

Sequential Correlations Change In-Context Learning: Effective Context Length and Architectural Mismatch

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2026

Abstract

Modern sequence models have a striking capacity for in-context learning (ICL); they can perform new tasks based only on examples given in the prompt. Understanding how this ability emerges requires theory that capture important properties of natural data. Linear regression has served as a useful sandbox for ICL theory, but existing work has largely focused on prompts with independent examples. In this work, we extend this setting to sequentially correlated data, a basic feature of real sequences, present a solvable model based on linear attention, and test our predictions on realistic transformer architectures. We identify two distinct effects: First, when the query token is independent of the context, within-context correlations induce an effective context length: correlated prompts behave like shorter i.i.d. prompts. Second, when the query is also correlated with its context, test error is reduced, particularly for softmax attention when compared to linear attention. These results suggest that correlated prompts alter not only the effective sample size of in-context learning, but also which attention architectures are best matched to the task.

1. Introduction

In-context learning (ICL) is an important and useful capability of modern sequence models (Brown et al., 2020; Von Oswald et al., 2023; Wei et al., 2022). In this setting, a model performs a task implicitly from few examples without parameter updates. Regression has been a useful toy setting for studying how sequence models can achieve this ability. In-context regression requires a model to predict the label of a final query based on a prompt of input-output examples. The appeal of this setup is that it is simple enough to admit explicit analyses, yet rich enough to distinguish between architectures and learning mechanisms (Akyürek et al., 2023; Letey et al., 2026; Lu et al., 2025; Oko et al., 2024; Von Oswald et al., 2023; Zhang et al., 2024b, 2023).

A standard simplifying assumption in the in-context regression literature is that the examples within a prompt are independent. This assumption is analytically convenient, but not reflective of real datasets. This has prompted the study of sequential correlations in other non-regression in-context data settings, considering for example Markov chains and dynamical systems (Cole et al., 2025; Edelman et al., 2024; Nichani et al., 2024; Park et al., 2025); for the regression setting, sequential correlations have yet to be analysed. We know that sequential correlations matter even for classical ridge regression, where they change the asymptotic risk and invalidate estimators inherited from the i.i.d. setting (Atanasov et al., 2025). Thus it is natural to ask what is the effect of sequential correlations in the

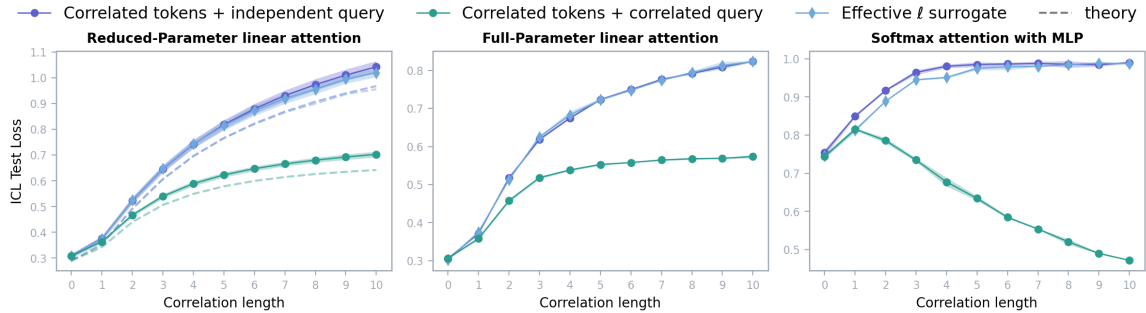


Figure 1: ICL test loss against token correlation strength for three attention architectures. We highlight two findings. (1) ICL performance on correlated data with an independent query (purple) is well approximated by uncorrelated data at an effective context length (blue) given by eq (1). (2) When the query token is equivalently correlated with its corresponding context (green), ICL error is reduced in all models as they utilise this correlation for inference. These predictions are made from our theory curves (dashed) derived for reduced-parameter linear attention. Experimental details given in Appendix A.

in-context regression setting, particularly as attention architectures are designed to operate on sequences in the first place.

In this work, we study a solvable model of in-context linear regression with sequentially correlated tokens. We begin from a reduced linear-attention model analysed in prior works (Bordelon et al., 2025; Letey et al., 2026; Lu et al., 2025; Wu et al., 2024; Zhang et al., 2024a), and perturb the usual i.i.d. setup by introducing a sequence correlation kernel for the context tokens, together with an optional correlation between the test query and the preceding context. This gives a controlled setting in which we can separate the effects of correlations *within* the context, and correlations *between* the query and the context.

Figure 1 summarises our findings. First, when correlations are confined to the context and the query remains independent, their effect is equivalent to an effective context-length reduction: correlated prompts behave like shorter i.i.d. prompts. Second, when the query token is correlated with its preceding context, ICL error decreases substantially, because the model can exploit this extra statistical structure to improve inference. Third, these correlated settings reveal an architectural mismatch: when additional query-dependent statistics enter, the performance gap between attention architectures (linear vs softmax in Figure 1) widens sharply.

An additional technical contribution of this work is distinguishing two cases of the asymptotic theory. Depending on the strength of the correlations (*i.e.*, the “length” over which they persist) relative to the sequence length, the high-dimensional treatment, as originally done by Lu et al. (2025), may fail: for correlations that persist on the same scale as the context length, moments of the reduced-linear attention estimator do not concentrate. We discuss this briefly in the main document and more thoroughly in Appendix D, and leave implications of this analysis to future work.

2. Theory and Data Setup

In-context linear regression. We study a correlated-token version of the solvable linear-regression ICL model induced by a reduced linear-attention block, as studied by [Lu et al. \(2025\)](#); [Zhang et al. \(2023\)](#). A single context consists of a sequence $\{\mathbf{x}_1, y_1, \dots, \mathbf{x}_\ell, y_\ell, \mathbf{x}_{\text{test}}\}$ of examples, where

$$y_a = \mathbf{x}_a^\top \mathbf{w} + \varepsilon_a, \quad y_{\text{test}} = \mathbf{x}_{\text{test}}^\top \mathbf{w} + \varepsilon_{\text{test}},$$

with task vector $\mathbf{w} \in \mathbb{R}^d$, and independent Gaussian noise $\varepsilon \sim \mathcal{N}(0, \rho I_\ell)$, $\varepsilon_{\text{test}} \sim \mathcal{N}(0, \rho)$. The goal is for the model’s output on this sequence to be close to y_{test} , thus performing the correct regression inference given examples in the context.

Sequential correlations. We model sequential structure within the context as linear dependence between the tokens, following previous works by [Atanasov et al. \(2025\)](#); [Moniri and Hassani \(2025\)](#) in the ridge regression setting. Under this model, the mean- $\mathbf{0}$ Gaussian tokens $\mathbf{x}_1, \dots, \mathbf{x}_\ell$ have second moments given by

$$\mathbb{E}[\mathbf{x}_a \mathbf{x}_b^\top] = K_{ab} I_d / d,$$

i.e., $K \in \mathbb{R}^{\ell \times \ell}$ controls correlations across context positions, and the feature-feature covariance is given by I_d/d . The case of independently-sampled tokens is described by $K = I_\ell$.

We will consider both cases where the mean- $\mathbf{0}$ test query \mathbf{x}_{test} is independent of X and correlated with X . We define the query correlations as above through the second moment

$$\mathbb{E}[\mathbf{x}_{\text{test}} \mathbf{x}_a^\top] = \mathbf{k}_{\text{test } a} I_d / d$$

where here $\mathbf{k}_{\text{test}} \in \mathbb{R}^\ell$ specifies how strongly the query is coupled to each preceding context token; for an independent query, $\mathbf{k}_{\text{test}} = \mathbf{0}$.

Correlation summaries and correlation length. Without loss of generality, we take $\text{tr}(K) = \ell$. The remaining correlation-dependent scalars that appear in the theory are

$$k_0 := \mathbf{k}_{\text{test}}^\top \mathbf{k}_{\text{test}}, \quad k_1 := \mathbf{k}_{\text{test}}^\top K \mathbf{k}_{\text{test}}, \quad k_2 := \text{tr}(K^2) / \ell.$$

These quantities act as effective summaries of the token correlations seen by the ICL estimator, arising from the structure of the covariates H .

For general K and \mathbf{k}_{test} , we can think of k_0, k_1, k_2 as measures of token correlation strength. The quantity k_2 is a bulk measure of context correlations: under our chosen normalisation $\text{tr}(K) = \ell$, it is minimised by the uncorrelated case $K = I_\ell$, where $k_2 = 1$, and increases as the spectrum of K becomes more anisotropic or concentrated. The query-dependent quantity k_0 measures the overall strength of the direct query-context coupling, while the mixed quantity k_1 measures how strongly that query-correlation profile is amplified by the correlated structure of the context itself.

These quantities become more intuitive if we take a specific kernel as an example: the exponential kernel $K_{ab} = \exp(-|a - b|/\xi)$ and $\mathbf{k}_{\text{test } a} = \exp(-(\ell + 1 - a)/\xi)$. Because K and \mathbf{k}_{test} only depend on *distances* between tokens, it makes sense to think of correlation “strength” in terms of correlation *length*; here, the correlation length is explicitly ξ . This case provides an intuitive sandbox, as the ICL summary statistics above are directly related to the correlation length by $k_2, k_0 \approx \xi$, $k_1 \approx \xi^2$. Given this clear scaling, we choose to use the exponential kernel for the experiments in [Figures 1 and 2](#). We save detailed interpretation of the summary statistics k_0, k_1, k_2 for non-exponential kernels to future work.

High-dimensional scaling. As standard for theoretical analysis of linear regression, we will work in a proportional regime of context length, *i.e.*, $d \rightarrow \infty$ with $\ell \propto d$. We assume the task signal remains identifiable, e.g. $\mathbf{w}^\top \mathbf{w} = \Theta(d)$.

We will analytically study the case of *weak correlations* with respect to context length ℓ , namely, for the computation of theory curves we will take $k_2, k_0, k_1 = \Theta_\ell(1)$. The case of *strong correlations* corresponds to $k_2, k_0 = \Theta(\ell), k_1 = \Theta(\ell^2)$. This case does not lend itself to an analytical solution in our asymptotic limit as the ICL error does not concentrate. Our theoretical results will thus focus on the former “weak” correlation case, with the “strong” correlation case discussed briefly in Section 3 and Appendices C and D.

Reduced linear-attention predictor. Following prior work on linear attention, the next-token prediction for y_{test} made by linear attention can be well-approximated by

$$\hat{y}_{\text{test}} := \text{tr}(\Gamma H^\top), \quad \Gamma \in \mathbb{R}^{d \times (d+1)}, \quad H := \mathbf{x}_{\text{test}} \left[\frac{d}{\ell} (X\mathbf{w} + \boldsymbol{\varepsilon})^\top X \quad \frac{1}{\ell} (X\mathbf{w} + \boldsymbol{\varepsilon})^\top (X\mathbf{w} + \boldsymbol{\varepsilon}) \right].$$

The parameter matrix Γ here is defined in terms of components of the value, key, and query matrices from full-parameter linear attention. This is the predictor that we will study in the theory, and is the “reduced model” referred to on the left-most panel of Figure 1.

Given a training batch of n sequences, corresponding to n such data-matrices H for the reduced model, the optimal parameters and final test loss can be written as

$$\Gamma^* = \left(\frac{n}{d} \lambda I_{d(d+1)} + \sum_{\mu=1}^n H^\mu \otimes H^\mu \right)^{-1} \sum_{\mu=1}^n y_{\text{test}}^\mu H^\mu, \quad \mathcal{E}_{\text{ICL}}(\Gamma^*) = \mathbb{E}[(y_{\text{test}}^{\text{new}} - \langle H^{\text{new}}, \Gamma^* \rangle)^2].$$

3. Effective in-context sample size

Working in the proportional limit described above, we can analyse the random parameter matrix Γ^* . The full result is given in Appendix E. We present its first implication.

Result: Effective-Sample Size

For prompts with weakly-correlated tokens $\mathbf{x}_1, \dots, \mathbf{x}_\ell$ and an *independent* query \mathbf{x}_{test} , the correlated-ICL performance at context length ℓ is equivalent to **uncorrelated**-ICL performance at lower context length given by

$$\ell_{\text{eff}} := \frac{1 + \rho}{k_2 + \rho} \ell, \tag{1}$$

where by Jensen’s inequality $k_2 = \text{tr}(K^2)/\ell \geq (\text{tr}(K)/\ell)^2 = 1$, so $\ell_{\text{eff}} \leq \ell$.

Eq (1) makes precise how sequential correlations reduce the usable information content of the prompt. Since k_2 increases as the context covariance becomes more anisotropic, stronger or more persistent correlations decrease the effective context length ℓ_{eff} . For stationary kernels, where k_2 grows with the correlation scale, this recovers the intuitive statement that longer-ranged correlations make fewer context examples effectively independent. The simplicity of the resulting formula is itself notable: the reduction depends linearly on the single scalar k_2 , rather than on the full spectrum of K . This contrasts with classical correlated linear regression, where test error typically depends more delicately on the spectrum and eigenvectors of the design covariance.

Eq (1) also clarifies why the strongly-correlated regime is analytically difficult. Our theory here assumes that the correlation summary k_2 remains $\Theta(1)$ as $\ell, d \rightarrow \infty$ with $\ell \propto d$. But once correlations persist across more and more of the context, we instead have $k_2 = \Theta(\ell)$. In this regime, ℓ_{eff} no longer scales proportionally with ℓ , but remains only $\Theta(1)$. Thus, even as the nominal context length increases, the model effectively sees only finitely many independent examples. This is precisely the regime in which the standard high-dimensional concentration underlying our asymptotic analysis breaks down, as discussed more thoroughly in Appendix D.

4. Gain from query correlations

We now turn to the second effect in Figure 1: correlations between the test query and its preceding context. To isolate this contribution, define the *query-correlation gain* as

$$\Delta_{\text{query}} := e_{\text{ICL}}(\ell, k_2; \mathbf{k}_{\text{test}} = 0) - e_{\text{ICL}}(\ell, k_2, k_0, k_1),$$

namely, the reduction in ICL error obtained by introducing query-context correlations while holding the bulk context kernel K fixed.

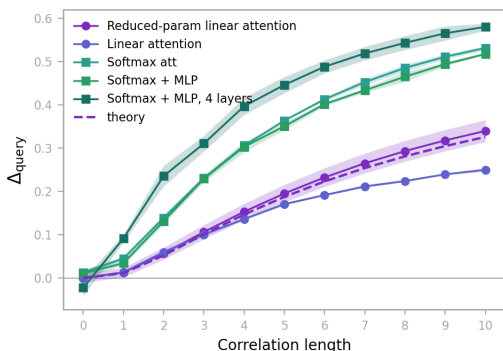


Figure 2: Correlated-query effects depend on architecture.

modest for linear attention architectures, but much larger for architectures that use softmax attention, indicating that the latter is better able to exploit informative query-context dependence.

5. Conclusion

Prompt correlations change in-context learning in two distinct ways: bulk context correlations reduce effective sample size, while query-context correlations provide useful signal for inference. In correlated linear regression, the first effect is captured by a simple effective-context-length law, while the second reveals a sharp architectural mismatch, with softmax attention benefiting much more than linear attention. These results suggest that understanding ICL on realistic structured data requires going beyond i.i.d. prompt models.

In Figure 2, this gain is almost always nonnegative, and vanishes only in the i.i.d. case. Thus, unlike bulk context correlations which only reduce effective sample size and hurt ICL performance, query correlations provide additional predictive signal that can in principle be exploited for inference.

In our theory, this effect is controlled not only by the bulk statistic k_2 , but also by the query-dependent summaries k_0 and k_1 , which increase with stronger correlations. The resulting closed-form expression for Δ_{query} is substantially more complicated than the effective-sample law in eq. (1), and we therefore omit it from the main text.

That said, we see that Δ_{query} is relatively

References

- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=OgOX4H8yN4I>.
- Alexander Atanasov, Jacob A Zavatore-Veth, and Cengiz Pehlevan. Risk and cross validation in ridge regression with correlated samples. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=GMwKpJ9TiR>.
- Blake Bordelon, Mary I. Letey, and Cengiz Pehlevan. Theory of scaling laws for in-context regression: Depth, width, context and time, 2025. URL <https://arxiv.org/abs/2510.01098>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Frank Cole, Yulong Lu, Tianhao Zhang, and Yuxuan Zhao. In-context learning of linear dynamical systems with transformers: Error bounds and depth-separation, 2025. URL <https://arxiv.org/abs/2502.08136>.
- Benjamin L. Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution of statistical induction heads: In-context learning markov chains, 2024. URL <https://arxiv.org/abs/2402.11004>.
- Mary Letey, Jacob A Zavatore-Veth, Yue M. Lu, and Cengiz Pehlevan. Pretrain-test task alignment governs generalization in in-context learning. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=KZLeg0MQ2r>.
- Yue M. Lu, Mary Letey, Jacob A. Zavatore-Veth, Anindita Maiti, and Cengiz Pehlevan. Asymptotic theory of in-context learning by linear attention. *Proceedings of the National Academy of Sciences*, 122(28):e2502599122, 2025. doi: 10.1073/pnas.2502599122. URL <https://www.pnas.org/doi/abs/10.1073/pnas.2502599122>.
- Behrad Moniri and Hamed Hassani. Asymptotics of linear regression with linearly dependent data. In Necmiye Ozay, Laura Balzano, Dimitra Panagou, and Alessandro Abate, editors, *Proceedings of the 7th Annual Learning for Dynamics & Control Conference*, volume 283 of *Proceedings of Machine Learning Research*, pages 72–85. PMLR, 04–06 Jun 2025. URL <https://proceedings.mlr.press/v283/moniri25a.html>.
- Eshaan Nichani, Alex Damian, and Jason D. Lee. How transformers learn causal structure with gradient descent. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of*

- Machine Learning Research*, pages 38018–38070. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/nichani24a.html>.
- Kazusato Oko, Yujin Song, Taiji Suzuki, and Denny Wu. Pretrained transformer efficiently learns low-dimensional target functions in-context, 2024. URL <https://arxiv.org/abs/2411.02544>.
- Core Francisco Park, Ekdeep Singh Lubana, Itamar Pres, and Hidenori Tanaka. Competition dynamics shape algorithmic phases of in-context learning, 2025. URL <https://arxiv.org/abs/2412.01003>.
- 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.
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 35151–35174. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/von-oswald23a.html>.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=yzkSU5zdwD>.
- Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter Bartlett. How many pretraining tasks are needed for in-context learning of linear regression? In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=vSh5ePa0ph>.
- Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024a. URL <http://jmlr.org/papers/v25/23-1042.html>.
- Ruiqi Zhang, Jingfeng Wu, and Peter L. Bartlett. In-context learning of a linear transformer block: Benefits of the mlp component and one-step gd initialization, 2024b.
- Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in-context learning learn? bayesian model averaging, parameterization, and generalization, 2023. URL <https://arxiv.org/abs/2305.19420>.

Appendix A. Experimental details

The data structure follows the Section 2 setup. For all experiments we choose a stationary kernel, *i.e.*, distance-dependent correlations, given by

$$K_{ab} := \exp(-|a - b|/\xi). \quad (2)$$

Here the correlation length ξ is manifestly part of the definition of K . This definition is convenient to choose due to the clarity of its correlation length parameter, as well as the fact that our Gaussian model with this kernel is equivalent to an AR(1) process which allows for efficient sampling. When query correlations are present, *i.e.*, $\mathbf{k}_{\text{test}} \neq \mathbf{0}$, we choose matching structure

$$(\mathbf{k}_{\text{test}})_a := \exp(-|\ell + 1 - a|/\xi). \quad (3)$$

This parameter ξ is what is referred to on the x -axes of Figures 1 and 2.

Data parameters. A batch of data contains sequences of length ℓ with one query each. We train in an offline manner with n total sequences. Each sequence is defined by a task vector \mathbf{w} . As in Letey et al. (2026); Lu et al. (2025); Raventós et al. (2023), the total number of task vectors in the pretraining batch may differ from n . The total number of unique task vectors in the batch is called k , and each unique vector is sampled i.i.d. from $\mathcal{N}(\mathbf{0}, I_d)$. For all experiments, we take

$$d = 32, \quad n = 4096 = 4d^2, \quad \ell = 128 = 4d, \quad k = 320 = 10d$$

matching the asymptotic scalings of these variables derived in Lu et al. (2025). We take label noise $\rho = 0.01$.

Architecture details. The experiments train sequence models on the above linear regression sequence data with sequentially correlated inputs. The architectures considered include softmax-only layers, softmax+MLP transformer-like architectures, and pure linear attention. Note that we distinguish between the reduced model (where the optimal solution is given directly by Γ^*) from *trainable* fully-parameterised linear attention. Softmax models embed inputs into width $d+1$ with a dedicated label channel, initialise input embeddings with scale proportional to $\sqrt{(d+1)}$, initialise query/key maps on the input channels, and initialise value/output maps to primarily use the label channel with small Gaussian noise. MLP blocks use GELU and are initialised close to zero. Linear-attention models operate directly on concatenated (\mathbf{x}_i, y_i) tokens of dimension $d+1$, with query/key/value matrices initialised near identity. Models are trained by minimising mean-squared error on the held-out query label only. Optimisation uses Optax AdamW-style updates with weight decay `lamb`; the total number of gradient steps is set to T (something large enough to reach training loss convergence), but we employ early stopping and report the test loss of the best checkpointed-model. We use a linear learning-rate warmup from 0 to `max_lr` over the first 10% of training, with learning rate remaining constant after warmup. We use exponential moving average weights for evaluation. Minibatching is used to divide the full dataset of n sequences; we find that smaller batch sizes perform better for softmax attention architectures. Figures 1 and 2 show mean best-test loss over different seeds, where seeds control train data batch, testing samples, and batch divisions. A summary of our training hyperparameters is

$$\text{max_lr} = 0.001, \quad \text{lamb} = 0.0001, \quad T = 10000, \quad \text{batch size} = 64 = 2d, \quad \# \text{ seeds} = 5.$$

Appendix B. Detailed setup

Here we re-iterate the data and model setup given in the main document with more details fleshed out. We leave the token design matrix general, *i.e.*,

$$\mathbf{x}_a \sim_K \mathcal{N}(0, \Sigma/d)$$

where the sample-sample correlations are given by positive semi-definite K as

$$\mathbb{E}_X[X_{ai}X_{bj}] = \frac{1}{d}K_{ab}\Sigma_{ij} \quad \text{for } a, b \in [\ell] \text{ and } i, j \in [d]. \quad (4)$$

The label noise we take to be sequentially uncorrelated

$$\varepsilon \sim \mathcal{N}(0, \rho I_{\ell+1}) \quad \text{independent of all } X \text{ and } \mathbf{x}_{\text{test}}.$$

The query \mathbf{x}_{test} may be correlated with the contexts as

$$\mathbb{E}[\mathbf{x}_{\text{test}}|X] = X^\top \mathbf{m}, \quad \text{Var}(\mathbf{x}_{\text{test}}|X) = \frac{1}{d}(1 - k_{-1})\Sigma$$

for

$$\mathbf{m} \equiv K^{-1}\mathbf{k}_{\text{test}} \in \mathbb{R}^\ell, \quad k_{-1} \equiv \mathbf{k}_{\text{test}}^\top K^{-1}\mathbf{k}_{\text{test}} \leq 1.$$

This is equivalent to sampling, as in eq (4), using $\ell+1 \times \ell+1$ positive-semidefinite correlation kernel

$$K_{\text{query}} = \begin{bmatrix} K & \mathbf{k}_{\text{test}} \\ \mathbf{k}_{\text{test}}^\top & 1 \end{bmatrix}.$$

As a final note, this entire data setup is equivalent to

$$X = \frac{1}{\sqrt{d}}\sqrt{K}Z\sqrt{\Sigma}, \quad Z_{si} \sim_{\text{i.i.d.}} \mathcal{N}(0, 1).$$

B.1. Notation and Assumptions

We will consider a standard high-dimensional proportional limit, as in (Letey et al., 2026; Lu et al., 2025). The context length scales as

$$\alpha \equiv \frac{\ell}{d} = \Theta_{d,\ell}(1).$$

The usual high dimensional assumptions on Σ are such that

$$\text{tr}(\Sigma), \text{tr}(\Sigma\mathbf{w}\mathbf{w}^\top) = \Theta(d).$$

These are a succinct summary of conditions on both Σ and \mathbf{w} . The first condition $\text{tr}(\Sigma) = \Theta(d)$ ensures that there is reasonable signal to estimate in the tokens; note that after this section we will always assume $\Sigma = I_d$. The second condition $\text{tr}(\Sigma\mathbf{w}\mathbf{w}^\top) = \Theta(d)$ is equivalent to $\|\mathbf{w}\|^2 = \Theta(d)$ ensuring that there is enough recoverable task signal. Given these assumptions, everything will be written in normalised $\Theta_d(1)$ quantities

$$\text{tr}[\Sigma] := \frac{1}{d} \text{tr}(\Sigma), \quad \text{tr}[\Sigma\mathbf{w}\mathbf{w}^\top] := \frac{1}{d} \text{tr}(\Sigma\mathbf{w}\mathbf{w}^\top).$$

These need to be extended to assumptions about K . We will take

$$\text{tr}[K] = \frac{1}{\ell} \text{tr}(K) = 1$$

The final thing to reason about is the contributions from K^2 and from \mathbf{k}_{test} . This is where the notion of ‘‘correlation length’’ captured in K and \mathbf{k}_{test} becomes important, as the terms

$$\text{tr}(K^2), \quad k_{-1} \equiv \mathbf{k}_{\text{test}}^\top K^{-1}\mathbf{k}_{\text{test}}, \quad k_0 \equiv \mathbf{k}_{\text{test}}^\top \mathbf{k}_{\text{test}}, \quad k_1 \equiv \mathbf{k}_{\text{test}}^\top K \mathbf{k}_{\text{test}}$$

will appear in the computation.

Working definition of correlation length Suppose we choose K to be PSD and Toeplitz, *i.e.*,

$$K_{st} = c(|s - t|).$$

This means the correlations between \mathbf{x}_s and \mathbf{x}_t are stationary: they only depend on their relative distance in the sequence. A decent starting definition for the “correlation length” defined by c would be the sum of one of its columns, or the sum over all the lags, *i.e.*,

$$\hat{\xi} = \sum_{\tau=1}^{\ell-1} c(\tau).$$

Example: the exponential kernel Here take $c(\tau) = \exp(-\tau/\xi) \equiv \phi^\tau$. Then $\hat{\xi}$ recovers at least the correct scaling of the correlation length.

- **Finite correlations:** $\xi = \Theta_\ell(1)$. Then $\hat{\xi} \approx \xi$ (for ξ not too close to 0).
- **Proportional correlations:** $\xi = \beta\ell$. Then $\hat{\xi} \approx (1 - \exp(-1/\beta))\xi$.

We can compute k_{-1}, k_0, k_1 explicitly in each of these regimes, giving us a surrogate for how these quantities should scale for a general K or $c(\tau)$. Have

$$\frac{1}{\ell} \text{tr}(K^2) \sim \xi, \quad k_{-1} \leq 1, \quad k_0 \sim \xi, \quad k_1 \sim \xi^2$$

and

$$\frac{1}{\ell} \text{tr}(K^2) \sim \beta\ell, \quad k_{-1} = 1, \quad k_0 \sim (1 - \exp(-2/\beta))\frac{\beta}{2}\ell, \quad k_1 \sim \beta^2\ell^2.$$

Final assumptions Given this example, we will distinguish the asymptotic treatment between the two regimes

Finite correlations: $k_2 = \frac{1}{\ell} \text{tr}(K^2), \quad k_0, \quad k_1 = \Theta(1)$

Proportional correlations: $\tilde{k}_2 = \frac{1}{\ell^2} \text{tr}(K^2), \quad \tilde{k}_0 = \frac{k_0}{\ell}, \quad \tilde{k}_1 = \frac{k_1}{\ell^2} = \Theta(1)$

Appendix C. Population losses

We will be considering linear features

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\varepsilon}, \quad y_{\text{test}} = \mathbf{x}_{\text{test}}^\top \mathbf{w} + \mathcal{N}(0, \rho).$$

For each such sequence of $X, \mathbf{y}, \mathbf{x}_{\text{test}}, y_{\text{test}}$, define a data matrix

$$H = \mathbf{x}_{\text{test}} \left[\frac{d}{\ell} (X\mathbf{w} + \boldsymbol{\varepsilon})^\top X \quad \frac{1}{\ell} (X\mathbf{w} + \boldsymbol{\varepsilon})^\top (X\mathbf{w} + \boldsymbol{\varepsilon}) \right] \in \mathbb{R}^{d \times (d+1)}.$$

The predictor we will use for y_{test} is given by

$$\hat{y}_{\text{test}} = \langle \Gamma, H \rangle = \text{tr}(\Gamma H^\top)$$

for parameters $\Gamma \in \mathbb{R}^{d \times (d+1)}$. Before we consider ICL error at optimal parameters Γ^* , we can compute a population formula for

$$\mathcal{E}_{\text{ICL}}(\Gamma) = \mathbb{E}_{\text{new data}} \left[(y_{\text{test}}^{\text{new}} - \langle \Gamma, H^{\text{new}} \rangle)^2 \right].$$

This formula will be different depending on which asymptotic treatment of k_2, k_0, k_1 we are considering.

Notation. We will use the convention of row-wise vectorisation, and so have

$$\text{Vec}(vu^\top) = v \otimes u \in \mathbb{R}^{\dim(v) \times \dim(u)}, \quad \text{Vec}(vu^\top)\text{Vec}(vu^\top)^\top = (vv^\top) \otimes (\mathbf{u}\mathbf{u}^\top)$$

Lemma 1 *Suppose we have weak-range correlations, i.e., $k_2, k_0, k_1 = \Theta(1)$. For $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, C)$, writing $\text{tr}[C] = \text{tr}(C)/d$, with*

$$\rho_1 \equiv \text{tr}[C] + \rho, \quad \rho_2 \equiv k_2 \text{tr}[C] + \rho$$

have **population ICL risk**

$$\mathcal{E}_{\text{ICL}}(\Gamma) = \rho + \text{tr}[C] - 2\frac{1}{d} \text{tr}(\Gamma A^\top) + \frac{1}{d} \text{tr}(\Gamma B \Gamma^\top) + \frac{1}{d^2} \text{tr}(\mathcal{T} \text{vec}(\Gamma) \text{vec}(\Gamma)^\top) \quad (5)$$

for

$$\begin{aligned} A &\equiv \begin{bmatrix} C + \frac{1}{\alpha} k_0 \text{tr}[C] I_d & \mathbf{0} \\ \mathbf{0} & \rho_1^2 \end{bmatrix}, & B &\equiv \begin{bmatrix} C + \frac{\rho_2}{\alpha} I_d & \mathbf{0} \\ \mathbf{0} & \rho_1^2 \end{bmatrix} \\ \mathcal{T} &\equiv \frac{k_1 \text{tr}[C] + \rho k_0}{\alpha^2} \text{vec}([I \ \mathbf{0}]) \text{vec}([I \ \mathbf{0}])^\top \\ &+ \frac{k_0}{\alpha} \left(\text{vec}([I \ \mathbf{0}]) \text{vec}([C \ \mathbf{0}])^\top + \text{vec}([C \ \mathbf{0}]) \text{vec}([I \ \mathbf{0}])^\top \right) \end{aligned}$$

Lemma 2 *Suppose we have strong-range correlations, i.e., $k_0, k_2 = \Theta(\ell)$ and $k_1 = \Theta(\ell^2)$. Then writing $\Gamma = [\Gamma_{\text{sq}} \ \gamma]$ we have*

$$\mathcal{E}_{\text{ICL}}(\Gamma) = \rho + \text{tr}[C] - 2 \text{tr}[C] \text{tr}(\Gamma_{\text{sq}}) + \text{tr}[C] \tilde{k}_2 \text{tr}(\Gamma_{\text{sq}} \Gamma_{\text{sq}}^\top) + \text{tr}[C] \tilde{k}_1 \text{tr}(\Gamma_{\text{sq}} \Gamma_{\text{sq}}) + \text{tr}[C] \tilde{k}_1 \text{tr}(\Gamma_{\text{sq}})^2$$

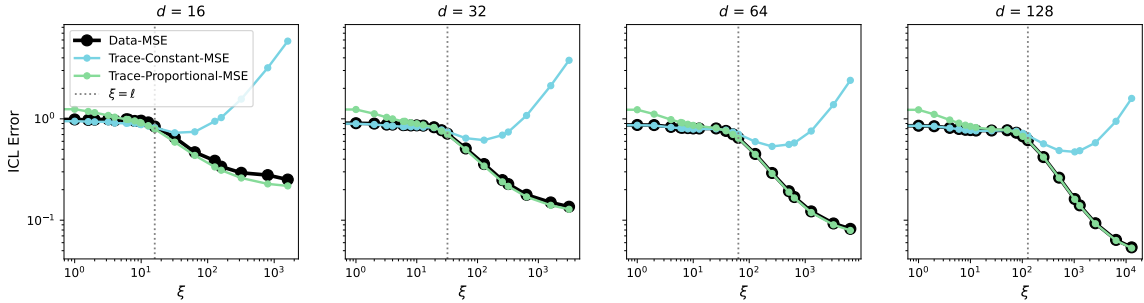


Figure 3: Illustration of transition between these two regimes. Quantities here are ICL MSE for Γ^* computed numerically from data, where the MSE is computed from (1) sampling the test distribution directly (2) using the population average in Lemma 1 (the weakly correlated case) and (3) using the population average in Lemma 2 (the strongly correlated case).

Figure 3 shows clearly the different regimes of correlation strength. Here we use the exponential kernel setting with correlated queries, i.e., K and \mathbf{k}_{test} given by eq.s (2) and (3). We see the Lemma 1 formula tracks the data-sampled MSE for low correlation lengths, and likewise Lemma 2 for proportional correlation lengths, as predicted.

Appendix D. Non-concentration of finite-sample ICL error for strong correlations

We can immediately see in Figure 3 the lack of concentration that appears from the proportional correlation case: as d increases, the MSE losses for $\xi \propto \ell$ keeps decreasing, while the $\xi \ll \ell$ losses show minimal variation across d .

In section we will derive an asymptotically-precise formula for the finite-sample ICL error, i.e. $\mathcal{E}_{\text{ICL}}(\Gamma^*)$ where Γ^* depends on a *finite* selection of training data sequences. Indeed $\mathcal{E}_{\text{ICL}}(\Gamma^*)$ is a random quantity, and the fact that we can write down a deterministic equivalent for it hinges on its concentration in high dimensions. This quantity concentrates for the weak-correlations case, i.e. the case covered in the population analysis by Lemmas 1 and 4, where we have $k_0, k_1, k_2 = \Theta(1)$ as $\ell \rightarrow \infty$. Here we will explain why $\mathcal{E}_{\text{ICL}}(\Gamma^*)$ does *not* concentrate in the strong-correlations case, where $k_0, k_2 \propto \ell$, $k_1 \propto \ell^2$.

The optimal parameters for our reduced linear attention model are given by

$$\text{vec}(\Gamma^*) = \left(\frac{n}{d} \lambda J_{d(d+1)} + \sum_{\mu=1}^n \text{vec}(H^\mu) \otimes \text{vec}(H^\mu) \right)^{-1} \sum_{\mu=1}^n y_{\text{test}}^\mu \text{vec}(H^\mu). \quad (6)$$

We can see from this form that a necessary condition of \mathcal{E}_{ICL} will be the concentration of the data vectors $\text{vec}(H^\mu)$. By concentration, we mean that its norm concentrates as $d, \ell \rightarrow \infty$ in the proportional limit. We will thus analyse the asymptotic behaviour of $\|\text{vec}(H)\|$ where we are suppressing the index $\mu \in [n]$.

Recall that we have

$$H = \mathbf{x}_{\text{test}} \left[\frac{d}{\ell} (X\mathbf{w} + \boldsymbol{\varepsilon})^\top X \quad \frac{1}{\ell} (X\mathbf{w} + \boldsymbol{\varepsilon})^\top (X\mathbf{w} + \boldsymbol{\varepsilon}) \right].$$

We will simplify our analysis here, as our goal is to focus on intuition, by ignoring the noise terms given by $\boldsymbol{\varepsilon}$; these do not affect the argument. Thus we write

$$\text{vec}(H) = \mathbf{v} \otimes \mathbf{x}_{\text{test}} + \text{ignored noise contribution}$$

where

$$\mathbf{v} = \begin{bmatrix} \frac{d}{\ell} X^\top X \mathbf{w} \\ \frac{1}{\ell} \mathbf{w}^\top X^\top X \mathbf{w} \end{bmatrix}.$$

We thus have that

$$\begin{aligned} \frac{1}{d^2} \|\text{vec}(H)\|^2 &= \frac{1}{d^2} \|\mathbf{x}_{\text{test}}\|^2 \|\mathbf{v}\|^2 \\ &= \frac{1}{d^2} \|\mathbf{x}_{\text{test}}\|^2 \left(\frac{d^2}{\ell^2} \mathbf{w}^\top X^\top X X^\top X \mathbf{w} + \frac{1}{\ell^2} (\mathbf{w}^\top X^\top X \mathbf{w})^2 \right) \\ &\stackrel{d}{=} \left(\frac{1}{d} \mathbf{m}^\top K^{1/2} Z Z^\top K^{1/2} \mathbf{m} + 2 \frac{\sqrt{1-k_{-1}}}{d} \mathbf{z}^\top Z^\top K^{1/2} \mathbf{m} + \frac{1-k_{-1}}{d} \|\mathbf{m}_z\|^2 \right) \\ &\quad \times \left(\frac{1}{d^2 \ell^2} \mathbf{w}^\top Z^\top K Z Z^\top K Z \mathbf{w} + \frac{1}{d^4 \ell^2} (\mathbf{w}^\top Z^\top K Z \mathbf{w})^2 \right), \end{aligned}$$

where $\stackrel{d}{=}$ means equivalent in distribution. Here Z is an $\ell \times d$ matrix with i.i.d. $\mathcal{N}(0, 1)$ entries and $\mathbf{z} \sim \mathcal{N}(0, I_d)$ for the query (independent of Z).

Upon analysis of each of these terms, which we omit for brevity, we find that the problematic term is

$$\begin{aligned} \frac{1}{d^2 \ell^2} \mathbf{w}^\top Z^\top K Z Z^\top K Z \mathbf{w} &\stackrel{d}{=} \|\mathbf{w}\|^2 (Z^\top K Z Z^\top K Z)_{11} \\ &\stackrel{d}{=} \frac{\|\mathbf{w}\|^2}{d} \frac{1}{\ell^2} \sum_{a,b=1}^{\ell} \kappa_a \kappa_b Z_{a1} Z_{b1} \frac{1}{d} \sum_{i=1}^d Z_{ai} Z_{bi} \\ &= \frac{\|\mathbf{w}\|^2}{d} \frac{1}{d} \left(\frac{1}{\ell} \sum_{a=1}^{\ell} \kappa_a Z_{a1}^2 \right)^2 + \frac{\|\mathbf{w}\|^2}{d} \frac{1}{\ell^2} \sum_{a,b=1}^{\ell} \kappa_a \kappa_b Z_{a1} Z_{b1} \frac{1}{d} \sum_{i=2}^d Z_{ai} Z_{bi}. \end{aligned}$$

We have arrived at this expression by exploiting the fact that the distribution of the matrix Z is invariant under both left and right rotation to perform two changes of basis: (1) we let K have orthogonal eigendecomposition $K = O \text{diag}(\kappa_1, \dots, \kappa_\ell) O^\top$, and (2) we choose a basis such that $\mathbf{w} = (\|\mathbf{w}\|, 0, \dots, 0)^\top$.

Now we must analyse this expression to provide a condition on K , specifically its eigenvalues $\kappa_1, \dots, \kappa_\ell$, for when this term does or does not concentrate. We focus on the second term in the sum

$$A := \frac{\|\mathbf{w}\|^2}{d} \frac{1}{\ell^2} \sum_{a,b=1}^{\ell} \kappa_a \kappa_b Z_{a1} Z_{b1} \frac{1}{d} \sum_{i=2}^d Z_{ai} Z_{bi}.$$

It has moments given by

$$\begin{aligned} \mathbb{E}[A] &= \frac{d-1}{d} \frac{1}{\ell^2} \sum_{a=1}^{\ell} \kappa_a^2, \\ \text{var}[A] &= \frac{2}{\ell^4 d^2} (d-1) \left(\sum_{a=1}^{\ell} \kappa_a^2 \right)^2 + \frac{2}{\ell^4 d^2} (d^2-1) \sum_{a=1}^{\ell} \kappa_a^4, \end{aligned}$$

and so the square of the coefficient of variation of A is therefore

$$\begin{aligned} \frac{\text{var}[A]}{\mathbb{E}[A]^2} &= 2 \left[\frac{1}{d-1} + \frac{d^2-1}{(d-1)^2} \frac{\sum_{a=1}^{\ell} \kappa_a^4}{\left(\sum_{a=1}^{\ell} \kappa_a^2 \right)^2} \right] \\ &= 2[1 + \mathcal{O}(d^{-1})] \frac{\sum_{a=1}^{\ell} \kappa_a^4}{\left(\sum_{a=1}^{\ell} \kappa_a^2 \right)^2} + \mathcal{O}(d^{-1}). \end{aligned}$$

We now notice that the K -dependence here

$$\frac{\sum_{a=1}^{\ell} \kappa_a^4}{\left(\sum_{a=1}^{\ell} \kappa_a^2 \right)^2} = \frac{\text{tr}(K^4)}{\text{tr}(K^2)^2}$$

looks like some sort of participation ratio. Indeed if $K = I_\ell$ (uncorrelated case, we know this concentrates) then

$$\frac{\text{tr}(K^4)}{\text{tr}(K^2)^2} = \frac{1}{\ell} \rightarrow 0,$$

while if $K = \mathbf{1}\mathbf{1}^\top$ (maximally correlated case, here $\mathbf{x}_1 = \dots = \mathbf{x}_\ell$) then

$$\frac{\text{tr}(K^4)}{\text{tr}(K^2)^2} = 1.$$

This gives us an intuitive heuristic for concentration of $\|\text{vec}(H)\|$.

Necessary condition for concentration

The random quantity $\mathcal{E}_{\text{ICL}}(\Gamma^*)$, where Γ^* is determined by data sampled with correlation kernel K , concentrates to a deterministic high-dimensional limit only if

$$\frac{\text{tr}(K^4)}{\text{tr}(K^2)^2} \rightarrow 0. \quad (7)$$

We can estimate this in terms of “correlation length” ξ for the exponential case

$$K_{ab} = \exp(-|a - b|/\xi)$$

finding that

$$\frac{\text{tr}(K^4)}{\text{tr}(K^2)^2} \sim \frac{\xi}{\ell}$$

which gives the condition we expected. We see that for correlations that persist on the same order as the context length, we do not have concentration of $\text{vec}(H)$, and thus we do not have concentration of $\mathcal{E}_{\text{ICL}}(\Gamma^*)$. An extension of great interest to us would be to analyse this nonconcentration for a wider range of K choices, *e.g.* how does this participation ratio behave for other non-exponential (or even non-Toeplitz) choices of K ?

Appendix E. Asymptotic formula for ICL error in the weak-correlations case

Here we present a deterministic formula for $\mathcal{E}_{\text{ICL}}(\Gamma^*)$ in the weak-correlations case. Again, the optimal parameters are given by eq. (6). The explicit ridge λ is a parameter that adds regularisation to this solution; it is common to take the ridgeless limit $\lambda \rightarrow 0$ in similar works.

Here we have n sample sequences, giving n different $H^\mu \in \mathbb{R}^{d \times (d+1)}$ data embeddings, where the sequence length is of course $\ell + 1$ (ℓ main tokens and one query). We will also introduce a new parameter k that controls the number of regression task vectors \mathbf{w} that define the labels in these sequences. The phenomenology surrounding k is not discussed in this work for space reasons, but we include it to best match previous work on this model for consistency (Letey et al., 2026; Lu et al., 2025; Raventós et al., 2023). We sample k task vectors

$$\mathbf{w}_j \sim_{\text{i.i.d.}} \mathcal{N}(\mathbf{0}, I_d)$$

and these k vectors are distributed evenly across the n sequences. This serves to limit the “task diversity” of the training batch, allowing us to judge how much of the true task distribution $\mathcal{N}(\mathbf{0}, I_d)$ the model is truly learning in-context.

We consider the same asymptotic treatment as derived in Lu et al. (2025), namely $d \rightarrow \infty$ with

$$\alpha := \frac{\ell}{d}, \quad \tau := \frac{n}{d^2}, \quad \kappa := \frac{k}{d}$$

all taken to be $\Theta(1)$ as $d \rightarrow \infty$.

An important quantity to define is the Stieltjes transform of the task sample covariance. This will be given by

$$\mathcal{M}_\kappa(z) := \lim_{d \rightarrow \infty, k = \kappa d} \frac{1}{d} \operatorname{tr} \left(\left(\frac{1}{k} \sum_{j=1}^k \mathbf{w}_j \mathbf{w}_j^\top + z I_d \right)^{-1} \right) \quad (8)$$

We will also use

$$\mathcal{M}'_\kappa(z) := \frac{d}{dz} \mathcal{M}_\kappa(z).$$

These quantities will appear throughout our formulas.

The explicit ridge parameter λ will be modulated by the fact that we have finite samples, and will appear in our formula through an *effective* ridge defined implicitly by

$$\tilde{\lambda} \mathcal{M}_\kappa \left(\tilde{\lambda} + \frac{\rho_2}{\alpha} \right) + \frac{\lambda \tau}{\tilde{\lambda}} = 1 - \tau$$

We will sometimes write $\mathcal{M} := \mathcal{M}_\kappa(\sigma)$, $\mathcal{M}' := \mathcal{M}'_\kappa(\sigma)$ for shorthand, where

$$\sigma = \tilde{\lambda} + \frac{\rho_2}{\alpha}, \quad \tilde{\sigma} = \sigma - \frac{k_0}{\alpha}.$$

Finally, we are only considering the case where the sequential correlations defined by $K, \mathbf{k}_{\text{test}}$ are sufficiently weak, *i.e.*,

$$k_2 := \frac{1}{\ell} \operatorname{tr}(K^2), \quad k_0 := \mathbf{k}_{\text{test}}^\top \mathbf{k}_{\text{test}}, \quad k_1 := \mathbf{k}_{\text{test}}^\top K \mathbf{k}_{\text{test}}$$

are all $\Theta(1)$ as $d, \ell \rightarrow \infty$. These terms will appear in the constants

$$\rho_1 \equiv 1 + \rho, \quad \rho_2 \equiv k_2 + \rho, \quad \phi_1 \equiv \frac{1}{\alpha^2} (k_1 + \rho k_0), \quad \phi_2 \equiv \frac{1}{\alpha} k_0.$$

We present the main formula below. All proofs are given in subsequent sections.

Proposition 3 *Consider isotropic tokens and tasks, *i.e.*, $\Sigma = I_d = C$, with sequential correlations in the tokens given by PSD $K \in \mathbb{R}^{\ell \times \ell}$ and $\mathbf{k}_{\text{test}} \in \mathbb{R}^\ell$ for $\mathbf{k}_{\text{test}}^\top K^{-1} \mathbf{k}_{\text{test}} \leq 1$. We then have that ICL error for a model defined by Γ^* concentrates as $d \rightarrow \infty$ in the above limit, *i.e.*,*

$$\mathcal{E}_{\text{ICL}}(\Gamma^*) \simeq e_{\text{ICL}}^{\text{corr}}(\alpha, \kappa, \tau, \rho, k_0, k_1, k_2)$$

for

$$e_{\text{ICL}}^{\text{corr}}(\alpha, \kappa, \tau, \rho, k_0, k_1, k_2) = e_{\text{ICL}}^{\text{uncorr}} \left(\frac{1 + \rho}{k_2 + \rho} \alpha, \kappa, \tau, \rho \right) + e_{\text{query}}(\alpha, \kappa, \tau, \rho, k_0, k_1, k_2)$$

where

$$\begin{aligned} e_{\text{query}} &= -2\phi_2 + q_{\text{query}}(\mathcal{M} + \tilde{\lambda} \mathcal{M}') + \phi_2(2\sigma - \phi_2) \mathcal{M}' + 2\phi_2 \tilde{\sigma} \mathcal{M} \\ &\quad + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2 - 2\mathbf{m}_2^\top S \mathbf{m}_3 + 2(1 + \phi_2) \mathbf{m}_1^\top S \mathbf{m}_2 + (\phi_1 + 2\phi_2) \left(1 - \tilde{\sigma} \mathcal{M} - \mathbf{m}_1^\top S \mathbf{m}_2 \right)^2 \\ &\quad + \frac{\rho_2}{\alpha} \left[q_{\text{query}}(\mathcal{M} + \tilde{\lambda} \mathcal{M}') + \phi_2(2\sigma - \phi_2) \mathcal{M}' + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2 - 2\mathbf{m}_2^\top S \mathbf{m}_3 + 2\phi_2 \mathcal{M} \right]. \end{aligned}$$

for

$$\begin{aligned} \mathbf{m}_1 &= \begin{bmatrix} \mathcal{M} \\ \phi_2(1 - \sigma \mathcal{M}) \end{bmatrix}, \quad \mathbf{m}_2 = \begin{bmatrix} 1 - \tilde{\sigma} \mathcal{M} \\ \phi_2(1 - \tilde{\sigma} + \sigma \tilde{\sigma} \mathcal{M}) \end{bmatrix}, \quad \mathbf{m}_3 = \begin{bmatrix} \mathcal{M} + \tilde{\sigma} \mathcal{M}' \\ \phi_2(1 - \sigma \mathcal{M} - \tilde{\sigma} \mathcal{M} - \sigma \tilde{\sigma} \mathcal{M}') \end{bmatrix} \\ M &= \begin{bmatrix} -\mathcal{M}' & \phi_2(\mathcal{M} + \sigma \mathcal{M}') \\ \phi_2(\mathcal{M} + \sigma \mathcal{M}') & \phi_2^2(1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}') \end{bmatrix}, \quad S = \begin{bmatrix} \mathcal{M} & 1 + \phi_2(1 - \sigma \mathcal{M}) \\ 1 + \phi_2(1 - \sigma \mathcal{M}) & -\phi_1 + \phi_2^2(1 - \sigma + \sigma^2 \mathcal{M}) \end{bmatrix}^{-1} \end{aligned}$$

$$q_{\text{query}} = \frac{1}{\tau - (1 - 2\tilde{\lambda}\mathcal{M} - \tilde{\lambda}^2\mathcal{M}')} \left((2\frac{k_0}{\alpha} (\tilde{\lambda}(\mathcal{M} + \sigma\mathcal{M}') + (1 - 2\sigma\mathcal{M})) + \frac{k_0^2}{\alpha^2} (\mathcal{M} - \tilde{\lambda}\mathcal{M}')) \right. \\ \left. + \mathbf{m}_2^\top S \mathbf{m}_2 - \tilde{\lambda} (-2\mathbf{m}_2^\top S \mathbf{m}_3 + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2) \right)$$

Appendix F. Proofs of Population Averages

We begin for necessary y_{test}, H averages needed for terms in the population loss formulas in Lemmas 1 and 2. *A language model was used at various stages of these proofs to help with the Wick expansions of the 6th order moments in terms of traces. All expressions were subsequently carefully verified by the authors both analytically and numerically.*

Lemma 4 *Suppose we have $k_2, k_0, k_1 = \Theta(1)$ i.e., weak correlations. Write $\mathbf{v} = \Sigma \mathbf{w}$. Then we have*

$$\mathbb{E}[y_{\text{test}}^2] = \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho \\ \mathbb{E}[y_{\text{test}} H] \approx \frac{1}{d} [\mathbf{v} \mathbf{v}^\top + \frac{1}{\alpha} k_0 \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] \Sigma \quad (\text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho) \mathbf{v}] \\ \mathbb{E}[\text{vec}(H) \text{vec}(H)^\top] \approx \frac{1}{d} \Sigma \otimes \begin{bmatrix} \mathbf{v} \mathbf{v}^\top + (k_2 \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho) \Sigma / \alpha & (\text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho) \mathbf{v} \\ (\text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho) \mathbf{v}^\top & (\text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho)^2 \end{bmatrix} + \frac{1}{d^2} \mathcal{X}$$

for $d(d+1) \times d(d+1)$ tensor \mathcal{X} given by

$$\mathcal{X} = \frac{1}{\alpha^2} (k_1 \text{tr}[\Sigma W] + \rho k_0) \text{vec}([\Sigma \mathbf{0}]) \text{vec}([\Sigma \mathbf{0}])^\top \\ + \frac{k_0}{\alpha} (\text{vec}([\Sigma \mathbf{0}]) \text{vec}([S (\text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho) \mathbf{v}])^\top \\ + \text{vec}([S (\text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho) \mathbf{v}]) \text{vec}([\Sigma \mathbf{0}])^\top).$$

The “ \approx ” means subleading terms in ℓ, d are neglected.

Proof Let’s go term by term.

Label-label term. Conditioning on X , we have

$$\mathbb{E}[\mathbf{x}_{\text{test}} \mathbf{x}_{\text{test}}^\top | X] = \frac{1}{d} (1 - k_{-1}) \Sigma + (X^\top \mathbf{m})(X^\top \mathbf{m})^\top$$

and so

$$\mathbb{E}[y_{\text{test}}^2] = \mathbb{E}[\text{tr}(\mathbf{x}_{\text{test}} \mathbf{x}_{\text{test}}^\top \mathbf{w} \mathbf{w}^\top) + \rho] \\ = \frac{1}{d} (1 - k_{-1}) \mathbb{E}_X [\text{tr}(\Sigma \mathbf{w} \mathbf{w}^\top)] + \mathbb{E}_X [\text{tr}((X^\top \mathbf{m})(X^\top \mathbf{m})^\top \mathbf{w} \mathbf{w}^\top)] + \rho \\ = (1 - k_{-1} + k_{-1}) \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho.$$

Label-feature term. For

$$\mathbf{b} = X^\top (X \mathbf{w} + \boldsymbol{\varepsilon}), \quad \mathbf{c} = (X \mathbf{w} + \boldsymbol{\varepsilon})^\top (X \mathbf{w} + \boldsymbol{\varepsilon})$$

we have

$$y_{\text{test}} H = \mathbf{x}_{\text{test}} (\mathbf{x}_{\text{test}}^\top \mathbf{w} + \boldsymbol{\varepsilon}_{\text{test}}) \left[\frac{d}{\ell} \mathbf{b}^\top \quad \frac{1}{\ell} \mathbf{c} \right]$$

and so

$$\mathbb{E}[y_{\text{test}} H] = \frac{1}{d} (1 - k_{-1}) \Sigma \mathbf{w} \left[\frac{d}{\ell} \mathbb{E}_X [\mathbf{b}] \right]^\top + \left[\frac{d}{\ell} \mathbb{E}_X [(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \mathbf{w} \mathbf{b}^\top] \quad \frac{1}{\ell} \mathbb{E}_X [(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \mathbf{w} \mathbf{c}] \right]$$

Have

$$\begin{aligned}
 \frac{d}{\ell} \mathbb{E}_X[\mathbf{b}] &= \text{tr}[K] \mathbf{v} \\
 \frac{1}{\ell} \mathbb{E}_X[c] &= \text{tr}[K] \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho \\
 \frac{d}{\ell} \mathbb{E}_X[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \mathbf{w} \mathbf{b}^\top] &= \frac{1}{d} \cdot \frac{1}{\ell} (k_{-1} \text{tr}(K) + k_0) \mathbf{v} \mathbf{v}^\top + \frac{1}{d} \cdot \frac{1}{\ell} k_0 \text{tr}(\Sigma \mathbf{w} \mathbf{w}^\top) \Sigma \\
 &\approx \frac{1}{d} \cdot \left(k_{-1} \text{tr}[K] \mathbf{v} \mathbf{v}^\top + \frac{1}{\alpha} k_0 \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] \Sigma \right) \\
 \frac{1}{\ell} \mathbb{E}_X[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \mathbf{w} c] &= \frac{1}{d} \cdot \rho k_{-1} \mathbf{v} + \frac{1}{d} \cdot k_{-1} \text{tr}[K] \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] \mathbf{v} + \frac{2}{d^2} \frac{k_0}{\ell} \mathbf{v} \mathbf{v}^\top \mathbf{w} \\
 &\approx \frac{1}{d} \cdot k_{-1} \left(\text{tr}[K] \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] + \rho \right) \mathbf{v}
 \end{aligned}$$

Combining gives

$$\mathbb{E}[y_{\text{test}} H] = \frac{1}{d} \left[\text{tr}[K] \mathbf{v} \mathbf{v}^\top + \frac{1}{\alpha} k_0 \text{tr}[\Sigma \mathbf{w} \mathbf{w}^\top] \Sigma - \rho \mathbf{1} \mathbf{v} \right]$$

Feature-feature term Using \mathbf{b}, c as above we can write more easily

$$\text{Vec}(H) \text{Vec}(H)^\top = (\mathbf{x}_{\text{test}} \mathbf{x}_{\text{test}}^\top) \otimes \left(\begin{bmatrix} \frac{d}{\ell} \mathbf{b} \\ \frac{1}{\ell} c \end{bmatrix} \begin{bmatrix} \frac{d}{\ell} \mathbf{b} \\ \frac{1}{\ell} c \end{bmatrix}^\top \right)$$

Taking conditional expectation over \mathbf{x}_{test} for fixed X , our expression simplifies as

$$\mathbb{E}[\text{Vec}(H) \text{Vec}(H)^\top] = \frac{1}{d} (1 - k_{-1}) \Sigma \otimes \mathbb{E}_{X, \epsilon}[\mathbf{u} \mathbf{u}^\top] + \mathbb{E}_{X, \epsilon} \left[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes (\mathbf{u} \mathbf{u}^\top) \right].$$

where I'm writing $\mathbf{u}^\top = \left[\frac{d}{\ell} \mathbf{b}^\top \quad \frac{1}{\ell} c \right]$ for convenience. Now let's start with the $\mathbb{E}[\mathbf{u} \mathbf{u}^\top]$ term.

We have

$$\begin{aligned}
 \mathbb{E}[\mathbf{b} \mathbf{b}^\top] &= \mathbb{E}[X^\top X \mathbf{w} \mathbf{w}^\top X^\top X] + \mathbb{E}[X^\top \epsilon \epsilon^\top X] \\
 &= \frac{1}{d^2} \left(\text{tr}(K^2) + \text{tr}(K)^2 \right) \Sigma \mathbf{w} \mathbf{w}^\top \Sigma + \frac{1}{d^2} \text{tr}(K^2) \text{tr}(\Sigma \mathbf{w} \mathbf{w}^\top) \Sigma + \frac{1}{d} \rho \text{tr}(K) \Sigma \\
 \mathbb{E}[c \mathbf{b}] &= \mathbb{E}[\mathbf{w}^\top X^\top X \mathbf{w} X^\top X \mathbf{w}] + 2 \mathbb{E}[\epsilon^\top X \mathbf{w} X^\top \epsilon] + \mathbb{E}[\epsilon^\top \epsilon X^\top X \mathbf{w}] \\
 &= \left(\frac{1}{d^2} \left(\text{tr}(K)^2 + 2 \text{tr}(K^2) \right) \text{tr}(\Sigma \mathbf{w} \mathbf{w}^\top) + \frac{1}{d} \rho (\ell + 2) \text{tr}(K) \right) \Sigma \mathbf{w} \\
 \mathbb{E}[c^2] &= \mathbb{E}[\mathbf{w}^\top X^\top X \mathbf{w} \mathbf{w}^\top X^\top X \mathbf{w}] + 2 \mathbb{E}[\mathbf{w}^\top X^\top X \mathbf{w} \epsilon^\top \epsilon] + 4 \mathbb{E}[\mathbf{w}^\top X^\top \epsilon \epsilon^\top X \mathbf{w}] + \mathbb{E}[\epsilon^\top \epsilon \epsilon^\top \epsilon] \\
 &= \frac{1}{d^2} \left(\text{tr}(K)^2 + 2 \text{tr}(K^2) \right) \text{tr}(\Sigma \mathbf{w} \mathbf{w}^\top)^2 + \frac{1}{d} \rho (2\ell + 4) \text{tr}(K) \text{tr}(\Sigma \mathbf{w} \mathbf{w}^\top) + \rho^2 (\ell^2 + 2\ell).
 \end{aligned}$$

This simplifies as

$$\begin{aligned}
 \frac{d^2}{\ell^2} \mathbb{E}[\mathbf{b} \mathbf{b}^\top] &\approx \text{tr}[K]^2 (\Sigma \mathbf{w})(\Sigma \mathbf{w})^\top + \frac{1}{\alpha} (\rho \text{tr}[K] + \text{tr}[K^2] \text{tr}[\Sigma W]) \Sigma \\
 \frac{d}{\ell^2} \mathbb{E}[c \mathbf{b}] &\approx \text{tr}[K] (\text{tr}[K] \text{tr}[\Sigma W] + \rho) \Sigma \mathbf{w} \\
 \frac{1}{\ell^2} \mathbb{E}[c^2] &\approx (\rho + \text{tr}[K] \text{tr}[\Sigma W])^2
 \end{aligned}$$

upon making the above high-dimensional assumptions.

For the higher order term,

$$\mathbb{E} \left[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes (\mathbf{u} \mathbf{u}^\top) \right] = \left[\frac{d^2}{\ell^2} \mathbb{E}[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes \mathbf{b} \mathbf{b}^\top] \quad \frac{d}{\ell^2} \mathbb{E}[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes c \mathbf{b}] \right] \\ \left[\frac{d}{\ell^2} \mathbb{E}[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes c \mathbf{b}^\top] \quad \frac{1}{\ell^2} \mathbb{E}[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes c^2] \right]$$

Write

$$\mathbf{v} := \Sigma \mathbf{w}, \quad S := \mathbf{v} \mathbf{v}^\top, \quad k_{-1} \equiv \text{tr}(KM) = \mathbf{k}_{\text{test}}^\top K^{-1} \mathbf{k}_{\text{test}}, \quad k_0 \equiv \mathbf{k}_{\text{test}}^\top \mathbf{k}_{\text{test}}, \quad k_1 \equiv \mathbf{k}_{\text{test}}^\top K \mathbf{k}_{\text{test}}$$

For large ℓ, d and the chosen high-dimensional assumptions, the Isserlis expansion gives

$$\begin{aligned} \frac{d^2}{\ell^2} \mathbb{E} \left[((X^\top m)(X^\top m)^\top) \otimes (\mathbf{b} \mathbf{b}^\top) \right]_{(I,i),(J,j)} &\approx \frac{1}{d} k_{-1} \left(\Sigma_{IJ} S_{ij} \text{tr}[K]^2 + \frac{1}{\alpha} (\text{tr}[\Sigma W] \text{tr}[K^2] + \rho \text{tr}[K]) \Sigma_{IJ} \Sigma_{ij} \right) \\ &+ \frac{1}{d^2} \mathcal{X}_{(I,i)(J,j)} \end{aligned}$$

with

$$\begin{aligned} \mathcal{X}_{(I,i)(J,j)} &\equiv \frac{1}{\alpha^2} k_1 \text{tr}[\Sigma W] \Sigma_{Ii} \Sigma_{Jj} + \frac{1}{\alpha} \text{tr}[K] k_0 (\Sigma_{Ii} S_{Jj} + \Sigma_{Jj} S_{Ii}) + \rho \frac{1}{\alpha^2} k_0 \Sigma_{Ii} \Sigma_{Jj} \\ &= \frac{1}{\alpha^2} (k_1 \text{tr}[\Sigma W] + \rho k_0) \Sigma_{Ii} \Sigma_{Jj} + \frac{1}{\alpha} k_0 \text{tr}[K] (\Sigma_{Ii} S_{Jj} + \Sigma_{Jj} S_{Ii}). \end{aligned}$$

Now for the second term,

$$\begin{aligned} \mathbb{E}[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes \mathbf{c} \mathbf{b}]_{(I,i),J} &= \mathbb{E}[c(X^\top \mathbf{m})_I (X^\top \mathbf{m})_J b_i] \\ &= \mathbb{E}[(X_{sk} w_k + \varepsilon_s)(X_{sn} w_n + \varepsilon_s) X_{tI} X_{pJ} M_{tp} X_{qi} (X_{qj} w_j + \varepsilon_q)] \\ &= M_{tp} W_{kn} w_j \mathbb{E}[X_{sk} X_{sn} X_{tI} X_{pJ} X_{qi} X_{qj}] + (\ell + 2) \rho M_{tp} w_j \mathbb{E}[X_{tI} X_{pJ} X_{qi} X_{qj}] \end{aligned}$$

Expanding all the Wick terms gives

$$\frac{d}{\ell^2} \mathbb{E} \left[((X^\top m)(X^\top m)^\top) \otimes (\mathbf{c} \mathbf{b}) \right]_{(I,i),J} \approx \frac{1}{d} k_{-1} (\text{tr}[\Sigma W] + \rho) \Sigma_{IJ} v_i + \frac{1}{d^2} \mathcal{X}_{(I,i)(J,d+1)}$$

where

$$\mathcal{X}_{(I,i)(J,d+1)} \equiv \frac{1}{\alpha} k_0 (\text{tr}[\Sigma W] + \rho) \Sigma_{Ii} v_J.$$

Finally, the last term gives

$$\frac{1}{\ell^2} \mathbb{E} \left[((X^\top m)(X^\top m)) \otimes c^2 \right]_{IJ} \approx \frac{1}{d} k_{-1} (\text{tr}[\Sigma W] + \rho)^2 \Sigma_{IJ}$$

Gathering everything together gives the required formula. \blacksquare

Lemma 5 *We have*

$$\mathbb{E}[y_{\text{test}}^2] = \text{tr}[\Sigma W] + \rho \tag{9}$$

$$\mathbb{E}[y_{\text{test}} H] \approx [\tilde{k}_0 \text{tr}[\Sigma W] \Sigma \quad \mathbf{0}] \tag{10}$$

$$\mathbb{E}[\text{vec}(H) \text{vec}(H)^\top] \approx \begin{bmatrix} \text{tr}[\Sigma W] \left(\tilde{k}_2 \Sigma \otimes \Sigma + \tilde{k}_1 \Sigma \tilde{\otimes} \Sigma + \tilde{k}_1 \text{vec}(\Sigma) \text{vec}(\Sigma)^\top \right) & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix} \tag{11}$$

where I'm using $\tilde{\otimes}$ to refer to the transposition+Kronecker operation, i.e., for $A, B \in \mathbb{R}^{d \times d}$, have

$$(A \otimes B)_{(I,i),(J,j)} = A_{IJ} B_{ij}, \quad (A \tilde{\otimes} B)_{(I,i),(J,j)} = A_{Ij} B_{Ji}.$$

Proof Let's go term by term.

Label-label term. As before $\mathbb{E}[y_{\text{test}}^2] = \text{tr}[\Sigma W] + \rho$.

Label-feature term. Have

$$\begin{aligned}\frac{d}{\ell}\mathbb{E}_X[\mathbf{b}] &= \mathbf{v} \\ \frac{1}{\ell}\mathbb{E}_X[c] &= \text{tr}[\Sigma W] + \rho \\ \frac{d}{\ell}\mathbb{E}_X[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \mathbf{w} \mathbf{b}^\top] &= \frac{1}{d} \left(k_{-1} + \tilde{k}_0 \right) S + \tilde{k}_0 \text{tr}[\Sigma W] \Sigma \\ \frac{1}{\ell}\mathbb{E}_X[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \mathbf{w} c] &= \frac{1}{d} \cdot k_{-1} (\text{tr}[\Sigma W] + \rho) \mathbf{v} + \frac{2}{d^2} \tilde{k}_0 S \mathbf{w}\end{aligned}$$

and so

$$\mathbb{E}[y_{\text{test}} H] \approx [\tilde{k}_0 \text{tr}[\Sigma W] \Sigma \quad \mathbf{0}]$$

Feature-feature term Again

$$\mathbb{E}[\text{Vec}(H)\text{Vec}(H)^\top] = \frac{1}{d}(1 - k_{-1})\Sigma \otimes \mathbb{E}_{X,\epsilon}[\mathbf{u}\mathbf{u}^\top] + \mathbb{E}_{X,\epsilon}[(X^\top \mathbf{m})(X^\top \mathbf{m})^\top \otimes (\mathbf{u}\mathbf{u}^\top)].$$

where $\mathbf{u}^\top = [\frac{d}{\ell}\mathbf{b}^\top \quad \frac{1}{\ell}c]$. The order-4 terms are

$$\begin{aligned}\frac{d^2}{\ell^2}\mathbb{E}[\mathbf{b}\mathbf{b}^\top] &= (1 + \tilde{k}_2)S + d\tilde{k}_2 \text{tr}[\Sigma W] \Sigma + \frac{1}{\alpha}\rho \Sigma \\ \frac{d}{\ell^2}\mathbb{E}[\mathbf{c}\mathbf{b}] &= \left((1 + 2\tilde{k}_2) \text{tr}[\Sigma W] + \rho \right) \mathbf{v} \\ \frac{1}{\ell^2}\mathbb{E}[c^2] &= (1 + 2\tilde{k}_2) \text{tr}[\Sigma W]^2 + 2\rho \text{tr}[\Sigma W] + \rho^2.\end{aligned}$$

Expanding the order-6 terms, the first is

$$\begin{aligned}\frac{d^2}{\ell^2}\mathbb{E}\left[\left((X^\top \mathbf{m})(X^\top \mathbf{m})^\top\right) \otimes (\mathbf{b}\mathbf{b}^\top)\right]_{(I,i),(J,j)} &= \frac{1}{d}\Sigma_{IJ} S_{ij} k_{-1} (1 + \tilde{k}_2) + \Sigma_{IJ} \Sigma_{ij} k_{-1} \tilde{k}_2 \text{tr}[\Sigma W] \\ &\quad + \tilde{k}_1 \text{tr}[\Sigma W] \left(\Sigma_{Ii} \Sigma_{Jj} + \Sigma_{Ij} \Sigma_{Ji} \right) \\ &\quad + \frac{1}{d}(\tilde{k}_0 + \tilde{k}_1) \left(\Sigma_{Ii} S_{Jj} + \Sigma_{Ji} S_{Ij} + \Sigma_{Jj} S_{Ii} + \Sigma_{Ij} S_{Ji} \right) \\ &\quad + 2\frac{1}{d}\tilde{k}_1 \Sigma_{ij} S_{IJ} \\ &\quad + \frac{1}{d}\frac{1}{\alpha}\rho \left(\Sigma_{IJ} \Sigma_{ij} k_{-1} + \tilde{k}_0 (\Sigma_{Ii} \Sigma_{Jj} + \Sigma_{Ij} \Sigma_{Ji}) \right)\end{aligned}$$

Now for the second term,

$$\begin{aligned}\frac{d}{\ell^2}\mathbb{E}\left[\left((X^\top \mathbf{m})(X^\top \mathbf{m})^\top\right) \otimes (\mathbf{c}\mathbf{b})\right]_{(I,i),J} &= \frac{1}{d} \cdot \text{tr}[\Sigma W] \Sigma_{IJ} v_i k_{-1} (1 + 2\tilde{k}_2) \\ &\quad + \frac{1}{d} \cdot \text{tr}[\Sigma W] (\Sigma_{Ii} v_J + \Sigma_{Ji} v_I) (\tilde{k}_0 + 2\tilde{k}_1) \\ &\quad + \frac{1}{d^2} (v_I v_J v_i) (2\tilde{k}_0 + 4\tilde{k}_1) \\ &\quad + \frac{1}{d} \cdot \rho \left(k_{-1} \Sigma_{IJ} v_i + \tilde{k}_0 (\Sigma_{Ii} v_J + \Sigma_{Ji} v_I) \right)\end{aligned}$$

Finally, the last term

$$\begin{aligned} \frac{1}{\ell^2} \mathbb{E} \left[((X^\top m)(X^\top m)) \otimes c^2 \right]_{IJ} &= \frac{1}{d} \cdot \left[\Sigma_{IJ} \operatorname{tr}[\Sigma W]^2 k_{-1}(1 + 2\tilde{k}_2) + \frac{1}{d} \operatorname{tr}[\Sigma W] v_I v_J (4\tilde{k}_0 + 8\tilde{k}_1) \right] \\ &\quad + \frac{2}{d} \rho \cdot \left[\Sigma_{IJ} \operatorname{tr}[\Sigma W] k_{-1} + \frac{2}{d} \tilde{k}_0 v_I v_J \right] + \frac{1}{d} \cdot \rho^2 \Sigma_{IJ} k_{-1}. \end{aligned}$$

Thus we have

$$\approx (1 - k_{-1}) \Sigma \otimes \begin{bmatrix} \tilde{k}_2 \operatorname{tr}[\Sigma W] \Sigma & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{bmatrix} + \operatorname{tr}[\Sigma W] \begin{bmatrix} k_{-1} \tilde{k}_2 \Sigma \otimes \Sigma + \tilde{k}_1 \Sigma \tilde{\otimes} \Sigma + \tilde{k}_1 \operatorname{vec}(\Sigma) \operatorname{vec}(\Sigma)^\top & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}$$

and the result follows. ■

Appendix G. Proof of Proposition 3

This entire section is an extended calculation that serves as a complete proof of the deterministic formula for $\mathcal{E}_{\text{ICL}}(\Gamma^*)$ given in Proposition 3. The methodology in this section is to analyse the random object Γ^* , which contains randomness through the particular training batch $(H^\mu, y_{\text{test}}^\mu)_{\mu=1}^n$, as a resolvent of a larger random matrix, and the necessary $\text{tr}(\Gamma A^\top)$ and $\text{tr}(\Gamma^\top B \Gamma)$ terms as traces of a resolvent against deterministic test matrices. Throughout this entire section we will take $\Sigma = I_d$.

This section will be relatively light on detail as the full calculation setup can be found in both Lu et al. (2025) and Letey et al. (2026). For a more rigorous analysis, with error terms bounded properly, see Lu et al. (2025); here we will simply write \approx when terms can be dropped due to the negligibility in high dimensions.

We care about the MSE loss (computed by Theorem 1) for the optimal parameters

$$\text{vec}(\Gamma^*) = \left(\frac{n}{d} \lambda I + \sum_{\mu=1}^n \text{vec}(H^\mu) \text{vec}(H^\mu)^\top \right)^{-1} \sum_{\mu=1}^n y_{\text{test}}^\mu \text{vec}(H^\mu) \quad (12)$$

for $\mu \in [n]$ denoting the index for the n training sample sequences. Define

$$\mathbf{z}_\mu = \begin{bmatrix} y_{\text{test}}^\mu / d \\ \text{vec}(H^\mu) / \sqrt{d} \end{bmatrix} \in \mathbb{R}^{d(d+1)+1}$$

and construct extended resolvent

$$G_{\text{ext}}(\pi) = \frac{1}{\sum_{\mu \in [n]} \mathbf{z}_\mu \mathbf{z}_\mu^\top + \pi B_{\text{ext}} + \tau \lambda I}. \quad (13)$$

for some test matrix $B_{\text{ext}} \in \mathbb{R}^{(d(d+1)+1) \times (d(d+1)+1)}$.

It's important to note that even though we're studying sequentially-correlated data, the vectors \mathbf{z}_μ s are still independent of each other. We can thus apply the ‘‘leave-one-out’’ or cavity method here over the μ index. We end up with

$$\sum_{\mu \in [n]} \frac{1}{1 + \mathbf{z}_\mu^\top G_{\text{ext}}^{[\mu]} \mathbf{z}_\mu} G_{\text{ext}}^{[\mu]} \mathbf{z}_\mu \mathbf{z}_\mu^\top + G_{\text{ext}}(\pi B_{\text{ext}} + \tau \lambda I) = I. \quad (14)$$

For fixed task \mathbf{w}_μ we can apply Lemma 4 to compute

$$\mathbb{E}_{X, \mathbf{x}_{\text{test}}, \varepsilon} [\mathbf{z} \mathbf{z}^\top] = \frac{1}{d^2} \Upsilon(\mathbf{w})$$

for

$$\Upsilon(\mathbf{w}) \equiv \begin{bmatrix} \text{tr}[\mathbf{w} \mathbf{w}^\top] + \rho & \frac{1}{\sqrt{d}} \text{vec}([\mathbf{w} \mathbf{w}^\top + \frac{1}{\alpha} k_0 \text{tr}[\mathbf{w} \mathbf{w}^\top] I_d \quad \rho_1(\mathbf{w}) \mathbf{w}]) \\ \frac{1}{\sqrt{d}} \text{vec}([\mathbf{w} \mathbf{w}^\top + \frac{1}{\alpha} k_0 \text{tr}[\mathbf{w} \mathbf{w}^\top] I_d \quad \rho_1(\mathbf{w}) \mathbf{w}])^\top & I_d \otimes E(\mathbf{w}) + \frac{1}{d} \mathcal{X}(\mathbf{w}) \end{bmatrix}$$

where

$$E(\mathbf{w}) \equiv \begin{bmatrix} \mathbf{w} \mathbf{w}^\top + \rho_2(\mathbf{w}) I_d / \alpha & \rho_1(\mathbf{w}) \mathbf{w} \\ \rho_1(\mathbf{w}) \mathbf{w}^\top & \rho_1(\mathbf{w})^2 \end{bmatrix}, \quad \rho_1(\mathbf{w}) = \text{tr}[\mathbf{w} \mathbf{w}^\top] + \rho, \quad \rho_2(\mathbf{w}) = \text{tr}[\mathbf{w} \mathbf{w}^\top] k_2 + \rho$$

and $d(d+1) \times d(d+1)$ tensor

$$\begin{aligned}\mathcal{X}(\mathbf{w}) &\equiv \frac{1}{\alpha^2} \left(\text{tr} \left[\mathbf{w} \mathbf{w}^\top \right] k_1 + \rho k_0 \right) \text{vec} \left([I_d \ \mathbf{0}] \right) \text{vec} \left([I_d \ \mathbf{0}] \right)^\top \\ &\quad + \frac{1}{\alpha} k_0 \left(\text{vec} \left([I_d \ \mathbf{0}] \right) \text{vec} \left([\mathbf{w} \mathbf{w}^\top \ \mathbf{0}] \right)^\top + \text{vec} \left([\mathbf{w} \mathbf{w}^\top \ \mathbf{0}] \right) \text{vec} \left([I_d \ \mathbf{0}] \right)^\top \right) \\ &\quad + \frac{1}{\alpha} \rho_1(\mathbf{w}) k_0 \left(\text{vec} \left([I_d \ \mathbf{0}] \right) \text{vec} \left(\mathbf{w} \mathbf{e}_{d+1}^\top \right)^\top + \text{vec} \left(\mathbf{w} \mathbf{e}_{d+1}^\top \right) \text{vec} \left([I_d \ \mathbf{0}] \right)^\top \right).\end{aligned}$$

The quadratic form $\mathbf{z}_\mu^\top G_{\text{ext}}^{[\mu]} \mathbf{z}_\mu$ concentrates for fixed tasks \mathbf{w}_μ by standard arguments, so we have

$$\mathbf{z}_\mu^\top G_{\text{ext}}^{[\mu]} \mathbf{z}_\mu \simeq \chi^\mu(\mathbf{w}_\mu) \quad (15)$$

where

$$\chi^\mu(\mathbf{w}_\mu) \equiv \frac{1}{d^2} \text{tr} \left([G_{\text{ext}}^\mu]_{\setminus 0} \cdot \left[I \otimes E(\mathbf{w}_\mu) + \frac{1}{d} \mathcal{X}(\mathbf{w}_\mu) \right] \right). \quad (16)$$

Replacing $\mathbf{z}_\mu^\top G_{\text{ext}}^{[\mu]} \mathbf{z}_\mu$ in (14) with $\chi^\mu(\mathbf{w}_\mu)$ gives

$$\sum_{\mu \in [n]} \frac{1}{1 + \chi^\mu(\mathbf{w}_\mu)} G_{\text{ext}}^{[\mu]} \mathbf{z}_\mu \mathbf{z}_\mu^\top + G_{\text{ext}}(\pi B_{\text{ext}} + \tau \lambda I) \simeq I. \quad (17)$$

In this equation we will also replace $\mathbf{z}_\mu \mathbf{z}_\mu^\top$ with its conditional expectation over $X, \mathbf{x}_{\text{test}}, \boldsymbol{\varepsilon}$, giving

$$\frac{\tau}{n} \sum_{\mu \in [n]} \frac{1}{1 + \chi^\mu(\mathbf{w}_\mu)} G_{\text{ext}}^{[\mu]} \Upsilon(\mathbf{w}_\mu) + G_{\text{ext}}(\pi B_{\text{ext}} + \tau \lambda I) \simeq I \quad (18)$$

where recall $\tau = n/d^2$.

In high dimensions and for large n , there is negligible difference between $\sum_{\nu \neq \mu}$ and \sum_μ . Thus, we replace G_{ext}^μ by G_{ext} , and $\chi^\mu(\mathbf{w}_\mu)$ by

$$\chi(\mathbf{w}_\mu) \equiv \frac{1}{d^2} \text{tr} \left([G_{\text{ext}}]_{\setminus 0} \cdot \left[I \otimes E(\mathbf{w}_\mu) + \frac{1}{d} \mathcal{X}(\mathbf{w}_\mu) \right] \right). \quad (19)$$

So finally we have the expression for G_{ext}

$$G_{\text{ext}} \left(\frac{\tau}{n} \sum_{\mu \in [n]} \frac{1}{1 + \chi(\mathbf{w}_\mu)} \Upsilon(\mathbf{w}_\mu) + \pi B_{\text{ext}} + \tau \lambda I \right) \simeq I. \quad (20)$$

Exploit finiteness of training task set. So far we are summing over n task vectors, but really only n/k of these are unique. Thus, we can simplify (20) as

$$G_{\text{ext}} \left(\frac{\tau}{k} \sum_{j \in [k]} \frac{1}{1 + \chi(\mathbf{w}_j)} \Upsilon(\mathbf{w}_j) + \pi B_{\text{ext}} + \tau \lambda I \right) \simeq I. \quad (21)$$

We replace $\chi(\mathbf{w}_j)$, which is self-averaging in \mathbf{w}_j , with its mean

$$\hat{\chi}_{\text{ave}} \equiv \frac{1}{k} \sum_{j \in [k]} \chi(\mathbf{w}_j). \quad (22)$$

To clean up the sums over the tasks $\mathbf{w}_1, \dots, \mathbf{w}_k$ we use the that

$$\begin{aligned} \frac{1}{k} \sum_{j \in [k]} \text{tr} [\mathbf{w}_j \mathbf{w}_j^\top] &= \text{tr}[R_{\text{tr}}] \\ \frac{1}{k} \sum_{j \in [k]} \left(\rho + \text{tr} [\mathbf{w}_j \mathbf{w}_j^\top] \right) \mathbf{w}_j &\simeq (\rho + \text{tr}[R_{\text{tr}}]) \mathbf{b}_{\text{tr}} \\ \frac{1}{k} \sum_{j \in [k]} \left(\rho + \text{tr} [\mathbf{w}_j \mathbf{w}_j^\top] \right)^2 &\simeq (\rho + \text{tr}[R_{\text{tr}}])^2 \end{aligned}$$

for

$$\mathbf{b}_{\text{tr}} \equiv \frac{1}{k} \sum_{j \in [k]} \mathbf{w}_j, \quad R_{\text{tr}} \equiv \frac{1}{k} \sum_{j \in [k]} \mathbf{w}_j \mathbf{w}_j^\top.$$

We thus have that

$$\begin{aligned} \frac{1}{k} \sum_{j \in [k]} E(\mathbf{w}_j) &\approx B_{\text{tr}} \equiv \begin{bmatrix} R_{\text{tr}} + \rho_2 I_d / \alpha & \rho_1 \mathbf{b}_{\text{tr}} \\ \rho_1 \mathbf{b}_{\text{tr}}^\top & \rho_1^2 \end{bmatrix} \\ \frac{1}{k} \sum_{j \in [k]} \mathcal{X}(\mathbf{w}_j) &\approx \Phi \end{aligned}$$

for

$$\begin{aligned} \Phi &\equiv \frac{1}{\alpha^2} (\text{tr}[R_{\text{tr}}] k_1 + \rho k_0) \text{vec} \left(\begin{bmatrix} I_d & \mathbf{0} \end{bmatrix} \right) \text{vec} \left(\begin{bmatrix} I_d & \mathbf{0} \end{bmatrix} \right)^\top \\ &\quad + \frac{1}{\alpha} k_0 \left(\text{vec} \left(\begin{bmatrix} I_d & \mathbf{0} \end{bmatrix} \right) \text{vec} \left(\begin{bmatrix} R_{\text{tr}} & \mathbf{0} \end{bmatrix} \right)^\top + \text{vec} \left(\begin{bmatrix} R_{\text{tr}} & \mathbf{0} \end{bmatrix} \right) \text{vec} \left(\begin{bmatrix} I_d & \mathbf{0} \end{bmatrix} \right)^\top \right) \\ &\quad + \frac{1}{\alpha} \rho_1 k_0 \left(\text{vec} \left(\begin{bmatrix} I_d & \mathbf{0} \end{bmatrix} \right) \text{vec} \left(\mathbf{b}_{\text{tr}} \mathbf{e}_{d+1}^\top \right)^\top + \text{vec} \left(\mathbf{b}_{\text{tr}} \mathbf{e}_{d+1}^\top \right) \text{vec} \left(\begin{bmatrix} I_d & \mathbf{0} \end{bmatrix} \right)^\top \right) \end{aligned}$$

where

$$\rho_1 = \text{tr}[R_{\text{tr}}] + \rho, \quad \rho_2 = \text{tr}[K^2] \text{tr}[R_{\text{tr}}] + \rho.$$

Finally, we have

$$\hat{\chi}_{\text{ave}} = \frac{1}{d^2} \text{tr} \left([G_{\text{ext}}]_{\setminus 0} \cdot \left[I \otimes B_{\text{tr}} + \frac{1}{d} \Phi \right] \right). \quad (23)$$

Thus, after averaging over $X, \mathbf{x}_{\text{test}}, \boldsymbol{\varepsilon}$ in the extended resolvent G_{ext} , we have a deterministic equivalent $G_{\text{ext}} \simeq \mathcal{G}_{\text{ext}}$ (still depending on random task quantities $R_{\text{tr}}, \mathbf{b}_{\text{tr}}$) defined by self-consistent equations

$$[\mathcal{G}_{\text{ext}}]_{\setminus 0} = \left(\frac{\tau}{1 + \chi_\pi} I_d \otimes B_{\text{tr}} + \frac{1}{d} \frac{\tau}{1 + \chi_\pi} \Phi + \pi \Pi + \tau \lambda I_d \otimes I_{d+1} \right)^{-1} \quad (24)$$

$$\chi_\pi = \frac{1}{d^2} \text{tr} \left([\mathcal{G}_{\text{ext}}]_{\setminus 0} \left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi \right) \right) \quad (25)$$

where the full matrix is given by

$$\begin{aligned} \mathcal{G}_{\text{ext}}(\pi)^{-1} &\equiv \frac{\tau}{1 + \chi_\pi} \begin{bmatrix} \rho_1 & & & \\ \frac{1}{\sqrt{d}} \text{vec} \left(\begin{bmatrix} R_{\text{tr}} + \frac{1}{\alpha} k_0 \text{tr}[R_{\text{tr}}] I_d & \rho_1 \mathbf{b}_{\text{tr}} \end{bmatrix} \right) & & & \\ & & & \\ & & & \frac{1}{\sqrt{d}} \text{vec} \left(\begin{bmatrix} R_{\text{tr}} + \frac{1}{\alpha} k_0 \text{tr}[R_{\text{tr}}] I_d & \rho_1 \mathbf{b}_{\text{tr}} \end{bmatrix} \right)^\top \\ & & & I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi \end{bmatrix} \\ &\quad + \pi B_{\text{ext}} + \tau \lambda I \end{aligned} \quad (26)$$

Note that up until this point, the results have matched the previous analysis from Lu et al. (2025) and Letey et al. (2026) with the exception of ρ_2 (depending on k_2) and the low-rank term Φ (depending on k_0 and k_1).

The low-rank term is given by

$$\begin{aligned}\Phi &= \phi_1 \mathbf{1}_+ \mathbf{1}_+^\top + \phi_2 (\mathbf{1}_+ r^\top + r \mathbf{1}_+^\top) + \phi_3 (\mathbf{1}_+ \mu^\top + \mu \mathbf{1}_+^\top) \\ &= [\mathbf{1}_+ \quad v] \begin{bmatrix} \phi_1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{1}_+^\top \\ v^\top \end{bmatrix} \\ &= T \phi T^\top\end{aligned}$$

for

$$\mathbf{1}_+ \equiv \text{vec} \left(\begin{bmatrix} I_d & \mathbf{0} \end{bmatrix} \right), \quad r \equiv \text{vec} \left(\begin{bmatrix} R_{\text{tr}} & \mathbf{0} \end{bmatrix} \right), \quad \mu \equiv \text{vec} \left(\mathbf{b}_{\text{tr}} \mathbf{e}_{d+1}^\top \right), \quad v \equiv \text{vec} \left(\begin{bmatrix} \phi_2 R_{\text{tr}} & \phi_3 \mathbf{b}_{\text{tr}} \end{bmatrix} \right)$$

and

$$\phi_1 \equiv \frac{1}{\alpha^2} (\text{tr}[R_{\text{tr}}] k_1 + \rho k_0), \quad \phi_2 \equiv \frac{1}{\alpha} k_0, \quad \phi_3 \equiv \frac{1}{\alpha} (\text{tr}[R_{\text{tr}}] + \rho) k_0.$$

G.1. Effective ridge

The effective ridge will be $\tilde{\lambda} = \lambda(1 + \chi_0)$ where χ_0 is defined as the solution to the implicit equations

$$[\mathcal{G}_{\text{ext}}]_{\setminus 0} = \left(\frac{\tau}{1 + \chi_0} I_d \otimes B_{\text{tr}} + \frac{1}{d} \frac{\tau}{1 + \chi_0} \Phi + \tau \lambda I_d \otimes I_{d+1} \right)^{-1} \quad (27)$$

$$\chi_0 = \frac{1}{d^2} \text{tr} \left([\mathcal{G}_{\text{ext}}]_{\setminus 0} \left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi \right) \right). \quad (28)$$

We can recast the above implicit equations using Woodbury, as

$$\frac{\tau \chi_0}{1 + \chi_0} = \frac{1}{d^2} \text{tr} \left(\left(I_d \otimes F_0 - \frac{1}{d} \Phi_C \right) \left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi \right) \right) \quad (29)$$

where

$$F_0 = (B_{\text{tr}} + \lambda(1 + \chi_0) I_{d+1})^{-1}$$

$$\Phi_C = (I_d \otimes F_0) [\mathbf{1}_+ \quad v] S \begin{bmatrix} \mathbf{1}_+^\top \\ v^\top \end{bmatrix} (I_d \otimes F_0) \quad \text{for some } 2 \times 2 \text{ mtx } S \text{ with } \Theta(1) \text{ elements.}$$

However, all the $1/d$ terms in (29) are subleading, e.g.

$$\frac{1}{d^3} \text{tr} \left((I_d \otimes F_0) r r^\top (I_d \otimes F_0 B_{\text{tr}}) \right) = \frac{1}{d^3} \text{tr} (R_{\text{tr}} F_0 B_{\text{tr}} F_0 R_{\text{tr}}) = \mathcal{O} \left(\frac{1}{d^2} \right)$$

and so we write

$$\frac{\tau \chi_0}{1 + \chi_0} \approx \frac{1}{d^2} \text{tr} (I_d \otimes F_0 B_{\text{tr}}) = \text{tr} \left[(B_{\text{tr}} + \tilde{\lambda} I_{d+1})^{-1} B_{\text{tr}} \right]. \quad (30)$$

This can be simplified as

$$\begin{aligned}\frac{\tau \chi_0}{1 + \chi_0} &= 1 - \lambda(1 + \chi_0) \text{tr} \left[(B_{\text{tr}} + \lambda(1 + \chi_0) I_{d+1})^{-1} \right] \\ &\approx 1 - \lambda(1 + \chi_0) \text{tr} \left[\left(R_{\text{tr}} + \left(\lambda(1 + \chi_0) + \frac{\rho_2}{\alpha} \right) I_d \right)^{-1} \right]\end{aligned}$$

after ignoring the \mathbf{b}_{tr} terms in B_{tr} . This is where the Stieltjes transform of R_{tr} is first introduced:

$$\mathcal{M}_\kappa(w) = \lim_{d \rightarrow \infty, k \rightarrow \infty, k/d = \kappa} \frac{1}{d} \text{tr} \left((R_{\text{tr}} + w I_d)^{-1} \right).$$

Using this, we find the effective ridge self-consistency equation

$$\tilde{\lambda} \mathcal{M}_\kappa \left(\tilde{\lambda} + \frac{\rho_2}{\alpha} \right) + \frac{\lambda \tau}{\tilde{\lambda}} = 1 - \tau \quad (31)$$

G.2. Relating Γ^* to G

Recall that

$$\text{vec}(\Gamma^*) = \left(\frac{n}{d} \lambda I + \sum_{\mu=1}^n \text{vec}(H^\mu) \text{vec}(H^\mu)^\top \right)^{-1} \sum_{\mu=1}^n y_{\text{test}}^\mu \text{vec}(H^\mu)$$

and so we have

$$G_{\text{ext}}(0) = d \begin{bmatrix} \frac{1}{d} \sum_{\mu} (y_{\text{test}}^\mu)^2 + \frac{n}{d} \lambda & \frac{1}{\sqrt{d}} \sum_{\mu} y_{\text{test}}^\mu \text{vec}(H^\mu)^\top \\ \frac{1}{\sqrt{d}} \sum_{\mu} y_{\text{test}}^\mu \text{vec}(H^\mu) & \sum_{\mu} \text{vec}(H^\mu) \text{vec}(H^\mu)^\top + \frac{n}{d} \lambda I \end{bmatrix}^{-1}$$

Using that

$$\begin{bmatrix} a & \mathbf{b}^\top \\ \mathbf{b} & D \end{bmatrix}^{-1} = \begin{bmatrix} c & -c\mathbf{q}^\top \\ -c\mathbf{q} & D^{-1} + c\mathbf{q}\mathbf{q}^\top \end{bmatrix}, \quad c = \frac{1}{a - \mathbf{b}^\top D^{-1} \mathbf{b}}, \quad \mathbf{q} = D^{-1} \mathbf{b}$$

we have a formula for Γ^* in terms of $G_{\text{ext}}(0)$ as

$$\frac{\sqrt{d}}{\mathbf{e}_1^\top G_{\text{ext}}(0) \mathbf{e}_1} G_{\text{ext}}(0) \mathbf{e}_1 = \begin{bmatrix} \sqrt{d} \\ \text{vec}(\Gamma^*) \end{bmatrix}$$

Using our above formulas, we also have

$$\frac{\sqrt{d}}{\mathbf{e}_1^\top \mathcal{G}_{\text{ext}}(0) \mathbf{e}_1} \mathcal{G}_{\text{ext}}(0) \mathbf{e}_1 = \left[\left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi + \tilde{\lambda} I_{d(d+1)} \right)^{-1} \text{vec} \left([R_{\text{tr}} + \frac{k_0}{\alpha} \text{tr}[R_{\text{tr}}] I_d \quad \rho_1 \mathbf{b}_{\text{tr}}] \right) \right]$$

and so

$$\text{vec}(\Gamma^*) \simeq \text{vec}(\Gamma_{\text{de}}) \equiv \left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi + \tilde{\lambda} I_{d(d+1)} \right)^{-1} \text{vec} \left([R_{\text{tr}} + \frac{k_0}{\alpha} \text{tr}[R_{\text{tr}}] I_d \quad \rho_1 \mathbf{b}_{\text{tr}}] \right) = F_\Phi \mathbf{g}$$

for

$$F_\Phi = \left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi + \tilde{\lambda} I_{d(d+1)} \right)^{-1} \\ \mathbf{g} = \text{vec} \left([R_{\text{tr}} + \frac{k_0}{\alpha} \text{tr}[R_{\text{tr}}] I_d \quad \rho_1 \mathbf{b}_{\text{tr}}] \right).$$

This immediately gives the linear term in the MSE error expression as

$$\text{tr}[\Gamma^* A^\top] = \frac{1}{d} \text{tr}(\Gamma^* A^\top) = \frac{1}{d} \text{tr}(\text{vec}(\Gamma^*) \text{vec}(A)^\top) \simeq \frac{1}{d} \text{tr}(\text{vec}(\Gamma_{\text{de}}) \text{vec}(A)^\top)$$

where we will specifically use $\text{vec}(A) = (1 + \phi_2) \mathbf{1}_+$ from Lemma 1.

For the quadratic terms we need to be a bit more careful, as the correct object to work with for high-dimensional equivalence is technically the components of \mathcal{G} and not this linear representation of $\text{vec}(\Gamma_{\text{de}})$. This is what the parameterisation πB_{ext} is for. Take

$$B_{\text{ext}} = \begin{bmatrix} 0 & 0 \\ 0 & \Pi \end{bmatrix}$$

for some $\Pi \in \mathbb{R}^{d(d+1) \times d(d+1)}$. Then

$$\frac{d}{d\pi} \frac{1}{c(\pi)} (\pi = 0) = \frac{1}{d} \text{vec}(\Gamma^*)^\top \Pi \text{vec}(\Gamma^*)$$

where $c(\pi) = \mathbf{e}_1^\top G_{\text{ext}}(\pi) \mathbf{e}_1$. We can safely replace

$$c(\pi) = \mathbf{e}_1^\top G_{\text{ext}}(\pi) \mathbf{e}_1 \simeq \mathbf{e}_1 \mathcal{G}_{\text{ext}}(\pi) \mathbf{e}_1.$$

By Schur complement on $\mathcal{G}_{\text{ext}}(\pi)$ we have

$$\begin{aligned} \frac{1}{\mathbf{e}_1 \mathcal{G}_{\text{ext}}(\pi) \mathbf{e}_1} &= \frac{\tau}{1 + \chi_\pi} \rho_1 + \tau \lambda - \frac{1}{d} \frac{\tau^2}{(1 + \chi_\pi)^2} \mathbf{g}^\top \left(\frac{\tau}{1 + \chi_\pi} \left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi \right) + \pi \Pi + \tau \lambda I \right)^{-1} \mathbf{g} \\ &= \frac{\tau}{1 + \chi_\pi} \rho_1 + \tau \lambda - \frac{1}{d} \frac{\tau}{1 + \chi_\pi} \mathbf{g}^\top \left(I_d \otimes B_{\text{tr}} + \frac{1}{d} \Phi + \pi \frac{1 + \chi_\pi}{\tau} \Pi + \lambda(1 + \chi_\pi) I \right)^{-1} \mathbf{g}. \end{aligned}$$

Given the eventual MSE term we want from Lemma 1, we choose

$$\Pi = I_d \otimes B + \frac{1}{d} \Psi$$

for

$$\Psi = [1_+ \quad \text{vec}([I_d \quad \mathbf{0}])] \begin{bmatrix} \phi_1 & \phi_2 \\ \phi_2 & 0 \end{bmatrix} \begin{bmatrix} 1_+^\top \\ \text{vec}([I_d \quad \mathbf{0}]^\top) \end{bmatrix} = (\phi_1 + 2\phi_2) 1_+ 1_+^\top.$$

Using the same approximation as before, we will take

$$[\mathcal{G}_{\text{ext}}]_{\setminus 0} = I_d \otimes \left(\frac{\tau}{1 + \chi_\pi} B_{\text{tr}} + \pi B + \tau \lambda I \right)^{-1} + \frac{1}{d} \text{low rank terms from } \Phi \text{ and } \Psi$$

and approximate χ_π as

$$\chi_\pi = \text{tr} \left[\left(\frac{\tau}{1 + \chi_\pi} B_{\text{tr}} + \pi B + \tau \lambda I \right)^{-1} B_{\text{tr}} \right].$$

As before, we find that

$$\frac{\tau \chi'_0}{(1 + \chi_0)^2} = \frac{\text{tr}[F_0 B F_0 B_{\text{tr}}]}{\text{tr}[F_0 B_{\text{tr}} F_0 B_{\text{tr}}] - \tau}, \quad F_0 = (B_{\text{tr}} + \tilde{\lambda} I)^{-1}.$$

Thus,

$$\frac{1}{d} \text{vec}(\Gamma^*)^\top \Pi \text{vec}(\Gamma^*) \simeq \frac{1}{d} \text{vec}(\Gamma_{\text{de}})^\top \Pi \text{vec}(\Gamma_{\text{de}}) - \frac{\tau \chi'_0}{(1 + \chi_0)^2} \left(\rho_1 - \frac{1}{d} \mathbf{g}^\top \text{vec}(\Gamma_{\text{de}}) - \tilde{\lambda} \frac{1}{d} \text{vec}(\Gamma_{\text{de}})^\top \text{vec}(\Gamma_{\text{de}}) \right).$$

Note that

$$\begin{aligned} \text{vec}(\Gamma_{\text{de}}) &= (I_d \otimes F_0) \mathbf{g} - \frac{1}{d} \Phi_C \mathbf{g} \\ \Pi &= I_d \otimes B_{\text{test}} + \frac{1}{d} \Psi \end{aligned}$$

where

$$\begin{aligned} \Phi_C &= ((I_d \otimes F_0) T) \begin{bmatrix} \mathcal{M} & 1 + \phi_2(1 - \sigma \mathcal{M}) \\ 1 + \phi_2(1 - \sigma \mathcal{M}) & -\phi_1 + \phi_2^2(1 - \sigma + \sigma^2 \mathcal{M}) \end{bmatrix}^{-1} ((I_d \otimes F_0) T)^\top \\ &= (I_d \otimes F_0) [1_+ \quad \phi_2 r + \phi_3 \mu] S \begin{bmatrix} 1_+^\top \\ \phi_2 r^\top + \phi_3 \mu^\top \end{bmatrix} (I_d \otimes F_0) \end{aligned}$$

This comes from expanding

$$\begin{aligned}
 F_0 &= \begin{bmatrix} F & \mathbf{f} \\ \mathbf{f}^\top & f \end{bmatrix} \\
 \frac{1}{d} \begin{bmatrix} 1_+^\top \\ r^\top \\ \mu^\top \end{bmatrix} [1_+ \quad r \quad \mu] &= \begin{bmatrix} \text{tr}[F] & \text{tr}[FR_{\text{tr}}] & \text{tr}[\mathbf{f}\mathbf{b}_{\text{tr}}^\top] \\ \text{tr}[FR_{\text{tr}}] & \text{tr}[FR_{\text{tr}}^2] & \text{tr}[F\mathbf{f}\mathbf{b}_{\text{tr}}^\top] \\ \text{tr}[\mathbf{f}\mathbf{b}_{\text{tr}}^\top] & \text{tr}[F\mathbf{f}\mathbf{b}_{\text{tr}}^\top] & f \text{tr}[\mathbf{b}_{\text{tr}}\mathbf{b}_{\text{tr}}^\top] \end{bmatrix} \\
 &\simeq \begin{bmatrix} \mathcal{M}_\kappa(\sigma) & 1 - \sigma\mathcal{M}_\kappa(\sigma) & 0 \\ 1 - \sigma\mathcal{M}_\kappa(\sigma) & 1 - \sigma + \sigma^2\mathcal{M}_\kappa(\sigma) & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad \sigma \equiv \frac{\rho_2}{\alpha} + \tilde{\lambda}.
 \end{aligned}$$

As we can see, all the $\mathbf{b}_{\text{tr}}, \mathbf{f}$ contributions are negligible compared to the resolvent $(R_{\text{tr}} + \sigma I_d)^{-1}$ contributions. Using similar intuition, in the expansions of the linear and quadratic error terms, we will neglect the $\rho_1 \mathbf{b}_{\text{tr}}$ component of \mathbf{g} , as well as $\mu = \text{vec}(\mathbf{b}_{\text{tr}} e_{d+1}^\top)$ in Φ_C .

Linear error term Following this, and writing the matrix component of \mathbf{g} as

$$\tilde{R} = R_{\text{tr}} + \frac{k_0}{\alpha} \text{tr}[R_{\text{tr}}] I_d$$

we have

$$\begin{aligned}
 \frac{1}{d} \text{vec}(A)^\top \text{vec}(\Gamma_{\text{de}}) &= (1 + \phi_2) \frac{1}{d} \left(1_+^\top (I \otimes F) \mathbf{g} - \frac{1}{d} 1_+^\top \Phi_C \mathbf{g} \right) \\
 &\simeq (1 + \phi_2) \left(\text{tr}[F\tilde{R}] - [\text{tr}[F] \quad \phi_2 \text{tr}[FR_{\text{tr}}]] S \begin{bmatrix} \text{tr}[F\tilde{R}] \\ \phi_2 \text{tr}[R_{\text{tr}}F\tilde{R}] \end{bmatrix} \right)
 \end{aligned}$$

Quadratic error term Recall

$$B = \begin{bmatrix} (1 + \frac{\rho_2}{\alpha}) I_d & \mathbf{0} \\ \mathbf{0}^\top & \rho_1^2 \end{bmatrix}, \quad \Psi = (\phi_1 + 2\phi_2) 1_+ 1_+^\top.$$

Have

$$\begin{aligned}
 \frac{1}{d} \mathbf{g}^\top \text{vec}(\Gamma_{\text{de}}) &\approx \text{tr}[\tilde{R}F\tilde{R}] - [\text{tr}[\tilde{R}F] \quad \phi_2 \text{tr}[\tilde{R}FR_{\text{tr}}]] S \begin{bmatrix} \text{tr}[F\tilde{R}] \\ \phi_2 \text{tr}[R_{\text{tr}}F\tilde{R}] \end{bmatrix} \\
 \frac{1}{d} \text{vec}(\Gamma_{\text{de}})^\top \text{vec}(\Gamma_{\text{de}}) &\approx \text{tr}[\tilde{R}F^2\tilde{R}] - 2 [\text{tr}[\tilde{R}F] \quad \phi_2 \text{tr}[\tilde{R}FR_{\text{tr}}]] S \begin{bmatrix} \text{tr}[F^2\tilde{R}] \\ \phi_2 \text{tr}[R_{\text{tr}}F^2\tilde{R}] \end{bmatrix} \\
 &\quad + [\text{tr}[\tilde{R}F] \quad \phi_2 \text{tr}[\tilde{R}FR_{\text{tr}}]] S^\top \begin{bmatrix} \text{tr}[F^2] & \phi_2 \text{tr}[F^2R_{\text{tr}}] \\ \phi_2 \text{tr}[F^2R_{\text{tr}}] & \phi_2^2 \text{tr}[R_{\text{tr}}F^2R_{\text{tr}}] \end{bmatrix} S \begin{bmatrix} \text{tr}[F\tilde{R}] \\ \phi_2 \text{tr}[R_{\text{tr}}F\tilde{R}] \end{bmatrix} \\
 \frac{1}{d} \text{vec}(\Gamma_{\text{de}})^\top (I_d \otimes B) \text{vec}(\Gamma_{\text{de}}) &\approx \left(1 + \frac{\rho_2}{\alpha}\right) \cdot \frac{1}{d} \text{vec}(\Gamma_{\text{de}})^\top \text{vec}(\Gamma_{\text{de}}) \\
 \frac{1}{\phi_1 + 2\phi_2} \frac{1}{d^2} \text{vec}(\Gamma_{\text{de}})^\top \Psi \text{vec}(\Gamma_{\text{de}}) &\approx \text{tr}[\tilde{R}F]^2 - 2 \text{tr}[\tilde{R}F] [\text{tr}[F] \quad \phi_2 \text{tr}[R_{\text{tr}}F]] S \begin{bmatrix} \text{tr}[F\tilde{R}] \\ \phi_2 \text{tr}[R_{\text{tr}}F\tilde{R}] \end{bmatrix} \\
 &\quad + \left([\text{tr}[F] \quad \phi_2 \text{tr}[FR_{\text{tr}}]] S \begin{bmatrix} \text{tr}[F\tilde{R}] \\ \phi_2 \text{tr}[R_{\text{tr}}F\tilde{R}] \end{bmatrix} \right)^2.
 \end{aligned}$$

The final step is to remember that $F = (R_{\text{tr}} + \sigma)^{-1}$ and $\tilde{R} = R_{\text{tr}} + k_0 I_d / \alpha$ and so we have the following dictionary of terms

$$\begin{aligned}
\text{tr}[F] &\simeq \mathcal{M}_\kappa(\sigma) \\
\text{tr}[FR_{\text{tr}}] &\simeq 1 - \sigma \mathcal{M}_\kappa(\sigma) \\
\text{tr}[F\tilde{R}] &\simeq 1 - \sigma \mathcal{M}_\kappa(\sigma) + \frac{k_0}{\alpha} \mathcal{M}_\kappa(\sigma) \\
\text{tr}[FR_{\text{tr}}\tilde{R}] &\simeq \frac{k_0}{\alpha} (1 - \sigma \mathcal{M}_\kappa(\sigma)) + 1 - \sigma + \sigma^2 \mathcal{M}_\kappa(\sigma) \\
\text{tr}[F\tilde{R}\tilde{R}] &\simeq \left(\frac{k_0}{\alpha}\right)^2 \mathcal{M}_\kappa(\sigma) + 2\frac{k_0}{\alpha} (1 - \sigma \mathcal{M}_\kappa(\sigma)) + 1 - \sigma + \sigma^2 \mathcal{M}_\kappa(\sigma) \\
\text{tr}[F^2] &\simeq -\mathcal{M}'_\kappa(\sigma) \\
\text{tr}[F^2 R_{\text{tr}}] &\simeq \mathcal{M}_\kappa(\sigma) + \sigma \mathcal{M}'_\kappa(\sigma) \\
\text{tr}[F^2 \tilde{R}] &\simeq \mathcal{M}_\kappa(\sigma) + \sigma \mathcal{M}'_\kappa(\sigma) - \frac{k_0}{\alpha} \mathcal{M}'_\kappa(\sigma) \\
\text{tr}[F^2 R_{\text{tr}} R_{\text{tr}}] &\simeq 1 - 2\sigma \mathcal{M}_\kappa(\sigma) - \sigma^2 \mathcal{M}'_\kappa(\sigma) \\
\text{tr}[F^2 R_{\text{tr}} \tilde{R}] &\simeq 1 - 2\sigma \mathcal{M}_\kappa(\sigma) - \sigma^2 \mathcal{M}'_\kappa(\sigma) + \frac{k_0}{\alpha} (\mathcal{M}_\kappa(\sigma) + \sigma \mathcal{M}'_\kappa(\sigma)) \\
\text{tr}[F^2 \tilde{R} \tilde{R}] &\simeq 1 - 2\sigma \mathcal{M}_\kappa(\sigma) - \sigma^2 \mathcal{M}'_\kappa(\sigma) + 2\frac{k_0}{\alpha} (\mathcal{M}_\kappa(\sigma) + \sigma \mathcal{M}'_\kappa(\sigma)) - \left(\frac{k_0}{\alpha}\right)^2 \mathcal{M}'_\kappa(\sigma)
\end{aligned}$$

So finally, if we use some shorthand

$$\tilde{\sigma} = \sigma - \frac{k_0}{\alpha}, \quad \mathbf{m}_1 = \begin{bmatrix} \mathcal{M} \\ \phi_2(1 - \sigma \mathcal{M}) \end{bmatrix}, \quad \mathbf{m}_2 = \begin{bmatrix} 1 - \tilde{\sigma} \mathcal{M} \\ \phi_2(1 - \tilde{\sigma} + \sigma \tilde{\sigma} \mathcal{M}) \end{bmatrix}, \quad \mathbf{m}_3 = \begin{bmatrix} \mathcal{M} + \tilde{\sigma} \mathcal{M}' \\ \phi_2(1 - \sigma \mathcal{M} - \tilde{\sigma} \mathcal{M} - \sigma \tilde{\sigma} \mathcal{M}') \end{bmatrix}$$

$$M = \begin{bmatrix} -\mathcal{M}' & \phi_2(\mathcal{M} + \sigma \mathcal{M}') \\ \phi_2(\mathcal{M} + \sigma \mathcal{M}') & \phi_2^2(1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}') \end{bmatrix}$$

and so

$$\begin{aligned}
\frac{1}{d} \text{vec}(A)^\top \text{vec}(\Gamma_{\text{de}}) &\simeq (1 + \phi_2) \left(1 - \tilde{\sigma} \mathcal{M} - \mathbf{m}_1^\top S \mathbf{m}_2\right) \\
\frac{1}{d} \mathbf{g}^\top \text{vec}(\Gamma_{\text{de}}) &\simeq 1 + \sigma - 2\tilde{\sigma} + \tilde{\sigma}^2 \mathcal{M} - \mathbf{m}_2^\top S \mathbf{m}_2 \\
\frac{1}{d} \text{vec}(\Gamma_{\text{de}})^\top \text{vec}(\Gamma_{\text{de}}) &\simeq 1 - 2\tilde{\sigma} \mathcal{M} - \tilde{\sigma}^2 \mathcal{M}' - 2\mathbf{m}_2^\top S \mathbf{m}_3 + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2 \\
\frac{1}{d^2} \text{vec}(\Gamma_{\text{de}})^\top \Psi \text{vec}(\Gamma_{\text{de}}) &\simeq (\phi_1 + 2\phi_2) \left((1 - \tilde{\sigma} \mathcal{M}) - \mathbf{m}_1^\top S \mathbf{m}_2 \right)^2
\end{aligned}$$

with

$$c = \frac{\tau \chi'_0}{(1 + \chi_0)^2} \simeq \left(1 + \frac{\rho_2}{\alpha}\right) \frac{\mathcal{M} + \tilde{\lambda} \mathcal{M}'}{1 - 2\tilde{\lambda} \mathcal{M} - \tilde{\lambda}^2 \mathcal{M}' - \tau}.$$

So finally, we have

$$\begin{aligned}
e_{\text{ICL}} &= (1 - c) \rho_1 - 2(1 + \phi_2) (1 - \tilde{\sigma} \mathcal{M} - \mathbf{m}_1^\top S \mathbf{m}_2) + (\phi_1 + 2\phi_2) (1 - \tilde{\sigma} \mathcal{M} - \mathbf{m}_1 S \mathbf{m}_2)^2 \\
&\quad + \left(1 + \frac{\rho_2}{\alpha} + c \tilde{\lambda}\right) (1 - 2\tilde{\sigma} \mathcal{M} - \tilde{\sigma}^2 \mathcal{M}' - 2\mathbf{m}_2^\top S \mathbf{m}_3 + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2) \\
&\quad + c(1 + \sigma - 2\tilde{\sigma} + \tilde{\sigma}^2 \mathcal{M} - \mathbf{m}_2^\top S \mathbf{m}_2)
\end{aligned}$$

for

$$\rho_1 \equiv 1 + \rho, \quad \rho_2 \equiv \text{tr}[K^2] + \rho$$

$$\begin{aligned}\sigma &= \tilde{\lambda} + \frac{\rho_2}{\alpha}, \quad \tilde{\sigma} = \sigma - \frac{k_0}{\alpha}, \quad \tilde{\lambda} \mathcal{M}_\kappa \left(\tilde{\lambda} + \frac{\rho_2}{\alpha} \right) + \frac{\lambda \tau}{\tilde{\lambda}} = 1 - \tau, \quad c = \left(1 + \frac{\rho_2}{\alpha} \right) \frac{\mathcal{M} + \tilde{\lambda} \mathcal{M}'}{1 - 2\tilde{\lambda} \mathcal{M} - \tilde{\lambda}^2 \mathcal{M}' - \tau} \\ \mathbf{m}_1 &= \begin{bmatrix} \mathcal{M} \\ \phi_2(1 - \sigma \mathcal{M}) \end{bmatrix}, \quad \mathbf{m}_2 = \begin{bmatrix} 1 - \tilde{\sigma} \mathcal{M} \\ \phi_2(1 - \tilde{\sigma} + \sigma \tilde{\sigma} \mathcal{M}) \end{bmatrix}, \quad \mathbf{m}_3 = \begin{bmatrix} \mathcal{M} + \tilde{\sigma} \mathcal{M}' \\ \phi_2(1 - \sigma \mathcal{M} - \tilde{\sigma} \mathcal{M} - \sigma \tilde{\sigma} \mathcal{M}') \end{bmatrix} \\ M &= \begin{bmatrix} -\mathcal{M}' & \phi_2(\mathcal{M} + \sigma \mathcal{M}') \\ \phi_2(\mathcal{M} + \sigma \mathcal{M}') & \phi_2^2(1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}') \end{bmatrix}, \quad S = \begin{bmatrix} \mathcal{M} & 1 + \phi_2(1 - \sigma \mathcal{M}) \\ 1 + \phi_2(1 - \sigma \mathcal{M}) & -\phi_1 + \phi_2^2(1 - \sigma + \sigma^2 \mathcal{M}) \end{bmatrix}^{-1} \\ \phi_1 &\equiv \frac{1}{\alpha^2}(k_1 + \rho k_0) = \phi_1, \quad \phi_2 \equiv \frac{1}{\alpha} k_0 = \phi_2\end{aligned}$$

Sanity check For $K = I_d$ have $\rho_1 = \rho_2$, $\sigma = \tilde{\sigma}$, $\phi_1 = \phi_2 = \phi_1 = \phi_2 = 0$, and all \mathbf{m}, S, M terms don't contribute. Get

$$\begin{aligned}e_{\text{ICL}} &= 1 + \rho - 2(1 - \sigma \mathcal{M}) + \left(1 + \frac{1 + \rho}{\alpha} \right) (1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}') \\ &\quad - c(\rho + \sigma - \sigma^2 \mathcal{M} - \tilde{\lambda}(1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}')) \\ &= 1 + \rho - 2(1 - \sigma \mathcal{M}) + \left(1 + \frac{1 + \rho}{\alpha} \right) (1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}') + \left(1 + \frac{1 + \rho}{\alpha} \right) (\mathcal{M} + \tilde{\lambda} \mathcal{M}') q \\ &= \rho + q \mathcal{M} + (\tilde{\lambda} q - \sigma^2) \mathcal{M}' + \frac{1 + \rho}{\alpha} (1 - (q - 2\sigma) \mathcal{M} + (\tilde{\lambda} q - \sigma^2) \mathcal{M}')\end{aligned}$$

as previously in [Letey et al. \(2026\)](#).

G.3. Dependence on query terms

This formula can be cleaned up a fair bit by gathering terms that depend on k_0, k_1 *i.e.*, the query correlation terms. First we write all the terms depending on c as

$$\begin{aligned}(-c) &\left(\rho_1 - \tilde{\lambda} \left(1 - 2\tilde{\sigma} \mathcal{M} - \tilde{\sigma}^2 \mathcal{M}' - 2\mathbf{m}_2^\top S \mathbf{m}_3 + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2 \right) - \left(1 + \sigma - 2\tilde{\sigma} + \tilde{\sigma}^2 \mathcal{M} - \mathbf{m}_2^\top S \mathbf{m}_2 \right) \right) \\ &= \left(1 + \frac{\rho_2}{\alpha} \right) \frac{\mathcal{M} + \tilde{\lambda} \mathcal{M}'}{\tau - (1 - 2\tilde{\lambda} \mathcal{M} - \tilde{\lambda}^2 \mathcal{M}')} \left(\rho + \sigma - \sigma^2 \mathcal{M} - \tilde{\lambda}(1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}') \right) \\ &\quad - \left(1 + \frac{\rho_2}{\alpha} \right) \frac{\mathcal{M} + \tilde{\lambda} \mathcal{M}'}{\tau - (1 - 2\tilde{\lambda} \mathcal{M} - \tilde{\lambda}^2 \mathcal{M}')} \left(2 \frac{k_0}{\alpha} \left(\tilde{\lambda}(\mathcal{M} + \sigma \mathcal{M}') + (1 - 2\sigma \mathcal{M}) \right) + \frac{k_0^2}{\alpha^2} (\mathcal{M} - \tilde{\lambda} \mathcal{M}') \right) \\ &\quad + \left(1 + \frac{\rho_2}{\alpha} \right) \frac{\mathcal{M} + \tilde{\lambda} \mathcal{M}'}{\tau - (1 - 2\tilde{\lambda} \mathcal{M} - \tilde{\lambda}^2 \mathcal{M}')} \left(\mathbf{m}_2^\top S \mathbf{m}_2 - \tilde{\lambda} \left(-2\mathbf{m}_2^\top S \mathbf{m}_3 + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2 \right) \right) \\ &= \left(1 + \frac{\rho_2}{\alpha} \right) (\mathcal{M} + \tilde{\lambda} \mathcal{M}') (q_{\text{old}} + q_{\text{query}})\end{aligned}$$

and the remaining ICL error terms can be expressed as

$$\begin{aligned}\rho_1 &- 2(1 - \sigma \mathcal{M}) + \left(1 + \frac{\rho_2}{\alpha} \right) (1 - 2\sigma \mathcal{M} - \sigma^2 \mathcal{M}') \\ &\quad - \frac{k_0}{\alpha} \left(2\mathcal{M} - 2 \left(1 + \frac{\rho_2}{\alpha} \right) (\mathcal{M} + \sigma \mathcal{M}') \right) - \frac{k_0^2}{\alpha^2} \left(1 + \frac{\rho_2}{\alpha} \right) \mathcal{M}' \\ &\quad - 2\phi_2(1 - \tilde{\sigma} \mathcal{M}) + 2(1 + \phi_2) \mathbf{m}_1^\top S \mathbf{m}_2 \\ &\quad + \left(1 + \frac{\rho_2}{\alpha} \right) (\mathbf{m}_2^\top S^\top M S \mathbf{m}_2 - 2\mathbf{m}_2^\top S \mathbf{m}_3) \\ &\quad + (\phi_1 + 2\phi_2) (1 - \tilde{\sigma} \mathcal{M} - \mathbf{m}_1^\top S \mathbf{m}_2)^2.\end{aligned}$$

Thus we can write

$$\begin{aligned}
e_{\text{ICL}}(k_0, k_1) &= e_{\text{ICL}}(\text{independent query}) - 2\phi_2 + q_{\text{query}}(\mathcal{M} + \tilde{\lambda}\mathcal{M}') + \phi_2(2\sigma - \phi_2)\mathcal{M}' \\
&\quad + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2 - 2\mathbf{m}_2^\top S \mathbf{m}_3 + 2\phi_2 \tilde{\sigma} \mathcal{M} + 2(1 + \phi_2) \mathbf{m}_1^\top S \mathbf{m}_2 \\
&\quad + (\phi_1 + 2\phi_2) \left(1 - \tilde{\sigma} \mathcal{M} - \mathbf{m}_1^\top S \mathbf{m}_2 \right)^2 \\
&\quad + \frac{\rho_2}{\alpha} \left[q_{\text{query}}(\mathcal{M} + \tilde{\lambda}\mathcal{M}') + \phi_2(2\sigma - \phi_2)\mathcal{M}' + \mathbf{m}_2^\top S^\top M S \mathbf{m}_2 - 2\mathbf{m}_2^\top S \mathbf{m}_3 + 2\phi_2 \mathcal{M} \right].
\end{aligned}$$