Never Lost in the Middle: Mastering Long-Context Question Answering with Position-Agnostic Decompositional Training

Anonymous ACL submission

Abstract

While large language models (LLMs) are equipped with longer text input capabilities than before, they are struggling to seek correct information in long contexts. The "lost in the middle" problem challenges most LLMs, referring to the dramatic decline in accuracy when correct information is located in the middle. To overcome this crucial issue, this paper proposes to enhance the information searching and reflection ability of LLMs in long contexts via specially designed tasks called Position-Agnostic Multi-step QA (PAM QA). Trained in this task, our model excels in focusing more precisely on the desired information. Experimental results show substantial improvement in Multi-doc QA and other benchmarks, superior to state-of-the-art models by 13.7% absolute gain in shuffled settings, by 21.5% in passage retrieval task. We release our model and code to promote related research in the community.¹

1 Introduction

002

005

007

009

011

012

016

017

019

022

027

038

Large Language Models (LLMs), renowned for their exceptional generative and zero-shot learning abilities across diverse natural language processing (NLP) fields, have found extensive downstream applications (OpenAI, 2023; Boiko et al., 2023; Cheng et al., 2023; Waisberg et al., 2023; Hu et al., 2023). However, LLMs suffer from severe hallucinations, significantly compromising their performance in knowledge-oriented QA, dialogue, and writing (Roberts et al., 2020; Agrawal et al., 2023). Retrieval Augmented Generation (RAG) is an effective solution to hallucinations and remarkable improvements have been achieved by incorporating supporting knowledge into the input of LLMs (Lewis et al., 2020b; Shuster et al., 2021; Thoppilan et al., 2022; Shi et al., 2023a). The most fundamental challenge to address in RAG is long context and Multi-document question answering (Multi-doc QA).

Thorough research has been conducted to deal with long context inputs, categorized into three mainstreams: The first one is to expand the context window using a sliding window (Dai et al., 2019). Other researchers proposed to enhance the extrapolation ability by improving the Relative Positional Encoding in Transformers, the backbone of most LLMs (Su et al., 2021; Press et al., 2021; Luo et al., 2022; Vaswani et al., 2017). These two kinds of modifications both show substantial improvement in language modelling (LM). The third category of studies focuses on the recurrent compression of memory for long-range sequence learning (Rae et al., 2019; Peng et al., 2023). This methodology effectively learns the comprehensive representation of context, demonstrating notable proficiency in rapid computation and cost-effectiveness during inference. Though the methods above show strong performance in specific tasks and support LLMs with extra-long context windows, i.e. GPT3.5-Turbo-16k, Claude-v1.3-100k, Longchat (Dacheng et al., 2023), LLMs fail to produce correct answers if related documents are located in the middle of the context, called "lost in the middle" (Liu et al.). It is fatal for Multi-doc QA. However, whether a similar deterioration exists in Chinese LLMs has been unexplored and solutions to this problem have rarely been researched.

041

042

043

044

045

047

050

051

055

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

We hypothesise that the scale of attention scores of the beginning context grows large after pretraining and instruction tuning while that of the middle context, whose position is less trained, remains small for a long distance to the current token. The small attention scores limit the contribution of related information to the answer and result in lower QA accuracy.

To overcome the pitfall, we proposed positionagnostic decompositional training to even up the attention scores over input context. Concretely, we designed a tailored Multi-doc QA task in which positive documents are located at various posi-

¹It is publicly available at https://xxx



Figure 1: The workflow of PAM QA. The blue dashed lines indicate information flows. The desired output of a sample is composed of three parts, corresponding to three steps: Question repetition, index prediction, and answer summarization.

tions in contexts among noisy documents. The task presents a significant challenge, compelling the models to extract and summarize information despite the interference of useless ones (Ye et al., 2022). As human beings routinely solve complex tasks by decomposition to obtain higher quality outcomes (Cheng et al., 2015; Correa et al., 2023), we modified the Multi-doc QA task as a multi-step reasoning task, the Position-Agnostic Multi-step QA (PAM QA), combining the Chain-of-Thought (COT, Wei et al.) and position-agnostic Multi-doc QA. Trained with explicit extraction of the question and the index of supporting documents before generating answers, models learn to distinguish correct information from noisy ones and attend to them. It also makes the attention to the question and supporting indexes stronger because the attention scale decays with increasing distance (Su et al., 2021). Empirical results on Multi-doc QA and other benchmarks show that, with only 1/2 or 1/4 context window size, our model improves upon state-of-the-art (SOTA) models by 7.0% in the topranked setting and by 13.7% in the shuffled setting. Competitive results are shown in other attentiondependent tasks including passage retrieval and summarization.

The contribution of this paper is threefold:

- This paper proposed a novel task named PAM QA to tackle the "lost in the middle" issue, which is fatal for knowledge-intensive scenarios. To our knowledge, it is the first attempt to solve the problem by training on special tasks.
- We investigate the model's behaviour in-depth, revealing that failing to focus on target information may be the cause of "lost in the middle".

• Comprehensive experiments have shown that the proposed PAM QA is effective in solving the "lost in the middle" problem. Our model surpasses SOTA in Multi-doc QA and other related tasks on renowned Chinese benchmarks. It is non-trivial that the general QA ability of the model is also strong and satisfying. The model is open-sourced to boost future research in the community. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

2 Position-Agnostic Multi-step QA

Multi-doc QA refers to a type of QA task where a model is presented with multiple documents and asked to answer questions based on these documents. It is difficult both for models and human beings, where we need to aggregate and comprehend information from noisy candidates to generate accurate answers while struggling with fading memory. In this situation, task decomposition, identifying subproblems and reasoning about them, becomes essential (Correa et al., 2023).

Therefore we decomposed the difficult Multidoc QA to PAM QA. This innovative task comprises three steps, as depicted in Figure 1.

2.1 Question repetition

The first step is question repetition (QR). The questions are placed at the front as a contextual-aware representation (Liu et al.). The subtask is started with $prefix_1$, "As for the question:" (or expressions with identical meaning) to prompt the model.

2.2 Index Prediction

The second step is index prediction (IP), namely to predict the indexes of the supported documents for the question as an MRC task, beginning with $prefix_2$: "Based on the information numbered".

117

Supporting evidence not only helps LLMs ver-152 ify themselves but also aids users in evaluating 153 responses (Menick et al., 2022). Remarkable re-154 sults have been shown in generating quotes and 155 citations (Thoppilan et al., 2022; Menick et al., 156 2022). Different from previous works, it is sug-157 gested to predict the indexes of corresponding evi-158 dence rather than a verbatim quote extracted from 159 a longer source retrieved from documents. In this 160 step, we hypothesize that the indicator helps to en-161 code and navigate the attention to corresponding documents. Considering the indexes in the sec-163 ond step only count for a small number of tokens 164 and are hard to emphasize in the sequential cross-165 entropy loss during training, a relevance MRC task 166 is added as a supplement. The task needs to predict indexes of correct documents only.

2.3 Answer Summarization

169

170

171

172

173

174

175

176

177

178

180

181

182

184

185

188

189

190

191

192

193

The third step is to generate the final answer for the Multi-doc QA task. Thanks to the previous two steps, this task can be simplified as a summarization task, called answer summarization (AS). The step begins with an answer indicator like "my answer is" as $prefix_3$.

The entire process of PAM QA unfolds as follows: when receiving a question, a set of candidate documents, and a specific instruction, the model initiates by generating $prefix_1$. It then proceeds to restate the question, predicting the indexes of related evidence after incorporating a connecting phrase, denoted as $prefix_2$. Finally, it formulates an answer to the question by aggregating previous information, following an answer indicator, $prefix_3$.

In line with the proverb *"the palest ink is better than the best memory,"* we teach the model to take notes, turning these annotations into a highway to the relevant knowledge. It can reduce the distraction of extraneous information and make the attention to the question and supporting index stronger because the attention scale decays with increasing distance.

3 Training Data Construction

We equipped our model with distinguishing ability
through instruction tuning. The training procedure
is composed of two stages. We expand the LLM's
context window to 8K in the first stage. In the
second stage, the model was trained with PAM
QA data to solve the attention (or memory) failure

called "lost in the middle".

3.1 Context Window Expansion

We used about 300k selected data for general supervised finetuning (SFT). The data cover various categories of tasks including QA, MRC, role-playing, writing, coding, translation, brainstorming, math, Language Modeling (LM), and other natural language understanding (NLU) tasks like text classification. Then the data are packed to 8k window size in a multi-turn conversation style except for the LM task, while the LM task calculates the cross-entropy loss on the whole sequence. 201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

224

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

3.2 PAM QA

Data are constructed according to the three steps in PAM QA. We collected Multi-doc QA data and adapted it to PAM QA data.

First, we filtered out 30k samples of the Fact category with a single answer from DuReader2.0 dataset (He et al., 2018) and 20k samples from WebCPM (Qin et al., 2023). DuReader2.0 is the largest Chinese MRC dataset from Web doc and community QA, containing 200K questions, 420K answers and 1M documents. To ensure the quality of data, we creatively utilize a reward model to score the samples and select the high-quality part of them with a certain threshold, inspired by (Li et al., 2023). The reward model is trained with 69k human-ranked samples for alignment in general tasks, following (Köpf et al., 2023; Ouyang et al., 2022). As both datasets exclusively consist of positive samples, negative samples are ingeniously generated and incorporated in a sophisticated manner.

For each sample, documents in the whole collection except the positive ones are regarded as negative samples. As collaborative learning is beneficial to RAG (Izacard et al., 2022), we build an embedding-based search engine with all the documents in the corresponding dataset. Subsequently, we retrieved negative documents with questions from the search engine for a partition comprising 70% of the data, while we randomly sampled from the original negative candidates for the remaining portion of data. The retrieved negative samples are more relevant to questions and harder to distinguish from the positive ones than random samples. Next, documents are shuffled within each sample in 50% of the data to prevent positive ones from consistently being positioned at the beginning of contexts. Next, 25k samples were sampled from a

retrieval benchmark, T2Rank (Xie et al., 2023) as
the relevance MRC data for a supplement for step 2.
The negative samples are randomly sampled from
the original hard negative collections. Indexes of
positive candidates are recorded and shuffled with
the negative ones. The max lengths are sampled
from 1k to 8k under the uniform distribution in all
samples. This ensures our model can deal with samples of various input lengths. In this way, correct
documents can be located at any position.

Additionally, we generate 5% of samples, where all documents are negative as the positive ones are excluded. The answer for these samples is a constant term indicating "I don't know." Termed "Synthetic Unknown," these samples are employed to train the model to recognize situations where the correct document is absent.

Finally, We replayed the same distribution of data in general SFT with a ratio of 20% in this stage to alleviate the catastrophic forgetting (McCloskey and Cohen, 1989; Rebuffi et al., 2017). The total training samples in this stage summed up to 90k.

3.3 Training

262

264

265

269

271

272

274

275

276

279

284

286

289

290

296

We trained our model based on a pre-trained LLM that adapted from LLaMA2 (Touvron et al., 2023; Gan et al., 2023; Zhang et al., 2023). We trained for 2 epochs on 16 A100 GPUs in both stages with constructed data. The learning rate began with 1e-5 then decayed to 1e-6 with a warmup for the first 0.05% steps in the first stage. The max learning rate for the second stage was 5e-6. Flash Attention (Dao et al., 2022) was utilized to accelerate the training procedure. Sampling is turned on for all models during testing in the benchmarks. The hyper-params for testing are listed in Appendix A.

4 Experiments

In this section, we evaluate the long-context QA abilities of our model and existing representative LLMs. By inspecting the performance, we can verify whether our model overcomes the so-called "lost in the middle" problem (Liu et al.).

4.1 Benchmarks

We conducted experiments on a long context benchmark, LongBench (Bai et al., 2023) and Retrieval-Augmented Generation Benchmark (RGB, Chen et al.). The benchmark measures various abilities of the models given long input contexts. Specifically, we tested models on three related tasks in Long-

Datasets	Avg length	Source	Metrics
Multi-doc.	15,768	DuReader	Rouge-L
Synt.	6,745	C4 Chinese	Accuracy
Summ.	15,380	VCSUM	Rouge-L
RGB NR.	1,105.7	Self Generated	EM

Table 1: The statistics of input lengths of the testing datasets. Multi-doc. is short for Multi-doc QA. Synt. and Summ. represent Synthetic Tasks and Summarization respectively while RGB NR is the abbreviation of RGB noise robustness task.

Bench: Chinese Multi-doc QA, Synthetic tasks, and summarization. We also used the noise robustness testbed in RGB to test the QA ability in short texts, which examines the information extraction ability given a certain ratio of noise documents. 299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

330

331

332

333

334

335

336

Synthetic task is a document retrieval task, where given a summary, the goal is to find the corresponding document from a large number of documents. This task evaluates the information retrieval ability of LLMs in long contexts. The summarization task gives extremely long meeting records from multiple speakers and asks for a summary. It assesses the model's memory and summarization capabilities. The context lengths and other statistics of the datasets are listed in Table 1. The evaluation scripts provided by the LongBench official website² and RGB official repository were used to compute the results.

We also re-constructed Synthetic Task to examine whether the models are "lost in the middle". The correct passages are relocated at the 1st, 5th, 10th, 15th and 20th with passages located beyond the 20th removed. The results are in Figure 2.

Considering that the documents in the samples of Multi-doc QA tasks are basically sorted by relevance, we shuffled the first 10 candidate documents in each sample to make the real performance exposed, called Multi-doc QA shuffled.

In addition, we conducted a comprehensive human evaluation of model capabilities to see if training on PAM QA harms the general abilities of LLM. The test set contains 200 questions from a wide range of categories.

4.2 Baselines

We compared the performance of the most popular LLMs with a long context window. These strong baselines include: GPT3.5T-turbo-16k extends the context window to 16K tokens, while

²https://github.com/THUDM/LongBench

Model	Multi-doc QA	Synthetic Tasks	Summarization	
(Baichuan2-Turbo-192k)	36.8	90.0	18.4	
Longchat-v1.5-7B-32k	19.5	7.6	9.9	
ChatGLM2-6B-32k	37.6	64.5	16.1	
(ChatGLM3-6B-32k)	44.8	94.0	17.8	
GPT3.5-Turbo-16k	28.7	77.5	16.0	
Vicuna-v1.5-7B-16k	19.3	5.0	15.1	
Xgen-7B-8k	11.0	3.5	2.2	
InternLM-7B-8k	16.3	0.9	12.4	
Our model	44.6	98.5	15.6	

Table 2: The results are Rouge-L percentage for Multi-doc QA and Summarization while Synthetic Tasks compute the accuracy (EM scores). Models are separated in lines by certain context window sizes. ChatGLM3-6B-32k and Baichuan2-turbo-192k are new models after our work.

both Longchat-v1.5-7B-32k (Dacheng et al., 2023) and ChatGLM2(3)-6B-32k (Du et al., 2022) further push the boundary to 32K tokens. Vicunav1.5-7B-16k (Zheng et al., 2023) and Xgen-7B-8k (Nijkamp et al., 2023) offer fine-tuned models on user-shared conversations and 8K sequences respectively. Baichuan2-13B-Chat (Yang et al., 2023) stands out in few-shot learning with a 4K token window, alongside a larger closed-source variant. Lastly, Qwen-14B-chat introduces a 14B parameter model with dynamic ntk (dyn, 2023), trained on an 8K token window size.

5 Results and Discussion

337

339

340

341

342

346

348

359

361

367

In this section, we analyze the experimental results of the LLMs and discuss the reason for the findings.

5.1 Longer window size does not guarantee better performance

As shown in Table 2, Our model has a Rouge-L of 44.6% in the Multi-doc QA task, 7.0% higher than ChatGLM2-6B-32k, which was the SOTA model. With only 1/4 window size, Our model can outperform ChatGLM2-6B-32k at this task. It reveals the strong attention ability of our model since it is an open-book QA task. This Chinese Multi-doc QA dataset does not need to consider all of the contexts, as the correct documents are located at the beginning of contexts.

In the Synthetic task, namely an abstract retrieval task, Our model achieves the highest result with an accuracy of 98.5%. This indicates that the "lost in the middle" issue is almost solved by the proposed method in this paper, as long as the average length is covered.

As for summarization, ChatGLM2-6B-32k and

GPT3.5-Turbo-16k have similar performance with different context window sizes, showing that longer context window sizes do not guarantee better performance. The Rouge-L of our model is only 0.5% lower than SOTA, without any summarization data in the PAM QA training. As the average length of the task is much longer than 8k, Our model with a longer context length will have a promising improvement.



Figure 2: The EM score on Synthetic (passage retrieval) Task from LongBench with correct document inserted to certain position ranging from 1st to 20th.

5.2 PAM QA alleviates lost in the middle (and tail) problem

Experiments on the Synthetic (passage retrieval) benchmark with modifications to display the performance of models on different positive document positions. Concretely, the correct passage for each sample is inserted into the 1st, 5th, 10th, 15th, and 20th locations respectively among the other documents in each experiment. Ideally, we should see a U-curve described in (Liu et al.), called "lost in the 381

382

383

384

385

386



Figure 3: The attention scores over the input tokens in the self-attention procedure within ChatGLM2-6B-32k and our model on a document repeated 20 times. Length differs with tokenizers.

middle". Results are displayed in Figure 2.

Figure 2 suggests that most open-source LLMs are lost not only in the middle but also in the tail. A significant decrease is observed when positive documents are placed at the 10th position. Despite the employment of techniques such as Al-ibi() or NTK to expand the context window (i.e. Baichuan2-13B-Chat and Qwen-14B-Chat), models still demonstrate low results in the absence of training. However, Our model can survive in different settings of positions, holding a record of 99%. It reveals the effects of PAM QA training.

5.3 Models defeated by shuffled dataset, attention failure being the culprit



Figure 4: Performance on Multi-doc QA before and after shuffling. ChatGLM2 is short for ChatGLM2-6B-32k, GPT3.5-Turbo is short for GPT3.5-Turbo-16k. Scores are in percentage.

Figure 4 demonstrates the models' performance on Multi-doc QA before and after shuffling. We can see a sharp decline in all three models except ours. The largest gap reaches 17.3%, from ChatGLM2-6B-32k. Meanwhile, Baichuan2-13B-Chat also has a 7% reduction although the shuffled documents are within its context length. Therefore, LLMs without extra long context windows also have difficulty dealing with the challenge. Our model is the most robust model with a 3.7% decrease.

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

To unearth the cause of the decline and examine the attention capabilities of models, we visualize the attention scores of the last layer for the identical input. We repeat a sentence including the correct answer 20 times as the context to find if all the answers will be highlighted in the selfattention procedure in models. Attention scores of ChatGLM2-6B-32k and Our model over the input are depicted in Figure 3.

We can see the attention scores on documents are fading away in ChatGLM2-6B-32k, as the context after the first 100 tokens is almost neglected. The situation in Our model is quite different. 20 peaks of attention scores are observed (the last one is next to the beginning of instruction), corresponding to the answers in sentences. It reveals that attention to related tokens is the key to the performance gap between models. The models struggle to precisely focus on the correct tokens, paying tremendous attention to the beginning and the ending tokens (where instruction and query are frequently located), which is the culprit of the "lost in the middle" problem.

438 439 440

441

442

443

444

445

446

447

448

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

5.4 Each step in PAM QA matters

Here we inspect the contribution of each step in PAM QA. The variants are evaluated on Multi-doc QA and Synthetic tasks. Results of this ablation study are listed in Table 3.

Variants	Multi-doc QA	Synt.	
Our model	44.6	98.5	
- QR	38.8	98.0	
- QR - IP	37.8	1.3	
Only-SFT	8.7	7.5	

Table 3: Synt. is short for Synthetic tasks. Results are in percentage. QR is short for question repetition. IP is short for index prediction. Only-SFT represents the finetuned model only with context window expansion.

When we remove the question repetition, the first step in PAM QA, a 5.8% decrease can be observed in Multi-doc QA, showing the inevitable contribution to high performance. It strengthens the attention of the question and constrains the following generation procedure explicitly. By repeating the question first, models can directly attend to the question in the subsequent steps without going through a long context, avoiding the disruption of context when performing self-attention.

When the index prediction(IP) step is also removed, a pronounced decrease in Synthetic tasks emerges, which emphasizes the importance of IP. It not only teaches LLMs to distinguish between related and irrelevant information but also frees the model from a prior experience (i.e., seeking information from the beginning and the end of context). Meanwhile, this step relieves models from scanning the long input again and makes them focus on the abstracted related information in the subsequent process. Since the scale of the attention scores decays as the distance grows (Su et al., 2021), models with rotary position embeddings (RoPE) struggle to remember the remote tokens without training. A slight drop of Rouge-L in addition also shows the benefit of the second step in QA. With the former two steps, the question and the potentially correct evidence are listed just a few tokens ahead. This reduces the probability of forgetting questions and context by decreasing the distance.

An enormous gap between the results of Our model and the model without QR and IP, which represents the performance of only trained with filtered Multi-doc QA data, indicates the substantial improvement from PAM QA training. We also

Noise Ratio	0	0.2	0.4	0.6	0.8
GPT3.5-Turbo	95.67	94.67	91.00	87.67	70.67
ChatGLM2-6B	86.67	82.33	76.67	72.33	54.00
(ChatGLM3-6B)	91.67	90.00	89.00	84.67	66.33
Baichuan2-13B-Chat	93.00	90.33	89.00	82.33	63.33
(Qwen-14B-Chat)	94.67	92.00	88.00	85.30	69.67
Our model	96.00	90.67	90.00	85.50	67.33

Table 4: Performance in RGB noise robustness testbed. EM scores are in percentage. ChatGLM3-6B and Qwen-14B-Chat are new models after our work.

visualize the attention scores when predicting the first token and discover the generated questions and indexes are highlighted, shown in Figure 7 in Appendix C. 477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

Compared with Only-SFT, the model only with the first stage training, the variant model without QR and IP steps also gains 29.1% improvement, which shows position-agnostic Multi-doc QA data benefit the related abilities. By simply transforming the Multi-doc QA into PAM QA, the same data can boost the performance by 16.8% in Multi-doc QA, and 97.2% in Synthetic Task.

5.5 Competitive results observed in short text Multi-doc QA

As reported in Table 4, Our model still has a competitive performance among open-source models on short-text multi-doc QA although not trained on any short texts. Even compared with the latest popular Chinese LLMs, Qwen-14B-chat and ChatGLM3-6B-32k, results from Our model are higher under the setting of noise rate in [0,0.4,0.6].

5.6 General ability is preserved with PAM QA Training

A side-by-side (SBS) comparison was performed by 3 human annotators to check the general ability of Our model. They are all master students. Capabilities including commonsense, math, reasoning, QA, writing, harmlessness, etc. are examined in the test, as shown in Figure 5. The annotators have to choose which one is better based on the answers unless the answers are both bad or the same, as in (Zheng et al., 2023). They are blind to the models and other information. Results compared with similar size models, Ziya-LLaMa-13B-v1.1³ and Baichuan2-13B-Chat respectively are illustrated in Figure 6.

³https://huggingface.co/IDEA-CCNL/Ziya-LLaMA-13Bv1.1



Figure 5: The distribution of tasks in the general ability test.



Figure 6: SBS results on general ability evaluation that contains a wide range of tasks. The results of our model are superior to Ziya-LLaMa-13B-v1.1 and close to Baichuan2-13B-Chat.

Figure 6 summarizes the human preference between Our model and other open-source LLMs. Although Our model was trained only with PAM QA data, its performance is only slightly inferior to Baichuan2-13B-chat and better than Ziya-LLaMa-13B-v1.1 significantly. Thus, the general capabilities are maintained after the PAM QA training.

6 Related Works

513

514

516

517

518

519

520

521

6.1 Retrieval-Augmented Language Models

522Retrieval-Augmented Language Models (RALMs)523mark notable progress in natural language process-524ing by merging the capabilities of expansive lan-525guage models with the precision and intricacy of-526fered by external knowledge sources. (Guu et al.,5272020; Lewis et al., 2020a; Izacard et al., 2022).528These models use a retriever to search through a529large body of evidence, like Wikipedia, to find a530specific set of documents related to the user's query.

Afterwards, a reader component is utilized to carefully examine these documents and generate a response. This two-step process guarantees both relevance and depth in the produced answers. Recent research efforts have concentrated on enhancing the performance of the retriever (Karpukhin et al., 2020; Sachan et al., 2023) or the reader(Izacard and Grave, 2020; Cheng et al., 2021), training the system end-to-end (Lewis et al., 2020a; Sachan et al., 2021), and integrating the retrieval systems with black-box large language models (Shi et al., 2023b; Yu et al., 2023; Trivedi et al., 2023) 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

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

6.2 RALMs Adapted to Long and Noisy Context

Recent research emphasizes the influence of contextual length and the position of related context on the performance of language models (Krishna et al., 2023; Bai et al., 2023; Liu et al.). The research closely aligned with ours is the study by (Yoran et al., 2023), focusing on training RALMs to disregard irrelevant contexts. However, it overlooked scenarios involving long contexts, specifically the "lost in the middle" issue, a key consideration in our work.

7 Conclusion

In this paper, we assume that the widely recognized "lost in the middle" phenomenon may caused by weak attention to target information. We found popular Chinese LLMs are "lost" not only in the middle but also tail. A novel approach is proposed to address the deficiency in LLMs by training models with Posistion-Agnostic Multi-step (PAM) QA. Experimental results show the superiority and effectiveness of our method, surpassing SOTA LLMs in Multi-doc QA and passage retrieval significantly, with only 1/4 context window size. By shuffling the candidate documents in open benchmarks, degraded performance is observed in all models, among which our model is the most robust one. The ablation study also reveals the significant effect of PAM QA and the positive contribution of its components. Our study also finds that LMs with extremely long context windows do not ensure better performance on Multi-doc QA and passage retrieval tasks. We hope our study provides profound insight into the "lost in the middle" problem at a broader scale and sheds light on developing more intelligent LLMs.

588

590

591

592

593

595

606

607

610

611

612

614

615

616

617

621

622

623

627

580 Our work only covers the "lost in the middle" issue 581 and experiments with Chinese Benchmarks using 582 popular Chinese and some English LLMs with long 583 context capability. More crucial further investiga-584 tions are welcome.

References

Limitations

- 2023. Dynamically scaled rope further increases performance of long context llama with zero fine-tuning.
- Ayush Agrawal, Lester Mackey, and Adam Tauman Kalai. 2023. Do language models know when they're hallucinating references? *arXiv preprint arXiv:2305.18248*.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.
- DaniilA. Boiko, Robert MacKnight, and Gabe Gomes. 2023. Emergent autonomous scientific research capabilities of large language models.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2023. Benchmarking large language models in retrieval-augmented generation.
- Hao Cheng, Yelong Shen, Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2021. UnitedQA:
 A hybrid approach for open domain question answering. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3080–3090, Online. Association for Computational Linguistics.
- Justin Cheng, Jaime Teevan, Shamsi T Iqbal, and Michael S Bernstein. 2015. Break it down: A comparison of macro-and microtasks. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, pages 4061–4064.
- Kunming Cheng, Qiang Guo, Yongbin He, Yanqiu Lu, Shuqin Gu, and Haiyang Wu. 2023. Exploring the potential of gpt-4 in biomedical engineering: the dawn of a new era. *Annals of Biomedical Engineering*, pages 1–9.
- Carlos G Correa, Mark K Ho, Frederick Callaway, Nathaniel D Daw, and Thomas L Griffiths. 2023. Humans decompose tasks by trading off utility and computational cost. *PLOS Computational Biology*, 19(6):e1011087.
- Li* Dacheng, Shao* Rulin, Xie Anze, Sheng Ying, Zheng Lianmin, Gonzalez Josep E., Stoica Ion, Ma Xuezhe, and Hao Zhang. 2023. How long can open-source llms truly promise on context length?

- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. arXiv preprint arXiv:1901.02860.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In Advances in Neural Information Processing Systems.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335.
- Ruyi Gan, Ziwei Wu, Renliang Sun, Junyu Lu, Xiaojun Wu, Dixiang Zhang, Kunhao Pan, Ping Yang, Qi Yang, Jiaxing Zhang, et al. 2023. Ziya2: Datacentric learning is all llms need. *arXiv preprint arXiv:2311.03301*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.
- Wei He, Kai Liu, Jing Liu, Yajuan Lyu, Shiqi Zhao, Xinyan Xiao, Yuan Liu, Yizhong Wang, Hua Wu, Qiaoqiao She, et al. 2018. Dureader: a chinese machine reading comprehension dataset from real-world applications. ACL 2018, page 37.
- Chenxu Hu, Jie Fu, Chenzhuang Du, Simian Luo, Junbo Zhao, and Hang Zhao. 2023. Chatdb: Augmenting llms with databases as their symbolic memory. *arXiv preprint arXiv:2306.03901*.
- Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, et al. 2023. Openassistant conversations–democratizing large language model alignment. *arXiv preprint arXiv:2304.07327*.

639

640

641

642

632

633

634

649

650

651

656

657

658

659

660

661

669

670

671

672

673

674

675

676

677

678

679

680

681

682

684

685

795

741

Kalpesh Krishna, Erin Bransom, Bailey Kuehl, Mohit Iyyer, Pradeep Dasigi, Arman Cohan, and Kyle Lo. 2023. Longeval: Guidelines for human evaluation of faithfulness in long-form summarization.

687

697

704

708

710

711

712

713

714

715

717

718

719

720

721

722

723

724

725

729

730

731 732

733

734

735

736

737

739

740

- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020a. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- PatrickS.H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Michael Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020b. Retrieval-augmented generation for knowledge-intensive nlp tasks. *arXiv: Computation and Language,arXiv: Computation and Language.*
 - Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston, and Mike Lewis. 2023. Self-alignment with instruction backtranslation. *arXiv preprint arXiv:2308.06259*.
- NelsonF Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts.
- Shengjie Luo, Shanda Li, Shuxin Zheng, Tie-Yan Liu, Liwei Wang, and Di He. 2022. Your transformer may not be as powerful as you expect. *Advances in Neural Information Processing Systems*, 35:4301–4315.
- Michael McCloskey and Neal J. Cohen. 1989. *Catastrophic interference in connectionist networks: the sequential learning problem*, page 109–165.
- Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, et al. 2022. Teaching language models to support answers with verified quotes. arXiv preprint arXiv:2203.11147.
- Erik Nijkamp, Tian Xie, Hiroaki Hayashi, Bo Pang, Congying Xia, Chen Xing, Jesse Vig, Semih Yavuz, Philippe Laban, Ben Krause, et al. 2023. Long sequence modeling with xgen: A 7b llm trained on 8k input sequence length. *Salesforce AI Research Blog*.
- OpenAI OpenAI. 2023. Gpt-4 technical report.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Przemyslaw Kazienko, Jan Kocon,

Jiaming Kong, Bartlomiej Koptyra, Hayden Lau, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Xiangru Tang, Bolun Wang, Johan S. Wind, Stansilaw Wozniak, Ruichong Zhang, Zhenyuan Zhang, Qihang Zhao, Peng Zhou, Jian Zhu, and Rui-Jie Zhu. 2023. Rwkv: Reinventing rnns for the transformer era.

- Ofir Press, Noah A Smith, and Mike Lewis. 2021. Train short, test long: Attention with linear biases enables input length extrapolation. *arXiv preprint arXiv:2108.12409*.
- Yujia Qin, Zihan Cai, Dian Jin, Lan Yan, Shihao Liang, Kunlun Zhu, Yankai Lin, Xu Han, Ning Ding, Huadong Wang, Ruobing Xie, Fanchao Qi, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2023. WebCPM: Interactive web search for Chinese long-form question answering. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8968–8988, Toronto, Canada. Association for Computational Linguistics.
- Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, and Timothy P. Lillicrap. 2019. Compressive transformers for long-range sequence modelling.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. 2017. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? *arXiv preprint arXiv:2002.08910.*
- Devendra Singh Sachan, Mike Lewis, Dani Yogatama, Luke Zettlemoyer, Joelle Pineau, and Manzil Zaheer. 2023. Questions are all you need to train a dense passage retriever. *Transactions of the Association for Computational Linguistics*, 11:600–616.
- Devendra Singh Sachan, Siva Reddy, William Hamilton, Chris Dyer, and Dani Yogatama. 2021. End-to-end training of multi-document reader and retriever for open-domain question answering.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and ScottWen-tau Yih. 2023a. Trusting your evidence: Hallucinate less with context-aware decoding.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023b. Replug: Retrieval-augmented black-box language models.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. In *Findings* of the Association for Computational Linguistics: *EMNLP 2021*, pages 3784–3803.

897

898

854

Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2021. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*.

796

797

799

810

811

812

813

814

815

816

817

818

819

820

821

822

823

825

829

830

831

836

837 838

839

841

847 848

849

851

852

- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. arXiv preprint arXiv:2201.08239.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledgeintensive multi-step questions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10014–10037, Toronto, Canada. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Ethan Waisberg, Joshua Ong, Mouayad Masalkhi, Sharif Amit Kamran, Nasif Zaman, Prithul Sarker, Andrew G Lee, and Alireza Tavakkoli. 2023. Gpt-4: a new era of artificial intelligence in medicine. *Irish Journal of Medical Science (1971-)*, pages 1–4.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Xiaohui Xie, Qian Dong, Bingning Wang, Feiyang Lv, Ting Yao, Weinan Gan, Zhijing Wu, Xiangsheng Li, Haitao Li, Yiqun Liu, et al. 2023. T2ranking: A large-scale chinese benchmark for passage ranking. *arXiv preprint arXiv:2304.03679*.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, JunTao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023. Baichuan 2: Open large-scale language models.

- Ziyi Ye, Xiaohui Xie, Yiqun Liu, Zhihong Wang, Xuesong Chen, Min Zhang, and Shaoping Ma. 2022. Towards a better understanding of human reading comprehension with brain signals. In *Proceedings of the ACM Web Conference 2022*, pages 380–391.
- Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2023. Making retrieval-augmented language models robust to irrelevant context. *arXiv preprint arXiv:2310.01558*.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023. Generate rather than retrieve: Large language models are strong context generators.
- Jiaxing Zhang, Ruyi Gan, Junjie Wang, Yuxiang Zhang, Lin Zhang, Ping Yang, Xinyu Gao, Ziwei Wu, Xiaoqun Dong, Junqing He, Jianheng Zhuo, Qi Yang, Yongfeng Huang, Xiayu Li, Yanghan Wu, Junyu Lu, Xinyu Zhu, Weifeng Chen, Ting Han, Kunhao Pan, Rui Wang, Hao Wang, Xiaojun Wu, Zhongshen Zeng, and Chongpei Chen. 2023. Fengshenbang 1.0: Being the foundation of chinese cognitive intelligence.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging Ilm-as-a-judge with mt-bench and chatbot arena.

A Hyper-parameters

We used the following settings in the inference phrase in the LongBench test: $do_sample =$ $True, top_p = 0.85, temperature = 0.8,$ $repetition_penalty = 1.0, early_stopping =$ True.

As for the RGB noise robustness testbed, the temperature was set to 0.2 for all models, with other parameters unchanged.

B Demonstration of PAM QA

A case of Multi-doc QA and the answer produced by Our model and ChatGPT3.5 are shown in Table 5.

C Visualization of attention scores on tokens in answer summarization step

The heatmap of attention scores of each token when generating the final answer is displayed in Figure 7.

Chinese:

<s><human> :给定问题:联合国2022年对人口展望中,预计世界人口峰值达到多少?</human></s>
文章:
[1] 联合国发布的《2022 年世界人口展望》报告中预测,印度将在明年超过中国,成为世界上人口最多的国家。
[2] 根据联合国经济和社会事务部发布的《2022 年世界人口展望》报告,全球人口预计将在 2022 年 11 月 15 日达到 80 亿。报告同时预
测,世界人口预计将在 2080 年代达到约 104 亿的峰值,并在 2100 年之前保持这一水平。
[3] 11月15日,根据统计网站Worldometers实时统计数据,全球人口已超80亿。
[4] 联合国(UN)核监督机构周二批准了日本的"核排污"计划,不顾当地渔民和该地区其他国家的反对,允许日本将受损的福岛第一核电站逾100万
吨核废水排入海洋。
[5] 据联合国官网消息,近日,联合国193个成员国一致通过了一项具有法律约束力的协定,旨在加强海洋生物多样性养护和可持续发展,尤其是
对公海的保护。
请阅读理解上面多个文章,正确地回答问题。若检索结果都不相关,请回答都不相关。
<bot> : 对于问题"联合国2022年对人口展望中,预计世界人口峰值达到多少?",根据上面的编号为1,2,3的信息,我的答案是根据联合国经济</bot>
和社会事务部发布的《2022 年世界人口展望》报告,全球人口预计将在 2022 年 11 月 15 日达到 80 亿。报告同时预测,世界人口预计将在
2080 年代达到约 104 亿的峰值,并在 2100 年之前保持这一水平。因此,联合国预计世界人口峰值将达到约 104 亿。

English:

<s><human>: Given question: What is the expected peak world population in the UN's 2022 population outlook?

Essays:

[1] The World Population Prospects 2022 report released by the United Nations predicts that India will surpass China as the world's most populous country next year.

[2] According to the World Population Prospects 2022 report published by the United Nations Department of Economic and Social Affairs, the global population is expected to reach 8 billion by 15 November 2022. The report also predicts that the world's population is expected to peak at about 10.4 billion in the 2080s and remain at that level until 2100.

[3] On November 15, according to real-time statistics from the statistics website Worldometers, the global population has exceeded 8 billion.

[4] The United Nations (UN) nuclear watchdog on Tuesday approved Japan's plan to "dismantle nuclear wastewater," allowing Japan to discharge more than 1 million tons of nuclear wastewater from the damaged Fukushima Daiichi nuclear power plant into the ocean, despite opposition from local fishermen and other countries in the region.

[5] According to the official website of the United Nations, the 193 member states of the United Nations recently unanimously adopted a legally binding agreement aimed at strengthening the conservation and sustainable development of marine biodiversity, especially the protection of the high seas.

Please read and understand many of the passages above and answer the questions correctly. If the search results are not relevant, please answer that they are not relevant.

<bot> : In response to the question "What is the expected peak world population in the United Nations 2022 population outlook?" Based on the information numbered 1,2,3 above, my answer is that according to the World Population Prospects 2022 report published by the United Nations Department of Economic and Social Affairs, the global population is expected to reach 8 billion on November 15, 2022. The report also predicts that the world's population is expected to peak at about 10.4 billion in the 2080s and remain at that level until 2100. As a result, the United Nations expects the world's population to peak at about 10.4 billion.

ChatGPT3.5: The expected peak world population in the UN's 2022 population outlook is not explicitly mentioned in the provided essay passages.

Table 5: A Multi-doc QA sample and the output of Our model, ChatGPT3.5 in the last two cells respectively. The answer of Our model starts after the <bot> :. Special tokens are removed when testing ChatGPT3.5. We can see that Our model produces a correct answer while ChatGPT3.5 fails to answer the question.



Figure 7: The heatmap of attention scores on the last 100 tokens when our model begins to generate an answer. The darker represents the higher score. The English version is translated from the Chinese heatmap token by token. The first token of the generated question and the first predicted index are attended to, indicating the utility of the question repetition and index prediction in PAM QA.