



Adaptive Activation Functions

Does choice of activation function matter in smaller Language Models?

Filip Ignjić

Supervisor(s): Aral de Moor, Responsible Professors: Maliheh Izadi, Arie van Deursen

EEMCS, Delft University of Technology, The Netherlands

A Thesis Submitted to EEMCS Faculty Delft University of Technology,
In Partial Fulfilment of the Requirements
For the Bachelor of Computer Science and Engineering
June 19, 2024

Name of the student: Filip Ignjić
Final project course: CSE3000 Research Project
Thesis committee: Thomas Abeel, Maliheh Izadi, Arie van Deursen, Aral de Moor

An electronic version of this thesis is available at <http://repository.tudelft.nl/>.

Abstract

The rapid expansion of large language models (LLMs) driven by the transformer architecture has introduced concerns about the lack of high-quality training data. This study investigates the role of activation functions in smaller-scale language models, specifically those with approximately 10million (M) parameters, to ensure sustained progress in LLM development despite data limitations. Activation functions, crucial for neural network performance, have evolved significantly, but comprehensive comparisons under consistent conditions remain scarce, especially for smaller parameter count models. This research systematically evaluates traditional and novel activation functions, including learnable variants, and introduces the Kolmogorov-Arnold Network (KAN) to language modeling. Using Hugging Face implementations of GPT-Neo and RoBERTa models, this study assesses performance impacts through the BabyLM evaluation pipeline. TODO: Add results and conclusions.

1 Introduction

The transformer architecture [21] has revolutionized the AI landscape and enabled the development of commercial large language models (LLMs) like ChatGPT. However, as these models continue to grow in size, current trends forecast running out of high-quality data required for optimal performance [22]. This limitation stimulates the initiative to improve sample efficiency motivated by the observation that LLMs are exposed to orders of magnitude more information than a human in their lifetime, yet are certainly not leveraging all this information nearly as efficiently. Therefore, this research aims to investigate the impact of architectural decisions on smaller models which are often left neglected. By understanding how to optimize smaller models, we can ensure the progress of LLMs even with the projected lack of high-quality data. Furthermore, smaller models can be deployed on edge devices can insure privacy and be used for spellcheck, predictive typing, conversational assistance etc.

The activation function in a neural network, determines if neuron should be activated or not. The choice of activation functions has historically played a crucial role in the advancement of neural networks, such as the shift from Sigmoid activation functions to RELU (Rectified Linear Unit), which simplified computation, improved feature learning and mitigated vanishing gradient problem [15]. As the era of LLMs unfolded, the development of activation functions continued to evolve, resulting in over 400 documented activation functions [12]. Despite this progress, literature typically compares new activation functions against their immediate predecessors and usually using all the latest state-of-the-art configuration (bigger and bigger models) leading to a gap in the literature: there is no comprehensive comparison of multiple activation functions under consistent model sizes, furthermore, with the trend of increasing model sizes, the investigation of the impact activation functions have on smaller models is neglected.

This research aims to address the aforementioned gap by exploring the impact of existing and novel activation functions on smaller-scale language models with approximately 10 million (10M) parameters.

For purposes of the research, we will modify Hugging face implementations of GPTNEO [7] and RoBERTA [8] with selected existing activation functions and a novel Learnable GELU activation function, set the hyperparameters to amount to a total size of approximately 10M trainable parameters and evaluate them on the BabyLM evaluation pipeline [25]. More details on the setup provided in section 4.

This paper provides the evolution and background of activation functions in section 2 and the approach of the investigation in section 3. It describes the experimental setup using in section 4. The results are presented and analysed in section 5, and the discussion addresses performance outcomes and future research directions in section 6. Conclusion and responsible research considerations are presented in sections 7 and 8, respectively.

2 Background and related work

Transformers comprise multiple layers, each crucial in processing input data and generating meaningful representations. Among these layers, the Feed-Forward Neural Network (FFN) layer typically consists of two linear transformations with an activation function in between, effectively forming a simple Multi-Layer Perceptron (MLP) [9].

Activation functions are used to introduce non-linearity into neural networks, allowing them to model complex relationships in the data. They are applied to the nodes of the network and are essential for enabling the network to learn and perform a wide range of tasks, beyond what linear models can achieve [4].

The famous paper “Attention is All You Need” [21] used the state-of-the-art activation function at the time, Rectified Linear Unit RELU. Since then, no significant improvements were made until the introduction of Gaussian Error Linear Unit GELU, which quickly became the default activation function in most of the state-of-the-art LLMs. The popularity of GELU stems from its ability to enhance model performance without introducing an efficiency overhead by combining the properties of RELU and SIGMOID function. It balances between linearity and non linearity, allowing it to handle both positive and negative inputs gracefully [10]. Despite these advancements, continuous innovation leads to alternatives like Gated GELU GEGLU, noted for its effectiveness [20], also used in last year’s winner of the BabyLM challenge [18]. However, a recently published paper suggests a return to ReLU [14], adding further confusion to the search for the optimal activation function. Fortunately, it also provides some clues that guide further exploration and motivate this research.

The paper suggests that suggests that the impact of activation functions diminishes as the model size increases, evident in models with over a billion parameters [14]. This also explains the initial move away from RELU, since the research on activation alternatives was done on models with the sizes of 100+ million parameters [20] [10]. Highlighting this find-

ing further motivates the need to investigate activation functions in smaller models. The impact of activation functions is expected to be more significant in smaller models, and until now, decisions have been made with the trend of increasing model size in mind.

Furthermore, another gap appears in the research on activation functions with trainable parameters (adaptable activation functions). The possible explanation for this could be the trade-off between additional trainable parameters and performance. However, this was primarily studied in larger models. To our knowledge the only paper on adaptive activation functions in LLMs is by Rajanand et al. [16], which found that adaptive activation functions outperform static ones in text-to-text machine translation, a task also performed by language models. This suggests that adaptive activation functions could be beneficial for smaller models, but further research is needed to confirm this hypothesis. Given these insights, this research will explore the impact of various activation functions on smaller-scale language models with around 10 million parameters. Hypothesizing that at smaller scales, the choice of activation function is crucial, having learnable parameters could be beneficial, while the added parameters from the function’s parametrization remain relatively minimal compared to the total model size.

Kolmogorov-Arnold Networks (KAN) represent a recent development in neural network architecture, where activation functions are applied on edges instead of nodes [13]. This approach has been shown to outperform traditional neural networks in some tasks, particularly in scientific applications such as solving partial differential equations. However, at the time of this study, it has yet to be tested on language models. The primary benefit of KAN is the optimization of activation on each edge using splines. A spline is a piece-wise-defined polynomial function used in interpolation and approximation to create smooth curves through a set of points [1]. With a spline on each edge (see figure 1), each edge can have its own custom activation function, trained separately and uniquely shaped. In contrast, adaptive activation functions have the same shape but different gradients. However, this comes with the drawback of an increased number of trainable parameters. This research will experiment with applying KAN to language modeling to assess its efficacy at smaller scales, addressing the gap in the current literature.

3 Approach

The transformer structure (as introduced in section 2) allows the default activation function to be switched out with different activation functions for testing, enabling a direct comparison of their performance while keeping the rest of the architecture the same. To explore the effectiveness of various activation functions, we will modify the existing Hugging face implementations of GPTNEO [7] and ROBERTA [8].

3.1 Choosing activation functions

The activation functions evaluated were Rectified Linear Unit (ReLU) ¹, Sigmoid Linear Unit (SiLU) [10], Gated

SiLU with learnable parameters (SWISH) [5], Parametric ReLU (PReLU) ², Gaussian Error Linear Unit (GELU) (baseline models), ADAPTABLE GELU and NO ACTIVATION FUNCTION. Additionally, the Kolmogorov-Arnold Networks KAN NETWORK [13] will be compared against all of these options.

GELU

Currently, the most popular activation function in LLMs is also used as the default activation function in the baseline models GPT-NEO and ROBERTA. It is a smooth approximation of ReLU, originally defined as $\text{GELU}(x) = x \cdot \Phi(x)$, where $\Phi(x)$ is the Cumulative Distribution Function for the Gaussian Distribution. For optimization purposes, since calculating $\Phi(x)$ is computationally expensive, it is instead calculated with the TANH approximation as:

$$\text{GELU}(x) = 0.5x \left(1 + \tanh \left(\sqrt{\frac{2}{\pi}} (x + 0.044715x^3) \right) \right) \quad [10].$$

This function will serve as the baseline for comparing all other activation functions, with a particular focus on evaluating it against the novel Learnable GELU.

ReLU and PReLU

ReLU was considered state-of-the-art at the time of the original transformer paper [21], but has since been surpassed by other activation functions.

$$\text{ReLU}(x) = \max(0, x)$$

ReLU will be used for comparison with PReLU and GELU. The comparison with GELU is motivated by the findings of *I. Mirzadeh et al.* [14], which suggest that the use of ReLU is acceptable as the impact of activation functions diminishes with increasing model size. The objective is to assess the extent to which ReLU underperforms compared to GELU when applied to smaller models.

PReLU is a variant of ReLU that incorporates learnable parameters, allowing the activation function to adaptably learn the optimal slope for negative values.

$$\text{PReLU}(x) = \max(0, x) + a \min(0, x)$$

where a is a learnable parameter. This introduces x learnable parameters where x is the FFNs’ intermediate size in the transformer. The objective is to evaluate whether adding a learnable parameter to ReLU can enhance performance or increase the training time.

SiLU and Swish

SiLU was evaluated on LLMs in the original GELU paper [10] but was found to perform worse than the GELU activation function. It is defined as

$$\text{SiLU}(x) = x \cdot \sigma(x), \text{ where } \sigma(x) \text{ is the logistic sigmoid.}$$

It will be used only as a baseline comparison for the *Swish* activation function, which is its adaptable counterpart with learnable parameters, aiding the objective of exploring the impact of adding learnable parameters to activation functions.

¹PyTorch Library, "torch.nn.ReLU," *PyTorch Documentation*, <https://pytorch.org/docs/stable/generated/torch.nn.ReLU.html>

²PyTorch Library, *torch.nn.PReLU*, *PyTorch Documentation*, <https://pytorch.org/docs/stable/generated/torch.nn.PReLU.html>

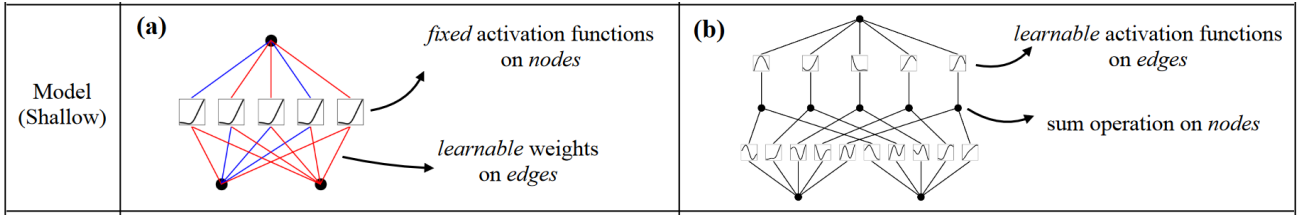


Figure 1: KAN vs MLP Z. Liu et al., “KAN: Kolmogorov-Arnold Networks.” arXiv, May 02, 2024. doi: 10.48550/arXiv.2404.19756.

Swish is a self-gated activation function that was proposed by Ramachandran et al. [17]. It was implemented as proposed in a paper

$$\text{swish}(x) = x \cdot \text{silu}(\alpha \cdot x)$$

where α is a learnable parameter. The objective is to evaluate whether the SWISH activation function outperforms the non adaptable SiLU to evaluate the impact of adding learnable parameters to activation functions.

Adaptable GELU

The GELU activation function will be parameterized with learnable parameters. The implementation will be based on the PyTorch GELU implementation[cite] with TANH approximation, which, out of the box, does not support learnable parameters. This approach appears to be novel, as no prior research has been found that explores the impact of adding learnable parameters to the GELU activation function. The new GELU activation function adds a learnable parameter as a scaling factor and is defined as follows:

$$0.5 \cdot \alpha \cdot x \cdot \left(1.0 + \tanh \left(\sqrt{\frac{2.0}{\pi}} \cdot (\alpha \cdot x + 0.044715 \cdot x^3) \right) \right)$$

where α is a learnable parameter.

This activation functions adds a total of 2048 learnable parameters. The objective is to evaluate whether the ADAPTABLE GELU activation function outperforms the standard GELU activation function to evaluate the impact of adding learnable parameters to activation functions.

No Activation Function

The NO ACTIVATION FUNCTION will be used as a baseline for comparison with all other activation functions. This baseline will help determine the impact of using an activation function compared to not using one.

KAN-Network

The KAN NETWORK is a novel activation function that has shown promising results in the literature. The implementation of the KAN NETWORK is available by original authors, but has shown to be problematic and slow during our research. Those issues were addressed in efficient-KAN³ implementation, which was the most popular optimization on GitHub of the original KAN paper at the time of this study.

³Blealtan, *Efficient-KAN, GitHub Repository*, 2024. Available at: <https://github.com/Blealtan/efficient-kan> (Accessed: 2024-06-03).

This approach requires careful selection of certain parameters, specifically the number of grid intervals and the order of piecewise polynomials. Based on recommendations from the paper, we will set these parameters to 3 and 5, respectively. The configuration aims to balance performance and computational efficiency, ensuring that the KAN implementation is both effective and practical for the experiments. FFN layers from GPT-NEO and roBERTa both use MLPs in their implementation, which can be directly replaced with efficient-KAN implementation. As KAN network increases the amount of total trainable parameters in FFN from 2M to slightly over 8M, the intermediate size for models with KAN NETWORK was decreased to 256 to keep the number of parameters in FFN around 2M. The objective is to evaluate the performance of the KAN network compared to the other activation functions and to determine whether it is a viable alternative for MLPs in LLMs. Additionally, the KAN NETWORK offers the unique capability to extract symbolic representations of the learned activation functions from splines. It also provides the ability to retrain or further train specific parts of the network if necessary. Unfortunately, due to time constraints and the current implementation, this aspect will not be explored in this study. For further discussion, see Section 6.

4 Experimental setup

4.1 Research questions

Through the experiment, we aimed to answer the following research questions:

- *Is the choice of activation function relevant to the performance of smaller models with 10M parameters?* We compared the baseline models (with GELU activation) and models with *ReLU* activation function.
- *How does the addition of learnable parameters to the activation function improve the performance of the model?* We modified static activation functions GELU, SiLU and RELU to include learnable parameters.
- *Do FFNs using KAN-networks outperform FFNs using MLP networks?* We used efficient-kan implementation³ of KAN-networks.

All of the setups above were evaluated on BabyLM evaluation pipeline [25].

4.2 Models^{4 5}

We used Hugging face implementations of GPTNeoForCausalLM [7] and RobertaForMaskedLM [8] models. Gpt-NEO is GPT2-based decoder model, while roBERTa is based on Google’s BERT model from 2018. We chose these models as they are well-established without and without the newer optimizations found in more recent other architectures, providing simplicity that is more suitable the scope of this project. Additionally, testing on encoder and decoder architectures increases the generalizability the research. GPT-Neo and Roberta also come in pre-trained form, however as we do not have the resources to train them to that extent, we pre-train on our own, smaller dataset so we have a baseline to compare to. The hyperparameters used for the models can be seen in table 1. The change in intermediate size for KAN-MLP implementation was made due to increased number of learnable parameters in the activation function and the need to keep the total number of parameters in the model approximately constant. The final parameter count for models using KAN was 11 million, whereas the parameter count for the other models was 9 million.

4.3 Activation functions

RELU and SiLU are provided by used hugging face implementations and can be set via configuration file. The rest of the activation functions was implemented as described in Section 3 and are available in the replication package. [todo: cite replication]

4.4 Dataset⁵

We used the TinyStoreis [6] dataset for pre-training. It is a dataset of short stories, with the total of 2.1 million samples with an average of 175.4 words per story, containing words that a typical 4-year-old would likely understand, generated by GPT-3.5 and GPT-4. We used them as it has been shown that training on reduced dataset complexity exhibits better natural language understanding than GPT-2 (125M) at the fraction of training cost.

4.5 Evaluation Setting and Metrics

Evaluation pipeline^{4 5}

We used the BabyLM evaluation pipeline [25] to evaluate the models. The pipeline consists of three components: BLiMP, GLUE and SuperGLUE. BLiMP is a benchmark used to evaluate the linguistic capabilities of language models. It consists of pairs of minimally different sentences, with one sentence in the pair containing a grammatical error. Performance is measured by the model’s ability to correctly assign a higher likelihood to the grammatically correct sentence. [26] [27]. GLUE and SuperGLUE are benchmarks that assess the performance of language models across a range of tasks focused on text understanding and reasoning. The BabyLM pipeline

⁴ The research group supervisor implemented the pre-training and evaluation for the baseline models, which were then individually modified by the research group members.

⁵ These sections were written based on the guidelines provided by the research group supervisor and discussions with the other group members.

integrates select tasks from both GLUE and SuperGLUE [24] [23].

Statistical significance testing

To ensure the results are statistically significant, we trained each model 6 times with different seeds. The following section justifies the choices made for the statistical significance testing, based and summarised from The Hitchhiker’s Guide to Testing Statistical Significance in Natural Language Processing [3].

Since the distribution of the test statistic is not known, we had to choose approaches from non parametric family of approaches. Which is split in two categories. Sampling based tests and sampling free tests. Sampling free tests hold less statistical power, therefore we decided to use sampling based tests.

Sampling based tests compensate for the lack of distribution information with resampling, which makes them computationally more expensive. We used bootstrapping as a sampling based test. It is a resampling method that involves drawing samples with replacement from the original results. These samples are used to approximate the distribution of the statistic. We use that to calculate the means and 95% confidence interval of BLiMP and GLUE scores for each of the models. We used 10 000 samples for bootstrapping in our experiments. The exact implementation used, available in the replication package [TODO: cite replication].

4.6 Hardware

All the models were trained and evaluated on a single NVIDIA A100 GPU with 4 CPUs and 24GB of memory on DelftBlue cluster.

5 Results

ToDo: Fix the postition of all those tables bellow. All results of evaluations on BLiMP and GLUE datasets are shown in Table ??.

After performing bootstrap resampling for 10 000 samples, we calculated the mean difference and confidence interval to make each of the comparisons mentioned in section 4.

Comparison of Baseline (GELU) model with Adaptive GELU model showed that the mean difference is negative for both BLiMP and GLUE benchmarks, showing that the Adaptive GELU model performed slightly better than the Baseline (GELU) model. However the confidence intervals for both benchmarks contain zero, which means that the difference is not statistically significant. The results are shown in Table ??.

Bootstrap comparing static activations ReLU and SiLU with their adaptive counterparts PReLU and Swish showed that the mean difference is positive for both BLiMP and GLUE benchmarks, showing that the static activation functions performed slightly better than the adaptive activation functions. However the confidence intervals for both benchmarks contain zero, which means that the difference is not statistically significant. The results are shown in Table ??.

Parameter	GPT Neo	RoBERTa
Embedding Parameters		
Vocab Size	10,000	10,000
Hidden Size	512	512
Max Position Embeddings	512	513
Blocks (Attention & FFN)		
Number of Layers	2	2
Attention Types	global and local	global
Number of Attention Heads	4	4
Window Size	256	N/A
Intermediate Size	1024 (256 for KAN-MLP)	1024 (256 for KAN-MLP)

Table 1: Comparison of Parameters for GPT Neo and RoBERTa

The highest GLUE score was achieved by the BERT-KAN model with 63.65, the highest BLiMP score was achieved by the GPT-KAN model with 63.43. **ToDo: Further discuss after all statistical tests are done.**

Across all models, the training times were relatively consistent, with two outliers: KAN network models and ReLU models. The prolonged training time for the ReLU models might be attributed to the busy DelftBlue nodes, whereas the extended time for the KAN models is expected. We will compute the standard deviations for all models once we obtain all the results. Excluding the two outliers, the standard deviation in training times was 3 minutes and 33 seconds.

6 Discussion

6.1 Implications

Our results show that choice of activation functions remains irrelevant even at smaller scales. All the comparisons in section 5 show that the differences are not statistically significant. This suggests that the choice of activation function does not have a significant impact on the performance of language models with 10 million parameters. This is in contrast to the findings of Mirzadeh et al. [14], who suggested that the impact of activation functions diminishes as the model size increases. Our results suggest that the impact of activation functions is negligible even at smaller scales. Furthermore, the training durations of models employing various activation functions showed minimal differences, indicating that the selection of one activation function over another does not significantly impact training efficiency.

The comparison of results between GELU and the novel Learnable GELU further corroborates the findings of Dror et al. [3], which emphasize the importance of statistical significance analysis in evaluating various architectural decisions in language models. Our initial findings indicated an improved performance of Learnable GELU, but subsequent analysis revealed these improvements to be statistically insignificant.

ToDo: Add more specific implications of other comparisons.

Models with KAN-Networks do seem to perform better than models with traditional MLPs, as shown by the highest scores achieved by the BERT-KAN and GPT-KAN models. This suggests that the KAN architecture could be beneficial for language models, even at smaller scales. However, the increased number of trainable parameters and increased training times is a significant drawback that needs to be considered.

6.2 Threats to validity

Internal validity examines certainty that the observed results are due to independent variable and not other factors. External validity concerns the extent to which the findings of the study can be generalized to contexts outside the study. Con-

Table 2: Generated by Spread-Latex

Model	Glue Mean	95% CI	Blimp Mean	95% CI Blimp	Avg. Train. Time
BERT Baseline (GELU)	56,04	[52.35, 58.74]	58,48	[56.68, 59.5]	1h 41m 27s \pm 32s
BERT Learnable GELU	57,69	[56.75, 58.73]	59,25	[58.78, 59.65]	1h 47m 56s \pm 313s
GPT Baseline (GELU)	58,95	[58.27, 59.7]	57,82	[56.97, 58.68]	1h 42m 41s \pm 103s
GPT Learnable GELU	59,3	[58.92, 59.73]	57,24	[56.47, 58.03]	1h 41m 59s \pm 102s
BERT ReLU	57,91	[56.98, 58.87]	58,4	[57.37, 59.23]	2h 25m 23s \pm 147s
BERT PReLU	N/A	N/A	N/A	N/A	N/A
GPT ReLU	59,41	[59.0, 59.82]	56,91	[55.93, 57.8]	2h 20m 27s \pm 203s
GPT PReLU	N/A	N/A	N/A	N/A	N/A
BERT SiLU	57,54	[57.15, 58.1]	59,16	[58.75, 59.53]	1h 43m 50s \pm 268s
BERT Swish	N/A	N/A	N/A	N/A	N/A
GPT SiLU	59,05	[58.4, 59.63]	57,25	[56.43, 58.2]	1h 39m 7s \pm 270s
GPT Swish	N/A	N/A	N/A	N/A	N/A
GPT KAN	N/A	N/A	N/A	N/A	N/A

struct validity refers to the extent to which the the mearument tools are appropriate for the study.

Internal validity

To ensure that our results were not due to chance, we trained models using six different random seeds. While this sample size is relatively small, it was a necessary compromise given our resource constraints. To address this limitation, we employed bootstrapping with 10,000 samples to calculate confidence intervals and mean differences between models. Additionally, to verify that our findings were not model-specific, we utilized two distinct types of models: an encoder (roBERTa) and a decoder (GPT-Neo). Nevertheless, it is important to note that all models were trained on the same dataset, which may have introduced some bias. The train times could have been affected by busy delft blue nodes, but across models trained on multiple seeds, the training times were relatively consistent with standard deviation of only 9 minutes and 52 seconds accorss all models including the ones with KAN - Network.

Another potential threat to validity is that the implementation of adaptable activation functions added 2048 learnable parameters to each model, which could have influenced the results. However, given that 2048 parameters constitute a relatively small proportion in models with 9M, and the results were not statistically significant, this impact is likely minimal.

External validity

All the models were trained for only one epoch due to computational constraints, resulting in none of the models fully converging. This may have impacted the performance of the models and the significance of the results, as differences might be more pronounced or minimized with additional epochs. In a real-world scenario, models would typically be trained for more epochs. Additionally, the models utilized in this study are not the latest state-of-the-art, which may affect the generalizability of the results. Furthermore, the TinyStories dataset used for pre-training, which comprises only short fictional stories, may not be representative of the datasets typically used for production-level language models. **ToDo: talk about kan implementation**

Construct validity

The evaluation pipeline used in this study, BabyLM, primarily focuses on grammatical tasks, which may not fully capture the comprehensive capabilities of language models. However, considering the scope of this research, the tasks evaluated by the BabyLM evaluation pipeline are suitable. Since BabyLM is specifically designed for the evaluation of smaller language models, it aligns well with our study's focus. Thus, it is an appropriate tool for assessing the impact of activation functions on these models.

6.3 Future work

I hope this research will inspire further studies to give accitivation functions less attention. Before publishing new activation functions, researchers should conduct multiple runs and perform statistical significance tests to ensure the robustness of their findings.

The KAN architecture is a relatively new concept, and the implementation used in this study may not have been optimal. Although we utilized more efficient and optimized implementation as suggested by the original paper [13], the training process was less stable compared to baseline models (see Figure 2), and the models did not converge as expected. This instability may have impacted the performance of the KAN MODELS, suggesting that future research should investigate this issue further. Despite its potential, this promising direction concerning activation functions and its interpretability benefits were not fully explored in this study due to time constraints.

7 Conclusions and Future Work

TODO:

- briefly repeat the RQs
- metion the the implications of results and tie them back to the RQs

8 Reponsibile research

To prevent test set contamination, we pre-trained our models on datasets that were separate from those used for evaluation. This addresses an issue highlighted in recent works,

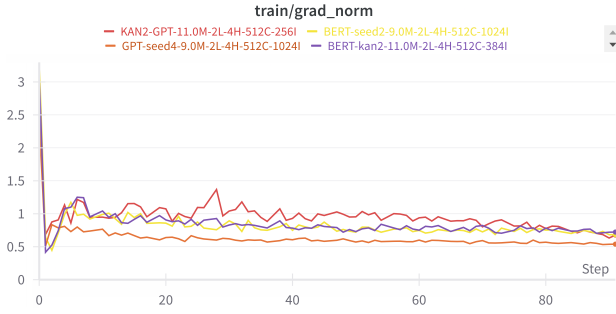


Figure 2: Gradient normalization during training of KAN models. The KAN models show higher variance in gradient norms compared to the baseline models. **ToDo: Fix readability**

where pre-training on test sets can artificially inflate performance metrics and call into question the validity of the results [19]. Ensuring distinct separation between training and evaluation datasets maintains the integrity of our findings and contributes to the reliability of our research.

In conducting this research, we ensured transparency and reproducibility by sharing the experimental setup under section 4. The datasets and models used are publicly available and can be found in the references. Furthermore, all the implemented code is available on GitHub. We adhered to scientific integrity by fabrication, and plagiarism, ensuring that everything reported accurately with proper citations.

Guided by the Netherlands Code of Conduct for Research Integrity, we incorporated principles of honesty, transparency, and responsibility, ensuring our research practices align with the highest standards. We followed the educational and normative framework from chapters 2 and 3 of the Code, emphasizing good research practices that promote a responsible research environment [11].

9 Acknowledgements

The research was conducted at TU Delft as the final thesis of the Bachelors of Computer Science and Engineering. We would like to thank for support of the project amazing project supervisor, supervising professors and the entire research research project group. Additionally, we thank the DHCP [2] for providing the necessary computational resources for the project.

References

- [1] Arindam Chaudhuri. B-splines. *arXiv:2108.06617*, 2021.
- [2] Delft High Performance Computing Centre (DHPC). DelftBlue Supercomputer (Phase 2). <https://www.tudelft.nl/dhpc/ark:/44463/DelftBluePhase2>, 2024.
- [3] Rotem Dror, Gil Baumer, Segev Shlomov, and Roi Reichart. The hitchhiker’s guide to testing statistical significance in natural language processing. In Iryna Gurevych and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1383–1392, Melbourne, Australia, Jul 2018. Association for Computational Linguistics, Association for Computational Linguistics.
- [4] Shiv Ram Dubey, Satish Kumar Singh, and B. B. Chaudhuri. Activation functions in deep learning: A comprehensive survey and benchmark. *arXiv preprint arXiv:2109.14545*, Jun 2022.
- [5] Steffen Eger, Paul Youssef, and Iryna Gurevych. Is it time to swish? comparing deep learning activation functions across NLP tasks. *arXiv:1901.02671*, 2019.
- [6] R. Eldan and Y. Li. Tinystories: How small can language models be and still speak coherent english? *arXiv 2305.07759*, May 2023. Accessed: Nov. 01, 2023.
- [7] Hugging Face. Gpt neo, 2024. Accessed: Jun. 03, 2024 https://huggingface.co/docs/transformers/en/model_doc/gpt_neo.
- [8] Hugging Face. Roberta, 2024. Accessed: Jun. 03, 2024 https://huggingface.co/docs/transformers/en/model_doc/roberta.
- [9] Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. *arXiv:2203.14680*, 2022.
- [10] D. Hendrycks and K. Gimpel. Gaussian error linear units (gelus). *arXiv*, arXiv:1606.08415, Jun. 2023.
- [11] KNAW, NFU, NWO, TO2-Federatie, Vereniging Hogescholen, and VSNU. Nederlandse gedragscode wetenschappelijke integriteit, 2018.
- [12] V. Kunc and J. Kléma. Three decades of activations: A comprehensive survey of 400 activation functions for neural networks. *arXiv*, arXiv:2402.09092, 2024.
- [13] Z. Liu and Others. Kan: Kolmogorov-arnold networks. *arXiv*, arXiv:2404.19756, May 2024.
- [14] I. Mirzadeh and Others. Relu strikes back: Exploiting activation sparsity in large language models. *arXiv*, arXiv:2310.04564, Oct 2023.
- [15] Vinod Nair and Geoffrey Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, volume 27, page 814, 2010.

- [16] A. Rajanand and P. Singh. Erfrelu: Adaptive activation function for deep neural network. *arXiv 2306.01822*, Jun 2023.
- [17] P. Ramachandran, B. Zoph, and Q. V. Le. Searching for activation functions. *arXiv*, arXiv:1710.05941, Oct. 2017.
- [18] D. Samuel, A. Kutuzov, L. Øvrelid, and E. Velldal. Trained on 100 million words and still in shape: Bert meets british national corpus. *arXiv*, arXiv:2303.09859, May 2023.
- [19] Rylan Schaeffer. Pretraining on the test set is all you need. *arXiv*, September 2023. doi: 10.48550/arXiv.2309.08632.
- [20] N. Shazeer. Glu variants improve transformer. *arXiv*, arXiv:2002.05202, Feb 2020.
- [21] A. Vaswani and Others. Attention is all you need. *arXiv*, arXiv:1706.03762v5, Dec 2017.
- [22] P. Villalobos, J. Sevilla, L. Heim, T. Besiroglu, M. Hobbhahn, and A. Ho. Will we run out of data? an analysis of the limits of scaling datasets in machine learning. *arXiv*, arXiv:2211.04325, Oct 2022.
- [23] A. Wang et al. Superglue: A stickier benchmark for general-purpose language understanding systems. *arXiv*, February 2020.
- [24] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*, February 2019. Accessed: Jan. 25, 2024.
- [25] A. Warstadt, L. Choshen, A. Mueller, A. Williams, E. Wilcox, and C. Zhuang. Call for papers – the babyLM challenge: Sample-efficient pretraining on a developmentally plausible corpus. *arXiv 2301.11796*, January 2023. Accessed: Apr. 15, 2024.
- [26] A. Warstadt et al. Blimp: The benchmark of linguistic minimal pairs for english. *arXiv*, February 2023. Accessed: Jan. 22, 2024.
- [27] Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, and Ryan Cotterell. Findings of the BabyLM challenge: Sample-efficient pretraining on developmentally plausible corpora. In Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, and Ryan Cotterell, editors, *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, pages 1–34, Singapore, December 2023. Association for Computational Linguistics.