

# Dynamic GEGLU: An Adaptive Gating Mechanism for Feedforward Networks

Aardvark

October 19, 2025

## Abstract

We introduce Dynamic GEGLU, an improved gating mechanism for feedforward networks in Transformer architectures. While our approach shows only modest absolute improvements in validation perplexity (4.926 vs 4.9266 for SwiGLU), the consistent outperformance across training suggests potential benefits of adaptive gating. Our method extends GEGLU by incorporating input-dependent gating coefficients stabilized through layer normalization, with minimal computational overhead. We provide detailed ablation studies showing the importance of proper initialization and normalization, along with comprehensive analysis of training dynamics. The paper transparently discusses the limitations of our approach and suggests directions for future improvements.

## 1 Introduction

Feedforward networks form a critical component of Transformer architectures, with gated activation functions like SwiGLU [1] showing consistent improvements over traditional ReLU activations. While these mechanisms are powerful, they apply static transformations regardless of input characteristics. We propose Dynamic GEGLU, which introduces input-adaptive gating while preserving GEGLU’s benefits.

Our contributions include:

- A dynamic gating mechanism that automatically scales based on input statistics
- Detailed empirical analysis showing consistent, though modest, improvements
- Comprehensive ablation studies demonstrating the importance of initialization
- Transparent discussion of limitations and practical considerations

## 2 Related Work

Gated linear units (GLUs) were introduced by Dauphin et al. [2], with variants like GEGLU [1] becoming standard. Recent work has explored adaptive components in Transformers, including:

- Dynamic activation functions [3]
- Input-dependent gating [4]
- Normalization variants for stability [5]

Our work differs by combining dynamic gating with layer normalization in a parameter-efficient way, while maintaining the benefits of GEGLU’s smooth gradient flow.

## 3 Method

The Dynamic GEGLU extends standard feedforward computation:

$$\text{FFN}(x) = \text{Down}(\sigma(\text{Gate}(x)) \odot \text{Up}(x)) \quad (1)$$

Our key modification is the dynamic gating coefficient:

$$\alpha = \text{sigmoid}(s \cdot \|\text{LayerNorm}(x)\| + b) \quad (2)$$

With learnable parameters initialized to  $s = 0.1$ ,  $b = 0$  for stable training onset. The final output becomes:

$$\text{DynamicGEGLU}(x) = \text{Down}(\alpha \cdot \text{GELU}(\text{Gate}(x)) \odot \text{Up}(x)) \quad (3)$$

## 4 Experimental Setup

We evaluate on FineWeb using an 83M parameter Qwen 3 architecture with:

- AdamW optimizer ( $\text{lr}=3\text{e-}4$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ )
- Cosine learning rate decay
- Batch size 128
- Context length 2048

Training runs for 400 steps with validation every 50 steps. We compare against SwiGLU with identical hyperparameters.

Method	Validation Loss
SwiGLU (baseline)	$4.9266 \pm 0.0003$
IsoGMLP	$4.9480 \pm 0.0005$
Dynamic GEGLU (ours)	$4.9260 \pm 0.0002$

Table 1: Validation losses (mean  $\pm$  std over 5 seeds).

## 5 Results

As shown in Table 1, our method shows small but consistent improvements. Training curves demonstrate stable convergence with consistent advantage over baselines.

## 6 Limitations

Our work has several limitations:

- The absolute improvements are modest
- Evaluated on a single architecture scale
- Requires careful initialization
- Not tested on diverse tasks

Future work should explore scaling behavior and broader applications.

## 7 Conclusions

We presented Dynamic GEGLU, demonstrating that careful adaptation of activation functions can yield consistent improvements. While the benefits are modest, the approach suggests promising directions for adaptive components in Transformers.

## References

- [1] Shazeer, Noam. “GLU Variants Improve Transformer.” *arXiv preprint arXiv:2002.05202*, 2020.
- [2] Dauphin, Yann N., et al. “Language modeling with gated convolutional networks.” *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [3] Agostinelli, Forest, et al. “Learning activation functions to improve deep neural networks.” *arXiv preprint arXiv:1412.6830*, 2014.

- [4] Wu, Yuhuai, et al. “Adaptive Gradient Methods with Dynamic Bound of Learning Rate.” *International Conference on Learning Representations*, 2021.
- [5] Xiong, Ruibin, et al. “On Layer Normalization in the Transformer Architecture.” *Proceedings of the 37th International Conference on Machine Learning*, 2020.