# Focus Directions Make Your Language Models Pay More Attention to Relevant Contexts

Youxiang Zhu, Ruochen Li, Danqing Wang, Daniel Haehn, Xiaohui Liang

## 1. LLMs are prone to distracted by irrelevant context. Why?
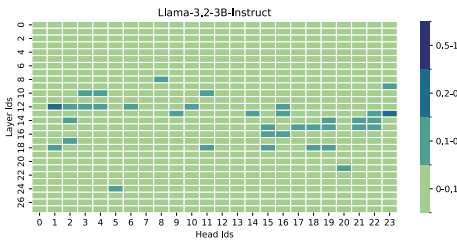
### 1.1 Identify contextual heads

**Contextual Scoring:** A metric that quantifies the degree of attention allocated to specific segments of the input (e.g., relevant contexts) during response generation.

**Contextual Heads:** The top-$k$ attention heads ranked by contextual score (to relevant contexts)
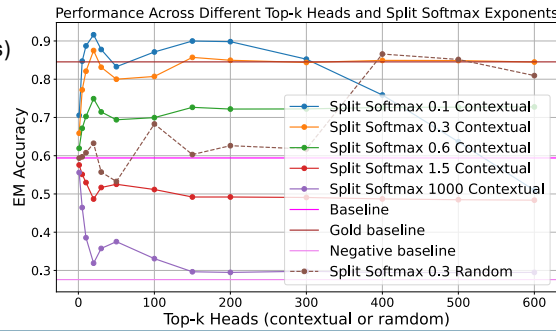
### 1.2 Properties of contextual heads

- Contextual heads are **sparse**, located in **middle and late layers**

| LLM Response | Correct | Wrong |
|---|---|---|
| Focus to relevant contexts | More | Less |


Llama-3.2-3B-Instruct

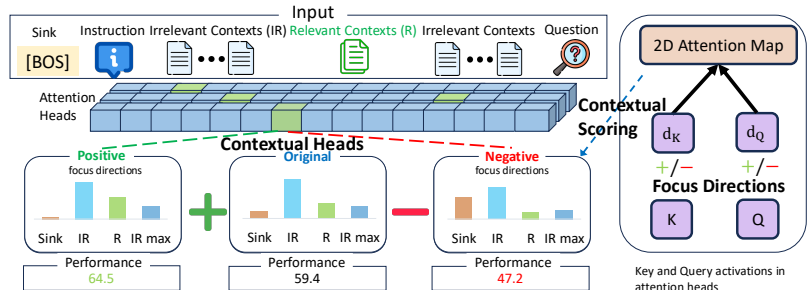### 1.3 Modifying attention on contextual heads

- **Method**: Split softmax (<1 increase attention, > 1 decrease)
- **Increase** attention to relevant contexts: performance ↑
- **Decrease** attention to relevant contexts: performance ↓
- Non contextual heads have minimum such effects


Performance Across Different Top-k Heads and Split Softmax Exponents

**Contextual heads controls the overall attention of LLMs**

Modifying attention on contextual heads could make performance better than gold baseline (relevant context only), or close to the negative baseline (irrelevant contexts only)

### 0. Overview



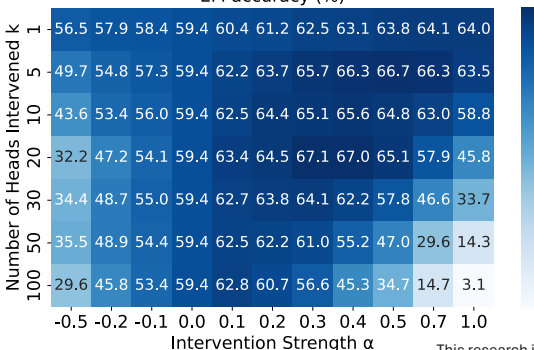## 2. Focus directions move attention from sink to relevant contexts

**Motivation:** Can contextual heads figure out the relevant contexts by themselves?

**Locating focus directions**: simply train $d_K$ and $d_Q$ to maximize the contextual score for the relevant contexts.

**Main findings:**

| Focus direction | Relevant contexts | Sink |
|---|---|---|
| Positive | More attention | Less attention |
| Negative | Less attention | More attention |

- Focus directions only help mitigate distraction on **contextual heads**.

EM accuracy (%)



## 3. Focus directions mitigate poor task alignment

**Benchmark:** HELMET (5 categories, 16 tasks used)

**Main findings:**

- Focus directions help for the **long context tasks** that LLM could do well in the **short context**.
- Most of the tasks could be improved by **either positive or negative** focus direction.
- Focus direction improves the overall performance of **poorly aligned LLMs**. (e.g., base vs. instruct, inconsistent sink score for the same context length)

| Model | Recall | RAG | Re-ranking | ICL | Long QA | Overall Average | Model | Recall | RAG | Re-ranking | ICL | Long QA | Overall Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Llama-3.2-3B** | | | | | | | **Llama-3.2-3B** | | | | | | |
| 20. -0.2 | 66.00 | 54.96 | 29.22 | 82.20 | - | 58.10 | 20. -0.2 | 55.50 | 50.38 | 6.83 | 85.20 | - | 49.48 |
| 10. -0.2 | 73.81 | 56.58 | 26.73 | 83.00 | - | 60.03 | 10. -0.2 | 64.56 | 53.96 | 6.24 | 86.20 | - | 52.74 |
| 20.0.2 | 81.50 | 58.75 | 25.37 | 80.20 | - | 61.46 | 20.0.2 | 66.31 | 56.46 | 7.08 | 85.40 | - | 53.81 |
| 10.0.2 | 82.00 | 58.54 | 26.16 | 80.60 | - | 61.83 | 10.0.2 | 65.69 | 55.83 | 9.27 | 85.40 | - | 54.05 |
| baseline | 78.88 | 58.83 | 26.10 | 82.20 | - | 61.50 | baseline | 65.50 | 54.83 | 7.29 | 86.20 | - | 53.46 |
| **Llama-3.2-3B-Instruct** | | | | | | | **Llama-3.2-3B-Instruct** | | | | | | |
| 20. -0.2 | 73.00 | 58.04 | 13.68 | 78.80 | 27.32 | 50.17 | 20. -0.2 | 56.38 | 56.75 | 3.77 | 83.80 | 28.64 | 45.87 |
| 10. -0.2 | 79.12 | 60.21 | 13.52 | 79.40 | 26.66 | 51.74 | 10. -0.2 | 61.06 | 58.33 | 2.44 | 85.00 | 30.38 | 47.44 |
| 20.0.2 | 83.50 | 60.25 | 20.58 | 80.60 | 26.09 | 54.20 | 20.0.2 | 61.81 | 59.21 | 2.72 | 83.40 | 28.23 | 47.87 |
| 10.0.2 | 83.69 | 62.05 | 20.77 | 80.40 | 25.94 | 54.58 | 10.0.2 | 64.25 | 59.79 | 3.10 | 84.20 | 26.80 | 47.63 |
| baseline | 84.38 | 63.00 | 17.13 | 80.20 | 26.78 | 54.30 | baseline | 64.12 | 59.96 | 3.77 | 85.00 | 31.13 | 48.80 |
| **Qwen2.5-7B** | | | | | | | **Qwen2.5-7B** | | | | | | |
| 20. -0.2 | 94.56 | 53.50 | 22.86 | 77.60 | - | 62.13 | 20. -0.2 | 42.06 | 41.96 | 1.30 | 77.40 | - | 40.68 |
| 10. -0.2 | 95.31 | 54.36 | 24.84 | 78.40 | - | 63.28 | 10. -0.2 | 45.00 | 43.42 | 2.64 | 77.40 | - | 42.11 |
| 20.0.2 | 95.88 | 54.04 | 23.29 | 79.40 | - | 63.14 | 20.0.2 | 46.56 | 43.08 | 1.19 | 77.60 | - | 42.11 |
| 10.0.2 | 95.50 | 54.08 | 23.11 | 80.00 | - | 63.17 | 10.0.2 | 46.56 | 43.62 | 1.07 | 78.80 | - | 42.51 |
| baseline | 96.00 | 54.21 | 23.15 | 79.60 | - | 63.24 | baseline | 45.19 | 44.12 | 1.88 | 78.00 | - | 42.30 |
| **Qwen2.5-7B-Instruct** | | | | | | | **Qwen2.5-7B-Instruct** | | | | | | |
| 20. -0.2 | 94.38 | 55.87 | 35.88 | 78.40 | 33.73 | 59.65 | 20. -0.2 | 46.31 | 43.96 | 11.92 | 78.80 | 22.07 | 40.61 |
| 10. -0.2 | 95.38 | 56.54 | 34.64 | 78.40 | 34.03 | 59.94 | 10. -0.2 | 45.88 | 44.42 | 12.21 | 78.60 | 21.66 | 40.65 |
| 20.0.2 | 95.44 | 58.50 | 36.75 | 78.40 | 33.48 | 60.51 | 20.0.2 | 51.38 | 46.88 | 10.28 | 78.20 | 22.95 | 41.94 |
| 10.0.2 | 95.50 | 54.08 | 35.85 | 78.40 | 32.64 | 60.18 | 10.0.2 | 49.25 | 47.79 | 11.66 | 78.20 | 23.94 | 42.17 |
| baseline | 95.25 | 57.71 | 36.56 | 77.40 | 31.92 | 59.77 | baseline | 47.88 | 46.54 | 11.88 | 78.00 | 22.43 | 41.46 |
| **Ministral-8B-Instruct-2410** | | | | | | | **Ministral-8B-Instruct-2410** | | | | | | |
| 20. -0.2 | 94.62 | 61.79 | 31.31 | 77.20 | 33.59 | 59.70 | 20. -0.2 | 30.56 | 46.17 | 0.00 | 80.60 | 21.41 | 35.75 |
| 10. -0.2 | 94.56 | 62.17 | 29.74 | 78.80 | 33.17 | 59.69 | 10. -0.2 | 30.06 | 46.04 | 0.00 | 80.00 | 20.62 | 35.34 |
| 20.0.2 | 93.81 | 63.46 | 38.86 | 79.40 | 35.08 | 60.91 | 20.0.2 | 30.88 | 47.12 | 0.00 | 81.80 | 19.98 | 35.96 |
| 10.0.2 | 93.81 | 63.67 | 36.69 | 79.60 | 28.74 | 60.54 | 10.0.2 | 31.19 | 46.79 | 0.00 | 81.80 | 19.49 | 36.05 |
| baseline | 94.75 | 63.58 | 33.68 | 79.00 | 31.56 | 60.51 | baseline | 30.62 | 47.17 | 0.00 | 81.40 | 21.40 | 36.12 |

Table 2: Results of HELMET benchmark under 32k (left) and 64k (right) context. Green indicates better than the baseline; red indicates worse than the baseline.