## **Query-Focused Retrieval Heads Improve Long-Context Reasoning and Re-ranking**

Anonymous ACL submission

#### Abstract

Recent work has identified retrieval heads (Wu et al., 2025b), a subset of attention heads responsible for retrieving salient information in long-context language models (LMs), as measured by their copy-paste behavior in Needlein-a-Haystack tasks. In this paper, we introduce **QRHEAD** (Query-Focused Retrieval Head), an improved set of attention heads that significantly enhance retrieval from long contexts. We identify QRHEAD by aggregating attention scores with respect to the input query, using real-world tasks such as long-context QA. We further introduce QRRETRIEVER, an efficient and effective retriever that uses the accumulated attention mass of QRHEAD as retrieval scores. We evaluate QRRETRIEVER as a reranker on the BEIR benchmark and find that it achieves strong zero-shot performance, outperforming other LLM-based re-rankers such as RankGPT. We also use QRRETRIEVER for long-context reasoning by selecting the most relevant parts with the highest retrieval scores. On long-context, multi-hop reasoning tasks LongMemEval and CLIPPER, this yields over 10% performance gains over full context and outperforms strong dense retrievers. Further analysis shows that both the query-context attention scoring and task difficulty are crucial for identifying QRHEAD with strong downstream utility. Overall, our work contributes a generalpurpose retriever and offers interpretability insights into the long-context capabilities of LMs.

#### 1 Introduction

004 005

007

012

015

017

027

028

034

Retrieving salient information from long contexts serves as a foundation for language models (LMs), enabling a wide range of downstream applications, such as long document understanding and passage re-ranking. Prior work has identified a subset of attention heads in transformers (Vaswani et al., 2017) that are responsible for retrieving relevant information, known as *retrieval heads* (Wu et al., 2025b).



Figure 1: **Top**: Masking the top 32 original retrieval heads (Wu et al., 2025b) of Llama-3.1-8B-Instruct. **Bot-tom**: Masking the top 32 QRHeads which has a more pronounced impact on Needle-in-a-Haystack.

043

044

047

053

055

056

060

061

062

063

064

However, these retrieval heads are identified based on the frequency of their copy-paste operations in a simple synthetic task—Needle-in-a-Haystack (NIAH; Kamradt, 2024). Although they exhibit some significance on downstream tasks, such as extractive question answering, we argue that the copy-paste objective and synthetic data used to identify them are misaligned with how language models retrieve pertinent information in realworld settings.

To this end, we propose a more effective approach for identifying retrieval heads and introduce QRHEAD, a distinct subset of attention heads whose attention mass plays a more critical role in retrieving relevant information from long contexts. Compared to original retrieval heads, our method incorporates two key changes: (1) a query-context scoring function that measures attention mass allocated to pertinent context spans with respect to an input query, and (2) the use of more natural data from real-world tasks, such as question answering over long texts. Our method only requires a small

1

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

113

114

115

116

065amount of data to be effective. As shown in Fig-066ure 1, we detect QRHEAD using 70 examples from067a natural long-context QA task, LongMemEval,068and find masking them out results in more severe069degradation in NIAH compared to original retrieval070heads detected from in-domain data.

072

094

Furthermore, we build QRRETRIEVER on top of QRHEAD as a general-purpose retriever for improving LMs on diverse long-context downstream applications. Given a query and a set of passages (e.g., a claim and a book consisting of multiple chapters), QRRETRIEVER scores each passage using the aggregated attention mass from the QR-HEAD of a language model, and returns the topranked passages. Using QRHEADS, we instantiate QRRETRIEVER with multiple LMs of different scales (3B–70B) and families (Llama-3.1, Llama-3.2, and Qwen).

We evaluate QRRETRIEVER as a re-ranker on the standard BEIR benchmark (Thakur et al., 2021). It exhibits strong zero-shot performance across diverse domains and outperforms other LLMbased re-rankers, such as RankGPT (Sun et al., 2024). We further evaluate QRRETRIEVER on two long-context, multi-hop reasoning tasks: Long-MemEval (Wu et al., 2025a) and CLIPPER (Pham et al., 2025). Using QRRETRIEVER to select top-ranked documents yields substantial improvements in retrieval recall and downstream task performance. For example, with Llama-3.1-8B-Instruct, QRRETRIEVER outperforms dense retrievers and improves performance by over 10% on both datasets, compared to full-context generation.

098 Finally, we provide extensive analyses of the effectiveness of QRHEAD. First, using QRHEAD outperforms both full attention heads and original retrieval heads. Second, QRHEAD generalizes 101 across input lengths—the heads identified at 32K 102 tokens transfer well to tasks with 128K context lengths. Lastly, we show that both key modifica-104 tions-our query-focused scoring objective and the 105 use of natural data-contribute to the improved 106 downstream performance of QRHEAD over origi-108 nal retrieval heads. Together, these findings highlight the practicality and robustness of QRHEAD as 109 a foundation for long-context retrieval and suggest 110 opportunities for further exploration of retrieval 111 mechanisms in language models. 112

### 2 Background: Retrieval Heads

Retrieval heads are a specialized subset of attention heads that are pivotal for extracting relevant information from long input contexts.

**Original retrieval heads.** Wu et al. (2025b) first discovered a set of retrieval heads that exhibit copypaste behavior during decoding-effectively copying tokens from the long context input context into the generated output. Their retrieval head detection method roots from the Needle-in-a-Haystack test (NIAH) with a triple (C, q, a) of context, question, and answer: the answer span a (the "needle") is embedded within a long context sequence  $C = d_1...d_N$  where  $d_1, ..., d_N$  are N irrelevant sequences (the "haystack"). The LM is tasked to generate an answer to q based on the provided context. Successful generation of a demonstrates effective copy-paste behavior by extracting a from the haystack and copying it over to the output. To quantify this behavior, the retrieval score of an attention head h is defined as the fraction of tokens copied from a by the head h during decoding:

Retrieval\_Score(h) = 
$$\frac{|g_h \cap a|}{|a|}$$
, (1)

where  $g_h$  denotes the set of tokens copied by head h to the output. Attention heads with the highest retrieval scores are selected as retrieval heads.

**Shortcomings.** The scoring mechanism described above focuses only on attention heads that perform strict copy-paste operations, potentially missing heads involved in semantic-based retrieval, such as paraphrasing or reasoning over relevant context. Moreover, recent work has shown that heads identified through copy-paste metrics exhibit limited cross-domain generalizability (Zhao et al., 2024). This suggests that the simplified formulation may not fully capture the complexity of incontext retrieval behavior in LLMs and has limited relevance for downstream applications.

### **3** QRHEAD: Identifying Query-Focused Retrieval Heads

In this section, we introduce a new approach for detecting retrieval heads that significantly improves upon prior retrieval head detection. For clarity, we refer to our heads as Query-Focused Retrieval Heads (QRHEAD) and the original retrieval heads as RETHEAD. Our approach introduces two key improvements. First, we propose a query-focused



Figure 2: Comparison between Retrieval Heads (Wu et al., 2025b) and QRHEADS (Ours).

retrieval score (QRscore), which captures querycontext attention rather than relying solely on copypaste behavior (§3.1). Second, we leverage realistic tasks that require in-context retrieval to identify
effective heads (§3.2). Lastly, we present a comparison between QRHEAD and RETHEAD (§3.3).

166

181

182

183

184

187

188

190

192

193

#### Task formulation: LMs for in-context retrieval.

Our study focuses on the task of in-context retrieval with LMs, i.e., identifying relevant information from given context. Formally, let  $\Sigma$  denote the vocabulary. Given an input query  $q \in \Sigma^*$  and a con-170 text  $D \in \Sigma^*$ , the objective is to retrieve the most 171 relevant information from the context with respect 172 to q, denoted as  $D_{[q]} \subseteq D$ . Typically, the context D consists of a sequence of passages (or chunks), 174 represented as  $D = \{d_1, d_2, \dots, d_N\}$ . With both q and D jointly fed into an LM as input, we assign a score  $\mathcal{R}(q, d_i)$  to each passage  $d_i$  with respect to 177 q. We measure the effectiveness of the retriever by evaluating whether the top-scored passages align with the ground-truth relevant documents  $D_{[a]}^{*-1}$ .

# 3.1 Scoring Heads with Query-Context Attention

Instead of scoring attention heads based on their activations in copy-paste operations, we propose to evaluate them based on their effectiveness in realistic in-context retrieval tasks. This offers a more general and realistic measure of retrieval capability, as it captures semantic relevance rather than relying solely on verbatim copying.

**Query-focused retrieval score (QRscore).** We use QRscore as a measure of the retrieval capability of an attention head in response to a specific query. Formally, let  $h \in \mathcal{H}$  be an attention head

within the language model, and let  $A_h$  denote the attention weights (post-softmax) of head h over a query prompt  $\{D, q\}$ , such as a prompt with a book followed by a question over its contents. The query-focused attention scores of head h towards a document  $d_i$  is calculated as follows:

194

195

196

197

198

199

201

202

203

204

205

206

207

210

211

212

213

214

215

216

217

218

219

220

$$QRscore_h(q, d_i) = \frac{1}{|q|} \sum_{t_q \in q} \sum_{t_d \in d_i} A_h^{t_q \to t_d}$$
(2)

where  $t_q$  denotes tokens in the query q,  $t_d$  represents tokens in the document  $d_i$ , and  $A_h^{t_q \to t_d}$  is the attention weight of h from  $t_q$  to  $t_d$ . This formulation quantifies the degree to which head h focuses on document  $d_i$  in response to q. Lastly, we aggregate the scores for all documents  $d_i$  within the gold document set  $D_{[q]}^*$ , resulting in the final QRscore for head h with respect to the query q:

$$QRscore_h(q) = \frac{1}{|q|} \sum_{d_i \in D^*} \sum_{t_q \in q} \sum_{t_d \in d_i} A_h^{t_q \to t_d}$$
(3)

#### 3.2 Detecting QRHEAD on Real-World Tasks

With the QRscore defined in Eq. 3, we can now quantify the retrieval capabilities of each attention head over a given set of documents in response to a query. To achieve this, we leverage a head detection dataset  $\mathcal{T} = \{(q, D, D_{[q]}^*)\}$ , which consists of a query q, a set of candidate documents D, and the corresponding gold documents  $D_{[q]}^*$ . Notably, our approach does not require explicit answers to the queries—only the annotations of the gold document. Using this detection dataset  $\mathcal{T}$ , we compute the empirical effectiveness of an attention head h for retrieval as follows:

$$\operatorname{QRscore}_{h,\mathcal{T}} = \frac{1}{|\mathcal{T}|} \sum_{(q,D,D^*)\in\mathcal{T}} \operatorname{QRscore}_h(q) \quad (4)$$

As shown in Figure 2, instead of synthetic needlein-a-haystack task (NIAH) (Kamradt, 2023), we 225

<sup>&</sup>lt;sup>1</sup>We note NIAH task can also be viewed as a special case of this formulation, where the gold document set only contains one document (the needle).

use more realistic in-context retrieval task for head detection (e.g., claim verification over books). We argue that more natural and realistic distractors provide more effective supervision that allows identifying heads that are better at differentiating relevant contexts from distracting context. We also note that even a small amount (< 100) of realistic data points can be sufficient, allowing us to find QRHEADS heads that contribute to improved downstream performance (see §6.1 for detailed results).

226

227

235

240

241

242

243

244

245

246

247

248

251

255

260

261

262

263

267

268

269

270

271

#### **3.3 Comparing QRHEAD and RETHEAD**

We have demonstrated our method for detecting QRHEAD, which improves upon RETHEAD. Here, we compare the two head sets within the same model, using Llama-3.1-8B-Instruct (Llama-3 Team, 2024) as a case study, where we identify the top QRHEAD and top RETHEAD

First, following the analysis setup of Wu et al. (2025b), we measure the impact of pruning by the performance drop on NIAH test. Specifically, we prune the top 32 heads (roughly 3% of all attention heads in LLaMA-3.1-8B), following the commonly reported 5% sparsity level of retrieval heads in Wu et al. (2025b); Zhao et al. (2024). As shown in Figure 1, pruning the top 32 QRHEAD results in near-complete failure on the NIAH performance, whereas pruning the top 32 RETHEAD yields a much smaller performance decline.<sup>2</sup> In addition, we find substantial divergence between the two sets. Among the top 32 and top 64 heads, only 8 and 32 overlap, respectively. This less than 25% overlap in the top 32 highlights the distinct roles of **ORHEAD and RETHEAD.** 

#### 4 Building General-Purpose Retriever with QRHEAD

In this section, we describe how the detected QRHEAD can be used in downstream applications. Specifically, we find the attention mass of QRHEAD provides highly reliable signals for incontext retrieval.

#### 4.1 The Method

Given a selected set of QRHEAD  $\mathcal{H}^{\text{select}}$ , a query q, and a collection of passages D, we compute the retrieval score for each passage  $d_i$  by aggregating the QRscore across all heads in  $\mathcal{H}^{\text{select}}$ :

$$\mathcal{R}(q, d_i) = \frac{1}{|\mathcal{H}^{\text{select}}|} \sum_{h \in \mathcal{H}^{\text{select}}} \text{QRscore}_h(q, d_i). \quad (5)$$

Passages are then weighed using their retrieval scores. We call our retrieval system QRRE-TRIEVER. It offers several advantages: (1) *Modelagnostic:* compatible with any transformer-based LMs without modification, (2) *Efficient:* leverages attention patterns to process long contexts simultaneously without expensive generation or pairwise comparisons, (3) *High-performing:* outperforms various baselines, as shown in §5.

272

273

274

275

276

277

278

279

281

282

283

284

287

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

**Implementation details.** To mitigate intrinsic biases in LMs' attention weights, we adopt the score calibration method proposed by Chen et al. (2025). Instead of directly using  $R(q, d_i)$  as the score, we additionally compute baseline score,  $R(q_{null}, d_i)$ , using a context-free null query  $q_{null}$  ("N/A"). We use calibrated the score  $R(q, d_i) - R(q_{null}, d_i)$  as the final retriever score.

#### 4.2 Applications

**Long-context reasoning.** Language models, including long-context language models, often struggle with performance degradation when processing long contexts (Yen et al., 2025; Ye et al., 2025; Liu et al., 2024). To address this, we integrate QRRE-TRIEVER within a retrieval-augmented generation (RAG) framework. Given a long-context input and a query, we segment the input into smaller chunks and use QRRETRIEVER to score and subsequently extract the most relevant ones. The extracted contexts are concatenated to create a reduced context that is then given to the LM for generation.

**Passage re-ranking.** Text retrieval powers many retrieval-augmented downstream applications (Lewis et al., 2020). A critical component in the retrieval pipeline is the re-ranker, which reorders the passages returned by a first-stage retriever to enhance top passage relevance (Nogueira and Cho, 2020; Ma et al., 2024). QRRETRIEVER can naturally be used as a re-ranker as part of any retrieval pipeline without any fine-tuning by simply concatenating the retrieved passages in the input and scoring their relevance directly.

#### **5** Experiments

We evaluate QRRETRIEVER on two tasks: longcontext reasoning (§5.2) and re-ranking (§5.3).

#### 5.1 Base Models and Baselines

Base LMs.We experiment with open-weight,317instruction-tuned LMs from two families across318

<sup>&</sup>lt;sup>2</sup>See Appendix B for results on Qwen-2.5-7B-Instruct.

	LongMemEval				CLIPPER			
Retriever	$\begin{array}{l} \mathbf{R} \mathbf{E} \mathbf{T} \mathbf{R} \\ \mathbf{R} \mathbf{E} \mathbf{C} \mathbf{A} \\ \mathbf{k} = 5 \end{array}$	AIEVAL LL@K k = 10	END-1 Perfoi Top-5	ro-End rmance Top-10	<b>RETR</b> <b>RECA</b> k = 3	IEVAL     LL@K     k = 5	END-T Perfoi Top-3	<b>CO-END</b> RMANCE Top-5
Base LM: Llama-3.2-3B-Instruct Full context BM25 Contriever Stella	57.5 62.7 63.9	67.5 79.2 77.6	2: 46.1 <b>48.6</b> 44.9	8.1 44.9 46.5 47.7	74.6 60.2 83.3	83.7 78.9 90.0	20.0 12.6 21.3	<b>5.2</b> 22.8 18.4 25.1
RankGPT RankGPT <sup>Bubble</sup> ICR QRRETRIEVER (Ours)	1.8 2.1 68.8 <b>76.5</b>	3.4 3.8 78.8 <b>86.1</b>	23.5 24.0 46.5 47.4	23.3 24.4 45.8 <b>48.6</b>	16.8 17.0 75.7 <b>87.9</b>	27.3 27.4 86.1 <b>94.6</b>	3.6 3.8 22.4 24.0	8.8 8.8 23.7 24.4
Base LM: Llama-3.1-8B-Instruct Full context BM25 Contriever Stella	57.5 62.7 63.9	67.5 79.2 77.6	48.8 52.6 50.9	6.5 50.9 55.4 58.4	74.6 60.2 83.3	83.7 78.9 90.0	37.9 28.2 38.8	1.3 37.9 31.1 39.6
RankGPT RankGPT <sup>Bubble</sup> ICR QRRETRIEVER (Ours)	2.1 8.3 78.2 <b>85.6</b>	4.0 9.0 85.3 <b>91.8</b>	26.7 28.1 58.4 <b>59.5</b>	24.2 27.0 58.1 <b>60.2</b>	30.0 36.7 89.9 <b>93.7</b>	39.4 44.3 95.4 <b>97.0</b>	15.9 19.7 43.8 <b>46.5</b>	19.4 20.4 42.3 <b>44.4</b>
Base LM: Llama-3.1-70B-Instruct Full context BM25 Contriever Stella	57.5 62.7 63.9	67.5 79.2 77.6	34 52.8 53.7 56.3	4.2 53.0 60.5 62.3	74.6 60.2 83.3	83.7 78.9 90.0	60.1 38.5 65.9	3.9 66.5 49.7 71.2
RankGPT RankGPT <sup>Bubble</sup> ICR QRRETRIEVER (Ours)	1.8 47.9 45.6 <b>77.5</b>	3.5 49.0 58.2 <b>88.3</b>	21.2 44.0 43.0 <b>64.2</b>	27.4 42.6 48.4 <b>63.3</b>	57.0 74.3 88.3 <b>95.5</b>	63.4 78.8 94.2 <b>98.2</b>	44.7 58.4 71.0 <b>76.7</b>	50.4 61.5 73.3 <b>74.1</b>

Table 1: Results on LongMemEval and CLIPPER. The base model denotes the LM used for both the retriever and end-to-end generation. QRHEADS used for CLIPPER are found through using LongMemEval.

different sizes, including Llama-3.2 (3B), Llama-3.1 (8B and 70B) of Llama family (Llama-3 Team, 2024), and Qwen2.5 (7B) of Qwen family (Yang et al., 2024). With QRRETRIEVER, we use 16 heads for models with fewer than 10B parameters, and 32 heads for LLaMA-3.1-70B. This corresponds to approximately 1–2% of the total attention heads, given the sparsity of retrieval heads.

319

320

321

322

326

327

331

332

335

Baselines. We compare our methods against several strong baselines. Following Wu et al. (2025a), we compare against dense retrievers, including Contriever (Izacard et al., 2022) and 1.5B Stella V5 (Zhang et al., 2025), two popular strong dense retrievers. For Contriever, we truncate the input to 512 tokens according to its maximum context length. We also compare against existing LLM-based re-rankers, including:

RankGPT (Sun et al., 2024) is a generative re-ranker that instructs LLMs to output the ranking order of a given set of documents based on a query. We experiment with two variants of RankGPT: (1) RankGPT without sliding window, which directly inputs all documents into the model prompt simultaneously, and (2) RankGPT with sliding window
(RankGPT<sup>Bubble</sup>), which leverages bubble sort to rank smaller subsets of documents incrementally.

• **In-Context-Reranker** (ICR; Chen et al., 2025) is a re-ranker that also leverages the attention for relevance scoring. ICR uses full attention heads for scoring relevace, whereas we only use the attention weights of selected QRHEADS.

346

347

348

350

351

352

353

354

355

356

357

359

360

361

362

363

365

366

367

368

369

370

371

#### 5.2 Long-Context Multi-Hop Reasoning

**Datasets.** We use 1) **LongMemEval** (Wu et al., 2025a), which evaluates the long-term memory capabilities of LLM-driven chat assistants, and 2) **CLIPPER** (Pham et al., 2025), which evaluate claim-verification over books. Both datasets feature long-contexts (90K to 120K) and require multihop reasoning over several pieces of evidences. We segment each dataset according to its natural structure (e.g., message in multi-turn conversation or chapters in a book). For evaluation, we measure retrieval performance using recall and assess downstream task performance with accuracy. Please refer to Appendix A for more details.

**Data for head detection.** We detect QRHEAD using a small subset of single-hop data from Long-MemEval, specifically the single-session-user subset consisting of 70 examples, which we exclude from downstream evaluation. We use the set of heads for both LongMemEval and CLIPPER, testing generalization to multi-hop reasoning.

	NQ	COVID	NFCorpus	FiQA	Scifact	Scidocs	FEVER	Climate	DBPedia	Robust04	News	Avg
BM25	30.5	59.5	32.2	23.6	67.9	14.9	65.1	16.5	31.8	40.7	39.5	38.4
	Base LM: Llama-3.2-3B-Instruct											
RankGPT	30.0	59.5	32.2	23.6	67.9	14.9	65.9	17.1	31.8	40.7	39.5	38.5
RankGPT <sup>Bubble</sup>	33.2	61.8	32.0	22.4	66.1	14.8	65.8	17.1	34.8	40.5	40.2	39.0
ICR	48.0	71.3	33.6	31.6	73.1	16.1	83.8	23.8	34.8	47.1	44.3	46.1
QRRETRIEVER (Ours)	<b>54.0</b>	77.2	<b>35.0</b>	<b>35.2</b>	<b>74.5</b>	<b>17.1</b>	<b>84.8</b>	<b>24.6</b>	<b>35.9</b>	<b>49.8</b>	<b>45.0</b>	<b>48.5</b>
					Base LM	l: Llama-3	.1-8B-Instr	uct				
RankGPT	30.0	59.5	32.2	23.6	67.9	14.9	65.9	16.8	31.8	40.7	39.5	38.4
RankGPT <sup>Bubble</sup>	53.7	75.5	34.3	31.4	69.3	17.4	67.5	23.8	<b>42.9</b>	47.8	<b>46.2</b>	46.3
ICR	53.7	73.3	34.8	36.1	75.5	17.4	<b>87.1</b> 86.3	<b>25.2</b>	36.9	49.1	44.4	48.5
QRRETRIEVER (Ours)	<b>57.2</b>	<b>76.7</b>	<b>35.3</b>	<b>39.5</b>	<b>76.0</b>	<b>17.9</b>		24.1	36.6	<b>50.7</b>	46.0	<b>49.7</b>
	Base LM: Qwen-2.5-7B-Instruct											
RankGPT	30.0	59.5	32.2	23.6	67.9	14.9	65.9	16.8	31.8	40.7	39.5	38.4
RankGPT <sup>Bubble</sup>	42.7	70.5	<b>34.1</b>	29.5	69.3	<b>16.6</b>	70.5	19.7	<b>37.1</b>	<b>46.4</b>	<b>43.6</b>	43.6
ICR	41.1	65.3	32.6	27.1	70.8	15.1	80.8	19.7	34.9	43.2	40.3	42.8
QRRETRIEVER (Ours)	<b>48.8</b>	<b>67.7</b>	33.1	<b>29.8</b>	<b>70.9</b>	14.2	<b>82.7</b>	<b>19.8</b>	35.5	43.7	40.5	<b>44.2</b>
	Base LM: Llama-3.1-70B-Instruct											
RankGPT	45.4	62.7	33.6	28.6	71.3	16.1	74.2	18.9	37.6	41.3	39.8	42.7
RankGPT <sup>Bubble</sup>	58.4	<b>81.2</b>	<b>36.1</b>	41.0	76.1	<b>20.2</b>	80.0	<b>25.1</b>	<b>45.5</b>	<b>59.0</b>	<b>48.5</b>	<b>51.9</b>
ICR	57.0	71.9	34.0	37.9	73.5	17.5	<b>87.5</b>	22.6	38.3	39.1	39.0	47.1
QRRETRIEVER (Ours)	<b>60.5</b>	74.8	34.7	<b>43.8</b>	<b>76.5</b>	18.5	86.7	23.2	35.9	51.8	44.1	50.1

Table 2: Performance comparison (nDCG@10) on BEIR benchmarks across LMs. QRRETRIEVER generally outperforms other baselines across all models. With Llama-3.1-70B, QRRETRIEVER underperforms RankGPT with (Bubble sort), which requires substantial amount of LLM generation calls.

**QRRETRIEVER achieves strong retrieval performance for long contexts, leading to improved end-to-end performance.** Table 1 demonstrates the strong performance of QRRETRIEVER on both LongMemEval and CLIPPER: it outperforms other baselines regarding both retrieval recall and endto-end performance. For instance, Llama-3.1-8B-Instruct as the base LM, we see end-to-end performance improvements of over 10% on both tasks with Llama-3.1-8B-Instruct.

372

373

374

375

378

379

382**QRRETRIEVER generalizes across domains.**383The fact that QRRETRIEVER outperforming off-384the-shelf dense retrievers (Contriever and Stella)385by a large margin on LongMemEval and CLIPPER.386In particular, none of these methods are trained or387calibrated on CLIPPER. The better performance of388QRRETRIEVER suggests its stronger cross-domain389generalization capabilities than dense retrievers.390Moreover, while QRHEADS are detected using391only the single-hop questions, it also performs well392on the multi-hop questions.

393**QRRETRIEVER scales with the model sizes.**394We note that LM-based re-rankers show incon-395sistent performance patterns across model scales:396RankGPT achieves near-zero retrieval recall with397small models, and retrieval performance of ICR398sees significant degradation when scaling up model399size from 8B to 70B. At the same time, the perfor-

mance of QRRETRIEVER generally improves as the model size scales up.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

#### 5.3 Passage Re-Ranking

To test the general applicability of QRRE-TRIEVER, we evaluate our method on BEIR benchmark (Thakur et al., 2021) consisting of diverse domains. We compare against zero-shot LLM-based re-rankers, RankGPT and ICR.

**Setting.** Our setting largely follows prior work (Chen et al., 2025). We re-rank 200 passages retrieved using BM25, resulting a overall context length ranging from 16K to 64K depending on the average document length of domains. We report the performance on the set of tasks used in Chen et al. (2025), we sub-sampled 512 random questions for each domain for evaluation.

**Data for head detection.** For BEIR, we utilize the 128 (held-out) data points from NQ and use them for on all other domains zero-shot.

Results.Table 2 summarizes the BEIR results,419demonstrating the strong effectiveness of QRRE-<br/>TRIEVER as a general-purpose retriever. For mod-<br/>els under 10B parameters, QRRETRIEVER consis-<br/>tently outperforms other baselines. With LLaMA-<br/>3.1-8B, it achieves an average score of 49.7, out-<br/>performing RankGPT by 3.4 points and ICR by 1.2419

		BEIR <sub>SHUFFLED</sub>	LONGMEM RECALL
Llama-8B	RANDOMHEADS	37.5	59.8
	FULLHEADS	42.8	73.2
	RETRIEVALHEADS	43.4	81.5
	QRHEADS	<b>47.5</b>	<b>85.6</b>
Qwen-7B	Random	19.9	57.2
	Full Heads	22.6	67.1
	Retrieval Heads	27.4	70.7
	QRHeads	<b>31.9</b>	<b>83.2</b>

Table 3: Comparison across head selection strategies. Using QRHEADS substantially outperforms using all heads or using original retrieval heads.

	Model: LLama-3.1-8B-Instruct					
	NQ+	Fever	LongMemEva			
	32K	128K	32K	128K		
ICR	66.7	56.5	85.2	78.2		
QRRETRIEVER <sup>32K</sup>	70.1	63.9	89.2	85.2		
QRRETRIEVER <sup>128K</sup>	68.8	67.2	89.2	85.6		
	Mo	del: Qwen-	2.5-7B-Instruct			
	NQ+	Fever	LongN	IemEval		
	32K	64K	32K	64K		
ICR	40.0	17.4	83.4	67.1		
QRRETRIEVER <sup>32K</sup>	51.9	25.3	90.2	77.9		
QRRETRIEVER <sup>64K</sup>	54.1	29.1	90.1	77.0		

Table 4: Results on short-to-long generalization of QR-HEADS. QRHEADS detected with relative short-context data can be used for retrieval on longer contexts.

points. For the larger LLaMA-3.1-70B model, QR-RETRIEVER significantly surpasses ICR, though it generally lags RankGPT<sup>Bubble</sup> (which require over 200 generation calls). Nevertheless, QRRE-TRIEVER achieves the best performance on several domains, suych as SciFact and FEVER.

#### 6 Analysis

426

497

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

#### 6.1 Impact of Head Selection

We provide further ablation on head selection, the core idea behind QRRETRIEVER. We experiment with different sets of heads, including (1) using our QRHEADS, (2) using all the attention heads (Full), (3) using original retrieval head (Retrieval), and (4) using randomly selected heads (Random). We use 16 heads for all settings. Table 3 presents the retrieval performance on LongMemEval re-ranking performance on BEIR (aggregated across tasks).<sup>3</sup> The performance gaps between different strategies demonstrate the importance of using the right heads for retrieval. Using original retrieval heads is effective, compared to using random heads or full heads. Using our improved QRHEADS consistently outperforms using original retrieval heads.

	Data	BEIR <sub>shuffled</sub> NDCG@10	LongMem Recall
	Ma	del: LLama-3.1	-8B-Inst
QRHEADS	NQ	47.5	83.9
QRHEADS	LMÈ	47.1	85.6
QRHEADS	NIAH	46.8	83.4
RETHEAD	NIAH	43.4	81.5
	M	odel: Qwen-2.5-	7B-Inst
QRHEADS	NQ	31.9	80.2
QRHEADS	LMÈ	32.1	83.2
QRHEADS	NIAH	30.9	79.7
RetHead	NIAH	27.4	70.7

Table 5: Analysis of factors contributing to improved head selection. Applying QRScore (§3.1) on NIAH results in more effective heads than the original retrieval heads. Using QRScore on realistic tasks yields the most effective head selection overall.

#### 6.2 Generalizability Across Lengths

We test the length generalization of QRHEADS: if we detect QRHEADS on relatively short context length (32K), can the heads generalize to longer context lengths (128K)? 449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

We test such short-to-long generalization by controlling the number of documents (messages). The resulting in datasets of different lengths ranging from 32K to 128K tokens. We detect QRHEADS from both short and long datasets and test their performance on re-ranking tasks (using two representative subsets: NQ and Fever) and Long-MemEval. For Qwen-2.5-7B, we set the longer context length to 64K due to its original 32K limit. As shown in Table 4, QRHEAD detected using short-context data can generalize to longer-context settings, though heads detected from longer data generally yield better long-context performance.

#### 6.3 What Contributes to Better Selection?

In §3, we describe two key factors for head detection: using query-context attention objective, and using realistic data. To assess the importance of these factors, we experiment with detecting heads on NIAH using QRScore (§3). As shown in Table 5, applying QRScore on NIAH leads to improved performance compared to using the original retrieval heads detected from the same task. However, using realistic tasks with QRScore yields the best overall performance. These results highlight the importance of both the scoring method and head detection data.

# 6.4 Sensitivity of QRHEADS Detection to Variation in Detection Data

In Section 5.3, we show using a small number of samples from NQ is sufficient to identify effective QRHEADS for BEIR re-ranking tasks. e as-

<sup>&</sup>lt;sup>3</sup>Here, we use BEIR where input documents are randomly shuffled rather than ranked by BM25. This setup allows uniform evaluation of retrieval across the full context.

	Ove Set0	rlap (Toj Set1	BEIR nDCG@10		
		Model:	-8B-Inst		
QRHEADS <sup>Set0</sup>	64	51	51	49.8	
QRHEADS <sup>Set1</sup>	51	64	53	49.7	
QRHEADS <sup>Set2</sup>	51	53	64	49.9	
		Model:	Qwen-2.5	-7B-Inst	
QRHEADS <sup>Set0</sup>	64	50	53	44.2	
QRHEADS <sup>Set1</sup>	50	64	57	44.4	
QRHEADS <sup>Set2</sup>	53	57	64	44.5	

Table 6: Left: Overlap in QRHEADS identified using three disjoint sets of 128 random samples from NQ. Right: BEIR performance (nDCG@10) using QR-HEADS detected from each sample set.

sess the robustness of this head detection process to different random samples of detection set, by experimenting with three disjoint random subsets of NQ, each containing 128 examples. Table 6 presents the overlap among the top-64 heads selected from these subsets and their performance on BEIR benchmark. Across two LLMs from different model families (Llama and Qwen), we observe a high degree of consistency with over 50 heads overlapping among the top 64 across subsets. Furthermore, the downstream performance remains stable across these variations. These results indicate that QRRETRIEVER can be reliably identified using a small sample of data.

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

503

504

505

506

507

510

511

512

513

514

515

516

#### 6.5 Discussion: Retrieval-Generation Gap

Interestingly, we observe that even compact LMs exhibit strong retrieval capabilities despite their limited generation abilities. As shown in Table 1, on CLIPPER, Llama-3.2-3B-Instruct achieves a Recall@10 of 86.1, closely matching the 88.3 score of the much larger LlamA-3.1-70B. However, Llama-3.2-3B only achieves a final end-to-end performance of 24.0, largely lagging 70B's performance of 76.7. We hypothesize that the long-context limitations of compact models stem more from their generation capabilities than from their retrieval abilities, revealing a significant retrieval-generation gap. These findings open up promising future directions. Compact LMs could serve as efficient long-context retrievers, paired with larger models for the actual generation.

#### 7 Related Work

517 LM-based retrieval and re-ranking. LMs are
518 widely used in retrieval, including embedding519 based methods (Muennighoff, 2022; Lee et al.,
520 2021) and generative approaches (Tay et al., 2022;

Cao et al., 2021; Sun et al., 2023). For re-ranking, instruction-tuned LMs been adapted as re-rankers in various ways (Sun et al., 2024; Drozdov et al., 2023; Sachan et al., 2023; Ma et al., 2023; Pradeep et al., 2023), leveraging their generation capabilities. Similar to our approach, recent work has explored using logits (Reddy et al., 2024) or aggregated attention scores (Chen et al., 2025) for re-ranking. In contrast, we identify a specialized set of attention heads responsible for retrieval, offering improved performance and interpretability. 521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

570

Localizing model behavior. Interpretability studies have shown that many core behaviors of LMs, including in-context learning (Olsson et al., 2022; Todd et al., 2024; McDougall et al., 2023) and retrieval (Wu et al., 2025b), can be traced to specialized transformer modules (Meng et al., 2022; Dai et al., 2022; Stolfo et al., 2024). Techniques have been proposed to localize such modules with a small amount of data (Meng et al., 2022; Geiger et al., 2024; Bhaskar et al., 2024), and to intervene on them for control (Li et al., 2023; Yin et al., 2024; Huang et al., 2025) or efficiency (Tang et al., 2025; Xiao et al., 2025). However, only a few works (Zhao et al., 2024) have examined attention head specialization in long-context settings, where attention is known to be not robust (Liu et al., 2024; Xiao et al., 2024), and it is an open question if intervening the localized modules is crucial in practical settings (Hase et al., 2023; Wang and Veitch, 2024). Our work contributes to this line of research by finding better specialized set of attention heads that explain the model behavior for query-focused longcontext retrieval, and that can be practically useful for zero-shot efficient retrieval.

#### 8 Conclusion

We introduced Query-Focused Retrieval Heads (QRHEADS), a set of attention heads specialized in identifying query-relevant information in longcontext inputs. Detected using query-context attention scores on realistic data, QRHEADS are better aligned with practical retrieval tasks than original retrieval heads. Built on top of QRHEADS, our retrieval method QRRETRIEVER achieves strong performance on both long-context reasoning and re-ranking tasks, outperforming dense retrievers and other LLM-based re-rankers in many settings. These findings highlight the practical utility of QR-HEADS and offer insights for further improving retrieval with LMs.

669

670

671

672

673

674

675

621

622

623

#### Limitations

571

573

574

576

577

578

581

582

586

588

594

598

601

610

611

612

613

614

615

616

617

618

619

620

Our work detects improved retrieval heads and builds general-purpose retrievers based on them. We do not explore techniques that involve updating model parameters, as our goal is to develop flexible methods that can directly use off-the-shelf models as retrievers. Consequently, we leave to future work the investigation of parameter-updating techniques that leverage insights from QRHEADS.

While our method finds that QRHEADS can enhance downstream performance, and shows the importance of two factors leading to selection of better heads. We lack a complete understanding of the internal mechanism accounting for QRHEADS's effectiveness. Future work could apply circuit analysis techniques (e.g., Bhaskar et al. (2024); Shi et al. (2024)) to dissect the fine-grained behaviors and roles of these heads.

Our evaluation primarily targets passage reranking and long-context multi-hop reasoning tasks. Although our approach is conceptually applicable to broader long-context tasks—such as longdocument summarization (Shaham et al., 2023; Laban et al., 2024)—it remains unclear whether it generalizes to such tasks without thorough empirical validation.

Finally, our experiments are limited to English datasets. As LMs may exhibit different behaviors across languages, the cross-lingual robustness of our approach remains an open question.

#### References

- Adithya Bhaskar, Alexander Wettig, Dan Friedman, and Danqi Chen. 2024. Finding transformer circuits with edge pruning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems.*
- Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2021. Autoregressive entity retrieval. *Preprint*, arXiv:2010.00904.
- Shijie Chen, Bernal Jiménez Gutiérrez, and Yu Su. 2025. Attention in large language models yields efficient zero-shot re-rankers. *Preprint*, arXiv:2410.02642.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8493– 8502, Dublin, Ireland. Association for Computational Linguistics.
- Andrew Drozdov, Honglei Zhuang, Zhuyun Dai, Zhen Qin, Razieh Rahimi, Xuanhui Wang, Dana Alon,

Mohit Iyyer, Andrew McCallum, Donald Metzler, and Kai Hui. 2023. Parade: Passage ranking using demonstrations with large language models. *Preprint*, arXiv:2310.14408.

- Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. 2024. Finding alignments between interpretable causal variables and distributed neural representations. In *Proceedings of the Third Conference on Causal Learning and Reasoning*, volume 236 of *Proceedings of Machine Learning Research*, pages 160–187. PMLR.
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2023. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. In *Thirtyseventh Conference on Neural Information Processing Systems*.
- Lei Huang, Xiaocheng Feng, Weitao Ma, Yuchun Fan, Xiachong Feng, Yangfan Ye, Weihong Zhong, Yuxuan Gu, Baoxin Wang, Dayong Wu, Guoping Hu, and Bing Qin. 2025. Improving contextual faithfulness of large language models via retrieval heads-induced optimization. *arXiv preprint arXiv:2501.13573*.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. *Preprint*, arXiv:2112.09118.
- Garrett Kamradt. 2024. Needle in a haystack pressure testing llms.
- Gregory Kamradt. 2023. Needle In A Haystack pressure testing LLMs.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Philippe Laban, Alexander Fabbri, Caiming Xiong, and Chien-Sheng Wu. 2024. Summary of a haystack: A challenge to long-context LLMs and RAG systems. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 9885–9903. Association for Computational Linguistics.
- Jinhyuk Lee, Alexander Wettig, and Danqi Chen. 2021. Phrase retrieval learns passage retrieval, too. *Preprint*, arXiv:2109.08133.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020.

783

784

Retrieval-augmented generation for knowledgeintensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

676

686

701

702

703

704

705

706

710

712

715

717

718

719

720

721

722

723

724

727

- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Inferencetime intervention: Eliciting truthful answers from a language model. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
  - Llama-3 Team. 2024. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2024. Fine-tuning llama for multi-stage text retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, page 2421–2425, New York, NY, USA. Association for Computing Machinery.
- Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. 2023. Zero-shot listwise document reranking with a large language model. *Preprint*, arXiv:2305.02156.
- Callum McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. 2023. Copy Suppression: Comprehensively Understanding an Attention Head. *arXiv preprint arXiv:2310.04625*.
- Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In Advances in Neural Information Processing Systems.
- Niklas Muennighoff. 2022. Sgpt: Gpt sentence embeddings for semantic search. *Preprint*, arXiv:2202.08904.
- Rodrigo Nogueira and Kyunghyun Cho. 2020. Passage re-ranking with bert. *Preprint*, arXiv:1901.04085.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, and 7 others. 2022. In-context Learning and Induction Heads. *arXiv preprint arXiv:2209.11895*.
- Chau Minh Pham, Yapei Chang, and Mohit Iyyer. 2025. Clipper: Compression enables long-context synthetic data generation. *Preprint*, arXiv:2502.14854.

- Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. 2023. Rankzephyr: Effective and robust zeroshot listwise reranking is a breeze! *Preprint*, arXiv:2312.02724.
- Revanth Gangi Reddy, JaeHyeok Doo, Yifei Xu, Md Arafat Sultan, Deevya Swain, Avirup Sil, and Heng Ji. 2024. First: Faster improved listwise reranking with single token decoding. *Preprint*, arXiv:2406.15657.
- Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2023. Improving passage retrieval with zero-shot question generation. *Preprint*, arXiv:2204.07496.
- Uri Shaham, Maor Ivgi, Avia Efrat, Jonathan Berant, and Omer Levy. 2023. ZeroSCROLLS: A zero-shot benchmark for long text understanding. In *Findings* of the Conference on Empirical Methods in Natural Language Processing (EMNLP Findings).
- Claudia Shi, Nicolas Beltran-Velez, Achille Nazaret, Carolina Zheng, Adrià Garriga-Alonso, Andrew Jesson, Maggie Makar, and David Blei. 2024. Hypothesis testing the circuit hypothesis in LLMs. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Alessandro Stolfo, Ben Peng Wu, Wes Gurnee, Yonatan Belinkov, Xingyi Song, Mrinmaya Sachan, and Neel Nanda. 2024. Confidence regulation neurons in language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems.*
- Weiwei Sun, Lingyong Yan, Zheng Chen, Shuaiqiang Wang, Haichao Zhu, Pengjie Ren, Zhumin Chen, Dawei Yin, Maarten de Rijke, and Zhaochun Ren. 2023. Learning to tokenize for generative retrieval. *Preprint*, arXiv:2304.04171.
- Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2024. Is chatgpt good at search? investigating large language models as re-ranking agents. *Preprint*, arXiv:2304.09542.
- Hanlin Tang, Yang Lin, Jing Lin, Qingsen Han, Danning Ke, Shikuan Hong, Yiwu Yao, and Gongyi Wang. 2025. Razorattention: Efficient KV cache compression through retrieval heads. In *The Thirteenth International Conference on Learning Representations*.
- Yi Tay, Vinh Q. Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao, Jai Gupta, Tal Schuster, William W. Cohen, and Donald Metzler. 2022. Transformer memory as a differentiable search index. *Preprint*, arXiv:2202.06991.
- Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2).*

880

881

882

883

884

885

886

887

888

841

842

Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. 2024. Function vectors in large language models. In *The Twelfth International Conference on Learning Representations.* 

788

789

790

791

792

793

794

795

797

799

801

803

808

809

810

811

812

813

814

815

816

817

818

819

820

821

824

825

826

832

833

834

836

837

838

839

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Zihao Wang and Victor Veitch. 2024. Does editing provide evidence for localization? In *ICML 2024 Workshop on Mechanistic Interpretability*.
- Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. 2025a. Longmemeval: Benchmarking chat assistants on long-term interactive memory. *Preprint*, arXiv:2410.10813.
- Wenhao Wu, Yizhong Wang, Guangxuan Xiao, Hao Peng, and Yao Fu. 2025b. Retrieval head mechanistically explains long-context factuality. In *The Thirteenth International Conference on Learning Representations*.
- Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, junxian guo, Shang Yang, Haotian Tang, Yao Fu, and Song Han. 2025. Duoattention: Efficient long-context LLM inference with retrieval and streaming heads. In *The Thirteenth International Conference on Learning Representations*.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 39 others. 2024. Qwen2 technical report. ArXiv, abs/2407.10671.
- Xi Ye, Fangcong Yin, Yinghui He, Joie Zhang, Howard Yen, Tianyu Gao, Greg Durrett, and Danqi Chen. 2025. Longproc: Benchmarking long-context language models on long procedural generation. In *arXiv*.
- Howard Yen, Tianyu Gao, Minmin Hou, Ke Ding, Daniel Fleischer, Peter Izsak, Moshe Wasserblat, and Danqi Chen. 2025. Helmet: How to evaluate longcontext language models effectively and thoroughly.
- Fangcong Yin, Xi Ye, and Greg Durrett. 2024. Lofit: Localized fine-tuning on LLM representations. In The Thirty-eighth Annual Conference on Neural Information Processing Systems.
- Dun Zhang, Jiacheng Li, Ziyang Zeng, and Fulong Wang. 2025. Jasper and stella: distillation of sota embedding models. *Preprint*, arXiv:2412.19048.

Xinyu Zhao, Fangcong Yin, and Greg Durrett. 2024. Understanding synthetic context extension via retrieval heads.

#### **A** Details about Evaluation Datasets

We use LongMemEval (Wu et al., 2025a) and CLIP-PER (Pham et al., 2025) for evaluating our systems on long-context reasoning.

**LongMemEval** evaluates the long-term memory capabilities of LLM-driven chat assistants across five fundamental abilities: information extraction, multi-session reasoning, temporal reasoning, knowledge updates, and abstention. We segment the LongMemEval-S dataset ( $\sim$ 115k tokens/question) at the round level, where each round is a document consisting of a single user message paired with the corresponding assistant response.

**CLIPPER** targets narrative claim verification—a challenging long-context reasoning task that requires verifying claims over entire books, with an average length of 90K tokens and 23 chapters. In CLIPPER, data is split at the chapter level, with each chapter treated as an individual document during retrieval.

**Evaluation Process** For each question, we first feed the entire context (e.g., all chapters or dialogue rounds) into the language model without using any first-stage retriever. We compute a retrieval score for each document or segment using our method described in §4. We then select the top-k documents based the scores, concatenate them, and feed them together with the query into the language model in a second pass to generate the final answer. We choose k = 5, 10 for LongMemEval and k = 3, 5 for Clipper. We report retrieval performance using recall and downstream task performance using accuracy.

#### **B** NIAH Test on Qwen-2.5-7B-Instruct

We evaluate Qwen-2.5-7B-Instruct on the NIAH test by masking selected attention heads. As shown in Figure 3 and Figure 4, pruning the top 16 QR-HEAD leads to a more substantial degradation in NIAH performance compared to pruning the top 16 RETHEAD, indicating the greater functional importance of QRHEAD. When pruning the top 32 heads, the performance gap between QRHEAD and RET-HEAD narrows, suggesting that QRHEAD achieves better efficiency and effectiveness with fewer heads for retrieval in NIAH task.



Random Heads 0.1 0.5 0.9 Original Retrieval Heads (Wu et al., 2025) 0.1 Depth 0.5 0.9 **QRHeads** (Ours) 0 1 0.5 0.9 100K 1K 50k Context Length

Figure 3: Top: Masking 16 random heads of Qwen2.5-7B-Instruct. Middle: Masking the top 16 original retrieval heads (Wu et al., 2025b). Bottom: Masking the top 16 QRHeads.

The licenses datasets used in our work include:

• LongMemEval (Wu et al., 2025a) under MIT

• Clipper (Pham et al., 2025) under Apache li-

• NQ (Kwiatkowski et al., 2019) under Creative

• BEIR (Thakur et al., 2021) under Creative

**Computational Resources and Model** 

We use Llama-3.2 (3B), Llama-3.1 (8B and

70B) (Llama-3 Team, 2024), and Qwen2.5

(7B) (Yang et al., 2024). 8B models were run using

Commons Attribution Share Alike 4.0 Read

Commons Attribution Share Alike 3.0.

**License of Datasets** 

on choosealicense.com

License.

cense 2.0.

890

С

D

Sizes

894

895

# 901

902

903

904

906

Figure 4: Top: Masking 32 random heads of Qwen2.5-7B-Instruct. Middle: Masking the top 32 original retrieval heads (Wu et al., 2025b). Bottom: Masking the top 32 QRHeads.

infrastructure.

#### Е **Potential Risks of Our Work**

N/A. Our work investigates the capabilities of existing language models, without proposing new model architectures or training procedures. While large language models pose well-known risks-including potential misuse, generation of harmful content, and encoding of societal biases-our study does not introduce new risks beyond those already covered in the broader literature. As such, we do not believe any specific risk mitigation measures are necessary for the scope of this work.

All experiments were conducted on A100-based

a single NVIDIA A100 GPU with 80GB of memory, and 70B models were run using 4 A100 GPUs.

907

908