

GENOMIC NEXT-TOKEN PREDICTORS ARE IN-CONTEXT LEARNERS

Anonymous authors

Paper under double-blind review

ABSTRACT

In-context learning (ICL) – the capacity of a model to infer and apply abstract patterns from examples provided within its input – has been extensively studied in large language models trained for next-token prediction on human text. In fact, prior work often attributes this emergent behavior to distinctive statistical properties in *human* language. This raises a fundamental question: can ICL arise *organically* in other sequence domains purely through large-scale predictive training?

To explore this, we turn to genomic sequences, an alternative symbolic domain rich in statistical structure. Specifically, we study the Evo2 genomic model, trained predominantly on next-nucleotide (A/T/C/G) prediction, at a scale comparable to mid-sized LLMs. We develop a controlled experimental framework comprising symbolic reasoning tasks instantiated in both linguistic and genomic forms, enabling direct comparison of ICL across genomic and linguistic models. Our results show that genomic models, like their linguistic counterparts, exhibit log-linear gains in pattern induction as the number of in-context demonstrations increases. To the best of our knowledge, this is the first evidence of organically emergent ICL in genomic sequences, supporting the hypothesis that ICL arises as a consequence of large-scale predictive modeling over rich data. These findings extend emergent meta-learning beyond language, pointing toward a unified, modality-agnostic view of in-context learning.

1 INTRODUCTION

Scaling Large Language Models (LLMs) has revealed an unexpected and powerful capacity: in-context learning (ICL) (Radford et al., 2018; Brown et al., 2020a), the ability to infer and apply abstract patterns purely from examples contained within their input. Almost all prominent evidence of emergent ICL so far (Srivastava et al., 2023, *inter alia*) comes from training on human language (e.g., English). This raises a fundamental question: *is there something inherently special about human language that enables ICL to emerge?* To investigate this, we turn to a radically different substrate: **the genome**. Genomic sequences can be viewed as a form of *natural* language that has evolved through nature. Like human language, they comprise complex symbol sequences that exhibit rich statistical regularities – motifs, repeats, dependencies (Benegas et al., 2025) that could, in principle, support pattern induction within context.

Testing for the existence of ICL requires large-scale models. We use the **Evo2** series (Brix et al., 2025): a large-scale model trained solely on ‘next-nucleotide’ prediction. With training scale comparable to mid-sized LLMs (e.g., Qwen3-14B (Yang et al., 2025)), these genomic models finally make it possible to systematically compare genomic and linguistic representations under large-scale autoregressive training. If ICL arises primarily from scale and predictive compression, then large-scale genomic models such as Evo2 may already exhibit it.

Our experiment’s results (§3) show that large genomic models do, in fact, exhibit ICL. **Both linguistic and genomic models exhibit log-linear improvements** in pattern inference as the number of demonstrations grows. We further validate these findings by demonstrating that Evo2 also exhibits robust few-shot learning on a real-world genomic promoter classification task (§3.1). These findings suggest *ICL is not an artifact that is specific to human language, but a broader consequence of*

LLMs were lightly used during the conception, implementation, and refinement of this paper. We assume full responsibility for all content.

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

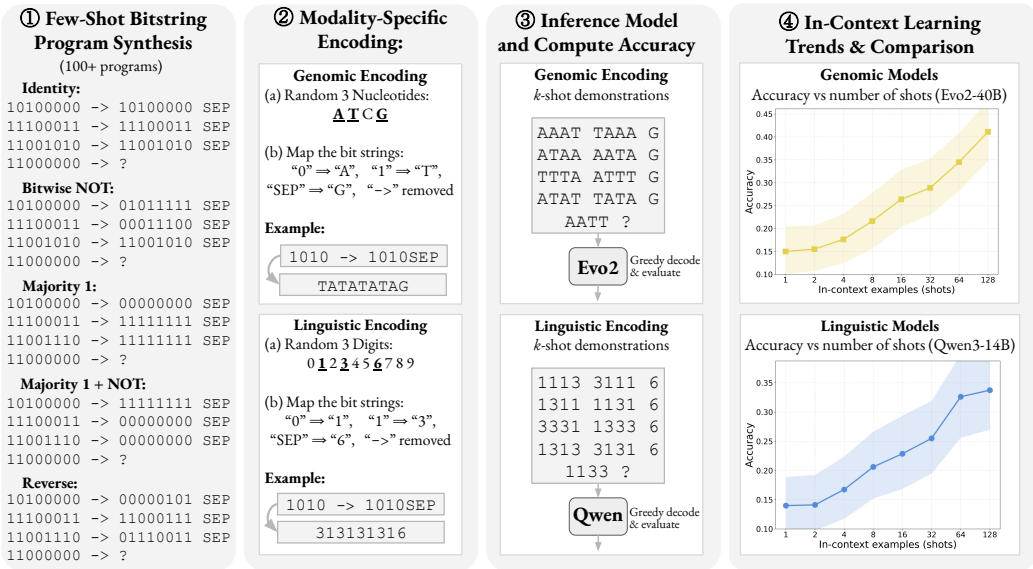


Figure 1: We design parallel symbolic reasoning tasks that allow direct comparison of ICL behavior across modalities. ① Few-shot bitstring program-synthesis tasks require models to infer mappings from examples. ② Each task is rendered in two modality-specific encodings. ③ Both genomic (Evo2) and linguistic (Qwen3) models receive k -shot demonstrations and are greedily decoded to compute exact-match accuracy. ④ Both models show log-linear accuracy gains with more demonstrations.

compression and predictive modeling in pattern-rich sequence spaces. This points toward a unified view of context-adaptive computation that spans different modalities of data.

2 FRAMEWORK FOR CROSS-DOMAIN IN-CONTEXT LEARNING

This section presents our framework for evaluating ICL in linguistic and genomic models. We outline experimental desiderata (§B), setup (§2.1), models (§C), and evaluation protocol (§2.2).

2.1 EXPERIMENTAL SETUP

Notation: We define in-context learning (ICL) as follows. Let S and O be the input and output domains, and let each task be a deterministic latent function $f : S \rightarrow O$. Given an ordered n -shot context $E = ((x_1, f(x_1)), \dots, (x_n, f(x_n)))$ with $x_i \in S$ and a query $x \in S \setminus \{x_1, \dots, x_n\}$, the model predicts $\hat{y} = M(E, x)$. A trial succeeds if $\hat{y} = f(x)$, i.e., $\mathbb{1}[\hat{y} = f(x)]$.

Model families and selection rationale. Details on the model families and our selection criteria are provided in Appendix C.

Task Formulation: in-context program induction over abstract symbols. Given our desiderata (§B), we construct symbolic induction tasks with a small lexicon: the model infers a hidden transformation from a few input–output demonstrations and applies it to a new query. Related setups appear in work on compositional reasoning (Brown et al., 2020b), Raven-style analogical reasoning (Raven et al., 1962; Webb et al., 2023), and ARC-AGI (Chollet, 2019).

Bitstring domain and function space. We use a shared symbolic domain $S = \{0, 1\}^k$ to ensure comparable token granularity across genomic and linguistic models. For genomic models, two nucleotides (e.g., A/C) encode bits while the remaining bases (G/T) serve as separators; for linguistic models, binary digits are typically single tokens (§G; Fig. 1). We set $k = 8$, so $|S| = 2^8 = 256$. Our task set is $F \subseteq \{f : S \rightarrow S\}$ with $|F| = 100$, where each f is either one primitive transformation or a composition of two primitives.

Generation of F . We build F from a library of 30 fundamental primitives (listed in §E.1) covering low-level logical operations and higher-order compositional structure. We construct the final function set, we use GPT-5-Codex to automatically generate 100 unique transformations by composing primitives into valid Python programs. The full list of transformations and corresponding source code is provided in §E.2.

2.2 EVALUATION PROTOCOL

Evaluation metric. For each transformation $f \in F$, sample i.i.d. trials $(E^{(t)}, x^{(t)})$ with $E^{(t)} \sim \text{Unif}(E_n)$ and $x^{(t)} \sim \text{Unif}(S \setminus E^{(t)})$, where $E_n = \{E \subset S : |E| = n\}$. Given $E^{(t)}$ and $x^{(t)}$, the model predicts $\hat{y}^{(t)} = M(E^{(t)}, x^{(t)})$ and a trial is correct if $\hat{y}^{(t)} = f(x^{(t)})$. Our Monte Carlo estimate of n -shot accuracy is

$$\hat{P}(M, n) = \frac{1}{|F|m} \sum_{f \in F} \sum_{t=1}^m \mathbb{1} \left[M(E^{(t)}, x^{(t)}) = f(x^{(t)}) \right]. \quad (1)$$

Model suites and sampling. We use $m = 8$ Monte Carlo trials and shot counts $\mathcal{N} = \{1, 2, 4, 8, 16, 32, 64, 128\}$. Genomic models: $\mathcal{M}_G = \{\text{evo2_1b_base}, \text{evo2_7b}, \text{evo2_40b}\}$; linguistic models: $\mathcal{M}_L = \{\text{Qwen3-0.6B-Base}, \text{Qwen3-1.7B-Base}, \text{Qwen3-4B-Base}, \text{Qwen3-8B-Base}, \text{Qwen3-14B-Base}\}$. For each $(M, n) \in (\mathcal{M}_L \cup \mathcal{M}_G) \times \mathcal{N}$, we report $\hat{P}(M, n)$. Standard errors are estimated by a two-stage cluster bootstrap (resample f , then its trials) with 5000 replicates per (M, n) .

Mode baseline. We define a mode baseline that always predicts the most frequent output observed in the context. We formally define this baseline in §H.

3 EMPIRICAL RESULTS

This section reports empirical findings on few-shot bitstring generalization. §3.1 presents the main accuracy trends with respect to model size and shot count. §D.1 examines sensitivity to task complexity using a BitLoad measure, and §D.2 contrasts the models’ qualitative competencies across individual transformations.

3.1 MAIN RESULTS

Across both model families, accuracy generally rises linearly with respect to $\log(\text{shots})$. A linear regression linking performance to $\log(\text{shots})$ yields highly significant slopes for all models (all $p \leq 10^{-3}$ via one-sided t-test on slope). As Fig. 3 shows, Evo2 models show cleanly monotonic gains, with a pronounced step from 1B to 7B, and near-indistinguishable curves for 7B and 40B; by 128 shots, both 7B and 40B surpass 40% accuracy (Fig. 3a). Qwen3 also trends upward overall but with non-monotonic patches—especially for smaller models in the 4–16 shot band—before resuming clear improvements from 32 to 128 shots; at 128 shots, the 14B model approaches 35% (Fig. 3b). A full table of all accuracies is in §I.

The mode baseline’s performance doesn’t improve with more shots, and lags Qwen3 and Evo2 at high shot counts (Fig. 3d). For Qwen3, advantages over the mode baseline become consistently significant ($p < 0.05$ via one-sided z-test on bootstrapped standard errors) at $n=128$ for all sizes greater than 0.6B. For Evo2, statistically significant advantages emerge slightly earlier: for 1B at $n=64$, the 7B at $n = 32$, and the 40B at $n = 16$. Regardless, all models surpass the naive baseline at high shot counts.

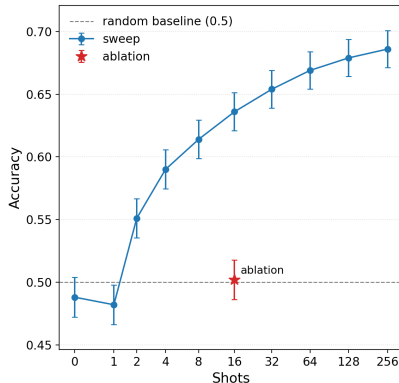


Figure 2: ICL Performance of Evo2 7B on human_nontata_promoters at varying shot counts. Error bars are \pm one standard error. The performance shows clear monotonic improvements with respect to shot count.

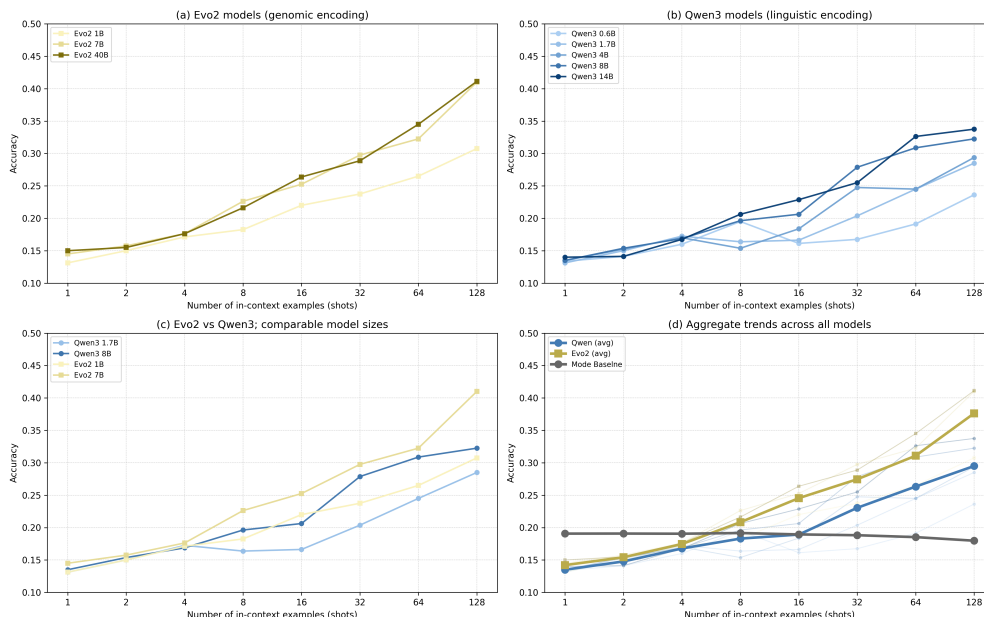


Figure 3: Few-shot performance of Qwen3 and Evo2 models. (a) Evo2 model performance with respect to $\log(\text{shots})$. All models monotonically improve. 7B and 40B beat 1B. (b) Qwen3 model performance with respect to $\log(\text{shots})$. All models improve, but not always monotonically. (c) At comparable sizes, Evo2 outperforms Qwen3. (d) Averaged performance across both model families. All models exceed the mode baseline shown in gray color. The exact accuracies and error bars (bootstrap-based standard errors) are included in §I.

Evo2 performs ICL on genomic tasks. In Fig. 2, we demonstrate that Evo2 7B can also perform few-shot in-context learning (ICL) on a relatively in-distribution binary genomic classification task: `human_nontata_promoters`, in which the model must predict whether a 251 nucleotide-long human DNA sequence is a promoter. Details of this experiment can be found in §K.

Further Analysis. We provide additional analyses of *BitLoad* (how many input bits influence the output) sensitivity in §D.1, *BitDiversity* (the number of minority bits in the output string) sensitivity in §L, and qualitative task differences in §D.2. Finally, we also provide analysis of how model scale impacts accuracy in §J.

4 CONCLUSION

We introduce a suite of bitstring reasoning tasks that can be encoded in both natural language and genomic sequences, showing that genomic models – like their linguistic counterparts – exhibit clear in-context learning. Across all Evo2 model sizes, we observe robust log-linear gains in accuracy with increasing demonstrations, paralleling the scaling trends of Qwen3 language models. These findings challenge the notion that ICL is unique to human language, suggesting it emerges whenever an expressive model is trained autoregressively on structured, pattern-rich data. See §A for a broader discussion of implications.

Potential future work: This work motivates searching for ICL in other non-linguistic modalities – time series (Das et al., 2024), system logs (Akhauri et al., 2025), physics simulations (Holzschuh et al., 2025), chess games (Ruoss et al., 2024), and climate projections (Duncan et al., 2025). Each offers a structured, patterned substrate that could support its own form of contextual reasoning. These diverse modalities, each with their unique structure and constraints, suggest a rich world of non-linguistic ICL capability waiting to be explored, and this work represents a maiden voyage into these extremely interesting waters.

REFERENCES

- 216
217
218 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transform-
219 ers learn to implement preconditioned gradient descent for in-context learn-
220 ing. In *Advances in Neural Information Processing Systems* (NeurIPS), 2023.
221 URL [http://papers.nips.cc/paper_files/paper/2023/hash/
222 8ed3d610ea4b68e7afb30ea7d01422c6-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/8ed3d610ea4b68e7afb30ea7d01422c6-Abstract-Conference.html).
- 223 Yash Akhauri, Bryan Lewandowski, Cheng-Hsi Lin, Adrian N. Reyes, Grant C. Forbes, Arissa
224 Wongpanich, Bangding Yang, Mohamed S. Abdelfattah, Sagi Perel, and Xingyou Song. Perform-
225 ance prediction for large systems via text-to-text regression, 2025. URL [https://arxiv.
226 org/abs/2506.21718](https://arxiv.org/abs/2506.21718).
- 227 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning al-
228 gorithm is in-context learning? investigations with linear models. In *International Conference on
229 Learning Representations* (ICLR), 2022. URL <https://arxiv.org/abs/2211.15661>.
- 230 Arc Institute. Evo 2 1b base. [https://huggingface.co/arcinstitute/evo2_1b_
231 base](https://huggingface.co/arcinstitute/evo2_1b_base), 2025a. Model card on Hugging Face.
- 232 Arc Institute. Evo 2 40b. https://huggingface.co/arcinstitute/evo2_40b, 2025b.
233 Model card on Hugging Face.
- 234 Arc Institute. Evo 2 7b. https://huggingface.co/arcinstitute/evo2_7b, 2025c.
235 Model card on Hugging Face.
- 236
237 Hritik Bansal, Karthik Gopalakrishnan, Saket Dingliwal, Sravan Bodapati, Katrin Kirchhoff, and
238 Dan Roth. Rethinking the role of scale for in-context learning: An interpretability-based case
239 study at 66 billion scale. In *Annual Meeting of the Association for Computational Linguistics*
240 (ACL), 2023. URL <https://aclanthology.org/2023.acl-long.660/>.
- 241 Gonzalo Benegas, Chengzhong Ye, Carlos Albors, Jianan Canal Li, and Yun S Song. Genomic
242 language models: opportunities and challenges. *Trends in Genetics*, 2025.
- 243
244 Garyk Brixi, Matthew G. Durrant, Jerome Ku, Michael Poli, Greg Brockman, Daniel Chang,
245 Gabriel A. Gonzalez, Samuel H. King, David B. Li, Aditi T. Merchant, Mohsen Naghipourfar,
246 Eric Nguyen, Chiara Ricci-Tam, David W. Romero, Gwanggyu Sun, Ali Taghibakshi, Anton
247 Vorontsov, Brandon Yang, Myra Deng, Liv Gorton, Nam Nguyen, Nicholas K. Wang, Etowah
248 Adams, Stephen A. Baccus, Steven Dillmann, Stefano Ermon, Daniel Guo, Rajesh Ilango, Ken
249 Janik, Amy X. Lu, Reshma Mehta, Mohammad R.K. Mofrad, Madelena Y. Ng, Jaspreet Pannu,
250 Christopher Ré, Jonathan C. Schmok, John St. John, Jeremy Sullivan, Kevin Zhu, Greg Zynda,
251 Daniel Balsam, Patrick Collison, Anthony B. Costa, Tina Hernandez-Boussard, Eric Ho, Ming-
252 Yu Liu, Thomas McGrath, Kimberly Powell, Dave P. Burke, Hani Goodarzi, Patrick D. Hsu, and
253 Brian L. Hie. Genome modeling and design across all domains of life with evo 2. *bioRxiv*, 2025.
254 doi: 10.1101/2025.02.18.638918. URL [https://www.biorxiv.org/content/early/
255 2025/02/21/2025.02.18.638918](https://www.biorxiv.org/content/early/2025/02/21/2025.02.18.638918).
- 256 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
257 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
258 few-shot learners. *Advances in Neural Information Processing Systems* (NeurIPS), 2020a. URL
259 <https://arxiv.org/abs/2005.14165>.
- 260 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
261 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
262 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
263 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
264 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
265 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*,
266 [abs/2005.14165](https://arxiv.org/abs/2005.14165), 2020b. URL <https://arxiv.org/abs/2005.14165>.
- 267 Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond,
268 James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learn-
269 ing in transformers. *Advances in Neural Information Processing Systems* (NeurIPS), 2022. URL
<https://arxiv.org/abs/2205.05055>.

- 270 Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. Parallel structures in pre-
271 training data yield in-context learning. *ArXiv preprint*, abs/2402.12530, 2024. URL <https://arxiv.org/abs/2402.12530>.
272
273
- 274 François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
275
- 276 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
277 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
278 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
279 2021. URL <https://arxiv.org/pdf/2110.14168>.
- 280 Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
281 time-series forecasting, 2024. URL <https://arxiv.org/abs/2310.10688>.
- 282 James P. C. Duncan, Elynn Wu, Surya Dheeshjith, Adam Subel, Troy Arcomano, Spencer K.
283 Clark, Brian Henn, Anna Kwa, Jeremy McGibbon, W. Andre Perkins, William Gregory, Carlos
284 Fernandez-Granda, Julius Busecke, Oliver Watt-Meyer, William J. Hurlin, Alistair Adcroft, Laure
285 Zanna, and Christopher Bretherton. Samudrace: Fast and accurate coupled climate modeling with
286 3d ocean and atmosphere emulators, 2025. URL <https://arxiv.org/abs/2509.12490>.
- 287 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,
288 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Gan-
289 guli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal
290 Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris
291 Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
292 URL <https://transformer-circuits.pub/2021/framework/index.html>.
- 293 Eric Elmoznino, Thomas Jiralerspong, Yoshua Bengio, and Guillaume Lajoie. A complexity-based
294 theory of compositionality. *ArXiv preprint*, abs/2410.14817, 2024a. URL <https://arxiv.org/abs/2410.14817>.
295
296
- 297 Eric Elmoznino, Tom Marty, Tejas Kasetty, Leo Gagnon, Sarthak Mittal, Mahan Fathi, Dhanya
298 Sridhar, and Guillaume Lajoie. In-context learning and occam’s razor. *ArXiv preprint*,
299 abs/2410.14086, 2024b. URL <https://arxiv.org/abs/2410.14086>.
- 300 Zhouxiang Fang, Aayush Mishra, Muhan Gao, Anqi Liu, and Daniel Khashabi. ICL Ciphers: Quan-
301 tifying “Learning” in In-Context Learning via Substitution Ciphers. In *Conference on Empiri-
302 cal Methods in Natural Language Processing (EMNLP)*, 2025. URL <https://arxiv.org/abs/2504.19395>.
303
304
- 305 Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle
306 use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions
307 of the Association for Computational Linguistics (TACL)*, 2021. URL <https://arxiv.org/abs/2101.02235>.
308
- 309 Riccardo Grazi, Julien Siems, Simon Schrod, Thomas Brox, and Frank Hutter. Is mamba capable
310 of in-context learning? *ArXiv preprint*, abs/2402.03170, 2024. URL <https://arxiv.org/abs/2402.03170>.
311
- 312 Katarína Grešová, Vlastimil Martinek, David Čechák, Petr Šimeček, and Panagiotis Alexiou. Ge-
313 nomic benchmarks: a collection of datasets for genomic sequence classification. *BMC Genomic
314 Data*, 24:25, May 2023. doi: 10.1186/s12863-023-01123-8.
315
- 316 Michael Hahn and Navin Goyal. A theory of emergent in-context learning as implicit structure
317 induction. *arXiv preprint arXiv:2303.07971*, 2023. URL <https://arxiv.org/abs/2303.07971>.
318
- 319 Damian Hodel and Jevin West. Response: Emergent analogical reasoning in large language models.
320 *arXiv preprint arXiv:2308.16118*, 2023.
321
322
323

- 324 Benjamin Holzs Schuh, Qiang Liu, Georg Kohl, and Nils Thuerey. PDE-transformer: Efficient and
325 versatile transformers for physics simulations. In Aarti Singh, Maryam Fazel, Daniel Hsu, Si-
326 mon Lacoste-Julien, Felix Berkenkamp, Tegan Maharaj, Kiri Wagstaff, and Jerry Zhu (eds.),
327 *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Pro-
328 ceedings of Machine Learning Research*, pp. 23562–23602. PMLR, 13–19 Jul 2025. URL <https://proceedings.mlr.press/v267/holzs Schuh25a.html>.
329
- 330 Zhongtao Jiang, Yuanzhe Zhang, Kun Luo, Xiaowei Yuan, Jun Zhao, and Kang Liu. On the in-
331 context generation of language models. In *Conference on Empirical Methods in Natural Lan-
332 guage Processing (EMNLP)*, 2024. doi: 10.18653/v1/2024.emnlp-main.568. URL <https://aclanthology.org/2024.emnlp-main.568/>.
333
- 334 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
335 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
336 models. *arXiv preprint arXiv:2001.08361*, 2020. URL <https://arxiv.org/abs/2001.08361>.
337
- 338 Ivan Lee, Nan Jiang, and Taylor Berg-Kirkpatrick. Exploring the relationship between model
339 architecture and in-context learning ability. *ArXiv preprint*, abs/2310.08049, 2023. URL <https://arxiv.org/abs/2310.08049>.
340
- 341 Martha Lewis and Melanie Mitchell. Evaluating the robustness of analogical reasoning in large
342 language models. *arXiv preprint arXiv:2411.14215*, 2024.
343
- 344 Shuai Li, Zhao Song, Yu Xia, Tong Yu, and Tianyi Zhou. The closeness of in-context learning and
345 weight shifting for softmax regression. *arXiv preprint arXiv:2304.13276*, 2023. URL <https://arxiv.org/abs/2304.13276>.
346
- 347 Ziqian Lin and Kangwook Lee. Dual operating modes of in-context learning. In *International
348 Conference on Machine Learning (ICML)*, 2024. URL <https://arxiv.org/pdf/2402.18819>.
349
- 350 Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is
351 provably the optimal in-context learner with one layer of linear self-attention. *ArXiv preprint*,
352 abs/2307.03576, 2023. URL <https://arxiv.org/abs/2307.03576>.
353
- 354 Ziv Nevo and Ran El-Yaniv. On online learning of decision lists. *J. Mach. Learn. Res.*, 3:271–301,
355 2002. URL <https://jmlr.org/papers/v3/nevo02a.html>.
356
- 357 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
358 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction
359 heads. *arXiv preprint arXiv:2209.11895*, 2022. URL <https://arxiv.org/abs/2209.11895>.
360
- 361 Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What in-context learning “learns” in-context:
362 Disentangling task recognition and task learning. In *Findings of the Association for Computa-
363 tional Linguistics: ACL 2023*, July 2023. URL <https://aclanthology.org/2023.findings-acl.527>.
364
- 365 Madhur Panwar, Kabir Ahuja, and Navin Goyal. In-context learning through the bayesian prism.
366 *ArXiv preprint*, abs/2306.04891, 2023. URL <https://arxiv.org/abs/2306.04891>.
367
- 368 Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kang-
369 wook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative
370 study on in-context learning tasks. *ArXiv preprint*, abs/2402.04248, 2024. URL <https://arxiv.org/abs/2402.04248>.
371
- 372 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language
373 understanding by generative pre-training. 2018. URL <https://openai.com/index/language-unsupervised/>.
374
- 375 John Carlyle Raven, John H. Court, and John Carlyle Raven. Manual for raven’s progressive
376 matrices and vocabulary scales. 1962. URL <https://api.semanticscholar.org/CorpusID:143337389>.
377

- 378 Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context
379 classification task. In *International Conference on Learning Representations (ICLR)*, 2023. URL
380 <https://arxiv.org/pdf/2312.03002>.
381
- 382 Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Quanshi Zhang, Xipeng Qiu, and Dahua
383 Lin. Identifying semantic induction heads to understand in-context learning. *arXiv preprint*
384 *arXiv:2402.13055*, 2024.
- 385 Anian Ruoss, Grégoire Delétang, Sourabh Medapati, Jordi Grau-Moya, Li Kevin Wenliang, Elliot
386 Catt, John Reid, Cannada A. Lewis, Joel Veness, and Tim Genewein. Amortized planning with
387 large-scale transformers: A case study on chess, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2402.04494)
388 [2402.04494](https://arxiv.org/abs/2402.04494).
- 389 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
390 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,
391 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W.
392 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain,
393 Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, An-
394 ders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, An-
395 drew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh
396 Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabas-
397 sum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Her-
398 rick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph,
399 Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin
400 Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron
401 Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh,
402 Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites,
403 Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera,
404 Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Gar-
405 rette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy,
406 Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito,
407 Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Den-
408 nis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta
409 Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Eka-
410 terina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Eliza-
411 beth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem,
412 Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Ev-
413 genii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé,
414 Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán
415 Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-
416 López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh
417 Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hong-
418 ming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson
419 Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel,
420 James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema
421 Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova,
422 Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Ji-
423 acheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis,
424 Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph
425 Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua,
426 Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja
427 Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chia-
428 fullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo
429 Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency,
430 Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón,
431 Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi,
Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy,

- 432 Michael Starritt, Michael Strube, Michał Śwędrowski, Michele Bevilacqua, Michihiro Yasunaga,
 433 Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal,
 434 Moin Aminnaseri, Mor Geva, Mozhdah Gheini, Mukund Varma T, Nanyun Peng, Nathan A.
 435 Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron,
 436 Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar,
 437 Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar El-
 438 baghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung,
 439 Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Pe-
 440 ter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour,
 441 Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer
 442 Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A.
 443 Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Ro-
 444 man Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov,
 445 Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Moham-
 446 mad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R.
 447 Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghaz-
 448 arian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schus-
 449 ter, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar
 450 Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upad-
 451 hyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy,
 452 Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene,
 453 Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Pianta-
 454 dosi, Stuart M. Shieber, Summer Mishnerghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen,
 455 Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore
 456 Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Ti-
 457 tus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz,
 458 Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh,
 459 Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders,
 460 William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen,
 461 Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang,
 462 Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid,
 463 Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game:
 Quantifying and extrapolating the capabilities of language models. *Transactions on Machine
 Learning Research* (TMLR), 2023. URL <https://arxiv.org/abs/2206.04615>.
- 464 Ramzan Kh. Umarov and Victor V. Solovyev. Recognition of prokaryotic and eukaryotic pro-
 465 moters using convolutional deep learning neural networks. *PLOS ONE*, 12(2):e0171410,
 466 2017. doi: 10.1371/journal.pone.0171410. URL <https://doi.org/10.1371/journal.pone.0171410>.
- 467
 468 Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun.
 469 Label words are anchors: An information flow perspective for understanding in-context learn-
 470 ing. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023a.
 471 doi: 10.18653/v1/2023.emnlp-main.609. URL <https://aclanthology.org/2023.emnlp-main.609>.
- 472
 473 Xiaolei Wang, Xinyu Tang, Wayne Xin Zhao, and Ji-Rong Wen. Investigating the pre-training dy-
 474 namics of in-context learning: Task recognition vs. task learning. *ArXiv preprint*, abs/2406.14022,
 475 2024. URL <https://arxiv.org/abs/2406.14022>.
- 476
 477 Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large
 478 language models are latent variable models: Explaining and finding good demonstra-
 479 tions for in-context learning. In *Advances in Neural Information Processing Systems*
 480 (NeurIPS), 2023b. URL [http://papers.nips.cc/paper_files/paper/2023/
 481 hash/3255a7554605a88800f4e120b3a929e1-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/3255a7554605a88800f4e120b3a929e1-Abstract-Conference.html).
- 482
 483 Taylor Webb, Keith J Holyoak, and Hongjing Lu. Emergent analogical reasoning in large lan-
 484 guage models. *Nature Human Behaviour*, 7(9):1526–1541, 2023. URL <https://arxiv.org/abs/2212.09196>.
- 485

486 Taylor W Webb, Keith J Holyoak, and Hongjing Lu. Evidence from counterfactual tasks supports
487 emergent analogical reasoning in large language models. *PNAS nexus*, 4(5):pgaf135, 2025.
488

489 Kevin Christian Wibisono and Yixin Wang. In-context learning from training on unstructured
490 data: The role of co-occurrence, positional information, and training data structure. In *ICML*
491 *2024 Workshop on Theoretical Foundations of Foundation Models*, 2024. URL [https://](https://openreview.net/forum?id=Zvwwnfwxa4)
492 openreview.net/forum?id=Zvwwnfwxa4.

493 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
494 learning as implicit bayesian inference. In *International Conference on Learning Representations*
495 (ICLR), 2021. URL <https://arxiv.org/abs/2111.02080>.

496 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
497 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
498 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
499 Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
500 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
501 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
502 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
503 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
504 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

A DISCUSSION

What are the implications of our findings on prior efforts to explain the emergence of ICL?

First, let us organize the existing frameworks for pinpointing the conditions under which ICL emerges:

- (E1) **ICL’s emergence is due to data distributional properties:** The distributional properties of data, such as “parallel structures” in human language pretraining data (Chen et al., 2024), its compositional structure (Hahn & Goyal, 2023), “burstiness” (Chan et al., 2022) and other such properties (Wibisono & Wang, 2024; Reddy, 2023) may be of importance (and perhaps necessary) for the emergence of ICL.
- (E2) **ICL’s emergence is due to a compression mechanism:** The large-scale compression mechanism during massive pretraining might drive ICL (Elmoznino et al., 2024a;b; Hahn & Goyal, 2023).
- (E3) **ICL’s emergence may require specific architectural properties:** While Transformers might be better suited for ICL than LSTMs (Xie et al., 2021), evidence is mixed (Lee et al., 2023), and non-Transformer models have also demonstrated ICL capabilities (Grazzi et al., 2024; Park et al., 2024).

Our findings refine existing hypotheses about ICL’s origins. The emergence of ICL in genomic models challenges accounts that rely solely on language-specific distributional structures (E1). The presence of ICL across both genomic and linguistic models supports the compression-based explanation (E2), suggesting that large-scale sequence compression and its induced inductive biases drive ICL across modalities. With respect to architecture (E3), Evo2 – an autoregressive hybrid combining convolutional and attention layers rather than a pure Transformer – exhibits similar scaling behavior, indicating that ICL does not depend on the pure Transformer form. Instead, architecture provides an expressive substrate that enables pattern induction once exposed to sufficiently large and structured data. Overall, these results position ICL as a modality-agnostic outcome of large-scale next-token prediction, rather than a phenomenon tied to linguistic statistics or a specific architecture.

What are the implications of our findings on the frameworks to explain how ICL operates?

We next consider the major perspectives that seek to explain how ICL operates. One view holds that ICL functions as a mix of *task learning* and *task retrieval*, with demonstrations serving either to recall pretrained capabilities or to enable learning on the fly (Pan et al., 2023; Lin & Lee, 2024; Wang et al., 2024; Fang et al., 2025). Our symbolic reasoning tasks, instantiated in both linguistic and genomic domains, provide direct evidence for this *task learning* mode, aligning with this hypothesis and prior work (Pan et al., 2023; Fang et al., 2025). Because these tasks do not depend on pre-trained semantic priors, they do *not* invoke *task retrieval*, offering limited insight into the Bayesian view that interprets ICL as implicit inference over latent concepts (Xie et al., 2021; Panwar et al., 2023; Wang et al., 2023b; Jiang et al., 2024). Meanwhile, our results remain agnostic toward the optimization-based hypothesis, which posits that ICL implements an implicit gradient-descent-like process (Akyürek et al., 2022; Ahn et al., 2023; Mahankali et al., 2023; Li et al., 2023), as well as the induction-based account, which attributes ICL to specialized “circuits” for performing inductive generalization (Elhage et al., 2021; Olsson et al., 2022; Wang et al., 2023a; Bansal et al., 2023; Ren et al., 2024). Together, our results most strongly support the presence of genuine *task learning* within ICL.

Does Evo2 have an innate advantage on these tasks? Possibly, for multiple reasons. First, though Evo2’s is trained on less tokens total, all of Evo2’s training tokens are long sequences of repeated nucleotides, and very little of Qwen3’s training tokens are long sequences of repeated digits. Second, Evo2’s StripedHyena2 architecture was found to significantly outperform a vanilla transformer on long DNA sequences (Brix et al., 2025). This could give Evo2 an innate advantage on long contexts containing the same few symbols vs Qwen3. Our result results do not imply Evo2’s ICL ability is superior to Qwen3’s, and we concede that our few-shot prompting setup may be biased toward Evo2. Attempts to make the task more legible to Qwen3, however, ran into the confound that Qwen3 has pretraining exposure to bitstring manipulation. Any prompting setup that included 0s and 1s immediately resulted in an extreme increase in Qwen3’s performance. Future work is needed to establish a method that controls for Evo2’s structural advantages without granting Qwen3 an unfair edge via its pretraining knowledge.

594 **Why not test the models on semantic tasks?** Semantic tasks – such as identifying the capitals
 595 associated with countries, classifying malformed proteins, etc. – require a fair amount of pretraining
 596 exposure to the concepts involved in the task as well as measuring ICL. While it’s undoubtedly ICL
 597 when a model infers a semantic transformation (for instance, country→capital, word→opposite), the
 598 pretraining knowledge necessary to manifest this ICL precludes its use in extreme cross-modality
 599 comparisons. Focusing on far simpler bitstring transformations that can be learned entirely in-
 600 context allows for an apples-to-apples comparison between Evo2 and Qwen3. To ensure that the
 601 ICL observed is not an artifact of this simplified domain, we additionally demonstrate in Appendix K
 602 that Evo2 exhibits robust, scaling ICL on a native genomic task (promoter classification), though we
 603 exclude this from the main comparison to maintain parity with the linguistic models.

605 B EXPERIMENTAL DESIDERATA

607 **Desideratum 1: Cross-domain comparability.** Our experimental framework requires that each
 608 task be performable by *both* language and genomic models. Thus, every task must be representable
 609 in both linguistic and nucleotide alphabets. This constraint excludes existing benchmarks that rely
 610 on domain-specific semantics – such as language reasoning datasets (e.g., BIG-Bench, StrategyQA,
 611 GSM8K (Srivastava et al., 2023; Geva et al., 2021; Cobbe et al., 2021)) or biological tasks (e.g.
 612 variant effect prediction, exon identification (Brixi et al., 2025)). While one could, in principle,
 613 translate domain-specific tasks into an alternate alphabet (e.g., mapping language tokens to base
 614 sequences via quaternary encoding), doing so inherently biases the evaluation: the source domain
 615 retains advantages, while the target domain must operate on representations that are unnatural to
 616 it. As a result, such cross-domain encodings confound the comparison, reflecting representational
 617 translation artifacts rather than genuine differences in ICL behavior.

618 **Desideratum 2: Limited vocabulary.** Since genomic models operate on only four nucleotides
 619 (A, T, C, G), tasks must be expressible within an equally compact alphabet. This rules out the
 620 existing symbolic reasoning and analogy benchmarks (Hodel & West, 2023; Lewis & Mitchell,
 621 2024; Webb et al., 2025), which rely on richer vocabularies (e.g., shapes, colors, or linguistic tokens
 622 with explicit semantic roles). Such tasks cannot be faithfully represented in a four-token regime
 623 without introducing artifacts or structural loss. Accordingly, our evaluation focuses on tasks that
 624 retain abstract reasoning structure while remaining compatible with the low-vocabulary symbolic
 625 space shared across linguistic and genomic models.

626 C MODEL FAMILIES AND SELECTION RATIONALE

628 For fair cross-domain comparison, we use two representative model families: Qwen3 for human
 629 language and Evo2 for genomics. The rationale for this selection is as follows:

- 631 (a) *Parameter scaling:* Both families span multiple orders of magnitude in parameter count,
 632 enabling systematic scaling analysis of ICL ability. The Qwen3 series ranges from 0.6B
 633 to 14B parameters, while Evo2 includes 1B, 7B, and 40B models (Yang et al., 2025; Brixi
 634 et al., 2025). This parallel scaling structure facilitates consistent measurement of how ICL
 635 performance evolves with model size across linguistic and genomic modalities.
- 636 (b) *Compute matching:* The largest models in each family are trained with comparable total
 637 compute, offering an opportunity for an approximately compute-matched cross-domain
 638 comparison. Using the standard $6ND$ estimate (Kaplan et al., 2020), Qwen3-14B-Base
 639 is trained with about 3.2×10^{24} FLOPs, while Evo2-40B is trained with 2.25×10^{24}
 640 FLOPs (Yang et al., 2025; Brixi et al., 2025), making Evo2 uniquely well-suited for com-
 641 parison with Qwen3 at scale.
- 642 (c) *Availability of base models:* The Qwen3 family releases base (*pre*-instruction-tuned) mod-
 643 els at all scales up to 14B parameters, enabling direct evaluation of the intrinsic inductive
 644 reasoning ability of pure next-token predictors, without instruction-tuning artifacts. The
 645 Evo2 base models have no additional instruction tuning applied—the only nuance is that
 646 Evo2 1B is trained at a context length of 8192 nucleotides, whereas the 7B/40B are ex-
 647 tended to a context length of one million.

- 648 (d) *Tokenizer*: Qwen3 uses a standard BPE tokenizer with a vocabulary size of 151,669. Evo2
649 uses a byte-level tokenizer, so individual nucleotides are mapped to individual tokens. No-
650 tably, Qwen’s tokenizer maps single digits to single tokens, allowing for parity in our ex-
651 periments (Yang et al., 2025; Brixi et al., 2025).
- 652 (e) *Context length*: Qwen3’s 0.6B and 1.7B models have a context length of 32K. The remain-
653 ing dense models have a context length of 128K. Evo2’s 1B model has a context length of
654 8K, and the remaining models have a context length of 1M. As none of our experiments
655 approach these limits, we can use all models in both families without fear of context length
656 as a confound (Yang et al., 2025; Brixi et al., 2025).
- 657 (f) *Training corpora*: All Qwen3 models are trained on 36 trillion text tokens covering a wide
658 variety of topics and languages. Evo1 1B is trained on 1 trillion tokens, Evo2 7B is trained
659 on 2.4 trillion tokens, and Evo2 40B is trained on 9.3 trillion tokens. These extensive
660 training corpora ensure that any potential ICL dynamics have thoroughly emerged (Yang
661 et al., 2025; Brixi et al., 2025).
- 662 (g) *Architecture*: Qwen3 is based on a conventional Llama-like transformer architecture, while
663 Evo2 uses the StripedHyena2 architecture which intersperses convolutional layers with
664 attention-based ones. While ideally both model families would use the same architecture,
665 there were no vanilla transformers at Evo2’s compute scale (Yang et al., 2025; Brixi et al.,
666 2025).
- 667 (h) *Licensing*: All models are released under Apache 2.0 (Yang et al., 2025; Arc Institute,
668 2025a;c;b), ensuring reproducibility.

670 D SUPPLEMENTARY ANALYSES OF TASK COMPLEXITY AND QUALITATIVE
671 DIFFERENCES

673 D.1 ICL SENSITIVITY TO TASK COMPLEXITY: BITLOAD ANALYSIS

674
675 To understand which transformations Qwen3 and Evo2 infer most effectively, we analyze their per-
676 formance across varying task complexities. We focus on the largest models in each family (Qwen3-
677 14B and Evo2-40B) under the ($n = 128$) shot regime, providing both models ample opportunity to
678 display their ICL abilities.

679 **Defining BitLoad.** We introduce *BitLoad*, a
680 measure of a function’s intrinsic complexity.
681 Informally, BitLoad quantifies how many input
682 bits influence the output. Formally, it is
683 defined as:

684
685
$$\text{BitLoad}(f) = \sum_{i=1}^k \mathbb{1}[\exists x, j : f_j(x) \neq f_j(x^{\oplus i})],$$

686
687 (2)

688 where $x^{\oplus i}$ denotes x with bit i flipped, and
689 $f_j(x)$ returns the j -th output bit. Intuitively, it
690 counts the number of bit positions whose per-
691 turbation changes the output. So, the larger the
692 BitLoad of a function, the harder it is since it
693 requires the model to attend to more bits. This
694 metric is similar to the existing statistical mea-
695 sure of "relevant features", just defined specifically for our bitstring manipulation tasks. (Nevo &
696 El-Yaniv, 2002).

697 Fig. 4 shows the mean accuracy of Qwen vs Evo with respect to the BitLoad of all of our tasks.
698 A full table of the BitLoad of every tested task is attached in the appendix E.3. We see that both
699 models achieve near-perfect accuracy on constant (0 BitLoad) tasks and remain similar at BitLoad
700 1. However, by BitLoad 2, Evo begins to outperform Qwen, and beyond that point Qwen’s accuracy
701 drops sharply – falling below 20% by BitLoad 4 – while Evo’s performance declines more gradually,
remaining above 40% before converging near 20% at BitLoad 8. This asymmetry suggests that Evo

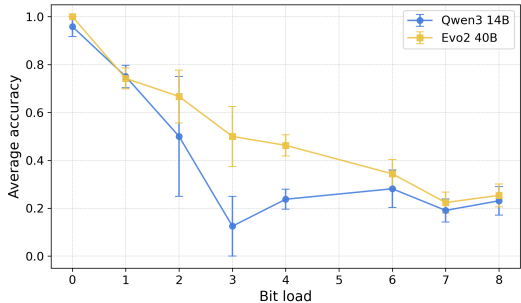


Figure 4: Accuracy vs. BitLoad averaged across all tasks (BitLoad; Eq.2). Qwen declines sharply with increasing BitLoad, while Evo degrades more gradually, indicating greater robustness. Details in §D.1.

702 maintains partial generalization as dependency depth increases (BitLoad value), whereas Qwen’s
703 ICL collapses more abruptly.

704 While BitLoad strongly correlates with task accuracy, Fig. 4 also shows notable deviations within
705 the error bars. Thus, while task complexity is a strong predictor of ICL accuracy, other factors
706 (e.g., transformation depth, pretraining exposure) likely play meaningful roles, which motivate our
707 qualitative study in §D.2.

709 D.2 QUALITATIVE ANALYSIS OF FUNCTIONAL AND BEHAVIORAL DIFFERENCES

710
711 To complement the quantitative BitLoad analysis (§D.1), we conduct a per-task comparison to iden-
712 tify qualitative differences between Qwen3 and Evo2. Specifically, we analyze which tasks each
713 model performs best on and how their inductive profiles diverge. We focus on the $n = 128$ shot
714 regime for the largest models – Qwen3-14B and Evo2-40B – where both have maximal opportunity
715 to exhibit ICL. We rank all tasks by model accuracy and examine the top ten for each model. Tasks
716 that appear in one model’s top ten but not the other are considered “exclusive.”

717 **Exclusive competencies.** Qwen’s exclusive tasks involve right-shift operations:
718 "spread_last_bit" \rightarrow "shift_right_zero" and "edge_mask" \rightarrow
719 "shift_right_zero". For instance, "shift_right_zero" pads the bitstring on the
720 left with a zero and truncates the last bit (e.g., 01010000 \rightarrow 00101000). Qwen achieves 100%
721 accuracy on "spread_last_bit" \rightarrow "shift_right_zero", whereas Evo2 achieves
722 only 50%. A similar but smaller gap appears for "edge_mask" \rightarrow "shift_right_zero"
723 (87.5% vs. 62.5%).

724 In contrast, Evo2’s exclusive strengths involve multi-bit transformations. A clear example is
725 "flip_bits" \rightarrow "right_half", which applies bitwise NOT followed by masking the first
726 half of the input (e.g., 01011100 \rightarrow 00000011). Evo2 achieves 87.5% accuracy, while Qwen only
727 25%. This task has a BitLoad of four, consistent with Evo’s superior performance on medium-
728 complexity (2–4 bit) transformations observed in Fig. 4.

729 **Shared strengths.** Despite these differences, 7 of the top 10 tasks are shared between Qwen and
730 Evo, yielding an intersection-over-union of 0.54. Both models excel at simple transformations such
731 as constant-outputs or single-bit dependencies, e.g., "spread_last_bit" which copies the final
732 bit to all positions.

733 **Differential skill profiles.** To sharpen contrasts, we identify the ten tasks most favoring each model.
734 Qwen’s advantages are concentrated in simple shifts and aggregation tasks. It outperforms Evo by
735 37.5% on the "minority" operation (output all 1s if zeros $>$ ones), and by a similar margin on
736 two parity-based tasks requiring counting the number of 1s. These trends suggest that Qwen may be
737 better at reasoning over global properties of bitstrings. Its superiority on simple shifts may also be
738 explained by the extreme rarity of frame shift mutations in DNA due to how catastrophic they are – a
739 single nucleotide offset can decimate an entire protein. This is empirically supported by the fact that
740 single-nucleotide deletions/shifts increase perplexity far more than other common mutations when
741 presented to Evo2 (Brix et al., 2025).

742 Evo2, by contrast, dominates tasks requiring full-bitstring manipulation. It achieves 62.5% on
743 bitwise NOT (vs. 0% for Qwen), 62.5% on identity (vs. 12.5%), and large margins on com-
744 positions such as "flip_bits" \rightarrow "right_half" (87.5% vs. 25%) and "rot11" \rightarrow
745 "flip_bits" (37.5% vs. 0%). This exposes perhaps the most important difference between
746 Qwen and Evo’s ICL in this specific context: Evo can learn simple full-bitstring operations in-
747 context, whereas Qwen cannot. Notably, Qwen3’s base models are trivially capable of learning the
748 identity in a more familiar few-shot context – when examples are presented with arrows and new-
749 lines separating them, instead of our intentionally unfamiliar encoding. Thus these results should
750 be taken as an existence proof of ICL in Evo2, not a definitive statement of Evo2 having more ICL
751 ability than Qwen. We leave a comparison of these models across broader tasks to future work.

E FURTHER DETAILS ON SYNTHETIC TASK DEFINITION

E.1 TABLE OF PRIMITIVES

The following table describes the thirty primitives used to construct the task space via composition – their use is described in §2.1.

Primitive	Description
alternating_start_one	Produce a mask marking positions that differ from the alternating 1010... pattern starting with 1 (1 = mismatch, 0 = match).
alternating_start_zero	Produce a mask marking positions that differ from the alternating 0101... pattern starting with 0 (1 = mismatch, 0 = match).
center_mask	Zero out the first and last bits while leaving the interior bits unchanged; strings of length ≤ 2 become all zeros.
double_rotl	Circularly rotate the bitstring two positions to the left.
double_rotr	Circularly rotate the bitstring two positions to the right.
edge_mask	Preserve the first and last bits and zero out every interior bit (length 0/1 strings pass through).
flip_bits	Invert every bit, swapping 0s for 1s and vice versa.
identity	Return the bitstring unchanged.
invert_prefix	Flip the bits in the left half of the string, keep the right half as is.
invert_suffix	Keep the left half as is and flip every bit in the right half of the string.
keep_even_positions	Keep bits at even indices (0-based) and zero out bits at odd indices.
keep_odd_positions	Keep bits at odd indices (0-based) and zero out bits at even indices.
left_half	Preserve the left half of the string and replace the right half with zeros.
majority	Fill the string with the majority bit from the input; ties resolve to all 1s.
meta_constant	Returns a random, pre-set constant.
minority	Fill the string with the minority bit from the input; ties resolve to all 0s.
mirror_half	Copy the left half of the string onto the right half in reverse order, keeping the center bit unchanged for odd lengths.
ones_if_palindrome	Output all 1s if the input is a palindrome; otherwise output all 0s.
parity_fill	Output all 1s when the input contains an odd number of 1s; otherwise output all 0s.
reverse_bits	Reverse the order of the bits in the string.
right_half	Zero out the left half and keep the right half unchanged.
rotll	Circularly rotate the bitstring one position to the left.
rotrl	Circularly rotate the bitstring one position to the right.
shift_left_zero	Shift the string left by one, dropping the first bit and appending a 0 on the right.
shift_right_zero	Shift the string right by one, inserting a 0 on the left and dropping the last bit.
spread_first_bit	Replace every position with the first bit of the input.
spread_last_bit	Replace every position with the last bit of the input.
swap_halves	Swap the left and right halves of the string.
swap_pairs	Swap each adjacent pair of bits (positions 0/1, 2/3, ...).
xor_with_s0	Can only be applied after another primitive. Computes the logical XOR between the original input s_0 and the output of the first primitive.

Table 1: The 30 unary primitives used to construct functions in F .

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

E.2 TABLE OF FUNCTIONS

The following table describes the one hundred specific compositions of primitives used to construct the evaluation suite in §??.

Function	Function	Function
identity	spread_last_bit	swap_halves → shift_left_zero
rotll	invert_prefix	swap_halves → shift_right_zero
reverse_bits	invert_suffix	shift_left_zero → swap_halves
flip_bits	meta_constant	shift_right_zero → swap_halves
swap_halves	flip_bits → reverse_bits	keep_even_positions → flip_bits
majority	rotll → reverse_bits	keep_odd_positions → flip_bits
minority	reverse_bits → rotll	flip_bits → keep_even_positions
parity_fill	rotll → flip_bits	flip_bits → keep_odd_positions
alternating_start_one	swap_halves → reverse_bits	edge_mask → flip_bits
alternating_start_zero	swap_halves → flip_bits	center_mask → flip_bits
left_half	double_rotl → flip_bits	shift_left_zero → keep_even_positions
right_half	rotl → flip_bits	shift_left_zero → keep_odd_positions
double_rotl	spread_first_bit → flip_bits	shift_right_zero → keep_even_positions
rotl	spread_last_bit → flip_bits	shift_right_zero → keep_odd_positions
double_rotl	left_half → flip_bits	keep_even_positions → reverse_bits
ones_if_palindrome	right_half → flip_bits	keep_odd_positions → reverse_bits
mirror_half	flip_bits → left_half	shift_left_zero → parity_fill
spread_first_bit	flip_bits → right_half	shift_right_zero → parity_fill
spread_last_bit	double_rotl → reverse_bits	parity_fill → shift_left_zero
invert_prefix	rotll → swap_halves	parity_fill → shift_right_zero
invert_suffix	xor_with_s0	spread_first_bit → shift_left_zero
meta_constant	flip_bits → xor_with_s0	spread_last_bit → shift_right_zero
shift_left_zero	ones_if_palindrome → flip_bits	spread_first_bit → keep_even_positions
shift_right_zero	flip_bits → mirror_half	spread_last_bit → keep_odd_positions
swap_pairs	invert_prefix → reverse_bits	spread_first_bit → edge_mask
keep_even_positions	left_half → reverse_bits	spread_last_bit → edge_mask
keep_odd_positions	right_half → reverse_bits	spread_first_bit → center_mask
edge_mask	parity_fill → flip_bits	spread_last_bit → center_mask
center_mask	rotll → spread_first_bit	rotll → shift_left_zero
xor_with_s0	shift_left_zero → flip_bits	rotll → shift_right_zero
flip_bits → reverse_bits	shift_right_zero → flip_bits	shift_left_zero → rotll
rotll → reverse_bits	flip_bits → shift_left_zero	shift_right_zero → rotll
reverse_bits → rotll	flip_bits → shift_right_zero	reverse_bits → edge_mask
rotll → flip_bits	swap_pairs → flip_bits	reverse_bits → center_mask
swap_halves → reverse_bits	shift_left_zero → reverse_bits	edge_mask → shift_left_zero
swap_halves → flip_bits	shift_right_zero → reverse_bits	edge_mask → shift_right_zero
double_rotl → flip_bits	spread_first_bit → flip_bits	shift_left_zero → shift_left_zero
rotl → flip_bits	shift_left_zero → edge_mask	shift_left_zero → swap_pairs

Table 2: The complete set of 100 functions in F , consisting of 30 single primitives and 70 composed functions ($(f \rightarrow g)(x) = g(f(x))$).

E.3 BITLOAD OF FUNCTIONS

Function / Composition	BitLoad	Function / Composition	BitLoad
identity	8	rotll	8
reverse_bits	8	flip_bits	8
swap_halves	8	majority	8
minority	8	parity_fill	8
alternating_start_one	8	alternating_start_zero	8
left_half	4	right_half	4
double_rotl	8	rotrr	8
double_rotr	8	ones_if_palindrome	8
mirror_half	4	spread_first_bit	1
spread_last_bit	1	invert_prefix	8
invert_suffix	8	meta_constant	0
flip_bits → reverse_bits	8	rotll → reverse_bits	8
reverse_bits → rotll	8	rotll → flip_bits	8
swap_halves → reverse_bits	8	swap_halves → flip_bits	8
double_rotl → flip_bits	8	rotrr → flip_bits	8
spread_first_bit → flip_bits	1	spread_last_bit → flip_bits	1
left_half → flip_bits	4	right_half → flip_bits	4
flip_bits → left_half	4	flip_bits → right_half	4
double_rotl → reverse_bits	8	rotll → swap_halves	8
xor_with_s0	0	flip_bits → xor_with_s0	0
ones_if_palindrome → flip_bits	8	flip_bits → mirror_half	4
invert_prefix → reverse_bits	8	left_half → reverse_bits	4
right_half → reverse_bits	4	parity_fill → flip_bits	8
rotll → spread_first_bit	1	shift_left_zero	7
shift_right_zero	7	swap_pairs	8
keep_even_positions	4	keep_odd_positions	4
edge_mask	2	center_mask	6
shift_left_zero → flip_bits	7	shift_right_zero → flip_bits	7
flip_bits → shift_left_zero	7	flip_bits → shift_right_zero	7
swap_pairs → flip_bits	8	shift_left_zero → reverse_bits	7
shift_right_zero → reverse_bits	7	swap_halves → shift_left_zero	7
swap_halves → shift_right_zero	7	shift_left_zero → swap_halves	7
shift_right_zero → swap_halves	7	keep_even_positions → flip_bits	4
keep_odd_positions → flip_bits	4	flip_bits → keep_even_positions	4
flip_bits → keep_odd_positions	4	edge_mask → flip_bits	2
center_mask → flip_bits	6	shift_left_zero → keep_even_positions	4
shift_left_zero → keep_odd_positions	3	shift_right_zero → keep_even_positions	3
shift_right_zero → keep_odd_positions	4	keep_even_positions → reverse_bits	4
keep_odd_positions → reverse_bits	4	shift_left_zero → parity_fill	7
shift_right_zero → parity_fill	7	parity_fill → shift_left_zero	8
parity_fill → shift_right_zero	8	spread_first_bit → shift_left_zero	1
spread_last_bit → shift_right_zero	1	spread_first_bit → keep_even_positions	1
spread_last_bit → keep_odd_positions	1	spread_first_bit → edge_mask	1
spread_last_bit → edge_mask	1	spread_first_bit → center_mask	1
spread_last_bit → center_mask	1	rotll → shift_left_zero	7
rotll → shift_right_zero	7	shift_left_zero → rotll	7
shift_right_zero → rotll	7	reverse_bits → edge_mask	2
reverse_bits → center_mask	6	edge_mask → shift_left_zero	1
edge_mask → shift_right_zero	1	shift_left_zero → shift_left_zero	6
shift_left_zero → swap_pairs	7	shift_left_zero → edge_mask	1

Table 3: BitLoad for every primitive and composed function in F . Used in §D.1.

F A GLOBAL METRIC OVER ALL PROGRAMS

For a given function f , model M , and number of in-context examples n , we define the average accuracy over all context sets $E_N = \{E \subset S : |E| = n\}$:

$$A_f(M, n) = \frac{1}{|E_N|(|S| - n)} \sum_{E \in E_N} \sum_{x \in S \setminus E} A(f, E, x, M).$$

The overall benchmark score across all functions F is then

$$P(M, n) = \frac{1}{|F|} \sum_{f \in F} A_f(M, n).$$

Because exact evaluation is intractable (for $n = 8$, $|E_N|(|S| - N) = \binom{|S|}{n}(|S| - n) = \binom{256}{8} \cdot 248 \approx 1.016 \times 10^{17}$) we estimate $A_f(M, n)$ via Monte Carlo, as discussed in §2.2.

G PROMPT ENCODING

Prompts are encoded using a unified symbolic scheme to enable evaluation across linguistic and genomic models. For linguistic models, bits (0, 1) are mapped to two random digits (0–9); for genomic models, to random nucleotides (A, T, C, G). The separator token (\rightarrow) is omitted, and in-context examples are separated by a randomly chosen unused token distinct from those representing 0 and 1. (See Fig. 1 for examples.) All mappings are randomized per trial to avoid memorization or positional bias.

H MODE BASELINE

We define a mode baseline that always predicts the most frequent output observed in the context. For a given function f , context set $E \subset S$, and query input $x \in S \setminus E$, the mode prediction is $\hat{y}_{\text{mode}}(f, E, x) = \arg \max_{y \in S} |\{e \in E : f(e) = y\}|$, where ties in the $\arg \max$ are broken randomly. This simply corresponds to guessing the most common output in the few-shot examples. The overall mode baseline with n shots and across the set of all functions F is:

$$\hat{P}_{\text{mode}}(n) = \frac{1}{m|F|} \sum_{f \in F} \sum_{t=1}^m \mathbb{1}[\hat{y}_{\text{mode}}(f, E^{(t)}, x^{(t)}) = f(x^{(t)})], \quad (3)$$

where $E^{(t)}$ and $x^{(t)}$ are sampled identically to the model evaluation in Eq.1. This baseline corresponds to making an educated guess based only on the overall distribution of function outputs that simply learns the majority statistics of the prompt without attempting to infer the underlying transformation.

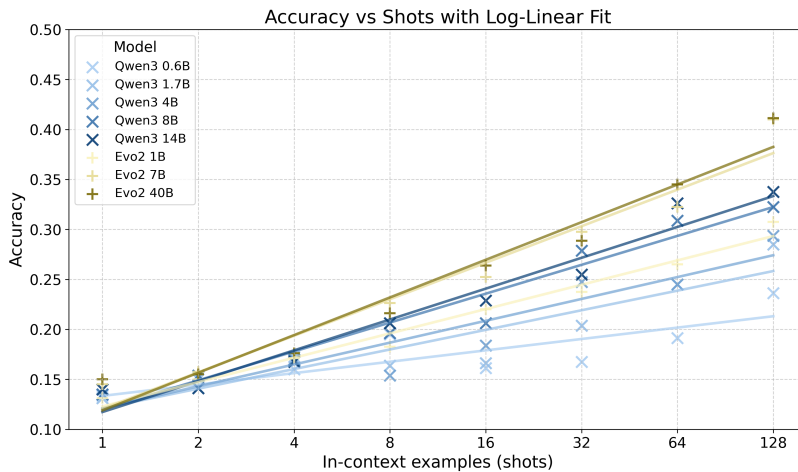
I FULL ACCURACY RESULTS

Here we show the full table of accuracies used in §3.1 to analyze the ICL capabilities of Evo2 and Qwen3 and perform the necessary statistical tests.

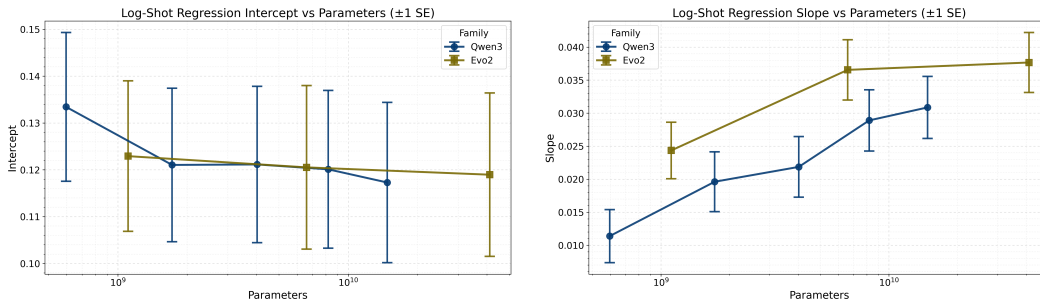
Model	1 Shot	2 Shots	4 Shots	8 Shots	16 Shots	32 Shots	64 Shots	128 Shots
Qwen3 0.6B	13.4 \pm 2.3	14.1 \pm 2.3	16.0 \pm 2.6	19.5 \pm 2.7	16.1 \pm 2.4	16.8 \pm 2.5	19.1 \pm 2.7	23.6 \pm 2.9
Qwen3 1.7B	13.1 \pm 2.2	15.0 \pm 2.7	17.2 \pm 2.8	16.4 \pm 2.6	16.6 \pm 2.7	20.4 \pm 2.9	24.5 \pm 3.4	28.5 \pm 3.5
Qwen3 4B	13.5 \pm 2.4	15.2 \pm 2.4	17.0 \pm 2.6	15.4 \pm 2.5	18.4 \pm 2.7	24.8 \pm 3.2	24.5 \pm 3.4	29.4 \pm 3.4
Qwen3 8B	13.5 \pm 2.3	15.4 \pm 2.5	16.9 \pm 2.7	19.6 \pm 2.7	20.6 \pm 2.8	27.9 \pm 3.3	30.9 \pm 3.5	32.2 \pm 3.5
Qwen3 14B	14.0 \pm 2.4	14.1 \pm 2.4	16.8 \pm 2.7	20.6 \pm 2.9	22.9 \pm 3.1	25.5 \pm 3.2	32.6 \pm 3.7	33.8 \pm 3.5
Evo2 1B	13.1 \pm 2.2	15.0 \pm 2.4	17.1 \pm 2.9	18.2 \pm 2.7	22.0 \pm 2.9	23.8 \pm 2.8	26.5 \pm 3.0	30.8 \pm 3.1
Evo2 7B	14.5 \pm 2.4	15.8 \pm 2.7	17.6 \pm 2.9	22.6 \pm 2.9	25.2 \pm 3.0	29.8 \pm 3.1	32.2 \pm 3.1	41.0 \pm 3.4
Evo2 40B	15.0 \pm 2.6	15.5 \pm 2.5	17.6 \pm 2.8	21.6 \pm 3.1	26.4 \pm 3.1	28.9 \pm 3.1	34.5 \pm 3.2	41.1 \pm 3.3

Table 4: In-context learning performance across model families and shot counts. Values show accuracy \pm standard error. All numbers are percentages. **Bold** numbers show the best performance within a model family. Models are ordered by parameter count within each family.

J META-REGRESSION FOR ICL EFFICACY WITH NUMBER OF DEMONSTRATIONS



(a) Few-shot performance of Qwen3 and Evo2 models. All models show consistent linear improvement with respect to $\log(\text{shots})$. In contrast, no such improvement occurs for the naive baseline.



(b) Baseline accuracy decreases slightly with scale as the model gains more parameters for both Evo2 and Qwen3. (c) ICL rate vs. model size: sharp gains up to 4B for Qwen3; mild boost from Evo2 1B to 7B; both plateau after 4–7B.

Figure 5: Few-shot behavior and scaling trends across Qwen3 and Evo2.

We perform linear regressions to predict accuracy from shot count with each model. For each model M , fit the linear regression: $\hat{P}(M, n) = \alpha_0(M) + \alpha_1(M) \log(n) + \varepsilon$. The raw regressions are shown in Fig. 5a. Predictably, all α_1 are positive as all models are capable of learning in-context.

We can interpret α_0 as representing the model’s base accuracy at the task, what it would logically achieve with only one shot to identify the task. We can then interpret α_1 as the model’s ICL efficacy: the speed at which it adapts to the task being presented and at which its accuracy improves. Analyzing how these values change across parameter values reveals insights into the ICL abilities of both the Qwen3 and Evo2 models.

First, we analyze how α_0 changes in each model family – this analysis can be seen in Fig. 5b. Both Evo2’s and Qwen3’s initial α_0 remains essentially constant model-to-model, indicating that all models have similar levels of few-shot baseline performance. Notably, Evo2 and Qwen3 have essentially identical intercepts at around 0.12. This implies that despite drastically different training data, the overall amount of prior knowledge the models have coming into this task is roughly similar. This rules out Qwen simply having ‘less experience’ with this sort of task.

If one looks at α_1 – ICL efficacy – in Fig. 5c, a dramatically different picture is painted. Here, both Qwen3 and Evo2 follow similar patterns with a significant difference. Qwen3’s ICL efficacy increases monotonically with respect to parameters, more than doubling from the 0.6B to the 14B. Evo2 follows suit (albeit less dramatically), with ICL efficacy monotonically increasing from the 1B to the 40B.

In absolute terms, however, Evo2, when parameter-matched, adapts in-context faster than Qwen3 does. Evo2 40B outperforms Qwen3 14B significantly, and it takes until Qwen3 8B to exceed the ICL ability of Evo2 1B.

Taken together, this data suggests that Qwen3 and Evo2 have similar amounts of pretraining exposure to be able to solve these tasks, and that Evo2 simply has better overall ICL capability (in this regime) – even though Qwen’s ICL ability increases more rapidly with respect to parameters.

K DEMONSTRATION OF EVO2’S ICL ABILITY ON A GENOMIC TASK

We demonstrate that Evo2 7B can perform few-shot in-context learning (ICL) on `human_nontata_promoters`, a binary genomic classification task in which the model must predict whether a 251 nucleotide-long human DNA sequence is a promoter. To make the task non-trivial, sequences in the positive class are restricted to non-TATA promoters, removing the “TATA box”, a pattern that serves as an obvious, exploitable marker (Umarov & Solovyev, 2017; Grešová et al., 2023).

For each test query, we construct a few-shot prompt by concatenating balanced examples from the training set (equal numbers of promoters and non-promoters), followed by a held-out test sequence. Each training example is immediately followed by an encoded label. The label consists of a 24 nucleotide buffer (a single nucleotide repeated 24 times) and then a 24 nucleotide label run (another nucleotide repeated 24 times). The buffer nucleotide and the two label nucleotides (for class 0 vs. class 1) are re-sampled uniformly at random without replacement from $\{A, C, G, T\}$ for every trial, so the mapping changes across evaluations. This prevents any nucleotide-specific biasing effects.

To predict the test label, we score two candidate prompts - one where the test sequence is followed by the class-0 label and one followed by the class-1 label - and choose the label with lower perplexity on the label run (i.e., computed only over the final 24 label nucleotides, excluding the buffer). This lets us measure whether Evo2 can infer the label mapping from the few-shot context and apply it to the query sequence.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

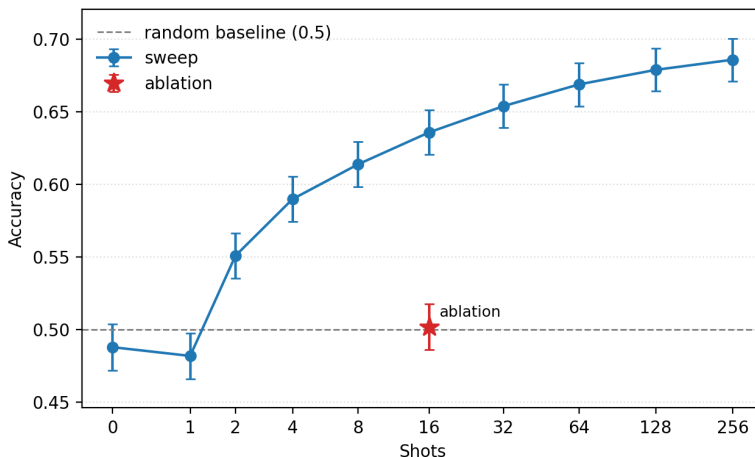


Figure 6: Performance of Evo2 7B on human_nontata_promoters at varying shot counts. Error bars are \pm one standard error. Evo2 7B’s performance shows clear monotonic improvements with respect to shot count.

In Figure 6, we observe monotonic improvement with respect to shot count, with accuracy rising from below random baseline at 1-shot (48.2%) to 63.6% at 16-shot. These gains further rise to 68.6% after 256 shots. A saturated power law achieves an R^2 of 0.998 in the 2-256 shot regime, which is characteristic of few-shot ICL. Furthermore, we show via ablation that the improvement is due to correctly-labeled examples presented in context—randomly permuting the labels in the 16-shot regime drops performance from 63.6% back to random chance (50.2%).

L ANALYSIS OF ICL WITH PROGRAM SYNTHESIS TASKS: BIT-DIVERSITY

BitLoad captures a measure of theoretical information bottleneck required to make predictions in program synthesis tasks. However, this does not always correlate with what models find easy to do. In Fig. 8 and Fig. 9, we plot the accuracy of Evo and Qwen models with respect to BitLoad. We notice that few high BitLoad tasks have high accuracy compared to all medium BitLoad tasks, illustrating that some high BitLoad tasks can be easy for these models to solve, probably due to the nature of patterns found during pre-training.

To further analyze the nature of ICL exhibited through these tasks, we define BitDiversity as *the number of minority bits in the output string*. In Fig. 10 and Fig. 11, we plot the accuracy with respect to BitDiversity. These plots try to estimate the effect of entropy in the output on model performance, and we see a more expected trend: models tend to perform better on low-entropy outputs, regardless of BitLoad. However, bin-wise performance trends are always increasing with shots, supporting the central hypothesis that ICL increases with increasing number of demos in these models.

Prevalence of 0-BitDiversity outputs. Looking at the expected and predicted outputs of trials from our tasks (Fig. 12), we made a few interesting observations about 0-BitDiversity (BD) outputs.

- Around 25% of true targets are 0-BD. This is a significantly high number as we only have 2 0-BD bit strings in all possible bit strings of any length K . It implies that many of our tasks create low BitDiversity outputs on random inputs, i.e., 0-BD outputs.
- Models tend to produce a lot of 0-BD outputs, much higher than the number of 0-BD true targets in the low-shot regime. But this number quickly drops to the expected number with higher shots.
- A large portion of the baseline (1-shot) performance can be explained by this prevalence of 0-BD outputs, but with more shots we get stronger evidence of ICL with increasing correctly predicted non 0-BD cases.

Understandable Mistakes. A potential confound is what we will call "understandable mistakes". These mistakes occur when the model outputs an incorrect answer that would be correct for some other programs given the few-shot context. Formally, a model's output y (see 2.2 for notation):

$$y = M(e_1 \rightarrow f(e_1), e_2 \rightarrow f(e_2), \dots, e_N \rightarrow f(e_N), x)$$

is an **understandable mistake** if $y \neq f(x)$ but there exists $f_2 \in F$ such that $\forall e_i \quad f(e_i) = f_2(e_i)$ and $y = f_2(x)$.

These understandable mistakes are an alarming confound at low shot counts, but their effect vanishes by $N = 16$. They occur most in tasks with low BitDiversity, which can often be confused with each other. Figure 7 below shows how understandable mistakes in Qwen3-4B's inferences decay exponentially as the number of shots increases.

Noise due to Monte Carlo Trials. We use $m = 8$ for our per-function Monte Carlo trials when we compute accuracy. This has little aggregate impact when comparing model performance across the entire suite, but drastically reduces the significance of results when comparing models and trends at the task-level. We only have eight discrete bands of accuracy at which we can estimate a model's per-task performance – which reduces the expressiveness of regression. Worse, this can lead to models getting a lucky prompt or two with a 0-BD output, which raises performance to 12.5% or 25% without the model truly understanding the task. Addressing these confounds would simply require increasing m by an order of magnitude or so. Alternatively, tasks of interest could be identified and m selectively increased for those tasks to enable more nuanced analysis.

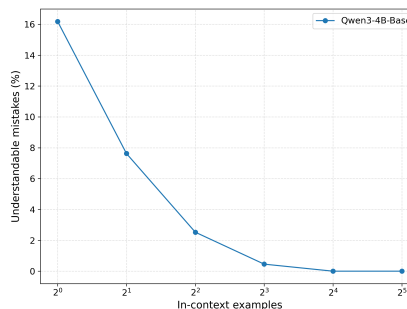


Figure 7: Qwen3-4B's rate of understandable mistakes with respect to the number of shots. Despite starting at 16% in the one-shot regime, they fall to less than 1% by 8 shots and vanish entirely at 32 shots. This underscores how understandable mistakes are only a confound at very low shot counts and can essentially be ignored past 4 shots.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

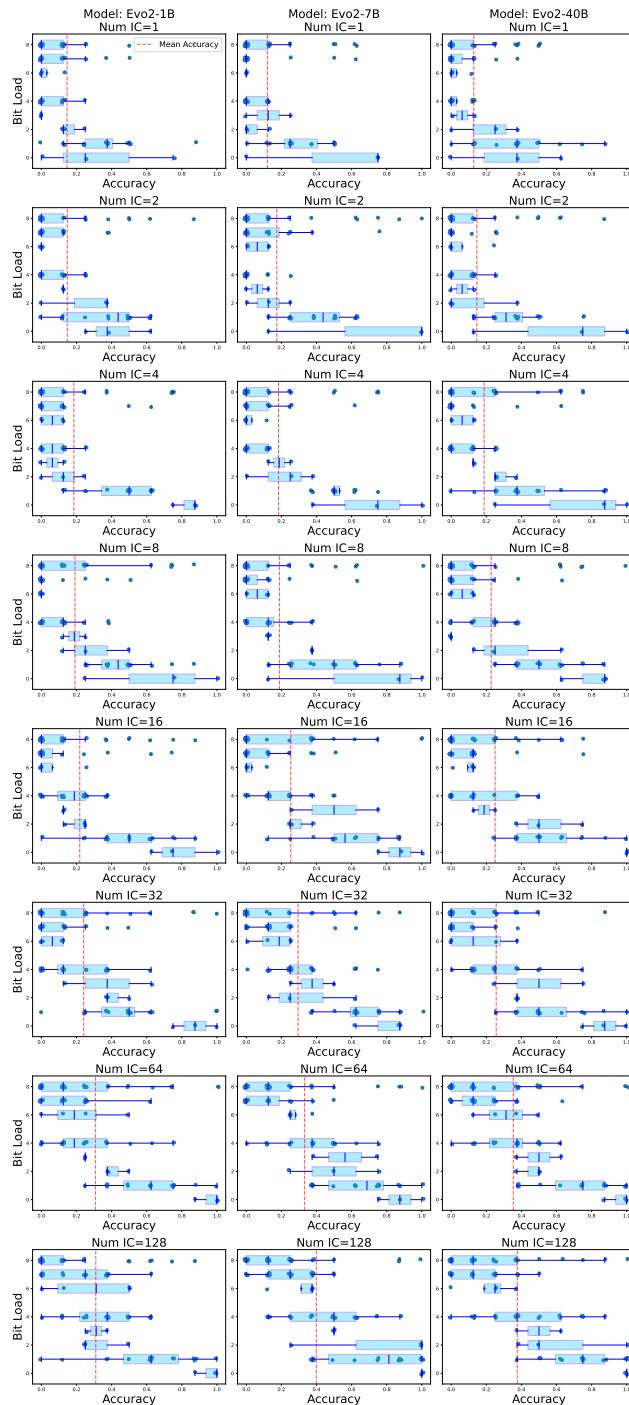


Figure 8: BitLoad vs Accuracy for all Evo models at all ICL shot-numbers. The per-BitLoad distribution of model performance at each shot level shows an uneven affinity of models for solving tasks at different BitLoads. However, all trends support our overall hypotheses.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295



Figure 9: BitLoad vs Accuracy for all Qwen models at all ICL shot-numbers.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

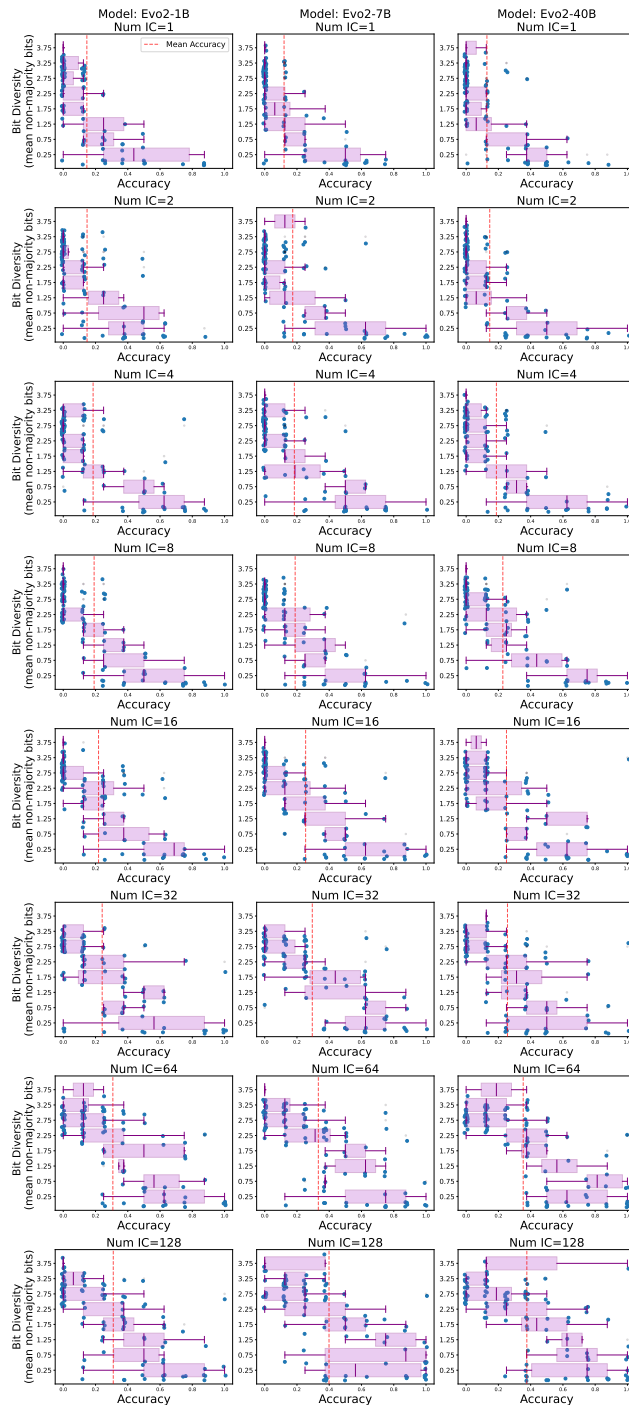


Figure 10: BitDiversity vs Accuracy for all Evo models at all ICL shot-numbers. Model performance follows a more natural pattern of increasing performance with decreasing output entropy. This highlights that it is difficult for models to decipher ICL patterns for high entropy outputs.

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

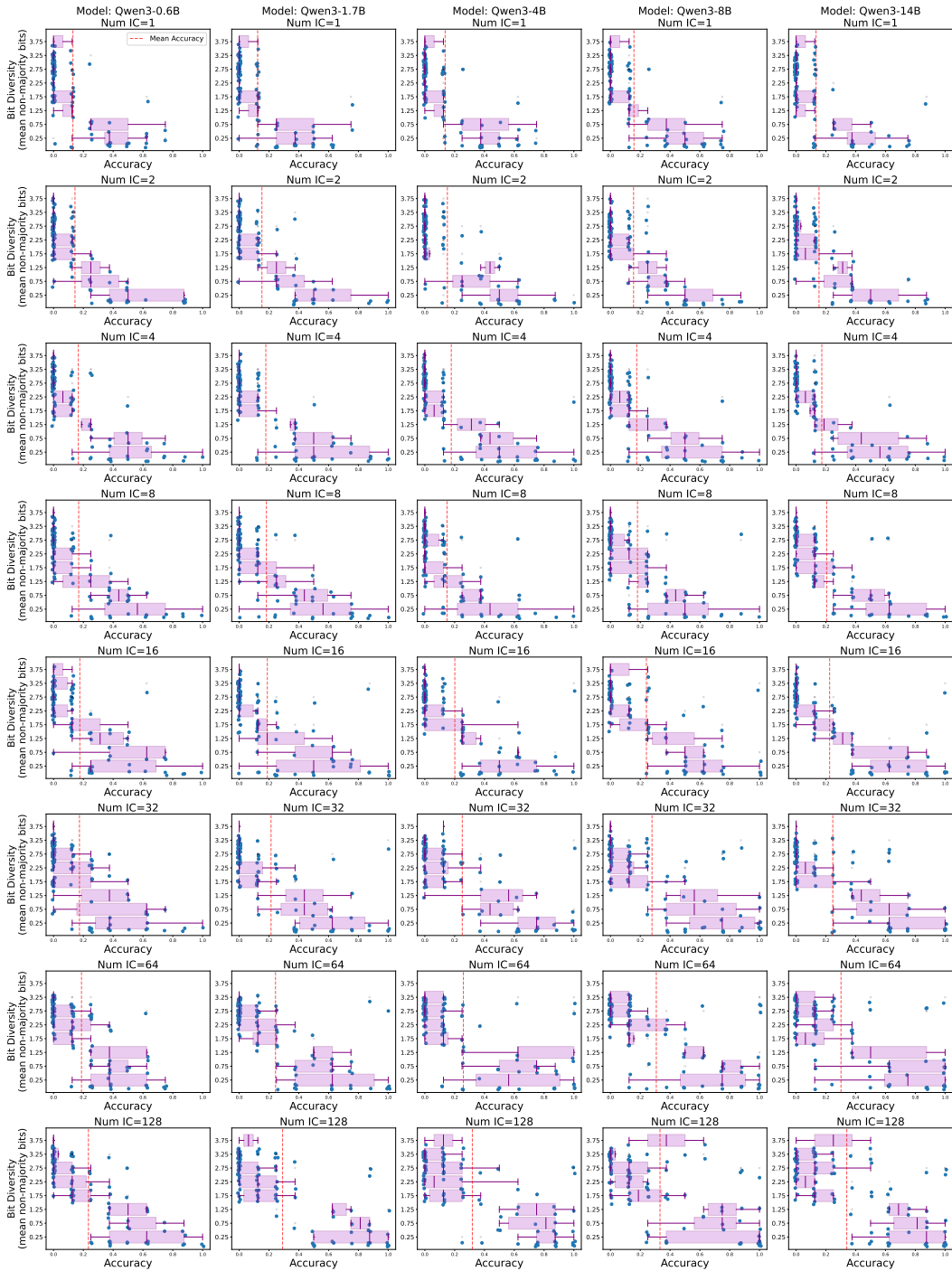


Figure 11: BitDiversity vs Accuracy for all Qwen models at all ICL shot-numbers.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

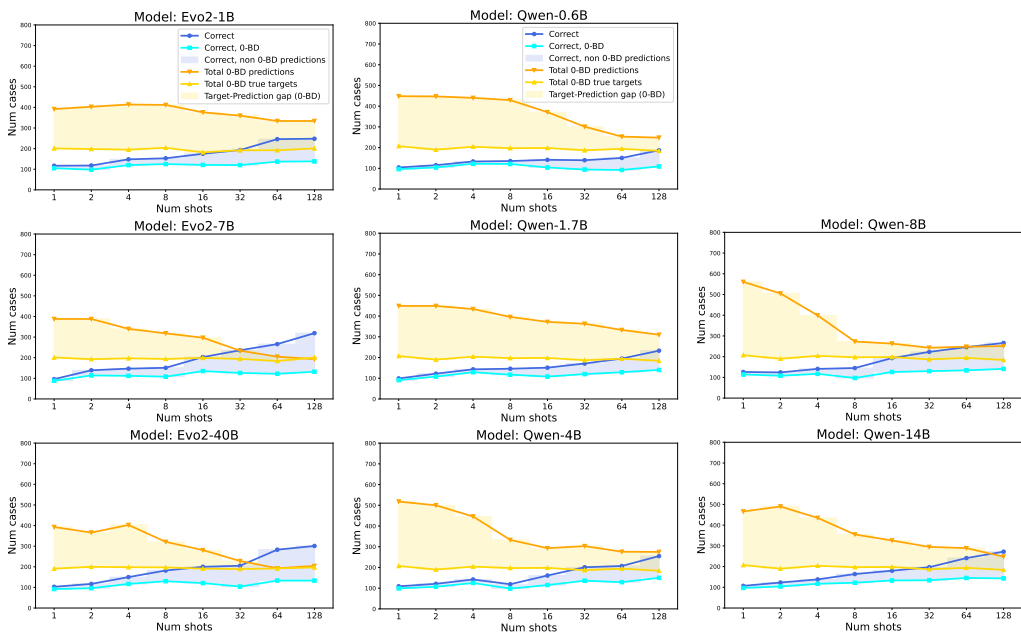


Figure 12: Enumerating cases of 0 BitDiversity. For our defined tasks, around 25% of the true targets have 0 BitDiversity4 (BD). With low-shots, models tend to produce a large number of 0-BD predictions. But this number decreases significantly with increasing shots and tends to match the actual prior value. Similarly, a majority of the baseline (1-shot) performance of models can be explained through 0-BD outputs. With a higher number of shots, the model starts learning higher entropy patterns and presents stronger evidence of ICL.