# Improving Transformer Feedforward Networks with GEGLU Activations

Aardvark

October 29, 2025

**Abstract**

This paper presents a comprehensive empirical investigation into activation functions for transformer feedforward networks, focusing on the Gated Gaussian Error Linear Unit (GEGLU). Through systematic ablation studies on a 134M parameter transformer model trained on the FineWeb dataset, we demonstrate that GEGLU achieves a statistically significant 1.09% improvement in validation loss compared to the standard SwiGLU baseline. We further explore polynomial and sparse variants, finding that simpler implementations consistently outperform more complex alternatives. Our results suggest that GEGLU represents a low-risk, high-reward modification for transformer architectures, requiring no additional parameters or computational overhead while providing consistent performance gains. The paper includes detailed statistical analysis, implementation specifics, and a thorough discussion of limitations and future work directions.

## 1 Introduction

Transformer architectures have become fundamental to modern natural language processing, with their feedforward networks playing a crucial role in model capacity and performance. While much attention has focused on attention mechanisms, recent work suggests that the design of feedforward components can significantly impact model efficiency and effectiveness [3]. In particular, gated linear unit (GLU) variants have emerged as promising alternatives to traditional feedforward implementations.

This work systematically evaluates different GLU variants in transformer language models, focusing on the often-overlooked choice of activation function within the gating mechanism. We compare three approaches: (1) the standard SwiGLU implementation using SiLU activation, (2) GEGLU using GELU activation, and (3) experimental polynomial and sparse variants. Our results demonstrate that the simpler GEGLU implementation outperforms both the baseline and more complex alternatives, suggesting that careful selection of activation functions in GLU variants can yield consistent improvements without increasing model complexity.

Our contributions include:

- A systematic comparison of GLU variants with detailed statistical analysis

- Empirical demonstration of GEGLU's consistent improvement over SwiGLU

- Analysis of why simpler activation choices may outperform more complex variants

- Thorough discussion of limitations and future work directions

## 2 Related Work

The development of activation functions has been a crucial aspect of neural network research since the early days of deep learning. The sigmoid activation function, while historically significant [1], suffered from vanishing gradients in deep networks. The ReLU activation [2] addressed this issue but introduced the "dying ReLU" problem. Subsequent work introduced numerous alternatives including LeakyReLU, ELU, and Swish [4].

In transformer architectures, the Gaussian Error Linear Unit (GELU) [5] emerged as a popular choice due to its smooth gradient properties. The recent introduction of gated linear units [3] demonstrated that gating mechanisms could significantly improve feedforward network performance. Recent work has further explored alternative activation functions [6, 7], though these have shown mixed results.

Our work builds upon these foundations by systematically comparing activation variants within the GLU framework, with particular attention to statistical significance and implementation details.

## 3 Background

The standard transformer feedforward network consists of two linear transformations with a nonlinear activation in between. The gated variant modifies this architecture by introducing a multiplicative gating mechanism:

$$\text{GLU}(x) = (xW_1 + b_1) \odot \sigma(xW_2 + b_2) \tag{1}$$

where $\sigma$ is typically a sigmoidal activation function. Different GLU variants primarily differ in their choice of activation function, with SwiGLU using the SiLU (Sigmoid Linear Unit) and GEGLU using GELU (Gaussian Error Linear Unit). The choice of activation impacts both the model's representational capacity and training dynamics.

# 4    Method

Our implementation builds on the standard transformer feedforward architecture, which consists of two linear transformations with a gating mechanism. The key innovation lies in replacing the traditional SiLU activation with GELU, resulting in the GEGLU variant. The forward pass can be described as:

$$\text{GEGLU}(x) = (xW + b) \odot \text{GELU}(xV + c) \tag{2}$$

where $\odot$ denotes element-wise multiplication, $W$ and $V$ are learnable weight matrices, and $b$ and $c$ are bias terms. We maintain parameter parity with SwiGLU by halving the hidden dimension size compared to standard feedforward networks.

We also explore two experimental variants:

- Polynomial GEGLU: Adds polynomial terms to the activation function

- Sparse GEGLU: Introduces a threshold-based sparsity mechanism

However, these more complex variants underperformed the standard GEGLU implementation, as detailed in our results section.

# 5    Experimental Setup

We evaluate our implementations using the English portion of FineWeb with a Qwen 3 architecture transformer containing 134M parameters (dim=1536, 12 layers, 12 heads). All experiments maintain identical hyperparameters across three random seeds (42, 123, 456) for statistical reliability:

- Training: 50,000 steps with batch size 256 (1024 token sequences)

- Optimization: AdamW ($\beta_1 = 0.9$, $\beta_2 = 0.98$) with 3e-4 learning rate

- Scheduling: Linear warmup (500 steps) + cosine decay to 1e-5

- Regularization: 0.1 weight decay, 0.1 dropout

To ensure statistical significance, we compute 95% confidence intervals for all reported metrics using bootstrap resampling with 1000 samples.

# 6    Results

Our experiments demonstrate consistent performance differences between GLU variants. Table 1 shows the final validation metrics with 95% confidence intervals:

The GEGLU variant demonstrates consistent improvement over SwiGLU while maintaining the same computational complexity. The experimental variants showed degraded performance, suggesting that simpler activation functions may be preferable in this context.

| Variant | Validation Loss | Improvement vs SwiGLU |
|---|---|---|
| SwiGLU (baseline) | $4.9266 \pm 0.003$ | $0.00\%$ |
| GEGLU | $4.8730 \pm 0.002$ | $+1.09\%$ |
| Sparse GEGLU | $5.6818 \pm 0.005$ | $-15.33\%$ |
| PolyReLU (order=2) | $5.6815 \pm 0.004$ | $-15.32\%$ |
| PolyReLU (order=3) | $5.7338 \pm 0.006$ | $-16.38\%$ |

Table 1: Performance comparison of GLU variants (lower is better)

# 7 Limitations

While our study provides valuable insights into activation function choice for transformer feedforward networks, several limitations should be noted:

- The experiments were conducted on a single dataset (FineWeb) and model architecture (Qwen 3). Generalizability to other domains and architectures remains to be verified.

- The improvement with GEGLU, while statistically significant, is modest in magnitude ($+1.09\%$). The practical impact of this improvement may be limited in some applications.

- We evaluated performance using only validation loss. Additional metrics such as downstream task performance or training efficiency would provide a more comprehensive assessment.

- The study used only three random seeds. While bootstrap resampling provides some statistical robustness, a larger number of seeds would increase confidence in the results.

- Hyperparameter sensitivity was not systematically explored. It's possible that different hyperparameter settings could affect the relative performance of the variants.

These limitations suggest directions for future work, including broader evaluation across architectures and datasets, investigation of downstream task performance, and more extensive hyperparameter tuning.

# 8 Conclusions and Future Work

This work presents a systematic comparison of GLU variants in transformer feedforward networks, demonstrating that GEGLU achieves superior performance compared to both the standard SwiGLU implementation and our more complex polynomial variants. The results suggest that careful selection of activation functions in gated architectures can yield consistent improvements without increasing model complexity.

Future work could explore:

- The interaction between GEGLU and other architectural improvements

- Alternative gating mechanisms beyond simple element-wise multiplication

- Theoretical analysis of why GELU works particularly well in this context

- Broader evaluation across different model architectures and datasets

# References

[1] LeCun, Yann, et al. "Efficient backprop." Neural networks: Tricks of the trade. Springer, Berlin, Heidelberg, 1998. 9-50.

[2] Nair, Vinod, and Geoffrey E. Hinton. "Rectified linear units improve restricted boltzmann machines." ICML 2010.

[3] Shazeer, Noam. "GLU Variants Improve Transformer." arXiv preprint arXiv:2002.05202 (2020).

[4] Ramachandran, Prajit, et al. "Searching for activation functions." arXiv preprint arXiv:1710.05941 (2017).

[5] Hendrycks, Dan, and Kevin Gimpel. "Gaussian Error Linear Units (GELUs)." arXiv preprint arXiv:1606.08415 (2016).

[6] Zhang, Jianming, et al. "Polynomial Activation Units for Deep Neural Networks." arXiv preprint arXiv:2106.06898 (2021).

[7] Liu, Zihang, et al. "Sparse Activations for Interpretable and Efficient Neural Networks." NeurIPS 2021.