How does Multi-Task Training Affect Transformer In-Context Capabilities? Investigations with Function Classes

Anonymous ACL submission

Abstract

Large language models (LLM) have recently 001 shown the extraordinary ability to perform unseen tasks based on few-shot examples pro-004 vided as text, also known as in-context learning (ICL). While recent work has attempted to understand the mechanisms driving ICL, few have explored training strategies that incen-007 800 tivize these models to generalize to multiple tasks. Multi-task learning (MTL) for generalist models is a promising direction that offers 011 transfer learning potential, enabling large parameterized models to be trained from simpler, 012 related tasks. In this work, we investigate the combination of MTL with ICL to build models that efficiently learn tasks while being robust to out-of-distribution examples. We propose several effective curriculum learning strategies that 017 018 allow ICL models to achieve higher data efficiency and more stable convergence. Our exper-019 iments ¹ reveal that ICL models can effectively learn difficult tasks by training on progressively harder tasks while mixing in prior tasks, denoted as mixed curriculum in this work. 023

1 Introduction

034

Recently, the emergence of in-context-learning capabilities in LLMs has revolutionized the field of NLP (Wei et al., 2022a). By pre-training with nextword predictions, these models can be prompted with few-shot examples and make accurate incontext predictions during inference (Brown et al., 2020). The ICL capability demonstrated even by smaller Transformer models presents an alternative way to understand LLMs (Dong et al., 2023; Li et al., 2023a; Lu et al., 2023). To empirically understand this phenomenon, Garg et al. (2022) focus on learning a single function class in-context by a Transformer model. Their model achieves competitive normalized mean-squared error (MSE) compared to the optimal ordinary least squares estimator when performing in-context linear regression. Nevertheless, these models sometimes fail to converge and often struggle to generalize to more challenging function classes. While the follow-up studies (Akyürek et al., 2023; Von Oswald et al., 2023; Yang et al., 2023) have extensively analyzed how these models conduct ICL, little work exists exploring how training on *multiple function classes* can enable transformer models to generalize and perform ICL more efficiently. As we believe that these generalist models are typically designed to perform multiple tasks, there is a pressing need to study the multi-task ICL capability of these models, which is still missing in the literature. 040

041

042

045

046

047

048

051

052

054

056

057

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

079

From prior multi-task training (MTR) studies (Zhang et al., 2023; Ruder, 2017; Weiss et al., 2016), models can be trained on multiple related tasks to improve their performance on individual tasks. Despite its popularity, multi-task learning has been difficult to understand in Transformer models when trained on natural language, most likely due to the difficulty of ranking and scheduling language tasks (Crawshaw, 2020). However, the newly introduced framework of learning function classes in context (Garg et al., 2022) provides an easier way to study this multi-task training paradigm. For example, for a polynomial function class, its difficulty can be scaled by changing its degree (e.g., linear to quadratic), or changing the input distribution (e.g., standard Gaussian to Gaussian distribution with decaying eigenvalues). This allows us to understand the ICL capabilities to transfer across similar function classes from multitask training. Motivated by this, we conduct a systematic analysis by training a Transformer on varying function class families and input distributions in a multi-task manner to examine if the same principles from MTR carry over into ICL.

During training, we explore different curriculum learning strategies to schedule the ICL tasks of multiple function classes: (*mixed*, *sequential*, *random*)

¹Code and models will be released on acceptance

(§2.3). For benchmarking, we train another set of models only on a single function class family following (Garg et al., 2022). We quantitatively and qualitatively compare our models trained with and without curriculum across all tasks and analyze the normalized mean-squared error (MSE) and attention matrices (§3). Our experiments show that curriculum learning is more data-efficient, achieving comparable performance to single-task models using only 1/9 of the training data. These curriculum models can also obtain an optimal MSE in function classes where none of the single-task models converge.

2 Methods

081

087

097

100

101

103

104

105

106

107

108

109

110

111

112 113

114

115

116

117

118

119

121

122

123

124

2.1 Problem Definition

Following Garg et al. (2022), we define the problem of ICL as passing in an n-shot sequence $(x_1, f(x_1), x_2, f(x_2), \ldots, x_n, f(x_n), x_{n+1})$ to the Transformer and generating an output for $f(x_{n+1})$, where the examples have not been seen during training. We refer to this n-shot prediction problem, where input is given in pairs, as in-context learning.

We consider a data-generating process where ddimensional covariates are drawn $x_i \sim \mathcal{D}_x$ and a function $f \sim \mathcal{F}$ where \mathcal{D}_x is any arbitrary distribution and \mathcal{F} is the class of functions related to single-index normalized Hermite polynomials.

2.2 Tasks

We explore two types of separate tasks: learning a function class and learning a data distribution (in Appendix). We consider a single-index function:

$$f(x) = \varphi(\langle x, w \rangle).$$

Function Class Learning We look at the class of functions derived from normalized probabilist's Hermite polynomial, $\frac{1}{\sqrt{n!}}He_n(x)$ which satisfies orthogonality. This is useful as it guarantees that the function values of all tasks are uncorrelated. For each task, we sample w uniformly from the unit sphere. We define K = 3 polynomial function classes as follows: denoting $t = \langle x, w \rangle$, we pick $\varphi \in \{\varphi_{\text{linear}}, \varphi_{\text{quadratic}}, \varphi_{\text{cubic}}\}$ for three function classes $\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3$ respectively.

$$\varphi_{\text{linear}}(t) = t,$$

25
$$\varphi_{\text{quadratic}}(t) = \frac{1}{\sqrt{2}}(t + \frac{1}{\sqrt{2}}(t^2 - 1))$$

26
$$\varphi_{\text{cubic}}(t) = \frac{1}{\sqrt{3}}(t + \frac{1}{\sqrt{2}}(t^2 - 1) + \frac{1}{\sqrt{6}}(t^3 - 3t))$$

2.3 Curriculum Learning

We define the total training steps T to be 500,000, where the *t*-th training step ranges from t = 1, 2, ..., T. Our curriculum learning strategy (Sequential, Mixed, Random) is used to allocate our Ktasks across training time. In this paper, we explore K = 3 function classes defined earlier.

Sequential Curriculum Given T total training steps, we split the training steps into K partitions. Within the k-th partition of training steps, we train the model on learning a function from the k-th function class, in order of increasing difficulty:

$$f \sim \begin{cases} \mathcal{F}_1 & 0 \le t \le \frac{T}{3} \\ \mathcal{F}_2 & \frac{n}{3} \le t \le \frac{2T}{3} \\ \mathcal{F}_3 & \frac{2T}{3} \le t \end{cases}$$
 139

127

128

129

130

131

132

134

136

137

138

147

148

149

151

Mixed CurriculumGiven T total training steps,140we again split the training steps into K partitions.141Let ξ be (uniformly) drawn from $\{1, 2\}$ and ζ be142(uniformly) drawn from $\{1, 2, 3\}$. We select tasks143from the previous k partitions with equal probabil-144ity (1 denotes the indicator function):145

$$f \sim \begin{cases} \mathcal{F}_{1} & 0 \le t \le \frac{T}{3} \\ \sum_{s=1}^{2} \mathbf{1}(\xi = s) \mathcal{F}_{s} & \frac{T}{3} \le t \le \frac{2T}{3} \\ \sum_{s=1}^{3} \mathbf{1}(\zeta = s) \mathcal{F}_{s} & \frac{2n}{3} \le t \end{cases}$$
 146

Random Curriculum At each training step t, randomly sample from the list of K tasks with distribution equal probability:

$$f \sim \sum_{s=1}^{3} \mathbf{1}(\zeta = s) \mathcal{F}_s, \quad 0 \le t \le T$$
 150

2.4 Attention Analysis

To understand how single and multi-task models 152 learn, we analyze the Transformer's self-attention 153 weights. Specifically, we mask out the attention 154 matrices for each head to filter out the self-attention 155 scores between each $f(x_i)$ token and its corre-156 sponding x_i token. To summarize the head's in-157 clination to attend to previous tokens, we aggregate these scores by taking the mean across all $f(x_i)$ 159 tokens. We then do this for all attention heads in 160 all layers and plot this as a head-by-layer heatmap. 161 We define a "retrospective head" as an attention 162 head that has a lighter value in the heatmap, indi-163 cating that this specific head learns to attend to the 164 previous input token when constructing a represen-165 tation for the current token, a natural pattern that encourages understanding of the input pairs. 167

3 Results



Figure 1: Comparison of moving average of all three curriculum learning strategies when evaluated on quadratic function class dataset during test time. The mixed curriculum is the only model that is able to achieve an accurate normalized MSE. The random curriculum performs comparatively worse, whereas the sequential curriculum performs substantially worse (y-axis is limited in order for mixed and random curricula to be differentiated).

Curriculum Strategy Comparison Figure 1 169 shows that the mixed curriculum outperforms both 170 the random and sequential curriculum when eval-171 uating all models on a quadratic function class 172 dataset during test time. We find that the mixed 173 curriculum strategy provides the most benefit to-174 wards learning multiple tasks. This is further vali-175 dated in Supplementary Figure 5, which shows that 176 the mixed curriculum is most stable over all tasks, achieving an accurate solution after sufficient few-178 shot examples (20/80/90-shot examples for Lin-179 ear/Quadratic/Cubic respectively). We hypothesize 180 this is due to stable periods of training, where the 181 model is able to adapt to the given function class, whereas the random curriculum does not have such 183 184 a schedule. Additionally, mixed curriculum likely outperforms sequential curriculum because includ-185 ing tasks from previous training blocks mitigates catastrophic forgetting (Zhai et al., 2023). There-187 fore, we stick with the mixed curriculum model in the following experiments.

Qualitative Attention Analysis Figure 2 dis-190 plays how masking 7 retrospective heads (as de-191 fined in §2.4) causes a significant increase in 192 normalized MSE compared to 7 non-retrospective 193 194 heads in the mixed curriculum model. Using our attention analysis in Supplementary Figure 4, we 195 identify retrospective heads as those with yellow 196 values, whereas non-retrospective heads are highlighted with dark purple values. This supports the 198

theory that specific heads may be reasonable for the ICL abilities of these models (Olsson et al., 2022). Additionally, we find that the same attention heads have high scores across related tasks, indicating that these models are conducting approximations, rather than learning the true tasks.

199

200

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

246

247

Curriculum Learning Convergence Figure 3 reveals 60% of mixed curriculum models converge, whereas 0% of the single-task models trained on quadratic function classes converge. Specifically, these models do not achieve optimal (below 1) normalized MSE during training time and at test time. We believe curriculum learning aids in this task, as we allow the model to warm up the training with the objective (calculate f(x) from x) on easier tasks. In contrast, the poor performance of the single-task models may be explained by their cryptic attention patterns (Supp Fig. 2). These findings help us understand how curriculum learning can be used to learn difficult function classes that are otherwise unlearable by single-task models.

Curriculum Learning Data Efficiency Figure 4 illustrates the performance of a single-task model and a mixed curriculum model during training when evaluated on a cubic function class validation dataset. Our experiments uncover that the mixed curriculum model can improve data efficiency, learning harder tasks with fewer examples. The mixed curriculum model is pre-trained on 1/9 of the training examples seen by the single-task cubic model, yet the mixed curriculum model has better performance on the validation set. Pulling from qualitative attention analysis, we hypothesize that the mixed curriculum model is able to use its approximate understanding of the linear and quadratic function classes to improve the initial normalized MSE of a cubic function class. This explains why the cubic model starts at 450 normalized MSE, whereas the mixed model starts at 200 normalized MSE. When analyzing both models at test time (Supp Fig. 3, 5) the mixed model has comparable performance to the single-task cubic model. These findings suggest that curriculum learning can assist data efficiency by making use of transfer learning from easier tasks.

4 Discussion

In this paper, we examine how different curriculum 245 learning strategies affect Transformer's in-context learning capability. We first introduce the three



Figure 2: Masking retrospective heads (*bottom row*) causes significant increase in normalized MSE compared to non-retrospective heads (*top row*) in the mixed curriculum model.



Figure 3: Comparison of the moving average of five different seeded single-task (*blue-purple*) and curriculum models (*orange-red*) evaluated on a quadratic function class dataset during test time. Mixed models are able to learn quadratic function classes whereas the single task models are unable to, indicated by the spikes and upward trend in normalized MSE.

types of curriculum (mixed, random, sequential) along with the two tasks (function class learning and distribution learning). We then compare these curriculum models against models that we only train on a single specific task and evaluate them across related tasks. This reveals that the mixed curriculum provides the best results, as well as increases data efficiency and model convergence.

251

255



Figure 4: Comparison of moving average normalized MSE of a single-task cubic model to a mixed curriculum model during training. Data points are generated by evaluating the model on a separate validation set of cubic function examples. The mixed curriculum model is initialized with a checkpoint trained on linear and quadratic function examples, while the single-task model is initialized with random weights.

Through our analysis of attention, we show that these multi-task models had high attention scores across related tasks in the same heads, and that if we mask these heads during test time, the accuracy of these models drop drastically, indicating that specific heads are responsible for ICL. These results provide an important insight into how we can better pre-train LLMs to in-context-learn efficiently.

4

264 Limitations

Our work investigates ICL on standard function classes which can be mathematically defined, however it may be difficult to extend our work to natural 267 language tasks as they are hard to define. The extensibility of our work to natural language tasks therefore remains an open question. We make use of three well-known scheduling methods, how-271 ever, more effective curriculum learning scheduling strategies should be investigated. We work with a relatively small model, thus our results may not 274 be transferable to larger models such as Llama-2 or GPT-4 and we work with noiseless data which may inflate the accuracy. Lastly, we acknowledge 277 that in-context learning can be inconsistent (models only learn approximations for tasks and have 279 varying performance across seeds) and should not be used in high-risk situations.

References

287

289

290

291

296

301

303

305

306

310

311

312

313

314

315 316

- Samira Abnar and Willem Zuidema. 2020. Quantifying attention flow in transformers. *arXiv preprint arXiv:2005.00928.*
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. 2023. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. 2023. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.
- Center for High Throughput Computing. 2006. Center for high throughput computing.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of BERT's attention. In Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pages 276–286, Florence, Italy. Association for Computational Linguistics.

Michael Crawshaw. 2020. Multi-task learning with deep neural networks: A survey.

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

339

340

341

342

343

345

346

347

348

349

350

351

352

353

354

355

356

357

360

361

362

363

364

365

366

367

368

369

370

371

372

373

- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. A survey on in-context learning.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*. Https://transformercircuits.pub/2021/framework/index.html.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. 2022. What can transformers learn incontext? a case study of simple function classes. In *Advances in Neural Information Processing Systems*, volume 35, pages 30583–30598. Curran Associates, Inc.
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. 2023a. Transformers as algorithms: Generalization and stability in in-context learning. In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages 19565–19594. PMLR.
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. 2023b. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, pages 19565–19594. PMLR.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pretrain, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *ACM Computing Surveys*, 55.
- Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych. 2023. Are emergent abilities in large language models just in-context learning?
- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2022a. MetaICL: Learning to learn in context. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2791–2809, Seattle, United States. Association for Computational Linguistics.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022b. Rethinking the role of demonstrations: What makes in-context learning work? In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2022. In-context learning and induction heads. *Transformer Circuits Thread*. Https://transformer-circuits.pub/2022/incontext-learning-and-induction-heads/index.html.

375

376

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420 421

422

423

424

425

426

427

428

429

- Sebastian Ruder. 2017. An overview of multi-task learning in deep neural networks.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
 - Jesse Vig and Yonatan Belinkov. 2019. Analyzing the structure of attention in a transformer language model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 63–76, Florence, Italy. Association for Computational Linguistics.
 - Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, Joao Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2023. Transformers learn in-context by gradient descent. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 35151–35174. PMLR.
 - Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. Emergent abilities of large language models.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. 2023. Larger language models do in-context learning differently.
- Karl Weiss, Taghi M. Khoshgoftaar, and DingDing Wang. 2016. A survey of transfer learning. *Journal of Big Data*, 3(1):9.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2022. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*.

Steve Yadlowsky, Lyric Doshi, and Nilesh Tripuraneni. 2023. Pretraining data mixtures enable narrow model selection capabilities in transformer models. 430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

- Jiaxi Yang, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. 2023. Iterative forward tuning boosts in-context learning in language models.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*.
- Fan Yin, Jesse Vig, Philippe Laban, Shafiq Joty, Caiming Xiong, and Chien-Sheng Wu. 2023. Did you read the instructions? rethinking the effectiveness of task definitions in instruction learning. In *Proceedings* of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3063–3079, Toronto, Canada. Association for Computational Linguistics.
- Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. 2023. Investigating the catastrophic forgetting in multimodal large language models.
- Zhihan Zhang, Wenhao Yu, Mengxia Yu, Zhichun Guo, and Meng Jiang. 2023. A Survey of Multi-task Learning in Natural Language Processing: Regarding Task Relatedness and Training Methods. In *Proceedings* of the 17th Conference of the European Chapter of the Association for Computational Linguistics.

A Experimental Settings

Appendix

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

487

488

489

490

491

492

493

494

495

496

497

498

We train on the GPT-2 model (22.4 million parameters) with 12 heads, 8 layers, and an embedding size of 256 over 500,000 steps, where each batch size is 64. Each batch consists of 100 $(x_i, f(x_i))$ pairs (we expect higher order polynomials to require more in-context examples to converge). During training, we evaluate each model every 2,000 steps on a validation set with size of 32,000. During test time evaluation, we evaluate each model on 64 randomly selected examples. We train GPT-2 using a A100-SXM4-80GB provided by the Center for High Throughput Computing (2006).

B Distribution Learning

474 In addition to different function classes, we explore training data generated from different distributions, 475 given that recent literature has shown that these 476 models do not perform well under distributional 477 shifts (Garg et al., 2022; Yadlowsky et al., 2023). 478 Particularly, we sample inputs x_i from (i) Gaussian 479 distributions, (ii) skewed Gaussian distributions 480 (decaying eigenvalues), and (iii) student-t distri-481 butions (df=4). Attention matrices (Supp Fig. 6 482 and 8) and normalized MSE (Supp Fig. 7 and 483 9) across tasks may be found for both single-task 484 485 and curriculum-based models in the Supplementary Materials. 486

C Instruction Prompting

We explore two sets of instruction prompting architectures: one-hot encoded vectors and preset instruction vectors. The goal of instruction prompting was to evaluate whether our objective could benefit from instruction prompting the way language translation or other NLP tasks do.

C.1 One Hot Encoded Vectors (OHEI)

After generating our $(x_i, f(x_i))$ pairs, we append a single one hot encoded vector, p, to the beginning of the sequence, with the one hot encoding corresponds to the "task":

499
$$p = \begin{cases} p_0 = 1 & \varphi = \varphi_1 \\ p_1 = 1 & \varphi = \varphi_2 \\ p_2 = 1 & \varphi = \varphi_3 \end{cases}$$

We then apply a linear transformation to transform the concatenation into the dimension of our transformer, 256.

C.2 Preset Instruction Vectors (PI)

After we use a linear transformation to transform our $(x_i, f(x_i))$ pairs to the input dimension of our transformer,256, we append a unique vector, $p \sim \mathcal{N}(0, I_d)$, that has been sampled from an isotropic Gaussian distribution. This "instruction vector" remains constant throughout the training of all models, but remains different for each of the different tasks.

C.3 Instruction prompting remains unclear

Supplementary figure 1 shows the comparison of a mixed curriculum model with no instruction prompting, to the two instruction prompting architectures listed above, evaluated over all function class tasks. Applications of the one hot encoded instruction (OHEI) vector to the mixed model causes minimal improvement, whereas application of the preset instruction (PI) vector to the mixed model worsens model performance in the quadratic and cubic function class evaluation during test time. We believe the former has minimal effect in performance as the one-hot encoded vectors may just be seen as noise, whereas the latter most likely worsens the ability of the model to learn the task as it may be seen as an extreme version of noise (i.e. it disrupts the flow of x_i , $f(x_i)$ confusing the model). Overall, we believe that instruction tokens may not be tractable in this setting due to the difficult of learning a 20-dimensional instruction.

D Related Work

In-context Learning In-context learning has been around for a few years now (Dong et al., 2023), and many papers have analyzed in-context learning with natural language (Min et al., 2022b; Xie et al., 2022; Min et al., 2022a). It wasn't until Garg et al. (2022) that papers started analyzing ICL through the paradigm of function class learning. Garg et al. (2022) released a paper that showed that transformers could learn linear regression close to the OLS estimator, and other more complex function classes with respectable accuracy. However, they found that some function classes were hard to learn or did not converge (e.g. skewed gaussian). Yadlowsky et al. (2023) looked at a framework similar to Garg et al. (2022), where they explored training models on a mixture of function classes

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

500 501

502

- however they do not explore curriculum strategies. 549 Many papers also explore how in-context learning 550 works, with current literature pointing to it being 551 a fuzzy gradient descent, (Akyürek et al., 2023; 552 Von Oswald et al., 2023; Yang et al., 2023). Ad-553 ditional theoretical work studies how transformers 554 can implement near-optimal regression algorithms 555 and describe stability conditions for ICL (Li et al., 2023b).
- Attention Analysis Until Vaswani et al. (2017), 558 attention was not widely used in neural networks, however Transformers have revolutionized our capabilities of performing tasks in a variety of fields. 561 Given the power behind attention, we wanted to 562 figure out a way to analyze it, similar to what has 563 been done in previous work (Clark et al., 2019). 564 (Olsson et al., 2022; Elhage et al., 2021) found that specific heads, specified as "induction heads", were 566 responsible for the in-context-learning ability of transformers, both in large and small transformers. To measure this, they created their own metric. 569 Interested in seeing if specific heads attended to 570 specific tasks in a multi-task framework, we de-571 cided to visualize the attention matrix. (Vig and Belinkov, 2019) showed a simple and easy way to 573 visualize attention, that was interpretable. We used 574 575 this as a proxy to develop our own analyses of the 576 attention matrices. Other recent work focuses on summarizing attention flow through transformer 577 models from input embeddings to later layers with attention rollout (Abnar and Zuidema, 2020).
- **Instruction prompting** Prompting has been 580 widely used in natural language tasks to improve 581 accuracy and tends to be robust to variations dur-582 ing test time (Liu et al., 2023). Wei et al. (2023) showed that models of different architectures responded differently to instruction tokens, with the 585 formatting of the instruction effecting multi-task 586 settings. Yin et al. (2023) showed that providing key information in tasks in a common format improved the ability of the model to learn the task. Recently frameworks have emerged which 590 prompt LLMs with intermediate reasoning steps 591 to elicit better reasoning capabilities (Wei et al., 592 2022b), known as Chain of Thought (CoT) prompt-593 ing. (Besta et al., 2023) and (Yao et al., 2023) ex-594 tend CoT prompting to consider multiple reasoning 595 paths to improve performance. Future work may 596 consider using these methods to improve in-context 597 learning in the multi-task setting.

Supplementary Materials



Figure 1: Normalized MSE over number of in-context examples for mixed curriculum model, mixed curriculum model with one hot encoded instruction (OHEI) vector and mixed curriculum model with preset instruction (PI) vector. Solid line represents the moving average (window=10) whereas the dashed line is the actual value. Scientific notation is used for the y-axis.



Figure 2: Attention analysis for single-task function learning models.



Figure 3: Normalized MSE over number of in-context examples for single-task function learning models. Solid line represents the moving average (window=10) whereas the dashed line is the actual value. Scientific notation is used for the y-axis.



Figure 4: Attention analysis for curriculum function learning model.



Figure 5: Normalized MSE over number of in-context examples for curriculum function learning models. Solid line represents the moving average (window = 10) whereas the dashed line is the actual value. Scientific notation is used for the y-axis.



Figure 6: Attention analysis for single-task distribution learning models.



Figure 7: Normalized MSE over number of in-context examples for single-task distribution learning models. Solid line represents the moving average (window = 10) whereas the dashed line is the actual value. Scientific notation is used for the y-axis.



Figure 8: Attention analysis for curriculum distribution learning.



Figure 9: Normalized MSE over number of in-context examples for curriculum distribution learning models. Solid line represents the moving average (window = 10) whereas the dashed line is the actual value. Scientific notation is used for the y-axis.