

---

# The Evolution of Statistical Induction Heads: In-Context Learning Markov Chains

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models have the ability to generate text that mimics patterns in  
2 their inputs. We introduce a simple Markov Chain (MC) sequence modeling task in  
3 order to study how this in-context learning (ICL) capability emerges. Transformers  
4 trained on this task (ICL-MC) form *statistical induction heads* which compute  
5 accurate next-token probabilities given the bigram statistics of the context. During  
6 the course of training, models pass through multiple phases: after an initial stage  
7 in which predictions are uniform, they learn to sub-optimally predict using in-  
8 context single-token statistics (unigrams); then, there is a rapid phase transition to  
9 the correct in-context bigram solution. We conduct an empirical and theoretical  
10 investigation of this multi-phase process, showing how successful learning results  
11 from the interaction between the transformer’s layers, and uncovering evidence that  
12 the presence of simpler solutions delays formation of the final optimal solutions.

## 13 1 Introduction

14 Large language models (LLMs) exhibit a remarkable ability to perform *in-context learning* (ICL)  
15 from patterns in their input context [10, 14]. The ability of LLMs to adaptively learn from context is  
16 profoundly useful, yet the underlying mechanisms of this emergent capability are not fully understood.

17 In an effort to better understand ICL, some recent works propose to study ICL in controlled synthetic  
18 settings—in particular, training transformers on mathematically defined tasks which require learning  
19 from the input context. For example, a recent line of works studies the ability of transformers to  
20 perform ICL of standard supervised learning problems such as linear regression [2, 18, 23, 35].  
21 Studying these well-understood synthetic learning tasks enables fine-grained control over the data  
22 distribution, allows for comparisons with established supervised learning algorithms, and facilitates  
23 the examination of the in-context “algorithm” implemented by the network.

24 The goal of this work is to propose and analyze a simple synthetic setting for studying ICL. To achieve  
25 this, we consider  $n$ -gram models [9, 13, 32], one of the simplest and oldest methods for language  
26 modeling. An  $n$ -gram language model predicts the probability of a token based on the preceding  $n - 1$   
27 tokens, using fixed-size chunks ( $n$ -grams) of text data to capture linguistic patterns. Our work studies  
28 ICL of  $n$ -gram models, where the network needs to compute the conditional probability of the next  
29 token based on the statistics of the tokens observed in the input context, rather than on the statistics  
30 of the entire training data. We mainly focus on the simple case of  $n = 2$ ; i.e., bigram models, which  
31 can be represented as Markov chains. We therefore consider ICL of Markov chains (ICL-MC): we  
32 train a two layer attention-only transformer on sequences of tokens, where each sequence is produced  
33 by a different Markov chain, generated using a different transition matrix (see Figure 1 (left)).

34 We summarize our key findings:

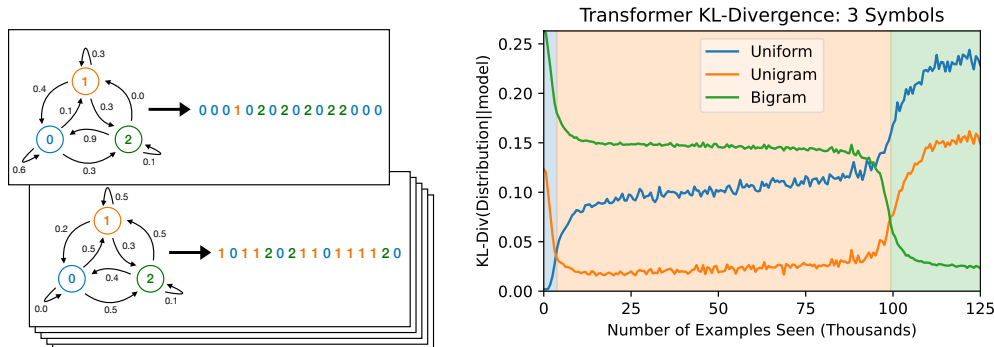


Figure 1: (*left*) We train small transformers to perform in-context learning of Markov chains (ICL-MC). Each training sequence is generated by sampling a transition matrix from a prior distribution, and then sampling a sequence from this Markov chain. (*right*) Distance of a transformer’s output distribution to several well-defined strategies over the course of training. The model passes through three stages: (1) predicting a uniform distribution (*blue* region), (2) predicting based on in-context unigram statistics (*orange* region), (3) predicting based on in-context bigram statistics (*green* region).

35 **(1) Transformers learn statistical induction heads to optimally solve ICL-MC.** We show that in  
 36 order to solve ICL-MC, transformers learn *statistical* induction heads [16] that are able to compute  
 37 the correct *conditional (posterior) probability* of the next token given all previous occurrences of the  
 38 prior token (see attention patterns in Figure 4). We show that these statistical induction heads lead to  
 39 the transformer achieving performance approaching that of the Bayes-optimal predictor.

40 **(2) Transformers learn predictors of increasing complexity and undergo a phase transition**  
 41 **when increasing complexity.** We observe that transformers display *phase transitions* when learning  
 42 Markov chains—learning appears to be separated into phases, with fast drops in loss between the  
 43 phases. We are able to show that different phases correspond to learning models of increased  
 44 complexity—unigrams, then bigrams (see Figure 1)—and characterize the transition between the  
 45 phases.

46 **(3) Simplicity bias may slow down learning.** We provide evidence that the model’s inherent bias  
 47 towards simpler solutions (in particular, in-context unigrams) causes learning of the optimal solution  
 48 to be delayed. Changing the distribution of the in-context examples to remove the usefulness of  
 49 in-context unigrams leads to faster convergence, even when evaluated on the original distribution.

50 **(4) Alignment of layers is crucial.** We show that the transition from a phase of learning the  
 51 simple-but-inadequate solution to the complex-and-correct solution happens due to an alignment  
 52 between the layers of the model: the learning signal for the first layer is tied to the extent to which  
 53 the second layer approaches its correct weights.

54 Finally, in Appendix E we provide experiments with higher order Markov Chains, where we also  
 55 observe a similar multi-stage learning process.

56 **Concurrent work.** In parallel to this work, there have been a number of papers devoted to the  
 57 study of similar questions regarding in-context learning of Markov chains [3, 19, 25]. Perhaps closest  
 58 to our work, [27] introduces a general family of in-context learning tasks with causal structure, a  
 59 special case of which is in-context Markov chains, and shows that simplified transformers (similar  
 60 to the ones we introduce in Section B.2) can learn to identify the causal relationships. The focus of  
 61 our work, instead, is on the different stages of training and how they relate to specific, well-defined,  
 62 strategies. See Section A for a detailed discussion on prior work.

## 63 2 Setup

64 **ICL-MC Task.** Our learning task consists of sequences generated from Markov Chains with  
 65 random transition matrices. The goal is to in-context estimate the transition probabilities from

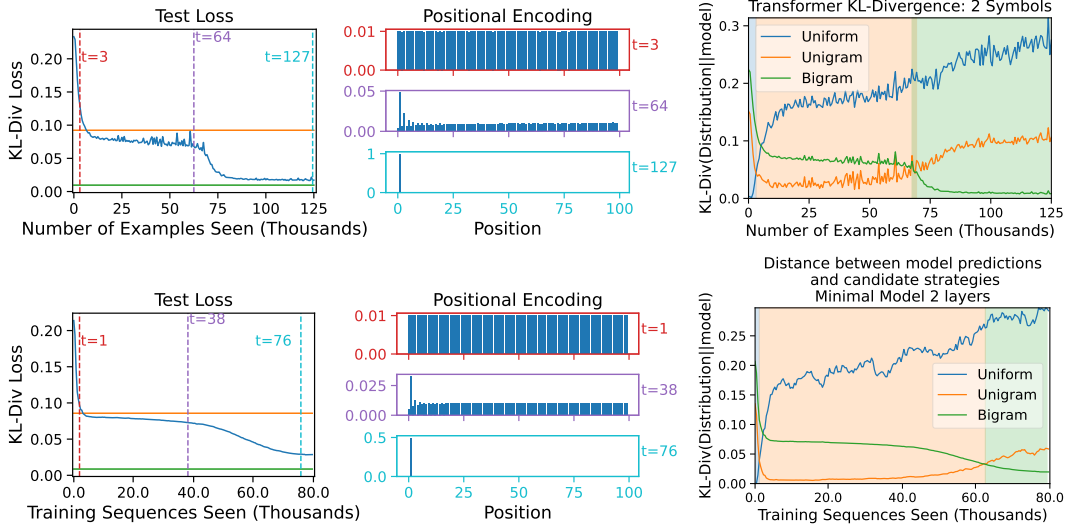


Figure 2: A two layer transformer (*top*) and a minimal model (*bottom*) trained on our in-context Markov Chain task. A comparison of the two layer attention-only transformer and minimal model (4). The graphs on the left are test loss measured by KL-Divergence from the underlying truth. The orange line shows the loss of the unigram strategy, and the green line shows the loss of the bigram strategy. The middle graph shows the effective positional encoding (for the transformer, these are for the first layer). The graph on the right shows the KL-divergence between the outputs of the models and three strategy. The lower the KL-divergence, the more similar the model is to that strategy.

66 sampled sequences, in order to predict the next state. Formally, each sample sequence is generated  
 67 by a Markov Chain with state space  $S = \{1, \dots, k\}$  and a transition matrix  $\mathcal{P}$  sampled from a prior  
 68 distribution, with  $x_1$  drawn from some other prior distribution (potentially dependent on  $\mathcal{P}$ ), and the  
 69 rest of  $\mathbf{x} = (x_1, \dots, x_t)$  drawn from the Markov Chain. We focus on the case where each row of the  
 70 matrix is sampled from the Dirichlet distribution with concentration parameter  $\alpha$ , i.e.  $\mathcal{P}_{i,:} \sim \text{Dir}(\alpha)$ .  
 71 We want to learn a predictor that, given context  $x_1, \dots, x_t$ , predicts the next token,  $x_{t+1}$ .

72 **Strategies.** We consider two particular strategies that can be employed to solve the above task: a  
 73 (suboptimal) *unigram* strategy which assumes tokens in each sequence are i.i.d. samples (and counts  
 74 the frequency of the states in the sequence so far), and the *bigram* strategy which correctly takes  
 75 into account dependencies among adjacent tokens (and counts frequency of pairs of tokens). See  
 76 Section B in the Appendix for a detailed description of our learning setup.

### 77 3 Empirical Findings and Theoretical Validation

78 In this section, we present our empirical findings on how transformers succeed in in-context learning  
 79 Markov Chains, we demonstrate the different learning stages during training and the sudden transitions  
 80 between them, and draw analytical and empirical insights from a minimal model that we believe  
 81 captures the behavior of transformers for this task.

#### 82 3.1 Transformers In-Context Learn Markov Chains Hierarchically

83 We focus on attention-only transformers with 2 layers with causal masking and relative positional  
 84 encodings and train them with the Adam optimizer on ICL-MC. As can be seen in Figure 2, all the  
 85 models converge near the Bayes optimal solution, suggesting that they learn to implement the bigram  
 86 strategy. Curiously, however, the learning seems to be happening in stages; there is an initial rapid  
 87 drop and the model quickly finds a better than random solution. Afterwards, there is a long period of  
 88 only slight improvement before a second rapid drop brings the model close to the Bayes optimal loss.

89 Interestingly, as can be seen from the horizontal lines in Figure 2, the intermediate plateau corresponds  
 90 to a phase when the model reaches the unigram baseline. We provide evidence that this is not a

91 coincidence, and that after the initial drop in loss, the model’s strategy is very similar to the unigram  
92 strategy, before eventually being overtaken by the bigram strategy (see Figure 2). This final drop is  
93 what has been associated to prior work with *induction heads* formation [28]; special dedicated heads  
94 inside a transformer are suddenly being formed to facilitate in-context learning.

95 **Mechanistic evidence for solutions found by transformer.** To confirm how the two layer attention-  
96 only transformer solves ICL-MC, we inspected the attention in each layer throughout training. Figure  
97 4 shows the attention for a particular input during different parts of training. We observe that, by the  
98 end of training, each token in the first layer is attending to the previous token. In the second layer, the  
99 last token, a “2”, is attending to tokens that followed “2”s, allowing bigram statistics to be calculated.  
100 In Proposition B.2, we show how this behavior can be implemented in the transformer architecture.

101 **Varying the data distribution - Unigrams slow down learning.** Given the previous findings, one  
102 can ask the question: *is the unigram solution helpful for the eventual convergence of the model, or*  
103 *is it perhaps just a by-product of the learning procedure?* To answer these questions, we define  
104 distributions over Markov chains that are in between the distribution where unigrams is Bayes optimal,  
105 and the distribution where unigrams is as good as uniform. As we see in Figure 3, the transformers  
106 that are being trained on the distribution where there is no unigrams “signal” train much faster. It  
107 appears that this simplicity bias towards the unigrams solution actually slows down learning. See also  
108 Figure 10 in the Appendix that displays how the models perform on different parts of the distribution  
109 during training.

### 110 3.2 Theoretical Insights from the Minimal Model

111 We now provide theoretical insights on how training progresses stage by stage and how this is  
112 achieved by the synergy between the two layers. For this, we analyze the training dynamics of a  
113 minimal model which can be seen as a simplified 2-layer attention only transformer. Section D  
114 contains our main theoretical result. Here, we summarize our theoretical findings:

115 **Learning occurs in two phases.** Both in the theoretical and experimental models, training has  
116 two phases that work at very different speeds. The first phase is fast in both cases; in the theoretical  
117 setting, even a  $O\left(\frac{1}{T}\right)$  step size is sufficient for learning the second layer. In the second phase, a much  
118 larger step size of  $O(1)$  is needed in order to learn the positional encodings.

119 **Second layer is learned first.** It has been observed before in a similar bigram learning setting with  
120 a two-layer transformer that the model might be learning first the second layer [6]. We also make  
121 similar observations in our experiments with the minimal model and the transformers (see Figure  
122 4). For the minimal model, the gradient calculations, clearly suggest that starting from a default  
123 initialization, it is only the second layer that quickly “picks up” the right solution.

124 **Even/odd pattern in positional encodings.** We notice in the experiments that the positional  
125 embeddings of the models displayed an intriguing even/odd oscillating pattern - see Figure 2 (*top,*  
126 *center*), Figure 3 (*right*). We believe that a careful analysis the gradient of  $v$  in the second step will  
127 recover this pattern, which is likely related to the moments of the eigenvalues of the transition matrix.

## 128 4 Conclusion

129 In this work, we have introduced a simple learning problem which serves as a controlled setting  
130 for understanding in-context learning and the emergence of (statistical) induction heads. Through a  
131 combination of empirical investigation and theoretical analysis, we identify different stages during  
132 learning which we were able to precisely characterize. These validate similar observations from  
133 training large-scale language models.

134 It would be worthwhile to understand similar stage-wise learning with natural language data, and use  
135 insights from our minimal model to improve formation of induction heads. In particular, it would be  
136 great to understand if better data curriculum could remove the undesirable simplicity bias we observe  
137 from unigrams. Such simple but incomplete solutions may be commonplace in language modeling  
138 and other rich learning settings; for any such solution, one can ask to what extent its presence speeds  
139 up or slows down the formation of more complex circuits with higher accuracy.

## References

- [1] Jacob D. Abernethy, Alekh Agarwal, Teodor V. Marinov, and Manfred K. Warmuth. A mechanism for sample-efficient in-context learning for sparse retrieval tasks. *CoRR*, abs/2305.17040, 2023.
- [2] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- [3] Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Architectures and algorithms. *CoRR*, abs/2401.12973, 2024. doi: 10.48550/ARXIV.2401.12973. URL <https://doi.org/10.48550/arXiv.2401.12973>.
- [4] Jimmy Ba, Murat A. Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- [5] Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hidden progress in deep learning: Sgd learns parities near the computational limit. *Advances in Neural Information Processing Systems*, 35:21750–21764, 2022.
- [6] Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint, 2023.
- [7] Garrett Birkhoff. Three observations on linear algebra. *Univ. Nac. Tacuman, Rev. Ser. A*, 5: 147–151, 1946.
- [8] Charles Bordenave, Pietro Caputo, and Djilil Chafai. Circular law theorem for random markov matrices. *Probability Theory and Related Fields*, 152, 08 2008.
- [9] Peter F. Brown, Peter V. deSouza, Robert L. Mercer, Vincent J. Della Pietra, and Jenifer C. Lai. Class-based n-gram models of natural language. *Comput. Linguist.*, 18(4):467–479, dec 1992. ISSN 0891-2017.
- [10] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc’ Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [11] Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891, 2022.
- [12] Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L. Leavitt, and Naomi Saphra. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in mlms, 2023.
- [13] Noam Chomsky. Three models for the description of language. *IRE Transactions on information theory*, 2(3):113–124, 1956.
- [14] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- [15] Benjamin L. Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms, 2022.
- [16] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1, 2021.

- 193 [17] Robert Gallager. *Discrete Stochastic Processes (Draft of 2nd Edition)*. MIT OpenCourseWare,  
194 2011.
- 195 [18] Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers  
196 learn in-context? A case study of simple function classes. In Sanmi Koyejo, S. Mohamed,  
197 A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information  
198 Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022,  
199 NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- 200 [19] Jesse Hoogland, George Wang, Matthew Farrugia-Roberts, Liam Carroll, Susan Wei, and Daniel  
201 Murfet. The developmental landscape of in-context learning. *CoRR*, abs/2402.02364, 2024. doi:  
202 10.48550/ARXIV.2402.02364. URL <https://doi.org/10.48550/arXiv.2402.02364>.
- 203 [20] Arthur Jacot, Clément Hongler, and Franck Gabriel. Neural tangent kernel: Conver-  
204 gence and generalization in neural networks. In Samy Bengio, Hanna M. Wallach, Hugo  
205 Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances  
206 in Neural Information Processing Systems 31: Annual Conference on Neural Informa-  
207 tion Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*,  
208 pages 8580–8589, 2018. URL [https://proceedings.neurips.cc/paper/2018/hash/  
209 5a4be1fa34e62bb8a6ec6b91d2462f5a-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/5a4be1fa34e62bb8a6ec6b91d2462f5a-Abstract.html).
- 210 [21] Andrej Karpathy. Mingpt. <https://github.com/karpathy/minGPT/tree/master>, 2023.
- 211 [22] Tanishq Kumar, Blake Bordelon, Samuel J. Gershman, and Cengiz Pehlevan. Grokking  
212 as the transition from lazy to rich training dynamics. *CoRR*, abs/2310.06110, 2023. doi:  
213 10.48550/ARXIV.2310.06110. URL <https://doi.org/10.48550/arXiv.2310.06110>.
- 214 [23] Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Trans-  
215 formers as algorithms: Generalization and stability in in-context learning. In *International  
216 Conference on Machine Learning*, pages 19565–19594. PMLR, 2023.
- 217 [24] Kaifeng Lyu, Jikai Jin, Zhiyuan Li, Simon S. Du, Jason D. Lee, and Wei Hu. Dichotomy of  
218 early and late phase implicit biases can provably induce grokking. *CoRR*, abs/2311.18817,  
219 2023.
- 220 [25] Ashok Vardhan Makuva, Marco Bondaschi, Adway Girish, Alliot Nagle, Martin Jaggi, Hyeji  
221 Kim, and Michael Gastpar. Attention with markov: A framework for principled analysis of  
222 transformers via markov chains. *CoRR*, abs/2402.04161, 2024. doi: 10.48550/ARXIV.2402.  
223 04161.
- 224 [26] William Merrill, Nikolaos Tsilivis, and Aman Shukla. A tale of two circuits: Grokking as  
225 competition of sparse and dense subnetworks. *CoRR*, abs/2303.11873, 2023.
- 226 [27] Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with  
227 gradient descent. *arXiv preprint arXiv:2402.14735*, 2024.
- 228 [28] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom  
229 Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain,  
230 Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson  
231 Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan,  
232 Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer  
233 Circuits Thread*, 2022. [https://transformer-circuits.pub/2022/in-context-learning-and-induction-  
234 heads/index.html](https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html).
- 235 [29] Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gen-  
236 eralization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*,  
237 2022.
- 238 [30] Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context  
239 classification task, 2023.
- 240 [31] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language  
241 models a mirage? *CoRR*, abs/2304.15004, 2023.
- 242 [32] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical  
243 journal*, 27(3):379–423, 1948.
- 244 [33] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position repre-  
245 sentations. In Marilyn A. Walker, Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018  
246 Conference of the North American Chapter of the Association for Computational Linguistics:*

- 247 *Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018,*  
248 *Volume 2 (Short Papers)*, pages 464–468. Association for Computational Linguistics, 2018.
- 249 [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
250 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von  
251 Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman  
252 Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference*  
253 *on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA,*  
254 *pages 5998–6008, 2017.*
- 255 [35] Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L  
256 Bartlett. How many pretraining tasks are needed for in-context learning of linear regression?  
257 *arXiv preprint arXiv:2310.08391, 2023.*
- 258 [36] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of  
259 in-context learning as implicit bayesian inference. In *The Tenth International Conference on*  
260 *Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022, 2022.*

## 261 A Related Work

262 **In-Context Learning.** In [11], the authors discuss how properties of the data distribution promote  
263 ICL. Xie et al. [36] suggest a Bayesian interpretation of ICL and studies how ICL emerges when the  
264 training distribution comes from a Hidden Markov Model (HMM). Abernethy et al. [1] study the  
265 ability of transformers to segment the context into pairs of examples and labels and provide learning  
266 guarantees when the labeling is of the form of a sparse function. Finally, the work of Bietti et al.  
267 [6] studies the dynamics of training transformers on a task that is reminiscent of our Markov chain  
268 setting but has additional complexities. Instead of drawing a fresh Markov chain for each sequence,  
269 in their task all sequences are sampled from the same Markov chain; after certain ‘trigger’ tokens, the  
270 following ‘output’ token is chosen deterministically within a sequence. Thus, successful prediction  
271 requires incorporating both global bigram statistics and in-context deterministic bigram copying,  
272 unlike in our setting where the patterns computed by *statistical* induction heads are necessary and  
273 sufficient. As in our work, the authors identify multiple distinct stages of training and show how  
274 multiple top-down gradient steps lead to a solution.

275 **Induction Heads.** Elhage et al. [16] relates ICL with the formation of induction heads, sub-  
276 components of transformers that match previous occurrences of the current token, retrieving the  
277 token that succeeds the most recent occurrence. Reddy [30] studies the formation of induction heads  
278 and their role in ICL, showing empirically that a three layer network exhibits a sudden formation of  
279 induction heads towards solving some ICL problem of interest. Bietti et al. [6] study the effect of  
280 specific trigger tokens on the formation of induction heads.

281 **Phase Transitions.** It has been observed in different contexts that neural networks and language  
282 models display a sudden drop in loss during their training process. This phase transition is often  
283 related to emergence of new capabilities in the network. The work of Power et al. [29] observed the  
284 “grokking” phenomena, where the test loss of neural networks sharply drops, long after the network  
285 overfits the training data. Chen et al. [12] shows another example of a phase transition in language  
286 model training, where the formation of specific attention mechanisms happen suddenly in training,  
287 causing the loss to quickly drop. Barak et al. [5] observe that neural networks trained on complex  
288 learning problems display a phase transition when converging to the correct solution. Several works  
289 [22, 24] attribute these phase transitions to rapid changes in the inductive bias of networks, while  
290 Merrill et al. [26] argue that the models are sparser after the phase change. Schaeffer et al. [31] warn  
291 that phenomena in deep learning that seem to be discontinuous can actually be understood to evolve  
292 continuously once seen through the right lens.

293 **Concurrent works.** In parallel to this work, there have been a number of papers devoted to the study  
294 of similar questions regarding in-context learning or Markov chains: Akyürek et al. [3] empirically  
295 compare the ability of different architectures to perform in-context learning of regular languages.  
296 Their experiments with synthetic languages motivate architectural changes which improve natural  
297 language modeling in large scale datasets. Hoogland et al. [19] observe similar stage-wise learning  
298 behaviors on transformers trained on language or synthetic linear regression tasks. Makkuva et al.  
299 [25] study the loss landscape of transformers trained on sequences sampled from a single Markov  
300 Chain. Perhaps closest to our work, Nichani et al. [27] introduces a general family of in-context  
301 learning tasks with causal structure, a special case of which is in-context Markov chains. The authors  
302 prove that a simplified transformer architecture (similar to the one we introduce in Section B.2) can  
303 learn to identify the causal relationships by training via gradient descent, and also characterize the  
304 ability of the trained models to adapt to out-of-distribution data. The focus of our work, instead, is on  
305 the different stages of training and how they relate to specific, well-defined, strategies.

## 306 B Setup

307 In this section, we provide further details on our learning problem and present the neural network  
308 architectures that we consider.

309 **Details on ICL-MC Task.** We focus on the case of the *flat* Dirichlet distribution, with  $\alpha =$   
310  $(1, \dots, 1)^\top$ , that corresponds to uniformly random transition probabilities between states. We draw  
311 the initial state  $x_1$  from the stationary distribution  $\pi$  of the chain (which exists almost surely). We  
312 primarily consider the case where the number of states  $k$  is 2 or 3. In subsection E, we consider the  
313 generalization of this setting to  $n$ -grams for  $n > 2$ . Instead of  $\Pr(x_t)$  being determined by  $x_{t-1}$ ,



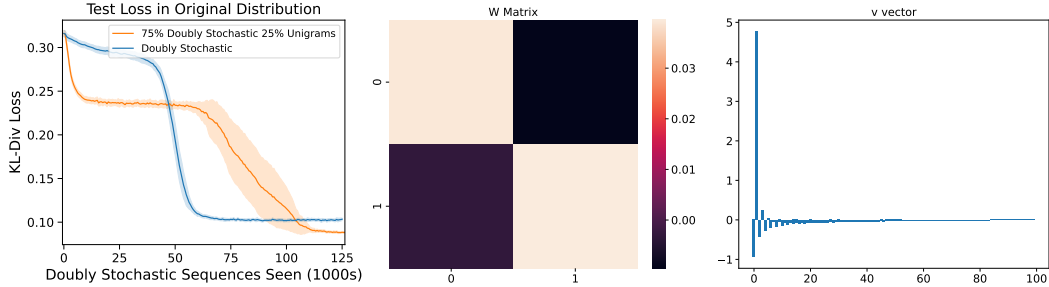


Figure 3: (left) Unigrams slow down optimization: Comparison of two-layer attention only transformers trained on two distributions; one with a uniformly random doubly stochastic transition matrix and another with a mixture of the doubly stochastic and unigrams distribution. We see that in absence of unigrams “signal” the model minimizes the loss (evaluated on the full distribution) much faster. (center, right) Training of the minimal model on ICL-MC with  $k = 2$  states: (center) The heatmap of the second layer ( $W$  matrix) that learns to be close to diagonal. (right) The values of the positional embeddings (1st layer) that display a curious even/odd pattern. This is before any softmax is applied to the positional embeddings.

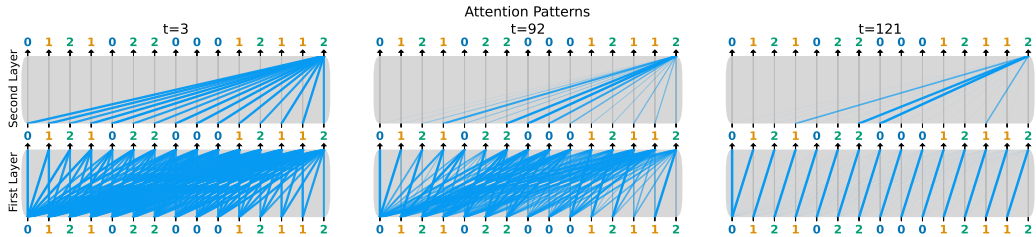


Figure 4: Attention patterns that correspond to the last token of the sequence for a transformer trained to perform ICL-MC. The intensity of each blue line signifies the strength of the corresponding attention value. As the model gets trained, we observe that the attention weights mimic the construction of Proposition B.2. Specifically, at the end of training (right), each token in the first layer is attending to the previous token. In the second layer, the last token, a “2”, is attending to tokens that followed “2”s, allowing bigram statistics to be calculated. See Figure 7 for full attention matrices

314 we let  $\Pr(x_t)$  be determined by  $x_{t-n+1}, \dots, x_{t-1}$ , according to a conditional distribution  $\mathcal{P}$  drawn  
 315 from some prior. In particular, for each tuple of  $n - 1$  tokens, we sample the vector of conditional  
 316 probabilities for the next state from a flat Dirichlet distribution.

### 317 B.1 Potential Strategies for (Partially) Solving ICL-MC

318 **1st strategy: Unigrams.** Since we let the Markov chain reach its stationary distribution (which  
 319 exists a.s.), the optimal strategy across unigrams is just to count frequency of states and form a  
 320 posterior belief about the stationary distribution. Unfortunately, the stationary distribution of this  
 321 random Markov chain does not admit a simple analytical characterization when there is a finite  
 322 number of states, but it can be estimated approximately. At the limit of  $k \rightarrow \infty$ , the stationary  
 323 distribution converges to the uniform distribution [8].

324 **2nd strategy: Bigrams.** For any pair of states  $i$  and  $j$ , let  $\mathcal{P}_{ij}$  be the probability of transition-  
 325 ing from  $i$  to  $j$ . On each sample  $\mathbf{x}$ , we can focus on the transitions from the  $i$ -th state, which  
 326 follow a categorical distribution with probabilities equal to  $(\mathcal{P}_{i1}, \dots, \mathcal{P}_{ik})$ . If we observe the in-  
 327 context empirical counts  $\{c_{ij}\}_{j=1}^k$  of the transitions, then  $\mathcal{P}_{ij}$  is given by:  $(\mathcal{P}_{i1}, \dots, \mathcal{P}_{ik}) | \mathbf{x} \sim$   
 328  $\text{Dir}(k, c_{i1} + \alpha_1, \dots, c_{ik} + \alpha_k)$ , where  $\alpha_1, \dots, \alpha_k$  are the Dirichlet concentration parameters of  
 329 the prior. Hence, each  $\mathcal{P}_{ij}$  has a (marginal) distribution that is actually a Beta distribution:  
 330  $\mathcal{P}_{ij} | \mathbf{x} \sim \text{Beta}(c_{ij} + \alpha_j, \sum_j \alpha_j + N_i - \alpha_j - c_{ij})$ , where  $N_i$  is the total number of observed transi-

331 tions from state  $i$ . As such, our best (point) estimate for each state  $j$  is given by:  $\mathbb{E}[\mathcal{P}_{ij}|\mathbf{x}] = \frac{c_{ij} + \alpha_j}{N + \sum_i \alpha_i}$ .  
 332 For the uniform Dirichlet,  $\alpha = (1, \dots, 1)^\top$ , it is  $\mathbb{E}[\mathcal{P}_{ij}|\mathbf{x}] = \frac{c_{ij} + 1}{N_i + k}$ .

333 *Remark B.1.* The bigram strategy implicitly assumes that the first token  $x_1$  is sampled uniformly,  
 334 as opposed to being sampled from the stationary distribution (which is used in our experiments and  
 335 theoretical results). As the context length grows, the bigram statistics approach the Bayes optimal  
 336 solution either way and this difference becomes negligible.

## 337 B.2 Architectures: Transformers and Simplifications

338 We are mainly interested in investigating how transformers [34] can succeed in in-context learning  
 339 this task. We focus on attention-only transformers with 2 layers with causal masking which is a  
 340 popular architecture for language modeling. Given an input sequence  $\mathbf{x}$ , the output of an  $n$ -layer  
 341 attention-only transformer<sup>1</sup> is:

$$TF(E) = P \circ (\text{Attn}_n + I) \cdots \circ (\text{Attn}_1 + I) \circ E. \quad (1)$$

342 Where  $E \in \mathbb{R}^{t \times d}$  is an embedding of  $\mathbf{x}$ ,  $P \in \mathbb{R}^{d \times k}$  is a linear projection to the output logits, and  
 343  $\text{Attn}(\mathbf{x})$  is masked self attention with relative position embeddings [33], which is parameterized by  
 344  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}, v \in \mathbb{R}^{t \times d}$ :

$$\text{Attn}(z) = \text{softmax}(\text{mask}(A))zW_V, \quad A_{i,j} = \frac{(z_i W_Q)(z_j W_K + v_{i-j+1})^\top}{\sqrt{d}}. \quad (2)$$

345 Transformers with more complicated components, such as MLPs, also display similar qualitative  
 346 behavior (see Figure 8). During training, we minimize this loss:

$$L(\theta) = \mathbb{E}_{\substack{\mathbf{x} \sim \mathcal{P} \\ \mathcal{P} \sim \text{Dir}(\alpha)^k}} \left[ \frac{1}{t} \sum_{p=1}^t l(TF(\mathbf{x}; \theta)_p, x_{p+1}) \right], \quad (3)$$

347 where  $\theta$  denotes the parameters of the model and  $l$  is the cross entropy loss.

348 We now show how a two-layer transformer can represent the optimal bigrams solution.

349 **Proposition B.2** (Transformer Construction). *A single-head two layer attention-only transformer*  
 350 *can find the bigram statistics in the in-context learning Markov chain task.*

351 Intuitively, the first layer of the transformer copies the previous token at each position, and in the  
 352 second layer each token sums the embeddings of all the tokens whose output from the first layer  
 353 matches itself. The full proof can be found in Appendix D.1.

354 **Simplified Transformer Architecture.** As we see from the construction, there are two main  
 355 ingredients in the solution realized by the transformer; (1st layer) the ability to look one token back  
 356 and (2nd layer) the ability to attend to itself. For this reason, we define a *minimal model* that is  
 357 expressive enough to be able to represent such a solution, but also simple enough to be amenable to  
 358 analysis. Let  $e_{x_i}$  denote the one-hot embedding that corresponds to the state at position  $i \in [T]$ , and  
 359 let  $E$  be the  $\mathbb{R}^{(T+1) \times k}$  one-hot embedding matrix. Then the model is parameterized by  $W \in \mathbb{R}^{k \times k}$   
 360 and  $v \in \mathbb{R}^{T+1}$  and defined as:

$$f(E) = \text{mask}(EW(\text{Softmax}(M)E)^\top)E, \quad M = \begin{pmatrix} v_0 & -\infty & \dots & -\infty \\ v_1 & v_0 & \dots & -\infty \\ \vdots & \vdots & \dots & \vdots \\ v_T & v_{T-1} & \dots & v_0 \end{pmatrix} \in \mathbb{R}^{(T+1) \times (T+1)}, \quad (4)$$

361 where  $\text{mask}(\cdot)$  is a causal mask, and  $\text{Softmax}(M)_{i,j} = \frac{\exp(M_{i,j})}{\sum_{\tilde{i}=0}^T \exp(M_{\tilde{i},j})}$ . Notice that the role of  $W$  is  
 362 to mimic the attention mechanism of the second layer and the role of  $v$  is that of the relative positional  
 363 embeddings. This model can be seen as a simplified version of a two-layer linear attention-only  
 364 transformer. See also Appendix D.2 for a discussion.

<sup>1</sup>For simplicity of notation we assume embedding dimension equals the hidden dimension, but in general they can be different.

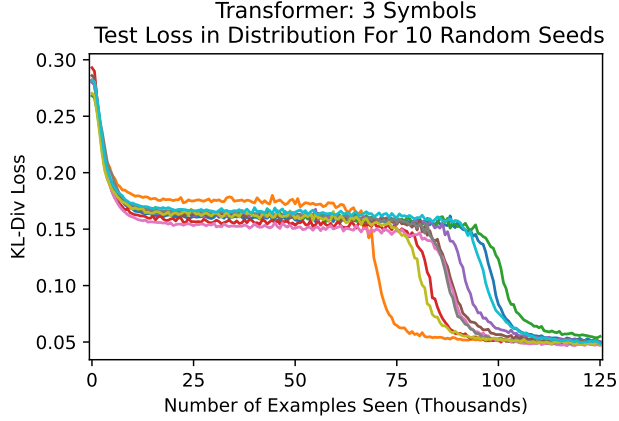


Figure 5: In distribution test loss for 10 two layer attention only transformers, with random seeds  $0, 1, \dots, 9$  (randomness affects initialization and the training data). The training dynamics are consistent for each model, though the exact position of the phase transitions changes.

365 *Fact B.3.* Both the bigrams strategy and the unigrams strategy can be expressed by the minimal  
 366 model with a simple choice of weights.

367 • **Bigrams:**  $v = (0, c, 0, \dots, 0)^\top$  and  $W = I_{k \times k}$ , then  $f(E)_{T,s} =$   
 368  $\sum_{t'=2}^T \mathbb{1}\{x_{t'} = s\} \mathbb{1}\{x_{t'-1} = x_T\} + O\left(\frac{kT^2}{\exp(c)}\right)$ .

369 • **Unigrams:** For  $v = (0, 0, 0, \dots, 0)^\top$ ,  $W = 11^\top$ , we have  $f(E)_{T,s} = \sum_{t'=1}^T \mathbb{1}\{x_{t'} = s\}$ .

## 370 C Experimental Details and Additional Experiments

371 **Note on KL-divergence** In our experiments, we used KL divergence to measure the difference  
 372 between the probabilities predicted by the model and other probability distributions. For test loss,  
 373 this other distribution was the appropriate rows of the transition matrices used to generate the test  
 374 examples.

Formally, let  $f(\mathbf{x}_{1:T-1})$  be the softmax distribution of the transformer’s output, given the input  
 sequence  $\mathbf{x}_{1:T-1}$ . In our standard setting, we measured

$$d_{KL}(\mathcal{P}_{\mathbf{x}_{T-1}} \| f(\mathbf{x}_{1:T-1}))$$

375 where  $\mathcal{P}_{\mathbf{x}_{T-1}}$  is the true distribution of the next state  $\mathbf{x}_T$  given the previous state, under the true  
 376 Markov chain  $\mathcal{P}$ . Note that  $\mathcal{P}$  varies from sequence to sequence (it is drawn from a prior over  
 377 transition matrices) and is not directly observable by the learner—this is what needs to be learned  
 378 in-context.

For measuring how close the model was to various strategies, we computed the predicted probabilities  
 given by said strategies, and used those as the base distribution. Note that the output of the bigrams  
 strategy (which is Bayes-optimal for our base setting) is different from the aforementioned ground-  
 truth  $\mathcal{P}_{\mathbf{x}_{T-1}}$ . Instead, as described in Section B, it is a Bayesian posterior distribution of the next  
 state given the observed sequence, with the prior determined by the prior distribution of transition  
 matrices. Formally:

$$\mathbb{E}[\mathcal{P}_{\mathbf{x}_{T-1}} | \mathbf{x}_{1:T-1}]$$

379 where the expectation is taken over the draw of Markov chain transition matrix.

380 **Experimental details** We train transformers of the form (1) with the AdamW optimizer with  
 381 learning rate  $3e - 5$  (for 3-grams a learning rate of  $3e - 2$  was used), batch size 64, and hidden  
 382 dimension 16. The sequence length of the examples is 100 tokens. The minimal model was trained  
 383 with SGD, with batch size 64, and learning rate  $3e - 4$ . We use PyTorch 2.1.2.

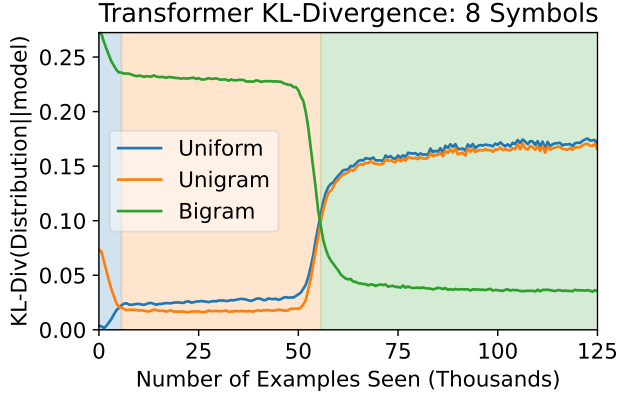


Figure 6: ICL-MC with  $k = 8$  states - KL-divergence between the transformer and the various strategies over training. This required a sequence length greater than 100 (200 in this case) for the difference between unigrams and bigrams to be large enough for the unigram phase to be visible (in either case there was a plateau before the final drop in test loss).

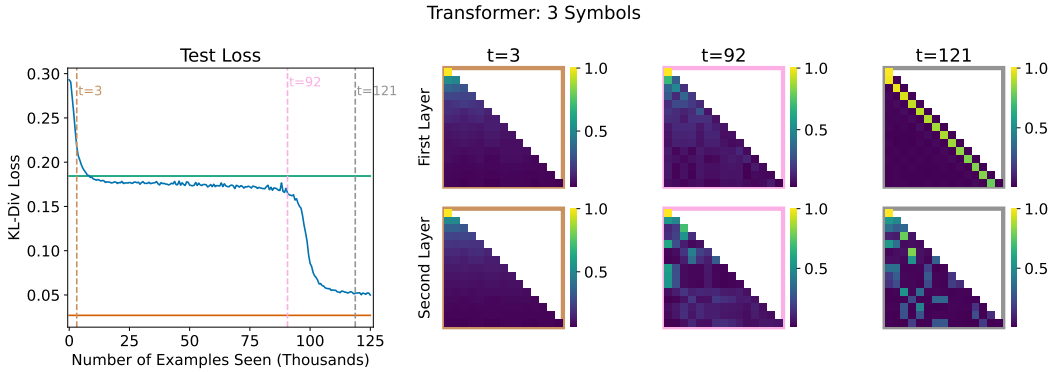


Figure 7: A two layer attention-only transformer trained with cross entropy loss on ICL-MC. The heatmaps on the right represent part of the attention for the transformer at various time steps, specifically the values of the matrix  $A$  from (2). The top row are showing  $A$  from the first layer, and the bottom row from the second layer.

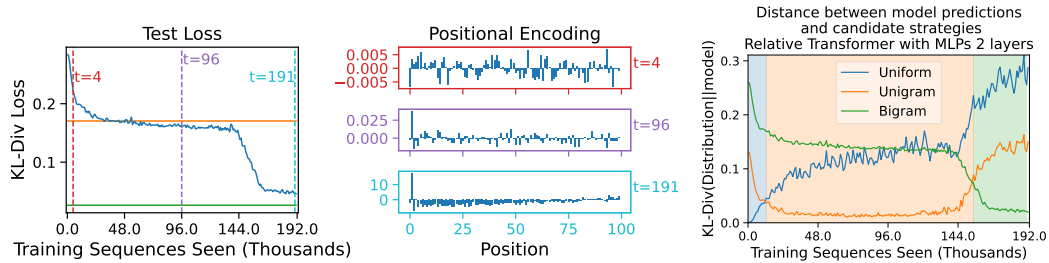


Figure 8: A two layer relative position encoding transformer with MLPs trained on ICL-MC with  $k=3$  symbols. Notice while slightly noisier, the overall trend and observations made regarding the attention only transformer still hold.

384 The data was generated in an online fashion, using `numpy.random.dirichlet` to generate each row  
 385 of the transition matrices. Both the model initialization (for the transformers) and the data were  
 386 randomized based on the seed (in a perfectly reproducible manner).

387 Some of the training and model code was based on minGPT [21]. The experiments all measure the  
 388 outputs of the models at the last token.

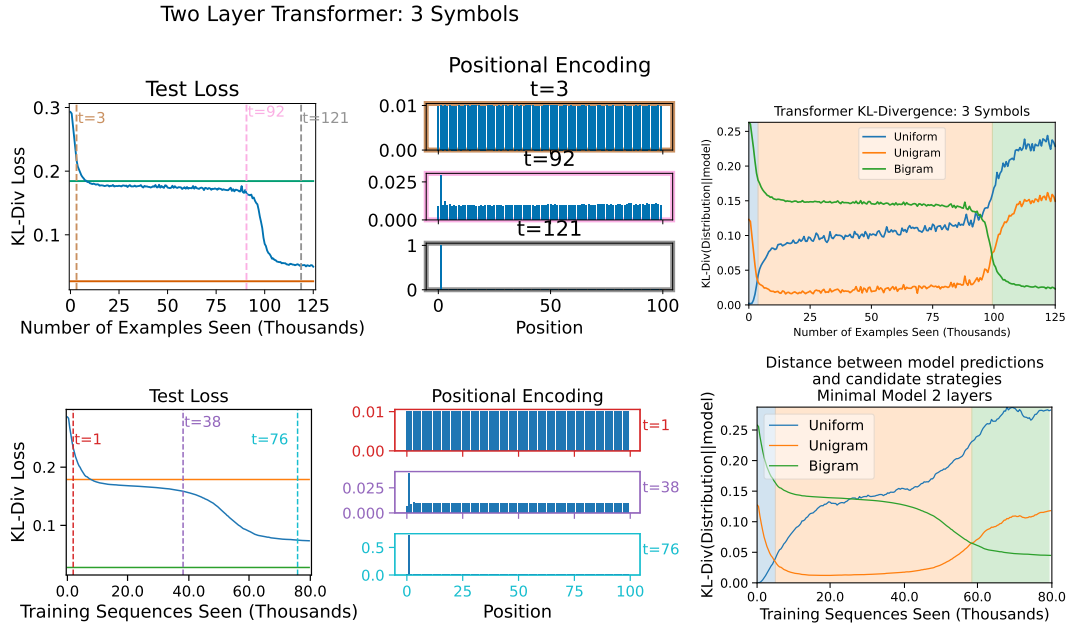


Figure 9: A comparison of the two layer attention only transformer and minimal model for  $k = 3$  symbols.

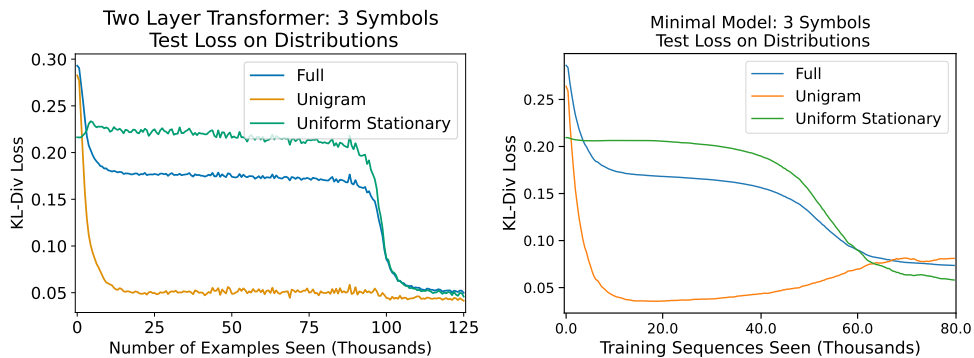


Figure 10: A two layer attention-only transformer (top) and minimal model (4) (bottom), trained on the main task with ICL-MC with cross entropy loss, test loss measured by KL-Divergence from the underlying truth (labels based on transition probabilities, not samples). The distributions test loss is measured in are (from left to right) in-distribution, a distribution where each token is sampled iid, and a distribution over uniformly random doubly stochastic transition matrices (equivalently, stationary distribution is identity, or unigram based guesses are as good as guessing uniform probability). For both models, the in distribution test loss quickly drops to the level of the unigram algorithm.

389 All of the experiments were performed with a single NVIDIA GeForce GTX 1650 Ti GPU with 4  
 390 gigabytes of vram with 32 gigabytes of system memory. Each training run took under ten minutes.

## 391 D Proofs

392 In this section, we present our theoretical results on in-context learning Markov Chains of Section  
 393 B.2.

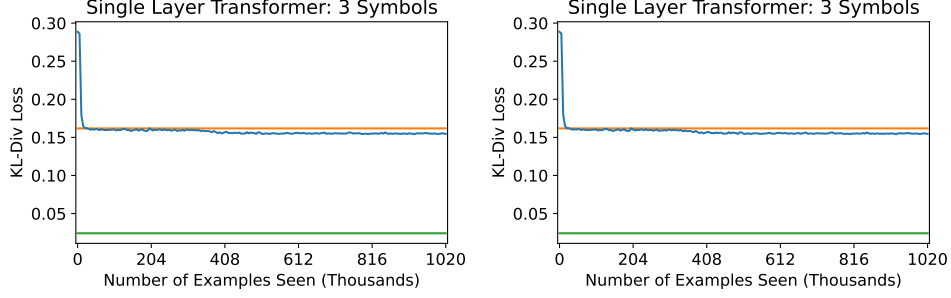


Figure 11: Graphs of test loss showing that a single layer transformer can not achieve good performance on ICL-MC. This result holds for transformers with or without MLPs, and with absolute or relative positional encodings. These graphs show that even trained 8 times longer, there is no notable increase in performance beyond the unigrams strategy (orange line).

### 394 D.1 Transformer Construction

395 *Proof of Proposition B.2.* Set the internal dimension  $d = 3k$ , and choose  $\mathbf{e}_x$  to be one-hot  
 396 embeddings—that is,  $\mathbf{e}_{x_i} = \delta_{x_i}$ , where  $\delta$  is the Kronecker delta. We will call the parameters  
 397 of attention layer  $i$ ,  $W_Q^{(i)}$ ,  $W_K^{(i)}$ ,  $W_V^{(i)}$ ,  $v^{(i)}$ . Let

$$v^{(1)} = \begin{pmatrix} \delta_2 \mathbf{1}_k^\top \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix} \quad W_Q^{(1)} = \begin{pmatrix} cI^{k \times k} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \quad W_K^{(1)} = \mathbf{0} \quad W_V^{(1)} = \begin{pmatrix} \mathbf{0} & I^{k \times k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}$$

So,

$$A_{i,j}^{(1)} = \frac{(e_i W_Q^{(1)})(v_{i-j+1}^{(1)})^\top}{\sqrt{d}}.$$

398 Notice that  $A_{i,j}^{(1)} = c\mathbf{1}[j = i - 1]$ . So,  $\text{softmax}(\text{mask}(A))_{i,j}^{(1)} \approx \mathbf{1}[j = i - 1]$  for large enough  $c$ .  
 399 So, for any  $2 \leq i < T$ ,  $1 \leq j < k$ ,  $\text{Attn}_1(e)_{i,j+2k} = e_{i-1,j}$ . Effectively, the first layer appends  
 400 the embedding of the previous token after the embedding of the current token, so that the output at  
 401 position  $i$  is approximately  $(e_{x_i} \ e_{x_{i-1}} \ \mathbf{0})$ .

402 The second layer is defined as follows:

$$v^{(2)} = \mathbf{0} \quad W_Q^{(2)} = \begin{pmatrix} cI^{k \times k} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \quad W_K^{(2)} = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ I^{k \times k} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \quad W_V^{(2)} = \begin{pmatrix} \mathbf{0} & \mathbf{0} & I^{k \times k} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}$$

Note that  $z = e + \text{Attn}_1(e)$ , then

$$A_{i,j}^{(2)} = \frac{(z_i W_Q^{(2)})(z_j W_K^{(2)})^\top}{\sqrt{d}} = \frac{ce_{x_i}(e_{x_{j-1}})^\top}{\sqrt{d}} = \frac{c}{\sqrt{d}} \mathbf{1}[x_{j-1} = x_i].$$

403 So, for all  $j < i$ ,  $\text{softmax}(\text{mask}(A))_{i,j} \approx \frac{\mathbf{1}[x_{j-1} = x_i]}{\sum_{h=1}^i \mathbf{1}[x_{h-1} = x_i]}$  for large enough  $c$ . For any  $2 \leq i <$   
 404  $T$ ,  $1 \leq j < k$ ,

$$\text{Attn}_2(e)_{i,j+2k} = \sum_{h=1}^{3k} \frac{\mathbf{1}[x_{h-1} = x_i]}{\sum_{g=1}^i \mathbf{1}[x_{g-1} = x_i]} (z W_V^{(2)})_{h,j} = \frac{\sum_{h=1}^k \mathbf{1}[x_{h-1} = x_i] \mathbf{1}[x_h = j]}{\sum_{g=1}^i \mathbf{1}[x_{g-1} = x_i]}.$$

405 Which is exactly the empirical bigram statistics (that is, the number of times  $x_i \rightarrow j$  appears before  
 406 position  $i$ ), so to make this the output,  $P = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ aI^{k \times k} \end{pmatrix}^2$   $\square$

## 407 D.2 ICL-MC with Minimal Model

408 To abstract away some of the many complicated components from the transformer architecture, we  
 409 focus our attention now to the minimal model of Section B.2. We train minimal models of eq. (4),  
 410 starting from a deterministic constant initialization, by minimizing the cross entropy loss with SGD.  
 411 Full experimental details can be found in the Appendix. Figure 2 (bottom) displays the training  
 412 curves for the minimal model.

413 **Lemma D.1.** *Let the model defined as in eq. (4) and initialized with  $W = \mathbf{0}, v = \mathbf{0}$ . If the random  
 414 transition matrices are either*

- 415 • *uniformly random from  $2 \times 2$  stochastic matrices*
- 416 • *With some constant probability  $0 < \alpha < 1$ , a uniformly random doubly stochastic matrices, and  
 417 otherwise  $\mathbf{1}^\top v$  where  $v$  is a uniformly random vector on the  $k$ -simplex.*

418 *after one step of full batch gradient descent with step size  $\eta$  we have:*

$$W = \eta(T+1)(AI + B\mathbf{1}^\top\mathbf{1}) + \eta O(\log T) \text{ and } v^{(1)} = \mathbf{0},$$

419 *where  $A, B \in \mathbb{R}^+$ .*

*Assuming in the first step  $\eta = O(\frac{1}{T^2})$ , after the second step of gradient descent, it holds:*

$$W = (\eta + \eta_W)(T+1)(AI + B\mathbf{1}^\top\mathbf{1}) + (\eta + \eta_W) O(\log T)$$

420 *where the step size on  $W$  in the second step is  $\eta_W$ . Furthermore,*

$$v_1 = \eta_v C \log T, \text{ and } v_1 - v_n = \eta_v \Omega(\log T) \forall n \neq 1,$$

421 *where  $\eta_v$  is the step size for  $v$  in the second step, and  $C \approx 0.0114$ .*

422 *If  $\eta_v = O(T)$ , and  $\eta_W = \frac{1}{T(A+B)}$  then the output of the model will be a weighted sum of bigrams  
 423 and unigrams. Formally,*

$$f(E)_{T,s} = \frac{A}{A+B} \sum_{i=1}^T \mathbb{1}[x_{i-1} = x_t, x_i = s] + \frac{B}{A+B} \sum_{i=1}^T \mathbb{1}[x_i = s] + O(\log T)$$

424 *Note that in the first distribution (uniformly random  $2 \times 2$ ) or the second distribution with  $k > 6$ ,  
 425  $A > B$ , so at the end of the two steps, the weight on bigrams is greater than that of the weight on  
 426 unigrams.*

427 *Proof Overview.* The idea of the proof is that a first step of gradient descent with a small learning rate  
 428 can align the second layer, while a second step can learn to identify the correct relative positional  
 429 embedding. The identity bias of  $W$  in the second layer ensures there is a strong signal in the gradient  
 430 to look back one in the first layer. Without a bias in  $W$ , the gradient for the positional encodings,  $v$ ,  
 431 turns out to be zero.

432 We get additional intuition from looking at the proof for just the second distribution: in the first step,  
 433 effectively all of the gradient comes from the examples where the unigram strategy is optimal, while  
 434 in the second step effectively all of the gradient comes from the examples where the bigram strategy  
 435 is optimal.

436 *Remark D.2.* It is worth noting that, while this is a simplified setting, the analysis goes beyond  
 437 NTK-based [20] analyses where the representations do not change much and it crucially involves  
 438 more than one step which has been a standard tool in the analysis of feature learning [4].

<sup>2</sup>Technically, the output of this construction is not the log probabilities as generally cross-entropy loss  
 assumes. These can be approximated linearly by setting  $P = \begin{pmatrix} b\mathbf{1}^\top\mathbf{1} \\ \mathbf{0} \\ aI^{k \times k} \end{pmatrix}$  to change the output from  $x$  to  $ax + b$ .  
 In practice, this approximation can achieve close to Bayes optimal loss.

439 **Setup and notation** Our data consists of sequences of length  $T + 1$ ,  $\mathbf{x} = (x_0, \dots, x_T)$ , drawn  
440 from a Markov Chain with state space  $S = \{1, \dots, k\}$  (i.e.,  $x_j \in \{1, \dots, k\}$  for all  $j \in [T]$ ), and  
441 a random transition matrix  $P$ . Each row of the matrix is sampled from a flat Dirichlet distribution,  
442 i.e.  $P_i \sim \text{Dir}(\mathbf{1})$ , corresponding drawing the row from a uniform distribution over the simplex. Let  
443  $E \in \{0, 1\}^{(T+1) \times k}$  be the one hot embedding matrix of  $x$ , that is,  $E_{i,x_i} = 1$  and for all  $s \neq x_i$   
444  $E_{i,s} = 0$ .

**Model** We define our model as a simplified sequence to sequence transformer  $f : \mathbb{R}^{T \times k} \rightarrow \mathbb{R}^{(T+1) \times k}$  with  $f(E) = \text{mask}(EW(\text{Softmax}(M)E)^\top)E$ . The trained parameters are  $W \in \mathbb{R}^{k \times k}$

and  $v \in \mathbb{R}^{T+1}$ . We define  $M \in \mathbb{R}^{(T+1) \times (T+1)}$  as  $M = \begin{pmatrix} v_0 & -\infty & \dots & -\infty \\ v_1 & v_0 & \dots & -\infty \\ \vdots & \vdots & \dots & \vdots \\ v_T & v_{T-1} & \dots & v_0 \end{pmatrix}$ , that is, for all

$T \geq i \geq j \geq 0$ ,  $M_{i,j} = v_{i-j}$  and if  $i > j$ ,  $M_{j,i} = -\infty$ . Furthermore,  $v = [v_0, v_2, \dots, v_T] \in \mathbb{R}^{T+1}$ . Softmax is defined as follows:

$$\text{Softmax}(M)_{i,j} = \frac{\exp(M_{i,j})}{\sum_{T=1}^T \exp(M_{i,j})}$$

445 The logit for symbol  $s$  at position  $T$  for our model is:

$$f(E)_{T,s} = \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i W_{x_T,u} \mathbf{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)}. \quad (5)$$

446 The model can represent the unigrams and bigrams solutions as following:

- 447 • Construction for bigrams:  $v = (0, c, 0, \dots, 0)^\top$  and  $W = I_{k \times k}$ , then  $f(E)_{T,s} =$   
448  $\sum_{i=0}^T \mathbf{1}[x_i = s \wedge x_{i-1} = x_T] + O\left(\frac{T^3}{\exp(c)}\right)$ . As  $c$  tends to infinity, this becomes bigrams.
- 449 • Construction for unigrams:  $v = \mathbf{0}$  and  $W = \mathbf{1}^\top \mathbf{1}$ , then  $f(E)_{T,s} = \sum_{i=0}^T \mathbf{1}[x_i = s]$ .

### 450 Proof of Fact B.3

451 *Proof.* We will first prove the unigrams construction.

$$\begin{aligned} f(E)_{T,s} &= \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i W_{x_T,u} \mathbf{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \\ &= \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i \mathbf{1}[x_{i-j} = u \wedge x_i = s] \frac{1}{i} \\ &= \sum_{i=0}^T \sum_{j=0}^i \mathbf{1}[x_i = s] \frac{k}{i} \\ &= k \sum_{i=0}^T \mathbf{1}[x_i = s] \end{aligned}$$

452 Which is exactly unigrams.

453 Now consider the bigrams construction. As  $c$  grows, the softmax of  $v$  very quickly becomes one hot.  
454 Formally, by lemma B.7 in [15], for any  $i > 0$ ,

$$\frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} = \mathbf{1}[j = 1] + O\left(\frac{T}{\exp(c)}\right)$$



455 So,

$$\begin{aligned} f(E)_{T,s} &= \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i W_{x_T,u} \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \\ &= \sum_{i=0}^T \sum_{j=0}^i \mathbb{1}[x_{i-j} = x_T \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \end{aligned}$$

456 We will take out the term  $i = 0$ ,

$$= \mathbb{1}[x_T = x_0 = s] + \sum_{i=1}^T \sum_{j=0}^i \mathbb{1}[x_T = x_{i-j} \wedge x_i = s] \frac{\exp(v_j)}{i - 1 + \exp(c)}$$

457 Then apply the softmax approximation mentioned earlier,

$$\begin{aligned} &= \mathbb{1}[x_T = x_0 = s] + \sum_{i=1}^T \sum_{j=0}^i \mathbb{1}[x_T = x_{i-j} \wedge x_i = s] \left( \mathbb{1}[j = 1] + O\left(\frac{2T}{\exp(c)}\right) \right) \\ &= \mathbb{1}[x_T = x_0 = s] + \sum_{i=1}^T \mathbb{1}[x_T = x_{i-1} \wedge x_i = s] + \sum_{i=1}^T \sum_{j=0}^i \mathbb{1}[x_T = x_{i-j} \wedge x_i = s] O\left(\frac{2T}{\exp(c)}\right) \\ &= \sum_{i=1}^T \mathbb{1}[x_T = x_{i-1} \wedge x_i = s] + \sum_{i=1}^T O\left(\frac{T^3}{\exp(c)}\right) \end{aligned}$$

458

□

459 This simplified model was constructed by taking a two layer transformer with relative positional  
 460 encodings and simplifying it. Our construction for how transformers would form induction heads  
 461 (corroborated with experiments such as the viewing of attention patterns in figure 4) implies that the  
 462 MLPs and the value matrices could just be identity functions, and the first layer query matrix, and the  
 463 second layer positional embeddings were zero matrices, so in the simplified model we froze these  
 464 parameters to their final states. We also remove the softmax on the attention in the first layer. Despite  
 465 these changes, the training dynamics, our main interest, stay remarkably similar.

466 **Training** We analyze gradient descent with the cross entropy loss  $L_T(f, E, x_{T+1}) =$   
 467  $-\sum_{s=1}^k \log \text{Softmax}(f(E))_{T,s} P_{X_T,s}^3$

### 468 D.3 Gradient Calculations

469 For use in the proofs, here we show the calculations of the gradients of the model with respect to the  
 470 parameters, and the loss with respect to the model.

$$\begin{aligned} \frac{\partial f(E)_{T,s}}{\partial W_{a,b}} &= \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i \mathbb{1}[x_T = a \wedge b = u] \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \\ &= \sum_{i=0}^T \sum_{j=0}^i \mathbb{1}[x_T = a] \mathbb{1}[x_{i-j} = b \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \end{aligned}$$

471

$$\begin{aligned} &\frac{\partial f(E)_{T,s}}{\partial v_a} \\ &= \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i W_{x_T,u} \mathbb{1}[x_{i-j} = u \wedge x_i = s] \left( \mathbb{1}[j = a] \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} - \mathbb{1}[a \leq i] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} \right) \end{aligned}$$

<sup>3</sup>In practice, one would often use the empirical value of  $x_{T+1}$  rather than its distribution  $P_{X_T,s}$ , but in full batch gradient descent this is in fact equivalent in our setting. This is because conditional on  $x_T$  and  $P$ ,  $x_{T+1}$  is independent of  $x_1, \dots, x_{T-1}$ .

$$\begin{aligned}
&= \sum_{u=1}^k \sum_{i=0}^T \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} \sum_{j=0}^i W_{x_T, u} \left( \mathbb{1}[x_{i-a} = u \wedge x_i = s] \mathbb{1}[j = a] - \mathbb{1}[x_{i-j} = u \wedge x_i = s] \mathbb{1}[a \leq i] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \right) \\
&= \sum_{u=1}^k \sum_{i=0}^T \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} W_{x_T, u} \left( \mathbb{1}[x_{i-a} = u \wedge x_i = s] \mathbb{1}[a \leq i] - \sum_{j=0}^i \mathbb{1}[x_{i-j} = u \wedge x_i = s] \mathbb{1}[a \leq i] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \right) \\
&= \sum_{u=1}^k \sum_{i=a}^T \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} W_{x_T, u} \left( \mathbb{1}[x_{i-a} = u \wedge x_i = s] - \sum_{j=0}^i \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \right) \\
&= \sum_{u=1}^k \sum_{i=a}^T \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} W_{x_T, u} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \mathbb{1}[x_{i-j} = u] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \right)
\end{aligned}$$

472

$$\frac{\partial L_T}{\partial f(E)_{T,s}} = \text{Softmax}(f(E))_{T,s} - P_{x_T, s}$$

#### 473 **D.4 Proof of lemma D.1**

474 *Proof.* Recall that at initialization,  $v = \mathbf{0}$  and  $W = \mathbf{0}$ , implying further that  $f(E) = \mathbf{0}$ .

475 **First step.**

476 First consider the gradient of the loss with respect to  $W$ . By chain rule,

$$\begin{aligned}
\frac{\partial L_T(E)}{\partial W_{a,b}} &= \sum_{s=1}^k \frac{\partial L_T}{\partial f(E)_{T,s}} \frac{\partial f(E)_{T,s}}{\partial W_{a,b}} \\
&= \sum_{s=1}^k (\text{Softmax}(f(E))_{T,s} - P_{x_T, s}) \sum_{i=0}^T \sum_{j=0}^i \mathbb{1}[x_T = a] \mathbb{1}[x_{i-j} = b \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \\
&= \sum_{s=1}^k \left( \frac{1}{k} - P_{x_T, s} \right) \mathbb{1}[x_T = a] \sum_{i=0}^T \sum_{j=0}^i \mathbb{1}[x_{i-j} = b \wedge x_i = s] \frac{1}{i+1} \\
&= \frac{1}{k} \mathbb{1}[x_T = a] \sum_{i=0}^T \left( \sum_{j=0}^i \mathbb{1}[x_{i-j} = b] \frac{1}{i+1} - \sum_{s=1}^k P_{a,s} \mathbb{1}[x_T = a] \sum_{j=0}^i \mathbb{1}[x_{i-j} = b \wedge x_i = s] \frac{1}{i+1} \right) \\
&= \frac{1}{k} \mathbb{1}[x_T = a] \sum_{i=0}^T \left( \mathbb{1}[x_i = b] - \sum_{s=1}^k P_{a,s} \mathbb{1}[x_T = a] \sum_{j=0}^i \mathbb{1}[x_{i-j} = b \wedge x_i = s] \frac{1}{i+1} \right) \\
&= \frac{1}{k} \sum_{i=0}^T \left( \mathbb{1}[x_i = b \wedge x_T = a] - \sum_{s=1}^k P_{a,s} \sum_{j=0}^i \mathbb{1}[x_{i-j} = b \wedge x_i = s \wedge x_T = a] \frac{1}{i+1} \right) \\
&= \frac{1}{k} \sum_{i=0}^T \left( \mathbb{1}[x_0 = b \wedge x_{T-i} = a] - \sum_{s=1}^k P_{a,s} \sum_{j=0}^i \mathbb{1}[x_0 = b \wedge x_j = s \wedge x_{T-i+j} = a] \frac{1}{i+1} \right)
\end{aligned}$$

477 Where the last line follows from the markov property.

478 Now we take the expectation over  $x$ ,  $x_{T+1}$  conditioned on the transition matrix  $P$ ,

$$\begin{aligned}
\mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial W_{a,b}} \right] &= \pi_b \sum_{i=0}^T \left( \frac{1}{k} (P^{T-i})_{b,a} - \sum_{s=1}^k P_{a,s} \frac{1}{i+1} (P^{T-i})_{s,a} \sum_{j=0}^i (P^j)_{b,s} \right) \\
&= \pi_b \sum_{i=0}^T \left( \frac{1}{k} (P^i)_{b,a} - \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} (P^i)_{s,a} \sum_{j=0}^{T-i} (P^j)_{b,s} \right)
\end{aligned}$$

$$= \pi_b \pi_a (T+1) \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) + O(\log T)$$

479 Where the last step follows from Lemma D.8. Then, by applying Lemma D.3 or lemmas D.5 and D.6  
 480 (depending on the distribution assumption on  $P$ ), there exist positive constants (potentially depending  
 481 on  $k$ , but not  $T$ )  $A, B$  such that for all  $a$

$$\mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial W_{a,a}} \right] = -(A+B)T + O(\log T)$$

482 and for all  $a \neq b$ ,

$$\mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial W_{a,b}} \right] = -BT + O(\log T)$$

483 The updated  $W_{a,b}$  after the gradient step is just  $-\eta \mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial W_{a,b}} \right]$  (because  $W$  is initialized at  $\mathbf{0}$ ).

484 Choose  $\eta = \Theta\left(\frac{1}{T}\right)$ , so that  $W$  will be  $O(1)$  with respect to  $T$  after the first step.

485 For the gradient with respect to  $v$ , since  $W = \mathbf{0}$ ,

$$\begin{aligned} \frac{\partial F(E)_{T,s}}{\partial v} &= \sum_{u=1}^k \sum_{i=a}^T \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} W_{x_T, u} \left( \mathbb{1}[x_{i-a} = u \wedge x_i = s] - \sum_{j=0}^i \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \right) \\ &= 0 \end{aligned}$$

486 So,

$$\frac{\partial L_T(E)}{\partial v} = \sum_{s=1}^k \frac{\partial L_T(E)}{\partial f(E)_{T,s}} \frac{\partial F(E)_{T,s}}{\partial v} = 0$$

487 Completing the first step calculations.

488 **Second step.**

489 After the first step,  $W = \eta (AI + B\mathbf{1}^\top \mathbf{1})$ . Now let us bound the output of the model,

$$\begin{aligned} |f(E)_{T,s}| &= \left| \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i W_{x_T, u} \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \right| \\ &= \left| \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i \eta (AI + B\mathbf{1}^\top \mathbf{1})_{x_T, u} \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{1}{i} \right| \\ &\leq \eta \left| \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i (A+B) \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{1}{i} \right| \\ &\leq \eta \left| \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i (A+B) \mathbb{1}[x_{i-j} = u] \frac{1}{i} \right| \\ &\leq \eta \left| \sum_{i=0}^T \sum_{j=0}^i (A+B) \frac{1}{i} \right| \\ &\leq \eta T |A+B| \end{aligned}$$

490 So, using the first order approximation of softmax,

$$\frac{\partial L_T(E)}{\partial f(E)_{T,s}} = \text{Softmax}(f(E))_{T,s} - \mathbb{1}[x_{T+1} = s]$$

$$\begin{aligned}
&= \frac{1}{k} + \frac{f(E)_{T,s}}{k} - \frac{\sum_{u=1}^k f(E)_{T,u}}{k^2} + O(f(E)_{T,s}^2) - \mathbb{1}[x_{T+1} = s] \\
&= \frac{1}{k} + O\left(\eta \frac{T}{k}(A+B)\right) + O(\eta^2 T^2(A+B)^2) - \mathbb{1}[x_{T+1} = s] \\
&= \frac{1}{k} + O\left(\eta \frac{T}{k}(A+B)\right) + O(\eta^2 T^2(A+B)^2) - \mathbb{1}[x_{T+1} = s] \\
&= \frac{1}{k} - \mathbb{1}[x_{T+1} = s] + O\left(\frac{1}{T}\right)
\end{aligned}$$

491 Where the last step follows since  $\eta = O\left(\frac{1}{T^2}\right)$ .

492 Now we can begin to analyze the gradients with respect to the parameters. For  $W$ , the gradient is  
493 approximately the same as in the last step. Notice that  $\frac{\partial f(E)_{T,s}}{\partial W_{a,b}}$  does not depend on  $W$ , and  $v$  is  
494 unchanged, so  $\frac{\partial f(E)_{T,s}}{\partial W_{a,b}}$  is unchanged. Furthermore,

$$\begin{aligned}
\frac{\partial f(E)_{T,s}}{\partial W_{a,b}} &= \sum_{s=1}^k \sum_{i=0}^T \sum_{j=0}^i \mathbb{1}[x_T = a] \mathbb{1}[x_{i-j} = b \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \\
&= \sum_{i=0}^T \sum_{j=0}^i \mathbb{1}[x_T = a] \mathbb{1}[x_{i-j} = b] \frac{1}{i} \\
&\leq \sum_{i=0}^T \sum_{j=0}^i \frac{1}{i} \\
&= T
\end{aligned}$$

495 We will now show that the gradient is approximately the same as in the first gradient step:

$$\begin{aligned}
\frac{\partial L_T(E)}{\partial W_{a,b}} &= \sum_{s=1}^k \frac{\partial L_T}{\partial f(E)_{T,s}} \frac{\partial f(E)_{T,s}}{\partial W_{a,b}} \\
&= \sum_{s=1}^k \left( \frac{1}{k} - \mathbb{1}[x_{T+1} = s] + O\left(\frac{1}{T}\right) \right) \frac{\partial f(E)_{T,s}}{\partial W_{a,b}} \\
&= \sum_{s=1}^k \left( \frac{1}{k} - \mathbb{1}[x_{T+1} = s] \right) \frac{\partial f(E)_{T,s}}{\partial W_{a,b}} + O\left(\frac{1}{T}\right) \frac{\partial f(E)_{T,s}}{\partial W_{a,b}} \\
&= \pi_b \pi_a (T+1) \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) + O(\log T)
\end{aligned}$$

496 Where the last lines follows from the gradient calculations in the first step.

497 Now we will consider the gradient with respect to  $v$ . First, notice that the uniform component of  $W$ ,  
498  $B\mathbf{1}^\top \mathbf{1}$ , has no affect on the gradient of  $v$ :

$$\begin{aligned}
\frac{\partial f(E)_{T,s}}{\partial v_a} &= \sum_{u=1}^k \sum_{i=a}^T W_{x_T, u} \frac{\exp(v_a)}{\sum_{\ell=0}^i \exp(v_\ell)} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \mathbb{1}[x_{i-j} = u] \right) \\
&= \sum_{u=1}^k \sum_{i=a}^T (mI + B\mathbf{1}^\top \mathbf{1})_{x_T, u} \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) \\
&= \sum_{u=1}^k \sum_{i=a}^T (A\mathbb{1}[x_T = u] + B) \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right)
\end{aligned}$$

$$\begin{aligned}
&= A \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) \\
&+ B \sum_{u=1}^k \sum_{i=a}^T \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) \\
&= A \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) \\
&+ B \sum_{i=a}^T \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \sum_{u=1}^k \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \sum_{u=1}^k \mathbb{1}[x_{i-j} = u] \right) \\
&= A \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) \\
&+ B \sum_{i=a}^T \frac{1}{i+1} \mathbb{1}[x_i = s] \left( 1 - \sum_{j=0}^i \frac{1}{i+1} \right) \\
&= A \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right)
\end{aligned}$$

499 By chain rule,

$$\begin{aligned}
\frac{\partial L_T}{\partial v_a} &= \sum_{s=1}^k \frac{\partial L_T}{\partial f(E)_{T,s}} \frac{\partial f(E)_{T,s}}{\partial v_a} \\
&= \sum_{s=1}^k \left( \frac{1}{k} - P_{x_T,s} + O\left(\frac{1}{T}\right) \right) \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) \\
&= \sum_{s=1}^k \left( \frac{1}{k} - P_{x_T,s} \right) \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) + O\left(\frac{\log T}{T}\right)
\end{aligned}$$

500 Where the last step follows because

$$\begin{aligned}
\left| \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \mathbb{1}[x_i = s] \left( \mathbb{1}[x_{i-a} = u] - \sum_{j=0}^i \frac{1}{i+1} \mathbb{1}[x_{i-j} = u] \right) \right| &\leq \left| \sum_{u=1}^k \sum_{i=a}^T \mathbb{1}[x_T = u] \frac{1}{i+1} \right| \\
&= \left| \sum_{i=a}^T \frac{1}{i+1} \right| \\
&\leq \log T
\end{aligned}$$

501 In expectation over the values of  $x$ , conditioned on the choice of  $P$ :

$$\begin{aligned}
\mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial v_a} \right] &= \sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \sum_{i=a}^T \frac{\pi_u}{i+1} (P^{T-i})_{s,u} \left( (P^a)_{u,s} - \frac{1}{i+1} \sum_{j=0}^i (P^{i-j})_{u,s} \right) + O\left(\frac{\log T}{T}\right) \\
&= \sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \sum_{i=a}^T \frac{\pi_u}{T-i+1} (P^i)_{s,u} \left( (P^a)_{u,s} - \frac{1}{T-i+1} \sum_{j=0}^{T-i} (P^j)_{u,s} \right) + O\left(\frac{\log T}{T}\right) \\
&= (\log(T+1) - \log(a+1)) \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} \left( \pi_s - (P^a)_{u,s} \right) + O(1)
\end{aligned}$$

502 Where the last step follows from lemma D.9. Then, by applying Lemma D.4 or lemmas D.5 and D.6  
 503 (depending on the distribution assumption on  $P$ ),

$$\begin{aligned}\mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial v_1} \right] &< \mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial v_a} \right] \\ \mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial v_1} \right] &< 0\end{aligned}$$

504 Therefore, after the step is taken,

$$\begin{aligned}v_1 &= \Theta(\eta_v \log T) \\ v_1 - v_n &= \eta_v \Omega(\log T)\end{aligned}$$

505 Finally, we can consider the state of the model after the second step. Assume that the step size for  $v$   
 506 in the second step is  $O(T)$ , and the step size for  $W$  is  $\frac{1}{T(A+B)}$

$$\begin{aligned}f(E)_{T,s} &= \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i W_{x_T,u} \mathbb{1}[x_{i-j} = u \wedge x_i = s] \frac{\exp(v_j)}{\sum_{\ell=0}^i \exp(v_\ell)} \\ &= \frac{1}{A+B} \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i \left( AI + B\mathbf{1}^\top \mathbf{1} + O\left(\frac{\log T}{T}\right) \right)_{x_T,u} \mathbb{1}[x_{i-j} = u \wedge x_i = s] \left( \mathbb{1}[j=1] + O\left(\frac{2T}{\exp(\log(T))}\right) \right) \\ &= \frac{1}{A+B} \sum_{u=1}^k \sum_{i=0}^T \sum_{j=0}^i \left( AI + B\mathbf{1}^\top \mathbf{1} + O\left(\frac{\log T}{T}\right) \right)_{x_T,u} \mathbb{1}[x_{i-j} = u \wedge x_i = s] \left( \mathbb{1}[j=1] + O\left(\frac{1}{T}\right) \right) \\ &= \frac{1}{A+B} \sum_{u=1}^k \sum_{i=0}^T (AI + B\mathbf{1}^\top \mathbf{1})_{x_T,u} \mathbb{1}[x_{i-1} = u \wedge x_i = s] + O(\log T) \\ &= \frac{A}{A+B} \sum_{i=0}^T \mathbb{1}[x_{i-1} = x_T \wedge x_i = s] + \frac{B}{A+B} \sum_{u=1}^k \sum_{i=0}^T \mathbb{1}[x_{i-1} = u \wedge x_i = s] + O(\log T) \\ &= \frac{A}{A+B} \sum_{i=0}^T \mathbb{1}[x_{i-1} = x_T \wedge x_i = s] + \frac{B}{A+B} \sum_{i=0}^T \mathbb{1}[x_i = s] + O(\log T)\end{aligned}$$

507 This completes the proof. □

### 508 D.5 Inequality lemmas for $k = 2$

**Lemma D.3.** *If  $P$  is a uniformly random stochastic  $2 \times 2$  matrix, and  $\pi$  is the stationary distribution of  $P$ , then*

$$\mathbb{E} \left[ \pi_a^2 \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] = \frac{5}{12} - \frac{2}{3} \log(2) \approx -0.045$$

and for any  $b \neq a$

$$\mathbb{E} \left[ \pi_a \pi_b \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] = -\frac{7}{6} + \frac{5}{3} \log(2) \approx -0.011$$

509 *Proof.* We have:

$$\begin{aligned}\mathbb{E} \left[ \pi_a^2 \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] &= \mathbb{E}_{a,b} \left[ \frac{(b-1)^2}{(a+b-2)^2} \left[ \frac{1}{2} - \frac{a(b-1)}{a+b-2} - \frac{(1-a)(a-1)}{a+b-2} \right] \right] \\ &= \frac{1}{2} \int_0^1 \int_0^1 \frac{(b-1)^2}{(a+b-2)^2} dadb - \int_0^1 \int_0^1 \frac{a(b-1)^3}{(a+b-2)^3} dadb + \int_0^1 \int_0^1 \frac{(b-1)^2(a-1)^2}{(a+b-2)^3} dadb \\ &= \frac{1}{2} (1 - \ln 2) - \frac{1}{2} (1 - \ln 2) + \frac{5}{12} (5 - 8 \ln 2) = \frac{5}{12} - \frac{2}{3} \ln 2.\end{aligned}$$

(6)

510 For the non-diagonal elements, it holds:

$$\begin{aligned}
& \mathbb{E} \left[ \pi_a \pi_b \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] \\
&= \frac{1}{2} \int_0^1 \int_0^1 \frac{(b-1)(a-1)}{(a+b-2)^2} da db - \int_0^1 \int_0^1 \frac{a(b-1)^2(a-1)}{(a+b-2)^3} da db + \int_0^1 \int_0^1 \frac{(b-1)(a-1)^3}{(a+b-2)^3} da db \\
&= \frac{1}{2} \left( \ln 2 - \frac{1}{2} \right) - \frac{1}{6} (1 - \ln 2) + \left( \ln 2 - \frac{3}{4} \right) = \frac{5}{3} \ln 2 - \frac{7}{6}.
\end{aligned} \tag{7}$$

511

□

**Lemma D.4.** *If  $P$  is a uniformly random stochastic  $2 \times 2$  matrix, and  $\pi$  is the stationary distribution of  $P$ , then,*

$$\mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - P_{u,s}) \right] = -7/2 + 5 \log(2) \approx -0.034$$

and for any  $n \neq 1$

$$\mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - P_{u,s}) \right] \leq \mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - (P^n)_{u,s}) \right]$$

512 *Proof.* We have:

$$\begin{aligned}
& \mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - P_{u,s}) \right] = \\
& \left[ \frac{1/6x^4 + (x+y)(6xy(-4x^2+2x+1) + 6y^4 + y^3(20-24x))}{12(x+y)} \right. \\
& \left. + \frac{y^2(12x^2 - 12x - 3) + \log((x+y)^{6x^2(4x^2+2x-1)}(x+y)^{6y^2(4y^2+2y-1)})}{12(x+y)} \right]_0^1 \\
& = -7/2 + 5 \log(2)
\end{aligned} \tag{8}$$

513 For the inequality, we have an intuition that doesn't depend on  $k$ , notice that:

$$\begin{aligned}
\sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - (P^n)_{u,s}) &\geq - \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} \left| \pi_s - (P^n)_{u,s} \right| \\
&\geq - \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} \alpha^n \\
&= - \sum_{u=1}^k \pi_u^2 \alpha^n \\
&\leq \alpha^n
\end{aligned}$$

514 As long as  $\alpha$  isn't concentrated around 1, then this shows that the magnitude of the RHS is bounded  
515 by a term that shrinks exponentially in  $n$ . For  $k = 2$ , we will find a similar bound, and then show  
516 separately that for all  $n$  for which the bound fails, the inequality still holds true.

$$\sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - (P^n)_{u,s}) = \frac{P_{1,2} P_{2,1} (4P_{1,2} P_{2,1} - P_{1,2} - P_{2,1})}{(P_{1,2} + P_{2,1})^3} (1 - P_{1,2} - P_{2,1})^n$$

517 We can show that for any choice of  $P_{1,2}$  and  $P_{2,1}$  on the unit square,

$$\left| \frac{P_{1,2} P_{2,1} (4P_{1,2} P_{2,1} - P_{1,2} - P_{2,1})}{(P_{1,2} + P_{2,1})^3} \right| \leq \frac{1}{4}$$

518 To see why this is true, observe that,

$$\begin{aligned}
& (4P_{1,2}P_{2,1} - P_{1,2} - P_{2,1})^2 \\
&= 16P_{1,2}^2P_{2,1}^2 + (P_{1,2} + P_{2,1})^2 - 8(P_{1,2} + P_{2,1})P_{1,2}P_{2,1} \\
&\leq 16P_{1,2}^2P_{2,1}^2 + (P_{1,2} + P_{2,1})^2 - 4(P_{1,2} + P_{2,1})^2P_{1,2}P_{2,1} && \text{since } P_{1,2} + P_{2,1} \leq 2 \\
&= 16P_{1,2}^2P_{2,1}^2 + (P_{1,2} + P_{2,1})^2 - 4P_{1,2}P_{2,1}((P_{1,2} + P_{2,1})^2 - 4P_{1,2}P_{2,1}) \\
&= (P_{1,2} + P_{2,1})^2 - 4P_{1,2}P_{2,1}(P_{1,2} - P_{2,1})^2 \\
&\leq (P_{1,2} + P_{2,1})^2
\end{aligned}$$

519 Using the above, we have

$$\begin{aligned}
\left( \frac{P_{1,2}P_{2,1}(4P_{1,2}P_{2,1} - P_{1,2} - P_{2,1})}{(P_{1,2} + P_{2,1})^3} \right)^2 &\leq \frac{P_{1,2}^2P_{2,1}^2(P_{1,2} + P_{2,1})^2}{(P_{1,2} + P_{2,1})^6} \\
&= \frac{P_{1,2}^2P_{2,1}^2}{(P_{1,2} + P_{2,1})^4} \\
&\leq \frac{P_{1,2}^2P_{2,1}^2}{16P_{1,2}^2P_{2,1}^2} && \text{using } (P_{1,2} + P_{2,1})^2 \geq 4P_{1,2}P_{2,1} \\
&= \frac{1}{16}.
\end{aligned}$$

520 So,

$$\left\| \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} \left( \pi_s - (P^n)_{u,s} \right) \right\| \leq \frac{1}{4} |1 - P_{1,2} - P_{2,1}|^n$$

521 Now,

$$\begin{aligned}
\mathbb{E} \left[ -1 \frac{1}{4} |1 - P_{1,2} - P_{2,1}|^n \right] &= -\frac{1}{4} \int_0^1 \int_0^1 |1 - x - y|^n \\
&= -\frac{1}{4} \frac{2}{(n+1)(n+2)} \\
&= -\frac{1}{2(n+1)(n+2)}
\end{aligned}$$

522 Notice that this decreases in  $n$ , and at  $n = 3$ ,  $\frac{1}{2(3+1)(3+2)} = \frac{1}{40} = 0.025$  which is less in magnitude  
523 than the value we proved at  $n = 1$ ,  $|-7/2 + 5 \log 2 \approx 0.034$ . So, solving for  $n = 2$  (verified by a  
524 symbolic algebra program)

$$\mathbb{E} \left[ \frac{P_{1,2}P_{2,1}(-P_{1,2} - P_{2,1} + 1)^2 \cdot (2P_{1,2}P_{2,1} + P_{1,2}(P_{2,1} - 1) + P_{2,1}(P_{1,2} - 1))}{(P_{1,2} + P_{2,1})^3} \right] = -\frac{413}{60} + \frac{149 \log(2)}{15} \approx 0.002$$

525 Which is not only greater than  $-7/2 + 5 \log 2$ , but positive. Lastly, we simply need to show that the  
526 inequality holds at  $n = 0$ , and we are done.

$$\begin{aligned}
& \mathbb{E} \left[ \frac{P_{1,2}P_{2,1}(-P_{1,2} - P_{2,1} + 1)^0 \cdot (2P_{1,2}P_{2,1} + P_{1,2}(P_{2,1} - 1) + P_{2,1}(P_{1,2} - 1))}{(P_{1,2} + P_{2,1})^3} \right] \\
&= -\mathbb{E} \left[ \frac{P_{1,2}P_{2,1} \cdot (2P_{1,2}P_{2,1} + P_{1,2}(P_{2,1} - 1) + P_{2,1}(P_{1,2} - 1))}{(P_{1,2} + P_{2,1})^3} \right] \\
&= -7/6 + 5 * \log(2)/3 \approx -0.0114
\end{aligned}$$

527 Which is greater than  $-7/2 + 5 \log 2$ , completing our proof.  $\square$



**Lemma D.5.** *If  $P$  is a uniformly random doubly stochastic matrix, then,*

$$\mathbb{E} \left[ \pi_a^2 \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] = \mathbb{E} \left[ \pi_a \pi_b \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right]$$

for all  $a, b$  and

$$\mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - P_{u,s}) \right] < \mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - (P^a)_{u,s}) \right]$$

For all non-negative  $a \neq 1$ . and,

$$\mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - P_{u,s}) \right] < 0$$

528 *Proof.* We will use the fact that for doubly stochastic matrices, the stationary distribution is the  
529 uniform vector  $\frac{1}{k} \mathbf{1}$ .

530 The first equality follows directly from  $\pi_a = \frac{1}{k} = \pi_b$ . Now we will prove the inequality.

$$\begin{aligned} \mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - (P^a)_{u,s}) \right] &= \mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \frac{1}{k^2} P_{u,s} \left( \frac{1}{k} - (P^a)_{u,s} \right) \right] \\ &= \frac{1}{k^2} - \frac{1}{k^2} \sum_{s=1}^k \sum_{u=1}^k \mathbb{E} \left[ P_{u,s} (P^a)_{u,s} \right] \end{aligned}$$

531 By Cauchy Schwartz,

$$\begin{aligned} &= \frac{1}{k^2} - \frac{1}{k^2} \mathbb{E} [\langle P, P^a \rangle_F] \\ &> \frac{1}{k^2} - \frac{1}{k^2} \mathbb{E} [\|P\|_F \|P^a\|_F] \end{aligned}$$

532 We can make the above inequality strict because Cauchy Schwartz is only tight when the vectorizations  
533 of  $P$  and  $P^a$  are linearly dependent, since both are still doubly stochastic, this can only happen when  
534  $P = P^a$ , which occurs only when each row of  $P$  is identical, which happens with probability zero.  
535 For now assume  $a > 0$ , then,

$$\begin{aligned} &= \frac{1}{k^2} - \frac{1}{k^2} \mathbb{E} [\|P\|_F \|P^a\|_F] \\ &= \frac{1}{k^2} - \frac{1}{k^2} \mathbb{E} [\|P\|_F \|P P^{a-1}\|_F] \\ &= \frac{1}{k^2} - \frac{1}{k^2} \mathbb{E} \left[ \|P\|_F \left\| \sum_i \alpha_i \Lambda_i P^{a-1} \right\|_F \right] \\ &\geq \frac{1}{k^2} - \frac{1}{k^2} \sum_i \alpha_i \mathbb{E} [\|P\|_F \|\Lambda_i P^{a-1}\|_F] \\ &= \frac{1}{k^2} - \frac{1}{k^2} \sum_i \alpha_i \mathbb{E} [\|P\|_F \|P^{a-1}\|_F] \\ &= \frac{1}{k^2} - \frac{1}{k^2} \mathbb{E} [\|P\|_F \|P^{a-1}\|_F] \end{aligned}$$

536 The third step used the well known Birkhoff-Von Nuemann Theorem [7] that any doubly stochastic  
537 matrix  $P$  is the convex combination of permutation matrices, so  $P = \sum_i \alpha_i \Lambda_i$  for some permutation  
538 matrices  $\Lambda_i$  and constants  $\alpha_i > 0$  with  $\sum_i \alpha_i = 1$ . The inequality step uses Jensen's inequality.  
539 Induction on positive  $a$  yields the desired inequality for positive  $a$ .

540 Now consider the remaining case,  $a = 0$ ,

$$\frac{1}{k^2} - \frac{1}{k^2} \sum_{s=1}^k \sum_{u=1}^k \mathbb{E} [P_{u,s} (P^0)_{u,s}] = \frac{1}{k^2} - \frac{1}{k^2} \sum_{s=1}^k \mathbb{E} [P_{s,s}]$$

541 Since cycling each row or column by 1 in a doubly stochastic matrix results in a doubly stochastic  
542 matrix, by symmetry the marginal distributions of any two entries in  $P$  are identical, so,

$$\frac{1}{k^2} - \frac{1}{k^2} \sum_{s=1}^k \mathbb{E} [P_{s,s}] = \frac{1}{k^2} - \frac{1}{k^2} = 0$$

543 While at  $a = 1$ , we have

$$\frac{1}{k^2} - \frac{1}{k^2} \sum_{s=1}^k \sum_{u=1}^k \mathbb{E} [P_{u,s}^2] < 0$$

544 Note that equality only occurs when  $P = \mathbb{1}\mathbb{1}^\top$ , which occurs with probability 0, hence why the  
545 inequality is strict.

546

□

**Lemma D.6.** *If  $P$  is a uniformly random  $k$  by  $k$  stochastic matrix subject to each row being the same, then,*

$$\mathbb{E} \left[ \pi_a^2 \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] < \mathbb{E} \left[ \pi_a \pi_b \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] < 0$$

547 and

$$\frac{\mathbb{E} \left[ \pi_a^2 \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right]}{\mathbb{E} \left[ \pi_a \pi_b \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right]} \geq \frac{8}{5}$$

for all  $a$  and  $b$  and

$$\mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - (P^a)_{u,s}) \right] = 0$$

548 For all  $a$ .

549 *Proof.* The equality statement follows from the facts that for such transition matrices,  $P^a = P$  for all  
550 natural  $a > 0$ , and that the stationary distribution matches the rows, that is, for any  $a, b$ ,  $\pi_b = P_{a,b}$ ,

$$\mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - (P^a)_{u,s}) \right] = \mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - P_{u,s}) \right] = \mathbb{E} \left[ \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} (\pi_s - \pi_s) \right] = 0$$

551 Now we will do the inequalities. We will also use the following facts derived from the moments of  
552 the Dirichlet distribution,

$$E [\|\pi\|_2^2] = \frac{2}{k+1}$$

$$E [\|\pi\|_2^4] = \frac{4(k+5)}{(k+1)(k+2)(k+3)}$$

553 So,

$$\begin{aligned} \mathbb{E} \left[ \pi_a^2 \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] &= \mathbb{E} \left[ \pi_a^2 \left( \frac{1}{k} - \sum_{s=1}^k \pi_s^2 \right) \right] \\ &= \frac{1}{k} \mathbb{E} \left[ \|\pi\|_2^2 \left( \frac{1}{k} - \|\pi\|_2^2 \right) \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{k^2} \mathbb{E} [\|\pi\|_2^2] - \frac{1}{k} \mathbb{E} [(\|\pi\|_2^4)] \\
&= \frac{2}{k^2(k+1)} - \frac{4(k+5)}{k(k+1)(k+2)(k+3)}
\end{aligned}$$

554 Which is negative for all  $k \geq 2$ . And,

$$\begin{aligned}
\mathbb{E} \left[ \pi_a \pi_b \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right] &= \mathbb{E} \left[ \pi_a \pi_b \left( \frac{1}{k} - \sum_{s=1}^k \pi_s^2 \right) \right] \\
&= \frac{1}{k^2} \mathbb{E} \left[ \left( \frac{1}{k} - \|\pi\|_2^2 \right) \right] \\
&= \frac{1}{k^3} - \frac{1}{k^2} \mathbb{E} [\|\pi\|_2^2] \\
&= \frac{1}{k^3} - \frac{2}{k^2(k+1)}
\end{aligned}$$

555 Which is also negative for all  $k \geq 2$ . Finally, notice that

$$\frac{\frac{2}{k^2(k+1)} - \frac{4(k+5)}{k(k+1)(k+2)(k+3)}}{\frac{1}{k^3} - \frac{2}{k^2(k+1)}} \geq \frac{8}{5}$$

556 For all  $k \geq 2$ . □

## 557 D.6 Approximation Lemmas

558 The following lemma is a well known property of stochastic matrices, (see Lemma 3.3.2 Gallager  
559 [17] for example).

**Lemma D.7.** *Let  $\alpha = 1 - 2 \min_{i,j} P_{i,j}$ . Then, for any  $i, j$*

$$\left| (P^n)_{i,j} - \pi_j \right| \leq \alpha^n$$

560 Lemma D.8 and lemma D.9 both share similar intuitions and proofs. They largely rely on lemma  
561 D.7, which shows that  $(P^n)_{i,j}$  approaches  $\pi_j$  exponentially fast with respect to  $n$ , to show that over  
562 the course of summations over  $n$  the stationary distribution dominates, allowing us to simplify the  
563 expressions.

564 **Lemma D.8.** *Let  $P$  be a stochastic matrix with all positive entries, and let  $a, b$  be states. Assume  
565 that  $\min_{i,j} P_{i,j}$  is positive and doesn't depend on  $T$ . Then,*

$$\begin{aligned}
&\pi_b \sum_{i=0}^T \left( \frac{1}{k} (P^i)_{b,a} - \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} (P^i)_{s,a} \sum_{j=0}^{T-i} (P^j)_{b,s} \right) \\
&= \pi_b \pi_a (T+1) \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) + O(\log T).
\end{aligned}$$

566 *Proof.* Let us bound the magnitude of the difference between the two expressions.

$$\begin{aligned}
&\left| \pi_b \sum_{i=0}^T \left( \frac{1}{k} (P^i)_{b,a} - \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} (P^i)_{s,a} \sum_{j=0}^{T-i} (P^j)_{b,s} \right) - \pi_b \pi_a (T+1) \left( \frac{1}{k} - \sum_{s=1}^k P_{a,s} \pi_s \right) \right| \\
&= \left| \pi_b \sum_{i=0}^T \left( \frac{1}{k} \left( (P^i)_{b,a} - \pi_a \right) - \sum_{s=1}^k P_{a,s} \left( \frac{1}{T-i+1} (P^i)_{s,a} \sum_{j=0}^{T-i} (P^j)_{b,s} - \pi_s \pi_a \right) \right) \right| \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \left| (P^i)_{b,a} - \pi_a \right| + \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \left| (P^i)_{s,a} (P^j)_{b,s} - \pi_s \pi_a \right| \right)
\end{aligned}$$

$$\begin{aligned}
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \left| (P^i)_{s,a} (P^j)_{b,s} - \pi_s \pi_a \right| \right) \\
&= \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \left| ((P^j)_{b,s} (P^i)_{s,a} - \pi_a) + \pi_a ((P^j)_{b,s} - \pi_s) \right| \right) \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \left( (P^j)_{b,s} \left| (P^i)_{s,a} - \pi_a \right| + \pi_a \left| (P^j)_{b,s} - \pi_s \right| \right) \right) \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \left( (P^j)_{b,s} \alpha^i + \pi_a \alpha^j \right) \right) \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \sum_{s=1}^k P_{a,s} \frac{1}{T-i+1} \sum_{j=0}^{T-i} (\alpha^i + \alpha^j) \right) \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \sum_{s=1}^k P_{a,s} \left( \alpha^i + \frac{1}{T-i+1} \frac{1 - \alpha^{T-i+1}}{1 - \alpha} \right) \right) \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \sum_{s=1}^k P_{a,s} \left( \alpha^i + \frac{1}{T-i+1} \frac{1}{1 - \alpha} \right) \right) \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \alpha^i + \frac{1}{T-i+1} \frac{1}{1 - \alpha} \right) \\
&\leq \pi_b \sum_{i=0}^T \left( \frac{1}{k} \alpha^i + \alpha^i + \frac{1}{T-i+1} \frac{1}{1 - \alpha} \right) \\
&\leq \pi_b \left( \left( 1 + \frac{1}{k} \right) \frac{1 - \alpha^{T+1}}{1 - \alpha} + \frac{\log(T+1) + 1}{1 - \alpha} \right) \\
&\leq \pi_b \frac{2 + \frac{1}{k} + \log(T+1)}{1 - \alpha} \\
&\leq \frac{2 \log T}{1 - \alpha} \\
&= \frac{\log T}{\min_{i,j} P_{i,j}} \\
&= O(\log T)
\end{aligned}$$

567 The last step follows from our assumption, completing the proof.  $\square$

568 **Lemma D.9.** *Let  $P$  be a stochastic matrix with all positive entries, and let  $a, b$  be states. Assume*  
569 *that  $\min_{i,j} P_{i,j}$  is positive and doesn't depend on  $T$ . Then,*

$$\begin{aligned}
&\sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \sum_{i=a}^T \frac{\pi_u}{T-i+1} (P^i)_{s,u} \left( (P^a)_{u,s} - \frac{1}{T-i+1} \sum_{j=0}^{T-i} (P^j)_{u,s} \right) \\
&= (\log(T+1) - \log(a+1)) \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 P_{u,s} \left( \pi_s - (P^a)_{u,s} \right) + O(1)
\end{aligned}$$

570 *Proof.* First notice that,

$$\sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \pi_u^2 \left( (P^a)_{u,s} - \pi_s \right) = \sum_{s=1}^k \sum_{u=1}^k \pi_u^2 \left( \frac{1}{k} \left( (P^a)_{u,s} - \pi_s \right) - P_{u,s} \left( (P^a)_{u,s} - \pi_s \right) \right)$$

$$\begin{aligned}
&= \sum_{u=1}^k \pi_u^2 \left( \sum_{s=1}^k \frac{1}{k} \left( (P^a)_{u,s} - \pi_s \right) - \sum_{s=1}^k P_{u,s} \left( (P^a)_{u,s} - \pi_s \right) \right) \\
&= \sum_{u=1}^k \pi_u^2 \left( \frac{1}{k} (1-1) - \sum_{s=1}^k P_{u,s} \left( (P^a)_{u,s} - \pi_s \right) \right) \\
&= \sum_{u=1}^k \pi_u^2 \sum_{s=1}^k P_{u,s} \left( \pi_s - (P^a)_{u,s} \right)
\end{aligned}$$

571 We will bound the distance between  $\sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \pi_u^2 \left( (P^a)_{u,s} - \pi_s \right)$  and  $\mathbb{E}_{x|P} \left[ \frac{\partial L_T}{\partial v_a} \right]$ .  
572 Define  $\alpha = 1 - 2 \min_{i,j} P_{i,j}$  as in lemma D.7.

$$\begin{aligned}
&= \left| \sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \sum_{i=a}^T \frac{\pi_u}{T-i+1} (P^i)_{s,u} \left( (P^a)_{u,s} - \frac{1}{T-i+1} \sum_{j=0}^{T-i} (P^j)_{u,s} \right) \right. \\
&\quad \left. - \sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \sum_{i=a}^T \pi_u^2 \frac{1}{T-i+1} \left( (P^a)_{u,s} - \pi_s \right) \right| \\
&= \left| \sum_{s=1}^k \sum_{u=1}^k \left( \frac{1}{k} - P_{u,s} \right) \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left( (P^i)_{s,u} \left( (P^a)_{u,s} - \frac{1}{T-i+1} \sum_{j=0}^{T-i} (P^j)_{u,s} \right) - \pi_u \left( (P^a)_{u,s} - \pi_s \right) \right) \right| \\
&\leq \sum_{s=1}^k \sum_{u=1}^k \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left| (P^i)_{s,u} \left( (P^a)_{u,s} - \frac{1}{T-i+1} \sum_{j=0}^{T-i} (P^j)_{u,s} \right) - \pi_u \left( (P^a)_{u,s} - \pi_s \right) \right| \\
&\leq \sum_{s=1}^k \sum_{u=1}^k \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left| (P^i)_{s,u} \left( \pi_s - \frac{1}{T-i+1} \sum_{j=0}^{T-i} (P^j)_{u,s} \right) - \left( (P^a)_{u,s} - \pi_s \right) \left( (P^i)_{s,u} - \pi_u \right) \right| \\
&\leq \sum_{s=1}^k \sum_{u=1}^k \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left| (P^i)_{s,u} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \left( \pi_s - (P^j)_{u,s} \right) - \left( (P^a)_{u,s} - \pi_s \right) \left( (P^i)_{s,u} - \pi_u \right) \right| \\
&\leq \sum_{s=1}^k \sum_{u=1}^k \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left( (P^i)_{s,u} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \left| \pi_s - (P^j)_{u,s} \right| + \left( (P^a)_{u,s} - \pi_s \right) \left| (P^i)_{s,u} - \pi_u \right| \right) \\
&\leq \sum_{s=1}^k \sum_{u=1}^k \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left( (P^i)_{s,u} \frac{1}{T-i+1} \sum_{j=0}^{T-i} \alpha^j + \left( (P^a)_{u,s} - \pi_s \right) \alpha^i \right) \quad \text{By lemma D.7} \\
&= \sum_{s=1}^k \sum_{u=1}^k \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left( (P^i)_{s,u} \frac{1}{T-i+1} \frac{1 - \alpha^{T-i+1}}{1 - \alpha} + \left( (P^a)_{u,s} - \pi_s \right) \alpha^i \right) \\
&\leq \sum_{s=1}^k \sum_{u=1}^k \sum_{i=a}^T \frac{\pi_u}{T-i+1} \left( \frac{1}{T-i+1} \frac{1}{1 - \alpha} + \left( (P^a)_{u,s} - \pi_s \right) \alpha^i \right) \\
&\leq \sum_{i=a}^T \frac{1}{T-i+1} \left( \sum_{s=1}^k \frac{1}{T-i+1} \frac{1}{1 - \alpha} + \sum_{s=1}^k \left( \sum_{u=1}^k \pi_u (P^a)_{u,s} - \pi_s \right) \alpha^i \right) \\
&\leq \sum_{i=a}^T \frac{1}{T-i+1} \left( \sum_{s=1}^k \frac{1}{T-i+1} \frac{1}{1 - \alpha} + \sum_{s=1}^k (\pi_s - \pi_s) \alpha^i \right) \\
&\leq \sum_{i=a}^T \frac{1}{T-i+1} \left( \frac{k}{T-i+1} \frac{1}{1 - \alpha} \right)
\end{aligned}$$

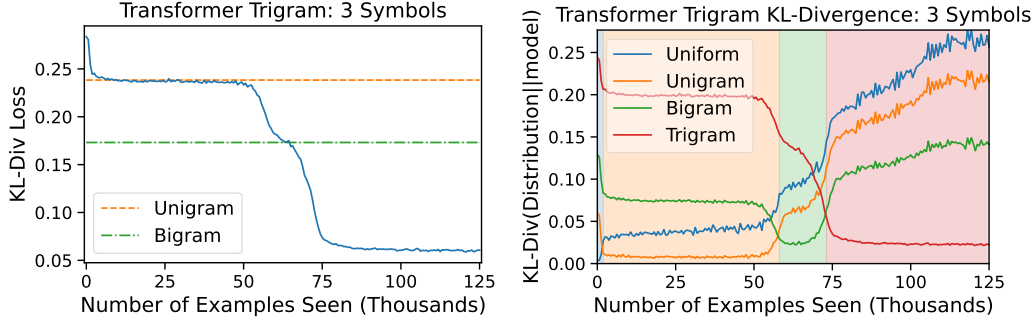


Figure 12: Three-headed transformer trained on In-Context Learning 3-grams (trigrams), with context length 200. (left) Loss during training. The model hierarchically converges close to the Bayes optimal solution. (right) KL divergence between the model and different strategies during training. As we observe, there are 4 stages of learning, each of them corresponding to a different algorithm implemented by the model.

$$\begin{aligned}
 &\leq \sum_{i=a}^T \frac{1}{(T-i+1)^2} \frac{k}{1-\alpha} \\
 &\leq \frac{2k}{1-\alpha} \\
 &= \frac{k}{\min_{i,j} P_{i,j}} \\
 &= O(1)
 \end{aligned}$$

573 The last step follows from our assumption, and the fact that  $k$  does not depend on  $T$ . □

## 574 E Beyond Bigrams: $n$ -gram Statistics

575 Finally, we investigate the performance of transformers on learning in-context  $n$ -grams for  $n > 2$ ; in  
 576 particular, 3-grams. We train attention-only transformers with three heads in each layer by minimizing  
 577 the in-context cross entropy loss with the Adam optimizer. As can be seen in Figure 12 (left), the  
 578 model eventually converges to the Bayes optimal solution. Interestingly, as in the case of Markov  
 579 Chains, the model displays a “hierarchical learning” behavior characterized by long plateaus and  
 580 sudden drops. In this setup, the different strategies correspond to unigrams, bigrams and trigrams,  
 581 respectively. This is presented clearly on the right of Figure 12, where we plot the similarity of the  
 582 model with the different strategies and it exhibits the same clear pattern as in the case of  $n = 2$ .  
 583 Curiously, single attention headed models could not achieve better performance than bigrams. We  
 584 leave a more thorough investigation of  $n$ -grams for future work. This behaviour is much less stable  
 585 for different number of heads and tokens. With two heads or four heads, there is sometimes no bigram  
 586 phase and faster convergence.

## 587 NeurIPS Paper Checklist

### 588 1. Claims

589 Question: Do the main claims made in the abstract and introduction accurately reflect the  
 590 paper’s contributions and scope?

591 Answer: [Yes]

592 Justification: All claims made in our paper are supported by theoretical proofs or rigorous  
 593 empirical evaluations.

594 Guidelines:

- 595 • The answer NA means that the abstract and introduction do not include the claims made  
 596 in the paper.

- 597
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- 598
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- 600
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.
- 602
- 603

## 604 2. Limitations

605 Question: Does the paper discuss the limitations of the work performed by the authors?

606 Answer: [Yes]

607 Justification: The purpose of this work is to introduce a novel task to study in-context learning  
608 in transformers. Since it focuses on a theoretical understanding, certain assumptions must be  
609 made. We mention in the introduction that our theoretical results apply to a simplified model  
610 of a transformer, while we precisely define our experimental setup.

611 Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

## 638 3. Theory Assumptions and Proofs

639 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

641 Answer: [Yes]

642 Justification: Our main paper may have some informal theorem statements for ease of read-  
643 ability, we give formal versions of all the statements with full proofs in the appendix. We also  
644 make sure to clarify the assumptions under which our results hold.

645 Guidelines:

- The answer NA means that the paper does not include theoretical results.
  - All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
  - All assumptions should be clearly stated or referenced in the statement of any theorems.
- 646
- 647
- 648
- 649

- 650 • The proofs can either appear in the main paper or the supplemental material, but if they  
651 appear in the supplemental material, the authors are encouraged to provide a short proof  
652 sketch to provide intuition.
- 653 • Inversely, any informal proof provided in the core of the paper should be complemented  
654 by formal proofs provided in appendix or supplemental material.
- 655 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 656 4. Experimental Result Reproducibility

657 Question: Does the paper fully disclose all the information needed to reproduce the main  
658 experimental results of the paper to the extent that it affects the main claims and/or conclusions  
659 of the paper (regardless of whether the code and data are provided or not)?

660 Answer: [Yes]

661 Justification: We provide all training details in Appendix C, and information on the data  
662 generating process and architecture in Appendix B.

663 Guidelines:

- 664 • The answer NA means that the paper does not include experiments.
- 665 • If the paper includes experiments, a No answer to this question will not be perceived  
666 well by the reviewers: Making the paper reproducible is important, regardless of whether  
667 the code and data are provided or not.
- 668 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
669 to make their results reproducible or verifiable.
- 670 • Depending on the contribution, reproducibility can be accomplished in various ways.  
671 For example, if the contribution is a novel architecture, describing the architecture fully  
672 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
673 be necessary to either make it possible for others to replicate the model with the same  
674 dataset, or provide access to the model. In general, releasing code and data is often  
675 one good way to accomplish this, but reproducibility can also be provided via detailed  
676 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
677 of a large language model), releasing of a model checkpoint, or other means that are  
678 appropriate to the research performed.
- 679 • While NeurIPS does not require releasing code, the conference does require all submissions  
680 to provide some reasonable avenue for reproducibility, which may depend on the  
681 nature of the contribution. For example
- 682 (a) If the contribution is primarily a new algorithm, the paper should make it clear  
683 how to reproduce that algorithm.
- 684 (b) If the contribution is primarily a new model architecture, the paper should describe  
685 the architecture clearly and fully.
- 686 (c) If the contribution is a new model (e.g., a large language model), then there  
687 should either be a way to access this model for reproducing the results or a way to  
688 reproduce the model (e.g., with an open-source dataset or instructions for how to  
689 construct the dataset).
- 690 (d) We recognize that reproducibility may be tricky in some cases, in which case  
691 authors are welcome to describe the particular way they provide for reproducibility.  
692 In the case of closed-source models, it may be that access to the model is limited in  
693 some way (e.g., to registered users), but it should be possible for other researchers  
694 to have some path to reproducing or verifying the results.

#### 695 5. Open access to data and code

696 Question: Does the paper provide open access to the data and code, with sufficient instructions  
697 to faithfully reproduce the main experimental results, as described in supplemental material?

698 Answer: [No]

699 Justification: We plan to release the code prior to publication. However, we believe our  
700 experiments are easy to reproduce with the provided information.

701 Guidelines:

- 702 • The answer NA means that paper does not include experiments requiring code.



- 703 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 704
- 705 • While we encourage the release of code and data, we understand that this might not be
- 706 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
- 707 including code, unless this is central to the contribution (e.g., for a new open-source
- 708 benchmark).
- 709 • The instructions should contain the exact command and environment needed to run to
- 710 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 711
- 712 • The authors should provide instructions on data access and preparation, including how
- 713 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 714 • The authors should provide scripts to reproduce all experimental results for the new
- 715 proposed method and baselines. If only a subset of experiments are reproducible, they
- 716 should state which ones are omitted from the script and why.
- 717 • At submission time, to preserve anonymity, the authors should release anonymized
- 718 versions (if applicable).
- 719 • Providing as much information as possible in supplemental material (appended to the
- 720 paper) is recommended, but including URLs to data and code is permitted.

## 721 6. Experimental Setting/Details

722 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparam-

723 eters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

724 Answer: [Yes]

725 Justification: The details are available in Appendix C.

726 Guidelines:

- 727 • The answer NA means that the paper does not include experiments.
- 728 • The experimental setting should be presented in the core of the paper to a level of detail
- 729 that is necessary to appreciate the results and make sense of them.
- 730 • The full details can be provided either with the code, in appendix, or as supplemental
- 731 material.

## 732 7. Experiment Statistical Significance

733 Question: Does the paper report error bars suitably and correctly defined or other appropriate

734 information about the statistical significance of the experiments?

735 Answer: [Yes]

736 Justification: Our main results regard the multiphase nature of training, and we include Figure 5

737 which shows that for the seeds 0, 1, . . . 9 the model each time has the same patterns in the

738 training curve. The curves were shown individually instead of through error bars since the

739 main purpose of the loss curves is to show the shape, but because the phase transition doesn't

740 occur at a consistent time, adding error bars smooths out the curve making the phase transition

741 look less sharp.

742 Guidelines:

- 743 • The answer NA means that the paper does not include experiments.
- 744 • The authors should answer "Yes" if the results are accompanied by error bars, confidence
- 745 intervals, or statistical significance tests, at least for the experiments that support the
- 746 main claims of the paper.
- 747 • The factors of variability that the error bars are capturing should be clearly stated (for
- 748 example, train/test split, initialization, random drawing of some parameter, or overall
- 749 run with given experimental conditions).
- 750 • The method for calculating the error bars should be explained (closed form formula, call
- 751 to a library function, bootstrap, etc.)
- 752 • The assumptions made should be given (e.g., Normally distributed errors).
- 753 • It should be clear whether the error bar is the standard deviation or the standard error of
- 754 the mean.

- 755
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- 756
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 758
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
- 759
- 760
- 761
- 762

## 763 8. Experiments Compute Resources

764 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

765

766

767 Answer: [Yes]

768 Justification: The resources are described in Appendix C. Only a single computer was used, and none of the training runs took longer than ten minutes.

769

770 Guidelines:

- The answer NA means that the paper does not include experiments.
  - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 771
- 772
- 773
- 774
- 775
- 776
- 777
- 778

## 779 9. Code Of Ethics

780 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

781

782 Answer: [Yes]

783 Justification: We read the code of ethics in its entirety and strongly believe that our research abides by the stated code.

784

785 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 786
- 787
- 788
- 789
- 790

## 791 10. Broader Impacts

792 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

793

794 Answer: [NA]

795 Justification: Our work focuses on understanding the internal mechanisms of Transformer models on synthetic tasks. We do not foresee any direct societal impact of this work.

796

797 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
  - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- 798
- 799
- 800
- 801
- 802
- 803
- 804

- 805 • The conference expects that many papers will be foundational research and not tied  
806 to particular applications, let alone deployments. However, if there is a direct path to  
807 any negative applications, the authors should point it out. For example, it is legitimate  
808 to point out that an improvement in the quality of generative models could be used to  
809 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
810 that a generic algorithm for optimizing neural networks could enable people to train  
811 models that generate Deepfakes faster.
- 812 • The authors should consider possible harms that could arise when the technology is being  
813 used as intended and functioning correctly, harms that could arise when the technology is  
814 being used as intended but gives incorrect results, and harms following from (intentional  
815 or unintentional) misuse of the technology.
- 816 • If there are negative societal impacts, the authors could also discuss possible mitigation  
817 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
818 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
819 feedback over time, improving the efficiency and accessibility of ML).

## 820 11. Safeguards

821 Question: Does the paper describe safeguards that have been put in place for responsible  
822 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
823 image generators, or scraped datasets)?

824 Answer: [NA]

825 Justification: Our work does not release any data or models that poses safety risks.

826 Guidelines:

- 827 • The answer NA means that the paper poses no such risks.
- 828 • Released models that have a high risk for misuse or dual-use should be released with  
829 necessary safeguards to allow for controlled use of the model, for example by requiring  
830 that users adhere to usage guidelines or restrictions to access the model or implementing  
831 safety filters.
- 832 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
833 should describe how they avoided releasing unsafe images.
- 834 • We recognize that providing effective safeguards is challenging, and many papers do not  
835 require this, but we encourage authors to take this into account and make a best faith  
836 effort.

## 837 12. Licenses for existing assets

838 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the  
839 paper, properly credited and are the license and terms of use explicitly mentioned and properly  
840 respected?

841 Answer: [Yes]

842 Justification: We cite the Github repository we use as the codebase for our Transformer models.

843 Guidelines:

- 844 • The answer NA means that the paper does not use existing assets.
- 845 • The authors should cite the original paper that produced the code package or dataset.
- 846 • The authors should state which version of the asset is used and, if possible, include a  
847 URL.
- 848 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 849 • For scraped data from a particular source (e.g., website), the copyright and terms of  
850 service of that source should be provided.
- 851 • If assets are released, the license, copyright information, and terms of use in the pack-  
852 age should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has  
853 curated licenses for some datasets. Their licensing guide can help determine the license  
854 of a dataset.
- 855 • For existing datasets that are re-packaged, both the original license and the license of the  
856 derived asset (if it has changed) should be provided.

857 • If this information is not available online, the authors are encouraged to reach out to the  
858 asset's creators.

### 859 13. New Assets

860 Question: Are new assets introduced in the paper well documented and is the documentation  
861 provided alongside the assets?

862 Answer: [NA]

863 Justification: We do not introduce any new assets.

864 Guidelines:

- 865 • The answer NA means that the paper does not release new assets.
- 866 • Researchers should communicate the details of the dataset/code/model as part of their  
867 submissions via structured templates. This includes details about training, license,  
868 limitations, etc.
- 869 • The paper should discuss whether and how consent was obtained from people whose  
870 asset is used.
- 871 • At submission time, remember to anonymize your assets (if applicable). You can either  
872 create an anonymized URL or include an anonymized zip file.

### 873 14. Crowdsourcing and Research with Human Subjects

874 Question: For crowdsourcing experiments and research with human subjects, does the paper  
875 include the full text of instructions given to participants and screenshots, if applicable, as well  
876 as details about compensation (if any)?

877 Answer: [NA]

878 Justification: Our work does not involve crowdsourcing nor research with human subjects.

879 Guidelines:

- 880 • The answer NA means that the paper does not involve crowdsourcing nor research with  
881 human subjects.
- 882 • Including this information in the supplemental material is fine, but if the main contri-  
883 bution of the paper involves human subjects, then as much detail as possible should be  
884 included in the main paper.
- 885 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
886 or other labor should be paid at least the minimum wage in the country of the data  
887 collector.

### 888 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 889 Subjects

890 Question: Does the paper describe potential risks incurred by study participants, whether such  
891 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals  
892 (or an equivalent approval/review based on the requirements of your country or institution)  
893 were obtained?

894 Answer: [NA]

895 Justification: Our work does not involve crowdsourcing nor research with human subjects.

896 Guidelines:

- 897 • The answer NA means that the paper does not involve crowdsourcing nor research with  
898 human subjects.
- 899 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
900 may be required for any human subjects research. If you obtained IRB approval, you  
901 should clearly state this in the paper.
- 902 • We recognize that the procedures for this may vary significantly between institutions  
903 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
904 guidelines for their institution.
- 905 • For initial submissions, do not include any information that would break anonymity (if  
906 applicable), such as the institution conducting the review.