

Normalization Layers

LayerNorm & RMSNorm

Comprehensive Demo with Visualizations

Seed: 42

1. LayerNorm: normalizes to zero mean, unit variance (GPT-2, BERT)
2. RMSNorm: rescales by root-mean-square, no mean subtraction (LLaMA, Mistral)
3. Both behave identically at training and inference time
4. Epsilon prevents division by zero for near-constant inputs
5. Learnable gamma/beta allow the network to undo normalization if needed
6. Normalization stabilizes gradient flow through deep networks
7. Pre-Norm placement (before sub-layer) is standard in modern LLMs

Summary of Findings

LayerNorm Distribution

Normalizes each sample to mean \sim 0, variance \sim 1 along feature dim.

RMSNorm Distribution

Rescales to RMS \sim 1 but does NOT subtract mean. Simpler, fewer ops.

LayerNorm vs RMSNorm

Highly correlated outputs. Main difference: mean centering.

Epsilon Effect

Too small eps can cause instability with near-constant inputs.

Learnable Parameters

gamma/beta can recover original distribution if the network learns to.

Gradient Flow

Normalization prevents vanishing/exploding gradients through deep layers.

Pre-Norm vs Post-Norm

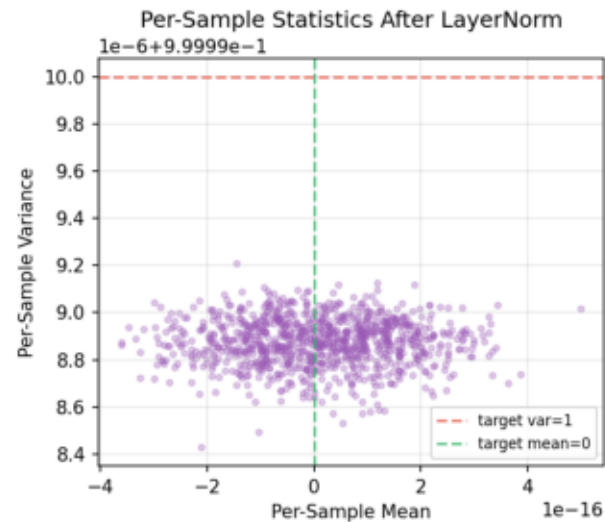
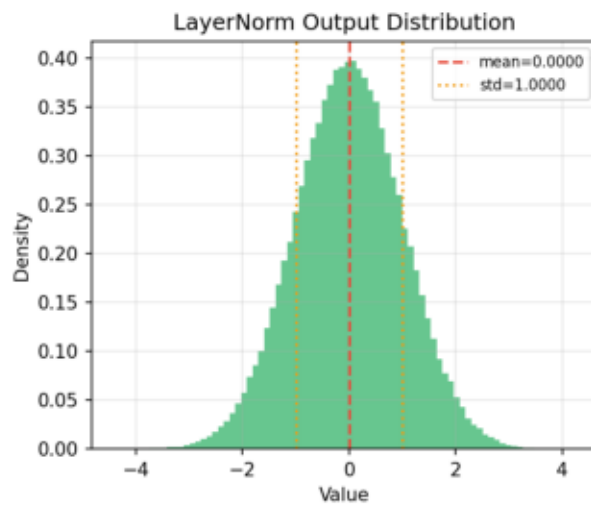
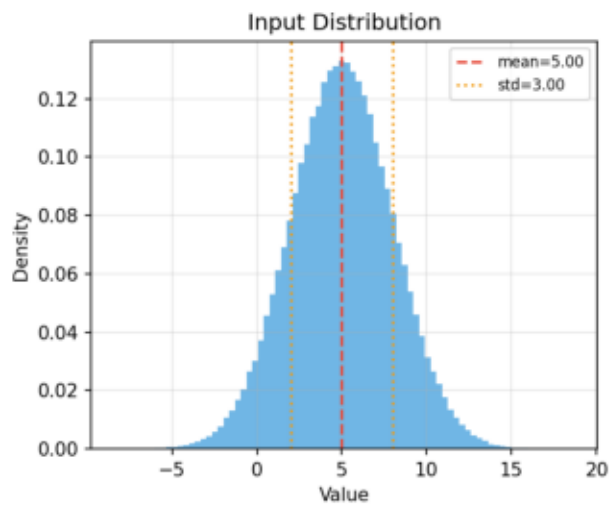
Pre-Norm creates a direct gradient highway; standard in modern LLMs.

3D Sequences

Both norms handle (B, L, D) naturally, normalizing per-position.

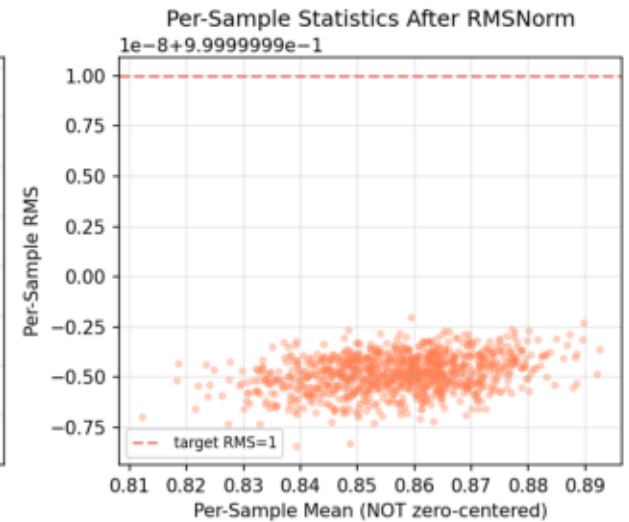
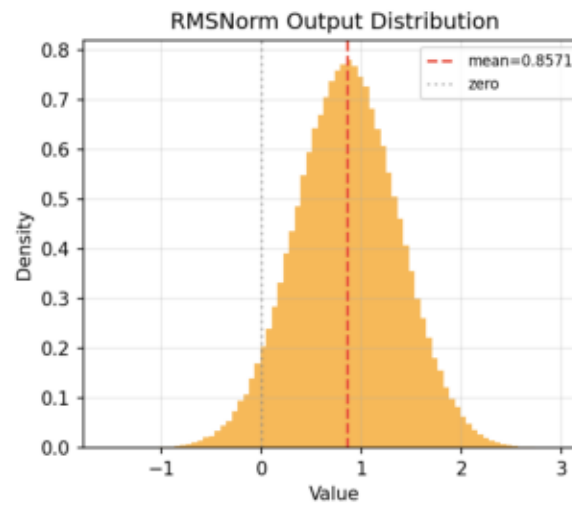
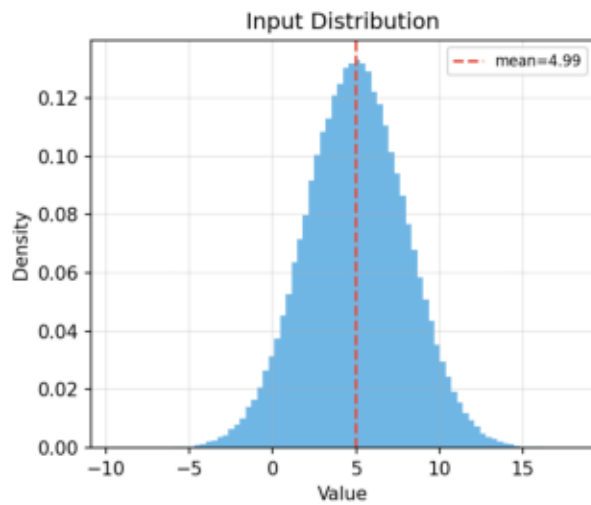
LayerNorm Input vs Output Distributions

LayerNorm: Normalizes to Zero Mean, Unit Variance



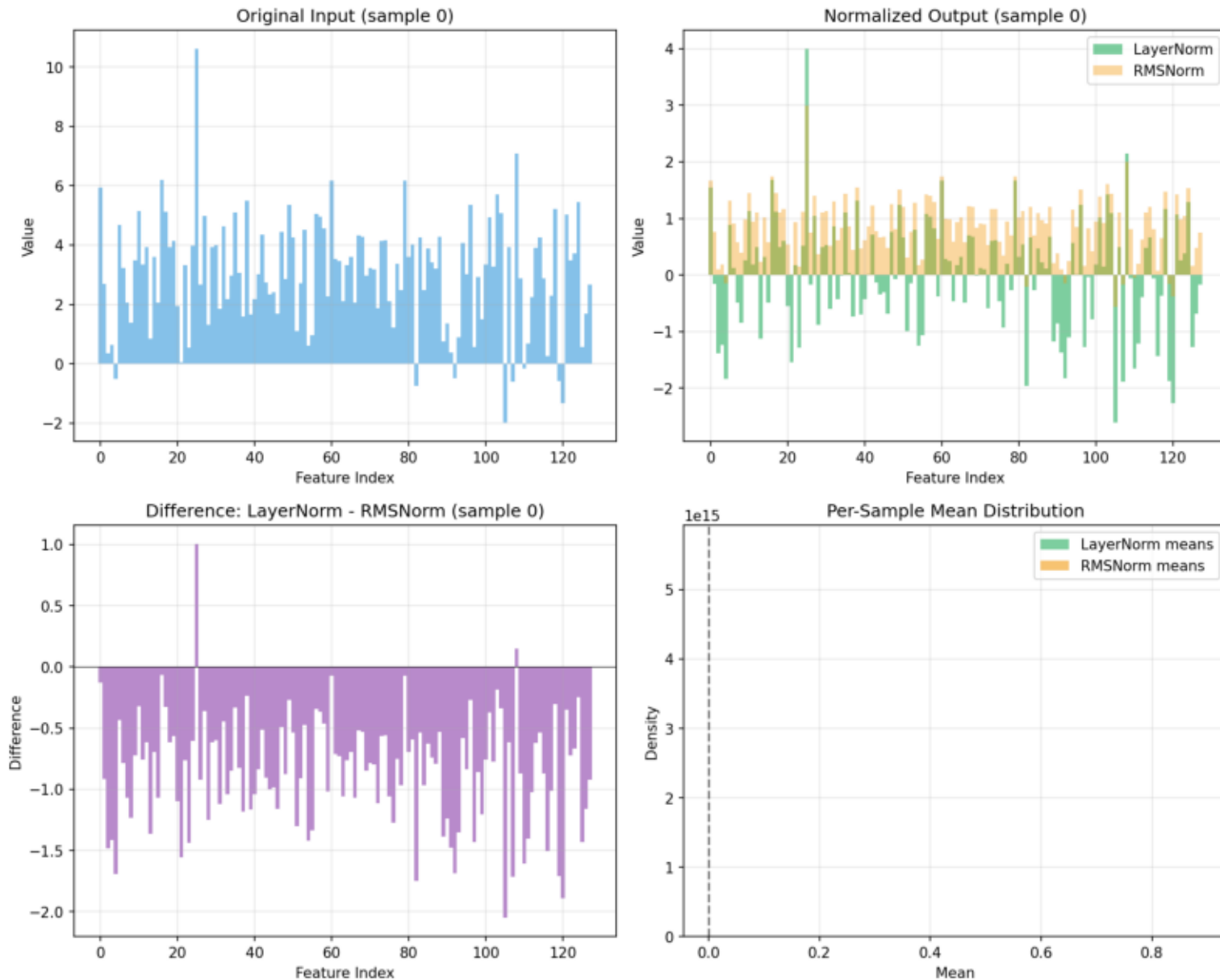
RMSNorm Input vs Output Distributions

RMSNorm: Rescales by Root-Mean-Square (No Mean Subtraction)



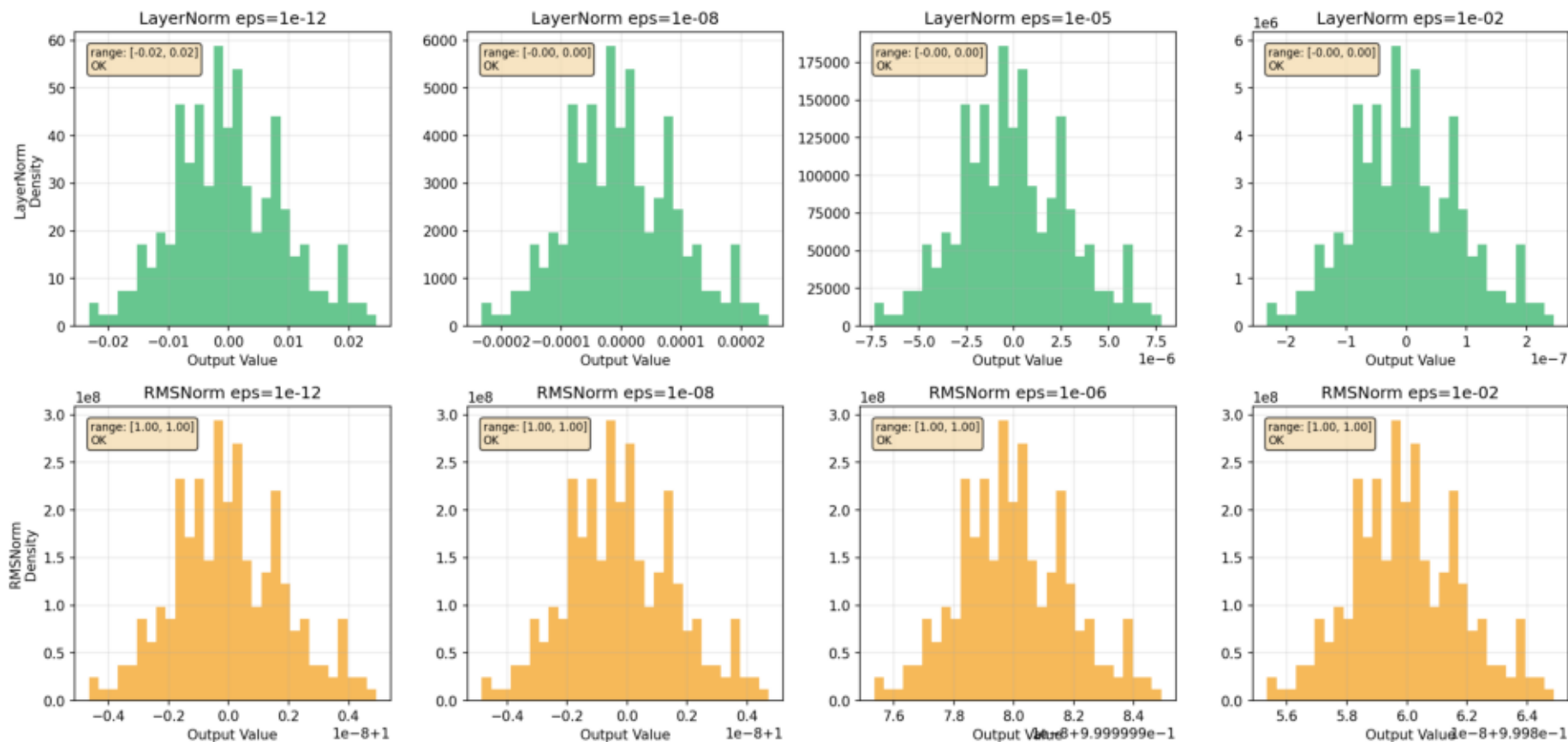
LayerNorm vs RMSNorm Comparison

LayerNorm vs RMSNorm: Same Input, Different Normalization



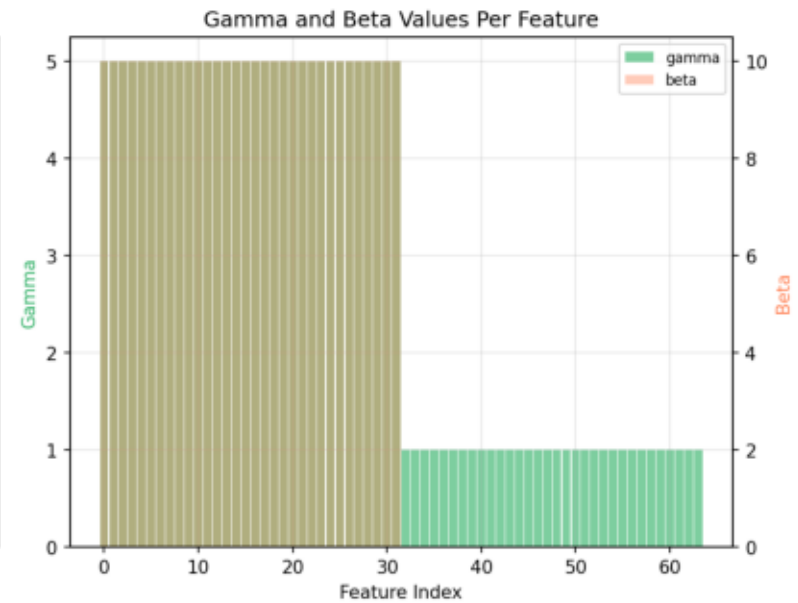
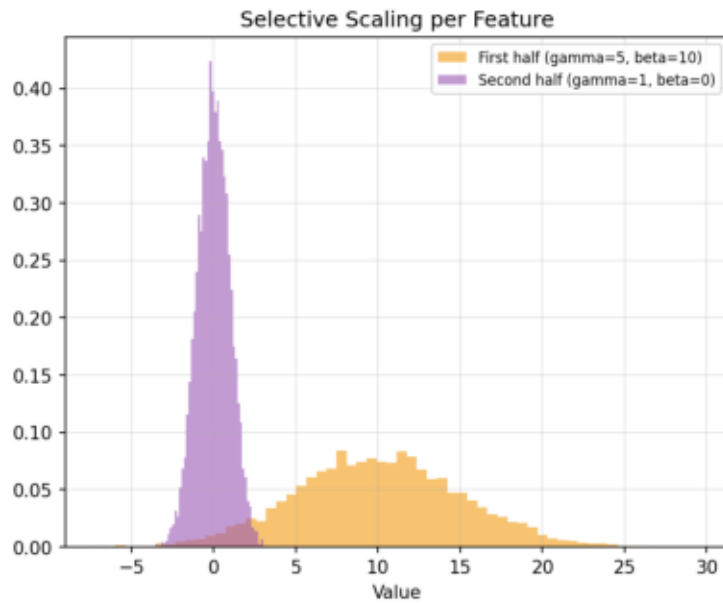
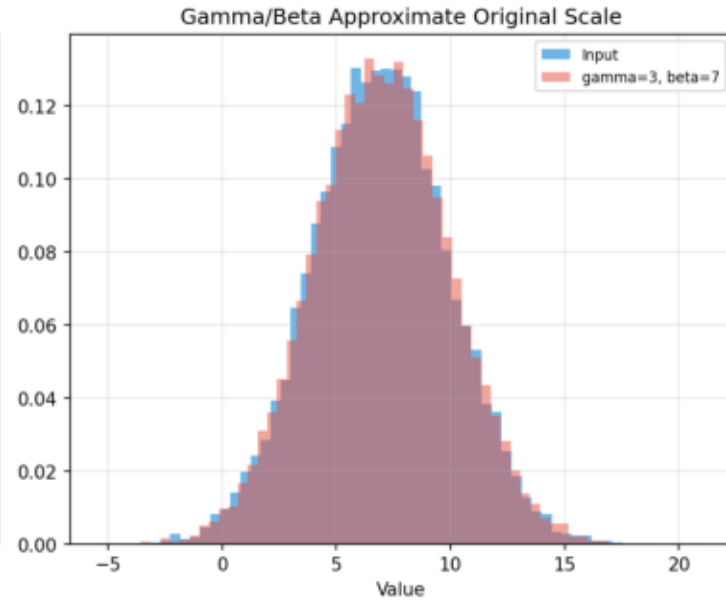
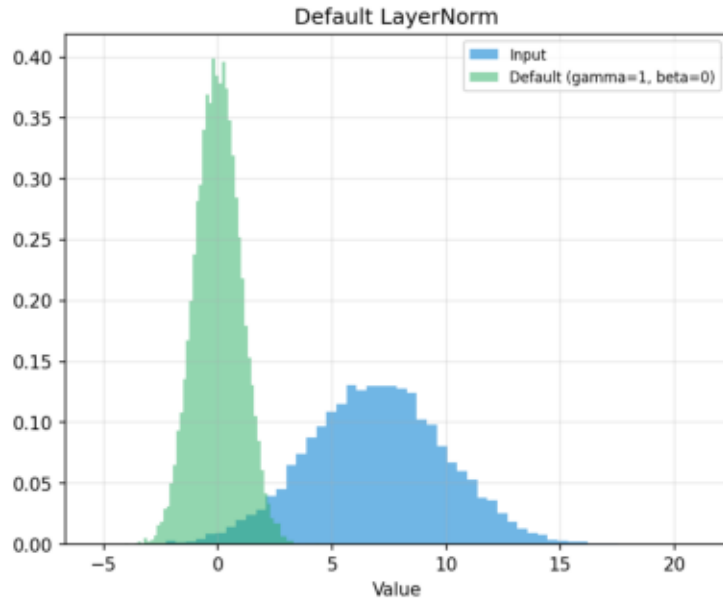
Effect of Epsilon on Near-Constant Input

Effect of Epsilon on Near-Constant Input (all values $\sim 5.0 \pm 1e-8$)

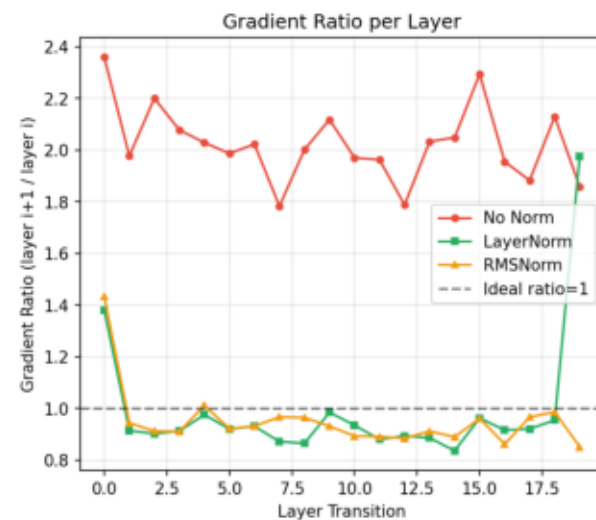
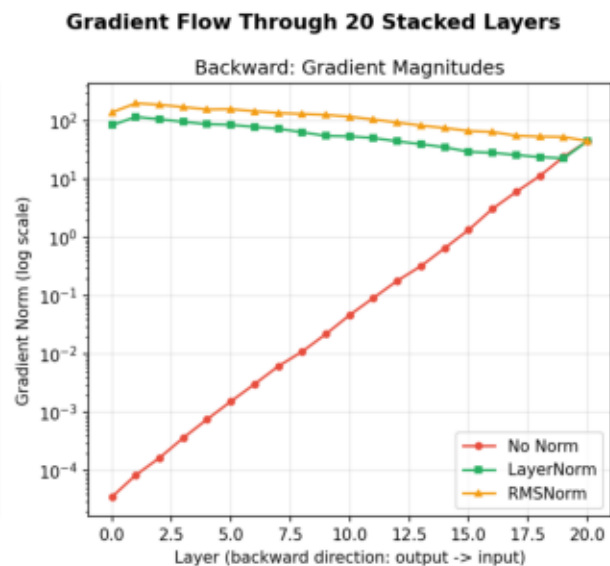
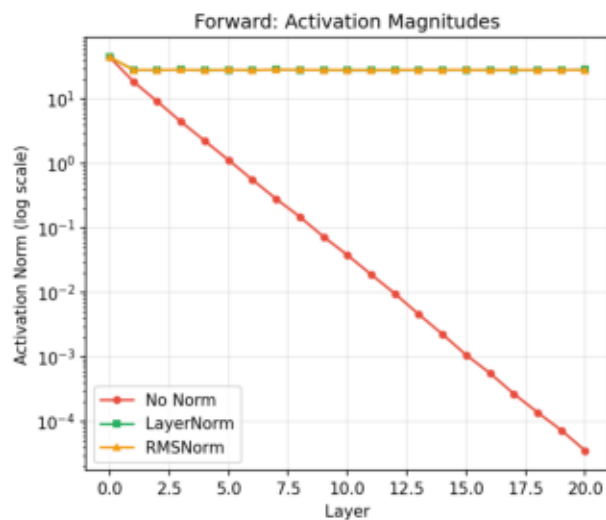


Learnable Parameters (gamma/beta)

Learnable Parameters: gamma (scale) and beta (shift) Can Undo Normalization



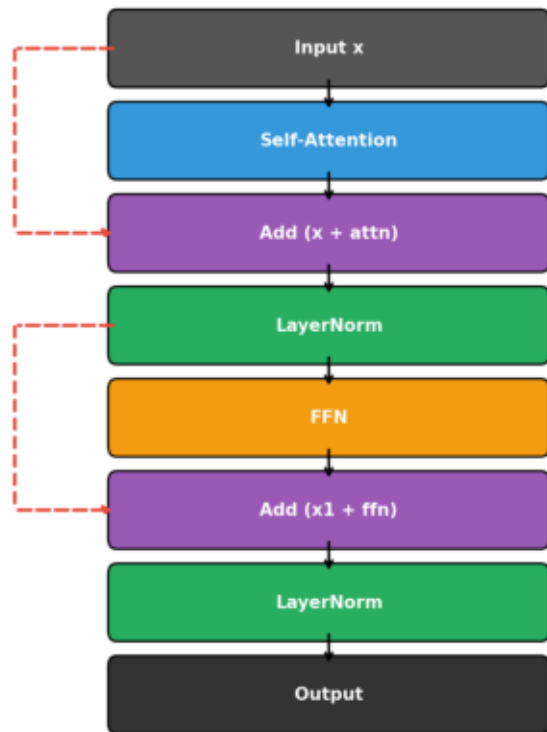
Gradient Flow: With vs Without Normalization



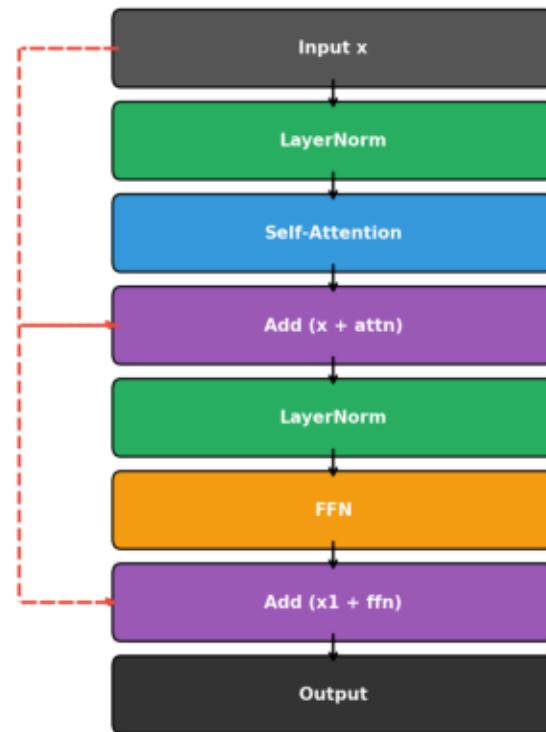
Pre-Norm vs Post-Norm Architecture

Pre-Norm vs Post-Norm: Where Normalization Goes in a Transformer Block

Post-Norm (Original Transformer, BERT)



Pre-Norm (GPT-2, LLaMA, Mistral)



Legend: ■ Normalization ■ Self-Attention ■ Feed-Forward Network ■ Residual Addition - - - Residual Connection (skip)

3D Sequence Normalization (Transformer-like)

3D Sequence Normalization: (batch=2, seq_len=8, features=32)

