

A Phase Transition between Positional and Semantic Learning in a Solvable Model of Dot-Product Attention

Hugo Cui

Statistical Physics Of Computation laboratory, EPFL, Switzerland

HUGO.CUI@EPFL.CH

Freya Behrens

Statistical Physics Of Computation laboratory, EPFL, Switzerland

FREYA.BEHRENS@EPFL.CH

Florent Krzakala

Information Learning & Physics laboratory, EPFL, Switzerland

FLORENT.KRZAKALA@EPFL.CH

Lenka Zdeborová

Statistical Physics Of Computation laboratory, EPFL, Switzerland

LENKA.ZDEBOROVA@EPFL.CH

Abstract

A theoretical understanding of how algorithmic abilities emerge in the learning of language models remains elusive. In this work, we provide a tight theoretical analysis of the emergence of semantic attention in a solvable model of dot-product attention and consider a non-linear self-attention layer with trainable tied and low-rank query and key matrices. In the asymptotic limit of high-dimensional data and a comparably large number of training samples we provide a tight closed-form characterization of the global minimum of the non-convex empirical loss landscape. We show that this minimum corresponds to either a positional attention mechanism (with tokens attending to each other based on their respective positions) or a semantic attention mechanism (with tokens attending to each other based on their meaning), and evidence an emergent phase transition from the former to the latter with increasing sample complexity.

1. Introduction

Self-attention layers [60] have been instrumental in advancing the abilities of language models, as they provide an efficient method of extracting information from sentences – both the information encoded in the ordering (i.e. *positions*) of the words, and that encoded in the meaning (i.e. *semantics*) of the words. In theory, attention layers can learn to leverage both types of information, by having tokens attend to each other based on their respective positions (a mechanism called *positional attention* in [27], through some form of positional encodings [24, 28, 46, 50, 52]) and/or respective meanings (henceforth referred to as *semantic attention*). A growing body of work on mechanistic interpretability aims to empirically understand which precise algorithmic mechanisms a neural network learns and has shown that attention layers are able to implement a wide range of different algorithms using both positional and semantic attributes of the inputs [1, 11, 38, 43, 44, 51, 65]. Simultaneously, empirical studies have provided evidence for the emergence of specific algorithmic mechanisms (abilities) in the learning of language models that lead to qualitative improvements of the model capabilities [51, 61, 62]. Despite these efforts, it remains an open question how to theoretically characterise the conditions under which such an ability emerges in the model. Even the nature of this algorithmic emergence is unclear, i.e. whether it constitutes a fast but smooth change in performance or it is due to a sharp boundary between fundamentally different regimes of learning [48]. While there is

a plethora of work on the theory of attention investigating various aspects such as their expressivity [16, 18, 23], inductive bias [47, 56, 57], training dynamics [9, 27, 29, 59], and in-context learning [3, 8, 21, 30, 64], these studies do not allow to capture sharp changes in the behaviour of attention mechanisms such as phase transitions or do not capture an emergent phenomenon [3, 18, 27, 45].

In our work, we take inspiration from physics, where a similar theoretical questions about the nature of phase transitions and emergent phenomena were posed a century ago for models of interacting particles, such as the famous Ising model [25, 39] and more recently in the theory of feed-forward fully connected neural networks, e.g. [2, 4, 5, 22, 33, 49, 53]. Here, for the first time, we bring this type of study to the analysis of neural networks with attention layers. We introduce and analyse a tractable model that permits a sharp high-dimensional characterisation for attention layers. In particular, we describe a model with a single self-attention layer with tied, low-rank query and key matrices, with Gaussian input data and realizable labels:

- (1) We show that this model exhibits a phase transition in terms of sample complexity between a semantic and a positional mechanism.
- (2) We analyse this model in the asymptotic limit where the embedding dimension d of the tokens and the number n of training samples grow proportionally to infinity and provide a *tight closed-form characterization* of the test error and training loss achieved at the minima of the non-convex empirical loss.
- (3) Using this high-dimensional characterization, we locate the positional-semantic phase transition, thus providing the first theoretical result about the *emergence of sharp phase transitions* in a model of dot-product attention.

2. Tied low-rank attention model

Input data model We consider a model of embedded sentences with uncorrelated (1-gram) words. More precisely, a sentence $\mathbf{x} \in \mathbb{R}^{L \times d}$, where L is the sentence length and d represents the embedding dimension, consists of L tokens $\{\mathbf{x}_\ell\}_{1 \leq \ell \leq L}$ independently drawn from a Gaussian distribution $\mathbf{x}_\ell \sim \mathcal{N}(0, \Sigma_\ell)$ with covariance $\Sigma_\ell \in \mathbb{R}^{d \times d}$. In the following, we denote the probability distribution of \mathbf{x} as p_x . Note that while this sentence model does not involve in itself statistical correlations between tokens, the task (target function) will entail interactions between different tokens.

Target function The target function (teacher) is assumed to be of the form

$$y(\mathbf{x}) = \mathbb{T} \left[\frac{1}{\sqrt{d}} \mathbf{x} \mathbf{Q}_* \right] \mathbf{x}, \quad (1)$$

for a function $\mathbb{T} : \mathbb{R}^{L \times r_t} \rightarrow \mathbb{R}^{L \times L}$. The term $\mathbb{T} \left[\frac{1}{\sqrt{d}} \mathbf{x} \mathbf{Q}_* \right] \in \mathbb{R}^{L \times L}$ in (1) should be interpreted as the target attention matrix, which mixes the tokens of the input \mathbf{x} . This attention matrix is parametrized by the target weights $\mathbf{Q}_* \in \mathbb{R}^{d \times r_t}$.

Tied attention We consider the learning of the target (1) by a single attention layer

$$f_{\mathbf{Q}}(\mathbf{x}) = \mathbb{S} \left[\frac{1}{\sqrt{d}} (\mathbf{x} + \mathbf{p}) \mathbf{Q} \right] (\mathbf{x} + \mathbf{p}). \quad (2)$$

In (2), $\mathbf{p} \in \mathbb{R}^{L \times d}$ is a *fixed* matrix, corresponding to positional encodings, and $\mathbf{Q} \in \mathbb{R}^{d \times r_s}$ is a trainable weight matrix. We denote subsequently $\mathbf{p}_\ell \in \mathbb{R}^d$ the ℓ -th row of \mathbf{p} . Like the target (1), the parametric function (2) takes the form of a data-dependent attention matrix $\mathbb{S} \left[\frac{1}{\sqrt{d}} (\mathbf{x} + \mathbf{p}) \mathbf{Q} \right] \in \mathbb{R}^{L \times L}$ mixing the tokens of the input \mathbf{x} .

Empirical risk minimization We study the learning of the attention layer (2), when a training set $\mathcal{D} = \{\mathbf{x}^\mu, y(\mathbf{x}^\mu)\}_{\mu=1}^n$ with n independently sampled sentences $\{\mathbf{x}^\mu\}_{\mu=1}^n$, and the associated labels $\{y(\mathbf{x}^\mu)\}_{\mu=1}^n$, is available. The target (1) can be learnt by carrying out an empirical risk minimization, with the generalization error measured at test time by the mean squared error (MSE)

$$\hat{\mathbf{Q}} = \operatorname{argmin}_{\mathbf{Q} \in \mathbb{R}^{d \times r}} \left[\sum_{\mu=1}^n \frac{1}{2d} \|y(\mathbf{x}^\mu) - f_{\mathbf{Q}}(\mathbf{x}^\mu)\|^2 + \frac{\lambda}{2} \|\mathbf{Q}\|^2 \right]; \quad \epsilon_g \equiv \frac{1}{dL} \mathbb{E}_{\mathbf{x} \sim p_x} \|y(\mathbf{x}) - f_{\hat{\mathbf{Q}}}(\mathbf{x})\|^2. \quad (3)$$

3. Closed-form characterization of the training loss

We analyze the learning problem (3) in the limit where the embedding dimension d and the number of training samples n jointly tend to infinity, while their ratio $\alpha = n/d$, the sample complexity, stays of order $\Theta_d(1)$. We further assume the sentence length L , the ranks r_s, r_t of the weights \mathbf{Q}, \mathbf{Q}_* , and the norm of the positional embeddings $\|\mathbf{p}\|$ to be $\Theta_d(1)$. This limit has been considered in a stream of previous works (e.g. [14, 32, 36]) and allows to derive closed-form characterization of the ERM problem (3). It also exhibits a particularly rich learning phenomenology which we further explore in Section 4.

The main technical result of the present work is a closed-form characterization of the test MSE and training loss (3) achieved in the high-dimensional limit for the model (2) trained via the empirical risk minimization of (3), and we state it stated in Appendix A. The derivation is exploiting a mapping of the model (2) to a (variant of) a Generalized Linear Model (GLM) [34, 37]. The summary statistics characterized by the self-consistent state evolution equations (8) [26] asymptotically describe the fixed points of a Generalized Approximate Message Passing (GAMP) algorithm [42] (A). The fixed points of GAMP in turn correspond to critical (zero-gradient) points of the non-convex empirical loss landscape (3). Therefore, while the technical Result 1 is stated as a characterization of the global minimum of (3), which is the main concern of the present work, solutions of (8) also describe local minima and saddles. Note that Appendix B provides an alternative derivation of the result for different losses using the replica method from statistical physics [40]. The methodology we use is similar to many recent work that study asymptotics of a large number of high-dimensional problems, e.g. [7, 15, 17, 19, 31].

4. Positional-to-semantic phase transition

Rank one dot-product attention In the following, we turn to a special case of tied low-rank attention (2) – namely a dot-product attention layer, with the student

$$\mathbf{S} \left[\frac{1}{\sqrt{d}} (\mathbf{x} + \mathbf{p}) \mathbf{Q} \right] = \operatorname{softmax} \left(\frac{1}{d} (\mathbf{x} + \mathbf{p}) \mathbf{Q} \mathbf{Q}^\top (\mathbf{x} + \mathbf{p})^\top \right), \quad (4)$$

and a specific case of target attention matrix (1) of the form

$$\mathbf{T} \left[\frac{1}{\sqrt{d}} \mathbf{x} \mathbf{Q}_* \right] = (1 - \omega) \operatorname{softmax} \left(\frac{1}{d} \mathbf{x} \mathbf{Q}_* \mathbf{Q}_*^\top \mathbf{x}^\top \right) + \omega A. \quad (5)$$

with $A \in \mathbb{R}^{L \times L}$ a fixed matrix. In (5), the parameter $\omega \in [0, 1]$ tunes the relative strength of the dot-product term and the fixed matrix term, and interpolates between a fully positional and a fully semantic task: With $\omega = 0$ we have a purely semantic target as the i, j -th element of the score matrix $\operatorname{softmax}(1/d \mathbf{x} \mathbf{Q}_* \mathbf{Q}_*^\top \mathbf{x}^\top)$ only depends on the tokens $\mathbf{x}_i, \mathbf{x}_j$ and not explicitly on their respective placements i, j inside the sentence. For $\omega = 1$, the attention matrix A associated thereto is purely

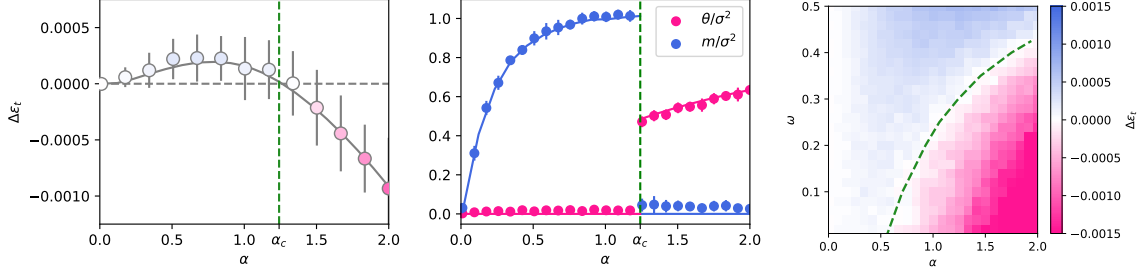


Figure 1: **Mixed positional/semantic teacher for $\omega = 0.3$.** Setting is $r_s = r_t = 1, L = 2, A = ((0.6, 0.4), (0.4, 0.6)), \Sigma_1 = \Sigma_2 = 0.25\mathbb{I}_d, \mathbf{p}_1 = \mathbf{1}_d = -\mathbf{p}_2$ and $\mathbf{Q}_\star \sim \mathcal{N}(0, \mathbb{I}_d)$. We compare the solutions of (8) in Result 1 (solid lines) with the models (2) trained using gradient descent initialized resp. at \mathbf{Q}_\star and at \mathbf{p}_1 (markers). The green dashed line is the theoretical prediction for the threshold $\alpha_c(\omega)$ above which the semantic solution of (8) in Result 1 has lower loss than the positional solution. **(left)** Difference in training loss $\Delta\epsilon_t$ for $\omega = 0.3$. **(center)** overlap θ between the learnt weights $\hat{\mathbf{Q}}$ and the target weights \mathbf{Q}_\star . **Overlap m** between the learnt weights $\hat{\mathbf{Q}}$ and the positional embedding \mathbf{p}_1 , where only the solution of (8) corresponding to the lowest found training loss is represented. **(right)** Empirical difference in training loss for a range of ω, α .

positional, in the sense that A_{ij} is a function of i, j but not of $\mathbf{x}_i, \mathbf{x}_j$. To complete the learning task, a *positional* mechanism then needs to be learnt. The parameter ω thus allows to tune the amount of semantic/positional content in the target (5), and thus the extent to which the task requires the model to implement semantic attention (small ω s) or rather positional attention (large ω s). In the following, for definiteness, we further assume $r_s = r_t = 1$ and set \mathbf{Q}_\star to be a fixed random Gaussian vector drawn from $\mathcal{N}(0, \mathbb{I}_d)$, and choose the positional encodings $\mathbf{p}_1 = -\mathbf{p}_2 = \mathbf{1}_d$. Finally, for simplicity, we consider sentences with two tokens $L = 2$ and isotropic token covariances $\Sigma_1 = \Sigma_2 = \sigma^2\mathbb{1}_d$.

Semantic and positional mechanisms The summary statistics θ_ℓ, m_ℓ describing the global minimizer of the empirical loss minimization (3) of the dot-product attention (4) on the target (5) are captured alongside the corresponding test error (3) and training loss (3), by Result 1. The solution of the system of equations (8) is not unique, and different stable fixed points describe different corresponding critical points of the non-convex empirical loss landscape (3). In practice, we find two solutions of (8), corresponding to two minima associated with different mechanisms implemented by the dot-product attention (4) when approximating the target (5):

–**Positional solution** One solution of (8) correspond to vanishing overlap $\theta = 0$ between the trained weights $\hat{\mathbf{Q}}$ and the semantic target weights \mathbf{Q}_\star , and non-zero $m > 0$ between the trained weights $\hat{\mathbf{Q}}$ and the positional embedding $\mathbf{p}_1 = -\mathbf{p}_2$. Consequently, the argument of the dot-product attention $\hat{\mathbf{Q}}(\mathbf{x} + \mathbf{p})$ has a sizeable token-independent –thus positional– contribution $\hat{\mathbf{Q}}\mathbf{p}$, alongside a token-dependent semantic part $\hat{\mathbf{Q}}\mathbf{x}$. Because of the positional terms, the resulting learnt attention matrix $\text{softmax}(1/d(\mathbf{x} + \mathbf{p})\hat{\mathbf{Q}}\hat{\mathbf{Q}}^\top(\mathbf{x} + \mathbf{p})^\top)$ implements a partly positional mechanism.

–**Semantic solution** Another solution of the system of equations (8) is associated with a vanishing overlap $m = 0$ between the learnt weights $\hat{\mathbf{Q}}$ and the positional embeddings, and a finite overlap $\theta > 0$ with the target weights \mathbf{Q}_\star . Therefore the resulting learnt attention matrix $\text{softmax}(1/d(\mathbf{x} + \mathbf{p})\hat{\mathbf{Q}}\hat{\mathbf{Q}}^\top(\mathbf{x} + \mathbf{p})^\top) \approx \text{softmax}(1/d\mathbf{x}\hat{\mathbf{Q}}\hat{\mathbf{Q}}^\top\mathbf{x}^\top)$ is largely semantic.

Positional-to-semantic phase transition While the system of self-consistent equations (8) may admit other solutions, we did not find solutions with lower training loss than the two aforesaid fixed points. Which of these solutions corresponds to the global minimum – and thus the solution of the optimization (3) – depends on the sample complexity α and the positional/semantic parameter ω (5). For a fixed parameter ω in (5), an analysis of equations (8) (Appendix B), reveals that for a sizeable range of ω there exists a threshold α_c for the sample complexity so that for $\alpha < \alpha_c$, the global minimum of (3) corresponds to a positional mechanism, and is described by the positional solution of (8) of Result 1 with $\theta = 0, m > 0$. For $\alpha > \alpha_c$, the global minimum of (3) corresponds to a semantic mechanism, and is described by the semantic solution of (8) of Result 1 with $\theta > 0, m = 0$.

The dot-product attention thus displays a *phase transition in sample complexity from a positional mechanism to a semantic mechanism*, implementing the simpler positional mechanism when having access to small amounts of data, and only learning the semantic content of the target (5) when presented sufficient data. The critical sample complexity α_c generically grows with the positionality ω of the target function (5), as the semantic content – i.e. the first term of (5) – is less apparent for larger ω , and thus requires larger amounts of data to be identified and approximated by the dot-product attention (4). In Fig. 1 (left) for $\omega = 0.3$ the difference in training loss $\Delta\epsilon_t$ between the positional and semantic solutions of (8) is represented, alongside the difference in training loss at convergence experimentally reached by gradient descent. For small (resp. large) sample complexity $\alpha < \alpha_c$ (resp. $\alpha > \alpha_c$), the training loss of the positional (resp. semantic) minimum is lower, and thus corresponds to the global minimum.

Experimentally, the positional minimum can be reached for $\alpha < \alpha_c$ via gradient descent by initializing the weights Q of the attention (4) close to the positional embedding p_1 . By the same means, the semantic minimum can be reached from an initialization at the teacher weights Q_* (5). Note that the semantic initialization is informed in nature, in that it necessitates the knowledge of the target parameters Q_* . We conduct numerical experiments from a random initialization of Q in Appendix C.3, and show that the dynamics may reach either of the local minima, or get stuck in a different one.

In Fig. 1 (center), we compare our analytical characterizations for different metrics at the global minimum – the summary statistics θ, m (middle), and the test MSE (right) –, with the corresponding experimental estimates, obtained by optimizing (3) with the `PyTorch` implementation of gradient descent, from a positional (resp. semantic) initialization for $\alpha < \alpha_c$ (resp. $\alpha > \alpha_c$), displaying overall good agreement. In Appendix C.1 we further verify that in the scaling limit of our analysis, namely $n, d \rightarrow \infty$ for $\alpha = O(1)$, the agreement improves with growing n, d . Overall, we observe the emergence of semantic learning as a function of the task and the sample complexity in Fig. 1.

Conclusion

We explored the interplay between positional and semantic attention, through the prism of tied low-rank self-attention in high dimensions. In a theoretically controlled setting, we characterized the global optimum of the empirical loss, when learning a target attention layer. This global optimum was found to correspond to either a positional or a semantic mechanism, with a phase transition between the two mechanisms occurring as the sample complexity increases. We believe the present asymptotic analysis of the inner workings of attention mechanisms opens up exciting research directions. Considering untied query and key matrices, appending a readout network after the attention layer, or addressing more practical training procedures such as masked language modelling, are some possible extensions which will hopefully pave the way towards a satisfactory theoretical comprehension of attention mechanisms.

References

- [1] Sanjeev Arora and Anirudh Goyal. A theory for emergence of complex skills in language models. *arXiv preprint arXiv:2307.15936*, 2023.
- [2] Benjamin Aubin, Antoine Maillard, Florent Krzakala, Nicolas Macris, Lenka Zdeborová, et al. The committee machine: Computational to statistical gaps in learning a two-layers neural network. *Advances in Neural Information Processing Systems*, 31, 2018.
- [3] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*, 2023.
- [4] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Optimal errors and phase transitions in high-dimensional generalized linear models. *Proceedings of the National Academy of Sciences*, 116(12):5451–5460, 2019.
- [5] Adriano Barra, Giuseppe Genovese, Peter Sollich, and Daniele Tantari. Phase transitions in restricted boltzmann machines with generic priors. *Physical Review E*, 96(4):042156, 2017.
- [6] Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2): 764–785, 2011.
- [7] Mohsen Bayati and Andrea Montanari. The lasso risk for gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2011.
- [8] Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36, 2024.
- [9] Enric Boix-Adsera, Etai Littwin, Emmanuel Abbe, Samy Bengio, and Joshua Susskind. Transformers learn through gradual rank increase. *arXiv preprint arXiv:2306.07042*, 2023.
- [10] Erwin Bolthausen. An iterative construction of solutions of the tap equations for the sherrington–kirkpatrick model. *Communications in Mathematical Physics*, 325(1):333–366, 2014.
- [11] Angelica Chen, Ravid Schwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in mlms. *arXiv preprint arXiv:2309.07311*, 2023.
- [12] Lucas Clarté, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Theoretical characterization of uncertainty in high-dimensional linear classification. *Machine Learning: Science and Technology*, 4(2):025029, 2023.
- [13] Elisabetta Cornacchia, Francesca Mignacco, Rodrigo Veiga, Cédric Gerbelot, Bruno Loureiro, and Lenka Zdeborová. Learning curves for the multi-class teacher–student perceptron. *Machine Learning: Science and Technology*, 4(1):015019, 2023.
- [14] Hugo Cui and Lenka Zdeborová. High-dimensional asymptotics of denoising autoencoders. *arXiv preprint arXiv:2305.11041*, 2023.

- [15] David Donoho and Andrea Montanari. High dimensional robust m-estimation: Asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166:935–969, 2016.
- [16] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages 5793–5831. PMLR, 2022.
- [17] Melikasadat Emami, Mojtaba Sahraee-Ardakan, Parthe Pandit, Sundeep Rangan, and Alyson Fletcher. Generalization error of generalized linear models in high dimensions. In *International Conference on Machine Learning*, pages 2892–2901. PMLR, 2020.
- [18] Hengyu Fu, Tianyu Guo, Yu Bai, and Song Mei. What can a single attention layer learn? a study through the random features lens. *arXiv preprint arXiv:2307.11353*, 2023.
- [19] Cedric Gerbelot, Alia Abbara, and Florent Krzakala. Asymptotic errors for teacher-student convex generalized linear models (or: How to prove kabashima’s replica formula). *IEEE Transactions on Information Theory*, 69(3):1824–1852, 2022.
- [20] Yehoram Gordon. On milman’s inequality and random subspaces which escape through a mesh in \mathbb{R}^n . In *Geometric Aspects of Functional Analysis: Israel Seminar (GAFA) 1986–87*, pages 84–106. Springer, 1988.
- [21] Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. How do transformers learn in-context beyond simple functions? a case study on learning with representations. *arXiv preprint arXiv:2310.10616*, 2023.
- [22] Géza Györgyi. First-order transition to perfect generalization in a neural network with binary synapses. *Physical Review A*, 41(12):7097, 1990.
- [23] Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions of the Association for Computational Linguistics*, 8:156–171, 2020.
- [24] Adi Haviv, Ori Ram, Ofir Press, Peter Izsak, and Omer Levy. Transformer language models without positional encodings still learn positional information. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1382–1390, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.99. URL <https://aclanthology.org/2022.findings-emnlp.99>.
- [25] Ernst Ising. Contribution to the theory of ferromagnetism. *Z. Phys*, 31(1):253–258, 1925.
- [26] Adel Javanmard and Andrea Montanari. State evolution for general approximate message passing algorithms, with applications to spatial coupling. *Information and Inference: A Journal of the IMA*, 2(2):115–144, 2013.
- [27] Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. *Advances in Neural Information Processing Systems*, 35:37822–37836, 2022.

- [28] Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan, Payel Das, and Siva Reddy. The impact of positional encoding on length generalization in transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Drrl2gcjzl>.
- [29] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. *arXiv preprint arXiv:2302.06015*, 2023.
- [30] Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, pages 19565–19594. PMLR, 2023.
- [31] Bruno Loureiro, Gabriele Sicuro, Cedric Gerbelot, Alessandro Pocco, Florent Krzakala, and Lenka Zdeborová. Learning gaussian mixtures with generalized linear models: Precise asymptotics in high-dimensions. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 10144–10157. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/543e83748234f7cbab21aa0ade66565f-Paper.pdf.
- [32] Bruno Loureiro, Gabriele Sicuro, Cédric Gerbelot, Alessandro Pocco, Florent Krzakala, and Lenka Zdeborová. Learning gaussian mixtures with generalized linear models: Precise asymptotics in high-dimensions. *Advances in Neural Information Processing Systems*, 34:10144–10157, 2021.
- [33] Antoine Maillard, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Phase retrieval in high dimensions: Statistical and computational phase transitions. *Advances in Neural Information Processing Systems*, 33:11071–11082, 2020.
- [34] Peter McCullagh. *Generalized linear models*. Routledge, 2019.
- [35] Marc Mézard, Giorgio Parisi, and Miguel Angel Virasoro. *Spin glass theory and beyond: An Introduction to the Replica Method and Its Applications*, volume 9. World Scientific Publishing Company, 1987.
- [36] Francesca Mignacco, Florent Krzakala, Yue Lu, Pierfrancesco Urbani, and Lenka Zdeborova. The role of regularization in classification of high-dimensional noisy gaussian mixture. In *International conference on machine learning*, pages 6874–6883. PMLR, 2020.
- [37] John Ashworth Nelder and Robert WM Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 135(3):370–384, 1972.
- [38] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer*

- Circuits Thread*, 2022. <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- [39] Lars Onsager. Crystal statistics. i. a two-dimensional model with an order-disorder transition. *Physical Review*, 65(3-4):117, 1944.
- [40] Giorgio Parisi. Toward a mean field theory for spin glasses. *Physics Letters A*, 73(3):203–205, 1979.
- [41] Giorgio Parisi. Order parameter for spin-glasses. *Physical Review Letters*, 50(24):1946, 1983.
- [42] Sundeep Rangan, Philip Schniter, Erwin Riegler, Alyson K Fletcher, and Volkan Cevher. Fixed points of generalized approximate message passing with arbitrary matrices. *IEEE Transactions on Information Theory*, 62(12):7464–7474, 2016.
- [43] Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 14228–14246. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/2e10b2c2e1aa4f8083c37dfe269873f8-Paper-Conference.pdf.
- [44] Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context classification task. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=aN4Jf6Cx69>.
- [45] Riccardo Rende, Federica Gerace, Alessandro Laio, and Sebastian Goldt. Optimal inference of a generalised Potts model by single-layer transformers with factored attention. *arXiv preprint arXiv:2304.07235*, 2023.
- [46] Anian Ruoss, Grégoire Delétang, Tim Genewein, Jordi Grau-Moya, Róbert Csordás, Mehdi Bennani, Shane Legg, and Joel Veness. Randomized positional encodings boost length generalization of transformers. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1889–1903, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-short.161. URL <https://aclanthology.org/2023.acl-short.161>.
- [47] Arda Sahiner, Tolga Ergen, Batu Ozturkler, John Pauly, Morteza Mardani, and Mert Pilanci. Unraveling attention via convex duality: Analysis and interpretations of vision transformers. In *International Conference on Machine Learning*, pages 19050–19088. PMLR, 2022.
- [48] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language models a mirage? *Advances in Neural Information Processing Systems*, 36, 2024.
- [49] Henry Schwarze. Learning a rule in a multilayer neural network. *Journal of Physics A: Mathematical and General*, 26(21):5781, 1993.

- [50] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-Attention with Relative Position Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 464–468. Association for Computational Linguistics, 2018. doi: 10.18653/v1/N18-2074. URL <https://aclanthology.org/N18-2074>.
- [51] Aaditya K Singh, Stephanie C.Y. Chan, Ted Moskovitz, Erin Grant, Andrew M Saxe, and Felix Hill. The transient nature of emergent in-context learning in transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Of0GBzow8P>.
- [52] Koustuv Sinha, Amirhossein Kazemnejad, Siva Reddy, Joelle Pineau, Dieuwke Hupkes, and Adina Williams. The curious case of absolute position embeddings. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4449–4472, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.326. URL <https://aclanthology.org/2022.findings-emnlp.326>.
- [53] Haim Sompolinsky, Naftali Tishby, and H Sebastian Seung. Learning from examples in large neural networks. *Physical Review Letters*, 65(13):1683, 1990.
- [54] Mihailo Stojnic. Meshes that trap random subspaces. *arXiv preprint arXiv:1304.0003*, 2013.
- [55] Mihailo Stojnic. Upper-bounding ℓ_1 -optimization weak thresholds. *arXiv preprint arXiv:1303.7289*, 2013.
- [56] Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines. *arXiv preprint arXiv:2308.16898*, 2023.
- [57] Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [58] Christos Thrampoulidis, Samet Oymak, and Babak Hassibi. The gaussian min-max theorem in the presence of convexity. *arXiv preprint arXiv:1408.4837*, 2014.
- [59] Yuandong Tian, Yiping Wang, Beidi Chen, and Simon Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. *arXiv preprint arXiv:2305.16380*, 2023.
- [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- [61] Taylor Webb, Keith J Holyoak, and Hongjing Lu. Emergent analogical reasoning in large language models. *Nature Human Behaviour*, 7(9):1526–1541, 2023.

- [62] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [63] Lenka Zdeborová and Florent Krzakala. Statistical physics of inference: Thresholds and algorithms. *Advances in Physics*, 65(5):453–552, 2016.
- [64] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*, 2023.
- [65] Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob Andreas. The clock and the pizza: Two stories in mechanistic explanation of neural networks, 2023.

Appendix A. Derivation of Main Result

In this Appendix, we provide a detailed derivation our works main technical result, as stated in subsection A.1. In subsection A.3 we introduce a Generalized Approximate Message Passing algorithm (GAMP) [42]. Subsection A.4 then establishes that equations (8) of Result 1 track the dynamics of summary statistics describing the GAMP algorithm. In particular, the equations (8) describe the fixed points of GAMP. Finally, subsection A.5 shows that fixed points of GAMP correspond to critical (zero-gradient) points of the empirical loss landscape (3), thus establishing that equations 8 of Result 1 describe fixed points of GD.

A.1. Main Result

Assumption 1 *The covariances $\{\Sigma_\ell\}_{\ell=1}^L$ admit a common set of eigenvectors $\{e_i\}_{i=1}^d$. We further note $\{\lambda_i^\ell\}_{i=1}^d$ the eigenvalues of Σ_ℓ . The eigenvalues $\{\lambda_i^\ell\}_{i=1}^d$ and the projection of the positional embedding $\{\mathbf{p}_\ell\}_{\ell=1}^L$ and the teacher columns $\{\mathbf{Q}_j^*\}_{j=1}^{r_t}$ on the eigenvectors $\{e_i^\top \mathbf{p}_\ell\}_{i,\ell}$, $\{e_i^\top \mathbf{Q}_j^*\}_{i,j}$ are assumed to admit a well-defined joint distribution ν as $d \rightarrow \infty$ – namely, for $\gamma = (\gamma_1, \dots, \gamma_L) \in \mathbb{R}^L, \pi = (\pi_1, \dots, \pi_{r_t}) \in \mathbb{R}^{r_t}$ and $\tau = (\tau_1, \dots, \tau_L) \in \mathbb{R}^L$:*

$$\frac{1}{d} \sum_{i=1}^d \prod_{\ell=1}^L \delta(\lambda_i^\ell - \gamma_\ell) \delta(\sqrt{d} e_i^\top \mathbf{p}_\ell - \tau_\ell) \prod_{j=1}^{r_t} \delta(e_i^\top \mathbf{Q}_j^* - \pi_j) \xrightarrow{d \rightarrow \infty} \nu(\gamma, \tau, \pi). \quad (6)$$

Result 1 *Under Assumption 1, in the limit $n, d \rightarrow \infty$, $\|\mathbf{p}\|, n/d, L, r_s, r_t = \Theta_d(1)$, the summary statistics*

$$\begin{aligned} \rho_\ell &\equiv \frac{\mathbf{Q}_*^\top \Sigma_\ell \mathbf{Q}_*}{d} \in \mathbb{R}^{r_t \times r_t}, & q_\ell &\equiv \frac{\hat{\mathbf{Q}}^\top \Sigma_\ell \hat{\mathbf{Q}}}{d} \in \mathbb{R}^{r_s \times r_s}, \\ m_\ell &\equiv \frac{\hat{\mathbf{Q}}^\top \mathbf{p}_\ell}{d} \in \mathbb{R}^{r_s}, & \theta_\ell &\equiv \frac{\hat{\mathbf{Q}}^\top \Sigma_\ell \mathbf{Q}_*}{d} \in \mathbb{R}^{r_s \times r_t} \end{aligned} \quad (7)$$

concentrate in probability, and are solutions of the set of finite-dimensional self-consistent equations

$$\begin{cases} q_\ell = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \left(\sum_{\kappa=1}^L \gamma_\kappa \hat{q}_\kappa + \left(\sum_{\kappa=1}^L \hat{m}_\kappa \tau_\kappa + \gamma_\kappa \hat{\theta}_\kappa \cdot \pi \right)^{\otimes 2} \right) \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \\ V_\ell = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \\ m_\ell = \int d\nu(\gamma, \tau, \pi) \tau_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \left(\sum_{\kappa=1}^L \hat{m}_\kappa \tau_\kappa + \gamma_\kappa \hat{\theta}_\kappa \cdot \pi \right) \\ \theta_\ell = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \left(\sum_{\kappa=1}^L \hat{m}_\kappa \tau_\kappa + \gamma_\kappa \hat{\theta}_\kappa \cdot \pi \right) \pi^\top. \end{cases} \quad (8)$$

$$\begin{cases} \hat{q}_\ell = \alpha \mathbb{E}_{\Xi, U} V_\ell^{-1} \left(\text{prox}(\Xi, U)_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right)^{\otimes 2} V_\ell^{-1} \\ \hat{V}_\ell = \hat{\theta}_\ell \theta_\ell^\top q_\ell^{-1} - \alpha \mathbb{E}_{\Xi, U} V_\ell^{-1} \left(\text{prox}(\Xi, U)_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \xi_\ell^\top q_\ell^{-\frac{1}{2}} \\ \hat{m}_\ell = \alpha \mathbb{E}_{\Xi, U} V_\ell^{-1} \left(\text{prox}(\Xi, U)_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \\ \hat{\theta}_\ell = \alpha \mathbb{E}_{\Xi, U} V_\ell^{-1} \left(\text{prox}(\Xi, U)_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \left(u_\ell - \xi_\ell^\top q_\ell^{-1/2} \theta_\ell \right)^\top \left(\rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell \right)^{-1} \end{cases} \quad (9)$$

In (8), $U = \{u_\ell\}_{\ell=1}^L$ and $\Xi = \{\xi_\ell\}_{\ell=1}^L$, with $u_\ell \sim \mathcal{N}(\xi_\ell^\top q_\ell^{-1/2} \theta_\ell, \rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell)$ and $\xi_\ell \sim \mathcal{N}(0, \mathbb{I}_{r_s})$, and $\cdot^{\otimes 2}$ denotes the outer product of a vector with itself. Finally, the resolvents $\{\text{prox}(\Xi, U)_\ell\}_{\ell=1}^L$ are defined as the minimizers of the Moreau envelope

$$\mathcal{M}(\Xi, U) = \inf_{z_1, \dots, z_L} \left\{ \sum_{\ell=1}^L \text{Tr} \left[V_\ell^{-1} \left(x_\ell - q_\ell^{1/2} \xi_\ell - m_\ell \right)^{\otimes 2} \right] + \text{Tr} \left[\mathbf{S}(Z) \rho_\Sigma \mathbf{S}(Z)^\top \right] - 2 \text{Tr} \left[\mathbf{T}(U) \rho_\Sigma \mathbf{S}(Z)^\top \right] \right\}.$$

We noted $Z \in \mathbb{R}^{L \times r_s}$ (resp. $U \in \mathbb{R}^{L \times r_t}$) the matrix whose rows are z_ℓ (resp. u_ℓ) and:

$$\rho_\Sigma \equiv \text{diag} \left[\left(\int d\nu(\gamma, \tau, \pi) \gamma_\ell \right)_{\ell=1}^L \right] \in \mathbb{R}^{L \times L}. \quad (10)$$

In the same limit, the test error (3) converges in probability to

$$\epsilon_g = \mathbb{E}_h \text{Tr} \left[\mathbf{S}[h] \rho_\Sigma \mathbf{S}[h]^\top \right] + \mathbb{E}_{h^*} \text{Tr} \left[\mathbf{T}[h^*] \rho_\Sigma \mathbf{T}[h^*]^\top \right] - 2 \mathbb{E}_{h, h^*} \text{Tr} \left[\mathbf{S}[h] \rho_\Sigma \mathbf{T}[h^*]^\top \right]. \quad (11)$$

where the average bears on $h \in \mathbb{R}^{L \times r_s}$, $h^* \in \mathbb{R}^{L \times r_t}$ with independent rows with statistics

$$(h_\ell, h_\ell^*) \sim \mathcal{N} \left[\begin{pmatrix} m_\ell \\ 0 \end{pmatrix}, \begin{pmatrix} q_\ell & \theta_\ell \\ \theta_\ell^\top & \rho_\ell \end{pmatrix} \right] \quad (12)$$

Finally, the training loss ϵ_t converges in probability to

$$\epsilon_t = \alpha \mathbb{E}_{Y, \Xi} \mathcal{M} - \frac{1}{2} \sum_{\ell=1}^L \text{Tr}[\hat{q}_\ell V_\ell] + \frac{\lambda}{2} \int d\nu(\gamma, \tau) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right]. \quad (13)$$

A.2. Notations

For simplicity, we place ourselves in the setting $r_s = 1$ explored in Section 4 of the main text, but allow the length L of the sentences to be arbitrary, and allow a generic learning model \mathbf{S} (2), i.e. not necessarily the dot-product attention model analyzed in Section 4. The case $r_s \geq 2$ follows identical derivation steps, modulo the replacement of all variables by tensor objects. We provide another alternative derivation of Result 1 in full generality in Appendix B, using the replica method from statistical physics. Let us note $\{X_\ell\}_{1 \leq \ell \leq L}$ a series of L $n \times d$ matrices, with X_ℓ corresponding to the ℓ -th rows (tokens) of each input sentence x^μ stacked vertically, and normalized by \sqrt{d} . We denote $\tilde{X}_\ell \equiv X_\ell + P_\ell$, where $P \in \mathbb{R}^{n \times d}$ is the matrix with all rows equal to the ℓ -th positional encoding p_ℓ . Let us further define $\rho \in \mathbb{R}^{n \times L \times L}$ the tensor corresponding to the sequence of n matrices $\frac{1}{d} x^\mu (x^\mu)^\top \in \mathbb{R}^{L \times L}$. Finally, let us denote $T \in \mathbb{R}^{n \times L \times L}$ the tensor so that the μ -th row of T satisfies $y(x^\mu) = T^\mu x^\mu$, see equation (1). In other words, T corresponds to the concatenation of the target attention matrices.

Before detailing the derivation, we first highlight a simplifying observation. Note that a loss item can be expanded as

$$\begin{aligned} \frac{1}{d} \left\| y(x^\mu) - \mathbf{S} \left[\frac{1}{\sqrt{d}} (x^\mu + p) Q \right] x^\mu \right\|^2 &= \|y(x^\mu)\|^2 + \text{Tr} \mathbf{S} \left[\frac{1}{\sqrt{d}} (x^\mu + p) Q \right] \rho_\Sigma \mathbf{S} \left[\frac{1}{\sqrt{d}} (x^\mu + p) Q \right]^\top \\ &\quad - 2 \text{Tr} \mathbf{T} \left[\frac{1}{\sqrt{d}} x^\mu Q_\star \right] \rho_\Sigma \mathbf{S} \left[\frac{1}{\sqrt{d}} (x^\mu + p) Q \right]^\top, \end{aligned} \quad (14)$$

where we used that with high probability in the considered asymptotic limit, for all $1 \leq \mu \leq n$,

$$xx^\top = (x+p)(x+p)^\top = x(x+p)^\top = \rho_\Sigma. \quad (15)$$

Since the first term of (14) does not depend on the weights Q , it can be without loss of generality subtracted from the loss. Without loss of generality, one can thus consider the equivalent empirical risk minimization problem

$$\hat{Q} = \operatorname{argmin}_{Q \in \mathbb{R}^{d \times r}} \sum_{\mu=1}^n \frac{1}{2d} \left[\operatorname{Tr} \mathbf{s} \left[\frac{1}{\sqrt{d}} (x^\mu + p) Q \right] \rho_\Sigma \mathbf{s} \left[\frac{1}{\sqrt{d}} (x^\mu + p) Q \right]^\top - 2 \operatorname{Tr} \mathbf{T} \left[\frac{1}{\sqrt{d}} x_\ell Q_* \right] \rho_\Sigma \mathbf{s} \left[\frac{1}{\sqrt{d}} (x_\ell^\mu + p_\ell) Q \right]^\top \right] + \frac{\lambda}{2} \|Q\|^2. \quad (16)$$

The risks (16) and (3) are equivalent, and we shall use the former in the following.

Finally, for arguments $T \in \mathbb{R}^{L \times L}$, $\rho \in \mathbb{R}^{L \times L}$, $\omega \in \mathbb{R}^L$, $V \in \mathbb{R}^{L \times L}$ we introduce the resolvent

$$\operatorname{prox}(T, \rho, \omega, V) \equiv \operatorname{arginf}_{x = \{x_\ell \in \mathbb{R}\}_{\ell=1}^L} \left\{ \sum_{\ell, \kappa=1}^L (x_\ell - \omega_\ell) (V^{-1})_{\ell\kappa} (x_\kappa - \omega_\kappa) - 2 \operatorname{Tr} \left[\mathbf{S}[x] \rho T^\top \right] + \operatorname{Tr} \left[\mathbf{S}[x] \rho \mathbf{S}[x]^\top \right] \right\} \quad (17)$$

Note that the latter part of the bracketed term corresponds to the simplified loss (14) derived in the beginning of Appendix B, which is the one we shall without loss of generality consider in the present appendix. For ease of presentation, we place ourselves under Assumption 1, where all the input covariances $\{\Sigma_\ell\}_\ell$ are codiagonalizable. In the following, without loss of generality, we thus assume them diagonal, by placing ourselves in the common basis $\{e_i\}_{1 \leq i \leq d}$ of Assumption 1.

A.3. AMP algorithm

We are now in a position to state the AMP algorithm:

Algorithm 1 GAMP

Inputs : $\{\tilde{X}_\ell \in \mathbb{R}^{n \times d}\}_{\ell=1}^L, T \in \mathbb{R}^{n \times L \times L}, \rho \in \mathbb{R}^{n \times L \times L}$
Initialize $\hat{Q}^0 = \sim \mathcal{N}(0, \mathbb{I}_d), \hat{c}^0 = \mathbb{I}_d, \{f_\ell^0 = 0_n\}_{\ell=1}^L$
for $t \leq t_{\max}$ **do**
 $\forall 1 \leq \ell, \kappa \leq L, V_{\ell\kappa}^t = (\tilde{X}_\ell \odot \tilde{X}_\kappa) \hat{c}^t$
 $\forall 1 \leq \ell \leq L, \omega_\ell^t = \tilde{X}_\ell \hat{Q}^t - \sum_{\kappa=1}^L V_{\ell\kappa}^t f_\kappa^{t-1}$
 $\forall 1 \leq \ell \leq L, f_\ell^t = \sum_{\kappa} (V^{-1})_{\ell\kappa} (\text{prox}(T, \rho, \omega^t, V^t)_\kappa - \omega_\kappa^t)$
 $\forall 1 \leq \ell, \kappa \leq L, g_{\ell\kappa}^t = \partial_{\omega_\ell} f_\kappa^t$
 $A^t = - \sum_{\ell, \kappa=1}^L (\tilde{X}_\ell \odot \tilde{X}_\kappa)^\top g_{\ell\kappa}^t$
 $b^t = \sum_{\ell=1}^L \tilde{X}_\ell^\top f_\ell^t + A^t \odot \hat{Q}^t$
 $\hat{Q}^{t+1} = (\lambda \mathbb{I}_d + A^t)^{-1} b^t$
 $\hat{c}^{t+1} = (\lambda \mathbb{I}_d + A^t)^{-1}$
end for
return Estimator \hat{Q}

The GAMP algorithm can be derived in standard fashion from the Belief Propagation (BP) algorithm, see e.g. [6, 42] or [63] for an overview. Compared to the standard GAMP iterations for Generalized linear models, one needs to account for the fact that there exist different sources of data X_ℓ (corresponding to the ℓ -th tokens of each input sentence), and for the fact that the output of the equivalent GLM are $\mathbb{R}^{L \times L}$ -valued attention matrices. In the following subsection, we show that the fixed points of GAMP 1 correspond to critical points of the empirical loss (3), i.e. fixed points of Gradient Descent (GD), allowing to connect Result 1 to our numerical experiments using GD.

A.4. State evolution

In this section we show that the dynamics of the GAMP Algorithm 1 are tracked by the summary statistics of Result 1. In particular, the equations (8) describe the statistics of the GAMP fixed points. To see this, it is convenient to take as a starting point the relaxed Belief Propagation (rBP) equations, which are a step upstream in the derivation of the GAMP iterations, and which are asymptotically equivalent— see e.g. [63] for a review or e.g. [12], Appendix A, for a detailed walkthrough. The rBP equations read

As conventional, we note \cdot_μ the version of a variable $\cdot_{\mu \rightarrow i}$ where the summation also encompasses the index i , and \cdot_i the version of a variable $\cdot_{i \rightarrow \mu}$ where the summation also encompasses the index μ . Note that in all cases above the two variables differ by at most $\Theta_d(1/\sqrt{d})$.

Concentration of $(V_{\mu \rightarrow i}^t)_{\ell\kappa}$ We first study the statistics of $V_{\mu \rightarrow i}^t, A_{i \rightarrow \mu}^t$, remembering that the data $\tilde{x}_{\ell i}^\mu \equiv (x_\ell^\mu)_i / \sqrt{d} + (p_\ell)_i / \sqrt{d}$ in the notation of the main text, with $(x_\ell^\mu)_i = \Theta_d(1), (p_\ell)_i = \Theta_d(1/\sqrt{d})$.

Algorithm 2 rBP

Inputs : $\{\tilde{X}_\ell \in \mathbb{R}^{n \times d}\}_{\ell=1}^L, T \in \mathbb{R}^{n \times L \times L}, \rho \in \mathbb{R}^{n \times L \times L}$
Initialize $\forall 1 \leq \mu \leq n, 1 \leq i \leq d, \hat{Q}_{i \rightarrow \mu}^0 = 0, \hat{c}_{i \rightarrow \mu}^0 = 1, \{f_{\ell \mu \rightarrow i}^0 = 0\}_{\ell=1}^L$
for $t \leq t_{\max}$ **do**

$$\forall 1 \leq \ell, \kappa \leq L, 1 \leq \mu \leq n, 1 \leq i \leq d, (V_{\mu \rightarrow i}^t)_{\ell \kappa} = \sum_{j \neq i} (\tilde{x}_{\ell j}^\mu)(\tilde{x}_{\kappa j}^\mu) \hat{c}_{j \rightarrow \mu}^t$$

$$\forall 1 \leq \ell, 1 \leq \mu \leq n, 1 \leq i \leq d, \omega_{\ell, \mu \rightarrow i}^t = \sum_{j \neq i} \tilde{x}_{\ell, i}^\mu \hat{Q}_{j \rightarrow \mu}^t$$

$$\forall 1 \leq \ell, 1 \leq \mu \leq n, 1 \leq i \leq d, f_{\ell, \mu \rightarrow i}^t = \sum_{\kappa} (V_{\mu \rightarrow i}^{-1})_{\ell \kappa} (\text{prox}(T_\mu, \rho_\mu, \omega_{\mu \rightarrow i}^t, V_{\mu \rightarrow i}^t)_\kappa - \omega_{\kappa, \mu \rightarrow i}^t)$$

$$\forall 1 \leq \ell, \kappa \leq L, 1 \leq \mu \leq n, 1 \leq i \leq d, g_{\ell \kappa, \mu \rightarrow i}^t = \partial_{\omega_\ell} f_{\kappa \mu \rightarrow i}^t$$

$$\forall 1 \leq \mu \leq n, 1 \leq i \leq d, A_{i \rightarrow \mu}^t = - \sum_{\ell, \kappa=1}^L \sum_{\nu \neq \mu} (\tilde{x}_{\ell i}^\nu)(\tilde{x}_{\kappa i}^\nu) g_{\ell \kappa, \nu \rightarrow i}^t$$

$$\forall 1 \leq \mu \leq n, 1 \leq i \leq d, b_{i \rightarrow \mu}^t = \sum_{\ell=1}^L \sum_{\nu \neq \mu} x_{\ell i}^\nu f_{\ell, \nu \rightarrow i}^t$$

$$\forall 1 \leq \mu \leq n, 1 \leq i \leq d, \hat{Q}_{i \rightarrow \mu}^{t+1} = (\lambda \mathbb{I}_d + A_{i \rightarrow \mu}^t)^{-1} b_{i \rightarrow \mu}^t$$

$$\forall 1 \leq \mu \leq n, 1 \leq i \leq d, \hat{c}_{i \rightarrow \mu}^{t+1} = (\lambda \mathbb{I}_d + A_{i \rightarrow \mu}^t)^{-1}$$

end for
return Estimator \hat{Q}

Replacing in the rBP updates:

$$\begin{aligned} (V_{\mu \rightarrow i}^t)_{\ell \kappa} &= \sum_{j \neq i} (\tilde{x}_{\ell j}^\mu)(\tilde{x}_{\kappa j}^\mu) \hat{c}_{j \rightarrow \mu}^t \\ &= \underbrace{\frac{1}{d} \sum_{j \neq i} (x_\ell^\mu)_j (x_\kappa^\mu)_j \hat{c}_{j \rightarrow \mu}^t}_{\delta_{\ell \kappa} \Theta_d(1) + (1 - \delta_{\ell \kappa}) \Theta_d(1/\sqrt{d})} + \underbrace{\frac{1}{d} \sum_{j \neq i} (x_\ell^\mu)_j (p_\kappa)_j \hat{c}_{j \rightarrow \mu}^t}_{\Theta_d(1/d)} + (\ell \leftrightarrow \kappa) + \underbrace{\frac{1}{d} \sum_{j \neq i} (p_\ell)_j (p_\kappa)_j \hat{c}_{j \rightarrow \mu}^t}_{\Theta_d(1/d)} \\ &= \delta_{\ell \kappa} \frac{1}{d} \sum_j (\Sigma_\ell)_{jj} \hat{c}_j^t \equiv V_\ell^t \end{aligned} \quad (18)$$

Distribution of $\omega_{\ell, \mu \rightarrow i}^t$ Let us first introduce the teacher local field

$$h_{\mu, \ell} = \sum_i (x_\ell^\mu)_i Q_i^*. \quad (19)$$

 Like e.g. [12], Appendix A, we first ascertain the joint distribution of $h_{\mu, \ell}, \omega_{\ell, \mu \rightarrow i}^t$ with respect to the data. These variables have mean

$$\mathbb{E}[\omega_{\ell, \mu \rightarrow i}^t] = \frac{p_\ell^\top \hat{Q}^t}{\sqrt{d}} \equiv m_\ell^t \quad (20)$$

and respective variance

$$\mathbb{V}[\omega_{\ell, \mu \rightarrow i}^t \omega_{\kappa, \nu \rightarrow j}^t] = \delta_{\mu\nu} \delta_{\ell\kappa} \frac{1}{d} \sum_{i,j} \hat{Q}_i^t(\Sigma_\ell)_{ij} \hat{Q}_j^t \equiv \delta_{\mu\nu} \delta_{\ell\kappa} q_\ell^t \quad (21)$$

$$\mathbb{E}[h_{\mu\ell} h_{\nu\kappa}] = \delta_{\mu\nu} \delta_{\ell\kappa} \frac{1}{d} \sum_{i,j} Q_i^*(\Sigma_\ell)_{ij} Q_j^* \equiv \delta_{\mu\nu} \delta_{\ell\kappa} \rho_\ell \quad (22)$$

$$\mathbb{E}[h_{\mu\ell} \omega_{\kappa, \nu \rightarrow j}^t] = \delta_{\mu\nu} \delta_{\ell\kappa} \frac{1}{d} \sum_{i,j} Q_i^*(\Sigma_\ell)_{ij} \hat{Q}_j^t \equiv \delta_{\mu\nu} \delta_{\ell\kappa} \theta_\ell^t \quad (23)$$

Distribution of $b_{i \rightarrow \mu}^t$ Let us ascertain the distribution of $b_{i \rightarrow \mu}^t$.

$$\begin{aligned} b_{i \rightarrow \mu}^t &= \sum_{\ell} \sum_{\nu \neq \mu} (\tilde{x}_\ell^\nu)_i (V_\ell^t)^{-1} \underbrace{\left(\text{prox}(T_\nu, \rho_\nu, \{\omega_{\kappa, \nu \rightarrow i}^t\}_\kappa, V_{\nu \rightarrow i}^t) \right)}_{\equiv \text{pr}\ddot{\text{ox}}(T_\nu, \rho_\nu, \{\omega_{\kappa, \nu \rightarrow i}^t\}_\kappa, V_{\nu \rightarrow i}^t)_\ell} - \omega_{\ell, \nu \rightarrow i}^t \\ &= \sum_{\ell} \sum_{\nu \neq \mu} 1/\sqrt{d} ((x_\ell^\nu)_i + (p_\ell)_i) \left[\text{pr}\ddot{\text{ox}}(\mathbb{T}[\{h_{\nu \rightarrow i, \kappa}\}_\kappa], \rho_\nu, \{\omega_{\kappa, \nu \rightarrow i}^t\}_\kappa, V_{\nu \rightarrow i}^t)_\ell \right. \\ &\quad \left. + 1/\sqrt{d} \sum_{\gamma} (x_\gamma^\nu)_i Q_i^* \partial_{h_\gamma} \text{pr}\ddot{\text{ox}}(\mathbb{T}[\{h_{\nu \rightarrow i, \kappa}\}_\kappa], \rho_\nu, \{\omega_{\kappa, \nu \rightarrow i}^t\}_\kappa, V_{\nu \rightarrow i}^t)_\ell \right], \end{aligned} \quad (24)$$

leading asymptotically to

$$\begin{aligned} \mathbb{E}[b_{i \rightarrow \mu}^t] &= \sum_{\ell} (\sqrt{d} p_\ell)_i \underbrace{\alpha \mathbb{E}_{H=\{h_\kappa\} \Xi = \{\xi_\kappa\}} \text{pr}\ddot{\text{ox}}(\mathbb{T}[H], \rho_\Sigma, \{m_\kappa^t + \sqrt{q_\kappa^t} \xi_\kappa\}_\kappa, \{V_\kappa^t\}_\kappa)_\ell}_{\equiv \hat{m}_\ell^t} \\ &\quad + Q_i^* \sum_{\ell} (\Sigma_\ell)_{ii} \underbrace{\alpha \mathbb{E}_{H, \Xi} \partial_{h_\ell} \text{pr}\ddot{\text{ox}}(\mathbb{T}[H], \rho_\Sigma, \{m_\kappa^t + \sqrt{q_\kappa^t} \xi_\kappa\}_\kappa, \{V_\kappa^t\}_\kappa)_\ell}_{\equiv \hat{\theta}_\ell^t} \end{aligned} \quad (25)$$

where the expectations bear over $\xi_\ell \sim \mathcal{N}(0, 1)$ and $h_\ell \sim \mathcal{N}(\xi_\ell \theta_\ell^t / \sqrt{q_\ell^t}, \rho_\ell - (\theta_\ell^t)^2 / q_\ell^t)$. The variance is given by

$$\mathbb{V}[b_i^t, b_j^t] = \delta_{ij} \sum_{\ell} (\Sigma_\ell)_{ii} \underbrace{\alpha \mathbb{E}_{H, \Xi} \text{pr}\ddot{\text{ox}}(\mathbb{T}[H], \rho_\Sigma, \{m_\kappa^t + \sqrt{q_\kappa^t} \xi_\kappa\}_\kappa, \{V_\kappa^t\}_\kappa)_\ell^2}_{\equiv \hat{q}_\ell^t} \quad (26)$$

Concentration of $A_{i \rightarrow \mu}^t$ Similarly to the derivation for $V_{\mu \rightarrow i}^t$, $A_{i \rightarrow \mu}^t$ concentrates to

$$A_{i \rightarrow \mu}^t = \sum_{\ell} -\alpha \frac{1}{V_\ell^t} \underbrace{\left(\mathbb{E}_{H, \Xi} \partial_{\omega_\ell} \text{pr}\ddot{\text{ox}}(\mathbb{T}[H], \rho_\Sigma, \{m_\kappa^t + \sqrt{q_\kappa^t} \xi_\kappa\}_\kappa, \{V_\kappa^t\}_\kappa)_\ell - 1 \right)}_{\equiv \hat{V}_\ell^t} (\Sigma_\ell)_{ii} \quad (27)$$

Recovering Result 1 Wrapping up, we now massage these equations to recover equations (8) from Result 1 of the main text. Starting from (18):

$$\begin{aligned} V_\ell^t &= \frac{1}{d} \sum_j (\Sigma_\ell)_{jj} \frac{1}{\lambda + \sum_{\kappa} \hat{V}_\kappa^{t-1}(\Sigma_\kappa)} \\ &= \int d\nu(\gamma, \tau) \gamma_\ell \left(\lambda + \sum_{\kappa} \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-1}. \end{aligned} \quad (28)$$

Next, for q_ℓ^t (21):

$$\begin{aligned} q_\ell^t &= \frac{1}{d} \sum_i (\Sigma_\ell)_{ii} \left(\left(\sum_\kappa (\sqrt{d}(p_\kappa)_i \hat{m}_\kappa^{t-1} + Q_i^*(\Sigma_\kappa)_{ii} \hat{\theta}_\kappa) \right)^2 + (\Sigma_\kappa)_{ii} \hat{q}_\kappa^{t-1} \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1}(\Sigma_\kappa) \right)^{-2} \\ &= \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\left(\sum_\kappa \hat{m}_\kappa^{t-1} \tau_\kappa + \hat{\theta}_\kappa \gamma_\kappa \pi_\kappa \right)^2 + \gamma_\kappa \hat{q}_\kappa^{t-1} \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-2} \end{aligned} \quad (29)$$

For θ_ℓ^t (21):

$$\begin{aligned} \theta_\ell^t &= \frac{1}{d} \sum_i (\Sigma_\ell)_{ii} Q_i^* \left(\sum_\kappa (\sqrt{d}(p_\kappa)_i \hat{m}_\kappa^{t-1} + Q_i^*(\Sigma_\kappa)_{ii} \hat{\theta}_\kappa) \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1}(\Sigma_\kappa) \right)^{-1} + o_d(1) \\ &= \int d\nu(\gamma, \tau, \pi) \gamma_\ell \pi_\ell \left(\sum_\kappa \hat{m}_\kappa^{t-1} \tau_\kappa + \hat{\theta}_\kappa \gamma_\kappa \pi_\kappa \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-1}. \end{aligned} \quad (30)$$

Finally for m_ℓ^t (20):

$$\begin{aligned} m_\ell^t &= \frac{1}{d} \sum_i (\sqrt{d} p_\ell)_i \left(\sum_\kappa (\sqrt{d}(p_\kappa)_i \hat{m}_\kappa^{t-1} + Q_i^*(\Sigma_\kappa)_{ii} \hat{\theta}_\kappa) \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1}(\Sigma_\kappa) \right)^{-1} + o_d(1) \\ &= \int d\nu(\gamma, \tau, \pi) \tau_\ell \left(\sum_\kappa \hat{m}_\kappa^{t-1} \tau_\kappa + \hat{\theta}_\kappa \gamma_\kappa \pi_\kappa \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-1}. \end{aligned} \quad (31)$$

For \hat{m}_ℓ^t (25):

$$\hat{m}_\ell^t = \alpha \mathbb{E}_{H, \Xi} \frac{1}{V_\ell^t} \left[\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell^t \right], \quad (32)$$

while for $\hat{\theta}_\ell^t$ (25):

$$\begin{aligned} \hat{\theta}_\ell^t &= \alpha \mathbb{E}_{H, \Xi} \frac{1}{V_\ell^t} \partial_{h_\ell} \left[\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell^t \right] \\ &= \alpha \mathbb{E}_{H, \Xi} \frac{1}{V_\ell^t} \frac{h_\ell - \theta_\ell^t / \sqrt{q_\ell^t} \xi_\ell}{\rho_\ell - (\theta_\ell^t)^2 / q_\ell^t} \left[\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell^t \right]. \end{aligned} \quad (33)$$

Now turning to \hat{q}_ℓ^t :

$$\hat{q}_\ell^t = \alpha \mathbb{E}_{H, \Xi} \left[\left(\frac{1}{V_\ell^t} \text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell^t \right)^2 \right]. \quad (34)$$

Finally, for \hat{V}_ℓ^t (27):

$$\begin{aligned}
 \hat{V}_\ell^t &= -\alpha \mathbb{E}_{H,\Xi} \frac{1}{V_\ell^t} [\partial_{\omega_\ell} \text{prox}_\ell - 1] \\
 &= -\alpha \mathbb{E}_{H,\Xi} \frac{1}{V_\ell^t} \left[\frac{1}{\sqrt{q_\ell^t}} \partial_\xi (\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell) \right] \\
 &= \alpha \mathbb{E}_{H,\Xi} \frac{1}{\sqrt{q_\ell^t} V_\ell^t} \left[\frac{\theta_\ell^t}{\sqrt{q_\ell^t} V_\ell^t} \left(\frac{h_\ell - \sqrt{q_\ell^t} \xi_\ell}{\rho_\ell - (\theta_\ell^t)^2 / q_\ell^t} - \xi \right) (\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell) \right] \\
 &= \frac{\theta_\ell^t \hat{\theta}_\ell^t}{q_\ell^t} - \alpha \mathbb{E}_{H,\Xi} \frac{1}{\sqrt{q_\ell^t} V_\ell^t} (\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell) \xi_\ell
 \end{aligned} \tag{35}$$

Summary : State evolution equations The state evolution equations asymptotically describing the dynamics of the GAMP algorithm 1 thus read

$$\begin{cases}
 V_\ell^t = \int d\nu(\gamma, \tau) \gamma_\ell \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-1} \\
 q_\ell^t = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\left(\sum_\kappa \hat{m}_\kappa^{t-1} \tau_\kappa + \hat{\theta}_\kappa \gamma_\kappa \pi_\kappa \right)^2 + \gamma_\kappa \hat{q}_\kappa^{t-1} \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-2} \\
 \theta_\ell^t = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \pi_\ell \left(\sum_\kappa \hat{m}_\kappa^{t-1} \tau_\kappa + \hat{\theta}_\kappa \gamma_\kappa \pi_\kappa \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-1} \\
 m_\ell^t = \int d\nu(\gamma, \tau, \pi) \tau_\ell \left(\sum_\kappa \hat{m}_\kappa^{t-1} \tau_\kappa + \hat{\theta}_\kappa \gamma_\kappa \pi_\kappa \right) \left(\lambda + \sum_\kappa \hat{V}_\kappa^{t-1} \gamma_\kappa \right)^{-1} . \\
 \hat{V}_\ell^t = \frac{\theta_\ell^t \hat{\theta}_\ell^t}{q_\ell^t} - \alpha \mathbb{E}_{H,\Xi} \frac{1}{\sqrt{q_\ell^t} V_\ell^t} (\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell) \xi_\ell \\
 \hat{q}_\ell^t = \alpha \mathbb{E}_{H,\Xi} \left[\left(\frac{1}{V_\ell^t} \text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell^t \right)^2 \right] \\
 \hat{\theta}_\ell^t = \alpha \mathbb{E}_{H,\Xi} \frac{1}{V_\ell^t} \frac{h_\ell - \theta_\ell^t / \sqrt{q_\ell^t} \xi_\ell}{\rho_\ell - (\theta_\ell^t)^2 / q_\ell^t} \left[\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell^t \right] \\
 \hat{m}_\ell^t = \alpha \mathbb{E}_{H,\Xi} \frac{1}{V_\ell^t} \left[\text{prox}_\ell - \sqrt{q_\ell^t} \xi_\ell - m_\ell^t \right]
 \end{cases} \tag{36}$$

which exactly recovers equations (8) of Result 1 of the main text, for the case $r_s = 1$ considered in the present Appendix. Again, we mention that the case $r_s \geq 2$ should follow straightforwardly with the exact same derivation steps, using tensor variables (see e.g. [13]). This subsection has thus established that the equations (8) (with time indices) describe the summary statistics capturing the dynamics of GAMP iterations 1. In particular, (8) describe the fixed points of GAMP. The next subsection further shows that the (stable) fixed points of GAMP correspond to critical (zero-gradient) points of the empirical landscape (3), i.e. fixed points of gradient descent. Finally, we provide in Appendix B an alternative derivation of the state evolution equations (8)(36), using the replica method from statistical physics [40, 41].

A.5. Fixed points of GAMP are fixed points of GD

In this subsection, we show that fixed points of GAMP 1, as asymptotically described by (8) in Result 1, correspond to critical (zero gradient) points of the empirical landscape (3). Again, we present the

result for $r_s = 1$ for clarity, the generalization to $r_s \geq 2$ being straightforward (see e.g. [13]). In the previous notations, let us denote the (simplified, see (14)) empirical loss as

$$L(\{\tilde{X}_\ell Q\}_\ell) + g(Q) \quad (38)$$

where we introduced the shorthands

$$L(\{h_\ell\} \in \mathbb{R}^n) \equiv \sum_{\mu=1}^n \left(-2 \operatorname{Tr} \left[\mathbf{S}[\{h_\ell^\mu\}_\ell] \rho_\mu T_\mu^\top \right] + \operatorname{Tr} \left[\mathbf{S}[\{h_\ell^\mu\}_\ell] \rho_\mu \mathbf{S}[\{h_\ell^\mu\}_\ell]^\top \right] \right) \quad (39)$$

$$g(Q) \equiv \frac{\lambda}{2} \|Q\|^2, \quad (40)$$

i.e. respectively the simplified empirical loss (3) and the regularization, as functions with matrix arguments. The empirical minimization problem (3) can thus be written compactly as

$$\hat{Q} = \operatorname{argmin}_{Q \in \mathbb{R}^d} \left\{ L(\{\tilde{X}_\ell Q\}_\ell) + r(Q) \right\} \quad (41)$$

with the critical (zero-gradient) condition being given by

$$\sum_{\ell=1}^L \tilde{X}_\ell^\top \partial_\ell L(\{\tilde{X}_\ell Q\}_\ell) + \partial g(Q) \stackrel{!}{=} 0. \quad (42)$$

Let us choose a diagonal definite $A \in \mathbb{R}^{d \times d}$, and a sequence $\{V_\mu\}_{1 \leq \mu \leq n}$ of symmetric definite $L \times L$ matrices. Group them into a block diagonal matrix $\check{V} \in \mathbb{R}^{Ln \times Ln}$, so that the μ -th block of \check{V} corresponds to V_μ . It shall prove useful to further introduce the matrices $\check{X} \in \mathbb{R}^{d \times nL}$ (resp. $\check{\partial} L(\check{X}) \in \mathbb{R}^{nL}$), defined as the concatenation of the matrices $\check{X}_1, \dots, \check{X}_L$ (resp. $\partial_1 L, \dots, \partial_L L$), viewed as n blocks of length L . Then without loss of generality the zero-gradient condition can be rewritten as

$$\check{X}^\top V^{-1} \left(V \check{\partial} L(\check{X}) - \check{X} Q \right) + A(A^{-1} \partial g(Q) + Q) \stackrel{!}{=} \check{X}^\top V^{-1} \check{X} Q + A Q. \quad (43)$$

Similarly to [13], let us introduce

$$\check{\omega} \equiv V \check{\partial} L(\check{X}) - \check{X} Q. \quad (44)$$

This can be written in terms of a resolvent as

$$\check{X} Q = \operatorname{prox}(\check{\omega}) \quad (45)$$

where

$$\operatorname{prox}(\check{\omega}) \in \mathbb{R}^{nL} = \operatorname{argmin}_{\check{x} \in \mathbb{R}^{nL}} \left\{ \frac{1}{2} \|\check{x} - \check{\omega}\|_V^2 + L(\check{x}) \right\} \quad (46)$$

which corresponds to (17). Similarly, we denote

$$b \equiv A^{-1} \partial g(Q) + Q \quad (47)$$

So that

$$Q = \text{prox}_g(b) = \underset{x \in \mathbb{R}^d}{\text{argmin}} \left\{ \frac{1}{2} \|x - b\|_{A^{-1}}^2 + g(x) \right\} \quad (48)$$

In the particular case of an ℓ_2 regularization $g(\cdot) = \lambda/2 \|\cdot\|^2$, note that

$$\text{prox}_g(b) = (\lambda \mathbb{I}_d + A)^{-1} Ab. \quad (49)$$

The zero-gradient condition can now be rewritten as

$$\begin{cases} \tilde{X}^\top V^{-1} (\text{prox}(\tilde{\omega}) - \tilde{\omega}) = A(b - \text{prox}_g(b)) \\ \tilde{X} \text{prox}_g(b) = \text{prox}(\tilde{\omega}) \end{cases} \quad (50)$$

One is now in a position to expand the concatenated variables $\tilde{\cdot}$ into a sequence of L n -dimensional parameters. For $u = \text{prox}(\tilde{\omega})$, $\tilde{\omega}$ let us denote $u_{\mu\ell}$ ($1 \leq \mu \leq n, 1 \leq \ell \leq L$) the ℓ -th component of the μ -th block. Introduce

$$f_{\mu\ell} \equiv \sum_{\kappa} (V_{\mu}^{-1})_{\ell\kappa} (\text{prox}(\tilde{\omega})_{\mu\kappa} - \tilde{\omega}_{\mu\kappa}). \quad (51)$$

Denote $f_{\ell} \equiv (f_{\mu\ell})_{1 \leq \mu \leq n} \in \mathbb{R}^n$, $\omega_{\ell} \equiv (\omega_{\mu\ell})_{1 \leq \mu \leq n} \in \mathbb{R}^n$. The system of equations (50) can then be rewritten as (further redefining $b \leftarrow Ab$):

$$\begin{cases} \sum_{\ell} \tilde{X}_{\ell}^\top f_{\ell} = b - A(\lambda \mathbb{I}_d + A)^{-1} b \\ \tilde{X}_{\ell} (\lambda \mathbb{I}_d + A)^{-1} b = \sum_{\kappa} V_{\ell\kappa} f_{\kappa} - \omega_{\ell} \end{cases} \quad (52)$$

We used the assumption that $g(\cdot)$ is an ℓ_2 regularization. Finally, introducing $\hat{Q} = \text{prox}_g(A^{-1}b) = (\lambda \mathbb{I}_d + A)^{-1} b$, one reaches

$$\begin{cases} \sum_{\ell} \tilde{X}_{\ell}^\top f_{\ell} = b - A\hat{Q} \\ \tilde{X}_{\ell} \hat{Q} = \sum_{\kappa} V_{\ell\kappa} f_{\kappa} - \omega_{\ell} \end{cases} \quad (53)$$

which corresponds to the fixed-point equations of GAMP (Algorithm 1). This finishes to show the correspondence between the fixed points of GAMP and the critical points of the empirical landscape (3). To summarize, we have shown that equations (8) describe the zero-gradient points of the empirical loss landscape (3), i.e. fixed points of GD.

A.6. Towards a rigorous proof of result 1

While the connection between the GAMP fixed point and the extrema of the loss is sound, and has been at the roots of many rigorous results for convex losses, see e.g. [7, 15, 17, 19, 31], there exist technical difficulties in adapting these rigorous arguments to the present setting, and a fully rigorous proof would warrant sizable work. While we leave this challenging task for future work, we wish to discuss how it can be potentially achieved. The first task would require the proof of point-wise

convergence of GAMP, as indeed, the identification of the GAMP estimates with the one of the extrema of the loss function requires to be at the fixed point of the iteration. This difficulty, discussed in detail in, e.g. [10, 19, 31], can be in principle addressed by computing the convergence criterion from the state evolution equations (see [10] the discussion in Lemma 7 in [19]), a criterion sometimes called the "replicon" in the context of replica theory [35].

Provided the replicon criterion is satisfied, all converging fixed point described by our theory thus correspond rigorously to fixed point of the loss. The last task would be to prove that the minimum of the loss is indeed the fixed point we found with minimum energy. A potential strategy to prove this would be to use the Gordon-Minimax approach of [20]. While it is used in many situations for convex problems (e.g. [54, 55, 58]), only one side would be required for our (non-convex) problem thanks to the GAMP matching bound. We hope that our results would provide inspiration for further research in this direction.

Appendix B. Derivation of the Main Result with the replica method

In the Appendix we provide an alternative derivation of Result 1, which sharply characterizes the global minimum of the empirical loss (3), using the heuristic replica method from statistical physics [40, 41] in its replica-symmetric formulation. First observe that for any test function $\phi(\hat{Q})$ of the minimizer \hat{Q} of (3),

$$\phi(\hat{Q}) = \lim_{\beta \rightarrow \infty} \mathbb{E}_{\mathcal{D}} \frac{1}{Z} \int dQ \phi(Q) e^{-\beta \mathcal{R}[Q]}, \quad (54)$$

where we denoted $R[Q]$ the empirical loss (3), and

$$Z \equiv \int dQ e^{-\beta \mathcal{R}[Q]} \quad (55)$$

the normalization factor, also known as the *partition function* in statistical physics. We remind that \mathcal{D} refers to the training set. In order to access key summary statistics and learning metrics associated to \hat{Q} , it is therefore reasonable to seek to compute the generating function associated to the measure (54), namely $\mathbb{E} \ln Z$. Such computations can be addressed using the *replica* method from statistical physics [40, 41], building on the identity

$$\ln Z = \lim_{s \rightarrow 0} \frac{Z^s - 1}{s}. \quad (56)$$

The backbone of the derivation thus lies in the computation of $\mathbb{E} Z^s$. Below, we detail the derivation for a generic convex regularizer $g : \mathbb{R}^d \rightarrow \mathbb{R}_+$ and later specialize to the case of ℓ_2 regularization. The replicated partition function thus reads

$$\begin{aligned} \mathbb{E} Z^s &= \int \prod_{a=1}^s dQ_a e^{-\beta \sum_{a=1}^s g(Q_a)} \\ &\quad \prod_{\mu=1}^n \mathbb{E}_x e^{-\beta \sum_{a=1}^s \left(\text{Tr} \mathbb{S} \left[\frac{1}{\sqrt{d}} (x+p) Q_a \right] \rho_{\Sigma} \mathbb{S} \left[\frac{1}{\sqrt{d}} (x+p) Q_a \right]^{\top} - 2 \text{Tr} \mathbb{T} \left[\frac{1}{\sqrt{d}} x_{\ell} Q_{\star} \right] \rho_{\Sigma} \mathbb{S} \left[\frac{1}{\sqrt{d}} (x+p) Q_a \right]^{\top} \right)}. \end{aligned} \quad (57)$$

Introduce the local fields

$$h^a \equiv \frac{xQ_a}{\sqrt{d}} \in \mathbb{R}^{L \times r}, \quad h^* \equiv \frac{xQ_\star}{\sqrt{d}} \in \mathbb{R}^{L \times t} \quad (58)$$

and the overlaps

$$m_a \equiv \frac{pQ_a}{\sqrt{d}} \in \mathbb{R}^{L \times r}, \quad (59)$$

with rows m_a^ℓ . These fields have statistics

$$\mathbb{E}_x[h_\ell^a(h_\kappa^b)^\top] = \delta_{\ell\kappa} \frac{Q_a^\top \Sigma_\ell Q_b}{d} \equiv q_{ab}^\ell \quad (60)$$

$$\mathbb{E}_x[h_\ell^*(h_\kappa^*)^\top] = \delta_{\ell\kappa} \frac{Q_\star^\top \Sigma_\ell Q_\star}{d} \equiv \rho_\ell \quad (61)$$

$$\mathbb{E}_x[h_\ell^a(h_\kappa^*)^\top] = \delta_{\ell\kappa} \frac{Q_a^\top \Sigma_\ell Q_\star}{d} \equiv \theta_a^\ell. \quad (62)$$

Thus

$$\begin{aligned} \mathbb{E}Z^s &= \int dm d\hat{m} d\theta d\hat{\theta} dq d\hat{q} e^{\underbrace{-d \sum_a \sum_\ell [\hat{m}_a^{\ell\top} m_a^\ell + \text{Tr}(\theta_a^\ell \hat{\theta}_a^{\ell\top})] - d \sum_\ell \sum_{1 \leq a \leq b \leq s} \text{Tr}(q_{ab}^\ell \hat{q}_{ab}^{\ell\top})}_{e^{sd\Psi_t}}} \\ &\quad \underbrace{\int \prod_{a=1}^s dQ_a e^{-\beta \sum_{a=1}^s g(Q_a) + \sum_a \sum_\ell (\sqrt{d} \hat{m}_a^{\ell\top} Q_a^\top p_\ell + \text{Tr}[\theta_a^\ell Q_\star^\top \Sigma_\ell Q_a]) + \sum_{1 \leq a \leq b \leq s} \sum_\ell \text{Tr}[q_{ab}^\ell Q_b^\top \Sigma_\ell Q_a]}}_{e^{sd\Psi_Q}} \\ &\quad \underbrace{\left[\mathbb{E}_{h^*, \{h_a\}_{a=1}^s} e^{-\beta \sum_{a=1}^s (\text{Tr} S[h^a + m^a] \rho_\Sigma S[h^a + m^a]^\top - 2 \text{Tr} T[h^*] \rho_\Sigma S[h^a + m^a]^\top)} \right]^{\alpha d}}_{e^{s\alpha d\Psi_y}}, \quad (63) \end{aligned}$$

where we decomposed the replicated free entropy into the trace, entropic and energetic potentials Ψ_t, Ψ_Q, Ψ_y . Note that all exponents are scaling with $d \rightarrow \infty$. Therefore the integral in (63) can be computed using a Laplace saddle-point approximation.

B.1. Replica-Symmetric ansatz

We have thus rephrased the analysis of the measure (54) as a optimization problem over the order parameters $\{q_{ab}^\ell, \theta_a^\ell, m_a\}$, and the associated conjugate variables. However, these still represent $2L(s^2 + 1) + 2s$ variables, and $s \rightarrow 0$. In order to make progress, we assume that the maximizer is of *replica-symmetric* (RS) form [40, 41]

$$q_{ab}^\ell = (r_\ell - q_\ell) \delta_{ab} + q_\ell \quad (64)$$

$$m_a^\ell = m_\ell \quad (65)$$

$$\theta_a^\ell = \theta_\ell \quad (66)$$

$$\hat{q}_{ab}^\ell = -(\hat{r}_\ell/2 + \hat{q}_\ell) + \hat{q}_\ell \quad (67)$$

$$\hat{m}_a^\ell = \hat{m}_\ell \quad (68)$$

$$\hat{\theta}_a^\ell = \hat{\theta}_\ell \quad (69)$$

The RS ansatz holds in a number of machine learning settings, notably for convex problems and Bayes-optimal settings, see e.g. [63] for a review. In the present setting, since the empirical loss (3) is non-convex, we emphasize that the RS ansatz constitutes a heuristic technical assumption of our analysis.

B.2. Entropic potential

We now turn to the entropic potential Ψ_w . It is convenient to introduce the variance order parameter

$$\hat{V}_\ell \equiv \hat{r}_\ell + \hat{q}_\ell. \quad (70)$$

The entropic potential can then be expressed as

$$\begin{aligned} & e^{\beta s d \Psi_Q} \\ &= \int \prod_{a=1}^s dQ_a e^{-\beta \sum_a g(Q_a) + \sum_{\ell=1}^L \sum_{a=1}^s (\sqrt{d} \hat{m}_\ell^\top Q_a^\top p_\ell + \text{Tr}[\hat{Q}_\star^\top \Sigma_\ell Q_a]) - \frac{1}{2} \sum_{\ell=1}^L \sum_{a=1}^s \text{Tr}[\hat{V}_\ell Q_a \Sigma_\ell Q_a^\top] + \frac{1}{2} \sum_{\ell=1}^L \sum_{a,b} \text{Tr}[\hat{q}_\ell Q_a \Sigma_\ell Q_b^\top]} \\ &= \int \prod_{\ell=1}^L D\Xi_\ell \\ & \quad \left[\int dQ e^{-\beta g(Q) - \frac{1}{2} \text{Tr} \left[\sum_{\ell=1}^L \hat{V}_\ell Q \Sigma_\ell Q^\top \right] + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell) + \sum_{\ell=1}^L \Xi_\ell \odot (\hat{q}_\ell \otimes \Sigma_\ell)^{\frac{1}{2}} \right) \odot Q} \right]^s \\ &= \mathbb{E}_\Xi \left[\int dQ e^{-\beta g(Q) - \frac{1}{2} Q \odot \left[\sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right] \odot Q + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell) + \sum_{\ell=1}^L \Xi_\ell \odot (\hat{q}_\ell \otimes \Sigma_\ell)^{\frac{1}{2}} \right) \odot Q} \right]^s. \quad (71) \end{aligned}$$

Therefore

$$\beta \Psi_w = \frac{1}{d} \int \mathbb{E}_\Xi \ln \left[\int dQ e^{-\beta g(Q) - \frac{1}{2} Q \odot \left[\sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right] \odot Q + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell) + \sum_{\ell=1}^L \Xi_\ell \odot (\hat{q}_\ell \otimes \Sigma_\ell)^{\frac{1}{2}} \right) \odot Q} \right]. \quad (72)$$

For a matrix $\Xi \in \mathbb{R}^{r \times d}$ and tensors $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{r \times d} \otimes \mathbb{R}^{r \times d}$, we denoted $(\Xi \odot \mathbf{A})_{kl} = \sum_{ij} \Xi^{ij} \mathbf{A}_{ij,kl}$ and $(\mathbf{A} \odot \mathbf{B})_{ij,kl} = \sum_{rs} \mathbf{A}_{ij,rs} \mathbf{B}_{rs,kl}$.

B.3. Energetic potential

The computation of the energetic potential Ψ_y is rather standard and follows the same lines as in e.g. [2], yielding

$$\begin{aligned}
 \beta\Psi_y &= \int_{\mathbb{R}^{L \times t}} dY DZ \int_{\mathbb{R}^{L \times r}} D\Xi \prod_{\ell=1}^L \delta \left[y_\ell - (\rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell)^{\frac{1}{2}} z_\ell - \theta_\ell^\top q_\ell^{\frac{1}{2}} \xi_\ell \right] \\
 &\quad \times \ln \left[\int_{\mathbb{R}^{L \times r}} dX \prod_{\ell=1}^L \frac{e^{-\frac{1}{2} \left(x_\ell - q_\ell^{\frac{1}{2}} \xi_\ell \right)^\top V_\ell^{-1} \left(x_\ell - q_\ell^{\frac{1}{2}} \xi_\ell \right)}}{\det(2\pi V_\ell)} e^{-\beta \operatorname{Tr} \mathbf{S}[x+m] \rho_\sigma \mathbf{S}[x+m]^\top - 2 \operatorname{Tr} \mathbf{T}[y] \rho_\Sigma \mathbf{S}[x+m]^\top} \right] \\
 &= \underbrace{\int_{\mathbb{R}^{L \times t}} dY \int_{\mathbb{R}^{L \times r}} D\Xi \prod_{\ell=1}^L \frac{e^{-\frac{1}{2} \left(y_\ell - \theta_\ell^\top q_\ell^{\frac{1}{2}} \xi_\ell \right)^\top (\rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell)^{-1} \left(y_\ell - \theta_\ell^\top q_\ell^{\frac{1}{2}} \xi_\ell \right)}}{\det[2\pi(\rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell)]}}_{\equiv \mathbb{E}_{Y, \Xi}} \\
 &\quad \times \ln \left[\int_{\mathbb{R}^{L \times r}} dX \prod_{\ell=1}^L \frac{e^{-\frac{1}{2} \left(x_\ell - q_\ell^{\frac{1}{2}} \xi_\ell \right)^\top V_\ell^{-1} \left(x_\ell - q_\ell^{\frac{1}{2}} \xi_\ell \right)}}{\det(2\pi V_\ell)} e^{-\beta \operatorname{Tr} \mathbf{S}[x+m] \rho_\sigma \mathbf{S}[x+m]^\top - 2 \operatorname{Tr} \mathbf{T}[y] \rho_\Sigma \mathbf{S}[x+m]^\top} \right]
 \end{aligned} \tag{73}$$

B.4. Zero-temperature limit

We now take the limit $\beta \rightarrow \infty$. Rescaling

$$\beta \hat{V}_\ell \leftarrow \hat{V}_\ell, \quad \frac{1}{\beta} V_\ell \leftarrow V_\ell, \quad \beta \hat{m}_\ell \leftarrow \hat{m}_\ell, \quad \beta \hat{\theta}_\ell \leftarrow \hat{\theta}_\ell, \quad \beta^2 \hat{q}_\ell \leftarrow \hat{q}_\ell \tag{74}$$

The entropic potential then reduces to

$$\begin{aligned}
 \Psi_w &= \frac{1}{2d} \mathbb{E}_\Xi \operatorname{Tr} \left[\left(\sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right) \odot \left(\sum_{\ell=1}^L \left(\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell \right) + \sum_{\ell=1}^L \Xi_\ell \odot (\hat{q}_\ell \otimes \Sigma_\ell)^{\frac{1}{2}} \right) \right]^{\otimes 2} \\
 &\quad - \frac{1}{d} \mathbb{E}_\Xi \mathcal{M}_g(\Xi)
 \end{aligned} \tag{75}$$

where we defined the entropic Moreau envelope

$$M_g(\Xi) \equiv \inf_Q \left\{ \frac{1}{2} \left\| \left(\sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right)^{1/2} \left(Q - \left(\sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right)^{-1} \left(\sum_{\ell=1}^L \left(\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell \right) + \sum_{\ell=1}^L \Xi_\ell \odot (\hat{q}_\ell \otimes \Sigma_\ell)^{\frac{1}{2}} \right) \right\|^2 + g(Q) \right\}. \tag{76}$$

The energetic potential can be similarly recast into a more compact form

$$\Psi_y = -\mathbb{E}_{Y, \Xi} \mathcal{M}(Y, \Xi) \tag{77}$$

where the Moreau envelope is defined as

$$\mathcal{M}(Y, \Xi) = \inf_X \frac{1}{2} \left\{ \sum_{\ell=1}^L \text{Tr} \left[V_\ell^{-1} \left(x_\ell - q_\ell^{1/2} \xi_\ell - m_\ell \right)^{\otimes 2} \right] + \text{Tr} \left[\mathbf{S}(X) \rho_\Sigma \mathbf{S}(X)^\top \right] - 2 \text{Tr} \left[\mathbf{T}(Y) \rho_\Sigma \mathbf{S}(X)^\top \right] \right\}. \quad (78)$$

B.5. Replica free entropy

One finally reaches an expression for the replica free entropy as

$$\begin{aligned} \Phi = & \frac{1}{2} \sum_{\ell=1}^L \left(\text{Tr} \hat{V}_\ell q_\ell - \text{Tr} \hat{q}_\ell V_\ell \right) - \sum_{\ell=1}^L \hat{m}_\ell^\top m_\ell - \sum_{\ell=1}^L \text{Tr} \hat{\theta}_\ell^\top \theta_\ell - \frac{1}{d} \mathbb{E}_\Xi \mathcal{M}_g(\Xi) \\ & + \frac{1}{2d} \mathbb{E}_\Xi \text{Tr} \left[\left(\sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right) \odot \left(\left(\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell \right) + \sum_{\ell=1}^L \Xi_\ell \odot \left(\hat{q}_\ell \otimes \Sigma_\ell \right)^{\frac{1}{2}} \right)^{\otimes 2} \right] - \alpha \mathbb{E}_{y, \xi} \mathcal{M}(y, \xi) \end{aligned} \quad (79)$$

B.6. Saddle-point equations : general regularizer

The extremization of the free entropy (79) yields, similarly to [14], the following system of self-consistent equations on the summary statistics:

$$\begin{cases} V_\ell = \frac{1}{d} \mathbb{E}_\Xi \left[\left(\text{prox}_g \odot \left(\hat{q}_\ell \otimes \Sigma_\ell \right)^{-\frac{1}{2}} \odot \left(\mathbb{I}_r \otimes \Sigma_\ell \right) \right) \Xi_\ell^\top \right] \\ q_\ell = \frac{1}{d} \mathbb{E}_\Xi \left[\text{prox}_g \Sigma_\ell \text{prox}_g^\top \right] \\ m_\ell = \frac{1}{\sqrt{d}} \mathbb{E}_\Xi \left[\text{prox}_g p_\ell \right] \\ \theta_\ell = \frac{1}{\sqrt{d}} \mathbb{E}_\Xi \left[\text{prox}_g \Sigma_\ell Q_\star \right] \\ \hat{q}_\ell = \alpha \mathbb{E}_{\Xi, Y} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right)^{\otimes 2} V_\ell^{-1} \\ \hat{V}_\ell = q_\ell^{-1} \hat{\theta}_\ell \theta_\ell^\top - \alpha q_\ell^{-\frac{1}{2}} \mathbb{E}_{\Xi, Y} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \xi_\ell^\top \\ \hat{m}_\ell = \alpha \mathbb{E}_{\xi, \eta} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \\ \hat{\theta}_\ell = \alpha \mathbb{E}_{\xi, \eta} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \left(y_\ell - \xi_\ell^\top q_\ell^{-1/2} \theta_\ell \right)^\top \left(\rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell \right)^{-1} \end{cases}, \quad (80)$$

$$\begin{cases} \hat{q}_\ell = \alpha \mathbb{E}_{\Xi, Y} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right)^{\otimes 2} V_\ell^{-1} \\ \hat{V}_\ell = q_\ell^{-1} \hat{\theta}_\ell \theta_\ell^\top - \alpha q_\ell^{-\frac{1}{2}} \mathbb{E}_{\Xi, Y} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \xi_\ell^\top \\ \hat{m}_\ell = \alpha \mathbb{E}_{\xi, \eta} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \\ \hat{\theta}_\ell = \alpha \mathbb{E}_{\xi, \eta} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \left(y_\ell - \xi_\ell^\top q_\ell^{-1/2} \theta_\ell \right)^\top \left(\rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell \right)^{-1} \end{cases}, \quad (81)$$

where the proximals prox_g and prox_ℓ respectively refer to the arginf in Q (resp. x_ℓ) of the envelopes \mathcal{M}_g (76) (resp. \mathcal{M} 78).

B.7. Saddle-point equations : ℓ_2

We now specialize the saddle-point equations (80) to the case of an ℓ_2 regularizer $g(\cdot) = 1/2\|\cdot\|$ the entropic potential admits the simple form

$$\begin{aligned} \Psi_Q &= \frac{1}{2d} \text{Tr} \left[\left(\lambda \mathbb{I}_r \odot \mathbb{I}_d + \sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right)^{-1} \odot \left(\sum_{\ell=1}^L \hat{q}_\ell \otimes \Sigma_\ell + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell) \right)^{\otimes 2} \right) \right] \\ &= \frac{1}{2d} \sum_{i=1}^d \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \lambda_i^\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \lambda_i^\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_\ell p_\ell^\top e_i + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell e_i) \right) \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_\ell p_\ell^\top e_i + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell e_i) \right)^\top \right) \right] \\ &\stackrel{d \rightarrow \infty}{=} \frac{1}{2} \int d\nu(\gamma, \tau, \pi) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \gamma_\ell \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right]. \end{aligned} \quad (82)$$

The replica free energy thus reads

$$\begin{aligned} \Phi &= \frac{1}{2} \sum_{\ell=1}^L \left(\text{Tr} \hat{V}_\ell q_\ell - \text{Tr} \hat{q}_\ell V_\ell \right) - \sum_{\ell=1}^L \hat{m}_\ell^\top m_\ell - \sum_{\ell=1}^L \text{Tr} \hat{\theta}_\ell^\top \theta_\ell - \alpha \mathbb{E}_{y, \xi} \mathcal{M}(y, \xi) \\ &\quad + \frac{1}{2} \int d\nu(\gamma, \tau, \pi) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \gamma_\ell \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right], \end{aligned} \quad (83)$$

leading to the saddle point equations

$$\begin{cases} \hat{q}_\ell = \alpha \mathbb{E}_{\Xi, Y} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right)^{\otimes 2} V_\ell^{-1} \\ \hat{V}_\ell = \hat{\theta}_\ell \theta_\ell^\top q_\ell^{-1} - \alpha \mathbb{E}_{\Xi, Y} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \xi_\ell^\top q_\ell^{-\frac{1}{2}} \\ \hat{m}_\ell = \alpha \mathbb{E}_{\xi, \eta} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \\ \hat{\theta}_\ell = \alpha \mathbb{E}_{\xi, \eta} V_\ell^{-1} \left(\text{prox}_\ell - q_\ell^{\frac{1}{2}} \xi_\ell - m_\ell \right) \left(y_\ell - \xi_\ell^\top q_\ell^{-1/2} \theta_\ell \right)^\top \left(\rho_\ell - \theta_\ell^\top q_\ell^{-1} \theta_\ell \right)^{-1} \\ q_\ell = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \left(\sum_{\kappa=1}^L \gamma_\kappa \hat{q}_\kappa + \left(\sum_{\kappa=1}^L \hat{m}_\kappa \tau_\kappa + \gamma_\kappa \hat{\theta}_\kappa \cdot \pi \right)^{\otimes 2} \right) \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \\ V_\ell = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \\ m_\ell = \int d\nu(\gamma, \tau, \pi) \tau_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \left(\sum_{\kappa=1}^L \hat{m}_\kappa \tau_\kappa + \gamma_\kappa \hat{\theta}_\kappa \cdot \pi \right) \\ \theta_\ell = \int d\nu(\gamma, \tau, \pi) \gamma_\ell \left(\lambda \mathbb{I}_r + \sum_{\kappa=1}^L \gamma_\kappa \hat{V}_\kappa \right)^{-1} \left(\sum_{\kappa=1}^L \hat{m}_\kappa \tau_\kappa + \gamma_\kappa \hat{\theta}_\kappa \cdot \pi \right) \pi^\top. \end{cases} \quad (84)$$

which finishes to recover (8). Let us finally mention that the update equations (8) for the summary statistics (7) do *not* describe the dynamics of gradient descent, but rather that of an Approximate Message Passing algorithm [6], which we elicit in Appendix A for completeness. q.e.d.

B.8. test MSE

The generalization performance is measured by the test error

$$\epsilon_g \equiv \mathbb{E}_{\mathcal{D}} \mathbb{E}_x \left\| \mathbf{T} \left[\frac{1}{\sqrt{d}} x Q_{\star} \right] x - \mathbf{S} \left[\frac{1}{\sqrt{d}} (x+p) \hat{Q} \right] (x+p) \right\|^2. \quad (85)$$

Expliciting this expression in terms of the correlated Gaussian variables xQ_{\star}, xQ allows to straightforwardly show that ϵ_g admits the sharp asymptotic characterization in terms of the summary statistics characterized by (84):

$$\epsilon_g = \mathbb{E}_X \text{Tr} \left[\mathbf{S}[X] \rho_{\Sigma} \mathbf{S}[X]^{\top} \right] + \mathbb{E}_Y \text{Tr} \left[\mathbf{T}[Y] \rho_{\Sigma} \mathbf{T}[Y]^{\top} \right] - 2 \mathbb{E}_{X,Y} \text{Tr} \left[\mathbf{S}[X] \rho_{\Sigma} \mathbf{T}[Y]^{\top} \right], \quad (86)$$

where the average bears on $X \in \mathbb{R}^{L \times r}, Y \in \mathbb{R}^{L \times t}$ with independent rows with statistics

$$(x_{\ell}, y_{\ell}) \sim \mathcal{N} \left[\begin{pmatrix} m \\ 0 \end{pmatrix}, \begin{pmatrix} q_{\ell} & \theta_{\ell} \\ \theta_{\ell}^{\top} & \rho_{\ell} \end{pmatrix} \right] \quad (87)$$

B.9. Training loss

We finally turn to the training loss. It is reasonable to expect, from statistical physics, that the training loss should be equal to the free energy $-\Phi$ at zero temperature. We provide below an alternative derivation, for simplicity in the case of ℓ_2 regularization $g = 1/2 \|\cdot\|^2$. First note that the training loss ϵ_t can be expressed as

$$\epsilon_t = - \lim_{\beta \rightarrow \infty} \partial_{\beta} \underbrace{\frac{1}{d} \ln Z(\beta)}_{\Phi(\beta)} \quad (88)$$

Where $\Phi(\beta)$ is the free entropy at finite temperature. The trace potential Ψ_t bears no explicit dependence on β . On the other hand,

$$\begin{aligned} \beta \Psi_Q &= -\frac{1}{2} \ln \det \left[\beta \lambda \mathbf{I}_r \otimes \mathbf{I}_d + \sum_{\ell=1}^L \hat{V}_{\ell} \otimes \Sigma_{\ell} \right] \\ &+ \frac{1}{2d} \text{Tr} \left[\left(\beta \lambda \mathbf{I}_r \odot \mathbf{I}_d + \sum_{\ell=1}^L \hat{V}_{\ell} \otimes \Sigma_{\ell} \right)^{-1} \odot \left(\sum_{\ell=1}^L \hat{q}_{\ell} \otimes \Sigma_{\ell} + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_{\ell} p_{\ell}^{\top} + \hat{\theta}_{\ell} Q_{\star}^{\top} \Sigma_{\ell}) \right)^{\otimes 2} \right) \right] \end{aligned} \quad (89)$$

Thus

$$\begin{aligned} \partial_{\beta}(\beta \Psi_Q) &= -\frac{\lambda}{2} \text{Tr} \left[\beta \lambda \mathbf{I}_r \otimes \mathbf{I}_d + \sum_{\ell=1}^L \hat{V}_{\ell} \otimes \Sigma_{\ell} \right]^{-1} \\ &- \frac{\lambda}{2d} \text{Tr} \left[\left(\beta \lambda \mathbf{I}_r \odot \mathbf{I}_d + \sum_{\ell=1}^L \hat{V}_{\ell} \otimes \Sigma_{\ell} \right)^{-2} \odot \left(\sum_{\ell=1}^L \hat{q}_{\ell} \otimes \Sigma_{\ell} + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_{\ell} p_{\ell}^{\top} + \hat{\theta}_{\ell} Q_{\star}^{\top} \Sigma_{\ell}) \right)^{\otimes 2} \right) \right] \end{aligned} \quad (90)$$

Finally, going through the same rescaling steps to take the $\beta \rightarrow \infty$ limit,

$$\lim_{\beta \rightarrow \infty} \partial_\beta (\beta \Psi_Q) = -\frac{\lambda}{2d} \text{Tr} \left[\left(\lambda \mathbb{I}_r \odot \mathbb{I}_d + \sum_{\ell=1}^L \hat{V}_\ell \otimes \Sigma_\ell \right)^{-2} \odot \left(\sum_{\ell=1}^L \hat{q}_\ell \otimes \Sigma_\ell + \left(\sum_{\ell=1}^L (\sqrt{d} \hat{m}_\ell p_\ell^\top + \hat{\theta}_\ell Q_\star^\top \Sigma_\ell) \right)^{\otimes 2} \right) \right] \quad (91)$$

By the same token, it is straightforward to see that

$$\lim_{\beta \rightarrow \infty} \partial_\beta (\beta \Psi_y) = -\mathbb{E}_{Y, \Xi} \left[\mathcal{M}(Y, \Xi) - \frac{1}{2} \sum_{\ell=1}^L \text{Tr} \left[\underbrace{V_\ell^{-1} \left(x_\ell - q_\ell^{1/2} \xi_\ell - m_\ell \right)^{\otimes 2}}_{\hat{q}_\ell V_\ell} \right] \right] \quad (92)$$

We used the self-consistent equations (8) to identify the term in underbrace. Putting everything together,

$$\begin{aligned} -\epsilon_t &= \lim_{\beta \rightarrow \infty} \partial_\beta \Psi(\beta) = -\frac{\lambda}{2} \int d\nu(\gamma, \tau) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right] \\ &\quad - \alpha \mathbb{E}_{Y, \Xi} [\mathcal{M}(Y, \Xi)] + \frac{1}{2} \sum_{\ell=1}^L \text{Tr}[\hat{q}_\ell V_\ell]. \end{aligned} \quad (93)$$

This constitutes a sharp asymptotic characterization of the training loss ϵ_t as a function of the summary statistics characterized in Result 1.

For completeness, we finally explicit the connection between ϵ_t and the negative free entropy (i.e. the *free energy* in statistical physics). We go back to massage the expression for the free entropy

$$\begin{aligned}
 \Phi &= \frac{1}{2} \sum_{\ell=1}^L \left(\text{Tr} \hat{V}_\ell q_\ell - \text{Tr} \hat{q}_\ell V_\ell \right) - \sum_{\ell=1}^L \hat{m}_\ell^\top m_\ell - \sum_{\ell=1}^L \text{Tr} \hat{\theta}_\ell^\top \theta_\ell - \alpha \mathbb{E}_{Y, \Xi} \mathcal{M}(Y, \Xi) \\
 &\quad + \frac{1}{2} \int d\nu(\gamma, \tau) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \gamma_\ell \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right] \\
 &= \frac{1}{2} \sum_{\ell=1}^L \text{Tr} [\hat{q}_\ell V_\ell + \hat{V}_\ell q_\ell] - \frac{1}{2} \int d\nu(\gamma, \tau) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \gamma_\ell \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right] \\
 &\quad - \alpha \mathbb{E}_{Y, \Xi} \mathcal{M}(Y, \Xi) \\
 &= \frac{1}{2} \sum_{\ell=1}^L \text{Tr} [\hat{q}_\ell V_\ell] + \frac{1}{2} \int d\nu(\gamma, \tau) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-2} \left(\sum_{\ell=1}^L \hat{V}_\ell \gamma_\ell - \lambda - \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right) \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \gamma_\ell \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right] \\
 &\quad - \alpha \mathbb{E}_{Y, \Xi} \mathcal{M}(Y, \Xi) \\
 &= -\frac{\lambda}{2} \int d\nu(\gamma, \tau) \text{Tr} \left[\left(\lambda + \sum_{\ell=1}^L \gamma_\ell \hat{V}_\ell \right)^{-1} \left(\sum_{\ell=1}^L \gamma_\ell \hat{q}_\ell + \left(\sum_{\ell=1}^L \tau_\ell \hat{m}_\ell + \hat{\theta}_\ell \cdot \pi \right)^{\otimes 2} \right) \right] \\
 &\quad - \alpha \mathbb{E}_{Y, \Xi} [\mathcal{M}(Y, \Xi)] + \frac{1}{2} \sum_{\ell=1}^L \text{Tr} [\hat{q}_\ell V_\ell] \\
 &= -\epsilon_t
 \end{aligned} \tag{94}$$

In other words, the training loss is equal to the zero-temperature free energy.

Appendix C. Supplementary Experiments

C.1. Empirical scaling of $\alpha = d/n$

In the following we verify that our experiments are consistent with the scaling behaviour predicted from the theory. We jointly increase d and n for a fixed value of α . In Fig. 2 we indeed observe the expected behaviour for an exemplary value of $\alpha = 2$. The same holds for the summary statistics θ and m , which concentrate as d and n jointly grow, shown in Fig. 1 (center) in the main text.

C.2. Alternative hyperparameters

We provide supplementary results for different parameter settings. Fig. 3 on the left shows more slices from the phase diagram that appears in the main Fig. 1. For the experimental section of the main text, we chose a specific A for definiteness. In the following, we present the same results for a different A with a stronger off-diagonal and a higher rank,

$$A = \begin{pmatrix} 0.3 & 0.7 \\ 0.8 & 0.2 \end{pmatrix}. \tag{95}$$

In Fig. 4 we present the analogous simulations to Fig. 1. While the global phenomena match the previous example, the details of the transitions location differ.

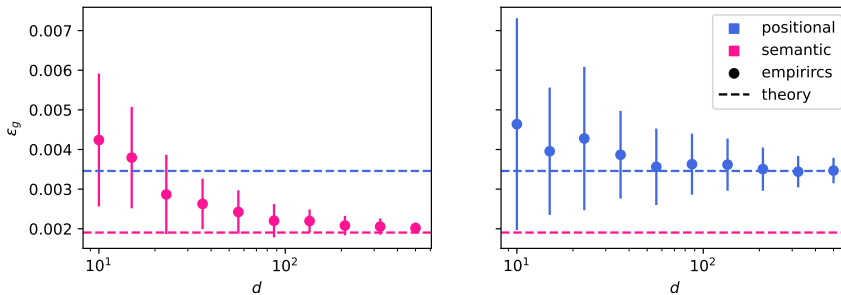


Figure 2: Scaling d and n jointly for $\alpha = 1.5$ approaches the theoretical prediction of the generalization error of the positional and semantic local minima. Experimental settings as in Fig. 1, with 70 runs per datapoint.

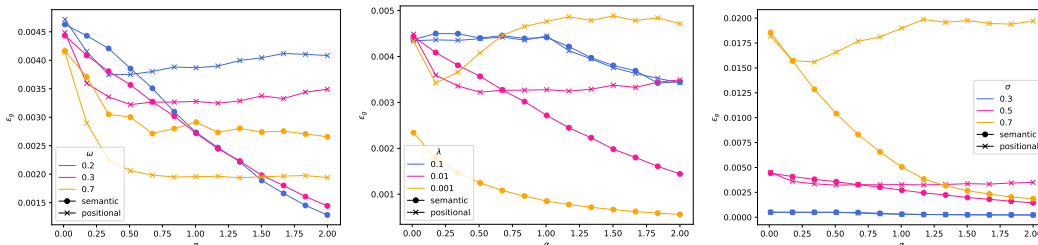


Figure 3: *Alternative Parameters*. Mixed positional/semantic teacher for $\omega = 0.3$. Settings is $r_s = r_t = 1, L = 2, A = ((0.6, 0.4), (0.4, 0.6)), \Sigma_1 = \Sigma_2 = 0.25\mathbb{I}_d, \mathbf{p}_1 = \mathbf{1}_d = -\mathbf{p}_2$ and $\mathbf{Q}_\star \sim \mathcal{N}(0, \mathbb{I}_d)$. While keeping all other settings the same, we vary from left to right: The target positionality ω , the student regularizer λ and the standard deviation σ (which is $0.5 = \sqrt{\Sigma_1}$). Experiment settings as in Fig. 1.

C.3. Uninformed initialization and training via Adam

In our experiments, to obtain the empirical results, we initialize the GD optimizer in an informed fashion, i.e. initializing \mathbf{Q}_\star of the student with $r = 1$ as either \mathbf{p}_1 (positional) or \mathbf{Q}_\star (semantics). GD then converges in the two local optima described by our theory.

Since our theory only ascertains that these solutions predicted are indeed fixed points of GD for large sizes, this does not have direct implications for other types of optimization algorithms. In Fig. 5 we show that indeed running the Adam optimizer from an *uninformed* initialization may lead one to either of the local minima for $d = 100$. For larger d we observe the semantic minimum is reached less often than the positional minimum, and a considerable number of times the algorithm simply does not find either of them.

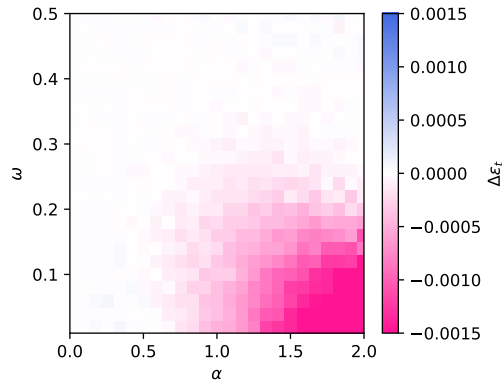


Figure 4: *Alternative Positional Matrix.* $r_s = r_t = 1, L = 2, \Sigma_1 = \Sigma_2 = 0.25\mathbb{I}_d, \mathbf{p}_1 = -\mathbf{p}_2$ and $\mathbf{p}_1, \mathbf{Q}_* \sim \mathcal{N}(0, \mathbb{I}_d)$ independently. Here, we use a definite matrix A from (95), which differs from the one used in the main text. Experiments were conducted as in Fig. 1.

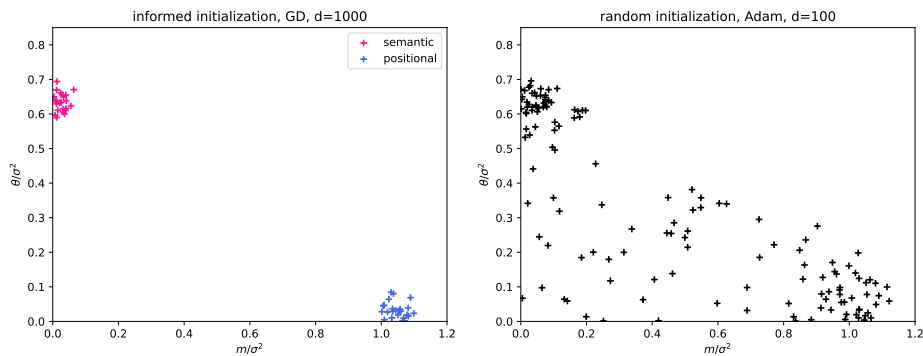


Figure 5: *Comparing GD and Adam.* Settings as in Fig. 1 for the sample complexity $\alpha = 2$. The student parameter \mathbf{Q} is obtained via either (left) positional and semantic informed initialization and (right) GD training from a random initialization are compared. Each point represents a single run. For the informed GD, we used the same optimization parameters as in Fig. 1 (24 runs per initialization). For Adam we trained on the same data, but for 2,500 epochs with learning rate $\eta = 0.01$ (showing 140 runs).