

Laying Anchors: Semantically Priming Numerals in Language Modeling

Mandar Sharma
Virginia Tech
mandarsharma@vt.edu

Rutuja Murlidhar Taware
Virginia Tech
trutujamurlidhar@vt.edu

Pravesh Koirala
Vanderbilt University
pravesh.koirala@vanderbilt.edu

Nikhil Muralidhar
Stevens Institute of Technology
nmurali1@stevens.edu

Naren Ramakrishnan
Virginia Tech
naren@cs.vt.edu

Abstract

Off-the-shelf pre-trained language models have become the de facto standard in NLP pipelines for a multitude of downstream tasks. However, the inability of these models to properly encode numerals limits their performance on tasks requiring numeric comprehension. We introduce strategies to semantically prime numerals in any corpus by generating anchors governed by the distribution of numerals in said corpus, thereby enabling mathematically grounded representations of these numeral tokens. We establish the superiority of our proposed techniques through evaluation on a range of numeracy tasks for both in-domain (seen) and out-domain (unseen) numerals. Further, we expand our empirical evaluations to numerals ranging from 1 to 10 billion, a significantly broader range compared to previous studies of the same nature, and we demonstrate significant improvements in the mathematical grounding of our learned embeddings.¹

1 Introduction

Numeracy, at its core, is the comprehension of numbers, akin to the comprehension of words in literacy. The magnitude of a number is especially tied to its meaning (Dehaene et al., 1998); as such, in developmental psychology, children able to distinguish numbers based on their magnitudes are said to possess the concept of numbers (Piaget, 1952). In the context of NLP, because numbers often grant objectivity to language (Porter, 1996), language models that can comprehend numeric magnitude and scales allow for better inference (Naik et al., 2018), information extraction (Madaan et al., 2016), and data-to-text generation (Sharma et al., 2021, 2022a).

Numeric comprehension can indeed be induced in language models through explicit supervision

¹Our codebase with the data and pre-trained models are hosted at <https://github.com/Mandar-Sharma/Laying-Anchors>

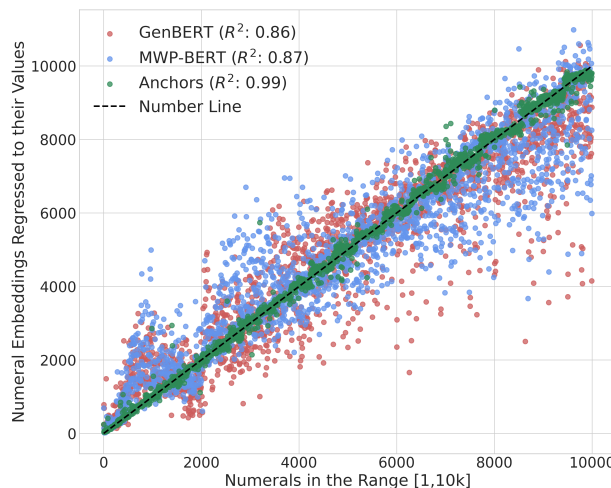


Figure 1: *Anchor-based embeddings correlate significantly better to the number line*: The plot above showcases how well the numeral embeddings from the baselines and our model (Anchors) correlate to the number line with their R^2 goodness-of-fit scores presented. The numeral range [1,10k] is employed for this plot as it contains a healthy mixture of both in-domain and out-domain numerals from our dataset.

(Vinyals et al., 2016); however, the inherent numeric capabilities of off-the-shelf language models induced from unsupervised training have been shown to be inadequate (Naik et al., 2018) and often fail to extrapolate to numerals *not seen* in the training set (Wallace et al., 2019; Razeghi et al., 2022) - referred to as *out-of-domain* (OOD) numerals. Approaches for numeracy induction to-date either involve strategies that learn representations for numerals separately from regular tokens (Spithourakis and Riedel, 2018; Jiang et al., 2020) or do so by training models on numeracy-specific tasks (Geva et al., 2020; Liang et al., 2022). In contrast, we *prime* (see §2) the numerals in the training corpus by laying anchors such that numeracy is induced via the unsupervised pre-training of the model itself without separately training numerical embeddings. As illustrated in Figure 1, our model shows substantial improvements in numeral representations for both numerals present in the training corpus (in-domain) as well as numerals

Laying Anchors

Training Sample: Kemp 's injury woes continued as he crashed into the outfield wall at Coors Field and injured both his knee and his shoulder . He continued to play despite the injuries . He finished the season batting 303 with 23 home runs and 69 RBIs in 106 games.

Training Sample (primed with **anchors**): Kemp 's injury woes continued as he crashed into the outfield wall at Coors Field and injured both his knee and his shoulder . He continued to play despite the injuries . He finished the season batting 303 **<ANC> 300** with 23 **<ANC> 23** home runs and 69 **<ANC> 68** RBIs in 106 **<ANC> 105** games.

Training Sample (primed with **directional anchors**): Kemp 's injury woes continued as he crashed into the outfield wall at Coors Field and injured both his knee and his shoulder . He continued to play despite the injuries . He finished the season batting 303 **<LA> 300** with 23 **<ANC> 23** home runs and 68 **<RA> 70** RBIs in 106 **<LA> 105** games.

Figure 2: *How are the numerals in the training corpus primed?* Showcasing samples from the training corpus - as-is, primed with simple anchors **<ANC>** where each numeral in the sample is augmented with the its closest anchor, and directional anchors **<LA>**/**<RA>** where the direction of the anchor with respect to the numeral (left or right in the number-line) is also embedded.

absent from the training corpus (out-domain) over the state-of-the-art baselines.

Further, the evaluation of numeracy in language models through their ability to predict numbers in a manner similar to textual tokens (Spithourakis and Riedel, 2018; Chen et al., 2019) omits the influence of rote-memorization (Zhang et al., 2020). In order to decouple the rote-memorization of numerals with respect to the linguistic context in which they appear, our study follows the evaluation protocols of Wallace et al. (2019) wherein the quality of learned representations are assessed through a set of numeric comprehension tasks. Our contributions can be summarized as:

- We develop new techniques for mathematical grounding of numerals in a corpus and quantitatively demonstrate significant improvements in model numeracy.
- We evaluate our models on a range of numerical tasks for numerals 1 to 10 billion (10^{10}), the largest analysis scope to the best of our knowledge, and evaluate its extrapolation capabilities to unseen (out-domain) numerals.
- Through rigorous evaluation, we demonstrate that the anchoring mechanisms lead to improved magnitude estimation (from *compressive representations*) and relative ordering (from *directional priming*) of numerals.

2 Priming Numerals with Anchors

How does one prime numerals? The *priming* effect is a temporary change in the perception of a target stimulus that frequently occurs in conjunction with a priming stimulus (Bargh and Chartrand, 2000). Similarly, semantic priming establishes the strength of relations among items belonging to the same or different categories (Zorzi et al., 2004).

Now, what does this mean in the context of numerals in a training corpus? Consider numerals 0

and 10 that are both equidistant to a supposed anchor numeral 5. If a language model has never seen the numerals 0 and 10 in its training corpus, the anchor numeral 5—that the model has seen during its training—can now be used to ground the magnitudes of these unseen numerals such that the model can now reason its magnitude. *Essentially, we intend to ground the magnitudes of numerals that the model rarely sees or has never seen based on the magnitudes of the numerals that it has frequently seen, known as the anchors.*

How are the anchors determined? First, we extract all numerals X from a training corpus C through which we intend to induce our anchors. The intuition that anchors should be numerals widely represented (frequent) in the corpus leads to the choice of Gaussian mixture models (GMMs) in contrast to clustering methods such as k-means that lack probabilistic cluster assignment. The set of anchors is induced from the means μ_k of each Gaussian $k \in K$ such that each numeral $n \in X$ can be tied to its closest anchor (1). Here, \mathcal{N} represents the probability density function and π_k, σ_k represent the mixing coefficient and standard deviation for the k -th Gaussian component. The initialization and the choice of K is described in §A.1.

$$p(n) = \sum_{k=1}^K \pi_k \mathcal{N}(n; \mu_k, \sigma_k^2) \quad (1)$$

Devising the four categories of anchors: Theories for mental representation of cardinality further divides our implementation of these anchors into two halves: a continuous linear representation (Dehaene, 2003) and a compressive representation where the difference between numerals n and $n + 1$ decreases as n increases (Dehaene et al., 1990). As such, for linear representation of the number line, we associate numerals with their closest anchor without alteration - giving us our first model *Anchors*. Similarly, for compressive representation, a given numeral n is anchored to m from a set of

log-normalized anchors such that $\ln(n) \approx m$ - our second model *ln Anchors*. In both these methods, the priming is implemented through a specialized token `<ANC>` added to the tokenizer.

Further, this priming effect is known to be symmetric with respect to the priming direction and additive to the effect of repetition priming (Reynvoet et al., 2002). This notion leads to our second category of models, viz. *directional anchors* represented with bi-directional arrows \rightleftarrows . Thus, in addition to attaching anchors to numerals in the corpus, we signify *where* the anchor lies in the number line with respect to the target numeral using specialized tokens `<LA>` (stating the anchor lies to the left of the target numeral in the number line) and `<RA>` (stating the anchor lies to the right of the target numeral in the number line). Training samples augmented with both `<ANC>` and `<LA>/<RA>` are depicted in Figure 2.

3 Experimentation and Results

As delineated in the previous section, we evaluate four configurations of our model pre-trained on the *anchor-augmented* WikiText-103 corpus (Merity et al., 2017): Anchors, *ln Anchors*, Anchors (\rightleftarrows), and *ln Anchors* (\rightleftarrows). The details of the datasets, pre-training and fine-tuning configurations, and embedding retrieval are described in §A.2.

3.1 Baselines

GenBERT (Geva et al., 2020): This model is based on the pre-trained BERT model and is additionally trained for quantitative reasoning (arithmetic, list minimum/maximum operations) with a corpus of 1 million synthetically generated quantitative reasoning prompts.

MWP-BERT (Liang et al., 2022): Also based on the pre-trained BERT model, MWP-BERT is trained for solving math word problems (MWP) through the injection of numerical properties via multiple numeracy grounded pre-training objectives that encourages contextual representations to capture numerical information.

3.2 Numeracy of Embeddings

In line with the premise set by Wallace et al. (2019), we evaluate the performance of the model embeddings on the tasks described below for different numerical ranges. The configurations for regressors and classifiers for the tasks mentioned below, are described in §A.3.

Decoding: Given embeddings for a set of numerals, the task is to regress them to their numerical values, thus assessing the fidelity of the numerical magnitudes captured by the embeddings.

Addition: Given sets of concatenated embeddings of two numerals, the task is to regress them to the numerical sum of the two numerals. In addition to assessing the magnitude fidelity, this task additionally requires number manipulation.

List Maximum-Minimum: While the first two tasks assess the magnitude captured by the embeddings, the task of predicting the maximum or minimum numeral in a set of randomly sampled numerals assesses whether the embeddings capture relative ordering.

4 Results

The results of above four tasks are illustrated in Table 1 for in-domain numerals, and similarly in Table 2 for out-of-domain numerals² (see §A.3). Our findings paint a consistent picture:

- For the lower numeral ranges $[1, 100]$ and $[100, 1k]$, all models do seemingly well. However, the performance of the baselines decreases sharply as the magnitude of numerals increase (for ranges $[1k, 10k]$ and $[10k, 10^{10}]$). However, *Anchors* and its variants have consistent performance across all the numeral ranges for both in-domain numerals and out-of-domain numerals.
- **Estimation of Numeral Magnitudes (I):** Within our models, the first notable phenomena we observe is that for the decoding and addition tasks designed to assess the fidelity of numerical magnitudes captured by the numeral embeddings, the logarithmic compression (*ln Anchors*) has a greater contribution to the model performance than directional anchors (Anchors (\rightleftarrows)).
- **Estimation of Numeral Magnitudes (II):** As the GMM-based anchors favor numerals frequent in the corpus, the anchors become sparse at higher numeral ranges - $[10k, 10^{10}]$. Thus, for this range specifically, we see that the model that strictly relies on directional

²Please note that as all numerals in range $[1, 100]$ and $[100, 1k]$ appear in the training corpus, only numeral ranges $[1k, 10k]$ and $[10k, 10^{10}]$ qualify for OOD evaluation.

Table 1: For in-domain numerals, Anchors consistently showcases enhanced numeracy across all numeral ranges while the baselines suffer significant degradation for larger numeral ranges: Performance of our model variants (Anchors) vs the baselines for in-domain numerals on four tasks evaluating the numeracy captured by model embeddings. The tasks are further sub-divided into number ranges and column $\forall Z \in C$ includes all numerals Z in corpus C .

Models	Decoding (Log-RMSE)					Addition (Log-RMSE)				
	[1,100]	[100, 1k]	[1k, 10k]	[10k, 10 ¹⁰]	$\forall Z \in C$	[1,100]	[100, 1k]	[1k, 10k]	[10k, 10 ¹⁰]	$\forall Z \in C$
GenBERT	0.0926	0.0301	0.0215	0.0639	0.0700	0.0250	0.0204	0.0237	0.0905	0.0752
MWP-BERT	0.0633	0.0213	0.0150	0.0540	0.0575	<u>0.0077</u>	0.0128	0.0200	0.0871	0.0533
Anchors	0.1279	0.0196	0.0074	0.0344	0.0424	0.0449	0.0172	0.0102	0.0442	0.0401
Anchors (\rightleftarrows)	0.1269	0.0123	0.0057	<u>0.0290</u>	0.0422	0.0180	0.0122	0.0089	<u>0.0426</u>	0.0378
<i>ln</i> Anchors	0.0279	0.0087	0.0049	0.0375	0.0304	0.0119	0.0067	0.0084	0.0572	0.0329
<i>ln</i> Anchors (\rightleftarrows)	0.1729	0.0109	0.0054	0.0375	0.0525	0.0157	0.0079	0.0106	0.0585	0.0443
	List Maximum (Accuracy)					List Minimum (Accuracy)				
	[1,100]	[100, 1k]	[1k, 10k]	[10k, 10 ¹⁰]	$\forall Z \in C$	[1,100]	[100, 1k]	[1k, 10k]	[10k, 10 ¹⁰]	$\forall Z \in C$
GenBERT	92.49%	91.49%	82.50%	82.50%	83.50%	94.99%	81.50%	83.50%	70.49%	86.00%
MWP-BERT	<u>93.00%</u>	91.50%	85.00%	79.00%	87.25%	<u>96.00%</u>	88.50%	88.50%	75.00%	87.00%
Anchors	92.50%	91.00%	63.00%	87.00%	87.75%	90.49%	88.99%	92.00%	86.00%	88.87%
Anchors (\rightleftarrows)	<u>93.00%</u>	83.00%	82.50%	83.00%	88.37%	92.50%	90.00%	86.50%	85.50%	91.00%
<i>ln</i> Anchors	92.00%	88.00%	88.50%	81.50%	89.37%	93.50%	92.00%	81.00%	85.00%	90.50%
<i>ln</i> Anchors (\rightleftarrows)	89.00%	<u>93.50%</u>	<u>90.50%</u>	88.00%	89.87%	94.00%	<u>93.50%</u>	<u>92.50%</u>	91.50%	92.50%

anchors outperforms the log-compressive anchors on magnitude estimation tasks. Essentially, when the anchors are further from each other, knowing which direction they reside in with respect to the target numeral aids the model in reasoning about that numeral.

- **Estimation of Relative Ordering:** The second phenomena we observe is that for the task of retrieving the maximum/minimum numeral from a list of numerals, designed to assess the relative ordering capabilities of the numeral embeddings, the model that leverages both compressive representations and directional priming (Reynvoet et al., 2002) (*ln* Anchors (\rightleftarrows)), has the best performance. Establishing that the incorporation of directional priming through the use of directional anchors further increases the relative ordering capabilities of the numeral embeddings.

For easier comparisons among models, the measure employed for the decoding and addition tasks is *log-RMSE*; as the error is log-compressed, seemingly small changes to the log-RMSE score translates to visible changes in numerical estimation through their embeddings, as depicted in Figure 1.

5 Conclusions

In this paper, we have presented a simple plug-and-play BERT variant with enhanced numerical capabilities. Through our rigorous interpolation (in-domain) and extrapolation (out-of-domain) analyses, we showcase the superiority of our model in

numeric comprehension while outlining the impact of logarithmic compression on magnitude estimation and the impact of directionality on relative ordering capabilities. Further, as a consequence of introducing anchors, we find the learning of niche pockets of similar embeddings for numerals closer in their magnitudes (§A.4).

6 Related Work

Although the majority of recent scholarly work in this domain revolves around training models to solve math problems (Wang et al., 2017; Nogueira et al., 2021; Liang et al., 2022) or strict arithmetic (Sharma et al., 2022b, 2023), several notable articles have looked exclusively into numeracy. Spithourakis and Riedel (2018) and Jiang et al. (2020) devise strategies with Gaussian mixture models to generate embeddings for out-of-vocabulary numeral tokens. Similarly, Razeghi et al. (2022) study the impact of numeral frequency in the pre-training corpus for few-shot arithmetic reasoning. Naik et al. (2018), Wallace et al. (2019), and Pal and Baral (2021) perform exploratory analysis of numeric comprehension through probing strategies.

Limitations

The restrictions from our in-house GPU resources do not allow scaling this study to more recent models that exceed 1 billion parameters. Nevertheless, recently published baselines that we evaluate against use the same underlying architecture that we employ, viz. the base BERT model. Given that larger models also depend on the base transformer

architecture (Vaswani et al., 2017) and use similar learning mechanisms, we believe that these observations will carry over to larger models as well.

Ethics Statement

The datasets we use in this study are established benchmark datasets from publicly accessible websites and do not contain any personally identifiable information. Our analyses does not constitute human subjects and thus do not fall within the purview of the IRB.

References

- John A. Bargh and Tanya L. Chartrand. 2000. The mind in the middle: A practical guide to priming and automaticity research.
- Johannes Blömer and Kathrin Bujna. 2013. Simple methods for initializing the em algorithm for gaussian mixture models. *CoRR*.
- Chung-Chi Chen, Hen-Hsen Huang, Hiroya Takamura, and Hsin-Hsi Chen. 2019. Numeracy-600k: Learning numeracy for detecting exaggerated information in market comments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6307–6313.
- Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794.
- Stanislas Dehaene. 2003. The neural basis of the weber–fechner law: a logarithmic mental number line. *Trends in cognitive sciences*, 7(4):145–147.
- Stanislas Dehaene, Ghislaine Dehaene-Lambertz, and Laurent Cohen. 1998. Abstract representations of numbers in the animal and human brain. *Trends in neurosciences*, 21(8):355–361.
- Stanislas Dehaene, Emmanuel Dupoux, and Jacques Mehler. 1990. Is numerical comparison digital? analogical and symbolic effects in two-digit number comparison. *Journal of experimental Psychology: Human Perception and performance*, 16(3):626.
- 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.
- Chengyue Jiang, Zhonglin Nian, Kaihao Guo, Shanbo Chu, Yinggong Zhao, Libin Shen, and Kewei Tu. 2020. Learning numeral embedding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2586–2599.
- Zhenwen Liang, Jipeng Zhang, Lei Wang, Wei Qin, Jie Shao, and Xiangliang Zhang. 2022. Mwp-bert: A numeracy-augmented pre-trained encoder for math word problems. *36th Conference on Neural Information Processing Systems (NeurIPS 2022) Workshop on Math-AI*.
- Aman Madaan, Ashish Mittal, Ganesh Ramakrishnan, Sunita Sarawagi, et al. 2016. Numerical relation extraction with minimal supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel mixture models. *International Conference on Learning Representations (ICLR)*.
- Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353.
- Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. 2021. Investigating the limitations of transformers with simple arithmetic tasks. *arXiv preprint arXiv:2102.13019*.
- Kuntal Kumar Pal and Chitta Baral. 2021. Investigating numeracy learning ability of a text-to-text transfer model. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3095–3101.
- Jean Piaget. 1952. *The Child’s Conception of Number*. London: Routledge and Kegan Paul.
- Theodore M Porter. 1996. *Trust in numbers*. Princeton University Press.
- Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. 2022. Impact of pretraining term frequencies on few-shot reasoning. *arXiv preprint arXiv:2202.07206*.
- Bert Reynvoet, Marc Brysbaert, and Wim Fias. 2002. Semantic priming in number naming. *The Quarterly Journal of Experimental Psychology: Section A*, 55(4):1127–1139.
- Mandar Sharma, John S Brownstein, and Naren Ramakrishnan. 2021. T 3: Domain-agnostic neural time-series narration. In *2021 IEEE International Conference on Data Mining (ICDM)*, pages 1324–1329. IEEE.
- Mandar Sharma, Ajay Gogineni, and Naren Ramakrishnan. 2022a. Innovations in neural data-to-text generation. *arXiv preprint arXiv:2207.12571*.
- Mandar Sharma, Nikhil Muralidhar, and Naren Ramakrishnan. 2022b. Overcoming barriers to skill injection in language modeling: Case study in arithmetic. *36th Conference on Neural Information Processing Systems (NeurIPS 2022) Workshop on Math-AI*.

Mandar Sharma, Nikhil Muralidhar, and Naren Ramakrishnan. 2023. [Learning non-linguistic skills without sacrificing linguistic proficiency](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6178–6191, 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.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 2016. Order matters: Sequence to sequence for sets. In *4th International Conference on Learning Representations, ICLR 2016*.

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*, pages 5307–5315.

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.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Michael C. Mozer, and Yoram Singer. 2020. Identity crisis: Memorization and generalization under extreme overparameterization. In *8th International Conference on Learning Representations, ICLR 2020*.

Marco Zorzi, Ivilin Peev Stoianov, and Carlo Umiltà. 2004. Computational modeling of numerical cognition.

A Appendix

A.1 Gaussian Mixture Models Initialization and Parameters

As Gaussian mixture models are sensitive to initialization methods (Blömer and Bujna, 2013), we initialize our models with random sampling from the dataset. The heterogeneous nature of the numeral distribution in the dataset lends this as the optimal initialization strategy. The models are trained to a

convergence tolerance of 0.001 with each component given its own general covariance matrix. The choice of $K = 1000$ Gaussian components was established stabilizing AIC and BIC values through a parameter sweep with K ranging from 10 to 5000.

A.2 Experimental Setup

A.2.1 Training Corpus

The WikiText-103 corpus (Merity et al., 2017) consists of 611,725 training instances (that includes over 100 million tokens) extracted from the set of verified *good* and *featured* articles on Wikipedia. Numeral tokens account for 2.4% of the corpus tokens with quadruple-digit numbers accounting for the greatest concentration of numerals - 41.8% .

A.2.2 Training Configurations

For both our baselines GenBERT (Geva et al., 2020) and MWP-BERT (Liang et al., 2022), the pre-trained models that the authors have provided are used as-is, thus ensuring no performance degradation as a consequence of in-house training/replication. For our Anchor models, the scheme for training follows BERT’s standard training protocol of using masked-language modeling. However, instead of randomly masking 15% of the tokens as done in BERT, we mask the anchor numeral as we intend to ground the learning of the target numerals based on their anchors. With the standard sequence size of 512 for BERT, the models were trained for 6 epochs each in a cluster of 4 Tesla P100 GPUs. The pre-trained BERT models are loaded from the Huggingface library (Wolf et al., 2019).

A.2.3 Embedding Retrieval

As recommended in the original BERT configuration, we tested hidden representations from the last hidden layer as well as from the sum of the last 4 hidden layers. We observed the best performance using a sum of the last 4 hidden layer representations, which we adopt for our experimentation.

A.2.4 Regressors and Classifiers

For consistency in our experimental results, we opted for Extreme Gradient Boosting (XGBoost) (Chen and Guestrin, 2016) for regression over standard neural networks for their robustness to parameterization. The regressors were initialized with 1000 components with each tree having a maximum depth of 5 and trained with a learning rate

Table 2: *Anchors generalize much better to unseen OOD numerals*: Performance of our model variants (Anchors) vs the baselines for out-of-domain numerals on four tasks evaluating the numeracy captured by model embeddings. The tasks are further sub-divided into number ranges and column $\forall Z \in C$ includes all numerals Z in corpus C .

Models	Decoding (Log-RMSE)		Addition (Log-RMSE)		
	Range	OOD [1k, 10k]	OOD [10k, 10 ¹⁰]	OOD [1k, 10k]	OOD [10k, 10 ¹⁰]
GenBERT		0.0132	0.0602	0.0130	0.0922
MWP-BERT		0.0097	0.0537	0.1205	0.0788
Anchors		0.0059	0.0328	0.0082	0.0419
Anchors (\Leftrightarrow)		0.0043	0.0278	0.0067	0.0409
\ln Anchors		0.0033	0.0338	0.0043	0.0557
\ln Anchors (\Leftrightarrow)		0.0029	0.0347	0.0033	0.0625
		List Maximum (Accuracy)		List Minimum (Accuracy)	
		OOD [1k, 10k]	OOD [10k, 10 ¹⁰]	OOD [1k, 10k]	OOD [10k, 10 ¹⁰]
GenBERT		86.50%	78.49%	90.00%	76.00%
MWP-BERT		87.00%	82.50%	88.50%	77.00%
Anchors		84.50%	83.50%	89.49%	83.50%
Anchors (\Leftrightarrow)		86.00%	88.50%	90.00%	81.50%
\ln Anchors		87.50%	86.99%	90.00%	83.50%
\ln Anchors (\Leftrightarrow)		88.00%	87.00%	91.50%	84.00%

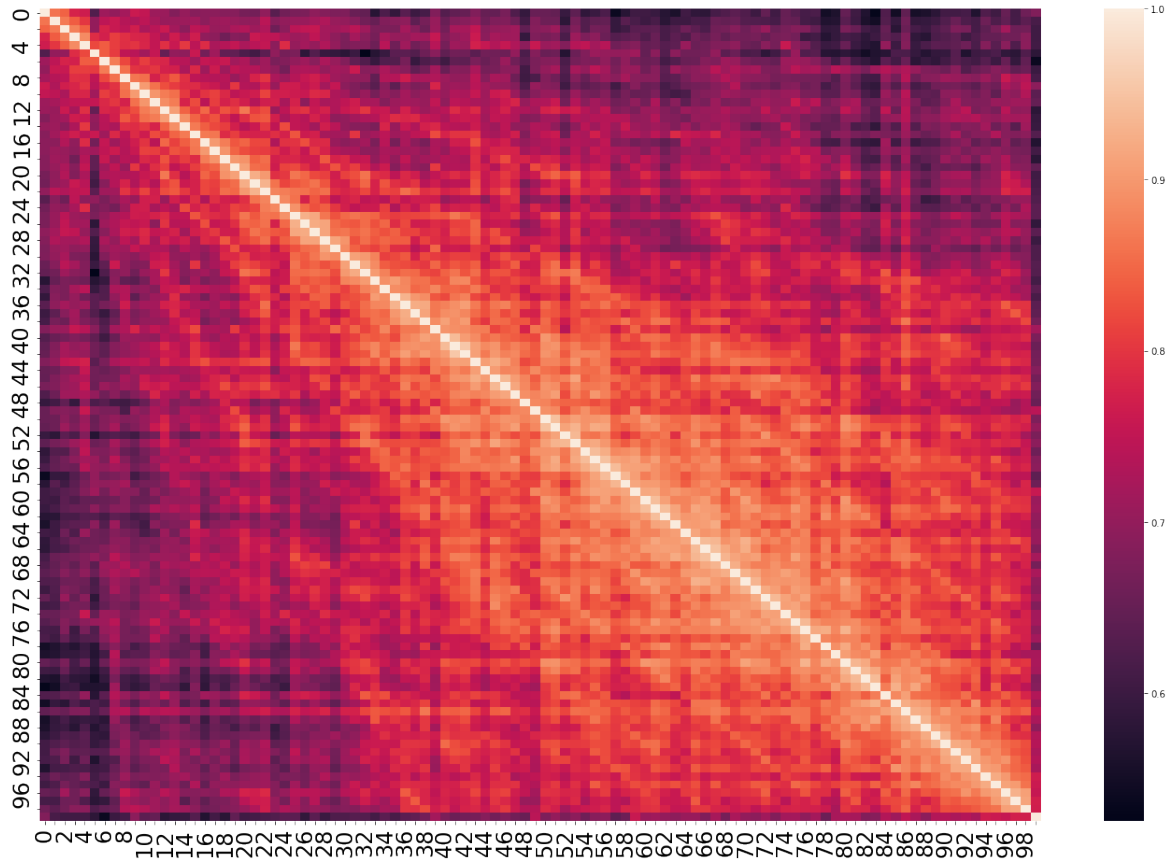
of 0.01. Similarly, a standard LSTM setup with 4 stacked LSTMs coupled with a sigmoid activation for the final linear layer was used as the classifier. Each classifier was trained for 150 epochs with a learning rate of 1e-4.

A.3 Extrapolation for Out-domain Numerals

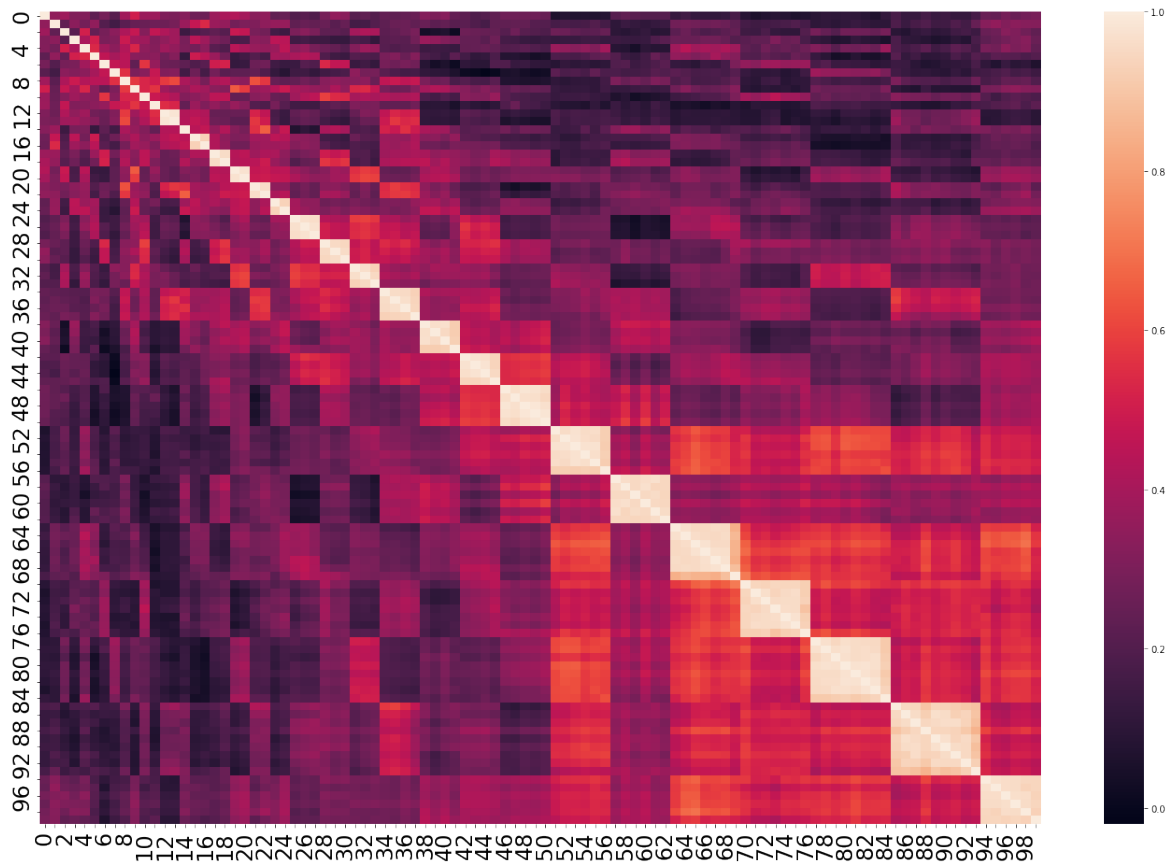
As depicted in Table 1 for in-domain numerals, we perform the same set of evaluations for out-of-domain (unseen) numerals in Table 2, corroborating the same performance gains that we observed for in-domain numerals. Please note that all numerals in range [1,100] and [100, 1k] appear in the training corpus, thus only the ranges [1k, 10k] and [10k, 10¹⁰] qualify for OOD evaluation.

A.4 Embedding Visualizations

As an alternative visualization tool, we contrast heatmaps generated through the cosine similarities of numeral embeddings for the base BERT model and our model. As illustrated in Figure 3, the heatmap for the base BERT model has uniformly low cosine similarity throughout, leading to little distinction between numeral embeddings. In contrast, the heatmap for our model demonstrates sophisticated patterns of similarity for proximal numerals along its diagonal. Also seen are sections of low similarity scores in the top right and bottom left - indicating the ability to discern numerical magnitudes of lower and higher number ranges.



(a) Base BERT model



(b) Our model

Figure 3: Heatmaps computed from cosine similarities of numeral embeddings in range [1,100].