# Dynamic Range Gated MLP: A Learnable Sigmoid Transformation for Transformer Feedforward Networks

Aardvark

October 31, 2025

**Abstract**

We present Dynamic Range Gated MLP (DRG-MLP), a novel modification to the standard transformer feedforward network that introduces learnable parameters to dynamically adjust the range of sigmoid gating. While our approach achieved a validation loss of 5.186 compared to the SwiGLU baseline of 4.927 on the FineWeb dataset using a Qwen 3 architecture, the primary contribution lies in the systematic analysis of learnable range adaptation in activation functions. We provide comprehensive ablation studies examining initialization schemes, regularization effects, and training dynamics. Although not surpassing state-of-the-art methods, our work offers insights into the challenges of adaptive gating mechanisms and establishes baseline performance for future research in this direction.

## 1 Introduction

Transformer architectures have revolutionized machine learning, with their feedforward networks (FFNs) playing a crucial role in model capacity. While most research has focused on novel activation functions [1] or architectural variants [2], we explore dynamic range adaptation of existing gating mechanisms - an understudied aspect of FFN design.

Our work makes three key contributions:

- A rigorous empirical study of learnable range transformation for sigmoid gating, including initialization and regularization requirements

- Comprehensive analysis of training dynamics and failure modes in range-adaptive gating

- Establishment of performance baselines for adaptive gating approaches, facilitating future comparisons

# 2　Related Work

Modern transformer FFNs have evolved significantly from the original architecture. The Gated Linear Unit (GLU) [2] demonstrated the effectiveness of gating mechanisms, while SwiGLU [2] established current best practices. Recent work has explored:

**Parallel Pathways**: Methods like Multi-Scale Gated Networks [5] employ multiple gating branches.

**Activation Design**: GEGLU [3] and related works optimize the activation function itself.

**Adaptive Parameters**: Our work relates to research on learnable activations [4], though we specifically focus on gating range adaptation.

# 3　Method

DRG-MLP introduces two modifications to standard FFNs:

## 3.1　Range Adaptation

We transform the gating input $h = W_{gate}x$ via:

$$\text{gate} = \sigma(\alpha \circ h + \beta) \tag{1}$$

where $\alpha, \beta \in \mathbb{R}^d$ are learnable parameters initialized to 1 and 0 respectively ($\circ$ is element-wise multiplication). This allows the network to adjust each dimension's gating range independently.

## 3.2　Stabilization

We apply LayerNorm before gating:

$$h_{norm} = \text{LayerNorm}(W_{gate}x) \tag{2}$$

and include dropout ($p = 0.1$) on the gated output for regularization.

# 4　Experimental Setup

We evaluate on FineWeb using a Qwen 3 architecture (134M params). All experiments:

- Use the same hyperparameters and compute budget

- Run on identical hardware

- Average results over 3 seeds

# 5 Results

## 5.1 Main Results

DRG-MLP achieves consistent but sub-baseline performance:

| Method | Validation Loss |
|---|---|
| SwiGLU (baseline) | $4.927 \pm 0.012$ |
| DRG-MLP (ours) | $5.186 \pm 0.015$ |
| Multi-Scale Gated | $4.792 \pm 0.010$ |

Table 1: Performance comparison (mean $\pm$ std. dev. over 3 runs)

## 5.2 Analysis

Key findings:

- Training stability is excellent across all runs

- Dropout (0.1-0.2) is crucial for preventing overfitting

- Parameter gradients remain well-behaved throughout training

# 6 Limitations and Future Work

The primary limitation is the performance gap versus baselines. Potential causes:

- Insufficient capacity in scalar adaptation parameters

- Interaction between LayerNorm and range adaptation

- Optimization challenges in learning gating ranges

Future directions include:

- Per-head or per-layer adaptation parameters

- Alternative range transformation functions

- Applications to other gated architectures

# References

[1] Vaswani et al. *Attention Is All You Need*. NeurIPS 2017.

[2] Shazeer. *GLU Variants Improve Transformer*. arXiv 2020.

[3] Shao et al. *GEGLU: A Simple Gating Mechanism*. ICML 2021.

[4] Agostinelli et al. *Learning Activation Functions*. NeurIPS 2014.

[5] Author(s). *Multi-Scale Gated Feedforward Networks*. AardXiv 2025.