

# The 63rd Annual Meeting of the Association for Computational Linguistics

Vienna, Austria

July 27–August 1st, 2025



## Efficient OpAmp Adaptation for Zoom Attention to Golden Contexts

Haoyuan Wu ♠† Rui Ming ♠† Haisheng Zheng ♡ Zhuolun He ♠♣ Bei Yu ♠

♠ The Chinese University of Hong Kong

♡ Shanghai Artificial Intelligent Laboratory

♣ ChatEDA Technology

Jul. 27, 2025



上海人工智能实验室  
Shanghai Artificial Intelligence Laboratory

# Outline

1 Introduction

2 Method

3 Experiments

# Introduction

# Background

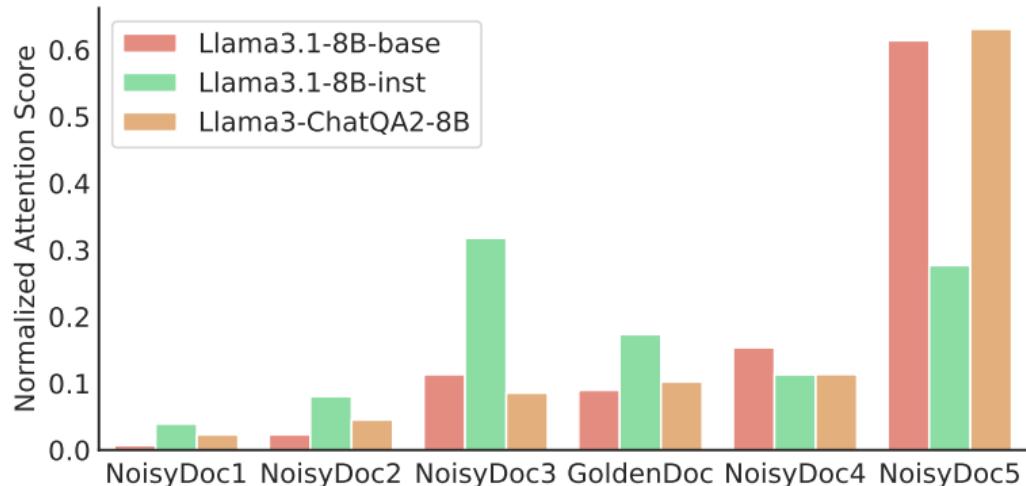


Figure: Normalized Attention Score.

- Transformers often miss the golden document in a noisy context.

# Contribution

- We introduce the *OpAmp Adaptation* for zoom attention to the most relevant context in noisy contexts.
- Implement *OpAmp Adaptation* with adapters, which are fine-tuned with our noisy context dataset, achieving significant improvements.
- Develop *OpAmp Models* with our *OpAmp Adaptation* method, surpassing current SOTA LLMs in various noisy-context benchmarks.

# Method

# Operational Amplifier

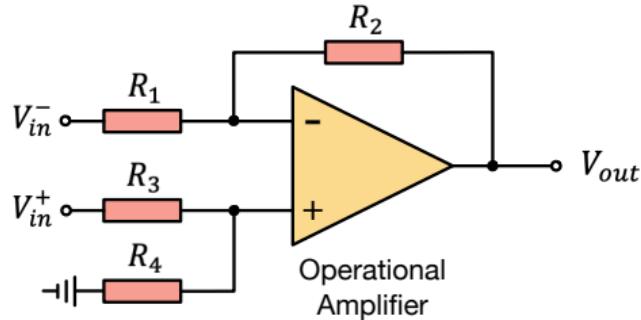


Figure: Visualization of *Operational Amplifier*.

- The *Operational Amplifier* with two input voltages  $V_{\text{in}}^+$  and  $V_{\text{in}}^-$ . The CMRR  $\mathcal{K}$  is controlled by resistances  $R_1, R_2, R_3, R_4$ .

$$\begin{aligned}V_{\text{out}} &= V_{\text{in}}^+ \cdot \left( \frac{R_4}{R_3 + R_4} \cdot \frac{R_1 + R_2}{R_1} \right) - V_{\text{in}}^- \cdot \frac{R_2}{R_1} \\&= A_d(V_{\text{in}}^+ - V_{\text{in}}^-) + \frac{A_c}{2}(V_{\text{in}}^+ + V_{\text{in}}^-).\end{aligned}\tag{1}$$

# Common-Mode Rejection Ratio

The *Common-Mode Rejection Ratio (CMRR)* is defined as the ratio of the differential gain to the common-mode gain:

$$\mathcal{K} = \frac{A_d}{A_c}. \quad (2)$$

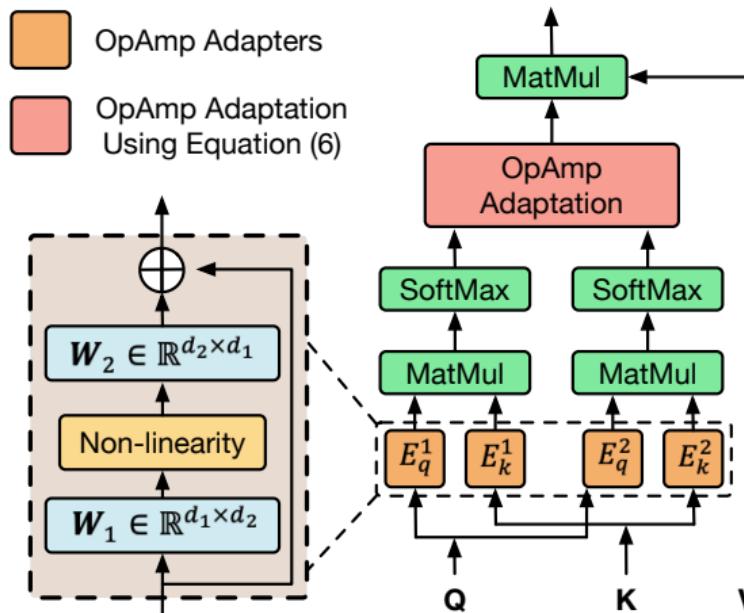
- CMRR is a very important parameter, which indicates the suppression and attenuation of the common-mode signal passing through the amplifier.

Inspired by the operational amplifier, we propose the *OpAmp Adaptation*, which modifies the original attention mechanism into the *OpAmp Attention* mechanism.

$$\bar{\mathbf{M}} = A_d(\mathbf{M}^+ - \mathbf{M}^-) + \frac{A_c}{2}(\mathbf{M}^+ + \mathbf{M}^-), \quad (3)$$

where  $\bar{\mathbf{M}}$  is the denoised attention matrix via OpAmp adaptation,  $\mathbf{M}^+$  and  $\mathbf{M}^-$  are formulated through adapters.

# Architecture Design



**Figure:** Overview of the OpAmp adaptation using Equation (3) with adapters.

At the onset of training, we employ zero initialization to promote identity mapping. Specifically,  $\mathbf{W}_2$  is initialized to zero to guarantee that  $E_j^i(\mathbf{x}) = \mathbf{x}$ .

Furthermore, to prevent any disruption to the original  $\mathbf{M}$  during the initial phase of training, we set  $A_c = 1$  and regulate  $\mathcal{K} = \frac{A_d}{A_c}$  by adjusting the values of  $A_d$ .

As a result, at the initial stage, Equation (3) reduces to:

$$\begin{aligned}\bar{\mathbf{M}} &= A_d \cdot (\mathbf{M} - \mathbf{M}) + \frac{A_c}{2} \cdot (\mathbf{M} + \mathbf{M}), \\ &= A_d \cdot 0 + \frac{A_c}{2} \cdot 2\mathbf{M} = \mathbf{M},\end{aligned}\tag{4}$$

# Experiments

# Experiments Setting

	LongCite-45k <sup>[1]</sup>	Neural-Bridge-RAG <sup>[2]</sup>	Tulu3-SFT-Mix <sup>[3]</sup>
NCFT	30k	20k	450k

**Table:** Training Dataset Composition and Proportions.

<sup>[1]</sup> Jiajie Zhang et al. (2024). “LongCite: Enabling LLMs to Generate Fine-grained Citations in Long-context QA”. In: *arXiv preprint*.

<sup>[2]</sup> Neural Bridge AI (2024). *Retrieval-Augmented Generation (RAG) Dataset 12000*. URL: <https://huggingface.co/datasets/neural-bridge/rag-dataset-12000>.

<sup>[3]</sup> Nathan Lambert et al. (2024). “Tülu 3: Pushing Frontiers in Open Language Model Post-Training”. In: *arXiv preprint*.

# Experiments Setting

Benchmark	Source	Max Length	Metric	# Data
<i>Long-Context QA</i>				
NarrativeQA	Literature, Film	64K	EM	1009
Qasper	Science	8K	PM	200
QuALITY	Literature	8K	Acc.	1065
LooGLE	Science	32K	EM	1427
<i>Multi-Hop QA</i>				
HotpotQA	Wikipedia	16K	EM	200
MuSiQue	Wikipedia	16K	EM	200
MultiHopRAG	News	8K	EM	2255
<i>Noisy-RAG QA</i>				
CoQA	Multi-field	4K	EM	500
QuAC	Wikipedia	4K	PM	996
QReCC	Multi-field	4K	PM	643

**Table:** An Overview of the Dataset Statistics for the *Noisy-Context* Benchmark.

# Evaluation on Noisy-Context Benchmarks

	Qwen2.5 OpAmp-72B	Llama3 ChatQA2-70B	Qwen2.5 72B inst	Llama3.3 70B inst	DeepSeek V3	GPT-4o 0806
LooGLE	<b>66.3</b>	59.1	64.9	63.0	63.4	62.7
NarrativeQA	<b>61.7</b>	59.8	60.2	61.5	60.5	61.5
MultiHopRAG	<b>89.6</b>	78.2	89.2	83.7	88.6	87.7
HotpotQA	<b>77.5</b>	70.5	76.0	74.5	77.0	<b>77.5</b>
MuSiQue	48.0	39.0	44.0	47.5	52.5	<b>53.0</b>
CoQA	<b>92.4</b>	80.2	85.8	88.2	88.4	88.6

**Table:** Performance of Qwen2.5-OpAmp-72B on various *Noisy Context Benchmarks*.

# Evaluation on Noisy-Context Benchmarks

	Llama3.1 OpAmp-8B	Llama3 ChatQA2-8B	Mistral 7B inst-v0.3	Llama3.1 8B inst	Qwen2.5 7B inst
LooGLE	<b>56.6</b>	50.7	51.6	56.1	53.8
NarrativeQA	<b>57.4</b>	53.1	44.7	55.9	47.7
MultiHopRAG	<b>70.5</b>	50.9	69.5	63.9	66.9
HotpotQA	<b>61.0</b>	56.5	58.0	58.5	59.5
MuSiQue	<b>35.0</b>	23.0	28.5	29.5	31.5
CoQA	<b>85.4</b>	78.2	70.6	82.2	84.2

**Table:** Performance of Llama3.1-OpAmp-8B on various *Noisy Context Benchmarks*.

# Ablation Studies: Common-Mode Rejection Ratio

Method	$\mathcal{K}$	Avg.	Qasper (PM)	LooGLE (EM)	NarrativeQA (EM)	QuALITY (Acc.)	MultiHopRAG (EM)	HotpotQA (EM)	MuSiQue (EM)	CoQA (EM)	QuAC (PM)	QReCC (PM)
QLoRA	-	52.4	38.9	53.1	55.7	76.1	68.4	56.5	31.5	83.6	25.2	35.4
	1	54.1 (+1.7)	40.8	56.0	56.4	79.2	68.5	57.5	32.5	<b>85.8</b>	26.1	38.3
OpAmp	5	54.3 (+1.9)	41.2	56.5	56.9	77.8	69.5	<b>62.0</b>	31.5	84.6	25.5	37.1
Adapter	10	<b>55.4 (+3.0)</b>	<b>43.1</b>	<b>56.6</b>	<b>57.4</b>	79.0	70.5	61.0	<b>35.0</b>	85.4	<b>26.5</b>	<b>39.8</b>
	20	54.4 (+2.0)	41.5	55.4	56.4	<b>79.3</b>	<b>71.4</b>	59.0	33.0	84.0	26.2	37.0

Table: Ablation Studies on various *Noisy Context Benchmarks* using Llama3.1-8B-Base as the Base Model.

# Ablation Studies: Noise Ratios

		CoQA (EM)			QuAC (PM)			QReCC (PM)		
Noise Ratio		0.0	0.8	0.9	0.0	0.8	0.9	0.0	0.8	0.9
QLoRA		89.8	85.4	83.6	27.5	26.1	25.2	36.5	36.4	35.4
OpAmp Adapter	1	90.4	85.6	<b>85.8</b>	28.5	26.2	26.1	39.4	39.1	38.3
	5	90.0	85.6	84.6	27.5	26.7	25.5	38.2	37.3	37.1
	10	91.2	<b>88.0</b>	85.4	28.5	26.5	<b>26.5</b>	<b>40.8</b>	<b>39.8</b>	<b>39.8</b>
	20	<b>91.8</b>	86.6	84.0	<b>28.6</b>	<b>28.0</b>	26.2	38.5	38.1	37.0

**Table:** Ablation Studies on *Different Noise Ratios* using Llama3.1-8B-Base as the Base Model.

# Ablation Studies: Hallucination

Method	$\kappa$	FaithEval				Avg.
		Inconsistent (EM)	Unanswerable (EM)	Counterfactual (EM)		
QLoRA	-	24.1	46.1	71.6		47.3
	1	<b>45.5</b>	53.1	<b>76.3</b>		<b>58.3 (+11.0)</b>
OpAmp	5	42.1	<b>53.7</b>	75.9		57.2 (+9.90)
Adapter	10	45.3	53.0	75.1		57.8 (+10.5)
	20	22.3	58.8	73.8		51.6 (+4.30)

Table: Ablation Studies on *FaithEval* using Llama3.1-8B-Base as the Base Model.

# Visualization of Attention

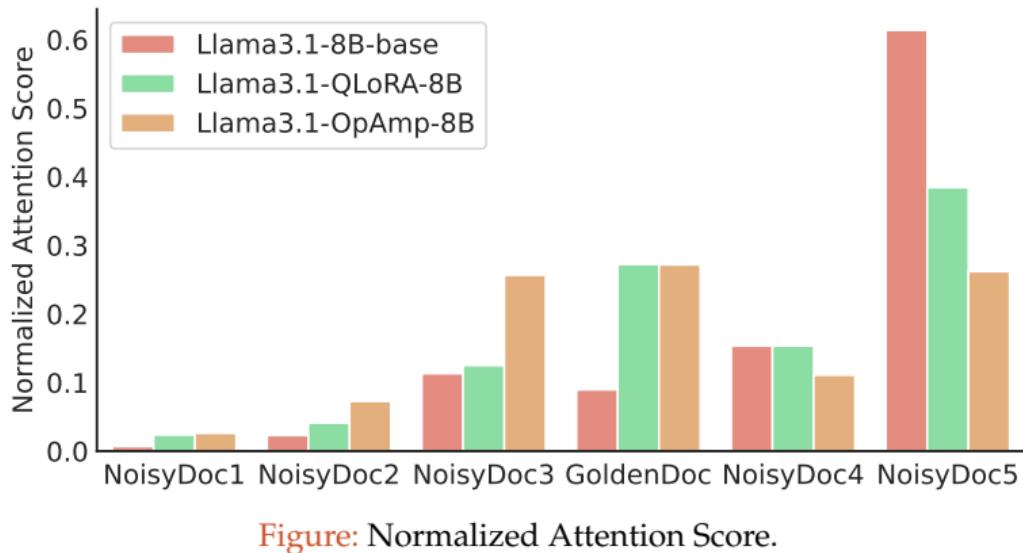


Figure: Normalized Attention Score.

- Our *OpAmp Model* demonstrates significant *Attention Denoise* capability compared to the base model and QLoRA model.

# Visualization of Attention

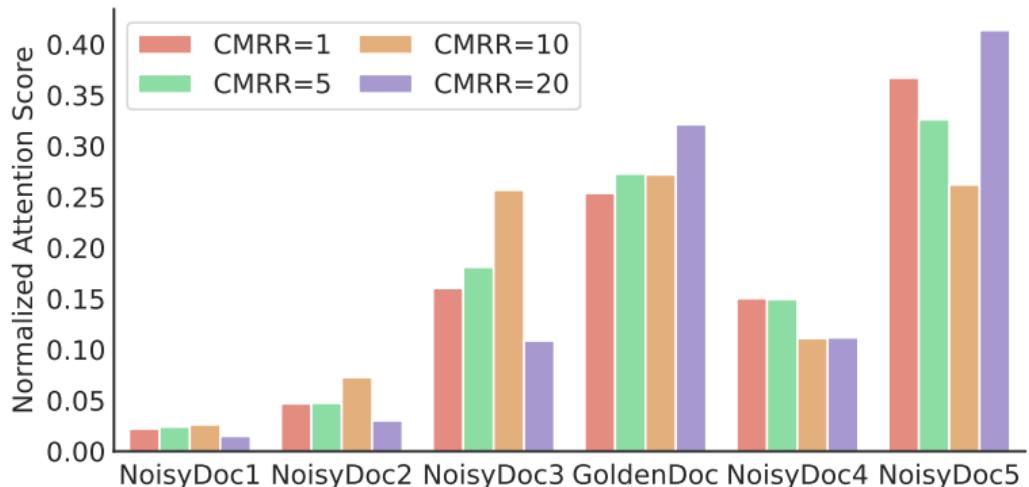


Figure: Normalized Attention Score with *Different Values of K Utilizing for OpAmp Adaptation.*

**THANK YOU!**