

図 3: 分布形状による丸め誤差シミュレーション

#### 4. 提案手法

本研究では、GELU 関数によるロングテール分布の課題を解決するために、アクティベーション分布を整形する新たなコネクションである Activation Shaping Connection (ASC) を提案する。ASC は、GELU 関数により生じるロングテール分布を整形することで、量子化時の丸め誤差の累積を抑制し、精度を維持することを目指す。

##### 4.1 Activation Shaping Connection

ASC は、アクティベーションの分布を整形することで、丸め誤差の累積を抑制し、学習の安定性を向上させるためのコネクションである。本研究では、ViT の MLP ブロックで使用される GELU 関数のロングテール分布を、正負の値がより均等に分布する形状に整形するように ASC を導入する。ASC を導入した MLP ブロックの構造を図 4 に示す。この整形処理では、BitLinear 層を用いて変換し、分布を整形する。ASC の処理を式 (4) に示す。

$$\mathbf{y} = \text{GELU}(\text{BitLinear}(\mathbf{x})) \odot \text{BitLinear}(\mathbf{x}) \quad (4)$$

ここで、 $\mathbf{x}$  はアクティベーション、 $\odot$  は要素積を示す。具体的には、BitLinear 層の後に GELU 関数を適用する通常の処理に加え、新たに BitLinear 層を通過するコネクションを導入する。この導入したコネクションの出力と、通常の処理の出力間で要素積を計算することで、整形された出力  $\mathbf{y}$  を得る。この結果、整形した  $\mathbf{y}$  を MLP ブロックの 2 層目の BitLinear 層に入力する。これにより、GELU 関数のロングテール分布に起因する課題が解決され、丸め誤差の累積を抑制できる。

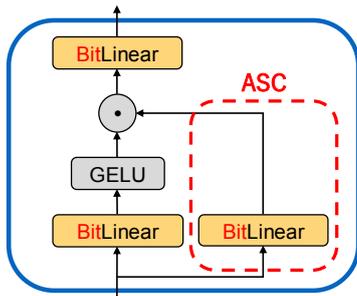


図 4: ASC を導入した MLP ブロック

##### 4.2 ガウス分布整形の Loss の導入

MLP ブロックの 2 層目の BitLinear 層への入力分布は、ASC により得られる出力分布に依存する。そのため、ASC の出力分布が極端に偏ると、意図した分布整形が機能しなくなるため、適切な対称性を持たせる必要がある。そこで、Loss 関数にガウス分布への近似を促す制約を追加する。具体的には、ASC により出力される BitLinear 層の分布をガウス分布に近似させるため、 $\lambda_{\text{gauss}}$  を追加して分布の整形を行う。Loss 関数の全体式を式 (5) に示す。

$$\lambda_{\text{all}} = \lambda_{\text{label}} + \alpha \cdot \lambda_{\text{gauss}} \quad (5)$$

ここで、 $\lambda_{\text{label}}$  は従来のラベルに基づく損失、 $\lambda_{\text{gauss}}$  は出力分布をガウス分布に近似させるための損失、 $\alpha$  は両損失項の重み付けバランスを調整するハイパーパラメータであ

る。 $\lambda_{\text{gauss}}$  は、各層の ASC におけるアクティベーションを標準化し、標準化したアクティベーションとガウス分布との間の KL ダイバージェンスを計算することで、アクティベーションをガウス分布に近似させる。

#### 5. 評価実験

本章では、提案手法の有効性を検証するために、ImageNet-1k データセットを用いて評価実験を実施する。

##### 5.1 Activation Shaping Connection の有効性

提案手法である ASC の有効性を評価するため、従来の MLP ブロックを使用した BitNet と、ASC を導入したモデルと比較する。表 1 に示すように、ASC の導入により BitNet と比較して精度が 7.22 ポイント向上した。この改善は、GELU 関数のロングテール分布を適切に整形し、丸め誤差の累積を抑制したことが主な要因であると考えられる。また、量子化なしの ViT と比較して、Linear 層の圧縮率が 63.59% に達した。一方で、ガウス分布に近似する Loss のハイパーパラメータ  $\alpha$  を 1.0 に設定すると精度が低下した。しかし、 $\alpha = 0.0009$  に設定した場合は、Loss を導入しない場合より精度が向上した。これは、ガウス分布の近似より分布整形が安定し、丸め誤差の抑制および学習全体の安定性が向上したと考える。

表 1: ASC を用いた評価結果

モデル	量子化	ASC	$\alpha$	サイズ [MB]	精度 [%]
ViT	-	-	-	2084.33	64.68
BitNet	✓	-	-	531.14	49.38
	✓	✓	0	758.76	56.60
	✓	✓	1.0	758.76	48.93
	✓	✓	0.9	758.76	53.56
	✓	✓	0.09	758.76	55.98
	✓	✓	0.009	758.76	56.61
	✓	✓	0.0009	758.76	56.97

##### 5.2 ASC によるアクティベーション分布の可視化

ASC の有効性を検証するため、アクティベーション分布の変化を可視化する。図 5 に示すように、ASC の各ルートで得られる分布形状を比較した結果、アクティベーションが最終的に正負の値がより均等に分布する形状に整形されることを確認した。この結果、ASC は GELU 関数によるロングテール分布の課題を抑制し、次の BitLinear 層での丸め誤差の累積を抑制できることを示す。

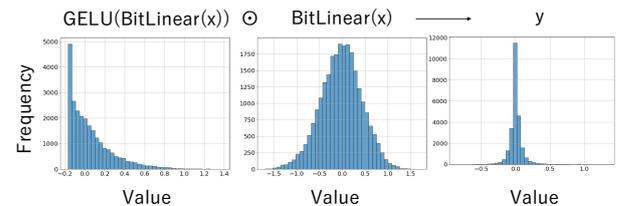


図 5: ASC アクティベーション分布の可視化

#### 6. おわりに

本研究では、小規模ハードウェア環境で動作可能な Transformer モデルを実現するため、GELU 関数のロングテール分布を整形し、丸め誤差を抑制する ASC を導入した量子化手法を提案した。評価実験の結果、提案手法が精度向上とモデル圧縮の両立に有効であることを確認した。また、小規模なハードウェア環境への適用可能性が示され、実用化への期待が高まる結果となった。今後は、より複雑なモデルや多様なタスクへの適用可能性についてさらなる検証を進める予定である。

#### 参考文献

- [1] L. Yijiang, *et al.*, “Noisyquant: Noisy bias-enhanced post-training activation quantization for vision transformers”, WACV, 2023.
- [2] M. Shuming, *et al.*, “The Era of 1-bit LLMs: All Large Language Models are in 1.58 Bits”, arXiv:2402.17764, 2024.

#### 研究業績

- [1] 若山浩之 等, “CNN と ViT を組み合わせたモデルのノイズへの頑健性”, 画像の認識・理解シンポジウム, 2023.