
Teaching Arithmetic to Small Transformers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models like GPT-4 exhibit emergent capabilities across general-
2 purpose tasks, such as basic arithmetic, when trained on extensive text data, even
3 though these tasks are not explicitly encoded by the unsupervised, next-token
4 prediction objective. This study investigates how even small transformers, trained
5 from random initialization, can efficiently learn arithmetic operations such as addi-
6 tion, multiplication, and elementary functions like square root, using the next-token
7 prediction objective. We first demonstrate that conventional training data is not the
8 most effective for arithmetic learning, and simple formatting changes can signif-
9 icantly improve accuracy. This leads to sharp phase transitions as a function of
10 training data scale, which, in some cases, can be explained through connections
11 to low-rank matrix completion. Building on prior work, we then train on chain-
12 of-thought style data that includes intermediate step results. Even in the complete
13 absence of pretraining, this approach significantly and simultaneously improves
14 accuracy, sample complexity, and convergence speed. We also study the interplay
15 between arithmetic and text data during training and examine the effects of few-shot
16 prompting, pretraining, and parameter scaling. Additionally, we discuss the chal-
17 lenges associated with length generalization. Our work highlights the importance
18 of high-quality, instructive data that considers the particular characteristics of the
19 next-word prediction loss for rapidly eliciting arithmetic capabilities.¹

20 1 Introduction

21 Large language models like GPT-3/4, PaLM, LaMDA (Brown et al., 2020; Chowdhery et al., 2022;
22 Thoppilan et al., 2022) have demonstrated general-purpose properties, often referred to as *emergent*
23 *abilities* (Wei et al., 2022a), for a wide range of downstream tasks like language and code translation,
24 compositional reasoning, and basic arithmetic operations (Webb et al., 2022; Nye et al., 2021; Wei
25 et al., 2022b; Shi et al., 2022; Wang et al., 2022; Srivastava et al., 2022; Chen et al., 2023). What is
26 perhaps surprising, is that these tasks are not explicitly encoded in the model’s training objective,
27 which typically is an auto-regressive, next-token-prediction loss.

28 Prior research (see Appendix 5) has delved into exploring these capabilities and how they emerge
29 as the scale and of training compute, type of data, and model size vary (Wei et al., 2022a; Chung
30 et al., 2022; Tay et al., 2022). Untangling the factors, however, remains challenging due to the data
31 complexity and the variety of tasks examined. Driven by the curiosity to understand the factors
32 that elicit these capabilities in next-token predictors, we set out to pinpoint the key contributors that
33 accelerate the emergence of such abilities. These contributors may include the format and scale of
34 data, model scale, the presence of pre-training, and the manner of prompting.

¹Our code is available at <https://anonymous.4open.science/r/nanoGPT-25D2>

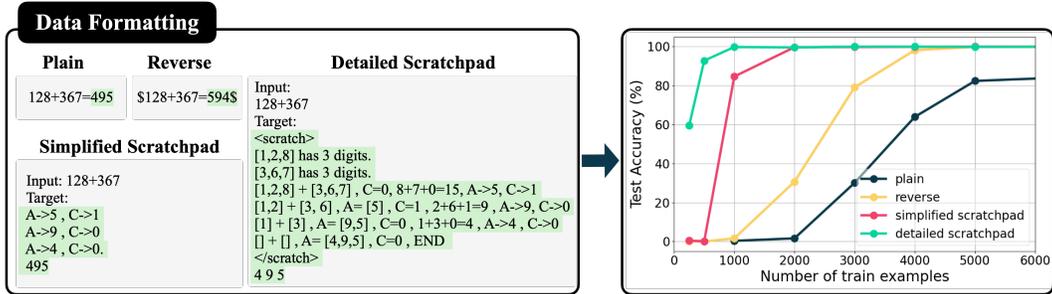


Figure 1: We investigate four data formatting approaches: (i) **Plain**: standard addition formatting (Appendix 6), (ii) **Reverse**: reversing the output (Appendix 6), (iii) **Simplified Scratchpad**: recording the digit-wise sum and carry-ons (Appendix 8), and (iv) **Detailed Scratchpad**: providing detailed intermediate steps (Appendix 8). We train small decoder-only transformers from scratch on addition data in these formats. The results (right) highlight the crucial role of data formatting in accuracy and sample efficiency. Plain never reaches 100% accuracy and the sample complexity for the remaining methods steadily improves with the level of details in the data format.

35 To provide a more precise examination of these factors, our study is conducted in a controlled setting:
 36 we first focus on teaching arithmetic to small decoder-only transformer models, such as NanoGPT
 37 and GPT-2, when trained from random initialization. Starting with a model of 10.6M parameters and
 38 scaling up to 124M parameters, we use the standard autoregressive next-token prediction loss. Our
 39 objective is to understand if and to what degree these models can efficiently learn basic arithmetic
 40 operations like addition, subtraction, multiplication, square root, and sine, thereby providing a clearer
 41 lens through which to view the elicitation of emergent abilities. Below, we summarize our findings.

42 **Data format and sampling plays a significant role.** We first observe that teaching a model addition
 43 (or any other operation) using standard addition samples, *i.e.*, ‘ $A_3A_2A_1 + B_3B_1B_1 = C_3C_2C_1$ ’, is
 44 suboptimal, as it requires the model to evaluate the most significant digit C_3 of the result first, which
 45 depends globally on all the digits of the two summands. By training on samples with reversed results,
 46 *i.e.*, ‘ $A_3A_2A_1 + B_3B_1B_1 = C_1C_2C_3$ ’, we enable the model to learn a simpler function, significantly
 47 improving sample complexity (Appendix 6). Additionally, balanced sampling of different “variants”
 48 of addition, based on the number of carries and digits involved, further enhances learning. Even
 49 in this simple setting, we observe relatively sharp phase transitions from 0 to 100% accuracy as a
 50 function of the size of the training data. Although this may seem surprising, we observe that learning
 51 an addition map on n digits from random samples is equivalent to completing a low-rank matrix.
 52 This connection allows us to offer a reasonable explanation for such phase transitions (Appendix 7).

53 **Chain-of-thought data during training.** Building on these findings, we then explore the potential
 54 benefits of chain-of-thought (CoT) data during training. This format includes step-by-step operations
 55 and intermediate results, allowing the model to learn the individual components of compositional
 56 tasks. This format is directly borrowed from related literature, *e.g.*, Ling et al. (2017); Wei et al.
 57 (2022b); Zhou et al. (2022a,b). We find that CoT-type training data significantly improved learning in
 58 terms of both sample complexity and accuracy in agreement with CoT fine-tuning literature (Nye
 59 et al., 2021; Chung et al., 2022), but *even in the complete absence of pretraining* (Appendix 8). We
 60 conjecture that this is because breaking down the required compositional function to be learned into
 61 individual components allows the model to learn a higher-dimensional but easier-to-learn function
 62 map, in agreement with recent theoretical findings (Li et al., 2023; Malach, 2023). In Figure 1, we
 63 provide examples of the data formatting methods explored in our work.

64 **Training on text and arithmetic mixtures and the role of few-shot prompting.** We also explore the
 65 interplay between arithmetic and text data during training, as LLMs are trained on massive amounts
 66 of data scraped from the internet (Bubeck et al., 2023; Peterson et al., 2019), where it is impractical to
 67 carefully separate different types of data. We observe how the model’s perplexity and accuracy vary
 68 with the ratio of text to arithmetic data. We find that jointly training on all the arithmetic operations
 69 discussed earlier can improve the individual performance of each task and that going from zero-shot
 70 to 1-shot prompting (showing one arithmetic example) yields a large accuracy improvement, but
 71 there is no significant improvement in accuracy by showing more examples (Appendix 9).

72 **The role of pre-training and model scale.** We further investigate the role of pretraining by fine-
 73 tuning pretrained models like GPT-2 and GPT-3 (davinci) and observe that while the zero-shot
 74 performance on arithmetic operations is poor, prior “skills” acquired during pretraining facilitate
 75 quick learning of some basic arithmetic tasks, even with a small number of finetuning samples.

76 However, finetuning on non-standard data, such as those that result from reverse formatting, can
77 interfere with the model’s performance when pretrained, leading to decreased accuracy. We finally
78 share our observations on how performance in arithmetic changes with scale, and although we find
79 that scale does aid when finetuning for these tasks, it is not a necessary trait (Appendix 10).

80 **Compositional and length generalization.** One might question if our trained models truly grasp
81 arithmetic. Our findings present a nuanced answer. We find that length generalization beyond trained
82 digit lengths is still challenging. For instance, if a model is trained on all n -digit lengths, excluding a
83 specific length, it still struggles to accurately calculate this missing digit length. Consequently, the
84 models achieve high accuracy within trained digit lengths but struggle significantly beyond this range.
85 This suggests that the models learn arithmetic not as a flexible algorithm, but as a mapping function
86 constrained to trained digit lengths. While this significantly surpasses memorization, it falls short of
87 comprehensive arithmetic “understanding” (Appendix 16).

88 **Novelty over prior work.** Our approach heavily builds upon prior work that uses reasoning-
89 augmented data to enhance model performance, and we do not purport originality in the types of
90 training data used, nor in achieving the highest performance with the smallest model parameters
91 possible. What sets our work apart is the primary focus on meticulously ablating our settings and
92 extensive studies on various sampling techniques, training data formats, data source mixing ratios,
93 and model scales. Our goal is to pinpoint the factors that contribute to the fast emergence of arithmetic
94 capabilities. In the process, we also provide several straightforward yet novel and insightful theoretical
95 explanations for some of the phase transition phenomena we observe. Our emphasis on arithmetic is
96 not due to its intrinsic significance — one can easily delegate calculations to external tools (Schick
97 et al., 2023; Gao et al., 2023). Instead, arithmetic serves as an emergent skill, easy to isolate and test,
98 facilitating a more precise exploration of emergent phenomena.

99 2 Preliminaries and Experimental Setup

100 In this section, we provide a detailed description of our experimental setup, including the model
101 architecture and an overview of the various data formatting and sampling techniques used.

102 **Model and Data.** To examine the individual factors at play, we use NanoGPT (Karpathy, 2022), a
103 lightweight implementation of the GPT family of models. NanoGPT is a decoder-only transformer
104 with six self-attention layers, six heads, and an embedding dimension of 384, resulting in approxi-
105 mately 10.6M parameters. Unless stated otherwise, we use character-level tokenization and absolute
106 position encoding. We train NanoGPT from random initialization, which we refer to as *training from*
107 *scratch*, using the conventional next-token prediction objective. To study the effect of scale, we extend
108 our experiments to GPT-2 and GPT-3 in Appendix 10. We investigate both training from scratch
109 as well as fine-tuning using a pretrained GPT-2, whereas, for GPT-3, we only consider fine-tuning
110 pretrained models. Refer to Appendix 19 for more details on the models and data used.

111 For arithmetic tasks like addition, subtraction, and multiplication, we define the training dataset for
112 a binary operator $f(\cdot)$ as $\mathcal{D}_{\text{train}} = \{(a_i, b_i), y_i\}_{i=1}^N$ where $y_i = f(a_i, b_i)$. For unary operations like
113 sine, the training dataset is formulated as $\mathcal{D}_{\text{train}} = \{a_i, y_i\}_{i=1}^N$, where $y_i = f(a_i)$. The test dataset
114 $\mathcal{D}_{\text{test}}$ is constructed by randomly sampling pairs of operands not included in $\mathcal{D}_{\text{train}}$. We then apply
115 different *data formatting* techniques on each data sample from the training dataset, creating the final
116 sequence that serves as the model’s input. Note that while we view a_i as a single integer, the model
117 will see it as a sequence of digits after character-level tokenization.

118 **Data Formatting.** In the following sections, we will delve into the four data formatting approaches
119 in our arithmetic experiments. See Figure 1 and Appendix 20 for examples. In Appendix 6, we explore
120 the limitations of the conventional plain-format data and demonstrate how a simple reversal of the
121 output order can lead to substantial performance improvements and enhanced sample efficiency. We
122 introduce two Lemmas to support and explain these findings. Additionally, in Appendix 8, we present
123 results on the simplified and detailed scratchpad formats, highlighting significant enhancements in
124 sample efficiency for learning addition. We also emphasize the importance of carefully designing the
125 intermediate steps in the detailed scratchpad method. Note that the scratchpad formats are largely
126 adopted from the literature of chain-of-thought (CoT) training (Nye et al., 2021; Zhou et al., 2022b).

127 **Structured Data Sampling.** While data formatting plays a crucial role, we also discover that
128 choosing the samples carefully is also essential. When sampling operands for n -digit addition
129 uniformly at random between 1 to $10^n - 1$, the dataset inevitably becomes highly skewed in terms

130 of the number of samples with (i) operands containing a certain number of digits and (ii) operands
 131 resulting in a certain number of *carry-on* operations. For instance, in the case of 3-digit addition,
 132 random sampling results in a meager 0.01% probability of selecting a 1-digit number. Additionally, 1
 133 or 2 carry-on operations are more likely to occur than 0 or 3. To address this *imbalance*, we employ a
 134 structured sampling approach. Specifically, we aim to **(i) balance digits** by assigning higher weights
 135 to lower-digit numbers during the sampling process and **(ii) balance carry-ons** by ensuring an equal
 136 distribution of examples with 0, 1, . . . , n carry-on operations.

137 When sampling 10,000 examples of 3-digit addition,
 138 we include all 100 1-digit additions, 900 2-digit sam-
 139 ples and 9000 3-digit samples. Note that while the
 140 number of samples increases, the fraction of all possi-
 141 ble k -digit additions that we sample for $k = 2, 3$
 142 decreases due to the inherent skew. The split was
 143 chosen to ensure we saw a “reasonable” fraction of
 144 all possible k -digit samples for all k . Similarly, we
 145 ensure that the number of samples with 0, 1, 2, or 3
 146 carry-ons are all approximately 2500.

147 Figure 2 reveals the importance of balancing. We
 148 observe improvements in accuracy across the board
 149 while using *balanced* data when compared to ran-
 150 dom sampling. Further, random sampling performs
 151 relatively poorly even for the simple task of 2-digit
 152 addition, possibly due to the fact that the model has
 153 not seen enough of these examples. For the remaining
 154 experiments, we set the default dataset for addition
 155 to be one that has both balanced digits and carry-ons.

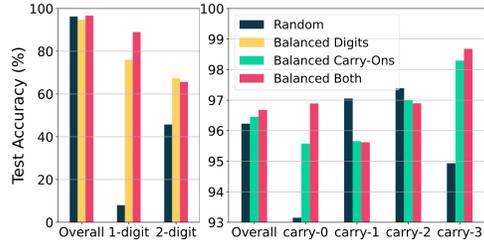


Figure 2: Performance of 3-digit addition on various data sampling methods used: **(i) Random**: uniform sampling of operands; **(ii) Balanced digits**: assigning higher sampling weights to operations involving 1 and 2-digit numbers; **(iii) Balanced carry**: balancing the dataset to contain an equal number of carry-on operations; **(iv) Balanced both**: balancing digits and carry-ons. We observe that *balanced data* improves accuracy compared to random sampling. Experiments on addition with the ‘\$’ symbol wrapped for each sample.

156 3 Limitations

157 **Length generalization.** In our experiments, we did not observe any instances where the model
 158 could predict beyond the number of digits it had been trained on (see Appendix 16). Shaw et al.
 159 (2018); Sun et al. (2022) reported similar difficulties and proposed approaches such as relative
 160 positional encodings. Anil et al. (2022) suggests that models can only perform out-of-distribution
 161 tasks by combining fine-tuning, prompting, and scratchpad techniques.

162 **Model/Data scale.** Due to the smaller scale of our experiments, we were able to thoroughly
 163 examine the impact of individual components on the model’s arithmetic learning capabilities. Our
 164 model was limited to decoder-only architectures, primarily focusing on character-level tokenization.
 165 Although we have some preliminary results on scaling up and incorporating BPE-based tokenization,
 166 it is not clear if all our findings can be generalized to the scale of LLMs being used in practice.

167 **Beyond elementary arithmetic.** We choose to analyze simple arithmetic operations in order
 168 to carefully isolate factors that contribute to emergence. While the existing literature has already
 169 demonstrated the emergence of complicated abilities in practice, our work seeks to provide a better
 170 understanding of this behavior by extensive ablations in a controlled setting.

171 4 Conclusion

172 In this study, we investigate teaching arithmetic operations to small randomly initialized transformers
 173 using the next-token prediction objective. We carefully ablate different aspects of the training setting
 174 so as to isolate the factors that contribute to the emergence of arithmetic capabilities. Our results reveal
 175 that traditional training data is sub-optimal for learning arithmetic, and training on data with detailed
 176 intermediate steps or even simply reversing the output improves accuracy and sample complexity.
 177 We also study the effects of few-shot prompting, pretraining, and model scale. Despite improvements
 178 from detailed data, length generalization remains a challenge, highlighting the need for better-curated
 179 training data to ensure successful learning of specific algorithms as opposed to just learning an
 180 approximate function map. The correct approach for learning multiple arithmetic operations, of
 181 different levels of complexity, is still unclear. We anticipate that this research will contribute to a
 182 more nuanced understanding of the mechanisms by which transformers (approximately) acquire
 183 algorithmic skills.

184 **References**

- 185 Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh,
186 Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization
187 in large language models. *arXiv preprint arXiv:2207.04901*, 2022.
- 188 Samuel R Bowman. Can recursive neural tensor networks learn logical reasoning? *arXiv preprint*
189 *arXiv:1312.6192*, 2013.
- 190 Samuel R Bowman, Christopher Potts, and Christopher D Manning. Recursive neural networks for
191 learning logical semantics. *CoRR, abs/1406.1827*, 5, 2014.
- 192 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
193 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
194 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 195 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,
196 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:
197 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 198 Jonathon Cai, Richard Shin, and Dawn Song. Making neural programming architectures generalize
199 via recursion. *arXiv preprint arXiv:1704.06611*, 2017.
- 200 François Charton. Linear algebra with transformers. *arXiv preprint arXiv:2112.01898*, 2021.
- 201 François Charton. What is my math transformer doing?—three results on interpretability and general-
202 ization. *arXiv preprint arXiv:2211.00170*, 2022.
- 203 Xinyun Chen, Chang Liu, and Dawn Song. Towards synthesizing complex programs from input-
204 output examples. *arXiv preprint arXiv:1706.01284*, 2017.
- 205 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to
206 self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
- 207 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
208 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
209 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 210 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi
211 Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models.
212 *arXiv preprint arXiv:2210.11416*, 2022.
- 213 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
214 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
215 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 216 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal
217 transformers. *arXiv preprint arXiv:1807.03819*, 2018.
- 218 Andrew Drozdov, Nathanael Schärli, Ekin Akyürek, Nathan Scales, Xinying Song, Xinyun Chen,
219 Olivier Bousquet, and Denny Zhou. Compositional semantic parsing with large language models.
220 *arXiv preprint arXiv:2209.15003*, 2022.
- 221 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Peter West,
222 Chandra Bhagavatula, Ronan Le Bras, Jena D. Hwang, Soumya Sanyal, Sean Welleck, Xiang
223 Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. Faith and fate: Limits of transformers on
224 compositionality, 2023.
- 225 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
226 Graham Neubig. Pal: Program-aided language models, 2023.
- 227 Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-than?: Inter-
228 preting mathematical abilities in a pre-trained language model. *arXiv preprint arXiv:2305.00586*,
229 2023.

- 230 Samy Jelassi, Stéphane d’Ascoli, Carles Domingo-Enrich, Yuhuai Wu, Yuanzhi Li, and François
231 Charton. Length generalization in arithmetic transformers, 2023.
- 232 Łukasz Kaiser and Ilya Sutskever. Neural gpu learn algorithms. *arXiv preprint arXiv:1511.08228*,
233 2015.
- 234 Andrej Karpathy. char-rnn. <https://github.com/karpathy/char-rnn>, 2015.
- 235 Andrej Karpathy. Andrej karpathy’s lightweight implementation of medium-sized gpts. *GitHub*,
236 2022. URL <https://github.com/karpathy/nanoGPT>.
- 237 Jeonghwan Kim, Giwon Hong, Kyung-min Kim, Junmo Kang, and Sung-Hyon Myaeng. Have you
238 seen that number? investigating extrapolation in question answering models. In *Proceedings of the*
239 *2021 Conference on Empirical Methods in Natural Language Processing*, pp. 7031–7037, 2021.
- 240 Franz J Király, Louis Theran, and Ryota Tomioka. The algebraic combinatorial approach for low-rank
241 matrix completion. *J. Mach. Learn. Res.*, 16(1):1391–1436, 2015.
- 242 Brenden Lake and Marco Baroni. Generalization without systematicity: On the compositional skills
243 of sequence-to-sequence recurrent networks. In *International conference on machine learning*, pp.
244 2873–2882. PMLR, 2018.
- 245 Yingcong Li, Kartik Sreenivasan, Angeliki Giannou, Dimitris Papailiopoulos, and Samet Oymak.
246 Dissecting chain-of-thought: A study on compositional in-context learning of mlps. *arXiv preprint*
247 *arXiv:2305.18869*, 2023.
- 248 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
249 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
250 *arXiv:2305.20050*, 2023.
- 251 Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale genera-
252 tion: Learning to solve and explain algebraic word problems. *arXiv preprint arXiv:1705.04146*,
253 2017.
- 254 Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Exposing attention
255 glitches with flip-flop language modeling. *arXiv preprint arXiv:2306.00946*, 2023.
- 256 Eran Malach. Auto-regressive next-token predictors are universal learners, 2023.
- 257 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
258 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv*
259 *preprint arXiv:2202.12837*, 2022.
- 260 MosaicML. Introducing mpt-7b: A new standard for open source, commercially usable llms, 2023.
261 URL www.mosaicml.com/blog/mpt-7b. Accessed: 2023-05-05.
- 262 Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. Investigating the limitations of transformers with
263 simple arithmetic tasks. *arXiv preprint arXiv:2102.13019*, 2021.
- 264 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David
265 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work:
266 Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*,
267 2021.
- 268 Santiago Ontanón, Joshua Ainslie, Vaclav Cvicek, and Zachary Fisher. Making transformers solve
269 compositional tasks. *arXiv preprint arXiv:2108.04378*, 2021.
- 270 OpenAI. OpenAI platform. URL <https://platform.openai.com/docs/models/gpt-3>. Ac-
271 cessed: 2023-09-28.
- 272 Joshua Peterson, Stephan Meylan, and David Bourgin. Open clone of openai’s unreleased webtext
273 dataset scraper. *GitHub*, 2019. URL <https://github.com/jcpeterson/openwebtext>.

- 274 Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Jonathan
275 Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. Efficiently scaling transformer inference.
276 *Proceedings of Machine Learning and Systems*, 5, 2023.
- 277 Jing Qian, Hong Wang, Zekun Li, Shiyang Li, and Xifeng Yan. Limitations of language models in
278 arithmetic and symbolic induction. *arXiv preprint arXiv:2208.05051*, 2022.
- 279 Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training.
280 2018.
- 281 Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. Impact of pretraining term
282 frequencies on few-shot reasoning. *arXiv preprint arXiv:2202.07206*, 2022.
- 283 Benjamin Recht. A simpler approach to matrix completion. *Journal of Machine Learning Research*,
284 12(12), 2011.
- 285 Scott Reed and Nando De Freitas. Neural programmer-interpreters. *arXiv preprint arXiv:1511.06279*,
286 2015.
- 287 Subhro Roy and Dan Roth. Solving general arithmetic word problems. *arXiv preprint*
288 *arXiv:1608.01413*, 2016.
- 289 Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer,
290 Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to
291 use tools, 2023.
- 292 Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations.
293 *arXiv preprint arXiv:1803.02155*, 2018.
- 294 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
295 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. Language models are mul-
296 tilingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057*, 2022.
- 297 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
298 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the
299 imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint*
300 *arXiv:2206.04615*, 2022.
- 301 Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary,
302 Xia Song, and Furu Wei. A length-extrapolatable transformer. *arXiv preprint arXiv:2212.10554*,
303 2022.
- 304 Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks.
305 *Advances in neural information processing systems*, 27, 2014.
- 306 Yi Tay, Jason Wei, Hyung Won Chung, Vinh Q Tran, David R So, Siamak Shakeri, Xavier Garcia,
307 Huaixiu Steven Zheng, Jinfeng Rao, Aakanksha Chowdhery, et al. Transcending scaling laws with
308 0.1% extra compute. *arXiv preprint arXiv:2210.11399*, 2022.
- 309 Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze
310 Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog
311 applications. *arXiv preprint arXiv:2201.08239*, 2022.
- 312 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
313 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
314 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
315 models, 2023.
- 316 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia
317 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and
318 outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- 319 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
320 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
321 *systems*, 30, 2017.

- 322 Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh, and Matt Gardner. Do nlp models know
323 numbers? probing numeracy in embeddings. *arXiv preprint arXiv:1909.07940*, 2019.
- 324 Cunxiang Wang, Boyuan Zheng, Yuchen Niu, and Yue Zhang. Exploring generalization ability of
325 pretrained language models on arithmetic and logical reasoning. In *Natural Language Processing*
326 *and Chinese Computing: 10th CCF International Conference, NLPCC 2021, Qingdao, China,*
327 *October 13–17, 2021, Proceedings, Part I 10*, pp. 758–769. Springer, 2021.
- 328 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency
329 improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- 330 Taylor Webb, Keith J Holyoak, and Hongjing Lu. Emergent analogical reasoning in large language
331 models. *arXiv preprint arXiv:2212.09196*, 2022.
- 332 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
333 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.
334 *arXiv preprint arXiv:2206.07682*, 2022a.
- 335 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny
336 Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint*
337 *arXiv:2201.11903*, 2022b.
- 338 Zhen Yang, Ming Ding, Qingsong Lv, Zhihuan Jiang, Zehai He, Yuyi Guo, Jinfeng Bai, and Jie Tang.
339 Gpt can solve mathematical problems without a calculator, 2023.
- 340 Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, and Songfang Huang. How well do large
341 language models perform in arithmetic tasks? *arXiv preprint arXiv:2304.02015*, 2023.
- 342 Wojciech Zaremba and Ilya Sutskever. Learning to execute. *arXiv preprint arXiv:1410.4615*, 2014.
- 343 Wojciech Zaremba, Karol Kurach, and Rob Fergus. Learning to discover efficient mathematical
344 identities. *Advances in Neural Information Processing Systems*, 27, 2014.
- 345 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,
346 Olivier Bousquet, Quoc Le, and Ed Chi. Least-to-most prompting enables complex reasoning in
347 large language models. *arXiv preprint arXiv:2205.10625*, 2022a.
- 348 Hattie Zhou, Azade Nova, Hugo Larochelle, Aaron Courville, Behnam Neyshabur, and Hanie Sedghi.
349 Teaching algorithmic reasoning via in-context learning. *arXiv preprint arXiv:2211.09066*, 2022b.

350

351 Appendix

352

353 Table of Contents

354	5 Related Works	10
355	6 Data Format Challenges and Arithmetic Emergence	10
356	7 Matrix Completion: an Incomplete Tale of Emergence	11
357	8 Training on Chain-of-Thought Data Expedites Emergence	12
358	9 Longer Digits, Varied Operations, and Blending Arithmetic with Shakespeare	13
359	10 Fine-tuning, Scaling, and Pretraining in Larger Models	14
360	11 Extending to Longer Digit Addition	15
361	11.1 Training from Random Initialization	15
362	11.2 Fine-Tuning from Pretrained Models	16
363	11.3 Impact of Formats on Fine-Tuning	17
364	12 Teaching Arithmetic Operations Beyond Addition	19
365	12.1 Extended Arithmetic Operations	19
366	12.2 Jointly Training on All Five Arithmetic Tasks	21
367	13 Mixing Text with Arithmetic Data	22
368	13.1 Few-Shot Prompting	22
369	13.2 Disentangling the effect of text on prompting	22
370	13.3 Prompting with Text	23
371	14 Fine-tuning, Scaling, and Pretraining in Larger Models	26
372	15 Token Efficiency Across Data Formats	28
373	16 Length Generalization	29
374	17 Proofs	32
375	18 Additional Experiments	33
376	18.1 Zero-Padding and Symbol Wrapping	33
377	18.2 Low-Rank Matrix Completion	34
378	18.3 The Importance of Intermediate Step Design	35
379	18.4 The Effect of Noisy Inputs on Accuracy	37
380	18.5 Analyzing the results on Sine/Sqrt	38
381	19 Experimental Setup	40
382	19.1 Dataset	40
383	19.2 Model	43
384	19.3 Hyperparameter Configurations	44
385	20 Prompt Examples	45
386	20.1 Addition	45
387	20.2 Subtraction	47
388	20.3 Multiplication	48
389	20.4 Sine	48

390	20.5 Square Root	49
391	20.6 Noisy Simple Scratchpad	49
392	20.7 Example data for GPT-3 fine-tuning	50

396 5 Related Works

397 **Instructional data/chain-of-thought.** Detailed reasoning in training data has roots predating
 398 Transformers (Vaswani et al., 2017). Ling et al. (2017); Cobbe et al. (2021) use natural language to
 399 generate reasoning steps while Roy & Roth (2016); Reed & De Freitas (2015); Chen et al. (2017);
 400 Cai et al. (2017); Nye et al. (2021) show that symbolic reasoning may suffice. Nogueira et al. (2021)
 401 stress the importance of large number of small-digit samples (Yuan et al., 2023). Razeghi et al.
 402 (2022) observe a correlation between the frequency of numbers in the dataset and the performance
 403 involving them. In contrast, we find that transformers can learn to add numbers that were not seen
 404 during training. Chain-of-thought (Wei et al., 2022b) refers to the model’s improved accuracy when
 405 prompted to produce intermediate reasoning steps. Zhou et al. (2022b) show that this can be achieved
 406 by providing sufficiently informative exemplars as a few-shot prompt (Brown et al., 2020). Zhou et al.
 407 (2022a) showed that *least-to-most* prompting can help GPT-3 solve problems decomposable into
 408 simpler sub-problems, by sequentially solving these subproblems. We extend this notion to simple
 409 addition and show that asking the model to output the least significant bit first has a similar effect.

410 **Arithmetic using Transformer models.** Our work focuses on decoder-only models as they are
 411 widely used in LLMs (Brown et al., 2020; Touvron et al., 2023; MosaicML, 2023). However,
 412 encoder-decoder models have also been extensively studied in the context of learning arithmetic
 413 operations (Kim et al., 2021; Wang et al., 2021). Wallace et al. (2019) on the other hand, focus
 414 on the impact of the learned embeddings. Ontanón et al. (2021) extensively study the problem of
 415 compositional generalization on benchmark datasets, such as SCAN (Lake & Baroni, 2018; Drozdov
 416 et al., 2022), and conclude that design choices, like relative position encoding (Shaw et al., 2018),
 417 can improve performance. Charton (2022, 2021) show that Transformers can learn linear algebra
 418 operations with carefully chosen encodings. Hanna et al. (2023) use mechanistic interpretability
 419 techniques to explain the limited numerical reasoning capabilities of GPT-2. Dziri et al. (2023);
 420 Jelassi et al. (2023); Yang et al. (2023) focus on the challenges of length generalization. A recent line
 421 of work explores finetuning techniques to improve arithmetic capabilities in pretrained models (Qian
 422 et al., 2022; Lightman et al., 2023; Uesato et al., 2022).

423 **Beyond Transformers.** While we focus our attention on GPT-like models, there is a rich literature
 424 studying other seq-to-seq models such as recurrent neural networks (RNNs) (Bowman, 2013; Bowman
 425 et al., 2014; Zaremba et al., 2014). Zaremba & Sutskever (2014) show that RNNs can learn how to
 426 execute simple programs with for-loops provided they are trained with curriculum learning. Sutskever
 427 et al. (2014) show that LSTMs show improved performance on text-based tasks such as translation
 428 when the source sentences are reversed, which is closely related to what we observe in addition.
 429 Kaiser & Sutskever (2015) propose Neural GPUs which outperform prior RNNs on binary arithmetic
 430 tasks and even show length generalization *i.e.*, they can perform arithmetic on inputs of lengths that
 431 were unseen during training. This is yet to be seen even in modern pre-trained models (Bubeck
 432 et al., 2023) and therefore it is interesting to see if we can leverage some of these techniques and
 433 apply them to existing modern architectures. Dehghani et al. (2018) propose Universal Transformers
 434 (UTs) which introduce a recurrent transition function to apply recurrence over revisions of the vector
 435 representation at each position as opposed to the different positions in the input. They show that on
 436 the tasks from Zaremba & Sutskever (2014), UTs outperform traditional Transformers and RNNs.

437 6 Data Format Challenges and Arithmetic Emergence

438 We start by examining *integer addition*. We first focus on 3-digit addition, *i.e.*, where the two
 439 summands have at most 3 digits (≤ 999). Later, in Section 9, we extend our findings to numbers with
 440 up to 10 digits. Surprisingly, teaching addition can be more complex than expected.

441 **Training on Conventional Data.** We start by training NanoGPT on standard addition data
 442 represented as ‘ $A_3A_2A_1 + B_3B_2B_1 = C_3C_2C_1$ ’, termed the *plain* format. However, as shown in
 443 Figure 1, this leads to fairly poor performance. We suspect that this is because the next-token

444 prediction objective outputs the most significant digit (MSB) first. The following lemma clarifies the
 445 necessity to access all operand digits for outputting the MSB first.

446 **Lemma 1.** *Let A and B be two n -digit numbers, and let $C = A + B$. Suppose an algorithm \mathcal{A}
 447 outputs the digits of C in decreasing order of significance, then \mathcal{A} must have access to all digits of A
 448 and B starting from the first digit that it outputs.*

449 The lemma suggests that to train a model for addition and to output the MSB first, it is necessary to
 450 emulate a “*global*” algorithm. Unlike the standard “*local*” algorithm for addition, which consists of
 451 computing digit-wise sums and carry-ons, approximating the global algorithm would require learning
 452 a more complicated function than necessary. The increased complexity results in decreased accuracy,
 453 as observed in our experiments. Liu et al. (2023) refer to this phenomenon as *attention glitches*.

454 **Reversing the Output.** We propose that the *reverse* format ‘ $\$A_3A_2A_1 + B_3B_2B_1 = C_1C_2C_3\$$ ’ is
 455 more suitable for next-word prediction models. The rationale behind this is that when generating the
 456 sum by starting with the least significant digit (LSB), the model only needs to learn a local function
 457 of three inputs per digit – the two relevant digits of the operands and the carry-on from the previous
 458 digit. This local operation simplifies the function to be learned. The following lemma formalizes this:

459 **Lemma 2.** *There exists an algorithm that computes $C = A + B$ for two n -digit numbers A and B
 460 and outputs its digits in increasing order of significance such that, at each position i , the algorithm
 461 only requires access to the i^{th} digits of A and B , as well as the carry-on from the previous position.*

462 Lemma 2 directly follows from the *standard* algorithm for addition, which performs the sum and
 463 carry-on operations digit by digit. The implications of these lemmata are evident in our experiments
 464 when comparing the accuracy of the *plain* and *reverse* formats. As shown in Figure 1, the accuracy
 465 of *plain* plateaus at around 85%, even with 10k addition examples. In contrast, training on reversed
 466 outputs significantly enhances accuracy. Moreover, *reverse* requires considerably fewer samples.
 467 What is particularly remarkable is the rapid emergence of addition and a *phase transition* occurring
 468 between 1k to 4k samples for reverse. During that, the model rapidly transitions from being unable to
 469 add two numbers to being capable of perfectly adding. This leads us to ask:

470 *Why does addition rapidly emerge as the number of training examples increases?*

471 7 Matrix Completion: an Incomplete Tale of Emergence

472 Although the rapid phase transition observed in the previous section may initially seem surprising,
 473 closer examination reveals a fascinating equivalence – learning an addition map on n digits from
 474 random samples can be considered as completing a rank-2 matrix. Establishing this connection
 475 with low-rank matrix completion (LRMC) provides meaningful insights into the observed phe-
 476 nomenon. However as we explain in this Section, this connection does not tell the complete story,
 477 and Transformers possess generalization capabilities far beyond what LRMC would predict.

478 **Learning addition tables is Matrix Comple-**
 479 **tion.** Learning addition from samples $i + j$
 480 can be formulated as a rank-2 Matrix Comple-
 481 tion (MC) problem, where we partially obser-
 482 ve an $n \times n$ matrix M , whose (i, j) -th
 483 entry represents $i + j$. M can be decom-
 484 posed into the sum of two rank-1 matrices,
 485 $\mathbf{N}\mathbf{1}^T + \mathbf{1}\mathbf{N}^T$, where \mathbf{N} is a vector with en-
 486 tries $\{1, \dots, n\}$, and $\mathbf{1}$ is the vector of ones.
 487 Recovering a rank-2 matrix, in the absence of
 488 noise, can be sample-optimally performed by
 489 a simple iterative algorithm from Király et al.
 490 (2015) (Algorithm 1 in Appendix 18.2). As
 491 depicted in Figure 3a, a sharp *phase transition*
 492 occurs at $\mathcal{O}(n)$, a well-known matrix recovery
 493 phenomenon (Recht, 2011).

494 We notice a similar phase transition in
 495 NanoGPT. To investigate it, we focus on 2-digit

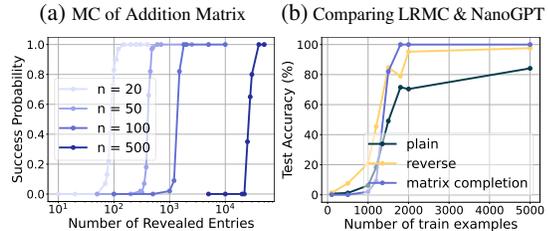


Figure 3: (a) We run Algorithm 1 (Király et al., 2015) on the addition matrix for $n = 20, 50, 100, 500$ and report the success probability while varying the number of revealed entries. As expected, a sharp phase transition occurs when approximately $\mathcal{O}(n)$ entries are revealed. (b) We compare the performance of NanoGPT trained on a dataset containing $n = 100$ samples (*i.e.*, 2-digit addition) to that of the corresponding LRMC problem using the same sample set. Remarkably, at ≈ 1500 samples, both NanoGPT and Algorithm 1 begin learning addition almost flawlessly.

of learning addition through NanoGPT and LRMC (Figure 3a) by constructing train data as the revealed entries of the M matrix. Note that the dataset is no longer *balanced*, as the revealed entries are randomly sampled for LRMC experiments, to match the standard MC probabilistic settings Recht (2011). In Figure 3b, both NanoGPT and LRMC exhibit phase transitions at approximately 1500 samples.

While the observed phase transition can be attributed to the principles of LRMC, shedding light on the emergent arithmetic skill in NanoGPT, this connection falls short of capturing the full generalization capabilities displayed by NanoGPT.

NanoGPT generalizes better than Matrix Completion solutions. Upon further investigation, we find that NanoGPT exhibits capabilities beyond LRMC. Notably, LRMC is constrained by its inability to generalize in the presence of missing rows or columns. In our context, this equates to certain numbers being omitted from the training data. To assess NanoGPT’s ability to overcome this, we deliberately exclude specific numbers or digits from our training data and assess the model’s ability in learning addition. Can the model still generalize to unseen numbers?

As shown in Table 1, the answer to this question is a resounding *Yes!* The model achieves almost perfect accuracy even when excluding half of all possible 3-digit numbers. NanoGPT can successfully learn 3-digit addition even when numbers or digits are intentionally excluded from the training data, thereby exhibiting generalization capabilities that far exceed what standard LRMC would predict.

Table 1: Impact of excluding numbers on addition task: NanoGPT models trained with 100/200/500 excluded operands show no significant drop in accuracy and in some cases, the performance even improves. Note that models trained with *reverse* data remain consistently at 100% accuracy.

	No Exclusion		Excluding 100 numbers		Excluding 200 numbers		Excluding 500 numbers	
	Plain	Rev	Plain	Rev	Plain	Rev	Plain	Rev
Overall Accuracy	87.18%	99.97%	87.94%	100.00%	87.24%	99.99%	88.15%	99.99%
Exclusion Accuracy	-	-	92.55%	100.00%	92.15%	99.95%	90.85%	100%

Specifically, we randomly choose 100/200/500 numbers and exclude them from the training data. We then evaluate the trained models using two metrics: (i) *Overall accuracy*: which measures the accuracy over a random set of 10,000 examples and (ii) *Exclusion accuracy*: which measures the accuracy only over the excluded set. Remarkably, excluding numbers from the training data sometimes leads to improved performance. We conjecture that this may be due to a regularization effect, similar to random masking or cropping images in vision tasks. In Appendix 18.2.1, we further find that NanoGPT models can even generalize to *unseen digits*.

8 Training on Chain-of-Thought Data Expedites Emergence

So far, we observed that simply reversing the output can result in remarkable performance, exceeding that of LRMC in learning addition. Here, we investigate if it is possible to expedite the emergence of addition by further enhancing the data format. As addition is a multi-step process, we explore the idea of incorporating additional information about each intermediate step. We adopt a CoT style approach, where we guide the model to learn addition step-by-step. We explore two levels of detail in the provided instruction steps, as shown in Figure 1: (i) **Simplified Scratchpad** with minimal information – the sum and carry information for each digit/step. (ii) **Detailed Scratchpad** with comprehensive information on detailed traces of execution for each intermediate step.

The results in Figure 4a show that the model trained on simplified scratchpad achieves 100% accuracy with only 2000 samples, whereas reverse requires more than twice as many. Detailed scratchpad, which provides even more fine grained information,

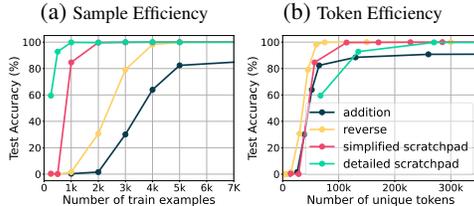


Figure 4: (a) Comparison of sample efficiency: evaluating performance on training datasets with different numbers of addition samples. While all variants other than plain achieve 100% accuracy, they differ in terms of sample complexity. (b) Number of unique tokens required by NanoGPT to learn addition. Reverse is the most efficient in terms of token usage for model training, as the scratchpad methods, although more sample-efficient, require more tokens per sample.

540 achieves perfect addition with just 1000 samples. This indicates a clear message: incorporating more
 541 information enables the model to learn addition with far fewer examples. We conjecture that this is
 542 because breaking down the required compositional function to be learned into individual, simpler
 543 components allows the model to learn a higher-dimensional but easier-to-learn function map, in
 544 agreement with recent theoretical work (Li et al., 2023; Malach, 2023).

545 We note that while CoT-style training enhances sample efficiency, it may not necessarily be the most
 546 “token-efficient” approach. To account for the cost associated with training and inference, we conduct
 547 a cost analysis based on the number of “unique” tokens (number of training samples \times number of
 548 tokens per sample – see Appendix 15 for details). encountered during training. The result in Figure 4b
 549 shows that reverse is the most efficient in terms of *token usage* for model training. The scratchpad
 550 methods, although more *sample-efficient*, require more tokens per sample.

551 In summary, incorporating scratchpad data and decomposing the addition task into steps offer a
 552 promising strategy to improve the performance and efficiency of small models in learning addition
 553 from scratch. Nevertheless, for practical usage, it is crucial to evaluate both the number of samples for
 554 achieving the desired performance and the actual token requirements during training and inference.

555 9 Longer Digits, Varied Operations, and Blending Arithmetic with 556 Shakespeare

557 In this section, we go beyond 3-digit addition to encompass
 558 a wider range of arithmetic tasks and longer digits to show
 559 that our insights on data sampling and formatting hold true
 560 even in this regime. We also explore the effect of mixing
 561 arithmetic with text data, and few-shot prompting.

562 **Extending to longer digit addition.** We repeat the exper-
 563 iment from Section 2 with up to 10 digit integers. Figure 5
 564 shows that the behavior of all data formats remains similar
 565 across varying number of digits. In fact, the performance
 566 gap between the modified formats and plain grows with
 567 longer digits. While plain requires an increasing number
 568 of samples to learn higher-digit additions, the reverse and
 569 scratchpad formats maintain a consistent sample complexity.

570 We also observe similar results in the *fine-tuning* setting, where we fine-tune a model initially trained
 571 on k -digits on $k + 1$ -digit data. See Appendix 11 for details on the experimental setup.

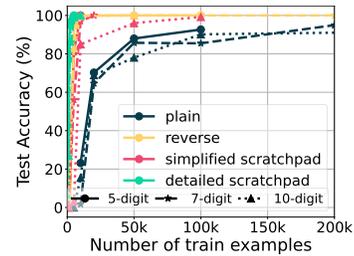


Figure 5: Comparison of sample efficiency for 5, 7, and 10-digit additions. Training on plain requires an increasing number of samples for higher digits, while the sample complexity for other methods remains relatively consistent.

572 **Teaching arithmetic operations beyond addition.** While our experiments so far were primarily
 573 focused on *addition*, we include other arithmetic operations to demonstrate the broader applicability
 574 of our insights. We consider a mix of arithmetic tasks – *subtraction*, *multiplication*, *sine*, and *square*
 575 *root*. Each operation entails its unique challenges and intricacies. For instance, subtraction introduces
 576 the concept of negative numbers, multiplication can generate significantly longer outputs, and sine
 577 and square root functions entail computations involving floating-point numbers.

578 The results depicted in Figure 6 indicate that similar
 579 to the findings of addition, the detailed scratchpad
 580 format significantly improves performance over
 581 *plain* or *reverse* formats and yields efficient results
 582 even with few samples for subtraction and multi-
 583 plication tasks. Interestingly, we find *reverse* is not
 584 particularly effective in multiplication. On the other
 585 hand, the detailed scratchpad format exhibits re-
 586 duced efficiency for *sin* and *sqrt* compared to other
 587 operations. This discrepancy can be traced back to
 588 the complexity of the intermediate steps involved
 589 in the detailed scratchpad. While addition, subtrac-
 590 tion, and multiplication are decomposed into sim-
 591 pler functions, sine and square root involve more in-
 592 tricate operations. See Appendix 18.5 for a broader
 593 analysis of the error profile, and Appendix 12 for

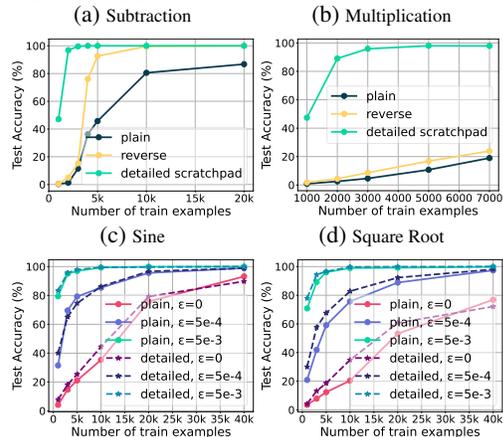


Figure 6: Performance of 3-digit subtraction, 2-digit multiplication, 4-digit precision sine and square root with varying data formats.

594 detailed experimental setup and results on jointly
 595 training on all five arithmetic tasks.

596 **Mixing Text with Arithmetic Data.** While
 597 the models so far were trained exclusively on
 598 arithmetic tasks, in practice, LLMs utilize a combination of arithmetic and *text* data for training.
 599 How does that affect the emergence of arithmetic skills? To explore that we incorporate
 600 both addition samples and text into our train data and evaluate the models with few-shot prompting (showing a few examples of addition in the prompt) to see if it is able to be effectively conditioned for the appropriate context (arithmetic/text generation). As we see in Figure 7, we
 601 find that few-shot prompting improves the performance of the model, allowing it to perform
 602 addition accurately even in the plain format.
 603
 604
 605
 606
 607
 608
 609
 610

611 Intriguingly, accuracy remains high using *plain* even with the inclusion of a text prompt preceding
 612 “A+B=”. This is likely due to the structure of our mixed dataset where addition examples are
 613 interspersed within Shakespeare text. With the incorporation of more addition examples, instances
 614 where addition follows Shakespeare text increases, leading to a decrease in potential inconsistencies
 615 when text content is present during addition test queries. We further analyze the effect of text on
 616 prompting for both cases with and without text in the training data in Appendix 13.

617 10 Fine-tuning, Scaling, and Pretraining in Larger Models

618 We extend our study from NanoGPT to larger
 619 models like GPT-2 and GPT-3 to explore the
 620 impact of pretraining and model size. Initially,
 621 we compare the performance of NanoGPT and
 622 GPT-2, both trained from scratch. This highlights the advantages of larger model scales,
 623 especially in zero-shot scenarios. We then fine-tune a pretrained GPT-3 on various arithmetic
 624 tasks using different data formats, reaffirming
 625 the importance of data formatting for larger
 626 pretrained models.
 627
 628

629 **Comparing NanoGPT and GPT-2.** We repeat our experiments on a GPT-2 model with 85M
 630 parameters, with twice as many layers, heads, and embedding size compared to NanoGPT. We
 631 train the model from scratch using character-level tokenization, jointly on text and addition tasks,
 632 adopting both plain and detailed scratchpad formats as in Section 9. The results depicted in Figure 8
 633 demonstrate that the larger model outperforms in both plain and detailed scratchpad evaluations.

634 **GPT-3 experiments.** We consider three GPT-3 variants: Ada, Curie, and Davinci (OpenAI). We
 635 fine-tune these models using the same four data formatting methods as our NanoGPT experiments
 636 except that we introduce spaces between numbers in *plain* and *reverse* formatting to ensure consistent
 637 tokenization of numbers.

Table 2: Evaluation of addition performance for fine-tuned GPT-3 models: Davinci, Curie, and Ada. In each case, the model is finetuned on 1000 samples of addition in the corresponding format.

GPT-3 Model	Zero-shot	Plain	Reverse	Simplified Scratchpad	Detailed Scratchpad
Davinci	2%	34%	80.9%	88.7%	99.5%
Curie	0.0%	1.4%	12.3%	10.7%	99.7%
Ada	0.0%	0.3%	6.3%	0.6%	99.8%

638 The results in Table 2 show that starting with pretrained GPT-3 significantly improves performance
 639 compared to training NanoGPT or GPT-2 from scratch with only 1000 examples (Figure 4a). Similar

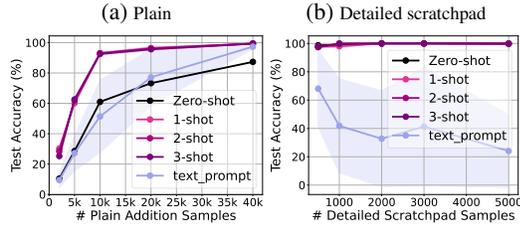


Figure 7: Performance of NanoGPT model trained with the Shakespeare dataset, addition dataset in plain, and detailed scratchpad format. The number of plain (left) and detailed scratchpad (right) formatted samples are varied. Performance is evaluated on zero-shot, few-shot, and text prompts, with the shaded area representing the standard deviation across various prompt exemplar sets.

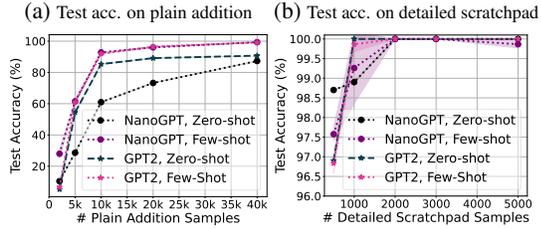


Figure 8: Comparing NanoGPT and GPT-2 trained jointly on the Shakespeare dataset and addition tasks using plain and detailed scratchpad formatting. Larger model scale and using few-shot prompting enhances performance.

640 to the result of training NanoGPT from scratch, the modified formats all outperform the plain format.
 641 Detailed scratchpad data achieves near-perfect accuracy, albeit with increased training and inference
 642 costs due to higher context length requirements. For our detailed experimental setup and further
 643 experiments on fine-tuning GPT-3 refer to Appendix 14.

644 11 Extending to Longer Digit Addition

645 In this section, we extend our experiments beyond 3-digit addition and explore longer-digit settings,
 646 ranging up to 10 digits. Our aim is to investigate whether our previous findings regarding the sample
 647 efficiency of reverse and scratchpad formats hold true for larger numbers of digits.

648 We begin by observing that the phase transition behavior observed in previous sections also applies to
 649 longer-digit addition. Furthermore, we discover that the advantages of using reverse and scratchpad
 650 formats become even more pronounced as the number of digits increases. Next, we examine the
 651 number of training samples required to learn $k + 1$ digit addition when fine-tuning a pretrained
 652 model trained on k digit addition. We find that while the number of samples needed to further learn
 653 $k + 1$ digit addition remains relatively consistent for reverse and scratchpad formats, the plain format
 654 requires an increasing number of samples.

655 **Experimental setup and data generation.** To explore the performance of the model in higher-digit
 656 addition scenarios, we extend the experimental setup described in Section 2. We adopt a balanced
 657 sampling approach for training data with D digits, ensuring an equal number d of all combinations of
 658 digits for both operands as follows:

659 We begin by sampling all 100-digit additions. For the remaining number of digits, ranging from
 660 2 to D , we generate addition examples of the form “ $A + B = C$ ”. The two operands, A and B ,
 661 are randomly sampled $d = \lfloor (N - 100) / (D(D + 1) / 2 - 1) \rfloor$ times for every D , where N is the
 662 total number of training examples. Operand A is sampled between $[10^{k_1 - 1}, 10^{k_1} - 1]$ and operand
 663 B is sampled between $[10^{k_2 - 1}, 10^{k_2} - 1]$, for all $1 \leq k_1 \leq k_2 \leq D$, excluding the case where
 664 $k_1 = k_2 = 1$. After sampling the two operands, we randomly interchange them to cover cases where
 665 A has fewer digits than B and vice versa.

666 11.1 Training from Random Initialization

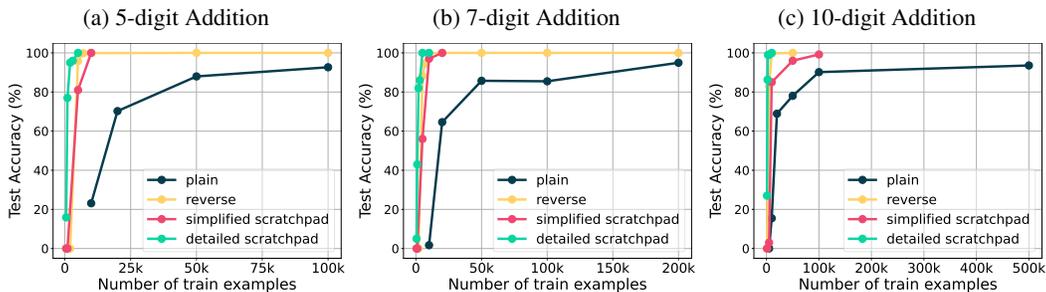


Figure 9: Comparison of sample efficiency for 5, 7 and 10-digit additions: performance of models trained with varying numbers of addition samples on each data format. The plain format data requires an increasing number of training examples for higher digits, while the number of samples required for other methods remains relatively consistent.

667 We repeat the experiment from Section 2 on nanoGPT with longer digits. The results shown in Figure 9
 668 demonstrate a similar behavior to the findings observed in Figure 4a for 3-digit addition. This indicates
 669 that our previous observations generalize to longer sequence lengths. Notably, the performance gap
 670 between the modified formats (reverse, simplified scratchpad, and detailed scratchpad) and the plain
 671 format becomes even more significant in the context of higher digits. While the plain format requires
 672 an increasing number of training examples to learn higher-digit additions, the reverse or scratchpad
 673 formats exhibit a more consistent requirement in terms of the number of training examples.

674 This prompts us to explore the differences between each format in a fine-tuning setting. Specifically,
 675 we ask whether a model trained on reverse or scratchpad-formatted k digit addition data would find it
 676 easier to learn $k + 1$ digit addition compared to a model trained with plain format addition.

677 **11.2 Fine-Tuning from Pretrained Models**

678 In this section, we investigate the generalization ability of transformer models, specifically focusing
 679 on their capacity to learn higher-digit additions based on their knowledge of lower-digit additions.
 680 Additionally, we explore how the choice of data format affects the number of samples required to
 681 learn higher-digit additions.

682 **Forgetting of k -digit addition when trained on $k + 1$ -digit addition.**

683 We begin by fine-tuning a model that was initially trained on 3-digit addition. We fine-tune this
 684 model using 4-digit addition training data, with each data format being used separately. To mitigate
 685 the ‘catastrophic forgetting’ phenomenon, we experiment with different learning rates, gradually
 686 reducing the magnitude. We continue this process until the learning rate becomes too small for the
 687 model to effectively learn 4-digit addition.

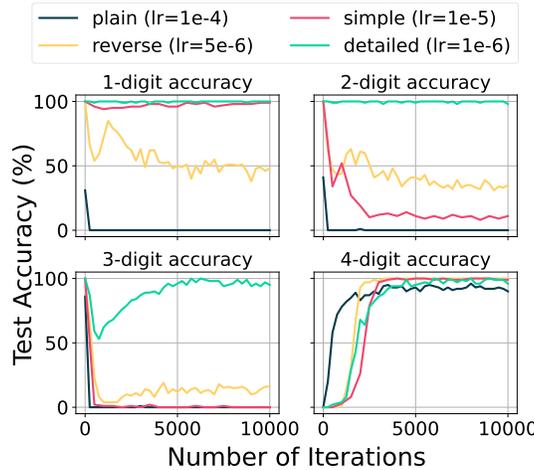


Figure 10: Accuracy of 1 to 4-digit additions during fine-tuning of a pretrained model on 3-digit additions using different data formats. The model is fine-tuned using only 4-digit addition data with corresponding formats. We observe that the plain format ‘forgets’ 1 to 3-digit additions entirely when learning 4-digit addition. In contrast, the detailed scratchpad method successfully learns 4-digit addition while maintaining high performance on 1 to 3-digit additions.

688 The results depicted in Figure 10 reveal interesting insights about the fine-tuning process. When
 689 training the model using the plain format with only 4-digit addition data, there is an immediate drop
 690 in accuracy for 1 to 3 digit additions. This indicates that the model experiences significant forgetting
 691 of previously learned additions. In contrast, the reverse and scratchpad methods exhibit a more
 692 favorable behavior. The model trained with these methods does not completely forget 1 or 2 digit
 693 additions while learning 4-digit addition. Remarkably, the detailed scratchpad method stands out by
 694 enabling the model to learn 4-digit addition without compromising its performance on 1 to 3 digit
 695 additions. Although there is a slight decrease in performance for 3-digit additions initially, the model
 696 quickly recovers and picks up the knowledge again as it trains on 4-digit additions.

697 This result can be explained by the hypothesis that learning a $k + 1$ digit addition from a k -digit
 698 model is an incremental process for the detailed scratchpad method. The model already has a solid
 699 foundation in understanding the intermediate steps involved in addition, so it only needs to adapt to
 700 longer sequences. In contrast, for the plain format, learning higher-digit additions requires the model
 701 to establish new mappings to generate correct outputs, which is a more challenging task.

702 **Sample efficiency of fine-tuning k -digit models with $k + 1$ -digit examples.** Building upon our
 703 previous findings that fine-tuning a model solely on $k + 1$ -digit addition leads to a loss in performance
 704 for k -digit addition, we modify our approach to prevent the loss of performance in the k -digit
 705 task. Instead of training solely on $k + 1$ -digit examples, we construct a dataset that includes all
 706 addition tasks from 1-digit to $k + 1$ -digit, with the method described in the previous section. By
 707 doing so, we aim to maintain the performance of 1 to k -digit addition while enabling the model to
 708 learn $k + 1$ -digit addition during fine-tuning.

709 In this experiment, we investigate the number of $k + 1$ -digit training examples required for the model
 710 to effectively learn $k + 1$ -digit addition when fine-tuning a pretrained model on k -digit addition. It is
 711 important to note that this setting differs from the previous section (Section 11.1), where we focused
 712 on training models from random initialization. Here, we specifically focus on the fine-tuning process.
 713 We fine-tune individual models pretrained on each data format (using k -digit addition) and further
 714 train them using the same data format on a new dataset that includes all addition examples from
 715 1-digit to $k + 1$ -digit.

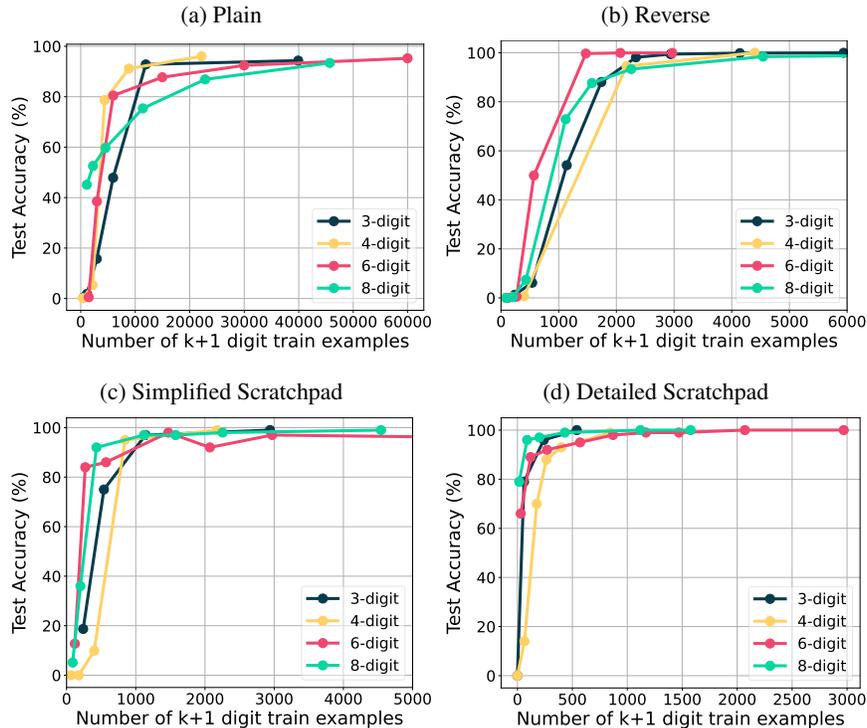


Figure 11: Fine-tuning performance of pretrained k -digit models using varying numbers of $k + 1$ -digit examples, with corresponding data formats. The plain format requires an increasing number of $k + 1$ -digit examples as the number of digits ($k + 1$) increases. In contrast, the modified formats (reverse, scratchpad) exhibit consistent performance across different numbers of digits, requiring a relatively consistent number of examples to learn the additional digit.

716 The results in Figure 11 demonstrate the number of $k + 1$ -digit addition samples required for a
 717 pretrained model capable of performing k -digit addition to learn the addition of $k + 1$ digits. The
 718 findings reveal that modified formats (reverse, scratchpad) require a relatively small number of
 719 samples (between 1000 and 5000) to learn the addition of an extra digit. In contrast, the plain format
 720 necessitates a significantly larger number of training examples, with the requirement increasing as
 721 the number of digits grows.

722 This observation aligns with our previously established Lemma 2 and Lemma 1, which suggest that
 723 learning higher-digit addition in the reverse format involves processing the i -th digit of the operands
 724 and carrying from the previous position. This operation *remains consistent regardless of the number*
 725 *of digits being added*. As a result, the model primarily needs to learn how to handle longer digits to
 726 perform addition effectively.

727 In contrast, the plain addition format requires the model to learn a more complex function that
 728 incorporates all digits from both operands. As the number of digits increases, the complexity of
 729 this function grows as well. This highlights the greater difficulty faced by the plain format in
 730 accommodating additions with a larger number of digits.

731 11.3 Impact of Formats on Fine-Tuning

732 We delve deeper into the impact of different formats on the fine-tuning process. Specifically, we
 733 investigate whether training a model in one format helps in learning addition in another format,

734 and vice versa. To conduct this analysis, we begin with a model trained on each data format using
 735 3-digit addition examples. We then individually fine-tune these pretrained models using different
 736 data formats, on 4-digit addition examples.

737 The results depicted in Figure 12 highlight some interesting findings. Firstly, we observe that a
 738 model trained with the same format as the fine-tuning format exhibits faster learning in terms of
 739 the number of iterations. For instance, training a model with the plain format outperforms
 740 training a model pretrained with scratchpad formats. For instance, training a model with the plain
 741 format outperforms training a model pretrained with scratchpad formats. This suggests that the model
 742 benefits from the consistency and familiarity provided by the same format throughout the training
 743 process. This suggests that the model benefits from the consistency and familiarity provided by the
 744 same format throughout the training process. This suggests that the model benefits from the
 745 consistency and familiarity provided by the same format throughout the training process.
 746

747 Additionally, we notice that fine-tuning a detailed scratchpad pretrained model on other
 748 formats proves to be more challenging. This observation can be attributed to the need for the
 749 model to “unlearn” the intricacies of the verbose detailed scratchpad format and adapt to
 750 the new format. For example, the plain format does not involve the use of alphabet characters
 751 in the data, so a model pretrained with the plain format would have a low probability of gener-
 752 ating alphabetic outputs. In contrast, a detailed scratchpad pretrained model would have encoun-
 753 tered various alphabets and may have a tendency to output them. Therefore, adjusting to a new
 754 format requires additional effort for the model to “unlearn” the patterns specific to the previous
 755 format and effectively learn the new format it is being trained on.

764 These findings highlight the importance of considering format consistency during the fine-tuning
 765 process, as it can impact the efficiency and effectiveness of the learning process. We will delve
 766 further into this topic in the upcoming section 10, where we fine-tune pretrained GPT-3 models.
 767 Notably, we observe that fine-tuning with reverse or simplified scratchpad formats actually yields
 768 worse results compared to fine-tuning with plain formats. For a detailed exploration of these
 769 observations, please refer to the forthcoming section.

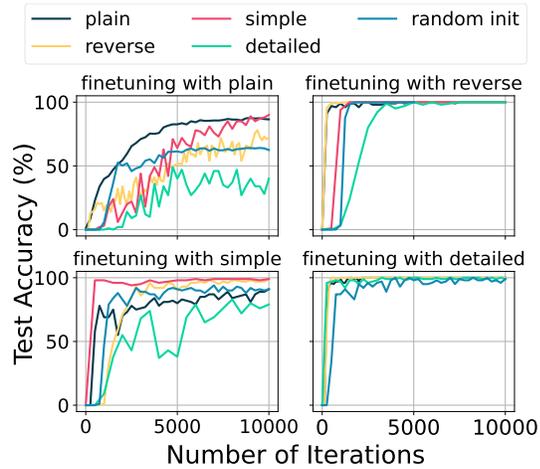


Figure 12: Performance of fine-tuning a 3-digit model trained on different data formats (plain, reverse, simple scratchpad, detailed scratchpad, and random initialization) individually with different data formats of 4-digit addition. The results demonstrate that fine-tuning yields the best performance when the pretrained model and the fine-tuning format are consistent. Notably, fine-tuning a detailed scratchpad format model shows suboptimal performance. We hypothesize that this is due to the need for the model to “unlearn” the rigid and verbose format and adapt to the new format.

770 12 Teaching Arithmetic Operations Beyond Addition

771 While this study has a primary focus on the *addition* operation and aims to comprehend the signifi-
772 cance of data sampling and formatting, its findings are applicable beyond the realm of addition alone.
773 In this section, we expand our examination to include other arithmetic operations, thus demonstrating
774 the broader applicability of our insights. We consider a mix of arithmetic tasks, including binary
775 operations like *subtraction* and *multiplication*, and unary operations such as *sine* and *square root*.
776 Each operation entails its unique challenges and intricacies. For instance, subtraction introduces the
777 concept of negative numbers, multiplication can generate significantly longer outputs, and sine and
778 square root functions entail computations involving floating-point numbers, which are considered up
779 to four digits of precision in our work.

780 We acknowledge that while our examination is detailed, it does not encompass all the fundamental
781 arithmetic operations or the entire scope of floating-point arithmetic. Specifically, our focus is primar-
782 ily on integer arithmetic for binary operations, considering a limited length of digits. Additionally,
783 for unary operations, we confine ourselves to a restricted number of digits below the decimal point.

784 In Section 12.1, we delve into each arithmetic operation individually, exploring the impact of
785 data formatting and determining the relevancy of our insights across disparate tasks. Further, in
786 Section 12.2, we perform an analysis of joint training across all five tasks, investigating the potential
787 performance implications for each individual task.

788 12.1 Extended Arithmetic Operations

789 In order to extend our analysis to arithmetic operations beyond addition, we consider the following
790 tasks:

791 **Subtraction** (−). We consider subtraction of positive numbers up to 3 digits, written as
792 $A_3A_2A_1 - B_3B_2B_1 = C_3C_2C_1$ in (i) plain formatting, and $\$A_3A_2A_1 - B_3B_2B_1 = C_1C_2C_3\$$ in (ii)
793 reverse formatting. As with addition, scratchpad-based methods (iii, iv), present the intermediate steps
794 of digit-wise subtraction and handling of carry-ons. These steps proceed from the least significant
795 bit (LSB) to the most significant bit (MSB). If the final result after computing all the digit-wise
796 subtractions is negative, we subtract the number in the most significant bit (MSB) position multiplied
797 by 10 to the power of (number of digits in the output - 1) from the remaining digits in the output. In
798 Section 18.3, we present an alternative version of the detailed scratchpad formatting for subtraction.

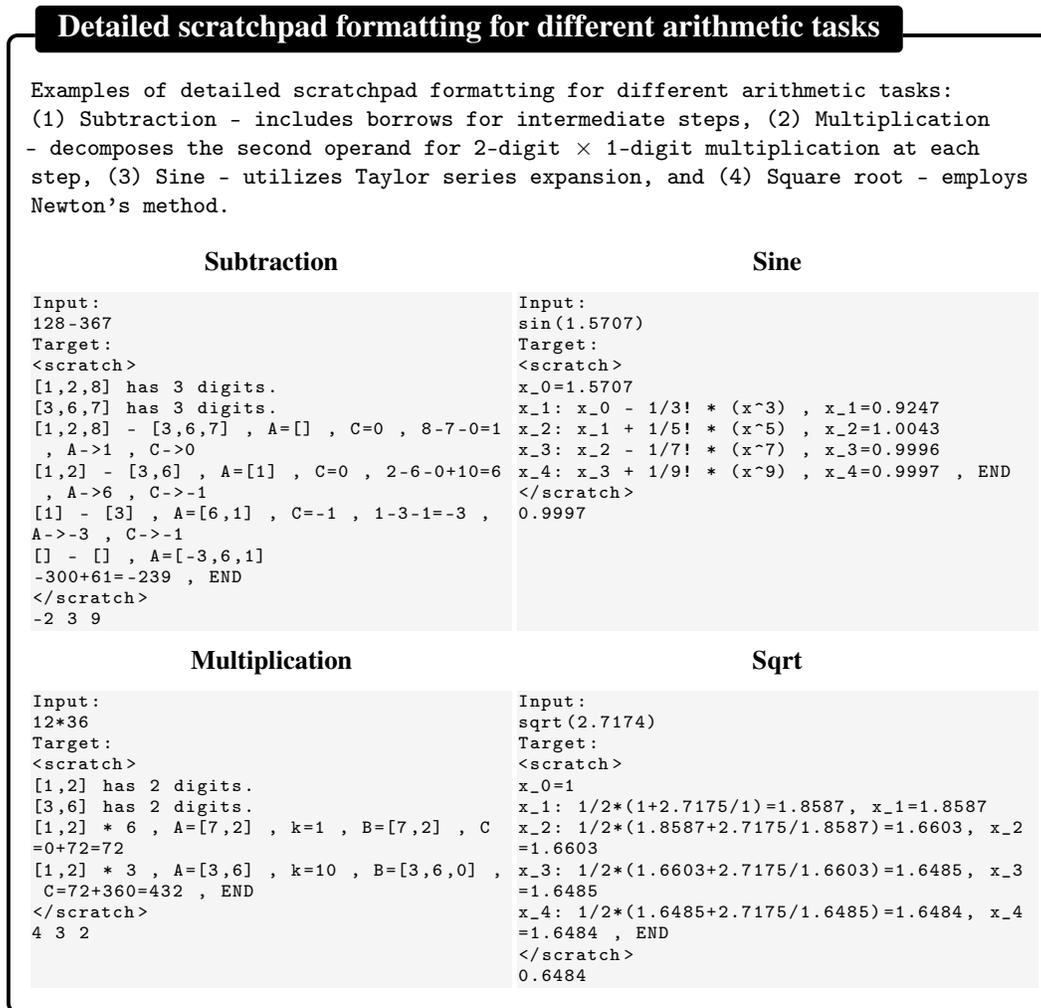
799 **Multiplication** (×). We consider multiplication of positive numbers up to 2-digits. (i) Plain
800 formatting examples are formatted as $A_2A_1 * B_2B_1 = C_4C_3C_2C_1$, while (ii) reverse formatting is
801 formatted as $\$A_2A_1 * B_2B_1 = C_1C_2C_3C_4\$$. The (iv) detailed scratchpad method simplifies each
802 intermediate step by conducting a series of multiplications between the first operand and each digit
803 of the second operand, starting from the least significant bit (LSB) and moving toward the most
804 significant bit (MSB). For each step, we multiply the result by an exponentiation of 10 corresponding
805 to the relative digit position.

806 **Sine** (sin). We consider decimal numbers within the range $[-\pi/2, \pi/2]$, truncated to 4-digit
807 precision. (i) Plain formatting examples are formatted as $\sin(A_0.A_1A_2A_3A_4) = B_0.B_1B_2B_3B_4$.
808 For (iv) detailed scratchpad method, we include the Taylor series expansion steps for sine, which
809 is represented as $\sin(x) = x - \frac{1}{3!}x^3 + \frac{1}{5!}x^5 - \frac{1}{7!}x^7 + \dots$. These intermediate steps involve
810 exponentiation, which may not be any easier to compute than the sine operation itself.

811 **Square Root** ($\sqrt{\cdot}$). We consider decimal numbers within $[1, 10)$, truncated to 4-digits of precision
812 with the format, written as $\text{sqrt}(A_0.A_1A_2A_3A_4) = B_0.B_1B_2B_3B_4$ for (i) plain formatting. For (iv)
813 detailed scratchpad method, we enumerate each step of Newton’s method to compute the square root
814 function. The iterative formula is given by $x_n = \frac{1}{2}(x_{n-1} + \frac{x}{x_{n-1}})$, where x_0 is initialized as the
815 floor of the square root value of the operand x . These intermediate steps involve a division operation,
816 which can be as complex as the square root operation itself.

817 For evaluation of sine and square root, we classify the result \hat{y}_i as correct if the absolute difference
818 between \hat{y}_i and the ground truth value y_i is less than or equal to a predefined threshold $\epsilon \geq 0$.

819 For each arithmetic task, we explore both the plain format and the detailed scratchpad format.
 820 The detailed scratchpad formatting for each task is illustrated in Figure 13 and Appendix 20. For
 821 subtraction, the process involves breaking down the operation into intermediate steps of digit-wise
 822 subtraction, including carry-ons when necessary. Unlike addition, subtraction requires an additional
 823 step to handle cases where the first operand is smaller than the second. Further details on the detailed
 824 scratchpad for subtraction can be found in Section 18.3. For multiplication, each intermediate step
 825 carries out a 2-digit \times 1-digit multiplication between the first operand and each separate digit of the
 826 second operand. For sine and square root, we utilize a sequence of *iterative approximations* instead
 827 of algorithmic explanations. Specifically, Taylor's series expansion steps for sine and Newton's
 828 method steps for square root are used. It is important to note that while addition, subtraction, and
 829 multiplication are broken down into simpler operations at each step, CoT for sine and square root
 830 functions requires intermediate steps involving operations like exponentiation or division, which
 831 might not be inherently simpler.



832 Figure 13: Examples of the detailed scratchpad format for different arithmetic tasks such as subtraction, sine,
 833 multiplication, and square root.

834 The results depicted in Figure 6 indicate that similar to the findings of addition, the detailed scratchpad
 835 format significantly improves performance over *plain* or *reverse* formats and yields efficient results
 836 even with few samples for subtraction and multiplication tasks. Interestingly, we find *reverse* is not
 837 particularly effective in multiplication. On the other hand, the detailed scratchpad format exhibits
 838 reduced efficiency for \sin and $\sqrt{\quad}$ compared to other operations ($+$, $-$, \times). This discrepancy can be
 839 traced back to the complexity of the intermediate steps involved in the detailed scratchpad. While
 840 addition, subtraction, and multiplication are decomposed into simpler functions, sine and square

841 root operations involve more intricate operations. For a broader analysis of the error profile, see
 842 Appendix 18.5.

843 **12.2 Jointly Training on All Five Arithmetic Tasks**

844 So far, we only considered the problem of learning different arithmetic operations individually. In
 845 this section, we study the effect of jointly training on all five arithmetic tasks - addition, subtraction,
 846 multiplication, sine, and square root. We construct a single train dataset incorporating all
 847 task $\mathcal{D}_{\text{train}} = \{\mathcal{D}_{\text{train}}^+, \mathcal{D}_{\text{train}}^-, \mathcal{D}_{\text{train}}^\times, \mathcal{D}_{\text{train}}^{\sin}, \mathcal{D}_{\text{train}}^\sqrt{\cdot}\}$, and randomize the sequence of tasks in our train
 848 samples. For example, a randomly chosen segment of the training data may exhibit a task order
 849 such as $(+, -, \sin, -, \times, \times, \sqrt{\cdot}, \dots)$. We consider 10,000 training examples for each task of addition,
 850 subtraction, sine, and square root and 3,000 for multiplication.

851 The model’s performance, after training on our joint dataset $\mathcal{D}_{\text{train}}$, is evaluated in both zero-shot and
 852 few-shot settings. These results are also compared with the performance of models that were trained
 853 separately on each dataset $(\mathcal{D}_{\text{train}}^+, \mathcal{D}_{\text{train}}^-, \mathcal{D}_{\text{train}}^\times, \mathcal{D}_{\text{train}}^{\sin}, \mathcal{D}_{\text{train}}^\sqrt{\cdot})$, identical to those used to construct
 854 $\mathcal{D}_{\text{train}}$. In the few-shot setting, each task is given examples from any of the five arithmetic tasks (not
 855 necessarily related to the test task under consideration) or prompt texts, followed by test queries
 856 specific to the task of interest. For further details on the few-shot prompting methods used, please
 857 refer to Section ??.

858 Table 3 shows that joint training significantly enhances the zero-shot performance for multiplication
 859 and square root tasks, yet it slightly reduces the performance for subtraction. Generally, few-shot
 860 prompting exhibits improved performance. Notably, the performance of few-shot prompting remains
 861 consistent regardless of whether the exemplars provided are from unrelated tasks or are task-specific.
 862 We propose that this consistency is due to our randomized task sequence during training, which
 863 presents the model with numerous instances where one task directly follows another, thus simulating
 864 few-shot prompting with different tasks. Furthermore, we observe that text prompting performs
 865 similar to zero-shot. We conjecture that this is because the training data does not include text data
 866 and the model has never encountered text and therefore, text prompting serves as a random prefix
 867 attached to our test query.

Table 3: Performance of models trained individually and jointly on five arithmetic tasks. The threshold ϵ for \sin and $\sqrt{\cdot}$ functions is set to 0. For the models trained jointly on all five tasks, we evaluate their performance in both a zero-shot setting and a few-shot setting. In the few-shot setting, each task is presented with exemplars from one of the five arithmetic tasks or prompted with text, followed by task-specific test queries. The results show that few-shot prompting with any arithmetic operators (even unrelated to the test task) generally improves performance. However, text prompting shows performance similar to the zero-shot setting.

	Trained on individual task	Trained jointly on all 5 tasks						
		Zero-shot	Few-shot exemplar format					
			+	-	\times	sin	sqrt	text
+	84.06	87.96	96.45	96.90	96.92	97.06	97.01	88.71
-	79.97	72.83	81.28	79.59	81.39	81.84	81.74	68.91
\times	4.58	14.28	18.86	18.96	15.43	19.20	19.59	15.48
sin	35.03	34.74	34.35	34.31	34.34	32.64	33.42	33.96
sqrt	19.85	27.37	26.65	26.74	26.70	25.60	25.61	26.02

868 **13 Mixing Text with Arithmetic Data**

869 Until now, our focus was primarily on models trained exclusively on arithmetic tasks. However,
 870 in practice, large language models (LLMs) utilize a combination of arithmetic and *text* data for
 871 training. In this section, we broaden our scope by incorporating both addition samples and text into
 872 our pretraining data. We then evaluate the trained models with various few-shot prompts to analyze if
 873 the model is able to effectively identify the correct context.

874 **Experimental Setup.** We mix addition and text data in our experiment using the Shakespeare
 875 dataset (Karpathy, 2015) that includes 1, 115, 394 tokens of text, 10, 000 plain addition examples
 876 (120, 027 tokens), and 3, 000 detailed scratchpad formatted addition examples (813, 510 tokens). We
 877 fix the number of detailed scratchpad examples and plain addition examples (3, 000 and 10, 000
 878 respectively) while varying the number of each example type in the training process. The Shakespeare
 879 text is segmented into dialogue chunks, with a random number of addition data inserted between
 880 them. We use a character-level tokenizer with a vocabulary size of 80, containing all characters
 881 present in the dataset, including alphabets, digits, and certain symbols like +, = and \n.

882 **13.1 Few-Shot Prompting**

883 Given the mixed nature (arithmetic and text) of our dataset, introducing relevant examples seems an
 884 effective strategy to prime the model to generate the desired type of output. To assess the performance
 885 of such few-shot (1/2/3-shot) prompting, we provide task-specific exemplars as illustrated in
 886 Figure 14. Plain addition formatted exemplars are used for testing plain addition inputs, while
 887 detailed scratchpad formatted exemplars are utilized for assessing performance on detailed scratchpad
 888 formatted inputs. Additionally, we experiment with demonstrating text (see Appendix 13.3. for
 889 details) before querying addition (which we denote, Text-prompt). For each 1/2/3-shot and text
 890 prompting, average performance is reported over a fixed set of exemplars. Standard deviations of
 891 these prompts are denoted by shaded areas in the plots. The term “few-shot” refers to the reported
 892 mean of all 1/2/3-shot prompting results.

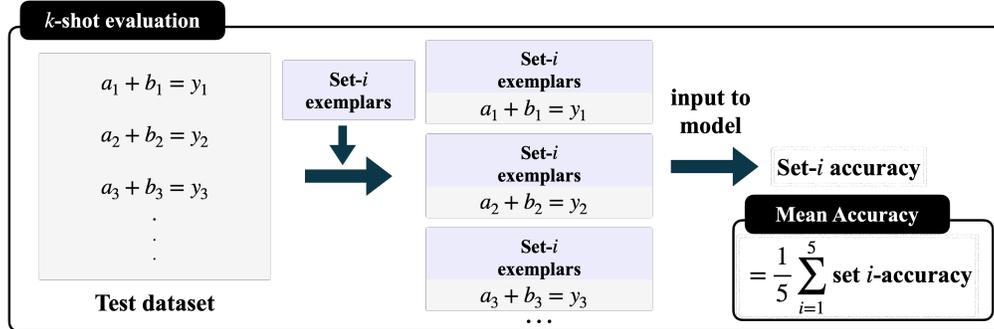


Figure 14: Few-shot prompting method. Few-shot prompting performance is evaluated by presenting relevant exemplars of addition and detailed scratchpad formatted inputs. Each 1/2/3-shot prompting is tested on a fixed five set of exemplars, and the accuracy is averaged over these evaluations.

893 Figure 15 shows that few-shot prompting directs the enhancement of performance, thereby allowing
 894 plain addition to perform almost perfectly with 40,000 train samples. Intriguingly, performance
 895 remains high on plain addition even with the inclusion of a text prompt, given a substantial number
 896 of addition examples. We hypothesize that this is due to the structure of our mixed dataset where
 897 addition examples are interspersed within Shakespeare data. With the incorporation of more addition
 898 examples, instances where addition examples directly follow Shakespeare text increase, leading to a
 899 decrease in potential inconsistencies when text content is present during addition test queries.

900 **13.2 Disentangling the effect of text on prompting**

901 To disentangle the effects of the textual content in the training data, we train a model strictly on plain
 902 addition, utilizing an enlarged vocabulary that also includes alphabet characters, thereby enabling text
 903 prompting. (Note that previous experimental settings on plain formatted additions used a vocabulary
 904 size of 13, which only includes 10 numerals and 3 symbols - “+”, “=”, “\n”). We introduce a variant

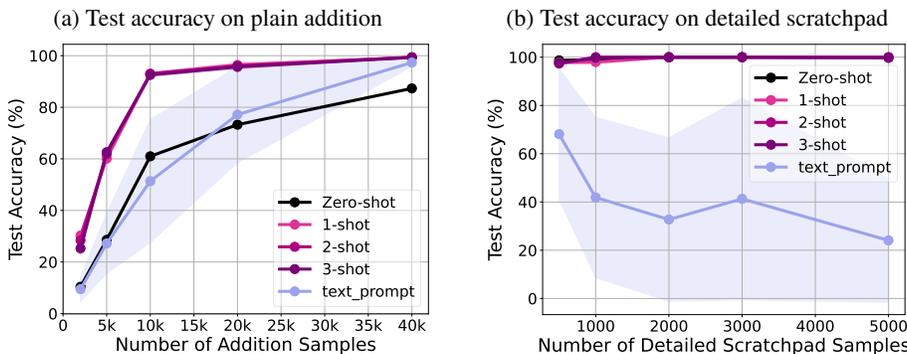


Figure 15: Performance of NanoGPT model trained with the Shakespeare dataset, addition dataset in plain, and detailed scratchpad format. The number of plain (left) and detailed scratchpad (right) formatted addition samples are varied. Performance is evaluated on zero-shot, few-shot, and text prompts, with the shaded area representing the standard deviation across various prompt exemplar sets. The results indicate a consistent enhancement in model performance using few-shot prompting.

905 of few-shot prompting, termed as **noisy-prompt**, which prompts the model with erroneous addition
 906 exemplars, *i.e.*, $A + B = C$, where $C \neq A + B$.

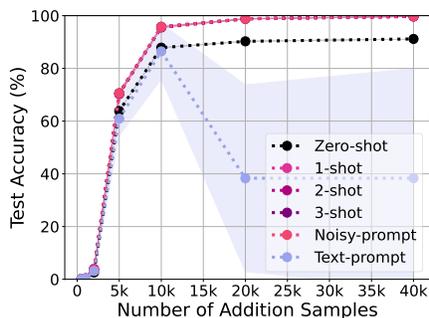


Figure 16: Performance of NanoGPT model trained exclusively on plain addition, but with an extended vocabulary including both addition and alphabets (vocabulary size = 80). Few-shot prompting, using both correct addition examples (1, 2, 3-shot) and incorrect addition examples (noisy-prompt) leads to enhanced performance, while the use of text prompts results in a degradation of performance when the model is trained solely on addition.

907 Figure 16 shows that few-shot prompting contributes to performance enhancement even when the
 908 model is confined to training on a single plain addition task. Even in the presence of noisy prompting,
 909 simply providing the model with the $A + B = C$ format yields performance nearly identical to few-
 910 shot prompting, aligning with the result observed by [Min et al. \(2022\)](#). Conversely, we notice that
 911 text prompts negatively influence performance when the model is trained only on addition. This
 912 finding reinforces our earlier observation in Figure 7 that the advantageous impact of text prompts
 913 originates from the combined text and addition data.

914 13.3 Prompting with Text

915 To extend on the few-shot prompting experiments from Section 12.2, we also evaluate the effect of
 916 prompting the model with pure-text prompts. If few-shot prompting with addition samples improves
 917 accuracy through in-context learning, we expect few-shot prompting with text to hurt accuracy since
 918 the text exemplars are *out-of-context*. We use five different types of text exemplars: (i) **Prompt1**: a
 919 short text prompt that is not present in the Shakespeare dataset, (ii) **Prompt2**: a short text prompt
 920 extracted from within Shakespeare dataset, (iii) **Prompt3**: a longer form text prompt extracted from
 921 within the Shakespeare dataset, (iv) **Prompt4**: a prompt that includes numbers, and (v) **Prompt5**: a
 922 long text prompt that is not present in the Shakespeare dataset. More details on the text prompts can
 923 be found in Figure 17.

Text prompts for few-shot experiments

Examples of the different text prompts used in the few-shot experiment. Each exemplar is separated by '---'.

Prompt 1. Short, ∉ Shakespeare

```
et tu brute
---
hello, world
---
how are you doing?
---
agi is coming
---
boom! stability
```

Prompt 2. Short, ∈ Shakespeare

```
JULIET:
Romeo!
---
All:
Resolved. resolved.
---
VOLUMNIA:
Why, I pray you?
---
CORIOLANUS:
May! prithee, woman,--
---
MENENIUS:
I mean, thy general.
```

Prompt 3. Long, ∈ Shakespeare

```
JULIET:
Romeo!
ROMEO:
My dear?
---
MENENIUS:
This is good news:
I will go meet the ladies. This Volumnia
Is worth of consuls, senators, patricians,
---
LADY ANNE:
Foul devil, for God's sake, hence, and trouble
us not;
For thou hast made the happy earth thy hell,
Fill'd it with cursing cries and deep exclaims
.
---
BUCKINGHAM:
I fear he will.
How now, Catesby, what says your lord?
---
CATESBY:
Bad news, my lord: Ely is fled to Richmond;
And Buckingham, back'd with the hardy Welshmen
,
Is in the field, and still his power
increaseth.
```

Prompt 4. Has number, ∉ Shakespeare

```
I go 16-12
That's the code to my heart, ah
I go 1-6-1-2
Star
---
Like a river flows 17-23
Surely to the sea 15-22
Darling, so it goes 46-92
Some things are meant to be
---
I got my first real 6-string
Bought it at the five and dime
Played it 'til my fingers bled
Was the summer of '69
---
I think someday I might just 5-3-2-1 get a real job
I spent half of my life 1-2-3 in a bus or on a flight
I'm getting off 17-36-8-2 the road and in a real job
---
Every time that 27-67-29 I look in the mirror
All these lines on my 1-3-92-5 face getting clearer
The past 45-5-3 is gone
```

Prompt 5. Long, ∉ Shakespeare

```
Is this the real life? Is this just fantasy? Caught in a landside, no escape from
reality.
Open your eyes, look up to the skies and see.
I'm just a poor boy, I need no sympathy. Because I'm easy come, easy go,
Little high, little low,
Any way the wind blows doesn't really matter to me, to me.
---
It's my life
And it's now or never
I ain't gonna live forever
I just want to live while I'm alive
My heart is like an open highway
Like Frankie said, I did it my way
---
Destruction leads to a very rough road but it also breeds creation
And earthquakes are to a girl's guitar, they're just another good vibration
And tidal waves couldn't save the world from Californication
```

```

---
I want to stay
But I need to go
I want to be the best for you
But I just don't know what to do
'Cause baby, say I've cried for you
The time we have spent together
Riding through this English whether
---
Lorem ipsum dolor sit amet, consectetur adipiscing elit. Vestibulum mattis in leo
vel gravida.
Pellentesque libero elit, scelerisque varius vehicula a, hendrerit et tellus.
Proin convallis neque nisl, nec lobortis est scelerisque tincidunt.
Nunc venenatis auctor urna.
Class aptent taciti sociosqu ad litora torquent per conubia nostra.

```

925

Figure 17: Text prompt exemplars for few-shot experiments.

926

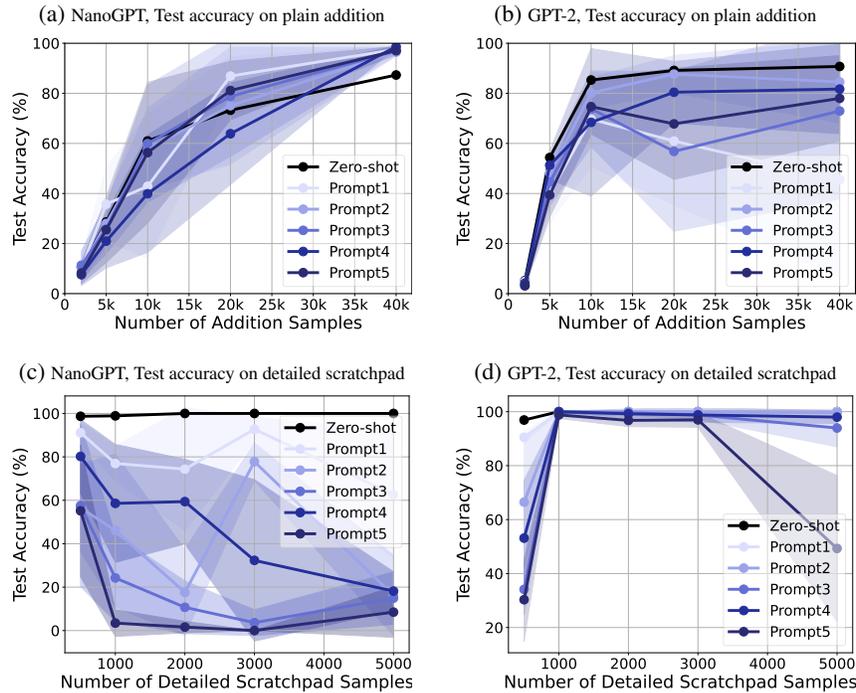


Figure 18: Experiments on few-shot prompting with different text prompts: (i) Prompt1: short text not in Shakespeare dataset (ii) Prompt2: short text within Shakespeare dataset (iii) Prompt3: long text within Shakespeare dataset (iv) Prompt4: text with numbers (v) Prompt5: long text not in the Shakespeare dataset. Each prompt (Prompt 1-5) consists of five distinct exemplars. The solid lines represent the mean performance across the five exemplars, while the shaded area indicates the standard deviation. We observe that the effectiveness of text prompts varies greatly depending on the exemplars used.

927 The results presented in Figure 18 show notable variations in evaluation accuracy for addition,
 928 depending on the chosen text prompts. Longer text prompts (Prompt 5) typically result in a more
 929 significant decline in performance. With the exception of NanoGPT trained on plain addition, the
 930 result in Figure 19 indicates that employing text prompts followed by test addition queries tends
 931 to have an adverse impact on the overall model performance, whereas incorporating relevant few-shot
 932 exemplars (1/2/3-shot) is beneficial. This aligns well with our intuition on the benefits on in-context
 933 learning.

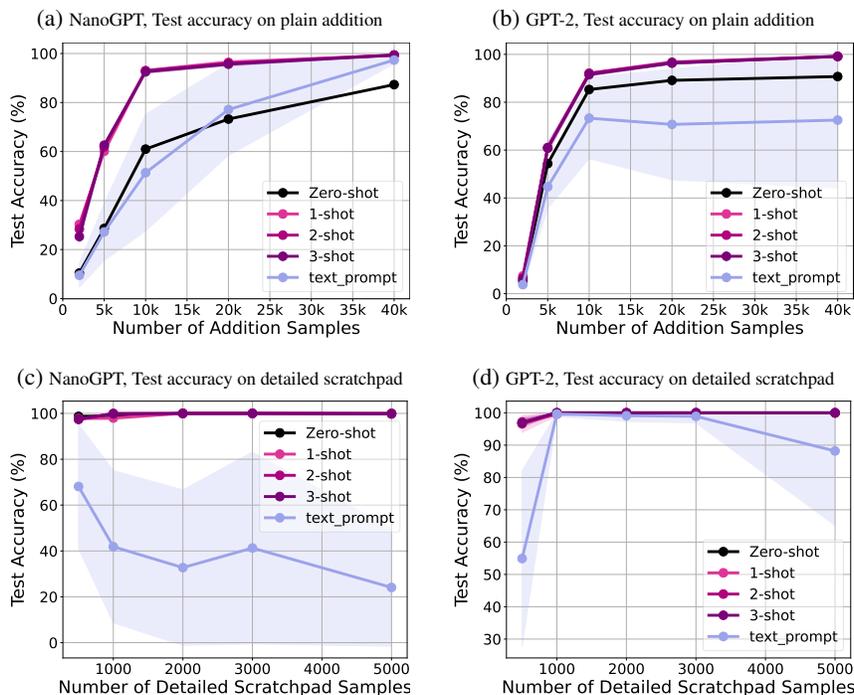


Figure 19: Performance of NanoGPT and GPT-2 model trained with entire Shakespeare dataset and a varying number of samples of plain addition, and addition with detailed scratchpad dataset. Performance is evaluated on test prompts formatted as plain addition and detailed scratchpad. Few-shot experiments are based on an average of 5 exemplars, while text prompts involve an average of 25 exemplars. The shaded area represents the standard deviation. Our observations indicate that few-shot prompting consistently improves performance, whereas text prompts generally have a negative impact.

934 14 Fine-tuning, Scaling, and Pretraining in Larger Models

935 This section focuses on bridging the gap between our experiments on NanoGPT and the more realistic
 936 setting of larger language models like GPT-2 and GPT-3. We begin by comparing the performance of
 937 NanoGPT and GPT-2 models when trained from random initialization. This comparison highlights
 938 the improved performance achieved with the larger model scale, especially in the zero-shot setting.
 939 Subsequently, we delve into the impact of tokenization methods and model pretraining in GPT-2
 940 models. Our exploration reveals the crucial role of pretrained models and the consistent tokenization
 941 of numbers (achieved by introducing spaces) during the training phase for arithmetic tasks. Building
 942 on these findings, we proceed to fine-tune a pretrained GPT-3 model on various arithmetic tasks,
 943 employing different data formats.

944 **Comparing NanoGPT and GPT-2.** To examine the impact of scale on arithmetic performance,
 945 we explore a larger GPT-2 model with 85 million parameters, featuring twice as many self-attention
 946 layers, heads, and embedding size compared to the previously used NanoGPT model. We train
 947 the GPT-2 model from scratch using character-level tokenization, jointly on text and addition tasks,
 948 adopting both plain and detailed scratchpad formats; an approach mirroring the setting in Section ??.
 949 The results depicted in Figure 8 demonstrate that the larger model outperforms in both plain and
 950 detailed scratchpad evaluations. For a comprehensive analysis of GPT-2, including few-shot learning
 951 and the influence of text prompts, refer to Figure 18 and Figure 19.

952 **Going from character-level tokenization to BPE.** The transition to a GPT-2 setup necessitates
 953 several modifications. Firstly, we shift to OpenAI’s Tiktoken BPE tokenizer, which is the default
 954 tokenizer for the pretrained GPT-2 model, featuring a vocabulary size of 50,257. We also examined
 955 two different training approaches: training the model from random initialization (scratch) and
 956 fine-tuning the pretrained model sourced from Huggingface. To ensure uniform digit tokenization,

957 alterations were made in data formatting to include spaces between numbers. This change aims to
958 circumvent potential inconsistent tokenization of numbers while utilizing the Tiktoken tokenizer.

959 Figure 20 shows that GPT-2 demonstrates high performance in addition tasks with both character-level
960 tokenization and Tiktoken with spaces between digits. This aligns with the results by Wallace et al.
961 (2019), suggesting that character-level tokenization exhibits stronger numeracy capabilities compared
962 to a word or sub-word methods. Furthermore, comparing the models trained from scratch and the
963 models trained from the pretrained model, we observe that fine-tuning a pretrained model results in
964 better performance compared to training a model from scratch.

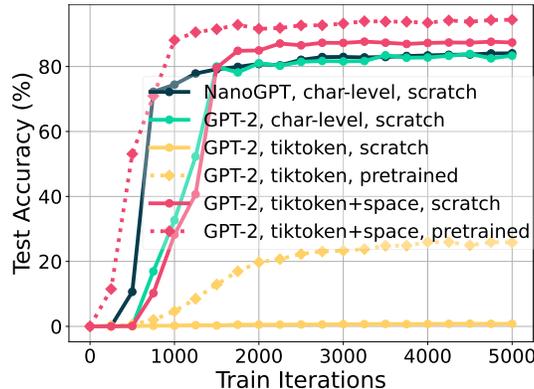


Figure 20: Performance of various configurations of the GPT-2 model on the addition task. We compare the effects of tokenization methods, specifically character-level tokenization versus Tiktoken (OpenAI’s BPE tokenizer), training initialization (training from scratch versus training from a pretrained GPT-2 model), and the inclusion or exclusion of spaces between numbers. The results highlight the significance of utilizing pretrained models and incorporating spaces for consistent tokenization of numbers when training a model for arithmetic tasks.

965 **GPT-3 experiments: Supervised fine-tuning.** We extend our experiments to verify if our observa-
966 tions hold while fine-tuning larger pre-trained models. In the following, we consider three GPT-3
967 variants: Ada, Curie, and Davinci. Note that since we perform fine-tuning using the OpenAI APIs, by
968 default *only the completions are loss generating tokens*. Therefore, these experiments are slightly
969 different when compared to the previous settings. We fine-tune these models using the same four
970 data formatting methods as our NanoGPT experiments: (i) *plain* formatting, (ii) *reverse* formatting,
971 (iii) *simplified scratchpad*, and (iv) *detailed scratchpad*. These formats are identical to those from our
972 NanoGPT experiments except for one aspect. We introduce spaces between numbers in *plain* and
973 *reverse* formatting to ensure consistent tokenization.

974 Due to budget constraints, all experiments were conducted using a fine-tuning dataset of 1,000
975 examples, and models were trained for 4 epochs. Performance evaluation was carried out on 1,000
976 examples that were disjoint from the training dataset. Note that this training scale is significantly
977 smaller than our experiments on NanoGPT, which employed 10,000 training examples for 5,000
978 iterations, with evaluations conducted on 10,000 test examples. However, given these models’
979 extensive pretraining on large data corpora, this scale can be deemed rational.

980 The results for addition and subtraction tasks are presented in Table 2 and Table 4, respectively. We
981 observed that initiating with a pretrained GPT-3 model significantly improves performance compared
982 to training NanoGPT or GPT-2 models from random initialization with only 1000 samples. This
983 indicates the utility of leveraging pretrained models for improved arithmetic performance. Interest-
984 ingly, while reverse formatting and simplified scratchpad formats improve addition performance,
985 they adversely affect subtraction performance. This observation is consistent with our earlier finding
986 depicted in Figure 12, wherein transitioning from one data format to another often results in lower
987 performance compared to initiating training from random initialization. We postulate that this discrep-
988 ancancy may be due to the pretrained GPT-3 model’s requirement to adapt to the reversed approach and
989 “unlearn” its knowledge of plain formatting arithmetic, thereby introducing additional complexity. On
990 the other hand, the detailed scratchpad method achieves excellent performance, albeit with increased
991 training and inference costs due to higher token requirements.

Table 4: Evaluation of subtraction performance for fine-tuned GPT-3 models: Davinci, Curie, and Ada. In each case, the model is finetuned on 1000 samples of addition in the corresponding format.

GPT-3 Model	Zero-shot	Plain	Reverse	Simplified Scratchpad	Detailed Scratchpad
Davinci	0.1%	84.8%	66.0%	15.4%	99.5%
Curie	0.1%	24.1%	6%	3.8%	92.5%
Ada	0.0%	3.7%	2.6%	3.4%	81.5%

Table 5: Evaluation of sine and square root performance for fine-tuned GPT-3 models: Davinci, Curie, and Ada. In each case, the model is finetuned on 1000 samples of addition in the corresponding format.

GPT-3 Model	eps	Sine			Square Root		
		Zero-shot	Plain	Detailed Scratchpad	Zero-shot	Plain	Detailed Scratchpad
Davinci	0	0%	11.0%	10.3%	0%	0.7%	4.6%
	5e-4	0%	35.9%	29.7%	0%	7.5%	17.2%
	5e-3	0.4%	85.5%	72.8%	0%	59%	60.5%
Curie	0	0.0%	8.6%	1.2%	0.0%	0.7%	2.1%
	5e-4	0.4%	32.7%	5.4%	0.1%	6.5%	6.0%
	5e-3	0.9%	80.8%	15%	0%	52.7%	30.2%
Ada	0	0.0%	5.8%	4.3%	0.0%	0.3%	2.7%
	5e-4	0.0%	21.4%	9.1%	0.0%	3.8%	11.9%
	5e-3	0.3%	67.8%	25.2%	0.0%	32.2%	45.8%

992 For the more complex sine and square root tasks as shown in Table 5, we found that training with
 993 only 1000 samples is insufficient to generate exact answers (eps=0). The GPT-3 model, fine-tuned
 994 with 1,000 samples, performs worse than the NanoGPT model trained with 10,000 samples. Further
 995 experiments with larger training datasets are necessary for deeper insights and improved performance
 996 on these tasks.

997 It is worth mentioning that while few-shot prompting notably improves the performance of all
 998 three GPT-3 models, their zero-shot performance is quite poor (as shown in the leftmost column
 999 of the tables). However, post-training, few-shot prompting becomes less effective as OpenAI’s
 1000 fine-tuning process trains the model on individual prompts and desired completions serially, rather
 1001 than in concatenation with multiple examples like in our NanoGPT experiments. Consequently, our
 1002 comparisons primarily focus on the **zero-shot performances** of each task.

1003 15 Token Efficiency Across Data Formats

1004 Figure 4a demonstrates that more detailed training data leads to improved sample efficiency. However,
 1005 this comparison does not account for the cost associated with training and inference. To address this,
 1006 we conduct a cost analysis based on the number of “unique” tokens encountered during training.
 1007 Each data sample is treated as a set of unique tokens, and the number of unique tokens is derived by
 1008 *multiplying the number of samples with the tokens per sample*. For instance, the mean token count for
 1009 a single training example in a 3-digit addition task is 13 for plain format, 15 for reverse format, 64 for
 1010 simplified scratchpad format, and 281 for detailed scratchpad format. Note that this calculation does
 1011 not evaluate uniqueness of tokens across samples *i.e.*, if the first sample is “112 + 129 = 241” and
 1012 the second sample is “112 + 128 = 240”, we will still consider that the model has seen 26 unique
 1013 tokens even though only two tokens differ across samples. This approach ensures our cost calculation
 1014 accounts for a vanilla implementation of attention with no additional optimizations (Pope et al., 2023).
 1015 Table 6 presents the number of tokens required for prompting and completion in each data format,
 1016 per example. Evidently, the detailed scratchpad method uses considerably more tokens compared to
 1017 other techniques.

1018 The result in Figure 4b indicates that reverse formatting is the most token-efficient approach. While
 1019 detailed scratchpad training is more sample efficient, it necessitates a larger number of tokens per
 1020 sample, both during training and inference. Given that the inference cost for commercial models is

1021 determined by the number of tokens utilized per inference call (sum of prompting and completion
 1022 tokens), abundant use of models trained on detailed scratchpad formats may escalate overall costs.
 1023 Furthermore, since the cost of a single forward pass is cubic in the number of tokens, this is important
 1024 to consider. Therefore, for practical usage, it is crucial to evaluate both the number of samples needed
 1025 for achieving the desired performance and the actual token demands during training and inference.

Table 6: Token requirements for prompting and completion per single example of 3-digit addition.

	Plain	Reverse	Simplified Scratchpad	Detailed Scratchpad
Prompt	8	9	23	23
Completion	5	6	41	258
Total	13	15	64	281

1026 **16 Length Generalization**

1027 In this section, we present results from experiments conducted to assess the model’s ability to gener-
 1028 alize across different digit lengths. Initially, we exclude training examples featuring 2-digit operands
 1029 from the 10,000-sample addition dataset, yielding a reduced dataset of 7,655 samples, consisting
 1030 solely of 1 or 3-digit operands. The model is trained with reverse format and its performance is evalu-
 1031 ated on test dataset containing 100 random samples of 1-digit, 2-digit, 3-digit, and 4-digit additions.
 1032 The results in Figure 21 demonstrate that the NanoGPT model is incapable of performing 2-digit
 1033 and 4-digit additions. This suggests an inherent necessity for exposure to all digit combinations to
 1034 perform accurate calculations and lacks generalization capabilities for unseen digit lengths.

1035 Additionally, we investigate the model’s ability to extrapolate over larger digit lengths. The model is
 1036 trained on 7-digit plain-formatted additions (each digit addition comprises 16650 samples, except
 1037 1-digit addition, which is trained on 100 samples). Its ability to add 8-digit numbers is then put to
 1038 test. The results in Figure 21 show that the model is unable to generalize to a greater number of digits
 1039 beyond what it has been trained on. Similarly, when training the model on 10-digit binary numbers, it
 1040 fails to generalize to 11-digit binary additions, further confirming its limited ability to handle unseen
 1041 digit combinations.

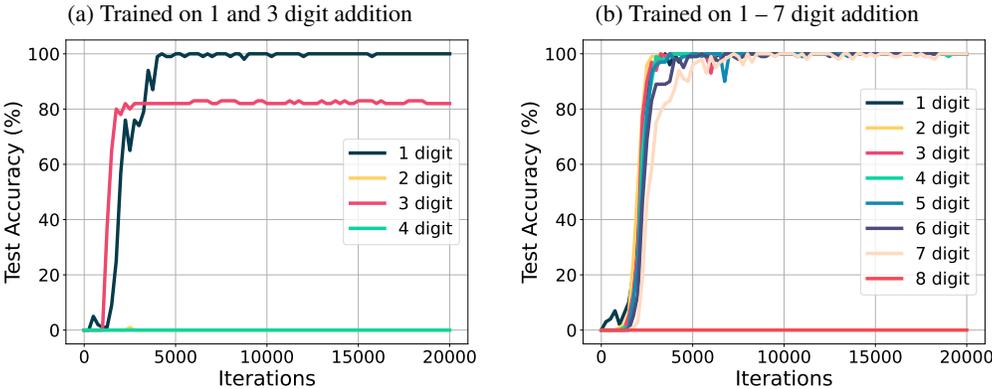


Figure 21: Generalization experiments testing NanoGPT’s performance on unseen numbers of digits in addition tasks. (Left): NanoGPT trained on reverse formatted addition with 1 and 3 digits, and tested on additions ranging from 1 to 4 digits. (Right): NanoGPT trained on up to 7-digit plain formatted addition and tested on additions ranging from 1 to 8 digits. In both cases, NanoGPT exhibits an inability to perform addition on digits it has not been exposed to.

1042 We further explore the impact of detailed scratchpad formatting. The model trained on additions of
 1043 up to 3 digits, struggles to generalize to 4-digit additions. Notably, it randomly drops a single digit
 1044 from the 4-digit number, erroneously perceiving it as a 3-digit number. We illustrate this difficulty
 1045 in Figure 22 through multiple detailed error cases, ranging from instances in which only the test
 1046 query is provided (Case 1) to scenarios where all intermediate steps are provided except only the final
 1047 answer (Case 5). The prompts are highlighted in light grey and the responses generated by our trained

1048 NanoGPT model are highlighted in light green. These cases emphasize the model's shortcomings in
1049 accurately managing larger digit lengths.

Examples for length generalization prompts

Results obtained by prompting the NanoGPT model with larger digits than those it was trained on. The model is trained using detailed scratchpad formats with 3-digit numbers. We evaluate its performance on 4-digit numbers, with varying levels of provided information. The prompt input is highlighted in a light blue box, while the model's output is highlighted in a light green box.

Case 1: Just asking the question

```
Input:
8465+3541
Target:
<scratch>
[8,4,6] has 3 digits.
[3,5,1] has 3 digits.
[8,4,6] + [3,5,1] , A=[], C=0 , 6+1+0=7 , A->7 , C->0
[8,4] + [3,5] , A=[7] , C=0 , 4+5+0=9 , A->9 , C->0
[8] + [3] , A=[9,7] , C=0 , 8+3+0=11 , A->1 , C->1
[] + [] , A=[1,9,7] C=1 , END
</scratch>
1 1 9 7
```

Case 2: Giving number of digits

```
Input:
8465+3541
Target:
<scratch>
[8,4,6,5] has 4 digits.
[3,5,4,1] has 4 digits.
[8,4,6] + [3,5,4] , A=[], C=0 , 6+4+0=10 , A->0 , C->1
[8,4] + [3,5] , A=[0] , C=1 , 4+5+1=10 , A->0 , C->1
[8] + [3] , A=[0,0] , C=1 , 8+3+1=12 , A->2 , C->1
[] + [] , A=[2,0,0] C=1 , END
</scratch>
1 2 0 0
```

Case 3: Giving 1 intermediate steps

```
Input:
8465+3541
Target:
<scratch>
[8,4,6,5] has 4 digits.
[3,5,4,1] has 4 digits.
[8,4,6,5] + [3,5,4,1] , A=[], C=0 , 5+1+0=6 , A->6 , C->0
[8,4] + [3,5] , A=[6] , C=0 , 4+5+0=9 , A->9 , C->0
[8] + [3] , A=[9,6] , C=0 , 8+3+0=11 , A->1 , C->1
[] + [] , A=[1,9,6] C=1 , END
</scratch>
1 1 9 6
```

Case 4: Giving all but one intermediate steps

```
Input:
8465+3541
Target:
<scratch>
[8,4,6,5] has 4 digits.
[3,5,4,1] has 4 digits.
[8,4,6,5] + [3,5,4,1] , A=[], C=0 , 5+1+0=6 , A->6 , C->0
[8,4,6] + [3,5,4] , A=[6] , C=0 , 6+4+0=10 , A->0 , C->1
[8,4] + [3,5] , A=[0,6] , C=1 , 4+5+1=10 , A->0 , C->1
[8] + [3] , A=[0,0,6] , C=1 , 8+3+1=12 , A->2 , C->1
[] + [] , A=[2,0,6] C=1 END
</scratch>
1 0 0 6
```

Case 5: Giving all intermediate steps, all it has to do is aggregate the answer

```
Input:
8465+3541
Target:
```

1050

```

<scratch>
[8,4,6,5] has 4 digits.
[3,5,4,1] has 4 digits.
[8,4,6,5] + [3,5,4,1] , A=[] , C=0 , 5+1+0=6 , A->6 , C->0
[8,4,6] + [3,5,4] , A=[6] , C=0 , 6+4+0=10 , A->0 , C->1
[8,4] + [3,5] , A=[0,6] , C=1 , 4+5+1=10 , A->0 , C->1
[8] + [3] , A=[0,0,6] , C=1 , 8+3+1=12 , A->2 , C->1
[] + [] , A=[2,0,0,6] C=1 END
</scratch>
1 2 0 6

```

Case 6: Giving 1 in-context example

```

Input:
8465+3541
Target:
<scratch>
[8,4,6,5] has 4 digits.
[3,5,4,1] has 4 digits.
[8,4,6,5] + [3,5,4,1] , A=[] , C=0 , 5+1+0=6 , A->6 , C->0
[8,4,6] + [3,5,4] , A=[6] , C=0 , 6+4+0=10 , A->0 , C->1
[8,4] + [3,5] , A=[0,6] , C=1 , 4+5+1=10 , A->0 , C->1
[8] + [3] , A=[0,0,6] , C=1 , 8+3+1=12 , A->2 , C->1
[] + [] , A=[2,0,0,6] C=1 , END
</scratch>
1 2 0 0 6
Input:
1946+3598
Target:
<scratch>
[1,9,4] has 3 digits.
[3,5,8] has 3 digits.
[1,9,4] + [3,5,8] , A=[] , C=0 , 4+8+0=12 , A->2 , C->1
[1,9] + [3,5] , A=[2] , C=1 , 9+5+1=15 , A->5 , C->1
[1] + [3] , A=[5,2] , C=1 , 1+3+1=5 , A->5 , C->0
[] + [] , A=[5,5,2] C=0 , END
</scratch>
5 5 2

```

Case 7: Giving 1 In-context example, and all intermediate steps

```

Input:
8465+3541
Target:
<scratch>
[8,4,6,5] has 4 digits.
[3,5,4,1] has 4 digits.
[8,4,6,5] + [3,5,4,1] , A=[] , C=0 , 5+1+0=6 , A->6 , C->0
[8,4,6] + [3,5,4] , A=[6] , C=0 , 6+4+0=10 , A->0 , C->1
[8,4] + [3,5] , A=[0,6] , C=1 , 4+5+1=10 , A->0 , C->1
[8] + [3] , A=[0,0,6] , C=1 , 8+3+1=12 , A->2 , C->1
[] + [] , A=[2,0,0,6] C=1 , END
</scratch>
1 2 0 0 6
Input:
1946+3598
Target:
<scratch>
[1,9,4,6] has 4 digits.
[3,5,9,8] has 4 digits.
[1,9,4,6] + [3,5,9,8] , A=[] , C=0 , 6+8+0=14 , A->4 , C->1
[1,9,4] + [3,5,9] , A=[4] , C=1 , 4+9+1=14 , A->4 , C->1
[1,9] + [3,5] , A=[4,4] , C=1 , 9+5+1=15 , A->5 , C->1
[1] + [3] , A=[5,4,4] , C=1 , 1+3+1=5 , A->5 , C->0
[] + [] , A=[5,5,4,4] C=0 , END
</scratch>
5 5 4

```

1051

Figure 22: Example results on the model's output when prompted with a larger number of digits than those it was trained on.

1052

1053 **17 Proofs**

1054 Here, we present the proofs of Lemma 1 and 2.

1055 **Lemma 1.** *Let A and B be two n -digit numbers, and let $C = A + B$. Suppose an algorithm A*
 1056 *outputs the digits of C in decreasing order of significance, then A must have access to all digits of A*
 1057 *and B starting from the first digit that it outputs.*

1058 *Proof.* We begin by assuming for contradiction that there does exist an algorithm Algo that does
 1059 not have access to all digits of A and B and still outputs $C = A + B$ correctly for all n -digit
 1060 numbers A, B . Without loss of generality, say Algo does not have access to the k -th digit of A
 1061 where $k \in [n]$ represents the position counting from the least significant digit. Then consider the
 1062 example $B = (10^n - 1)$ and $(A = 000 \dots A_k 00 \dots 0)$ where B is just the integer with n 9's and A
 1063 is just 0's with A_k in the k th position. If $A_k = 0$, then $C_{n+1} = 0$, but if $A_k = 1$, then $C_{n+1} = 1$.
 1064 Therefore, without access to the k -th digit of A , there exist examples where the algorithm will surely
 1065 make a mistake. Therefore, by contradiction such an Algo cannot exist. \square

1066 **Lemma 2.** *There exists an algorithm that computes $C = A + B$ for two n -digit numbers A and B*
 1067 *and outputs its digits in increasing order of significance such that, at each position i , the algorithm*
 1068 *only requires access to the i^{th} digits of A and B , as well as the carry-on from the previous position.*

1069 *Proof.* First note that the trivial algorithm for addition is exactly the proof of this Lemma. However,
 1070 we present a more formal argument below for completeness. Let A, B be n -digit numbers and
 1071 $C = A + B$ be at most an $(n + 1)$ digit number. Define the digits of A, B , and C as A_i, B_i , and C_i ,
 1072 respectively, for $i \in [n]$ counting from the least significant digit once again. Then, the addition can be
 1073 performed using the following steps. First, $C_i = (A_i + B_i + \text{carry}_i) \bmod 10$ where carry_i is the
 1074 carry-on from the addition of digits at position $i - 1$. If there is no carry from the previous position,
 1075 then $\text{carry}_i = 0$. The carry for the next position is then calculated as $\text{carry}_{i+1} = \left\lfloor \frac{A_i + B_i + \text{carry}_i}{10} \right\rfloor$.

1076 Putting this together, the algorithm for addition can be described as follows:

1077 **Step 1:** Set $\text{carry}_1 = 0$. **Repeat** for $i = \{1, \dots, n\}$: **Step 2:** Compute $C_i = (A_i + B_i + \text{carry}_i)$
 1078 $\bmod 10$ and $\text{carry}_{i+1} = \left\lfloor \frac{A_i + B_i + \text{carry}_i}{10} \right\rfloor$, **Step 3:** Output C_i .

1079 It is easy to see that this algorithm computes the digits of the sum C correctly and requires only the
 1080 individual digits at position i and the carry from the previous position. Therefore, this algorithm
 1081 satisfies the conditions of the lemma. \square

1082 **18 Additional Experiments**

1083 **18.1 Zero-Padding and Symbol Wrapping**

1084 As discussed briefly in Section 2, we found a significant benefit to using padding for multi-digit
 1085 addition. Throughout our experiments, we use the plain format without any such padding (denoted as
 1086 “vanilla” below) as the default baseline representing the conventional data format used in training.
 1087 Nonetheless, we explore modifications to this plain format to enhance performance; zero-padding,
 1088 and wrapping with a single symbol. Zero-padding ensures a fixed length for operands and the
 1089 output. In the case of 3-digit addition, this means 3-digit operands and a 4-digit output. For
 1090 example, ‘112 + 29 = 141’ becomes ‘112 + 029 = 0141’. As shown in Table 7, this modification
 1091 significantly improves model performance. Next, we wrap each sample using the ‘\$’ symbol as in
 1092 ‘\$112 + 29 = 141\$’. We found this performs on par with zero-padding.

1093 As a result, we adopt the ‘\$’ symbol for efficient data delimiter, extending its use to the reverse format.
 1094 Figure 23 shows ‘\$’-wrapping also enhances the performance of the reverse format. Despite the
 1095 plain format being improved with the ‘\$’ delimiter, it remains short of the reverse format’s accuracy
 1096 and sample efficiency. We continue to maintain the original plain format as a baseline since it not
 1097 only exemplifies conventional data but further emphasizes the need for improved data formatting to
 1098 ensure efficient training. As such, for the reverse format, we have incorporated the ‘\$’ delimiter in
 1099 our formatting modifications.

Table 7: Test accuracy of NanoGPT model on 3-digit addition trained on 10,000 samples of plain format data, comparing (i) vanilla format without modifications, (ii) Zero-padding format, and (iii) ‘\$’-wrapped format. The results show significant performance enhancement through zero-padding for fixed length and similar improvements when deploying a single-symbol wrapping.

Vanilla	Zero-pad	‘\$’-Wrapped
88.17%	97.74%	97.76%

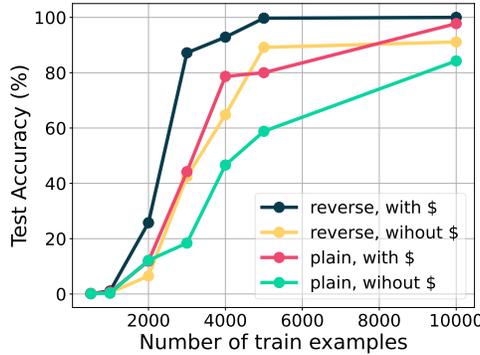


Figure 23: Performance of NanoGPT model on 3-digit addition using plain and reverse format, both with and without ‘\$’ delimiter. The addition of the ‘\$’ symbol noticeably enhances performance in both formats. Nevertheless, the plain format underperforms compared to the reverse format, particularly in terms of sample efficiency. While we maintain the original plain format as a baseline – emphasizing the necessity for improved data formatting for efficient emergence – we incorporate the ‘\$’ wrapping in our modified reverse format.

1100 18.2 Low-Rank Matrix Completion

1101 In our Low-Rank Matrix Completion experiment for the addition matrix (which is of rank-2), we
 1102 employ an iterative algorithm proposed by Király et al. (2015). This algorithm systematically searches
 1103 for a 2×2 submatrix in which three entries are known and one entry is unknown. It then fills the
 1104 unknown entry to ensure that the determinant of the 2×2 submatrix becomes zero, where the solution
 1105 is known to be optimal. We present the full pseudo-code in Algorithm 1.

1106 To assess the performance of the algorithm, we generate $n \times n$ addition matrices for various values
 1107 of n (e.g., 20, 50, 100, 500). We vary the number of revealed entries, randomly sampling a sparse
 1108 matrix where only a specified number of entries between n and $n \times n$ are known, while the remaining
 1109 entries are set to zero. We repeat this process 100 times for each number of revealed entries, tracking
 1110 the algorithm's success or failure in finding the solution. We calculate the average success rate across
 1111 the trials and present the success probabilities in Figure 3a, where we observe a sharp phase transition
 1112 when $\mathcal{O}(n)$ entries are observed, as expected.

Algorithm 1: Iterative 2×2 Matrix Completion Algorithm

Data: Data Matrix $M \in \mathbb{R}^{n \times n}$ with partially revealed entries. Assumed to be of Rank 2.

Result: $\widehat{M} \in \mathbb{R}^{n \times n}$, Success/Fail.

```

1   $n_1 \leftarrow 1$  represents number of resolved submatrices.
2   $n_2 \leftarrow 0$  represents number of unresolved submatrices.
3   $\widehat{M} \leftarrow M$ 
4  while  $n_1 \geq 1$  do
5      /* As long as we resolved at least one submatrix in the previous iteration */
6       $n_1 \leftarrow 0$ 
7       $n_2 \leftarrow 0$ 
8      for  $i = 1$  to  $n$  do
9          for  $j = 1$  to  $n$  do
10             /* do something */
11             if  $\widehat{M}_{i,j}$  is not revealed and all its neighbors are revealed then
12                  $\widehat{M}_{i,j} = \frac{\widehat{M}_{i+1,j} \times \widehat{M}_{i,j+1}}{\widehat{M}_{i+1,j+1}}$ 
13                  $n_1 \leftarrow n_1 + 1$ 
14             if  $\widehat{M}_{i+1,j}$  is not revealed and all its neighbors are revealed then
15                  $\widehat{M}_{i+1,j} = \frac{\widehat{M}_{i,j} \times \widehat{M}_{i+1,j+1}}{\widehat{M}_{i+1,j}}$ 
16                  $n_1 \leftarrow n_1 + 1$ 
17             if  $\widehat{M}_{i+1,j+1}$  is not revealed and all its neighbors are revealed then
18                  $\widehat{M}_{i+1,j+1} = \frac{\widehat{M}_{i+1,j} \times \widehat{M}_{i,j+1}}{\widehat{M}_{i,j}}$ 
19                  $n_1 \leftarrow n_1 + 1$ 
20             if  $\widehat{M}_{i,j}, \widehat{M}_{i+1,j}, \widehat{M}_{i,j+1}, \widehat{M}_{i+1,j+1}$  are all revealed then
21                 continue
22             else
23                  $n_2 \leftarrow n_2 + 1$ 
24 if  $n_2 > 0$  then
25     return  $\widehat{M}$ , Fail
26 else
27     return  $\widehat{M}$ , Success

```

1113 18.2.1 Generalizing to unseen digits

1114 Building upon the model's robustness to excluded numbers, we further investigate its ability to handle
 1115 excluded digits. Intuitively, this should be even more challenging since excluding a digit means the
 1116 model cannot learn directly how to operate in that position. Instead, it would have to generalize and

1142 task into *simpler intermediate steps*. Unlike addition, subtraction behaves differently depending on
 1143 whether the first operand (a) is greater than the second operand (b) or vice versa.

1144 The first strategy (Version 1 in Figure 25) involves performing digit-wise subtraction starting from the
 1145 least significant bit (LSB) and considering borrows when necessary. However, this strategy produces
 1146 incorrect results when the first operand is smaller than the second operand. In such cases, we subtract
 1147 the number in the most significant bit (MSB) position multiplied by 10 to the power of (number
 1148 of digits in the output - 1) from the remaining digits in the output. An example illustrating this
 1149 approach is shown in Version 1, Case 2. Alternatively, we can adopt a more familiar strategy. If
 1150 the first operand is smaller than the second, we swap the operands and compute the negation of the
 1151 subtraction of the swapped operands: $a - b = -(b - a)$ (referred to as Version 2).

1152 The results in Figure 26 indicate that Version 2, which involves comparing two operands, performs
 1153 considerably worse than Version 1. In Version 1, each intermediate step only requires the simpler
 1154 1-digit subtraction, along with addition in the final result processing step. Upon analyzing the failure
 1155 cases of Version 2, we observe that the majority of errors stem from incorrectly identifying which of
 1156 the two operands is larger, while the intermediate steps are handled correctly. This finding underscores
 1157 the significance of breaking down arithmetic operations into *simpler intermediate steps*. Unless
 1158 otherwise specified, we use Version 1 in all detailed scratchpad experiments.

Detailed scratchpad formatting for different arithmetic tasks

Examples of two variations of detailed scratchpad formatting for subtraction, considering the scenario where the first operand a is greater than the second operand b , and vice versa. In Version 1, a result processing step is included in the final stage to handle negative outputs. In Version 2, the operands are compared at the beginning, and if b is larger, their order is reversed.

Prompt (Case 1. $a - b \geq 0$):

Input:
 367-128
 Target:

Version 1.

```
...
<scratch>
[3,6,7] has 3 digits.
[1,2,8] has 3 digits.
[3,6,7] - [1,2,8] , A=[], C=0 ,
7-8-0+10=9 , A->9 , C->-1
[3,6] - [1,2] , A=[9] , C=-1 , 6-2-1=3 , A
->3 , C->0
[3] - [1] , A=[3,9] , C=0 , 3-1-0=2 , A->2
, C->0
[] - [] , A=[2,3,9]
200+39=239 , END # result processing
</scratch>
2 3 9
```

Version 2.

```
...
<scratch>
[3,6,7] has 3 digits.
[1,2,8] has 3 digits.
367>=128 # comparison of two operands
[3,6,7] - [1,2,8] , A=[], C=0 ,
7-8-0+10=9 , A->9 , C->-1
[3,6] - [1,2] , A=[9] , C=-1 , 6-2-1=3 , A
->3 , C->0
[3] - [1] , A=[3,9] , C=0 , 3-1-0=2 , A->2
, C->0
[] - [] , A=[2,3,9] , END
</scratch>
2 3 9
```

Prompt (Case 2. $a - b < 0$):

Input:
 128-367
 Target:

Version 1.

```
...
<scratch>
[1,2,8] has 3 digits.
[3,6,7] has 3 digits.
[1,2,8] - [3,6,7] , A=[], C=0 , 8-7-0=1 ,
A->1 , C->0
[1,2] - [3,6] , A=[1] , C=0 , 2-6-0+10=6 ,
A->6 , C->-1
[1] - [3] , A=[6,1] , C=-1 , 1-3-1=-3 , A
->-3 , C->-1
[] - [] , A=[-3,6,1]
-300+61=-239 , END # result processing
</scratch>
-2 3 9
```

Version 2.

```
...
<scratch>
[1,2,8] has 3 digits.
[3,6,7] has 3 digits.
128<367 : 128-367=- (367-128) # comparison
[3,6,7] - [1,2,8] , A=[], C=0 ,
7-8-0+10=9 , A->9 , C->-1
[3,6] - [1,2] , A=[9] , C=-1 , 6-2-1=3 , A
->3 , C->0
[3] - [1] , A=[3,9] , C=0 , 3-1-0=2 , A->2
, C->0
[] - [] , A=[2,3,9] , END
</scratch>
-2 3 9
```

1159

1160

Figure 25: Two versions of detailed scratchpad formatting for subtraction.

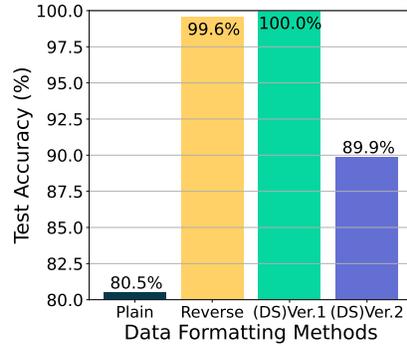


Figure 26: Comparison of performance among various data formatting approaches (plain, reverse, and two versions of detailed scratchpad (DS)) for the subtraction task. The experiments were conducted on a NanoGPT model trained on a dataset of 10,000 examples. Version 2, which incorporates operand comparison, exhibits significantly lower performance compared to Version 1. This observation highlights the substantial impact of the construction of intermediate steps on the model’s performance.

1161 **18.4 The Effect of Noisy Inputs on Accuracy**

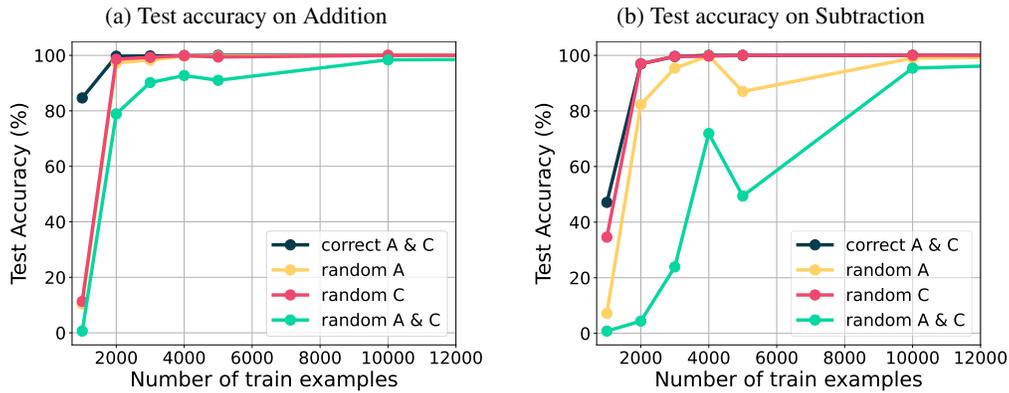


Figure 27: Comparison of training with simplified scratchpad formatting using correct A and C information with formatting using random A/C and their effect on sample efficiency and accuracy. Results show that noisy labels degrade sample efficiency, but with sufficient training data, the model eventually reaches full accuracy.

1162 **Noisy intermediate steps in the scratchpad data.** We further investigate the significance of
 1163 providing accurate intermediate steps in the scratchpad during the training process. While this was
 1164 inspired by the findings of Min et al. (2022), it is inherently different. Min et al. (2022) show that
 1165 using random labels in ICL demonstrations caused minimal degradation when compared to the gold
 1166 labels. However, those models were trained on gold labels and then evaluated on multiple downstream
 1167 tasks. In our setting, the model is trained and evaluated on a single arithmetic task. Further, the
 1168 final result(or label) is left untouched as the correct answer to the arithmetic operation. We only
 1169 replace the intermediate steps. The goal of this study is to verify whether the model actually learns to
 1170 reason using the given intermediate steps or merely uses the scratchpad to improve its expressivity.
 1171 We compare the performance of training with our simplified scratchpad formatting, which includes
 1172 accurate *A* (digit sum) and *C* (carry) information, with formatting that includes random *A*, random
 1173 *C*, or random *A* and *C* for each intermediate step, as depicted in Figure 1.

1174 The results in Figure 27, demonstrate that the inclusion of noisy labels can impede sample efficiency.
 1175 However, with enough samples, the model ultimately achieves full accuracy. This suggests that while
 1176 the model is capable of leveraging the information contained in the intermediate steps, it can also
 1177 gradually learn how to perform addition while disregarding the presence of noisy intermediate steps.

1178 **Model robustness to noise in the auto-regressive output.** In this analysis, we explore the
 1179 robustness of models trained on plain or reverse formatted data (without noise) when exposed to
 1180 noise during an auto-regressive generation process. In particular, we aim to unravel how much the
 1181 learned mapping of the *i*-th output relies on the operands and preceding tokens in the addition result,

1182 given that transformer models generate tokens sequentially in an autoregressive manner, making them
 1183 prone to error propagation.

1184 For this experiment, we focus on 3-digit addition. We train models on either plain or reverse format
 1185 data and evaluate the accuracy of next-token predictions when the output sequence contains noise.
 1186 Specifically, in the plain format setting, we expect a well-performing model to generate the correct
 1187 output tokens O_3, O_2, O_1 sequentially, where $O_3 = C_3, O_2 = C_2, O_1 = C_1$, and $C_3C_2C_1$ represents
 1188 the correct answer. We consider two types of perturbation: (i) **random** perturbation, where we
 1189 modify the first two output tokens O_3O_2 to random numbers different from C_3C_2 , and (ii) **precise**
 1190 perturbation, where we perturb only the second output token O_2 by 1. The second case is particularly
 1191 relevant since a common error case is where the model misses a digit by 1. We provide the model with
 1192 an expression of the form “ $A_3A_2A_1 + B_3B_2B_1 = O_3O_2$ ”, where O_3O_2 can be either (i) a random
 1193 incorrect number, *i.e.*, $O_3O_2 \neq C_3C_2$, or (ii) $O_2 = C_2 \pm 1 \pmod{10}$, and observe the next token
 1194 generated by the model. A corresponding process is deployed for the reverse format, introducing a
 1195 noisy sequence to models trained on reverse format data.

1196 To evaluate the performance, we define two accuracy criteria for O_1 : **exact accuracy**, reckoning
 1197 O_1 as accurate only when $O_1 = C_1$, and **relaxed accuracy**, considering O_1 correct if it deviates
 1198 from the original output C_1 by at most 1. In other words, $C_1 = O_1, C_1 = O_1 + 1 \pmod{10}$ or
 1199 $C_1 = O_1 - 1 \pmod{10}$.

Table 9: Prediction accuracy for the third digit output under different types of noise in the preceding output tokens. **Random** perturbation, applies random flips whereas **precise** perturbation shifts the preceding output tokens by 1. **Relaxed accuracy**, allows for a ± 1 deviation from the true output whereas **Exact accuracy** is strict. Reverse consistently outputs a number that is at most 1 different from the true output, even in the presence of noise. The plain format has high exact accuracy in the presence of precise perturbation, as the noise in the output token has a lower impact on predicting the next token, which is of lower significance. However, with completely random noise, the plain format shows poor performance, suggesting a strong dependence on all digits. (See Lemma 1 and 2).

Perturbation Type	Random		Precise	
	Plain	Reverse	Plain	Reverse
Exact Acc	49.88%	81.26%	99.85%	90.47%
Relaxed Acc	61.55%	100%	100%	100%

1200 The results presented in Table 9 reveal intriguing findings. We observe that the reverse format
 1201 consistently outputs a result that *deviates by no more than 1 from the true answer*, regardless
 1202 of whether the preceding outputs O_3O_2 are subjected to random or precise perturbation. This
 1203 consistency can be explained by Lemma 2, indicating that the reverse format only requires learning
 1204 a straightforward function of digit-wise addition for each corresponding position, along with the
 1205 carry-on (0 or 1). Therefore, even with noise in the preceding tokens, the model accurately performs
 1206 digit-wise addition, albeit with occasional carry-on prediction errors. With an exact accuracy of
 1207 81.26% even in the presence of random perturbation, the reverse format demonstrates the model’s
 1208 ability to rely less on the preceding output tokens, indicating a robust learned output mapping.

1209 On the contrary, models using the plain format have to decipher a more intricate function drawing
 1210 from all digits within the sequence, as described by Lemma 1. Given that in addition, carry operations
 1211 transition from right to left (*i.e.*, least to most significant digit), the introduction of precise perturbation
 1212 on preceding output tokens, which possess higher significance, has a minor impact on the output
 1213 (which has less significance). As a result, models trained using the plain format attain an exact
 1214 accuracy rate of 99.85% and a relaxed accuracy of 100% for cases involving precise perturbation.
 1215 Interestingly, under purely random perturbation, the plain format struggles, leading to a reduced
 1216 relaxed accuracy of 61.55% and exact accuracy of 49.88%. This suggests that the output mapping
 1217 learned by the plain format is not merely a function of the two operands but rather enmeshed in
 1218 complex dependencies on preceding output tokens.

1219 18.5 Analyzing the results on Sine/Sqrt

1220 Since sine and sqrt are arguably more complicated functions than the remaining arithmetic tasks, we
 1221 decided to more carefully analyze their performance. As shown in Figure 28, sin shows excellent
 1222 performance across all data formats around $\sin(x) = 0$. We conjecture that this is because $\sin(x) \approx x$

1223 for $x \approx 0$, which is easy to learn. We also note that accuracy once again improves close to ± 1
 1224 potentially for similar reasons.

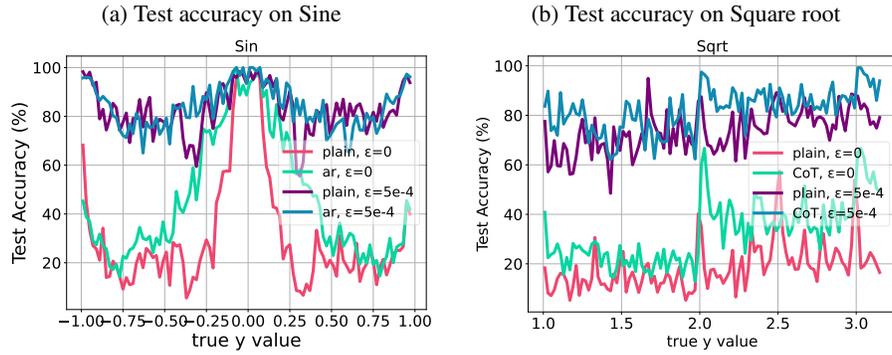


Figure 28: Error analysis of sine and square root functions, considering varying error tolerance (eps) thresholds to determine correct output. The sine function demonstrates excellent performance across all data formats, particularly around $\sin(x) = 0$, where $\sin(x) \approx x$ for $x \approx 0$. Additionally, we observe improved accuracy near ± 1 .

1225 19 Experimental Setup

1226 In this section, we summarize the datasets, models and hyperparameters used for experiments. All
1227 of our experiments on NanoGPT and GPT-2 models are run using PyTorch 2.1 and CUDA 11.7 on
1228 Nvidia 2808 TIs and NVIDIA 3090s. Detailed dependencies are provided on our github repository².

1229 19.1 Dataset

1230 In this section, we explain the details of the datasets used for our experiments. For arithmetic tasks,
1231 we construct our own datasets as described below while we use the standard shakespeare (Karpathy,
1232 2015) dataset for text.

1233 **Arithmetic Tasks** As mentioned above, for all arithmetic tasks, we prepare our own datasets.
1234 We refer to the training dataset for a binary operator $f(\cdot)$ as $\mathcal{D}_{\text{train}} = \{(x_i^1, x_i^2), y_i\}_{i=1}^N$ where
1235 $y_i = f(x_i^1, x_i^2)$. Similarly, the test dataset $\mathcal{D}_{\text{test}}$ is constructed by randomly sampling pairs of
1236 operands that do not appear in $\mathcal{D}_{\text{train}}$. During both training and inference, we then apply different
1237 formatting techniques (see Section 2), to construct the final sequence that is input to the model. We
1238 would like to repeat that both the careful choice of samples in the training dataset as well as their
1239 formatting play a crucial role in the final performance of the model.

1240 **Text** For text data, we use the Shakespeare dataset which was introduced by Karpathy (2015)
1241 originally featured in the blog post “The Unreasonable Effectiveness of Recurrent Neural Networks”.
1242 It consists of 40,000 lines of dialogue carefully curated from William Shakespeare’s plays. The dataset
1243 comprises of a total of 1,115,394 characters and 64 unique tokens(when using the character-level
1244 tokenizer that we employed in all NanoGPT experiments).

1245 19.1.1 Data Balancing

1246 As mentioned in Section 2, we carefully sample our data to ensure that they are “*balanced*” with
1247 respect to the number of carries and number of digits. As mentioned earlier, sampling the operands
1248 uniformly at random would lead to an extremely skewed dataset. To avoid this, we try to **(i) Balance**
1249 **digits** by sampling lower-digit numbers with higher weights and **(ii) Balance carry-ons** by sampling
1250 such that we have equal number of examples with 0, 1, 2 and 3 carry-on operations³.

1251 Specifically, we create a balanced dataset of 10,000 samples. This dataset includes all 100 1-digit
1252 additions and a random sampling of 900 2-digit additions (including both $(2 + 1)$ and $(1 + 2)$
1253 digit additions) and 9,000 3-digit additions. For the 3-digit addition samples, we employ *rejection*
1254 *sampling* to ensure an equal distribution of carry-ons (0, 1, 2, or 3). For the test dataset, we uniformly
1255 sample 10,000 addition examples that do not overlap with the train dataset. Results in Figure 2 and
1256 Table 10 demonstrate a clear advantage of the employed data balancing methods.

1257 For the train dataset, we follow a specific approach based on the number of examples. For sample
1258 sizes smaller than 10,000 (*e.g.*, 500, 1,000, 2,000, 3,000, 4,000, 5,000), we include all 1-digit
1259 additions and a proportionate number of 2-digit samples (*e.g.*, for a total of 5,000 samples, we
1260 include $900 \times 5,000/10,000 = 450$ two-digit additions). The remaining samples are filled with
1261 3-digit additions from the constructed train dataset of 10,000 samples. For sample sizes larger than
1262 10,000 (*e.g.*, 20,000, 40,000), we include all examples from the 10,000-sample train dataset and then
1263 add additional samples as needed. Similar to before, we perform rejection sampling to maintain an
1264 equal number of carry operations. Table 11. provides detailed information on the number of samples
1265 with 1-digit, 2-digit, and 3-digit additions, as well as the number of carry-ons.

1266 For the other arithmetic operations (subtraction, multiplication, sine, and square root), we construct
1267 the train dataset using the following approach: (i) For subtraction, we use the same pairs of operands
1268 that were used for addition. (ii) For multiplication, we include all 100 cases of a 1-digit number
1269 multiplied by a 1-digit number. Additionally, we randomly sample multiplications involving operands
1270 of up to 2 digits. (iii) For sine, we sample a random number in $[\pi/2, \pi/2]$ and truncate it to 4 decimal

²<https://anonymous.4open.science/r/nanoGPT-25D2>

³In this paper, we adopt the definition that a carry-on operation involves transferring information from one digit position to another position of higher significance. Therefore, we refer to the “borrow” operation in subtraction as a carry operation.

1271 places. (iv) For square root, we sample a random number between $[1, 10]$ and truncate it to 4 decimal
 1272 places. For the test dataset, we sample 10,000 data points (7,000 for multiplication) that do not
 1273 overlap with the train dataset.

Table 10: Performance of addition on various data sampling methods used: (i) Random - uniform sampling of operands; (ii) Balanced digits - sampling more 1 and 2-digit operations ; (iii) Balanced carry - balancing the dataset to contain an equal number of carry-on operations. Experiments on addition with zero-padding each operand and output to have 3 and 4 digits, respectively. We observe that balancing the dataset can significantly improve the performance or arithmetic operations.

Data Sampling	Overall	1-digit	2-digit	Carry-0	Carry-1	Carry-2	Carry-3
Random	97.74	98.00	96.20	95.88	98.61	98.74	94.98
Balanced Digits	98.13	100.00	99.70	98.87	98.64	98.13	95.93
Balanced Carry-Ons	98.29	100.00	99.70	98.38	97.56	99.02	98.22

Table 11: Number of examples of digit 1/2/3 and 0/1/2/3 carry-ons for NanoGPT experiments on addition for different number of samples varying from 500 to 40,000.

Total number	1-digit	2-digit	3-digit	0-carry-ons	1-carry-ons	2-carry-ons	3-carry-ons
500	100	45	355	163	141	97	99
1000	100	90	810	283	268	236	213
2000	100	180	1720	535	502	481	482
3000	100	270	2630	781	782	748	689
4000	100	360	3540	1020	1016	958	1006
5000	100	450	4450	1279	1271	1229	1221
10000	100	900	9000	2500	2500	2500	2500
20000	121	1937	17942	5000	5000	5000	5000
40000	132	3939	35929	10000	10000	10000	10000

1274 19.1.2 Data Formatting

1275 For each of the four formatting techniques, as applied to each arithmetic operation we provide the
 1276 details below. (i) **Plain** refers to the simplest formatting where we simply create a sequence as the
 1277 mathematical representation of the corresponding operation (e.g., $A_3A_2A_1 + B_3B_1B_1 = C_3C_2C_1$).
 1278 For (ii) **Reverse**, we simply reverse the digits of the output so that they appear in increasing order
 1279 from LSB to MSB (e.g., $\$A_3A_2A_1 + B_3B_1B_1 = C_1C_2C_3\$$). (iii) **Simplified Scratchpad** and (iv)
 1280 **Detailed Scratchpad** provide algorithmic reasoning steps like (Nye et al., 2021; Zhou et al., 2022b)
 1281 so as to help the model get more “information” per sample. Our intuition is that this approach nudges
 1282 the model towards actually learning the algorithm of addition or subtraction rather than merely trying
 1283 to fit the training examples. Refer to Appendix 20 for detailed examples of data formatting for each
 1284 arithmetic operation.

Different data formatting methods for addition

Four input formatting methods used for the addition task:

- (i) **Plain**: standard formatting of addition
- (ii) **Reverse**: flips the order of the output and encapsulates each data sample with the '\$' symbol at the start and end.
- (iii) **Simplified Scratchpad**: provides carry and digit-sum information for each step of addition, from the LSB to the MSB⁴.
- (iv) **Detailed Scratchpad**: provides explicit details of intermediate steps of addition.

Plain	Detailed Scratchpad
128+367=495	Input: 128+367
Reverse	Target: <scratch>
\$128+367=594\$	[1,2,8] has 3 digits. [3,6,7] has 3 digits.
Simplified Scratchpad	[1,2,8] + [3,6,7] , C=0, 8+7+0=15, A->5, C->1 [1,2] + [3, 6] , A= [5], 2+6+1=9 , A->9, C->0 [1] + [3] , A= [9,5] , C=0 , 1+3+0=4 , A->4 , C->0 [] + [] , A= [4,9,5] , C=0 , END
Input: 128+367 Target: A->5 , C->1 A->9 , C->0 A->4 , C->0. 495	</scratch> 4 9 5

1285

Figure 29: The four input formatting methods used for the addition task. We progressively increase the amount of detail with each format.

1286

1287

1288 Note that we wrap each data sample in the reverse format with the '\$' symbol at the beginning and end
 1289 as a delimiter. We originally observed improved performance in both the plain and reverse formats
 1290 when the operands and outputs were zero-padded to a fixed length (e.g., 3 and 4 digits, respectively,
 1291 for 3-digit addition). But later realized that a single symbol can effectively replace zero-padding.
 1292 While we maintain the original plain format without padding as a baseline – emphasizing the necessity
 1293 for improved data formatting for efficient emergence – we incorporate the '\$'-encapsulation in our
 1294 modified reverse format. For further details, refer to Appendix 18.1.

1295 **Addition (+)**. We focus on additions of positive numbers up to 3-digits, in which the plain
 1296 formatting would look like $A_3A_2A_1 + B_3B_1B_1 = C_3C_2C_1$. For experiments on comparing data
 1297 sampling presented in Figure 2, we pad the two operands and the output with zero, to be of length 3
 1298 and 4 respectively. For all other experiments, we **do not utilize zero-padding**. For Scratchpad-based
 1299 methods (iii, iv), we provide the digit-wise addition (denoted as A) and carry-on (denoted as C)
 1300 information for intermediate steps from the least significant bit (LSB) to the most significant bit
 1301 (MSB).

1302 **Subtraction (-)**. We consider subtraction of positive numbers up to 3 digits, written as
 1303 $A_3A_2A_1 - B_3B_2B_1 = C_3C_2C_1$ in (i) plain formatting, and $\$A_3A_2A_1 - B_3B_1B_1 = C_1C_2C_3\$$ in (ii)
 1304 reverse formatting. As with addition, scratchpad-based methods (iii, iv), present the intermediate steps
 1305 of digit-wise subtraction and handling of carry-ons. These steps proceed from the least significant
 1306 bit (LSB) to the most significant bit (MSB). If the final result after computing all the digit-wise
 1307 subtractions is negative, we subtract the number in the most significant bit (MSB) position multiplied
 1308 by 10 to the power of (number of digits in the output - 1) from the remaining digits in the output. In
 1309 Section 18.3, we present an alternative version of the detailed scratchpad formatting for subtraction.

1310 **Multiplication (x)**. We consider multiplication of positive numbers up to 2-digits. (i) Plain
 1311 formatting examples are formatted as $A_2A_1 * B_2B_1 = C_4C_3C_2C_1$, while (ii) reverse formatting is
 1312 formatted as $\$A_2A_1 * B_2B_1 = C_1C_2C_3C_4\$$. The (iv) detailed scratchpad method simplifies each

⁴We deviate from the strict definition of “most significant bit” (MSB) and “least significant bit” (LSB), typically associated with binary numbers, and reinterpret them for the purpose of this paper as the most significant “digit” and least significant “digit”, respectively.

1313 intermediate step by conducting a series of multiplications between the first operand and each digit
 1314 of the second operand, starting from the least significant bit (LSB) and moving toward the most
 1315 significant bit (MSB). For each step, we multiply the result by an exponentiation of 10 corresponding
 1316 to the relative digit position.

1317 **Sine (sin).** We consider decimal numbers within the range $[-\pi/2, \pi/2]$, truncated to 4-digit
 1318 precision. (i) Plain formatting examples are formatted as $\text{sin}(A_0.A_1A_2A_3A_4) = B_0.B_1B_2B_3B_4$.
 1319 For (iv) detailed scratchpad method, we include the Taylor series expansion steps for sine, which
 1320 is represented as $\text{sin}(x) = x - \frac{1}{3!}x^3 + \frac{1}{5!}x^5 - \frac{1}{7!}x^7 + \dots$. These intermediate steps involve
 1321 exponentiation, which may not be any easier to compute than the sine operation itself.

1322 **Square Root ($\sqrt{\cdot}$).** We consider decimal numbers within $[1, 10)$, truncated to 4-digits of precision
 1323 with the format, written as $\text{sqrt}(A_0.A_1A_2A_3A_4) = B_0.B_1B_2B_3B_4$ for (i) plain formatting. For (iv)
 1324 detailed scratchpad method, we enumerate each step of Newton’s method to compute the square root
 1325 function. The iterative formula is given by $x_n = \frac{1}{2}(x_{n-1} + \frac{x}{x_{n-1}})$, where x_0 is initialized as the
 1326 floor of the square root value of the operand x . These intermediate steps involve a division operation,
 1327 which can be as complex as the square root operation itself.

1328 For evaluation of sine and square root, we classify the result \hat{y}_i as correct if the absolute difference
 1329 between \hat{y}_i and the ground truth value y_i is less than or equal to a predefined threshold $\epsilon \geq 0$.

1330 **19.2 Model**

1331 For all experiments, we use a Decoder-only Transformer architecture. Specifically, we primarily use
 1332 the NanoGPT model, a scaled-down variant of the GPT-2 model with half the number of self-attention
 1333 layers, heads, and embedding dimension. Note that we use character-level tokenization instead of
 1334 using the OpenAI’s BPE tokenizer (Tiktoken) of vocabulary size 50257, making the vocabulary
 1335 size significantly smaller. We use a learnable absolute positional embedding initialized randomly,
 1336 following the GPT-2 model. Are results are generated using a temperature of 0.8.

1337 In the case of arithmetic tasks performed on plain and reverse formatting, we set a context length of
 1338 256 for NanoGPT experiments. The length of a single train example falls within the range of 13 to
 1339 15, approximately. However, when conducting experiments on scratchpad formatting, we increase
 1340 the context length to 1024. This adjustment allows us to accommodate more examples per batch. In
 1341 the case of simplified scratchpad, the length of each train example is approximately 64, while the
 1342 detailed scratchpad has a length of approximately 281. For GPT-2 experiments we fix the context
 1343 length to 1024 for all experiments. See Table 12 for details on model configuration.

1344 For experiments on fine-tuning a pretrained large language model, we use OpenAI’s GPT-3 model -
 1345 Ada, Curie, and Davinci.

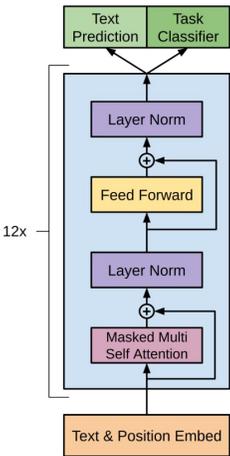


Figure 30: The GPT-2 Architecture. Image from (Radford & Narasimhan, 2018). NanoGPT model is a smaller model with half the number of self-attention layers, multi-heads, and embedding dimensions.

Table 12: NanoGPT and GPT-2 model configuration

Model	Input Formatting	Context Length	Self-Attn Layers	Num Heads	Embedding Dim
NanoGPT	Plain, Reverse	256	6	6	384
	Scratchpad	1024	6	6	384
GPT-2	Plain, Reverse	1024	12	12	768
	Scratchpad	1024	12	12	768

1346 19.3 Hyperparameter Configurations

1347 In this section, we provide a detailed overview of the hyperparameter configuration used in our
 1348 experiments in Table 13 and 14. To enhance memory efficiency and training speed, we employ flash
 1349 attention. For most experiments, we utilize the bfloat16 data type. However, when working with
 1350 Nvidia 2080 GPUs, which do not support bfloat16, we switch to float16. It is worth noting that we
 1351 did not observe significant differences in training and evaluation performance between the two data
 1352 types.

1353 For the GPT-2 experimentation, we reduced the batch size to 8 to accommodate the GPU memory
 1354 limitations. However, to mitigate the impact of the smaller batch size, we employed gradient accu-
 1355 mulation steps. This approach involves taking multiple steps between gradient updates, effectively
 1356 increasing the *effective* batch size to 64. For specific hyperparameter details, please refer to Table 14.

Table 13: Hyper Parameters used for NanoGPT experiments on arithmetic tasks

Input Format	Batch Size	Optimizer	LR	Betas	Iterations	Warmup Iter	Wt decay	Dropout
Plain, Reverse	256	AdamW	0.001	(0.9, 0.99)	5000	100	0.1	0.2
Scratchpad	16	AdamW	0.001	(0.9, 0.99)	50000	0	0.1	0.2

Table 14: Hyper Parameters used for GPT-2 experiments on arithmetic tasks

Input Format	Batch Size	Optimizer	LR	Betas	Iterations	Warmup Iter	Wt decay	Dropout
Plain, Reverse	64	AdamW	0.0005	(0.9, 0.99)	5000	100	0.1	0.2
Scratchpad	64	AdamW	0.0005	(0.9, 0.99)	20000	0	0.1	0.2

Table 15: Hyper Parameters used for tandem training experiments in Section ??.

Model	Batch Size	Optimizer	LR	Betas	Iterations	Warmup Iter	Wt decay	Dropout
NanoGPT	16	AdamW	0.001	(0.9, 0.99)	5000	0	0.1	0.2
GPT-2	40	AdamW	0.0006	(0.9, 0.95)	50000	2000	0.1	0.2

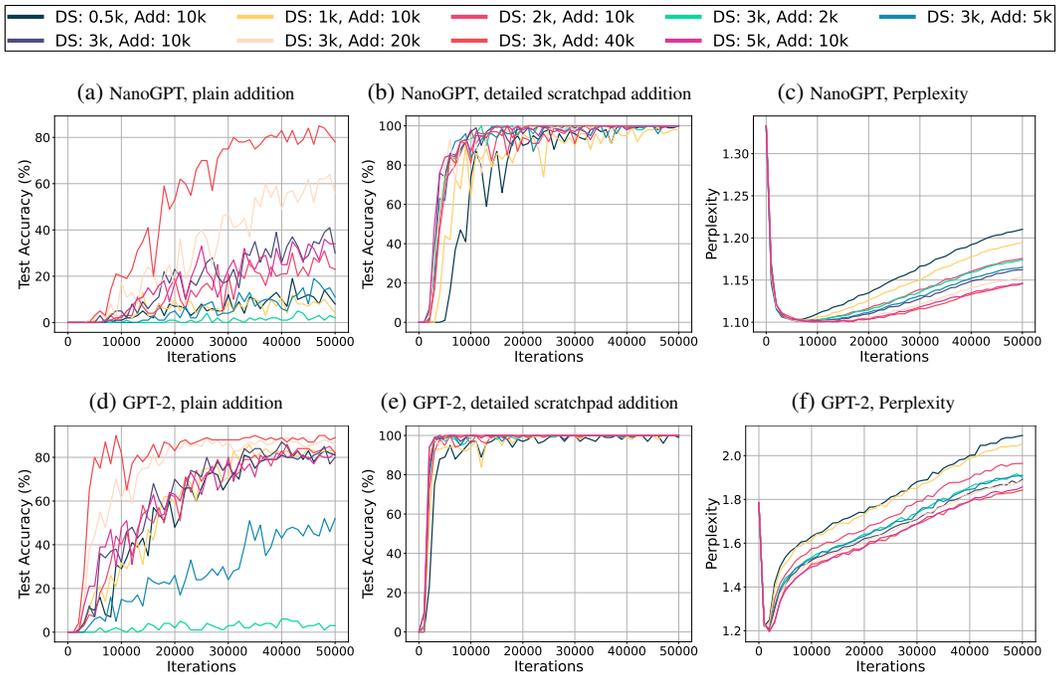


Figure 31: Training loss curves for NanoGPT and GPT-2 trained with varying numbers of plain (Add) and detailed scratchpad (DS) samples as well as the shakespeare dataset as described in Section ???. As we can see, the model continues to improve in addition accuracy as the number of iterations increases. However, the training perplexity on Shakespeare also tends to increase, which indicates some overfitting. However, we note that the model still outputs “reasonable” text when prompted with shakespeare text.

1357 **20 Prompt Examples**

1358 In this section, we provide three examples of each formatting (plain, reverse, simplified scratchpad,
 1359 detailed scratchpad) of arithmetic operations (+, −, ×, sin, √).

1360 **20.1 Addition**

Addition Examples

Plain

266+738=1004
 980+743=1723
 41+34=75

Reverse

\$913+524=1437\$
 \$226+598=824\$
 \$35+58=93\$

Simplified Scratchpad

Input:
 922+244
 Target:
 A->6 , C->0
 A->6 , C->0
 A->1 , C->1.
 1166
 Input:
 285+43
 Target:
 A->8 , C->0
 A->2 , C->1
 A->3 , C->0.
 328
 Input:
 993+849
 Target:
 A->2 , C->1
 A->4 , C->1
 A->8 , C->1.
 1842

Detailed Scratchpad

Input:
 396+262
 Target:
 <scratch>
 [3,9,6] has 3 digits.
 [2,6,2] has 3 digits.
 [3,9,6] + [2,6,2] , A=[], C=0 , 6+2+0=8 , A->8 , C->0
 [3,9] + [2,6] , A=[8] , C=0 , 9+6+0=15 , A->5 , C->1
 [3] + [2] , A=[5,8] , C=1 , 3+2+1=6 , A->6 , C->0
 [] + [] , A=[6,5,8] C=0 , END
 </scratch>
 6 5 8
 Input:
 796+890
 Target:
 <scratch>
 [7,9,6] has 3 digits.
 [8,9,0] has 3 digits.
 [7,9,6] + [8,9,0] , A=[], C=0 , 6+0+0=6 , A->6 , C->0
 [7,9] + [8,9] , A=[6] , C=0 , 9+9+0=18 , A->8 , C->1
 [7] + [8] , A=[8,6] , C=1 , 7+8+1=16 , A->6 , C->1
 [] + [] , A=[6,8,6] C=1 , END
 </scratch>
 1 6 8 6
 Input:
 788+989
 Target:
 <scratch>
 [7,8,8] has 3 digits.
 [9,8,9] has 3 digits.
 [7,8,8] + [9,8,9] , A=[], C=0 , 8+9+0=17 , A->7 , C->1
 [7,8] + [9,8] , A=[7] , C=1 , 8+8+1=17 , A->7 , C->1
 [7] + [9] , A=[7,7] , C=1 , 7+9+1=17 , A->7 , C->1
 [] + [] , A=[7,7,7] C=1 , END
 </scratch>
 1 7 7 7

Subtraction Examples	
Plain	Detailed Scratchpad
<p>266 - 738 = -472 980 - 743 = 237 41 - 34 = 7</p>	<p>Input: 396 - 262 Target: <scratch> [3,9,6] has 3 digits. [2,6,2] has 3 digits. [3,9,6] - [2,6,2], A=[], C=0, 6-2-0=4, A->4, C->0 [3,9] - [2,6], A=[4], C=0, 9-6-0=3, A->3, C->0 [3] - [2], A=[3,4], C=0, 3-2-0=1, A->1, C->0 [] - [], A=[1,3,4] 100+34=134, END </scratch></p>
Reverse	
<p>\$913 - 524 = 983\$ \$226 - 598 = 273-\$ \$35 - 58 = 32-\$</p>	<p>1 3 4 Input: 796 - 890 Target: <scratch> [7,9,6] has 3 digits. [8,9,0] has 3 digits. [7,9,6] - [8,9,0], A=[], C=0, 6-0-0=6, A->6, C->0 [7,9] - [8,9], A=[6], C=0, 9-9-0=0, A->0, C->0 [7] - [8], A=[0,6], C=0, 7-8-0=-1, A->-1, C->-1 [] - [], A=[-1,0,6] </scratch> -9 4</p>
Simplified Scratchpad	
<p>Input: 396 - 262 Target: A->4, C->0 A->3, C->0 A->1, C->0 100+34=134. 134 Input: 796 - 890 Target: A->6, C->0 A->0, C->0 A->-1, C->-1 -100+6=-94. -94 Input: 788 - 989 Target: A->9, C->-1 A->9, C->-1 A->-3, C->-1 -300+99=-201. -201</p>	<p>Input: 796 - 890 Target: <scratch> [7,9,6] has 3 digits. [8,9,0] has 3 digits. [7,9,6] - [8,9,0], A=[], C=0, 6-0-0=6, A->6, C->0 [7,9] - [8,9], A=[6], C=0, 9-9-0=0, A->0, C->0 [7] - [8], A=[0,6], C=0, 7-8-0=-1, A->-1, C->-1 [] - [], A=[-1,0,6] </scratch> -9 4 Input: 788 - 989 Target: <scratch> [7,8,8] has 3 digits. [9,8,9] has 3 digits. [7,8,8] - [9,8,9], A=[], C=0, 8-9-0+10=9, A->9, C->-1 [7,8] - [9,8], A=[9], C=-1, 8-8-1+10=9, A->9, C->-1 [7] - [9], A=[9,9], C=-1, 7-9-1=-3, A->-3, C->-1 [] - [], A=[-3,9,9] -300+99=-201, END </scratch> -2 0 1</p>

Multiplication Examples

Plain	Detailed Scratchpad
$5 * 32 = 160$ $66 * 76 = 5016$ $67 * 74 = 4958$	Input : 22*52 Target : <scratch> [2,2] has 2 digits. [5,2] has 2 digits. [2,2] * 2 , A=[4,4] , k=1 , B=[4,4] , C=0+44=44 [2,2] * 5 , A=[1,1,0] , k=10 , B=[1,1,0,0] , C=44+1100=1144 , END </scratch> 1 1 4 4 Input : 8*69 Target : <scratch> [8] has 1 digits. [6,9] has 2 digits. [8] * 9 , A=[7,2] , k=1 , B=[7,2] , C=0+72=72 [8] * 6 , A=[4,8] , k=10 , B=[4,8,0] , C=72+480=552 , END </scratch> 5 5 2 Input : 52*34 Target : <scratch> [5,2] has 2 digits. [3,4] has 2 digits. [5,2] * 4 , A=[2,0,8] , k=1 , B=[2,0,8] , C=0+208=208 [5,2] * 3 , A=[1,5,6] , k=10 , B=[1,5,6,0] , C=208+1560=1768 , END </scratch> 1 7 6 8
Reverse $\$5 * 32 = 061\$$ $\$66 * 76 = 6105\$$ $\$67 * 74 = 8594\$$	

1365

Sine Examples

Plain	Detailed Scratchpad
$\sin(1.0313) = 0.8579$ $\sin(-0.6909) = -0.6373$ $\sin(-0.5719) = -0.5413$	Input : sin(1.0313) Target : <scratch> x_0=1.0313 x_1: x_0 - 1/3! * (x^3) , x_1=0.8484 x_2: x_1 + 1/5! * (x^5) , x_2=0.8581 x_3: x_2 - 1/7! * (x^7) , x_3=0.8578 x_4: x_3 + 1/9! * (x^9) , x_4=0.8578 , END </scratch> 0.8578 Input : sin(-0.6909) Target : <scratch> x_0=-0.6909 x_1: x_0 - 1/3! * (x^3) , x_1=-0.636 x_2: x_1 + 1/5! * (x^5) , x_2=-0.6374 x_3: x_2 - 1/7! * (x^7) , x_3=-0.6374 x_4: x_3 + 1/9! * (x^9) , x_4=-0.6375 , END </scratch> -0.6375 Input : sin(-0.5719) Target : <scratch> x_0=-0.5719 x_1: x_0 - 1/3! * (x^3) , x_1=-0.5408 x_2: x_1 + 1/5! * (x^5) , x_2=-0.5414 x_3: x_2 - 1/7! * (x^7) , x_3=-0.5414 x_4: x_3 + 1/9! * (x^9) , x_4=-0.5415 , END </scratch> -0.5415

1367

Square Root Examples

Plain

```
sqrt(7.2726)=2.6967
sqrt(3.6224)=1.9032
sqrt(1.0895)=1.0437
```

Detailed Scratchpad

```
Input :
sqrt(7.1042)
Target :
<scratch>
x_0=2
x_1: 1/2*(2+7.1042/2)=2.776, x_1=2.776
x_2: 1/2*(2.776+7.1042/2.776)=2.6675, x_2=2.6675
x_3: 1/2*(2.6675+7.1042/2.6675)=2.6653, x_3=2.6653
x_4: 1/2*(2.6653+7.1042/2.6653)=2.6653, x_4=2.6653 , END
</scratch>
2.6653
Input :
sqrt(6.2668)
Target :
<scratch>
x_0=2
x_1: 1/2*(2+6.2668/2)=2.5667, x_1=2.5667
x_2: 1/2*(2.5667+6.2668/2.5667)=2.5041, x_2=2.5041
x_3: 1/2*(2.5041+6.2668/2.5041)=2.5033, x_3=2.5033
x_4: 1/2*(2.5033+6.2668/2.5033)=2.5033, x_4=2.5033 , END
</scratch>
2.5033
Input :
sqrt(8.3216)
Target :
<scratch>
x_0=2
x_1: 1/2*(2+8.3216/2)=3.0804, x_1=3.0804
x_2: 1/2*(3.0804+8.3216/3.0804)=2.8909, x_2=2.8909
x_3: 1/2*(2.8909+8.3216/2.8909)=2.8847, x_3=2.8847
x_4: 1/2*(2.8847+8.3216/2.8847)=2.8847, x_4=2.8847 , END
</scratch>
2.8847
```

1369

20.6 Noisy Simple Scratchpad

1371 We provide one example for each case of adding noise in the simplified scratchpad experiments
 1372 discussed in Section 18.4.

Noisy Simple Scratchpad Examples

We provide one example for each case of adding noise in the simplified scratchpad experiments discussed in Section 18.4. The input prompt is highlighted in light blue, while the remaining part is highlighted in light green. We construct the dataset to have either correct or random digit-sum A and carry information C. For all cases, the final answer remains accurate.

Prompt:

```
Input :
686+886
Target :
```

Correct A & C

```
A->2 , C->1
A->7 , C->1
A->5 , C->1.
1572
```

Random C

```
A->2 , C->0
A->7 , C->0
A->5 , C->1.
1572
```

Random A

```
A->0 , C->1
A->9 , C->1
A->9 , C->1.
1572
```

Random A & C

```
A->8 , C->1
A->1 , C->0
A->2 , C->1.
1572
```

1373

1374 **20.7 Example data for GPT-3 fine-tuning**

1375 We provide an example from the training dataset consisting of one prompt-completion pair used for
 1376 fine-tuning the GPT-3 model using OpenAI’s API. The prompt is highlighted in light grey, while
 1377 the completion is highlighted in light green. Note that for plain and reverse formatting, we include
 1378 spacing between digits to ensure consistent tokenization of numbers. “###” is used as the stop
 1379 sequence for generation.

1380 **20.7.1 Addition**

Addition Examples

Plain	Detailed Scratchpad
<p>6 7 7 + 8 9 8 = 1 5 7 5###</p>	<p>Input: 356+787 Target: <scratch></p>
<p style="text-align: center;">Reverse</p> <p>7 4 9 + 7 8 5 = 4 3 5 1###</p>	<p>[3,5,6] has 3 digits. [7,8,7] has 3 digits. [3,5,6] + [7,8,7] , A=[], C=0 , 6+7+0=13 , A->3 , C->1 [3,5] + [7,8] , A=[3] , C=1 , 5+8+1=14 , A->4 , C->1 [3] + [7] , A=[4,3] , C=1 , 3+7+1=11 , A->1 , C->1 [] + [] , A=[1,4,3] C=1 , END</p>
<p style="text-align: center;">Simplified Scratchpad</p> <p>Input: 32+981 Target: A->3 , C->0 A->2 , C->1 A->0 , C->1. 1013###</p>	<p></scratch> 1 1 4 3###</p>

1381

1382 **20.7.2 Subtraction**

Subtraction Examples

Plain	Detailed Scratchpad
<p>2 0 4 - 5 0 1 = - 2 9 7###</p>	<p>Input: 848-367 Target: <scratch></p>
<p style="text-align: center;">Reverse</p> <p>7 3 4 - 9 6 7 = 3 3 2 -###</p>	<p>[8,4,8] has 3 digits.[3,6,7] has 3 digits. [8,4,8] - [3,6,7] , A=[], C=0 , 8-7-0=1 , A->1 , C->0 [8,4] - [3,6] , A=[1] , C=0 , 4-6-0+10=8 , A->8 , C->-1 [8] - [3] , A=[8,1] , C=-1 , 8-3-1=4 , A->4 , C->0 [] - [] , A=[4,8,1] 400+81=481 , END</p>
<p style="text-align: center;">Simplified Scratchpad</p> <p>Input: 695-489 Target: A->6 , C->-1 A->0 , C->0 A->2 , C->0 200+6=206. 206###</p>	<p></scratch> 4 8 1###</p>

1383

1384 **20.7.3 Sine**

Sine Examples	
Plain	Detailed Scratchpad
<pre>sin(-0.8649) -0.7611###</pre>	<pre>Input: sin(-1.3516) Target: x_0=-1.3516 x_1: -1.3516 - 1/3! * (x*x*x), x_1=-0.9401 x_2: -0.9401 + 1/5! * (x*x*x*x*x), x_2=-0.9777 x_3: -0.9777 - 1/7! * (x*x*x*x*x*x*x), x_3=-0.9761 x_4: -0.9761 + 1/9! * (x*x*x*x*x*x*x*x*x), x_4=-0.9762, END </scratch> -0.9762###</pre>

1385

1386 **20.7.4 Square Root**

Square Root Examples	
Plain	Detailed Scratchpad
<pre>sqrt(1.2178) 1.1035###</pre>	<pre>Input: sqrt(5.5808) Target: <scratch> x_0=2 x_1: 1/2*(2+5.5808/2)=2.3952, x_1=2.3952 x_2: 1/2*(2.3952+5.5808/2.3952)=2.3625, x_2=2.3625 x_3: 1/2*(2.3625+5.5808/2.3625)=2.3623, x_3=2.3623 x_4: 1/2*(2.3623+5.5808/2.3623)=2.3623, x_4=2.3623, END </scratch> 2.3623###</pre>

1387

1388