

# Optimal Attention Temperature Improves the Robustness of In-Context Learning under Distribution Shift in High Dimensions

author names withheld

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

## Abstract

Pretrained Transformers can perform in-context learning (ICL) from a few demonstrations, but this ability can fail sharply when the test distribution differs from pretraining—a common deployment setting. We study attention temperature as a simple inference-time control for improving ICL robustness under such shifts. In a high-dimensional linear-regression framework, we analyze a Transformer with “approximate softmax” attention, which preserves softmax’s normalization and temperature-dependent selectivity while remaining tractable. We derive a closed-form expression for the ICL generalization error under distribution shift, and show that it is minimized by an explicit optimal attention temperature. This characterization yields interpretable guidance by linking the best temperature to moments of the pre-softmax attention scores, and predicts when temperature adjustment can recover near Bayes-optimal performance. We validate the theory with extensive simulations, and further demonstrate gains on pretrained LLMs (GPT-2 and Llama2-7B) on question-answering benchmarks under distribution shift induced by noisy in-context demonstrations. Overall, attention temperature emerges as a principled, lightweight knob for improving the robustness of ICL in pretrained Transformers.

## 1. Introduction

Transformers [28] underpin contemporary AI systems, and a key driver of their success is *in-context learning* (ICL): adaptation to novel tasks directly from prompts, without gradient updates [4]. ICL has motivated extensive efforts to understand its mechanisms and scaling behavior [2, 30, 31, 34]. Yet ICL is fragile: even mild distributional shifts between pretraining and downstream tasks can sharply degrade performance [36], raising practical concerns about robustness and adaptability.

Within softmax self-attention, the *attention temperature*<sup>1</sup>  $\tau > 0$  rescales the pre-softmax scores, controlling their variance and hence the selectivity of attention. Although the original Transformer fixes  $\tau = \sqrt{d_k}$  (key dimension) [28], empirical work shows that adjusting  $\tau$  can substantially improve performance across NLP and vision benchmarks [5, 13, 16, 21, 37, 38]. To the best of our knowledge, however, its role in ICL under distributional shift remains largely unexplored<sup>2</sup>.

**This work** — We give a unified theoretical and empirical study of attention temperature for ICL under distribution shift, in the linear-regression framework that has been productive for analyzing ICL [9, 36]. Departing from prior work on linear attention, we analyze a Transformer with *ap-*

---

1. We distinguish this attention temperature from the LLM *sampling temperature* [25] that is used to adjust the output distribution of generative models.  
2. Related work is deferred to Appendix A.

*proximate softmax* attention — which preserves softmax’s temperature-dependent selectivity while remaining tractable — and derive a closed-form *optimal temperature*  $\tau_{\text{opt}}$  that minimizes the ICL generalization error. The expression depends explicitly on the type of shift, and setting the temperature  $\tau = \tau_{\text{opt}}$  recovers near-Bayes-optimal performance in our experiments. We also derive a moment-ratio heuristic that transfers the insight to standard softmax attention, which we evaluate on Llama2-7B.

### Contributions —

1. We derive, to our knowledge, the first closed-form expression for the *optimal attention temperature*  $\tau_{\text{opt}}$  that minimizes the ICL generalization error of a pretrained approximate-softmax Transformer (Theorems 2–3).
2. We characterize the effect of input-, task-, and noise-distribution shifts on ICL: the model is invariant to input-mean shifts but sensitive to input-covariance, task, and noise shifts, all of which  $\tau_{\text{opt}}$  partially mitigates.
3. We distill the analysis into a *moment-ratio heuristic* (9) that extends the optimal-temperature insight to standard softmax attention, and empirically validate it with Llama2-7B on SCIQ.

**Notation** — We follow Goodfellow et al. [10] for the notation. The spectral norm of a matrix  $M$  is  $\|M\|$  and the trace is  $\text{Tr}(M)$ . Matrix entries and slices are denoted by  $M_{i,j}$ ,  $M_{:,j}$ ,  $M_{i,:}$ .

## 2. Setting

We study ICL on linear-regression tasks: from context  $\{(\mathbf{x}_i, y_i)\}_{i=1}^{l-1}$  and query  $\mathbf{x}_l \in \mathbb{R}^d$ , the model predicts  $y_l$ , with each pair drawn i.i.d. as

$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x), y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma^2), \quad (1)$$

with task vector  $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$  fixed within a context and varying across tasks. We impose two assumptions; the first is weaker than the corresponding setup of Zhang et al. [36].

**Assumption 1 (Well-behaved data distributions)** *There exist  $c_1, c_2, c_3 > 0$  such that  $\|\boldsymbol{\mu}_x\|, \|\boldsymbol{\mu}_w\| \leq c_1$ ,  $\lambda_{\min}(\boldsymbol{\Sigma}_x), \lambda_{\min}(\boldsymbol{\Sigma}_w) \geq c_2$ , and  $\lambda_{\max}(\boldsymbol{\Sigma}_x), \lambda_{\max}(\boldsymbol{\Sigma}_w) \leq c_3$ .*

**Assumption 2 (High-dimensional regime)** *The context length  $l$  and input dimension  $d$  diverge jointly:  $l, d \rightarrow \infty$ .*

**Definition 1 (In-context learning)** *A model succeeds at ICL if its generalization error nearly matches that of the Bayes-optimal linear model (Appendix B).*

**Approximate softmax attention** — Following Zhang et al. [36], we encode the context as the embedding matrix

$$\mathbf{Z} := \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_{l-1} & \mathbf{x}_l \\ y_1 & \cdots & y_{l-1} & 0 \end{bmatrix} \in \mathbb{R}^{(d+1) \times l}, \quad (2)$$

whose last column is the query. We then analyze the approximate-softmax attention

$$\mathbf{E} := \mathbf{Z} + \mathbf{V} \mathbf{Z} \cdot \widehat{\text{softmax}} \left( \frac{(\mathbf{K} \mathbf{Z})^T (\mathbf{Q} \mathbf{Z})}{\tau} \right), \quad (3)$$

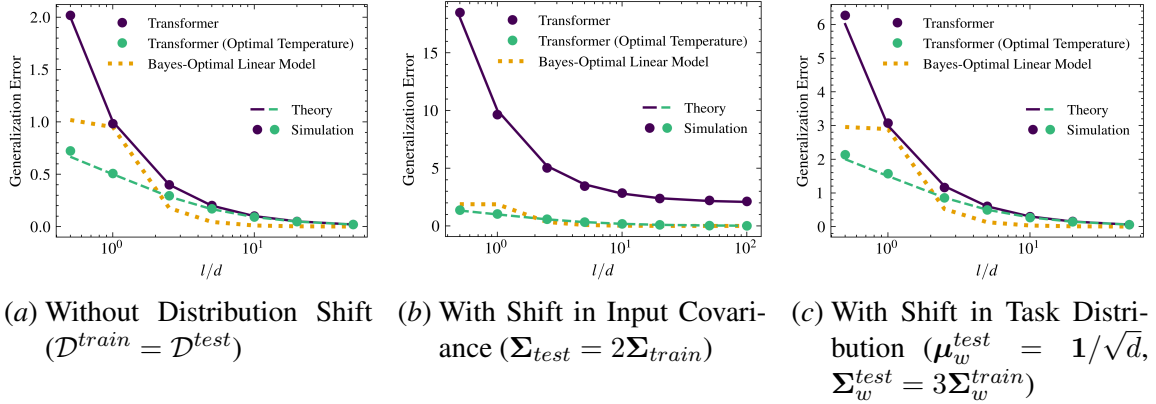


Figure 1: Experiments with Transformer (3) on ICL under distribution shifts. Parameters are set using (52) while the optimal temperature is by Theorem 3. Here,  $d = 50$ ,  $m = 5000$  (with a new task per sample),  $\sigma = 0.1$ ,  $\mu_x^{train} = \mu_w^{train} = \mathbf{0}$ , and  $\Sigma_x^{train} = \Sigma_w^{train} = I$ .

where  $\mathbf{K}, \mathbf{Q}, \mathbf{V}$  are key/query/value matrices,  $\tau$  is the temperature, and the explicit softmax form is derived in Appendix C. Unlike linear attention [36], softmax preserves row-wise normalization (so the model is robust to mean shifts; Remark 4) and mirrors softmax’s temperature behavior (Figure 3; Appendix D).

**Reparametrization and generalization error** — With  $\mathbf{M} := \mathbf{K}^\top \mathbf{Q}$ , the prediction  $\hat{y}(\mathbf{Z}; \mathbf{V}, \mathbf{M}) := E_{d+1,l}$  depends on  $(\mathbf{V}, \mathbf{M})$  only through their lower blocks:

$$\mathbf{V} = \begin{bmatrix} * & * \\ \mathbf{v}_{21}^\top & v_{22} \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} \mathbf{M}_{11} & * \\ \mathbf{m}_{21}^\top & * \end{bmatrix}, \quad (4)$$

where the entries marked  $*$  are irrelevant for predicting  $y_l$ . The ICL (generalization) error is

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) := \mathbb{E}_{(\mathbf{Z}, y_l) \sim \mathcal{D}^{test}} [(y_l - \hat{y}(\mathbf{Z}; \mathbf{V}, \mathbf{M}))^2], \quad (5)$$

evaluated on a test distribution  $\mathcal{D}^{test}$  whose task vectors were not seen during training.

### 3. Theoretical results

We characterize the ICL generalization error of the Transformer model (3) and identify the temperature that minimizes it. Throughout, Assumptions 1–2 are in force, together with a boundedness condition on the pretrained  $(\mathbf{M}, \mathbf{V})$  (Assumption 3, Appendix G).

**Theorem 2 (Generalization error for ICL)** *At test time, let the input, task, and noise distributions be  $\mathcal{N}(\mu_x, \Sigma_x)$ ,  $\mathcal{N}(\mu_w, \Sigma_w)$ , and  $\mathcal{N}(0, \sigma^2)$ ; let  $\mathbf{A} := \Sigma_x + \mu_x \mu_x^\top$  and  $\mathbf{B} := \Sigma_w + \mu_w \mu_w^\top$  denote the corresponding second-moment matrices. Then*

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) = \frac{1}{\tau^2} \text{Tr}(\mathbf{A} \mathbf{M}_{11}^\top \mathbf{F}_1 \mathbf{M}_{11}) - \frac{1}{\tau} \text{Tr}(\mathbf{A} (\mathbf{F}_2 \mathbf{M}_{11} + \mathbf{M}_{11}^\top \mathbf{F}_2^\top)) + \text{Tr}(\mathbf{A} \mathbf{B}) + \sigma^2, \quad (6)$$

where  $\mathbf{F}_1, \mathbf{F}_2$  depend on  $(\mathbf{V}, \Sigma_x, \Sigma_w, \sigma^2)$ ; their explicit forms and the proof (via Isserlis’ theorem [12]) are in Appendix G.

The temperature enters  $\mathcal{G}$  as  $1/\tau^2$  and  $1/\tau$ , so a single scalar adjustment can rebalance both terms. Minimizing  $\mathcal{G}$  over  $\tau$  yields:

**Theorem 3 (Optimal attention temperature)** *The generalization error is minimized at*

$$\tau_{\text{opt}} = \frac{2 \operatorname{Tr}(\mathbf{A} \mathbf{M}_{11}^T \mathbf{F}_1 \mathbf{M}_{11})}{\operatorname{Tr}(\mathbf{A}(\mathbf{F}_2 \mathbf{M}_{11} + \mathbf{M}_{11}^T \mathbf{F}_2^T))}, \quad (7)$$

provided both traces are positive (Appendix I).

When  $\tau_{\text{opt}} \neq 1$ , an unadjusted temperature is suboptimal. Without distribution shift, the Transformer with the pretrained parameters (Proposition 5 in Appendix F) emulates the Bayes-optimal linear model at  $\tau_{\text{opt}} = 1$ .

**Effect of distribution shift** — Under distribution shift ( $\mathcal{D}^{\text{test}} \neq \mathcal{D}^{\text{train}}$ ), the pretrained  $(\mathbf{M}, \mathbf{V})$  no longer match the test-time data, and three cases arise (Appendix I). (i) *Input shift*: centering renders the approximate model invariant to mean shifts, but a covariance mismatch ( $\Sigma_x^{\text{train}} \neq \Sigma_x^{\text{test}}$ ) drives the resulting estimator away from Bayes-optimality at every  $l$ , echoing prior results on linear attention [36]. (ii) *Task shift*:  $\mu_w^{\text{train}}$  and  $\Sigma_w^{\text{train}}$  are encoded in  $\mathbf{M}_{11}, \mathbf{v}_{21}$  via Proposition 5, but the model’s dependence on the task distribution decays with  $l$ , so task-shift error is most pronounced at small context lengths. (iii) *Noise shift*: a mismatch  $\sigma_{\text{train}}^2 \neq \sigma_{\text{test}}^2$  similarly moves the resulting estimator away from Bayes-optimality at small  $l$ , with the effect vanishing as  $l \rightarrow \infty$ . In the three sensitive cases (input-covariance, task, and noise shifts), applying  $\tau_{\text{opt}}$  from (7) minimizes the resulting generalization error and empirically recovers near-Bayes-optimal performance (Figure 1).

**Closed form under isotropic shift** — Consider isotropic training  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\epsilon_i \sim \mathcal{N}(0, \hat{\sigma}^2)$  and test  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, a\mathbf{I})$ ,  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, b\mathbf{I})$ ,  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  with shift parameters  $a, b, \sigma > 0$ . Then (7) reduces to

$$\tau_{\text{opt}} = \left( a + \frac{1}{l} \left( \frac{\sigma^2}{b} + ad \right) \right) \frac{\operatorname{Tr}(\mathbf{M}_{11}^T \mathbf{M}_{11})}{\operatorname{Tr}(\mathbf{M}_{11})}. \quad (8)$$

In particular,  $\tau_{\text{opt}} \rightarrow a$  as  $l \rightarrow \infty$ , since  $\operatorname{Tr}(\mathbf{M}_{11}^T \mathbf{M}_{11}) / \operatorname{Tr}(\mathbf{M}_{11}) \approx 1$  under the training distribution.

**Moment-ratio heuristic** — Rewriting (8) in terms of the moments of pre-softmax scores  $\mathbf{z}_i^\top \mathbf{M} \mathbf{z}_j$  (Appendix J) gives, for  $i \neq j$ ,

$$\tau_{\text{opt}} \approx \underbrace{\frac{\mathbb{E}[(\mathbf{z}_i^\top \mathbf{M} \mathbf{z}_j)^2]}{\mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \mathbf{z}_i]}}_{\text{moment-ratio}} + \underbrace{\frac{1}{l} \left( \frac{\sigma^2}{b} + ad \right)}_{\text{small-}l \text{ correction}}. \quad (9)$$

The moment-ratio is computable for any softmax attention layer, providing a recipe to transfer the optimal-temperature insight beyond approximate softmax — in particular, to large language models (cf. Figure 2).

## 4. Experimental results

We empirically validate our theory on (i) linear-regression ICL with the approximate-attention Transformer (3) and GPT-2 [23] (the latter in Appendix K), and (ii) Llama2-7B [27] on SCIQ [32].

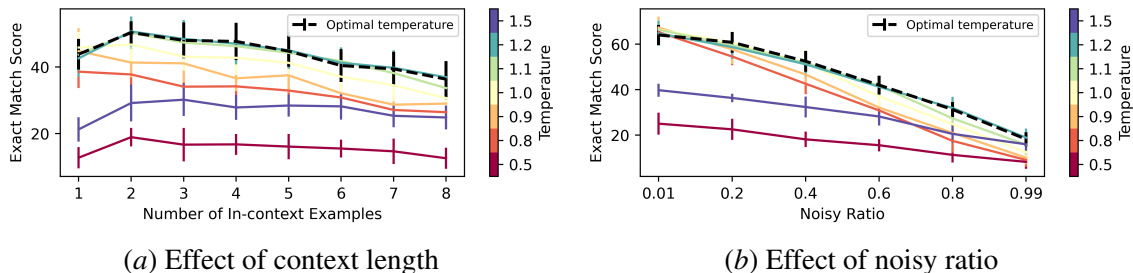


Figure 2: Attention-temperature effect on Llama2-7B [27] ICL on SCIQ [32] under noisy-label distribution shift [8]. Panel (a) fixes the noisy ratio at 0.6; panel (b) fixes the number of in-context examples at 6. The dashed black line marks the moment-ratio temperature (Eq. (9) and Appendix J) derived from Theorem 3. Error bars show one standard deviation over 12 Monte Carlo runs. The experimental details are deferred to Appendix K.

#### 4.1. Experiments on linear regression tasks

Figure 1 compares the approximate-attention Transformer to the Bayes-optimal linear model on the linear-regression setup (1). The pretrained model converges to Bayes-optimal as  $l$  grows (panel a) and recovers Bayes-optimal performance under input-covariance shift once  $\tau_{\text{opt}}$  is applied (panel b); task-shift error decays with  $l$  (panel c). Robustness to label-noise shift is reported in Appendix I (Figure 4), where temperature adjustment is critical at small  $l$  and  $\tau_{\text{opt}}$  increases with noise level. GPT-2 results (Appendix K) corroborate that the optimal-temperature insight transfers from approximate softmax to a full multi-head softmax Transformer with MLP layers.

#### 4.2. Experiments with LLMs for in-context Q&A

Since Theorem 3’s  $\tau_{\text{opt}}$  is specific to approximate softmax attention, we use the moment-ratio heuristic (9) as a transferable proxy for standard softmax. Following Gao et al. [8], we generate SCIQ ICL tasks with distribution shift induced by noisy labels and evaluate Llama2-7B via exact-match score (Appendix K). Figure 2 shows that the moment-ratio temperature consistently improves ICL across context lengths and noise ratios; higher noise ratios push the optimal temperature upward, matching our theoretical prediction.

### 5. Conclusion

This work provides a unified theoretical and empirical account of how attention temperature governs the in-context learning (ICL) performance of pretrained Transformers under distribution shift. Using a simplified yet expressive framework based on *approximate softmax attention*, we analytically show how shifts in input covariance and label noise degrade ICL and derive an *optimal temperature* that provably minimizes generalization error. Extensive experiments on synthetic regression tasks, GPT-2, and Llama-2 validate our predictions, demonstrating that temperature selection is not a mere heuristic but a principled mechanism for improving robustness. Taken together, our results advance the theoretical understanding of Transformer behavior under distribution shift and establish attention temperature as a powerful, practical lever for building more adaptive and generalizable foundation models.

## References

- [1] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. In *Advances in Neural Information Processing Systems*, 2023.
- [2] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *International Conference on Learning Representations*, 2023.
- [3] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. In *Advances in Neural Information Processing Systems*, 2023.
- [4] 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*, 2020.
- [5] Xiangyu Chen, Qinghao Hu, Kaidong Li, Cuncong Zhong, and Guanghui Wang. Accumulated trivial attention matters in vision transformers on small datasets. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023.
- [6] Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin Belanger, Lucy J Colwell, and Adrian Weller. Rethinking attention with performers. In *International Conference on Learning Representations*, 2021.
- [7] Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn to achieve second-order convergence rates for in-context linear regression. In *Advances in Neural Information Processing Systems*, 2024.
- [8] Hongfu Gao, Feipeng Zhang, Wenyu Jiang, Jun Shu, Feng Zheng, and Hongxin Wei. On the noise robustness of in-context learning for text generation. In *Advances in Neural Information Processing Systems*, 2024.
- [9] Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 2022.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [11] Dongchen Han, Yifan Pu, Zhuofan Xia, Yizeng Han, Xuran Pan, Xiu Li, Jiwen Lu, Shiji Song, and Gao Huang. Bridging the divide: Reconsidering softmax and linear attention. In *Advances in Neural Information Processing Systems*, 2024.
- [12] L. Isserlis. On a formula for the product-moment coefficient of any order of a normal frequency distribution in any number of variables. *Biometrika*, 12(1/2):134–139, 1918.

- [13] Seung Hoon Lee, Seunghyun Lee, and Byung Cheol Song. Vision transformer for small-size datasets. *arXiv preprint arXiv:2112.13492*, 2021.
- [14] Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, 2023.
- [15] Yingcong Li, Ankit Singh Rawat, and Samet Oymak. Fine-grained analysis of in-context linear estimation: Data, architecture, and beyond. In *Advances in Neural Information Processing Systems*, 2024.
- [16] Junyang Lin, Xu Sun, Xuancheng Ren, Muyu Li, and Qi Su. Learning when to concentrate or divert attention: Self-adaptive attention temperature for neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2985–2990, 2018.
- [17] Jiachang Liu, Dinghan Shen, Yizhe Zhang, William B Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for gpt-3? In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114, 2022.
- [18] Arvind V. Mahankali, Tatsunori Hashimoto, and Tengyu Ma. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. In *International Conference on Learning Representations*, 2024.
- [19] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- [20] Core Francisco Park, Ekdeep Singh Lubana, Itamar Pres, and Hidenori Tanaka. Competition dynamics shape algorithmic phases of in-context learning. *arXiv preprint arXiv:2412.01003*, 2024.
- [21] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. YaRN: Efficient context window extension of large language models. In *International Conference on Learning Representations*, 2024.
- [22] Zhen Qin, Weixuan Sun, Hui Deng, Dongxu Li, Yunshen Wei, Baohong Lv, Junjie Yan, Lingpeng Kong, and Yiran Zhong. cosformer: Rethinking softmax in attention. In *International Conference on Learning Representations*, 2022.
- [23] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI*, 2019.
- [24] Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. *Advances in Neural Information Processing Systems*, 2023.

- [25] Matthew Renze and Erhan Guven. The effect of sampling temperature on problem solving in large language models. In *Findings of the association for computational linguistics: EMNLP*, 2024.
- [26] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language models a mirage? In *Advances in Neural Information Processing Systems*, 2023.
- [27] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- [29] Petar Veličković, Christos Perivolaropoulos, Federico Barbero, and Razvan Pascanu. Softmax is not enough (for sharp size generalisation). *International Conference on Machine Learning*, 2025.
- [30] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, Joao Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, 2023.
- [31] 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.
- [32] Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, 2017.
- [33] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, 2020.
- [34] 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 *International Conference on Learning Representations*, 2024.
- [35] Zhenyu Wu, YaoXiang Wang, Jiacheng Ye, Jiangtao Feng, Jingjing Xu, Yu Qiao, and Zhiyong Wu. Openicl: An open-source framework for in-context learning. *arXiv preprint arXiv:2303.02913*, 2023.
- [36] Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024.

- [37] Shengqiang Zhang, Xingxing Zhang, Hangbo Bao, and Furu Wei. Attention temperature matters in abstractive summarization distillation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 127–141, 2022.
- [38] Yixiong Zou, Ran Ma, Yuhua Li, and Ruixuan Li. Attention temperature matters in vit-based cross-domain few-shot learning. In *Advances in Neural Information Processing Systems*, 2024.

## Appendix A. Related work

**Theory of in-context learning** — Simplified Transformer variants—particularly those using linear attention—have proven useful for gaining analytical insights about ICL [9, 24, 36]. Notably, Zhang et al. [36] showed that linear Transformers approximate Bayes-optimal inference in linear regression tasks, even under distribution shift. We build on this line of research but focus explicitly on the role of the *attention temperature*. In contrast to Zhang et al. [36], we (i) employ *approximate softmax attention* to isolate the effect of temperature, (ii) study how temperature adjustments can mitigate the impact of distribution shifts, and (iii) derive and empirically evaluate the *optimal temperature* for improving ICL performance. These advances extend prior analyses and yield a deeper theoretical and empirical understanding of how principled temperature selection enhances the robustness of Transformers under distributional shift.

**Linear vs. softmax attention** — Although linear attention has gained traction for its computational efficiency, it typically lags behind softmax-based counterparts in predictive performance, spurring efforts to narrow this gap [6, 22]. A pivotal advance in this direction is due to Han et al. [11], who showed that an *approximate variant of softmax attention* can closely approximate the performance of standard softmax attention. Building on this insight, we adopt the *approximate softmax* formulation, which preserves the essential temperature-dependent behavior of standard attention while enabling tractable theoretical analysis. This choice provides a principled framework for investigating how attention-temperature selection shapes ICL performance in pretrained Transformers.

**Attention temperature** — Research on attention temperature remains limited. Veličković et al. [29] recently proposed an adaptive temperature scheme to sharpen softmax outputs, and several empirical studies in natural language processing and computer vision [5, 13, 16, 21, 37, 38] suggest that adjusting the attention temperature can enhance Transformer performance. However, these works do not examine ICL under distributional shift. To our knowledge, no prior study has systematically analyzed how attention temperature influences ICL in such settings—a gap our work directly addresses.

**ICL by Transformers** — The ICL capability of Transformers was first brought to prominence by [4], leading to a surge of empirical and theoretical investigations. Several works have demonstrated that ICL performance improves with model scale [19, 26, 31], underscoring its importance in modern AI systems. To better understand this phenomenon, synthetic tasks such as linear regression have served as controlled testbeds for analyzing ICL in Transformers [9, 24, 36]. A prevailing hypothesis in recent theoretical work is that Transformers implicitly learn algorithms during pre-training, which they subsequently execute during inference [1–3, 7, 14, 15, 18, 20, 30, 36]. There remains ongoing debate over the precise nature of these learned procedures. However, our work focuses on a fundamentally different question, which is how attention temperature affects the ICL performance of pretrained Transformers under distribution shifts.

## Appendix B. Derivation of Bayes-optimal ridge estimator for $w$

We derive the Bayes-optimal ridge estimator for  $w$  given a set of context samples. We place a Gaussian prior on  $w$ , assumed to be a random vector  $w \sim \mathcal{N}(\mu_0, \Sigma_0)$  with prior mean  $\mu_0$  and covariance  $\Sigma_0$ . Let the observed (centered) inputs and labels be

$$\bar{X} = [\bar{x}_1, \dots, \bar{x}_{l-1}]^T, \quad \bar{y} = [\bar{y}_1, \dots, \bar{y}_{l-1}]^T,$$

and assume i.i.d. Gaussian noise  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ . The likelihood of  $\bar{y}$  given  $w$  is

$$p(\bar{y} | \bar{X}, w) = \prod_{i=1}^{l-1} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(\bar{y}_i - w^T \bar{x}_i)^2}{2\sigma^2}\right] \quad (10)$$

$$\propto \exp\left[-\frac{1}{2\sigma^2}(\bar{y} - \bar{X}w)^T(\bar{y} - \bar{X}w)\right], \quad (11)$$

where  $\propto$  denotes proportionality.

By Bayes' rule, the posterior of  $w$  is proportional to the product of likelihood and prior:

$$p(w | \bar{y}, \bar{X}) \propto p(\bar{y} | \bar{X}, w) p(w). \quad (12)$$

Substituting the Gaussian prior yields

$$p(w | \bar{y}, \bar{X}) \propto \exp\left[-\frac{1}{2\sigma^2}(\bar{y} - \bar{X}w)^T(\bar{y} - \bar{X}w)\right] \exp\left[-\frac{1}{2}(w - \mu_0)^T \Sigma_0^{-1}(w - \mu_0)\right]. \quad (13)$$

To determine the form of the posterior distribution, we complete the square in the exponent by collecting all terms involving  $w$ . Expanding the exponent in the joint expression from above, we obtain:

$$-\frac{1}{2\sigma^2}(\bar{y}^T \bar{y} - 2\bar{y}^T \bar{X}w + w^T \bar{X}^T \bar{X}w) - \frac{1}{2}(w^T \Sigma_0^{-1}w - 2\mu_0^T \Sigma_0^{-1}w + \mu_0^T \Sigma_0^{-1}\mu_0). \quad (14)$$

Grouping the quadratic and linear terms in  $w$ , we arrive at:

$$-\frac{1}{2}w^T \left(\frac{\bar{X}^T \bar{X}}{\sigma^2} + \Sigma_0^{-1}\right) w + w^T \left(\frac{\bar{X}^T \bar{y}}{\sigma^2} + \Sigma_0^{-1}\mu_0\right) + \text{terms independent of } w. \quad (15)$$

Defining the posterior precision and linear coefficient terms as  $\Sigma_l^{-1} = \frac{\bar{X}^T \bar{X}}{\sigma^2} + \Sigma_0^{-1}$  and  $b_l = \frac{\bar{X}^T \bar{y}}{\sigma^2} + \Sigma_0^{-1}\mu_0$ , the exponent can be rewritten as

$$-\frac{1}{2}w^T \Sigma_l^{-1}w + w^T b_l = -\frac{1}{2}(w - \mu_l)^T \Sigma_l^{-1}(w - \mu_l) + \text{const}, \quad (16)$$

where  $\mu_l = \Sigma_l b_l$  denotes the posterior mean. Expanding this expression gives:

$$\mu_l = \left(\frac{\bar{X}^T \bar{X}}{\sigma^2} + \Sigma_0^{-1}\right)^{-1} \left(\frac{\bar{X}^T \bar{y}}{\sigma^2} + \Sigma_0^{-1}\mu_0\right). \quad (17)$$

Hence, the posterior distribution of  $w$  given the observed data is Gaussian:

$$w | \bar{y}, \bar{X} \sim \mathcal{N}(\mu_l, \Sigma_l), \quad (18)$$

where  $\boldsymbol{\mu}_l$  is the posterior mean and  $\boldsymbol{\Sigma}_l$  is the posterior covariance matrix.

Under squared-error loss, the Bayes-optimal estimator coincides with the posterior mean, yielding the Bayes-optimal ridge estimator:

$$\hat{\boldsymbol{w}}_{\text{Ridge}} = \mathbb{E}[\boldsymbol{w} \mid \bar{\boldsymbol{y}}, \bar{\boldsymbol{X}}] = \boldsymbol{\mu}_l = \left( \frac{\bar{\boldsymbol{X}}^T \bar{\boldsymbol{X}}}{\sigma^2} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1} \left( \frac{\bar{\boldsymbol{X}}^T \bar{\boldsymbol{y}}}{\sigma^2} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right). \quad (19)$$

This expression provides the Bayes-optimal ridge estimate of  $\boldsymbol{w}$  under a Gaussian prior and additive Gaussian noise—minimizing expected squared error with respect to the posterior.

### Appendix C. Derivation of approximate softmax

The function softmax :  $\mathbb{R}^l \rightarrow \mathbb{R}^l$  is defined component-wise as

$$\text{softmax}(\boldsymbol{z})_i := \frac{e^{z_i}}{\sum_{j=1}^l e^{z_j}} \quad \forall i \in \{1, \dots, l\}. \quad (20)$$

To obtain an approximation, we expand around the origin  $\boldsymbol{z} = \mathbf{0}$  using a first-order Taylor series:

$$\text{softmax}(\boldsymbol{z}) \approx \text{softmax}(\mathbf{0}) + J_{\text{softmax}}(\mathbf{0})\boldsymbol{z}, \quad (21)$$

where  $J_{\text{softmax}}(\mathbf{0})$  is the Jacobian matrix of the softmax function evaluated at  $\boldsymbol{z} = \mathbf{0}$ .

We first compute the zeroth-order term:

$$\text{softmax}(\mathbf{0}) = \frac{e^0}{\sum_{j=1}^l e^0} \mathbf{1} = \frac{1}{l} \mathbf{1}. \quad (22)$$

Next, we evaluate the Jacobian entries at  $\boldsymbol{z} = \mathbf{0}$ :

$$J_{\text{softmax}}(\mathbf{0})_{ii} = \text{softmax}(\mathbf{0})_i (1 - \text{softmax}(\mathbf{0})_i) = \frac{l-1}{l^2}, \quad \forall i, \quad (23)$$

$$J_{\text{softmax}}(\mathbf{0})_{ij} = -\text{softmax}(\mathbf{0})_i \cdot \text{softmax}(\mathbf{0})_j = -\frac{1}{l^2}, \quad \forall i \neq j. \quad (24)$$

This yields the compact matrix form:

$$J_{\text{softmax}}(\mathbf{0}) = \frac{1}{l} \mathbf{I} - \frac{1}{l^2} \mathbf{1}\mathbf{1}^T. \quad (25)$$

Substituting back, we obtain the approximate softmax:

$$\text{softmax}(\boldsymbol{z}) \approx \frac{1}{l} \mathbf{1} + \left( \frac{1}{l} \mathbf{I} - \frac{1}{l^2} \mathbf{1}\mathbf{1}^T \right) \boldsymbol{z}, \quad (26)$$

$$= \left( \frac{1}{l} - \frac{1}{l^2} \sum_{j=1}^l z_j \right) \mathbf{1} + \frac{1}{l} \boldsymbol{z}, \quad (27)$$

$$=: \widehat{\text{softmax}}(\boldsymbol{z}). \quad (28)$$

This derivation yields the approximate softmax attention formulation in (3). From a practical standpoint, approximate softmax attention mechanisms have been empirically evaluated and shown to achieve performance comparable to standard softmax attention [11].

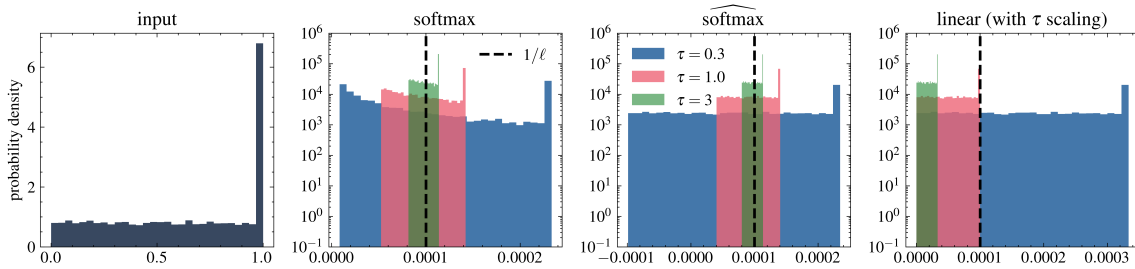


Figure 3: Comparison of temperature effects in softmax, approximate softmax, and linear (with temperature scaling) cases. We consider an input vector  $\mathbf{x} \in \mathbb{R}^l$  whose histogram is illustrated on the left-most plot. Rest of the plots illustrates histograms of the elements of  $\text{softmax}(\mathbf{x}/\tau)$ , approximation  $\widehat{\text{softmax}}(\mathbf{x}/\tau)$  derived in Appendix C and  $\mathbf{x}/(l\tau)$  from left to right, respectively.

#### Appendix D. Temperature effects for softmax and approximate softmax

In contrast to linear attention [36],  $\mathbf{Z} + \frac{1}{l}\mathbf{V}\mathbf{Z}(\mathbf{K}\mathbf{Z})^T(\mathbf{Q}\mathbf{Z})$ , the approximate softmax formulation in (3) explicitly preserves normalization, which is essential for both interpretability and robustness. This difference is described in the following remark.

**Remark 4 (Linear vs. approximate softmax attention)** *Approximate softmax attention maintains row-wise normalization, making it inherently more robust to shifts in input means — a critical failure mode of linear attention in ICL.*

The temperature parameter in softmax directly controls the variance of the output distribution. At higher temperatures, the variance across components decreases, and in the limit  $\tau \rightarrow \infty$ , all elements converge to  $1/l$  with zero variance. Conversely, lower temperatures increase variance, and as  $\tau \rightarrow 0^+$ , the output approaches a one-hot vector, achieving maximal variance.

In the approximate case, temperature similarly acts as an inverse scaling of the variance of the output components, capturing the limit  $\tau \rightarrow \infty$  (all elements equal to  $1/l$ ). For  $\tau \rightarrow 0^+$ , approximate softmax also reflects the maximal variance, but it does not produce a true one-hot distribution. Thus, approximate softmax closely mirrors the temperature behavior of softmax, except in the degenerate limit  $\tau \rightarrow 0^+$ , which is not of practical relevance in this work.

To further illustrate these effects, Figure 3 compares softmax and approximate softmax across different temperatures. The figure demonstrates that approximate softmax faithfully captures the variance effect of temperature: the variance of the output components is inversely proportional to  $\tau$ . Moreover, as  $\tau \rightarrow \infty$ , both softmax and approximate softmax concentrate around  $1/l$ , whereas linear attention with temperature scaling does not. Overall, the output distributions of softmax and approximate softmax are highly similar, except at very small values of  $\tau$ , where approximate softmax may yield negative components while softmax tends toward sparsity with many zeros. By contrast, linear attention with temperature scaling produces qualitatively different distributions. This comparison highlights the advantage of approximate softmax as a faithful surrogate for analyzing temperature effects relevant to softmax.

### Appendix E. Expanded form of approximate softmax attention

Using block matrix notation, the prediction from the approximate softmax attention model can be expanded as:

$$\hat{y}(\mathbf{Z}; \mathbf{V}, \mathbf{M}) = A_{d+1,l}, \quad (29)$$

$$= \frac{1}{l} \mathbf{V}_{d+1,:} \mathbf{Z} \left( \frac{(\mathbf{KZ})^T (\mathbf{QZ}_{:,l})}{\tau} - \frac{1}{l} \sum_{j=1}^l \frac{(\mathbf{KZ}_{:,j})^T (\mathbf{QZ}_{:,l})}{\tau} + \mathbf{1} \right), \quad (30)$$

$$= \frac{1}{l} [\mathbf{v}_{21}^T \ v_{22}] \mathbf{Z} \left( \frac{\mathbf{Z}^T \mathbf{M} \mathbf{Z}_{:,l}}{\tau} - \frac{1}{l} \sum_{j=1}^l \frac{(\mathbf{Z}_{:,j})^T \mathbf{M} \mathbf{Z}_{:,l}}{\tau} + \mathbf{1} \right), \quad (31)$$

$$= \frac{1}{l} [\mathbf{v}_{21}^T \ v_{22}] [\mathbf{X} \ \mathbf{y}]^T \left( \frac{[\mathbf{X} \ \mathbf{y}] \mathbf{M} [\mathbf{x}_l^T \ 0]^T}{\tau} - \frac{1}{l} \sum_{j=1}^l \frac{[\mathbf{x}_j^T \ y_j] \mathbf{M} [\mathbf{x}_l^T \ 0]^T}{\tau} + \mathbf{1} \right), \quad (32)$$

$$= \frac{1}{l} [\mathbf{v}_{21}^T \ v_{22}] [\mathbf{X} \ \mathbf{y}]^T \left( \frac{1}{\tau} [\mathbf{X} - \mathbf{1} \mathbf{s}_x^T \ \mathbf{y} - s_y \mathbf{1}] \begin{bmatrix} \mathbf{M}_{11} & * \\ \mathbf{m}_{21}^T & * \end{bmatrix} [\mathbf{x}_l^T \ 0]^T + \mathbf{1} \right), \quad (33)$$

$$= \frac{1}{l} [\mathbf{v}_{21}^T \ v_{22}] [\mathbf{X} \ \mathbf{y}]^T \left( \frac{1}{\tau} (\mathbf{X} - \mathbf{1} \mathbf{s}_x^T) \mathbf{M}_{11} \mathbf{x}_l + \frac{1}{\tau} (\mathbf{y} - s_y \mathbf{1}) \mathbf{m}_{21}^T \mathbf{x}_l + \mathbf{1} \right), \quad (34)$$

$$= \frac{1}{l} (\mathbf{v}_{21}^T \mathbf{X}^T + v_{22} \mathbf{y}^T) \left( \frac{1}{\tau} (\mathbf{X} - \mathbf{1} \mathbf{s}_x^T) \mathbf{M}_{11} \mathbf{x}_l + \frac{1}{\tau} (\mathbf{y} - s_y \mathbf{1}) \mathbf{m}_{21}^T \mathbf{x}_l + \mathbf{1} \right), \quad (35)$$

$$= \frac{1}{\tau} \left( \mathbf{v}_{21}^T \left( \frac{\mathbf{X}^T \mathbf{X}}{l} - \mathbf{s}_x \mathbf{s}_x^T \right) + v_{22} \left( \frac{\mathbf{y}^T \mathbf{X}}{l} - s_y \mathbf{s}_x^T \right) \right) \mathbf{M}_{11} \mathbf{x}_l, \\ + \frac{1}{\tau} \left( \mathbf{v}_{21}^T \left( \frac{\mathbf{X}^T \mathbf{y}}{l} - s_y \mathbf{s}_x \right) + v_{22} \left( \frac{\mathbf{y}^T \mathbf{y}}{l} - s_y^2 \right) \right) \mathbf{m}_{21}^T \mathbf{x}_l + \mathbf{v}_{21}^T \mathbf{s}_x + v_{22} s_y, \quad (36)$$

$$= \frac{1}{\tau} (\mathbf{v}_{21}^T \mathbf{C}_{xx} + v_{22} \mathbf{C}_{xy}^T) \mathbf{M}_{11} \mathbf{x}_l + \frac{1}{\tau} (\mathbf{v}_{21}^T \mathbf{C}_{xy} + v_{22} \mathbf{C}_{yy}) \mathbf{m}_{21}^T \mathbf{x}_l + \mathbf{v}_{21}^T \mathbf{s}_x + v_{22} s_y, \quad (37)$$

$$= \frac{1}{\tau} ((\mathbf{v}_{21}^T \mathbf{C}_{xx} + v_{22} \mathbf{C}_{xy}^T) \mathbf{M}_{11} + (\mathbf{v}_{21}^T \mathbf{C}_{xy} + v_{22} \mathbf{C}_{yy}) \mathbf{m}_{21}^T) \mathbf{x}_l + \mathbf{v}_{21}^T \mathbf{s}_x + v_{22} s_y, \quad (38)$$

where the summary statistics are defined as:

$$\mathbf{s}_x := \frac{1}{l} \sum_{i=1}^l \mathbf{x}_i, \quad s_y := \frac{1}{l} \sum_{i=1}^{l-1} y_i, \\ \mathbf{C}_{xx} := \frac{1}{l} \sum_{i=1}^l \mathbf{x}_i \mathbf{x}_i^T - \mathbf{s}_x \mathbf{s}_x^T, \quad \mathbf{C}_{xy} := \frac{1}{l} \sum_{i=1}^{l-1} y_i \mathbf{x}_i - s_y \mathbf{s}_x, \quad \mathbf{C}_{yy} := \frac{1}{l} \sum_{i=1}^{l-1} y_i^2 - s_y^2.$$

Then, we define

$$\hat{\mathbf{w}}_{Att}(\mathbf{C}_{xx}, \mathbf{C}_{xy}, \mathbf{C}_{yy}; \mathbf{M}, \mathbf{V}) = \mathbf{M}_{11}^T (\mathbf{C}_{xx} \mathbf{v}_{21} + v_{22} \mathbf{C}_{xy}) + (\mathbf{v}_{21}^T \mathbf{C}_{xy} + v_{22} \mathbf{C}_{yy}) \mathbf{m}_{21}, \quad (39)$$

$$b_{Att}(\mathbf{s}_x, s_y; \mathbf{V}) = \mathbf{v}_{21}^T \mathbf{s}_x + v_{22} s_y, \quad (40)$$

which allows us to write

$$\hat{y}(\mathbf{Z}; \mathbf{V}, \mathbf{M}) = \frac{1}{\tau} \hat{\mathbf{w}}_{Att}(\mathbf{C}_{xx}, \mathbf{C}_{xy}, \mathbf{C}_{yy}; \mathbf{M}, \mathbf{V})^T \mathbf{x}_l + b_{Att}(\mathbf{s}_x, s_y; \mathbf{V}). \quad (41)$$

## Appendix F. Derivation of the pretraining for ICL by mimicking the Bayes-optimal estimator

Here, we derive the pretraining of the approximate softmax attention model by mimicking the Bayes-optimal ridge estimator (Appendix B). The prediction of the approximate softmax attention model can be written as

$$\hat{y}(\mathbf{Z}; \mathbf{V}, \mathbf{M}) = \frac{1}{\tau} \hat{\mathbf{w}}_{Att}(\mathbf{C}_{xx}, \mathbf{C}_{xy}, \mathbf{C}_{yy}; \mathbf{M}, \mathbf{V})^T \mathbf{x}_l + b_{Att}(\mathbf{s}_x, \mathbf{s}_y; \mathbf{V}), \quad (42)$$

which is derived in Appendix E. Furthermore, the Bayes-optimal ridge regression model's prediction is

$$\hat{y}_{Bayes} = \hat{\mathbf{w}}_{Bayes}^T \mathbf{x}_l. \quad (43)$$

Therefore, we select the parameters  $\mathbf{M}$  and  $\mathbf{V}$  such that

$$\hat{\mathbf{w}}_{Att}(\mathbf{C}_{xx}, \mathbf{C}_{xy}, \mathbf{C}_{yy}; \mathbf{M}, \mathbf{V}) \approx \hat{\mathbf{w}}_{Bayes}, \quad b_{Att}(\mathbf{s}_x, \mathbf{s}_y; \mathbf{V}) \approx 0, \quad (44)$$

which makes the prediction of the approximate softmax attention model approximately equal to that of the Bayes-optimal regression. Furthermore, we consider  $\tau = 1$  for the pretraining. Let's first focus on  $\hat{\mathbf{w}}_{Att}(\mathbf{C}_{xx}, \mathbf{C}_{xy}, \mathbf{C}_{yy}; \mathbf{M}, \mathbf{V})$  as follows

$$\begin{aligned} \hat{\mathbf{w}}_{Att}(\mathbf{C}_{xx}, \mathbf{C}_{xy}, \mathbf{C}_{yy}; \mathbf{M}, \mathbf{V}) \\ = (\mathbf{M}_{11}^T (\mathbf{C}_{xx} \mathbf{v}_{21} + v_{22} \mathbf{C}_{xy}) + (\mathbf{v}_{21}^T \mathbf{C}_{xy} + v_{22} \mathbf{C}_{yy}) \mathbf{m}_{21}), \end{aligned} \quad (45)$$

$$= \left( \mathbf{M}_{11}^T \left( \frac{\bar{\mathbf{X}}^T \bar{\mathbf{X}}}{l} \mathbf{v}_{21} + v_{22} \frac{\bar{\mathbf{X}}^T \bar{\mathbf{y}}}{l} \right) + \left( \mathbf{v}_{21}^T \frac{\bar{\mathbf{X}}^T \bar{\mathbf{y}}}{l} + v_{22} \frac{\bar{\mathbf{y}}^T \bar{\mathbf{y}}}{l} \right) \mathbf{m}_{21} \right). \quad (46)$$

To reach the last line, we use the fact that  $\mathbf{C}_{xx} := \mathbf{X}^T \mathbf{X} / l - \mathbf{s}_x \mathbf{s}_x^T = \bar{\mathbf{X}}^T \bar{\mathbf{X}} / l$ ,  $\mathbf{C}_{xy} := \mathbf{X}^T \mathbf{y} / l - \mathbf{s}_y \mathbf{s}_x^T = \bar{\mathbf{X}}^T \bar{\mathbf{y}} / l$  and  $\mathbf{C}_{yy} := \mathbf{y}^T \mathbf{y} / l$ , where  $\bar{\mathbf{X}} := \mathbf{X} - \mathbf{s}_x^T$  and  $\bar{\mathbf{y}} := \mathbf{y} - \mathbf{s}_y$  denote centered input matrix and centered label vector. Now, recall that the Bayes-optimal ridge estimator is

$$\hat{\mathbf{w}}_{Bayes} = \left( \frac{\bar{\mathbf{X}}^T \bar{\mathbf{X}}}{\sigma^2} + \Sigma_w^{-1} \right)^{-1} \left( \frac{\bar{\mathbf{X}}^T \bar{\mathbf{y}}}{\sigma^2} + \Sigma_w^{-1} \boldsymbol{\mu}_w \right), \quad (47)$$

as derived in Appendix B. Looking at equations (47) and (46) together, we can see that setting the parameters as follows would make  $\hat{\mathbf{w}}_{Att} = \hat{\mathbf{w}}_{Bayes}$  hold

$$\mathbf{M}_{11} = \frac{l}{\sigma^2} \left( \frac{\bar{\mathbf{X}}^T \bar{\mathbf{X}}}{\sigma^2} + \Sigma_w^{-1} \right)^{-1}, \quad \mathbf{v}_{21} = \frac{\sigma^2}{l} \left( \frac{\bar{\mathbf{X}}^T \bar{\mathbf{X}}}{l} \right)^{-1} \Sigma_w^{-1} \boldsymbol{\mu}_w, \quad \mathbf{m}_{21} = \mathbf{0}, \quad v_{22} = 1. \quad (48)$$

However, while Bayes-optimal estimator  $\hat{\mathbf{w}}_{Bayes}$  is different for each sample, the attention model should be pretrained and fixed. Thus, we replace  $\bar{\mathbf{X}}^T \bar{\mathbf{X}}$  in (48) with  $\hat{\mathbf{X}}^T \hat{\mathbf{X}} / m$  as follows, where  $\hat{\mathbf{X}} \in \mathbb{R}^{ml \times d}$  is the centred input matrix including all the (pre)training data consisting of  $ml$  samples.

$$\mathbf{M}_{11} = \frac{l}{\sigma^2} \left( \frac{\hat{\mathbf{X}}^T \hat{\mathbf{X}}}{m\sigma^2} + \Sigma_w^{-1} \right)^{-1}, \quad \mathbf{v}_{21} = \frac{\sigma^2}{l} \left( \frac{\hat{\mathbf{X}}^T \hat{\mathbf{X}}}{ml} \right)^{-1} \Sigma_w^{-1} \boldsymbol{\mu}_w, \quad \mathbf{m}_{21} = \mathbf{0}, \quad v_{22} = 1. \quad (49)$$

In practice, the variance of noise  $\sigma^2$ , the mean  $\boldsymbol{\mu}_w$ , and covariance  $\boldsymbol{\Sigma}_w$  of the task vectors are unknown. Yet, we can use their estimates based on the (pre)training data.

Now, we can focus on making  $b_{Att}(\mathbf{s}_x, \mathbf{s}_y; \mathbf{V}) \approx 0$  hold as follows

$$b_{Att}(\mathbf{s}_x, \mathbf{s}_y; \mathbf{V}) = \mathbf{v}_{21}^T \mathbf{s}_x + v_{22} s_y, \quad (50)$$

where  $\mathbf{s}_x$  and  $s_y$  are based on data so we have no control over them. Instead, by using Assumptions 1 and 2, we can choose  $\mathbf{v}_{21}$  and  $v_{22}$  such that  $b_{Att} \rightarrow 0$  as  $l, d \rightarrow \infty$ . Note that Assumption 1 makes  $\mathbf{v}_{21}^T \mathbf{s}_x + v_{22} s_y$  bounded with high probability for  $\mathbf{v}_{21}$  and  $v_{22}$  given in (49). Therefore, multiplying  $\mathbf{v}_{21}, v_{22}$  given in (49) with  $1/d$  would make  $b_{Att} \rightarrow 0$  as  $d \rightarrow \infty$ . To fix the impact of the multiplication for  $\hat{\mathbf{w}}_{Att}$ , we can multiply  $\mathbf{M}_{11}$  with  $d$  as well. So, by applying the mentioned multiplications, we reach the following pretrained parameters mimicking the Bayes-optimal regression model

$$\mathbf{M}_{11} = \frac{dl}{\sigma^2} \left( \frac{\hat{\mathbf{X}}^T \hat{\mathbf{X}}}{m\sigma^2} + \boldsymbol{\Sigma}_w^{-1} \right)^{-1}, \quad \mathbf{v}_{21} = \frac{\sigma^2}{dl} \left( \frac{\hat{\mathbf{X}}^T \hat{\mathbf{X}}}{ml} \right)^{-1} \boldsymbol{\Sigma}_w^{-1} \boldsymbol{\mu}_w, \quad \mathbf{m}_{21} = \mathbf{0}, \quad v_{22} = \frac{1}{d}. \quad (51)$$

We summarize the resulting pretrained-parameters configuration as the following proposition.

**Proposition 5 (Pretrained parameters)** *When the temperature parameter is set to  $\tau = 1$  during pretraining, the following parameter configuration approximates the Bayes-optimal estimator (Appendix B):*

$$\begin{aligned} \mathbf{M}_{11} &= d \left( \frac{\hat{\mathbf{X}}^T \hat{\mathbf{X}}}{ml} + \frac{\sigma^2}{l} \boldsymbol{\Sigma}_w^{-1} \right)^{-1}, & \mathbf{m}_{21} &= \mathbf{0}, \\ \mathbf{v}_{21} &= \frac{\sigma^2}{dl} \left( \frac{\hat{\mathbf{X}}^T \hat{\mathbf{X}}}{ml} \right)^{-1} \boldsymbol{\Sigma}_w^{-1} \boldsymbol{\mu}_w, & v_{22} &= \frac{1}{d}, \end{aligned} \quad (52)$$

where  $\hat{\mathbf{X}} \in \mathbb{R}^{ml \times d}$  is the centered input matrix formed from  $ml$  samples of  $\mathbf{x}$ . This configuration aligns our model with Bayes-optimal ridge regression. The quantities  $\boldsymbol{\mu}_w$  and  $\boldsymbol{\Sigma}_w$  can be estimated from the pretraining data.

**Remark 6** *While the pretrained parameters specified in Proposition 5 are not guaranteed to be optimal in all settings, they are analytically useful for examining the effects of distribution shifts.*

Based on Proposition 5, we arrive at the following:

**Corollary 7** *Suppose there is no distribution shift between training and inference. Then, under the parameter configuration of Proposition 5, the Transformer model (3) emulates the Bayes-optimal linear model, implying that it is capable of in-context learning according to Definition 1.*

## Appendix G. Characterization of generalization error for ICL under distribution shift

We first state the boundedness assumption on the pretrained parameters used throughout Theorems 2 and 3, and give the explicit forms of the matrices appearing in those statements.

**Assumption 3** *There exists a constant  $c > 0$  such that*

$$\|\mathbf{M}_{11}\| \leq cd, \quad \|\mathbf{m}_{21}\| = 0, \quad \|\mathbf{v}_{21}\| \leq \frac{c}{dl}, \quad |v_{22}| \leq \frac{c}{d}.$$

The matrices  $\mathbf{A}, \mathbf{B}, \mathbf{F}_1, \mathbf{F}_2, \hat{\mathbf{B}}$  appearing in Theorem 2 and Theorem 3 are defined as

$$\mathbf{A} := \boldsymbol{\Sigma}_x + \boldsymbol{\mu}_x \boldsymbol{\mu}_x^T, \quad \mathbf{B} := \boldsymbol{\Sigma}_w + \boldsymbol{\mu}_w \boldsymbol{\mu}_w^T, \quad (53)$$

$$\mathbf{F}_1 := \left( \boldsymbol{\Sigma}_x \hat{\mathbf{B}} + \frac{1}{l} \left( v_{22}^2 \sigma^2 + \text{Tr}(\hat{\mathbf{B}} \boldsymbol{\Sigma}_x) \right) \mathbf{I} \right) \boldsymbol{\Sigma}_x, \quad (54)$$

$$\mathbf{F}_2 := (\boldsymbol{\mu}_w \mathbf{v}_{21}^T + v_{22} \mathbf{B}) \boldsymbol{\Sigma}_x, \quad (55)$$

$$\hat{\mathbf{B}} := v_{22} \boldsymbol{\mu}_w \mathbf{v}_{21}^T + v_{22} \mathbf{v}_{21} \boldsymbol{\mu}_w^T + v_{22}^2 \mathbf{B}. \quad (56)$$

Here, we characterize the generalization error for in-context learning under distribution shift, given that  $\mathbf{M}$  and  $\mathbf{V}$  are pretrained and fixed. So, the impact of pretraining distribution  $\mathcal{D}^{train}$  is captured by  $\mathbf{M}$  and  $\mathbf{V}$ . Suppose that  $\mathcal{D}^{test}$  denotes the test distribution. To avoid additional notations, here, we again use  $\boldsymbol{\mu}_x, \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_x, \boldsymbol{\Sigma}_w, \sigma^2$  to denote means and covariances for input and task vectors and noise variance for the inference (test). However, note that these can be different from those used for pretraining. We begin studying the generalization error defined in (5) as follows

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) := \mathbb{E}_{(\mathbf{Z}, y_l) \sim \mathcal{D}^{test}} \left[ (y_l - \hat{y}(\mathbf{Z}; \mathbf{V}, \mathbf{M}))^2 \right], \quad (57)$$

$$= \mathbb{E}_{(\mathbf{Z}, y_l) \sim \mathcal{D}^{test}} \left[ \left( \frac{1}{\tau} \hat{\mathbf{w}}_{Att}(\mathbf{C}_{xx}, \mathbf{C}_{xy}, \mathbf{C}_{yy}; \mathbf{M}, \mathbf{V})^T \mathbf{x}_l + b_{Att}(\mathbf{s}_x, \mathbf{s}_y; \mathbf{V}) - y_l \right)^2 \right], \quad (58)$$

$$= \mathbb{E}_{(\mathbf{Z}, y_l) \sim \mathcal{D}^{test}} \left[ \left( \frac{1}{\tau} (\mathbf{M}_{11}^T (\mathbf{C}_{xx} \mathbf{v}_{21} + v_{22} \mathbf{C}_{xy}))^T \mathbf{x}_l - y_l \right)^2 \right], \quad (59)$$

where we use the parameters from pretraining (51) together with Assumptions 1 and 2 to reach the last line. Then,

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) = \mathbb{E}_{(\mathbf{Z}, y_l) \sim \mathcal{D}^{test}} \left[ \left( \frac{1}{\tau} (\mathbf{M}_{11}^T (\mathbf{C}_{xx} \mathbf{v}_{21} + v_{22} \mathbf{C}_{xy}))^T \mathbf{x}_l - y_l \right)^2 \right], \quad (60)$$

$$= \mathbb{E} \left[ \left( \frac{1}{\tau} \left( \mathbf{M}_{11}^T \left( \frac{1}{l} \sum_{i \leq l} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{v}_{21} + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i (\bar{\mathbf{x}}_i^T \mathbf{w} + \epsilon_i) \right) \right)^T \mathbf{x}_l - \mathbf{w}^T \mathbf{x}_l - \epsilon_l \right)^2 \right], \quad (61)$$

$$= \mathbb{E} \left[ \left( \frac{1}{\tau} \left( \mathbf{M}_{11}^T \left( \frac{1}{l} \sum_{i \leq l} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{v}_{21} + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i (\bar{\mathbf{x}}_i^T \mathbf{w} + \epsilon_i) \right) \right)^T \mathbf{x}_l - \mathbf{w}^T \mathbf{x}_l \right)^2 \right] + \sigma^2 \quad (62)$$

where  $\bar{\mathbf{x}}_i := \mathbf{x}_i - \mathbf{s}_x = \mathbf{x}_i - \frac{1}{l} \sum_{i \leq l} \mathbf{x}_i$  and we use  $\epsilon_l \sim \mathcal{N}(0, \sigma^2)$  to reach the final line. We continue by defining

$$\mathbf{w}_{diff} := \frac{1}{\tau} \mathbf{M}_{11}^T \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{v}_{21} + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i (\bar{\mathbf{x}}_i^T \mathbf{w} + \epsilon_i) \right) - \mathbf{w}, \quad (63)$$

which allows us to write

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) = \mathbb{E} \left[ (\mathbf{w}_{diff}^T \mathbf{x}_l)^2 \right] + \sigma^2, \quad (64)$$

$$= \mathbb{E} \left[ \mathbf{w}_{diff}^T \mathbb{E}_{\mathbf{x}_l} [\mathbf{x}_l \mathbf{x}_l^T] \mathbf{w}_{diff} \right] + \sigma^2, \quad (65)$$

$$= \mathbb{E} \left[ \mathbf{w}_{diff}^T (\boldsymbol{\mu}_x \boldsymbol{\mu}_x^T + \boldsymbol{\Sigma}_x) \mathbf{w}_{diff} \right] + \sigma^2, \quad (66)$$

by the law of total expectation since  $\mathbf{w}_{diff}$  is independent of  $\mathbf{x}_l$ . Note that when writing (64), we safely ignore terms with  $(1/l)\bar{\mathbf{x}}_l \bar{\mathbf{x}}_l^T \mathbf{v}_{21}$  in (62) since they vanish as  $l \rightarrow \infty$  by Assumptions 1-2 and 3. Letting  $\mathbf{A} := \boldsymbol{\mu}_x \boldsymbol{\mu}_x^T + \boldsymbol{\Sigma}_x$ , we write

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) = \mathbb{E} \left[ \mathbf{w}_{diff}^T \mathbf{A} \mathbf{w}_{diff} \right] + \sigma^2, \quad (67)$$

$$= \mathbb{E} \left[ \text{Tr} (\mathbf{w}_{diff}^T \mathbf{A} \mathbf{w}_{diff}) \right] + \sigma^2, \quad (68)$$

$$= \mathbb{E} \left[ \text{Tr} (\mathbf{A} \mathbf{w}_{diff} \mathbf{w}_{diff}^T) \right] + \sigma^2, \quad (69)$$

$$= \text{Tr} (\mathbf{A} \mathbb{E} [\mathbf{w}_{diff} \mathbf{w}_{diff}^T]) + \sigma^2, \quad (70)$$

where we first apply the cyclic property of trace and then use the linearity of expectation and trace to reach the last line. Now, we need to calculate  $\mathbb{E}[\mathbf{w}_{diff} \mathbf{w}_{diff}^T]$ , for which we first take the expectation over  $\mathbf{w}$ . To do so, we rewrite  $\mathbf{w}_{diff}$  as

$$\mathbf{w}_{diff} = \underbrace{\frac{1}{\tau} \mathbf{M}_{11}^T \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{v}_{21} + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \epsilon_i \right)}_{\mathbf{e}} + \underbrace{\left( \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T - \mathbf{I} \right)}_{\mathbf{D}} \mathbf{w}, \quad (71)$$

$$= \mathbf{e} + \mathbf{D} \mathbf{w}, \quad (72)$$

where we define

$$\mathbf{e} := \frac{1}{\tau} \mathbf{M}_{11}^T \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{v}_{21} + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \epsilon_i \right), \quad (73)$$

$$\mathbf{D} := \left( \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T - \mathbf{I} \right). \quad (74)$$

Since  $\mathbf{e}$  and  $\mathbf{D}$  are independent of  $\mathbf{w}$ , we can easily calculate  $\mathbb{E}_{\mathbf{w}}[\mathbf{w}_{diff} \mathbf{w}_{diff}^T]$  as follows

$$\mathbb{E} \left[ \mathbb{E}_{\mathbf{w}} [\mathbf{w}_{diff} \mathbf{w}_{diff}^T] \right] = \mathbb{E} \left[ \mathbb{E}_{\mathbf{w}} [(\mathbf{e} + \mathbf{D} \mathbf{w})(\mathbf{e} + \mathbf{D} \mathbf{w})^T] \right], \quad (75)$$

$$= \mathbb{E} [\mathbf{e} \mathbf{e}^T] + \mathbb{E} [\mathbf{e} \boldsymbol{\mu}_w^T \mathbf{D}^T] + \mathbb{E} [\mathbf{D} \boldsymbol{\mu}_w \mathbf{e}^T] + \mathbb{E} [\mathbf{D} (\boldsymbol{\mu}_x \boldsymbol{\mu}_x^T + \boldsymbol{\Sigma}_w) \mathbf{D}^T], \quad (76)$$

$$= \mathbb{E} [\mathbf{e} \mathbf{e}^T] + \mathbb{E} [\mathbf{D} \boldsymbol{\mu}_w \mathbf{e}^T]^T + \mathbb{E} [\mathbf{D} \boldsymbol{\mu}_w \mathbf{e}^T] + \mathbb{E} [\mathbf{D} \mathbf{B} \mathbf{D}^T], \quad (77)$$

where we first apply the law of total expectation, then take the expectation over  $\mathbf{w}$  and finally, we define  $\mathbf{B} := \boldsymbol{\mu}_x \boldsymbol{\mu}_x^T + \boldsymbol{\Sigma}_w$  to reach the last line. Note that  $\boldsymbol{\mu}_w$  and  $\mathbf{B}$  are fixed while  $\mathbf{e}$  and  $\mathbf{D}$  are random in the last line. Therefore, we are required to calculate the three expectations that appeared in (77).

Before getting into the calculations of the aforementioned expectations, we provide the following lemma that is useful for the calculation of the expectations.

**Lemma 8** Let  $\bar{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , where  $\bar{\mathbf{x}} \in \mathbb{R}^d$ . Let  $\bar{\mathbf{x}}_i$  be  $l-1$  independent samples of  $\bar{\mathbf{x}}$  for  $i = 1, \dots, l-1$ . Furthermore, let  $\mathbf{A}$  be a fixed  $d \times d$  matrix. Then, the following holds

$$\mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \mathbf{A} \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \right] = \frac{l-1}{l} \Sigma \mathbf{A} \Sigma + \frac{1}{l} \Sigma \mathbf{A}^T \Sigma + \frac{1}{l} \text{Tr}(\mathbf{A} \Sigma) \Sigma. \quad (78)$$

**Proof** This is proven by using Isserlis' theorem [12] in Appendix H. ■

Note that our inputs  $\bar{\mathbf{x}}_i$  are centered, i.e.,  $\bar{\mathbf{x}}_i = \mathbf{x}_i - \frac{1}{l} \sum_{i \leq l} \mathbf{x}_i$ , so their distribution is  $\mathcal{N}(\mathbf{0}, \Sigma_x)$  as  $l \rightarrow \infty$ . Therefore, Lemma 8 is directly applicable in our setting.

Next, we start the calculations of the expectations in (77) with  $\mathbb{E}[\mathbf{e}\mathbf{e}^T]$  as follows

$$\mathbb{E}[\mathbf{e}\mathbf{e}^T] = \frac{1}{\tau^2} \mathbf{M}_{11}^T \mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{v}_{21} + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \epsilon_i \right) \cdot \left( \frac{1}{l} \sum_{i \leq l-1} \mathbf{v}_{21}^T \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i^T \epsilon_i \right) \right] \mathbf{M}_{11}, \quad (79)$$

$$= \frac{1}{\tau^2} \mathbf{M}_{11}^T \left( \mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{v}_{21} \right) \left( \frac{1}{l} \sum_{i \leq l-1} \mathbf{v}_{21}^T \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) + \left( v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \epsilon_i \right) \left( v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i^T \epsilon_i \right) \right] \right) \mathbf{M}_{11}, \quad (80)$$

$$= \frac{1}{\tau^2} \mathbf{M}_{11}^T \left( \mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \mathbf{v}_{21} \mathbf{v}_{21}^T \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) + \left( v_{22}^2 \frac{\sigma^2}{l^2} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \right] \right) \mathbf{M}_{11}, \quad (81)$$

$$= \frac{1}{\tau^2} \mathbf{M}_{11}^T \left( \Sigma_x \mathbf{v}_{21} \mathbf{v}_{21}^T \Sigma_x + \frac{1}{l} \text{Tr}(\mathbf{v}_{21} \mathbf{v}_{21}^T \Sigma_x) \Sigma_x + v_{22}^2 \frac{\sigma^2 (l-1)}{l^2} \Sigma_x \right) \mathbf{M}_{11}, \quad (82)$$

$$= \frac{1}{\tau^2} \mathbf{M}_{11}^T \left( v_{22}^2 \frac{\sigma^2}{l} \Sigma_x \right) \mathbf{M}_{11}, \quad (83)$$

where we first use the independence of the random variables and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  to simplify the equation. Then, we apply Lemma 8 and use the fact that  $\mathbb{E}[\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T] = \Sigma_x$  to get the penultimate line. Finally, we drop the vanishing terms and simplify the result using Assumptions 1-2 and 3 in order to reach the last line.

We continue with the calculation of  $\mathbb{E} [\mathbf{D}\boldsymbol{\mu}_w \mathbf{e}^T]$  as

$$\begin{aligned} & \mathbb{E} [\mathbf{D}\boldsymbol{\mu}_w \mathbf{e}^T] \\ &= \frac{1}{\tau} \mathbb{E} \left[ \left( \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T - \mathbf{I} \right) \boldsymbol{\mu}_w \left( \frac{1}{l} \sum_{i \leq l-1} \mathbf{v}_{21}^T \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T + v_{22} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i^T \epsilon_i \right) \right] \mathbf{M}_{11}, \end{aligned} \quad (84)$$

$$= \frac{1}{\tau} \mathbb{E} \left[ \left( \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T - \mathbf{I} \right) \boldsymbol{\mu}_w \left( \frac{1}{l} \sum_{i \leq l-1} \mathbf{v}_{21}^T \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \right] \mathbf{M}_{11}, \quad (85)$$

$$\begin{aligned} &= \frac{1}{\tau} \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \boldsymbol{\mu}_w \mathbf{v}_{21}^T \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \right] \mathbf{M}_{11} \\ &\quad - \boldsymbol{\mu}_w \mathbf{v}_{21}^T \mathbb{E} \left[ \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right] \mathbf{M}_{11}, \end{aligned} \quad (86)$$

$$\begin{aligned} &= \frac{1}{\tau} \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \left( \boldsymbol{\Sigma}_x \boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x + \frac{1}{l} \boldsymbol{\Sigma}_x \mathbf{v}_{21} \boldsymbol{\mu}_w^T \boldsymbol{\Sigma}_x + \frac{1}{l} \text{Tr}(\boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right) \mathbf{M}_{11} \\ &\quad - \frac{1}{\tau} \frac{l-1}{l} \boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x \mathbf{M}_{11}, \end{aligned} \quad (87)$$

$$= \frac{v_{22}}{\tau^2} \mathbf{M}_{11}^T \left( \boldsymbol{\Sigma}_x \boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x + \frac{1}{l} \text{Tr}(\boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right) \mathbf{M}_{11} - \frac{1}{\tau} \boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x \mathbf{M}_{11}, \quad (88)$$

where we again first use the independence of the random variables and  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ . Then, we apply basic algebraic manipulations. To reach the penultimate line, we utilize Lemma 8 together with the fact that  $\mathbb{E}[\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T] = \boldsymbol{\Sigma}_x$ . Using Assumptions 1-2 and 3, we reach the last line.

Finally, we calculate  $\mathbb{E} [DBD^T]$  as follows

$$\begin{aligned} & \mathbb{E} [DBD^T] \\ &= \mathbb{E} \left[ \left( \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T - \mathbf{I} \right) \mathbf{B} \left( \frac{v_{22}}{\tau} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{M}_{11} - \mathbf{I} \right) \right], \end{aligned} \quad (89)$$

$$\begin{aligned} &= \mathbb{E} \left[ \left( \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \mathbf{B} \left( \frac{v_{22}}{\tau} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{M}_{11} \right) \right] \\ &\quad - \mathbb{E} \left[ \left( \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \mathbf{B} \right] - \mathbb{E} \left[ \mathbf{B} \left( \frac{v_{22}}{\tau} \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{M}_{11} \right) \right] + \mathbf{B}, \end{aligned} \quad (90)$$

$$\begin{aligned} &= \frac{v_{22}^2}{\tau^2} \mathbf{M}_{11}^T \mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \mathbf{B} \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \right] \mathbf{M}_{11} \\ &\quad - \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \right] \mathbf{B} - \frac{v_{22}}{\tau} \mathbf{B} \mathbb{E} \left[ \left( \frac{1}{l} \sum_{i \leq l-1} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right) \right] \mathbf{M}_{11} + \mathbf{B}, \end{aligned} \quad (91)$$

$$\begin{aligned} &= \frac{v_{22}^2}{\tau^2} \mathbf{M}_{11}^T \left( \boldsymbol{\Sigma}_x \mathbf{B} \boldsymbol{\Sigma}_x + \frac{1}{l} \text{Tr}(\mathbf{B} \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right) \mathbf{M}_{11} - \frac{v_{22}}{\tau} \frac{l-1}{l} \mathbf{M}_{11}^T \boldsymbol{\Sigma}_x \mathbf{B} \\ &\quad - \frac{v_{22}}{\tau} \frac{l-1}{l} \mathbf{B} \boldsymbol{\Sigma}_x \mathbf{M}_{11} + \mathbf{B}, \end{aligned} \quad (92)$$

$$= \frac{v_{22}^2}{\tau^2} \mathbf{M}_{11}^T \left( \boldsymbol{\Sigma}_x \mathbf{B} \boldsymbol{\Sigma}_x + \frac{1}{l} \text{Tr}(\mathbf{B} \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right) \mathbf{M}_{11} - \frac{v_{22}}{\tau} \mathbf{M}_{11}^T \boldsymbol{\Sigma}_x \mathbf{B} - \frac{v_{22}}{\tau} \mathbf{B} \boldsymbol{\Sigma}_x \mathbf{M}_{11} + \mathbf{B}, \quad (93)$$

where we first do basic algebraic manipulations. Then, we use Lemma 8 and  $\mathbb{E}[\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T] = \boldsymbol{\Sigma}_x$  to get the penultimate line. For the final line, we utilize  $l \rightarrow \infty$  by Assumption 2.

Putting the found expectation results into (77), we get

$$\mathbb{E} [\mathbb{E}_w[\mathbf{w}_{diff} \mathbf{w}_{diff}^T]] = \mathbb{E} [\mathbf{e} \mathbf{e}^T] + \mathbb{E} [\mathbf{D} \boldsymbol{\mu}_w \mathbf{e}^T]^T + \mathbb{E} [\mathbf{D} \boldsymbol{\mu}_w \mathbf{e}^T] + \mathbb{E} [DBD^T], \quad (94)$$

$$= \frac{1}{\tau^2} \mathbf{M}_{11}^T \mathbf{F}_1 \mathbf{M}_{11} - \frac{1}{\tau} \mathbf{F}_2 \mathbf{M}_{11} + \frac{1}{\tau} \mathbf{M}_{11}^T \mathbf{F}_2^T + \mathbf{B}. \quad (95)$$

where matrices  $\mathbf{F}_1$  and  $\mathbf{F}_2$  are defined as

$$\mathbf{F}_1 := v_{22}^2 \frac{\sigma^2}{l} \boldsymbol{\Sigma}_x + v_{22} \left( \boldsymbol{\Sigma}_x \boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x + \frac{1}{l} \text{Tr}(\boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right) \quad (96)$$

$$\begin{aligned} &+ v_{22} \left( \boldsymbol{\Sigma}_x \boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x + \frac{1}{l} \text{Tr}(\boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right)^T + v_{22}^2 \left( \boldsymbol{\Sigma}_x \mathbf{B} \boldsymbol{\Sigma}_x + \frac{1}{l} \text{Tr}(\mathbf{B} \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right), \\ &= \left( \boldsymbol{\Sigma}_x \hat{\mathbf{B}} + \left( v_{22}^2 \frac{\sigma^2}{l} + \frac{1}{l} \text{Tr}(\hat{\mathbf{B}} \boldsymbol{\Sigma}_x) \right) \mathbf{I} \right) \boldsymbol{\Sigma}_x, \end{aligned} \quad (97)$$

$$\mathbf{F}_2 := \boldsymbol{\mu}_w \mathbf{v}_{21}^T \boldsymbol{\Sigma}_x + v_{22} \mathbf{B} \boldsymbol{\Sigma}_x = (\boldsymbol{\mu}_w \mathbf{v}_{21}^T + v_{22} \mathbf{B}) \boldsymbol{\Sigma}_x, \quad (98)$$

with  $\hat{\mathbf{B}} := v_{22} \boldsymbol{\mu}_w \mathbf{v}_{21}^T + v_{22} \mathbf{v}_{21} \boldsymbol{\mu}_w^T + v_{22}^2 \mathbf{B}$ .

Going back to generalization error in (70), we have

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) = \text{Tr}(\mathbf{A}\mathbb{E}[\mathbf{w}_{diff}\mathbf{w}_{diff}^T]) + \sigma^2, \quad (99)$$

$$= \text{Tr}\left(\mathbf{A}\left(\frac{1}{\tau^2}\mathbf{M}_{11}^T\mathbf{F}_1\mathbf{M}_{11} - \frac{1}{\tau}\mathbf{F}_2\mathbf{M}_{11} + \frac{1}{\tau}\mathbf{M}_{11}^T\mathbf{F}_2^T + \mathbf{B}\right)\right) + \sigma^2, \quad (100)$$

where  $\mathbf{F}_1 = \left(\boldsymbol{\Sigma}_x\hat{\mathbf{B}} + \frac{1}{l}\left(v_{22}^2\sigma^2 + \text{Tr}(\hat{\mathbf{B}}\boldsymbol{\Sigma}_x)\right)\mathbf{I}\right)\boldsymbol{\Sigma}_x$ , and  $\mathbf{F}_2 = (\boldsymbol{\mu}_w\mathbf{v}_{21}^T + v_{22}\mathbf{B})\boldsymbol{\Sigma}_x$ . Furthermore,  $\hat{\mathbf{B}}$  is defined as  $\hat{\mathbf{B}} := v_{22}\boldsymbol{\mu}_w\mathbf{v}_{21}^T + v_{22}\mathbf{v}_{21}\boldsymbol{\mu}_w^T + v_{22}^2\mathbf{B}$ .

## Appendix H. Proof of Lemma 8

We first restate the lemma as follows.

Let  $\bar{\mathbf{x}} \sim \mathcal{N}(0, \boldsymbol{\Sigma})$ , where  $\bar{\mathbf{x}} \in \mathbb{R}^d$ . Let  $\bar{\mathbf{x}}_i$  be  $l$  independent samples of  $\bar{\mathbf{x}}$  for  $i = 1, \dots, l$ . Let  $\mathbf{A}$  be a fixed  $d \times d$  matrix. Then, the following holds

$$\mathbb{E}\left[\left(\frac{1}{l}\sum_{i \leq l}\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T\right)\mathbf{A}\left(\frac{1}{l}\sum_{i \leq l}\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T\right)\right] = \boldsymbol{\Sigma}\mathbf{A}\boldsymbol{\Sigma} + \frac{1}{l}\boldsymbol{\Sigma}\mathbf{A}^T\boldsymbol{\Sigma} + \frac{1}{l}\text{Tr}(\mathbf{A}\boldsymbol{\Sigma})\boldsymbol{\Sigma}. \quad (101)$$

**Proof** Let  $\mathbf{S}_x = \frac{1}{l}\sum_{i=1}^l\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T$ . First, note that  $E[\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T] = \boldsymbol{\Sigma}$  since  $\bar{\mathbf{x}}_i \sim \mathcal{N}(0, \boldsymbol{\Sigma})$ .

Thus,  $\mathbb{E}[\mathbf{S}_x] = \frac{1}{l}\sum_{i=1}^l\mathbb{E}[\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T] = \frac{1}{l}\sum_{i=1}^l\boldsymbol{\Sigma} = \boldsymbol{\Sigma}$ . We have

$$\mathbf{S}_x\mathbf{A}\mathbf{S}_x = \frac{1}{l^2}\sum_{i=1}^l\sum_{j=1}^l\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T\mathbf{A}\bar{\mathbf{x}}_j\bar{\mathbf{x}}_j^T \quad (102)$$

Taking the expectation, we get

$$\mathbb{E}[\mathbf{S}_x\mathbf{A}\mathbf{S}_x] = \frac{1}{l^2}\sum_{i=1}^l\sum_{j=1}^l\mathbb{E}[\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T\mathbf{A}\bar{\mathbf{x}}_j\bar{\mathbf{x}}_j^T] \quad (103)$$

When  $i \neq j$ ,  $\bar{\mathbf{x}}_i$  and  $\bar{\mathbf{x}}_j$  are independent, so

$$\mathbb{E}[\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T\mathbf{A}\bar{\mathbf{x}}_j\bar{\mathbf{x}}_j^T] = \mathbb{E}[\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T]\mathbf{A}\mathbb{E}[\bar{\mathbf{x}}_j\bar{\mathbf{x}}_j^T] = \boldsymbol{\Sigma}\mathbf{A}\boldsymbol{\Sigma} \quad (104)$$

When  $i = j$ ,

$$\mathbb{E}[\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T\mathbf{A}\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T] = \mathbb{E}[\bar{\mathbf{x}}\bar{\mathbf{x}}^T\mathbf{A}\bar{\mathbf{x}}\bar{\mathbf{x}}^T] \quad (105)$$

Let  $\bar{\mathbf{x}} = [x_1, x_2, \dots, x_d]^T$ . Then, from Isserlis' theorem [12], we have

$$\mathbb{E}[x_i x_j x_k x_l] = \boldsymbol{\Sigma}_{ij}\boldsymbol{\Sigma}_{kl} + \boldsymbol{\Sigma}_{ik}\boldsymbol{\Sigma}_{jl} + \boldsymbol{\Sigma}_{il}\boldsymbol{\Sigma}_{jk} \quad (106)$$

Let  $\mathbf{A} = [a_{ij}]$ . Then,  $\bar{\mathbf{x}}^T\mathbf{A}\bar{\mathbf{x}} = \sum_{i,j}a_{ij}x_i x_j$ . Thus, we reach

$$\bar{\mathbf{x}}\bar{\mathbf{x}}^T\mathbf{A}\bar{\mathbf{x}}\bar{\mathbf{x}}^T = \bar{\mathbf{x}}\bar{\mathbf{x}}^T\sum_{i,j}a_{ij}x_i x_j, \quad (107)$$

$$\mathbb{E}[\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T\mathbf{A}\bar{\mathbf{x}}_i\bar{\mathbf{x}}_i^T] = \text{Tr}(\mathbf{A}\boldsymbol{\Sigma})\boldsymbol{\Sigma} + \boldsymbol{\Sigma}\mathbf{A}\boldsymbol{\Sigma} + \boldsymbol{\Sigma}\mathbf{A}^T\boldsymbol{\Sigma}. \quad (108)$$

There are  $l^2$  terms in the double sum.  $l$  terms are of the form  $\mathbb{E}[\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \mathbf{A} \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T]$  and  $l^2 - l$  terms are of the form  $\Sigma \mathbf{A} \Sigma$ . Therefore, we can write

$$\mathbb{E}[\mathbf{S}_x \mathbf{A} \mathbf{S}_x] = \frac{1}{l^2} [l(\text{Tr}(\mathbf{A} \Sigma) \Sigma + \Sigma \mathbf{A} \Sigma + \Sigma \mathbf{A}^T \Sigma) + l(l-1) \Sigma \mathbf{A} \Sigma], \quad (109)$$

$$= \frac{1}{l} (\text{Tr}(\mathbf{A} \Sigma) \Sigma + \Sigma \mathbf{A} \Sigma + \Sigma \mathbf{A}^T \Sigma) + \frac{l-1}{l} \Sigma \mathbf{A} \Sigma, \quad (110)$$

$$= \Sigma \mathbf{A} \Sigma + \frac{1}{l} \Sigma \mathbf{A}^T \Sigma + \frac{1}{l} \text{Tr}(\mathbf{A} \Sigma) \Sigma, \quad (111)$$

which completes the proof. ■

### Appendix I. Analysis of optimal temperature for ICL under distribution shift

Here, we find the optimal temperature minimizing the generalization error. First, recall that we have the following generalization error.

$$\mathcal{G}(\mathbf{V}, \mathbf{M}) = \frac{1}{\tau^2} \text{Tr}(\mathbf{A} \mathbf{M}_{11}^T \mathbf{F}_1 \mathbf{M}_{11}) - \frac{1}{\tau} \text{Tr}(\mathbf{A} (\mathbf{F}_2 \mathbf{M}_{11} + \mathbf{M}_{11}^T \mathbf{F}_2^T)) + \text{Tr}(\mathbf{A} \mathbf{B}) + \sigma^2, \quad (112)$$

as specified in Theorem 2. So, we can express the generalization error as,

$$\mathcal{G}(\tau; \mathbf{V}, \mathbf{M}) = \frac{a}{\tau^2} - \frac{b}{\tau} + c, \quad (113)$$

where  $a := \text{Tr}(\mathbf{A} \mathbf{M}_{11}^T \mathbf{F}_1 \mathbf{M}_{11})$ ,  $b := \text{Tr}(\mathbf{A} (\mathbf{F}_2 \mathbf{M}_{11} + \mathbf{M}_{11}^T \mathbf{F}_2^T))$ , and  $c = \text{Tr}(\mathbf{A} \mathbf{B}) + \sigma^2$ . Therefore, we have the following optimization problem

$$\tau_{\text{opt}} := \arg \min_{\tau} \mathcal{G}(\tau; \mathbf{V}, \mathbf{M}), \quad (114)$$

$$= \arg \min_{\tau} \left\{ \frac{a}{\tau^2} - \frac{b}{\tau} + c \right\}. \quad (115)$$

To find the optimal value of  $\tau$  that minimizes the given function, we can take the derivative of the expression with respect to  $\tau$  and set it to zero. From now on, we consider generalization error as a function of  $\tau$ , written as  $\mathcal{G}(\tau)$ .

Next, find the derivative of  $\mathcal{G}(\tau)$  with respect to  $\tau$  as

$$\mathcal{G}'(\tau) = -2a\tau^{-3} + b\tau^{-2}. \quad (116)$$

To find the critical points, set  $\mathcal{G}'(\tau) = 0$  as follows

$$\mathcal{G}'(\tau) = -2a\tau^{-3} + b\tau^{-2} = 0, \quad (117)$$

Solving this equation for  $\tau$ , we reach the following critical point

$$\tau = \frac{2a}{b}. \quad (118)$$

Now, we need to check if this is a minimum by taking the second derivative, which is

$$\mathcal{G}''(\tau) = 6a\tau^{-4} - 2b\tau^{-3}. \quad (119)$$

Evaluate  $\mathcal{G}''(\tau)$  at  $\tau = \frac{2a}{b}$  as follows

$$\mathcal{G}''\left(\frac{2a}{b}\right) = 6a\left(\frac{2a}{b}\right)^{-4} - 2b\left(\frac{2a}{b}\right)^{-3} = 6a\left(\frac{b^4}{16a^4}\right) - 2b\left(\frac{b^3}{8a^3}\right) = \frac{b^4}{8a^3}. \quad (120)$$

Since  $a, b > 0$ , we reach  $\mathcal{G}''\left(\frac{2a}{b}\right) = \frac{b^4}{8a^3} > 0$ , which means the function has a minimum at  $\tau = \frac{2a}{b}$ . Therefore,  $\tau_{\text{opt}} = \frac{2a}{b}$  is the solution minimizing the generalization error  $\mathcal{G}(\tau)$ . Writing  $a, b$  back into the optimal solution, we get

$$\tau_{\text{opt}} = \frac{2\text{Tr}(\mathbf{A}\mathbf{M}_{11}^T\mathbf{F}_1\mathbf{M}_{11})}{\text{Tr}(\mathbf{A}(\mathbf{F}_2\mathbf{M}_{11} + \mathbf{M}_{11}^T\mathbf{F}_2^T))}, \quad (121)$$

which concludes our derivation of the optimal temperature  $\tau_{\text{opt}}$ .

### I.1. Effect of distribution shift

In this section, we explore scenarios where  $\mathcal{D}^{\text{test}} \neq \mathcal{D}^{\text{train}}$ , indicating a shift in the input, task, or noise distribution after pretraining the model. We consider three cases: (1) a shift in the input distribution (altered mean or covariance), (2) a shift in the task distribution, and (3) a change in the noise levels.

**Pretraining for explaining distribution shift effects** — To streamline our explanations below, for pretraining, we optimize the parameters  $\mathbf{V}$  and  $\mathbf{M}$  using  $m$  samples of  $(\mathbf{Z}, y_l)$  drawn from the distribution  $\mathcal{D}^{\text{train}}$ , where each  $\mathbf{Z}$  contains  $l - 1$   $(x, y)$  pairs intended for ICL. Building upon prior work that connects ICL in linear regression to the Bayes-optimal ridge estimator [24, 36], we configure  $\mathbf{M}$  and  $\mathbf{V}$  to emulate Bayes-optimal ridge regression (see Proposition 5).

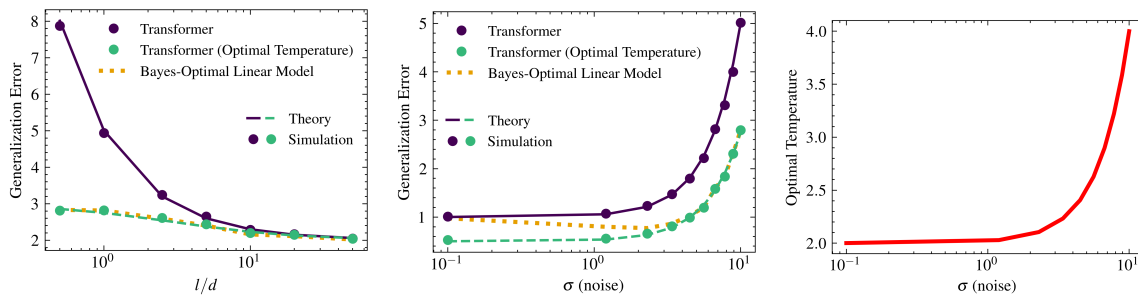
Since the pretrained model succeeds in ICL for  $\mathcal{D}^{\text{test}} = \mathcal{D}^{\text{train}}$ , we next investigate how distribution shifts affect its ICL performance.

**ICL under distribution shift** — To evaluate the impact of these distribution shifts on ICL performance, we assess whether adjustments to  $\mathbf{M}$  and/or  $\mathbf{V}$  are necessary to match the Bayes-optimal linear model under the new distribution. If so, the model is considered sensitive to the shift. Otherwise, it is deemed robust.

**Case I: Shift in input distribution** — Recall that inputs are drawn as  $x_i \sim \mathcal{N}(\mu_x, \Sigma_x)$ , as defined in (1). Let  $\mu_x^{\text{train}}, \Sigma_x^{\text{train}}$  and  $\mu_x^{\text{test}}, \Sigma_x^{\text{test}}$  denote the input means and covariances for pretraining and testing, respectively. We consider two subcases:

- (i) Mean shift ( $\mu_x^{\text{train}} \neq \mu_x^{\text{test}}$ ): Centering renders the approximate model invariant to mean shifts, but the uncentered linear attention model remains sensitive, as noted in Remark 4.
- (ii) Covariance shift ( $\Sigma_x^{\text{train}} \neq \Sigma_x^{\text{test}}$ ): Since  $\mathbf{M}_{11}$  is fitted to the pretraining covariance, a mismatch drives the estimator away from Bayes-optimality, echoing prior results on linear attention [36].

**Case II: Shift in task distribution** — The task vectors follow  $w \sim \mathcal{N}(\mu_w, \Sigma_w)$ . Let  $\mu_w^{\text{train}}, \Sigma_w^{\text{train}}$  and  $\mu_w^{\text{test}}, \Sigma_w^{\text{test}}$  be the mean and covariance of the task distribution during pretraining and testing, respectively. The Transformer model can incorporate  $\mu_w^{\text{train}}$  and  $\Sigma_w^{\text{train}}$  via the pretrained parameters  $\mathbf{M}_{11}$  and  $\mathbf{v}_{21}$  (see Proposition 5). However, as the context length  $l$  increases, the model's dependence on the task distribution diminishes. Thus, shifts in the task distribution primarily affect ICL performance for small  $l$ .



(a) Effect of  $l/d$  when  $\sigma_{test} = 10$       (b) Effect of  $\sigma_{test}$  when  $l/d = 1$       (c) Optimal temperature

Figure 4: Effect of noise shift on Transformer (3). The pretraining noise is  $\sigma_{train} = 0.1$ , while  $\sigma_{test}$  varies across plots. The optimal temperature is set by Theorem 3. This setting matches Figure 1a, except for changes in test-time noise  $\sigma_{test}$ .

**Case III: Shift in noise distribution** — Finally, consider a change in the noise distribution:  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ , with  $\sigma_{train}^2$  and  $\sigma_{test}^2$  denoting pretraining and testing noise variances. If  $\sigma_{train}^2 \neq \sigma_{test}^2$ , the parameters  $M_{11}$  and  $v_{21}$  become suboptimal relative to the Bayes-optimal linear model. However, as with the task distribution, the impact of noise shift diminishes as  $l \rightarrow \infty$ .

**Summary** — The Transformer model is robust to shifts in input mean but sensitive to input covariance changes. Shifts in task or noise distribution reduce ICL performance at small  $l$ , though increasing  $l$  mitigates these effects.

We further evaluate robustness to label noise in Figure 4. In Figure 4a, we observe that noise effects diminish as the context length increases, consistent with our theoretical predictions. However, at small  $l$ , temperature adjustment becomes critical. In Figure 4b (for  $l = d$ ), the Transformer increasingly diverges from the Bayes-optimal model as noise grows, yet optimal temperature correction closes this gap. Figure 4c shows that the optimal temperature increases with noise level, indicating a principled relationship between noise and temperature under limited context.

## Appendix J. An insight driven from optimal temperature for other settings

In this section, we extract a mathematical heuristic from the optimal temperature in Theorem 3 that can be applied to ICL settings beyond our existing setting involving approximate softmax attention and regression tasks. Specifically, we consider Transformers employing standard softmax attention. Recall that the attention temperature scales the pre-softmax scores  $(KZ)^\top(QZ)$ , thereby controlling the variance of the final scores. Since the optimal temperature depends on the distribution of these scores, it can be naturally characterized by the moments of that distribution. Our central intuition is that the optimal temperature identified in Theorem 3 relates directly to the first two moments of the pre-softmax scores. Although this optimal temperature was derived for *approximate softmax* attention, the insight remains relevant for softmax attention because the two mechanisms behave similarly in the regime considered (see Appendix D).

We now illustrate how the optimal temperature in Theorem 3 can be related to the first two moments of the pre-softmax scores. For simplicity, we consider the case  $\mu_x = \mu_w = m_{21} = 0$  and

$\Sigma_w = I$ , under which the optimal temperature reduces to

$$\tau_{\text{opt}} = \frac{v_{22} \text{Tr}(\Sigma_x M_{11} \Sigma_x M_{11}^\top \Sigma_x)}{\frac{1}{2} \text{Tr}(\Sigma_x (\Sigma_x M_{11} + M_{11}^\top \Sigma_x))}. \quad (122)$$

We next show how this expression connects to the first two moments of  $(KZ)^\top(QZ)$ . Let  $z_i$  denote the  $i$ -th column of  $Z$  from (2) and recall  $K^\top Q = M$ . We therefore compute  $\mathbb{E}[z_i^\top M z_j]$  and  $\mathbb{E}[(z_i^\top M z_j)^2]$  for  $i, j \in \{1, \dots, l\}$ . Starting with the first moment for  $i = j$ :

$$\mathbb{E}[z_i^\top M z_i] = \text{Tr}(\mathbb{E}[z_i z_i^\top] M) = \text{Tr}(\Sigma_x M_{11}), \quad (123)$$

where the block structure (and zero entries) of  $M$  is used in the last step. For  $i \neq j$ ,

$$\mathbb{E}[z_i^\top M z_j] = \text{Tr}(M \mathbb{E}[z_i z_j^\top]) = 0, \quad (124)$$

by independence of  $z_i$  and  $z_j$ . For the second moment with  $i \neq j$ :

$$\mathbb{E}[(z_i^\top M z_j)^2] = \mathbb{E}[z_i^\top M z_j z_j^\top M^\top z_i] \quad (125)$$

$$= \mathbb{E}[x_i^\top M_{11} x_j x_j^\top M_{11}^\top x_i] \quad (126)$$

$$= \mathbb{E}_{x_i} [x_i^\top M_{11} \mathbb{E}_{x_j} [x_j x_j^\top] M_{11}^\top x_i] \quad (127)$$

$$= \mathbb{E}_{x_i} [x_i^\top M_{11} \Sigma_x M_{11}^\top x_i] \quad (128)$$

$$= \text{Tr}(M_{11} \Sigma_x M_{11}^\top \mathbb{E}_{x_i} [x_i x_i^\top]) \quad (129)$$

$$= \text{Tr}(M_{11} \Sigma_x M_{11}^\top \Sigma_x), \quad (130)$$

where we again exploit the block structure of  $M$  and apply straightforward manipulations.

We observe a parallel between the numerator of (122) and the computed second moment (for  $i \neq j$ ), and between the denominator and the first moment (for  $i = j$ ). This motivates the heuristic that the optimal temperature should be roughly proportional to the ratio of the second moment (for  $i \neq j$ ) to the first moment (for  $i = j$ ). Accordingly, in our LLM experiments (Figure 2), we select the temperature proportional to this ratio while taking care to avoid numerical issues.

Finally, we note an important caveat: in order to obtain an insight of practical relevance, we intentionally relaxed the rigor applied in our main theoretical results. Consequently, the heuristic derived here—and the accompanying empirical findings—should be viewed as preliminary, intended to inspire future work on principled selection of attention temperature in practice.

## Appendix K. Experimental details and GPT-2 experiments

This section describes our experimental setups for GPT-2 and large language models (LLMs), including the motivation for our distribution-shift scenarios.

### K.1. GPT-2: Transformer with MLP layers

Building on the approximate-attention experiments, we investigate whether the optimal temperature also benefits more complex Transformer models on linear regression tasks. We evaluate GPT-2 [23]

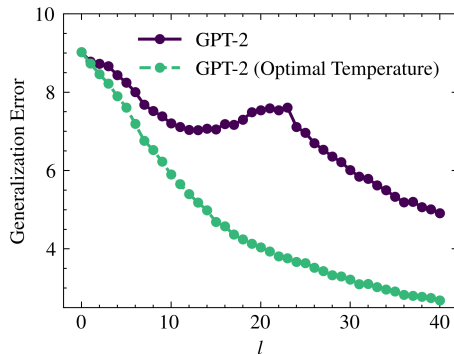


Figure 5: GPT-2 [23] under an input-covariance shift. GPT-2 exemplifies the Transformer architecture [28], combining multi-layer perceptrons with multi-head softmax self-attention. The model here is pretrained by Garg et al. [9] on the linear regression tasks defined in (1). We consider a shift from  $\Sigma_x^{\text{train}} = \mathbf{I}$  to  $\Sigma_x^{\text{test}} = 3\mathbf{I}$ . The attention temperature at each layer is scaled as  $\tau\sqrt{d_k}$  (where  $d_k$  is the key dimension) to ensure dimension-independent  $\tau$  values.

under a shift in input covariance (Figure 5). Consistent with prior work [9, 36], such shifts substantially degrade performance and can even induce nonmonotonic generalization error with respect to context length  $l$ . Remarkably, applying the optimal temperature mitigates this nonmonotonicity and improves in-context generalization.

### K.2. Details of the GPT-2 experiments in Figure 5

We use the standard GPT-2 architecture [23] as implemented in HuggingFace [33], leveraging the pretrained model of Garg et al. [9]. Training data match ours, while their training procedure differs slightly: the loss is auto-regressive, i.e., the average over the entire context sequence of length  $l = 40$ . We adopt the same embedding method as in Garg et al. [9]. The input dimension is  $d = 20$ , with 12 layers and 8 heads. All GPT-2 experiments run on an NVIDIA Tesla V100 GPU and complete in approximately 10 minutes.

### K.3. Details of the LLM experiments in Figure 2

For our large language model experiments, we employ Llama2-7B [27] and the SCIQ dataset [32], which contains science questions with supporting information. We generate ICL problems following Gao et al. [8], selecting in-context demonstrations using the TopK retrieval technique [17] to ensure relevance. An example ICL sample from SCIQ appears in Table 1. To simulate distribution shift, we follow Gao et al. [8] and introduce noisy labels—incorrect but semantically related—to the in-context demonstrations (Appendix K.4). Table 2 gives an example. The noisy ratio denotes the fraction of demonstrations with noisy labels (e.g., 0.6 means 60% noisy). We modify and use the codebase of Gao et al. [8], built on HuggingFace [33] and OpenICL [35]. The attention temperature at each layer is scaled as  $\tau\sqrt{d_k}$  (where  $d_k$  is the key dimension) so that the reported  $\tau$  values are dimension-independent. All LLM experiments run on an NVIDIA A40 GPU; a single Monte Carlo run per plot in Figure 2 takes a few hours.

**K.4. Why in-context demonstrations with noisy labels as an example of distribution shift?**

The link between noisy labels in demonstrations and distribution shift may not be immediately obvious. Quantifying pretraining–test shifts for pretrained LLMs is inherently difficult because their pretraining data are complex mixtures of sources [27]. However, we hypothesize—following Gao et al. [8]—that high perplexity can serve as an empirical indicator of distribution shift. Inputs aligned with the training distribution tend to yield low perplexity (high-confidence generation), whereas contradictory or out-of-distribution inputs induce high perplexity. Since noisy demonstrations are expected to contradict training-set patterns, they yield high perplexity and thereby act as a proxy for distribution shift. Consequently, introducing noisy labels into in-context demonstrations constitutes a principled way to test the robustness of in-context learning under distribution shift.

**In-context demonstration 1**


---

**Support:** Cells are organized into tissues, tissues are organized into organs.  
**Question:** What is considered the smallest unit of the organ?  
**Answer:** Cells

---

**In-context demonstration 2**


---

**Support:** ... four basic types of tissue: connective, muscle, **nervous**, and epithelial.  
**Question:** The four basic types of tissue are epithelial, muscle, connective, and what?  
**Answer:** **nervous**

---

⋮

---

**Test example**


---

**Support:** All forms of life are built of at least one cell. A cell is the basic unit of life.  
**Question:** What are the smallest structural and functional units of all living organisms?  
**Output:** ???

---

Table 1: A sample illustration of in-context learning on the SCIQ dataset.

Setting	In-context demonstration
True Label	<p><b>Support:</b> Cells are organized into tissues, tissues are organized into organs.</p> <p><b>Question:</b> What is considered the smallest unit of the organ?</p> <p><b>Label:</b> Cells</p>
Noisy Label	<p><b>Support:</b> Cells are organized into tissues, tissues are organized into organs.</p> <p><b>Question:</b> What is considered the smallest unit of the organ?</p> <p><b>Label:</b> tissues</p>

Table 2: An example of a true label vs. a relevant but noisy label. A relevant label is related to the question but is not necessarily true. Therefore, relevant labels can be considered noisy labels.

### Appendix L. Interpretation of the optimal attention temperature: full derivation

To obtain a simplified expression, we consider a simple but representative family of shifts as follows. The training distributions are

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \epsilon_i \sim \mathcal{N}(0, \hat{\sigma}^2).$$

We then introduce three independent shift parameters for the test distribution:

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, a\mathbf{I}), \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, b\mathbf{I}), \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2),$$

where  $a > 0$  controls the input variance shift,  $b > 0$  controls the task-parameter variance shift, and  $\sigma > 0$  controls the noise-variance shift. This setting preserves isotropy, which makes it possible to derive a clean closed-form expression while still connecting directly to realistic distribution shifts.

Substituting these shifted distributions into the optimal-temperature expression in (7) yields

$$\tau_{\text{opt}} = \frac{2\text{Tr}(a\mathbf{I}\mathbf{M}_{11}^T(ab\mathbf{I} + \frac{1}{l}(\sigma^2 + abd)\mathbf{I})a\mathbf{I}\mathbf{M}_{11})}{\text{Tr}(a\mathbf{I}(ab\mathbf{I}(\mathbf{M}_{11} + \mathbf{M}_{11}^T)))}, \quad (131)$$

$$= \left(a + \frac{1}{l} \left(\frac{\sigma^2}{b} + ad\right)\right) \frac{\text{Tr}(\mathbf{M}_{11}^T\mathbf{M}_{11})}{\text{Tr}(\mathbf{M}_{11})}, \quad (132)$$

recovering (8).

This concrete formula makes several effects fully explicit:

- *Input shift* — Increasing input variance  $a$  directly scales  $\tau_{\text{opt}}$  upward. This aligns with our earlier results (Figure 1) and the heuristic derived in Appendix J, which suggests that a greater variance of pre-softmax scores requires a higher temperature to maintain robustness to input shifts.
- *Noise shift* — Increasing noise variance  $\sigma^2$  also increases  $\tau_{\text{opt}}$ , but only through the  $\frac{1}{l}$  term, reflecting the diminishing effect of noise when more in-context examples are available.

- *Task shift* — Increasing task variance  $b$  reduces the effect of noise (via  $\sigma^2/b$ ), slightly lowering the optimal temperature.
- *Context length* — As  $l \rightarrow \infty$ , the  $\frac{1}{l}$  term vanishes, giving a simplified asymptotic rule:  $\tau_{\text{opt}} \rightarrow a \cdot \frac{\text{Tr}(\mathbf{M}_{11}^\top \mathbf{M}_{11})}{\text{Tr}(\mathbf{M}_{11})}$ .

Note that under the considered training distribution, we have  $\text{Tr}(\mathbf{M}_{11}^\top \mathbf{M}_{11})/\text{Tr}(\mathbf{M}_{11}) \approx 1$ , which implies that  $\tau_{\text{optimal}} \rightarrow a$  as  $l \rightarrow \infty$ .

**Moment-ratio heuristic** — We propose a practical heuristic (derived in Appendix J): the optimal attention temperature is roughly proportional to the ratio of the second and first moments of pre-softmax attention scores. The expression in (8) provides further theoretical justification for that heuristic. Indeed, for some  $i \neq j$ ,

$$\tau_{\text{opt}} = a \frac{\text{Tr}(\mathbf{M}_{11} \mathbf{M}_{11}^\top)}{\text{Tr}(\mathbf{M}_{11})} + \frac{1}{l} \left( \frac{\sigma^2}{b} + ad \right) \frac{\text{Tr}(\mathbf{M}_{11} \mathbf{M}_{11}^\top)}{\text{Tr}(\mathbf{M}_{11})}, \quad (133)$$

$$= \frac{\mathbb{E}[(\mathbf{z}_i^\top \mathbf{M} \mathbf{z}_j)^2]}{\mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \mathbf{z}_i]} + \frac{1}{l} \left( \frac{\sigma^2}{b} + ad \right) \frac{\text{Tr}(\mathbf{M}_{11} \mathbf{M}_{11}^\top)}{\text{Tr}(\mathbf{M}_{11})}, \quad (134)$$

where we used the moments calculated in Appendix J to reach the final line. Since  $\text{Tr}(\mathbf{M}_{11}^\top \mathbf{M}_{11})/\text{Tr}(\mathbf{M}_{11}) \approx 1$  for the considered training distribution, this gives the approximation:

$$\tau_{\text{opt}} \approx \underbrace{\frac{\mathbb{E}[(\mathbf{z}_i^\top \mathbf{M} \mathbf{z}_j)^2]}{\mathbb{E}[\mathbf{z}_i^\top \mathbf{M} \mathbf{z}_i]}}_{\text{moment-ratio}} + \underbrace{\frac{1}{l} \left( \frac{\sigma^2}{b} + ad \right)}_{\text{correction for small } l}, \quad (135)$$

recovering (9). This demonstrates that the moment-ratio heuristic is not an ad-hoc rule, but a theoretically grounded approximation of the exact closed-form optimal temperature.

## Appendix M. Numerical illustration of the optimal temperature and its generalization behavior

To complement the analytical reductions in the setting of Section 3, we now present numerical experiments illustrating how the optimal attention temperature varies under different types of distribution shift. These simulations closely follow the structure predicted by the closed-form expression in Theorem 3 and its simplified forms in (8) and (9).

Figure 6 shows the optimal temperature as a function of (i) input covariance shift, (ii) task covariance shift, and (iii) noise-variance shift. In each subplot, a single shift parameter ( $a$ ,  $b$ , or  $\sigma$ ) is varied while the others remain fixed. The closed-form optimal temperatures (7) align closely with the moment-ratio heuristic with correction (9). As anticipated:

- higher input variance  $a$  or noise level  $\sigma$  increases the optimal temperature, while
- task variance  $b$  does not significantly change the optimal temperature.

Figure 7 presents the corresponding generalization errors. These results demonstrate that the closed-form characterization accurately captures the key dependencies of the optimal temperature under a range of distribution shifts.

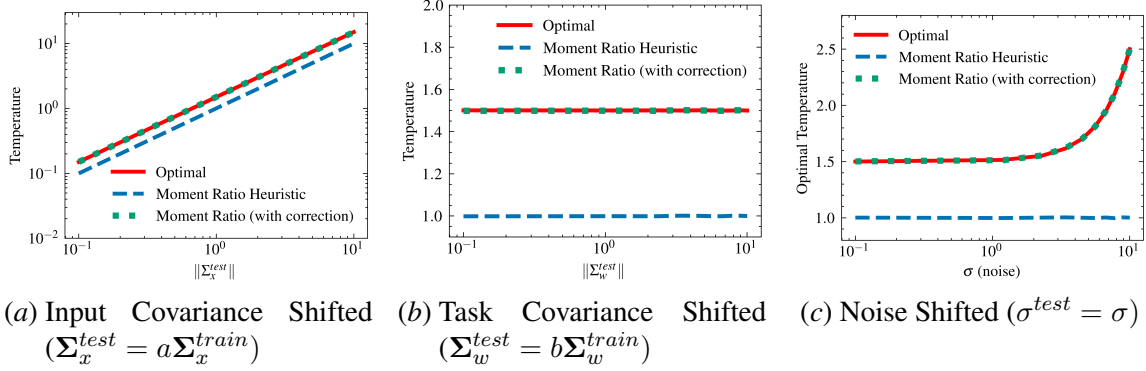


Figure 6: Optimal temperature under different types of distribution shift. The moment-ratio heuristic (Appendix J) is derived from the closed-form optimal temperature, with its corrected form given in (9). During training, we use  $m = 5000$  tasks, noise level  $\sigma^{train} = 0.1$ , means  $\mu_x^{train} = \mu_w^{train} = \mathbf{0}$ , and covariances  $\Sigma_x^{train} = \Sigma_w^{train} = \mathbf{I}$ . At test time, we set  $\mu_x^{test} = \mu_w^{test} = \mathbf{0}$ ,  $\Sigma_x^{test} = a\mathbf{I}$ ,  $\Sigma_w^{test} = b\mathbf{I}$ , and  $\sigma^{test} = \sigma$ . In each subplot, exactly one of  $a$ ,  $b$ , or  $\sigma$  is varied, while the other two remain fixed at their training-distribution values to isolate the effect of a single shift dimension. The dimension and context length are set to  $d = 50$  and  $l = 2d$ .

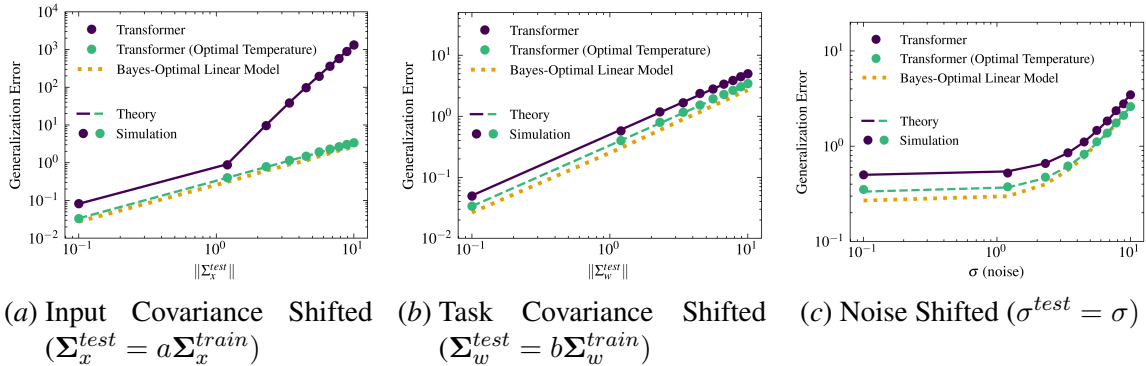


Figure 7: Generalization errors corresponding to Figure 6.