

# Reverse That Number! Decoding Order Matters in Arithmetic Learning

Anonymous ACL submission

## Abstract

Recent advancements in pretraining have demonstrated that modern Large Language Models (LLMs) possess the capability to effectively learn arithmetic operations. However, despite acknowledging the significance of digit order in arithmetic computation, current methodologies predominantly rely on sequential, step-by-step approaches for teaching LLMs arithmetic, resulting in a conclusion where obtaining better performance involves fine-grained step-by-step. Diverging from this conventional path, our work introduces a novel strategy that not only reevaluates the digit order by prioritizing output from the least significant digit but also incorporates a step-by-step methodology to substantially reduce complexity. We have developed and applied this method in a comprehensive set of experiments. Compared to the previous state-of-the-art (SOTA) method, our findings reveal an overall improvement of 11.1% in accuracy while requiring only a third of the tokens typically used during training. For the purpose of facilitating replication and further research, we have made our code and dataset publicly available at <https://anonymous.4open.science/r/RAIT-9FB7/>.

## 1 Introduction

Large language models (LLMs), though proficient in a range of tasks (Ouyang et al., 2022; Achiam et al., 2023; Anil et al., 2023), encounter challenges in arithmetic operations due to their inherent design limitations, such as reliance on next-token prediction methods and limited working memory (Bubeck et al., 2023). Despite their capability to utilize external tools for circumventing direct arithmetic computations during inference (Gao et al., 2023; Imani et al., 2023; Schick et al., 2023), efficiently and effectively incorporating arithmetic proficiency within LLMs is an unresolved issue. However, previous studies have demonstrated that LLMs can learn arithmetic effectively through pretraining (Yang et al., 2023). This suggests that it might

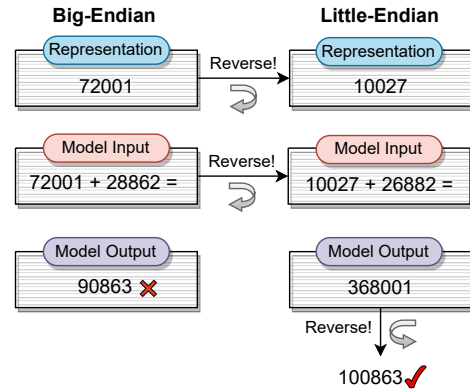


Figure 1: Reversing the numbers in training enables models to better learn to do arithmetic operations.

be feasible to efficiently teach LLMs arithmetic operations through fine-tuning alone, without the need external tool such as calculators.

The prevailing challenge in employing Large Language Models for arithmetic tasks is intricately linked to their next-token prediction mechanism. This mechanism often leads to a reversed computation order, where more significant digits are calculated before less significant ones, a flaw attributed to LLMs' inherent limitation in forward planning (Bubeck et al., 2023). This characteristic has led to the perception that arithmetic in LLMs is akin to other complex symbolic and logical tasks, necessitating a similar approach (Nye et al., 2021). Consequently, prior research has predominantly focused on the necessity of a step-by-step methodology, breaking down arithmetic into a series of sub-steps, as a critical strategy for addressing these challenges (Wei et al., 2022; Lee et al., 2023).

Such a technique achieves significant gains in performance but introduces a trade-off between efficiency and effectiveness, necessitating a balance between the number of tokens per training case and the total number of training cases. To enhance both efficiency and effectiveness without resorting to a brute-force integration of step-by-step processes,

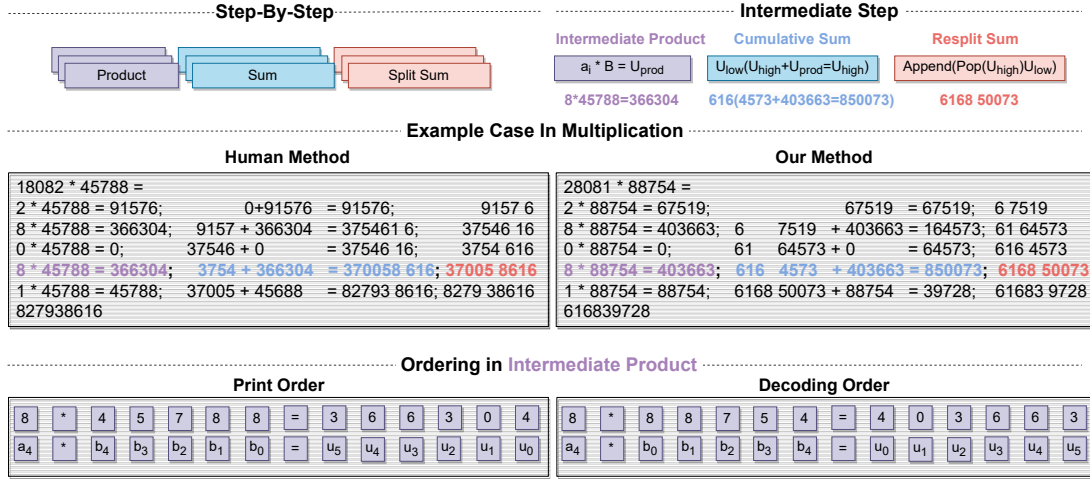


Figure 2: Example training data for Multiplication. Where the task is solved using a step-by-step process. During the  $i$ th intermediate step, the intermediate product is first computed. Then, inspired by the human process, we set the least significant digits( $U_{high}$ ) unchanged and directly added the product to the remaining digits( $U_{low}$ ) of the cumulative sum. Finally, we pop the least significant digit from the updated  $U_{high}$  and append it into  $U_{low}$  as it will not be added with non-zero digits in later steps. During decoding, we express all numbers in **Little-Endian**, where the least significant digit goes first. We convert all the numbers back to **Big-Endian** before printing.

we adopt a novel approach termed **LEFT (Little-Endian Fine-Tuning)**. Rather than incrementally integrating step-by-step mechanisms, we employ a strategy that reverses the number representation, prioritizing the computation of less significant digits. This approach utilizes the concept of **Little-Endian**, where numbers are represented with the least significant digits first, while maintaining the position of any negative signs. In contrast, the standard numeral representation is referred to as **Big-Endian**. Figure 1 demonstrates that initiating output generation with the most significant digit may result in carry-related errors. In contrast, employing a Little-Endian format, where the model produces the number 100863 as 368001, simplifies carry operations resulting in a correct solution. We present experimental results (Sec. 5) showcasing that **LEFT** not only improves accuracy by 11.1% against the current state-of-the-art (SOTA) for large digit inputs but also demonstrates efficiency by utilizing just 5.2% of the training tokens required by the previous SOTA for addition and subtraction tasks. Specifically, in multiplication, **LEFT** records a 35.7% performance gain while consuming only 56.6% of the training tokens used by prior SOTA. The key contributions of this paper include:

- We proposed a novel method, **LEFT**, leveraging Little-Endian to reduce the complexity of learning arithmetic operations.
- We conduct detailed evaluation and demon-

strate **LEFT** achieves better performance with lesser token used during training.

- Observations from our experiments indicate that, by reversing digit order, LLMs are capable of solving addition in human alike manner.

## 2 Problem Formulation

Consider the simple case where the input ( $\mathcal{I}$ ) consists of two numbers, **A** and **B**, combined with an operator  $op$ . We denote the digits of **A** as  $\mathbf{A} = \sum_{i=0}^{m-1} 10^i \mathbf{a}_i$ , where each  $\mathbf{a}_i$  is a single-digit integer ( $0 \leq \mathbf{a}_i \leq 9$ ), and  $\mathbf{a}_{m-1} \neq 0$  to ensure no leading zeros. Similarly, for **B**, we express its digits as  $\mathbf{B} = \sum_{i=0}^{n-1} 10^i \mathbf{b}_i$ , where each  $\mathbf{b}_i$  is a single-digit integer ( $0 \leq \mathbf{b}_i \leq 9$ ), and  $\mathbf{b}_{n-1} \neq 0$ .

We assume the ground truth output is a  $k$ -digit number,  $\mathbf{C} = \sum_{i=0}^{k-1} 10^i \mathbf{c}_i$  (for  $\mathbf{C} < 0$ , we use  $\mathbf{c}_{-1}$  to represent the negative sign). The trained LLM outputs an ordered sequence  $\mathcal{O} = \{\mathbf{o}_1, \mathbf{o}_2, \dots\}$ , which includes the output number  $\mathbf{C} \subseteq \mathcal{O}$ .

As step-by-step designs often incorporate intermediate results, we denote the  $i$ th intermediate result as  $\mathbf{U}^i$ . Finally, we define the remaining output as auxiliary tokens ( $\mathbf{X} = \mathcal{O} \setminus \{\mathbf{U}^i \mid \forall i\} \cup \{\mathbf{C}\}$ ).

## 3 Little-Endian Fine-Tuning

In order to effectively and efficiently teach LLMs arithmetic, we need to address three crucial questions: 1. What is the complexity in standard Big-Endian training(where no step-by-step is applied)?

2. Are there spaces for optimizing the standard method? 3. How to optimize cases when step-by-step is required? In the remaining parts of this section, we tackle such questions one by one.

### 3.1 Learning Complexity of Arithmetic

Autoregressive LLMs are interpreted as probabilistic models that predict output sequences by maximizing the likelihood of generating the correct output. In operations such as addition, this process of prediction can be formalized as follows:

$$\arg \max_{c_i} P(c_i | a_{0 \sim n-1}, b_{0 \sim m-1}, c_{i+1 \sim k})$$

Considering the specific nature of addition, where the outcome of each digit is influenced only by digits of equal or lesser significance, the process is refined to concentrate on pertinent inputs:

$$\arg \max_{c_i} P(c_i | a_{0 \sim i}, b_{0 \sim i}) \quad (1)$$

Assuming that all numbers involved possess an identical number of digits simplifies the analysis. Under this assumption, during the generation of each digit, there exist 10 potential inputs from each of the two numbers, resulting in  $10^{2i+2}$  possible input combinations. Given that the output digit can assume 10 possible values, the complexity of predicting a single digit’s value transitions from  $10^{2i+2}$  input conditions to 10 output conditions.

The overall learning complexity is quantified by summing the probabilities of accurately predicting each digit, based on the inputs up to that digit:

$$\mathcal{L}_{Big} = - \sum_{i=0}^n \log P(c_i | a_{0 \sim i}, b_{0 \sim i}) \quad (2)$$

Accordingly, the cumulative learning complexity, denoted as  $\mathcal{C}_{Big}$ , is conceptualized as the aggregate of complexities across all digits, with the input variations providing a lower bound:

$$\mathcal{C}_{Big} = \sum_{i=0}^n 10^{2i+2} \geq 10^{2n+2} \quad (3)$$

This model illustrates the exponential increase in learning complexity with the increment of digit count  $n$ , presenting a significant scalability challenge in teaching arithmetic to LLMs.

### 3.2 Optimizing Complexity via Little-Endian

In addressing the complexity of arithmetic operations, it is noted that the output token with the

greatest complexity is typically the most significant digit. Interestingly, unlike computational models, humans often do not consider all input digits simultaneously. Instead, they start from the least significant digit, using any carry-over to simplify the computation. Assuming the model can similarly infer the carry from the previous digit ( $a_{i-1}, b_{i-1}, c_{i-1}$ ), we can streamline the optimization target by focusing on this simplified context:

$$\arg \max_{c_i} P(c_i | a_i, a_{i-1}, b_i, b_{i-1}, c_{i-1})$$

Such adjustment leads to a significant reduction in input complexity, now quantified as  $10^5$ . By adopting this revised generating order, the task becomes markedly less challenging:

$$\mathcal{C}_{Little} = \sum_{i=0}^n 10^5 \leq n \cdot 10^5$$

For cases where  $n \geq 2$ , this model showcases a substantial decrease in learning complexity compared to the conventional approach ( $\mathcal{C}_{Little} \leq n \cdot 10^5 < 10^{2n+2} \leq \mathcal{C}_{Big}$ ). Such findings illuminate the potential benefits of inverting the decoding order to mitigate complexity. Motivated by this insight, we propose abandoning the classic, step-by-step design prevalent in previous methodologies in favor of revising addition and subtraction training to leverage this more efficient strategy.

**Addition.** In addressing addition within *LEFT*, the traditional approach of processing numbers from the most significant digit to the least significant is reimaged. By reversing both the input and output numbers, the calculation aligns with the Little-Endian format, where operations commence from the least significant digit and progress towards the most significant. Such conversion simplifies the decoding order, making it more intuitive and akin to human arithmetic practices. We hypothesized that the model can autonomously recompute the necessary carry for the subsequent significant digit. This method eliminates the need for a step-by-step design or the introduction of auxiliary tokens, streamlining the addition process without necessitating any extra tokens beyond the sum itself.

**Subtraction.** For subtraction, the model simplifies the process by first determining if the result will be negative, then applying the operation in Little-Endian order. This approach, which keeps the negative sign’s position unchanged (e.g., -256

214 becomes -652), enhances efficiency by eliminat- 260  
215 ing the need for intermediate results that assume 261  
216 a non-negative outcome. This streamlined method 262  
217 contrasts with traditional digit-wise subtraction, 263  
218 offering a more straightforward computation strategy. 264

### 219 3.3 Augmenting Step-by-Step 265

220 The application of Little-Endian formatting extends 266  
221 beyond the realms of addition and subtraction, of- 267  
222 fering substantial benefits in operations that inher- 268  
223 ently require a step-by-step approach due to their 269  
224 complexity. One prime example of such an opera- 270  
225 tion is multiplication, where the intricacies of the 271  
226 computation process are significantly amplified. 272

227 **Multiplication.** Traditional methods often in- 273  
228 volve breaking down the solving process into man- 274  
229 ageable chunks, typically computing the product of 275  
230 a single digit with a multi-digit number, and then 276  
231 summing these intermediate products. This con- 277  
232 ventional approach, however, often operates under 278  
233 the Big-Endian framework, starting with the most 279  
234 significant digits and potentially complicating the 280  
235 computation of intermediate products. 281

236 In contrast, the use of Little-Endian proposes a 282  
237 significant optimization. By reversing the order 283  
238 of digits—starting from the least significant—this 284  
239 method aligns with the natural flow of human com- 285  
240 putation, simplifying both the computation of inter- 286  
241 mediate product and subsequent sums. 287

## 242 4 Implementation 288

243 In this section, we delve into the detailed imple- 289  
244 mentation of *LEFT* and explore the methodologies 290  
245 applied in our experiments, along with the base- 291  
246 lines for comparison. Our discussion spans from 292  
247 the step-by-step design utilized in the experiments 293  
248 (Sec. 4.1) to dataset generation (Sec. 4.2) and other 294  
249 settings for the experiments (Sec. 4.3). 295

### 250 4.1 Step-By-Step Design 296

251 **Addition/Subtraction.** While our hypothesis 297  
252 posits that the step-by-step process might not be es- 298  
253 sential for efficiently learning addition and subtrac- 299  
254 tion, we incorporate it as a comparative measure 300  
255 to validate our assumption. We adopt the step-by- 301  
256 step design from the chain-of-thought methodol- 302  
257 ogy (Wei et al., 2022), as reproduced in previous 303  
258 studies (Zhou et al., 2022), for *LEFT*'s addition and 304  
259 subtraction tasks when necessary for evaluation. 305

**Addition/Subtraction.** Contrary to our initial hy- 260  
511 pothesis that a step-by-step process may not be 261  
512 crucial for efficiently mastering addition and sub- 262  
513 traction, we included it for comparative analysis 263  
514 to test our theory. Thus, we utilized the *Chain-Of-* 264  
515 *Thought* approach (Wei et al., 2022), as previously 265  
516 replicated (Zhou et al., 2022), in evaluating *LEFT* 266  
517 joined with step-by-step on addition/subtraction. 267

**Multiplication.** We previously outlined the key 268  
518 features of the step-by-step approach for multipli- 269  
519 cation within *LEFT*, yet a direct implementation 270  
520 was not provided. As shown in Figure 2, with the 271  
521 reversal of all numbers, the task is divided into nu- 272  
522 merous substeps. Each substep iterates over the 273  
523 digits of the first input number,  $a_i \in \mathbf{A}$ , starting 274  
524 from the least significant digit. In each iteration, 275  
525 the process begins by multiplying the current digit 276  
526 with the second input number to generate an in- 277  
527 termediate product. This intermediate product is 278  
528 then added to the cumulative sum of products from 279  
529 previous iterations. Since the lower  $i$  digits of the 280  
530 product are always zero, these are not explicitly 281  
531 represented; instead, the product is directly added 282  
532 to the higher section of the cumulative sum. The 283  
533 higher section is defined as the part of the cumula- 284  
534 tive sum obtained in the last step of the previous 285  
535 iteration, which considers the lower  $i$ -digits as a 286  
536 fixed result and defines the remaining digits as the 287  
537 higher section of the cumulative sum. 288

289 This refined step-by-step design for multiplica- 290  
538 tion highlights the efficiency and adaptability of 291  
539 the Little-Endian approach in managing complex 292  
540 arithmetic operations. By streamlining the inte- 293  
541 gration of intermediate products into a simplified 294  
542 cumulative sum, this method not only improves 295  
543 the performance and clarity of the model but also 296  
544 showcases the extensive utility of Little-Endian for- 297  
545 mating in enhancing computational processes. 298

### 299 4.2 Dataset 300

301 The inherent characteristics of arithmetic calcula- 302  
546 tions, which do not necessitate human-generated 303  
547 labels, enable the automated generation of training 304  
548 and testing sets in our study. Our primary objective 305  
549 is to create a dataset that is fair, isolated, and bal- 306  
550 anced, facilitating a comprehensive evaluation of 307  
551 the *LEFT*'s effectiveness and efficiency. 308

**Fairness.** Given that different methods may op- 309  
552 erate on varied data inputs, we aim to minimize 310  
553 the variance in performance attributable to differ- 311  
554 ent inputs as much as possible. To achieve this, 312

310	we initiate the process by generating a set of <i>meta</i>	359
311	data during the data generation phase. Each piece	360
312	of meta data is conceptualized as a triplet in the	361
313	form $(\mathbf{A}, op, \mathbf{B})$ . This triplet serves as a unified	362
314	seed for generating training and testing data for	
315	each method, ensuring that the same set of input	
316	is utilized across methods. Then, each triplet is	
317	expanded and formatted to suit the specific require-	
318	ments of each method’s data format.	
319	<b>Isolation.</b> Recognizing the critical importance	
320	of preventing data leakage, we take meticulous	
321	steps to ensure the uniqueness of input number sets,	
322	denoted by $\{\mathbf{A}, \mathbf{B}\}$ . This strategy guarantees that	
323	the test set contains no identical input number pairs	
324	as found in the training set, thereby also ensuring	
325	the uniqueness of each training and testing set.	
326	<b>Digit Distribution Balancing.</b> Echoing previous	
327	methods that have highlighted the importance of	
328	balanced data distribution (Lee et al., 2023), we en-	
329	sure that both the training and test sets are balanced	
330	such that the maximum quantity of any single num-	
331	ber in each data slice falls within the digit range of	
332	$[5, 12]$ . Specifically, we generate in total of $15K$	
333	training data and $3K$ test data, with $5K$ points	
334	for each operation, accompanied by $1K$ test data	
335	points for each operation, to maintain this balance.	
336		
337	<b>4.3 Experiment Setup</b>	
338	<b>Baseline.</b> We first include <i>End-To-End</i> training	
339	used in during pretraining methods (Yang et al.,	
340	2023) as a ground to compare performance in pre-	
341	vious methods. We then include <i>Scratchpad</i> (Nye	
342	et al., 2021), one of the early founders in using	
343	step-by-step approaches to break down arithmetic	
344	into multiple steps. We also include <i>Chain-Of-</i>	
345	<i>Thought</i> (Wei et al., 2022) which provided a gen-	
346	eral approach of breaking step-by-step to a wide	
347	range of complex tasks. In addition, we include the	
348	<i>Detailed-Scratchpad</i> method introduced in (Zhou	
349	et al., 2022). (Zhou et al., 2022) also introduces	
350	<i>Algorithmic-Prompting</i> technique but as it requires	
351	too many auxiliary tokens making it hard to fit 12-	
352	digit training into the context length. As a result,	
353	we exclude it during our evaluation.	
354		
355	<b>Metric.</b> As arithmetic reasoning is strongly af-	
356	ected by error propagation, solutions with interme-	
357	di-ate errors are almost impossible to provide the	
358	correct solution. As a result, we directly use the	
	accuracy (ACC) of the predicted output to evalu-	
	ate the effectiveness of the methods. As the dis-	
	ussion for efficiency is aimed at training better-	
	performed models using fewer resources, we record	
	the amount of tokens used for training and observe	
	the change in accuracy as more tokens are used.	
	<b>Backbone Model.</b> The base checkpoint for our	
	experimental framework is Llama2-13B (Touvron	
	et al., 2023), chosen for its status as a well-regarded	
	and openly accessible LLM. To address the need	
	for processing longer sequences, the model’s con-	
	text length has been extended to 4, 096 tokens.	
	<b>5 Experiments</b>	
	We now turn to a systematic evaluation of the pro-	
	posed method. Specifically, we design and conduct	
	a series of comprehensive analysis which seeks to	
	answer the following research questions:	
	<b>Q1</b> <i>Is LEFT effective and efficient?(Sec. 5.1)</i>	
	<b>Q2</b> <i>What grants LEFT the ability to effectively</i>	
	<i>tackles the provided task?(Sec. 5.2)</i>	
	<b>Q3</b> <i>What can be further done on LEFT?(Sec. 5.3)</i>	
	<b>5.1 Direct Evaluation Over Performance</b>	
	We began our analysis with the overall perfor-	
	mance of <i>LEFT</i> against previous methods for	
	jointly trained and evaluated addition, subtraction,	
	and multiplication performance. We then conduct	
	operation-by-operation analysis to observe the re-	
	sults of training when jointly training is opt-out.	
	<b>Observation 1: LEFT Learns Faster Than Base-</b>	
	<b>lines.</b> Table 1 shows the resulting performance of	
	each method after training. We order the baselines	
	according to token used during training. <i>LEFT</i>	
	used the least amount of training token among all	
	the step-by-step methods, yet achieving 11.1% per-	
	formance improvement over previous SOTA.	
	Specifically, <i>LEFT</i> ’s accuracy on addition and	
	subtraction is slightly below <i>Scratchpad-Detailed</i> .	
	However, <i>LEFT</i> only used $160K$ and $161K$ to-	
	kens for learning addition and subtraction. But	
	<i>Scratchpad-Detailed</i> used $2,936K$ and $3,254K$	
	for training. This means <i>LEFT</i> uses only $1/20$ of	
	training data yet still achieves similar performance.	
	<i>LEFT</i> also achieved 35.7% accuracy improvement	
	over previous SOTA on multiplication, further high-	
	lighting <i>LEFT</i> ’s effectiveness and efficiency.	
	<b>Observation 2: Using Little-Endian Alone Ob-</b>	
	<b>tains Better Efficiency On Addition/Subtraction.</b>	
	During method design(Sec. 3.2), we proposed that	
	Little-Endian is a better substitute than existing	

Method	Endian	StepByStep	+	-	×	Overall	Token Usage
<i>End-To-End</i>	Big	No	63.3	32.3	00.0	31.9	494,815
<i>Chain-of-Thought</i>	Big	Yes	88.0	83.5	08.2	59.9	4,938,148
<i>Scratchpad</i>	Big	Yes	94.8	73.1	00.0	56.0	5,747,670
<i>Scratchpad-Detailed</i>	Big	Yes	<b>99.8</b>	<b>97.3</b>	<u>52.8</u>	<u>83.3</u>	10,995,191
<i>LEFT (Our)</i>	Little	Mix	<u>98.8</u>	<u>95.9</u>	<b>88.5</b>	<b>94.4</b>	3,040,616

Table 1: Performance comparison between methods, trained with 5K data for each operation with randomly generated data. The maximum digits of input numbers for each data are equally distributed in the range of [5, 12] for each operation. The test set is generated in a similar manner but with only 1K data per operation. *LEFT* uses Little-Endian to represent all numbers and excludes the step-by-step process for addition and subtraction.

406 methods, which leverage step-by-step to reduce  
407 the complexity required for arithmetic. However,  
408 we have not yet examined such a statement. This  
409 raised two major questions: (1) Would it be better  
410 to contain step-by-step? (2) How does step-by-step  
411 itself perform? As a result, we apply step-by-step  
412 for closer observation. We scale down the training  
413 data to half and a quarter of training cases than the  
414 joint evaluation and observe the change in perform-  
415 ance. To omit influences caused by joint training,  
416 we train addition and subtraction separately.

417 As shown in Figure 3, we observe that the use  
418 of Little-Endian outperforms other settings in both  
419 operations, despite the use of fewer tokens when  
420 compared to the step-by-step settings.

421 Moreover, we observe that the conventional  
422 *Chain-Of-Thought* approach, which does not incor-  
423 porate Little-Endian formatting, also significantly  
424 lags behind the *LEFT* configuration. This outcome  
425 suggests that employing a step-by-step methodol-  
426 ogy does not invariably enhance performance. Par-  
427 ticularly in addition, both the presence and absence  
428 of Little-Endian in the settings lead to inferior re-  
429 sults compared to employing Little-Endian without  
430 a step-by-step approach. This implies that reversing  
431 the endian inherently captures critical information,  
432 which the step-by-step process aimed to convey in  
433 digit generation. Consequently, not only does the  
434 step-by-step application decrease efficiency, but it  
435 also deteriorates model performance by introduc-  
436 ing additional chance of error propagation.

437 On the other hand, by taking a closer observation  
438 of subtraction, we see whether the use of step-by-  
439 step is integrated or not, the integration of Little-  
440 Endian brings much better performance. However,  
441 the learning curve of Little-Endian without step-by-  
442 step is smoother than in addition. We believe this  
443 could be related to the pretraining setting, where  
444 the model is trained with **Big-Endian**. On addition,  
445 when the carry is not occurring, knowing what en-

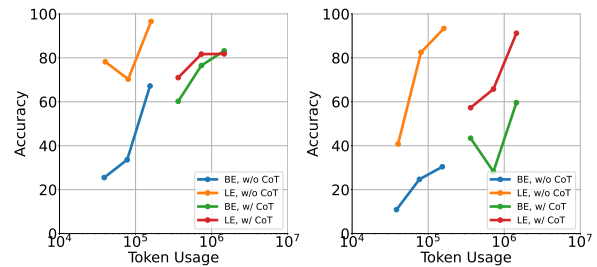


Figure 3: Performance when integrating step-by-step. BE stands for Big-Endian and LE stands for Little-Endian. The graph on the left shows the results after training on addition. The the right figure shows results for trained and evaluated on subtraction.

446 dian is involved doesn't have a strong effect on the  
447 result, the model could falsely interpret the task as  
448 aligning the numbers with the leftmost digit and  
449 still achieve some level of performance. However,  
450 on subtraction, the endian greatly affects the result,  
451 as whether the result is negative is affected by the  
452 most significant digit, which is strongly related to  
453 the endian. Such difference resulted in poor per-  
454 formance in the beginning, as the model will have  
455 a great chance of failing unless it actually under-  
456 stands the task. But it also brings faster learning as  
457 the chance for the model to falsely understand the  
458 task reduces. We believe such case highlights that  
459 the arithmetic ability of a fine-tuned model could  
460 be further improved with a backbone model that is  
461 pretrained with Little-Endian representation.

462 **Observation 3: Little-Endian And Step-by-Step**  
463 **Are Both Crucial For Multiplication.** We now  
464 conduct a detailed examination for multiplication.  
465 We re-evaluate our backbone model to examine our  
466 designs on multiplication. For better comparison,  
467 we include two additional settings other than the  
468 standard *End-To-End*. We first include a similar  
469 design as we proposed for solving addition and sub-  
470 traction, where the model directly outputs the result

Method	# of Epochs			Token Usage
	1	2	3	
<i>End-To-End</i>	-	-	-	186K
<i>Detailed-Scratchpad</i>	24.9	32.6	39.3	4,805K
<i>LEFT</i>	61.1	89.1	91.6	2,719K
<i>w/o Step-by-Step</i>	-	-	-	186K
<i>w/ Big-Endian</i>	24.2	42.8	52.7	2,719K

Table 2: Multiplication scores by different epochs and token usage. We observe settings without step-by-step solution failed to learn the task.

but the input and output are both in Little-Endian. We then include *LEFT*'s step-by-step design but convert the numbers into Big-Endian. We also measure the different performances after different epochs of training to observe the convergence for the same amount of training cases.

The results are shown in Table 2. We first observe that when the use of step-by-step is removed, it becomes impossible to learn multiplication. This demonstrates the need for step-by-step to break down the complexity in solving multiplication is still needed when only 5K of training data is available. We also observe that when Little-Endian is removed, the performance further improves over the step-by-step setting. The model also converges much faster, as the performance after 2 epochs of training is already close to the performance of the last epoch, an accuracy of 91.6%. We are amazed that *LEFT* achieves better performance when the model is trained only on multiplication, suggesting the potential for further optimization.

We also observe the number of tokens used during *LEFT*'s training in multiplication is approximately half of the tokens used by *Scratchpad-Detailed*. In addition and subtraction training, tokens are better off with a factor of 20. This shows that *LEFT* with better performance achieves even greater improvement in token efficiency.

## 5.2 Case Studies

We now conduct a detailed study of the results obtained in the previous section, seeking to discover findings that can help future studies.

**Finding 1: Little-Endian Reduces Step-By-Step Errors.** In this section, we conduct an error analysis for the errors in our main experiment in order to find an explanation of the performance gain caused by changing the endian. To do so, we first selected the place where the first error occurred as

an indication of the error of each falsely inferred test case. This is because error propagation is critical in arithmetic. We then focused on two crucial parts during each inference step, calculating the intermediate 1-by- $n$  product and the cumulative sum. As a result, we find that among the 417 errors that occurred during intermediate calculations in *Scratchpad-Detailed*: 1. 140 errors occurred during calculating the intermediate product; 2. 236 errors occurred during accumulating sum. Both operations had much better performance in *LEFT*, where only 77 errors were observed during computing the intermediate product and only 22 errors were observed when updating the cumulative sum. The error occurrence is decreased by a factor of 10 for summation and by a factor of 2 for the intermediate product. We believe this is because the carry is easier than to compute when the less significant digits are already shown, which possibly could reduce the complexity in computing the result for the current digit. The error for the intermediate sum is reduced by a greater factor as the addition training is transferable when accumulating sum on *LEFT*, whereas in *Scratchpad-Detailed*, the addition task stands more on its own. Despite slightly better performing while evaluated on addition, it cannot transfer its ability to other tasks like multiplication.

**Finding 2: LEFT Conducts Addition Just Like Humans** We now take a closer observation of how *LEFT* conducts addition. By logging the attention (Vaswani et al., 2017) scores in the model, we observe a correlation between the output digit and related digits from the input numbers, as shown in Figure 4. We observe that the input digits are recognized when computing the corresponding output during generation in some attention heads. We also observed that, in the 22th layer, shown traits suggest the fine-tuned LLM has learned to re-compute the carry from the previous digits. Addressing our hypothesized during the method design, this proves the assumption that the model can recover the carry when it's used (Sec. 3.2). This is an interesting indication because it suggests Little-Endian might be conducting training in a manner similar to how humans conduct addition without a draft paper.

## 5.3 Additional Error Analysis

Finally, we look at the errors occurred in *LEFT*'s joint experiment in the perspective of different maximum amount of input digits. As shown in Table 3, *LEFT* is able to perform well in lower digits, but

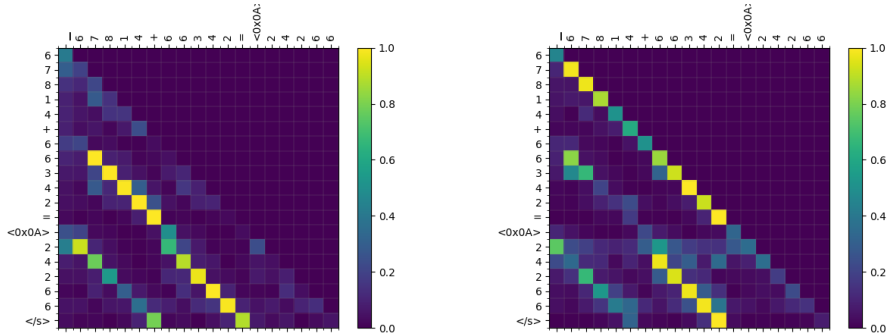


Figure 4: Visualization of attention weights during inference, with rows representing output tokens and columns indicating input tokens involved in generation. Attention weights are square-root transformed for enhanced visibility of correlations. The attention on the left(layer 14) reveals output digits are correlate with their inputs, while attention(right) from layer 22 suggests carry information reconstruction.

Max Digit	5	6	7	8	9	10	11	12
+	100.0	98.4	100.0	99.2	97.6	97.6	98.4	99.2
-	92.0	96.8	93.6	96.8	100.0	100.0	93.6	94.4
×	93.6	96.0	86.4	96.0	88.0	86.4	84.8	76.8

Table 3: Accuracy trends with increasing max input digits. We observe a steeper decline in multiplication’s performance compared to other operations.

when it is challenged towards higher digits of inputs, it loses part of its performance. Such a drop in performance is mostly significant when it comes to higher-digit multiplications, the digits being operated become much more complicated comparing to addition and subtraction. This stated that, despite well in performance, *LEFT* still faces challenges when inputted with larger digits, highlighting the need for future studies to not only focus on effectiveness and efficiency but also continue to narrow the gap for the LLMs’ inability to scale towards larger inputs and the amazing capability in humans.

## 6 Related Works

Previous methods that seek to teach LLMs to learn arithmetic mainly focus on the use of step-by-step processes. *Scratchpad* (Nye et al., 2021) was one of the early founders that recognized the use of step-by-step arithmetic solving. Zhou et al. focused on in-context learning and showed that a detailed version of *Scratchpad* could significantly improve the accuracy. Qian et al. recognized the challenger where LLM performance drops as repeated symbols increase. Goat (Liu and Low, 2023) classified tasks discussed the learnability of different operations and conducted supervised fine-tuning. Lee

and Kim proposed the Recursion of Thought to divide the solving process into short contexts.

On the other hand, some works also focus on analyzing arithmetic learning. Yuan et al. proposed MATH 401 to evaluate LLM’s arithmetic ability. Jelassi et al. discussed the length generalization ability in arithmetic. Muffo et al. evaluated the ability of Transformer to perform arithmetic operations following a pipeline that decomposes numbers in decimal before performing computations and demonstrated that this method was 60% more accurate than GPT-3 on 5-digit addition and subtraction tasks, but was inferior to GPT-3 on 2-digit multiplication tasks. Lee et al. conducted a compressive analysis on training strategies and discussed that reversing the output of addition can speed up the learning process.

## 7 Conclusion

In this study, we introduced a novel approach for teaching arithmetic to LLMs by reversing the number order to emphasize the least significant digit. This strategy, which aligns with human arithmetic practices, significantly reduces computational complexity and training data requirements, demonstrating an 11.1% increase in overall accuracy over previous SOTA and showcasing efficiency in token usage during training. The success of our method suggests the potential for broader applications in mathematical problem-solving and in environments with limited resources. We hope this study of ours paves the way for future investigations into optimizing LLM training techniques for numerical reasoning and arithmetic precision.



## 617 Limitations

618 Our study introduces a novel approach to arith-  
619 metic learning in LLMs but is not without limita-  
620 tions. Firstly, our focus on basic arithmetic opera-  
621 tions such as addition, subtraction, and multiplica-  
622 tion leaves unexplored territories in more complex  
623 arithmetic and mathematical problem-solving areas.  
624 Secondly, the generalizability of our method to do-  
625 mains beyond arithmetic is yet to be determined. A  
626 critical consideration is the reliance on LLMs pre-  
627 trained with standard numeral expressions; our ex-  
628 periments did not explore the potential benefits of  
629 pretraining models directly with reversed numeral  
630 expressions. Addressing these limitations could  
631 further enhance the applicability and efficiency of  
632 LLMs in numerical reasoning and arithmetic pre-  
633 cision, suggesting a promising direction for future  
634 research to broaden the scope of operations cover-  
635 ed and to investigate the impact of pretraining  
636 strategies.

## 637 Ethics Statement

638 Our research contributes to the field of artificial  
639 intelligence by proposing an innovative approach  
640 to improve the efficiency and accuracy of LLMs  
641 in performing arithmetic operations. This advance-  
642 ment has the potential to positively impact areas  
643 where numerical understanding is crucial, includ-  
644 ing but not limited to, educational technologies,  
645 data analysis, and automated reasoning systems.  
646 By improving the capability of LLMs to process  
647 and understand arithmetic, our work aims to sup-  
648 port further developments in technology that can  
649 assist in educational settings, enhance scientific re-  
650 search, and provide more reliable computational  
651 tools for industries relying on accurate numerical  
652 data processing.

653 We are mindful of the importance of conducting  
654 our research with a commitment to ethical princi-  
655 ples, ensuring that our methodologies and results  
656 are transparent, reproducible, and contribute con-  
657 structively to the academic community and society  
658 at large. While our work primarily focuses on the  
659 technical aspects of improving LLMs’ arithmetic  
660 abilities, we recognize the broader implications of  
661 AI and machine learning advancements. Therefore,  
662 we encourage the responsible use and continuous  
663 ethical evaluation of AI technologies, emphasizing  
664 the importance of using such advancements to  
665 foster positive societal outcomes.

## References

- 666 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
667 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
668 Diogo Almeida, Janko Altschmidt, Sam Altman,  
669 Shyamal Anadkat, et al. 2023. [GPT-4 technical re-  
670 port](#). *CoRR*, abs/2303.08774. 671
- 672 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-  
673 Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan  
674 Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Mil-  
675 lican, David Silver, Slav Petrov, Melvin Johnson,  
676 Ioannis Antonoglou, Julian Schrittwieser, Amelia  
677 Glaese, Jilin Chen, Emily Pitler, Timothy P. Lilli-  
678 crap, Angeliki Lazaridou, Orhan Firat, James Molloy,  
679 Michael Isard, Paul Ronald Barham, Tom Henni-  
680 gan, Benjamin Lee, Fabio Viola, Malcolm Reynolds,  
681 Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens  
682 Meyer, Eliza Rutherford, Erica Moreira, Kareem  
683 Ayoub, Megha Goel, George Tucker, Enrique Pi-  
684 queras, Maxim Krikun, Iain Barr, Nikolay Savinov,  
685 Ivo Danihelka, Becca Roelofs, Anaïs White, Anders  
686 Andreassen, Tamara von Glehn, Lakshman Yagati,  
687 Mehran Kazemi, Lucas Gonzalez, Misha Khalman,  
688 Jakub Sygnowski, and et al. 2023. [Gemini: A fam-  
689 ily of highly capable multimodal models](#). *CoRR*,  
690 abs/2312.11805. 691
- 692 Sébastien Bubeck, Varun Chandrasekaran, Ronen El-  
693 dan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Pe-  
694 ter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg,  
695 Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro,  
696 and Yi Zhang. 2023. [Sparks of Artificial General  
697 Intelligence: Early experiments with GPT-4](#). 698
- 699 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon,  
700 Pengfei Liu, Yiming Yang, Jamie Callan, and Gra-  
701 ham Neubig. 2023. [PAL: program-aided language  
702 models](#). In *International Conference on Machine  
703 Learning, ICML 2023, 23-29 July 2023, Honolulu,  
704 Hawaii, USA*, pages 10764–10799. 705
- 706 Shima Imani, Liang Du, and Harsh Shrivastava. 2023.  
707 [MathPrompter: Mathematical reasoning using large  
708 language models](#). In *Proceedings of the 61st An-  
709 nual Meeting of the Association for Computational  
710 Linguistics (Volume 5: Industry Track)*, pages 37–42. 711
- 712 Samy Jelassi, Stéphane d’Ascoli, Carles Domingo-  
713 Enrich, Yuhuai Wu, Yuanzhi Li, and François Char-  
714 ton. 2023. [Length generalization in arithmetic trans-  
715 formers](#). *CoRR*, abs/2306.15400. 716
- 717 Nayoung Lee, Kartik Sreenivasan, Jason D. Lee,  
718 Kangwook Lee, and Dimitris Papailiopoulos. 2023.  
719 [Teaching arithmetic to small transformers](#). *CoRR*,  
720 abs/2307.03381. 721
- 722 Soochan Lee and Gunhee Kim. 2023. [Recursion of  
723 thought: A divide-and-conquer approach to multi-  
724 context reasoning with language models](#). In *Find-  
725 ings of the Association for Computational Linguis-  
726 tics: ACL 2023*, pages 623–658. 727
- 728 Tiedong Liu and Bryan Kian Hsiang Low. 2023. [Goat:  
729 Fine-tuned llama outperforms gpt-4 on arithmetic  
730 tasks](#). 731

