Phi}_{\text{RS}}$ admits a unique maximizer m over $\mathbb{R}_+ \times \mathbb{R}_+$ and

$$\text{T-MMSE}_N \xrightarrow{N \rightarrow \infty} \Sigma_X^p - (m^*)^p,$$

$$\text{M-MMSE}_N \xrightarrow{N \rightarrow \infty} \Sigma_X^2 - (m^*)^2.$$

Here, we have defined the tensor-MMSE T-MMSE_N by the error in reconstructing the tensor:

$$\text{T-MMSE}_N(\Delta_2, \Delta_p) = \inf_{\hat{\theta}} \left\{ \frac{p!}{N^p} \sum_{i_1 < \dots < i_p} \left(x_{i_1}^0 \dots x_{i_p}^0 - \hat{\theta}(Y)_{i_1 \dots i_p} \right)^2 \right\},$$

and the matrix-MMSE M-MMSE_N by the error in reconstructing the matrix:

$$\text{M-MMSE}_N(\Delta_2, \Delta_p) = \inf_{\hat{\theta}} \left\{ \frac{2}{N^2} \sum_{i < j} \left(x_i^0 x_j^0 - \hat{\theta}(Y)_{i,j} \right)^2 \right\},$$

where in both cases the infimum is taken over all measurable functions $\hat{\theta}$ of the observations Y .

The result concerning the MMSE is a simple application of the I-MMSE theorem [52], that relates the derivative of the free energy with respect to the noise variances and the MMSE. The details of the arguments are the same than in the matrix ($p = 2$) case ([31], corollary 17) and the tensor one ([10], theorem 2). Indeed, as discussed in [31, 10], these M-MMSE and T-MMSE results implies the vector MMSE result of Eq. (41) when p is odd, and thus in particular for the $p = 3$ case discussed in the main text.

Sketch of proof In this section we give a detailed sketch of the proof theorem 1. Following the techniques used in many recent works [46, 47, 9, 31, 10, 49, 50, 8], we shall make few technical remarks:

- We will consider only priors with bounded support, $\text{supp}(P_X) = S \subset [-K; K]$. This allows to switch integrals and derivatives without worries. This condition can then be relaxed to unbounded distributions with bounded second moment using the same techniques as the ones that we are going to present, and the proof is therefore valid in this case. This is detailed for instance in [31] sec. 6.2.2.
- Another key ingredient is the introduction of a small perturbation in the model that takes the form of a small amount of side information. This kind of techniques are frequently used in statistical physics, where a small “magnetic field” forces the Gibbs measure to be in a single pure state [53]. It has also been used in the context of coding theory [54] for the same reason. In the context of Bayesian inference, we follow the generic scheme proposed by Montanari in [55] (see also [56]) and add a small additional source of information that allows the system to be in a single pure state so that the overlap concentrates on a single value. This source depends on Bernoulli random variables $L_i \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(\epsilon)$, $i \in [N]$; if $L_i = 1$, the channel, call it A , transmits the correct information. We can then consider the posterior of this new problem, $P(X|A, Y, T)$, and focus on the associated free energy density $F_{N,\epsilon}$ defined as the expected value of the average of the logarithm of normalization constant divided by the number of spins. Then we can immediately prove that for all $N \geq 1$ and $\epsilon, \epsilon' \in [0; 1]$ it follows: $|F_{N,\epsilon} - F_{N,\epsilon'}| \leq \left(\frac{K^{2p}}{\Delta_p} + \frac{K^4}{\Delta_2}\right) |\epsilon - \epsilon'|$. This allows (see for instance [10]) to obtain the concentration of the posterior distribution around the replica parameter ($q = \frac{1}{N} \langle x^{(1)} \cdot x^{(2)} \rangle$)

$$\mathbb{E} \left\langle \left(\frac{x^{(1)} \cdot x^{(2)}}{N} - q \right)^2 \right\rangle \xrightarrow{N \rightarrow \infty} 0 ; \quad (43)$$

$$\mathbb{E} \left\langle \left(\frac{x^* \cdot x}{N} - q \right)^2 \right\rangle \xrightarrow{N \rightarrow \infty} 0 , \quad (44)$$

where $x, x^{(1)}, x^{(2)}$ are sampled from the posterior distribution and the averages $\langle \cdot \rangle$ and $\mathbb{E}[\cdot]$ are respectively the average over the posterior measure and the remaining random variables.

- Finally, a fundamental property of inference problems which is a direct consequence of the Bayes theorem², is the so-called Nishimori symmetry [57, 6]: Let (X, Y) be a couple of random variables on a polish space. Let $k \geq 1$ and let $X^{(1)}, \dots, X^{(k)}$ be k i.i.d. samples (given Y) from the distribution $P(X = \cdot | Y)$, independently of every other random variables. Let us denote $\langle \cdot \rangle$ the expectation with respect to $P(X = \cdot | Y)$ and \mathbb{E} the expectation with respect to (X, Y) . Then, for all continuous bounded function f

$$\mathbb{E} \langle f(Y, X^{(1)}, \dots, X^{(k)}) \rangle = \mathbb{E} \langle f(Y, X^{(1)}, \dots, X^{(k-1)}, X) \rangle .$$

While the consequences of this identity are important, the proof is rather simple: It is equivalent to sample the couple (X, Y) according to its joint distribution or to sample first Y according to its marginal distribution and then to sample X conditionally to Y from its conditional distribution $P(X = \cdot | Y)$. Thus the $(k + 1)$ -tuple $(Y, X^{(1)}, \dots, X^{(k)})$ is equal in law to $(Y, X^{(1)}, \dots, X^{(k-1)}, X)$.

²And the fact that we are in the Bayes optimal setting where we know the statistical properties of the signal, namely the prior, and the statistical properties of the channels, namely Δ_2 and Δ_p .

The proof of Theorem 1 is obtained by using the Guerra interpolation technique to prove a lower bound for the free entropy and then by applying the Aizenman-Sims-Star scheme to get a matching upper bound.

Lower bound: Guerra interpolation We now move to the core of the proof. The first part combines the Guerra interpolation method [58] developed for matrices in [47] and tensors in [10].

Consider the interpolating Hamiltonian depending of $t \in [0, 1]$

$$\begin{aligned} \mathcal{H}_{N,t} = & - \sum_{i < j} \left[\frac{\sqrt{t}}{\Delta_2 \sqrt{N}} Y_{ij} x_i x_j + \frac{t}{2\Delta_2 N} (x_i x_j)^2 \right] + \\ & - \sum_{i_1 < \dots < i_p} \left[\frac{\sqrt{t(p-1)!}}{\Delta_p N^{\frac{p-1}{2}}} T_{i_1 \dots i_p} x_{i_1} \dots x_{i_p} + \frac{t(p-1)!}{2\Delta_p N^{p-1}} (x_{i_1} \dots x_{i_p})^2 \right] + \\ & - \sum_j \left[\sqrt{1-t} \sqrt{\frac{m^{p-1}}{\Delta_p} + \frac{m}{\Delta_2}} W_j x_j + (1-t) \left(\frac{m^{p-1}}{\Delta_p} + \frac{m}{\Delta_2} \right) x_j^* x_j + \frac{1-t}{2} \left(\frac{m^{p-1}}{\Delta_p} + \frac{m}{\Delta_2} \right) x_j^2 \right], \end{aligned} \quad (45)$$

where we have for $t = 1$ the regular Hamiltonian and for $t = 0$ the first term of Eq. (36) where W_j are i.i.d. canonical Gaussian variables. More importantly, for all $t \in [0, 1]$ we can show that the Hamiltonian above can be seen as the one emerging for an appropriate inference problem, so that the Nishimori property is kept valid for generic $t \in [0, 1]$ [47].

Given the interpolating Hamiltonian we can write the corresponding Gibbs measure,

$$P(x|W, Y, T) = \frac{1}{\mathcal{Z}_{N,t}} P_X(x) e^{H_{N,t}(x)}, \quad (46)$$

and the interpolating free entropy

$$\psi_N(t) \doteq \frac{1}{N} \mathbb{E} [\log \mathcal{Z}_{N,t}], \quad (47)$$

whose boundaries are $\psi_N(1) = \frac{1}{N} \mathcal{F}_N$ (our target) and $\psi_N(0) = \frac{1}{N} \tilde{\Phi}_{\text{RS}} + \frac{1}{4\Delta_2} m^2 + \frac{p-1}{2p\Delta_p} m^p$. We then use the fundamental theorem of calculus to write

$$\mathcal{F}_N = \psi_N(1) = \psi_N(0) + \underbrace{\frac{1}{N} \mathbb{E} \int_0^1 \left(-\frac{\partial \log \mathcal{Z}_{N,t}}{\partial t} \right) dt}_{\doteq \mathcal{R}}. \quad (48)$$

We work with the second term and use Stein's lemma which, given a well behaving function g , provides the useful relation for a canonical Gaussian variable Z : $\mathbb{E}_Z[Zg(Z)] = \mathbb{E}_Z[g'(Z)]$. This yields

$$\begin{aligned} \mathcal{R} = & -\mathbb{E} \int_0^1 \left[\frac{1}{\mathcal{Z}_{N,t}} \int dx^N \frac{\partial \mathcal{H}_{N,t}(x)}{\partial t} P_X(x) e^{\mathcal{H}_{N,t}(x)} \right] dt = -\mathbb{E} \int_0^1 \left\langle \frac{\partial \mathcal{H}_{N,t}(x)}{\partial t} \right\rangle dt \\ = & -\mathbb{E} \int_0^1 \left\langle \sum_{i < j} \frac{1}{\Delta_2 N} (x_i^* x_i x_j^* x_j) + \sum_{i_1 < \dots < i_p} \frac{(p-1)!}{\Delta_2 N^{p-1}} (x_{i_1}^* x_{i_1} \dots x_{i_p}^* x_{i_p}) - \sum_i \left(\frac{m}{2\Delta_2} + \frac{m^{p-1}}{2\Delta_p} \right) x_i^* x_i \right\rangle dt \\ = & \mathbb{E} \int_0^1 \left[\frac{1}{4\Delta_2} \left\langle \left(\frac{x \cdot x^*}{N} \right)^2 - 2m \left(\frac{x \cdot x^*}{N} \right) \right\rangle + \frac{1}{2p\Delta_p} \left\langle \left(\frac{x \cdot x^*}{N} \right)^p - pm^{p-1} \left(\frac{x \cdot x^*}{N} \right) \right\rangle \right] dt. \end{aligned}$$

where we have used the Nishimori property to replace terms such as $\langle x \rangle^2$ by $\langle x x^* \rangle$. At this point, we can write

$$\mathcal{R} = \mathbb{E} \int_0^1 \left[\frac{1}{4\Delta_2} \left\langle \left(\frac{x \cdot x^*}{N} \right)^2 - 2m \left(\frac{x \cdot x^*}{N} \right) \right\rangle \right] dt + \mathbb{E} \int_0^1 \left[\frac{1}{2p\Delta_p} \left\langle \left(\frac{x \cdot x^*}{N} \right)^p - pm^{p-1} \left(\frac{x \cdot x^*}{N} \right) \right\rangle \right] dt$$

$$= -\frac{m^2}{4\Delta_2} + \frac{1}{4\Delta_2} \mathbb{E} \int_0^1 \frac{1}{4\Delta_2} \left\langle \left(\frac{x \cdot x^*}{N} - m \right)^2 \right\rangle dt + \frac{1}{2p\Delta_p} \mathbb{E} \int_0^1 \left\langle \left(\frac{x \cdot x^*}{N} \right)^p - pm^{p-1} \left(\frac{x \cdot x^*}{N} \right) \right\rangle dt. \quad (49)$$

The first integral is clearly positive. The second one, however, seems harder to estimate. We may, however, use a simple convexity argument on the function $f(x) = x^k$. Indeed observe that $\forall a, b \geq 0$ and $p \geq 1$: $a^p - pb^{p-1}a \geq (1-p)b^p$. We would like to use this property but there is the subtlety that we need $x \cdot x^*$ to be non-negative. To bypass this problem we can add again a small perturbation that forces $x \cdot x^*$ to concentrate around a non-negative value, without affecting the ‘‘interpolating free entropy’’ $\psi_N(t)$ in the $N \rightarrow \infty$ limit. This is, again, the argument used in [10] and originally in [46]. In this way we can write

$$\begin{aligned} \mathcal{R} &\geq -\frac{m^2}{4\Delta_2} + \mathbb{E} \int_0^1 \left[\frac{1}{4\Delta_2} \left\langle \left(\frac{x \cdot x^*}{N} \right)^2 - 2m \left(\frac{x \cdot x^*}{N} \right) \right\rangle \right] dt + \frac{(1-p)m^p}{4\Delta_2} \\ &\geq -\frac{m^2}{4\Delta_2} - \frac{(p-1)m^p}{4\Delta_2}. \end{aligned} \quad (50)$$

This concludes the proof and yields the lower bound:

$$\mathcal{F}_N \geq \psi_N(0) - \frac{1}{4\Delta_2} m^2 - \frac{p-1}{2p\Delta_p} m^p = \frac{1}{N} \tilde{\Phi}_{\text{RS}}(m), \quad (51)$$

so that for all $m \geq 0$

$$\liminf_{N \rightarrow \infty} \mathcal{F}_N = \liminf_{N \rightarrow \infty} \psi_N(1) = \liminf_{N \rightarrow \infty} \left[\psi_N(0) + \int_0^1 \psi'_N(t) dt \right] \geq \tilde{\Phi}_{\text{RS}}(m).$$

Upper bound: Aizenman-Sims-Starr scheme. The matching upper bound is obtained using the Aizenman-Sims-Starr scheme [48]. This is a particularly effective tool that has been already used for these problems, see for example [31, 56, 10]. The method goes as follows. Consider the original system with N variables, \mathcal{H}_N and add an new variable x_0 so that we get an Hamiltonian \mathcal{H}_{N+1} . Define the Gibbs measures of the two systems, the first with N variables and the second with $N+1$ variables, and consider the two relative free entropies. Call $A_N = \mathbb{E} [\log \mathcal{Z}_{N+1}] - \mathbb{E} [\log \mathcal{Z}_N]$ their difference. First, we notice that we have $\limsup_N \mathcal{F}_N \leq \limsup_N A_N$ because

$$\mathcal{F}_N = \mathbb{E} \frac{1}{N} \log \mathcal{Z}_N = \frac{1}{N} \mathbb{E} \log \left(\frac{\mathcal{Z}_N}{\mathcal{Z}_{N-1}} \frac{\mathcal{Z}_{N-1}}{\mathcal{Z}_{N-2}} \dots \frac{\mathcal{Z}_1}{\mathcal{Z}_0} \right) = \frac{1}{N} \sum_i A_i \leq \sup_i A_i.$$

Moreover, we can separate the contribution of the additional variable in the Hamiltonian \mathcal{H}_{N+1} so that $\mathcal{H}_{N+1} = \tilde{\mathcal{H}}_N + x_0 z(x) + x_0^2 s(x)$, with $x = (x_1, \dots, x_N)$, and

$$\begin{aligned} z(x) &= \frac{1}{\sqrt{\Delta_2(N+1)}} \sum_{i=1}^N Z_{0i} x_i + \frac{\sqrt{(p-1)!}}{\sqrt{\Delta_p(N+1)^{(p-1)/2}}} \sum_{1 \leq i_1 < \dots < i_{p-1} \leq N} Z_{0i_1 \dots i_{p-1}} x_{i_1} \dots x_{i_{p-1}} + \\ &+ \frac{1}{\Delta_2(N+1)} \sum_{i=1}^N x_0^* x_i^* x_i + \frac{(p-1)!}{\Delta_p(N+1)^{p-1}} \sum_{1 \leq i_1 < \dots < i_{p-1} \leq N} x_0^* x_{i_1}^* x_{i_1} \dots x_{i_{p-1}}^* x_{i_{p-1}} \\ s(x) &= -\frac{1}{2\Delta_2(N+1)} \sum_{i=1}^N x_i^2 - \frac{(p-1)!}{2\Delta_p(N+1)^{p-1}} \sum_{1 \leq i_1 < \dots < i_{p-1} \leq N} (x_{i_1} \dots x_{i_{p-1}})^2 \end{aligned}$$

and \mathcal{H}_{N+1} is the same expression as Eq. (17) where the N in the denominators are replaced by $N+1$. We rewrite also $\mathcal{H}_N(x)$ as a perturbation of $\tilde{\mathcal{H}}_N$: $\mathcal{H}_N(x) = \tilde{\mathcal{H}}_N(x) + y(x) + O(1)$ with

$$\begin{aligned}
y(x) &= \frac{1}{\sqrt{\Delta_2 N}} \sum_{i < j} V_{ij} x_i x_j + \sqrt{p-1} \frac{\sqrt{(p-1)!}}{\sqrt{\Delta_p N^{p/2}}} \sum_{i_1 < \dots < i_p} V_{i_1 \dots i_p} x_{i_1} \dots x_{i_p} + \\
&+ \frac{1}{N^2} \sum_{i < j} \left[x_i^* x_i x_j^* x_j - \frac{1}{2} (x_i x_j)^2 \right] + (p-1)! \frac{p-1}{N^p} \sum_{i_1 < \dots < i_p} \left[x_{i_1}^* x_{i_1} \dots x_{i_p}^* x_{i_p} - \frac{1}{2} (x_{i_1} \dots x_{i_p})^2 \right],
\end{aligned}$$

where the Z s and the V s are standard Gaussian random variables.

Finally we can observe the partition functions Z_N can be interpreted as ensemble averages with respect to $\tilde{\mathcal{H}}_N$. Thus $A_N = \mathbb{E} \log \left\langle \int P_X(x_0) e^{x_0 z(x) + x_0^2 s(x)} dx_0 \right\rangle_{\tilde{\mathcal{H}}_N} - \mathbb{E} \log \left\langle e^{y(x)} \right\rangle_{\tilde{\mathcal{H}}_N}$. Now, using the Nishimori property and the concentration of the overlap around a non-negative value—that we denote $m(Y, T)$ since it depends explicitly on the disorder—it yields (see [31], see section 4.3 for details) eq. (36) in the thermodynamic limit, with $m(Y, T)$ instead of m . From this, we can now obtain the upper bound that concludes the proof:

$$\limsup_N \mathcal{F}_N \leq \limsup_N A_N \leq \limsup_N \mathbb{E}_{Y, T} \tilde{\Phi}_{\text{RS}}[m(Y, T)] \leq \limsup_N \sup_m \Phi_{\text{RS}}(m) \leq \tilde{\Phi}_{\text{RS}}. \quad (52)$$

B.3 State evolution of AMP and its analysis

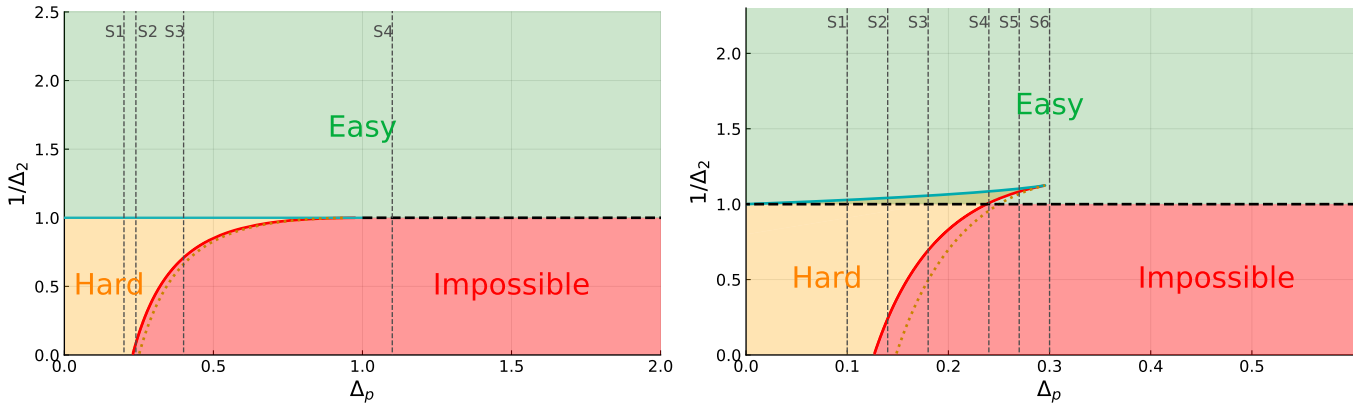


Figure 7: *On the left:* Phase diagram of the spiked matrix-tensor model for $p = 3$. The phase diagram identifies four regions: easy (green), impossible (red), and hard (orange). The lines correspond to different phase transitions namely the stability threshold (dashed black), the information theoretic threshold (solid red), the algorithmic threshold (solid cyan), and the dynamical threshold (dotted orange). The vertical cuts represent the section along which the magnetization is plotted in Fig. 9. *On the right:* Phase diagram of the spiked matrix-tensor model for $p = 4$. The main difference with respect to case $p = 3$, is that the algorithmic spinodal (solid cyan) is strictly above the stability threshold (dashed black). The hybrid-hard phase appears between these two lines (combined green and orange color). The vertical cuts represent the section along which the magnetization is plotted in Fig. 10.

The dynamical evolution of the AMP algorithm in the large N limit is described by the so-called State Evolution (SE) equations. The derivation of these equations can be straightforwardly done using the same techniques developed in [30]. They can be written in terms of two dynamical order parameters namely $m^t = \sum_i \hat{x}_i^t x_i^*/N$, which encodes for the alignment of the current estimation \hat{x}_i^t of the components of the signal with the signal itself at time t and $q^t = \sum_i \hat{x}_i^t \hat{x}_i^t/N$. Finally, using the Nishimori symmetry it can be shown that $m^t = q^t$ at all times³, see e.g. [6], and therefore the evolution of the algorithm is characterized by a single order

³Note that AMP satisfies the Nishimori property at all times while this condition is violated on the run by the Langevin dynamics. In that case the Nishimori symmetry is recovered only when equilibrium is reached and therefore it is violated when the Langevin algorithm gets trapped in the glass phase, see below.

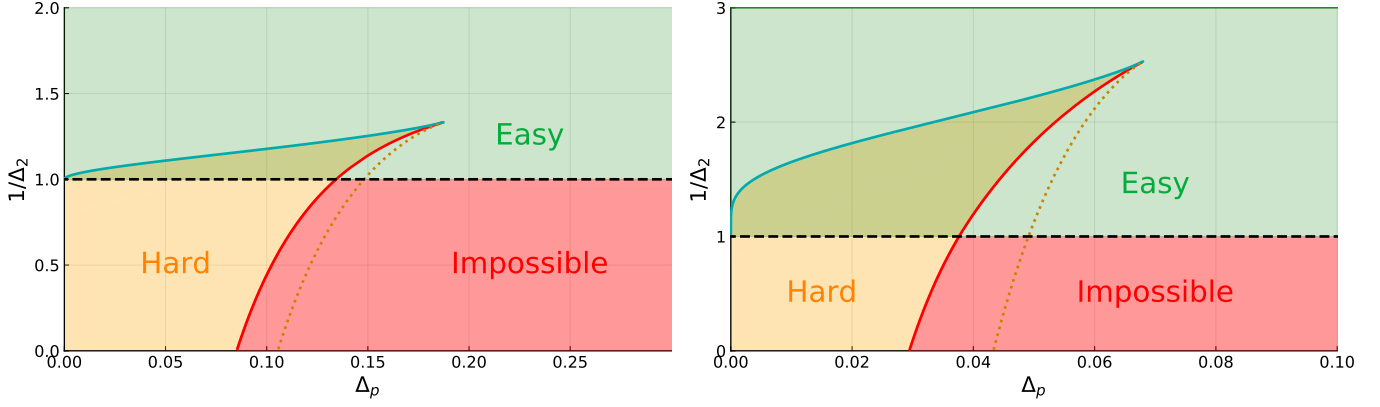


Figure 8: *On the left*: Phase diagram of the spiked matrix-tensor model for $p = 5$. *On the right*: Phase diagram of the spiked matrix-tensor model for $p = 10$. In both cases we observe qualitatively the same scenario found in the right panel of Fig. 7.

parameter m^t whose dynamical evolution is given by

$$m^{t+1} = 1 - \frac{1}{1 + \frac{m^t}{\Delta_2} + \frac{(m^t)^{p-1}}{\Delta_p}}. \quad (53)$$

If we initialize the configuration of the estimator \hat{x} at random, the initial value of m will be equal to zero on average. However finite size fluctuations will produce by chance a small bias towards the signal and therefore it is more meaningful to consider the initialization to be $m^{t=0} = \epsilon$ being ϵ an arbitrarily small positive number. We will call m_{AMP} the fixed point of Eq. (53) reached from this infinitesimal initialization. The mean-square-error (MSE) reached by AMP after convergence is then given by $\text{MSE}_{\text{AMP}} = 1 - m_{\text{AMP}}$.

We underline that Eq. (53) can be proven rigorously following [33, 17]. Finally we note that the fixed point of the SE satisfies the very same Eq. (38) that gives the replica free entropy. In the rest of this section we will study the fixed points of Eq. (53). This will allow too determine the phase diagram of the spiked matrix-tensor model.

We start by observing that $m = 0$ is a fixed point of Eq. (53). However in order to understand whether it is a possible attractor of the AMP dynamics we need to understand its local stability. This can be obtained perturbatively by expanding Eq. (53) around $m = 0$

$$m^{t+1} = \frac{m^t}{\Delta_2} + \left(\frac{m^t}{\Delta_2}\right)^2 - \frac{(m^t)^{p-1}}{\Delta_p} + O\left((m^t)^3\right). \quad (54)$$

It is clear that the non-informative fixed point $m = 0$ is stable as long as $\Delta_2 > 1$. We will call $\Delta_2 = 1$ the *stability threshold*.

When $p = 3$ the SE equations are particularly simple and the fixed points are written explicitly as

$$m_0 = 0; \quad m_{\pm} = \frac{1}{2} \left[1 - \frac{\Delta_3}{\Delta_2} \pm \sqrt{\left(1 + \frac{\Delta_3}{\Delta_2}\right)^2 - 4\Delta_3} \right]. \quad (55)$$

In the regime where $\Delta_2 > 1$, m_0 and m_+ are stable while m_- is unstable. When Δ_2 becomes smaller than one, m_+ becomes the only non-negative stable solution and therefore $\Delta_2 = 1$ is also known as the *algorithmic spinodal* since it corresponds to the point where the AMP algorithm converges to the informative fixed point. The informative solution m_+ exists as long as $\Delta_2 \leq \Delta_2^{\text{dyn}}$, where we have defined the *dynamical spinodal* by

$$\Delta_2^{\text{dyn}} = \frac{\Delta_3}{2\sqrt{\Delta_3} - 1}. \quad (56)$$

For a generic p we cannot determine the values of the informative fixed points explicitly but we can easily study Eq. (53) numerically to get the full phase diagram.

Furthermore we can obtain the spinodal transition lines as follows. The key observation is that the two spinodals are critical points of the equation $\Delta_p(m; \Delta_2)$ where Δ_2 is fixed, or analogously $\Delta_2(m; \Delta_p)$ where Δ_p is fixed (to have a pictorial representation of the idea you can see Fig. 10). We call $x = m/\Delta_2 + m^{p-1}/\Delta_p$, and $f_{\text{SE}}(x) \equiv 1 - \frac{1}{1+x}$, then

$$\Delta_p \equiv \Delta_p(x; \Delta_2) = \frac{(f_{\text{SE}}(x))^{p-1}}{x - \frac{f_{\text{SE}}(x)}{\Delta_2}}. \quad (57)$$

Then the stationary points are implicitly defined by

$$0 = \frac{d \log \Delta_p}{dm} = \frac{\partial \log \Delta_p}{\partial x} (1+x)^2 \propto (p-1) \frac{f'_{\text{SE}}(x)}{f_{\text{SE}}(x)} - \frac{1 - \frac{f'_{\text{SE}}(x)}{\Delta_2}}{x - \frac{f_{\text{SE}}(x)}{\Delta_2}} = \frac{\frac{2-p}{\Delta_2} + (1+x)(p-x-2)}{x(1+x) \left[x + 1 - \frac{1}{\Delta_p} \right]},$$

giving

$$x_{\pm}(\Delta_2) = \frac{1}{2} \left[p - 3 \pm \sqrt{(p-1)^2 - \frac{4}{\Delta_2}(p-2)} \right]. \quad (58)$$

Finally $\Delta_p(x_{\pm}(\Delta_2); \Delta_2)$ describes the two spinodals. We can also derive the tri-critical point, when the two spinodals meet, which is given by the zero discriminant condition on eq. (58)

$$\left(\Delta_p^{\text{tri}}; 1/\Delta_2^{\text{tri}} \right) = \left(\frac{4(p-2) \left(\frac{p-3}{p-1} \right)^{p-1}}{(p-3)^2}; \frac{(p-1)^2}{4(p-2)} \right). \quad (59)$$

B.4 Phase diagrams of spiked matrix-tensor model

In this section we present the phase diagrams for the spiked matrix-tensor model as a function of the two noise levels Δ_2 and Δ_p and for several values of p . These phase diagrams are plotted in Figs. 7 and 8.

Generically we can have four regions:

- **Easy phase** (green), where the MSE obtained through AMP coincides with the MMSE which is better than random sampling of the prior.
- **Impossible phase** (red), where the MMSE and MSE of AMP coincide and are equal to 1 (meaning that $m^* = m_{\text{AMP}}=0$).
- **Hard phase** (orange), where the MMSE is smaller than the MSE obtained from AMP and $m^* > m_{\text{AMP}} \geq 0$.
- **Hybrid-hard phase** [59] (mix of green and orange), is a part of the hard phase where the AMP performance is strictly better than random sampling from the prior, but still the MSE obtained this way does not match the MMSE, i.e. $m^* > m_{\text{AMP}} > 0$. The hybrid-hard phase can be found for $p \geq 4$.

All these phases are separated by the following transition lines:

- The **stability threshold** (dashed black line) at $\Delta_2 = 1$ for all p . This corresponds to the point where the uninformative fixed point $m = 0$ loses its local stability.
- The **information theoretic threshold** (solid red line). Here $m^* > 0$ and the MMSE jumps to a value strictly smaller than one.

- The **algorithmic threshold** (solid cyan line). This is where the fixed point of AMP jumps to the $\text{MMSE} < 1$. For $p = 3$ this line coincides with a segment of the stability threshold while for $p \geq 4$ it is strictly above.
- The **dynamic threshold** (dotted orange line). Here the most informative fixed point (the one with largest m_{AMP}) disappears.

In Figs. 9, and 10 we plot the evolution of the magnetization m , as found through the fixed points of the SE equation, for several fixed values of Δ_p and $p = 3$ and $p = 4$, respectively. The values of Δ_p are identified by the vertical cuts in the phase diagrams of Fig. 7.

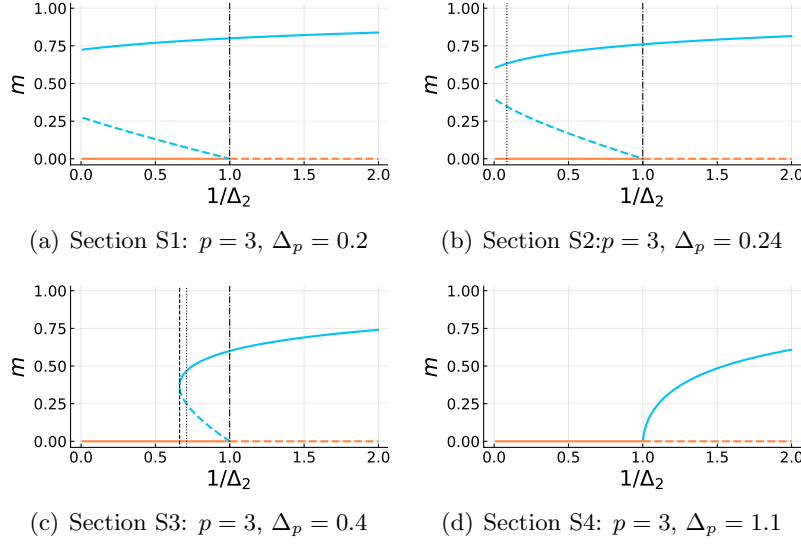


Figure 9: Fixed points of Eq. (53) as a function of Δ_2 for $p = 3$ and several fixed values of Δ_p . The values of Δ_p correspond to the vertical cuts in the left panel of Fig. 7. Solid lines are stable fixed point, dashed lines are unstable fixed points. The blue line represent informative fixed points with positive overlap with the signal while the orange line represent a uninformative fixed points with no overlap with the signal. Starting from high Δ_2 an informative fixed point appears at the dynamical threshold (vertical dashed line) but is energetically disfavored until the information theoretic threshold (vertical dotted line) and finally it becomes the only stable solution crossing the algorithmic threshold (vertical dotted-dashed line). When the transition is continuous the three vertical threshold lines merge and we have a single second order phase transition, which here occurs at $\Delta_p \geq 1$.

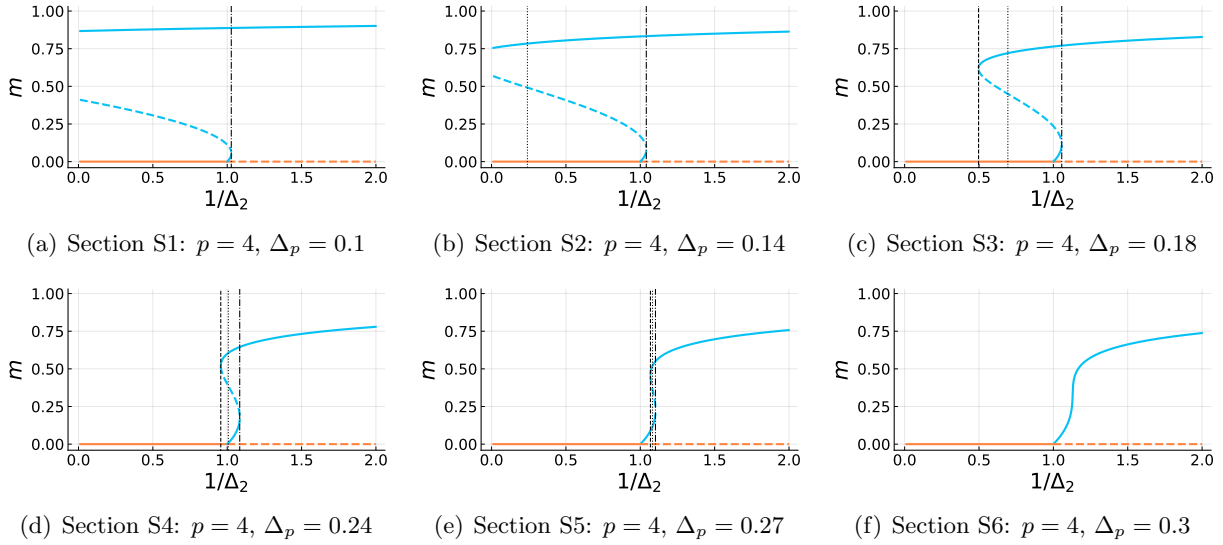


Figure 10: Fixed points Eq. (53) as a function of Δ_2 for $p = 4$ and several fixed values of Δ_p . The values of Δ_p correspond to the vertical cuts of the right panel of Fig. 7. The situation is qualitatively similar to Fig. 9, the difference being only the presence of the hybrid-hard phase. We can observe that when the transition is discontinuous, figure from (a) to (e), for $1/\Delta_2 > 1.0$ the uninformative solution becomes unstable and continuously goes to a stable-informative solution which is not the optimal one.

C Langevin Algorithm and its state evolution

The main goal of our analysis is to compare AMP with the performance of the Langevin dynamics. The advantage of the spiked matrix-tensor model is that in this case the Langevin dynamics can be studied in the large N limit through integro-differential equations for the correlation function, $C(t, t') = \lim_{N \rightarrow \infty} \sum_i \langle x_i(t) x_i(t') \rangle / N$, the response function $R(t, t') = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{d \langle x_i(t) \rangle}{d \eta_i(t')}$ and the magnetization $\bar{C}(t) = \lim_{N \rightarrow \infty} \sum_i \langle x_i(t) x_i^* \rangle / N$.

To obtain these equations we use the techniques developed in the context of mean-field spin glass systems [37, 60]. We call $\eta_i(t)$ a time dependent noise and we indicate with $\langle \cdot \rangle$ the average with respect to it. The noise is Gaussian and characterized by $\langle \eta_i(t) \rangle = 0$ for all t and $i = 1, \dots, N$ and $\langle \eta_i(t) \eta_j(t') \rangle = 2\delta_{ij} \delta(t - t')$. As before we will denote by $\mathbb{E}[\dots]$ the average with respect to the realization of disorder that in this case goes back to the specific realization of the signal.

Before proceeding, it is useful to introduce a set of auxiliary variables that will help in the following. For $k \in \{2, p\}$ we define $r_k \equiv r_k(t) = 2 / (k T_k(t) \Delta_k)$, $f_k(x) = x^k / 2$ and $m(t) \doteq \frac{1}{N} \sum_i x_i(t) x_i^*$, and the random variable $\tilde{\xi}_{i_1 \dots i_k} \equiv \frac{1}{\Delta_k} \xi_{i_1 \dots i_k} \sim \mathcal{N}(0, 1/\Delta_k)$. The time dependence in T_k , will be used in the smart annealing protocol that will be used to avoid part of the Langevin hard phase. We introduce a time dependent Hamiltonian

$$\begin{aligned} \mathcal{H}(t) = & -\frac{1}{T_2(t)\sqrt{N}} \sum_{i < j} \tilde{\xi}_{ij} x_i(t) x_j(t) - \frac{\sqrt{(p-1)!}}{T_p(t) N^{\frac{p-1}{2}}} \sum_{i_1 < \dots < i_p} \tilde{\xi}_{i_1 \dots i_p} x_{i_1}(t) \dots x_{i_p}(t) \\ & - N r_2(t) f_2(m(t)) - N r_p(t) f_p(m(t)), \end{aligned}$$

and the associated Langevin dynamics

$$\begin{aligned} \dot{x}_i(t) = & -\mu(t) x_i(t) - \frac{\partial \mathcal{H}}{\partial x_i}(t) - \eta_i(t) = -\mu(t) x_i(t) - \frac{1}{T_2(t)\sqrt{N}} \sum_{j(\neq i)} \tilde{\xi}_{ij} x_j(t) + \\ & + r_2(t) f_2'(m(t)) - \frac{\sqrt{(p-1)!}}{T_p(t) N^{\frac{p-1}{2}}} \sum_{(i, i_1, \dots, i_{p-1}) \setminus i} \tilde{\xi}_{i i_1 \dots i_{p-1}} x_{i_1}(t) \dots x_{i_{p-1}}(t) + r_p(t) f_p'(m(t)) - \eta_i(t), \end{aligned} \quad (60)$$

with μ a Langrange multiplier that enforces the spherical constraint $\sum_{i=1}^N x_i^2(t) = N$. If $T_k(t) = 1$ for all $k = 2, p$, the stationary equilibrium distribution for the Langevin dynamics is given by the posterior measure. Using Ito's lemma one finds

$$\frac{1}{N} \frac{d}{dt} \sum_i x_i^2(t) = \frac{2}{N} \sum_i x_i(t) \dot{x}_i(t) + 2.$$

Since the spherical constraint imposes the left-hand-side to be zero, one obtains a condition on the right-hand-side. By plugging the expression (60) in it, one gets that in the large N limit

$$\mu(t) = 1 - 2\mathcal{H}_2(t) - p\mathcal{H}_p(t) \quad (61)$$

where

$$\mathcal{H}_k = -\frac{\sqrt{(k-1)!}}{T_k(t)N^{\frac{k-1}{2}}} \sum_{i_1 < \dots < i_k} \tilde{\xi}_{i_1 \dots i_k} x_{i_1}(t) \dots x_{i_k}(t) - Nr_k(t)f_k(m(t)) \quad k = 2, p \quad (62)$$

are the parts of the Hamiltonian defined in Eq. (C) relative to the matrix ($k = 2$) and to the tensor ($k = p$).

Note that we have not specified any initial condition for the variables $x_i(t = 0)$. Therefore, since we always employ the spherical constraint, the initial condition for the dynamics is a point on the N dimensional hypersphere $|x|^2 = N$ extracted with the flat measure.

In order to analyze the Langevin dynamics in the large N limit, we will use the dynamical cavity method [37, 61, 62]. We will consider a system of N variables, with $N \gg 1$, and add a new one. This new variable will be considered as a small perturbation to the original system but at the same time will be treated self consistently.

C.1 Dynamical Mean-Field Equations

In the following we will drop the time dependence for simplicity restoring it only when it is needed. Given the system with N variables $i = 1 \dots N$, we add a new one, say $i = 0$, and define $\tilde{m} = \frac{1}{N+1} \sum_{i=0}^N x_i x_i^* \simeq \frac{1}{N} \sum_{i=0}^N x_i x_i^*$ (henceforth we use the symbol \simeq to denote two quantities that are equal up to terms that vanish in the large- N limit). The Langevin equation associated to the new variable is

$$\dot{x}_0 = -\mu x_0 - \frac{1}{T_2(t)\sqrt{N}} \sum_{j(\neq 0)} \tilde{\xi}_{0j} x_j + r_2 f_2'(\tilde{m}) - \frac{\sqrt{(p-1)!}}{T_p(t)N^{\frac{p-1}{2}}} \sum_{(0, i_1, \dots, i_{p-1}) \setminus 0} \tilde{\xi}_{0i_1 \dots i_{p-1}} x_{i_1} \dots x_{i_{p-1}} + r_p f_p'(\tilde{m}) - \eta_0, \quad (63)$$

where we used that $N \simeq N+1$ for $N \gg 1$. We will consider the contribution of the new variable on the others in perturbation theory. In the dynamical equations for the variables $i = 1, \dots, N$ we can isolate the variable $i = 0$ and write

$$\begin{aligned} \dot{x}_i = & -\mu x_i - \frac{1}{T_2(t)\sqrt{N}} \sum_{j(\neq i, 0)} \tilde{\xi}_{ij} x_j + r_2 f_2'(m) - \frac{\sqrt{(p-1)!}}{T_p(t)N^{\frac{p-1}{2}}} \sum_{(i, i_1, \dots, i_{p-1}) \setminus i, 0} \tilde{\xi}_{i i_1 \dots i_{p-1}} x_{i_1} \dots x_{i_{p-1}} \\ & + r_p f_p'(m) - \eta_i + H_i, \end{aligned} \quad (64)$$

with

$$H_i(t) = \left(r_2 f_2''(m) + r_p f_p''(m) \right) \frac{1}{N} x_0 - \frac{1}{T_2(t)\sqrt{N}} \tilde{\xi}_{0i} x_0 - \frac{\sqrt{(p-1)!}}{T_p(t)N^{\frac{p-1}{2}}} \sum_{(i, 0, i_1, \dots, i_{p-2}) \setminus i, 0} \tilde{\xi}_{i 0 i_1 \dots i_{p-2}} x_0 x_{i_1} \dots x_{i_{p-2}}. \quad (65)$$

Consider the unperturbed variables $x_i^0 = x_i|_{H_i=0}$. At leading order in N we can write

$$x_i \simeq x_i^0 + \int_{t_0}^t dt' \left. \frac{\delta x_i(t)}{\delta H_i(t')} \right|_{H_i=0} H_i(t'). \quad (66)$$

In the dynamical equation for the variable 0 we can identify a piece associated to the unperturbed variables x_i^0 . This term can be thought of collectively as a stochastic term $\Xi(t)$

$$\begin{aligned} \dot{x}_0 = & \overbrace{-\mu x_0 - \frac{1}{T_2(t)\sqrt{N}} \sum_{j(\neq 0)} \tilde{\xi}_{0j} x_j^0 - \frac{\sqrt{(p-1)!}}{T_p(t)N^{\frac{p-1}{2}}} \sum_{(0,i_1,\dots,i_{p-1})\setminus 0} \tilde{\xi}_{0i_1\dots i_{p-1}} x_{i_1}^0 \cdots x_{i_{p-1}}^0 - \eta_0 +}^{\doteq \Xi(t)} \\ & + r_2 f_2'(m) + r_p f_p'(m) + \left(r_2 f_2''(m) + r_p f_p''(m) \right) \frac{1}{N} x_0 - \frac{1}{T_2(t)\sqrt{N}} \sum_{j(\neq 0)} \tilde{\xi}_{0j} \int_{t_0}^t dt' \frac{\delta x_j(t)}{\delta H_j(t')} \Big|_{H_j=0} H_j(t') + \\ & - \left[\frac{\sqrt{(p-1)!}}{T_p(t)N^{\frac{p-1}{2}}} \sum_{(0,i_1,\dots,i_{p-1})\setminus 0} \tilde{\xi}_{0i_1\dots i_{p-1}} \int_{t_0}^t dt' \frac{\delta x_{i_1}(t)}{\delta H_{i_1}(t')} \Big|_{H_{i_1}=0} H_{i_1}(t') x_{i_2}^0 \cdots x_{i_{p-1}}^0 + \text{permutations} \right]. \end{aligned} \quad (67)$$

Indeed $\Xi(t)$ encodes the effect of a kind of bath made by of the unperturbed variables $i = 1, \dots, N$ to the new one. We can show that at leading order in N , $\Xi(t)$ is a Gaussian noise with zero mean and variance given by

$$\begin{aligned} \mathbb{E}\langle \Xi(t)\Xi(t') \rangle = & 2\delta(t-t') - \mathbb{E} \left[\frac{1}{T_2(t)T_2(t')N} \sum_{j(\neq 0)} \sum_{l(\neq 0)} \tilde{\xi}_{0j} \tilde{\xi}_{0l} x_j^0(t) x_l^0(t') \right] + \\ & - \mathbb{E} \left[\frac{(p-1)!}{T_p(t)T_p(t')N^{p-1}} \sum_{(0,i_1,\dots,i_{p-1})\setminus 0} \sum_{(0,j_1,\dots,j_{p-1})\setminus 0} \tilde{\xi}_{0i_1\dots i_{p-1}} \tilde{\xi}_{0j_1\dots j_{p-1}} x_{i_1}^0 \cdots x_{i_{p-1}}^0 x_{j_1}^0 \cdots x_{j_{p-1}}^0 \right] \end{aligned}$$

and the second term can be simplified as

$$\begin{aligned} \mathbb{E} \left[\frac{(p-1)!}{T_p(t)T_p(t')N^{p-1}} \sum_{(0,i_1,\dots,i_{p-1})\setminus 0} \sum_{(0,j_1,\dots,j_{p-1})\setminus 0} \tilde{\xi}_{0i_1\dots i_{p-1}} \tilde{\xi}_{0j_1\dots j_{p-1}} x_{i_1}^0 \cdots x_{i_{p-1}}^0 x_{j_1}^0 \cdots x_{j_{p-1}}^0 \right] = \\ \simeq \frac{(p-1)!}{N^{p-1}} \frac{1}{T_p(t)T_p(t')\Delta_p} \sum_{(0,i_1,\dots,i_{p-1})\setminus 0} \langle x_{i_1}^0(t) x_{i_1}^0(t') \cdots x_{i_{p-1}}^0(t) x_{i_{p-1}}^0(t') \rangle = \frac{1}{T_p(t)T_p(t')\Delta_p} C^{p-1}(t, t'), \end{aligned}$$

where we used $\sum_{(i_1,\dots,i_k)} = \frac{1}{k!} \sum_{1 \leq i_1, \dots, i_k \leq N}$, we neglected terms sub-leading in N , and we used the definition of the dynamical correlation function

$$C(t, t') = \frac{1}{N} \sum_{i=1}^N \langle x_i(t) x_i(t') \rangle.$$

Therefore we have

$$\mathbb{E}\langle \Xi(t) \rangle = 0; \quad (68)$$

$$\mathbb{E}\langle \Xi(t)\Xi(t') \rangle = 2\delta(t-t') + \frac{1}{T_2(t)T_2(t')} C(t, t') + \frac{1}{T_p(t)T_p(t')\Delta_p} C^{p-1}(t, t'). \quad (69)$$

Now we can focus of the deterministic term coming from the first order perturbation in eq. (67). Consider just the integral for the p -body term, the other will be given by setting $p = 2$

$$\begin{aligned} \frac{\sqrt{(p-1)!}}{T_p(t)N^{\frac{p-1}{2}}} \sum_{(0,i_1,\dots,i_{p-1})\setminus 0} \tilde{\xi}_{0i_1\dots i_{p-1}} \int_{t_0}^t dt' \frac{\delta x_{i_1}(t)}{\delta H_{i_1}(t')} \Big|_{H_{i_1}=0} H_{i_1}(t') x_{i_2}^0 \cdots x_{i_{p-1}}^0 + \text{permutations} = \\ \simeq \frac{(p-1)!}{T_p(t)N^{p-1}} \sum_{(0,i_1,\dots,i_{p-1})\setminus 0} \tilde{\xi}_{0i_1\dots i_{p-1}}^2 \int_{t_0}^t dt' \frac{1}{T_p(t')} \frac{\delta x_{i_1}(t)}{\delta H_{i_1}(t')} \Big|_{H_{i_1}=0} x_{i_1}^0(t) x_{i_1}^0(t') \cdots x_{i_{p-2}}^0(t) x_{i_{p-2}}^0(t') x_0(t') + \\ + \text{permutations} \simeq -p(p-1) \frac{1}{2T_p(t)\Delta_p} \int_{t_0}^t dt' \frac{1}{T_p(t')} R(t, t') C^{p-2}(t, t') x_0(t') \end{aligned} \quad (70)$$

where we have used the definition of the response function

$$R(t, t') = \frac{1}{N} \sum_{i=1}^N \left\langle \frac{\delta x_i(t)}{\delta H_i(t')} \right\rangle.$$

Plugging eq. (70) into eq. (67) we obtain an effective dynamical equation for the new variable in terms of the correlation and response function of the system with N variables

$$\begin{aligned} \dot{x}_0(t) = & -\mu(t)x_0(t) + \Xi(t) + r_p f'_p(\bar{C}(t)) + r_2 f'_2(\bar{C}(t)) + \\ & + (p-1) \frac{1}{T_p(t)\Delta_p} \int_{t_0}^t dt'' \frac{1}{T_p(t'')} R(t, t'') C^{p-2}(t, t'') x_0(t'') + \frac{1}{T_2(t)\Delta_2} \int_{t_0}^t dt'' \frac{1}{T_p(t'')} R(t, t'') x_0(t''). \end{aligned} \quad (71)$$

In order to close Eq. (71) we need to give the recipe to compute the correlation and response function.

C.2 Integro-differential equations

In order to obtain the final equations for dynamical order parameters we will assume that the new variable x_0 is a typical one, namely it has the same statistical nature of all the others. Therefore we can assume that

$$\begin{aligned} C(t, t') & \doteq \mathbb{E}\langle x_0(t)x_0(t') \rangle \\ R(t, t') & \doteq \mathbb{E} \left\langle \frac{\delta x_0(t)}{\delta \Xi(t')} \right\rangle \\ \bar{C}(t) & \doteq \mathbb{E}\langle x_0(t)x_0^* \rangle. \end{aligned} \quad (72)$$

Eqs. (72) give a way to obtain the equation for all the correlation functions. Indeed we can consider Eq. (71), multiply it by $x_0(t')$, or differentiate it with respect to an external field $h_0(t')$, or multiply it by x_0^* and we can average the results over the disorder and thermal noise. Using the following identity

$$\begin{aligned} \mathbb{E}\langle \Xi(t)x_0(t') \rangle & = \int \mathcal{D}\Xi(t) \Xi(t)x_0(t') e^{-\int d\bar{t} d\bar{t}' \Xi(\bar{t}) \mathbb{K}^{-1}(\bar{t}, \bar{t}') \Xi(\bar{t}')} = \\ & = - \int dt'' \int \mathcal{D}\Xi(t) x_0(t') \frac{\delta}{\delta \Xi(t'')} e^{-\int d\bar{t} d\bar{t}' \Xi(\bar{t}) \mathbb{K}^{-1}(\bar{t}, \bar{t}') \Xi(\bar{t}')} \mathbb{K}(t, t'') = \\ & = \int dt'' \mathbb{E} \left\langle \frac{\delta x_0(t')}{\delta \Xi(t'')} \mathbb{K}(t, t'') \right\rangle = \int dt'' R(t', t'') \mathbb{K}(t, t'') = \\ & = 2R(t', t) + \frac{1}{T_p(t)\Delta_p} \int_{t_0}^{t'} dt'' \frac{1}{T_p(t'')} R(t', t'') C^{p-1}(t, t'') + \frac{1}{T_2(t)\Delta_2} \int_{t_0}^{t'} dt'' \frac{1}{T_p(t'')} R(t', t'') C(t, t'') \end{aligned} \quad (73)$$

we get the following Langevin State Evolution (LSE) equations

$$\begin{aligned} \frac{\partial}{\partial t} C(t, t') & = \mathbb{E}\langle \dot{x}_0(t)x_0(t') \rangle = 2R(t', t) - \mu(t)C(t, t') + r_p(t) f'_p(\bar{C}(t)) \bar{C}(t') + r_2(t) f'_2(\bar{C}(t)) \bar{C}(t') + \\ & + (p-1) \frac{1}{T_p(t)\Delta_p} \int_{t_0}^t dt'' \frac{1}{T_p(t'')} R(t, t'') C^{p-2}(t, t'') C(t', t'') + \\ & + \frac{1}{T_p(t)\Delta_p} \int_{t_0}^{t'} dt'' \frac{1}{T_p(t'')} R(t', t'') C^{p-1}(t, t'') + \\ & + \frac{1}{T_2(t)\Delta_2} \int_{t_0}^t dt'' \frac{1}{T_2(t'')} R(t, t'') C(t', t'') + \frac{1}{T_2(t)\Delta_2} \int_{t_0}^{t'} dt'' \frac{1}{T_p(t'')} R(t', t'') C(t, t''); \\ \frac{\partial}{\partial t} R(t, t') & = \mathbb{E} \left\langle \frac{\delta \dot{x}_0(t)}{\delta \Xi(t')} \right\rangle = \end{aligned} \quad (74)$$

$$\begin{aligned} & = \delta(t-t') - \mu(t)R(t, t') + (p-1) \frac{1}{T_p(t)\Delta_p} \int_{t'}^t dt'' \frac{1}{T_p(t'')} R(t, t'') R(t'', t') C^{p-2}(t, t'') + \\ & + \frac{1}{T_2(t)\Delta_2} \int_{t'}^t dt'' \frac{1}{T_2(t'')} R(t, t'') R(t'', t'); \end{aligned} \quad (75)$$

$$\begin{aligned} \frac{\partial}{\partial t} \bar{C}(t) &= \mathbb{E} \langle \dot{x}_0(t) x_0^* \rangle = \\ &= -\mu(t) \bar{C}(t) + r_p(t) f'_p(\bar{C}(t)) + r_2(t) f'_2(\bar{C}(t)) + \end{aligned} \quad (76)$$

$$\begin{aligned} &+ (p-1) \frac{1}{T_p(t) \Delta_p} \int_{t_o}^t dt'' \frac{1}{T_p(t'')} R(t, t'') C^{p-2}(t, t'') \bar{C}(t'') + \frac{1}{T_2(t) \Delta_2} \int_{t_o}^t dt'' \frac{1}{T_2(t'')} R(t, t'') \bar{C}(t''); \\ \mu(t) &= 1 + r_p(t) f'_p(\bar{C}(t)) \bar{C}(t) + r_2(t) f'_2(\bar{C}(t)) \bar{C}(t) + \\ &+ p \frac{1}{T_p(t) \Delta_p} \int_{t_o}^t dt'' \frac{1}{T_p(t'')} R(t, t'') C^{p-1}(t, t'') + 2 \frac{1}{T_2(t) \Delta_2} \int_{t_o}^t dt'' \frac{1}{T_2(t'')} R(t, t'') C(t, t''). \end{aligned} \quad (77)$$

Note that the last equation for $\mu(t)$ is obtained by imposing the spherical constraint $C(t, t) = 1 \forall t$ using the fact that $0 = \frac{dC(t, t)}{dt} = \left. \frac{\partial C(t, t')}{\partial t} \right|_{t'=t} + \left. \frac{\partial C(t', t)}{\partial t} \right|_{t'=t}$. The boundary conditions of this equations are: $C(t, t) = 1$ the spherical constrain, $R(t, t) = 0$ which comes from causality in the Itô approach and $R(t, t' \rightarrow t^-) = 1$. The initial condition for $\bar{C}(0) = \bar{C}_0$ is the overlap with the initial configuration with the true signal. If the initial configuration is random, $\bar{C}_0 = 0$ but will have finite size fluctuations, as in the case of AMP. Therefore we can think that $\bar{C}_0 = \epsilon$ being ϵ an arbitrary small positive number.

D Numerical solution of the LSE equations

The dynamical equations (74-75-76-77) were integrated numerically using two schemes:

- fixed time-grid: the derivatives were discretized and integrated according to their causal structure. This method is suited only for short times (up to 500 time units);
- dynamic time-grid: the step size is doubled after a given number of steps and the equations are solved self-consistently for every waiting-time. This is the approach proposed in [38] and described in Appendix C of [63]. It allows integration up to very large times (up to 10^6 time units).

The results of these algorithms are concisely reported in the phase diagram shown in the main paper. In what follows we will present the algorithms and a series of investigations that we carried out to check their stability, we will explain the procedure followed to delimit the Langevin hard region, and we will discuss how we can enter into part of that region by choosing a proper annealing protocol. The codes are available online [40].

D.1 Fixed time-grid $(2 + p)$ -spin

In this approach time-derivatives and integrals were discretized using $\frac{\partial}{\partial t} f(t, t') \simeq \frac{1}{\Delta t} [f(t + \Delta t, t') - f(t, t')]$, and the trapezoidal rule for integration $\int_0^t f(t) dt \simeq \frac{\Delta t}{2} \sum_{l=0}^{t/\Delta t-1} [f(l\Delta t) + f((l+1)\Delta t)]$. For instance we defined a function for computing the update in the the response function, eq. (75) as follows

$$\begin{aligned} R(t + \Delta t, t') &= R(t, t') - \Delta t \mu(t) R(t, t') + \frac{1}{2} \frac{\Delta t^2}{\Delta_2} \sum_{l=t'/\Delta t}^{t/\Delta t-1} [R(t, l\Delta t) R(l\Delta t, t') + R(t, (l+1)\Delta t) R((l+1)\Delta t, t')] + \\ &+ (p-1) \frac{\Delta t^2}{\Delta_p} \sum_{l=t'/\Delta t}^{t/\Delta t-1} \left[C^{p-2}(t, l\Delta t) R(t, l\Delta t) R(l\Delta t, t') + C^{p-2}(t, (l+1)\Delta t) R(t, (l+1)\Delta t) R((l+1)\Delta t, t') \right]. \end{aligned}$$

Analogously we defined the other integrators. A simple causal integration scheme, being careful with the Itô prescription, gives the pseudo-code below.

```

C(0, 0) ← 1; R(0, 0) ← 0;  $\bar{C}(0) \leftarrow \bar{C}_0$ ;
for  $t \leq t_{\max}$  do

```

```

 $C(t + \Delta t, t + \Delta t) \leftarrow 1; R(t + \Delta t, t + \Delta t) \leftarrow 0;$ 
 $\mu(t) \leftarrow \text{compute\_mu}(C, R, \bar{C}, t);$ 
 $\bar{C}(t + \Delta t) \leftarrow \text{compute\_mag}(\mu, C, R, \bar{C}, t);$ 
for  $t' \leq t$  do
   $C(t + \Delta t, t') \leftarrow \text{compute\_C}(\mu, C, R, \bar{C}, t);$ 
   $R(t + \Delta t, t') \leftarrow \text{compute\_R}(\mu, C, R, \bar{C}, t);$ 
end for
 $R(t + \Delta t, t) \leftarrow 1;$ 
end for

```

D.2 Dynamical time-grid $(2 + p)$ -spin

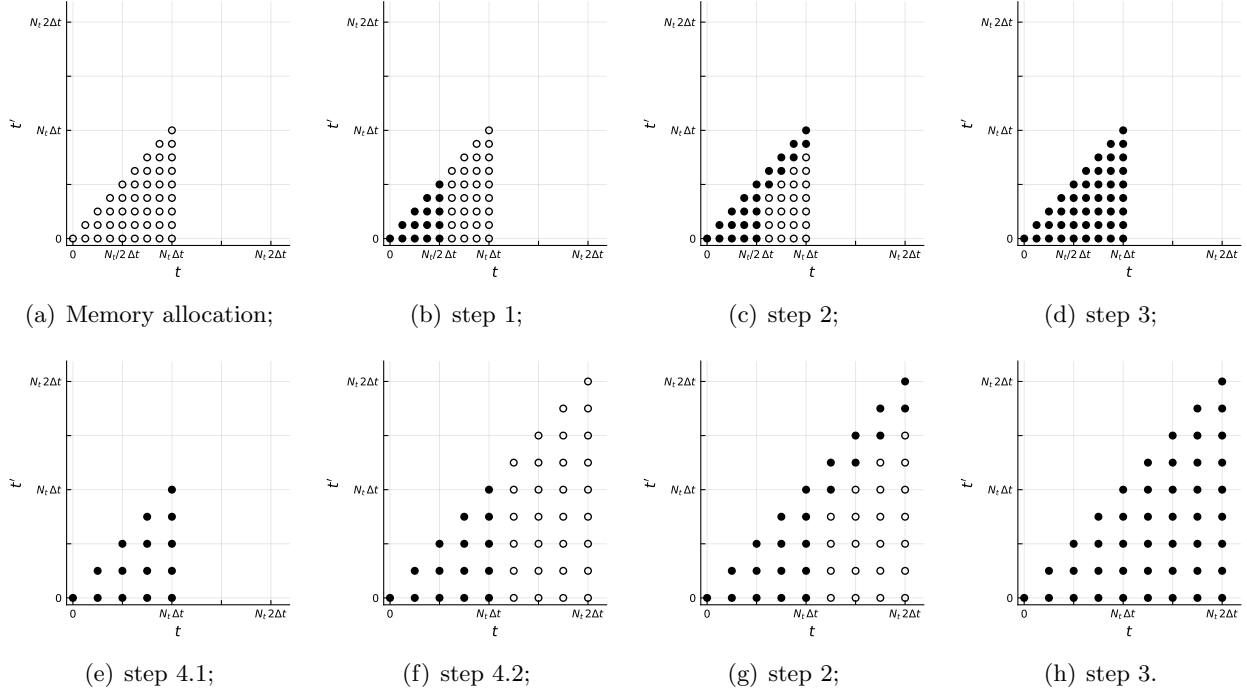


Figure 11: Representation of the initialization and the first two iterations for the evaluation of a two-times observable using the dynamic-grid algorithm. The empty circles represent slots allocated in memory but not associated to any specific value, while the full circles are memory slots already associated. For any two time function, it first allocates the memory (a), then it fills half of the grid by linear propagation (b). Still using linear propagation it fills the slots with $t - t' \ll 1$ (c), and it sets the other values by imposing self-consistency (d). Finally it halves the grid (e), doubles time step and it allocates the memory (f). Then the algorithm loops following the same scheme as in (b-c-d-e).

The numerical scheme we are going to discuss is presented in the Bayes-Optimal case where $T_2(t) \equiv T_p(t) \equiv 1$. However the derivation that we propose can be easily generalized to the case where the T s assume different values, but are constants⁴. It is convenient to manipulate the equations to obtain an equivalent set of equations for the functions $C(t, t')$, $Q(t, t') \doteq 1 - C(t, t') - \int_{t'}^t R(t, t'') dt''$, $\bar{C}(t)$, where $Q(t, t')$ represents the deviation from Fluctuation Dissipation Theorem (FDT) at time t starting from time t' . Indeed when the FDT theorem holds, it states that $R(t, t') = -\partial_t C(t, t')$.

⁴Therefore we do not employ this algorithm to solve the LSE equations in the smart annealing protocol for which instead we use the fixed time-grid algorithm.

We briefly anticipate the strategy that the algorithm uses to solve the equations. The algorithm discretizes the times into N_t intervals, first starting from the boundary conditions, $C(t, t) = 1$, $Q(t, t) = 0$ and $\bar{C} = \bar{C}_0 \in [0, 1]$, it fills the grid for small times (or small time differences $\tau = t - t' \ll 1$) using linear propagation. Given a time t and the initial guess for the Lagrange multiplier obtained by the linear propagator, the integrals are discretized and evaluated, then the results is used to update the value of the Lagrange multiplier. This procedure is repeated iteratively until convergence. Once that the first grid is filled, it follows a coarse-graining procedure where the sizes of the time intervals is doubled and only half of the information is retained. This procedure is repeated a fixed number of doubling of the original grid. The doubling scheme allows to explore exponentially long times at the cost of loosing part of the information, the direct consequence of this is the loss of stability for very large times (especially when the functions $C(t, t')$, $R(t, t')$, $\bar{C}(t)$ undergo fast changes at large times).

Dynamical equations in the algorithm. We recall the function $f_k(x) = \frac{x^k}{2}$ and its derivatives, $f'_k(x) = \frac{kx^{k-1}}{2}$ and $f''_k(x) = \frac{k(k-1)x^{k-2}}{2}$. For simplicity in the notation, we introduce also $f_k(t, t') \doteq f_k(C(t, t'))$

$$\begin{aligned}
(\partial_t + \mu(t))C(t, t') &= 2R(t', t) + r_2\bar{C}(t')f'_2(\bar{C}(t)) + r_p\bar{C}(t')f'_p(\bar{C}(t)) + \\
&\quad + \frac{1}{\Delta_2} \int_0^{t'} dt'' f'_2(t, t'')R(t', t'') + \frac{1}{\Delta_2} \int_0^t dt'' f''_2(t, t'')R(t, t'')C(t', t'') + \\
&\quad + \frac{2}{p\Delta_p} \int_0^{t'} dt'' f'_p(t, t'')R(t', t'') + \frac{2}{p\Delta_p} \int_0^t dt'' f''_p(t, t'')R(t, t'')C(t', t''), \\
(\partial_t + \mu(t))R(t, t') &= \delta(t - t') + \frac{1}{\Delta_2} \int_{t'}^t dt'' f''_2(t, t'')R(t, t'')R(t'', t') + \\
&\quad + \frac{2}{p\Delta_p} \int_{t'}^t dt'' f''_p(t, t'')R(t, t'')R(t'', t'), \\
(\partial_t + \mu(t))\bar{C}(t) &= r_2f'_2(\bar{C}(t)) + r_pf'_p(\bar{C}(t)) + \\
&\quad + \frac{1}{\Delta_2} \int_0^t dt'' f''_2(t, t'')R(t, t'')\bar{C}(t'') + \frac{2}{p\Delta_p} \int_0^t dt'' f''_p(t, t'')R(t, t'')\bar{C}(t''), \\
\mu(t) &= 1 + r_2\bar{C}(t)f'_2(\bar{C}(t)) + r_p\bar{C}(t)f'_p(\bar{C}(t)) + \\
&\quad + \frac{2}{\Delta_2} \int_0^t dt'' f'_2(t, t'')R(t, t'') + \frac{2}{\Delta_p} \int_0^t dt'' f'_p(t, t'')R(t, t'').
\end{aligned}$$

Following the lines of [63], we introduce the FDT violation function, $Q(t, t')$, and after some manipulation the systems becomes

$$\begin{aligned}
(\partial_t + \mu(t))C(t, t') &= \bar{C}(t') \left[r_2f'_2(\bar{C}(t)) + r_pf'_p(\bar{C}(t)) \right] + \\
&\quad + \frac{1}{\Delta_2} \left\{ \int_0^{t'} dt'' \left[f'_2(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f''_2(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] + \right. \\
&\quad \left. - \int_{t'}^t dt'' \left[f'_2(t, t'') \frac{\partial C(t'', t')}{\partial t''} - f''_2(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t'', t') \right] + f'_2(1)C(t, t') - f'_2(t, 0)C(t', 0) \right\} + \\
&\quad + \frac{2}{p\Delta_p} \left\{ \int_0^{t'} dt'' \left[f'_p(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f''_p(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] + \right. \\
&\quad \left. - \int_{t'}^t dt'' \left[f'_p(t, t'') \frac{\partial C(t'', t')}{\partial t''} - f''_p(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t'', t') \right] + f'_p(1)C(t, t') - f'_p(t, 0)C(t', 0) \right\}, \tag{78}
\end{aligned}$$

$$\begin{aligned}
(\partial_t + \mu(t))Q(t, t') &= \mu(t) - 1 + \frac{1}{\Delta_2} \left\{ - \int_{t'}^t dt'' f_2'(t, t'') \frac{\partial Q(t'', t')}{\partial t''} + \int_{t'}^t dt'' f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} [Q(t'', t') - 1] + \right. \\
&+ f_2'(1)[Q(t, t') - 1] + f_2'(t, 0)C(t', 0) - \int_0^{t'} dt'' \left[f_2'(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] \left. \right\} + \\
&+ \frac{2}{p\Delta_p} \left\{ - \int_{t'}^t dt'' f_p'(t, t'') \frac{\partial Q(t'', t')}{\partial t''} + \int_{t'}^t dt'' f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} [Q(t'', t') - 1] + \right. \\
&+ f_p'(1)[Q(t, t') - 1] + f_p'(t, 0)C(t', 0) - \int_0^{t'} dt'' \left[f_p'(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] \left. \right\} + \\
&- \bar{C}(t') \left[r_2 f_2'(\bar{C}(t)) + r_p f_p'(\bar{C}(t)) \right], \tag{79}
\end{aligned}$$

$$\begin{aligned}
(\partial_t + \mu(t))\bar{C}(t) &= r_2 f_2'(\bar{C}(t)) + r_p f_p'(\bar{C}(t)) + \frac{1}{\Delta_2} \left\{ f_2'(1)\bar{C}(t) - f_2'(t, 0)\bar{C}(0) - \int_0^t dt'' f_2'(t, t'') \frac{d}{dt''} \bar{C}(t'') + \right. \\
&+ \int_0^t dt'' f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} \bar{C}(t'') \left. \right\} + \frac{2}{p\Delta_p} \left\{ f_p'(1)\bar{C}(t) - f_p'(t, 0)\bar{C}(0) - \int_0^t dt'' f_p'(t, t'') \frac{d}{dt''} \bar{C}(t'') + \right. \\
&+ \int_0^t dt'' f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} \bar{C}(t'') \left. \right\}, \tag{80}
\end{aligned}$$

$$\begin{aligned}
\mu(t) &= 1 + r_2 \bar{C}(t) f_2'(\bar{C}(t)) + r_p \bar{C}(t) f_p'(\bar{C}(t)) + \frac{2}{\Delta_2} [f_2(1) - f_2(t, 0)] + \frac{2}{\Delta_p} [f_p(1) - f_p(t, 0)] + \\
&+ \int_0^t dt'' \left[\frac{2}{\Delta_2} f_2'(t, t'') + \frac{2}{\Delta_p} f_p'(t, t'') \right] \frac{\partial Q(t, t'')}{\partial t''}, \tag{81}
\end{aligned}$$

further simplifications can be obtained introducing $\mu'(t) = \mu(t) - \frac{2}{\Delta_2} f_2(1) - \frac{2}{\Delta_p} f_p(1)$

$$\begin{aligned}
\mu'(t) &= 1 + r_2 \bar{C}(t) f_2'(\bar{C}(t)) + r_p \bar{C}(t) f_p'(\bar{C}(t)) - \frac{2}{\Delta_2} f_2(t, 0) - \frac{2}{\Delta_p} f_p(t, 0) + \\
&+ \int_0^t dt'' \left[\frac{2}{\Delta_2} f_2'(t, t'') + \frac{2}{\Delta_p} f_p'(t, t'') \right] \frac{\partial Q(t, t'')}{\partial t''}, \tag{82}
\end{aligned}$$

$$\begin{aligned}
(\partial_t + \mu'(t))C(t, t') &= \bar{C}(t') \left[r_2 f_2'(\bar{C}(t)) + r_p f_p'(\bar{C}(t)) \right] + \\
&+ \frac{1}{\Delta_2} \left\{ \int_0^{t'} dt'' \left[f_2'(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] + \right. \\
&- \int_{t'}^t dt'' \left[f_2'(t, t'') \frac{\partial C(t'', t')}{\partial t''} - f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t'', t') \right] - f_2'(t, 0)C(t', 0) \left. \right\} + \\
&+ \frac{2}{p\Delta_p} \left\{ \int_0^{t'} dt'' \left[f_p'(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] + \right. \\
&- \int_{t'}^t dt'' \left[f_p'(t, t'') \frac{\partial C(t'', t')}{\partial t''} - f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t'', t') \right] - f_p'(t, 0)C(t', 0) \left. \right\}, \tag{83}
\end{aligned}$$

$$\begin{aligned}
(\partial_t + \mu'(t))Q(t, t') &= \mu'(t) - 1 - \bar{C}(t') \left[r_2 f_2'(\bar{C}(t)) + r_p f_p'(\bar{C}(t)) \right] + \\
&+ \frac{1}{\Delta_2} \left\{ - \int_{t'}^t dt'' f_2''(t, t'') \frac{\partial Q(t'', t')}{\partial t''} + \int_{t'}^t dt'' f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} [Q(t'', t') - 1] + \right. \\
&+ f_2'(t, 0)C(t', 0) - \int_0^{t'} dt'' \left[f_2'(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] \left. \right\} + \\
&+ \frac{2}{p\Delta_p} \left\{ - \int_{t'}^t dt'' f_p'(t, t'') \frac{\partial Q(t'', t')}{\partial t''} + \int_{t'}^t dt'' f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} [Q(t'', t') - 1] + \right. \\
&+ f_p'(t, 0)C(t', 0) - \int_0^{t'} dt'' \left[f_p'(t, t'') \frac{\partial Q(t', t'')}{\partial t''} + f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} C(t', t'') \right] \left. \right\}, \tag{84}
\end{aligned}$$

$$\begin{aligned}
(\partial_t + \mu'(t))\bar{C}(t) &= r_2 f_2'(\bar{C}(t)) + r_p f_p'(\bar{C}(t)) + \frac{1}{\Delta_2} \left\{ - f_2'(t, 0)\bar{C}(0) + \right. \\
&- \int_0^t dt'' f_2'(t, t'') \frac{d}{dt''} \bar{C}(t'') + \int_0^t dt'' f_2''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} \bar{C}(t'') \left. \right\} + \\
&+ \frac{2}{p\Delta_p} \left\{ - f_p'(t, 0)\bar{C}(0) + \right. \\
&- \int_0^t dt'' f_p'(t, t'') \frac{d}{dt''} \bar{C}(t'') + \int_0^t dt'' f_p''(t, t'') \frac{\partial Q(t, t'')}{\partial t''} \bar{C}(t'') \left. \right\}. \tag{85}
\end{aligned}$$

First order expansion coefficients. In the numerics we will initialize the grid by a linear propagation of the initial conditions. To determine the coefficients to use we can expand the functions up the second term for small values of τ (and in the last equation of t)

$$\begin{aligned}
C(t' + \tau, t') &= C(t', t') + C^{(1,0)}(t', t')\tau + \frac{1}{2}C^{(2,0)}(t', t') + O(\tau^3), \\
Q(t' + \tau, t') &= Q(t', t') + Q^{(1,0)}(t', t')\tau + \frac{1}{2}Q^{(2,0)}(t', t') + O(\tau^3), \\
\bar{C}(t) &= \bar{C}(0) + \bar{C}^{(1)}(0)\tau + \frac{1}{2}\bar{C}^{(2)}(0) + O(\tau^3). \tag{86}
\end{aligned}$$

This gives the following coefficients: $C(t, t) = 1$, $C^{(1,0)}(t, t) = -1$, $Q(t, t) = 0$, $Q^{(1,0)}(t, t) = 0$, $\bar{C}(0) = \bar{C}_0$ and $\bar{C}^{(1)}(0) = \left[r_2 f_2'(\bar{C}_0) + r_p f_p'(\bar{C}_0) \right] (1 - (\bar{C}_0)^2) - \bar{C}_0$, where \bar{C}_0 is the initial value of the overlap with the signal.

Numerical integration and derivation. The set of equations derived above presents six types of integrals

$$\begin{aligned}
I_{ij}^{(1AB)} &= \int_{t_j}^{t_i} dt'' A(t_i, t'') \frac{\partial B(t'', t_j)}{\partial t''}; \\
I_{ij}^{(2ABC)} &= \int_{t_j}^{t_i} dt'' A(t_i, t'') \frac{\partial B(t_i, t'')}{\partial t''} C(t'', t_j); \\
I_{ij}^{(3AB)} &= \int_0^{t_j} dt'' A(t_i, t'') \frac{\partial B(t_i, t'')}{\partial t''}; \\
I_{ij}^{(4ABC)} &= \int_0^{t_j} dt'' A(t_i, t'') \frac{\partial B(t_i, t'')}{\partial t''} C(t_j, t''); \\
I_i^{(5AB)} &= \int_0^{t_i} dt'' A(t_i, t'') \frac{\partial B(t'')}{\partial t''}; \\
I_i^{(6ABC)} &= \int_0^{t_i} dt'' A(t_i, t'') \frac{\partial B(t_i, t'')}{\partial t''} C(t'').
\end{aligned}$$

The integrals can be easily discretized

$$\begin{aligned}
I_{ij}^{(2ABC)} &= \sum_{t_l=t_j+\delta t}^{t_i} \int_{t_l-\delta t}^{t_l} dt'' A(t_i, t'') \frac{\partial B(t_i, t'')}{\partial t''} C(t'', t_j) \simeq \\
&\simeq \sum_{t_l=t_j+\delta t}^{t_i} \int_{t_l-\delta t}^{t_l} dt_1 A(t_i, t_1) \int_{t_l-\delta t}^{t_l} dt_2 \frac{\partial B(t_i, t_2)}{\partial t_2} \int_{t_l-\delta t}^{t_l} dt_3 C(t_3, t_j) \simeq \\
&\simeq \sum_{t_l=t_j+\delta t}^{t_i} \frac{1}{2} [A(t_i, t_l) + A(t_i, t_l - \delta t)] [B(t_i, t_l) - B(t_i, t_l - \delta t)] \frac{1}{2} [C(t_l, t_j) + C(t_l - \delta t, t_j)].
\end{aligned}$$

In particular the 6 integrals become

$$\begin{aligned}
I_{ij}^{(1AB)} &= A_{im} B_{mj} - A_{ij} B_{jj} + \sum_{l=m+1}^i \frac{1}{2} (A_{il} + A_{i(l-1)}) (B_{lj} - B_{(l-1)j}) + \\
&- \sum_{l=j+1}^m \frac{1}{2} (B_{lj} + B_{(l-1)j}) (A_{il} - A_{i(l-1)}) = \tag{87}
\end{aligned}$$

$$= A_{im} B_{mj} - A_{ij} B_{jj} + \sum_{l=m+1}^i dA_{il}^{(v)} (B_{lj} - B_{(l-1)j}) - \sum_{l=j+1}^m (A_{il} - A_{i(l-1)}) dB_{lj}^{(h)};$$

$$\begin{aligned}
I_{ij}^{(2ABC)} &= \sum_{l=j+1}^i \frac{1}{2} (A_{il} + A_{i(l-1)}) (B_{lj} - B_{(l-1)j}) \frac{1}{2} (C_{lj} + C_{(l-1)j}) = \\
&= \sum_{l=m+1}^i dA_{il}^{(h)} (B_{il} - B_{i(l-1)}) \frac{1}{2} (C_{lj} + C_{(l-1)j}) + \sum_{l=j+1}^m \frac{1}{2} (A_{il} + A_{i(l-1)}) (B_{il} - B_{i(l-1)}) dC_{lj}^{(v)}; \tag{88}
\end{aligned}$$

$$I_{ij}^{(3AB)} = A_{ij} B_{jj} - A_{i0} B_{j0} - \sum_{l=1}^j (A_{il} - A_{i(l-1)}) dB_{jl}^{(v)}; \tag{89}$$

$$I_{ij}^{(4ABC)} = \sum_{l=1}^j \frac{1}{2} (A_{il} + A_{i(l-1)}) (B_{il} - B_{i(l-1)}) dC_{jl}^{(v)}; \tag{90}$$

$$I_i^{(5AB)} = \sum_{l=1}^i dA_{il}^{(v)} (B_l - B_{l-1}); \tag{91}$$

$$I_i^{(6ABC)} = \sum_{l=1}^i \frac{1}{2} (A_{il} + A_{i(l-1)}) (B_{il} - B_{i(l-1)}) dC_l, \tag{92}$$

where the superscript (v) and (h) represent the vertical (t') and horizontal (t) derivatives in the discretized times, see Fig. 11 for an intuitive understanding.

We also discretized the derivative using the last two time steps

$$\frac{d}{dt} g(t) = \frac{3}{2\delta t} g(t) - \frac{2}{\delta t} g(t - \delta t) + \frac{1}{2\delta t} g(t - 2\delta t) + O(\delta t^3). \tag{93}$$

Given the time indices i and j , we will define and evaluate the following quantities

$$\begin{aligned}
&\{C_{ij}, Q_{ij}, M2_{ij}, N2_{ij}, Mp_{ij}, Np_{ij}, Cbar_i, P2_i, Pp_i, mu_i\} = \\
&= \{C(t_i, t_j), Q(t_i, t_j), f_2'(C(t_i, t_j)), f_2''(C(t_i, t_j)), f_p'(C(t_i, t_j)), f_p''(C(t_i, t_j)), \bar{C}(t_i), f_2'(\bar{C}(t_i)), f_p'(\bar{C}(t_i)), \mu(t_i)\}
\end{aligned}$$

plus the respective vertical and horizontal derivatives.

Calling $D_i = \frac{3}{2dt} + \mu'_i - \frac{1}{\Delta_2} M2_{ii} - \frac{2}{p\Delta_p} Mp_{ii}$, the original dynamical equations are integrated as follow

$$\begin{aligned}
C_{ij}D_i &= \frac{2}{dt}C_{(i-1)j} - \frac{1}{2dt}C_{(i-2)j} + \text{Cbar}_j(r_2P2_i + r_pPp_i) + \\
&+ \frac{1}{\Delta_2} \left(-\dot{I}_{ij}^{(1f'_2C)} + I_{ij}^{(2f'_2QC)} + I_{ij}^{(3f'_2Q)} + I_{ij}^{(4f''_2QC)} - M2_{i0}C_{j0} \right) + \\
&+ \frac{2}{p\Delta_p} \left(-\dot{I}_{ij}^{(1f'_pC)} + I_{ij}^{(2f''_pQC)} + I_{ij}^{(3f'_pQ)} + I_{ij}^{(4f''_pQC)} - Mp_{i0}C_{j0} \right), \tag{94}
\end{aligned}$$

$$\begin{aligned}
Q_{ij}D_i &= \mu'_i - 1 + \frac{2}{dt}Q_{(i-1)j} - \frac{1}{2dt}Q_{(i-2)j} + \text{Cbar}_j(r_2P2_i + r_pPp_i) + \\
&+ \frac{1}{\Delta_2} \left(-\dot{I}_{ij}^{(1f'_2Q)} + I_{ij}^{(2f'_2Q(Q-1))} - I_{ij}^{(3f'_2Q)} - I_{ij}^{(4f''_2QC)} - M2_{i0}C_{i0} \right) + \\
&+ \frac{2}{p\Delta_p} \left(-\dot{I}_{ij}^{(1f'_pQ)} + I_{ij}^{(2f''_pQ(Q-1))} - I_{ij}^{(3f'_pQ)} - I_{ij}^{(4f''_pQC)} - Mp_{i0}C_{i0} \right), \tag{95}
\end{aligned}$$

$$\begin{aligned}
\text{Cbar}_iD_i &= \frac{2}{dt}\text{Cbar}_{i-1} - \frac{1}{2dt}\text{Cbar}_{i-2} + r_2P2_i + r_pPp_i + \\
&+ \frac{1}{\Delta_2} \left(-\dot{I}_i^{(5f'_2\text{Cbar})} + I_i^{(6f''_2Q\text{Cbar})} - M2_{i0}\text{Cbar}_0 \right) + \\
&+ \frac{2}{p\Delta_p} \left(-\dot{I}_i^{(5f'_p\text{Cbar})} + I_i^{(6f''_pQ\text{Cbar})} - Mp_{i0}\text{Cbar}_0 \right). \tag{96}
\end{aligned}$$

In the systems we used \dot{I} to characterize the integrals where we remove from the sum the term present in the left-hand side (e.g. for C_{ij} eq. 94). Using Simpson's integration formula we define the increments

$$\begin{aligned}
\Delta_{il} &= \frac{1}{12}(Q_{il} - Q_{i(l-1)})\{W_2^2[-(M2_{i(l+1)} + N2_{i(l+1)}C_{i(l+1)}) + 8(M2_{il} + N2_{il}C_{il}) + 5(M2_{i(l-1)} + N2_{i(l-1)}C_{i(l-1)})] + \\
&+ W_p^2[-(Mp_{i(l+1)} + Np_{i(l+1)}C_{i(l+1)}) + 8(Mp_{il} + Np_{il}C_{il}) + 5(Mp_{i(l-1)} + Np_{i(l-1)}C_{i(l-1)})]\}
\end{aligned}$$

and we determine μ' as

$$\mu' = 1 + r_2P2_i + r_pPp_i + \delta\mu' + \sum_{l=1}^{i-N_t/4} \Delta_{il} - (W_2^2M2_{i0} + W_p^2Mp_{i0})C_{i0}, \tag{97}$$

with $\delta\mu'$ initially set to 0.

Algorithm: Here we describe the main steps of the algorithm, pictorially represented Fig. 11.

Discretize the time (t, t') in N_t (even) intervals, the results shown use $N_t = 1024$.

1. **Initialization.** Fill the first $N_t/2$ times by linear propagation of the value obtained from the perturbative analysis

$$C_{ij} = 1 - (i - j)dt; \tag{98}$$

$$Q_{ij} = 0; \tag{99}$$

$$\text{Cbar}_i = \bar{C}_0 + \left\{ \left[r_2f'_2(\bar{C}_0) + r_pf'_p(\bar{C}_0) \right] (1 + (\bar{C}_0)^2) - \bar{C}_0 \right\} dt; \tag{100}$$

$$M2_{ij} = f'_2(C_{ij}); \tag{101}$$

$$N2_{ij} = f''_2(C_{ij}); \tag{102}$$

$$Mp_{ij} = f'_p(C_{ij}); \tag{103}$$

$$Np_{ij} = f''_p(C_{ij}). \tag{104}$$

2. **Fill the grid (small τ).** Continue to propagate the values for small time differences $\tau = t - t' \ll 1$. In terms of the algorithm it means that we have some elements of the grid, N_c of them, close to the diagonal that will be updated by linear propagation because the approximation of small τ is still valid. In our simulation the first Δt is of the order 10^{-7} and $N_c = 2$.
3. **Fill the grid (larger τ).** The rest of the values will be copied from the previous t ($A_{t+\Delta t, t'} = A_{t, t'}$). These values are the initial guess for solving the self-consistent equations (94-95-96) and (97), in this procedure the derivatives are updated using the 2nd order discretization.
4. **Half the grid and expand.** The grid is decimated which means that each observable is contracted $A_{i,j} \leftarrow A_{2i, 2j}$ and the derivate are updated as follows $dA_{i,j}^{(h)} \leftarrow \frac{1}{2}(dA_{2i, 2j}^{(h)} + dA_{2i-1, 2j}^{(h)})$, $dA_{i,j}^{(v)} \leftarrow \frac{1}{2}(dA_{2i, 2j}^{(v)} + dA_{2i, 2j-1}^{(v)})$. The new time step is now: $\Delta t \leftarrow 2\Delta t$.
5. **Start over from step 2.**

D.3 Numerical checks on the dynamical algorithm

The dynamic-grid algorithm has been checked in a variety of ways.

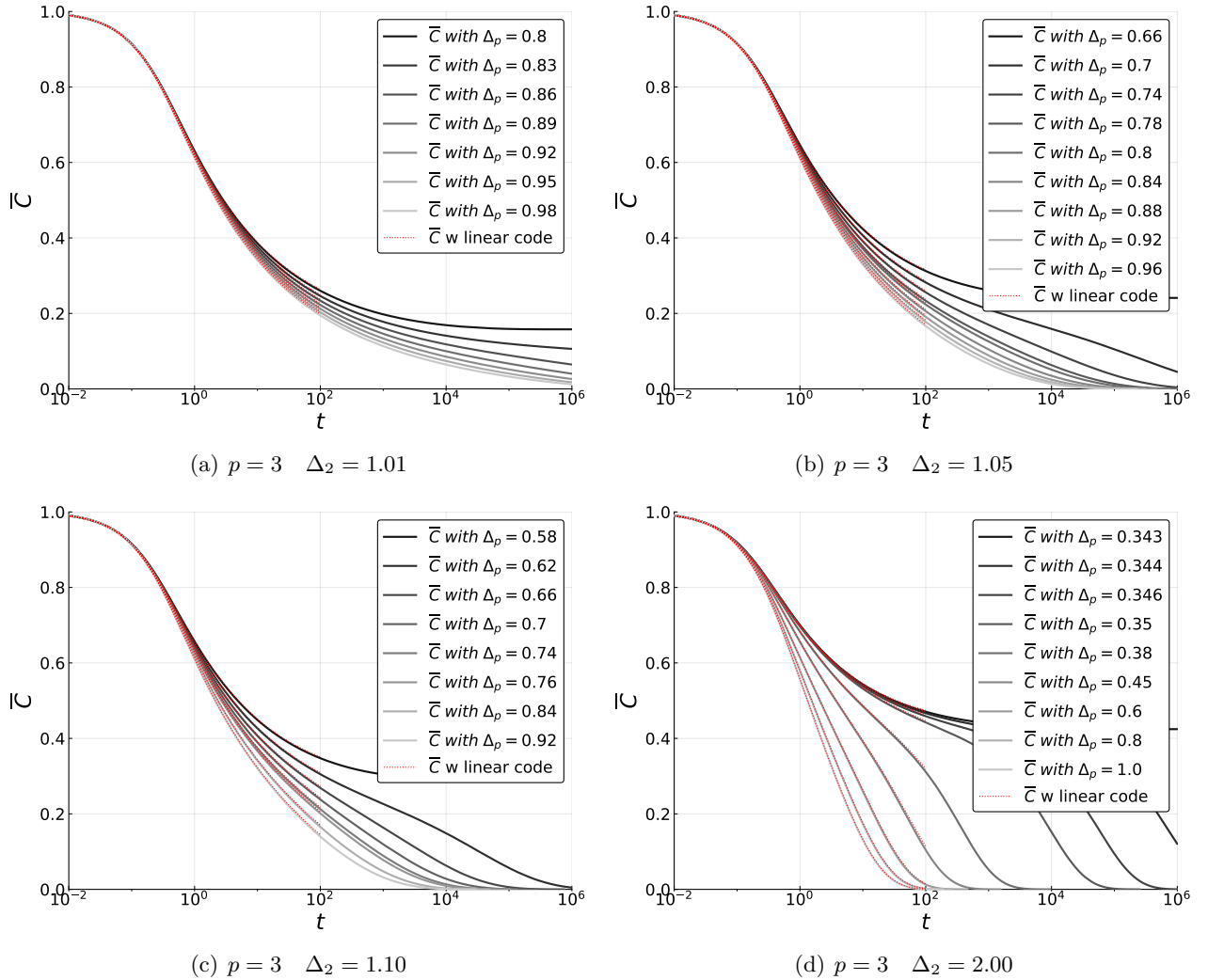


Figure 12: Evolution of the correlation with the signal at fixed Δ_2 for different Δ_p . The dotted red line overlapping with other lines, is the same quantity evaluated using the fixed grid algorithm up to time 100. We have started the LSE from an informative initial condition.

Cross-checking using the fixed-grid algorithm. For short times the dynamical equations were solved using the fixed-grid algorithm and compared with the outcome of the dynamic-grid algorithm, obtaining the same results, see Fig. 12. In the figure we used the fixed-grid with $t_{\max} = 100$ and the $\Delta t = 6.25 * 10^{-3}$.

Same magnetization in the easy region. In the impossible and easy regions, the overlap with the signal of both AMP and dynamic-grid integration, converges to the same value. In Fig. 13 we show the overlap obtained with AMP, black dashed line, and the overlap achieved by the integration scheme at a given time. We can see that the overlap with the signal as obtained solving the LSE equations converges to the same value of the fixed point of AMP. Given a fixed Δ_2 we can observe that the time to convergence increase very rapidly as we decrease Δ_p . We fitted this increase of the *relaxation time* to get the boundary of the Langevin hard region.

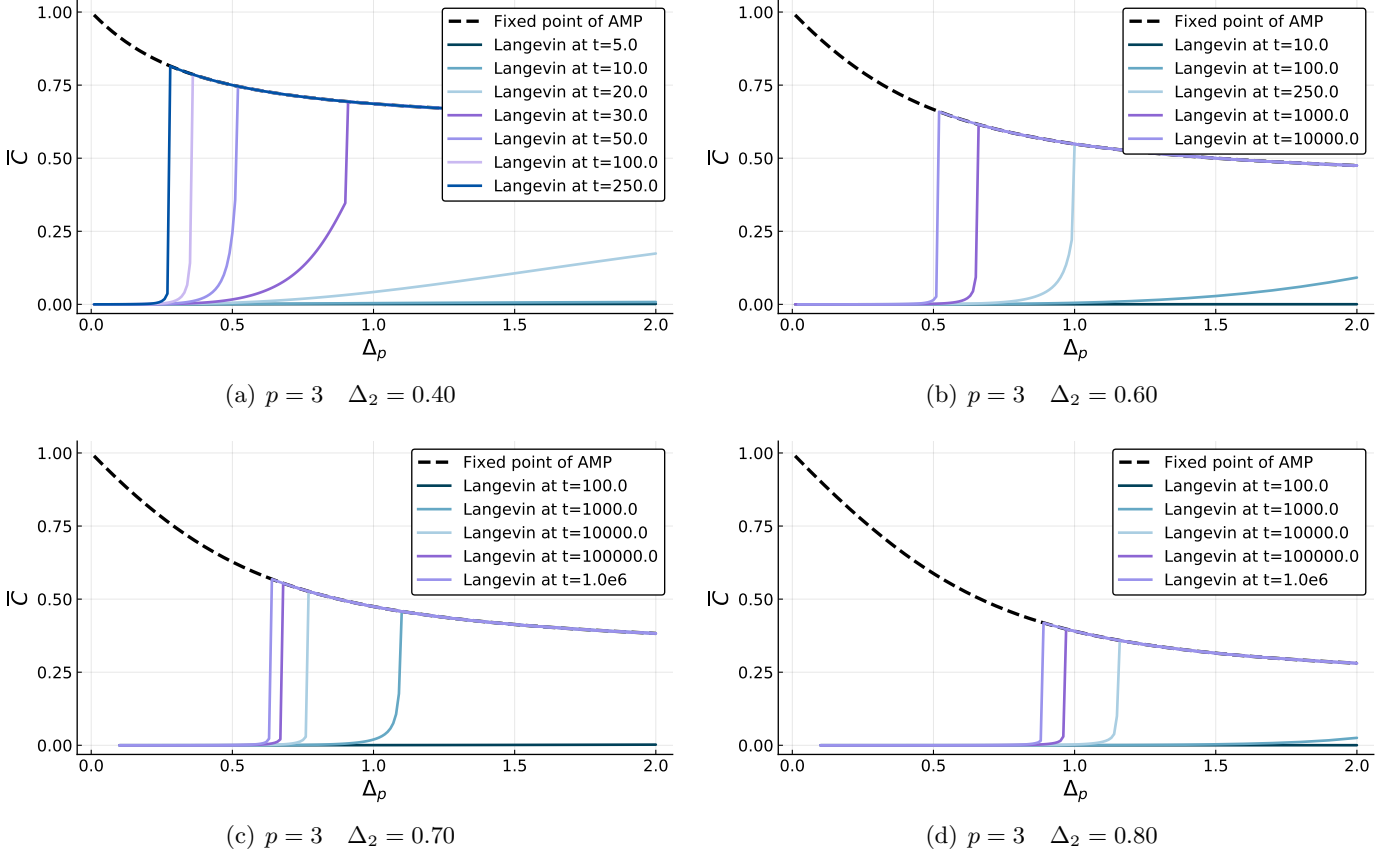


Figure 13: Correlation with the signal of AMP (dotted lines) and Langevin (solid lines) at k th iteration and t time respectively. The black dashed line is the asymptotic value predicted with AMP. In the easy region, provided enough running time, Langevin dynamics finds the same alignment as AMP. The figures show qualitatively the same behaviour for different values of Δ_2 .

Dynamical transition. The dynamical transition where the finite magnetization fixed point disappears can be regarded as a clustering or dynamical glass transition. Indeed coming from the impossible phase, going towards the hard phase, at the dynamical transition the free energy landscape changes and the unique ergodic paramagnetic minimum of the impossible phase gets clustered into an exponential number of metastable glassy states (see Sec. E). Correspondingly the relaxation time of the Langevin algorithm diverges. Fitting this divergence with a power law we obtain an alternative estimation of the dynamical line. In the right panel of Fig. 14 we plot with yellow points the dynamical transition line as extracted from the fit of the relaxation time of the Langevin algorithm extracted coming from the impossible phase and entering in the hard phase.

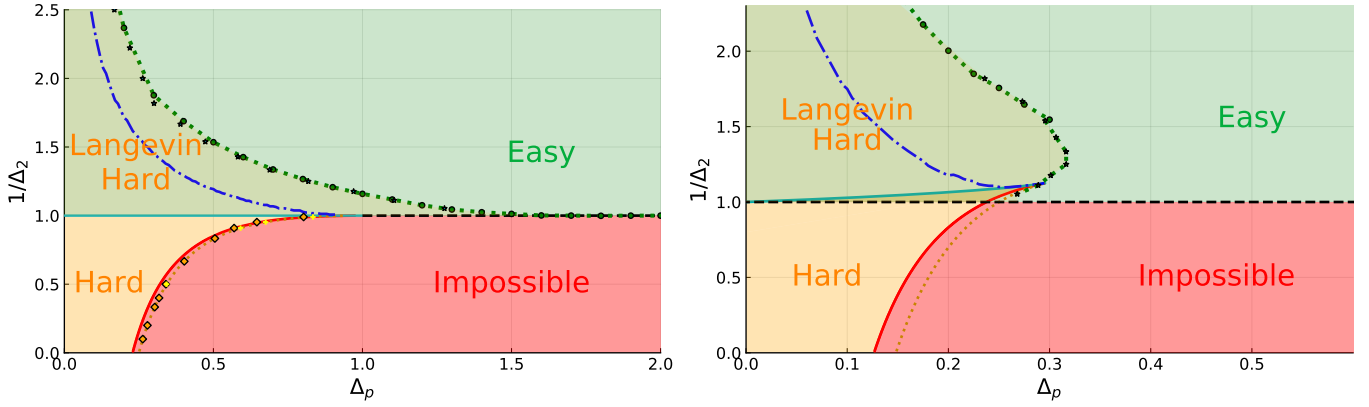


Figure 14: On the left: phase diagram of the spiked matrix-tensor model for $p = 3$ as presented in the left panel of Fig. 7 with the additional boundary of the Langevin hard phase (green circles and stars). The boundary of the Langevin hard phase has been obtained by fitting the relaxation time at fixed Δ_p and increasing Δ_2 (circles) and conversely, at fixed Δ_2 and decreasing Δ_p (stars). The blue dashed line marks a region above which we do not observe anymore a stable positive 1RSB complexity. Finally we plot with orange and yellow dots the dynamical transition line as extracted from the relaxation time of the Langevin algorithm coming from respectively the hard and impossible phase. On the right: phase diagram of the spiked matrix-tensor model for $p = 4$ as presented in the right panel of Fig. 7 with the additional Langevin hard phase boundary. Also in this case we observe that Langevin hard phase extends in the AMP easy phase. Interestingly the Langevin hard phase here folds and closes near the tri-critical point. The blue dashed line marks a region above which we do not observe anymore a stable positive 1RSB complexity.

D.4 Extrapolation procedure

In order to determine the Langevin Hard region, given a fixed value of Δ_p (Δ_2), we measure the relaxation time that it takes to relax to equilibrium. On approaching the Langevin hard region, this relaxation time increases and we extrapolate the divergence to obtain the critical Δ_p^* (Δ_2^* respectively) where the relaxation time diverges. The extrapolation is done assuming a power law divergence. Fig. 14 shows the results of this procedure for the cases $2 + 3$ and $2 + 4$.

Numerical checks on the extrapolation procedure. To test the quality of the fits we use a similar numerical procedure to locate the spinodal of the informative solution, which is given by the points where the informative solution ceases to exist. This spinodal must be the same for both the AMP and the Langevin algorithm [6].

Since we aim at studying the spinodal of the informative solution, we initialize the LSE with $\bar{C}_0 = 1$ and let it relax, measuring the time it takes to equilibrate at the value of \bar{C} given by the informative fixed point of AMP. We do this fixing Δ_2 and changing Δ_p . As we approach the critical $\Delta_{p,dyn}$, the relaxation time will diverge and we can fit this divergence with a power law. The dynamic threshold extracted in this way is finally compared with the one obtained from AMP. In Fig. 15 we show how this scheme has been applied for $\Delta_2 \in \{1.01, 1.05, 1.10, 2.00\}$. As we get closer to the critical line $\Delta_2 = 1$ the relaxation time increases (and the height of the plateau decreases), making the fit harder. All in all, we observe a very good agreement between the points found with these extrapolation procedures and the prediction obtained with AMP, as shown in Fig. 14.

D.5 Annealing protocol

In this section we show that using specific protocols we are able to enter in the Langevin hard region. A generic annealing scheme would lower the noises of both the channels simultaneously, which, as we tested and discussed

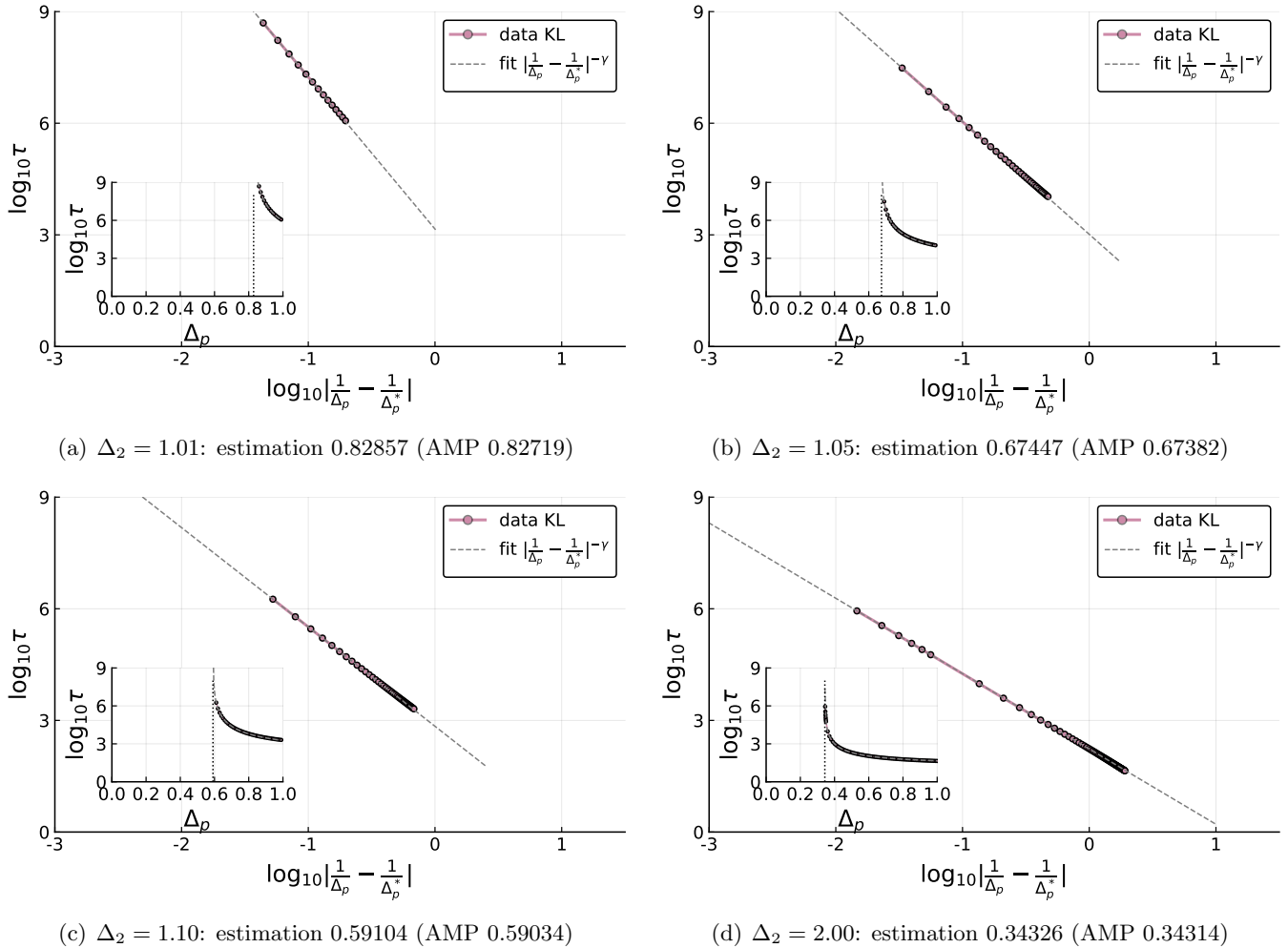


Figure 15: Relaxation time obtained from the LSE starting from an informative initial condition $\bar{C}_0 = 1$. The four cases refers to the 2 + 3 model and are fitted with a power law and the relaxation time appears to diverge very close to point predicted by AMP (the AMP prediction is given in the captions).

in the main text, will not be able to avoid the Langevin hard region. Instead we can use the following protocol

$$\begin{aligned}
 T_2 &\equiv 1, \\
 T_p &\equiv T_p(t) = 1 + \frac{C}{\Delta_p} e^{-\frac{t}{\tau_{\text{ann}}}}.
 \end{aligned}
 \tag{105}$$

The constant C allows to select at the initial time the desired effective $\Delta_{p,\text{eff}} = \Delta_p + C e^{-\frac{t}{\tau_{\text{ann}}}}$ far from (and much larger than) the original one. Instead τ_{ann} chooses the speed of the annealing protocol. Fig. 16 shows that using this protocol we are able to enter in Langevin Hard region even with Δ_2 close to the AMP threshold. To this purpose we initiated the effective $\Delta_{p,\text{eff}}$ close to 100 (i.e. $C = 100$), very far from the Langevin hard region, and we used different speeds for the annealing of Δ_p (different colors in the figures). In the figures we can observe that approaching the $\Delta_2 = 1$ we need slower and slower protocols (larger and larger τ_{ann}). The reason for this behavior is due to the fact that approaching $\Delta_2 = 1$ with $\Delta_p = 100$ a longer time is required to gain a non trivial overlap with the solution. Evidence of this growing timescale at $\Delta_p = 100$ is given in Fig. 17 where we show the relaxation time for magnetizing to the solution varying Δ_2 . In particular, we can observe that at $\Delta_2 = 0.70$ the relaxation time is of the order of 100 time units.

For the protocol to be successful it is therefore crucial that the annealing time τ_{ann} is large enough to give

the possibility of magnetizing the solution before $\Delta_{p,\text{eff}}$ has significantly decreased towards Δ_p . According to this analysis, it is not surprising that in Fig. 16 for $\Delta_2 = 0.90$ the proposed protocol seems not to be successful. For this value of the parameter Δ_2 , the time to find a solution even with $\Delta_p = 100$ should be larger than 1000 time units, which is much larger than the used τ_{ann} and anyway out of the time window of our numerical solution. However, with an annealing time large enough it would be in principle possible to recover exactly the same boundaries of the AMP easy region.

E Glassy nature of the Langevin hard phase: the replica approach

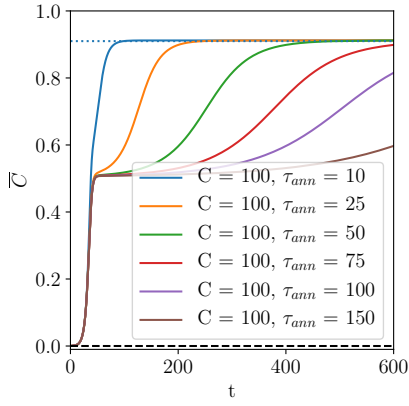
In this section we study the landscape of the spiked matrix-tensor problem following the approach of [25]. We underline here that we are interested in studying the free energy landscape problem rather than the energy landscape since the former is the relevant quantity for finite temperatures ($\beta = 1$ in our case, as discussed in Sec. A). The results of [25] suggest that the AMP-hard phase and part of the AMP-easy phase are glassy. Therefore we could expect that low magnetization glassy states trap the Langevin algorithm and forbid the relaxation to the equilibrium configurations that surrounds the signal. This may happen also in a region where AMP instead is perfectly fine in producing configurations strongly correlated with the signal. In order to check this hypothesis we compute the logarithm of the number of glassy states, called the complexity by using the replica method [64, 25]. The goal of this analysis is to trace an additional line in the phase diagram that delimits the region where stable one step replica symmetry breaking (1RSB) metastable states exist. We conjecture that this provides a physical lower bound to the Langevin hard phase in the $(\Delta_p, 1/\Delta_2)$ phase diagram.

E.1 Computation of the complexity through the replica method

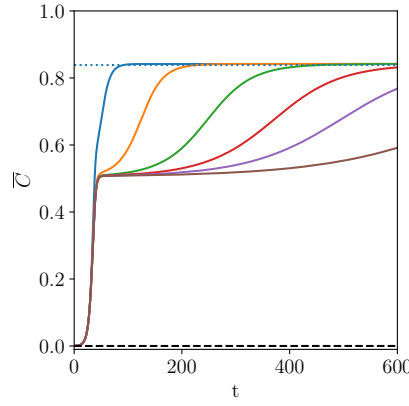
The replica trick is based on the simple identity: $\mathbb{E} \log x = \lim_{n \rightarrow 0} \frac{\partial}{\partial n} \mathbb{E} x^n$. Using this observation we can compute the expected value of the free energy, $\Phi = -(\log Z)/N$, averaging the Z^n and taking the limit $n \rightarrow 0$. This is in general as difficult as the initial problem, however, if we consider only integer n and extrapolate to 0, the computation becomes much less involved due to the fact that for integer n the average $\mathbb{E} x^n$ can be sometimes performed analytically. Indeed in this case the replicated partition function Z^n can be regarded as the partition function of n identical uncoupled systems or replicas. Averaging over the disorder we obtain a clean system of interacting replicas. The Hamiltonian of this system displays an emerging *replica symmetry* since it is left unchanged by a permutation of replicas. This symmetry can be spontaneously broken in certain disordered models where frustration is sufficiently strong [37].

In mean field models characterized by fully connected factor graphs, the resulting Hamiltonian of interacting replicas depends on the configuration of the system only through a simple order parameter, the overlap \tilde{Q} between them, which is a $n \times n$ matrix that describes the similarities of the configurations of different replicas in phase space. Furthermore the Hamiltonian is proportional to N which means that in the thermodynamic limit $N \rightarrow \infty$, the model can be solved using the saddle point method. In this case one needs to consider a simple ansatz for the saddle point structure of the matrix \tilde{Q} that allows to take the analytic continuation for $n \rightarrow 0$. The solution to this problem comes from spin glass theory and general details can be found in [37]. The saddle point solutions for \tilde{Q} can be classified according to the replica symmetry breaking level going from the replica symmetric solution where replica symmetry is not spontaneously broken to various degree of spontaneous replica symmetry breaking (including full-replica symmetry breaking). Here we will not review this subject but the interested reader can find details in [37]. The model we are analyzing can be studied in full generality at any degree of RSB (see for example [34, 35, 43] where the same models have been studied in absence of a signal). However here we will limit ourselves to consider saddle point solutions up to a 1RSB level.

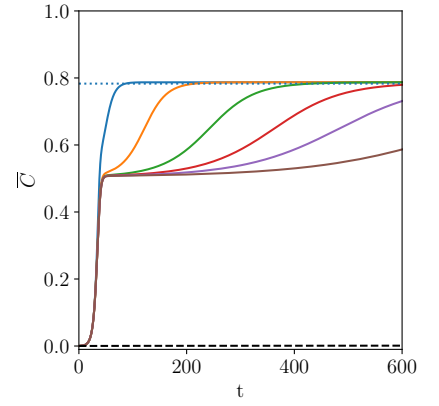
The complexity of the landscape can be directly related to replica symmetry breaking. A replica symmetric solution implies an ergodic free energy landscape characterized by a single pure state. When replica symmetry is broken instead, a large number of pure states arises and the phase space gets clustered in a hierarchical way [37]. Making a 1RSB approximation means to look for a situation in which the hierarchical organization



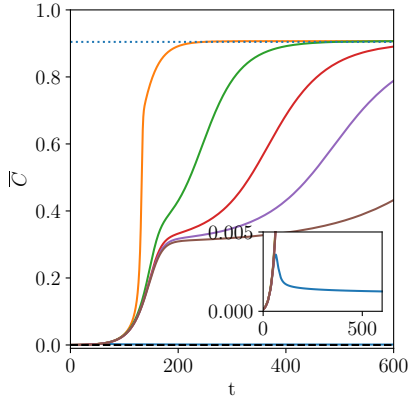
(a) $\Delta_2 = 0.50, \Delta_p = 0.10$



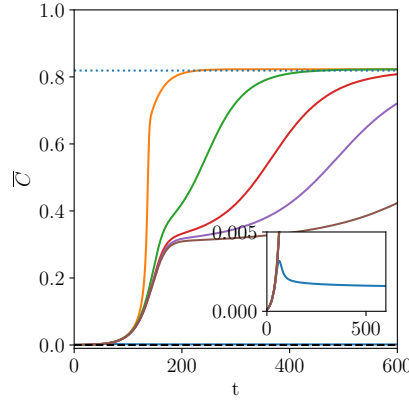
(b) $\Delta_2 = 0.50, \Delta_p = 0.20$



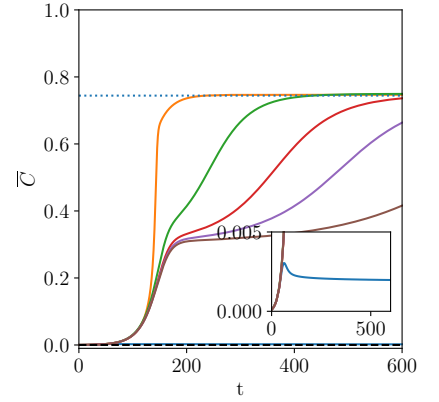
(c) $\Delta_2 = 0.50, \Delta_p = 0.30$



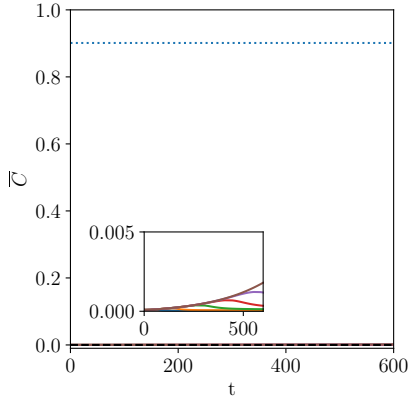
(d) $\Delta_2 = 0.70, \Delta_p = 0.10$



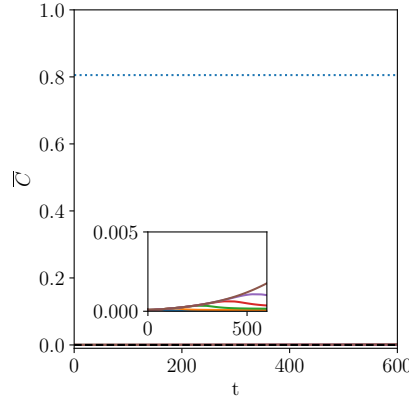
(e) $\Delta_2 = 0.70, \Delta_p = 0.20$



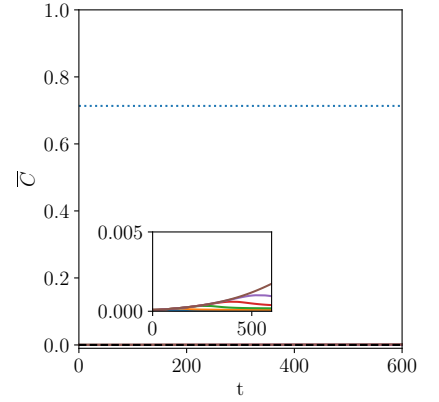
(f) $\Delta_2 = 0.70, \Delta_p = 0.30$



(g) $\Delta_2 = 0.90, \Delta_p = 0.10$

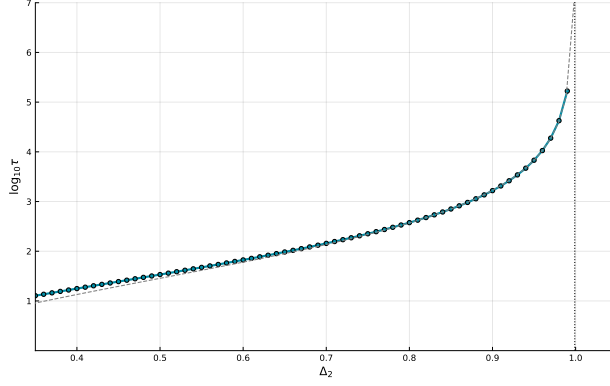


(h) $\Delta_2 = 0.90, \Delta_p = 0.20$



(i) $\Delta_2 = 0.90, \Delta_p = 0.30$

Figure 16: The figures show the correlation with the signal in time obtained using different annealing protocols, whose details are reported in the legend of the first figure. All the protocols have $C = 100$ which means that all the dynamics start with close effective $\Delta_p \sim 100, T_p(0)\Delta_p \simeq 100$. What changes among the different lines is the relaxation speed, from the fastest, drawn in blue, to the slowest, drawn in brown. They are compared with the asymptotic value of AMP, dotted line, and the Langevin dynamics without smart annealing, dashed line.



(a)

Figure 17: Relaxation times of Langevin dynamics at $\Delta_p = 100$ without using protocols.

contains just one level: the phase space gets clustered into an exponential number of pure states with no further internal structure.

If we assume a 1RSB glassy landscape, we can compute the complexity of metastable states using a recipe due to Monasson [64] (see also [65] for a pedagogical introduction). The argument goes as follows.

Let us consider system with x real replicas infinitesimally coupled. If the free energy landscape is clustered into an exponential number of metastable states, the replicated partition function, namely the partition function of the system of x real replicas, can be written as

$$Z^x \simeq e^{N[\Sigma(f^*) - x\beta f^*]}$$

where f^* is the internal free energy of the dominant metastable states that is determined by the saddle point condition $\frac{d\Sigma}{df}(f^*) = \beta x$ and β the inverse temperature. Note that since we are interested in the Bayes optimal case, this corresponds to set $\beta = 1$. The function $\Sigma(f)$ is the complexity of metastable states having internal entropy f . Therefore, using the free parameter x we can reconstruct the form of $\Sigma(f)$ from the replicated free energy. In order to compute the replicated free energy we need to apply the replica trick on the replicated system, $\overline{\log Z^x} = \lim_{n \rightarrow 0} \frac{\partial}{\partial n} (\overline{Z^x})^n$. Calling the replicated free energy $\Phi = -\frac{1}{N} \overline{\log Z^x}$, we get the complexity as $\Sigma = x \frac{\partial \Phi}{\partial x} - \Phi$.

We can now specify the computation to our case where the partition function is the normalization of the posterior measure. With simple manipulations of the equations [61], the partition function can be expressed as the integral over the overlap matrix

$$\overline{(Z^x)^n} = \overline{Z_x^n} \propto \int \prod_{ab} dQ_{ab} e^{NnxS(Q)} \simeq \limsup_Q e^{NnxS(Q)}; \quad (106)$$

where the overlap Q is a $(nx + 1) \times (nx + 1)$ matrix

$$Q = \left(\begin{array}{c|c} 1 & m \cdots m \\ \hline m & \tilde{Q} \\ \vdots & \\ m & \end{array} \right)$$

that contains a special row and column that encodes the overlap between different replicas with the signal and therefore the corresponding overlap is the *magnetization* m .

The 1RSB structure for the matrix \tilde{Q} can be obtained by defining the following $nx \times nx$ matrices: the identity matrix $\mathbb{1}_{ij} = \delta_{ij}$, the full matrix $\mathcal{J}_{nx,ij}^{(0)} = 1$, and $\mathcal{J}_{nx}^{(1)} = \text{diag}(J_x^{(0)}, \dots, J_x^{(0)})$ a block diagonal matrix

where the diagonal blocks $J_x^{(0)}$ have size $x \times x$ and are matrices full of 1. In this case the 1RSB ansatz for \tilde{Q} reads

$$\tilde{Q} = (1 - q_M)\mathbb{1}_{nx} + (q_M - q_m)\mathbb{J}_{nx}^{(1)} + q_m\mathbb{J}_{nx}^{(0)}.$$

Using this ansatz we can compute $S(Q)$ that is given by

$$\begin{aligned} S(Q) &= \frac{1}{nx} \left[\frac{1}{2} \log \det Q + \frac{1}{2p\Delta_p} \sum_{ab} Q_{ab}^p + \frac{1}{4\Delta_2} \sum_{ab} Q_{ab}^2 \right] = \\ &= \frac{1}{2} \log(1 - q_M) + \frac{1}{2x} \log \frac{1 - q_M + x(q_M - q_m)}{1 - q_M} + \frac{1}{2} \frac{q_m - m^2}{1 - q_M + x(q_M - q_m)} + \\ &+ \frac{1}{2p\Delta_p} (1 - q_M^p + x(q_M^p - q_m^p) + 2m^p) + \frac{1}{4\Delta_2} (1 - q_M^2 + x(q_M^2 - q_m^2) + 2m^2). \end{aligned} \quad (107)$$

We can observe that starting from this expression we can derive the RS free energy, eq. (36), taking $q_M = q_m$ or equivalently in the limit $x \rightarrow 1$. From eq. (107) we obtain the saddle point equations

$$\begin{aligned} 0 &= 2 \frac{\partial S}{\partial q_M} = (x - 1) \left[\frac{1}{x} \left(\frac{1}{1 - q_M + x(q_M - q_m)} - \frac{1}{1 - q_M} \right) - \frac{q_m - m^2}{[1 - q_M + x(q_M - q_m)]^2} + \frac{q_M^{p-1}}{\Delta_p} + \frac{q_M}{\Delta_2} \right]; \\ 0 &= 2 \frac{\partial S}{\partial q_m} = x \left[\frac{q_m - m^2}{[1 - q_M + x(q_M - q_m)]^2} - \frac{q_m^{p-1}}{\Delta_p} - \frac{q_m}{\Delta_2} \right]; \\ 0 &= \frac{\partial S}{\partial m} = \frac{-m}{1 - q_M + x(q_M - q_m)} + \frac{m^{p-1}}{\Delta_p} + \frac{m}{\Delta_2}. \end{aligned} \quad (108)$$

The low magnetization solution to these equations gives the complexity of the metastable branch of the posterior measure which is given by

$$\begin{aligned} -\Sigma(x; Q^*) &= -\log \frac{1 - q_M + x(q_M - q_m)}{1 - q_M} + x \left(\frac{q_m^{p-1}}{\Delta_p} + \frac{q_m}{\Delta_2} \right) - x^2 \frac{(q_m - m^2)(q_M - q_m)}{[1 - q_M + x(q_M - q_m)]^2} + x^2 \frac{q_M^p - q_m^p}{p\Delta_p} + \\ &+ x^2 \frac{q_M^2 - q_m^2}{2\Delta_2}. \end{aligned} \quad (109)$$

The free parameter x allows us to tune the free energy of the states of which we compute the complexity. Thus we can characterize the part of the phase diagram where an exponential number of states is present.

To complete the 1RSB analysis we must compute the stability of the 1RSB saddle point solution for Q . In particular it is important to compute the so called *replicon* eigenvalues that characterizes the instability of this solution towards further replica symmetry breaking. Following [66, 43] we have two replicon eigenvalues given by

$$\lambda_I = 1 - (1 - q_M + x(q_M - q_m))^2 \left[(p - 1) \frac{q_m^{p-2}}{\Delta_p} + \frac{1}{\Delta_2} \right], \quad (110)$$

$$\lambda_{II} = 1 - (1 - q_M)^2 \left[(p - 1) \frac{q_M^{p-2}}{\Delta_p} + \frac{1}{\Delta_2} \right]. \quad (111)$$

We can analyze what happens to the landscape when we fix $\Delta_p < 1$ and we start from a large value of $\Delta_2 < \Delta_{2,dyn}(\Delta_p)$ and we decrease Δ_2 . In this case for sufficiently high Δ_2 and large enough Δ_p the system is in a paramagnetic phase and no glassy states are present. At the dynamical transition line instead we find

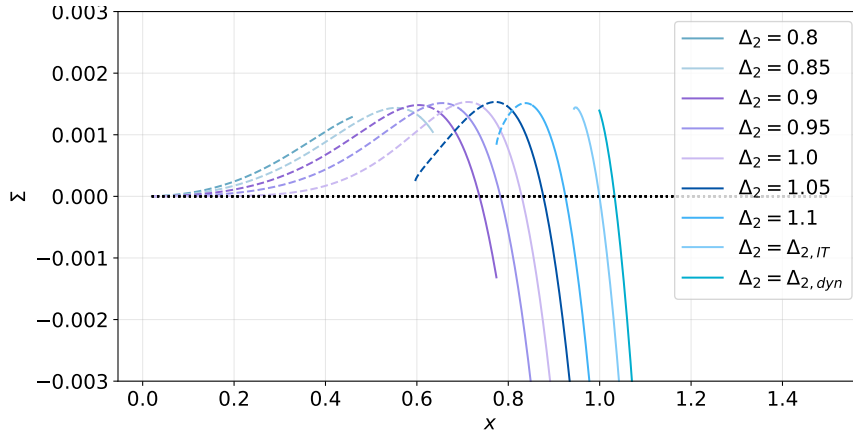
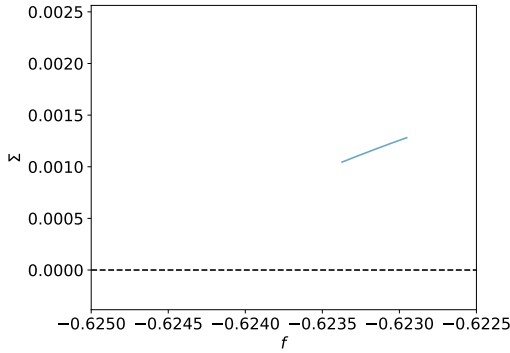
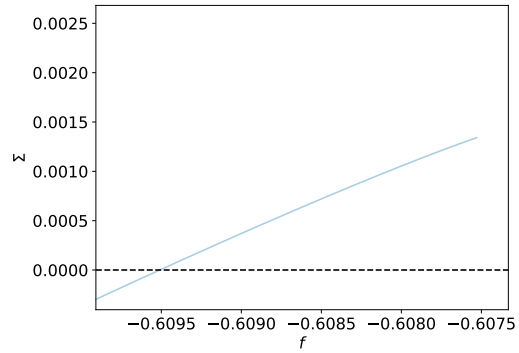


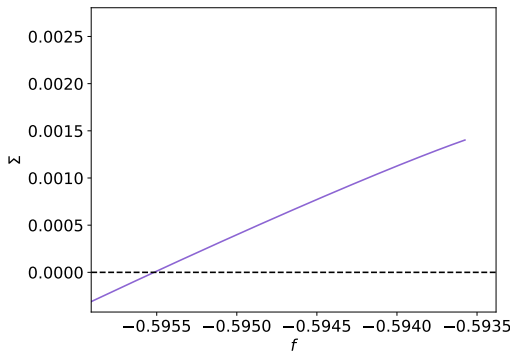
Figure 18: Complexity with as a function of the Parisi parameter x for $p = 3$ on the line $\Delta_p = 0.5$. The solid line characterizes the stable part of the complexity while the dashed line the unstable one.



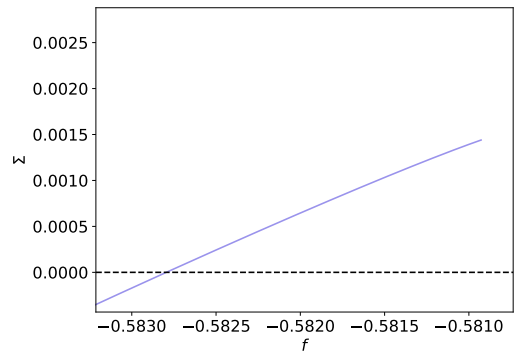
(a) $p = 3$; $\Delta_p = 0.50$; $\Delta_2 = 0.85$



(b) $p = 3$; $\Delta_p = 0.50$; $\Delta_2 = 0.90$



(c) $p = 3$; $\Delta_p = 0.50$; $\Delta_2 = 0.95$



(d) $p = 3$; $\Delta_p = 0.50$; $\Delta_2 = 1.00$

Figure 19: The stable part of the 1RSB complexity as a function of the free energy for $p = 3$ and $\Delta_p = 0.5$.

a positive complexity as plotted in Fig. 18. At this point the equilibrium states that dominate the posterior measure are the so called *threshold* states for which the complexity is maximal. For those states the eigenvalue $\lambda_{II} = 0$ which confirms that these states are marginally stable [27]. Decreasing Δ_2 one crosses the information theoretic phase transition where the relevant metastable states that dominate the posterior measure have zero complexity. This corresponds to a freezing/condensation/Kauzmann transition. Below the information theoretic phase transition the thermodynamics of the posterior measure is dominated by the state containing the signal. However one can neglect the high magnetization solution of the 1RSB equations to get the properties of the metastable branch and computing the complexity of states that have zero overlap with the signal. The complexity curves as a function of the Parisi parameter x for decreasing values of Δ_2 are plotted in Fig. 18 for fixed $\Delta_p = 0.5$ and several Δ_2 . The curves contain a stable 1RSB part and an unstable one where λ_{II} is negative. The 1RSB line shown in Figs. 14 is obtained by looking at when the states with positive complexity and $\lambda_{II} = 0$ disappear. This means that it gives the point where the 1RSB marginally stable states disappear and therefore it is expected to be a lower bound for the disappearance of glassiness in the phase diagram. The important outcome of this analysis is that for $\Delta_2 < 1$ but not sufficiently small, namely in part of the AMP-easy phase, the replica analysis predicts the existence of 1RSB marginally stable glassy states that may trap the Langevin algorithm from relaxing towards the signal [25] and therefore supports the existence of the Langevin hard phase.

Finally in Fig. 19 we plot the complexity as a function of the internal free energy of the metastable states for some values of Δ_2 and Δ_p .

E.2 Breakdown of the fluctuation-dissipation theorem in the Langevin hard phase

When the Langevin algorithm is able to reach equilibrium, being it the signal or the paramagnetic state, it should satisfy the Fluctuation-Dissipation Theorem (FDT) according to which the response function is related to the correlation function through $R(t, t') = -\frac{\partial C(t, t')}{\partial t}$. Furthermore, time translational invariance (TTI) should arise implying that both correlation and response functions should be functions of only the time difference meaning that $R(t, t') = R(t - t')$ and $C(t, t') = C(t - t')$ ⁵. When the dynamics is run in the glass phase, metastable states may forbid equilibration. In this case time translational invariance is never reached at long times⁶ and the dynamics displays aging violating at the same time the FDT relation. The analysis of the asymptotic aging dynamics has been cracked by Cugliandolo and Kurchan in [27, 67] (see also [60] for a pedagogical review) in the simplest spin glass model (see also [68] for a much more complex situation) where no signal is present. The outcome of this work is that when the dynamics started from a random initial conditions is run in the glass phase, it drives the system to surf on the *threshold* states. In the model analyzed in [27] these states correspond to the 1RSB marginally stable glassy states that maximize the complexity. In this section we analyze the Cugliandolo-Kurchan scenario by contrasting the numerical solution of the dynamical equations with the replica analysis of the complexity. According to [27], the long time Langevin dynamics⁷ can be characterized by two time regimes. For short times differences $t - t' \sim \mathcal{O}(1)$ and $t' \rightarrow \infty$, the system obeys the FDT theorem and TTI; this regime can be understood as a first fast local equilibration in the nearest metastable state available. On a longer timescale $t - t' \rightarrow \infty$ and $t/t' < \infty$, the dynamics surfs on threshold states and FDT and TTI are both violated. In this time window, both the response and correlation functions become functions of $\lambda = h(t)/h(t')$ being $h(t)$ an arbitrary reparametrization of the time variable⁸. By defining $\mathcal{C}(\lambda) = C(t, t')$ and $\mathcal{R}(\lambda) = tR(t, t')$ the Cugliandolo-Kurchan solution implies that in this aging regime the FDT relation can be generalized to

$$\mathcal{R}(\lambda) = x \mathcal{C}'(\lambda) \quad (112)$$

with x an *effective* FDT ratio that controls how much the FDT is violated. In the scenario of [27], the value of x coincides with the 1RSB Parisi parameter that corresponds to threshold states computed within the replica

⁵All one time quantities are constant in equilibrium.

⁶It is supposed to be reached only on exponential timescales in the system size.

⁷But still for times that are not exponentially large in the system size N .

⁸The function $h(t)$ must be a monotonously increasing function. The asymptotic reparametrization invariance is a key property of the dynamical equations [27].

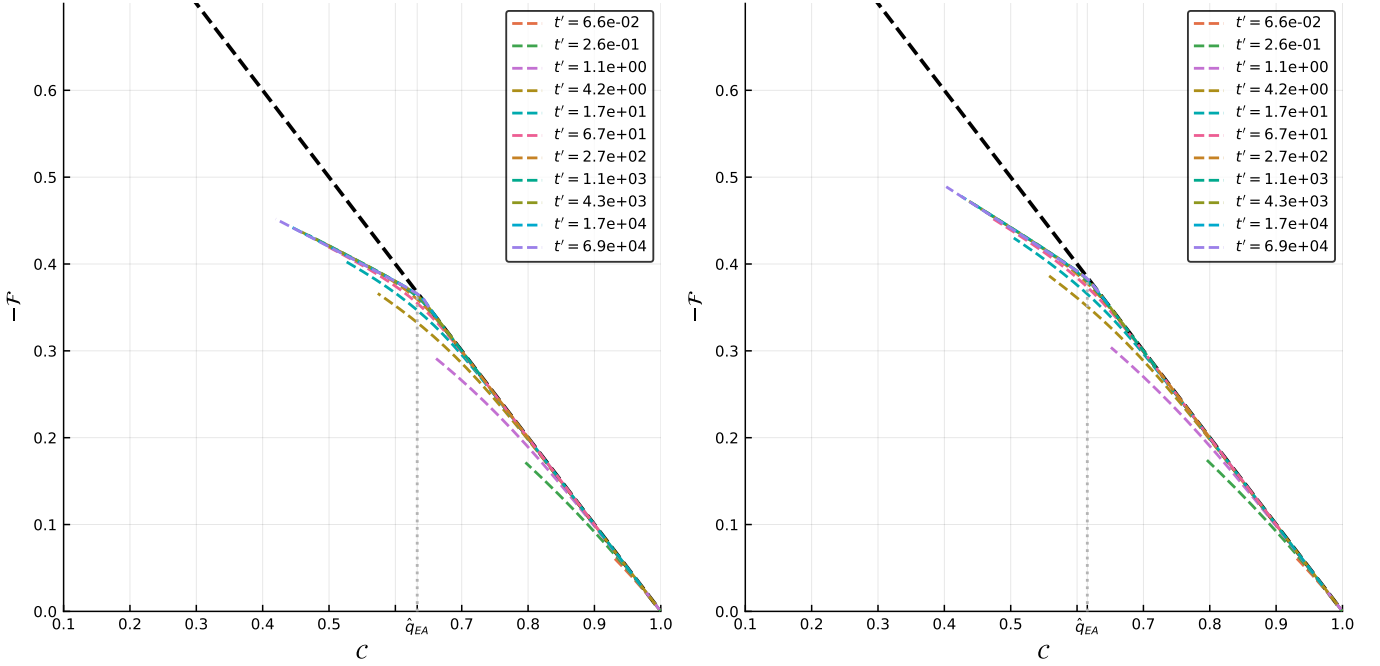


Figure 20: Left panel: parametric plot of integrate response function with respect to correlation function for $p = 3$, $\Delta_2 = 0.8$ and $\Delta_p = 0.2$. The different lines represent different waiting time, t' . The black dashed line correspond to the FDT prediction $x = 1$. The vertical dotted line is the point where we observe a kink, which we denote by $C = \hat{q}_{EA}$ and should be equal to the saddle point value of q_M as extracted from the 1RSB threshold states in the replica computation [27]: $\hat{q}_{EA} = 0.633$ and $q_M = 0.638$. For C smaller than q_{EA} the FDT is violated and is replaced by a generalized version as in Eq. (112). We can obtain the value of the FDT ratio from a fit of the slope of the asymptotic curves for $C < \hat{q}_{EA}$. We obtain $\hat{x} = 0.397$ which should be compared with the Parisi parameter that corresponds to 1RSB marginally stable states obtained from the replica computation that is $x = 0.408$. Right panel: parametric plot of the integrated response as a function of the correlation for $p = 3$ and $\Delta_2 = 1.4$ and $\Delta_p = 0.2$. In this case the value of the FDT ratio extracted from fitting the data is $\hat{x} = 0.397$ to be compared with the value of the Parisi parameter for the 1RSB threshold states that is $x = 0.408$. At the same time data gives $\hat{q}_{EA} = 0.633$ while the replica computation gives $q_M = 0.638$.

approach. In order to test this picture we follow Cugliandolo and Kurchan [69] and we plot the integrated response $\mathcal{F}(t, t') = -\int_{t'}^t R(t, t'') dt''$ as a function of $C(t, t')$ in a parametric way. This is done in Fig. 20.

If FDT holds at all timescales, one should see a straight line with slope -1 . Instead what we see in the Langevin hard phase is that for large values of t' the curves approach asymptotically for $t' \gg 1$ two straight lines. For high values of C , meaning for short time differences $t - t' \sim \mathcal{O}(1)$, the slope of the straight line is -1 which means that $\mathcal{F} = 1 - C$ as implied by the short time FDT relation. On longer timescales FDT is violated, confirming the glassiness of the Langevin hard phase. By doing a linear fit we can use the data plotted in Fig. 20 to estimate the FDT ratio x appearing in Eq. (112). This can be compared with the Parisi parameter x for which we have marginally stable 1RSB states. We find an overall very good agreement (data coming from the fit is reported in the caption of Fig. 20). The small discrepancy between the two values of x can be either due to the numerical accuracy in solving the dynamical equations as well as the possibility that the 1RSB threshold is not exactly the one that characterizes the long time dynamics. Further investigations are needed to clarify this point. Finally, according to [27] the value of C at which the two straight line cross should coincide with the value of q_M computed for the threshold states within the 1RSB solution. Again we find a very good agreement.