

# How to Leverage Digit Embeddings to Represent Numbers?

Jasivan Alex Sivakumar and Nafise Sadat Moosavi

School of Computer Science

The University of Sheffield

United Kingdom

{jasivakumar1|n.s.moosavi}@sheffield.ac.uk

## Abstract

Within numerical reasoning, understanding numbers themselves is still a challenge for existing language models. Simple generalisations, such as solving  $100+200$  instead of  $1+2$ , can substantially affect model performance (Sivakumar and Moosavi, 2023). Among various techniques, character-level embeddings of numbers have emerged as a promising approach to improve number representation. However, this method has limitations as it leaves the task of aggregating digit representations to the model, which lacks direct supervision for this process. In this paper, we explore the use of mathematical priors to compute aggregated digit embeddings and explicitly incorporate these aggregates into transformer models. This can be achieved either by adding a special token to the input embeddings or by introducing an additional loss function to enhance correct predictions. We evaluate the effectiveness of incorporating this explicit aggregation, analysing its strengths and shortcomings, and discuss future directions to better benefit from this approach. Our methods, while simple, are compatible with any pretrained model, easy to implement, and have been made publicly available.<sup>1</sup>

## 1 Introduction

Numbers play an integral role in language (Thawani et al., 2021), and they are crucial across various domains such as finance (Chen et al., 2018), medicine (Jullien et al., 2023) or even sarcasm (Dubey et al., 2019). Despite large language models improving their capacity in many tasks, numerical reasoning still poses a challenge (Hong et al., 2024). Recent advancements in enhancing numerical reasoning within language models have predominantly stemmed from using more extensive or higher-quality training datasets (Li et al., 2022a; Yu et al., 2024), scaling up models (Lewkowycz et al.,

2022; Kojima et al., 2022), or integrating prompt-based strategies such as chain-of-thought reasoning (Wei et al., 2022b; Yue et al., 2024). The effectiveness of such methods is significantly amplified when applied in conjunction with larger model architectures. With smaller models, the improvement shown is often minimal, for example, Wei et al. (2022b) use of chain-of-thought on a 20B parameter model only showed a 2.5% improvement on the MAWPS (Koncel-Kedziorski et al., 2016) dataset whereas it jumps to 14.7% with a 137B parameter model. In addition, many of these solutions are computationally expensive or inaccessible; we seek a low-cost approach that may have minimal impact on small-scale models but greater effects on larger models.

A key challenge for number understanding is that widely used tokenisation methods, like Byte-Pair Encoding (BPE) (Sennrich et al., 2016), work well for common words but not for numbers. Specifically, rarer numbers might be broken down into random, meaningless pieces. In light of this, digit tokenisation (Spithourakis and Riedel, 2018) stands out for its simplicity and efficacy at representing numbers. This technique involves breaking down numbers into their individual digits, reducing vocabulary size and ensuring all decimal numbers can be accurately represented enhancing numerical reasoning abilities across various model architectures, tasks, and datasets (Geva et al., 2020; Petrak et al., 2023; Sivakumar and Moosavi, 2023). However, the aggregation of digit embeddings into a complete number representation is implicitly handled by the model, which raises the question: **Can explicit aggregation using mathematical priors improve numerical understanding?**

In this paper, we investigate this question by integrating a mathematically grounded aggregation of digit embeddings explicitly, rather than relying solely on the model’s inherent capabilities. Our findings show that this aggregated digit embedding

<sup>1</sup><https://github.com/jasivan/Number-Embeddings>

enhances performance on small-scale models by up to 6.17% compared to our baseline without it, potentially leading to even greater improvements in larger models. However, the effectiveness of our integration strategy depends on the size, the architecture of the model used, and how these priors are integrated in the model. Our main contributions are as follows:

- We propose a novel approach to number embedding that requires no changes to the model’s architecture or additional pretraining, by showing that an aggregation is effective if it meets two criteria: (1) it distinguishes between distinct numbers, ensuring unique representations for each value, and (2) the aggregated embedding reflects natural numerical proximity.
- We explore two approaches for integrating our aggregation: adding a special token before the representation of individual digits to enhance input number representations, and incorporating an additional loss function to improve the representation of outputted digits.

## 2 Related Work

Numerical reasoning is the ability to interact with numbers using fundamental mathematical properties and model an area of human cognitive thinking (Saxton et al., 2019). Given a maths worded problem, for example “Sarah has 5 apples and eats 2. How many apples does she have left?”, the model needs to interpret the relation between both the numbers and the text to then solve the problem by means of arithmetic operations (Ahn et al., 2024). Therefore, an accurate number representation is primordial to distinguish between different numbers and also predict an accurate answer. The literature focuses on five different areas to better represent numbers.

### 2.1 Scaling

Increasing the number of parameters of pretrained models has improved their numerical reasoning but it is still nowhere near perfect. For example, Minerva (540B) (Lewkowycz et al., 2022) continues to struggle with greater than seven digit multiplication. Moreover, Frieder et al. (2023) found that very large models like ChatGPT and GPT4 are inconsistent when answering mathematical questions ranging from arithmetic problems to symbolic maths.

This suggests that the models lack a fundamental understanding of numbers and thus mathematics. One approach to improve number representation is to scale up the vocabulary by having more individual number tokens. For example, GPT3 (Brown et al., 2020) has unique tokens from the numbers 0-520, whereas GPT4 (OpenAI et al., 2023) has them up to 999. Despite general better performance of GPT4, it is not feasible to represent infinitely many numbers in finite memory capacity. Making the vocabulary larger also increases the computational costs.

### 2.2 Tokenisation

A more practical approach for representing all numbers is digit tokenisation (Spithourakis and Riedel, 2018; Geva et al., 2020) which separates numbers into a sequence of individual digits. This method improves upon conventional wordpiece tokenisation as shown with GenBERT (Geva et al., 2020) and Mistral-7B (Jiang et al., 2023) by reducing vocabulary size and ensuring precise representation of all numbers. Despite its advantages over conventional tokenisation algorithms, digit tokenisation has limitations. It relies on the model to aggregate digit embeddings into complete number representations. During pretraining, models typically learn to aggregate subword tokens effectively for common words. However, not all numbers are encountered frequently enough during pretraining for the model to learn accurate aggregation. As an example, when the same question is posed with numbers represented differently (once as an integer and once scaled to the thousands), FLAN large with digit tokenisation shows a performance drop of 10% (Sivakumar and Moosavi, 2023). This indicates that the model struggles with numerical consistency and accurate aggregation of digit embeddings.

### 2.3 Architectural level

Change in model architecture also aids numerical reasoning as shown by NumNET (Ran et al., 2019) and xVAL (Golkar et al., 2023). NumNET extracts the numbers from the input question and passage to create a directed graph with relative magnitude information about each number present, e.g. which is greater than the others. After encoding the input question, this comparative information is also passed to the model to also be leveraged in answering the query. Alternatively, xVAL generates two input encodings, a text only encoding where

numbers are replaced by [NUM], and a number encoding with empty space for the text but the actual value of the number in each number’s corresponding position. From the number preserving encoding, each number is converted to a vector with each entry as the number itself. The product of this vector with the embedding of [NUM] is then injected into the first layer of the transformer for each number in the input sequence. For decoding, a bespoke process is created to extract the predicted number instead of outputting the [NUM] token. Despite the positive contributions of these papers, their methods lack versatility as they are not adaptable off-the-shelf for any pretrained model.

## 2.4 Loss Functions

Another approach to improve numerical reasoning is for models to intrinsically learn better representation by introducing an inductive bias in the loss function. A simple approach is Wallace et al. (2019)’s use of the mean squared error (MSE) loss across the batch to directly predict floats on a subset of DROP (Dua et al., 2019) which consists of numerical answers. Unfortunately, this method is limited to datasets that only predict numbers. Contrastive loss has also been used to manipulate the representation of numbers, for instance, Petrak et al. (2023) draws nearer the representation generated by BPE and digit tokenisation through an auxiliary loss during extended pretraining to improve arithmetic reasoning in worded problems like DROP but also tables like SciGen (Moosavi et al., 2021). Similarly, Li et al. (2022b) use contrastive learning but on computation trees. They first generate computation trees for the mathematical operations. Then they use contrastive loss to pull nearer the graph representing the same operation, e.g. additions, and push other ones further. This is then integrated in the main loss and improves performance on two maths worded problem datasets, MathQA (Amini et al., 2019) and Math23K (Wang et al., 2017). While these loss functions are adaptable with different models, contrastive training is computationally expensive and requires annotated data.

## 2.5 Input Representation

The most model agnostic method involves changing the representation of the numbers in the input text. Wallace et al. (2019) explore worded forms of numbers, but this approach overly relies on the tokeniser which splits them into subwords. Muffo et al. (2022) decompose numbers into place val-

ues in reverse order, e.g.  $123 = 3 \text{ units}, 2 \text{ tens}, 1 \text{ hundreds}$  which helps when working with carry-on, e.g. when adding. However, this introduces many more tokens in the input which is undesirable, and also requires new vocabulary for each place value name. Zhang et al. (2020) convert all numbers into scientific notation, e.g. 314.1 is represented as  $3141[\text{EXP}]2$ , improving models’ ability to identify the magnitude of a number. Despite providing magnitudinal information, the number before [EXP] still needs to be represented faithfully. In fact, all the above strategies require the model to implicitly compute an overall aggregation for the numbers based on their individual components generated by the tokeniser of the model, whether these are digits or subwords. A simple, yet effective method is to introduce pause tokens before predicting the answer (Goyal et al., 2024). This is evaluated by training a 1B parameter transformer model on C4 using [PAUSE] tokens and a 1% improvement is shown on the numerical reasoning dataset, GSM8K (Cobbe et al., 2021). While this method can be used for inference only, they conclude that pretraining is recommended, therefore less applicable to existing models.

Within this line of research, our work is more versatile. Unlike previous methods that rely on the model to implicitly learn aggregation, we focus on the explicit aggregation of digit embeddings using mathematical priors. This provides direct supervision for the aggregation process, improving the accuracy of number representation. Furthermore, our method ensures that the embedding for a given number aligns with its numerical neighbours, enhancing the model’s numerical reasoning capabilities without altering the model architecture or requiring extensive retraining.

## 3 Aggregation of Digit Embeddings

Digit tokenisation has demonstrated its efficacy in enhancing numerical reasoning compared to BPE tokenisation. This improvement can be attributed to digit tokenisation’s utilisation of pretrained embeddings for individual digits, allowing the model to learn the overall representation through contextualised embeddings. In contrast, BPE tends to fragment longer and less frequent numbers into random subsequences, resulting in less meaningful aggregations than those achieved through digit tokenisation. However, the implicit aggregation process employed by digit tokenisation remains

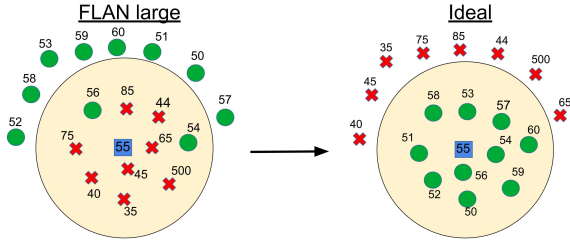


Figure 1: A 2D projection of the neighbourhood of the number token “55” in FLAN large is represented on the left. Ideally, number embeddings should reflect natural numerical proximity. In other words, the embedding for any given number should closely align with those of its immediate numerical neighbours, depicted on the right.

unclear, particularly how the model aggregates a number’s overall representation from its digit embedding.

In this paper, we investigate a natural continuation of digit tokenisation, a mathematically motivated aggregation that takes into account the relative position of each digit within a number. Our approach generates an overall embedding for the number by considering the positional weight of each individual digit in that number. For example, given “123”, the common understanding of numbers as base-10 is “ $1 \times 100 + 2 \times 10 + 3 \times 1$ ”, so the left most digits are weighted higher as they have a larger impact on the value of the number.

We design our weighted scheme such that: (1) the embeddings of single-digit numbers remain intact, as these embeddings are effectively learned during pretraining, evidenced by the high performance of models on single-digit operations (Sivakumar and Moosavi, 2023), (2) the weights of consecutive place values increase exponentially to reflect base-10, and (3) the weights do not sum to 1, meaning that it does not normalise the sum, allowing for numbers composed of the same digits, e.g. “111” and “11”, to be represented differently. These properties would introduce a bias towards an accurate length of numbers and the correct digits from left to right as the left most digits are amplified, hence preserving natural numerical order.

We propose to calculate the weighted aggregated embedding  $\mathbf{a}$  with  $a_i = \sum w_i \cdot d_i$  for  $1 \leq i \leq N$  where  $N$  is the number of digits, and the weights  $w_i$  are defined as:

$$w_i = 2^{N-i} \times \frac{3(N+1-i)(N+2-i)}{N(N+1)(N+2)}. \quad (1)$$

These weights are designed to satisfy three key

properties. **(1) Alignment with single-digit representations:** when  $N = 1$ ,  $w_1 = 1$ , ensuring compatibility with the model’s pretraining on single digits. **(2) Exponential growth:** the exponential component  $2^{N-i}$  mimics the base-10 system, providing an appropriate scale without causing the weights to grow too rapidly. This also ensures that the weights are not normalised. **(3) Regularisation Term:** the fractional component acts as a regularisation term, forming a normalised triangular number sequence. For instance, for a 3-digit number, the triangular sequence is 1,3,6, normalised to 0.1,0.3,0.6. This ensures that the difference between consecutive digit weights increases proportionally, i.e.,  $w_i - w_{i-1} = w_0 \times i$ , replicating the exponential ratio between digit positions in a logarithmic space.

To validate the ability of an aggregated embedding to accurately represent numerical relationships, we use the F1-score to compare natural k-Nearest Neighbours ( $nkNN$ ) with embedding k-Nearest Neighbours ( $ekNN$ ). This comparison serves two purposes: firstly, to assess the embeddings’ capacity to distinguish between distinct numbers, and secondly, to evaluate how well these embeddings mirror the natural numerical order. By defining  $nkNN$  as the set of mathematically adjacent numbers to a given integer  $n$ , and  $ekNN$  as the set of its closest neighbours in the embedding space, we create a direct measure of the embedding’s effectiveness in preserving numerical proximity. The F1-score evaluates the alignment between  $nkNN$  and  $ekNN$ , penalising both the inclusion of incorrect neighbours and the omission of correct ones. A strong correlation between  $nkNN$  and  $ekNN$ , as reflected in a high F1-score, indicates that the embeddings faithfully capture the essence of numerical data as illustrated in Figure 1.

We compare our bespoke weighted aggregation function to more standard aggregation functions such as sum. For a set of digit embeddings, we apply these functions along each dimension to generate a unique embedding for the number represented by these digits. Figure 2 graphs the F1-score for our weighted function and sum over different digit lengths, i.e. 2-digit would be the numbers 10 to 99. Appendix A has the results for other aggregation functions: max, min, mean and median; these have the lowest alignment with natural order with an F1-score below 5%. These functions all have a normalising property meaning that the length of the number has no bearing on the aggregated embed-

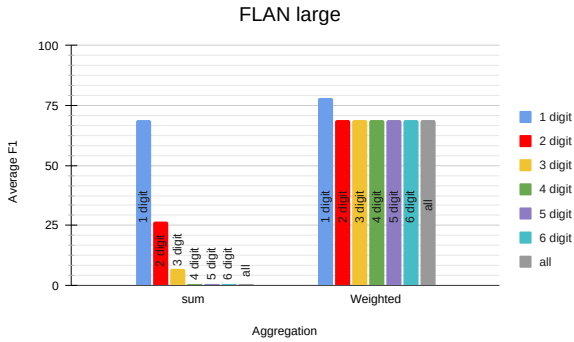


Figure 2: Average F1-score of FLAN large layer 1 numbers using sum and our weighted aggregation function with neighbourhood of 10.

ding, as the functions only retrieve one entry for each dimension therefore cases like “1111” would be equivalent to both “11” and “1”. Contrastingly, sum has better F1-scores for up to 3 digits as it possesses magnitudinal information since all the entries are summed up for each dimension distinguishing, for instance, a 2-digit set from a 3-digit set as it simply adds more numbers. However, it is position agnostic - it assigns equal weight to all the digits irrespective of their relative positions. Therefore, the embeddings generated from permutations of the same digits will always be equivalent, e.g. “85” and “58”. Since larger digit numbers have more such permutations, the F1-score reduces as the number of digits increases. Using this metric, the best aggregation is our weighted sum, the average F1-score rounds to 69% for 2 digits onwards suggesting that our weighted sum is closer to the ideal depiction in Figure 1. Undoubtedly, 1-digit F1-score is better as these embeddings are generated from pretraining, but also because the weighted scheme ensures that they are separated from the other number embeddings.

Despite this weighted scheme aligning the number embeddings with their natural order, the weights generated by Equation 1 can surpass the precision used making it too large after a certain point. However, this behaviour is attenuated by the regularisation term which maintains the high F1-score of 69% for, at least, up to 6-digit long numbers as shown in Figure 2. Theoretically, the newly formed number representation should contain numbers that beginning with the same digits and only vary at the unit level. For example, the neighbourhood of 4523 should contain all numbers of the form 452X where X is a digit from 0 to 9, therefore eight of these coincide in the 10-nearest

neighbour, namely 4520, 4521, 4522, 4524, 4525, 4526, 4527, and 4528.

## 4 Integrating Aggregated Embeddings

Given the construction of our mathematically grounded aggregation, we explore two distinct methodologies for enhancing numerical understanding in models, each targeting different aspects of number representation. The first method focuses on enriching the input data by integrating a mathematical aggregation directly into the input embedding as a special token. This approach requires no changes to the model’s architecture, making it a flexible solution compatible with various models and suitable for a broad spectrum of tasks.

In contrast, the second approach aims to refine the model’s output by improving how numbers are represented in the learned outcomes. This is achieved by incorporating the aggregation in the loss function, encouraging the model to generate number embeddings that align more closely to the correct numerical values. Specifically, this method includes an additional term in the loss calculation, which accounts for the distance between the aggregated embedding of the predicted number and that of the true number. This targeted intervention is particularly effective in tasks requiring precise numerical predictions, helping the model develop a more nuanced and accurate representation of numbers.

The baseline implementation for both methods is the same as Petrak et al. (2023) with digit tokenisation surrounded by [F] and [/F] tokens to mark the start and end of the number identified using the regular expression “(d\*\.)?d+”.

### 4.1 Aggregation in Input Embeddings

In our first approach, we enhance the input embedding by incorporating the computed aggregation directly. This is achieved by first digitising numbers and delineating them with special tokens as done by Petrak et al. (2023). Additionally, we introduce a special token, [AGG], positioned as follows where  $d_i$  represent the digit tokens: [F] [AGG] [ $d_1$ ] ... [ $d_n$ ] [/F]. The embedding for this [AGG] token is initialised with the aggregation of the digit embeddings based on Equation 1.

### 4.2 Aggregation in Loss Function

Language generation models typically use a cross-entropy loss function ( $\mathcal{L}_{CE}$ ) (Lewis et al., 2020;

Raffel et al., 2020). To improve the model’s ability to predict numbers accurately, we introduce an auxiliary loss ( $\mathcal{L}_{AUX}$ ) to calculate the mean squared error between the aggregate embedding of the gold and predicted numbers. Understanding and predicting numbers is inherently more complex than predicting a single word or sub-word because they consist of multiple digits, each carrying different significance. For example, in answering the question “Mary’s salary is £900 a month, but she pays £579 in rent. How much salary does she have left?”, the answers 320, 230, 32, or 456 are all incorrect. However, 320 is more accurate compared to others because its magnitude is closer to the correct answer, 321. Incorporating this new auxiliary loss would help the model predict digits that are closer to the gold answer, enhancing its precision in numerical predictions by recognising the relative significance of each digit within a number.

Given a prediction  $p$  and the gold label  $l$ , we compute the weighted sum of the digits<sup>2</sup> for both  $p$  and  $l$ . This process generates two single embedding representations:  $W(p)$  for the prediction, and  $W(l)$  for the gold label. The distance between these two embeddings is then calculated using the  $\log^3$  mean squared error (equivalent to the euclidean distance):

$$\mathcal{L}_{AUX} = \log_2 (\|W(p) - W(l)\|_2) \quad (2)$$

Finally, the two losses are linearly interpolated by a hyperparameter,  $\lambda$ :

$$\mathcal{L} = \lambda \times \mathcal{L}_{CE} + (1 - \lambda) \times \mathcal{L}_{AUX} \quad (3)$$

## 5 Experimental Setup

Both methods are evaluated on two different pre-trained models, BART base (140M) (Lewis et al., 2020) and FLAN base (250M) (Wei et al., 2022a). Additionally, we evaluate on FLAN large (780M) to explore the effect of model size. BART is an encoder-decoder pre-trained on five corrupted document tasks from books and Wikipedia data. FLAN is an instruction-finetuned version of T5 (Raffel et al., 2020) which is trained on C4 using transfer learning.

We evaluate our proposed methods on two different test sets: FERMAT (Sivakumar and Moosavi, 2023), and MAWPS (Koncel-Kedziorski et al., 2016). Both FERMAT and MAWPS consist of

<sup>2</sup>Should the answers not be numerical, the model is penalised by arbitrarily setting  $\mathcal{L}_{AUX}$  to 20.

<sup>3</sup>Log base 2 is used to regularise the auxiliary loss.

English maths worded problem that can be tackled by BART and FLAN, as shown by Sivakumar and Moosavi (2023) and where the answer is a single number. This enables us to evaluate our method strictly on numerical outputs reducing the interference of other difficulties such as predicting words and units, or extracting spans. FERMAT is a multi-view evaluation set which has different test sets with different number representations while keeping the maths problem fixed. The different test sets distinguish different number types of which we select the ones that separate integers into digit length (2-digit, 3-digit, 4-digit), contain a mixture of integers less than 1000, contain a mixture of integers greater than 1000, the sets of one and two decimal place numbers, and a test set that takes the original set and scales the number to more than 4-digit numbers; these allow us to evaluate which number representation the models support better. FERMAT’s training set is augmented from different templates making it independent to its test sets. MAWPS, on the other hand, has the same domain for both training and testing. It is a widely used dataset to evaluate numerical reasoning, chiefly because it is small and easy to train with small models. We finetune the models on each dataset’s respective training data (see Appendix B) using the hyperparameters described in Appendix C.

Accuracy is the general metric used to evaluate these datasets, however, since it is sometimes too stringent and neglects to reflect some improvements of the model, we also use a variation of edit distance (Levenshtein, 1966) as a supplementary metric. Edit distance helps see improvement in the predictions despite being incorrect; it calculates how many insertions, deletions or substitutions is required for the prediction to be transformed into the gold label number on a string level. In this paper, we will use Character Error Rate (CER) which is a character level (digit level) edit distance as a percentage over the string length of the target. The lower the CER, the closer the prediction is to the gold label.

## 6 Impact of Integrating Aggregations

Table 1 presents the results of our exploration into the effects of integrating mathematical aggregation into the three models across two distinct settings. The bold values indicate the stronger improvement between the two incorporation strategies. For the majority of the test splits, the strongest perfor-

Incorporating Weights (Accuracy %)		MAWPS	FERMAT													
			Original	Commuted	Integers 0 to 1000	2-digit integers	3-digit integers	4-digit integers	1000+	1000+ same	1dp random	2dp random	a+b	a-b	a*b	a/b
BART base (140M)	Digits	19.20	16.65	8.73	10.26	13.41	10.89	7.74	5.58	10.89	17.82	8.37	40.91	10.62	9.56	11.76
	[AGG] + Digits	<b>+2.00</b>	+0.63	+1.53	<b>-1.17</b>	-0.90	<b>-2.16</b>	-0.27	+0.09	<b>+0.09</b>	<b>+1.08</b>	-0.27	<b>-3.90</b>	-0.74	+1.77	0.00
	Digits + Aux Loss	+1.40	<b>+1.89</b>	<b>+1.80</b>	<b>+0.54</b>	<b>+0.81</b>	0.00	<b>+0.81</b>	<b>+1.17</b>	<b>-1.26</b>	+0.18	<b>+0.63</b>	<b>+2.01</b>	<b>+0.19</b>	<b>+4.25</b>	<b>-1.27</b>
FLAN base (250M)	Digits	23.00	28.35	17.82	17.10	22.86	17.37	13.77	10.35	18.72	25.83	18.45	63.38	19.57	12.92	11.27
	[AGG] + Digits	+0.80	<b>+2.79</b>	+0.27	+2.52	+0.81	<b>+1.80</b>	<b>+2.79</b>	<b>+1.80</b>	+0.90	+0.45	-0.09	<b>+4.48</b>	+3.21	-0.27	+1.08
	Digits + Aux Loss	<b>+1.80</b>	+2.25	<b>+0.36</b>	<b>+3.15</b>	<b>+2.16</b>	+1.71	<b>+2.79</b>	+0.81	<b>+3.87</b>	<b>+1.89</b>	-0.18	+3.90	<b>+5.80</b>	<b>+0.27</b>	<b>+1.57</b>
FLAN large (780M)	Digits	28.80	42.39	21.06	25.65	31.32	24.30	21.87	16.47	23.31	36.36	25.83	63.12	39.88	18.23	18.14
	[AGG] + Digits	<b>+1.20</b>	+0.45	<b>+0.45</b>	+0.81	+2.07	<b>+2.79</b>	<b>+0.99</b>	+1.35	<b>+2.88</b>	+0.27	+0.54	<b>+6.17</b>	<b>+3.83</b>	<b>+0.53</b>	<b>+1.47</b>
	Digits + Aux Loss	+1.00	<b>+0.99</b>	-0.18	<b>+1.62</b>	<b>+2.88</b>	<b>+2.79</b>	+0.72	<b>+1.53</b>	+1.26	<b>+1.26</b>	<b>+0.63</b>	-0.39	+1.79	+0.18	<b>-1.08</b>

Table 1: Results change in Accuracy from baseline after including aggregate embeddings in input embedding ([AGG] + Digits) and auxiliary loss (Digits + Aux Loss) for BART base, FLAN base and FLAN large. Darker shades of green and red indicate an absolute change greater than 1%.

Incorporating Weights (CER %)		MAWPS	FERMAT													
			Original	Commuted	Integers 0 to 1000	2-digit integers	3-digit integers	4-digit integers	1000+	1000+ same	1dp random	2dp random	a+b	a-b	a*b	a/b
BART base (140M)	Digits	77.73	89.59	90.32	72.87	71.93	72.25	74.04	77.01	50.29	54.42	62.23	50.31	74.12	60.73	75.51
	[AGG] + Digits	<b>-1.79</b>	<b>-12.40</b>	<b>-0.83</b>	+0.46	+0.51	<b>+1.19</b>	-0.16	-0.44	+0.94	-1.38	-1.28	<b>+3.08</b>	<b>-1.58</b>	<b>+1.21</b>	<b>-2.22</b>
	Digits + Aux Loss	+0.76	-1.88	-0.53	+0.17	+0.20	+0.34	<b>-1.06</b>	<b>-0.53</b>	<b>-1.89</b>	<b>-1.59</b>	<b>-1.78</b>	<b>-2.45</b>	-0.23	<b>-2.75</b>	+0.26
FLAN base (250M)	Digits	67.71	75.32	169.52	67.37	67.68	67.94	67.86	68.86	50.95	43.77	47.80	39.84	87.81	60.96	91.52
	[AGG] + Digits	-0.98	<b>-1.40</b>	-0.29	<b>-1.11</b>	<b>-1.41</b>	<b>-1.19</b>	<b>-1.67</b>	-0.96	<b>+1.26</b>	-1.33	<b>-0.39</b>	<b>-1.64</b>	-1.94	-0.17	-0.50
	Digits + Aux Loss	<b>-1.54</b>	-0.83	<b>-1.09</b>	-1.09	-1.15	-0.80	-1.39	<b>-1.23</b>	<b>-2.09</b>	<b>-1.82</b>	-0.30	-1.25	<b>-3.15</b>	<b>-0.72</b>	<b>-0.93</b>
FLAN large (780M)	Digits	63.13	69.71	76.46	63.02	62.69	63.53	63.96	66.67	49.90	37.63	42.31	39.00	58.84	52.84	70.49
	[AGG] + Digits	-2.57	<b>-44.77</b>	<b>-10.81</b>	-1.02	-0.10	<b>-1.63</b>	-0.65	-0.89	<b>+1.78</b>	-0.93	-1.23	<b>-6.16</b>	<b>-7.80</b>	<b>-5.49</b>	<b>-7.19</b>
	Digits + Aux Loss	<b>-3.45</b>	-45.42	-2.72	<b>-1.20</b>	<b>-0.24</b>	-1.09	<b>-1.23</b>	<b>-1.31</b>	<b>-2.57</b>	<b>-1.11</b>	<b>-1.27</b>	-3.47	-6.14	-2.93	-4.74

Table 2: Results in Character Error Rate (CER) as a percentage over the target string with change from baseline after including aggregate embeddings in input embedding ([AGG] + Digits) and auxiliary loss (Digits + Aux Loss) for BART base, FLAN base and FLAN large. With CER, a lower value indicates a better performance. Green highlight reduced CER (negative change), while red indicates the opposite. Darker shades of green and red indicate an absolute change greater than 1%.

mance of the examined models is observed when the aggregation is incorporated into the auxiliary loss. This suggests that incorporating aggregation at the output level is more effective than incorporating it in the input embedding. However, this may be due to the fact that adding a new token in the input might require more than just fine-tuning, such as an extended pretraining phase. This aligns with the observations made by Goyal et al. (2024), who found that the addition of the pause token only became effective from pretraining.

FLAN large, on the other hand, has a more balanced performance but an overall higher improvement when the aggregation is incorporated in the input as shown particularly from all the green cells in the row [AGG] + Digits. Therefore, a certain model size may be required to learn a new token and leverage the information it provides. This reinforces that an aggregated embedding provides useful signal to improve number understanding but

how it is integrated is also crucial.

For the operations, the improvements is generally positive across all of them, however, evidently greater for addition and subtraction than multiplication and division. This resonates with the fact that digits positions are more informative for the first two operations, especially when, for instance, aligning them to perform calculations.

When focusing on smaller integers (columns “Integers 0 to 1000” to “4-digit integers”), incorporating the weighted embedding in the auxiliary loss consistently yields better performance, with all cells being green and showing the highest scores. For smaller integers, models likely already possess a strong implicit representation, making the explicit [AGG] token less impactful. However, at the decoding stage, the auxiliary loss enhances precision by penalising incorrect predictions.

For the 1000+ columns, using accuracy, the pattern is not evident, however, Table 2 presents the

Aggregated Embedding (Accuracy %)		MAWPS	FERMAT													
			Original	Commuted	Integers 0 to 1000	2-digit integers	3-digit integers	4-digit integers	100+	100+ same	1dp random	2dp random	a+b	a-b	a*b	a/b
BART base (140M)	Digits	19.20	16.65	8.73	10.26	13.41	10.89	7.74	5.58	10.89	17.82	8.37	40.91	10.62	9.56	11.76
	Digits + [AGG]	-1.40	-14.76	-7.74	-8.82	-10.98	-8.73	-6.75	-5.58	-10.35	-14.76	-7.83	-36.82	-9.38	-8.94	-9.51
	[AGG] + Digits	<b>+2.00</b>	<b>+0.63</b>	<b>+1.53</b>	-1.17	-0.90	-2.16	0.27	<b>+0.09</b>	<b>+0.09</b>	<b>+1.08</b>	0.27	3.90	0.74	<b>+1.77</b>	0.00
	[PAUSE] + Digits	-1.40	+0.18	-0.45	-0.18	-0.63	-0.90	-0.36	-0.27	-3.87	-0.90	0.00	-8.51	-0.31	+1.68	-2.06
FLAN base (250M)	Digits	23.00	28.35	17.82	17.10	22.86	17.37	13.77	10.35	18.72	25.83	18.45	63.38	19.57	12.92	11.27
	Digits + [AGG]	<b>+1.80</b>	-1.53	-2.07	+0.99	-1.89	-0.36	+0.63	+1.35	-0.63	-1.98	-0.99	+0.45	+3.89	-2.39	-0.10
	[AGG] + Digits	+0.80	<b>+2.79</b>	<b>+0.27</b>	<b>+2.52</b>	+0.81	<b>+1.80</b>	<b>+2.79</b>	+1.80	+0.90	+0.45	-0.09	<b>+4.48</b>	+3.21	-0.27	+1.08
	[PAUSE] + Digits	+1.00	+2.07	-0.54	+1.98	<b>+1.44</b>	<b>+1.80</b>	+2.61	<b>+2.52</b>	<b>+2.16</b>	<b>+2.61</b>	<b>+1.71</b>	+3.18	<b>+5.99</b>	<b>1.95</b>	<b>+3.43</b>
FLAN large (780M)	Digits	28.80	42.39	21.06	25.65	31.32	24.30	21.87	16.47	23.31	36.36	25.83	63.12	39.88	18.23	18.14
	Digits + [AGG]	-2.80	-2.16	<b>+1.35</b>	<b>+1.89</b>	+1.08	+1.44	+1.62	+2.16	<b>+5.40</b>	-1.17	+0.54	<b>+8.57</b>	-8.15	-0.97	+1.18
	[AGG] + Digits	<b>+1.20</b>	<b>+0.45</b>	+0.45	+0.81	+2.07	+2.79	+0.99	+1.35	+2.88	+0.27	+0.54	+6.17	<b>+3.83</b>	<b>+0.53</b>	+1.47
	[PAUSE] + Digits	-1.40	-0.45	-0.45	<b>+1.89</b>	<b>+3.69</b>	<b>+2.88</b>	<b>+3.06</b>	<b>+2.25</b>	+5.04	<b>+1.17</b>	<b>+2.61</b>	+6.17	+1.17	-1.77	<b>+3.53</b>

Table 3: Comparing the aggregated embedding at the input level with a pause token and positioning the token after the digits. Darker shades of green and red indicate an absolute change greater than 1%.

character error rate (CER) comparing both incorporating strategies for all three models, and highlights that using the auxiliary loss clearly reduces the CER more than explicitly using the aggregation in the input. The auxiliary loss encourages the model to predict the correct answer as the CER is lower. However, since the weights assigned to each digit position is lower as it gets closer to the units, the auxiliary accounts less for it, reducing precision. As a consequence, despite the CER reducing, since the entire number is not predicted correctly, improvement fails to be reflected in the accuracy.

## 7 Analysis of Aggregation Embedding in the Input

The first integration method relies on prepending the aggregated embedding token, [AGG], before the digits. The position of the token is before what it represents, similar in nature to BERT’s (Devlin et al., 2019) [CLS] token, which is an aggregation token of the entire input. However, Goyal et al. (2024) use a [PAUSE] token posteriori to the digit tokens to act as processing time after concluding that prepending it had less impact. Consequently, we also evaluate our proposed method by appending the aggregation token, i.e. Digits + [AGG]. Table 3 clearly shows that this configuration for both base models underperforms compared to [AGG] + Digit as rows have more red entries. In fact, it performs worse than the baseline with only digit tokenisation. For FLAN large, the results between [AGG] prepended and appended are closer to one another, but prepended, the impact is positive for each test set and on average better by 1% than [AGG] used posteriori. Seeing the token before

the digits might provide magnitude information of the overall number which would indicate the importance of each digit to come, whereas having it after might interfere with the representation that the model has already started to create implicitly from seeing the digits first.

Additionally, we test the impact of providing the aggregated token by replacing it with a randomly initialised [PAUSE] token akin to Goyal et al. (2024). From Table 3, we observe that for BART, neither [AGG], nor [PAUSE] have a great positive impact on the performance. This confirms that BART struggles to learn new tokens from fine-tuning alone. The FLAN models are more adaptable to the new tokens as seen by the greener rows. However, the overwhelming bold entries with the [PAUSE] token indicate that both FLAN base and large perform better with a [PAUSE] token acting as a blank space for the model to process the information. It is possible that the model uses this token to create an implicit representation of the number. Nevertheless, the average improvement between the [PAUSE] and [AGG] differs by less than 0.5% implying that a different aggregation function or a full hyperparameter search could reverse the trend.

## 8 Future Work

Our proposed aggregation strategy has shown encouraging steps towards better number representation. However, as with observations made in previous work, the effect of new strategies report minimal improvement on smaller models but greater impact on larger models (Cobbe et al., 2021; Wei et al., 2022b). Therefore, an evaluation of our pro-



posed method on larger scale models would verify the scalability of this approach.

The weighting scheme, presented in Equation 1, offers a straightforward method for aggregating digit embeddings. However, as numbers increase in length, their aggregated embeddings tend to drift away from the original numerical embedding subspace. This divergence could be addressed by enabling the model to adapt to this new embedding space by exploring extended pretraining, or constructing weighting schemes that remain closer to the numerical subspace while satisfying the criteria outlined in Section 3.

Our auxiliary loss, grounded in Mean Squared Error, shows promising results for penalising the model’s erroneous predictions and nudging it towards more accurate outcomes. Given that the values resulting from standard cross-entropy and the MSE of the aggregated embeddings may span vastly different value ranges, crafting a loss function that aligns more closely in magnitude with the output of cross-entropy could mitigate the risk of exerting excessive regularisation pressure.

## 9 Conclusion

Improving numerical reasoning is a challenging task, increasing model sizes or focusing on data augmentation helps but at the cost of a substantial additional training time or computations. Digit tokenisation has been a pioneering in improving how models encode and decode numbers; however the aggregation of the digit is performed implicitly. We advance this idea by explicitly providing an aggregated number embedding that is more mathematically sound. These embeddings are generated as weighted sums of the digit embeddings by accounting for the digits relative position in the number. We then incorporate them in two model agnostic forms: in the input level as an additional token, and in an auxiliary MSE loss. Our promising results demonstrate that, as a proof-of-concept, even a straightforward aggregation with simple incorporation techniques can positively impact number understanding. Therefore, testing it at a larger scale, developing sophisticated aggregation functions, and refining the integration of the auxiliary loss presents valuable avenues for future research.

## 10 Limitations

Some of the limitations of this work are discussed in the Future Work section. However, we give de-

tails of further limitations relating to the size of the models used, and the compatibility and growth of our proposed weighted aggregation function.

Due to financial and resource constraints the hypothesis that the methods for incorporating the aggregated embedding in larger architectures would lead to greater performance based on the improvement observed on smaller model is not verified.

In addition, while the weighted scheme is designed using mathematical priors, it is specifically created for integers, therefore it may not be compatible with decimals or alternative representation of numbers such as 01 for 1. Nonetheless, from Table 2, we note that CER reduces for both 1dp and 2dp; therefore our aggregated embedding method has promising scope for all numbers.

Furthermore, the weight function described in Equation 1 does not converge, therefore for a sufficiently large number of digit it would grow beyond the accuracy provided by the model. However, we explain in Section 3 with the aid of Figure 2 that, for up to 6-digits, the weighted scheme functions well with no signs of deterioration. Moreover, in natural text, very large numbers tend to be shorten using a more appropriate unit, for example, the world population of 8114693010 is more often expressed as 8 billion reducing the numbers of digits needed considerably. But this raises the question of predicting the correct unit which would lead to future work.

Nonetheless, our weighting scheme leverages digit embedding, therefore it is heavily dependent on them, particularly on the relative distance of the digit embedding to one another. In FLAN large, the embedding of the digit 0 is more distant from the embeddings of other digits, which causes it to frequently include numbers ending in 0 when the target contains a 0, or to exclude them otherwise. As explained in Section 3, the optimal neighbour would include all numbers with different unit value, and this alone would achieve an F1-score of at least 70%.

Lastly, the experiments were conducted using a single random seed. While this ensures consistency and reproducibility, having access to better resources would have enabled us to run the experiments with multiple seeds. This would have allowed us to calculate the average improvement achieved by using aggregated digit embeddings to represent numbers.

## Acknowledgements

This work was supported by the Centre for Doctoral Training in Speech and Language Technologies (SLT) and their Applications funded by UK Research and Innovation [grant number EP/S023062/1]. We also acknowledge IT Services at The University of Sheffield for the provision of services for High Performance Computing. Additional thanks to the reviewers for their encouraging comments and discussion, and particularly to Danae Sanchez Villegas, Mugdha Pandya, Valeria Pastorino, Huiyin Xue and Constantinos Karouzos for their continued feedback throughout the research.

## References

- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. [Large language models for mathematical reasoning: Progresses and challenges](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 225–237, St. Julian’s, Malta. Association for Computational Linguistics.
- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [MathQA: Towards interpretable math word problem solving with operation-based formalisms](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Chung-Chi Chen, Hen-Hsen Huang, Yow-Ting Shiue, and Hsin-Hsi Chen. 2018. [Numeral understanding in financial tweets for fine-grained crowd-based forecasting](#). In *2018 IEEE/WIC/ACM International Conference on Web Intelligence (WI)*, pages 136–143.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. [Training verifiers to solve math word problems](#). *arXiv preprint arXiv:2110.14168*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.
- Abhijeet Dubey, Lakshya Kumar, Arpan Somani, Aditya Joshi, and Pushpak Bhattacharyya. 2019. [“when numbers matter!”: Detecting sarcasm in numerical portions of text](#). In *Proceedings of the Tenth Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 72–80, Minneapolis, USA. Association for Computational Linguistics.
- Simon Frieder, Luca Pinchetti, Alexis Chevalier, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Christian Petersen, and Julius Berner. 2023. [Mathematical capabilities of chatGPT](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Mor Geva, Ankit Gupta, and Jonathan Berant. 2020. [Injecting numerical reasoning skills into language models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 946–958, Online. Association for Computational Linguistics.
- Siavash Golkar, Mariel Pettee, Michael Eickenberg, Alberto Bietti, Miles Cranmer, Geraud Krawezik, Francois Lanusse, Michael McCabe, Ruben Ohana, Liam Parker, Bruno Régaldo-Saint Blancard, Tiberiu Tesileanu, Kyunghyun Cho, and Shirley Ho. 2023. [xval: A continuous number encoding for large language models](#). *arXiv preprint arXiv:2310.02989*.
- Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2024. [Think before you speak: Training language models with pause tokens](#). In *The Twelfth International Conference on Learning Representations*.
- Pengfei Hong, Navonil Majumder, Deepanway Ghosal, Somak Aditya, Rada Mihalcea, and Soujanya Poria.

2024. [Evaluating llms' mathematical and coding competency through ontology-guided interventions.](#) *arXiv preprint arXiv:2401.09395*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. [Mistral 7b.](#) *arXiv preprint arXiv:2310.06825*.
- Maël Jullien, Marco Valentino, Hannah Frost, Paul O'regan, Donal Landers, and André Freitas. 2023. [SemEval-2023 task 7: Multi-evidence natural language inference for clinical trial data.](#) In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 2216–2226, Toronto, Canada. Association for Computational Linguistics.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners.](#) In *ICML 2022 Workshop on Knowledge Retrieval and Language Models*.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. [MAWPS: A math word problem repository.](#) In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1152–1157, San Diego, California. Association for Computational Linguistics.
- Vladimir I. Levenshtein. 1966. [Binary codes capable of correcting deletions, insertions, and reversals.](#) *Soviet physics. Doklady*, 10:707–710.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Aitor Lewkowycz, Anders Johan Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Venkatesh Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. [Solving quantitative reasoning problems with language models.](#) In *Advances in Neural Information Processing Systems*.
- Ailisi Li, Yanghua Xiao, Jiaqing Liang, and Yunwen Chen. 2022a. [Semantic-based data augmentation for math word problems.](#) In *Database Systems for Advanced Applications: 27th International Conference, DASFAA 2022, Virtual Event, April 11–14, 2022, Proceedings, Part III*, page 36–51, Berlin, Heidelberg. Springer-Verlag.
- Zhongli Li, Wenxuan Zhang, Chao Yan, Qingyu Zhou, Chao Li, Hongzhi Liu, and Yunbo Cao. 2022b. [Seeking patterns, not just memorizing procedures: Contrastive learning for solving math word problems.](#) In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2486–2496, Dublin, Ireland. Association for Computational Linguistics.
- Nafise Sadat Moosavi, Andreas Rücklé, Dan Roth, and Iryna Gurevych. 2021. [Scigen: a dataset for reasoning-aware text generation from scientific tables.](#) In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Matteo Muffo, Aldo Cocco, and Enrico Bertino. 2022. [Evaluating transformer language models on arithmetic operations using number decomposition.](#) In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 291–297, Marseille, France. European Language Resources Association.
- R OpenAI et al. 2023. [Gpt-4 technical report.](#) *ArXiv*, 2303:08774.
- Dominic Petrak, Nafise Sadat Moosavi, and Iryna Gurevych. 2023. [Arithmetic-based pretraining improving numeracy of pretrained language models.](#) In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (\*SEM 2023)*, pages 477–493, Toronto, Canada. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer.](#) *J. Mach. Learn. Res.*, 21(1).
- Qiu Ran, Yankai Lin, Peng Li, Jie Zhou, and Zhiyuan Liu. 2019. [NumNet: Machine reading comprehension with numerical reasoning.](#) In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2474–2484, Hong Kong, China. Association for Computational Linguistics.
- David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. 2019. [Analysing mathematical reasoning abilities of neural models.](#) In *International Conference on Learning Representations*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units.](#) In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Jasivan Sivakumar and Nafise Sadat Moosavi. 2023. [FERMAT: An alternative to accuracy for numerical](#)

reasoning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15026–15043, Toronto, Canada. Association for Computational Linguistics.

Georgios Spithourakis and Sebastian Riedel. 2018. Numeracy for language models: Evaluating and improving their ability to predict numbers. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2104–2115, Melbourne, Australia. Association for Computational Linguistics.

Avijit Thawani, Jay Pujara, Filip Ilievski, and Pedro Szekely. 2021. Representing numbers in NLP: a survey and a vision. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–656, Online. Association for Computational Linguistics.

Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh, and Matt Gardner. 2019. Do NLP models know numbers? probing numeracy in embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5307–5315, Hong Kong, China. Association for Computational Linguistics.

Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. Deep neural solver for math word problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 845–854, Copenhagen, Denmark. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022b. Chain of thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.

Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. Meta-math: Bootstrap your own mathematical questions for large language models. In *The Twelfth International Conference on Learning Representations*.

Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2024. Mammoth: Building math generalist models through hybrid instruction tuning. In *The Twelfth International Conference on Learning Representations*.

Xikun Zhang, Deepak Ramachandran, Ian Tenney, Yanai Elazar, and Dan Roth. 2020. Do language

embeddings capture scales? In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4889–4896, Online. Association for Computational Linguistics.

## Appendix

### A Aggregation functions

Figure 3 shows that F1-score for numbers with up to 6-digits across six different aggregation functions. The F1-score for max, min, mean and median are all below 5%.

### B Datasets

The datasets’ split is given in Table 4. MAWPS is a dataset generated by combining different ones ranging from addition and subtraction to simultaneous equations. The collation of questions is split to create the train, development and test set. FERMAT is a large dataset which has a training and development set automatically generated from 100 templates using different numbers from the following four categories: small integers (less than 1000), large integers (between 1000 and 100000), 1 decimal place and 2 decimal place numbers. The test set is independently generated from two maths worded problem datasets, and then augmented to create 21 test sets of which we use 11.

Datasets	Train	Dev	Test
MAWPS	1500	373	500
FERMAT	200000	1000	1111x11

Table 4: Train, development, and test splits of MAWPS and FERMAT.

### C Hyperparameters

All experiments were conducted using an Nvidia Tesla A100 with 80G and with a weight decay of 0.005, warm-up of 100, float32 and 3 generation beams, max input length = 128, max target length=16, and seed=42. Due to limited computational resources, a full grid search of hyperparameter was impossible, however, we do a lambda search in the range 0.4 to 0.8 in 0.05 increments. Specific hyperparameters as well as computation time for dataset and model combinations can be found in Table 5.

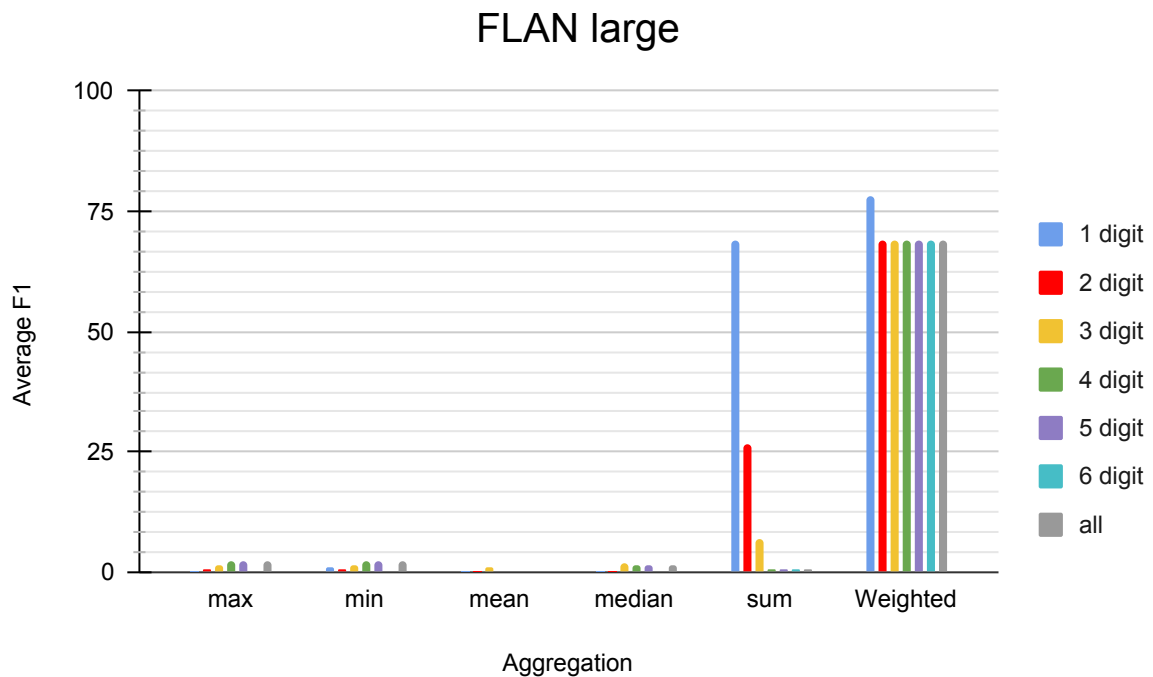


Figure 3: Average F1-score of FLAN large layer 1 numbers using max, min, median, mean sum and our weighted aggregation function with neighbourhood of 10. The bars are in the order of the legend top to bottom, reflected left to right.

Datasets	Models	Learning Rate	Epochs	Batch Size	Lambda	Training Time
MAWPS	BART base	1.00E-04	150	128	0.6	1h
	FLAN base		150	64	0.6	1h
	FLAN large		100	16	0.65	1.5h
FERMAT	BART base	1.00E-05	50	128	0.6	37h
	FLAN base		50	64	0.65	48h
	FLAN large		50	16	0.4	87h

Table 5: Specific hyperparameters for MAWPS and FERMAT based on the models trained. Training time is also provided as a rounded figure.