

Dense Hopfield Networks in the Teacher-Student Setting

Robin Thériault^{1*} and Daniele Tantari²

¹ Scuola Normale Superiore di Pisa,
Piazza dei Cavalieri 7, 56126, Pisa (PI), Italy

² Department of Mathematics, University of Bologna,
Piazza di Porta San Donato 5, 40126, Bologna (BO), Italy

* robin.theriault@sns.it

Abstract

Dense Hopfield networks with p -body interactions are known for their feature to prototype transition and adversarial robustness. However, theoretical studies have been mostly concerned with their storage capacity. We derive the phase diagram of pattern retrieval in the teacher-student setting of p -body networks, finding ferromagnetic phases reminiscent of the prototype and feature learning regimes. On the Nishimori line, we find the critical amount of data necessary for pattern retrieval, and we show that the corresponding ferromagnetic transition coincides with the paramagnetic to spin-glass transition of p -body networks with random memories. Outside of the Nishimori line, we find that the student can tolerate extensive noise when it has a larger p than the teacher. We derive a formula for the adversarial robustness of such a student at zero temperature, corroborating the positive correlation between number of parameters and robustness in large neural networks. Our model also clarifies why the prototype phase of p -body networks is adversarially robust.

Copyright attribution to authors.

This work is a submission to SciPost Physics.

License information to appear upon publication.

Publication information to appear upon publication.

Received Date

Accepted Date

Published Date

1

2 Contents

| | | |
|----|---|----|
| 3 | 1 Introduction | 2 |
| 4 | 2 Overview of Gardner's results | 4 |
| 5 | 3 Teacher-student setting | 6 |
| 6 | 3.1 Matched interaction orders | 7 |
| 7 | 3.2 Mismatched interaction orders | 8 |
| 8 | 4 Results and Discussion | 9 |
| 9 | 4.1 Retrieval transition at large interaction order | 9 |
| 10 | 4.2 Transition to the ordered phases: universality | 10 |
| 11 | 4.3 Phase diagram on the Nishimori line | 11 |
| 12 | 4.4 Inference temperature vs dataset noise | 13 |
| 13 | 4.5 Interaction order and noise tolerance | 15 |
| 14 | 4.5.1 Large noise scaling | 15 |

| | | |
|----|---|-----------|
| 15 | 4.5.2 Finite noise scaling | 17 |
| 16 | 4.6 Robustness against adversarial attacks | 17 |
| 17 | 5 Conclusion | 19 |
| 18 | References | 20 |
| 19 | A Gardner’s Hamiltonian vs K & H’s Hamiltonian | 26 |
| 20 | B Direct model cumulant expansions | 28 |
| 21 | C Teacher-student replicated partition function | 31 |
| 22 | D Teacher-student free entropy | 33 |
| 23 | E Direct model RSB ansatz | 41 |
| 24 | F Monte-Carlo simulations for various system sizes | 43 |

25
26

27 1 Introduction

28 Hopfield networks are artificial neural networks that model associative memory [1]. In the
29 Hopfield model, examples $\sigma \in \{-1, 1\}^N$ of memories $\xi^\mu \in \{-1, 1\}^N$, $\mu = 1, \dots, M$, are
30 retrieved by sampling the Gibbs distribution of a 2-body Hamiltonian $H[\sigma|\xi]$ at a given
31 temperature T [2]. Hopfield networks can be trained in a biologically plausible way using
32 Hebb’s rule [1, 3], which leads to $H[\sigma|\xi] = -\frac{1}{N} \sum_{\mu=1}^M \left(\sum_{i=1}^N \xi_i^\mu \sigma_i \right)^2$. However, they can only
33 store up to $M \sim \mathcal{O}(N)$ i.i.d. random memories in the limit of large N [1, 4, 5]. One way to
34 find this scaling is to study the phase diagram of $H[\sigma|\xi]$ as a function of the temperature T
35 and load $\alpha = \frac{M}{N}$ [5], where the so-called ferromagnetic phase, which extends up to $\alpha \approx 0.14$,
36 corresponds to accurate retrieval.

37 Since Hopfield’s seminal work, several generalizations have been investigated in relation
38 to their critical storage capacity and retrieval capabilities. For example, parallel retrieval
39 has been studied in relation to pattern sparsity [6–10] or hierarchical interactions [11–15],
40 and non-universality has been shown with respect to more general pattern entries and unit
41 priors [16–22]. Efforts to overcome the $\mathcal{O}(N)$ limitation of the capacity led to the development
42 of a novel class of modern Hopfield networks [23–25], which are sometimes called dense due to
43 their faculty to store much more memories than the original Hopfield model [26]. These neural
44 networks surpass $\mathcal{O}(N)$ storage capacity by using higher-order interactions instead of the
45 original 2-body couplings [27–32]. In particular, Gardner [30] calculated the replica-symmetric
46 (RS) phase diagram of the Hamiltonian $H[\sigma|\xi] = -\sum_{i_1 < \dots < i_p=1}^N J_{i_1 \dots i_p} \sigma_{i_1} \dots \sigma_{i_p}$ with p -body
47 interactions $J_{i_1 \dots i_p} = \frac{p!}{N^{p-1}} \sum_{\mu=1}^M \xi_{i_1}^\mu \dots \xi_{i_p}^\mu$ conditioned on i.i.d. random memories $\xi^\mu \in \{-1, 1\}^N$,
48 finding a $M = \mathcal{O}(N^{p-1})$ storage capacity. These calculations were later extended to include
49 the effects of one-step replica symmetry breaking (1RSB) [33].

50 Although they draw a rather detailed picture of the retrieval of individual i.i.d. random
51 memories, these results are not the end of the story. First of all, 1RSB calculations allegedly
52 struggle to find the paramagnetic to spin-glass phase transition accurately at large p because
53 of numerical instability issues [33]. Second of all, dense Hopfield networks have been rapidly

54 gaining a renewed attention for reasons other than their storage capacity since a recent
 55 paper [26] by Krotov and Hopfield (K & H), where they were used as a trainable machine
 56 learning architecture. For instance, they have been related to transformers [23,34] and diffusion
 57 models [35,36], and they were found to be significantly more explainable and adversarially
 58 robust than feedforward neural networks with ReLU activation functions [26,37].

59 One such aspect of dense Hopfield networks that is still poorly understood is their per-
 60 formance as generative models for unsupervised learning, where they are trained over some
 61 given dataset to reproduce its probability distribution. As far as we are aware, this problem has
 62 not yet been studied theoretically for p -body models with $p \geq 3$. However, it was studied for
 63 the original 2-body Hopfield network by using the teacher-student setting [38] first described
 64 in [16,17,39]. In the teacher-student setting, which is also called inverse problem in opposition
 65 to the direct problem of random pattern retrieval, a student model $H[\xi|\sigma]$ is trained with M
 66 teacher examples $\sigma^a \sim H[\sigma^a|\xi^*]$ conditioned on the planted pattern ξ^* . In other words, the
 67 student tries to infer the pattern ξ^* of the teacher using a structured set of examples σ^a .

68 At finite load $\alpha = \frac{M}{N}$, two regimes of pattern retrieval were found: example retrieval
 69 (eR) and signal retrieval (sR). In the eR phase, the student tries to reconstruct ξ^* by directly
 70 retrieving the examples σ^a , which is a good strategy provided that they are strongly correlated
 71 with ξ^* . In the sR phase, on the other hand, retrieval is done by extracting subtle cues from
 72 weakly correlated examples. The two types of examples used in these two retrieval strategies
 73 are respectively called prototypes and features of ξ^* [26]. Interestingly, a prototype regime
 74 and a feature regime were also observed by K & H in dense Hopfield networks trained to
 75 classify real data [26], where it was found that the prototype regime is significantly more
 76 adversarially robust than the feature regime. In other words, the prototype regime is more
 77 resistant than the feature regime to small data perturbations that are specifically designed to
 78 cause incorrect classification [40,41]. This prototype approach is arguably a big step towards
 79 designing adversarially robust neural networks, a long-standing problem that still lacks a fully
 80 satisfying solution [42–44].

81 In this work, we study the performance of p -body Hopfield networks in the teacher-student
 82 setting, revealing a prototype regime and a feature regime as in the 2-body model. In Section
 83 2, we review Gardner’s main results in studying p -body Hopfield models and summarize
 84 what the rest of the literature on spin-glass models with p -body interactions tell us about
 85 the paramagnetic to spin-glass phase transition in p -body Hopfield models. In Section 3, we
 86 compute the phase diagram of these p -body models in the teacher-student setting. In Section
 87 4.1, we discuss the transition to the retrieval phase in the inverse problem. In Section 4.2,
 88 we compare this retrieval transition against the transition to the spin-glass phase in the direct
 89 problem. Despite their different nature, we show that these two transitions are equivalent
 90 on the Nishimori line where the teacher and the student have the same p and T [45–48]. In
 91 Section 4.3, we discuss the phase diagram on the Nishimori line in more details. In Section
 92 4.4 and Section 4.5, we discuss the phase diagram outside of the Nishimori line. First of all,
 93 we investigate the effect of using an inference temperature different from the dataset noise.
 94 Second of all, we reveal that using a larger p for the student than the teacher gives the student
 95 an extensive tolerance against both teacher noise and pattern interference. Finally, in Section
 96 4.6, we derive a closed-form expression that measures the adversarial robustness of the student
 97 at zero temperature and explain what our results reveal about the nature of adversarial attacks.

98 2 Overview of Gardner's results

99 Consider the p -body Hamiltonian

$$H[\sigma|\xi] = - \sum_{i_1 < \dots < i_p = 1}^N J_{i_1 \dots i_p} \sigma_{i_1} \dots \sigma_{i_p} = - \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p = 1}^N \sum_{\mu=1}^M \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1} \dots \sigma_{i_p} \quad (1)$$

100 conditioned on a set of $M = \frac{\alpha N^{p-1}}{p!}$ quenched memories $\xi^\mu \in \{-1, 1\}^N$, $\mu = 1, \dots, M$, sampled
 101 i.i.d. from the Rademacher distribution $\frac{1}{2} [\delta(\xi_i^\mu - 1) + \delta(\xi_i^\mu + 1)]$. In the *direct model*,
 102 patterns σ are in turn sampled from the equilibrium Gibbs distribution $P(\sigma|\xi) = Z^{-1} e^{-\beta H[\sigma|\xi]}$,
 103 where $\beta \geq 0$ is the inverse temperature and $Z = \sum_{\sigma} e^{-\beta H[\sigma|\xi]}$ is the system's partition function.
 104 The so-called *direct problem* studied by Gardner [30] consists of quantifying the performance
 105 of this model as a method of memory retrieval. In that context, the overlap $\frac{1}{N} \sum_i \xi_i^\mu \sigma_i$ is a
 106 good measure of retrieval accuracy, and its expected value can be derived from the quenched
 107 free entropy $f = \frac{1}{N} \langle \log Z \rangle_\xi$ in the thermodynamic limit $N \rightarrow \infty$. At finite p , Gardner used
 108 the (non-rigorous) replica trick [49] to evaluate the RS approximation of f (see also Appendix
 109 B) in terms of a variational principle of the form

$$f = \lim_{N \rightarrow \infty} \frac{1}{N} \langle \log Z \rangle_\xi = \lim_{N \rightarrow \infty, L \rightarrow 0} \left(\frac{\partial}{\partial L} \left[\frac{1}{N} \log \langle Z^L \rangle_\xi \right] \right) = \text{Extr}_{m, k, q, r} f(m, k, q, r),$$

110 whose solution is

$$\begin{aligned} q &= \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) \tanh^2(\beta[\sqrt{\alpha r}x + k]) \\ m &= \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) \tanh(\beta[\sqrt{\alpha r}x + k]) \\ r &= pq^{p-1} \\ k &= pm^{p-1}, \end{aligned} \quad (2)$$

111 and the order parameters m and q are to be interpreted as expected overlaps. To be more
 112 precise, m can be shown to be the expected overlap of a retrieval attempt σ against one memory
 113 in the thermodynamic limit, i.e. $m = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \xi_i^\mu \sigma_i \right\rangle_{\xi, \sigma}$. Similarly, q is the expected
 114 overlap between two retrieval attempts σ^1 and σ^2 , i.e. $q = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \sigma_i^1 \sigma_i^2 \right\rangle_{\xi, \sigma}$ or
 115 equivalently $q = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \langle \sigma_i \rangle_\sigma^2 \right\rangle_\xi$. Intuitively, q measures the tendency of the system
 116 to stay frozen in specific configurations rather than visiting all possible values of σ .

117 The resulting RS phase diagram (see Fig. 1) are derived from the value of the order
 118 parameters as a function of three *hyperparameters*: the interaction order p , temperature
 119 $T = 1/\beta$ and load $\alpha = \frac{Mp!}{N^{p-1}}$. There are four different phases:

- 120 • In the Paramagnetic phase (**P**), the overlaps m and q both vanish. The network does not
 121 retrieve any specific pattern: sampled configurations are completely random.
- 122 • In the Spin-Glass phase (**SG**), m vanishes but $q > 0$. In other terms, the network does not
 123 retrieve individual stored memories but rather converges to spurious patterns depending
 124 on all the memories in a non-trivial way.
- 125 • In the signal Retrieval phases (**lR** and **gR**), $m \neq 0$ and $q > 0$, which means that the
 126 network is able to retrieve the stored memories. **lR** and **gR** are respectively locally

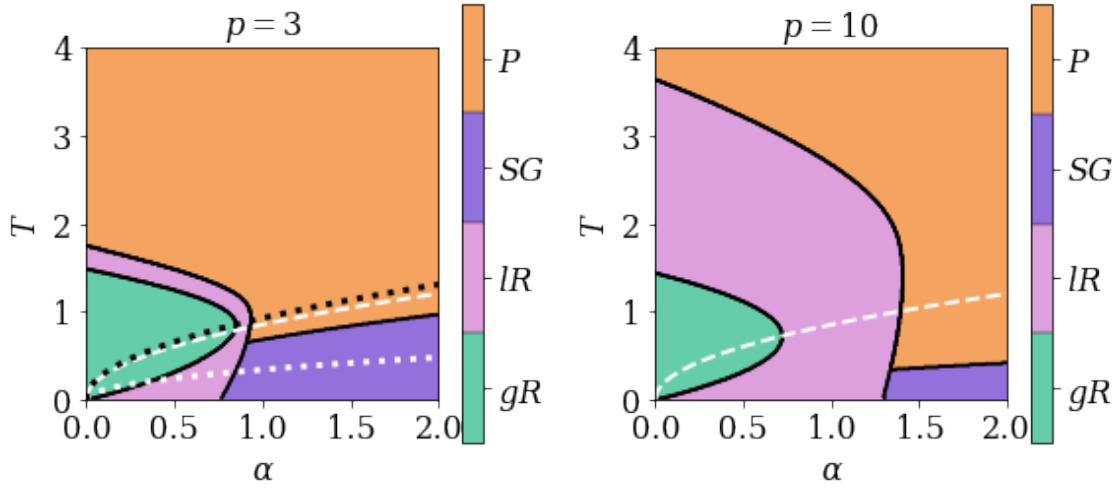


Figure 1: RS phase diagrams of the direct models with $p = 3$ on the left and $p = 10$ on the right. Accurate pattern retrieval is not possible in the paramagnetic phase (P) or in the spin-glass phase (SG), but it is possible in the local retrieval phase (IR) and in the global retrieval phase (gR). The ferromagnetic fixed point corresponding to accurate pattern retrieval is globally stable in the gR phase, but locally stable in the IR phase. The phase diagrams are inexact below the white dashed line where the total entropy of the paramagnetic phase becomes negative. The black dotted line overlaying the $p = 3$ diagram is the (exact) 1RSB P - SG transition temperature $T_s(\alpha, 3)$, which is obtained by rescaling by $\sqrt{2\alpha}$ the corresponding transition temperature of the spin-glass model with p -body Gaussian interactions. The d1RSB transition $T_d(\alpha, 3)$ is very close to $T_s(\alpha, 3)$ throughout the displayed range of α . The white dotted line in the $p = 3$ plot is the temperature $T_G(\alpha, 3)$ below which multiple steps of RSB are required to compute the free entropy. It is also obtained by rescaling by $\sqrt{2\alpha}$ the corresponding transition temperature of the Gaussian spin-glass model.

127 stable and globally stable. In other words, local retrieval IR is only attainable from
 128 initial conditions in a limited neighborhood of a memory ξ^μ , while global retrieval gR
 129 is accessible from any initial conditions given enough time. These two phases are said to
 130 be ferromagnetic.

131 Gardner also calculated the exact $p \rightarrow \infty$ phase diagram without making any assumptions
 132 about replica symmetry [30]. In this limit, the resulting paramagnetic to spin-glass (P - SG)
 133 phase transition occurs at a temperature $T_E(\alpha)$ that coincides with the boundary of the region
 134 where the total entropy of the paramagnetic phase becomes negative, given by $\beta^2 \alpha = 2 \log 2$
 135 (white dashed line in Fig. 1).

136 At finite p , Gardner's results only tell us that the model cannot be in the paramag-
 137 netic phase below $T_E(\alpha)$. Therefore, a spin-glass transition should occur at a temperature
 138 $T_s(\alpha, p) \geq T_E(\alpha)$. Since the RS spin glass solution of Eqs. (2) exists only below $T_E(\alpha)$ (violet
 139 region in Fig. 1), the spin-glass transition must be towards a RSB spin-glass phase.

140 Outside of the signal retrieval phases, the free entropy of the direct model is the same as
 141 for the spin-glass model with p -body Gaussian interactions where the temperature is rescaled
 142 by a factor of $\sqrt{2\alpha}$ [50, 51]. Therefore, the spin-glass and paramagnetic solutions are the
 143 same in the direct model as in this Gaussian spin-glass model, and we expect the exact phase
 144 diagrams of both models to be identical when the direct model is not in its signal retrieval
 145 phases. According to previous work on the Gaussian model with finite p [51], a 1RSB solution
 146 with $m = k = 0$ exists and is globally stable throughout a whole phase below $T_s(\alpha, p) \geq T_E(\alpha)$

147 (see Fig. 1). This solution becomes unstable at a lower transition temperature $T_G(\alpha, p)$
 148 (see Fig. 2), below which multiple steps of RSB are required. In the limit of $p \rightarrow \infty$, it
 149 holds that $T_s(\alpha, p) \rightarrow T_E(\alpha)$ and $T_G(\alpha, p) \rightarrow 0$. In other terms, the direct model becomes
 150 1RSB, which is consistent with the fact that it is converging to a random energy model with
 151 temperature rescaled by $\sqrt{2\alpha}$ [30, 50, 52]. Finally, we mention that this type of models
 152 exhibits a random first order transition phenomenology [53–56]: there is in fact a range of
 153 temperatures $T_s(\alpha, p) \leq T \leq T_d(\alpha, p)$ where the dynamics get trapped in an exponential
 154 number of metastable clusters, with an emerging RSB structure that does not affect the free
 155 energy (see Fig. 2). This range of temperatures thus defines a so-called dynamical 1RSB
 156 (d1RSB) phase. Below $T_s(\alpha, p)$, the number of clusters is no longer exponential, and the
 157 system undergoes the thermodynamic 1RSB phase transition that we mentioned previously.
 158 The critical temperatures $T_G(\alpha, p)$, $T_s(\alpha, p)$ and $T_d(\alpha, p)$ can all be obtained by standard RSB
 159 methods, but the resulting saddle-point equations can be prone to numerical instability at large
 160 p [33]. In Sections 4.2 and 4.3, we discuss an alternative way to obtain $T_s(\alpha, p)$ and $T_d(\alpha, p)$.

161 3 Teacher-student setting

162 On our end, we study a dense Hopfield network with Hamiltonian (1) as a generative model
 163 for unsupervised learning. In that context, the memories ξ are model parameters that have to
 164 be trained in such a way that the examples of a given dataset $\{\sigma^a\}_{a=1}^M$ result as typical network
 165 configurations.

166 In particular, we study a controlled teacher-student setting in which the examples are
 167 sampled from the probability distribution $P(\sigma^a|\xi^*)$ of a so-called *teacher* dense Hopfield
 168 network conditioned on a single *planted* pattern $\xi^* \in \{-1, 1\}^N$ whose entries are quenched
 169 Rademacher random variables. A *student* dense Hopfield network, also known as the *inverse*
 170 *model*, then samples its own student pattern ξ from the posterior distribution

$$P(\xi|\sigma) = \frac{P(\xi) \prod_{a=1}^M P(\sigma^a|\xi)}{P(\sigma)} = \frac{P(\xi)}{P(\sigma)} \prod_{a=1}^M Z^{-1} \exp(-\beta H[\sigma^a|\xi]),$$

171 where $P(\sigma^a|\xi)$ is the Gibbs distribution of the direct model with a single memory ξ , and $P(\xi)$
 172 is the prior on ξ that is chosen to be uniform. Since the direct model has only a single pattern,
 173 Z does not depend on ξ (see Appendix C), and the posterior simplifies to

$$P(\xi|\sigma) = Z^{-1}(\sigma) \exp(-\beta H[\xi|\sigma]).$$

174 In sum, the student posterior distribution is that of a dense Hopfield network where ξ plays the
 175 role of the sampled pattern and the examples σ act like the M quenched memories. Our task,
 176 called the *inverse problem*, consists of quantifying the student’s capability to infer the teacher
 177 pattern, which we will also call the *signal*. Like Gardner, we calculate a free entropy of the
 178 form $f = \frac{1}{N} \langle \log Z \rangle_\sigma$ in the thermodynamic limit $N \rightarrow \infty$. This time, however, the average
 179 $\langle \cdot \rangle_\sigma$ is over a structured set of examples σ . In fact, we recall that, unlike the i.i.d. memories
 180 studied by Gardner, the examples σ^a are sampled from the teacher distribution $P(\sigma^a|\xi^*)$.

181 In general, the student does not have access to the teacher generative model. In our
 182 controlled teacher-student setting, the student knows that the correct model for $P(\sigma^a|\xi)$ is a
 183 dense Hopfield network. Nevertheless, it does not necessarily have access to the interaction
 184 order p^* and inverse temperature β^* used by the teacher. Therefore, we denote the student
 185 hyperparameters by p and β and emphasize that they are not necessarily equal to p^* and β^* .

186 As previously stated, we calculate the free entropy

$$f = \frac{1}{N} \langle \log \mathcal{Z} \rangle_\sigma = 2^{-N} \sum_{\xi^*} \sum_{\sigma} [Z^*]^{-M} \exp \left(\beta^* \frac{p^*!}{N^{p^*-1}} \sum_{a=1}^M \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \right) \\ \times \log \sum_{\xi} \exp \left(\beta \frac{p!}{N^{p-1}} \sum_{a=1}^M \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right) \quad (3)$$

187 in the thermodynamic limit $N \rightarrow \infty$. We then draw phase diagrams of the inverse problem as
 188 a function of p^* , $T^* = 1/\beta^*$, p , $T = 1/\beta$ and α , where α is M normalized to $\mathcal{O}(1)$. Unless
 189 explicitly specified otherwise, we use $\alpha = \frac{Mp!}{N^{p-1}}$.

190 3.1 Matched interaction orders

191 We first consider the case where $p^* = p$ and the only possible mismatch between the teacher
 192 and student networks is in the inverse temperature, i.e. $\beta^* \neq \beta$. At low T^* , the student's task
 193 is easy. In fact, below the critical temperature T_{crit} of the direct problem with one pattern (see
 194 Fig. 1, $\alpha = 0$ axis), the teacher produces examples σ^a that cluster around ξ^* . Therefore, the
 195 student can infer ξ^* by aligning its pattern ξ with the examples σ^a . This retrieval strategy
 196 works even when using a very small amount of examples (see [38]). Since the size of our
 197 dataset is extensive, the retrieval accuracy is maximum in the thermodynamic limit. We call
 198 this region the (accurate) example Retrieval phase (**eR**).

199 Conversely, when T^* is above T_{crit} , the examples in the training set are very noisy and we
 200 do not observe a finite overlap between σ^a and ξ^* (see Fig. 1, $\alpha = 0$ axis). In this regime, we
 201 find that the RS approximation of the $p^* = p$ free entropy can be computed (see Appendix D)
 202 in terms of the variational principle

$$f = \underset{m, k, q, r, q^*, r^*}{\text{Extr}} \left\{ \beta^* \beta \alpha [q^*]^p - \frac{1}{2} \beta^2 \alpha q^p + \beta m^p - \beta^* \beta \alpha r^* q^* \right. \\ \left. + \frac{1}{2} \beta^2 \alpha r q - \frac{1}{2} \beta^2 \alpha r - \beta m k + \frac{1}{2} \beta^2 \alpha + \log 2 \right. \\ \left. + \int dx \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} x^2 \right\} \left\langle \log [\cosh (\beta [\sqrt{\alpha r} x + \beta^* \alpha r^* + k z])] \right\rangle_z \right\}, \quad (4)$$

203 whose solution is the saddle-point equations

$$q^* = \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} x^2 \right) \left\langle \tanh (\beta [\sqrt{\alpha r} x + \beta^* \alpha r^* + k z]) \right\rangle_z \\ q = \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} x^2 \right) \left\langle \tanh^2 (\beta [\sqrt{\alpha r} x + \beta^* \alpha r^* + k z]) \right\rangle_z \\ m = \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} x^2 \right) \left\langle z \tanh (\beta [\sqrt{\alpha r} x + \beta^* \alpha r^* + k z]) \right\rangle_z \quad (5) \\ r^* = p [q^*]^{p-1} \\ r = p q^{p-1} \\ k = p m^{p-1},$$

204 where z is a Rademacher random variable and $\alpha = \frac{Mp!}{N^{p-1}}$. As in the direct model described
 205 in Section 2, the order parameters m and q have a clear interpretation in terms of expected

206 overlaps. $m = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \xi_i \sigma_i^a \right\rangle_{\xi^*, \sigma, \xi}$ is the expected overlap of a retrieval attempt with
 207 an example σ^a , and $q = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \langle \xi_i \rangle_{\xi}^2 \right\rangle_{\xi^*, \sigma}$ is the expected overlap between two
 208 retrieval attempts. Similarly, q^* is the expected overlap between the teacher and student
 209 patterns, i.e. $q^* = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \xi_i^* \xi_i \right\rangle_{\xi^*, \sigma, \xi}$. Therefore, it is a good measure of inference
 210 performance. The free entropy (Eq. 4) is expected to be exact in absence of mismatch between
 211 the teacher and the student, i.e. $\beta^* = \beta$. This condition is known as the Nishimori line [45–48].
 212 Outside of the Nishimori region, RSB corrections are expected. Like the direct problem, the
 213 inverse problem with $T^* > T_{\text{crit}}$ has different phases characterized by the values of the order
 214 parameters:

- 215 • In the Paramagnetic phase (P), the overlaps m , q^* and q all vanish.
- 216 • In the signal Retrieval phases (lR and gR), $m = 0$ but $q^* \neq 0$ and $q > 0$. lR and gR are
 217 respectively locally stable and globally stable. In other words, local retrieval lR is only
 218 attainable from initial conditions in a limited neighborhood of ξ^* , while global retrieval
 219 gR is accessible from any initial conditions given enough time. These two phases are
 220 also said to be ferromagnetic.
- 221 • In the (inaccurate) example Retrieval phase (eR), $m \neq 0$ and $q > 0$ but $q^* = 0$.
- 222 • In the Spin-Glass phase (SG), $q > 0$ but q^* and m vanish.

223 In sum, when T^* is above T_{crit} , the student can only learn the teacher pattern in the signal
 224 retrieval phases. In all the other phases, the student pattern is uncorrelated with the signal, being
 225 either a random guess (P phase), aligned with a noisy example (inaccurate eR phase), or aligned
 226 with a spurious low energy state (SG phase). We stress that we cannot have $m \neq 0$ and $q^* \neq 0$
 227 at the same time (accurate eR phase) when $T^* > T_{\text{crit}}$ because $\lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \xi_i^* \sigma_i^a \right\rangle_{\xi^*, \sigma} = 0$
 228 in that regime (see Fig. 1, $\alpha = 0$ axis).

229 3.2 Mismatched interaction orders

230 We also investigate the $T^* > T_{\text{crit}}$ regime in the presence of a mismatch between the interaction
 231 orders of the teacher and student networks, i.e. $p^* \neq p$. We focus on the case of $p^* = 2$ and
 232 even $p \geq 3$ to study the consequences of fitting the teacher of [38] using a student with higher
 233 order interactions. We find two different scaling regimes of the training set size M and inverse
 234 temperature β^* that make retrieval possible (see Appendix D):

- 235 • a large-noise scaling where $\beta^* \sim \mathcal{O}(N^{2/p-1})$ and $M \sim \mathcal{O}(N^{p-1})$, such that $\alpha = \frac{Mp!}{N^{p-1}}$
 236 and $\lambda = \frac{[\beta^*]^{p/2}}{(p/2)!} N^{p/2-1}$ are finite;
- 237 • a finite-noise scaling where $\beta^* \sim \mathcal{O}(1)$ and $M \sim \mathcal{O}(N^{p/2})$, such that $\alpha = \frac{M(p/2+1)!}{N^{p/2}}$ is
 238 finite.

239 In the large-noise scaling, we obtain saddle point equations similar to Eqs. (5) but with β^*
 240 replaced by λ (see Appendix D). Conversely, the finite noise scaling leads to

$$\begin{aligned}
 q^* &= \left\langle \tanh(\beta [\eta \alpha r^* + k z]) \right\rangle_z \\
 m &= \left\langle z \tanh(\beta [\eta \alpha r^* + k z]) \right\rangle_z \\
 r^* &= p [q^*]^{p-1} \\
 k &= p m^{p-1},
 \end{aligned} \tag{6}$$

241 where η generally depends on β^* and p in a non-trivial way, but we find that $\eta = \frac{2[\beta^*]^2}{(1-2\beta^*)^2}$ when
 242 $p = 4$ (see Appendix D). These equations can also be derived by extrapolating the large-noise
 243 equations to $\alpha_{\text{large noise}} \rightarrow 0$ and $\lambda \rightarrow \infty$ with fixed $\lambda \alpha_{\text{large noise}} = \eta \alpha_{\text{finite noise}}$.

244 4 Results and Discussion

245 4.1 Retrieval transition at large interaction order

246 The paramagnetic solution of Eqs. (5) always exists and is globally stable in the part of the
 247 phase diagram where the temperature T is relatively large and $\alpha = \frac{Mp!}{N^{p-1}}$ is relatively small.
 248 On the other hand, the gR phase exists when $\beta^2 \alpha p$ and $\beta^* \beta \alpha p$ are both large. In fact, in
 249 that limit, $q^* = q = 1$ is a fixed point of Eqs. (5). The critical line where gR becomes globally
 250 stable instead of P is not clear from this analysis alone, but we can at least find it analytically
 251 in the limit of infinite p . As for the direct model, the free entropy and the total entropy of the
 252 paramagnetic phase are respectively $\frac{1}{2}\beta^2 \alpha + \log 2$ and $-\frac{1}{2}\beta^2 \alpha + \log 2$ [30]. At the same time,
 253 the $p \rightarrow \infty$ free entropy takes the form

$$f = \text{Extr} \left\{ \beta^* \beta \alpha \theta(q^* - 1) - \frac{1}{2} \beta^2 \alpha \theta(q - 1) - \beta^* \beta \alpha r^* q^* + \frac{1}{2} \beta^2 \alpha r q - \frac{1}{2} \beta^2 \alpha r + \frac{1}{2} \beta^2 \alpha \right. \\ \left. + \log 2 + \int dx \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2\right\} \log \left[\cosh\left(\sqrt{\beta^2 \alpha r x} + \beta^* \beta \alpha r^*\right) \right] \right\},$$

254 where $\theta(q - 1) := \lim_{p \rightarrow \infty} q^p$, $q \in [0, 1]$, is the Heaviside step function jumping at $q = 1$,
 255 i.e. $\theta(1) = 1$ and $\theta(q) = 0 \forall q \in [0, 1)$. In this limit, the ferromagnetic phase is characterized
 256 by $q = q^* = 1$, and its free entropy is then

$$f = \beta^* \beta \alpha - \beta^* \beta \alpha p + \int dx \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2\right\} \log \left[2 \cosh\left(\sqrt{\beta^2 \alpha p x} + \beta^* \beta \alpha p\right) \right] \\ \approx \beta^* \beta \alpha - \beta^* \beta \alpha p + \int dx \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2\right\} \left(\sqrt{\beta^2 \alpha p x} + \beta^* \beta \alpha p\right) \\ = \beta^* \beta \alpha.$$

257 The corresponding total entropy is $s = f - \beta \frac{\partial f}{\partial \beta} = 0$, as expected from a ferromagnetic phase
 258 with $q^* = q = 1$. On the Nishimori line, $f = \beta^* \beta \alpha$ becomes larger than the free entropy of
 259 the paramagnetic phase, which triggers a phase transition, if and only if

$$T < \sqrt{\frac{\alpha}{2 \log 2}}, \quad (7)$$

260 where $T_E = \sqrt{\frac{\alpha}{2 \log 2}}$ is also the temperature below which the total entropy of the paramag-
 261 netic phase becomes negative. Outside of the Nishimori line, this inequality generalizes to
 262 $\beta^* \beta \alpha > \frac{1}{2} \beta^2 \alpha + \log 2$, leading to

$$\beta^* - \sqrt{[\beta^*]^2 - \frac{2 \log 2}{\alpha}} < \beta < \beta^* + \sqrt{[\beta^*]^2 - \frac{2 \log 2}{\alpha}},$$

263 while the temperature where the paramagnetic total entropy becomes negative stays the same.

264 4.2 Transition to the ordered phases: universality

265 In the $p \rightarrow \infty$ limit, the transition towards \mathbf{gR} of the inverse model on the Nishimori line
 266 is identical to the exact $\mathbf{P-SG}$ transition of the direct model [30]. We claim that these two
 267 critical lines are actually closely related for any p . In the Hopfield model with $p = 2$, they
 268 were already shown to be identical [38]. We will now argue that they overlap for any p and
 269 β such that $T > T_{\text{crit}}$ (see Figs. 2 and 1). In the case of $p = 2$, both lines can be obtained
 270 exactly from the RS approximation of either the direct model or the inverse model, so there is
 271 no obvious advantage to using this equivalence in calculations. In general, while the inverse
 272 problem on the Nishimori line is replica symmetric, the direct problem is not, and the $p \geq 3$
 273 replica symmetric $\mathbf{P-SG}$ transition is not exact. Moreover, even the critical line calculated using
 274 1RSB may be inaccurate due to numerical instability [33]. In this situation, the knowledge
 275 of the \mathbf{gR} transition in the replica-symmetric inverse problem can be used to locate the exact
 276 $\mathbf{P-SG}$ transition of the direct problem, where symmetry breaking occurs.

277 For that purpose, we will argue that, given $T > T_{\text{crit}}$, *the direct model is in the paramagnetic*
 278 *phase if and only if the inverse model is in the paramagnetic phase.*

279 The converse implication comes from the fact that since (see Appendix C)

$$P(\boldsymbol{\sigma}) = \frac{1}{2^{MN}} \frac{\mathcal{Z}(\boldsymbol{\sigma})}{\langle \mathcal{Z} \rangle}, \quad (8)$$

280 the example distribution $P(\boldsymbol{\sigma})$ of the inverse problem is contiguous [57] to the uniform
 281 distribution, i.e. the memory distribution of the direct problem, when

$$\lim_{N \rightarrow \infty} \left\{ \frac{\log \mathcal{Z} - \log \langle \mathcal{Z} \rangle}{N} \right\} = 0. \quad (9)$$

282 As determined in Appendix C and D, the annealed expression $\frac{1}{N} \log \langle \mathcal{Z} \rangle$ is equal to the free
 283 entropy of the paramagnetic phase. Therefore, when the inverse model is in the paramagnetic
 284 phase, $P(\boldsymbol{\sigma})$ is contiguous to the uniform distribution. This property is called quiet planting
 285 and is known to occur more generally in mean-field paramagnets [58–61]. In our problem
 286 setting, it means that if the inverse model is in the paramagnetic phase, then it is equivalent to
 287 the direct model. In particular, if the inverse model is in the paramagnetic phase, then so is the
 288 direct model. In more intuitive terms, the \mathbf{gR} transition temperature of the inverse model must
 289 be greater than or equal to the $\mathbf{P-SG}$ transition temperature of the direct model because the
 290 ensemble of examples $\boldsymbol{\sigma}^a$ generated by the teacher model is on average at least as structured
 291 as the set of i.i.d. random memories stored in the direct model.

292 For the direct implication, notice that the average replicated partition function of the direct
 293 model in the paramagnetic phase can be approximated as (see Appendix E)

$$\begin{aligned} \langle Z^L \rangle \approx \frac{1}{\langle Z \rangle} & \left\langle \sum_{\boldsymbol{\sigma}} \exp \left(\beta N \sum_{\gamma} \sum_{\mu \in \Gamma_{\gamma}} \left[\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i^{\gamma} \right]^p \right. \right. \\ & \left. \left. + \beta \sum_{\gamma} \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^{\gamma} \dots \sigma_{i_p}^{\gamma} \right) \right. \\ & \left. \sum_{\boldsymbol{\sigma}_0} \exp \left(\beta \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^0 \dots \sigma_{i_p}^0 \right) \right\rangle. \end{aligned}$$

294 This expression is identical to the replicated partition function of the inverse model with
 295 $T > T_{\text{crit}}$, which therefore must also be in the paramagnetic phase.

296 As a consequence, when $T > T_{\text{crit}}$, the $\mathbf{P-SG}$ transition line of the direct model must be
 297 identical to the \mathbf{gR} transition line of the inverse model on the Nishimori line.

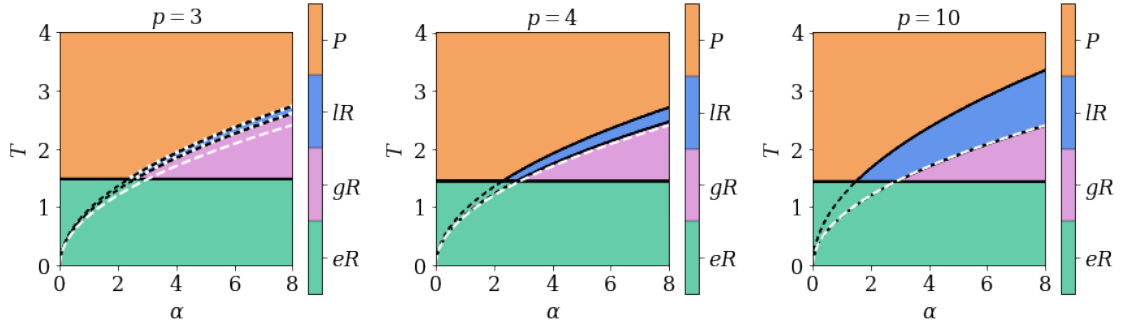


Figure 2: Exact RS phase diagrams of inverse models on the Nishimori line, i.e. $p^* = p$ and $\beta^* = \beta$. The left, center and right plots respectively have $p = 3$, $p = 4$ and $p = 10$. Accurate pattern retrieval is not possible in the paramagnetic phase (P), but it is possible in the local retrieval phase (IR), in the global retrieval phase (gR) and in the example retrieval phase (eR). The ferromagnetic fixed point corresponding to accurate pattern retrieval is globally stable in the gR phase, but locally stable in the IR phase. The critical temperature of the eR phase is the critical temperature T_{crit} of the direct problem with one pattern (see Fig. 1, $\alpha = 0$ axis). The black dashed lines mark the spurious continuation of the IR and gR phase boundaries through the eR phase. The white dashed line is the $p \rightarrow \infty$ gR critical line calculated analytically in Section 4.1. It matches the corresponding numerical phase boundary increasingly well as p grows larger. The white dotted lines on the $p = 3$ plot mark the 1RSB and d1RSB critical temperatures $T_s(\alpha, 3)$ and $T_d(\alpha, 3)$ of the direct model (see Section 2). We truncated them below T_{crit} for improved visibility. $T_s(\alpha, 3)$ and $T_d(\alpha, 3)$ are obtained by rescaling the corresponding critical temperatures found in [54] by $\sqrt{2\alpha}$.

298 4.3 Phase diagram on the Nishimori line

299 On the Nishimori line, the student is fully informed about the teacher generative model and
 300 uses $\beta = \beta^*$ and $p = p^*$. In this scenario, thanks to the Nishimori identities [46], it is well
 301 known that ξ^* and ξ play symmetric roles and that $q^* = q$. For the same reason, the overlaps
 302 $\frac{1}{N} \sum_i \xi_i^* \xi_i$ and $\frac{1}{N} \sum_i \xi_i^1 \xi_i^2$ have the same distribution. From the self-averaging of $\frac{1}{N} \sum_i \xi_i^* \xi_i$, it
 303 follows that the system is expected to be replica symmetric, and Eqs. (4) and (5) are expected
 304 to hold. Fig. (2) shows the phase diagrams obtained by solving the saddle-point equations
 305 numerically on the Nishimori line. Both $q^* = q$ and the replica symmetry condition are verified.
 306 In particular, numerical solutions of a few values of $p \geq 3$ show that the gR transition occurs
 307 at a higher T than the line $\beta^2 \alpha = 2 \log 2$ where the total entropy of the paramagnetic phase
 308 becomes negative. In other terms, the phase transition towards gR prevents the total entropy
 309 from becoming negative when T decreases below $\sqrt{\frac{\alpha}{2 \log 2}}$, which is consistent with the RS
 310 solution being exact on the Nishimori line.

311 At low T , the student can learn efficiently within the accurate eR regime. In this phase,
 312 learning is possible ($q^* \neq 0$) because the examples are correlated with the signal and the
 313 student can retrieve it by simply being aligned with them ($m \neq 0$).

314 At high T , learning is possible only if the amount of examples, i.e. the size of the dataset, is
 315 sufficiently large. When α is too small, Eqs. (5) have only a paramagnetic fixed point because
 316 the amount of information carried by the dataset is not large enough. Numerical solutions
 317 suggest that the paramagnetic fixed point always exist and it is actually locally stable in the
 318 whole high-temperature regime. When α is sufficiently large, the signal retrieval fixed point
 319 appears as a locally stable attractor (IR phase). It becomes globally stable (gR phase) as the
 320 size of the dataset is increased further or the student temperature decreases.

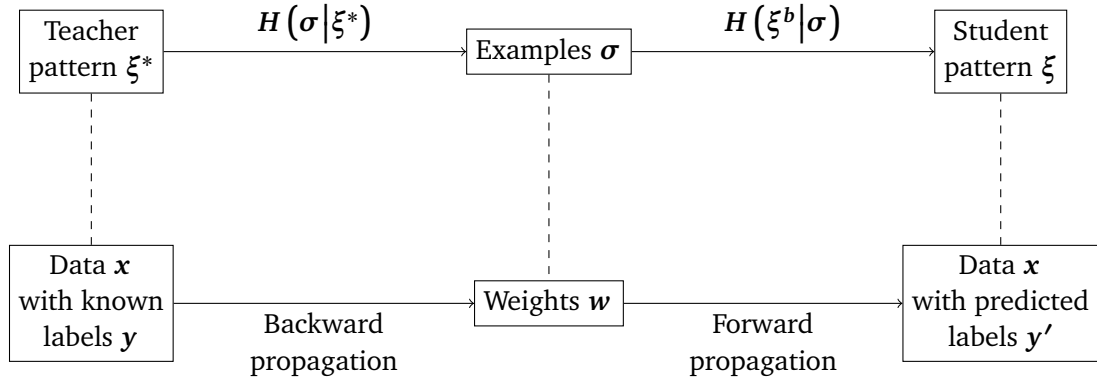


Figure 3: The first row of this diagram sketches how a p -body Hopfield network in the teacher-student setting can reconstruct an incomplete pattern ξ^b to match the teacher pattern ξ^* by relying on the examples σ obtained from ξ^* . The second row summarizes how a dense neural network trained by K & H can recover the labels y' of the data x given the weights w learned from x [26]. Both models tackle similar tasks using an approach where σ and ξ^b respectively play the same roles as w and (x, y') . The forward propagation algorithm used to generate y' is similar to the update rule of the student (see [26] and Appendix A), but the backpropagation algorithm used to learn w is very different from the update rule of the teacher.

321 As per the previous Section, the critical boundary of the gR phase obtained by solving Eqs.
 322 5 is identical to the 1RSB P -SG transition temperature $T_s(\alpha, p)$ of the direct model. Similarly,
 323 we observe that the metastable lR phase coincides with the d1RSB phase of the direct model
 324 (see Fig. 2). Our results are also consistent with the fact that $T_s(\alpha, p) \rightarrow T_E(\alpha)$ in the $p \rightarrow \infty$
 325 limit. In fact, we find that the analytical limit boundary closely agrees with the numerical
 326 solution of the saddle-point equations with $p^* = p = 10$ and remains a good approximation
 327 even down to $p^* = p = 4$.

328 In the student model, σ plays a similar role as the weights of the trainable dense Hopfield
 329 network model that K & H designed for classification of data [26]. In that context, ξ is analogous
 330 to the test data whose labels are being predicted (see Fig. 3). In fact, the computation performed
 331 by K & H's model to recover labels is similar to the update rule used by the student to infer the
 332 teacher pattern (see Appendix A). Moreover, the eR and gR phases are respectively reminiscent
 333 of the prototype and feature regimes of K & H's networks. Therefore, we believe that the
 334 student can act as a toy model of label prediction in these two regimes.

335 Comparing instead the phase diagrams of our inverse model with that of the inverse 2-body
 336 Hopfield model, we see that the eR and gR phases of the inverse p -body model with $p \geq 3$ are
 337 respectively analogous to the eR and sR (signal Retrieval) phases presented in [38]. One of
 338 the key differences between $p = 2$ and $p \geq 3$ is that the paramagnetic to signal retrieval phase
 339 transition of the p -body model is second order for $p = 2$ but first order for $p \geq 3$. On the one
 340 hand, the second order phase transition of $p = 2$ indicates that its paramagnetic fixed point is
 341 never locally stable and sets an unambiguous boundary between the sR phase where ξ^* can
 342 be recovered starting from any initial conditions and the paramagnetic phase where pattern
 343 retrieval is impossible [61]. On the other hand, the first order phase transition of $p \geq 3$ allows
 344 the retrieval and paramagnetic regimes to coexist. The lR phase is locally stable precisely
 345 because it coexists with the paramagnetic phase and has a lower free entropy. Meanwhile,
 346 the gR phase also coexists with the paramagnetic phase, but has a larger free entropy. In the
 347 presence of phase coexistence, an algorithm trying to retrieve ξ^* starting from random initial
 348 conditions can get stuck in the paramagnetic phase instead. In fact, it has been conjectured

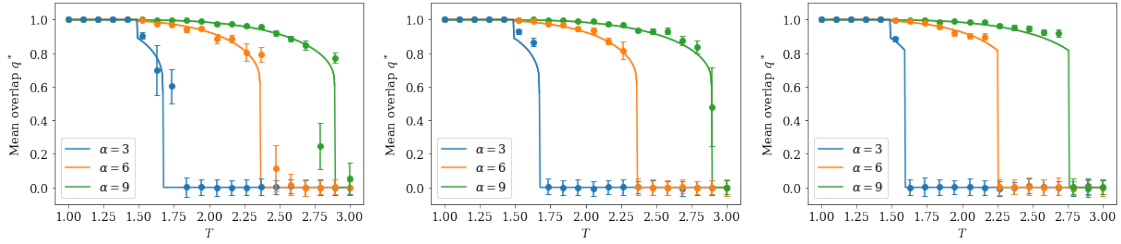


Figure 4: Monte-Carlo simulations of the $p = 3$ inverse model compared against RS saddle-point solutions. The IR phase is included on the left and central plots, but not on the right one. The left plot has $\varepsilon = 0$, and the two other ones have a handpicked ε such that the simulations are initialized near the saddle-point solutions. The dots are simulation data at a few values of α , and the lines are slices of the saddle-point solutions at the same α . The teacher generates $M = \frac{\alpha N^{p-1}}{p!}$ examples σ^a with $N = 512$ components each, and the simulation results are then averaged over $L = 100$ student patterns. The simulation data is sometimes systematically shifted up with respect to the saddle-point solution. This difference is notably visible on the central plot, right after the fall from eR to gR when $\alpha = 3$.

349 that there is no algorithm with random initial conditions that can find such a ferromagnetic
 350 fixed point in a tractable amount of time [61, 62]. That kind of metastable region was thus
 351 given the name *hard phase* [61, 63]. In summary, we expect that $p \geq 3$ models in the gR phase
 352 can only recover partially corrupted patterns whereas $p = 2$ can recover them entirely.

353 Fig. (4) shows results from Monte Carlo simulations with $p = 3$, where L replicas of
 354 the student pattern $\{\xi^b\}_{b=1}^L$ are initialized to the teacher pattern ξ^* corrupted by some
 355 Rademacher noise ε . In other words, the initial values of ξ_i^b are sampled from the distri-
 356 bution $(1 - \varepsilon) \delta(\xi_i - \xi_i^*) + \frac{\varepsilon}{2} [\delta(\xi_i + 1) + \delta(\xi_i - 1)]$ with $\varepsilon \in [0, 1]$. The value of ε is tuned
 357 so that the simulations start relatively close to the saddle-point solutions. As explained pre-
 358 viously, gR is a hard phase, so this initialization is necessary to make ξ^b converge to gR in a
 359 reasonable amount of time. Additionally, it is also used to make ξ^b converge to the IR phase
 360 rather than the P phase when desired. Once the simulations are over, the overlaps are averaged
 361 over all L replicas. If we fix $\varepsilon = 0$, then the simulations generally converge to the IR phase when
 362 it is a fixed point. If instead we initialize them to the saddle-point solutions by handpicking ε ,
 363 then they stay near the initial overlaps. In either case, the simulations converge to eR when it is
 364 globally stable. Some simulation data points might be systematically shifted up with respect to
 365 the saddle-point solutions. However, this difference decreases with the system size N , so finite
 366 size effects seem sufficient to explain it (see Fig. 9 in Appendix F). Overall, the Monte-Carlo
 367 simulations are in very good agreement with the $p = 3$ overlap landscape obtained by solving
 368 the saddle-point equations numerically.

369 4.4 Inference temperature vs dataset noise

370 In the two next Sections, we will discuss the phase diagram when the student is only partially
 371 informed about the teacher generative model, i.e. when the Nishimori conditions do not hold.
 372 We start with the case where $p = p^*$ but $\beta \neq \beta^*$, i.e. the inference temperature T is different
 373 from the dataset noise T^* . As we argued in Section 3.1, the student accurately retrieves ξ^*
 374 when $T^* < T_{\text{crit}}$. On the other hand, we must solve the saddle-points equations (see Eqs. 5) to
 375 study $T^* > T_{\text{crit}}$.

376 We show the phase diagram of this region on Fig. (5). At high inference temperature T , the
 377 situation is similar to Fig. (2): retrieval is possible if the data load α is sufficiently large, but

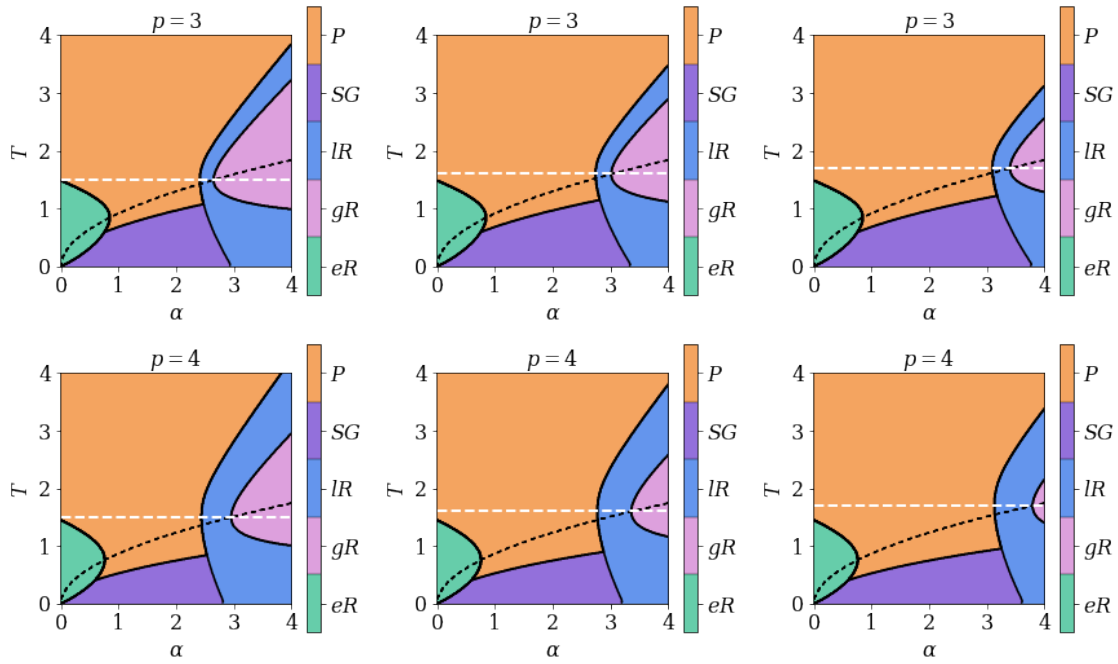


Figure 5: RS phase diagrams of inverse models with $p^* = p$ and fixed β^* . The top and bottom rows of plots respectively have $p^* = p = 3$ and $p^* = p = 4$. In the same way, the left, central and right columns correspond to $T^* = 1.5$, $T^* = 1.6$ and $T^* = 1.7$. Accurate pattern retrieval is not possible in the paramagnetic phase (P), in the spin-glass phase (SG) or in the example retrieval phase (eR), but it is possible in the local retrieval phase (LR) and in the global retrieval phase (gR). The ferromagnetic fixed point corresponding to accurate pattern retrieval is globally stable in the gR phase, but locally stable in the LR phase. Conversely, the SG fixed point is always locally stable and leads the student to a frozen spurious signal. The white dashed line indicates the Nishimori line $\beta^* = \beta$. The black dashed lined is the gR phase boundary on the Nishimori line. As explained in Section 4.3, we expect it to overlap the exact SG phase transition.

378 the paramagnetic phase is always locally stable. The situation is different when the inference
 379 temperature is low. In that case, there are two phases that we did not see for $\beta = \beta^*$: the
 380 inaccurate eR phase and the SG phase. When α is relatively small, the student falls in the
 381 inaccurate eR phase. In this regime, it has finite overlap with one of the noisy examples and
 382 cannot retrieve the signal ξ^* . When α is larger, the interference among the noisy examples
 383 prevents the student to be aligned with them. In this regime, the SG phase, the student locally
 384 converge to spurious patterns that are uncorrelated with the signal.

385 Accurate pattern retrieval is only possible in the LR and gR phases where α is so large that
 386 the student can gather enough information from the dataset to become very close to ξ^* . The
 387 phase diagrams indicate that pattern retrieval is optimal on the Nishimori line in the sense that
 388 $\beta = \beta^*$ is the inverse temperature where the student needs the least examples to recover ξ^* .
 389 In other words, the student's performance is non-monotonic in T and peaks at $T = T^*$. These
 390 properties were also observed in the teacher-student setting of the $p = 2$ Hopfield network [38].

391 Contrary to what one would expect to see on the exact phase diagram [45, 46], the Nishimori
 392 line $T = T^*$ does not cross a triple point on the RS phase diagram. The issue is that the RS
 393 phase diagram is not exact outside of the Nishimori line. In particular, the SG phase boundary
 394 is not exact. Outside of the retrieval regime, the free entropy of the inverse model is the same

395 as the direct model. Since the transition towards \mathbf{gR} of the inverse model on the Nishimori
 396 line overlaps the exact $\mathbf{P-SG}$ transition of the direct model (see Section 4.3), we deduce that
 397 it must also overlap the exact $\mathbf{P-SG}$ transition of the *inverse* model outside of the \mathbf{gR} phase.
 398 Plotting it on the RS phase diagrams, we see that it indeed crosses the Nishimori line and the
 399 \mathbf{gR} phase boundary at the same point, which therefore becomes a triple point, as expected.

400 4.5 Interaction order and noise tolerance

401 So far, we assumed that the student is informed about the interaction order used by the teacher,
 402 i.e. $\mathbf{p} = \mathbf{p}^*$. In this Section, we investigate the role of the student's choice of \mathbf{p} when the task
 403 is to learn from a dataset sampled by a 2-body Hopfield network, i.e. $\mathbf{p}^* = 2$. We study two
 404 different non trivial scalings regimes of M and β^* that make pattern inference possible (see
 405 Appendix D).

406 4.5.1 Large noise scaling

We first consider a large noise scaling where $\beta^* \sim \mathcal{O}(N^{2/p-1})$ and $M \sim \mathcal{O}(N^{p-1})$, such that

$$\alpha = \frac{Mp!}{N^{p-1}} \quad \text{and} \quad \lambda = \frac{[\beta^*]^{p/2}}{(p/2)!} N^{p/2-1}$$

407 are finite. In this scaling, a $\mathbf{p} \geq 3$ network requires $\mathcal{O}(N^{p-2})$ more training examples than a
 408 $\mathbf{p} = 2$ network with finite load $\gamma = \frac{M}{N}$, but also has a higher tolerance to teacher noise. For
 409 instance, a student with $\mathbf{p} = 4$ interactions is able to retrieve the pattern of a teacher with
 410 $T^* \sim \mathcal{O}(N^{1/2})$ noise when it is shown enough examples $M \sim \mathcal{O}(N^3)$ to be in the \mathbf{gR} phase
 411 (see Fig. 6).

412 $\mathcal{O}(N^{1/2})$ noise tolerance was also observed in the $\mathbf{p} = 4$ direct model, where it is a
 413 consequence of the redundancy stemming from storing $\mathcal{O}(N)$ memories rather than the $\mathcal{O}(N^3)$
 414 needed to saturate the storage capacity [64]. Our $\mathbf{p} = 4$ inverse model exploits a different
 415 kind of redundancy by learning from $\mathcal{O}(N^3)$ examples whereas $\mathbf{p} = 2$ only needs $\mathcal{O}(N)$. In
 416 other terms, both storing extensively less memories than the maximum allowed amount and
 417 generating extensively more examples than the minimum required amount provide enough
 418 redundancy to recover a pattern muddled in an extensive amount of noise. In both cases, there is
 419 an $\mathcal{O}(N^2)$ gap between the number of patterns used in the noise-tolerant and noise-susceptible
 420 regimes. Going beyond $\mathbf{p} = 4$, the inverse model has $\mathcal{O}(N^{1-2/p})$ noise tolerance as a function
 421 of \mathbf{p} . In particular, our theory predicts that the tolerance saturates at $T^* \sim \mathcal{O}(N)$ as $\mathbf{p} \rightarrow \infty$,
 422 but at the cost of using an intractable number of examples. This behavior is different from
 423 the $\mathcal{O}(N^{1/2-p/4})$ tolerance of the direct \mathbf{p} -body model in the noisy-learning regime studied
 424 in [65]. In other terms, the dataset noise that we are facing is of a different nature than the
 425 learning noise of [65]. In any case, it is interesting that both the direct and inverse models
 426 are able to tolerate an extensive amount of noise. Overall, our results suggest that it could be
 427 advantageous to use a student network with a relatively large \mathbf{p} to learn from a large but noisy
 428 dataset when the \mathbf{p}^* of the teacher generative model is unknown.

429 An unavoidable drawback of large teacher noise is that it always lead to uncorrelated
 430 examples, which makes accurate example retrieval impossible. Instead, it is replaced by the
 431 inaccurate example retrieval phase where the student has finite overlap \mathbf{m} with a noisy example
 432 generated by the teacher but no overlap with the signal (see Fig. 6). Depending on T and α ,
 433 this phase can be either globally stable or locally stable. For the sake of clarity, we plot only the
 434 globally stable phase on our phase diagram in Fig. (6). The locally stable phase is arguably less
 435 important to plot because it is identical to the locally stable ferromagnetic phase previously
 436 reported in the direct model when assuming replica symmetry (see [33] and Fig. 1).

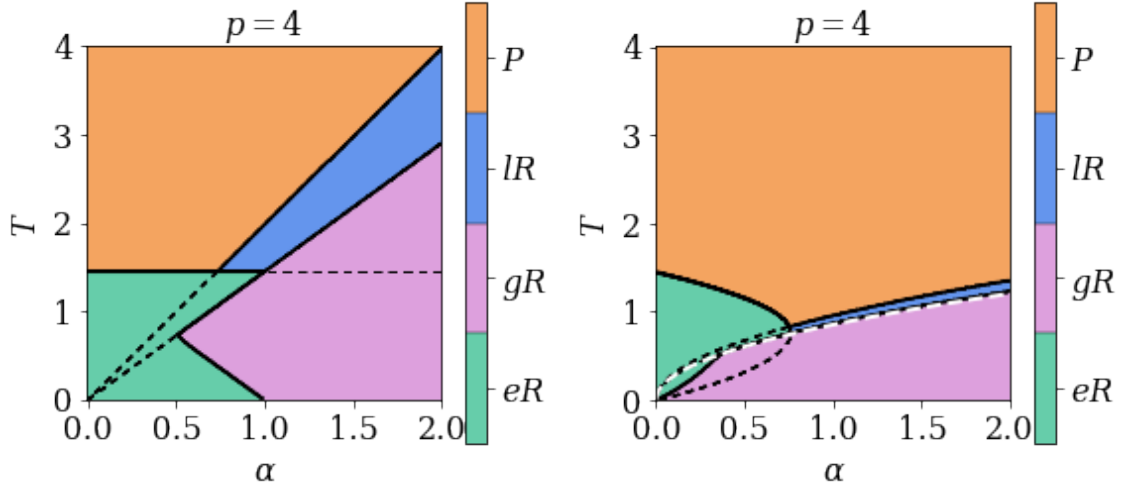


Figure 6: RS phase diagrams of inverse models with $p^* = 2$ and $p = 4$. The left plot is for $\alpha = \frac{M(p/2+1)!}{N^{p/2}}$, and $\beta^* = 1 - \frac{1}{\sqrt{2}}$ such that $\eta = 1$ and the right plot is for $\alpha = \frac{Mp!}{N^{p-1}}$ and $\beta^* = \sqrt{\frac{2\lambda}{N}}$ with $\lambda = \beta$. Accurate pattern retrieval is not possible in the paramagnetic phase (P) or in the example retrieval phase (eR), but it is possible in the local retrieval phase (IR) and in the global retrieval phase (gR). The ferromagnetic fixed point corresponding to accurate pattern retrieval is globally stable in the gR phase, but locally stable in the IR phase. The black dashed lines mark the metastable continuation of the eR , IR and gR phase boundaries through neighboring phases with a larger free entropy. The paramagnetic total entropy becomes negative below the white dashed line drawn on the right plot. However, the paramagnetic phase is no longer globally stable at that temperature.

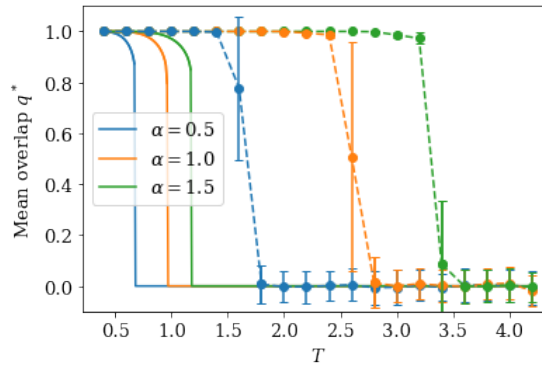


Figure 7: Monte-Carlo simulations (dashed lines) and RS saddle-point solutions (full lines) of the inverse model in the large-noise scaling with $p^* = 2$ and $p = 4$. The teacher generates $M = \frac{\alpha N^{p-1}}{p!}$ examples σ^a with $N = 256$ components each, and the simulation results are then averaged over $L = 100$ student patterns. The student patterns are all initialized to ξ^* .

437 Given $m = 0$, the free entropy of the inverse model with $p \geq 3$, $p^* = 2$ and $\beta = \lambda$ is
 438 the same as on the Nishimori line (see Eq. 5 and Appendix D). As a direct consequence, the
 439 total entropy is positive outside of the eR phase (see Fig. 6). Additionally, the $p^* = 2$, $p \geq 3$
 440 phase diagrams with $\beta \neq \lambda$ are identical to the $p = p^*$ phase diagrams with $\beta \neq \beta^*$, which
 441 suggests that $\beta = \lambda$ is optimal for $p^* = 2$, $p \geq 3$ in the same sense as $\beta = \beta^*$ is optimal for
 442 $p = p^*$ (see Fig. 5). Monte-Carlo simulations confirm that a student with $p \geq 3$ is able to
 443 retrieve the pattern of a teacher with $p = 2$ and $T^* \sim \mathcal{O}(N^{1/2})$ (see Fig. 7). However, the lR
 444 phase transition is at a higher T in the simulations than on the $\beta = \lambda$ RS phase diagram (see
 445 Fig. 5), which means that RSB is necessary to describe it accurately. One could check where
 446 replica symmetry holds by evaluating the stability of the RS saddle point throughout the phase
 447 diagram.

448 4.5.2 Finite noise scaling

We also consider a different scaling regime where $\beta^* \sim \mathcal{O}(1)$ and $M \sim \mathcal{O}(N^{p/2})$, such that

$$\alpha = \frac{M(p/2 + 1)!}{N^{p/2}}$$

449 is finite. In this finite-noise scaling, $p \geq 3$ requires $\mathcal{O}(N^{p/2-1})$ more training examples than
 450 $p = 2$, which is a lot less than the first scaling. For instance, a student with $p = 4$ needs $\mathcal{O}(N^2)$
 451 examples to retrieve ξ^* . As before, the phase transitions are all first order, the overlap q^* stays
 452 high throughout the gR and lR phase of $p = 4$ and gR is a hard phase. The saddle-point
 453 equations (see Eqs. 6) are free from the pattern interference term $\sqrt{ar}x$ present in their
 454 $p^* = p$ counterparts (see Eqs. 5) until β^* becomes so small that it approaches $\mathcal{O}(N^{2/p-1})$.
 455 Therefore, contrary to $p^* = p = 2$, the network is never in the SG phase. Practically, it means
 456 that $p \geq 3$ gives more freedom than $p = 2$ for tuning β and α . The only remaining restriction
 457 is that choosing α and T too small puts the network into the inaccurate eR phase resulting
 458 from the kz term (see Fig. 6). The saddle point equations can be derived without the RS
 459 ansatz because they do not involve q and r . Consequently, we expect them to yield an exact
 460 solution. Like on the Nishimori line, the total entropy of the paramagnetic phase is always
 461 positive, which is consistent with the solution being exact.

462 4.6 Robustness against adversarial attacks

463 Inverse models with $p^* = 2$ and $p \geq 3$ offer an opportunity to study adversarial attacks in a
 464 simple setting because their phase diagrams have regions where the signal retrieval phases (gR
 465 and lR) overlap with the inaccurate eR phase. Recall that, in the lR phase, a noisy student
 466 pattern ξ either converges to ξ^* or falls in the paramagnetic phase, depending on the amount
 467 of noise that ξ contains initially. The quantity of noise needed to prevent pattern retrieval
 468 becomes smaller as one approaches the lR to P phase transition and the basin of attraction
 469 of lR shrinks. Similarly, in the region of inaccurate eR where signal retrieval is metastable,
 470 patterns ξ that are corrupted by replacing some of their entries ξ_i by the components σ_i^a of an
 471 example σ^a may converge to σ^a when enough entries are replaced. The fraction ε of entries
 472 that need to be replaced becomes smaller as the basin of attraction of inaccurate eR expands
 473 and overtakes that of signal retrieval. In practice, an adversary can use this strategy to trick the
 474 student into converging to a pattern other than ξ^* . This scenario is similar to an adversarial
 475 attack targeting the input of K & H's dense Hopfield network model because the student pattern
 476 ξ plays a similar role in the inverse model as the test data in K & H's dense Hopfield networks
 477 (see Fig. 3, Section 4.3 and Appendix A). In that analogy, the examples σ are acting like the
 478 neural network weights rather than taking the role of the training data.

479 We will now investigate what values of the perturbation size ε are a threat by deriving a
 480 formula for the largest ε such that the student converges to the signal at zero temperature. This
 481 largest ε will be denoted ε^* , and we expect it to be a good measure of adversarial robustness.
 482 The saddle-point equations with $T = \mathbf{0}$ indicate that the student converges to one of the signal
 483 retrieval phases if and only if $k < \eta\alpha r^*$ (see Eqs. 6). Sampling the initial conditions of ξ_i
 484 from $(1 - \varepsilon) \delta(\xi_i - \xi_i^*) + \varepsilon \delta(\xi_i - \sigma_i^a)$ with $\varepsilon \in [0, 1]$, we get

$$r^* = p \left[\frac{1}{N} \sum_{i=1}^{(1-\varepsilon)N} \xi_i^* \xi_i^* + \frac{1}{N} \sum_{i=1}^{\varepsilon N} \xi_i^* \sigma_i^a \right]^{p-1},$$

$$k = p \left[\frac{1}{N} \sum_{i=1}^{(1-\varepsilon)N} \xi_i^* \sigma_i^a + \frac{1}{N} \sum_{i=1}^{\varepsilon N} \sigma_i^a \sigma_i^a \right]^{p-1}.$$

485 By the law of large numbers, $\frac{1}{\varepsilon N} \sum_{i=1}^{\varepsilon N} \xi_i^* \sigma_i^a$ and $\frac{1}{(1-\varepsilon)N} \sum_{i=1}^{(1-\varepsilon)N} \xi_i^* \sigma_i^a$ are both typically close
 486 to $m^* = \frac{1}{N} \sum_i \xi_i^* \sigma_i^a \approx \mathbf{0}$ as $N \rightarrow \infty$. If we take σ^a to be a typical example, then r^* and k
 487 reduce to

$$r^* \approx p (1 - \varepsilon)^{p-1}$$

$$k \approx p \varepsilon^{p-1}.$$

488 Substituting these expressions back in $k < \eta\alpha r^*$ yields

$$\varepsilon^{p-1} < \eta\alpha (1 - \varepsilon)^{p-1}$$

$$\varepsilon < \frac{[\eta\alpha]^{\frac{1}{p-1}}}{[\eta\alpha]^{\frac{1}{p-1}} + 1}.$$

489 In other terms, the inverse model with $p^* = 2$ and even $p \geq 3$ is resistant to adversarial attacks
 490 of size $\varepsilon^* = \frac{[\eta\alpha]^{\frac{1}{p-1}}}{[\eta\alpha]^{\frac{1}{p-1}} + 1}$ and smaller. For $p = 4$, ε^* is in good agreement with Monte-Carlo
 491 simulations of the inverse model corrupted by a typical example (see Fig. 8). This comparison
 492 is good evidence that our solution of the finite-noise scaling is indeed exact. Additionally, ε^* is
 493 a decent approximation of empirical robustness even when the inverse model is corrupted by
 494 the example that has the largest overlap with ξ^* . A similar construction with the perturbation
 495 sampled uniformly at random gives $k \sim \mathcal{O}(N^{1/2-p/2}) \approx \mathbf{0}$, so adversarial attacks are much
 496 more efficient at fooling the model than random noise. Just like adversarial attacks targeting
 497 more complicated neural networks [40, 41], our example-based attack can be hard to detect
 498 at low ε because a few adversarially perturbed entries ξ_i do not look very different from
 499 a low amount of meaningless noise. Moreover, ε^* grows monotonically with α , which is
 500 consistent with the common observation that larger neural networks are also more adversarially
 501 robust [43, 66–71]. At first glance, this effect can be counter-intuitive because adversarial
 502 vulnerability looks like a form of overfitting [42]. In our model, however, all examples work
 503 together to stabilize the \mathbf{IR} phase, and the best way to push the student into the \mathbf{eR} phase is
 504 to perturb it with a single example. Therefore, it is not surprising that increasing α makes
 505 the student more robust. We recall that the examples σ are a feature-based representation
 506 of ξ^* . Interestingly, it means that the underlying mechanism of our example-based attack
 507 is conceptually similar to gradient-based attacks targeting many common types of neural
 508 networks [42]. In fact, gradient-based attacks find features stored in neural network weights

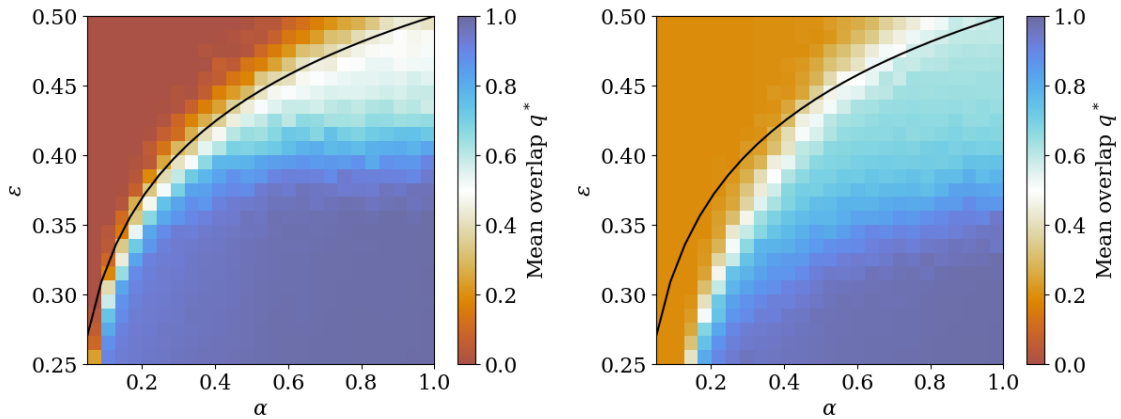


Figure 8: Monte-Carlo simulations of the overlap q^* as a function of α and adversarial attack size ϵ in the inverse model with $p^* = 2$, $\beta^* = 1 - \frac{1}{\sqrt{2}}$, $p = 4$, $\beta = \infty$ and $N = 1024$. The simulation results are averaged over $L = 100$ student patterns. On the left plot, the inverse model is corrupted by an example σ^a that has a small overlap with ξ^* in absolute value. On the right plot, it is corrupted by the example that has the largest overlap with ξ^* . The black line $\epsilon^* = \frac{\alpha^{1/3}}{\alpha^{1/3}+1}$ is our analytical formula for the largest adversarial perturbation ϵ such that the student retrieves ξ^* rather than the example σ^a .

509 and add them to the data in order to fool the network [42, 72–74]. It would be interesting to
 510 investigate, both empirically and theoretically, if only a small number of weights are involved
 511 in constructing these adversarial attacks. If it is the case, it could explain why larger neural
 512 networks are often more robust. In general, we expect this kind of one-example attack to be
 513 possible in any region of signal retrieval that overlaps with the inaccurate eR phase. Using
 514 $p \neq p^*$ may not be a necessary ingredient of adversarial vulnerability in more general models
 515 with other sources of mismatch, but in our case it ensures that the signal retrieval phases
 516 intersect the inaccurate eR phase. Conversely, the accurate eR phase is by definition robust to
 517 adversarial attacks since retrieving an example σ^a is the same as recovering ξ^* . This distinction
 518 clarifies why the dense Hopfield networks designed by K & H are adversarially robust in the
 519 prototype phase despite being adversarially vulnerable in the feature phase. In fact, K & H
 520 observed that adversarial attacks are unsuccessful in the prototype phase specifically because
 521 they retrieve stored examples that are semantically meaningful [37]. In summary, our model
 522 yields two main results concerning adversarial examples. First of all, it suggests a reason why
 523 large feature-based neural networks are more adversarially robust than smaller ones. Second
 524 of all, it clarifies why dense Hopfield networks are much more robust in the prototype phase
 525 than in the feature phase.

526 5 Conclusion

527 In this work, we derive the exact phase diagram of the p -dense networks in the teacher-
 528 student setting [16, 17, 30, 38]. On the Nishimori line, we find an example retrieval phase (eR)
 529 and a global retrieval phase (gR) reminiscent of the prototype and feature regimes observed
 530 empirically in dense Hopfield networks [26]. We show that the phase transition towards gR of
 531 the inverse model overlaps the paramagnetic to spin-glass (P - SG) transition of the direct model,
 532 which allows us to locate the P - SG transition much more precisely than before [30, 33]. On
 533 the other hand, we discover that inverse models outside of the Nishimori line are able to resist

534 an extensive amount of noise. In fact, a student with $p \geq 3$ is able to learn from a teacher with
535 $p^* = 2$ even when the teacher's inverse temperature β^* is as low as $\mathcal{O}(N^{2/p-1})$. Moreover, such
536 a student is immune to pattern interference until β^* reaches $\mathcal{O}(N^{2/p-1})$. In this setting, we
537 derive a formula measuring the adversarial robustness of the student with $p \geq 3$ and $T = 0$. We
538 then use this formula to describe how making a neural network larger can potentially increase
539 its robustness to adversarial attacks constructed with only a few learned weights [43, 66–71].
540 Our model also clarifies why the prototype phase of dense Hopfield networks is adversarially
541 robust [37]. We compare our key results against Monte-Carlo simulations.

542 Dense networks with exponential interactions have been argued to be the $p \rightarrow \infty$ limit of
543 the p -body models [75]. It would be interesting to see if they can achieve $\mathcal{O}(N)$ noise tolerance
544 at the cost of an exponential number of training examples. More generally, studying exponential
545 models in the teacher-student setting would be an interesting extension of this work and could
546 be used to complement existing studies of the direct model [75, 76]. A caveat of our model is
547 that the teacher has only one pattern. In fact, we would need to use a teacher with at least
548 two patterns to describe more completely the kind of adversarial attack aiming to misclassify
549 data. It should be possible to study this kind of teacher by using an approach similar to [77].
550 In particular, [16] and [77] argue that the performance of restricted Boltzmann machines with
551 a finite number P of i.i.d. planted patterns is independent of P in the teacher-student setting.
552 It would be interesting to investigate whether this characteristic also holds for p -body dense
553 networks. On the practical side, we highlight the untapped benefits of using p -body models to
554 either resist an extensive amount of noise in the feature phase or improve adversarial robustness
555 in the prototype phase. Overall, we stress that further investigations of dense Hopfield networks
556 could unlock their true potential.

557 **Funding information** This work was partially supported by project SERICS (PE00000014)
558 under the MUR National Recovery and Resilience Plan funded by the European Union - NextGen-
559 erationEU. The work was also supported by the project PRIN22TANTARI "Statistical Mechanics
560 of Learning Machines: from algorithmic and information-theoretical limits to new biolog-
561 ically inspired paradigms" 20229T9EAT – CUP J53D23003640001. DT also acknowledges
562 GNFM-Indam.

563 **Code availability** The figures can be reproduced using the code available on [this public](#)
564 [Github repository](#).

565 References

- 566 [1] J. J. Hopfield, *Neural networks and physical systems with emergent collective computa-*
567 *tional abilities.*, Proceedings of the National Academy of Sciences **79**(8), 2554 (1982),
568 doi:[10.1073/pnas.79.8.2554](https://doi.org/10.1073/pnas.79.8.2554), <https://www.pnas.org/doi/pdf/10.1073/pnas.79.8.2554>.
- 569 [2] D. J. Amit, H. Gutfreund and H. Sompolinsky, *Spin-glass models of neural networks*, Phys.
570 Rev. A **32**, 1007 (1985), doi:[10.1103/PhysRevA.32.1007](https://doi.org/10.1103/PhysRevA.32.1007).
- 571 [3] D. O. Hebb, *The organization of behavior: A neuropsychological theory*, Psychology press,
572 ISBN 9781410612403, doi:<https://doi.org/10.4324/9781410612403> (2005).
- 573 [4] T. M. Cover, *Geometrical and statistical properties of systems of linear inequalities with*
574 *applications in pattern recognition*, IEEE Transactions on Electronic Computers **EC-14**(3),
575 326 (1965), doi:[10.1109/PGEC.1965.264137](https://doi.org/10.1109/PGEC.1965.264137).

- 576 [5] D. J. Amit, H. Gutfreund and H. Sompolinsky, *Storing infinite numbers of pat-*
577 *terns in a spin-glass model of neural networks*, Phys. Rev. Lett. **55**, 1530 (1985),
578 doi:[10.1103/PhysRevLett.55.1530](https://doi.org/10.1103/PhysRevLett.55.1530).
- 579 [6] E. Agliari, A. Barra, A. Galluzzi, F. Guerra and F. Moauro, *Multitasking associative networks*,
580 Phys. Rev. Lett. **109**, 268101 (2012), doi:[10.1103/PhysRevLett.109.268101](https://doi.org/10.1103/PhysRevLett.109.268101), [1111.5191](https://arxiv.org/abs/1111.5191).
- 581 [7] E. Agliari, A. Annibale, A. Barra, A. C. C. Coolen and D. Tantari, *Immune networks: multi-*
582 *tasking capabilities at medium load*, Journal of Physics A: Mathematical and Theoretical
583 **46**(33), 335101 (2013), doi:[10.1088/1751-8113/46/33/335101](https://doi.org/10.1088/1751-8113/46/33/335101), [1302.7259](https://arxiv.org/abs/1302.7259).
- 584 [8] E. Agliari, A. Annibale, A. Barra, A. C. C. Coolen and D. Tantari, *Immune networks:*
585 *multitasking capabilities near saturation*, Journal of Physics A Mathematical General
586 **46**(41), 415003 (2013), doi:[10.1088/1751-8113/46/41/415003](https://doi.org/10.1088/1751-8113/46/41/415003), [1305.5936](https://arxiv.org/abs/1305.5936).
- 587 [9] P. Sollich, D. Tantari, A. Annibale and A. Barra, *Extensive parallel processing on scale-free*
588 *networks*, Phys. Rev. Lett. **113**, 238106 (2014), doi:[10.1103/PhysRevLett.113.238106](https://doi.org/10.1103/PhysRevLett.113.238106),
589 [1404.3654](https://arxiv.org/abs/1404.3654).
- 590 [10] E. Agliari, A. Annibale, A. Barra, A. C. C. Coolen and D. Tantari, *Retrieving infinite numbers*
591 *of patterns in a spin-glass model of immune networks*, Europhysics Letters **117**(2), 28003
592 (2017), doi:[10.1209/0295-5075/117/28003](https://doi.org/10.1209/0295-5075/117/28003), [1305.2076](https://arxiv.org/abs/1305.2076).
- 593 [11] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Retrieval capabilities*
594 *of hierarchical networks: From dyson to hopfield*, Phys. Rev. Lett. **114**, 028103 (2015),
595 doi:[10.1103/PhysRevLett.114.028103](https://doi.org/10.1103/PhysRevLett.114.028103), [1407.5019](https://arxiv.org/abs/1407.5019).
- 596 [12] E. Agliari, D. Migliozi and D. Tantari, *Non-convex Multi-species Hopfield Models*, Journal of
597 Statistical Physics **172**(5), 1247 (2018), doi:[10.1007/s10955-018-2098-6](https://doi.org/10.1007/s10955-018-2098-6), [1807.03609](https://arxiv.org/abs/1807.03609).
- 598 [13] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Hierarchical neural*
599 *networks perform both serial and parallel processing*, Neural Networks **66**, 22 (2015),
600 doi:<https://doi.org/10.1016/j.neunet.2015.02.010>, [1409.0227](https://arxiv.org/abs/1409.0227).
- 601 [14] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Metastable states in*
602 *the hierarchical Dyson model drive parallel processing in the hierarchical Hopfield network*,
603 Journal of Physics A Mathematical General **48**(1), 015001 (2015), doi:[10.1088/1751-](https://doi.org/10.1088/1751-8113/48/1/015001)
604 [8113/48/1/015001](https://doi.org/10.1088/1751-8113/48/1/015001), [1407.5176](https://arxiv.org/abs/1407.5176).
- 605 [15] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Topological properties of*
606 *hierarchical networks*, Phys. Rev. E **91**, 062807 (2015), doi:[10.1103/PhysRevE.91.062807](https://doi.org/10.1103/PhysRevE.91.062807),
607 [1412.5918](https://arxiv.org/abs/1412.5918).
- 608 [16] A. Barra, G. Genovese, P. Sollich and D. Tantari, *Phase transitions in re-*
609 *stricted boltzmann machines with generic priors*, Phys. Rev. E **96**, 042156 (2017),
610 doi:[10.1103/PhysRevE.96.042156](https://doi.org/10.1103/PhysRevE.96.042156), [1612.03132](https://arxiv.org/abs/1612.03132).
- 611 [17] A. Barra, G. Genovese, P. Sollich and D. Tantari, *Phase diagram of restricted boltzmann*
612 *machines and generalized hopfield networks with arbitrary priors*, Phys. Rev. E **97**, 022310
613 (2018), doi:[10.1103/PhysRevE.97.022310](https://doi.org/10.1103/PhysRevE.97.022310), [1702.05882](https://arxiv.org/abs/1702.05882).
- 614 [18] A. Barra, P. Contucci, E. Mingione and D. Tantari, *Multi-Species Mean Field Spin Glasses.*
615 *Rigorous Results*, Annales Henri Poincaré; **16**(3), 691 (2015), doi:[10.1007/s00023-](https://doi.org/10.1007/s00023-014-0341-5)
616 [014-0341-5](https://doi.org/10.1007/s00023-014-0341-5), [1307.5154](https://arxiv.org/abs/1307.5154).

- 617 [19] E. Agliari, A. Barra, C. Longo and D. Tantari, *Neural Networks Retrieving Boolean Pat-*
618 *terns in a Sea of Gaussian Ones*, Journal of Statistical Physics **168**(5), 1085 (2017),
619 doi:[10.1007/s10955-017-1840-9](https://doi.org/10.1007/s10955-017-1840-9), [1703.05210](https://arxiv.org/abs/1703.05210).
- 620 [20] A. Barra, G. Genovese, F. Guerra and D. Tantari, *How glassy are neural networks?*, Journal of
621 Statistical Mechanics: Theory and Experiment **2012**(7), 07009 (2012), doi:[10.1088/1742-](https://doi.org/10.1088/1742-5468/2012/07/P07009)
622 [5468/2012/07/P07009](https://doi.org/10.1088/1742-5468/2012/07/P07009), [1205.3900](https://arxiv.org/abs/1205.3900).
- 623 [21] G. Genovese and D. Tantari, *Legendre equivalences of spherical Boltzmann machines*,
624 Journal of Physics A Mathematical General **53**(9), 094001 (2020), doi:[10.1088/1751-](https://doi.org/10.1088/1751-8121/ab6b92)
625 [8121/ab6b92](https://doi.org/10.1088/1751-8121/ab6b92), [1910.14559](https://arxiv.org/abs/1910.14559).
- 626 [22] J. Rocchi, D. Saad and D. Tantari, *High storage capacity in the Hopfield model with*
627 *auto-interactions—stability analysis*, Journal of Physics A Mathematical General **50**(46),
628 465001 (2017), doi:[10.1088/1751-8121/aa8fd7](https://doi.org/10.1088/1751-8121/aa8fd7), [1704.07741](https://arxiv.org/abs/1704.07741).
- 629 [23] H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, L. Gruber, M. Holzleitner, T. Adler,
630 D. Kreil, M. K. Kopp, G. Klambauer, J. Brandstetter *et al.*, *Hopfield networks is all you need*,
631 In *International Conference on Learning Representations*, doi:[10.48550/arXiv.2008.02217](https://doi.org/10.48550/arXiv.2008.02217)
632 (2021), [2008.02217](https://arxiv.org/abs/2008.02217).
- 633 [24] M. Widrich, B. Schäfl, M. Pavlović, H. Ramsauer, L. Gruber, M. Holzleitner, J. Brandstetter,
634 G. K. Sandve, V. Greiff, S. Hochreiter and G. Klambauer, *Modern hopfield networks and*
635 *attention for immune repertoire classification*, In H. Larochelle, M. Ranzato, R. Hadsell,
636 M. Balcan and H. Lin, eds., *Advances in Neural Information Processing Systems*, vol. 33,
637 pp. 18832–18845. Curran Associates, Inc., doi:[10.48550/arXiv.2007.13505](https://doi.org/10.48550/arXiv.2007.13505) (2020),
638 [2007.13505](https://arxiv.org/abs/2007.13505).
- 639 [25] D. Krotov and J. J. Hopfield, *Large associative memory problem in neurobiol-*
640 *ogy and machine learning*, In *International Conference on Learning Representations*,
641 doi:[10.48550/arXiv.2008.06996](https://doi.org/10.48550/arXiv.2008.06996) (2021), [2008.06996](https://arxiv.org/abs/2008.06996).
- 642 [26] D. Krotov and J. J. Hopfield, *Dense associative memory for pattern recognition*, In D. Lee,
643 M. Sugiyama, U. Luxburg, I. Guyon and R. Garnett, eds., *Advances in Neural Information*
644 *Processing Systems*, vol. 29. Curran Associates, Inc., doi:[10.48550/arXiv.1606.01164](https://doi.org/10.48550/arXiv.1606.01164)
645 (2016), [1606.01164](https://arxiv.org/abs/1606.01164).
- 646 [27] H. H. Chen, Y. C. Lee, G. Z. Sun, H. Y. Lee, T. Maxwell and C. L. Giles, *High order*
647 *correlation model for associative memory*, AIP Conference Proceedings **151**(1), 86 (1986),
648 doi:[10.1063/1.36224](https://doi.org/10.1063/1.36224), [https://pubs.aip.org/aip/acp/article-pdf/151/1/86/12091820/](https://pubs.aip.org/aip/acp/article-pdf/151/1/86/12091820/86_1_online.pdf)
649 [86_1_online.pdf](https://pubs.aip.org/aip/acp/article-pdf/151/1/86/12091820/86_1_online.pdf).
- 650 [28] D. Psaltis and C. H. Park, *Nonlinear discriminant functions and associative memories*, AIP
651 Conference Proceedings **151**(1), 370 (1986), doi:[10.1063/1.36241](https://doi.org/10.1063/1.36241), [https://pubs.aip.](https://pubs.aip.org/aip/acp/article-pdf/151/1/370/12091772/370_1_online.pdf)
652 [org/aip/acp/article-pdf/151/1/370/12091772/370_1_online.pdf](https://pubs.aip.org/aip/acp/article-pdf/151/1/370/12091772/370_1_online.pdf).
- 653 [29] P. Baldi and S. S. Venkatesh, *Number of stable points for spin-glasses and neural networks*
654 *of higher orders*, Phys. Rev. Lett. **58**, 913 (1987), doi:[10.1103/PhysRevLett.58.913](https://doi.org/10.1103/PhysRevLett.58.913).
- 655 [30] E. Gardner, *Multiconnected neural network models*, Journal of Physics A: Mathematical
656 and General **20**(11), 3453 (1987), doi:[10.1088/0305-4470/20/11/046](https://doi.org/10.1088/0305-4470/20/11/046).
- 657 [31] L. F. Abbott and Y. Arian, *Storage capacity of generalized networks*, Phys. Rev. A **36**, 5091
658 (1987), doi:[10.1103/PhysRevA.36.5091](https://doi.org/10.1103/PhysRevA.36.5091).

- 659 [32] Horn, D. and Usher, M., *Capacities of multiconnected memory models*, J. Phys. France
660 **49**(3), 389 (1988), doi:[10.1051/jphys:01988004903038900](https://doi.org/10.1051/jphys:01988004903038900).
- 661 [33] L. Albanese, F. Alemanno, A. Alessandrelli and A. Barra, *Replica Symmetry Breaking*
662 *in Dense Hebbian Neural Networks*, Journal of Statistical Physics **189**(2), 24 (2022),
663 doi:[10.1007/s10955-022-02966-8](https://doi.org/10.1007/s10955-022-02966-8), [2111.12997](https://arxiv.org/abs/2111.12997).
- 664 [34] B. Hoover, Y. Liang, B. Pham, R. Panda, H. Strobelt, D. H. Chau, M. J. Zaki
665 and D. Krotov, *Energy Transformer*, arXiv e-prints arXiv:2302.07253 (2023),
666 doi:[10.48550/arXiv.2302.07253](https://doi.org/10.48550/arXiv.2302.07253), [2302.07253](https://arxiv.org/abs/2302.07253).
- 667 [35] B. Hoover, H. Strobelt, D. Krotov, J. Hoffman, Z. Kira and D. H. Chau, *Memory in Plain*
668 *Sight: A Survey of the Uncanny Resemblances between Diffusion Models and Associative*
669 *Memories*, arXiv e-prints arXiv:2309.16750 (2023), doi:[10.48550/arXiv.2309.16750](https://doi.org/10.48550/arXiv.2309.16750),
670 [2309.16750](https://arxiv.org/abs/2309.16750).
- 671 [36] L. Ambrogioni, *In search of dispersed memories: Generative diffusion mod-*
672 *els are associative memory networks*, arXiv e-prints arXiv:2309.17290 (2023),
673 doi:[10.48550/arXiv.2309.17290](https://doi.org/10.48550/arXiv.2309.17290), [2309.17290](https://arxiv.org/abs/2309.17290).
- 674 [37] D. Krotov and J. Hopfield, *Dense Associative Memory Is Robust to Adversarial Inputs*, Neural
675 Computation **30**(12), 3151 (2018), doi:[10.1162/neco_a_01143](https://doi.org/10.1162/neco_a_01143), [1701.00939](https://arxiv.org/abs/1701.00939).
- 676 [38] F. Alemanno, L. Camanzi, G. Manzan and D. Tantari, *Hopfield model with planted patterns:*
677 *A teacher-student self-supervised learning model*, Applied Mathematics and Computation
678 **458**, 128253 (2023), doi:<https://doi.org/10.1016/j.amc.2023.128253>, [2304.13710](https://arxiv.org/abs/2304.13710).
- 679 [39] A. Decelle, S. Hwang, J. Rocchi and D. Tantari, *Inverse problems for structured datasets*
680 *using parallel TAP equations and restricted Boltzmann machines*, Scientific Reports **11**,
681 19990 (2021), doi:[10.1038/s41598-021-99353-2](https://doi.org/10.1038/s41598-021-99353-2), [1906.11988](https://arxiv.org/abs/1906.11988).
- 682 [40] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrndić, P. Laskov, G. Giacinto and F. Roli,
683 *Evasion attacks against machine learning at test time*, In H. Blockeel, K. Kersting, S. Ni-
684 jssen and F. Železný, eds., *Machine Learning and Knowledge Discovery in Databases*, pp.
685 387–402. Springer Berlin Heidelberg, Berlin, Heidelberg, ISBN 978-3-642-40994-3,
686 doi:https://doi.org/10.1007/978-3-642-40994-3_25 (2013).
- 687 [41] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow and R. Fer-
688 gus, *Intriguing properties of neural networks*, arXiv e-prints arXiv:1312.6199 (2013),
689 doi:[10.48550/arXiv.1312.6199](https://doi.org/10.48550/arXiv.1312.6199), [1312.6199](https://arxiv.org/abs/1312.6199).
- 690 [42] I. J. Goodfellow, J. Shlens and C. Szegedy, *Explaining and harnessing adversarial examples*,
691 stat **1050**, arXiv:1412.6572 (2015), doi:[10.48550/arXiv.1412.6572](https://doi.org/10.48550/arXiv.1412.6572), [1412.6572](https://arxiv.org/abs/1412.6572).
- 692 [43] A. Madry, A. Makelov, L. Schmidt, D. Tsipras and A. Vladu, *Towards deep learning models*
693 *resistant to adversarial attacks*, In *International Conference on Learning Representations*,
694 doi:[10.48550/arXiv.1706.06083](https://doi.org/10.48550/arXiv.1706.06083) (2018), [1706.06083](https://arxiv.org/abs/1706.06083).
- 695 [44] A. Muhammad and S.-H. Bae, *A Survey on Efficient Methods for Adversarial Robustness*,
696 IEEE Access **10**, 118815 (2022), doi:[10.1109/ACCESS.2022.3216291](https://doi.org/10.1109/ACCESS.2022.3216291).
- 697 [45] H. Nishimori, *Exact results and critical properties of the ising model with competing interac-*
698 *tions*, Journal of Physics C: Solid State Physics **13**(21), 4071 (1980), doi:[10.1088/0022-](https://doi.org/10.1088/0022-3719/13/21/012)
699 [3719/13/21/012](https://arxiv.org/abs/3719/13/21/012).

- 700 [46] H. Nishimori, *Statistical Physics of Spin Glasses and Information Process-*
701 *ing: An Introduction*, Oxford University Press, ISBN 9780198509417,
702 doi:[10.1093/acprof:oso/9780198509417.001.0001](https://doi.org/10.1093/acprof:oso/9780198509417.001.0001) (2001), [https://academic.oup.com/](https://academic.oup.com/book/5185/book-pdf/54038185/acprof-9780198509400.pdf)
703 [book/5185/book-pdf/54038185/acprof-9780198509400.pdf](https://academic.oup.com/book/5185/book-pdf/54038185/acprof-9780198509400.pdf).
- 704 [47] P. Contucci, C. Giardinà and H. Nishimori, *Spin glass identities and the nishimori line*, In
705 A. B. de Monvel and A. Bovier, eds., *Spin Glasses: Statics and Dynamics*, pp. 103–121.
706 Birkhäuser Basel, Basel, doi:https://doi.org/10.1007/978-3-7643-9891-0_4 (2009),
707 [0805.0754](https://doi.org/10.1007/978-3-7643-9891-0_4).
- 708 [48] Y. Iba, *The Nishimori line and Bayesian statistics*, *Journal of Physics A Mathematical*
709 *General* **32**(21), 3875 (1999), doi:[10.1088/0305-4470/32/21/302](https://doi.org/10.1088/0305-4470/32/21/302), [cond-mat/9809190](https://arxiv.org/abs/cond-mat/9809190).
- 710 [49] P. Charbonneau, *From the replica trick to the replica symmetry breaking technique*, arXiv
711 e-prints arXiv:2211.01802 (2022), doi:[10.48550/arXiv.2211.01802](https://doi.org/10.48550/arXiv.2211.01802), [2211.01802](https://arxiv.org/abs/2211.01802).
- 712 [50] D. Gross and M. Mezard, *The simplest spin glass*, *Nuclear Physics B* **240**(4), 431 (1984),
713 doi:[https://doi.org/10.1016/0550-3213\(84\)90237-2](https://doi.org/10.1016/0550-3213(84)90237-2).
- 714 [51] E. Gardner, *Spin glasses with p-spin interactions*, *Nuclear Physics B* **257**, 747 (1985),
715 doi:[https://doi.org/10.1016/0550-3213\(85\)90374-8](https://doi.org/10.1016/0550-3213(85)90374-8).
- 716 [52] B. Derrida, *Random-energy model: An exactly solvable model of disordered systems*, *Phys.*
717 *Rev. B* **24**, 2613 (1981), doi:[10.1103/PhysRevB.24.2613](https://doi.org/10.1103/PhysRevB.24.2613).
- 718 [53] R. Monasson, *Structural glass transition and the entropy of the metastable states*, *Phys.*
719 *Rev. Lett.* **75**, 2847 (1995), doi:[10.1103/PhysRevLett.75.2847](https://doi.org/10.1103/PhysRevLett.75.2847).
- 720 [54] A. Montanari and F. Ricci-Tersenghi, *On the nature of the low-temperature phase in*
721 *discontinuous mean-field spin glasses*, *The European Physical Journal B - Condensed*
722 *Matter and Complex Systems* **33**(3), 339 (2003), doi:[10.1140/epjb/e2003-00174-7](https://doi.org/10.1140/epjb/e2003-00174-7).
- 723 [55] A. Crisanti, L. Leuzzi and T. Rizzo, *Complexity in mean-field spin-glass models: Ising p-spin*,
724 *Phys. Rev. B* **71**, 094202 (2005), doi:[10.1103/PhysRevB.71.094202](https://doi.org/10.1103/PhysRevB.71.094202).
- 725 [56] S. Franz, G. Parisi, M. Sevelev, P. Urbani and F. Zamponi, *Universality of the SAT-UNSAT*
726 *(jamming) threshold in non-convex continuous constraint satisfaction problems*, *SciPost*
727 *Phys.* **2**, 019 (2017), doi:[10.21468/SciPostPhys.2.3.019](https://doi.org/10.21468/SciPostPhys.2.3.019).
- 728 [57] G. G. Roussas, *Contiguity of Probability Measures: Some Applications in*
729 *Statistics*, Cambridge Tracts in Mathematics. Cambridge University Press,
730 doi:[10.1017/CBO9780511804373](https://doi.org/10.1017/CBO9780511804373) (1972).
- 731 [58] D. Achlioptas and A. Coja-Oghlan, *Algorithmic barriers from phase transitions*, In
732 *2008 49th Annual IEEE Symposium on Foundations of Computer Science*, pp. 793–802,
733 doi:[10.1109/FOCS.2008.11](https://doi.org/10.1109/FOCS.2008.11) (2008), [0803.2122](https://arxiv.org/abs/0803.2122).
- 734 [59] F. Krzakala and L. Zdeborová, *Hiding quiet solutions in random constraint satisfaction*
735 *problems*, *Phys. Rev. Lett.* **102**, 238701 (2009), doi:[10.1103/PhysRevLett.102.238701](https://doi.org/10.1103/PhysRevLett.102.238701),
736 [0901.2130](https://arxiv.org/abs/0901.2130).
- 737 [60] L. Zdeborová and F. Krzakala, *Quiet planting in the locked constraint satisfaction problems*,
738 *SIAM Journal on Discrete Mathematics* **25**(2), 750 (2011), doi:[10.1137/090750755](https://doi.org/10.1137/090750755),
739 [0902.4185](https://arxiv.org/abs/0902.4185).

- 740 [61] L. Zdeborová and F. Krzakala, *Statistical physics of inference: thresholds and algorithms*,
741 *Advances in Physics* **65**(5), 453 (2016), doi:[10.1080/00018732.2016.1211393](https://doi.org/10.1080/00018732.2016.1211393), [1511.02476](https://doi.org/10.1080/00018732.2016.1211393).
742
- 743 [62] F. Antenucci, S. Franz, P. Urbani and L. Zdeborová, *Glassy nature of the hard phase in*
744 *inference problems*, *Phys. Rev. X* **9**, 011020 (2019), doi:[10.1103/PhysRevX.9.011020](https://doi.org/10.1103/PhysRevX.9.011020),
745 [1805.05857](https://doi.org/10.1103/PhysRevX.9.011020).
- 746 [63] L. Zdeborová and F. Krzakala, *Phase transitions in the coloring of random graphs*, *Phys.*
747 *Rev. E* **76**, 031131 (2007), doi:[10.1103/PhysRevE.76.031131](https://doi.org/10.1103/PhysRevE.76.031131), [0704.1269](https://doi.org/10.1103/PhysRevE.76.031131).
- 748 [64] E. Agliari, F. Alemanno, A. Barra, M. Centonze and A. Fachechi, *Neural networks with a*
749 *redundant representation: Detecting the undetectable*, *Phys. Rev. Lett.* **124**, 028301 (2020),
750 doi:[10.1103/PhysRevLett.124.028301](https://doi.org/10.1103/PhysRevLett.124.028301), [1911.12689](https://doi.org/10.1103/PhysRevLett.124.028301).
- 751 [65] E. Agliari and G. De Marzo, *Tolerance versus synaptic noise in dense associative memories*,
752 *European Physical Journal Plus* **135**(11), 883 (2020), doi:[10.1140/epjp/s13360-020-](https://doi.org/10.1140/epjp/s13360-020-00894-8)
753 [00894-8](https://doi.org/10.1140/epjp/s13360-020-00894-8), [2007.02849](https://doi.org/10.1140/epjp/s13360-020-00894-8).
- 754 [66] S. Gowal, C. Qin, J. Uesato, T. Mann and P. Kohli, *Uncovering the Limits of Adversarial*
755 *Training against Norm-Bounded Adversarial Examples*, arXiv e-prints arXiv:2010.03593
756 (2020), doi:[10.48550/arXiv.2010.03593](https://doi.org/10.48550/arXiv.2010.03593), [2010.03593](https://doi.org/10.48550/arXiv.2010.03593).
- 757 [67] H. Huang, Y. Wang, S. Erfani, Q. Gu, J. Bailey and X. Ma, *Exploring architectural ingredients*
758 *of adversarially robust deep neural networks*, In M. Ranzato, A. Beygelzimer, Y. Dauphin,
759 P. Liang and J. W. Vaughan, eds., *Advances in Neural Information Processing Systems*,
760 vol. 34, pp. 5545–5559. Curran Associates, Inc., doi:[10.48550/arXiv.2110.03825](https://doi.org/10.48550/arXiv.2110.03825) (2021),
761 [2110.03825](https://doi.org/10.48550/arXiv.2110.03825).
- 762 [68] S. Bubeck, Y. Li and D. M. Nagaraj, *A law of robustness for two-layers neural networks*,
763 In M. Belkin and S. Kpotufe, eds., *Proceedings of Thirty Fourth Conference on Learn-*
764 *ing Theory*, vol. 134 of *Proceedings of Machine Learning Research*, pp. 804–820. PMLR,
765 doi:[10.48550/arXiv.2009.14444](https://doi.org/10.48550/arXiv.2009.14444) (2021), [2009.14444](https://doi.org/10.48550/arXiv.2009.14444).
- 766 [69] S. Bubeck and M. Sellke, *A universal law of robustness via isoperimetry*, In M. Ran-
767 zato, A. Beygelzimer, Y. Dauphin, P. Liang and J. W. Vaughan, eds., *Advances in Neu-*
768 *ral Information Processing Systems*, vol. 34, pp. 28811–28822. Curran Associates, Inc.,
769 doi:[10.48550/arXiv.2105.12806](https://doi.org/10.48550/arXiv.2105.12806) (2021), [2105.12806](https://doi.org/10.48550/arXiv.2105.12806).
- 770 [70] J. Puigcerver, R. Jenatton, C. Riquelme, P. Awasthi and S. Bhojanapalli, *On the adversarial*
771 *robustness of mixture of experts*, In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho
772 and A. Oh, eds., *Advances in Neural Information Processing Systems*, vol. 35, pp. 9660–9671.
773 Curran Associates, Inc., doi:[10.48550/arXiv.2210.10253](https://doi.org/10.48550/arXiv.2210.10253) (2022), [2210.10253](https://doi.org/10.48550/arXiv.2210.10253).
- 774 [71] A. H. Ribeiro and T. B. Schön, *Overparameterized Linear Regression Under*
775 *Adversarial Attacks*, *IEEE Transactions on Signal Processing* **71**, 601 (2023),
776 doi:[10.1109/TSP.2023.3246228](https://doi.org/10.1109/TSP.2023.3246228), [2204.06274](https://doi.org/10.1109/TSP.2023.3246228).
- 777 [72] S. Jetley, N. Lord and P. Torr, *With friends like these, who needs adversaries?*, In
778 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett,
779 eds., *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc.,
780 doi:[10.48550/arXiv.1807.04200](https://doi.org/10.48550/arXiv.1807.04200) (2018), [1807.04200](https://doi.org/10.48550/arXiv.1807.04200).

- 781 [73] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran and A. Madry, *Adversarial examples*
 782 *are not bugs, they are features*, In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-
 783 Buc, E. Fox and R. Garnett, eds., *Advances in Neural Information Processing Systems*,
 784 vol. 32. Curran Associates, Inc., doi:<https://doi.org/10.48550/arXiv.1905.02175> (2019),
 785 [1905.02175](https://doi.org/10.48550/arXiv.1905.02175).
- 786 [74] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner and A. Madry, *Robustness may*
 787 *be at odds with accuracy*, In *International Conference on Learning Representations*,
 788 doi:[10.48550/arXiv.1805.12152](https://doi.org/10.48550/arXiv.1805.12152) (2019), [1805.12152](https://doi.org/10.48550/arXiv.1805.12152).
- 789 [75] M. Demircigil, J. Heusel, M. Löwe, S. Uppgang and F. Vermet, *On a Model of Associative*
 790 *Memory with Huge Storage Capacity*, *Journal of Statistical Physics* **168**(2), 288 (2017),
 791 doi:[10.1007/s10955-017-1806-y](https://doi.org/10.1007/s10955-017-1806-y), [1702.01929](https://doi.org/10.1007/s10955-017-1806-y).
- 792 [76] C. Lucibello and M. Mézard, *The Exponential Capacity of Dense Associative Memories*, arXiv
 793 e-prints arXiv:2304.14964 (2023), doi:[10.48550/arXiv.2304.14964](https://doi.org/10.48550/arXiv.2304.14964), [2304.14964](https://doi.org/10.48550/arXiv.2304.14964).
- 794 [77] T. Hou, K. Y. M. Wong and H. Huang, *Minimal model of permutation symmetry in unsu-*
 795 *perervised learning*, *Journal of Physics A: Mathematical and Theoretical* **52**(41), 414001
 796 (2019), doi:[10.1088/1751-8121/ab3f3f](https://doi.org/10.1088/1751-8121/ab3f3f), [1904.13052](https://doi.org/10.1088/1751-8121/ab3f3f).
- 797 [78] N. Eddine Boukacem, A. Leary, R. Thériault, F. Gottlieb, M. Mani and P. François, *A*
 798 *Waddington landscape for prototype learning in generalized Hopfield networks*, arXiv
 799 e-prints arXiv:2312.03012 (2023), doi:[10.48550/arXiv.2312.03012](https://doi.org/10.48550/arXiv.2312.03012), [2312.03012](https://doi.org/10.48550/arXiv.2312.03012).

800 A Gardner's Hamiltonian vs K & H's Hamiltonian

801 Consider the generalized Hopfield Hamiltonian $H[\sigma|\xi] = -\sum_{i_1 < \dots < i_p=1}^N J_{i_1 \dots i_p} \sigma_{i_1} \dots \sigma_{i_p}$ with
 802 p -body interactions $J_{i_1 \dots i_p} = \frac{p!}{N^{p-1}} \sum_{\mu=1}^M \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu}$ described by Gardner [30], where M indicates
 803 the number of patterns ξ^{μ} used to construct J , and N denotes the number of components of
 804 each pattern ξ^{μ} and example σ . In this Section, we will omit ξ in the argument of $H[\sigma|\xi]$
 805 and write $H[\sigma]$ instead for notational simplicity. Unless indicated otherwise, we will assume
 806 a large number number of components $N \gg 1$ and patterns $M \sim \mathcal{O}(N^{p-1})$. We will start
 807 by comparing it to the dense Hopfield network Hamiltonian $\mathcal{H}[\sigma] = -\frac{1}{N^{p-1}} \sum_{\mu} \left(\sum_i \xi_i^{\mu} \sigma_i \right)^p$
 808 studied by K & H [26].

809 For that purpose, we rewrite H in the form $H[\sigma] = -\frac{1}{p!} \sum_{i_1 \neq \dots \neq i_p} J_{i_1 \dots i_p} \sigma_{i_1} \dots \sigma_{i_p}$ by sum-
 810 ming over all permutations of $\{i_1 \dots i_p\}$ in place of the restricted set $i_1 < \dots < i_p$ and compen-
 811 sating for double counting with the prefactor $\frac{1}{p!}$. This manipulation leads to

$$\begin{aligned} H[\sigma] &= -\frac{1}{p!} \sum_{i_1 \neq \dots \neq i_p} J_{i_1 \dots i_p} \sigma_{i_1} \dots \sigma_{i_p} \\ &= -\frac{1}{N^{p-1}} \sum_{\mu} \sum_{i_1 \neq \dots \neq i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1} \dots \sigma_{i_p}. \end{aligned}$$

812 On the other hand, K & H's Hamiltonian may be rewritten

$$\begin{aligned}\mathcal{H}[\sigma] &= -\frac{1}{N^{p-1}} \sum_{\mu} \left(\sum_i \xi_i^{\mu} \sigma_i \right)^p \\ &= -\frac{1}{N^{p-1}} \sum_{\mu} \left(\sum_{i_1} \xi_{i_1}^{\mu} \sigma_{i_1} \right) \dots \left(\sum_{i_p} \xi_{i_p}^{\mu} \sigma_{i_p} \right) \\ &= -\frac{1}{N^{p-1}} \sum_{\mu} \sum_{i_1 \dots i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1} \dots \sigma_{i_p},\end{aligned}$$

813 where the sum over $i_1 \dots i_p$ includes both the set of indices $i_1 \neq \dots \neq i_p$ found in $H[\sigma]$ and other
814 configurations where some indices are equal. For example, the configuration $i_1 \neq \dots \neq i_{p-1} = i_p$
815 contains the fewest equal indices after $i_1 \neq \dots \neq i_p$. In other words, $\mathcal{H}[\sigma]$ can be expressed as
816 an expansion around $H[\sigma]$, and the two Hamiltonians are equivalent when the normalized
817 residuals $\frac{\mathcal{H}[\sigma] - H[\sigma]}{N}$ vanish in the limit of large N . In this study, we encounter two cases which
818 bring different results.

819 **1** The Hamiltonians $\mathcal{H}[\sigma]$ and $H[\sigma]$ are dominated by a few closely packed configurations
820 ξ^{μ} that have finite overlap $\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i \sim \mathcal{O}(1)$ with σ . We say that they are aligned with
821 σ .

822 **2** The Hamiltonians $\mathcal{H}[\sigma]$ and $H[\sigma]$ are dominated by many spread out configurations
823 ξ^{μ} that have microscopic overlap $\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i \sim \mathcal{O}(N^{-1/2})$ with σ . We say that they are
824 misaligned with σ

825 We use the expansion of $\mathcal{H}[\sigma]$ to discuss both the aligned case and the misaligned case. We
826 start by writing the $i_1 \neq \dots \neq i_p$ and $i_1 \neq \dots \neq i_{p-1} = i_p$ terms explicitly, which leads to the
827 form

$$\begin{aligned}\mathcal{H}[\sigma] &= -\frac{1}{N^{p-1}} \sum_{\mu} \sum_{i_1 \neq \dots \neq i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1} \dots \sigma_{i_p} \\ &\quad - \frac{1}{2} \frac{p(p-1)}{N^{p-1}} \sum_{\mu} \sum_{i_1 \neq \dots \neq i_{p-1}} \xi_{i_1}^{\mu} \dots \left(\xi_{i_{p-1}}^{\mu} \right)^2 \sigma_{i_1} \dots \left(\sigma_{i_{p-1}} \right)^2 + \dots\end{aligned}$$

828 because there are $\binom{p}{2} = \frac{p(p-1)}{2}$ ways for the indices i_{p-1} and i_p to be equal. This expression
829 can be summarized by $\mathcal{H}[\sigma] = H[\sigma] + H'[\sigma] + \dots$, where $H'[\sigma]$ simplifies to

$$\begin{aligned}H'[\sigma] &= -\frac{1}{2} \frac{p(p-1)}{N^{p-1}} \sum_{\mu} \sum_{i_1 \neq \dots \neq i_{p-1}} \xi_{i_1}^{\mu} \dots \left(\xi_{i_{p-1}}^{\mu} \right)^2 \sigma_{i_1} \dots \left(\sigma_{i_{p-1}} \right)^2 \\ &= -\frac{1}{2} \frac{p(p-1)}{N^{p-2}} \sum_{\mu} \sum_{i_1 \neq \dots \neq i_{p-2}} \xi_{i_1}^{\mu} \dots \xi_{i_{p-2}}^{\mu} \sigma_{i_1} \dots \sigma_{i_{p-2}} \\ &= -\frac{1}{2} \frac{p!}{N^{p-2}} \sum_{\mu} \sum_{i_1 < \dots < i_{p-2}} \xi_{i_1}^{\mu} \dots \xi_{i_{p-2}}^{\mu} \sigma_{i_1} \dots \sigma_{i_{p-2}}.\end{aligned}$$

830 In the aligned case, $H'[\sigma]$ is $\mathcal{O}(1)$ in N because the sum over $i_1 < \dots < i_{p-2}$ is $\mathcal{O}(N^{p-2})$.
831 The terms implied by the ellipsis are even smaller because their sums are restricted by more
832 equality constraints. Therefore, the residuals $\frac{\mathcal{H}[\sigma] - H[\sigma]}{N}$ vanish in the limit of large N , and
833 the two Hamiltonians are equivalent. Conversely, we find that $\mathcal{H}[\sigma]$ and $H[\sigma]$ differ from

each other in the misaligned case (see Appendix B for more details). Therefore, although the phases of $H[\sigma]$ that we obtain in this study are qualitatively similar to the ones observed by K & H [26, 37], the phase diagram of $H[\sigma]$ must be compared against a simulation of $H[\sigma]$ rather than $\mathcal{H}[\sigma]$ in order to test our theory quantitatively.

To understand how to sample σ in both models, consider a Monte-Carlo simulation used to find the statistical equilibrium of a spin ensemble σ with Hamiltonian $G[\sigma]$. To be more specific, suppose σ is updated to a new state σ' with a randomly selected spin σ_i flipped with acceptance probability $P_i = \frac{1}{1+\exp[\beta(G[\sigma']-G[\sigma])]}$ for a large number of time-steps. This approach works well for $G[\sigma] = \mathcal{H}[\sigma]$. However, in the case of $H[\sigma]$, we find that the simulation only converges when we use the local field $h_i = \frac{p!}{N^{p-1}} \sum_{\mu} \xi_i^{\mu} \sum_{i_1 < \dots < i_{p-1}} \xi_{i_1}^{\mu} \dots \xi_{i_{p-1}}^{\mu} \sigma_{i_1} \dots \sigma_{i_{p-1}}$ mentioned by Gardner [30] to approximate $\frac{H[\sigma']-H[\sigma]}{2\sigma_i}$ at large N . In other words, we iteratively flip randomly chosen spins σ_i with acceptance probability $P_i = \frac{1}{1+\exp(2\beta h_i \sigma_i)}$ for a large number of time steps. For arbitrary p , it is not obvious how to compute h_i quickly as a sub-routine of the Monte-Carlo simulation. However, we find that both $p = 3$ and $p = 4$ have closed-formed expressions that are easy to evaluate numerically in an efficient way. To be more precise,

- $p = 3$ leads to $h_i = 3 \sum_{\mu} \xi_i^{\mu} \left[\left(\frac{1}{N} \sum_j \xi_j^{\mu} \sigma_j \right)^2 - \frac{1}{N} \right]$,
- and $p = 4$ leads to $h_i = 4 \sum_{\mu} \xi_i^{\mu} \left(\frac{1}{N} \sum_j \xi_j^{\mu} \sigma_j \right) \left[\left(\frac{1}{N} \sum_j \xi_j^{\mu} \sigma_j \right)^2 - \frac{3}{N} \right]$.

For this reason and also because the number $M \sim \mathcal{O}(N^{p-1})$ of patterns ξ^{μ} used in a Monte-Carlo simulations increases exponentially with p , we choose to simulate only $p = 3$ and $p = 4$.

The output of the neural network model that K & H designed for classification of data is $c_j = \tanh \left[\frac{1}{2} \beta (\mathcal{H}[\sigma'] - \mathcal{H}[\sigma]) \right] \approx \tanh \left[\beta p \sum_{\mu} \xi_j^{\mu} \left(\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i \right)^{p-1} \right]$. We omit the linear rectifier present in the original paper [26] because the overlaps $\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i$ are almost always positive (see for example the Supplement of [78]). The predicted class is then $j' = \text{argmax}_j \{c_j\}$. Using $1 - P_j = \frac{1}{1+\exp[\beta(\mathcal{H}[\sigma']-\mathcal{H}[\sigma])]}$ instead of c_j does not change j' because $1 - P_j$ and c_j are related by $1 - P_j = \frac{1}{2} [c_j + 1]$. When we evaluate P_i using H instead of \mathcal{H} , this relation does not always hold exactly. Rather, it should be considered an approximation.

B Direct model cumulant expansions

In the direct model, the average replicated partition function $\langle Z^L \rangle$ takes the form:

$$\langle Z^L \rangle = \left\langle \sum_{\sigma} \exp \left(-\beta \sum_{\gamma=1}^L H[\sigma^{\gamma} | \xi] \right) \right\rangle$$

with $\sigma = \{\sigma^1 \dots \sigma^L\}$. Gardner simplifies it to

$$\begin{aligned} \langle Z^L \rangle \approx & \left\langle \sum_{\sigma} \exp \left(\beta N \sum_{\gamma} \sum_{\mu \in \mathbb{I}_{\gamma}} \left[\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i^{\gamma} \right]^p \right. \right. \\ & \left. \left. + \beta \sum_{\gamma} \sum_{\mu \in \mathbb{I}_{\gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^{\gamma} \dots \sigma_{i_p}^{\gamma} \right) \right\rangle, \end{aligned} \quad (10)$$

864 where the sets Γ_γ contain the patterns ξ^μ that have macroscopic overlap with σ_γ , and their
 865 complement $\bar{\Gamma} = \cap_\gamma \bar{\Gamma}_\gamma$ consists of the remaining patterns. Two approximations are used to
 866 obtain this expression:

- 867 • $\sum_{\mu \in \Gamma_\gamma} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \approx N \sum_{\mu \in \Gamma_\gamma} \left[\frac{1}{N} \sum_i \xi_i^\mu \sigma_i^\gamma \right]^p$ because this part of
 868 $H[\sigma^\gamma | \xi]$ is aligned with σ (see Case 1 of Appendix A).
- 869 • $\sum_{\mu \in \bar{\Gamma}_\gamma} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \approx \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma$ since $\bar{\Gamma}$
 870 contains almost all of the elements in each $\bar{\Gamma}_\gamma$ when N is large.

871 Gardner evaluates the contribution of the $\mu \in \bar{\Gamma}$ terms via a cumulant expansion, resulting in:

$$\begin{aligned} & \log \left\langle \exp \left(\beta \sum_\gamma \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \right) \right\rangle \\ & \approx \beta \left\langle \sum_\gamma \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \right\rangle \\ & \quad + \frac{1}{2} \beta^2 \left\langle \left[\sum_\gamma \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \right]^2 \right\rangle \\ & \approx \frac{1}{2} \beta^2 \left\langle \left[\sum_\gamma \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \right] \left[\sum_{\delta} \frac{p!}{N^{p-1}} \sum_{j_1 < \dots < j_p} \xi_{j_1}^\mu \dots \xi_{j_p}^\mu \sigma_{j_1}^\delta \dots \sigma_{j_p}^\delta \right] \right\rangle \end{aligned}$$

872 because the product of independent spins $\xi_{i_1}^\mu \dots \xi_{i_p}^\mu$ averages to $\mathbf{0}$. The sums are then regrouped
 873 to get

$$\begin{aligned} & \log \left\langle \exp \left(\beta \sum_\gamma \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \right) \right\rangle \\ & = \frac{1}{2} \beta^2 \left[\frac{p!}{N^{p-1}} \right]^2 \left\langle \sum_\gamma \sum_{\delta} \sum_{i_1 < \dots < i_p} \sum_{j_1 < \dots < j_p} \xi_{i_1}^\mu \xi_{j_1}^\mu \dots \xi_{i_p}^\mu \xi_{j_p}^\mu \sigma_{i_1}^\gamma \sigma_{j_1}^\delta \dots \sigma_{i_p}^\gamma \sigma_{j_p}^\delta \right\rangle. \end{aligned}$$

874 Consider $\xi_i^\mu \xi_j^\mu$ for an arbitrary pair of indices i and j . There are two cases.

- 875 • If $i = j$, then $\xi_i^\mu \xi_j^\mu$ is deterministic and equal to 1.
- 876 • If $i \neq j$, then $\xi_i^\mu \xi_j^\mu$ can be either +1 and -1 with equal probabilities.

877 On the one hand, if $i_n = j_n$ for all n , then $\left\langle \xi_{i_1}^\mu \xi_{j_1}^\mu \dots \xi_{i_p}^\mu \xi_{j_p}^\mu \right\rangle = 1$. On the other hand, if $i_n \neq j_n$
 878 for some n , then $\left\langle \xi_{i_1}^\mu \xi_{j_1}^\mu \dots \xi_{i_p}^\mu \xi_{j_p}^\mu \right\rangle = 0$ because $\xi_{i_1}^\mu \xi_{j_1}^\mu \dots \xi_{i_p}^\mu \xi_{j_p}^\mu$ is still a product of independent
 879 random spins once the deterministic variables are removed. These two cases can be summarized

880 by $\left\langle \xi_{i_1}^\mu \xi_{j_1}^\mu \dots \xi_{i_p}^\mu \xi_{j_p}^\mu \right\rangle = \delta_{i_1 j_1} \dots \delta_{i_p j_p}$, which then gives

$$\begin{aligned}
& \log \left\langle \exp \left(\beta \sum_{\gamma} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \right) \right\rangle \\
&= \frac{1}{2} \beta^2 \left[\frac{p!}{N^{p-1}} \right]^2 \sum_{\gamma} \sum_{\delta} \sum_{i_1 < \dots < i_p} \sum_{j_1 < \dots < j_p} \delta_{i_1 j_1} \dots \delta_{i_p j_p} \sigma_{i_1}^\gamma \sigma_{j_1}^\delta \dots \sigma_{i_p}^\gamma \sigma_{j_p}^\delta \\
&= \frac{1}{2} \beta^2 \left[\frac{p!}{N^{p-1}} \right]^2 \sum_{\gamma} \sum_{\delta} \sum_{i_1 < \dots < i_p} \sigma_{i_1}^\gamma \sigma_{i_1}^\delta \dots \sigma_{i_p}^\gamma \sigma_{i_p}^\delta \\
&\approx \frac{1}{2} \beta^2 \frac{p!}{N^{p-1}} \frac{1}{N^{p-1}} \sum_{\gamma \delta} \left[\sum_i \sigma_i^\gamma \sigma_i^\delta \right]^p \\
&= \beta^2 \frac{p!}{N^{p-1}} N \sum_{\gamma < \delta} \left[\frac{1}{N} \sum_i \sigma_i^\gamma \sigma_i^\delta \right]^p + \frac{1}{2} \beta^2 \frac{p!}{N^{p-1}} LN.
\end{aligned}$$

881 The order $n > 2$ terms are subdominant in N and can be neglected when $p \geq 3$ [30]. The RS
882 free entropy is then obtained through a standard approach to the replica method. Note that
883 Gardner's Hamiltonian is misaligned with σ when the free entropy is dominated by this cumulant
884 expansion (see Case 2 of Appendix A). In the case of K & H's Hamiltonian, we must also take
885 into account the correction $H'[\sigma] = \frac{1}{2} \frac{p!}{N^{p-2}} \sum_{\gamma} \sum_{i_1 < \dots < i_{p-2}} \xi_{i_1}^\mu \dots \xi_{i_{p-2}}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_{p-2}}^\gamma$ introduced in
886 appendix A by imposing $i_{p-1} = i_p$. In fact, a cumulant expansion of this expression gives

$$\begin{aligned}
& \log \left\langle \exp \left(\beta p \sum_{\gamma} \frac{1}{2} \frac{p!}{N^{p-2}} \sum_{i_1 < \dots < i_{p-2}} \xi_{i_1}^\mu \dots \xi_{i_{p-2}}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_{p-2}}^\gamma \right) \right\rangle \\
&\approx \frac{1}{4} \beta^2 \frac{p!}{N^{p-2}} \frac{p(p-1)}{N^{p-2}} \sum_{\gamma < \delta} \left[\sum_i \sigma_i^\gamma \sigma_i^\delta \right]^{p-2} + \frac{1}{8} \beta^2 \frac{p!}{N^{p-2}} L \\
&= \frac{1}{4} p(p-1) \beta^2 \frac{p!}{N^{p-1}} N \sum_{\gamma < \delta} \left[\frac{1}{N} \sum_i \sigma_i^\gamma \sigma_i^\delta \right]^{p-2} + \frac{1}{8} \beta^2 \frac{p!}{N^{p-1}} LN,
\end{aligned}$$

887 which contributes to the free energy on the same order in N as Gardner's Hamiltonian. Therefore,
888 K & H's Hamiltonian is not equivalent to Gardner's Hamiltonian when the latter is misaligned
889 with σ (see Case 2). The index configurations with more equality constraints also contribute
890 to the free entropy on the same order in N because the factors of N that are lost to equality
891 constraints are restored when the sums get squared in the cumulant expansion.

892 $p = 2$ is the only positive integer such that Gardner's Hamiltonian and Krotov's Hamiltonian
893 are equivalent [5, 30]. In the misaligned case with a single stored pattern ξ^* (see Case 2), the
894 free entropy of $p = 2$ simplifies to

$$\begin{aligned}
\frac{\log(Z)}{N} &= \frac{1}{N} \log \left\langle \exp \left\{ \beta \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^\gamma \sigma_{i_2}^\gamma \right\} \right\rangle + \log 2 \\
&= \frac{1}{N} \log \left\langle \exp(-\beta) \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} x^2 + x \sqrt{\beta \frac{2}{N}} \sum_i \xi_i^* \sigma_i^\gamma \right\} \right\rangle + \log 2 \\
&= \frac{1}{N} \log \left[\int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} x^2 \right\} \cosh^N \left(x \sqrt{\beta \frac{2}{N}} \right) \right] - \beta \frac{1}{N} + \log 2,
\end{aligned}$$

895 by using the Hubbard-Stratonovich transformation. At large N , it approximates to:

$$\begin{aligned} \frac{\log(Z)}{N} &\approx \frac{1}{N} \log \left[\int dx \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}x^2 \right\} \left(1 + \beta \frac{1}{N}x^2 \right)^N \right] - \beta \frac{1}{N} + \log 2 \\ &\approx \frac{1}{N} \log \left[\int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}x^2 \right\} \exp(\beta x^2) \right] - \beta \frac{1}{N} + \log 2 \\ &= \left(-\frac{1}{2} \log(1 - 2\beta) - \beta \right) \frac{1}{N} + \log 2, \end{aligned}$$

896 thanks to the well-known limit $\lim_{N \rightarrow \infty} \left(1 + \frac{1}{N}z \right)^N = \exp(z)$. This free entropy is consistent
897 with the one found in literature when $\alpha = \frac{1}{N}$ [5].

898 C Teacher-student replicated partition function

899 Recall that the student samples its pattern from the posterior $P(\xi|\sigma) = \frac{P(\xi)\prod_a P(\sigma^a|\xi)}{P(\sigma)}$ (see
900 Section 3). Given $P(\xi)$ uniform, it can be rewritten as $P(\xi|\sigma) = \frac{\prod_a P(\sigma^a|\xi)}{\sum_{\xi} \prod_a P(\sigma^a|\xi)}$, where $P(\sigma^a|\xi)$
901 is the distribution of the direct model with a single pattern ξ . To simplify $P(\xi|\sigma)$ further,
902 we need to manipulate the partition function $Z = \sum_{\sigma^a} \exp(-\beta H[\sigma^a|\xi])$ of $P(\sigma^a|\xi)$ (see
903 Appendix A for the definition of $H[\sigma|\xi]$). Under the gauge transformation $\sigma_i^a \rightarrow \xi_i \sigma_i^a$, we
904 may write

$$Z = \sum_{\sigma^a} \exp \left(\beta \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \sigma_{i_1}^a \dots \sigma_{i_p}^a \right)$$

905 without changing the configurations of σ^a that we are summing over. Therefore, Z does not
906 depend on ξ , and we can factor it out of the sum \sum_{ξ} , which yields

$$\begin{aligned} P(\xi|\sigma) &= \frac{\prod_a \frac{1}{Z} \exp \left(\beta \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1} \dots \xi_{i_p} \sigma_{i_1}^a \dots \sigma_{i_p}^a \right)}{\sum_{\xi} \prod_a \frac{1}{Z} \exp \left(\beta \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1} \dots \xi_{i_p} \sigma_{i_1}^a \dots \sigma_{i_p}^a \right)} \\ &= \frac{\exp \left(\beta \frac{p!}{N^{p-1}} \sum_a \sum_{i_1 < \dots < i_p} \xi_{i_1} \dots \xi_{i_p} \sigma_{i_1}^a \dots \sigma_{i_p}^a \right)}{\sum_{\xi} \exp \left(\beta \frac{p!}{N^{p-1}} \sum_a \sum_{i_1 < \dots < i_p} \xi_{i_1} \dots \xi_{i_p} \sigma_{i_1}^a \dots \sigma_{i_p}^a \right)}. \end{aligned}$$

907 Therefore, we define the partition function of the inverse model to be $\mathcal{Z} = \sum_{\xi} \exp(-\beta H[\xi|\sigma])$
908 (again, see Appendix A for the definition of $H[\xi|\sigma]$). The L^{th} power of \mathcal{Z} and its average then
909 take the form

$$\begin{aligned} \mathcal{Z}^L &= \sum_{\xi} \prod_b \exp \left(\beta \frac{p!}{N^{p-1}} \sum_a \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right) \\ \langle \mathcal{Z}^L \rangle &= \sum_{\sigma} P(\sigma) \sum_{\xi} \exp \left(\beta \frac{p!}{N^{p-1}} \sum_{ab} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right), \end{aligned}$$

910 where $\mathbf{b} \in \{1 \dots L\}$ label replicas in the set of patterns $\xi = \{\xi^1 \dots \xi^L\}$ inferred by the
 911 student. Using the definition of conditional probability, we rewrite $P(\sigma)$ as

$$\begin{aligned} P(\sigma) &= \sum_{\xi^*} P(\sigma|\xi^*) P(\xi^*) \\ &= \frac{1}{2^N} \sum_{\xi^*} P(\sigma|\xi^*) \\ &= \frac{1}{2^N} \sum_{\xi^*} \prod_a P(\sigma^a|\xi^*), \end{aligned}$$

912 where $P(\sigma|\xi^*)$ has the same functional form as $P(\sigma|\xi^b)$, but has hyperparameters \mathbf{p}^* and β^*
 913 in place of \mathbf{p} and β . As we did for Z , we factor the partition function Z^* of $P(\sigma^a|\xi^*)$ out of
 914 the sum, which yields

$$\begin{aligned} P(\sigma) &= \frac{1}{2^N} \frac{1}{[Z^*]^M} \sum_{\xi^*} \prod_a \exp\left(\beta^* \frac{p^*!}{N^{p^*-1}} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a\right) \\ &= \frac{1}{2^N} \frac{Z^*}{[Z^*]^M} = \frac{1}{2^{MN}} \frac{Z^*}{[2^{N/M-N} Z^*]^M}, \end{aligned}$$

915 where $Z^* = \sum_{\xi^*} \exp(-\beta^* H[\xi^*|\sigma])$ is the partition function of the inverse model with in-
 916 teraction order p^* . Using $\sum_{\sigma} P(\sigma) = 1$, we immediately deduce that $[2^{N/M-N} Z^*]^M = \langle Z^* \rangle$.
 917 Plugging $P(\sigma) = \frac{1}{2^{MN}} \frac{Z^*}{\langle Z^* \rangle}$ back in $\langle Z^L \rangle$ then gives

$$\begin{aligned} \langle Z^L \rangle &= \frac{1}{2^{MN}} \frac{1}{\langle Z^* \rangle} \sum_{\xi^*} \sum_{\sigma} \exp\left(\beta^* \frac{p^*!}{N^{p^*-1}} \sum_a \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a\right) \\ &\quad \sum_{\xi} \exp\left(\beta \frac{p!}{N^{p-1}} \sum_{ab} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a\right). \end{aligned}$$

918 We simplify this expression to:

$$\begin{aligned}
\langle Z^L \rangle &= \frac{1}{2^{MN}} \frac{1}{\langle Z^* \rangle} \sum_{\xi^*} \sum_{\sigma} \exp \left(\beta^* \frac{p^*!}{N^{p^*-1}} \sum_{a \in \bar{\Gamma}_*} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \right. \\
&\quad \left. + \beta^* \frac{p^*!}{N^{p^*-1}} \sum_{a \in \bar{\Gamma}_*} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \right) \\
&\quad \sum_{\xi} \exp \left(\beta \frac{p!}{N^{p-1}} \sum_b \sum_{a \in \Gamma_b} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right. \\
&\quad \left. + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{a \in \bar{\Gamma}_b} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right) \\
&\approx \frac{1}{2^{MN}} \frac{1}{\langle Z^* \rangle} \sum_{\xi^*} \sum_{\sigma} \exp \left(\beta^* N \sum_{a \in \bar{\Gamma}_*} \left[\frac{1}{N} \sum_i \xi_i^* \sigma_i^a \right]^{p^*} \right. \\
&\quad \left. + \beta^* \frac{p^*!}{N^{p^*-1}} \sum_{a \in \bar{\Gamma}_*} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \right) \\
&\quad \sum_{\xi} \exp \left(\beta N \sum_b \sum_{a \in \Gamma_b} \left[\frac{1}{N} \sum_i \xi_i^b \sigma_i^a \right]^p \right. \\
&\quad \left. + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{a \in \bar{\Gamma}_b} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right),
\end{aligned}$$

919 where Γ_b represents the set of inputs σ^a which have macroscopic overlap with the pattern ξ^b ,
920 and $\bar{\Gamma} = [\cap_b \bar{\Gamma}_b] \cap \bar{\Gamma}_*$ contains almost all of the elements in each $\bar{\Gamma}_b$ and $\bar{\Gamma}_*$ for $N \rightarrow \infty$. The
921 reasoning used to build the sets Γ_* , Γ_b and $\bar{\Gamma}$ is the same as outlined at the start of appendix B.

922 D Teacher-student free entropy

923 Assuming that the teacher is misaligned with σ (see Case 2 of Appendix A), the form of $\langle Z^L \rangle$
924 obtained in appendix C simplifies to

$$\begin{aligned}
\langle Z^L \rangle &\approx \frac{1}{2^{MN}} \frac{1}{\langle Z^* \rangle} \sum_{\xi^*} \sum_{\sigma} \exp \left(\beta^* \frac{p^*!}{N^{p^*-1}} \sum_{a \in \bar{\Gamma}_*} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \right) \\
&\quad \exp \left(\beta N \sum_b \sum_{a \in \Gamma_b} \left[\frac{1}{N} \sum_i \xi_i^b \sigma_i^a \right]^p + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{a \in \bar{\Gamma}_b} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right).
\end{aligned}$$

925 In order to evaluate $\langle Z^* \rangle = [2^{N/M-N} Z^*]^M$, we recall that the teacher is a special case of the
926 direct model with a single memory (see Section 3). Since the teacher is in the misaligned case,
927 its free entropy is

$$\frac{\log(Z^*)}{N} = \begin{cases} \left(-\frac{1}{2} \log(1 - 2\beta^*) - \beta^* \right) \frac{1}{N} + \log 2 & p^* = 2 \\ \frac{1}{2} [\beta^*]^2 \frac{p^*!}{N^{p^*-1}} + \log 2 + \mathcal{O}\left(\frac{1}{N^{3p^*/2-2}}\right) & p^* \geq 3, \end{cases}$$

928 as derived in Appendix B. Given $\alpha^* = \frac{Mp^*!}{N^{p^*-1}}$, we use it to simplify $\frac{\log\langle Z^* \rangle}{N}$ to

$$\begin{aligned} \frac{\log\langle Z^* \rangle}{N} &= \frac{M \log[2^{N/M-N} Z^*]}{N} \\ &= \begin{cases} \frac{1}{2} \left(-\frac{1}{2} \log(1-2\beta^*) - \beta^* \right) \alpha^* + \log 2 & p^* = 2 \\ \frac{1}{2} [\beta^*]^2 \alpha^* + \log 2 + \mathcal{O}\left(\frac{1}{N^{p^*/2-1}}\right) & p^* \geq 3, \end{cases} \end{aligned}$$

929 which is the paramagnetic free entropy of a p^* -body Hopfield network [5, 30]. Coming back to
930 $\langle Z^L \rangle$, we fix order parameters q^{*b} , q^{bc} and m_a^b using the delta functions $\delta\left(Nq^{*b} - \sum_i \xi_i^* \xi_i^b\right)$,
931 $\delta\left(Nq^{bc} - \sum_i \xi_i^b \xi_i^c\right)$ and $\delta\left(Nm_a^b - \sum_i \xi_i^b \sigma_i^a\right)$, which results in

$$\begin{aligned} \langle Z^L \rangle &= \frac{1}{2^{MN}} \frac{1}{\langle Z^* \rangle} \sum_{\xi^* \xi} \sum_{\sigma} \int_{\mathbb{R}} \prod_b dq^{*b} \prod_{b<c} dq^{bc} \prod_b \prod_{a \in \Gamma_b} dm_a^b \\ &\quad \delta\left(Nq^{*b} - \sum_i \xi_i^* \xi_i^b\right) \delta\left(Nq^{bc} - \sum_i \xi_i^b \xi_i^c\right) \delta\left(Nm_a^b - \sum_i \xi_i^b \sigma_i^a\right) \\ &\quad \exp\left(\beta N \sum_b \sum_{a \in \Gamma_b} \left[\frac{1}{N} \sum_i \xi_i^b \sigma_i^a \right]^p\right) \\ &\quad + \beta^* \frac{p^*!}{N^{p^*-1}} \sum_{a \in \Gamma} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \\ &\quad + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{a \in \Gamma} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \Big). \end{aligned}$$

932 In Fourier space, this expression takes the form

$$\begin{aligned} \langle Z^L \rangle &= \frac{1}{\langle Z^* \rangle} \sum_{\xi^* \xi} \left\langle \int \prod_b dq^{*b} dr^{*b} \prod_{b<c} dq^{bc} dr^{bc} \prod_b \prod_{a \in \Gamma_b} dm_a^b dk_a^b \right. \\ &\quad \exp\left\{ \beta^* \beta \alpha \sum_b \left(\sum_i \xi_i^* \xi_i^b - Nq^{*b} \right) r^{*b} + \beta^2 \alpha \sum_{b<c} \left(\sum_i \xi_i^b \xi_i^c - Nq^{bc} \right) r^{bc} \right\} \\ &\quad \exp\left\{ \beta \sum_b \sum_{a \in \Gamma_b} \left(\sum_i \xi_i^b \sigma_i^a - Nm_a^b \right) k_a^b + \beta N \sum_b \sum_{a \in \Gamma_b} \left[\frac{1}{N} \sum_i \xi_i^b \sigma_i^a \right]^p \right. \\ &\quad + \beta^* \frac{p^*!}{N^{p^*-1}} \sum_{a \in \Gamma} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \\ &\quad \left. + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{a \in \Gamma} \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right\} \Bigg\rangle_{\sigma}, \end{aligned}$$

933 where the sum over σ with a pre-factor of $\frac{1}{2^{MN}}$ was replaced by the uniform average $\langle \rangle_{\sigma}$.
934 Following the same reasoning as in appendix B, a second order cumulant expansion of the last

935 two terms for any $\mathbf{a} \in \bar{\Gamma}$ yields

$$\begin{aligned}
& \log \left\langle \exp \left\{ \beta^* \frac{p^*!}{N^{p^*-1}} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \right. \right. \\
& \quad \left. \left. + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right\} \right\rangle \\
& \approx \frac{1}{2} \beta^2 \left[\frac{p!}{N^{p-1}} \right]^2 \sum_{b \neq c} \sum_{i_1 < \dots < i_p} \sum_{j_1 < \dots < j_p} \xi_{i_1}^b \xi_{j_1}^c \dots \xi_{i_p}^b \xi_{j_p}^c \left\langle \sigma_{i_1}^a \sigma_{j_1}^a \dots \sigma_{i_p}^a \sigma_{j_p}^a \right\rangle \\
& \quad + \beta^* \beta \frac{p^*!}{N^{p^*-1}} \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_{p^*}} \sum_{j_1 < \dots < j_p} \left\langle \xi_{i_1}^* \sigma_{i_1}^a \dots \xi_{i_{p^*}}^* \sigma_{i_{p^*}}^a \xi_{j_1}^b \sigma_{j_1}^a \dots \xi_{j_p}^b \sigma_{j_p}^a \right\rangle \\
& \quad + \frac{1}{2} \beta^2 \frac{p!}{N^{p-1}} LN + \frac{1}{2} [\beta^*]^2 \frac{p^*!}{N^{p^*-1}} N.
\end{aligned}$$

936 When $p^* = p$, it reduces to

$$\begin{aligned}
& \log \left\langle \exp \left\{ \beta^* \frac{p^*!}{N^{p^*-1}} \sum_{i_1 < \dots < i_{p^*}} \xi_{i_1}^* \dots \xi_{i_{p^*}}^* \sigma_{i_1}^a \dots \sigma_{i_{p^*}}^a \right. \right. \\
& \quad \left. \left. + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right\} \right\rangle \\
& = \beta^2 \frac{p!}{N^{p-1}} N \sum_{b < c} \left[\frac{1}{N} \sum_i \xi_i^b \xi_i^c \right]^p + \beta^* \beta \frac{p!}{N^{p-1}} N \sum_b \left[\frac{1}{N} \sum_i \xi_i^* \xi_i^b \right]^p \\
& \quad + \frac{1}{2} \beta^2 \frac{p!}{N^{p-1}} LN + \frac{1}{2} [\beta^*]^2 \frac{p!}{N^{p-1}} N,
\end{aligned}$$

937 because $\langle \sigma_{i_n}^a \sigma_{j_n}^a \rangle = \delta_{i_n j_n}$ (see Appendix B for more details). On the contrary, the second order
938 expectation $\langle \xi_{i_1}^* \sigma_{i_1}^a \dots \xi_{i_{p^*}}^* \sigma_{i_{p^*}}^a \xi_{j_1}^b \sigma_{j_1}^a \dots \xi_{j_p}^b \sigma_{j_p}^a \rangle$ vanishes when $p^* \neq p$. In fact, spins come in
939 pairs $\langle \sigma_{i_n}^a \sigma_{j_n}^a \rangle = \delta_{i_n j_n}$ only up to $n \leq \min\{p^*, p\}$, and the remaining single-spin averages
940 $\langle \sigma_{i_n}^a \rangle = 0$ make the second order expectation vanish.

941 We need to go beyond second order to treat $p^* \neq p$. We will focus on $p^* = 2$ and $p \geq 3$
942 to investigate the consequences of using a p -body model to learn examples generated by the
943 original 2-body Hopfield model. For simplicity, we take p even so that the spins of both terms
944 can be grouped in pairs at order $\frac{p}{2} + 1$, when the teacher term $\beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a$ is
945 raised to the power of $\frac{p}{2}$ and the student term is raised to the power of 1. This restriction will
946 simplify some of the incoming calculations. To leading order in N , the cumulant generating

947 function reduces to

$$\begin{aligned} & \log \left\langle \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right\} \right\rangle \\ & \approx \log \left[\left\langle \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \right. \\ & \quad \left. \left\langle \exp \left\{ \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right\} \right\rangle \right. \\ & \quad \left. + \left\langle \beta \frac{p!}{N^{p-1}} \sum_b \sum_{j_1 < \dots < j_p} \xi_{j_1}^b \dots \xi_{j_p}^b \sigma_{j_1}^a \dots \sigma_{j_p}^a \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \right], \end{aligned}$$

948 where the last term encompasses the teacher-student coupling that allows retrieval to take
949 place. The teacher term

$$\begin{aligned} & \log \left\langle \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \\ & \approx -\frac{1}{2} \log(1 - 2\beta^*) - \beta^* \end{aligned}$$

950 and the student term

$$\begin{aligned} & \log \left\langle \exp \left\{ \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right\} \right\rangle \\ & \approx \beta^2 \frac{p!}{N^{p-1}} N \sum_{b < c} \left[\frac{1}{N} \sum_i \xi_i^b \xi_i^c \right]^p + \frac{1}{2} \beta^2 \frac{p!}{N^{p-1}} LN \end{aligned}$$

951 are both known from Appendix B. Later on, we will use $\log(\mathbf{z}^*)$ and \mathbf{z}^* as shorthands for
952 $-\frac{1}{2} \log(1 - 2\beta^*) - \beta^*$ and $\exp\left(-\frac{1}{2} \log(1 - 2\beta^*) - \beta^*\right)$, respectively. The coupling between
953 the teacher and the student can be rewritten as

$$\begin{aligned} & \left\langle \beta \frac{p!}{N^{p-1}} \sum_b \sum_{j_1 < \dots < j_p} \xi_{j_1}^b \dots \xi_{j_p}^b \sigma_{j_1}^a \dots \sigma_{j_p}^a \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \\ & = \left\langle \beta \frac{p!}{N^{p-1}} \sum_b \sum_{j_1 < \dots < j_p} \xi_{j_1}^* \dots \xi_{j_p}^* \xi_{j_1}^b \dots \xi_{j_p}^b \xi_{j_1}^* \dots \xi_{j_p}^* \sigma_{j_1}^a \dots \sigma_{j_p}^a \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \\ & = \beta \frac{p!}{N^{p-1}} \sum_b \sum_{j_1 < \dots < j_p} \xi_{j_1}^* \dots \xi_{j_p}^* \xi_{j_1}^b \dots \xi_{j_p}^b \left\langle \xi_{j_1}^* \dots \xi_{j_p}^* \sigma_{j_1}^a \dots \sigma_{j_p}^a \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \end{aligned}$$

954 because $[\xi_{j_n}^*]^2 = 1$ for every index j_n . All interacting spin tuples of the form $\xi_{j_1}^* \dots \xi_{j_p}^* \sigma_{j_1}^a \dots \sigma_{j_p}^a$

955 are statistically equivalent as long as $j_1 < \dots < j_p$, so the teacher-student coupling simplifies to

$$\begin{aligned}
& \left\langle \beta \frac{p!}{N^{p-1}} \sum_b \sum_{j_1 < \dots < j_p} \xi_{j_1}^b \dots \xi_{j_p}^b \sigma_{j_1}^a \dots \sigma_{j_p}^a \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \\
&= \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^* \dots \xi_{i_p}^* \xi_{i_1}^b \dots \xi_{i_p}^b \\
& \left\langle \frac{p!}{N^p} \sum_{j_1 < \dots < j_p} \xi_{j_1}^* \dots \xi_{j_p}^* \sigma_{j_1}^a \dots \sigma_{j_p}^a \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \\
&= V(\beta^*, p) \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^* \dots \xi_{i_p}^* \xi_{i_1}^b \dots \xi_{i_p}^b,
\end{aligned}$$

956 where $V(\beta^*, p) = \left\langle \frac{p!}{N^p} \sum_{j_1 < \dots < j_p} \xi_{j_1}^* \dots \xi_{j_p}^* \sigma_{j_1}^a \dots \sigma_{j_p}^a \exp \left(\beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right) \right\rangle$ does not
957 depend on the microscopic details of the system. In fact, it can be expressed as a combination of
958 the moments of \mathbf{z}^* , which can all be derived from $\log(\mathbf{z}^*)$. To leading order in N , the cumulant
959 generating function expands to

$$\begin{aligned}
& \log \left\langle \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a + \beta \frac{p!}{N^{p-1}} \sum_b \sum_{i_1 < \dots < i_p} \xi_{i_1}^b \dots \xi_{i_p}^b \sigma_{i_1}^a \dots \sigma_{i_p}^a \right\} \right\rangle \\
& \approx -\frac{1}{2} \log(1 - 2\beta^*) - \beta^* + \beta^2 \frac{p!}{N^{p-1}} N \sum_{b < c} \left[\frac{1}{N} \sum_i \xi_i^b \xi_i^c \right]^p + \frac{1}{2} \beta^2 \frac{p!}{N^{p-1}} LN \\
& + [1 - 2\beta^*]^{1/2} \exp(\beta^*) V(\beta^*, p) \beta N \sum_b \left[\frac{1}{N} \sum_i \xi_i^* \xi_i^b \right]^p.
\end{aligned}$$

960 At this stage, we only need to find $V(\beta^*, p)$ in order to solve the system. We focus on two
961 different scalings of M and β^* that make the teacher-student coupling leading order in N :

962 **1** $M \sim \mathcal{O}(N^{p-1})$ and $\beta^* \sim \mathcal{O}(N^{2/p-1})$ will be called the large-noise scaling.

963 **2** $M \sim \mathcal{O}(N^{p/2})$ and $\beta^* \sim \mathcal{O}(1)$ will be called the finite-noise scaling.

964 The student term vanishes in the first scenario but is leading order in the second one. The case
965 of the teacher-student coupling is more subtle. When β^* is small, we may keep only the first
966 non-vanishing order of the exponential function present in the definition of $V(\beta^*, p)$. Since p
967 is even, it leads to

$$\begin{aligned}
V(\beta^*, p) & \approx \frac{1}{(p/2)!} \left\langle \frac{p!}{N^p} \sum_{j_1 < \dots < j_p} \xi_{j_1}^* \dots \xi_{j_p}^* \sigma_{j_1}^a \dots \sigma_{j_p}^a \left(\beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right)^{p/2} \right\rangle \quad (11) \\
& = \frac{[\beta^*]^{p/2} 2^{p/2} p!}{(p/2)! N^{p/2} 2^{p/2}} \\
& = \frac{[\beta^*]^{p/2} p!}{(p/2)! N^{p/2}}
\end{aligned}$$

968 because there are $\prod_{n=1}^{p/2} \binom{2n}{2} = \frac{p!}{2^{p/2}}$ spin pairings with non-zero expectation that satisfy the
969 inequality constraints. In the large-noise scaling, we set

$$\lambda = \frac{[\beta^*]^{p/2}}{(p/2)!} N^{p/2-1} \sim \mathcal{O}(1)$$

970 to get the asymptotically exact expression $V([(p/2)!]^{2/p} N^{1-2/p}, \mathbf{p}) = \lambda \frac{p!}{N^{p-1}}$. In the finite-
 971 noise scaling, this expansion is only an order of magnitude approximation. However, it still
 972 indicates that $V(\beta^*, \mathbf{p})$ is $\mathcal{O}(N^{-p/2})$ when β^* is $\mathcal{O}(1)$ in N . In other words, it shows that
 973 there is an $\mathcal{O}(1)$ parameter η such that $V(\beta^*(\eta, \mathbf{p}), \mathbf{p}) = \eta \frac{(p/2+1)!}{N^{p/2}}$. We will now use the
 974 cumulants $\frac{\partial \log(z^*)}{\partial \beta^*}$ and $\frac{\partial \log(z^*)}{\partial \beta^{*2}}$ of \mathbf{z}^* to derive the value of η corresponding to $\mathbf{p} = 4$. First of
 975 all, note that $\frac{4!}{N^4} \sum_{j_1 < \dots < j_4} \xi_{j_1}^* \dots \xi_{j_4}^* \sigma_{j_1}^a \dots \sigma_{j_4}^a$ can be expressed as:

$$\begin{aligned} & \frac{24}{N^4} \sum_{j_1 < j_2 < j_3 < j_4} \xi_{j_1}^* \xi_{j_2}^* \xi_{j_3}^* \xi_{j_4}^* \sigma_{j_1}^a \sigma_{j_2}^a \sigma_{j_3}^a \sigma_{j_4}^a \\ &= \frac{1}{N^4} \sum_{j_1 \neq j_2 \neq j_3 \neq j_4} \xi_{j_1}^* \xi_{j_2}^* \xi_{j_3}^* \xi_{j_4}^* \sigma_{j_1}^a \sigma_{j_2}^a \sigma_{j_3}^a \sigma_{j_4}^a \\ &= \frac{1}{N^4} \left[\sum_{j_1 \neq j_2} \xi_{j_1}^* \xi_{j_2}^* \sigma_{j_1}^a \sigma_{j_2}^a \right] \left[\sum_{j_3 \neq j_4} \xi_{j_3}^* \xi_{j_4}^* \sigma_{j_3}^a \sigma_{j_4}^a \right] - \frac{4}{N^3} \left[\sum_{j_1 \neq j_2} \xi_{j_1}^* \xi_{j_2}^* \sigma_{j_1}^a \sigma_{j_2}^a \right] - \frac{2}{N^2} \\ &= \frac{1}{N^2} \left[\frac{2}{N} \sum_{j_1 < j_2} \xi_{j_1}^* \xi_{j_2}^* \sigma_{j_1}^a \sigma_{j_2}^a \right]^2 - \frac{4}{N^2} \left[\frac{2}{N} \sum_{j_1 < j_2} \xi_{j_1}^* \xi_{j_2}^* \sigma_{j_1}^a \sigma_{j_2}^a \right] - \frac{2}{N^2} \end{aligned}$$

976 by subtracting the diagonals where pairs of indices are equal. Therefore, $\frac{1}{z^*} V(\beta^*, \mathbf{p})$ reduces to

$$\begin{aligned} \frac{1}{z^*} V(\beta^*, \mathbf{p}) &= \left\langle \frac{24}{N^4} \sum_{j_1 < j_2 < j_3 < j_4} \xi_{j_1}^* \xi_{j_2}^* \xi_{j_3}^* \xi_{j_4}^* \sigma_{i_1}^a \sigma_{i_2}^a \sigma_{i_3}^a \sigma_{i_4}^a \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \\ &= \frac{1}{z^*} \frac{1}{N^2} \left[\left\langle \left[\frac{2}{N} \sum_{j_1 < j_2} \xi_{j_1}^* \xi_{j_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right]^2 \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \right. \\ &\quad - 4 \left\langle \left[\frac{2}{N} \sum_{j_1 < j_2} \xi_{j_1}^* \xi_{j_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right] \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \\ &\quad \left. - 2 \left\langle \exp \left\{ \beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right\} \right\rangle \right] \\ &= \frac{1}{N^2} \left[\frac{\partial \log(z^*)}{\partial \beta^{*2}} + \left[\frac{\partial \log(z^*)}{\partial \beta^*} \right]^2 - 4 \frac{\partial \log(z^*)}{\partial \beta^*} - 2 \right]. \end{aligned}$$

977 The cumulants evaluate to

$$\begin{aligned} \frac{\partial \log(z^*)}{\partial \beta^*} &= \frac{\partial}{\partial \beta^*} \left[-\frac{1}{2} \log(1-2\beta^*) - \beta^* \right] = \frac{2\beta^*}{1-2\beta^*} \\ \frac{\partial \log(z^*)}{\partial \beta^{*2}} &= \frac{\partial}{\partial \beta^{*2}} \left[-\frac{1}{2} \log(1-2\beta^*) - \beta^* \right] = \frac{2}{(1-2\beta^*)^2}, \end{aligned}$$

978 so we obtain

$$\begin{aligned} \frac{1}{z^*} V(\beta^*, \mathbf{p}) &= \frac{1}{N^2} \left[\frac{2}{(1-2\beta^*)^2} + \frac{4[\beta^*]^2}{(1-2\beta^*)^2} - \frac{8\beta^*}{1-2\beta^*} - 2 \right] \\ &= \frac{6}{N^2} \frac{2[\beta^*]^2}{(1-2\beta^*)^2}. \end{aligned}$$

979 In other terms, we find $\eta = \frac{2[\beta^*]^2}{(1-2\beta^*)^2}$ when $p = 4$. In summary, depending on the scaling, the
 980 teacher student coupling either simplifies to

981 **1** $\beta\lambda\alpha\frac{N}{M}\sum_b\left[\frac{1}{N}\sum_i\xi_i^*\xi_i^b\right]^p$ where $\alpha = \frac{Mp!}{N^{p-1}}$ and $\lambda = \frac{[\beta^*]^{p/2}}{(p/2)!}N^{p/2-1}$ are finite,

982 **2** or $\beta\eta\alpha\frac{N}{M}\sum_b\left[\frac{1}{N}\sum_i\xi_i^*\xi_i^b\right]^p$ where $\alpha = \frac{M(p/2+1)!}{N^{p/2}}$ and η are finite.

983 In either case, the result is similar to $p^* = p$ except for its pre-factor. We describe the rest of
 984 the derivation only for $p^* = p$ because the $p^* = 2$ and $p \geq 3$ calculations are almost identical.
 985 Putting the result of the $p^* = p$ cumulant expansion back in $\langle Z^L \rangle$, we get:

$$\begin{aligned} \langle Z^L \rangle &\approx \frac{1}{\langle Z^* \rangle} \sum_{\xi^* \xi} \left\langle \int \prod_b dq^{*b} dr^{*b} \prod_{b<c} dq^{bc} dr^{bc} \prod_b \prod_{a \in \Gamma_b} dm_a^b dk_a^b \right. \\ &\quad \exp \left\{ \beta^* \beta \alpha \sum_b \left(\sum_i \xi_i^* \xi_i^b - Nq^{*b} \right) r^{*b} + \beta^2 \alpha \sum_{b<c} \left(\sum_i \xi_i^b \xi_i^c - Nq^{bc} \right) r^{bc} \right\} \\ &\quad \exp \left\{ \beta \sum_b \sum_{a \in \Gamma_b} \left(\sum_i \xi_i^b \sigma_i^a - Nm_a^b \right) k_a^b + \beta N \sum_b \sum_{a \in \Gamma_b} [m_a^b]^p \right\} \\ &\quad \left. \exp \left\{ \beta^* \beta \alpha N \sum_b [q^{*b}]^p + \beta^2 \alpha N \sum_{b<c} [q^{bc}]^p + \frac{1}{2} \beta^2 \alpha L N + \frac{1}{2} [\beta^*]^2 \alpha N \right\} \right\rangle, \end{aligned}$$

986 where $\alpha = \frac{Mp!}{N^{p-1}}$. The saddle point of $\langle Z^L \rangle$ then evaluates to

$$\begin{aligned} \frac{\log \langle Z^L \rangle}{N} &\approx \text{Extr}_{m,k,q,r,q^*,r^*} \left[\beta^* \beta \alpha \sum_b [q^{*b}]^p + \beta^2 \alpha \sum_{b<c} [q^{bc}]^p + \beta \sum_b \sum_{a \in \Gamma_b} [m_a^b]^p \right. \\ &\quad - \beta^* \beta \alpha \sum_b r^{*b} q^{*b} - \beta^2 \alpha \sum_{b<c} r^{bc} q^{bc} - \beta \sum_b \sum_{a \in \Gamma_b} m_a^b k_a^b \\ &\quad + \frac{1}{2} \beta^2 \alpha L + \frac{1}{2} [\beta^*]^2 \alpha - \frac{\log \langle Z \rangle}{N} + \log 2 \\ &\quad + \frac{1}{N} \log \left\langle \sum_{\xi} \exp \left\{ \beta \sum_b \sum_{a \in \Gamma_b} k_a^b \sum_i \xi_i^b \sigma_i^a \right. \right. \\ &\quad \left. \left. + \beta^* \beta \alpha \sum_b r^{*b} \sum_i \xi_i^* \xi_i^b + \beta^2 \alpha \sum_{b<c} r^{bc} \sum_i \xi_i^b \xi_i^c \right\} \right\rangle_{\xi^* \sigma} \left. \right] \\ &= \text{Extr} \left[\beta^* \beta \alpha \sum_b [q^{*b}]^p + \beta^2 \alpha \sum_{b<c} [q^{bc}]^p + \beta \sum_b \sum_{a \in \Gamma_b} [m_a^b]^p \right. \\ &\quad - \beta^* \beta \alpha \sum_b r^{*b} q^{*b} - \beta^2 \alpha \sum_{b<c} r^{bc} q^{bc} - \beta \sum_b \sum_{a \in \Gamma_b} m_a^b k_a^b \\ &\quad + \frac{1}{2} \beta^2 \alpha L + \frac{1}{2} [\beta^*]^2 \alpha - \frac{\log \langle Z \rangle}{N} + \log 2 \\ &\quad + \frac{1}{N} \sum_i \log \left\langle \sum_{\xi_i} \exp \left\{ \beta \sum_b \sum_{a \in \Gamma_b} k_a^b \xi_i^b \sigma_i^a \right. \right. \\ &\quad \left. \left. + \beta^* \beta \alpha \sum_b r^{*b} \xi_i^* \xi_i^b + \beta^2 \alpha \sum_{b<c} r^{bc} \xi_i^b \xi_i^c \right\} \right\rangle_{\xi_i^* \sigma_i} \left. \right], \end{aligned}$$

987 where the average over ξ^* and σ is uniform. We use $\frac{\log\langle Z^* \rangle}{N} = \frac{1}{2}[\beta^*]^2 \alpha + \log 2$ to simplify
 988 $\frac{\log\langle Z^L \rangle}{N}$ to

$$\begin{aligned} \frac{\log\langle Z^L \rangle}{N} &\approx \text{Extr} \left[\beta^* \beta \alpha \sum_b [q^{*b}]^p + \beta^2 \alpha \sum_{b<c} [q^{bc}]^p + \beta \sum_b \sum_{a \in \Gamma_b} [m_a^b]^p \right. \\ &\quad - \beta^* \beta \alpha \sum_b r^{*b} q^{*b} - \beta^2 \alpha \sum_{b<c} r^{bc} q^{bc} - \beta \sum_b \sum_{a \in \Gamma_b} m_a^b k_a^b \\ &\quad + \frac{1}{2} \beta^2 \alpha L + \frac{1}{N} \sum_i \log \left\langle \sum_{\xi_i} \exp \left\{ \beta \sum_b \sum_{a \in \Gamma_b} k_a^b \xi_i^b \sigma_i^a \right. \right. \\ &\quad \left. \left. + \beta^* \beta \alpha \sum_b r^{*b} \xi_i^* \xi_i^b + \beta^2 \alpha \sum_{b<c} r^{bc} \xi_i^b \xi_i^c \right\} \right\rangle_{\xi_i^* \sigma_i} \left. \right]. \end{aligned}$$

989 Assuming each ξ^b has macroscopic overlap with at most one pattern σ^a and using the replica-
 990 symmetric ansatz $q^{*b} = q^*$, $q^{bc} = q$, $r^{*b} = r^*$, $r^{bc} = r$, $m_a^b = m$, $k_a^b = k$, the free entropy
 991 approximates to

$$\begin{aligned} f &= \lim_{N \rightarrow \infty, L \rightarrow 0} \left(\frac{\partial}{\partial L} \left[\frac{1}{N} \log \langle Z^L \rangle \right] \right) \\ &\approx \text{Extr} \left[\beta^* \beta \alpha [q^*]^p - \frac{1}{2} \beta^2 \alpha q^p + \beta m^p - \beta^* \beta \alpha r^* q^* + \frac{1}{2} \beta^2 \alpha r q \right. \\ &\quad - \beta m k + \frac{1}{2} \beta^2 \alpha + \lim_{L \rightarrow 0} \left(\frac{\partial}{\partial L} \left[\log \left\langle \sum_{\xi_i} \exp \left\{ \beta k \sum_b \xi_i^b \sigma_i^a \right. \right. \right. \right. \\ &\quad \left. \left. \left. + \beta^* \beta \alpha r^* \sum_b \xi_i^* \xi_i^b + \beta^2 \alpha r \sum_{b<c} \xi_i^b \xi_i^c \right\} \right] \right) \left. \right]. \end{aligned}$$

992 Furthermore, the Hubbard-Stratonovich transformation gives

$$\exp \left\{ \beta^2 \alpha r \sum_{b<c} \xi_i^b \xi_i^c \right\} \propto \exp \left\{ -\frac{1}{2} \beta^2 \alpha r L \right\} \int_{\mathbb{R}} dx \exp \left\{ -\frac{1}{2} x^2 + x \beta \sqrt{\alpha r} \sum_b \xi_i^b \right\}.$$

993 We can then use the factorization

$$\begin{aligned} &\sum_{\xi_i} \exp \left\{ \beta k \sum_b \xi_i^b \sigma_i^a + \beta^* \beta \alpha r^* \sum_b \xi_i^* \xi_i^b + x \beta \sqrt{\alpha r} \sum_b \xi_i^b \right\} \\ &= \prod_b \sum_{\xi_i^b} \exp \left\{ \beta k \xi_i^b \sigma_i^a + \beta^* \beta \alpha r^* \xi_i^* \xi_i^b + x \beta \sqrt{\alpha r} \xi_i^b \right\} \\ &= \prod_b [2 \cosh(\beta k \sigma_i^a + \beta^* \beta \alpha r^* \xi_i^* + x \beta \sqrt{\alpha r})] \\ &= 2^L \cosh^L(\beta k \sigma_i^a + \beta^* \beta \alpha r^* \xi_i^* + x \beta \sqrt{\alpha r}) \end{aligned}$$

994 in order to simplify the free energy to

$$\begin{aligned}
f &= \text{Extr} \left\{ \beta^* \beta \alpha [q^*]^p - \frac{1}{2} \beta^2 \alpha q^p + \beta m^p - \beta^* \beta \alpha r^* q^* + \frac{1}{2} \beta^2 \alpha r q - \beta m k \right. \\
&\quad \left. - \frac{1}{2} \beta^2 \alpha r + \frac{1}{2} \beta^2 \alpha + \lim_{L \rightarrow 0} \left(\frac{\partial}{\partial L} \left[\log \left\langle \sum_{\xi_i} \int_{\mathbb{R}} dx \exp \left\{ -\frac{1}{x} x^2 \right\} \right. \right. \right. \right. \\
&\quad \left. \left. \left. \exp \left\{ \beta k \sum_b \xi_i^b \sigma_i^a + \beta^* \beta \alpha r^* \sum_b \xi_i^* \xi_i^b + x \beta \sqrt{\alpha r} \sum_b \xi_i^b \right\} \right] \right) \right\} \\
&= \text{Extr} \left\{ \beta^* \beta \alpha [q^*]^p - \frac{1}{2} \beta^2 \alpha q^p + \beta m^p - \beta^* \beta \alpha r^* q^* + \frac{1}{2} \beta^2 \alpha r q - \beta m k \right. \\
&\quad \left. - \frac{1}{2} \beta^2 \alpha r + \frac{1}{2} \beta^2 \alpha + \log 2 + \lim_{L \rightarrow 0} \left(\frac{\partial}{\partial L} \left[\log \left\langle \int_{\mathbb{R}} dx \exp \left\{ -\frac{1}{x} x^2 \right\} \right. \right. \right. \right. \\
&\quad \left. \left. \left. \cosh^L (\beta k \sigma_i^a + \beta^* \beta \alpha r^* \xi_i^* + x \beta \sqrt{\alpha r}) \right] \right) \right\}.
\end{aligned}$$

995 After differentiating and taking the limit, we get

$$\begin{aligned}
f &= \text{Extr}_{m,k,q,r,q^*,r^*} \left\{ \beta^* \beta \alpha [q^*]^p - \frac{1}{2} \beta^2 \alpha q^p + \beta m^p - \beta^* \beta \alpha r^* q^* \right. \\
&\quad \left. + \frac{1}{2} \beta^2 \alpha r q - \frac{1}{2} \beta^2 \alpha r - \beta m k + \frac{1}{2} \beta^2 \alpha + \log 2 \right. \\
&\quad \left. + \int dx \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} x^2 \right\} \left\langle \log [\cosh (\beta [\sqrt{\alpha r} x + \beta^* \alpha r^* + k z])] \right\rangle_z \right\}.
\end{aligned}$$

996 In the case of $p^* = 2$ and $p \geq 3$ with finite $\alpha = \frac{Mp!}{N^{p-1}}$ and $\lambda = \frac{[\beta^*]^{p/2}}{(p/2)!} N^{p/2-1}$, the free energy
997 has the same form but with β^* replaced by λ . On the other other hand, the free energy with
998 finite $\alpha = \frac{M(p/2+1)!}{N^{p/2}}$ and η evaluates to:

$$\begin{aligned}
f &= \text{Extr}_{m,k,q^*,r^*} \left\{ \beta \eta \alpha [q^*]^p - \beta m^p - \beta \eta \alpha r^* q^* - \beta m k + \log 2 \right. \\
&\quad \left. + \left\langle \log [\cosh (\beta [\eta \alpha r^* + k z])] \right\rangle_z \right\}.
\end{aligned}$$

999 E Direct model RSB ansatz

1000 Recall that the average replicated partition function of the direct model (see Eq. 10) takes the
1001 form

$$\begin{aligned}
\langle Z^L \rangle &\approx \left\langle \sum_{\sigma} \exp \left(\beta N \sum_{\gamma} \sum_{\mu \in \Gamma_{\gamma}} \left[\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i^{\gamma} \right]^p \right. \right. \\
&\quad \left. \left. + \beta \sum_{\gamma} \sum_{\mu \in \Gamma} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^{\gamma} \dots \sigma_{i_p}^{\gamma} \right) \right\rangle.
\end{aligned}$$

1002 Introducing a new replica σ^0 , we rewrite it as

$$\begin{aligned} \langle Z^L \rangle &= \left\langle \sum_{\sigma} \exp \left(\beta N \sum_{\gamma} \sum_{\mu \in \Gamma_{\gamma}} \left[\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i^{\gamma} \right]^p \right. \right. \\ &\quad \left. \left. + \beta \sum_{\gamma} \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^{\gamma} \dots \sigma_{i_p}^{\gamma} \right) \frac{\langle Z \rangle}{\langle Z \rangle} \right\rangle, \end{aligned}$$

1003 where $Z = \sum_{\sigma_0} \exp \left(\beta \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^0 \dots \sigma_{i_p}^0 \right)$. Recall that, in the paramag-
1004 netic phase, we have (see [30] and also Appendix B)

$$\begin{aligned} \langle Z \rangle &= \exp \left(\frac{1}{2} \beta^2 \alpha + \log 2 + \mathcal{O} \left(\frac{1}{N^{p/2-2}} \right) \right) \\ &= Z \exp \left(\mathcal{O} \left(\frac{1}{N^{p/2-2}} \right) \right), \end{aligned}$$

1005 so $\langle Z^L \rangle$ can be expressed as

$$\begin{aligned} \langle Z^L \rangle &= \frac{1}{\langle Z \rangle} \left\langle \sum_{\sigma} \exp \left(\beta N \sum_{\gamma} \sum_{\mu \in \Gamma_{\gamma}} \left[\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i^{\gamma} \right]^p \right. \right. \\ &\quad \left. \left. + \beta \sum_{\gamma} \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^{\gamma} \dots \sigma_{i_p}^{\gamma} \right) \right. \\ &\quad \left. \sum_{\sigma_0} \exp \left(\beta \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^0 \dots \sigma_{i_p}^0 + \mathcal{O} \left(\frac{1}{N^{p/2-2}} \right) \right) \right\rangle. \end{aligned}$$

1006 The $\mathcal{O} \left(\frac{1}{N^{p/2-2}} \right)$ corrections vanish to leading order in N when we calculate the free entropy.

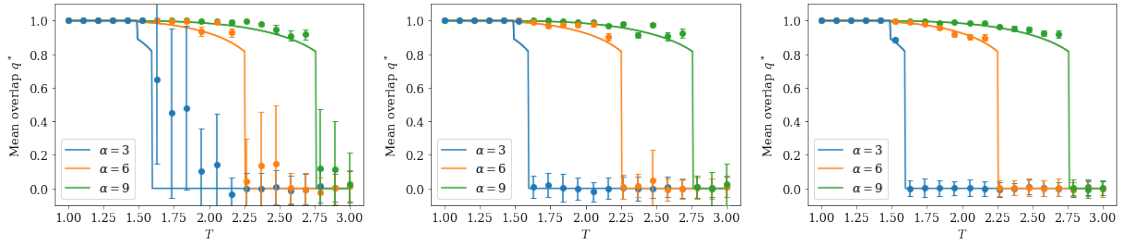
1007 **F Monte-Carlo simulations for various system sizes**


Figure 9: Monte-Carlo simulations of the $p = 3$ inverse model compared against saddle-point solutions for different values of N . The lR phase is not included in these plots. The left plot has $N = 128$, the center plot has $N = 256$, and the right plot has $N = 512$. The dots are simulation data at a few values of α , and the lines are slices of the saddle-point solutions at the same α . There are $M = \frac{\alpha N^{p-1}}{p!}$ examples σ^α , and simulation results are averaged over $L = 100$ student patterns. The simulation data is sometimes systematically shifted up with respect to the saddle-point solution, but the size of the difference tends to decrease with N . The shift is the most visible when $\alpha = 6$ and right after the fall from eR to gR when $\alpha = 3$. As expected, the fluctuations of the paramagnetic phase also decrease with N .