

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	3
5	3 Teacher-student setting	5
6	4 Results and Discussion	8
7	4.1 Transition to the ordered phases: universality	8
8	4.2 Phase diagram on the Nishimori line	10
9	4.3 Inference temperature vs dataset noise	13
10	4.4 Interaction order and noise tolerance	13
11	4.5 Robustness against adversarial attacks	17
12	5 Conclusion	19

13	References	20
14	A Gardner’s Hamiltonian vs Krotov’s Hamiltonian	25
15	B Direct model cumulant expansions	27
16	C Teacher-student replicated partition function	30
17	D Teacher-student free entropy	32
18	E Direct model RSB ansatz	40
19	F Monte-Carlo simulations for various system sizes	41

20

21

22 1 Introduction

23 Hopfield networks are artificial neural networks that model associative memory [1]. In the
 24 Hopfield model, examples $\sigma \in \{-1, 1\}^N$ of memories $\xi^\mu \in \{-1, 1\}^N$, $\mu = 1, \dots, M$, are
 25 retrieved by sampling the Gibbs distribution of a 2-body Hamiltonian $H[\sigma|\xi]$ at a given
 26 temperature T [2]. Hopfield networks can be trained in a biologically plausible way using
 27 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
 28 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
 29 find this scaling is to study the phase diagram of $H[\sigma|\xi]$ as a function of the temperature T
 30 and load $\alpha = \frac{M}{N}$ [5], where the so-called ferromagnetic phase, which extends up to $\alpha \approx 0.14$,
 31 corresponds to accurate retrieval.

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

45 Although they draw a rather detailed picture of the retrieval of individual i.i.d. random
 46 memories, these results are not the end of the story. First of all, there is still significant
 47 uncertainty on the location of the paramagnetic to spin-glass phase transition because of
 48 numerical instability. Second of all, dense Hopfield networks have been rapidly gaining a
 49 renewed attention for reasons other than their storage capacity since Krotov’s recent paper [26],
 50 where they were used as a trainable machine learning architecture. For instance, they have
 51 been related to transformers [23, 34] and diffusion models [35, 36], and they were found to be
 52 significantly more explainable and adversarially robust than feedforward neural networks with

53 ReLU activation functions [26, 37].

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

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

76 In this work, we study the performance of p -body Hopfield networks in the teacher-student
77 setting, revealing a prototype regime and a feature regime as in the 2-body model. In Section
78 2, we review Gardner’s main results in studying p -body Hopfield models. In Section 3, we
79 compute the phase diagram of these p -body models in the teacher-student setting. In Section
80 4.1, we discuss the transition to the retrieval phase in the inverse problem and compare it
81 against the transition to the spin-glass phase in the direct problem. Despite their different
82 nature, we show that these two transitions are equivalent on the Nishimori line where the
83 teacher and the student have the same p and T [45–48]. In Section 4.2, we discuss the phase
84 diagram on the Nishimori line in more details. In Section 4.3 and Section 4.4, we discuss the
85 phase diagram outside of the Nishimori line. First of all, we investigate the effect of using an
86 inference temperature different from the dataset noise. Second of all, we reveal that using
87 a larger p for the student than the teacher gives the student an extensive tolerance against
88 both teacher noise and pattern interference. Finally, in Section 4.5, we derive a closed-form
89 expression that measures the adversarial robustness of the student at zero temperature and
90 explain what our results reveal about the nature of adversarial attacks.

91 2 Overview of Gardner’s results

92 Consider the p -body Hamiltonian

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

93 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
94 i.i.d. from the Rademacher distribution $\frac{1}{2} [\delta(\xi_i^\mu - 1) + \delta(\xi_i^\mu + 1)]$. In the *direct model*,

95 patterns σ are in turn sampled from the equilibrium Gibbs distribution $P(\sigma|\xi) = Z^{-1}e^{-\beta H[\sigma|\xi]}$,
 96 where $\beta \geq 0$ is the inverse temperature and $Z = \sum_{\sigma} e^{-\beta H[\sigma|\xi]}$ is the system's partition function.
 97 The so-called *direct problem* studied by Gardner [30] consists of quantifying the performance
 98 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
 99 good measure of retrieval accuracy, and its expected value can be derived from the quenched
 100 free entropy $f = \frac{1}{N} \langle \log Z \rangle_{\xi}$ in the thermodynamic limit $N \rightarrow \infty$. At finite p , Gardner used
 101 the (non-rigorous) replica trick [49] to evaluate the RS approximation of f (see also Appendix
 102 B) in terms of a variational principle of the form

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

103 whose solution is

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

104 and the order parameters m and q are to be interpreted as expected overlaps. To be more
 105 precise, m can be shown to be the expected overlap of a retrieval attempt σ against one memory
 106 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
 107 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
 108 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
 109 to stay frozen in specific configurations rather than visiting all possible values of σ .

110 The resulting RS phase diagram (see Fig. 1) can be derived from the value of the order
 111 parameters as a function of three *hyperparameters*: the interaction order p , temperature
 112 $T = 1/\beta$ and load $\alpha = \frac{Mp^1}{N^{p-1}}$. There are four different phases:

- 113 • In the Paramagnetic phase (**P**), the overlaps m and q both vanish. The network does not
 114 retrieve any specific pattern: sampled configurations are completely random.
- 115 • In the Spin-Glass phase (**SG**), m vanishes but $q > 0$. In other terms, the network does not
 116 retrieve individual stored memories but rather converges to spurious patterns depending
 117 on all the memories in a non-trivial way.
- 118 • In the signal Retrieval phases (**lR** and **gR**), $m \neq 0$ and $q > 0$, which means that the
 119 network is able to retrieve the stored memories. **lR** and **gR** are respectively locally
 120 stable and globally stable. In other words, local retrieval **lR** is only attainable from
 121 initial conditions in a limited neighborhood of a memory ξ^{μ} , while global retrieval **gR** is
 122 accessible from any initial conditions given enough time. These two phases are said to
 123 be ferromagnetic.

124 Gardner also calculated the exact $p \rightarrow \infty$ phase diagram without making any assumptions
 125 about replica symmetry [30]. The resulting paramagnetic to spin-glass (**P-SG**) phase transition
 126 line, given by $\beta^2 \alpha = 2 \log 2$, coincides with the boundary of the region where the total
 127 entropy of the paramagnetic phase becomes negative. Therefore, the total entropy of the
 128 system is always positive, as expected. Conversely, the finite p phase diagrams obtained under

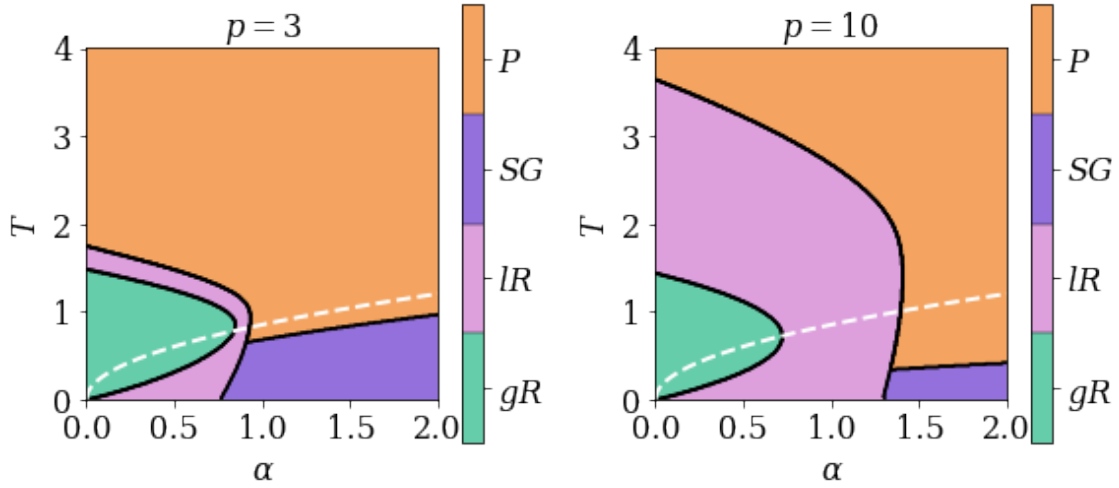


Figure 1: RS phase diagrams of the direct models with $p = 3$ on the left and $p = 10$ on the right. Accurate pattern retrieval is not possible in the paramagnetic phase (P) or in the spin-glass phase (SG), but it is possible in the local retrieval phase (lR) and in the global retrieval phase (gR). The ferromagnetic fixed point corresponding to accurate pattern retrieval is globally stable in the gR phase, but locally stable in the lR phase. The phase diagrams are inexact below the white dashed line where the total entropy of the paramagnetic phase becomes negative.

129 the RS approximation have a limited scope because RSB must be taken into account in the
 130 $\beta^2 \alpha > 2 \log 2$ region where the RS approximation of the total entropy becomes negative [30,50].
 131 Although these phase diagrams were recently extended by including the effects of 1RSB [33],
 132 there is still significant uncertainty on the location of the large- p P - SG transition because of
 133 numerical instability issues. In fact, the finite- p transition evaluated in [33] does not seem to
 134 get closer to the exact $p \rightarrow \infty$ transition as p gets larger.

135 3 Teacher-student setting

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

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

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

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

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

148 In sum, the student posterior distribution is that of a dense Hopfield network where ξ plays the
 149 role of the sampled pattern and the examples σ act like the M quenched memories. Our task,
 150 called the *inverse problem*, consists of quantifying the student's capability to infer the teacher
 151 pattern, which we will also call the *signal*. Like Gardner, we calculate a free entropy of the
 152 form $f = \frac{1}{N} \langle \log Z \rangle_{\sigma}$ in the thermodynamic limit $N \rightarrow \infty$. This time, however, the average
 153 $\langle \cdot \rangle_{\sigma}$ is over a structured set of examples σ . In fact, we recall that, unlike the i.i.d. memories
 154 studied by Gardner, the examples σ^a are sampled from the teacher distribution $P(\sigma^a | \xi^*)$.

155 In general, the student does not have access to the teacher generative model. In our
 156 controlled teacher-student setting, the student knows that the correct model for $P(\sigma^a | \xi)$ is a
 157 dense Hopfield network. Nevertheless, it does not necessarily have access to the interaction
 158 order p^* and inverse temperature β^* used by the teacher. Therefore, we denote the student
 159 hyperparameters by p and β and emphasize that they are not necessarily equal to p^* and β^* .
 160 As previously stated, we calculate the free entropy

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

161 in the thermodynamic limit $N \rightarrow \infty$. We then draw phase diagrams of the inverse problem as
 162 a function of p^* , $T^* = 1/\beta^*$, p , $T = 1/\beta$ and α .

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

171 Conversely, when T^* is above T_{crit} , the examples in the training set are very noisy and
 172 we do not observe a finite overlap between σ^a and ξ^* . In this regime, we find that the RS
 173 approximation of the free entropy can be computed (see Appendix D) in terms of the variational
 174 principle

$$f = \text{Extr}_{m,k,q,r,q^*,r^*} \left\{ \beta^* \beta \alpha [q^*]^p - \frac{1}{2} \beta^2 \alpha q^p + \beta m^p - \beta^* \beta \alpha r^* q^* \right. \\ \left. + \frac{1}{2} \beta^2 \alpha r q - \frac{1}{2} \beta^2 \alpha r - \beta m k + \frac{1}{2} \beta^2 \alpha + \log 2 \right. \\ \left. + \int dx \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} x^2 \right\} \left\langle \log [\cosh (\beta [\sqrt{\alpha r} x + \beta^* \alpha r^* + k z])] \right\rangle_z \right\}, \quad (3)$$

175 whose solution is the saddle-point equations

$$\begin{aligned}
q^* &= \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) \left\langle \tanh(\beta [\sqrt{ar}x + \beta^* ar^* + kz]) \right\rangle_z \\
q &= \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) \left\langle \tanh^2(\beta [\sqrt{ar}x + \beta^* ar^* + kz]) \right\rangle_z \\
m &= \int_{\mathbb{R}} dx \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) \left\langle z \tanh(\beta [\sqrt{ar}x + \beta^* ar^* + kz]) \right\rangle_z \\
r^* &= p [q^*]^{p-1} \\
r &= p q^{p-1} \\
k &= p m^{p-1},
\end{aligned} \tag{4}$$

176 where z is a Rademacher random variable. As in the direct model described in Section 2,
177 the order parameters m and q have a clear interpretation in terms of expected overlaps.
178 $m = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \xi_i \sigma_i^a \right\rangle_{\xi^*, \sigma, \xi}$ is the expected overlap of a retrieval attempt with an ex-
179 ample 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 retrieval
180 attempts. Similarly, q^* is the expected overlap between the teacher and student patterns, i.e.
181 $q^* = \lim_{N \rightarrow \infty} \left\langle \frac{1}{N} \sum_i \xi_i^* \xi_i \right\rangle_{\xi^*, \sigma, \xi}$. Therefore, it is a good measure of inference performance.
182 The free entropy (Eq. 3) is expected to be exact in absence of mismatch between the teacher
183 and the student, i.e. $\beta^* = \beta$. This condition is known as the Nishimori line [45–48]. Outside
184 of the Nishimori region, RSB corrections are expected. Like in the direct problem, there are
185 different phases characterized by the values of the order parameters:

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

194 In sum, when T^* is above T_{crit} , the student can only learn the teacher pattern in the signal
195 retrieval phases. In all the other phases, the student pattern is uncorrelated with the signal,
196 being either a random guess (P phase), aligned with a noisy example (inaccurate eR phase),
197 or aligned with a spurious low energy state (SG phase).

198 We also investigate the $T^* > T_{\text{crit}}$ regime in the presence of a mismatch between the
199 interaction orders of the teacher and student networks, i.e. $p^* \neq p$. We focus on the case of
200 $p^* = 2$ and even $p \geq 3$ to study the consequences of fitting the teacher of [38] using a student
201 with higher order interactions. We find two different scaling regimes of the training set size M
202 and inverse temperature β^* that make retrieval possible (see Appendix D):

- 203 • a large-noise scaling where $\beta^* \sim \mathcal{O}(N^{2/p-1})$ and $M \sim \mathcal{O}(N^{p-1})$, such that $\alpha = \frac{Mp!}{N^{p-1}}$
204 and $\lambda = \frac{[\beta^*]^{p/2}}{(p/2)!} N^{p/2-1}$ are finite;
- 205 • 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
206 finite.

207 In the large-noise scaling, we obtain saddle point equations similar to Eqs. (4) but with β^*
 208 replaced by λ (see Appendix D). Conversely, the finite noise scaling leads to

$$\begin{aligned} q^* &= \left\langle \tanh(\beta [\eta \alpha r^* + k z]) \right\rangle_z \\ m &= \left\langle z \tanh(\beta [\eta \alpha r^* + k z]) \right\rangle_z \\ r^* &= p [q^*]^{p-1} \\ k &= p m^{p-1}, \end{aligned} \quad (5)$$

209 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
 210 $p = 4$ (see Appendix D). These equations can also be derived by extrapolating the large-noise
 211 equations equations to $\alpha_{\text{large noise}} \rightarrow 0$ and $\lambda \rightarrow \infty$ with fixed $\lambda \alpha_{\text{large noise}} = \eta \alpha_{\text{finite noise}}$.

212 4 Results and Discussion

213 4.1 Transition to the ordered phases: universality

214 The paramagnetic solution of Eqs. (4) always exists and is globally stable in the part of the
 215 phase diagram where the T is relatively large and α is relatively small. On the other hand, the
 216 gR phase exists when $\beta^2 \alpha p$ and $\beta^* \beta \alpha p$ are both large. In fact, in that limit, $q^* = q = 1$ is
 217 a fixed point of Eqs. (4). The critical line where gR becomes globally stable instead of P is
 218 not clear from this analysis alone, but we can at least find it analytically in the limit of infinite
 219 p . As for the direct model, the free entropy and the total entropy of the paramagnetic phase
 220 are respectively $\frac{1}{2} \beta^2 \alpha + \log 2$ and $-\frac{1}{2} \beta^2 \alpha + \log 2$ [30]. At the same time, the $p \rightarrow \infty$ free
 221 entropy takes the form

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

222 where $\theta[q] := \lim_{p \rightarrow \infty} q^p$, $q \in [0, 1]$, is the Heaviside step function jumping at $q = 1$, i.e.
 223 $\theta(1) = 1$ and $\theta(q) = 0 \forall q \in [0, 1)$. In this limit, the ferromagnetic phase is characterized by
 224 $q = q^* = 1$, and its free entropy is then

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

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

$$T < \sqrt{\frac{\alpha}{2 \log 2}}, \quad (6)$$

228 where $T = \sqrt{\frac{\alpha}{2 \log 2}}$ is also the temperature below which the total entropy of the paramag-
 229 netic phase becomes negative. Outside of the Nishimori line, this inequality generalizes to
 230 $\beta^* \beta \alpha > \frac{1}{2} \beta^2 \alpha + \log 2$, leading to

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

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

232 In the $p \rightarrow \infty$ limit, the transition towards \mathbf{gR} of the inverse model on the Nishimori line
 233 is identical to the exact $\mathbf{P-SG}$ transition of the direct model [30]. We claim that these two
 234 critical lines are actually closely related for any p . In the Hopfield model with $p = 2$, they
 235 were already shown to be identical [38]. We will now argue that they overlap for any p and
 236 β such that the inverse model is outside of the \mathbf{eR} phase. In the case of $p = 2$, both lines
 237 can be obtained exactly from the RS approximation of either the direct model or the inverse
 238 model, so there is no obvious advantage to using this equivalence in calculations. In general,
 239 while the inverse problem on the Nishimori line is replica symmetric, the direct problem is not,
 240 and the $p \geq 3$ replica symmetric $\mathbf{P-SG}$ transition is not exact. Moreover, even the critical line
 241 calculated using 1RSB may be inaccurate due to numerical instability [33]. In this situation,
 242 the knowledge of the \mathbf{gR} transition in the replica-symmetric inverse problem can be used to
 243 locate the exact $\mathbf{P-SG}$ transition of the direct problem, where symmetry breaking occurs.

244 For that purpose, we will argue that *the direct model is in the paramagnetic phase if and only*
 245 *if the inverse model is in the paramagnetic phase.*

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

$$P(\sigma) = \frac{1}{2^{MN}} \frac{Z(\sigma)}{\langle Z \rangle}, \quad (7)$$

247 the example distribution $P(\sigma)$ of the inverse problem is contiguous [51] to the uniform
 248 distribution, i.e. the memory distribution of the direct problem, when

$$\lim_{N \rightarrow \infty} \left\{ \frac{\log Z - \log \langle Z \rangle}{N} \right\} = 0. \quad (8)$$

249 As determined in Appendix C and D, the annealed expression $\frac{1}{N} \log \langle Z \rangle$ is equal to the free
 250 entropy of the paramagnetic phase. Therefore, when the inverse model is in the paramagnetic
 251 phase, $P(\sigma)$ is contiguous to the uniform distribution. This property is called quiet planting
 252 and is known to occur more generally in mean-field paramagnets [52–55]. In our problem
 253 setting, it means that if the inverse model is in the paramagnetic phase, then it is equivalent to
 254 the direct model. In particular, if the inverse model is in the paramagnetic phase, then so is the
 255 direct model. In more intuitive terms, the \mathbf{gR} transition temperature of the inverse model must
 256 be greater than or equal to the $\mathbf{P-SG}$ transition temperature of the direct model because the
 257 ensemble of examples σ^a generated by the teacher model is on average at least as structured
 258 as the set of i.i.d. random memories stored in the direct model.

259 For the direct implication, notice that the average replicated partition function of the direct

260 model in the paramagnetic phase can be approximated as (see Appendix E)

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

261 This expression is identical to the replicated partition function of the inverse model outside of
262 the eR phase, which therefore must also be in the paramagnetic phase.

263 As a consequence, the P - SG transition line of the direct model must be identical to the gR
264 transition line of the inverse model on the Nishimori line.

265 4.2 Phase diagram on the Nishimori line

266 On the Nishimori line, the student is fully informed about the teacher generative model and
267 uses $\beta = \beta^*$ and $p = p^*$. In this scenario, thanks to the Nishimori identities [46], it is well
268 known that ξ^* and ξ play symmetric roles and that $q^* = q$. For the same reason, the overlaps
269 $\frac{1}{N} \sum_i \xi_i^* \xi_i$ 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
270 follows that the system is expected to be replica symmetric, and Eqs. (3) and (4) are expected
271 to hold. Fig. (2) shows the phase diagrams obtained by solving the saddle-point equations
272 numerically on the Nishimori line. Both $q^* = q$ and the replica symmetry condition are verified.
273 In particular, numerical solutions of a few values of $p \geq 3$ show that the gR transition occurs
274 at a higher T than the line $\beta^2 \alpha = 2 \log 2$ where the total entropy of the paramagnetic phase
275 becomes negative. In other terms, the phase transition towards gR prevents the total entropy
276 from becoming negative when T decreases below $\sqrt{\frac{\alpha}{2 \log 2}}$, which is consistent with the RS
277 solution being exact on the Nishimori line.

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

281 At high T , learning is possible only if the amount of examples, i.e. the size of the dataset, is
282 sufficiently large. When α is too small, Eqs. (4) have only a paramagnetic fixed point because
283 the amount of information carried by the dataset is not large enough. Numerical solutions
284 suggest that the paramagnetic fixed point always exist and it is actually locally stable in the
285 whole high-temperature regime. When α is sufficiently large, the signal retrieval fixed point
286 appears as a locally stable attractor. It becomes globally stable as the size of the dataset is
287 increased further.

288 The critical boundary of the gR phase can be obtained by solving the saddle-point equations
289 numerically (Eqs. 4), and the result is consistent with the analytical $p \rightarrow \infty$ gR boundary of
290 Eq. (6). In fact, we find that the analytical boundary closely agrees with the numerical solution
291 of the saddle-point equations with $p^* = p = 10$ and remains a good approximation even down
292 to $p^* = p = 4$.

293 In the student model, σ plays a similar role as the weights of the trainable dense Hopfield
294 network model that Krotov designed for classification of data [26]. In that context, ξ is
295 analogous to the test data whose labels are being predicted (see Fig. 3). In fact, the computation
296 performed by Krotov's model to recover labels is similar to the update rule used by the student
297 to infer the teacher pattern (see Appendix A). Moreover, the eR and gR phases are respectively

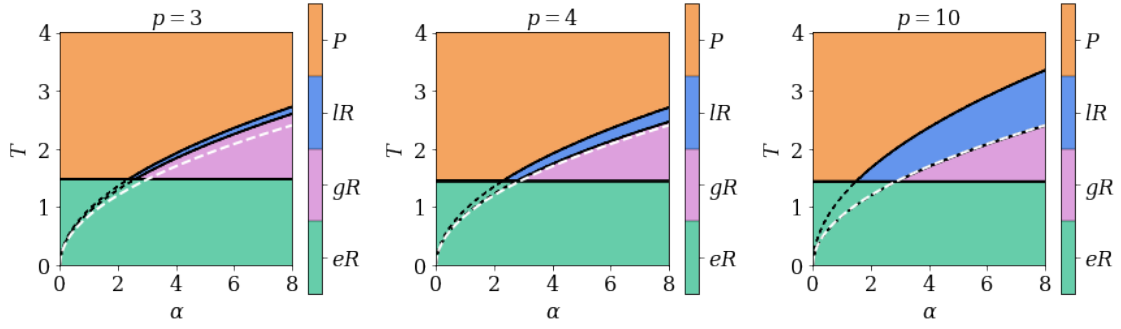


Figure 2: Exact replica-symmetric phase diagrams of the inverse models with $p = 3$ on the left, $p = 4$ in the center and $p = 10$ on the right. Accurate pattern retrieval is not possible in the paramagnetic phase (P), but it is possible in the local retrieval phase (IR), in the global retrieval phase (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 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 at the end of Section 3. It matches the corresponding numerical phase boundary increasingly well as p grows larger.

298 reminiscent of the prototype and feature regimes of Krotov’s networks. Therefore, we believe
 299 that the student can act as a toy model of label prediction in these two regimes.

300 Comparing instead the phase diagrams of our inverse model with that of the inverse 2-body
 301 Hopfield model, we see that the eR and gR phases of the inverse p -body model with $p \geq 3$ are
 302 respectively analogous to the eR and sR (signal Retrieval) phases presented in [38]. One of
 303 the key differences between $p = 2$ and $p \geq 3$ is that the paramagnetic to signal retrieval phase
 304 transition of the p -body model is second order for $p = 2$ but first order for $p \geq 3$. On the one
 305 hand, the second order phase transition of $p = 2$ indicates that its paramagnetic fixed point is
 306 never locally stable and sets an unambiguous boundary between the sR phase where ξ^* can
 307 be recovered starting from any initial conditions and the paramagnetic phase where pattern
 308 retrieval is impossible [55]. On the other hand, the first order phase transition of $p \geq 3$ allows
 309 the retrieval and paramagnetic regimes to coexist. The IR phase is locally stable precisely
 310 because it coexists with the paramagnetic phase and has a lower free entropy. Meanwhile,
 311 the gR phase also coexists with the paramagnetic phase, but has a larger free entropy. In the
 312 presence of phase coexistence, an algorithm trying to retrieve ξ^* starting from random initial
 313 conditions can get stuck in the paramagnetic phase instead. In fact, it has been conjectured
 314 that there is no algorithm with random initial conditions that can find such a ferromagnetic
 315 fixed point in a tractable amount of time [55, 56]. That kind of metastable region was thus
 316 given the name *hard phase* [55, 57]. In summary, we expect that $p \geq 3$ models in the gR phase
 317 can only recover partially corrupted patterns whereas $p = 2$ can recover them entirely.

318 Fig. (4) shows results from Monte Carlo simulations with $p = 3$, where L replicas of
 319 the student pattern $\{\xi^b\}_{b=1}^L$ are initialized to the teacher pattern ξ^* corrupted by some
 320 Rademacher noise ε . In other words, the initial values of ξ_i^b are sampled from the distri-
 321 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
 322 so that the simulations start relatively close to the saddle-point solutions. As explained pre-
 323 viously, gR is a hard phase, so this initialization is necessary to make ξ^b converge to gR in a
 324 reasonable amount of time. Additionally, it is also used to make ξ^b converge to the IR phase
 325 rather than the P phase when desired. Once the simulations are over, the overlaps are averaged

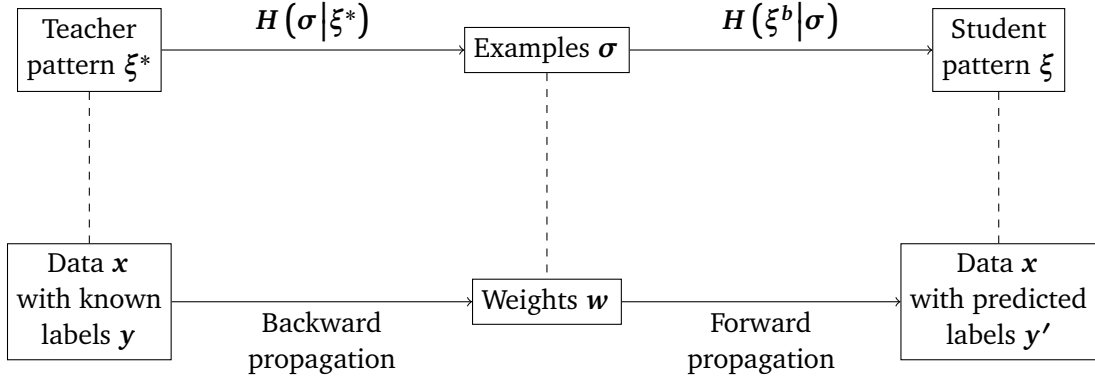


Figure 3: The first row of this diagram sketches how a p -body Hopfield network in the teacher-student setting can reconstruct an incomplete pattern ξ^b to match the teacher pattern ξ^* by relying on the examples σ obtained from ξ^* . The second row summarizes how a dense neural network trained by Krotov 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.

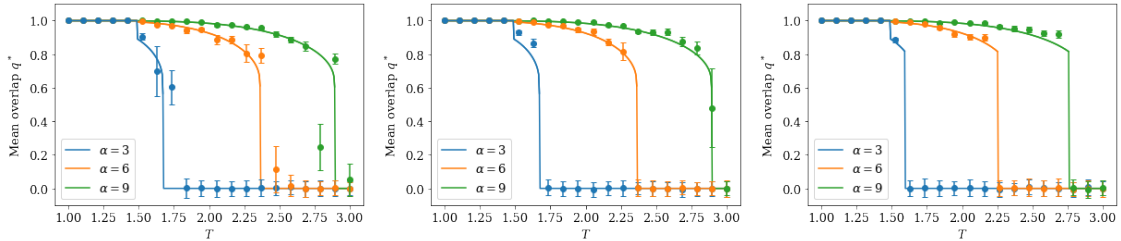


Figure 4: Monte-Carlo simulations of the $p = 3$ inverse model compared against RS saddle-point solutions. The lR phase is included on the left and central plots, but not on the right one. The left plot has $\varepsilon = 0$, and the two other ones have a handpicked ε such that the simulations are initialized 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 $\alpha = 3$.

326 over all L replicas. If we fix $\boldsymbol{\varepsilon} = \mathbf{0}$, then the simulations generally converge to the \mathbf{IR} phase when
 327 it is a fixed point. If instead we initialize them to the saddle-point solutions by handpicking $\boldsymbol{\varepsilon}$,
 328 then they stay near the initial overlaps. In either case, the simulations converge to \mathbf{eR} when it is
 329 globally stable. Some simulation data points might be systematically shifted up with respect to
 330 the saddle-point solutions. However, this difference decreases with the system size N , so finite
 331 size effects seem sufficient to explain it (see Fig. 9 in Appendix F). Overall, the Monte-Carlo
 332 simulations are in very good agreement with the $p = 3$ overlap landscape obtained by solving
 333 the saddle-point equations numerically.

334 4.3 Inference temperature vs dataset noise

335 In the two next Sections, we will discuss the phase diagram when the student is only partially
 336 informed about the teacher generative model, i.e. when the Nishimori conditions do not hold.
 337 We start with the case where $p = p^*$ but $\beta \neq \beta^*$, i.e. the inference temperature T is different
 338 from the dataset noise T^* . As we argued in Section 3, the student accurately retrieves ξ^* when
 339 $T^* < T_{\text{crit}}$. On the other hand, we must solve the saddle-points equations (see Eqs. 4) to study
 340 $T^* > T_{\text{crit}}$.

341 We show the phase diagram of this region on Fig. (5). At high inference temperature T , the
 342 situation is similar to Fig. (2): retrieval is possible if the data load α is sufficiently large, but
 343 the paramagnetic phase is always locally stable. The situation is different when the inference
 344 temperature is low. In that case, there are two phases that we did not see for $\beta = \beta^*$: the
 345 inaccurate \mathbf{eR} phase and the \mathbf{SG} phase. When α is relatively small, the student falls in the
 346 inaccurate \mathbf{eR} phase. In this regime, it has finite overlap with one of the noisy examples and
 347 cannot retrieve the signal ξ^* . When α is larger, the interference among the noisy examples
 348 prevents the student to be aligned with them. In this regime, the \mathbf{SG} phase, the student locally
 349 converge to spurious patterns that are uncorrelated with the signal.

350 Accurate pattern retrieval is only possible in the \mathbf{IR} and \mathbf{gR} phases where α is so large that
 351 the student can gather enough information from the dataset to become very close to ξ^* . The
 352 phase diagrams indicate that pattern retrieval is optimal on the Nishimori line in the sense that
 353 $\beta = \beta^*$ is the inverse temperature where the student needs the least examples to recover ξ^* .
 354 In other words, the student's performance is non-monotonic in T and peaks at $T = T^*$. These
 355 properties were also observed in the teacher-student setting of the $p = 2$ Hopfield network [38].

356 Contrary to what one would expect to see on the exact phase diagram [45,46], the Nishimori
 357 line $T = T^*$ does not cross a triple point on the RS phase diagram. The issue is that the RS
 358 phase diagram is not exact outside of the Nishimori line. In particular, the \mathbf{SG} phase boundary
 359 is not exact. Outside of the retrieval regime, the free entropy of the inverse model is the same
 360 as the direct model. Since the transition towards \mathbf{gR} of the inverse model on the Nishimori
 361 line overlaps the exact $\mathbf{P-SG}$ transition of the direct model (see Section 4.2), we deduce that
 362 it must also overlap the exact $\mathbf{P-SG}$ transition of the *inverse* model outside of the \mathbf{gR} phase.
 363 Plotting it on the RS phase diagrams, we see that it indeed crosses the Nishimori line and the
 364 \mathbf{gR} phase boundary at the same point, which therefore becomes a triple point, as expected.

365 4.4 Interaction order and noise tolerance

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

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

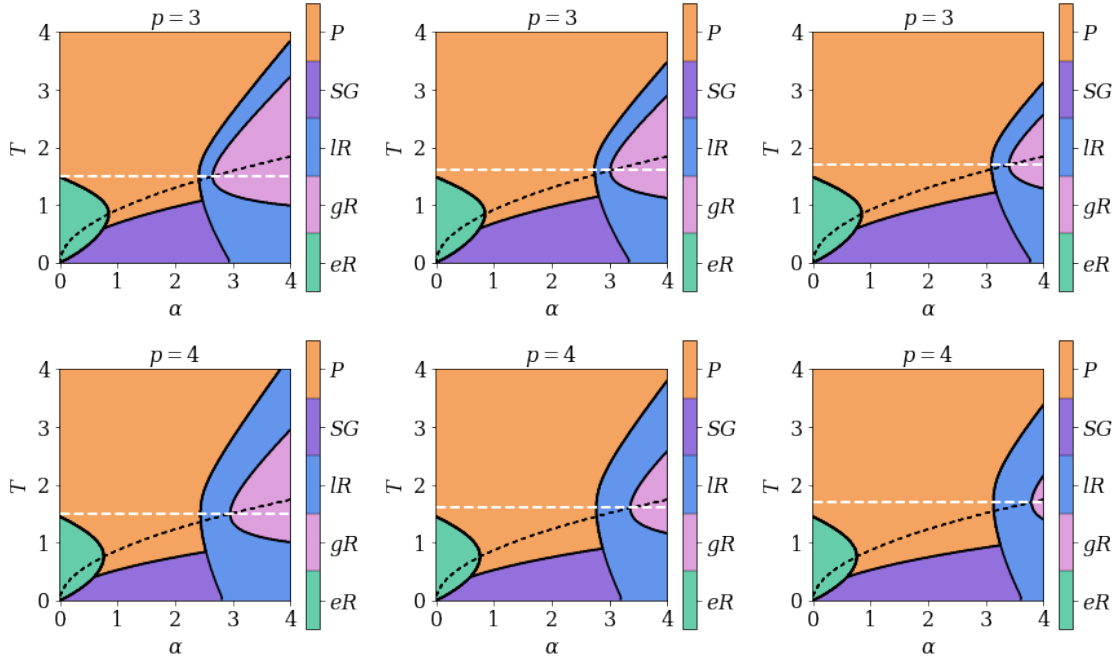


Figure 5: RS phase diagrams of inverse models with $p^* = p$ and fixed β^* . The top and bottom rows of plots respectively have $p^* = p = 3$ and $p^* = p = 4$. In the same way, the left, central and right columns correspond to $T^* = 1.5$, $T^* = 1.6$ and $T^* = 1.7$. Accurate pattern retrieval is not 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 line is the gR phase boundary on the Nishimori line. As explained in Section 4.2, we expect it to overlap the exact SG phase transition.

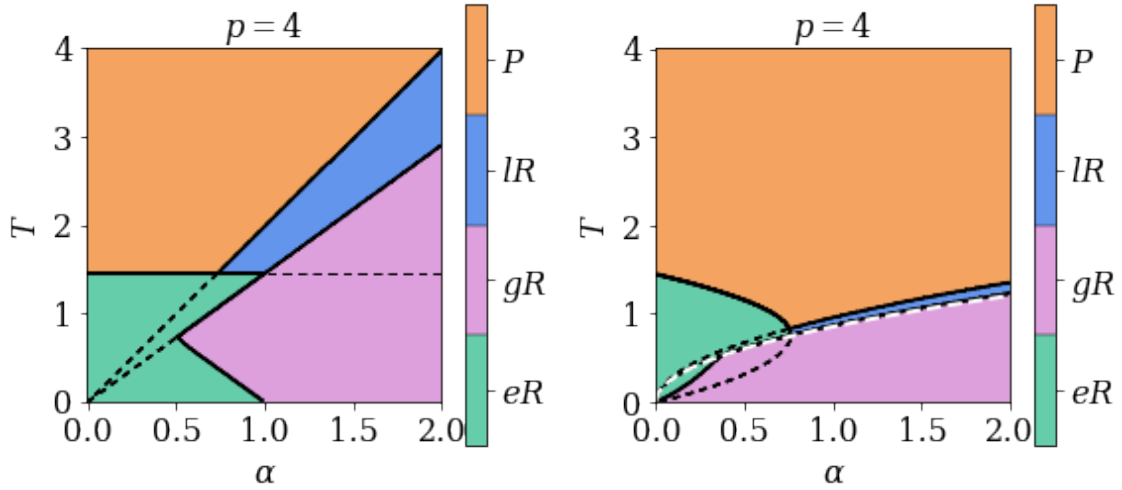


Figure 6: RS phase diagrams of inverse models with $p^* = 2$ and $p = 4$. The left plot is for $\alpha = \frac{M(p/2+1)!}{N^{p/2}}$, and $\beta^* = 1 - \frac{1}{\sqrt{2}}$ such that $\eta = 1$ and the right plot is for $\alpha = \frac{Mp!}{N^{p-1}}$ and $\beta^* = \sqrt{\frac{2\lambda}{N}}$ with $\lambda = \beta$. 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.

that

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

371 are finite. In this scaling, a $p \geq 3$ network requires $\mathcal{O}(N^{p-2})$ more training examples than a
 372 $p = 2$ network with finite load $\gamma = \frac{M}{N}$, but also has a higher tolerance to teacher noise. For
 373 instance, a student with $p = 4$ interactions is able to retrieve the pattern of a teacher with
 374 $T^* \sim \mathcal{O}(N^{1/2})$ noise when it is shown enough examples $M \sim \mathcal{O}(N^3)$ to be in the gR phase
 375 (see Fig. 6).

376 $\mathcal{O}(N^{1/2})$ noise tolerance was also observed in the $p = 4$ direct model, where it is a
 377 consequence of the redundancy stemming from storing $\mathcal{O}(N)$ memories rather than the $\mathcal{O}(N^3)$
 378 needed to saturate the storage capacity [58]. Our $p = 4$ inverse model exploits a different
 379 kind of redundancy by learning from $\mathcal{O}(N^3)$ examples whereas $p = 2$ only needs $\mathcal{O}(N)$. In
 380 other terms, both storing extensively less memories than the maximum allowed amount and
 381 generating extensively more examples than the minimum required amount provide enough
 382 redundancy to recover a pattern muddled in an extensive amount of noise. In both cases, there is
 383 an $\mathcal{O}(N^2)$ gap between the number of patterns used in the noise-tolerant and noise-susceptible
 384 regimes. Going beyond $p = 4$, the inverse model has $\mathcal{O}(N^{1-2/p})$ noise tolerance as a function
 385 of p . In particular, our theory predicts that the tolerance saturates at $T^* \sim \mathcal{O}(N)$ as $p \rightarrow \infty$,
 386 but at the cost of using an intractable number of examples. This behavior is different from
 387 the $\mathcal{O}(N^{1/2-p/4})$ tolerance of the direct p -body model in the noisy-learning regime studied
 388 in [59]. In other terms, the dataset noise that we are facing is of a different nature than the

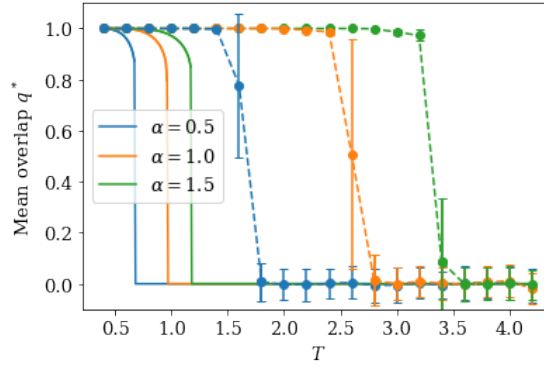


Figure 7: Monte-Carlo simulations (dashed lines) and RS saddle-point solutions (full lines) of the inverse model in the large-noise scaling with $p^* = 2$ and $p = 4$. The teacher generates $M = \frac{\alpha N^{p-1}}{p!}$ examples σ^α with $N = 256$ components each, and the simulation results are then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

389 learning noise of [59]. In any case, it is interesting that both the direct and inverse models
 390 are able to tolerate an extensive amount of noise. Overall, our results suggest that it could be
 391 advantageous to use a student network with a relatively large p to learn from a large but noisy
 392 dataset when the p^* of the teacher generative model is unknown.

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

401 Given $m = 0$, the free entropy of the inverse model with $p \geq 3$, $p^* = 2$ and $\beta = \lambda$ is
 402 the same as on the Nishimori line (see Eq. 4 and Appendix D). As a direct consequence, the
 403 total entropy is positive outside of the eR phase (see Fig. 6). Additionally, the $p^* = 2$, $p \geq 3$
 404 phase diagrams with $\beta \neq \lambda$ are identical to the $p = p^*$ phase diagrams with $\beta \neq \beta^*$, which
 405 suggests that $\beta = \lambda$ is optimal for $p^* = 2$, $p \geq 3$ in the same sense as $\beta = \beta^*$ is optimal for
 406 $p = p^*$ (see Fig. 5). Monte-Carlo simulations confirm that a student with $p \geq 3$ is able to
 407 retrieve the pattern of a teacher with $p = 2$ and $T^* \sim \mathcal{O}(N^{1/2})$ (see Fig. 7). However, the IR
 408 phase transition is at a higher T in the simulations than on the $\beta = \lambda$ RS phase diagram (see
 409 Fig. 5), which means that RSB is necessary to describe it accurately. One could check where
 410 replica symmetry holds by evaluating the stability of the RS saddle point throughout the phase
 411 diagram.

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

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

412 is finite. In this finite-noise scaling, $p \geq 3$ requires $\mathcal{O}(N^{p/2-1})$ more training examples than
 413 $p = 2$, which is a lot less than the first scaling. For instance, a student with $p = 4$ needs $\mathcal{O}(N^2)$
 414 examples to retrieve ξ^* . As before, the phase transitions are all first order, the overlap q^* stays
 415 high throughout the gR and IR phase of $p = 4$ and gR is a hard phase. The saddle-point
 416 equations (see Eqs. 5) are free from the pattern interference term $\sqrt{ar}x$ present in their

417 $p^* = p$ counterparts (see Eqs. 4) until β^* becomes so small that it approaches $\mathcal{O}(N^{2/p-1})$.
 418 Therefore, contrary to $p^* = p = 2$, the network is never in the **SG** phase. Practically, it means
 419 that $p \geq 3$ gives more freedom than $p = 2$ for tuning β and α . The only remaining restriction
 420 is that choosing α and T too small puts the network into the inaccurate **eR** phase resulting
 421 from the kz term (see Fig. 6). The saddle point equations can be derived without the RS
 422 ansatz because they do not involve q and r . Consequently, we expect them to yield an exact
 423 solution. Like on the Nishimori line, the total entropy of the paramagnetic phase is always
 424 positive, which is consistent with the solution being exact.

425 4.5 Robustness against adversarial attacks

426 Inverse models with $p^* = 2$ and $p \geq 3$ offer an opportunity to study adversarial attacks in a
 427 simple setting because their phase diagrams have regions where the signal retrieval phases (**gR**
 428 and **lR**) overlap with the inaccurate **eR** phase. Recall that, in the **lR** phase, a noisy student
 429 pattern ξ either converges to ξ^* or falls in the paramagnetic phase, depending on the amount
 430 of noise that ξ contains initially. The quantity of noise needed to prevent pattern retrieval
 431 becomes smaller as one approaches the **lR** to **P** phase transition and the basin of attraction
 432 of **lR** shrinks. Similarly, in the region of inaccurate **eR** where signal retrieval is metastable,
 433 patterns ξ that are corrupted by replacing some of their entries ξ_i by the components σ_i^a
 434 of an example σ^a may converge to σ^a when enough entries are replaced. The fraction ε of
 435 entries that need to be replaced becomes smaller as the basin of attraction of inaccurate **eR**
 436 expands and overtakes that of signal retrieval. In practice, an adversary can use this strategy
 437 to trick the student into converging to a pattern other than ξ^* . This scenario is similar to an
 438 adversarial attack targeting the input of Krotov's dense Hopfield network model because the
 439 student pattern ξ plays a similar role in the inverse model as the test data in Krotov's dense
 440 Hopfield networks (see Fig. 3, Section 4.2 and Appendix A). In that analogy, the examples σ
 441 are acting like the neural network weights rather than taking the role of the training data.

442 We will now investigate what values of the perturbation size ε are a threat by deriving a
 443 formula for the largest ε such that the student converges to the signal at zero temperature. This
 444 largest ε will be denoted ε^* , and we expect it to be a good measure of adversarial robustness.
 445 The saddle-point equations with $T = 0$ indicate that the student converges to one of the signal
 446 retrieval phases if and only if $k < \eta \alpha r^*$ (see Eqs. 5). Sampling the initial conditions of ξ_i
 447 from $(1 - \varepsilon) \delta(\xi_i - \xi_i^*) + \varepsilon \delta(\xi_i - \sigma_i^a)$ with $\varepsilon \in [0, 1]$, we get

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

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

448 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
 449 to $m^* = \frac{1}{N} \sum_i \xi_i^* \sigma_i^a \approx 0$ as $N \rightarrow \infty$. If we take σ^a to be a typical example, then r^* and k
 450 reduce to

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

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

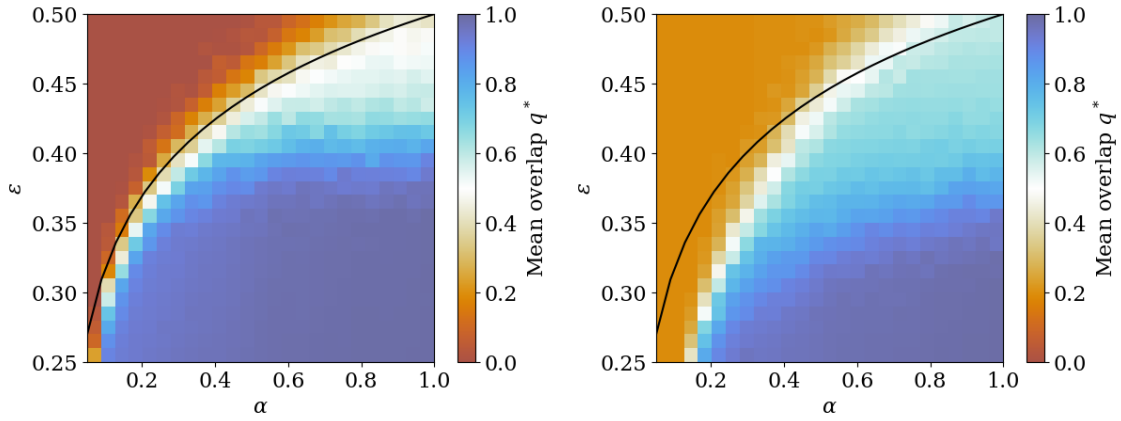


Figure 8: Monte-Carlo simulations of the overlap q^* as a function of α and ϵ in the inverse model with $p^* = 2$, $\beta^* = 1 - \frac{1}{\sqrt{2}}$, $p = 4$, $\beta = \infty$ and $N = 1024$. The simulation results are averaged over $L = 100$ student patterns. On the left plot, the inverse model is 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 .

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

$$\begin{aligned} \epsilon^{p-1} &< \eta \alpha (1 - \epsilon)^{p-1} \\ \epsilon &< \frac{[\eta \alpha]^{\frac{1}{p-1}}}{[\eta \alpha]^{\frac{1}{p-1}} + 1}. \end{aligned}$$

452 In other terms, the inverse model with $p^* = 2$ and even $p \geq 3$ is resistant to adversarial attacks
 453 of size $\epsilon^* = \frac{[\eta \alpha]^{\frac{1}{p-1}}}{[\eta \alpha]^{\frac{1}{p-1}} + 1}$ and smaller. For $p = 4$, ϵ^* is in good agreement with Monte-Carlo
 454 simulations of the inverse model corrupted by a typical example (see Fig. 8). This comparison
 455 is good evidence that our solution of the finite-noise scaling is indeed exact. Additionally, ϵ^* is
 456 a decent approximation of empirical robustness even when the inverse model is corrupted by
 457 the example that has the largest overlap with ξ^* . Just like adversarial attacks targeting more
 458 complicated neural networks [40, 41], our example-based attack can be hard to detect at low ϵ
 459 because a few adversarially perturbed entries ξ_i do not look very different from a low amount
 460 of meaningless noise. Interestingly, ϵ^* grows monotonically with α , which is consistent with the
 461 common observation that larger neural networks are also more adversarially robust [43, 60–65].
 462 At first glance, this effect can be counter-intuitive because adversarial vulnerability looks like a
 463 form of overfitting [42]. In our model, however, all examples work together to stabilize the \mathbf{IR}
 464 phase, and the best way to push the student into the \mathbf{eR} phase is to perturb it with a single
 465 example. Therefore, it is not surprising that increasing α makes the student more robust. We
 466 recall that the examples σ are a feature-based representation of ξ^* . Interestingly, it means
 467 that the underlying mechanism of our example-based attack is conceptually similar to gradient-
 468 based attacks targeting many common types of neural networks [42]. In fact, gradient-based
 469 attacks find features stored in neural network weights and add them to the data in order to
 470 fool the network [42, 66–68]. It would be interesting to investigate, both empirically and
 471 theoretically, if only a small number of weights are involved in constructing these adversarial

472 attacks. If it is the case, it could explain why larger neural networks are often more robust.
 473 In general, we expect this kind of one-example attack to be possible in any region of signal
 474 retrieval that overlaps with the inaccurate eR phase. Using $p \neq p^*$ may not be a necessary
 475 ingredient of adversarial vulnerability in more general models with other sources of mismatch,
 476 but in our case it ensures that the signal retrieval phases intersect the inaccurate eR phase.
 477 Conversely, the accurate eR phase is by definition robust to adversarial attacks since retrieving
 478 an example σ^a is the same as recovering ξ^* . This distinction clarifies why the dense Hopfield
 479 networks designed by Krotov are adversarially robust in the prototype phase despite being
 480 adversarially vulnerable in the feature phase. In fact, Krotov observed that adversarial attacks
 481 are unsuccessful in the prototype phase specifically because they retrieve stored examples that
 482 are semantically meaningful [37]. In summary, our model yields two main results concerning
 483 adversarial examples. First of all, it suggests a reason why large feature-based neural networks
 484 are more adversarially robust than smaller ones. Second of all, it clarifies why dense Hopfield
 485 networks are much more robust in the prototype phase than in the feature phase.

486 5 Conclusion

487 In this work, we derive the exact phase diagram of the p -dense networks in the teacher-
 488 student setting [16, 17, 30, 38]. On the Nishimori line, we find an example retrieval phase (eR)
 489 and a global retrieval phase (gR) reminiscent of the prototype and feature regimes observed
 490 empirically in dense Hopfield networks [26]. We show that the phase transition towards gR of
 491 the inverse model overlaps the paramagnetic to spin-glass ($P-SG$) transition of the direct model,
 492 which allows us to locate the $P-SG$ transition much more precisely than before [30, 33]. On
 493 the other hand, we discover that inverse models outside of the Nishimori line are able to resist
 494 an extensive amount of noise. In fact, a student with $p \geq 3$ is able to learn from a teacher with
 495 $p^* = 2$ even when the teacher's inverse temperature β^* is as low as $\mathcal{O}(N^{2/p-1})$. Moreover, such
 496 a student is immune to pattern interference until β^* reaches $\mathcal{O}(N^{2/p-1})$. In this setting, we
 497 derive a formula measuring the adversarial robustness of the student with $p \geq 3$ and $T = 0$. We
 498 then use this formula to describe how making a neural network larger can potentially increase
 499 its robustness to adversarial attacks constructed with only a few learned weights [43, 60–65].
 500 Our model also clarifies why the prototype phase of dense Hopfield networks is adversarially
 501 robust [37]. We compare our key results against Monte-Carlo simulations.

502 Dense networks with exponential interactions have been argued to be the $p \rightarrow \infty$ limit of
 503 the p -body models [69]. It would be interesting to see if they can achieve $\mathcal{O}(N)$ noise tolerance
 504 at the cost of an exponential number of training examples. More generally, studying exponential
 505 models in the teacher-student setting would be an interesting extension of this work and could
 506 be used to complement existing studies of the direct model [69, 70]. A caveat of our model is
 507 that the teacher has only one pattern. In fact, we would need to use a teacher with at least
 508 two patterns to describe more completely the kind of adversarial attack aiming to misclassify
 509 data. It should be possible to study this kind of teacher by using an approach similar to [71].
 510 On the practical side, we highlight the untapped benefits of using p -body models to either
 511 resist an extensive amount of noise in the feature phase or improve adversarial robustness in
 512 the prototype phase. Overall, we stress that further investigations of dense Hopfield networks
 513 could unlock their true potential.

514 **Funding information** This work was partially supported by project SERICS (PE00000014)
 515 under the MUR National Recovery and Resilience Plan funded by the European Union - NextGen-
 516 erationEU. The work was also supported by the project PRIN22TANTARI "Statistical Mechanics
 517 of Learning Machines: from algorithmic and information-theoretical limits to new biolog-

518 ically inspired paradigms" 20229T9EAT – CUP J53D23003640001. DT also acknowledges
519 GNFM-Indam.

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

522 References

- 523 [1] J. J. Hopfield, *Neural networks and physical systems with emergent collective computa-*
524 *tional abilities.*, Proceedings of the National Academy of Sciences **79**(8), 2554 (1982),
525 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>.
- 526 [2] D. J. Amit, H. Gutfreund and H. Sompolinsky, *Spin-glass models of neural networks*, Phys.
527 Rev. A **32**, 1007 (1985), doi:[10.1103/PhysRevA.32.1007](https://doi.org/10.1103/PhysRevA.32.1007).
- 528 [3] D. O. Hebb, *The organization of behavior: A neuropsychological theory*, Psychology press,
529 ISBN 9781410612403, doi:<https://doi.org/10.4324/9781410612403> (2005).
- 530 [4] T. M. Cover, *Geometrical and statistical properties of systems of linear inequalities with*
531 *applications in pattern recognition*, IEEE Transactions on Electronic Computers **EC-14**(3),
532 326 (1965), doi:[10.1109/PGEC.1965.264137](https://doi.org/10.1109/PGEC.1965.264137).
- 533 [5] D. J. Amit, H. Gutfreund and H. Sompolinsky, *Storing infinite numbers of pat-*
534 *terns in a spin-glass model of neural networks*, Phys. Rev. Lett. **55**, 1530 (1985),
535 doi:[10.1103/PhysRevLett.55.1530](https://doi.org/10.1103/PhysRevLett.55.1530).
- 536 [6] E. Agliari, A. Barra, A. Galluzzi, F. Guerra and F. Moauro, *Multitasking associative networks*,
537 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).
- 538 [7] E. Agliari, A. Annibale, A. Barra, A. C. C. Coolen and D. Tantari, *Immune networks: multi-*
539 *tasking capabilities at medium load*, Journal of Physics A: Mathematical and Theoretical
540 **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).
- 541 [8] E. Agliari, A. Annibale, A. Barra, A. C. C. Coolen and D. Tantari, *Immune networks:*
542 *multitasking capabilities near saturation*, Journal of Physics A Mathematical General
543 **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).
- 544 [9] P. Sollich, D. Tantari, A. Annibale and A. Barra, *Extensive parallel processing on scale-free*
545 *networks*, Phys. Rev. Lett. **113**, 238106 (2014), doi:[10.1103/PhysRevLett.113.238106](https://doi.org/10.1103/PhysRevLett.113.238106),
546 [1404.3654](https://arxiv.org/abs/1404.3654).
- 547 [10] E. Agliari, A. Annibale, A. Barra, A. C. C. Coolen and D. Tantari, *Retrieving infinite numbers*
548 *of patterns in a spin-glass model of immune networks*, Europhysics Letters **117**(2), 28003
549 (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).
- 550 [11] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Retrieval capabilities*
551 *of hierarchical networks: From dyson to hopfield*, Phys. Rev. Lett. **114**, 028103 (2015),
552 doi:[10.1103/PhysRevLett.114.028103](https://doi.org/10.1103/PhysRevLett.114.028103), [1407.5019](https://arxiv.org/abs/1407.5019).
- 553 [12] E. Agliari, D. Migliozi and D. Tantari, *Non-convex Multi-species Hopfield Models*, Journal of
554 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).

- 555 [13] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Hierarchical neural*
556 *networks perform both serial and parallel processing*, *Neural Networks* **66**, 22 (2015),
557 doi:<https://doi.org/10.1016/j.neunet.2015.02.010>, [1409.0227](https://arxiv.org/abs/1409.0227).
- 558 [14] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Metastable states in*
559 *the hierarchical Dyson model drive parallel processing in the hierarchical Hopfield network*,
560 *Journal of Physics A Mathematical General* **48**(1), 015001 (2015), doi:[10.1088/1751-](https://doi.org/10.1088/1751-8113/48/1/015001)
561 [8113/48/1/015001](https://doi.org/10.1088/1751-8113/48/1/015001), [1407.5176](https://arxiv.org/abs/1407.5176).
- 562 [15] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, D. Tantari and F. Tavani, *Topological properties of*
563 *hierarchical networks*, *Phys. Rev. E* **91**, 062807 (2015), doi:[10.1103/PhysRevE.91.062807](https://doi.org/10.1103/PhysRevE.91.062807),
564 [1412.5918](https://arxiv.org/abs/1412.5918).
- 565 [16] A. Barra, G. Genovese, P. Sollich and D. Tantari, *Phase transitions in re-*
566 *stricted boltzmann machines with generic priors*, *Phys. Rev. E* **96**, 042156 (2017),
567 doi:[10.1103/PhysRevE.96.042156](https://doi.org/10.1103/PhysRevE.96.042156), [1612.03132](https://arxiv.org/abs/1612.03132).
- 568 [17] A. Barra, G. Genovese, P. Sollich and D. Tantari, *Phase diagram of restricted boltzmann*
569 *machines and generalized hopfield networks with arbitrary priors*, *Phys. Rev. E* **97**, 022310
570 (2018), doi:[10.1103/PhysRevE.97.022310](https://doi.org/10.1103/PhysRevE.97.022310), [1702.05882](https://arxiv.org/abs/1702.05882).
- 571 [18] A. Barra, P. Contucci, E. Mingione and D. Tantari, *Multi-Species Mean Field Spin Glasses.*
572 *Rigorous Results*, *Annales Henri Poincaré*; **16**(3), 691 (2015), doi:[10.1007/s00023-](https://doi.org/10.1007/s00023-014-0341-5)
573 [014-0341-5](https://doi.org/10.1007/s00023-014-0341-5), [1307.5154](https://arxiv.org/abs/1307.5154).
- 574 [19] E. Agliari, A. Barra, C. Longo and D. Tantari, *Neural Networks Retrieving Boolean Pat-*
575 *terns in a Sea of Gaussian Ones*, *Journal of Statistical Physics* **168**(5), 1085 (2017),
576 doi:[10.1007/s10955-017-1840-9](https://doi.org/10.1007/s10955-017-1840-9), [1703.05210](https://arxiv.org/abs/1703.05210).
- 577 [20] A. Barra, G. Genovese, F. Guerra and D. Tantari, *How glassy are neural networks?*, *Journal of*
578 *Statistical Mechanics: Theory and Experiment* **2012**(7), 07009 (2012), doi:[10.1088/1742-](https://doi.org/10.1088/1742-5468/2012/07/P07009)
579 [5468/2012/07/P07009](https://doi.org/10.1088/1742-5468/2012/07/P07009), [1205.3900](https://arxiv.org/abs/1205.3900).
- 580 [21] G. Genovese and D. Tantari, *Legendre equivalences of spherical Boltzmann machines*,
581 *Journal of Physics A Mathematical General* **53**(9), 094001 (2020), doi:[10.1088/1751-](https://doi.org/10.1088/1751-8121/ab6b92)
582 [8121/ab6b92](https://doi.org/10.1088/1751-8121/ab6b92), [1910.14559](https://arxiv.org/abs/1910.14559).
- 583 [22] J. Rocchi, D. Saad and D. Tantari, *High storage capacity in the Hopfield model with*
584 *auto-interactions—stability analysis*, *Journal of Physics A Mathematical General* **50**(46),
585 465001 (2017), doi:[10.1088/1751-8121/aa8fd7](https://doi.org/10.1088/1751-8121/aa8fd7), [1704.07741](https://arxiv.org/abs/1704.07741).
- 586 [23] H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, L. Gruber, M. Holzleitner, T. Adler,
587 D. Kreil, M. K. Kopp, G. Klambauer, J. Brandstetter *et al.*, *Hopfield networks is all you need*,
588 *In International Conference on Learning Representations*, doi:[10.48550/arXiv.2008.02217](https://doi.org/10.48550/arXiv.2008.02217)
589 (2021), [2008.02217](https://arxiv.org/abs/2008.02217).
- 590 [24] M. Widrich, B. Schäfl, M. Pavlović, H. Ramsauer, L. Gruber, M. Holzleitner, J. Brandstetter,
591 G. K. Sandve, V. Greiff, S. Hochreiter and G. Klambauer, *Modern hopfield networks and*
592 *attention for immune repertoire classification*, In H. Larochelle, M. Ranzato, R. Hadsell,
593 M. Balcan and H. Lin, eds., *Advances in Neural Information Processing Systems*, vol. 33,
594 pp. 18832–18845. Curran Associates, Inc., doi:[10.48550/arXiv.2007.13505](https://doi.org/10.48550/arXiv.2007.13505) (2020),
595 [2007.13505](https://arxiv.org/abs/2007.13505).
- 596 [25] D. Krotov and J. J. Hopfield, *Large associative memory problem in neurobiol-*
597 *ogy and machine learning*, *In International Conference on Learning Representations*,
598 doi:[10.48550/arXiv.2008.06996](https://doi.org/10.48550/arXiv.2008.06996) (2021), [2008.06996](https://arxiv.org/abs/2008.06996).

- 599 [26] D. Krotov and J. J. Hopfield, *Dense associative memory for pattern recognition*, In D. Lee,
600 M. Sugiyama, U. Luxburg, I. Guyon and R. Garnett, eds., *Advances in Neural Information*
601 *Processing Systems*, vol. 29. Curran Associates, Inc., doi:[10.48550/arXiv.1606.01164](https://doi.org/10.48550/arXiv.1606.01164)
602 (2016), [1606.01164](https://doi.org/10.48550/arXiv.1606.01164).
- 603 [27] H. H. Chen, Y. C. Lee, G. Z. Sun, H. Y. Lee, T. Maxwell and C. L. Giles, *High order*
604 *correlation model for associative memory*, *AIP Conference Proceedings* **151**(1), 86 (1986),
605 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)
606 [86_1_online.pdf](https://pubs.aip.org/aip/acp/article-pdf/151/1/86/12091820/86_1_online.pdf).
- 607 [28] D. Psaltis and C. H. Park, *Nonlinear discriminant functions and associative memories*, *AIP*
608 *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)
609 [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).
- 610 [29] P. Baldi and S. S. Venkatesh, *Number of stable points for spin-glasses and neural networks*
611 *of higher orders*, *Phys. Rev. Lett.* **58**, 913 (1987), doi:[10.1103/PhysRevLett.58.913](https://doi.org/10.1103/PhysRevLett.58.913).
- 612 [30] E. Gardner, *Multiconnected neural network models*, *Journal of Physics A: Mathematical*
613 *and General* **20**(11), 3453 (1987), doi:[10.1088/0305-4470/20/11/046](https://doi.org/10.1088/0305-4470/20/11/046).
- 614 [31] L. F. Abbott and Y. Arian, *Storage capacity of generalized networks*, *Phys. Rev. A* **36**, 5091
615 (1987), doi:[10.1103/PhysRevA.36.5091](https://doi.org/10.1103/PhysRevA.36.5091).
- 616 [32] Horn, D. and Usher, M., *Capacities of multiconnected memory models*, *J. Phys. France*
617 **49**(3), 389 (1988), doi:[10.1051/jphys:01988004903038900](https://doi.org/10.1051/jphys:01988004903038900).
- 618 [33] L. Albanese, F. Alemanno, A. Alessandrelli and A. Barra, *Replica Symmetry Breaking*
619 *in Dense Hebbian Neural Networks*, *Journal of Statistical Physics* **189**(2), 24 (2022),
620 doi:[10.1007/s10955-022-02966-8](https://doi.org/10.1007/s10955-022-02966-8), [2111.12997](https://doi.org/10.1007/s10955-022-02966-8).
- 621 [34] B. Hoover, Y. Liang, B. Pham, R. Panda, H. Strobelt, D. H. Chau, M. J. Zaki
622 and D. Krotov, *Energy Transformer*, arXiv e-prints arXiv:2302.07253 (2023),
623 doi:[10.48550/arXiv.2302.07253](https://doi.org/10.48550/arXiv.2302.07253), [2302.07253](https://doi.org/10.48550/arXiv.2302.07253).
- 624 [35] B. Hoover, H. Strobelt, D. Krotov, J. Hoffman, Z. Kira and D. H. Chau, *Memory in Plain*
625 *Sight: A Survey of the Uncanny Resemblances between Diffusion Models and Associative*
626 *Memories*, arXiv e-prints arXiv:2309.16750 (2023), doi:[10.48550/arXiv.2309.16750](https://doi.org/10.48550/arXiv.2309.16750),
627 [2309.16750](https://doi.org/10.48550/arXiv.2309.16750).
- 628 [36] L. Ambrogioni, *In search of dispersed memories: Generative diffusion mod-*
629 *els are associative memory networks*, arXiv e-prints arXiv:2309.17290 (2023),
630 doi:[10.48550/arXiv.2309.17290](https://doi.org/10.48550/arXiv.2309.17290), [2309.17290](https://doi.org/10.48550/arXiv.2309.17290).
- 631 [37] D. Krotov and J. Hopfield, *Dense Associative Memory Is Robust to Adversarial Inputs*, *Neural*
632 *Computation* **30**(12), 3151 (2018), doi:[10.1162/neco_a_01143](https://doi.org/10.1162/neco_a_01143), [1701.00939](https://doi.org/10.1162/neco_a_01143).
- 633 [38] F. Alemanno, L. Camanzi, G. Manzan and D. Tantari, *Hopfield model with planted patterns:*
634 *A teacher-student self-supervised learning model*, *Applied Mathematics and Computation*
635 **458**, 128253 (2023), doi:<https://doi.org/10.1016/j.amc.2023.128253>, [2304.13710](https://doi.org/10.1016/j.amc.2023.128253).
- 636 [39] A. Decelle, S. Hwang, J. Rocchi and D. Tantari, *Inverse problems for structured datasets*
637 *using parallel TAP equations and restricted Boltzmann machines*, *Scientific Reports* **11**,
638 19990 (2021), doi:[10.1038/s41598-021-99353-2](https://doi.org/10.1038/s41598-021-99353-2), [1906.11988](https://doi.org/10.1038/s41598-021-99353-2).

- 639 [40] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrndić, P. Laskov, G. Giacinto and F. Roli,
640 *Evasion attacks against machine learning at test time*, In H. Blockeel, K. Kersting, S. Ni-
641 jssen and F. Železný, eds., *Machine Learning and Knowledge Discovery in Databases*, pp.
642 387–402. Springer Berlin Heidelberg, Berlin, Heidelberg, ISBN 978-3-642-40994-3,
643 doi:https://doi.org/10.1007/978-3-642-40994-3_25 (2013).
- 644 [41] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow and R. Fer-
645 gus, *Intriguing properties of neural networks*, arXiv e-prints arXiv:1312.6199 (2013),
646 doi:[10.48550/arXiv.1312.6199](https://doi.org/10.48550/arXiv.1312.6199), [1312.6199](https://arxiv.org/abs/1312.6199).
- 647 [42] I. J. Goodfellow, J. Shlens and C. Szegedy, *Explaining and harnessing adversarial examples*,
648 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).
- 649 [43] A. Madry, A. Makelov, L. Schmidt, D. Tsipras and A. Vladu, *Towards deep learning models*
650 *resistant to adversarial attacks*, In *International Conference on Learning Representations*,
651 doi:[10.48550/arXiv.1706.06083](https://doi.org/10.48550/arXiv.1706.06083) (2018), [1706.06083](https://arxiv.org/abs/1706.06083).
- 652 [44] A. Muhammad and S.-H. Bae, *A Survey on Efficient Methods for Adversarial Robustness*,
653 IEEE Access **10**, 118815 (2022), doi:[10.1109/ACCESS.2022.3216291](https://doi.org/10.1109/ACCESS.2022.3216291).
- 654 [45] H. Nishimori, *Exact results and critical properties of the ising model with competing interac-*
655 *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)
656 [3719/13/21/012](https://doi.org/10.1088/0022-3719/13/21/012).
- 657 [46] H. Nishimori, *Statistical Physics of Spin Glasses and Information Processing: An Introduction*,
658 Oxford University Press, ISBN 9780198509417,
659 doi:[10.1093/acprof:oso/9780198509417.001.0001](https://doi.org/10.1093/acprof:oso/9780198509417.001.0001) (2001), [https://academic.oup.com/
660 book/5185/book-pdf/54038185/acprof-9780198509400.pdf](https://academic.oup.com/book/5185/book-pdf/54038185/acprof-9780198509400.pdf).
- 661 [47] P. Contucci, C. Giardinà and H. Nishimori, *Spin glass identities and the nishimori line*, In
662 A. B. de Monvel and A. Bovier, eds., *Spin Glasses: Statics and Dynamics*, pp. 103–121.
663 Birkhäuser Basel, Basel, doi:https://doi.org/10.1007/978-3-7643-9891-0_4 (2009),
664 [0805.0754](https://doi.org/10.1007/978-3-7643-9891-0_4).
- 665 [48] Y. Iba, *The Nishimori line and Bayesian statistics*, Journal of Physics A Mathematical
666 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).
- 667 [49] P. Charbonneau, *From the replica trick to the replica symmetry breaking technique*, arXiv
668 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).
- 669 [50] D. Sherrington and S. Kirkpatrick, *Solvable model of a spin-glass*, Phys. Rev. Lett. **35**, 1792
670 (1975), doi:[10.1103/PhysRevLett.35.1792](https://doi.org/10.1103/PhysRevLett.35.1792).
- 671 [51] G. G. Roussas, *Contiguity of Probability Measures: Some Applications in*
672 *Statistics*, Cambridge Tracts in Mathematics. Cambridge University Press,
673 doi:[10.1017/CBO9780511804373](https://doi.org/10.1017/CBO9780511804373) (1972).
- 674 [52] D. Achlioptas and A. Coja-Oghlan, *Algorithmic barriers from phase transitions*, In
675 *2008 49th Annual IEEE Symposium on Foundations of Computer Science*, pp. 793–802,
676 doi:[10.1109/FOCS.2008.11](https://doi.org/10.1109/FOCS.2008.11) (2008), [0803.2122](https://arxiv.org/abs/0803.2122).
- 677 [53] F. Krzakala and L. Zdeborová, *Hiding quiet solutions in random constraint satisfaction*
678 *problems*, Phys. Rev. Lett. **102**, 238701 (2009), doi:[10.1103/PhysRevLett.102.238701](https://doi.org/10.1103/PhysRevLett.102.238701),
679 [0901.2130](https://arxiv.org/abs/0901.2130).

- 680 [54] L. Zdeborová and F. Krzakala, *Quiet planting in the locked constraint satisfaction problems*,
681 SIAM Journal on Discrete Mathematics **25**(2), 750 (2011), doi:[10.1137/090750755](https://doi.org/10.1137/090750755),
682 [0902.4185](https://doi.org/10.1137/090750755).
- 683 [55] L. Zdeborová and F. Krzakala, *Statistical physics of inference: thresholds and algorithms*,
684 Advances in Physics **65**(5), 453 (2016), doi:[10.1080/00018732.2016.1211393](https://doi.org/10.1080/00018732.2016.1211393), [1511.](https://doi.org/10.1080/00018732.2016.1211393)
685 [02476](https://doi.org/10.1080/00018732.2016.1211393).
- 686 [56] F. Antenucci, S. Franz, P. Urbani and L. Zdeborová, *Glassy nature of the hard phase in*
687 *inference problems*, Phys. Rev. X **9**, 011020 (2019), doi:[10.1103/PhysRevX.9.011020](https://doi.org/10.1103/PhysRevX.9.011020),
688 [1805.05857](https://doi.org/10.1103/PhysRevX.9.011020).
- 689 [57] L. Zdeborová and F. Krzakala, *Phase transitions in the coloring of random graphs*, Phys.
690 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).
- 691 [58] E. Agliari, F. Alemanno, A. Barra, M. Centonze and A. Fachechi, *Neural networks with a*
692 *redundant representation: Detecting the undetectable*, Phys. Rev. Lett. **124**, 028301 (2020),
693 doi:[10.1103/PhysRevLett.124.028301](https://doi.org/10.1103/PhysRevLett.124.028301), [1911.12689](https://doi.org/10.1103/PhysRevLett.124.028301).
- 694 [59] E. Agliari and G. De Marzo, *Tolerance versus synaptic noise in dense associative memories*,
695 European Physical Journal Plus **135**(11), 883 (2020), doi:[10.1140/epjp/s13360-020-](https://doi.org/10.1140/epjp/s13360-020-00894-8)
696 [00894-8](https://doi.org/10.1140/epjp/s13360-020-00894-8), [2007.02849](https://doi.org/10.1140/epjp/s13360-020-00894-8).
- 697 [60] S. Goyal, C. Qin, J. Uesato, T. Mann and P. Kohli, *Uncovering the Limits of Adversarial*
698 *Training against Norm-Bounded Adversarial Examples*, arXiv e-prints arXiv:2010.03593
699 (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).
- 700 [61] H. Huang, Y. Wang, S. Erfani, Q. Gu, J. Bailey and X. Ma, *Exploring architectural ingredients*
701 *of adversarially robust deep neural networks*, In M. Ranzato, A. Beygelzimer, Y. Dauphin,
702 P. Liang and J. W. Vaughan, eds., *Advances in Neural Information Processing Systems*,
703 vol. 34, pp. 5545–5559. Curran Associates, Inc., doi:[10.48550/arXiv.2110.03825](https://doi.org/10.48550/arXiv.2110.03825) (2021),
704 [2110.03825](https://doi.org/10.48550/arXiv.2110.03825).
- 705 [62] S. Bubeck, Y. Li and D. M. Nagaraj, *A law of robustness for two-layers neural networks*,
706 In M. Belkin and S. Kpotufe, eds., *Proceedings of Thirty Fourth Conference on Learn-*
707 *ing Theory*, vol. 134 of *Proceedings of Machine Learning Research*, pp. 804–820. PMLR,
708 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).
- 709 [63] S. Bubeck and M. Sellke, *A universal law of robustness via isoperimetry*, In M. Ran-
710 zato, A. Beygelzimer, Y. Dauphin, P. Liang and J. W. Vaughan, eds., *Advances in Neu-*
711 *ral Information Processing Systems*, vol. 34, pp. 28811–28822. Curran Associates, Inc.,
712 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).
- 713 [64] J. Puigcerver, R. Jenatton, C. Riquelme, P. Awasthi and S. Bhojanapalli, *On the adversarial*
714 *robustness of mixture of experts*, In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho
715 and A. Oh, eds., *Advances in Neural Information Processing Systems*, vol. 35, pp. 9660–9671.
716 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).
- 717 [65] A. H. Ribeiro and T. B. Schön, *Overparameterized Linear Regression Under*
718 *Adversarial Attacks*, IEEE Transactions on Signal Processing **71**, 601 (2023),
719 doi:[10.1109/TSP.2023.3246228](https://doi.org/10.1109/TSP.2023.3246228), [2204.06274](https://doi.org/10.1109/TSP.2023.3246228).
- 720 [66] S. Jetley, N. Lord and P. Torr, *With friends like these, who needs adversaries?*, In
721 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett,
722 eds., *Advances in Neural Information Processing Systems*, vol. 31. Curran Associates, Inc.,
723 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).

- 724 [67] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran and A. Madry, *Adversarial examples*
725 *are not bugs, they are features*, In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-
726 Buc, E. Fox and R. Garnett, eds., *Advances in Neural Information Processing Systems*,
727 vol. 32. Curran Associates, Inc., doi:<https://doi.org/10.48550/arXiv.1905.02175> (2019),
728 [1905.02175](https://doi.org/10.48550/arXiv.1905.02175).
- 729 [68] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner and A. Madry, *Robustness may*
730 *be at odds with accuracy*, In *International Conference on Learning Representations*,
731 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).
- 732 [69] M. Demircigil, J. Heusel, M. Löwe, S. Upgang and F. Vermet, *On a Model of Associative*
733 *Memory with Huge Storage Capacity*, *Journal of Statistical Physics* **168**(2), 288 (2017),
734 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).
- 735 [70] C. Lucibello and M. Mézard, *The Exponential Capacity of Dense Associative Memories*, arXiv
736 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).
- 737 [71] T. Hou, K. Y. M. Wong and H. Huang, *Minimal model of permutation symmetry in unsu-*
738 *perervised learning*, *Journal of Physics A: Mathematical and Theoretical* **52**(41), 414001
739 (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).
- 740 [72] N. Eddine Boukacem, A. Leary, R. Thériault, F. Gottlieb, M. Mani and P. François, *A*
741 *Waddington landscape for prototype learning in generalized Hopfield networks*, arXiv
742 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).

743 A Gardner's Hamiltonian vs Krotov's Hamiltonian

744 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
745 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
746 the number of patterns ξ^{μ} used to construct J , and N denotes the number of components of
747 each pattern ξ^{μ} and example σ . In this Section, we will omit ξ in the argument of $H[\sigma|\xi]$
748 and write $H[\sigma]$ instead for notational simplicity. Unless indicated otherwise, we will assume
749 a large number number of components $N \gg 1$ and patterns $M \sim \mathcal{O}(N^{p-1})$. We will start
750 by comparing it to the dense Hopfield network Hamiltonian $\mathcal{H}[\sigma] = -\frac{1}{N^{p-1}} \sum_{\mu} \left(\sum_i \xi_i^{\mu} \sigma_i \right)^p$
751 studied by Krotov [26].

752 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-
753 ming over all permutations of $\{i_1 \dots i_p\}$ in place of the restricted set $i_1 < \dots < i_p$ and compen-
754 sating for double counting with the prefactor $\frac{1}{p!}$. This manipulation leads to

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

755 On the other hand, Krotov's Hamiltonian may be rewritten

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

756 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
757 configurations where some indices are equal. For example, the configuration $i_1 \neq \dots \neq i_{p-1} = i_p$
758 contains the fewest equal indices after $i_1 \neq \dots \neq i_p$. In other words, $\mathcal{H}[\sigma]$ can be expressed as
759 an expansion around $H[\sigma]$, and the two Hamiltonians are equivalent when the normalized
760 residuals $\frac{\mathcal{H}[\sigma] - H[\sigma]}{N}$ vanish in the limit of large N . In this study, we encounter two cases which
761 bring different results.

- 762 **1** The Hamiltonians $\mathcal{H}[\sigma]$ and $H[\sigma]$ are dominated by a few closely packed configurations
763 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
764 σ .
- 765 **2** The Hamiltonians $\mathcal{H}[\sigma]$ and $H[\sigma]$ are dominated by many spread out configurations
766 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
767 misaligned with σ

768 We use the expansion of $\mathcal{H}[\sigma]$ to discuss both the aligned case and the misaligned case. We
769 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
770 form

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

771 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
772 can be summarized by $\mathcal{H}[\sigma] = H[\sigma] + H'[\sigma] + \dots$, where $H'[\sigma]$ simplifies to

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

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

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

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

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

794 For this reason and also because the number $M \sim \mathcal{O}(N^{p-1})$ of patterns ξ^{μ} used in a Monte-
 795 Carlo simulations increases exponentially with p , we choose to simulate only $p = 3$ and
 796 $p = 4$.

797 The output of the neural network model that Krotov designed for classification of data
 798 is $\mathbf{c}_j = \tanh \left[\frac{1}{2} \beta (\mathcal{H}[\boldsymbol{\sigma}'] - \mathcal{H}[\boldsymbol{\sigma}]) \right] \approx \tanh \left[\beta p \sum_{\mu} \xi_j^{\mu} \left(\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i \right)^{p-1} \right]$. We omit the linear
 799 rectifier present in the original paper [26] because the overlaps $\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i$ are almost always
 800 positive (see for example the Supplement of [72]). The predicted class is then $\mathbf{j}' = \mathbf{argmax}_j \{\mathbf{c}_j\}$.
 801 Using $\mathbf{1} - P_j = \frac{1}{1+\exp[\beta(\mathcal{H}[\boldsymbol{\sigma}']-\mathcal{H}[\boldsymbol{\sigma}])]}$ instead of \mathbf{c}_j does not change \mathbf{j}' because $\mathbf{1} - P_j$ and \mathbf{c}_j are
 802 related by $\mathbf{1} - P_j = \frac{1}{2} [\mathbf{c}_j + 1]$. When we evaluate P_i using \mathbf{H} instead of \mathcal{H} , this relation does
 803 not always hold exactly. Rather, it should be considered an approximation.

804 B Direct model cumulant expansions

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

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

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

$$\begin{aligned} \langle Z^L \rangle \approx & \left\langle \sum_{\boldsymbol{\sigma}} \exp \left(\beta N \sum_{\gamma} \sum_{\mu \in \mathbb{I}_{\gamma}} \left[\frac{1}{N} \sum_i \xi_i^{\mu} \sigma_i^{\gamma} \right]^p \right. \right. \\ & \left. \left. + \beta \sum_{\gamma} \sum_{\mu \in \mathbb{I}_{\gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^{\mu} \dots \xi_{i_p}^{\mu} \sigma_{i_1}^{\gamma} \dots \sigma_{i_p}^{\gamma} \right) \right\rangle, \end{aligned} \quad (9)$$

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

- 810 • $\sum_{\mu \in \Gamma_\gamma} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \approx N \sum_{\mu \in \Gamma_\gamma} \left[\frac{1}{N} \sum_i \xi_i^\mu \sigma_i^\gamma \right]^p$ because this part of
 811 $H[\sigma^\gamma | \xi]$ is aligned with σ (see Case 1 of Appendix A).
- 812 • $\sum_{\mu \in \bar{\Gamma}_\gamma} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma \approx \sum_{\mu \in \bar{\Gamma}} \frac{p!}{N^{p-1}} \sum_{i_1 < \dots < i_p} \xi_{i_1}^\mu \dots \xi_{i_p}^\mu \sigma_{i_1}^\gamma \dots \sigma_{i_p}^\gamma$ since $\bar{\Gamma}$
 813 contains almost all of the elements in each $\bar{\Gamma}_\gamma$ when N is large.

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

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

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

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

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

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

820 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 $i_n \neq j_n$
 821 for some n , then $\left\langle \xi_{i_1}^\mu \xi_{j_1}^\mu \dots \xi_{i_p}^\mu \xi_{j_p}^\mu \right\rangle = \mathbf{0}$ because $\xi_{i_1}^\mu \xi_{j_1}^\mu \dots \xi_{i_p}^\mu \xi_{j_p}^\mu$ is still a product of independent
 822 random spins once the deterministic variables are removed. These two cases can be summarized

823 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

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

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

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

831 which contributes to the free energy on the same order in N as Gardner's Hamiltonian. Therefore,
832 Krotov's Hamiltonian is not equivalent to Gardner's Hamiltonian when the latter is misaligned
833 with σ (see Case 2). The index configurations with more equality constraints also contribute
834 to the free entropy on the same order in N because the factors of N that are lost to equality
835 constraints are restored when the sums get squared in the cumulant expansion.

836 $p = 2$ is the only positive integer such that Gardner's Hamiltonian and Krotov's Hamiltonian
837 are equivalent [5, 30]. In the misaligned case with a single stored pattern ξ^* (see Case 2), the
838 free entropy of $p = 2$ simplifies to

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

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

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

840 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
841 with the one found in literature when $\alpha = \frac{1}{N}$ [5].

842 C Teacher-student replicated partition function

843 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
844 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 $P(\sigma^a|\xi)$
845 is the distribution of the direct model with a single pattern ξ . To simplify $P(\xi|\sigma)$ further,
846 we need to manipulate the partition function $Z = \sum_{\sigma^a} \exp(-\beta H[\sigma^a|\xi])$ of $P(\sigma^a|\xi)$ (see
847 Appendix A for the definition of $H[\sigma|\xi]$). Under the gauge transformation $\sigma_i^a \rightarrow \xi_i \sigma_i^a$, we
848 may write

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

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

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

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

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

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

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

856 where $P(\sigma|\xi^*)$ has the same functional form as $P(\sigma|\xi^b)$, but has hyperparameters \mathbf{p}^* and β^*
 857 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
 858 the sum, which yields

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

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

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

862 We simplify this expression to:

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

863 where Γ_b represents the set of inputs σ^a which have macroscopic overlap with the pattern ξ^b ,
864 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
865 reasoning used to build the sets Γ_* , Γ_b and $\bar{\Gamma}$ is the same as outlined at the start of appendix B.

866 D Teacher-student free entropy

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

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

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

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

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

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

873 which is the paramagnetic free entropy of a p^* -body Hopfield network [5, 30]. Coming back to
874 $\langle Z^L \rangle$, we fix order parameters q^{*b} , q^{bc} and m_a^b using the delta functions $\delta\left(Nq^{*b} - \sum_i \xi_i^* \xi_i^b\right)$,
875 $\delta\left(Nq^{bc} - \sum_i \xi_i^b \xi_i^c\right)$ and $\delta\left(Nm_a^b - \sum_i \xi_i^b \sigma_i^a\right)$, which results in

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

876 In Fourier space, this expression takes the form

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

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

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

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

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

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

881 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
882 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
883 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
884 $\langle \sigma_{i_n}^a \rangle = 0$ make the second order expectation vanish.

885 We need to go beyond second order to treat $p^* \neq p$. We will focus on $p^* = 2$ and $p \geq 3$
886 to investigate the consequences of using a p -body model to learn examples generated by the
887 original 2-body Hopfield model. For simplicity, we take p even so that the spins of both terms
888 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
889 raised to the power of $\frac{p}{2}$ and the student term is raised to the power of 1. This restriction will
890 simplify some of the incoming calculations. To leading order in N , the cumulant generating

891 function reduces to

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

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

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

894 and the student term

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

895 are both known from Appendix B. Later on, we will use $\log(\mathbf{z}^*)$ and \mathbf{z}^* as shorthands for
896 $-\frac{1}{2} \log(1 - 2\beta^*) - \beta^*$ and $\exp\left(-\frac{1}{2} \log(1 - 2\beta^*) - \beta^*\right)$, respectively. The coupling between
897 the teacher and the student can be rewritten as

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

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

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

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

900 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
901 depend on the microscopic details of the system. In fact, it can be expressed as a combination of
902 the moments of \mathbf{z}^* , which can all be derived from $\log(\mathbf{z}^*)$. To leading order in N , the cumulant
903 generating function expands to

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

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

906 **1** $M \sim \mathcal{O}(N^{p-1})$ and $\beta^* \sim \mathcal{O}(N^{2/p-1})$ will be called the large-noise scaling.

907 **2** $M \sim \mathcal{O}(N^{p/2})$ and $\beta^* \sim \mathcal{O}(1)$ will be called the finite-noise scaling.

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

$$\begin{aligned}
V(\beta^*, p) & \approx \frac{1}{(p/2)!} \left\langle \frac{p!}{N^p} \sum_{j_1 < \dots < j_p} \xi_{j_1}^* \dots \xi_{j_p}^* \sigma_{j_1}^a \dots \sigma_{j_p}^a \left(\beta^* \frac{2}{N} \sum_{i_1 < i_2} \xi_{i_1}^* \xi_{i_2}^* \sigma_{i_1}^a \sigma_{i_2}^a \right)^{p/2} \right\rangle \quad (10) \\
& = \frac{[\beta^*]^{p/2} 2^{p/2} p!}{(p/2)! N^{p/2} 2^{p/2}} \\
& = \frac{[\beta^*]^{p/2} p!}{(p/2)! N^{p/2}}
\end{aligned}$$

912 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
913 inequality constraints. In the large-noise scaling, we set

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

914 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-
 915 noise scaling, this expansion is only an order of magnitude approximation. However, it still
 916 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
 917 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
 918 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
 919 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:

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

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

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

921 The cumulants evaluate to

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

922 so we obtain

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

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

925 **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,

926 **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.

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

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

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

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

931 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
 932 $\frac{\log\langle Z^L \rangle}{N}$ to

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

933 Assuming each ξ^b has macroscopic overlap with at most one pattern σ^a and using the replica-
 934 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
 935 approximates to

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

936 Furthermore, the Hubbard-Stratonovich transformation gives

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

937 We can then use the factorization

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

938 in order to simplify the free energy to

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

939 After differentiating and taking the limit, we get

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

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

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

943 E Direct model RSB ansatz

944 Recall that the average replicated partition function of the direct model (see Eq. 9) takes the
945 form

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

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

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

947 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-
948 netic phase, we have (see [30] and also Appendix B)

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

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

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

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

951 F Monte-Carlo simulations for various system sizes

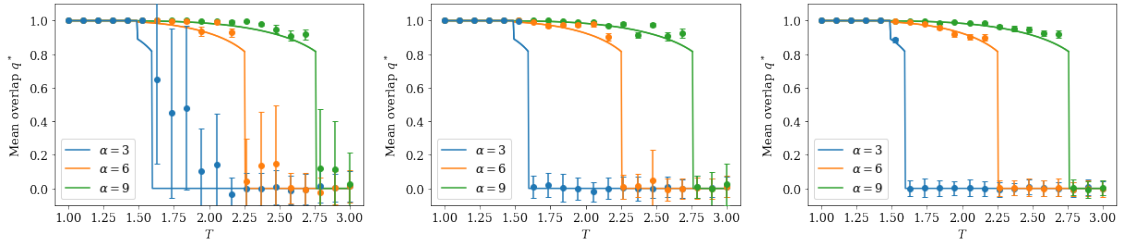


Figure 9: Monte-Carlo simulations of the $p = 3$ inverse model compared against saddle-point solutions for different values of N . The lR phase is not included in these plots. The left plot has $N = 128$, the center plot has $N = 256$, and the right plot has $N = 512$. The dots are simulation data at 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 .