

ノイズに着目している Attention Head の同定とスケーリングによる  
大規模言語モデル出力の正確性と多様性のトレードオフ制御  
Controlling the Trade-off between Accuracy and Diversity in LLM Outputs  
by Identifying and Scaling Noise-Focusing Attention Heads

下村 晃生<sup>†</sup> 花沢 明俊<sup>†</sup>  
Teruki Shimomura Akitoshi Hanazawa

### 概要

大規模言語モデル(LLM)は、汎用的なテキスト生成能力から、様々な分野への活用が広がってきている。しかし、不必要な情報が含まれる場合、推論能力が大幅に低下することが知られている。[1]本論文では、そのような現象を起こす原因として、LLM の中に無数にある Head の中に、文章中のノイズに対して過度に Attention を向けている Head があるからであるという仮説を立て、検証を行った。その結果、そのような Head は LLM の中間層に多く局在し、この Head を調整することによりモデルの出力を操作ができることを実証した。

## 1. 研究背景

大規模言語モデル(LLM)は、小学生レベルの算数の文章問題データセットである GSM8K[2]をはじめとした各種タスクにおいてその数理推論能力を評価されてきたが、GSM-Symbolic[1]では、回答に無関係な人物名や数量、背景説明といったノイズを付与するだけで、モデルの精度が著しく低下する現象が報告されている。このように、LLM は不必要な情報が含まれると出力品質が劣化するプロンプトロバスタ性の低さが問題となっているが、モデル内部でノイズに反応している要素を特定し、制御する手法は未だ確立されていない。

一方、Transformer の Attention Head 単位での役割分析は、Olsson et al.[3]や Wu et al.[4]らによって進められており、それぞれ特定のタスクや構文的特徴を捉えるヘッドの存在が示されている。本研究では、「文脈中のノイズに対して Attention を向けている Head も存在するのではないか」という仮説をたて、その有無を検証し、特定した Head が持つ役割を検証するために、Head の調整による出力特性への影響を検証することを目的とした。

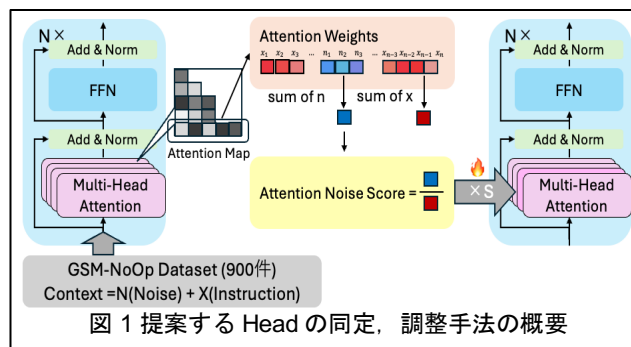
## 2. 提案手法

### 2.1 Attention Head の同定

本研究で行った実験には meta 社の Llama-3.1-8B-Instruct モデルを使用した。

#### 2.1.1 GSM-NoOp データセット

まず、GSM-NoOp データセットを独自に再現した。このデータセットは、Mirzadeh らによって提案[1]された小学生レベルの算数の文章問題に対して、回答に無関係な人物名や数量、背景説明を加えたものである。入力小学生レベルの算数文章問題 $x$ であり、文章中には解答に無関係なノイズ $N$ が混入している。モデル $f_{\theta}$ は $\hat{y} = f_{\theta}(x, N(x))$ を生成するが、理想的には $N(x)$ の有無で $\hat{y}$ は不変となるべきである。また、評価指標は正答率である。



#### 2.1.2 GSM-NoOp データセットの再現

OpenAI 社の GPT-4o に Few-Shot を提示して 1300 件作成し、混入させたノイズが問題とは無関係であり、回答に影響を与えていないことは人手によって判断を行った。900 件を Head 同定用、400 問を性能評価用に分割した。混入させたノイズ部の平均トークンは 29.4、割合は 43.6%であった。

## 2.2 Noise-Focusing Attention Head の同定

図1にパイプラインを示した。

### 2.2.1 Noise Focusing Attention Head の同定

入力するトークン集合をノイズ部分 $N$ と非ノイズ部分 $C$ に分割し、レイヤ $l$ 、ヘッド $h$ の Noise-Focusing Attention Score (NF-Score)を式(1)で定義する。

$$NF-Score_{l,h} = \frac{\sum_{i \in N} \alpha_{l,h}(i)}{\sum_{j \in C} \alpha_{l,h}(j)} \quad (1)$$

ここで $\alpha_{l,h}(i)$ はトークン $i$ へ向けられた注意重みである。同定用の GSM-NoOp データセット 900 問全体で平均をとり、上位 $k$  % を NF-Head と呼ぶ。

### 2.3 Head スケーリングによる重み調整

特定した Head について、式(2)のように NF-Head が持つ役割を検証するために再学習を行わずに value 行列にスカラー $S$  (Scaling Factor)を掛けるだけで NF-Head の調整を行う。

$$H_{NF} = \text{softmax} \left( \frac{W_q^i \times W_k^{iT}}{\sqrt{d_k}} \right) (W_v^i \cdot S) \quad (2)$$

- $S = 1$ : 元モデル
- $0 < S < 1$ : NF-Head 抑制
- $S = 0$ : NF-Head 無効化
- $S > 1$ : NF-Head 増幅

## 2.4 評価指標

### 2.4.1 推論能力の評価

推論能力は、同定に使用していない GSM-NoOp データセット 400 件に対する正答率で評価を行った。

