
UNDERSTANDING CONTEXTUAL RECALL IN TRANSFORMERS: HOW FINETUNING ENABLES IN-CONTEXT REASONING OVER PRETRAINING KNOWLEDGE

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformer-based language models excel at in-context learning (ICL), where they can adapt to new tasks based on contextual examples, without parameter updates. In a specific form of ICL, which we refer to as *contextual recall*, models pretrained on open-ended text leverage pairwise examples to recall specific facts in novel prompt formats. We investigate whether contextual recall emerges from pretraining alone, what finetuning is required, and what mechanisms drive the necessary representations. For this, we introduce a controlled synthetic framework where pretraining sequences consist of subject-grammar-attribute tuples, with attribute types tied to grammar statistics. We demonstrate that while such pretraining successfully yields factual knowledge, it is insufficient for contextual recall: models fail to implicitly infer attribute types when the grammar statistics are removed in ICL prompts. However, we show that finetuning on tasks requiring implicit inference, distinct from the ICL evaluation, using a subset of subjects, triggers the emergence of contextual recall across all subjects. This transition is accompanied by the formation of low-dimensional latent encodings of the shared attribute type. For mechanistic insight, we derive a construction for an attention-only transformer that replicates the transition from factual to contextual recall, corroborated by empirical validation.

1 INTRODUCTION

Transformer-based large language models (LLMs) exhibit remarkable abilities to extrapolate far beyond tasks seen during training. A notable instance of this extrapolation is in-context learning (ICL) (Brown et al., 2020), where models can adapt to new tasks based on contextual examples, without parameter updates.

In this paper, we investigate a specific form of ICL, which we refer to as *contextual recall*. Here, a model trained on open-ended text acquires factual knowledge and is later able to recall specific facts when prompted with example pairs in an unseen format. To illustrate, pretraining data might include descriptions of various landmarks, from which the model learns multiple attributes for each—such as the country where it is located, the year it was built, or its architectural style. At test time, the model receives a prompt of the form [Niagara Falls, Canada. Colosseum, Italy. Parthenon,] and must generate [Greece]. The prompt contains no explicit indication that the relevant attribute is “country” rather than, say, “year built”; the model must infer the attribute type from the in-context examples and recall the corresponding fact for the query subject.

The ability to succeed at contextual recall therefore requires both the acquisition of factual knowledge and adaptability to novel prompt formats, since the model may never have seen these facts presented as implicit subject-attribute pairs during pretraining (see Appendix A for a detailed discussion on related work on understanding ICL and factual recall using controlled, synthetic settings). In this work, we investigate the origins of this capability:

Does the ability to do contextual recall emerge naturally from pretraining, or does it necessitate specific finetuning? Furthermore, what mechanisms within the Transformer architecture enable the emergence of this ability?

Contributions. We summarize our contributions below.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

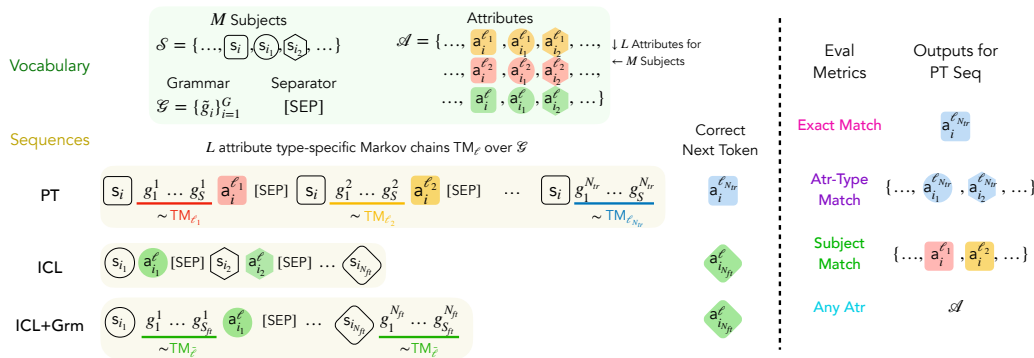


Figure 1: Illustration of the **data generation process** and the **evaluation metrics**. The vocabulary consists of M subjects with L attributes each (where the attribute tokens may be shared across the subjects), G grammar tokens and a separator token $.$ We sample L Markov chains, one for each attribute type. We consider three types of sequences (see Appendix B for details): PT sequences contain subject-grammar-attribute tuples, with attribute type information encoded in the grammar sequence statistics, for one subject, ICL sequences contain pairwise subject-attribute examples for a shared attribute type across subjects, but no grammar, and ICL+Grm sequences are analogous to ICL sequences, but contain subject-grammar-attribute tuples. We evaluate whether the model’s prediction (on PT sequences) matches the ground truth attribute (*Exact Match*), any attribute of the same type as the last subsequence (*Attribute-type Match*), any attribute that belongs to the same subject as the sequence (*Subject Match*), or any attribute token (*Any Attribute*).

First, we introduce a controlled synthetic framework to study the emergence of contextual recall in transformers. The pretraining (PT) sequences, used to instill factual knowledge in the model, contain multiple attributes of a single subject, interspersed with grammar tokens that encode information about each attribute type. The ICL sequences, used to evaluate contextual recall ability, consist of subject-attribute pairs sharing a common attribute type across different subjects, without any grammar. See Fig. 1 for an illustration of the data generation process and Appendix B for a detailed description.

In Section 2, we show that transformers trained with PT sequences succeed on factual recall, but fail to generalize to ICL sequences, which require the model to implicitly infer the attribute type from the in-context examples. However, we find that finetuning the model on sequences that are distinct from ICL sequences, but require implicit inference, using a subset of subjects, enables out-of-distribution generalization on ICL sequences on the held-out subjects (Section 3). We also probe the effect of some key dataset parameters on the model’s performance on PT and ICL sequences in Appendix D.1.

In Section 4, we analyze model representations and find that finetuning induces the formation of low-dimensional encodings of the shared attribute type based on the in-context examples. These representations become more disentangled as the number of in-context demonstrations increases.

Finally, in Appendix C, for mechanistic insight, we consider a simpler synthetic task and present constructions for an attention-only model that succeeds on factual recall after pretraining and contextual recall after finetuning, and corroborate them with empirical validation.

2 PRETRAINING ON PT SEQUENCES

We train a two-layer, single-head, decoder-only transformer on PT sequences to minimize the standard next-token prediction objective. We use an online training setup: at each iteration, we generate a fresh batch of PT sequences using the process described in App. B. Unless stated otherwise, we fix $M = 256$, $L = 8$, $S = 80$. We evaluate the model on two distinct held-out sets: i) PT sequences, to evaluate factual recall, and ii) ICL sequences, to evaluate contextual recall capabilities. For both, we measure the accuracy of predicting the final token given the preceding context, using four metrics, as illustrated in Fig. 1. See App. D for further details.

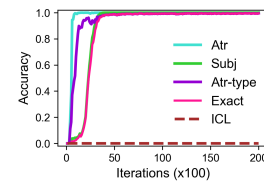


Figure 2: Transformer trained on PT sequences performs well on PT sequences, but does not generalize to ICL sequences.

Fig. 2 illustrates the performance as pretraining progresses. On the PT sequences (solid lines), we observe a stage-wise learning process: the model first learns to predict the attribute type, and then the exact attribute token. Crucially, however, we find that high performance on PT sequences does not transfer to ICL sequences, with the model achieving near-zero accuracy on ICL sequences. This shows that despite learning the factual associations, the model relies on *explicit* grammar statistics for retrieval and cannot *implicitly* infer the attribute type from in-context examples alone.

Finding 1: Pretraining on PT sequences does *not* suffice for good ICL performance.

3 FINETUNING EXPERIMENTS

Since pretraining on PT sequences alone does not suffice for good ICL performance, we investigate whether finetuning can bridge this gap. Specifically, we ask: can a model originally trained to rely on explicit grammar-based cues be adapted to perform implicit inference from in-context examples? To answer this, we finetune the pretrained model using the standard next-token prediction objective on a new data distribution (detailed below). Crucially, to test for generalization, we finetune on only a subset of subjects, reserving the remaining subjects as a held-out set. We then evaluate the model on ICL sequences (as defined in Appendix B) where the query subject belongs to this held-out set. We set the number of demonstrations $N_{ft} = 16$, and use 50% of the total subjects for finetuning.

FT on ICL using a subset of subjects.

We first consider the most direct approach: finetuning the pretrained model directly on ICL sequences. This serves as an upper bound, since the finetuning and evaluation formats are identical, *i.e.*, there is no distribution shift in prompt format. Fig. 3 shows performance on ICL sequences with seen and held-out query subjects. As finetuning progresses, the model successfully learns to perform contextual recall, generalizing to held-out subjects.

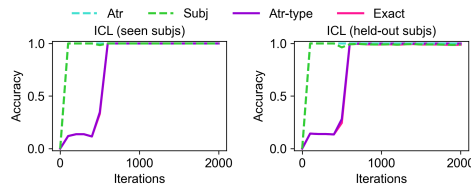


Figure 3: Finetuning the model pretrained on PT sequences, using ICL sequences with a subset of subjects, leads to good performance on ICL sequences with held-out query subjects.

Finding 2: Finetuning on ICL sequences with a subset of subjects leads to generalization on ICL sequences with held-out subjects.

Mechanism of Transfer. Why does finetuning on a subset of subjects enable generalization to the rest? We posit that pretraining and finetuning play distinct but complementary roles. During pretraining, the model *acquires factual knowledge*: given a PT subsequence containing s_i , the model learns to decode the attribute type ℓ from the grammar statistics and predict the corresponding attribute a_i^ℓ . During finetuning, the model does not learn new subject-attribute associations; rather, it learns a new *access mechanism*: inferring the attribute type implicitly from the attribute tokens in the context, rather than from explicit grammar cues. Because there is a shared structure between every s_i and its type- ℓ attribute a_i^ℓ (across $i \in [M]$), the model can be finetuned on a subset of subjects to learn this implicit inference mechanism and generalize to held-out subjects.

FT on ICL+Grm using a subset of subjects.

We next ask whether finetuning on ICL sequences is necessary to induce implicit attribute-type inference, or whether other distributions can achieve the same effect. To investigate this, we finetune the pretrained model on ICL+Grm sequences with short, variable grammar length S_{ft} . Specifically, we use a randomly sampled $S_{ft} \in \{1, \dots, 5\}$ to generate each ICL+Grm sequence. Due to the short grammar length, the sequence statistics are insufficient to reliably encode the attribute type, encouraging the model to instead infer it implicitly from the attribute tokens in the context. In this sense, these sequences serve the same purpose as ICL sequences.

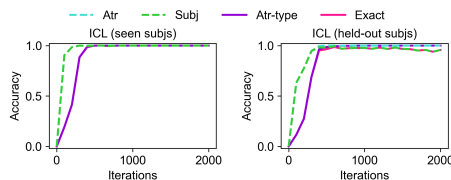


Figure 4: Finetuning the model pretrained on PT sequences, using ICL+Grm sequences with short, variable grammar length with a subset of subjects, leads to good performance on ICL sequences with held-out query subjects.

However, there remains a format distribution shift: during finetuning, attributes always follow a grammar token, whereas in evaluation ICL sequences, attributes directly follow subjects. The model must therefore learn to bridge this format gap at test time. Using variable grammar lengths prevents the model from overfitting to a fixed positional offset, encouraging it to rely on the context to infer the position of the next attribute token. Fig. 4 shows that the model successfully generalizes to ICL sequences for both seen and held-out subjects, as it is finetuned on ICL+Grm sequences.

Finding 3: Finetuning on ICL+Grm sequences with short, variable grammar length using a subset of subjects leads to out-of-distribution generalization on ICL sequences with held-out subjects.

4 REPRESENTATIONAL ANALYSIS

We now analyze the model’s internal representations to gain insight into the mechanism underlying contextual recall. Specifically, we investigate whether finetuning on ICL+Grm sequences induces the formation of a low-dimensional representation that encodes the shared attribute type ℓ from the in-context examples.

Let $X_{\ell,t} := [s_{i_0}, a_{i_0}^\ell, [\text{sep}], s_{i_1}, a_{i_1}^\ell, [\text{sep}], \dots, s_{i_{t+1}}]$ denote an ICL sequence for an attribute type ℓ , with $t \in [N]$ demonstrations. For each $\ell \in [L]$, we sample K such sequences, denoted X_ℓ^k for $k \in [K]$. Let $f_j(\cdot)$ denote the model’s representation at layer j at the last token position. For fixed j and t , we measure the cosine similarity for inter- and intra-task representations, averaged across the K contexts for each pair $\ell, \ell' \in [L]$. Additionally, we quantify the representation clustering strength, $\bar{S}_j^t \in [-1, 1]$, in terms of a clustering metric using $1 - \cos(\cdot, \cdot)$ as the distance and attribute-type ℓ as the cluster label (see App. D for details). A high \bar{S}^t indicates that representations of sequences with shared attribute type are tightly clustered and well-separated from those of different attribute types.

Fig. 5 shows the model’s performance on ICL sequences with held-out subjects, alongside the representation clustering strength \bar{S}^t computed from layer-2 attention representations, as the number of demonstrations t increases. (Results for other layers are included in App. D.) We consider two settings: $\text{Div} \approx 0.2$ (top) and $\text{Div} \approx 0.5$ (bottom). In both cases, accuracy and clustering strength initially improve as t is increased, and eventually saturate. This is also corroborated by the inset figures, which visualize the averaged cosine similarity for inter- and intra-task representations $\bar{C}^t(\ell, \ell')$ across attribute types $\ell, \ell' \in [L]$, after $t \in \{0, 2, 10\}$ demonstrations. This confirms that the finetuned model aggregates information from multiple examples to form a stable representation of the attribute type.

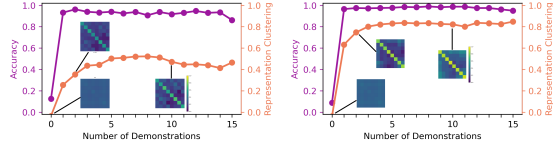


Figure 5: Comparison of accuracy on ICL sequences with held-out subjects, and layer-2 attention representation clustering strength (see text for details) for a finetuned model, as the number of demonstrations is increased, using $\text{Div} \approx 0.2$ (top) and $\text{Div} \approx 0.5$ (bottom). Inset figures visualize the averaged cosine similarity for inter- and intra-task representations. In both cases, both performance and clustering strength improve with number of demonstrations.

Interestingly, while accuracy is comparable in the two settings, higher Div during pretraining leads to stronger representation clustering. This suggests that the separation between the Markov chains (determined by Div) that encode attribute type information during pretraining, is reflected in the task vector separability after finetuning.

Finding 4: The finetuned model encodes attribute type information from in-context examples in layer-2 attn. representations. Clustering strength increases with both the number of in-context examples and the separation between the attribute-specific Markov chains during pretraining.

5 CONCLUSION

We studied contextual recall, a form of ICL that requires models pretrained to acquire knowledge about various subjects and associated attributes, to recall a specific attribute for a query subject by implicitly inferring the relevant attribute type based on in-context examples. Our results give insights into the complementary roles of pretraining and finetuning in enabling in-context reasoning involving learned knowledge. Important directions for future work include characterizing the finetuning dynamics in the minimal setting we study in Appendix C, and investigating how incorporating new knowledge via further finetuning might impact the learned contextual recall capability.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

REFERENCES

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=LziniAXEI9>.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? Investigations with linear models. In *Int. Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=0g0X4H8yN4I>.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- Tina Behnia, Puneesh Deora, and Christos Thrampoulidis. Facts in stats: Impacts of pretraining diversity on language model generalization, 2025. URL <https://arxiv.org/abs/2510.16096>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D 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 Ziegler, Jeffrey Wu, Clemens Winter, Chris 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 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.
- Liam Carroll, Jesse Hoogland, Matthew Farrugia-Roberts, and Daniel Murfet. Dynamics of transient structure in in-context linear regression transformers, 2025. URL <https://arxiv.org/abs/2501.17745>.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality, 2024a. URL <https://arxiv.org/abs/2402.19442>.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Unveiling induction heads: Provable training dynamics and feature learning in transformers. *arXiv preprint arXiv:2409.10559*, 2024b.
- Puneesh Deora, Bhavya Vasudeva, Tina Behnia, and Christos Thrampoulidis. In-context occam’s razor: How transformers prefer simpler hypotheses on the fly. In *Second Conference on Language Modeling*, 2025. URL <https://openreview.net/forum?id=ZSMnX3LBva>.
- Ezra Edelman, Nikolaos Tsilivis, Benjamin L. Edelman, Eran Malach, and Surbhi Goel. The evolution of statistical induction heads: In-context learning markov chains. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=qaRT6QTIqJ>.
- Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn to achieve second-order convergence rates for in-context linear regression. *Advances in Neural Information Processing Systems*, 37:98675–98716, 2024.
- Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=flNZJ2eOet>.
- Andrej Karpathy. `mingpt`. <https://github.com/karpathy/minGPT?tab=readme-ov-file> [Accessed: March 4, 2024].
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.

270 Eshaan Nichani, Jason D Lee, and Alberto Bietti. Understanding factual recall in transformers via
271 associative memories. *arXiv preprint arXiv:2412.06538*, 2024.
272

273 Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What in-context learning “learns” in-context:
274 Disentangling task recognition and task learning. In Anna Rogers, Jordan Boyd-Graber, and Naoaki
275 Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 8298–
276 8319, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/
277 v1/2023.findings-acl.527. URL [https://aclanthology.org/2023.findings-acl.
278 527/](https://aclanthology.org/2023.findings-acl.527/).

279 Core Francisco Park, Ekdeep Singh Lubana, and Hidenori Tanaka. Algorithmic phases of in-context
280 learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
281 <https://openreview.net/forum?id=XgH1wfHSX8>.

282 Nived Rajaraman, Marco Bondaschi, Ashok Vardhan Makkuva, Kannan Ramchandran, and Michael
283 Gastpar. Transformers on markov data: Constant depth suffices. In *The Thirty-eighth Annual
284 Conference on Neural Information Processing Systems*, 2024. URL [https://openreview.
285 net/forum?id=5uG9tp3v2q](https://openreview.net/forum?id=5uG9tp3v2q).

286 Allan Raventos, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the
287 emergence of non-bayesian in-context learning for regression. In *Thirty-seventh Conference on
288 Neural Information Processing Systems*, 2023. URL [https://openreview.net/forum?
289 id=BtAz4a5xDg](https://openreview.net/forum?id=BtAz4a5xDg).

290 Aaditya K. Singh, Stephanie C. Y. Chan, Ted Moskowitz, Erin Grant, Andrew M. Saxe, and Felix
291 Hill. The transient nature of emergent in-context learning in transformers, 2023. URL [https:
292 //arxiv.org/abs/2311.08360](https://arxiv.org/abs/2311.08360).

293 Aaditya K. Singh, Ted Moskowitz, Sara Dragutinovic, Felix Hill, Stephanie C. Y. Chan, and Andrew M.
294 Saxe. Strategy coepetition explains the emergence and transience of in-context learning, 2025.
295 URL <https://arxiv.org/abs/2503.05631>.
296

297 Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
298 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent.
299 *arXiv preprint arXiv:2212.07677*, 2022.
300

301 Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett.
302 How many pretraining tasks are needed for in-context learning of linear regression? *arXiv preprint
303 arXiv:2310.08391*, 2023.
304

305 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
306 learning as implicit bayesian inference. In *International Conference on Learning Representations*,
307 2022. URL <https://openreview.net/forum?id=RdJVFCHjUMI>.

308 Liu Yang, Ziqian Lin, Kangwook Lee, Dimitris Papailiopoulos, and Robert Nowak. Task vectors in
309 in-context learning: Emergence, formation, and benefit, 2025. URL [https://arxiv.org/
310 abs/2501.09240](https://arxiv.org/abs/2501.09240).

311 Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-context,
312 2023.
313

314 Ruiqi Zhang, Jingfeng Wu, and Peter Bartlett. In-context learning of a linear transformer block:
315 Benefits of the mlp component and one-step gd initialization. *Advances in Neural Information
316 Processing Systems*, 37:18310–18361, 2024.

317 Yedi Zhang, Aaditya K. Singh, Peter E. Latham, and Andrew Saxe. Training dynamics of in-context
318 learning in linear attention, 2025. URL <https://arxiv.org/abs/2501.16265>.
319

320 Nicolas Zucchet, Jörg Bornschein, Stephanie Chan, Andrew Lampinen, Razvan Pascanu, and Soham
321 De. How do language models learn facts? dynamics, curricula and hallucinations. *arXiv preprint
322 arXiv:2503.21676*, 2025.
323

324	CONTENTS	
325		
326		
327	1 Introduction	1
328		
329	2 Pretraining on PT Sequences	2
330		
331		
332	3 Finetuning Experiments	3
333		
334	4 Representational Analysis	4
335		
336		
337	5 Conclusion	4
338		
339	A Related Work	7
340		
341	B Data Generation Process	9
342		
343		
344	C Mechanistic Analysis in a Simpler Setting	10
345		
346	D Additional Results and Details of Experimental Settings	13
347		
348	D.1 Effect of data distribution	13
349	D.2 Details of Experimental Settings	14
350	D.3 Details of Representational Analysis	14
351		
352		
353	E Additional Details for Appendix C	16
354	E.1 Constructions for PT and ICL Sequences	16
355	E.2 Experimental Validation	18
356		
357		
358	A RELATED WORK	
359		
360	Prior work has leveraged controlled, synthetic settings to understand both ICL and factual recall in	
361	transformers trained from scratch, which we discuss below.	
362		
363	Factual Recall. Recent studies have utilized controlled synthetic setups to systematically explore	
364	how transformers acquire and recall factual associations. Allen-Zhu & Li (2023) and Zucchet	
365	et al. (2025) employ synthetic biography datasets, where each entry contains information about an	
366	individual or subject in the form of multi-sentence paragraphs. Each sentence associates some fact or	
367	attribute with an individual or subject, via a sentence template that contains information about the	
368	relation or attribute type, such as a person’s birthplace, in the linguistic structure. They consider a	
369	fixed set of templates for each attribute type. Zucchet et al. (2025) evaluate factual knowledge (<i>i.e.</i> ,	
370	correct attribute token prediction) using similar biography entries that are generated with templates	
371	from a held-out set, and identify a stage-wise learning dynamic where the model first learns to predict	
372	some attribute token, and then the correct attribute. We also observe similar stage-wise learning	
373	in our setup during pretraining in Section 2. In Allen-Zhu & Li (2023), the model is evaluated on	
374	question-answering (QA) formats, which constitutes a distribution shift from the biography-style	
375	format. Similar to our results in Section 3, they find that while pretraining is insufficient, finetuning	
376	(on QA) with a subset of subjects enables generalization on held-out subjects. Crucially, however,	
377	their QA prompts still contain explicit relation or attribute type information (e.g., "What is [subject]’s	
	city of birth?"), whereas our contextual recall task constitutes a more substantial distribution shift as	
	the context does not contain explicit relation cues. Investigating whether we see similar results as in	

378 Section 3 for our synthetic setup with synthetic biography data used in these studies would be an
379 interesting direction for future work.

380 A key differentiator across these frameworks is how "relation" information (the attribute type) is
381 represented. Nichani et al. (2024) adopts an abstracted triplet-based approach where the relation
382 is a single, dedicated token. In contrast, Behnia et al. (2025) omits the relation token entirely,
383 using a Markov chain-based grammar setup to study how sequence statistics affect generalization;
384 importantly, in their work, the facts are independent of the grammar. Our synthetic framework bridges
385 these approaches: similar to Behnia et al. (2025), we model templates using Markov chains, but in
386 contrast to their work, and similar to the other prior studies, we retain the relation information by
387 associating specific Markov chains with specific attribute types.

388 For mechanistic analysis, we adopt an abstracted setup similar to Nichani et al. (2024) and compare
389 our attention-only construction for factual recall (Proposition 1) with theirs. In both constructions,
390 the relation heads perform a similar role: they attend to the (most recent) relation token in the context
391 to output the sum of all attributes of the relevant type. However, the subject heads behave differently.
392 While the construction in Nichani et al. (2024) uses subject heads to boost the attributes associated
393 with the subject present in the context, our construction utilizes subject heads to suppress all attributes
394 not associated with the subject. We also present experimental validation for our construction.

395 Notably, in prior works focusing on factual recall, the context the context contains explicit information
396 (e.g., specific words, or linguistic patterns) required to recall specific learned facts, and does not
397 require the model to perform any implicit inference (e.g., what fact is relevant based on the in-context
398 examples) that is characteristic of ICL. In contrast, our work investigates contextual recall, a more
399 challenging task compared to factual recall, that requires the model to infer the attribute type implicitly
400 from given in-context examples.

401
402
403
404
405 **In-Context Learning.** Several recent studies have leveraged controlled synthetic environments to
406 analyze how Transformers learn in-context when trained from scratch. A common approach involves
407 training models on well-defined function classes, most notably linear regression (Garg et al., 2022;
408 Raventos et al., 2023; Akyürek et al., 2023; von Oswald et al., 2022; Ahn et al., 2023; Zhang et al.,
409 2025; Wu et al., 2023; Yang et al., 2025) and Markov chains (Park et al., 2025; Edelman et al.,
410 2024; Rajaraman et al., 2024; Deora et al., 2025). These works often explore whether Transformers
411 implement specific algorithms or functionalities—such as gradient descent (von Oswald et al., 2022;
412 Ahn et al., 2023) or higher-order-algorithms (Fu et al., 2024) for linear regression, and induction
413 heads for Markov chains (Edelman et al., 2024; Rajaraman et al., 2024; Chen et al., 2024b). Further
414 research (Raventos et al., 2023; Park et al., 2025) has investigated task diversity, comparing *task*
415 *retrieval* and *task learning* modes of ICL (Pan et al., 2023). Additionally, research has explored
416 transient dynamics to understand how these two modes evolve over the course of training (Singh
417 et al., 2023; Carroll et al., 2025; Singh et al., 2025). Finally, a growing body of work also explores the
418 training dynamics of in-context learning by examining the optimization dynamics of linear regression
419 (Zhang et al., 2025; 2023; 2024) in both one-layer linear attention and softmax attention models
420 (Chen et al., 2024a).

421 A defining characteristic of the aforementioned studies is that the training and inference formats
422 are identical; the model is evaluated on the same sequence structures it encountered during training.
423 In contrast, our work on contextual recall introduces a significant prompt-distribution shift. While
424 pretraining focuses on instilling factual knowledge in the model via explicit grammar-based cues,
425 the ICL evaluation requires the model to transition a novel format necessitating inference of the
426 relevant attribute type implicitly from the in-context demonstrations. Therefore, unlike standard ICL
427 setups that focus on function induction, our framework requires the model to bridge the gap between
428 structured knowledge acquisition and implicit in-context reasoning.

429 In contrast to our work, Xie et al. (2022) investigate a distribution shift that is primarily composi-
430 tional rather than structural. In their framework, the shift occurs when a model pretrained on long,
431 continuous documents is prompted with a sequence of independent, i.i.d. examples. While they
frame ICL as a statistical process of implicit Bayesian inference, our work provides a mechanistic
perspective discussed in Section 4 and Appendix C.

B DATA GENERATION PROCESS

To study contextual recall in transformers, we define three different types of sequences for pretraining, finetuning, and/or evaluation. The first, `PT`, is designed to instill factual knowledge, while the other two, `ICL` and `ICL+Grm`, are used to evaluate or facilitate how the model leverages such knowledge for in-context reasoning tasks.

We first introduce some useful notation. Let $\mathcal{S} = \{s_j\}_{j=1}^M$ denote the set of M subjects (e.g., landmarks such as `Parthenon` ($j = 1$), `Colosseum` ($j = 2$)). Let the set of unique attributes be denoted as $\mathcal{U} = \cup_{\ell=1}^L \mathcal{U}_\ell$, where $\mathcal{U}_\ell = \{u_i^\ell\}_{i=1}^{M_\ell}$, and ℓ indexes the attribute type (e.g., $\ell = 1$ for “country”, $\ell = 2$ for “year built”), and M_ℓ denotes the number of unique values for index ℓ (e.g., `{Greece, Italy, Canada, ...}` for “country”). Next, let $\mathcal{A} = \cup_{j=1}^M \{a_j^1, \dots, a_j^L\}$ denote the set of attributes assigned to the subjects, where a_j^ℓ is the type- ℓ attribute of subject j (e.g., $a_2^1 = \text{Italy}$). Each subject s_j has one attribute per type, giving L attributes per subject. Let $\mathcal{G} = \{\tilde{g}_1, \dots, \tilde{g}_G\}$ denote a set of G grammar tokens, and `[sep]` denote the separator token. The full vocabulary is $\mathcal{V} = \mathcal{S} \cup \mathcal{U} \cup \mathcal{G} \cup \{\text{[sep]}\}$, with size $V = M + \sum_{\ell=1}^L M_\ell + G + 1$. Let $\text{Unif}(\cdot)$ denote the uniform distribution, and `||` denote concatenation.

With this notation established, we now describe the data generation process. The key building block is a *subject-grammar-attribute subsequence*. For subject s_j , attribute type ℓ , and grammar sequence length S , this subsequence is denoted as $X_j^\ell = [s_j, g_1, \dots, g_S, a_j^\ell]$ (e.g., `[Parthenon, one, of, the, wonders, ..., is in, Greece]`). The grammar sequence $g_{1:S}$ is generated using a first-order Markov chain specific to attribute type ℓ . That is, $p(g_t = \tilde{g} | g_{t-1}, \dots, g_1) = p(g_t = \tilde{g} | g_{t-1})$ for all $\tilde{g} \in \mathcal{G}$, with the first token drawn uniformly at random. For each attribute type $\ell \in [L]$, we sample a row-stochastic transition matrix $\text{TM}_\ell \in \mathbb{R}^{G \times G}$, where each row is drawn independently from a Dirichlet prior with parameter α . Note that since each attribute type has its own Markov chain, the bigram statistics of the grammar sequence $g_{1:S}$ implicitly encode information about the attribute type ℓ .

We now define the three types of sequences used for pretraining, finetuning, and/or evaluation; see Fig. 1 for an illustration.

PT Sequences (Pretraining). To instill factual knowledge in the model, we use `PT` sequences which contain information about a specific subject and its associated attributes (analogous to a short encyclopedia entry for a given landmark). Let N_{tr} denote the number of subsequences (subject-grammar-attribute tuples) in each sequence. To generate a `PT` sequence, we first sample a subject $j \sim \text{Unif}([M])$, then draw N_{tr} attribute types $\ell_{1:N_{\text{tr}}} \sim \text{Unif}([L])^{N_{\text{tr}}}$, and sample each grammar subsequence $g_{1:S}^i$ from the corresponding Markov chain TM_{ℓ_i} . The resulting sequence is $\tilde{X} = X_j^{\ell_1} || [\text{sep}] || X_j^{\ell_2} \dots X_j^{\ell_{N_{\text{tr}}}}$, with sequence length $T = (S + 3)N_{\text{tr}} - 1$. To model generic text that does not convey factual information, we also include grammar-only subsequences: with probability p_G , each subject-grammar-attribute subsequence is replaced by a grammar-only subsequence $g_{1:S+2}$, generated from a separate fixed Markov chain (with transition matrix sampled from a Dirichlet prior). This gives the final sequence X .

ICL Sequences (Finetuning and Evaluation). To test for contextual recall, we use `ICL` sequences that contain subject-attribute pairs with a shared attribute type across different subjects—mirroring the evaluation format from the introduction (e.g., `[Niagara Falls, Canada, [sep], Colosseum, Italy, [sep], Parthenon,]`). Let N denote the number of in-context demonstrations. First, we sample an attribute type $\ell \sim \text{Unif}([L])$, and N subjects $j_{1:N+1} \sim \text{Unif}([M])^{N+1}$. The resulting `ICL` sequence is $X = [s_{j_1}, a_{j_1}^\ell, [\text{sep}], \dots, s_{j_{N+1}}, a_{j_{N+1}}^\ell]$. At evaluation, the model observes N subject-attribute pairs and must predict the attribute $a_{j_{N+1}}^\ell$ for the final subject $s_{j_{N+1}}$, without any grammar cues indicating the attribute type.

ICL+Grm Sequences (Finetuning). We also use `ICL+Grm` sequences for finetuning. As the name suggests, these sequences are similar to `ICL` sequences (subject-attribute pairs sharing a common attribute type across different subjects), but they also contain grammar tokens between the subject and attribute. (As we saw in Section 3, these sequences help instill the implicit inference ability required to succeed at contextual recall on `ICL` sequences.) Let S_{ft} denote the grammar sequence length used in each subsequence, and N_{ft} denote the number of subsequences. We sample attribute type $\ell \sim \text{Unif}([L])$ and subjects $j_{1:N_{\text{ft}}+1} \sim \text{Unif}([M])^{N_{\text{ft}}+1}$, then generate grammar

subsequences $g_{1:S_{\text{ft}}}^i$ from the corresponding Markov chain TM_ℓ . The resulting ICL+Grm sequence is $X = X_{j_1}^\ell ||[\text{sep}]|| X_{j_2}^\ell \dots X_{j_{N_{\text{tr}}}}^\ell$. Note that ICL sequences are a special case of ICL+Grm with $S_{\text{ft}} = 0$.

C MECHANISTIC ANALYSIS IN A SIMPLER SETTING

To gain mechanistic insight into how pretraining and finetuning contribute to contextual recall, we study an analytically tractable setting that preserves the qualitative behavior observed in Findings 1-2 of Section 2. Specifically, we simplify in two ways: i) we encode attribute type information in explicit “relation” tokens rather than grammar sequence statistics, and ii) we use a one-layer attention-only model. Below, we first describe our setup, then present constructions for accurately predicting the final attribute token on both the PT and ICL sequences, and then present experimental validation for the constructions.

Data Setting. Let $\mathcal{R} := \{r_1, \dots, r_L\}$ denote a set of relation tokens, one per attribute type. The PT sequences are generated in a similar manner to Appendix B, with a modification to the grammar subsequences. For each subsequence $g_{1:S}$, we first sample a position $i \sim \text{Unif}([S])$ and a relation token $r \sim \text{Unif}(\mathcal{R})$, and set $g_i = r$. Then, we sample the remaining entries independently as $g_j \sim \text{Unif}(\mathcal{G})$ for all $j \neq i$. This setup is similar to the synthetic setup used to study factual recall in Nichani et al. (2024). The ICL sequences are generated as in Appendix B.

Model. We consider a single-layer attention-only model with fixed relative positional encoding added to the key inputs. The output at position t is defined as follows:

$$\begin{aligned} g_h^t(\mathbf{X}) &= \mathbf{X}_{1:t}^\top \varphi((\mathbf{X}_{1:t} + \mathbf{P}_{\text{T}-t+1:\text{T}})^\top \mathbf{W}_{KQ}^h \mathbf{x}_t) \\ f_\Theta^t(\mathbf{X}) &= \sum_{h=1}^H \underbrace{\mathbf{W}_{OV}^h}_{f_h^t(\mathbf{X})} g_h^t(\mathbf{X}). \end{aligned} \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{\text{T} \times d}$ denotes the input sequence, $\mathbf{x}_t \in \mathbb{R}^d$ denotes the t^{th} token, $\mathbf{P} = [\mathbf{h}(p_{-\text{T}+1}), \dots, \mathbf{h}(p_0)]$ denotes the positional encodings, $\varphi(\cdot)$ denotes the softmax, and for convenience, we define $\mathbf{W}_{OV}^h = (\mathbf{W}_O^h)^\top \mathbf{W}_V^h$, $\mathbf{W}_{KQ}^h = (\mathbf{W}_K^h)^\top \mathbf{W}_Q^h$, for $h \in [H]$ using output, value, key, query weight matrices $\mathbf{W}_O^h, \mathbf{W}_V^h, \mathbf{W}_K^h, \mathbf{W}_Q^h \in \mathbb{R}^{d_h \times d}$. We use the shorthand $f(\cdot) = f_\Theta^{\text{T}}(\cdot)$ and $g_h(\cdot) = g_h^{\text{T}}(\cdot)$ to denote the outputs at the last token position. The model prediction is

$$v^* = \arg \max_{v \in \mathcal{V}} \mathbf{h}(v)^\top f(\mathbf{X}),$$

where $\mathbf{h}(v) \in \mathbb{R}^d$ denotes the embedding for token v . We use fixed embeddings and tied unembeddings. We consider one-hot embeddings and positional encodings and set $d = V + \text{T}$, so that these subspaces are orthogonal.

Construction for PT Sequences. We first investigate whether an attention-only model is expressive enough to perform factual recall on PT sequences of the form

$$\begin{aligned} \mathbf{X} &= \mathbf{X}_j^{\ell_1} ||[\mathbf{h}(a_j^{\ell_1})]|| \mathbf{X}_j^{\ell_2} \dots \mathbf{X}_j^{\ell_{N_{\text{tr}}}}, \text{ where} \\ \mathbf{X}_j^{\ell_1} &= [\mathbf{h}(s_j), \mathbf{h}(g_1), \dots, \mathbf{h}(r_{\ell_1}), \dots, \mathbf{h}(g_S), \mathbf{h}([\text{sep}])]. \end{aligned} \quad (2)$$

Here, $\mathbf{x}_{\text{T}} = \mathbf{h}([\text{sep}])$, and the correct last-token prediction is $\mathbf{h}(a_j^{\ell_{N_{\text{tr}}}})$. We first show that there exists an attention-only model capable of perfectly predicting the next token for such sequences. For simplicity, we focus on last-token prediction here; the construction extends to predictions at any position (see App. E.1).

Proposition 1 (Informal). *Consider the input \mathbf{X} in Eq. (2). There exists a single-layer attention-only model such that when $\|\mathbf{W}_{KQ}^h\| \rightarrow \infty$, correctly predicts the last token—returning the attribute $a_j^{\ell_{N_{\text{tr}}}}$ corresponding to the sequence’s subject s_j and attribute type $\ell_{N_{\text{tr}}}$.*

Proof Sketch. We present a construction with a 3-head model. At a high level, two heads, which we call the **relation** and **subject** heads, are responsible for the prediction at the target position (following $\mathbf{h}([\text{sep}])$), and the third **grammar** head for other cases. For simplicity, we present the construction

here assuming that $f_{\text{grm}}(\mathbf{X}) = 0$, since the outputs from this head do not affect the conclusions in this case, as shown in the full proof in Appendix E.1.

First, the **relation head** attends to the most recent relation token $r_{\ell_{N_r}}$ and maps it to the sum of all attributes of type ℓ_{N_r} . Specifically,

$$\mathbf{W}_{KQ}^{\text{rel}} = \beta \left(\sum_{\ell} \mathbf{h}(r_{\ell}) + \mathbf{p} \right) \mathbf{h}([\text{sep}])^{\top}, \quad \mathbf{W}_{OV}^{\text{rel}} = \sum_{\ell} \sum_j \mathbf{h}(u_j^{\ell}) \mathbf{h}(r_{\ell})^{\top},$$

where $\mathbf{p} := \sum_{i=1}^{S+2} \mathbf{h}(p_{-i})$. With $\beta \rightarrow \infty$, the head’s outputs become $g_{\text{rel}}(\mathbf{X}) = \mathbf{h}(r_{\ell_{N_r}})$, i.e., the most recent relation token in the sequence, and $f_{\text{rel}}(\mathbf{X}) = \sum_j \mathbf{h}(u_j^{\ell_{N_r}})$, i.e., all attributes of type ℓ_{N_r} .

On the other hand, the **subject head** attends to the subject token $s_{\bar{j}}$ and filters out irrelevant attributes (those not associated with $s_{\bar{j}}$). Specifically,

$$\mathbf{W}_{KQ}^{\text{subj}} = \beta \left(\sum_j \mathbf{h}(s_j) \right) \mathbf{h}([\text{sep}])^{\top}, \quad \mathbf{W}_{OV}^{\text{subj}} = - \sum_j \left(\sum_{j' \neq j, \ell} \mathbf{h}(a_{j'}^{\ell}) \right) \mathbf{h}(s_j)^{\top}.$$

Then, with $\beta \rightarrow \infty$, this head outputs $g_{\text{subj}}(\mathbf{X}) = \mathbf{h}(s_{\bar{j}})$, the subject token in the sequence, and $f_{\text{subj}}(\mathbf{X}) = - \left(\sum_u \mathbf{h}(u) - \sum_{\ell} \mathbf{h}(a_{\bar{j}}^{\ell}) \right)$, i.e., the negative of all attributes that are not associated with subject $s_{\bar{j}}$.

Combining these outputs, the model isolates the specific attribute $a_{\bar{j}}^{\ell_{N_r}}$, i.e., the attribute reinforced by both heads: $f(\mathbf{X}) = \sum_j \mathbf{h}(u_j^{\ell_{N_r}}) + \sum_{\ell} \mathbf{h}(a_{\bar{j}}^{\ell}) - \sum_u \mathbf{h}(u)$, yielding the correct prediction $v^* = a_{\bar{j}}^{\ell_{N_r}}$.

Construction for ICL Sequences. Next, consider ICL sequences of the form

$$\mathbf{X} = [\mathbf{h}(s_{j_1}), \mathbf{h}([\text{sep}]), \mathbf{h}(a_{j_1}^{\bar{\ell}}), \dots, \mathbf{h}(s_{j_{N_r+1}}), \mathbf{h}([\text{sep}])], \quad (3)$$

with correct last token prediction $\mathbf{h}(a_{j_{N_r+1}}^{\bar{\ell}})$. We show that there exists an attention-only model that perfectly predicts the last token for such sequences.

Proposition 2 (Informal). *Consider the input \mathbf{X} in Eq. (3). There exists single-layer attention-only model such that when $\|\mathbf{W}_{KQ}^h\| \rightarrow \infty$, it correctly predicts the last token, returning the attribute $a_{j_{N_r+1}}^{\bar{\ell}}$ corresponding to the query subject $s_{j_{N_r+1}}$ and the shared attribute type $\bar{\ell}$.*

Proof Sketch. We adapt the construction from Proposition 1 with minimal changes, but with one crucial modification that reflects the role of finetuning. The **subject** head operates similarly, attending to the query subject and filtering out attributes not associated with it. The key difference lies in the **relation** head: since ICL sequences contain no explicit relation tokens, this head must now *infer the attribute type implicitly* by attending to the attribute tokens in the context and mapping them to all attributes of the same type $\bar{\ell}$.

Specifically, for the **subject** head, we set

$$\mathbf{W}_{KQ}^{\text{subj}} = \beta \left(\sum_j \mathbf{h}(s_j) + \mathbf{h}(p_{-1}) \right) \mathbf{h}([\text{sep}])^{\top}, \quad \mathbf{W}_{OV}^{\text{subj}} = - \sum_j \left(\sum_{j' \neq j, \ell} \mathbf{h}(a_{j'}^{\ell}) \right) \mathbf{h}(s_j)^{\top}.$$

Note that as compared to the construction for PT sequences (Proposition 1), $\mathbf{W}_{OV}^{\text{subj}}$, which contains subject-attribute information, is unchanged. The only difference here is that $\mathbf{W}_{KQ}^{\text{subj}}$ now contains $\mathbf{h}(p_{-1})$, so that when $\beta \rightarrow \infty$, $g_{\text{subj}}(\mathbf{X}) = \mathbf{h}(s_{j_{N_r+1}})$, i.e., the query subject, and hence, $f_{\text{subj}}(\mathbf{X}) = - \left(\sum_u \mathbf{h}(u) - \sum_{\ell} \mathbf{h}(a_{j_{N_r+1}}^{\ell}) \right)$.

We now discuss the **relation** head. Since in our experiments, we finetune with a subset of subjects, let $\mathcal{S}' \subset \mathcal{S}$ denote a subset of subjects. Next, for each attribute type ℓ , let $\mathcal{U}'_{\ell} := \cup_{j \in \mathcal{S}'} \{a_j^{\ell}\}$ denote the set of unique attributes seen during finetuning. With this notation established, we set

$$\mathbf{W}_{KQ}^{\text{rel}} = \beta \sum_{\ell} \sum_{u \in \mathcal{U}'_{\ell}} \mathbf{h}(u) \mathbf{h}([\text{sep}])^{\top}, \quad \mathbf{W}_{OV}^{\text{rel}} = \sum_{\ell} \left(\sum_j \mathbf{h}(u_j^{\ell}) \right) \left(\sum_{u \in \mathcal{U}'_{\ell}} \mathbf{h}(u) + \mathbf{h}(r_{\ell}) \right)^{\top}.$$

In contrast to the construction for PT sequences, where both $\mathbf{W}_{KQ}^{\text{rel}}$ and $\mathbf{W}_{OV}^{\text{rel}}$ rely on the relation tokens $\mathbf{h}(r_\ell)$, in this case, they rely on the attributes seen during finetuning, *i.e.*, $\sum_{u \in \mathcal{U}_\ell} \mathbf{h}(u)$. In this case, when $\beta \rightarrow \infty$, $g_{\text{rel}}(\mathbf{X}) = \frac{1}{N_{\text{r}}^{\ell}} \sum_i \mathbf{h}(a_{j_i}^{\ell})$, *i.e.*, the average of the attributes that appear in the context, and $f_{\text{rel}}(\mathbf{X}) = \sum_j \mathbf{h}(u_j^{\ell})$, retrieving the sum of all attributes of the shared type $\bar{\ell}$.

Finally, by combining the outputs from the two heads, the model isolates the specific attribute $a_{j_{N_{\text{r}}+1}}^{\ell}$ that is encouraged by both, *i.e.* $f(\mathbf{X}) = \sum_j \mathbf{h}(u_j^{\ell}) + \sum_{\ell} \mathbf{h}(a_{j_{N_{\text{r}}+1}}^{\ell}) - \sum_u \mathbf{h}(u)$. Then, the final prediction $v^* = a_{j_{N_{\text{r}}+1}}^{\ell}$.

Note that, for the **relation** head to output all attributes of the same type as the *attributes* that appear in the context, for any type ℓ , the set of attributes seen during finetuning (\mathcal{U}_ℓ) can be much smaller than the full set \mathcal{U}_ℓ . This helps explain why finetuning on a subset of subjects can enable generalization on held-out subjects.

Experimental Validation. We present experimental evidence to corroborate our theoretical constructions, using a 1-layer 3-head attention-only model (see App. D for details).

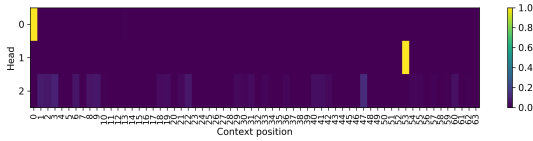


Figure 6: Attention scores for each head across a PT sequence at the end of pretraining. Head 0 attends to the first subject, while head 1 attends to the most recent relation token, as predicted by Proposition 1.

Validation of Pretraining Mechanism. We first train the model using PT sequences (Eq. (2)) with the next-token prediction objective. In Fig. 6, we visualize the attention scores for each head across a PT sequence, and find that the heads specialize into distinct roles consistent with Proposition 1. We find that head 0 attention score concentrates on the first subject token, and it outputs $\mathbf{h}(s_{\bar{j}})$, matching the role of $g_{\text{subj}}(\mathbf{X})$ in our construction, while head 1 attends to the most recent relation token, *i.e.*, it outputs $\mathbf{h}(r_{\ell_{N_{\text{r}}}})$, consistent with $g_{\text{rel}}(\mathbf{X})$. Further, in Fig. 7, we compute the cosine similarity between the head outputs and theoretical outputs of the subject and relation heads specified by our construction. Specifically, we report the cosine similarity with $f_{\text{subj}}(\mathbf{X})$, *i.e.*, negative sum of all attributes not associated with the subject $s_{\bar{j}}$ (left subplot) and $f_{\text{rel}}(\mathbf{X})$, *i.e.*, sum of all attributes of type $\ell_{N_{\text{r}}}$ (right), averaged across several sequences. We find that head 0 output closely matches $f_{\text{subj}}(\mathbf{X})$, while head 1 matches $f_{\text{rel}}(\mathbf{X})$. We designate these heads as subject and relation heads, respectively. Together, these results present experimental validation of our construction for PT sequences.

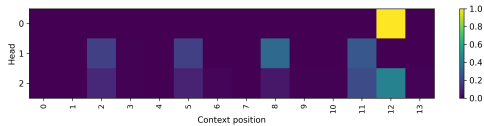


Figure 8: Attention scores for each head across an ICL sequence at the end of finetuning. Head 0 attends to the most recent subject, while head 1 attends to the attribute tokens in the sequence, as predicted by Proposition 2.

Validation of ICL Mechanism. We finetune the model on ICL sequences (Eq. (3)) with the last-token prediction objective, using 50% of the total subjects (Fig. 16 in the App. shows that the model generalizes on held-out subjects). Visualizing the attention scores on an ICL sequence in Fig. 8

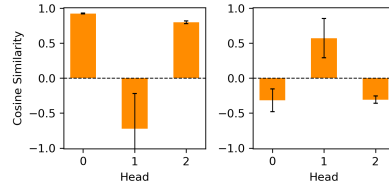


Figure 7: Cosine similarity between the actual outputs of each head and the outputs $f_{\text{subj}}(\mathbf{X})$ (left) and $f_{\text{rel}}(\mathbf{X})$ (right) from our construction in Proposition 1. Head 0 output closely matches $f_{\text{subj}}(\mathbf{X})$, while head 1 matches $f_{\text{rel}}(\mathbf{X})$.

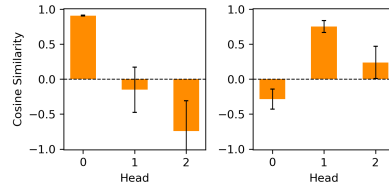


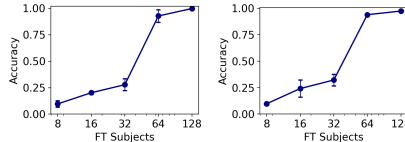
Figure 9: Cosine similarity between the actual outputs of each head and the outputs $f_{\text{subj}}(\mathbf{X})$ (left) and $f_{\text{rel}}(\mathbf{X})$ (right) from our construction in Proposition 2. Head 0 output closely matches $f_{\text{subj}}(\mathbf{X})$, while head 1 matches $f_{\text{rel}}(\mathbf{X})$.

648 reveals that the model repurposes its heads for the new task, consistent with Prop. 2. We find that the
 649 head 0 (subject head) attends to the most recent/query subject token, while head 1 (relation head)
 650 attends to the attribute tokens in the context, *i.e.*, it outputs a combination of attributes of type $\bar{\ell}$.
 651 Further, Fig. 9 confirms that the outputs of these heads closely match the theoretical outputs specified
 652 by our construction for ICL sequences: head 0 (subject head) outputs negative sum of all attributes
 653 not associated with the subject $s_{j_{N_R}}$ (left subplot), while head 1 (relation head) outputs the sum of all
 654 attributes of type $\bar{\ell}$ (right).

656 D ADDITIONAL RESULTS AND DETAILS OF EXPERIMENTAL SETTINGS

659 D.1 EFFECT OF DATA DISTRIBUTION

661 Here, we first investigate whether the format
 662 distribution shift inherent in finetuning with
 663 ICL+Grm sequences incurs a cost in terms of sam-
 664 ple efficiency compared to direct finetuning with
 665 ICL sequences. In Fig. 10, we compare the
 666 ICL performance on held-out subjects as we vary
 667 the number of finetuning subjects. Interestingly, we
 668 find that the sample complexity is comparable across
 669 both settings, with performance improving monotonically
 670 as the number of finetuning subjects increases.

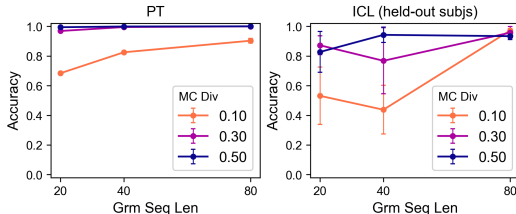


661 Figure 10: Increasing the number of
 662 finetuning subjects with ICL (left) or
 663 ICL+Grm (right) sequences improves the
 664 model’s performance on ICL sequences with
 665 held-out query subjects.

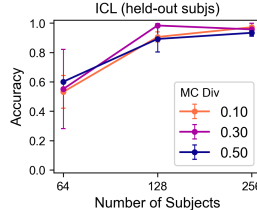
671 Having established (in Section 3) that finetuning on ICL+Grm enables contextual recall, we now
 672 investigate how the properties of the pretraining data influence this capability. Unless stated otherwise,
 673 we use the same parameters as in Sections 2 and 3 and use ICL+Grm sequences with short, variable
 674 grammar length for finetuning.

675 Keeping other parameters fixed, we first consider the effect of varying the sequence length S
 676 and the diversity between the Markov chains for different attribute types, quantified as $\text{Div} :=$
 677 $\min_{\ell, \ell'} \|\text{TM}_\ell - \text{TM}_{\ell'}\|_1$, where $\|\cdot\|_1$ denotes the ℓ_1 -norm. Intuitively, a higher Div for a fixed
 678 sequence length S and/or the Markov chain diversity Div used while pretraining, improves the
 679 model’s performance both on the PT sequences, as well as on ICL sequences with held-out subjects
 680 after finetuning on ICL+Grm sequences with short, variable grammar length.

682 Next, in Fig. 12, we probe the effect of increasing the number of subjects M used for pretraining.
 683 Increasing M improves the performance on ICL sequences with held-out subjects after fine-tuning.
 684 Here, we use grammar sequence length $S = 80$, and consistent with Fig. 11, observe that using a
 685 sufficiently large S for pretraining leads to comparable the final ICL performance across different
 686 levels of Div . We summarize these results in our next finding.



687 Figure 11: Increasing the sequence length S and/or
 688 Markov chain diversity Div used while pretraining
 689 improves the model’s performance on PT sequences
 690 as well as on ICL sequences with held-out subjects
 691 after finetuning on ICL+Grm sequences with short,
 692 variable grammar length.



693 Figure 12: Increasing the number of subjects
 694 M used while pretraining improves the
 695 model’s performance on ICL sequences
 696 with held-out subjects after finetuning on
 697 ICL+Grm sequences with short, variable
 698 grammar length.

Finding 5: Increasing the number of subjects, the grammar sequence length, or the separation between the Markov chains for different attribute types used while pretraining improves the final performance on ICL sequences.

D.2 DETAILS OF EXPERIMENTAL SETTINGS

We use a 2-layer 1-head GPT-2 type decoder-only transformer model (Karpathy) with embedding dimension 256. We train the model with AdamW optimizer (Loshchilov & Hutter, 2019) with learning rate 10^{-4} , weight decay 0.001, and batch size 64, for both pretraining and finetuning.

For the experiments in Section 2, we set $M = 256, L = 8, N_{\text{tr}} = 5, M_1 = 256, M_2 = \dots = M_8 = 32, S = 80$. We set the grammar-only subsequence probability $p_G = 0.2$, and separation between Markov chains $\text{Div} \approx 0.5$. To control for Div , we first randomly generate a large pool of transition matrices, and then use greedy selection to curate a subset of transition matrices that are assigned to each attribute type. We pretrain the model for $20k$ iterations.

For the experiments in Section 3, we use the same pretraining setting as in Section 2, and set $N = N_{\text{ft}} = 16$. We use 128 subjects for finetuning, unless stated otherwise. For experiments with ICL+Grm sequences, $S_{\text{ft}} \sim \text{Unif}(\{1, \dots, 5\})$. In all cases, we finetune for $2k$ iterations. In Fig. 10, we compare the best performance across finetuning iterations.

For Figs. 10 to 12, the results are reported after averaging across two random initialization seeds.

For the experiments in Appendix D.1, the details are as follows. All results are reported at end of pretraining/finetuning. We consider $\text{Div} \approx 0.1, 0.3, 0.5$. In Fig. 11, we use fixed $M = 256$ and for $S = 20$, we use $S_{\text{ft}} \sim \text{Unif}(\{1, \dots, 4\})$, while for $S = 40$ or $S = 80$, we use $S_{\text{ft}} \sim \text{Unif}(\{1, \dots, 5\})$. In Fig. 12, we use fixed $S = 80, S_{\text{ft}} \sim \text{Unif}(\{1, \dots, 5\})$.

The results in Fig. 5 (bottom) are reported with the same setting as Fig. 4. For the top plot, we only change $\text{Div} \approx 0.2$. We use 50 sequences for each attribute type.

The settings for the experiments in Appendix C are as follows. We train an attention-only model with $d_h = 256$ using AdamW optimizer with learning rate 0.001, weight decay 0.001, and batch size 64, for both pretraining ($20k$ iterations) and finetuning ($2k$ iterations). We use a cosine learning rate scheduler for pretraining. We set $M = 256, L = 8, N_{\text{tr}} = 5, M_1 = 256, M_2 = \dots = M_8 = 32, S = 10, N = N_{\text{ft}} = 5$.

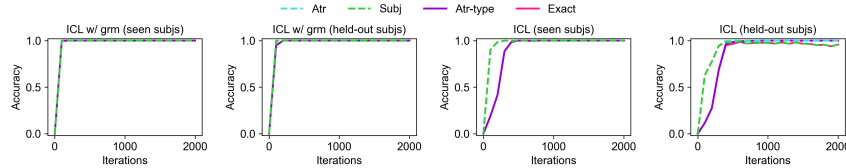


Figure 13: Performance of the model when finetuning with ICL+Grm sequences with a subset of subjects (same setting as Fig. 4), on ICL+Grm sequences with $S_{\text{ft}} = 1$ (left) and ICL sequences (right) with seen or held-out subjects. Finetuning with ICL+Grm sequences enables out-of-distribution generalization on ICL sequences with held-out subjects. Performance improves first on ICL+Grm and later on ICL sequences.

D.3 DETAILS OF REPRESENTATIONAL ANALYSIS

Consider ICL sequences, which are of the form

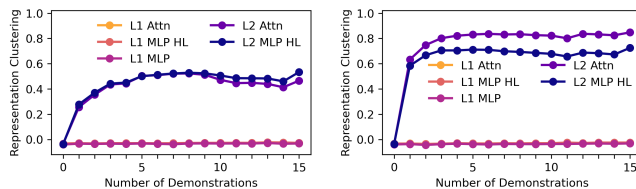
$$X_\ell = [s_{i_0}, a_{i_0}^\ell, [\text{sep}], s_{i_1}, a_{i_1}^\ell, [\text{sep}], \dots, s_{i_{N+1}}],$$

for a fixed attribute type ℓ . We sample K such sequences for each $\ell \in [L]$, denoted by X_ℓ^k . Also, define $X_{\ell,t} := [s_{i_0}, a_{i_0}^\ell, [\text{sep}], s_{i_1}, a_{i_1}^\ell, [\text{sep}], \dots, s_{i_{t+1}}]$, where $t \in [N]$ denotes the number of demonstrations.

Let $f_j(\cdot)$ denote the model’s representation at layer- j . Consider fixed j, t , and for $\ell, \ell' \in [L]$, define

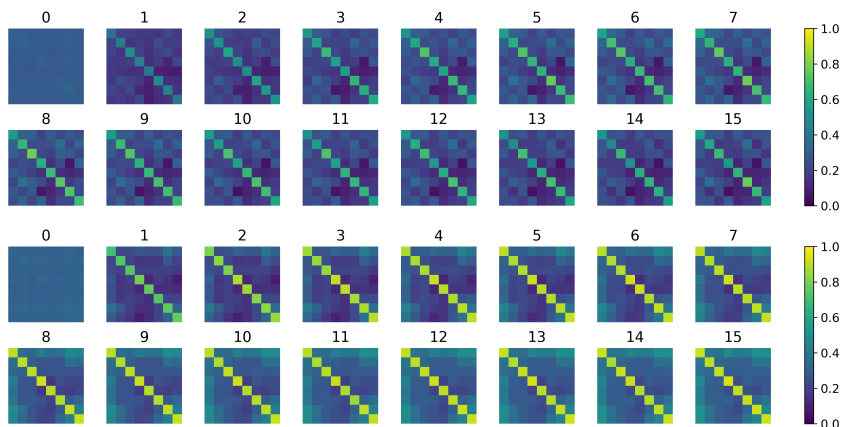
$$\bar{C}_j^t(\ell, \ell') = \frac{1}{K^2} \sum_{k,k'} \cos(f_j(X_{\ell,t}^k), f_j(X_{\ell',t}^{k'})).$$

756
757
758
759
760
761
762



763 Figure 14: Comparison of clustering strength (see Section 4 for details) using representations
764 from different layers of the finetuned 2-layer 1-head transformer with ICL sequences as inputs,
765 as the number of demonstrations is increased (same setting as Fig. 5) with $\text{Div} \approx 0.2$ (left) and
766 $\text{Div} \approx 0.5$ (right). We find that layer-2 attention representations cluster most strongly, while layer-1
767 representations exhibit no clustering based on the attribute type information in the context. Further,
768 using higher Div while pretraining leads to stronger representation clustering after finetuning.

769
770
771
772
773
774
775
776
777
778
779
780
781
782



783 Figure 15: Each subfigure visualizes the cosine similarity for inter- and intra-task representations
784 (from layer-2 attention layer of the finetuned model), $\bar{C}^t(\ell, \ell')$ (see Section 4 for details) across
785 attribute types $\ell, \ell' \in [L]$, averaged over 50 sequences for each attribute type (same setting as Fig. 5;
786 top: $\text{Div} \approx 0.2$, bottom: $\text{Div} \approx 0.5$). We find that as the number of demonstrations is increased
787 (from 0 to 15), the representations of ICL sequences with the same attribute type get clustered
788 together, and the clustering is stronger for higher Div (*i.e.*, more separated attribute-specific Markov
789 chains).

790
791
792
793
794
795

If the model perfectly disentangles attribute types in its representations, then, for some layer j and
number of demonstrations t , we would have $\bar{C}_j^t(\ell, \ell') = \mathbb{1}[\ell = \ell']$. Next, the clustering strength is
quantified as follows. For fixed j, t , we subsume these and define $v_\ell^k := f_j^t(X_\ell^k)$. Each representation
is assigned cluster label based on the attribute type ℓ . Define the cosine distance $d(v, w) :=$
 $1 - \cos(v, w)$. For each point v_ℓ^k , define the intra-cluster dissimilarity

796
797
798

$$a(k, \ell) = \frac{1}{K-1} \sum_{k' \neq k} d(v_\ell^k, v_\ell^{k'}),$$

and the nearest other-cluster dissimilarity

800
801
802

$$b(k, \ell) = \min_{\ell' \neq \ell} \left(\frac{1}{K} \sum_{k'} d(v_\ell^k, v_{\ell'}^{k'}) \right).$$

The silhouette value for a sample is

804
805
806

$$s(k, \ell) = \frac{b(k, \ell) - a(k, \ell)}{\max\{a(k, \ell), b(k, \ell)\}}.$$

Finally, the silhouette score for layer j after t demonstrations are seen is

807
808
809

$$\bar{S}_j^t = \frac{1}{KL} \sum_{\ell} \sum_k s(k, \ell).$$

E ADDITIONAL DETAILS FOR APPENDIX C

E.1 CONSTRUCTIONS FOR PT AND ICL SEQUENCES

Consider PT sequences of the form

$$\mathbf{X} = [\mathbf{h}(s_{\bar{i}}), \mathbf{h}(g_1), \dots, \mathbf{h}(r_{\ell_1}), \dots, \mathbf{h}([\text{sep}]), \mathbf{h}(a_{\bar{i}}^{\ell_1}), \mathbf{h}(s_{\bar{i}}), \dots, \mathbf{h}([\text{sep}])], \quad (4)$$

where the correct last token prediction is $\mathbf{h}(a_{\bar{i}}^{\ell_{N_{\text{tr}}}})$. The following result shows that there exists an attention-only model that always gives correct predictions at any position.

Proposition 3. *Consider the input \mathbf{X} in Eq. (4). There exists a one-layer attention-only model such that when $\|\mathbf{W}_{KQ}^h\| \rightarrow \infty$, it always gives the correct prediction for any token position.*

Proof. We present a construction for a 3-head model. For simplicity, we present the proof for the last subsequence, i.e., $t \in \{\mathbb{T} - S - 2, \dots, \mathbb{T}\}$, but it can be easily extended to other cases where $t < \mathbb{T} - S - 2$ as well. Hereafter, we assume that $t \geq \mathbb{T} - S - 2$.

At a high level, we use two heads, which we call the **relation** and **subject** heads, to get the output following $\mathbf{h}([\text{sep}])$, i.e., $t = \mathbb{T}$, and the third **grammar** head for other cases $t < \mathbb{T}$. Let $\mathbf{p} := \sum_{i=1}^{S+2} \mathbf{h}(p_{-i})$, $\mathcal{T} := \{\mathbb{T} - S, \dots, \mathbb{T} - 2\}$.

For the **relation** head, we set

$$\mathbf{W}_{KQ}^{\text{rel}} = \beta \left(\sum_{\ell} \mathbf{h}(r_{\ell}) + \mathbf{p} \right) \mathbf{h}([\text{sep}])^{\top}, \quad \mathbf{W}_{OV}^{\text{rel}} = \sum_{\ell} \sum_j \mathbf{h}(u_j^{\ell}) \mathbf{h}(r_{\ell})^{\top}.$$

Then, when $\beta \rightarrow \infty$, the outputs of this head are as follows. When $t = \mathbb{T}$, $g_{\text{rel}}(\mathbf{X}) = \mathbf{h}(r_{\ell_{N_{\text{tr}}}})$, i.e., the most recent relation token in the sequence, and $f_{\text{rel}}(\mathbf{X}) = \sum_{j, \ell_{N_{\text{tr}}}} \mathbf{h}(u_j^{\ell_{N_{\text{tr}}}})$, i.e., all attributes of type $\ell_{N_{\text{tr}}}$.

On the other hand, when $t < \mathbb{T}$, $f_{\text{rel}}^t(\mathbf{X}) = \frac{1}{t} \sum_{j, \ell_{N_{\text{tr}}}} \mathbf{h}(u_j^{\ell_{N_{\text{tr}}}})$.

For the **subject** head, we set

$$\mathbf{W}_{KQ}^{\text{subj}} = \beta \left(\sum_j \mathbf{h}(s_j) \right) \mathbf{h}([\text{sep}])^{\top}, \quad \mathbf{W}_{OV}^{\text{subj}} = - \sum_j \left(\sum_{j' \neq j, \ell} \mathbf{h}(a_{j'}^{\ell}) \right) \mathbf{h}(s_j)^{\top}$$

Then, when $\beta \rightarrow \infty$, the outputs of this head are as follows. When $t = \mathbb{T}$, $g_{\text{subj}}(\mathbf{X}) = \mathbf{h}(s_{\bar{i}})$, the subject token in the sequence, and $f_{\text{subj}}(\mathbf{X}) = - \left(\sum_u \mathbf{h}(u) - \sum_{\ell} \mathbf{h}(a_{\bar{i}}^{\ell}) \right)$, i.e., negative of all attributes that are not associated with subject $s_{\bar{i}}$.

On the other hand, when $t < \mathbb{T}$, $f_{\text{subj}}^t(\mathbf{X}) = -\frac{1}{t} \left(\sum_u \mathbf{h}(u) - \sum_{\ell} \mathbf{h}(a_{\bar{i}}^{\ell}) \right)$.

For the **grammar** head, we set

$$\begin{aligned} \mathbf{W}_{KQ}^{\text{grm}} &= \beta \left(\sum_{\ell, j} \mathbf{h}(u_j^{\ell}) \mathbf{h}(u_j^{\ell})^{\top} + \sum_j \mathbf{h}(s_j) \mathbf{h}(s_j)^{\top} + \sum_{\ell} (\mathbf{h}(r_{\ell}) + \mathbf{h}([\text{sep}]) + \mathbf{p}) \mathbf{h}(r_{\ell})^{\top} \right. \\ &\quad \left. + \left(\sum_{\ell} \mathbf{h}(r_{\ell}) + \mathbf{h}([\text{sep}]) + \mathbf{p} \right) \sum_{\tilde{g}} \mathbf{h}(\tilde{g})^{\top} \right) \\ \mathbf{W}_{OV}^{\text{grm}} &= \sum_j \sum_{\ell} \mathbf{h}(s_j) \mathbf{h}(a_j^{\ell})^{\top} + \left(\sum_{\tilde{g}} \mathbf{h}(\tilde{g}) + \sum_{\ell} \mathbf{h}(r_{\ell}) \right) \left(\sum_j \mathbf{h}(s_j)^{\top} + \mathbf{h}([\text{sep}])^{\top} \right) \\ &\quad + \left(0.5 \sum_{\tilde{g}} \mathbf{h}(\tilde{g}) + \mathbf{h}([\text{sep}]) \right) \sum_{\ell} \mathbf{h}(r_{\ell})^{\top}. \end{aligned}$$

Let $g_h^t = g_h^t(\mathbf{X})$ and similarly $f_h^t := f_h^t(\mathbf{X})$. The outputs in this case, when $\beta \rightarrow \infty$, are as follows:

- $t = \mathbb{T} - S - 2$: $g_{\text{grm}} = \mathbf{h}(a_{\bar{i}}^{\ell_{N_{\text{tr}}-1}})$, $f_{\text{grm}} = \mathbf{h}(s_{\bar{i}})$

- 864 • $t = T - S - 1$: $g_{\text{grm}} = \mathbf{h}(s_{\bar{i}})$, $f_{\text{grm}} = \sum_{v \in \mathcal{G} \cup \mathcal{R}} \mathbf{h}(v)$
865
866 • $t \in \mathcal{T}$: $g_{\text{grm}} = \begin{cases} 0.5(\mathbf{h}([\text{sep}]) + \mathbf{h}(r_{\ell})), & \text{if } \mathbf{x}_{t-S-2:t} \in \mathcal{R} \\ \mathbf{h}([\text{sep}]), & \text{if } \mathbf{x}_{t-S-2:t} \notin \mathcal{R} \end{cases}$,
867
868 $f_{\text{grm}} = \begin{cases} \sum_{v \in \mathcal{G} \cup \mathcal{R}} \mathbf{h}(v), & \text{if } g_{\text{grm}} = \mathbf{h}([\text{sep}]) \\ 0.75 \sum_{v \in \mathcal{G}} \mathbf{h}(v) + 0.5 \sum_{v \in \mathcal{R}} \mathbf{h}(v) + 0.5 \mathbf{h}([\text{sep}]), & \text{if } g_{\text{grm}} = 0.5(\mathbf{h}([\text{sep}]) + \mathbf{h}(r_{\ell_{N_r}})) \end{cases}$
869
870
871 • $t = T - 1$: $g_{\text{grm}} = \mathbf{h}(r_{\ell_{N_r}})$, $f_{\text{grm}} = \mathbf{h}([\text{sep}]) + 0.5 \sum_{v \in \mathcal{G}} \mathbf{h}(\tilde{g})$
872
873 • $t = T$: $g_{\text{grm}} = \frac{1}{T} \sum_t \mathbf{x}_t$, $f_{\text{grm}} = \frac{N_r - 1}{T} \mathbf{h}(s_{\bar{i}}) + \frac{N_r}{T} \left(1.5 \sum_{v \in \mathcal{G}} \mathbf{h}(v) + \sum_{\ell} \mathbf{h}(r_{\ell}) + \mathbf{h}([\text{sep}]) \right)$
874
875

876 Combining the outputs of the individual heads, the final output in different cases is as follows.

$$877 v^* = \begin{cases} a_{\bar{i}}^{\ell_{N_r}}, & \text{if } t = T, \\ [\text{sep}], & \text{if } t = T - 1, \\ v \sim \text{Unif}(\mathcal{G} \cup \mathcal{R}), & \text{if } t \in \mathcal{T}, \nexists j \in \{t - S - 2, \dots, t\}, \mathbf{x}_j \in \mathcal{R} \\ v \sim \text{Unif}(\mathcal{G}), & \text{if } t \in \mathcal{T}, \exists j \in \{t - S - 2, \dots, t\}, \mathbf{x}_j \in \mathcal{R} \\ v \sim \text{Unif}(\mathcal{G} \cup \mathcal{R}), & \text{if } t = T - S - 1, \\ s_{\bar{i}}, & \text{if } t = T - S - 2. \end{cases}$$

885 □

886
887
888 Next, consider ICL sequences of the form

$$889 \mathbf{X} = [\mathbf{h}(s_{j_1}), \mathbf{h}([\text{sep}]), \mathbf{h}(a_{j_1}^{\bar{\ell}}), \dots, \mathbf{h}(s_{j_{N_r+1}}), \mathbf{h}([\text{sep}])], \quad (5)$$

890
891 where the correct last token prediction is $\mathbf{h}(a_{j_{N_r}}^{\bar{\ell}})$. The following result shows that there exists an
892 attention-only model that always gives correct predictions for these sequences.

893 **Proposition 4.** *Consider the input \mathbf{X} in Eq. (5). There exists a one-layer attention-only model such*
894 *that when $\|\mathbf{W}_{KQ}^h\| \rightarrow \infty$, it always gives the correct prediction for the last token position.*

895
896 *Proof.* We present a construction for a 3-head model, with minimal changes to the construction
897 in the proof of Proposition 1. Specifically, the **grammar** head is unchanged. At a high level, the
898 **subject** head attends to the last/query subject and maps it to the negative of the sum of attributes not
899 associated with it, while the **relation** head attends to the attributes in the context, and maps them to
900 all attributes of the same type $\bar{\ell}$.

901 For the **subject** head, we set

$$902 \mathbf{W}_{KQ}^{\text{subj}} = \beta \left(\sum_j \mathbf{h}(s_j) + \mathbf{h}(p_{-1}) \right) \mathbf{h}([\text{sep}])^\top, \quad \mathbf{W}_{OV}^{\text{subj}} = - \sum_j \left(\sum_{j' \neq j, \ell} \mathbf{h}(a_{j'}^{\ell}) \right) \mathbf{h}(s_j)^\top$$

903
904 Then, when $\beta \rightarrow \infty$, $g_{\text{subj}}(\mathbf{X}) = \mathbf{h}(s_{j_{N_r}})$, i.e., the query subject, and $f_{\text{subj}}(\mathbf{X}) = - \left(\sum_u \mathbf{h}(u) - \sum_{\ell} \mathbf{h}(a_{j_{N_r}}^{\ell}) \right)$.

905
906 Let $S' \subset S$ denote a subset of subjects. We define the set of unique attributes for the subjects in S' :
907 for each attribute type ℓ , let $\mathcal{U}_{\ell}^i := \cup_{j \in S'} \{a_j^{\ell}\}$. For the **relation** head, we set

$$908 \mathbf{W}_{KQ}^{\text{rel}} = \beta \sum_{\ell} \sum_{u \in \mathcal{U}_{\ell}^i} \mathbf{h}(u) \mathbf{h}([\text{sep}])^\top, \quad \mathbf{W}_{OV}^{\text{rel}} = \sum_{\ell} \left(\sum_j \mathbf{h}(u_j^{\ell}) \right) \left(\sum_{u \in \mathcal{U}_{\ell}^i} \mathbf{h}(u) + \mathbf{h}(r_{\ell}) \right)^\top.$$

909
910 Then, when $\beta \rightarrow \infty$, $g_{\text{rel}}(\mathbf{X}) = \frac{1}{N_r} \sum_i \mathbf{h}(a_{j_i}^{\bar{\ell}})$, i.e., the average of the attributes that appear in the
911 context, and $f_{\text{rel}}(\mathbf{X}) = \sum_j \mathbf{h}(u_j^{\bar{\ell}})$.

Combining the outputs of the individual heads, the model output is

$$f(\mathbf{X}) = \sum_j \mathbf{h}(u_j^{\bar{\ell}}) - \left(\sum_u \mathbf{h}(u) - \sum_{\ell} \mathbf{h}(a_{j_{N_{\text{fit}}}^{\ell}}) \right) + \frac{1}{T} \sum_{i \in \mathcal{M}''} \mathbf{h}(s_i) + \frac{N_{\text{fit}}}{T} \sum_{v \in \mathcal{G} \cup \mathcal{R}} \mathbf{h}(v),$$

where $\mathcal{M}'' \subseteq [M]$. Then, the final prediction $v^* = a_{j_{N_{\text{fit}}}^{\bar{\ell}}}$.

□

E.2 EXPERIMENTAL VALIDATION

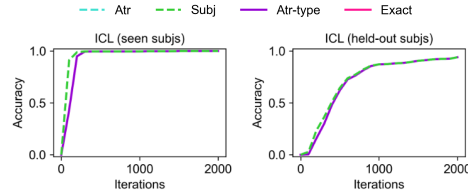


Figure 16: Finetuning the single-layer attention-only model pretrained on PT sequences using ICL sequences with a subset of subjects enables generalization on ICL sequences with held-out query subjects.

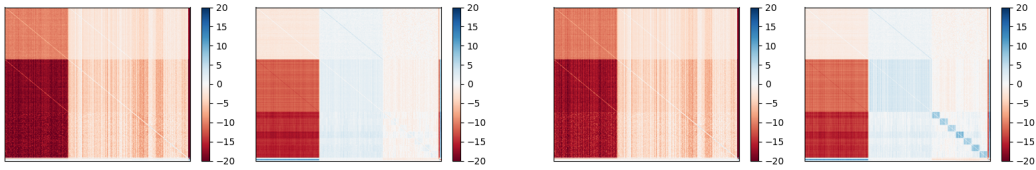


Figure 17: Visualization of the weight matrices $\mathbf{W}_{OV}^{\text{subj}}$ (left) and $\mathbf{W}_{OV}^{\text{rel}}$ (right) after pretraining the single-layer attention-only model on PT sequences.

Figure 18: Visualization of the weight matrices $\mathbf{W}_{OV}^{\text{subj}}$ (left) and $\mathbf{W}_{OV}^{\text{rel}}$ (right) after finetuning the single-layer attention-only model on ICL sequences.