
An empirical study on the limitation of Transformers in program trace generation

Simeng Sun
simengs@nvidia.com

Abstract

We study Transformers on the task *program trace generation* (PTG), where models produce step-by-step execution traces for synthetic programs. Unlike existing algorithmic problems, PTG externalizes reasoning through long traces where each step is trivial. We train small Transformers with diverse modifications, including alternative position encodings, softmax replacements, hybrid model, and short convolutions. While these models achieve strong in-distribution accuracy, they exhibit systematic failures when generalizing to various factors (e.g., program length, trace steps), though some designs significantly improve generalization.

1 Introduction

We study Transformers on the task *program trace generation* (PTG): given a synthetic program and inputs, models must output the complete step-by-step execution trace. This setup offloads the *internal* reasoning (e.g., sorting, tracking permutations) onto long, structured traces/scratchpads that explicitly store the states of intermediate steps. Following prior work (Sun et al., 2025), we generate simple, executable programs with additional constraints to reduce difficulty at each step. Complex execution thus emerges from composition of simple steps, which enable evaluation of models’ contextual flexibility at executing arbitrary programs, contrasting with fixed-procedure tasks (e.g., large number multiplication) whose complexity is tightly coupled with input size. With synthetic data, we train small Transformers ($\sim 154\text{M}$ parameters) with diverse modifications. Evaluation is limited below training sequence length to avoid confounding of length extrapolation. While all models achieve strong in-distribution accuracy, they struggle to reliably generalize to various factors. Experimental results reveal substantial impacts of architectural choices on out-of-distribution performance, with NAPE (mixing NoPE and ALiBi heads) on average outperforming other modifications. Our results expose limitations in existing Transformers at executing programs step by step, an ability crucial to reasoning tasks and for precise instruction-following in large language models.

2 Program trace generation

The synthetic programs are standalone functions varying code segments among assignment, binary operations,¹ if block, while block, list append or pop. We generate these examples following the methodology of L0-Bench (Sun et al., 2025). More details are provided in Appendix A. With constraints on the simplicity at each step, execution of complex programs becomes sequentially composing basic operations, and the trace steps grow at least linearly with the input size. Therefore, compared to tasks requiring direct production of final answers, we focus on the case where models leverage computations that change with problem complexity. The research question thus shifts from “Can Transformers execute a fixed program *internally* given unseen/longer input?” to “Can Transformers execute simple operations consistently over long sequence given arbitrary programs?”

¹We further restrict arithmetic ops to be only $x + 1$ and $x - 1$, comparison ops to only $==$ and $!=$.

3 Experimental setup

Data & Evaluation metrics. Training and evaluation data are generated while controlling four factors: program length, trace steps, number of variables, and input size. Training data contains programs of 5 to 40 lines, 5 to 103 trace steps, maximum 10 variables, and input list of 5 to 10 integers. Short examples are concatenated together, while longer examples do not exceed 4096 tokens. For out-of-distribution evaluation, we vary each factor independently while keeping others fixed: we test on *longer programs*, *longer traces*, *more variables*, and *longer lists* (i.e., larger input size). Motivation of studying these factors is in Appendix. We report whole-trace accuracy: a trace is correct only if every step matches the ground truth exactly. More details are provided in Appendix A.

Model & Training. We train small Transformers of size $\sim 154\text{M}$ parameters on a mixture of short and medium length programs. In addition to standard RoPE-based Transformers, we also experiment with other position encodings (NoPE, ALiBi, NaPE, Fox, PaTH), softmax replacements (STB, α -entmax), short 1d-convolution (Canon) and hybrid model (SWAN). Description of these methods are provided in Appendix B. We perform fixed-budget model comparison with 20k training steps, which is sufficient for the baseline (NoPE) to achieve near perfect in-distribution performance. We focus on OOD performance to explore the compatibility of various inductive biases with PTG. All models are trained with 4k-token sequences and vocabulary size of 176. We report the average and standard error of each model across 10 random seeds. More details are provided in Appendix C.

4 Results

Table 1: Whole-trace accuracy (%) across test conditions. ID: in-distribution test splits.

	ID (short)	ID (mid)	Longer lists	More variables	Longer traces - 1	Longer traces - 2	Longer programs
NoPE	99.8 \pm 0.2	99.0 \pm 0.4	90.3 \pm 2.5	94.1 \pm 2.0	91.6 \pm 2.7	14.2 \pm 2.5	46.5 \pm 2.2
RoPE	98.0 \pm 1.1	92.6 \pm 3.1	65.5 \pm 4.2	92.9 \pm 3.5	78.5 \pm 5.2	3.9 \pm 1.0	40.8 \pm 3.3
ALiBi	99.4 \pm 0.4	96.8 \pm 1.9	86.2 \pm 3.5	92.1 \pm 1.7	93.1 \pm 1.5	39.5 \pm 4.7	50.8 \pm 2.2
NaPE	99.7 \pm 0.3	98.3 \pm 1.6	93.4 \pm 1.6	98.4 \pm 0.9	97.9 \pm 1.3	93.5 \pm 1.4	57.4 \pm 2.1
Fox	99.2 \pm 0.4	95.8 \pm 2.1	90.0 \pm 3.3	95.8 \pm 1.8	77.9 \pm 1.5	19.2 \pm 1.7	42.9 \pm 1.1
PaTH	100.0 \pm 0.0	99.2 \pm 0.2	89.7 \pm 3.8	90.7 \pm 3.1	88.8 \pm 2.4	30.9 \pm 5.3	52.0 \pm 4.5
SWAN	97.9 \pm 1.1	85.3 \pm 3.6	29.8 \pm 3.9	80.5 \pm 2.6	49.5 \pm 3.9	0.5 \pm 0.2	21.9 \pm 2.5
NoPE+Canon	100.0 \pm 0.0	100.0 \pm 0.0	58.0 \pm 8.1	98.4 \pm 0.5	95.3 \pm 1.4	35.3 \pm 6.9	64.8 \pm 2.4
STB	99.1 \pm 0.5	82.2 \pm 4.4	81.9 \pm 4.4	50.7 \pm 3.8	66.4 \pm 2.7	11.2 \pm 1.5	21.6 \pm 3.7
NoPE+ α -Entmax	100.0 \pm 0.0	99.1 \pm 0.1	41.6 \pm 2.2	98.2 \pm 0.3	98.0 \pm 0.4	27.1 \pm 1.8	49.5 \pm 0.7

Table 1 shows whole-trace accuracy across architectures and test conditions. All models achieve strong in-distribution accuracy but vary significantly in generalization performance. Reliably generating long traces remains challenging across models with large degradation in whole-trace accuracy when increasing output trace steps (*longer traces* - 1 vs. - 2). Models also struggle with more complex/longer instructions in the context: none reaches above 70% accuracy. Dealing with larger entity size (*longer lists*), a factor often studied by existing work, is heavily impacted by positional encoding, with standard RoPE falling behind recent variants, such as Fox and PaTH. Despite its simplicity, NaPE is surprisingly powerful at generating longer traces, being the only variant achieving over 90% accuracy while others are below 40%. Finally, Canon layers improve performance on longer programs (64.8% for NoPE+Canon) but lead to degradation on longer lists (90.3% vs. 58.0%).

5 Discussion

We empirically studied Transformers on program trace generation, a new task requiring precise rule-following over long sequences of simple computational steps. We compare multiple modifications to standard Transformers, including alternative position encoding methods, softmax replacements and novel layer designs. Experimental results show that while these models achieve near-perfect in-distribution performance, they struggle to maintain global accuracy, especially for longer programs and longer traces unseen during training. Architectural choices also significantly impact performance, with NaPE outperforming others. In the future, we plan to extend our work to latest architectures potentially including non-Transformers and MoE models. Finally, we also plan to explore the transduction setting, which may require additional modeling capabilities atop contextual flexibility.

References

- Zeyuan Allen-Zhu. 2025. Physics of Language Models: Part 4.1, Architecture Design and the Magic of Canon Layers. *SSRN Electronic Journal*. <https://ssrn.com/abstract=5240330>.
- Zeyuan Allen-Zhu and Yuanzhi Li. 2025. Physics of language models: Part 1, learning hierarchical language structures.
- Gregoire Deletang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, and Pedro A Ortega. 2023. Neural networks and the chomsky hierarchy. In *The Eleventh International Conference on Learning Representations*.
- Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan, Payel Das, and Siva Reddy. 2023. The impact of positional encoding on length generalization in transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Zhixuan Lin, Evgenii Nikishin, Xu He, and Aaron Courville. 2025. Forgetting transformer: Softmax attention with a forget gate. In *The Thirteenth International Conference on Learning Representations*.
- Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. 2023. Exposing attention glitches with flip-flop language modeling. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- William Merrill, Jackson Petty, and Ashish Sabharwal. 2024. The illusion of state in state-space models. In *Forty-first International Conference on Machine Learning*.
- William Merrill and Ashish Sabharwal. 2024. The expressive power of transformers with chain of thought. In *The Twelfth International Conference on Learning Representations*.
- Ofir Press, Noah Smith, and Mike Lewis. 2022. Train short, test long: Attention with linear biases enables input length extrapolation. In *International Conference on Learning Representations*.
- Krishna C. Puvvada, Faisal Ladhak, Santiago Akle Serrano, Cheng-Ping Hsieh, Shantanu Acharya, Somshubra Majumdar, Fei Jia, Samuel Kriman, Simeng Sun, Dima Rekesh, and Boris Ginsburg. 2025. Swan-gpt: An efficient and scalable approach for long-context language modeling.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2023. Roformer: Enhanced transformer with rotary position embedding.
- Simeng Sun, Cheng-Ping Hsieh, Faisal Ladhak, Erik Arakelyan, Santiago Akle Serano, and Boris Ginsburg. 2025. L0-reasoning bench: Evaluating procedural correctness in language models via simple program execution.
- Shawn Tan, Songlin Yang, Aaron Courville, Rameswar Panda, and Yikang Shen. 2025. Scaling stick-breaking attention: An efficient implementation and in-depth study. In *The Thirteenth International Conference on Learning Representations*.
- Pavlo Vasylenko, Marcos Treviso, and André F. T. Martins. 2025. Long-context generalization with sparse attention.
- Jie Wang, Tao Ji, Yuanbin Wu, Hang Yan, Tao Gui, Qi Zhang, Xuanjing Huang, and Xiaoling Wang. 2024. Length generalization of causal transformers without position encoding.
- Songlin Yang, Yikang Shen, Kaiyue Wen, Shawn Tan, Mayank Mishra, Liliang Ren, Rameswar Panda, and Yoon Kim. 2025. Path attention: Position encoding via accumulating householder transformations.
- Wojciech Zaremba and Ilya Sutskever. 2015. Learning to execute.
- Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Josh Susskind, Samy Bengio, and Preetum Nakkiran. 2023. What algorithms can transformers learn? a study in length generalization.
- Łukasz Kaiser and Ilya Sutskever. 2016. Neural gpus learn algorithms.

A Program Trace Generation

Program generation & Trace format. We follow L0-Bench to generate simple synthetic Python programs using generative grammar. Each program is a standalone Python function with coding segments varying among assignment, if block, while block and list operations (append/pop). We further reduce the difficulty of each line of code by enforcing the binary arithmetic operations to be only increment or decrement by 1; the scope depth is set to 1 to avoid nested loops; list indexing is also excluded. An example of PTG instance is shown in Figure 1. The traces are generated by executing the programs given sampled inputs and post-processing Python execution bytecodes. Specifically, each trace is of format `line_number, var_name: var_value`. For code lines that do not have value updates, only `line_number,` is required.

Data generation & statistics. We list the distribution of training and test data configs in Table 2. Each program is unique. Training and in-distribution (ID) test splits do not have any overlap. During training, with $p = 0.5$, we shift the first line of the program (and correspondingly the trace) by an offset drawn between 1 and 100. Without doing so models fail to achieve non-zero accuracy for the *longer programs* split due to unseen line number in the programs and traces. Preliminary experiments also show the importance of including examples of diverse program/trace lengths to avoid overfitting to certain trace steps; training data are generated by binning short (5-20 lines of code) programs and mid (21-40 lines of code) programs by their trace steps, details listed in Table 3.

Task comparison & Four factors. The task program trace generation essentially forces models to perform reliable sequential operations like a computer (Łukasz Kaiser and Sutskever, 2016; Zaremba and Sutskever, 2015) with *guaranteed correctness*. However, PTG differs from existing algorithmic tasks (such as string reversal, modular addition, permutation composition, or in general recognizing regular languages) (Zhou et al., 2023; Kazemnejad et al., 2023; Deletang et al., 2023; Liu et al., 2023; Merrill et al., 2024; Allen-Zhu and Li, 2025) in two key ways. First, PTG provides the algorithm *explicitly* in the context, while existing tasks require models internalizing the fixed algorithms into parameters during training. In other words, the context window now serves as memory: the beginning part hosting the “instruction memory” and the trace/scratchpad/chain-of-thought (Merrill and Sabharwal, 2024) part “data memory.” Second, the inputs and algorithms are decoupled, i.e., programs now have additional complexity terms independent of input size. We further refine the setup by Sun et al. (2025) to encompass four factors for more targeted evaluation of program execution. Specifically, we focus on:

- **Program length:** Given a fixed max window size (4096 in our experiments), the split *longer programs* tests how models behave when the size of “instruction memory” increases. Longer programs also correspond to more complex procedures which often have higher branching factors.
- **Trace steps:** Given a fixed max window size (4096 in our experiments), the split *longer traces* tests how models behave when the size of “data memory” increases. This can be achieved by increasing the while block size and/or the while loop iteration numbers without increasing the program length.
- **Number of variables:** The split *more variables* tests whether models can handle programs with larger number of input variables, some of which may serve as distractors.
- **Input size:** The split *longer lists* tests if models can apply in-distribution algorithm to larger entity size. This is commonly studied in existing works, such as training models up to n -digit addition and testing on inputs with $[n + 1, k]$ digit inputs.

Scratchpad size. Give a list of T elements and a program of N code lines with k while loop (scope-depth 1) iterations, PTG requires generation of $O(kNT)$ tokens, similar to executing a linear-time algorithm with explicit scratchpad generation. This contrasts with trace-free tasks that test whether models can perform equivalent computation entirely within hidden states.

Table 2: Distribution of training and eval data configs. OOD dimensions are highlighted in blue. Program length in number of lines; Trace length in number of steps; input size in the length of the list.

Split	Program length	Trace steps	Num of variables	Input size	Avg. Out Tokens
TRAIN					
Short program	[5,20]	[5,103]	[1,10]	[5,10]	-
Mid program	[21,40]	[5,103]	[1,10]	[5,10]	-
EVAL					
ID (short)	[5,20]	[59,79]	[1,10]	[5,10]	713.2
ID (mid)	(20,40]	[90,103]	[1,10]	[5,10]	1052.3
Longer lists	(20,40]	[5,103]	[1,10]	[10,15]	828.1
More variables	(20,40]	[5,103]	[10,50]	[5,10]	644.8
Longer traces - 1	(20,40]	[121,142]	[1,10]	[5,10]	1580.5
Longer traces - 2	(20,40]	[172,300]	[1,10]	[5,10]	2402.3
Longer programs	[40,60]	[5,103]	[1,10]	[5,10]	1279.3

Table 3: Data split statistics. **pl**, program length in number of code liens; **tr** trace size in number of trace steps.

Train split	Avg. Token per example	Total num. tokens
pl [5, 20] tr (5, 7]	204.06	8.14E+07
pl [5, 20] tr (7, 8]	208.19	8.31E+07
pl [5, 20] tr (8, 10]	268.01	1.07E+08
pl [5, 20] tr (10, 15]	341.37	1.36E+08
pl [5, 20] tr (15, 40]	548.64	2.19E+08
pl [5, 20] tr (40, 59]	749.28	2.99E+08
pl [5, 20] tr (59, 79]	986.17	3.93E+08
pl [20, 40] tr (14, 25]	615.42	2.46E+08
pl [20, 40] tr (25, 48]	860.15	3.43E+08
pl [20, 40] tr (48, 63]	1057.39	4.22E+08
pl [20, 40] tr (63, 77]	1244.97	4.97E+08
pl [20, 40] tr (77, 90]	1419.38	5.66E+08
pl [20, 40] tr (90, 103]	1600.17	6.38E+08
Sum		4.03E+09

B Model description

C Experimental configurations

Detailed training configurations are provided in Table 5. Transformer variants evaluated in this work are listed in Table 4. For STB and ASentmax, we experiments with both RoPE and NoPE, and report the best for each (STB+RoPE, ASentmax+NoPE). All models are trained with A100 GPUs.

Figure 1: An example of program trace generation ask.

Program:	Output:	
11 def function(a, ai, lst_q, lst_v, lst_na, lst_p, cond_ae, cond_b):	L2, cond_nv: False	L25, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5]
12 cond_nv = 1 == 7	L3	L26, cnter_0: 4
13 if cond_nv:	L6	L27, cond_ai: True
14 lst_v.append(0)	L7, cond_f: True	L28, cnter_0: 5
15 lst_p.append(ai)	L8, lst_q: [4, 6, 4, 6, 9]	L29, lst_q: [4, 6, 4, 6, 9, 1, 5, 5, 5, 5]
16 if cond_b:	L9, lst_v: [0, 3, 8, 3]	L30, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5]
17 cond_f = 5 == 5	L10, lst_p: [6, 7, 1, 4, 1, 7, 7]	L31, cnter_0: 5
18 lst_q.pop()	L11	L32, cond_ai: True
19 lst_v.pop()	L12, lst_v: [0, 3, 8, 3]	L33, cnter_0: 5
20 lst_p.append(7)	L13, ad: 1	L34, cnter_0: 5
21 if cond_f:	L14, lst_v: [0, 3, 8, 3, 5]	L35, cond_nv: False
22 lst_v.pop()	L15, lst_v: [0, 3, 8, 3]	L36, lst_q: [4, 6, 4, 6, 9, 1, 5, 5, 5, 5]
23 ad = 2 + 1	L16, lst_na: [4, 7, 6, 7, 2, 2]	L37, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5]
24 lst_v.append(ai)	L17, lst_q: [4, 6, 4, 6, 9, 1]	L38, cnter_0: 6
25 lst_v.pop()	L18, ae: 1	L39, cond_ai: True
26 lst_na.append(a)	L19, cnter_0: 0	L40, cnter_0: 6
27 lst_q.append(ad)	L20, cond_ai: True	L41, cnter_0: 6
28 ae = ad	L21	L42, cnter_0: 6
29 cnter_0 = 0	L22, L17	L43, cond_nv: False
30 cond_ai = cnter_0 != 9	L23, cond_nv: False	L44, lst_q: [4, 6, 4, 6, 9, 1, 5, 5, 5, 5]
31 while cond_ai:	L24, lst_q: [4, 6, 4, 6, 9, 1, 5]	L45, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5, 5]
32 L1 = 7	L25, lst_p: [6, 7, 1, 4, 1, 7, 7, 5]	L46, cnter_0: 7
33 cond_nv = ae == 9	L26, cnter_0: 1	L47, cond_ai: True
34 lst_v.append(ai)	L27, cond_ai: True	L48, cnter_0: 7
35 lst_p.append(ai)	L28	L49, cnter_0: 7
36 cond_ai = cnter_0 + 1	L29, L17	L50, cond_nv: False
37 lst_p.append(ae)	L30, cond_nv: False	L51, lst_q: [4, 6, 4, 6, 9, 1, 5, 5, 5, 5]
38 na = 6 + 1	L31, lst_q: [4, 6, 4, 6, 9, 1, 5, 5]	L52, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5, 5]
39 if cond_ae:	L32, cnter_0: 2	L53, cnter_0: 8
40 lst_q.pop()	L33, cond_ai: True	L54, cond_ai: True
41 return	L34	L55, L17
42	L35, L17	L56, cond_nv: False
43	L36, cond_nv: False	L57, lst_q: [4, 6, 4, 6, 9, 1, 5, 5, 5, 5]
44	L37, lst_q: [4, 6, 4, 6, 9, 1, 5, 5]	L58, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5, 5]
45	L38, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5]	L59, cnter_0: 9
46	L39, cnter_0: 2	L60, cond_ai: False
47	L40, cond_ai: True	L61, cnter_0: 9
48	L41	L62, cond_ai: False
49	L42, L17	L63, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5, 5, 5]
50	L43, cond_nv: False	L64, na: 7
51	L44, lst_q: [4, 6, 4, 6, 9, 1, 5, 5, 5, 5]	L65
52	L45, lst_p: [6, 7, 1, 4, 1, 7, 7, 5, 5, 5, 5]	L66
53	L46, cnter_0: 3	L67
54	L47, cond_ai: True	L68
55	L48	L69
56	L49	L70
57	L50	L71
58	L51	L72
59	L52	L73
60	L53	L74
61	L54	L75
62	L55	L76
63	L56	L77
64	L57	L78
65	L58	L79
66	L59	L80
67	L60	L81
68	L61	L82
69	L62	L83
70	L63	L84
71	L64	L85
72	L65	L86
73	L66	L87
74	L67	L88
75	L68	L89
76	L69	L90
77	L70	L91
78	L71	L92
79	L72	L93
80	L73	L94
81	L74	L95
82	L75	L96
83	L76	L97
84	L77	L98
85	L78	L99
86	L79	L100
87	L80	L101
88	L81	L102
89	L82	L103
90	L83	L104
91	L84	L105
92	L85	L106
93	L86	L107
94	L87	L108
95	L88	L109
96	L89	L110
97	L90	L111
98	L91	L112
99	L92	L113
100	L93	L114
101	L94	L115
102	L95	L116
103	L96	L117
104	L97	L118
105	L98	L119
106	L99	L120
107	L100	L121
108	L101	L122
109	L102	L123
110	L103	L124
111	L104	L125
112	L105	L126
113	L106	L127
114	L107	L128
115	L108	L129
116	L109	L130
117	L110	L131
118	L111	L132
119	L112	L133
120	L113	L134
121	L114	L135
122	L115	L136
123	L116	L137
124	L117	L138
125	L118	L139
126	L119	L140
127	L120	L141
128	L121	L142
129	L122	L143
130	L123	L144
131	L124	L145
132	L125	L146
133	L126	L147
134	L127	L148
135	L128	L149
136	L129	L150
137	L130	L151
138	L131	L152
139	L132	L153
140	L133	L154
141	L134	L155
142	L135	L156
143	L136	L157
144	L137	L158
145	L138	L159
146	L139	L160
147	L140	L161
148	L141	L162
149	L142	L163
150	L143	L164
151	L144	L165
152	L145	L166
153	L146	L167
154	L147	L168
155	L148	L169
156	L149	L170
157	L150	L171
158	L151	L172
159	L152	L173
160	L153	L174
161	L154	L175
162	L155	L176
163	L156	L177
164	L157	L178
165	L158	L179
166	L159	L180
167	L160	L181
168	L161	L182
169	L162	L183
170	L163	L184
171	L164	L185
172	L165	L186
173	L166	L187
174	L167	L188
175	L168	L189
176	L169	L190
177	L170	L191
178	L171	L192
179	L172	L193
180	L173	L194
181	L174	L195
182	L175	L196
183	L176	L197
184	L177	L198
185	L178	L199
186	L179	L200
187	L180	L201
188	L181	L202
189	L182	L203
190	L183	L204
191	L184	L205
192	L185	L206
193	L186	L207
194	L187	L208
195	L188	L209
196	L189	L210
197	L190	L211
198	L191	L212
199	L192	L213
200	L193	L214
201	L194	L215
202	L195	L216
203	L196	L217
204	L197	L218
205	L198	L219
206	L199	L220
207	L200	L221
208	L201	L222
209	L202	L223
210	L203	L224
211	L204	L225
212	L205	L226
213	L206	L227
214	L207	L228
215	L208	L229
216	L209	L230
217	L210	L231
218	L211	L232
219	L212	L233
220	L213	L234
221	L214	L235
222	L215	L236
223	L216	L237
224	L217	L238
225	L218	L239
226	L219	L240
227	L220	L241
228	L221	L242
229	L222	L243
230	L223	L244
231	L224	L245
232	L225	L246
233	L226	L247
234	L227	L248
235	L228	L249
236	L229	L250
237	L230	L251
238	L231	L252
239	L232	L253
240	L233	L254
241	L234	L255
242	L235	L256
243	L236	L257
244	L237	L258
245	L238	L259
246	L239	L260
247	L240	L261
248	L241	L262
249	L242	L263
250	L243	L264
251	L244	L265
252	L245	L266
253	L246	L267
254	L247	L268
255	L248	L269
256	L249	L270
257	L250	L271
258	L251	L272
259	L252	L273
260	L253	L274
261	L254	L275
262	L255	L276
263	L256	L277
264	L257	L278
265	L258	L279
266	L259	L280
267	L260	L281
268	L261	L282
269	L262	L283
270	L263	L284
271	L264	L285
272	L265	L286
273	L266	L287
274	L267	L288
275	L268	L289
276	L269	L290
277	L270	L291
278	L271	L292
279	L272	L293
280	L273	L294
281	L274	L295
282	L275	L296
283	L276	L297
284	L277	L298
285	L278	L299
286	L279	L300
287	L280	L301
288	L281	L302
289	L282	L303
290	L283	L304
291	L284	L305
292	L285	L306
293	L286	L307
294	L287	L308
295	L288	L309
296	L289	L310
297	L290	L311
298	L291	L312
299	L292	L313
300	L293	L314
301	L294	L315
302	L295	L316
303	L296	L317
304	L297	L318
305	L298	L319
306	L299	L320
307	L300	L321
308	L301	L322
309	L302	L323
310	L303	L324
311	L304	L325
312	L305	L326
313	L306	L327
314	L307	L328
315	L308	L329
316	L309	L330
317	L310	L331
318	L311	L332
319	L312	L333
320	L313	L334
321	L314	L335
322	L315	L336
323	L316	L337
324	L317	L338
325	L318	L339
326	L319	L340
327	L320	L341
328	L321	L342
329	L322	L343
330	L323	L344
331	L324	L345
332	L325	L346
333	L326	L347
334	L327	L348
335	L328	L349
336	L329	L350
337	L330	L351
338	L331	L352
339	L332	L353
340	L333	L354
341	L334	L355
342	L335	L356
343	L336	L357
344	L337	L358
345	L338	L359
346	L339	L360
347	L340	L361
348	L341	L362
349	L342	L363
350	L343	L364
351	L344	L365
352	L345	L366
353	L346	L367
354	L347	L368
355	L348	L369
356	L349	L370
357	L350	L371
358	L351	L372
359	L352	L373
360	L353	L374
361	L354	L375
362	L355	L376
363	L356	L377
364	L357	L378
365	L358	L379
366	L359	L380
367	L360	L381
368	L361	L382
369	L362	L383
370	L363	L384
371	L364	L385
372	L365	L386
373	L366	L387
374	L367	L388
375	L368	L389
376	L369	L390
377	L370	L391
378	L371	L392
379	L372	L393
380	L373	L394
381	L374	L395
382	L375	L396
383	L376	L397
384	L377	L398
385	L378	L399
386	L379	L400
387	L380	L401
388		

Figure 2: NoPE + Canon accuracy across evaluation splits throughout training (average over 10 random seeds). In-distribution performance saturates much earlier than OOD splits. The learning speed for longer programs is much slower than other splits.

