000	WHY ARE SMALL TRANSFORMERS NOT BETTER?
001	
002	
003	Anonymous authors
004	Paper under double-blind review
005	
006	
007	ABSTRACT
800	
009	
010	The attention mechanism is powering rapid progress in terms of large-scale gener-
011	ative AI. It is conversely exceedingly difficult to find small-scale applications for
012	which allention-based models outperform traditional approaches, such as multi-
013	of 'task switching'. In this framework models work on ongoing taken sequences
014	with the current task being determined by stochastically interseeded control to-
015	kens. We show that standard transformers cannot solve a basic reference model.
016	IARC, which is based on finite-domain arithmetics. The model contains a trivial
017	unary operation, (I: increment the current input), a likewise trivial binary opera-
018	tion, (A: add last two inputs), and reverse copy, (R), a standard memory task. A
019	fourth control token, (C), adds recursive context dependency by modifying cur-
020	rent tasks. Tasks are maintained as long as no new control tokens appears in the
021	prompt, which happens stochastially every 3-9 steps. We show that transformers,
022	LSTM recurrent networks and plain MLPs of similar sizes ( $\sim$ 1.5M parameters)
023	trained transformers containing a modified attention mechanism expressive at
024	tention finding performance levels of around 95%. Our results indicate that the
025	workings of attention can be understood better, and even improved, when compar-
026	ing qualitatively different formulations is a task-switching setting.
027	
028	
029	
030	1 INTRODUCTION
031	
032	An undisputed advantage of the transformer architecture is that memory requirements scale only
033	linearly with context length Vaswani et al. (2017). Compute scales however quadratically, a feature
034	shared with fully-connected multi-layer perceptrons and recurrent networks. It remains an open
035	question whether the success of transformers is due to particular properties of the underlying atten-
036	tion mechanism, or a consequence of the resulting improved size scaling. Alternative models with
037	favorable scaling may be just as good in the later case Gu & Dao (2023). This question received
038	further urgency by the recent observation that MLPs learn in-context a la par with transformers when
039	given the same compute budget long & Pehlevan (2024). It is hence important to study to which
040	extend transformers excel or fail for small-sized applications, viz in a regime where scaling is not yet
041	relevant.
042	Here we work below the scaling regime, typically for a context length $N_c = 24$ . In this regime,
043	in which transformers and classical models have about the same number of adjustable parameters,

Here we work below the scaling regime, typically for a context length  $N_c = 24$ . In this regime, in which transformers and classical models have about the same number of adjustable parameters, we evaluated to which extend several standard models are able to switch task upon the appearance of a suitable control token. Individual tasks work on sequences of encoded numbers, skipping the interseeded control tokens, as explained further in Sect. 3. The resulting evaluation protocol can be varied and/or extended by selecting appropriate tasks, making it a versatile tool for comparative performance tests. For real-world applications, task switching protocols are relevant, e.g., for steering robots by switching between motor primitives Saveriano et al. (2023).

Our results support the notions that small transformers are not generically better than MLPs or recurrent networks. This holds, however, only when attention weights are  $\sim \exp(z_{ij}/\sqrt{N_c})$ , the standard formulation, where  $z_{ij}$  is the dot-product between query and key. We also tested transformers based on expressive attention Gros (2024), for which attention weights are proportional to  $z_{ij}^2/(1 + z_{ij}^2)$ , finding substantial performance leaps.



Figure 1: The IARC task switching evaluation framework, for details see Sect. 3. Shown are results for a LSTM recurrent network (black), a MLP (green), a standard transformer (blue), and a transformer with expressive attention (red). Left: As a function of training epochs, the prediction accuracy (performance). Right: Testing with various combinations of the fundamental tasks, including (I), incrementing the current input, (A), adding the last two inputs and (R), reverse copy. Recursive context dependency is encoded by (C).

## 2 **RELATED WORK**

The task-switching protocol used here is a specification of multi task learning Zhang & Yang (2021); 074 Chen et al. (2024), in the form of ever ongoing sequences of concatenated tasks. Task switching is 075 in particular important for reset-free robotic applications Gupta et al. (2021), e.g., in the context 076 of embodied robotics Kumar et al. (2024), or for large language models Knight & Duan (2024). 077 Related to our approach are testing procedures involving synthetic reasoning tasks Zhang et al. (2022), which have been applied to small models in the form of in-context and global bigrams 079 Bietti et al. (2024). Another example are toy models of superposition, which can be used to study polysemanticity Elhage et al. (2022). It is likewise important to find out which formal languages 081 transformers can express Strobl et al. (2024); Deletang et al. (2022). Equivalently the question arises 082 to which extend large-model scaling Kaplan et al. (2020); Hoffmann et al. (2022), and variantes 083 thereof Naveed et al. (2023); Shen et al. (2024), is retained when models are small Ivgi et al. (2022).

084 085

065

066

067

068

069

070 071

072 073

## 3 IARC TASK SWITCHING FRAMEWORK

087 A considerable number of benchmark tasks for the evaluation of transformer variantes have been 088 developed Liu et al. (2024), however nearly exclusively for large model sizes Tay et al. (2020). The 089 benschmark task introduced here is meant in contrast for the evaluation of small models, typically 090 with a few million parameters or less.

091 The vocabulary consists of a finite set of N numbers, plus a limited number of control tokens, 092 denoted here (I), (A), (R), and (C). Control tokens interseed the sequence of symbols,  $S_t = \{x_{t'} | t' \leq t\}$ 093 t}. The task is to predict the next symbol,  $x_{t+1}$ , but not the occurance of future control tokens. The 094 last control token, or the history of previous control tokens, determines the dependence of  $x_{t+1}$  on  $S_t$ . One has 096

$$x_{t+1}|_{I} = (x_t + 1)\% N, \qquad x_{t+1}|_{A} = (x_t + x_{t-1})\% N$$
 (1)

for increment (I) and addition (A). For N = 10, an example would be

$$2 | 3 | 4 | 7 | 1 | 8 | 9 | 7 | 8 | 9 | 0 | 1$$
(2)

100 Here with two interseeded control tokens, as indicated by the respective superscripts. Including 101 reverse copy (R), one could have

102 103

107

097

098

099

$$2 | 3 | 4 | 7 | 1 | 8 | 8 | 1 | 7 | 4 | 3 | 3 | 4 | 7$$
(3)

104 Tasks remain when (I/A) tokens are followed by (I/A) tokens, with the reverse copy process being 105 restarted by subsequent (R) tokens. The action of the context token (C) depends on the current task, consecutively increasing the increment by one for (I) tasks: 106

When the current task is (A/R), the context token (C) just acts as an additional (A/R) token, as illustrated in (3) for the case of two consecutive (R) token. As a basic protocol regulating the frequency of task switching, we use a  $6\pm3$  setup, which means that the distance between subsequent control tokens is drawn from a flat distribution out of [3,9]. The four control tokens, (A/I/R/C), appear with equal probabilities.

113 114

4 Results

115 116

For all models the context length is  $N_c = 24$ , together with one-hot embedding and a fixed embed-117 ding dimension d = 20. We did set d = N + S, where S is the number of control symbols, compare 118 (1). This implies N = 16 for IARC and N = 17 for an ablated version like IAR. Transformers 119 have L = 12 layers, causal self-attention with four heads and ALiBi positional encoding Press et al. 120 (2021), altogether 1.3M parameters. The multi-layer perceptron (MLP) had L = 16 layers with 121 causal connections, resulting in 1.9M parameters, the LSTM recurrent net was somewhat larger, 122 L = 2 with a size of 3.7M. For all models we did additional runs with increased numbers of layers, 123 finding only marginal improvements. All runs used identical training protocolls, with a batch size of 500. 124

125 Our results are summarized in Fig. 1. The three basic models, LSTM, MLP and the standard trans-126 former show roughly similar performances, both as a function of training compute and when ablating 127 the IARC model. Given the same resources, small standard transformers are not better. A massive 128 outperformance is seen however by transformers based on expressive attention Gros (2024). This is an interesting result, given that expressive attention is intrinsically quadratic in the scalar product be-129 tween queries and keys, making it bi-quadratic in token activities. It also implies that there is room 130 for gains when reformulating the core of the attention mechanism. Finally we want to point out 131 that task-switching frameworks paralleling the one used here are relevant for real world applications 132 taking place in dynamic environments, such as autonomous driving and robot control. 133

- 134 135 ACKNOWLEDGMENTS
- 136 to be added –
- 137 138
- 139 REFERENCES
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36, 2024.
- Shijie Chen, Yu Zhang, and Qiang Yang. Multi-task learning in natural language processing: An overview. ACM Computing Surveys, 56(12):1–32, 2024.
- Gregoire Deletang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, et al. Neural networks and the chomsky hierarchy. In *The Eleventh International Conference on Learning Representations*, 2022.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,
   Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
- Claudius Gros. Reorganizing attention-space geometry with expressive attention. arXiv preprint arXiv:2407.18601, 2024.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv* preprint arXiv:2312.00752, 2023.
- Abhishek Gupta, Justin Yu, Tony Z Zhao, Vikash Kumar, Aaron Rovinsky, Kelvin Xu, Thomas Devlin, and Sergey Levine. Reset-free reinforcement learning via multi-task learning: Learning dexterous manipulation behaviors without human intervention. In 2021 IEEE International Conference on Robotics and Automation (ICRA), pp. 6664–6671. IEEE, 2021.

162 163	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Train-
164	ing compute-optimal large language models. arXiv preprint arXiv:2203.15556, 2022.
165	Maor Ivgi Yair Carmon and Ionathan Berant Scaling laws under the microscope: Predicting
165	transformer performance from small scale experiments. In <i>Findings of the Association for Com</i> -
168	putational Linguistics: EMNLP 2022, pp. 7354–7371, 2022.
169	Jored Kaplan, Sam McCandlish, Tom Hanighan, Tom P. Brown, Banjamin Chass, Dawon Child
170	Scott Grav Alec Radford Jeffrey Wu and Dario Amodei Scaling laws for neural language
171	models. arXiv preprint arXiv:2001.08361, 2020.
172	Parker Knight and Rui Duan. Multi-task learning with summary statistics. Advances in Neural
173 174	Information Processing Systems, 36, 2024.
175	Vikash Kumar, Rutav Shah, Gaoyue Zhou, Vincent Moens, Vittorio Caggiano, Abhishek Gupta, and
176	Aravind Rajeswaran. Robohive: A unified framework for robot learning. Advances in Neural
177	Information Processing Systems, 36, 2024.
178	Nalson E Liu, Kovin Lin, John Howitt, Ashwin Dereniana, Michala Davilazaya, Eshia Detroni, and
179	Percy Liang Lost in the middle: How language models use long contexts. Transactions of the
180	Association for Computational Linguistics, 12:157–173, 2024.
181	
182	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman,
183	Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language
184	models. arxiv preprint arxiv:2307.00433, 2023.
185	Ofir Press, Noah A Smith, and Mike Lewis. Train short, test long: Attention with linear biases
186	enables input length extrapolation. arXiv preprint arXiv:2108.12409, 2021.
187	Matteo Saveriano, Fares I Abu-Dakka, Aliaž Kramberger, and Luka Peternel. Dynamic movement
100	primitives in robotics: A tutorial survey. <i>The International Journal of Robotics Research</i> , 42(13) 1133–1184, 2023.
190	
191	Xuyang Shen, Dong Li, Ruitao Leng, Zhen Qin, Weigao Sun, and Yiran Zhong. Scaling laws for
192 193	linear complexity language models. arXiv preprint arXiv:2406.16690, 2024.
194	Lena Strobl, William Merrill, Gail Weiss, David Chiang, and Dana Angluin. What formal lan-
195	guages can transformers express? a survey. <i>Transactions of the Association for Computational Linguistics</i> , 12:543–561, 2024.
107	
198	Yi Tay, Mostafa Denghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sabastian Budar, and Danald Matzlar. Long range arange to handburght for afficient
199	transformers, arXiv preprint arXiv:2011.04006, 2020.
200	
201	William L Tong and Cengiz Pehlevan. Mlps learn in-context. <i>arXiv preprint arXiv:2405.15618</i> 2024.
202	
203	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
204	Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
205	
206	Vi Thang Arturs Backurs Schastian Ruhack Bonan Eldon Suriva Gunasakar and Tal Was
207	Unveiling transformers with lego: a synthetic reasoning task. arXiv preprint arXiv:2206.04301
208	2022.
209	
210	Yu Zhang and Qiang Yang. A survey on multi-task learning. <i>IEEE transactions on knowledge and</i>
211	<i>uuu engineering</i> , 54(12):5560–5009, 2021.
212	
213	
214	
213	