

RA-DIT: Retrieval-Augmented Dual Instruction Tuning

Victoria Lin*, Xilun Chen*, Mingda Chen*

Weijia Shi, Maria Lomeli, Rich James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvassy,
Mike Lewis, Luke Zettlemoyer and Scott Wen-tau Yih

{victorialin, xilun, mingdachen, scottyih}@meta.com

 **Meta** @ICLR 2024

LLM with Access to Non-parametric Knowledge

Retrieval-**A**ugmented **D**ual **I**nstruction **T**uning

Can we improve RALMs via instruction tuning s.t.

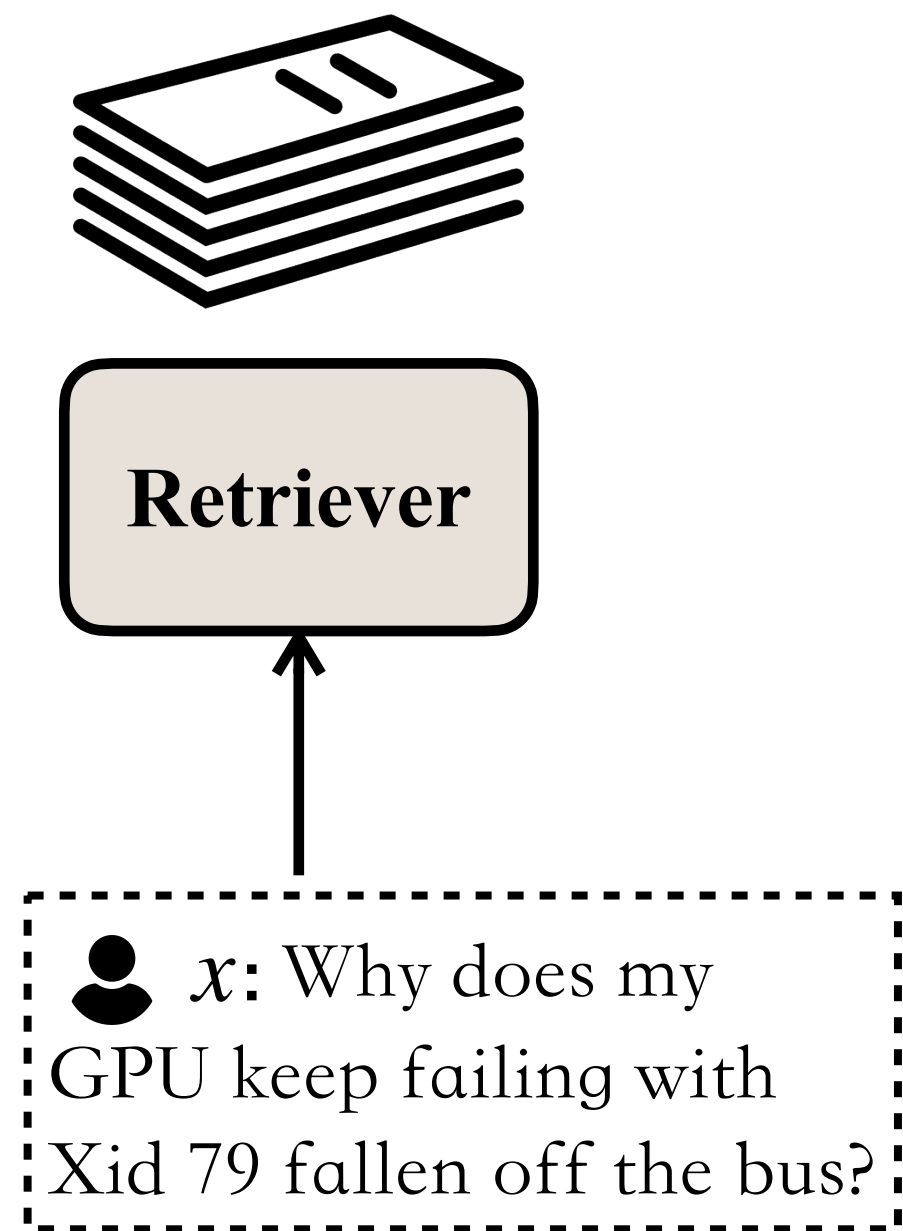
- A. The LLM can learn to better utilize the retrieved content in context
- B. The retriever can find more information relevant to the LLM

Concurrent related work:

Leo et al. 2023. SAIL: Search-Augmented Instruction Learning.

Asai et al. 2023. Self-RAG: Learning to Retrieve, Generate and Critique through Self-Reflection.

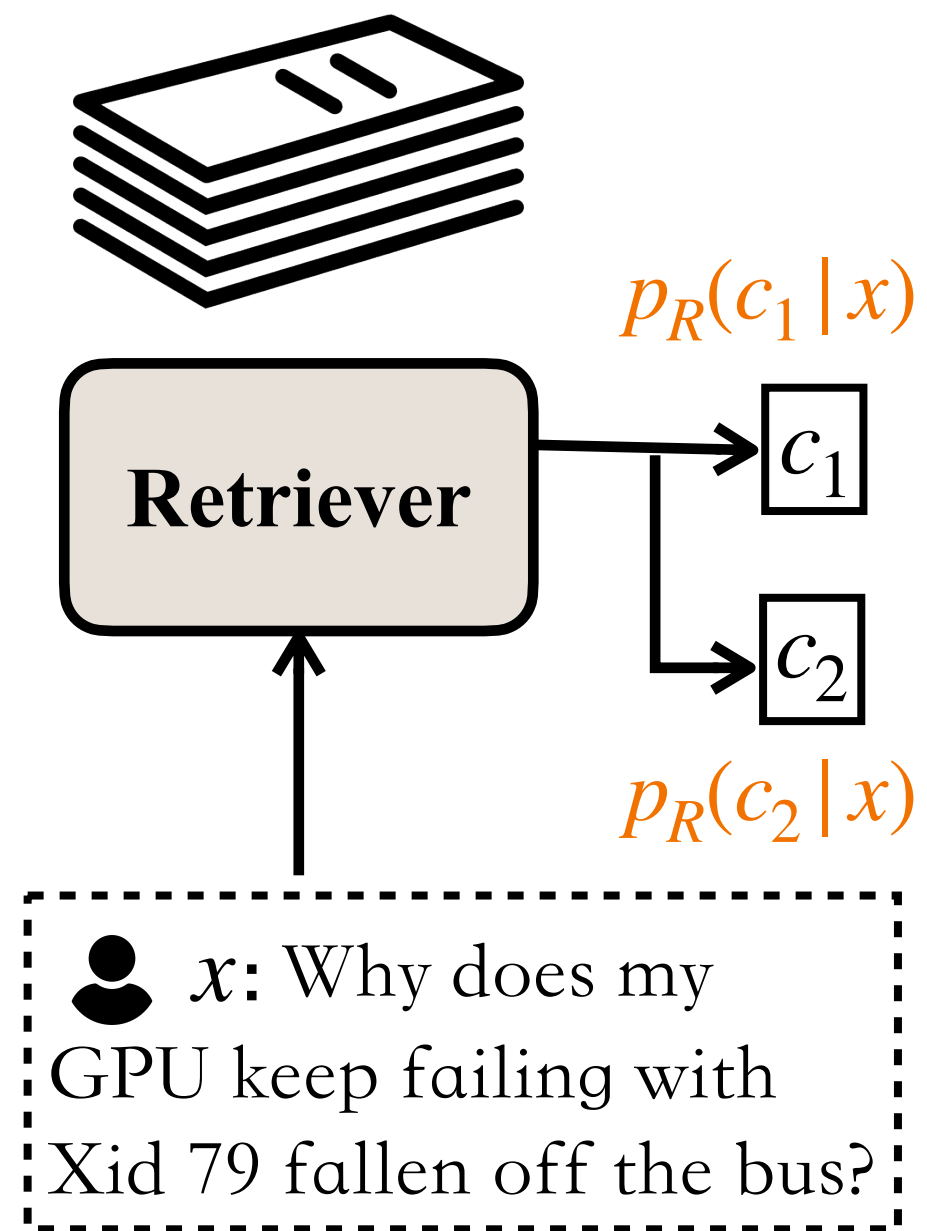
Retrieval-Augmented Dual Instruction Tuning



Q: Why does my GPU keep failing with Xid 79 fallen off the bus? **A:**

LLM

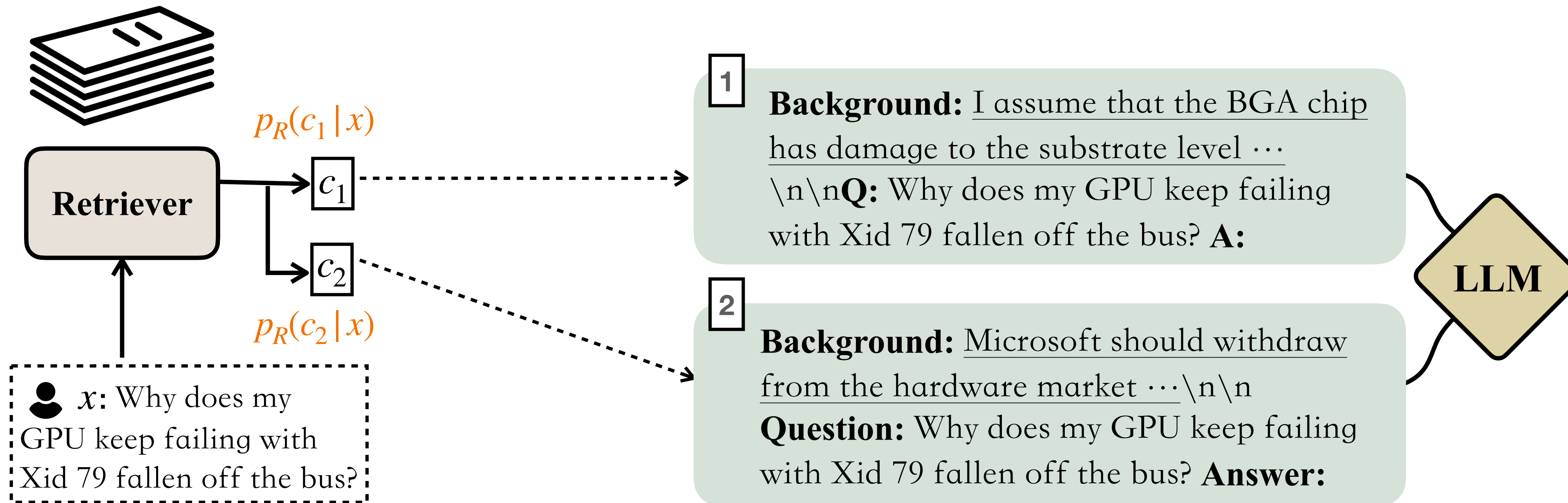
Retrieval-Augmented Dual Instruction Tuning



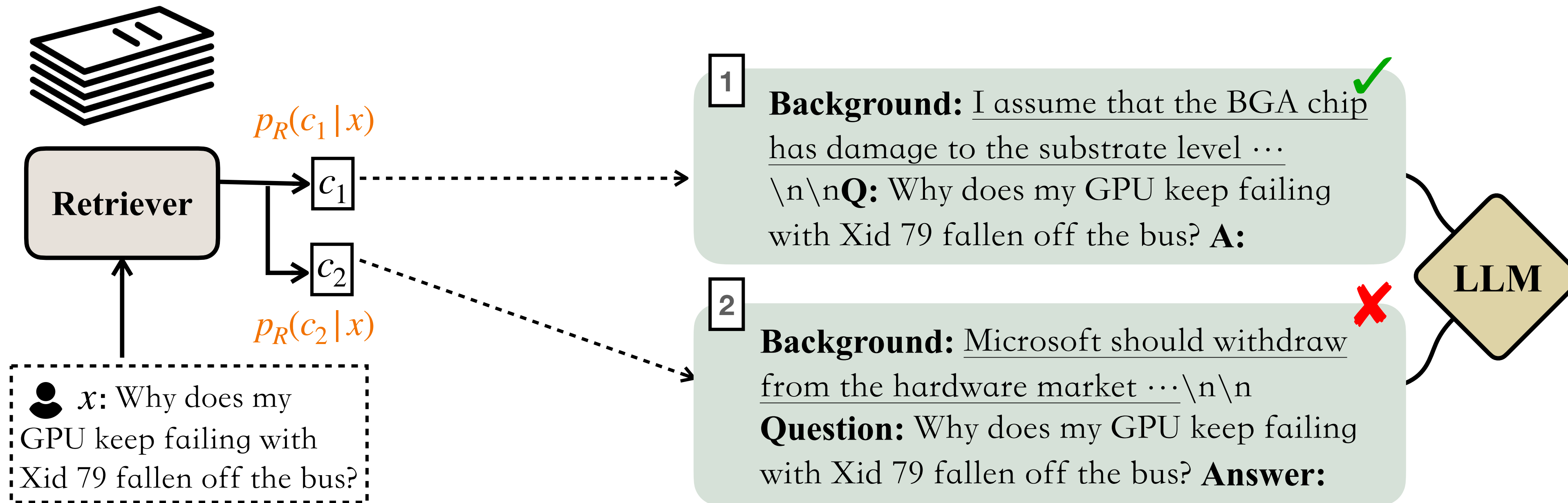
Q: Why does my GPU keep failing with Xid 79 fallen off the bus? **A:**

LLM

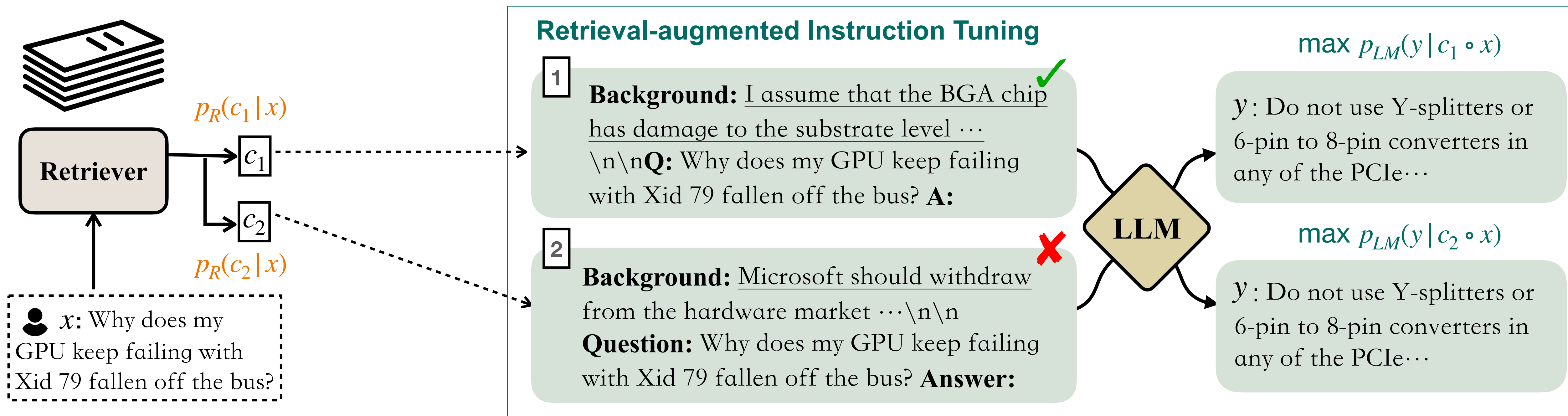
Retrieval-Augmented Dual Instruction Tuning



Retrieval-Augmented Dual Instruction Tuning



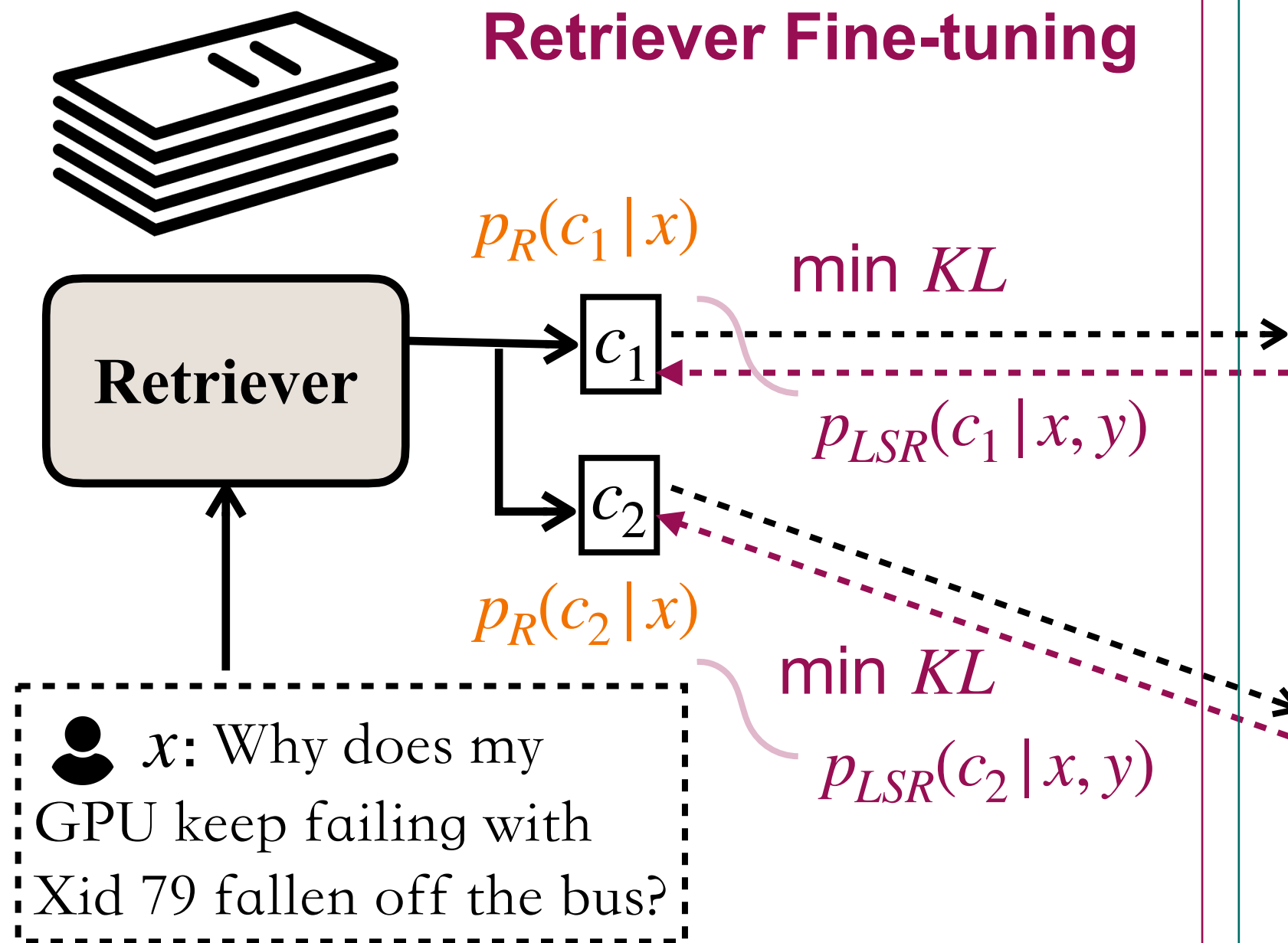
Retrieval-Augmented Dual Instruction Tuning



Retrieval-Augmented Dual Instruction Tuning

Independent Steps

Retriever Fine-tuning



Retrieval-augmented Instruction Tuning

- 1** **Background:** I assume that the BGA chip has damage to the substrate level ...
Q: Why does my GPU keep failing with Xid 79 fallen off the bus? **A:**
- 2** **Background:** Microsoft should withdraw from the hardware market ...
Question: Why does my GPU keep failing with Xid 79 fallen off the bus? **Answer:**



$$\max p_{LM}(y|c_1 \circ x)$$

y: Do not use Y-splitters or 6-pin to 8-pin converters in any of the PCIe...

$$\max p_{LM}(y|c_2 \circ x)$$

y: Do not use Y-splitters or 6-pin to 8-pin converters in any of the PCIe...

Experiment Setup

- LM initialization: Llama 65B
- Retriever initialization: DRAGON+ (Lin et al. 2023)
- Retrieval Corpus: **399M** text chunks
 - **37M** chunks from Wikipedia 2021 (Izacard et al. 2022)
 - **362M** chunks sampled from the 2017-2020 CommonCrawl dumps
- RA-DIT training data:

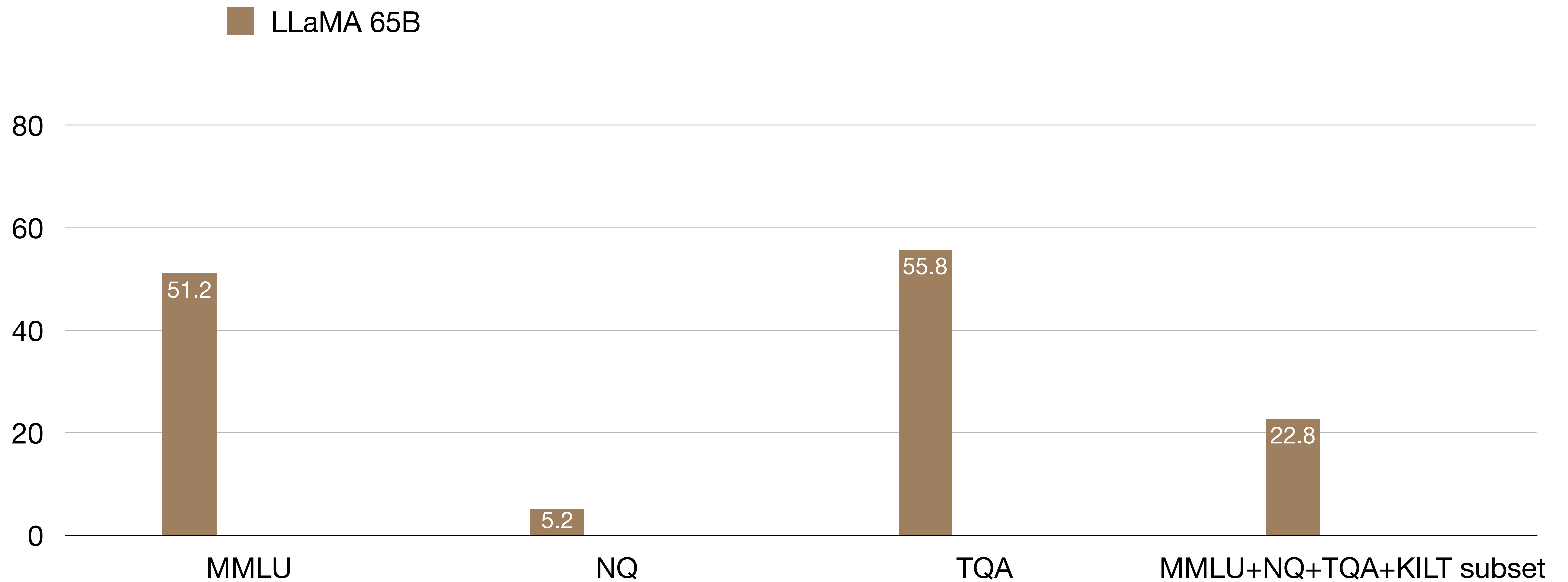
Task	HF identifier	Dataset name	\mathcal{D}_L	\mathcal{D}_R	#Train
Dialogue	oasst1	OpenAssistant Conversations Dataset (Köpf et al., 2023)	✓	✓	31,598
	commonsense_qa	CommonsenseQA (Talmor et al., 2019)	✓	✓	9,741
Open-Domain QA	math_qa	MathQA (Amini et al., 2019)	✓	✓	29,837
	web_questions	Web Questions (Berant et al., 2013)	✓	✓	3,778
	wiki_qa	Wiki Question Answering (Yang et al., 2015)	✓	✓	20,360
	yahoo_answers_qa	Yahoo! Answers QA	✓	✓	87,362
	freebase_qa	FreebaseQA (Jiang et al., 2019)		✓	20,358
	ms_marco*	MS MARCO (Nguyen et al., 2016)		✓	80,143
	coqa	Conversational Question Answering (Reddy et al., 2019)	✓		108,647
Reading Comprehension	drop	Discrete Reasoning Over Paragraphs (Dua et al., 2019)	✓		77,400
	narrativeqa	NarrativeQA (Kočíský et al., 2018)	✓		32,747
	newsqa	NewsQA (Trischler et al., 2017)	✓		74,160
	pubmed_qa	PubMedQA (Jin et al., 2019)	✓	✓	1,000
	quail	QA for Artificial Intelligence (Rogers et al., 2020)	✓		10,246
	quarel	QuaRel (Tafjord et al., 2019)	✓	✓	1,941
	squad_v2	SQuAD v2 (Rajpurkar et al., 2018)	✓		130,319
Summarization	cnn_dailymail	CNN / DailyMail (Hermann et al., 2015)	✓		287,113
	aqua_rat [‡]	Algebra QA with Rationales (Ling et al., 2017)	✓		97,467
Chain-of-thought	ecqa [‡]	Explanations for CommonsenseQ (Aggarwal et al., 2021)	✓		7,598
	gsm8k [‡]	Grade School Math 8K (Cobbe et al., 2021)	✓		7,473
Reasoning	math [‡]	MATH (Hendrycks et al., 2021c)	✓		7,500
	strategyqa [‡]	StrategyQA (Geva et al., 2021)	✓		2,290

* We only used the question-and-answer pairs in the MS MARCO dataset.

Knowledge Utilization

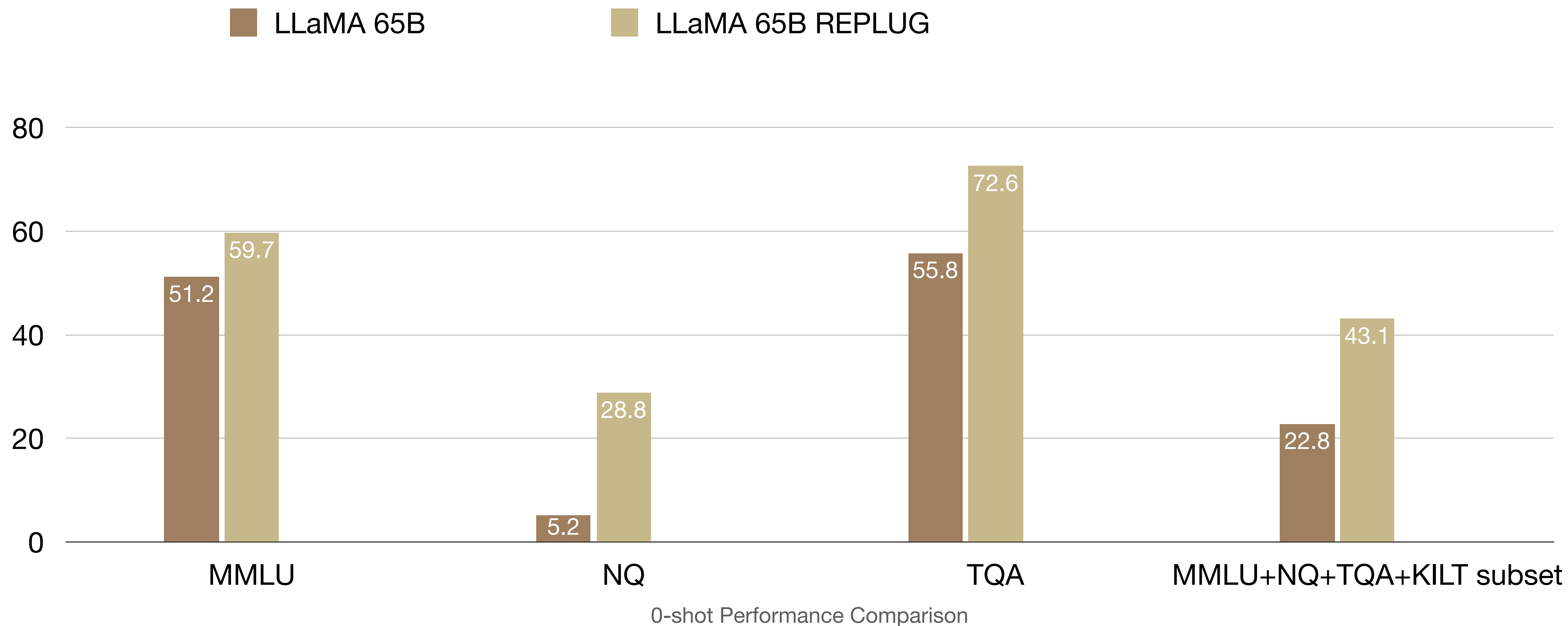
Contextual Awareness

Performance on Knowledge Intensive Tasks

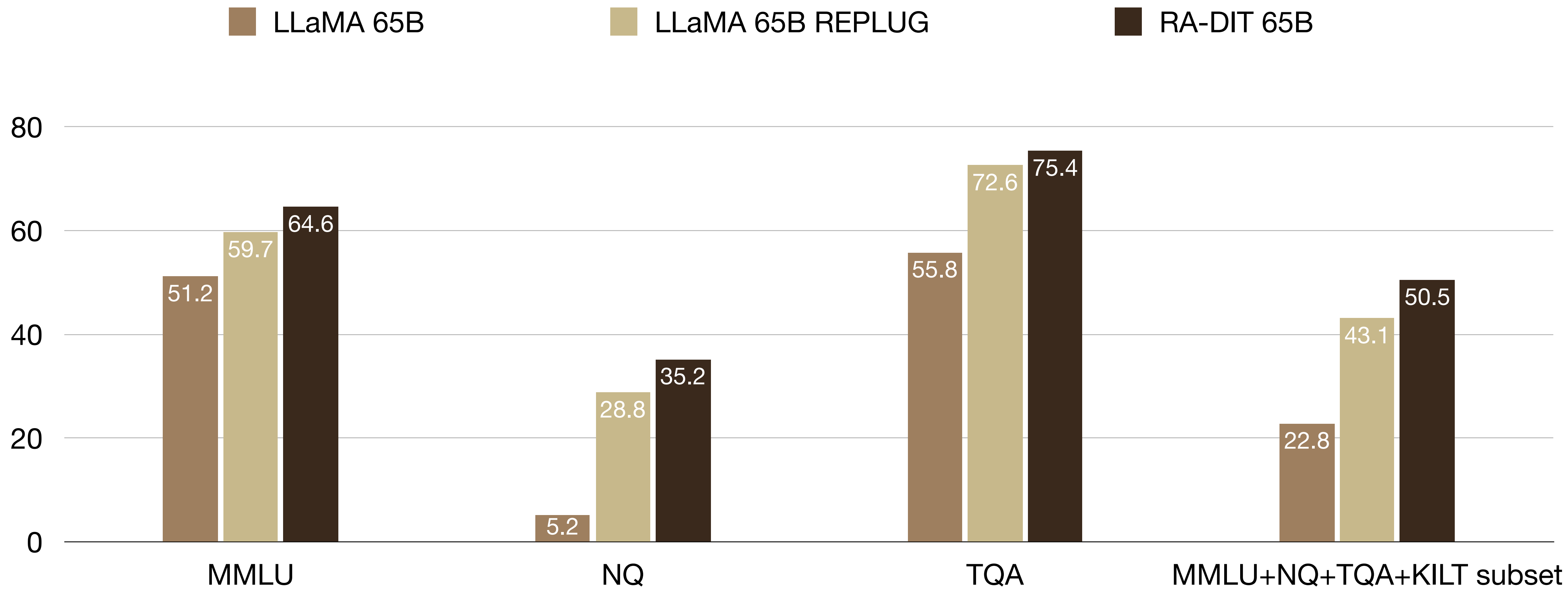


0-shot Performance Comparison

Performance on Knowledge Intensive Tasks



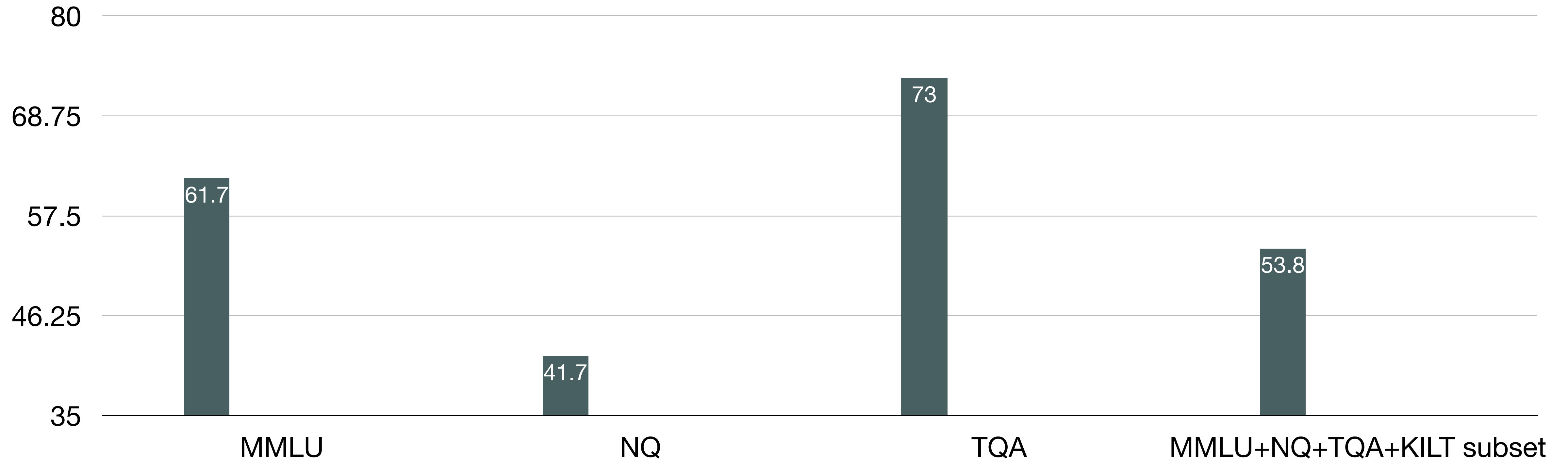
Performance on Knowledge Intensive Tasks



0-shot Performance Comparison

Dual Instruction Tuning Ablation

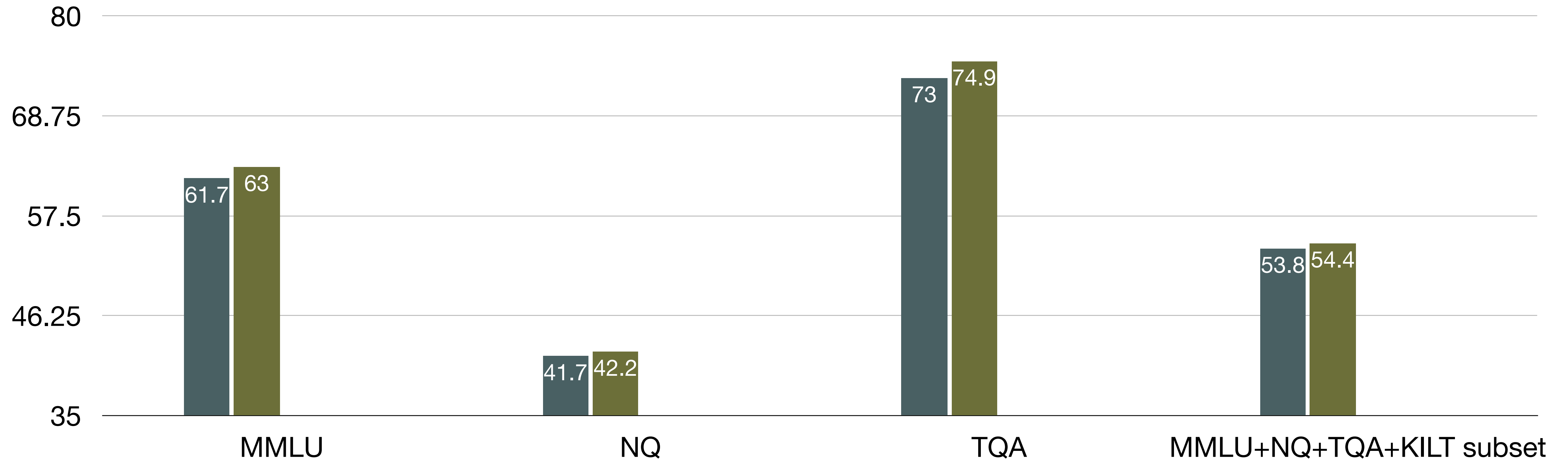
(RePlug) LLaMA 65B + DRAGON LLaMA 65B + FTed DRAGON RIT 65B + DRAGON
(RA-DIT) RIT 65B + FTed DRAGON



5-shot Performance Comparison

Dual Instruction Tuning Ablation

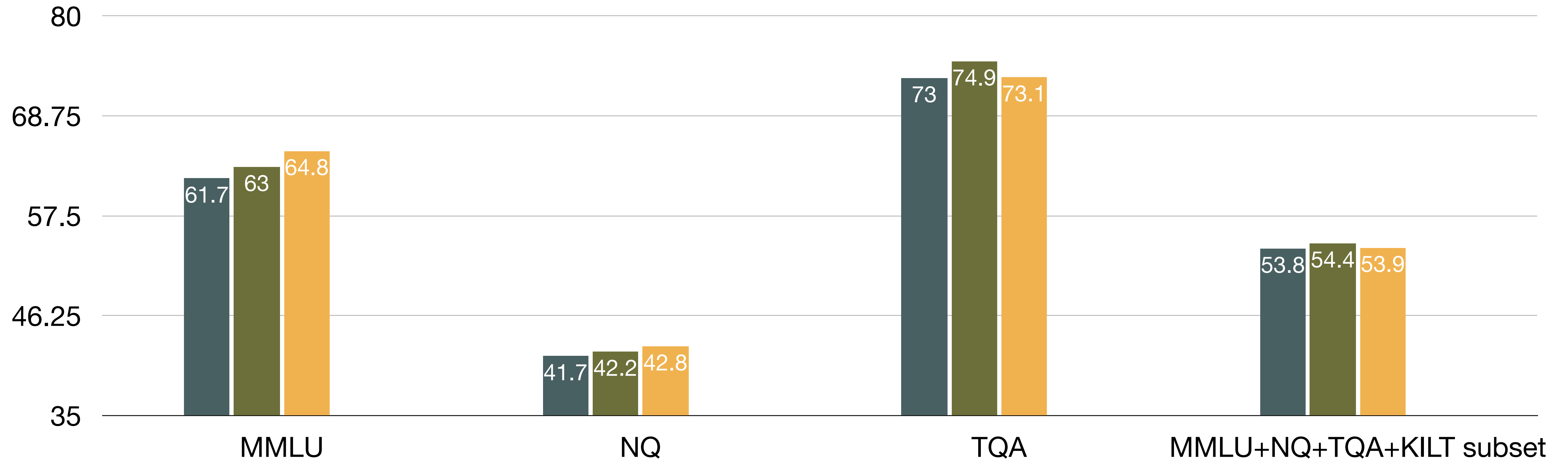
(RePlug) LLaMA 65B + DRAGON LLaMA 65B + FTed DRAGON RIT 65B + DRAGON
(RA-DIT) RIT 65B + FTed DRAGON



5-shot Performance Comparison

Dual Instruction Tuning Ablation

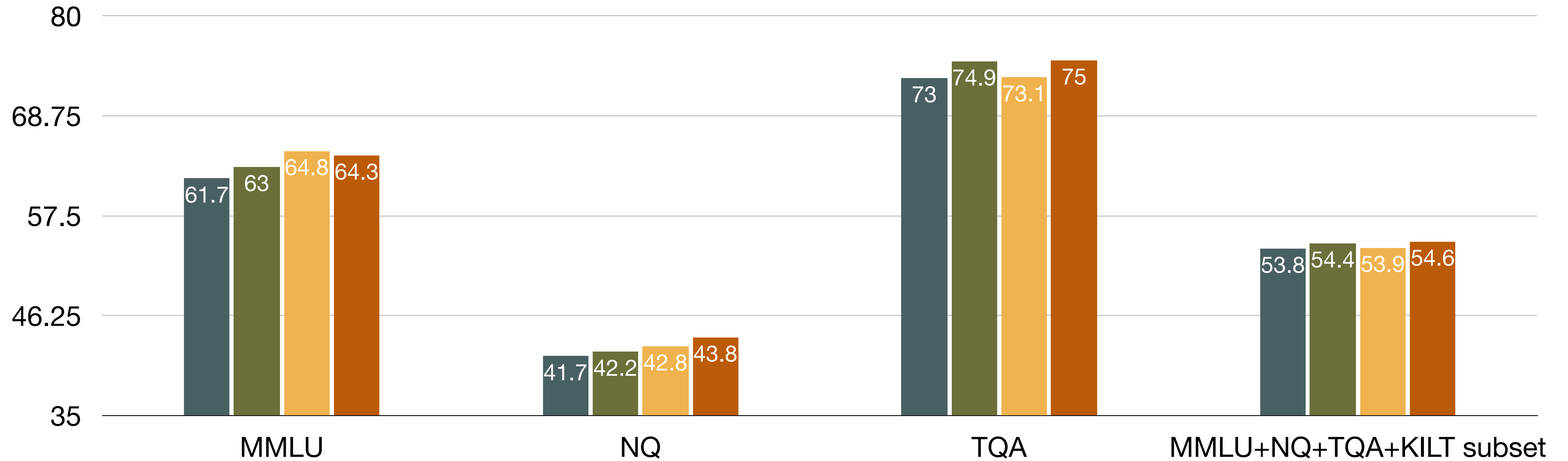
(RePlug) LLaMA 65B + DRAGON LLaMA 65B + FTed DRAGON RIT 65B + DRAGON
(RA-DIT) RIT 65B + FTed DRAGON



5-shot Performance Comparison

Dual Instruction Tuning Ablation

(RePlug) LLaMA 65B + DRAGON LLaMA 65B + FTed DRAGON RIT 65B + DRAGON
(RA-DIT) RIT 65B + FTed DRAGON



5-shot Performance Comparison

Conclusion

- **Fine-tuning with retrieval augmentation** is an effective approach that can improve **the LLM, the retriever** as well as **their integration**.
- See the paper for more ablations and discussions!