

# Understanding the Transient Nature of In-Context Learning: The Window of Generalization

Anonymous Author(s)  
Affiliation  
Address  
email

## Abstract

1 In-Context Learning (ICL) is one of the main mechanisms driving few shot learning  
2 capabilities of large language models (LLMs). A rich literatures explores the causal  
3 factors giving rise to this mechanism, while recent studies has pointed out that  
4 this mechanism can be transient. In this work, we study ICL on a synthetic task  
5 consisting of a probabilistic mixture of Markov chains which is simple enough  
6 to allow theoretical analysis yet rich enough to reproduce multiple phenomena  
7 discussed in the ICL literature. Here, we focus on analyzing the transient nature  
8 of ICL using this setup and elucidate the role of data and model training using a  
9 mechanistic phase diagram. Our findings conclude to: 1) A certain data diversity  
10 is required for ICL 2) a non-generalizing Bayesian solution might arise later  
11 in training if its circuit complexity is higher. We conclude: *ICL, or any other*  
12 *generalizing solution, is subject to transience if there exists a better solution*  
13 *narrowly fitting the training distribution accessible by gradient descent.*

## 14 1 Introduction

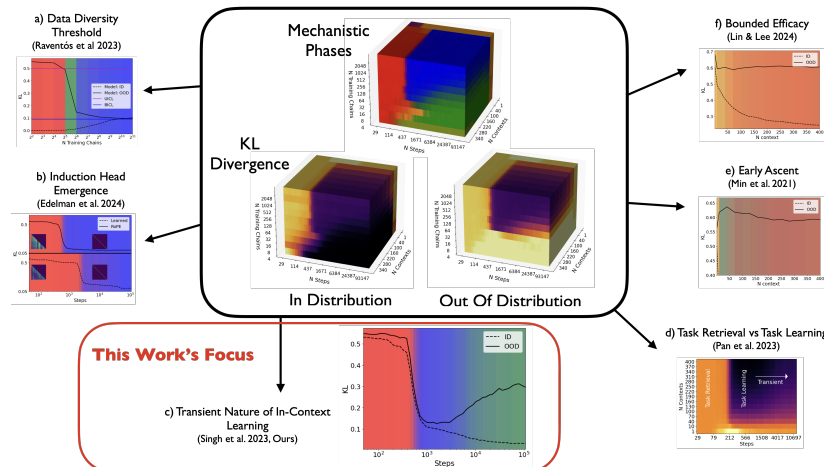


Figure 1: **Markov Mixtures data explains multiple In-Context Learning Phenomena** Autoregressive transformers trained on *Markov Mixtures* exhibits multiple phenomena of In-Context Learning. a) Data Diversity threshold for ICL [1] b) Emergence of Induction Heads for ICL [2, 3, 4, 5] c) **This work's focus: Transient Nature of ICL** [6, 7] d) Task Retrieval and Task Learning Phases of ICL [8] e) Early Ascent of ICL Accuracy [9] f) Bounded Efficacy of ICL [10]

15 In-Context Learning (ICL) is often attributed to be at the core of the impressive capabilities of  
 16 LLMs[11, 12, 13]. Motivated by this, many past work has explored different aspects of ICL,  
 17 including algorithms [14, 15, 16, 17, 18, 19], data dependence [1, 5], training dynamics [2, 3, 20, 21],  
 18 test time behavior [9, 8, 22, 23, 10]. Recently, it has been proposed that ICL can also be transient  
 19 [6, 7]. In this work, we develop a synthetic dataset which is defined in discrete token space unlike  
 20 the common linear regression setup[14, 16, 21, 1] called *Markov Mixtures*. This simple dataset  
 21 reproduces multiple phenomenologies explored in the literature as seen in Fig. 1. We then focus on  
 22 one specific phenomena, **the transient nature of ICL**, and explain why and when ICL is transient.

23 **Our main contributions is as follows:**

- 24 1. **Markov Mixtures Setup as a dataset to study ICL:** We introduce a discrete sequence  
 25 dataset, *Markov Mixture*, simple enough to postulate analytic solutions yet which reproduces  
 26 many ICL phenomena.
- 27 2. **Discovering control variables for ICL:** We discover data diversity and training time as  
 28 two sharp requirements for ICL to emerge.
- 29 3. **Explaining the transient nature of ICL** We provide a simple explanation of why ICL  
 30 emerges: “There exists a non-generalizing solution which performs better on the training  
 31 set”.
- 32 4. **Elucidating the role of data and model training for generalization** We run extensive  
 33 experiments to reveal the role of data and model training in controlling generalization  
 34 strategies.

## 35 2 Setup

36 **Markov Mixtures Dataset** Our data generating process (DGP) consists of a probabilistic mixture  
 37 of Markov chains.  $n$ ,  $k$  and  $l$  are hyperparameters of the DGP each controlling: how many transition  
 38 matrices can we draw from, how big is the state space and how long is the sequence. The dataset is  
 39 described in further detail in Appendix A

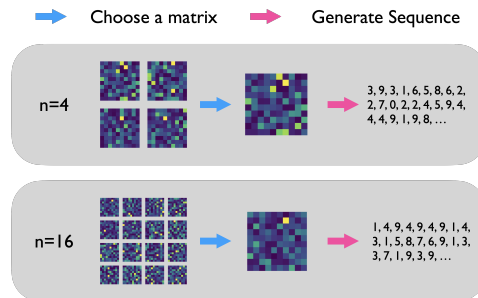


Figure 2: **Markov Mixtures Dataset** Our dataset consists of  $n$  predefined matrices from which a single matrix is chosen every time a datapoint(sequence) is drawn from the dataset.

40 **Bayesian Solutions** First, we introduce two simple Bayesian solutions which can be used to predict  
 41 the next token on the *Markov Mixtures* dataset. Since the training transition matrices are needed to  
 42 compute the solution, the model needs to internalize(memorize) the training matrices to implement  
 43 this solution. The first solution is **Unigram Likelihood Bayesian Averaging(ULBA)**, which observes  
 44 the stationary distribution of the given context and computes the likelihood over the different transition  
 45 matrices which can explain the given distribution. The second solution, conceptually very similar is  
 46 the **Bigram Likelihood Bayesian Averaging(BLBA)** solution, which is similar to ULBA, but uses a  
 47 bigram based likelihood. I.e., this solution computes the likelihood over *transitions* instead of states.  
 48 Please refer to Appendix B for the mathematical formulation of these solutions.

49 **Purely In-Context Solutions** We discuss two purely In-Context solutions. We add the prefix  
 50 “purely” to illustrate that these solutions are generalizing solutions with no performance gap between  
 51 an In-Distribution (ID) evaluation and an Out-Of-Distribution (OOD) evaluation. **Unigram In-  
 52 Context Learning (UICL)** infers the stationary token distribution from the context and draws new

53 token form that distribution. **Bigram In-Context Learning (BICL)** generates the transition matrix  
 54 by counting transitions from the context. This solution is the generalizing “world model” solution.

### 55 3 Results

56 We train a transformer[24] with a autoregressive next token cross entropy loss on our *Markov Mixtures*  
 57 dataset with  $k = 10$ ,  $l = 512$  and  $n \in 2^{\{2,3,4,5,6,7,8,9,10,11\}}$ . Please see Appendix A for more details.

#### 58 3.1 Requirements for ICL

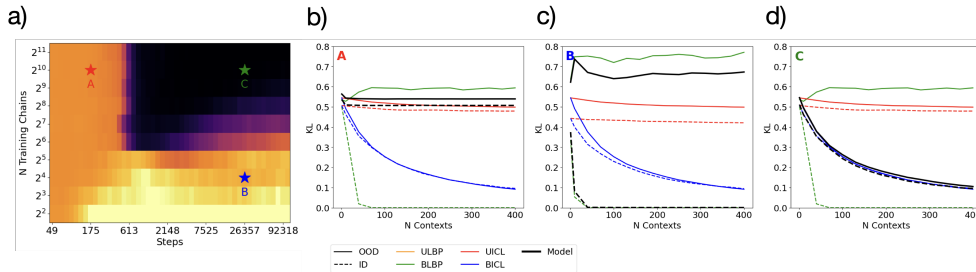


Figure 3: **Data Diversity and Enough Training is Required for ICL** a) OOD KL divergence depending on  $n$  and optimization steps. (Brighter is Higher, Worse) b,c,d) OOD and ID KL divergence of difference solutions and the model, respectively at the checkpoint corresponding to A,B,C in subpanel a).

59 Fig. 3 clearly demonstrates the joint data diversity and optimization step requirements for ICL. ICL  
 60 emerges when, in this case, at least 600 gradient steps has been applied and when there are more  
 61 than  $2^6$  latent transitions. It is notable that the compute requirement is extremely sharp, e.g. highly  
 62 emergent behavior shows up as seen in Edelman et al. [4]. Subpanel b suggestst that checkpoint A  
 63 is implementing UICL and checkpoint C is implementing BLBA, but only robustly for the ID data.  
 64 Subpanel d shows that BICL has emerged and both ID and OOD data has a lower KL with more  
 65 exemplars.

#### 66 3.2 Why Can ICL be Transient

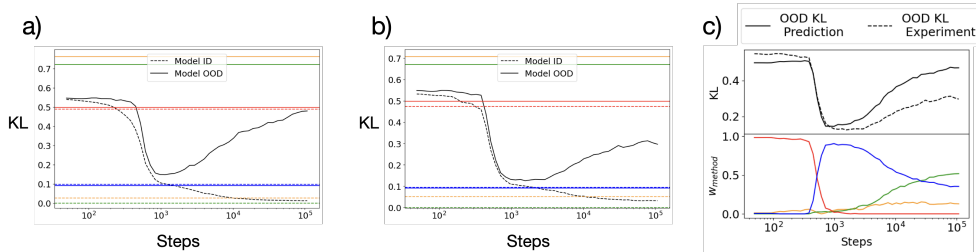


Figure 4: **The implementation BLBP drives the transience of ICL** a) ID, OOD model KL and different solution KLs.  $n=64$  b) same for  $n=128$  c) Mechanistic Decomposition of model outputs and corresponding test performance prediction.

67 Fig. 4 demonstrates why ICL can be transient. Fig. 4(a,b) shows that across different data diversity,  
 68 the BICL solution starts fading away and the OOD KL divergence starts increasing when BLBP is imple-  
 69 mented. The model does not optimize for test data, so the very mechanistic change which reduces  
 70 the training data KL harms the OOD performance driving the transience of ICL. Fig. 4(c) shows  
 71 that using a *mechanistic decomposition approach*(Explained in Appendix C), we can decompose the  
 72 model into the 4 suggested mechanisms and get their weights, only using the training distribution.  
 73 These weights are then used to predict the OOD performance as seen in the upper panel. This  
 74 suggests that mechanistic decomposition on ID data and proper modelling of different mechanisms  
 75 can explain/predict non-monotonic OOD performance.

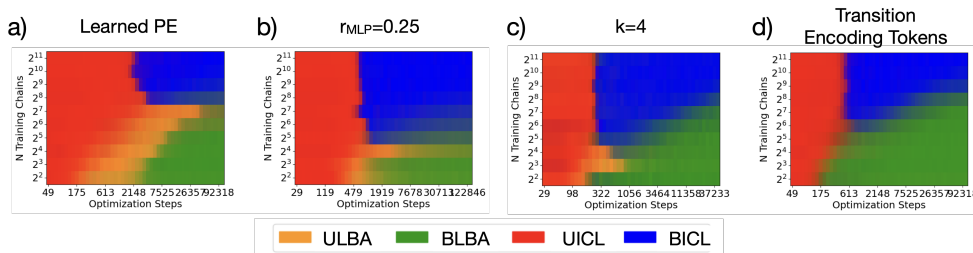


Figure 5: **Model Variants and Data Variants induces a phase diagram Change.** a) Using learned positional encodings b) Using thinner MLP layers c) using a smaller state space. d) Using transition encoding tokens.

77 Using the *mechanistic decomposition* seen above, we can draw a phase diagram of what solution the  
 78 model implements. Here, we show that we can analyze the role of data and model training using these  
 79 phase diagrams. Over different data and architectural changes, the general shape of the phase diagram  
 80 remains. Notably we find that a thinner MLP suppresses the BLBP solution which intuitively makes  
 81 sense as BLBP requires memorization of the training matrices. We also find that using tokens which  
 82 embed states AND transitions, the BLBP solution emerges faster and to higher data diversity. We  
 83 conclude: *Data defines the optimality of different solution mechanisms and model training determines*  
 84 *how these solutions are navigated.*

## 85 4 Discussion

86 **World Models v.s. Stochastic Parrots** Recently, researchers have been interested in taking a  
 87 position on the question of whether neural networks acquire world models or learn surface statistics  
 88 (commonly said stochastic parrots).<sup>1</sup> Some researchers have pointed out the possibility of a parrot[25,  
 89 26] while others found evidence towards world models[27, 28, 29]. This work suggests a simple  
 90 answer: “You get a parrot if the parrot works better, you get a world model if the world model works  
 91 better”, however, depending on your training setup, your world model *or your parrot* might not have  
 92 arrived yet.

93 **Optimization and Simplicity Bias** Past work has explored competing solutions to a task with a  
 94 simplicity perspective[30, 31]. From the simplicity perspective, we expect a transient generalizing  
 95 solution when there exists a solution which is 1) better performing with respect to the training  
 96 distribution and 2) complex enough to take a lot of optimization to learn. Our BLBA solution is  
 97 indeed such case, as one can see from its attention matrix, which requires additional diagonal weights  
 98 on top of an induction head. Our experiments with transition encoding tokens reveals that making  
 99 this solution more accessible accelerates(in training steps) and severses(in data diversity) the transient  
 100 nature of the BICL solution.

101 **Conclusion** We conclude our work with a hypothesis which we name *The Obvious Hypothesis*, as  
 102 it might be indeed obvious for some readers:

103 **The Obvious Hypothesis:** *Model training simply optimizes it on the training data, generalization*  
 104 *happens when the mechanism it currently implements on the training data turns out to be a generaliz-*  
 105 *ing one. A discovery of a better performing algorithm on the training set could lead to a deactivation*  
 106 *of the generalizing mechanism.*

<sup>1</sup>The authors believe this is not a question one can answer with rigor and generality, but more of a nu-  
 ance/position survey.

## References

- 107
- 108 [1] Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity  
109 and the emergence of non-bayesian in-context learning for regression, 2023. URL <https://arxiv.org/abs/2306.15063>.  
110
- 111 [2] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,  
112 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep  
113 Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt,  
114 Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and  
115 Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*,  
116 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- 117 [3] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom  
118 Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain,  
119 Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson  
120 Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Ka-  
121 plan, Sam McCandlish, and Chris Olah. In-context learning and induction heads, 2022. URL  
122 <https://arxiv.org/abs/2209.11895>.
- 123 [4] Benjamin L. Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The  
124 evolution of statistical induction heads: In-context learning markov chains, 2024. URL <https://arxiv.org/abs/2402.11004>.  
125
- 126 [5] Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning:  
127 Architectures and algorithms, 2024. URL <https://arxiv.org/abs/2401.12973>.
- 128 [6] Aaditya K. Singh, Stephanie C. Y. Chan, Ted Moskovitz, Erin Grant, Andrew M. Saxe, and  
129 Felix Hill. The transient nature of emergent in-context learning in transformers, 2023. URL  
130 <https://arxiv.org/abs/2311.08360>.
- 131 [7] Jesse Hoogland, George Wang, Matthew Farrugia-Roberts, Liam Carroll, Susan Wei, and  
132 Daniel Murfet. The developmental landscape of in-context learning, 2024. URL <https://arxiv.org/abs/2402.02364>.  
133
- 134 [8] Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What in-context learning "learns"  
135 in-context: Disentangling task recognition and task learning, 2023. URL <https://arxiv.org/abs/2305.09731>.  
136
- 137 [9] Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and  
138 Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning  
139 work? *arXiv preprint arXiv:2202.12837*, 2022.
- 140 [10] Ziqian Lin and Kangwook Lee. Dual operating modes of in-context learning, 2024. URL  
141 <https://arxiv.org/abs/2402.18819>.
- 142 [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
143 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
144 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 145 [12] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny  
146 Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint*  
147 *arXiv:2201.11903*, 2022.
- 148 [13] OpenAI. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- 149 [14] Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn  
150 in-context? a case study of simple function classes, 2023. URL <https://arxiv.org/abs/2208.01066>.  
151
- 152 [15] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning  
153 algorithm is in-context learning? investigations with linear models, 2023. URL <https://arxiv.org/abs/2211.15661>.  
154

- 155 [16] Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mord-  
156 vintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient  
157 descent, 2023. URL <https://arxiv.org/abs/2212.07677>.
- 158 [17] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:  
159 Provable in-context learning with in-context algorithm selection, 2023. URL <https://arxiv.org/abs/2306.04637>.  
160
- 161 [18] Yingcong Li, M. Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as  
162 algorithms: Generalization and stability in in-context learning, 2023. URL <https://arxiv.org/abs/2301.07067>.  
163
- 164 [19] Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai.  
165 How do transformers learn in-context beyond simple functions? a case study on learning with  
166 representations, 2023. URL <https://arxiv.org/abs/2310.10616>.
- 167 [20] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and vari-  
168 able creation in self-attention mechanisms. In *International Conference on Machine Learning*,  
169 pages 5793–5831. PMLR, 2022.
- 170 [21] Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models  
171 in-context, 2023. URL <https://arxiv.org/abs/2306.09927>.
- 172 [22] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan  
173 Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv*  
174 *preprint arXiv:2109.01652*, 2021.
- 175 [23] Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao  
176 Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning  
177 differently, 2023. URL <https://arxiv.org/abs/2303.03846>.
- 178 [24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
179 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- 180 [25] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On  
181 the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021*  
182 *ACM conference on fairness, accountability, and transparency*, pages 610–623, 2021.
- 183 [26] Keyon Vafa, Justin Y. Chen, Jon Kleinberg, Sendhil Mullainathan, and Ashesh Rambachan.  
184 Evaluating the world model implicit in a generative model, 2024. URL <https://arxiv.org/abs/2406.03689>.  
185
- 186 [27] Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin  
187 Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic  
188 task. In *The Eleventh International Conference on Learning Representations*, 2023. URL  
189 [https://openreview.net/forum?id=DeG07\\_TcZvT](https://openreview.net/forum?id=DeG07_TcZvT).
- 190 [28] Yida Chen, Fernanda Viégas, and Martin Wattenberg. Beyond surface statistics: Scene repre-  
191 sentations in a latent diffusion model. *arXiv preprint arXiv:2306.05720*, 2023.
- 192 [29] Adam Karvonen. Emergent world models and latent variable estimation in chess-playing  
193 language models, 2024. URL <https://arxiv.org/abs/2403.15498>.
- 194 [30] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The  
195 pitfalls of simplicity bias in neural networks, 2020. URL <https://arxiv.org/abs/2006.07710>.  
196
- 197 [31] Ekdeep Singh Lubana, Eric J Bigelow, Robert P Dick, David Krueger, and Hidenori Tanaka.  
198 Mechanistic mode connectivity. In *International Conference on Machine Learning*, pages  
199 22965–23004. PMLR, 2023.
- 200 [32] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer:  
201 Enhanced transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>.  
202
- 203 [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.

## 204 A Experimental Details

### 205 A.1 Data

206  $n$  determines the number of pre-drawn transition matrices of the training set.  $k$  is the number of  
207 states, here equivalent to the number of tokens. The transition matrix is thus  $\mathbb{R}^{(k,k)}$ . When drawing a  
208 new data point from the dataset, a transition matrix is chosen from the  $n$  matrices, say  $T_i$ . Then a  
209 start token is drawn from  $\pi_i$ , the stationary distribution corresponding to  $T_i$ . The Markov chain is  
210 continued using the same transition matrix, until a length of  $l$  is reached.  $k$  is fixed to 10 except in  
211 Fig 5(c), and  $l$  is always fixed to 512 during training.

### 212 A.2 Model

213 We use a standard transformer[24] architecture with a Rotational Positional Embedding(RoPE)[32]  
214 unless otherwise specified. The transformer has 2 layers, 4 heads and an embedding dimension of 64.  
215 The MLP layers upsample to 256 dimensions(a ratio of 4) except in Fig 5(b) where the ratio is set to  
216 0.25. The number of steps are adjusted to match the same total compute FLOPs.

### 217 A.3 Training

218 We train the model with the DGP using a batch size of 128. We have tried a batch size of 64 and  
219 256 and observed the same results. We use the AdamW[33] optimizer with a learning rate of  $1e-3$ .  
220 Model parameters are initialized to  $\mathcal{N}(\mathbf{0}, 0.02\Sigma)$ .

## 221 B Additional Details for Solutions

222 Eq. 1 and Eq. 2 respectively defines the unigram and bigram likelihood over the training set transition  
223 matrices.  $\pi_i$  is the stationary distribution corresponding to  $T_i$ , which satisfies  $\pi T = \pi$ .

$$\text{Unigram Likelihood: } \mathcal{L}_U(T_i|\mathbf{x}_{0:t}) = \prod_{j=1}^t \pi_{i[x_j]} \quad (1)$$

$$\text{Bigram Likelihood: } \mathcal{L}_B(T_i|\mathbf{x}_{0:t}) = \prod_{j=1}^{t-1} T_{i[x_{j-1}, x_j]} \quad (2)$$

$$\text{Bayesian Solution: } p(x_t|\mathbf{x}_{0:t-1}) = \sum_i p_i \mathcal{L}(T_i|\mathbf{x}_{0:t-1}) T_{i[x_{t-1}, x_t]} \quad (3)$$

224 Eq. 3 describes the Bayesian averaging process, which takes either of the unigram or bigram likelihood  
225 and computes the next token probability. The next token probability can be simply understood as a  
226 weighted average over hypotheses of what transition matrix might be underlying the data.

227 Eq. 4 and Eq. 5 respectively defines the transition matrix inferred by the UICL solution and the BICL  
228 solution. Note that the transition matrix elements inferred from UICL does not depend on the indice  
229  $i$ , i.e. it does not depend on the current state as it is a unigram solution.

$$\text{Unigram ICL: } \hat{T}_{[i,j]}^U(\mathbf{x}_t) = \hat{\pi}_{[j]}(\mathbf{x}_t) = \frac{\sum_{k=1}^t \delta_{x_k,j}}{t} \quad (4)$$

$$\text{Bigram ICL: } \hat{T}_{[i,j]}^B(\mathbf{x}_t) = \frac{1 + \sum_{k=1}^{t-1} \delta_{x_k,i} \delta_{x_{k+1},j}}{k + \sum_{k=1}^{t-1} \delta_{x_k,i}} \quad (5)$$

## 230 C Mechanistic Decomposition

231 We decompose the model output logits linearly as:

$$M(x_t|\mathbf{x}_{0:t-1}) = w_{ULBA}ULBA(x_t|\mathbf{x}_{0:t-1}) + w_{BLBA}BLBA(x_t|\mathbf{x}_{0:t-1}) \quad (6)$$

$$+ w_{UICL}UICL(x_t|\mathbf{x}_{0:t-1}) + w_{BICL}BICL(x_t|\mathbf{x}_{0:t-1}) \quad (7)$$

232 This system might seem under determined just from this equation, however  
233 *ULBA, BLBA, UICL, BICL* has no hyperparameters and there is an infinite amount of  
234  $\mathbf{x}_{0:t}$  we can sample, so there exists a unique optimal solution. In practice, we draw 300 chains to  
235 determine the weights at each checkpoint. The phase diagrams in Fig 5 are constructed by assigning  
236 each weight a color.