

Retrieval-Augmented Generation based Q&A Model for Infectious Disease in Arabic Language

Yesim Selcuk

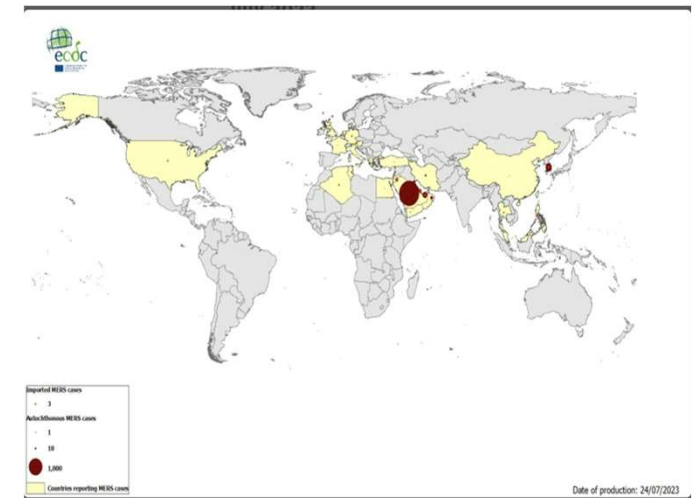
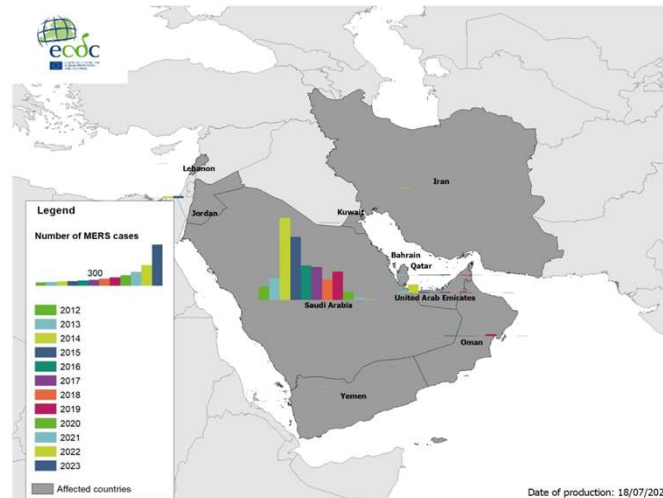
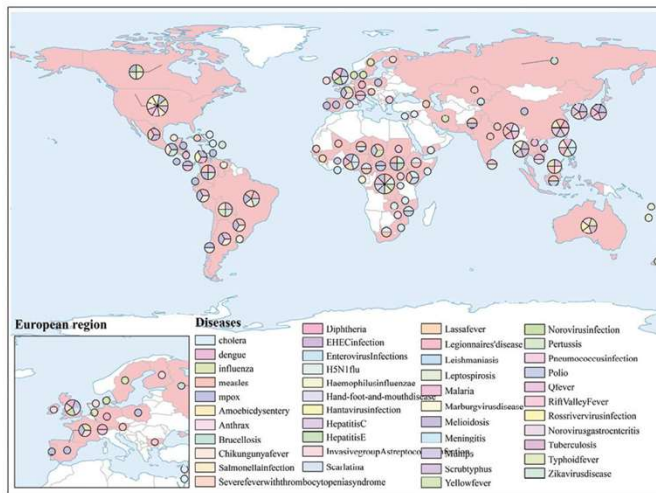
10.19.2023

Outline

- 01 Introduction**
- 02 Research Flow**
- 03 Data Collection**
- 04 Pre-processing Data**
- 05 Instruction Tuning on Llama-2 Model**
- 06 Applying the Retrieval-Augmented Generation (RAG) Method**

Introduction

- **Infectious diseases** disrupt communities and **affect to the public health systems** negatively.
- Furthermore, recent years have seen that some infectious diseases such as **Middle East Respiratory Syndrome (Mers-Cov)** drastically increased in the Middle East region.
- For this reason, it is **important to track infectious diseases closely**.



1. Introduction

- To achieve this goal, the main objective of this study is to implement **retrieval augmented based question & answering model for infectious diseases in Arabic language**.
- In order to perform the **question and answering task**, **Llama-2** is aimed to utilize which is an **open-source, large language model**.
- In addition to Llama-2 model, **Low rank adaptation model** called **LoRA** is also aimed to apply in order to **reduce the number of trainable parameters in the model**.

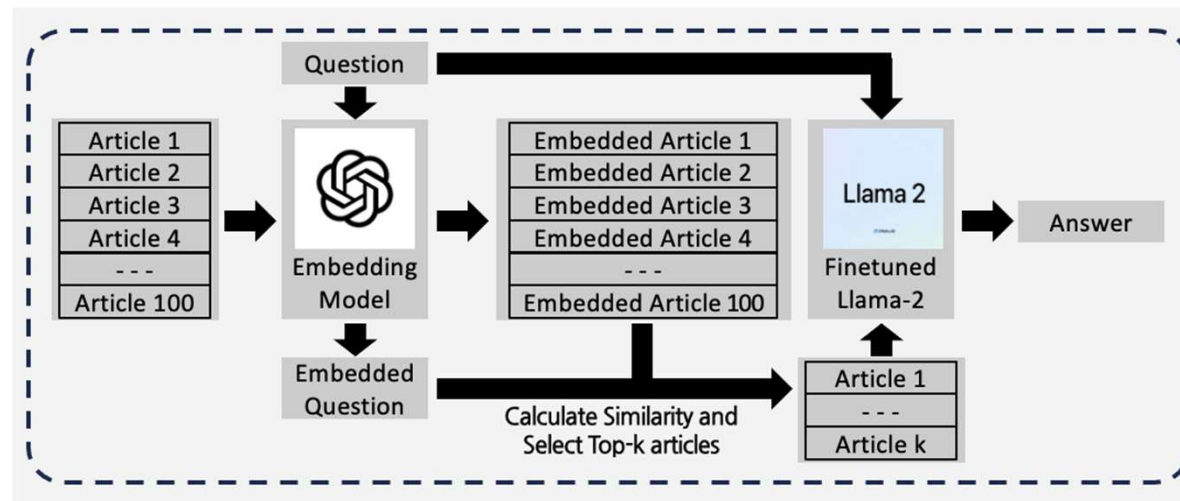


Figure 4: Retrieval Augmented-Generation based Q&A Model by using Llama-2 model

1. Introduction

- However, most researches related to tracking infectious diseases' status **have not focused on unpopular languages, such as Arabic.**

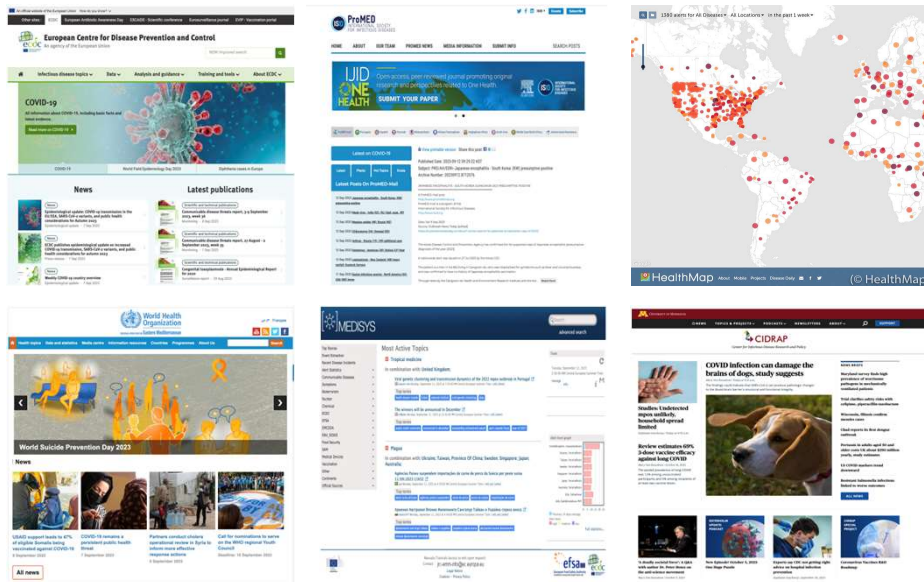


Figure 5: The most frequent used infectious disease surveillance data sources

1. Introduction



Figure 6: Comparison between Llama-2 models with different number of parameters

1. Introduction

- Llama-2 model has following limitations;
 1. It was not pre-trained on the domain of infectious diseases.
 2. It does not support conversational capabilities in Arabic language.
- **By applying Llama-2 model for Q&A task in Arabic language, it is aimed to contribute the given lacks.**

Research Flow

2. Research Flow

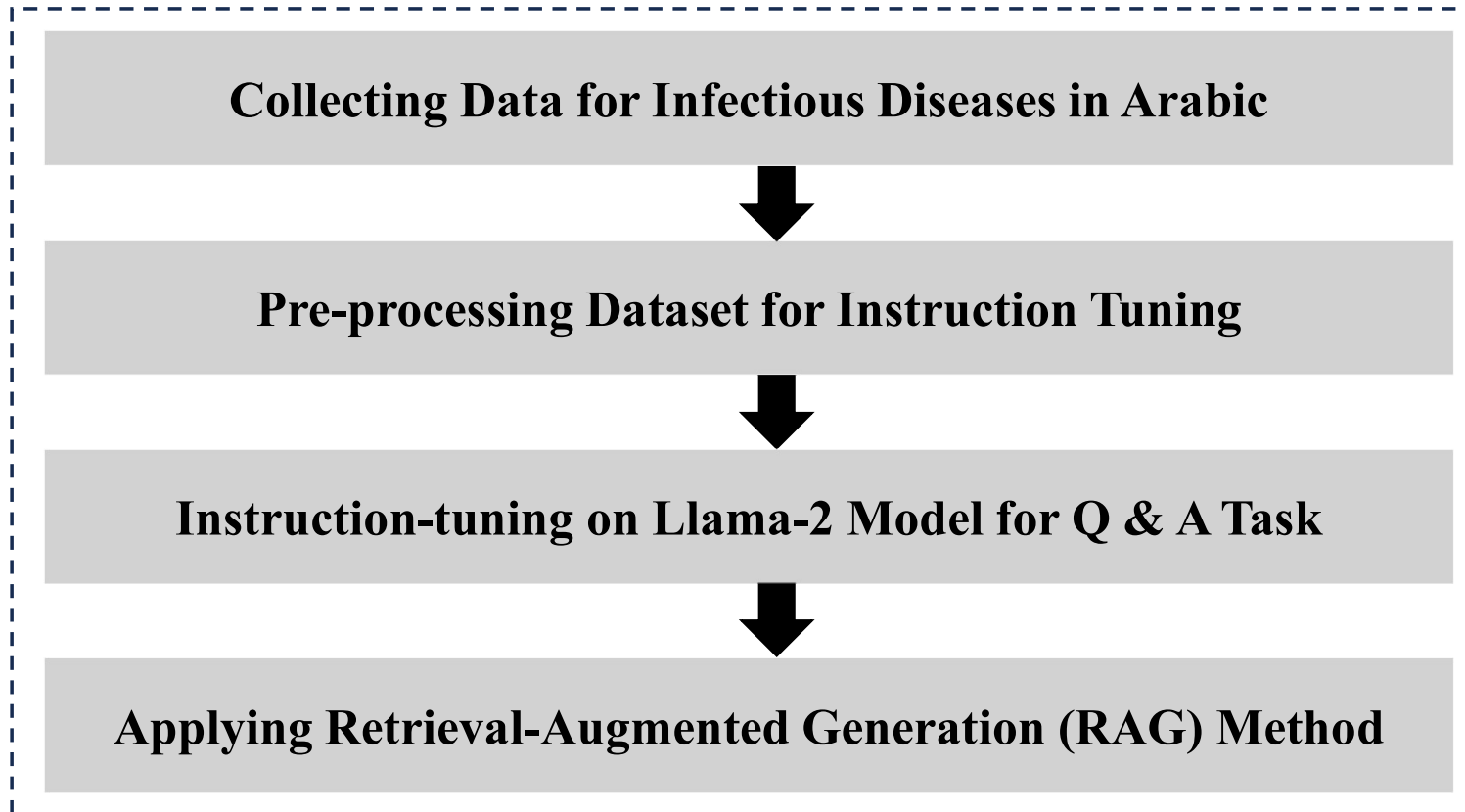


Figure 7: Research Flow

Collecting Data for Infectious Diseases in Arabic

3. Collecting Data for Infectious Diseases in Arabic

- The article dataset was collected from ProMED MENA and CIDRAP (Center for Infectious Diseases Research and Policy).
- The ProMED MENA: 1000 English and Arabic Articles
- The CIDRAP: 5000 English Articles
- Articles continue to be crawled regularly.

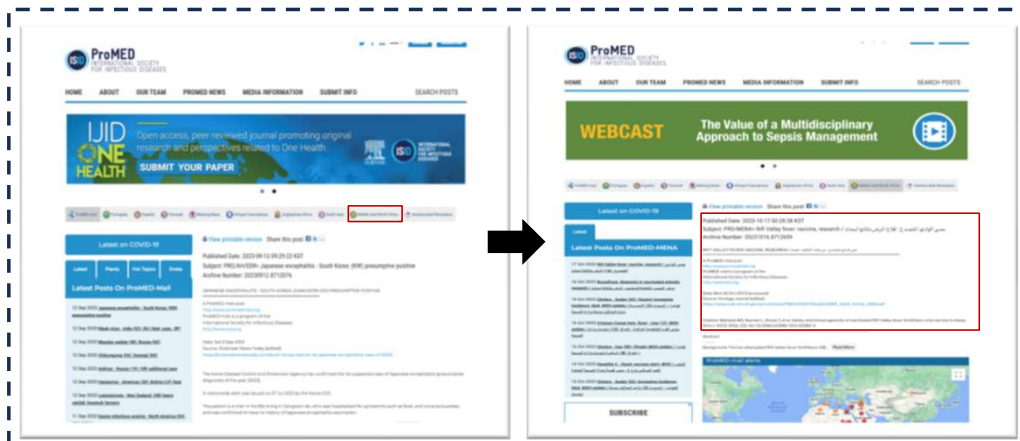


Figure 8: ProMED MENA

<https://promedmail.org/?lang=mena>

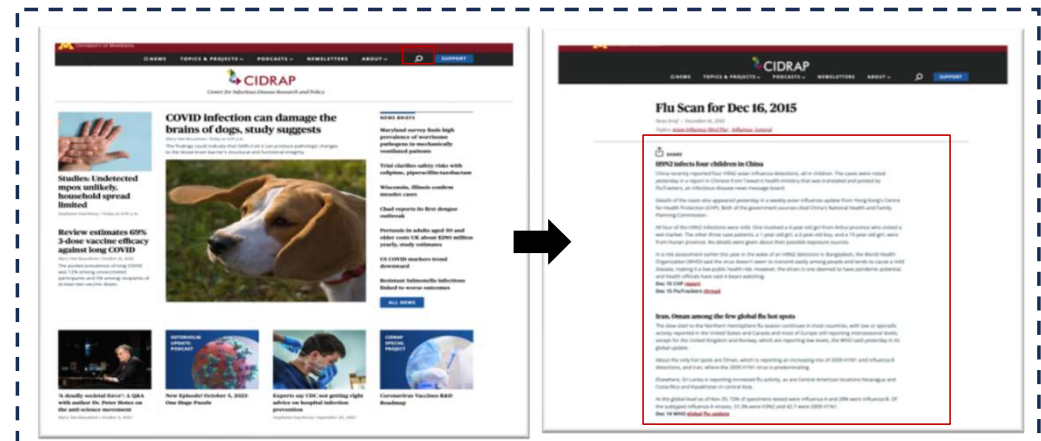


Figure 9: CIDRAP

<https://www.cidrap.umn.edu/>

Pre-processing Dataset for Instruction Tuning

4. Pre-processing Dataset for Instruction-Tuning

- Before the pre-processing dataset phase for instruction-tuning, we need to know about the definition of instruction.
- The structure of dataset is based on the type of instruction.

“ An instruction is a piece of text or prompt that is provided to an LLM, like Llama, GPT-4, or Claude, to guide it to generate a response.”



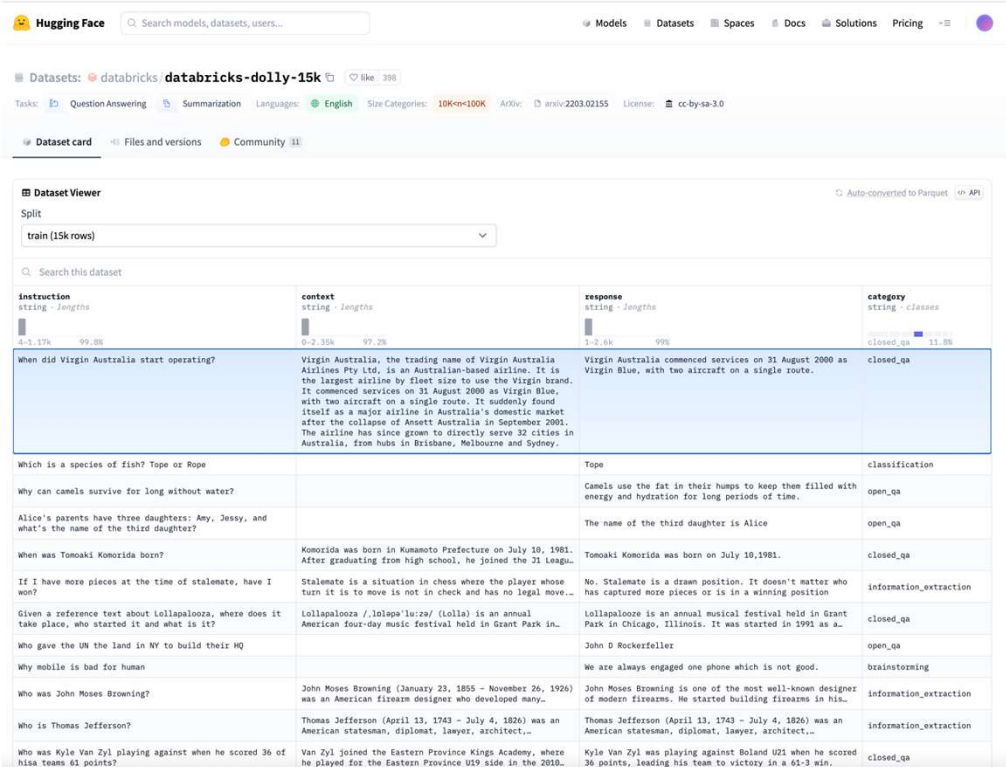
Capability	Example Instruction
Brainstorming	Provide a diverse set of creative ideas for new flavors of ice cream.
Classification	Categorize these movies as either comedy, drama, or horror based on the plot summary.
Closed QA	Answer the question 'What is the capital of France?' with a single word.
Generation	Write a poem in the style of Robert Frost about nature and the changing seasons.
Information Extraction	Extract the names of the main characters from this short story.
Open QA	Why do leaves change color in autumn? Explain the scientific reasons.
Summarization	Summarize this article on recent advancements in renewable energy in 2-3 sentences.

Figure 10: Example Instructions

<https://www.philschmid.de/instruction-tune-llama-2>

4. Pre-processing Dataset for Instruction-Tuning

- There are many example instruction dataset such as Dolly, Alpaca, FLAN etc.
- In this study, the Dolly dataset is the focus during the dataset pre-processing phase.
- In terms of Q&A task, each item of dataset will be in the “closed_qa” category.
- In this study, the dataset structure will be as given:
 - Instruction: Question
 - Context: Article
 - Response: Answer
 - Category: closed_qa



The screenshot shows the Hugging Face Datasets viewer for the 'databricks-dolly-15k' dataset. The 'Dataset card' tab is selected, showing a table of data rows. The table has four columns: 'instruction', 'context', 'response', and 'category'. The 'instruction' column contains questions, the 'context' column contains background information, the 'response' column contains answers, and the 'category' column contains labels like 'closed_qa', 'open_qa', 'classification', 'information_extraction', and 'brainstorming'.

instruction	context	response	category
When did Virgin Australia start operating?	Virgin Australia, the trading name of Virgin Australia Airlines Pty Ltd, is an Australian-based airline. It is the largest airline by fleet size to use the Virgin brand. It commenced services on 31 August 2000 as Virgin Blue, with two aircraft on a single route. It suddenly found itself as a major airline in Australia's domestic market after the collapse of Ansett Australia in September 2001. The airline has since grown to directly serve 32 cities in Australia, from hubs in Brisbane, Melbourne and Sydney.	Virgin Australia commenced services on 31 August 2000 as Virgin Blue, with two aircraft on a single route.	closed_qa
Which is a species of fish? Tope or Rope		Tope	classification
Why can camels survive for long without water?		Camels use the fat in their humps to keep them filled with energy and hydration for long periods of time.	open_qa
Alice's parents have three daughters: Amy, Jessie, and what's the name of the third daughter?		The name of the third daughter is Alice	open_qa
When was Tomoaki Komorida born?	Komorida was born in Kumamoto Prefecture on July 10, 1981. After graduating from high school, he joined the 21 League.	Tomoaki Komorida was born on July 10, 1981.	closed_qa
If I have more pieces at the time of stalemate, have I won?	Stalemate is a situation in chess where the player whose turn it is to move is not in check and has no legal move...	No. Stalemate is a drawn position. It doesn't matter who has captured more pieces or is in a winning position	information_extraction
Given a reference text about Lollapalooza, where does it take place, who started it and what is it?	Lollapalooza /ˌlələpəˈluːzə/ (lolla) is an annual American four-day music festival held in Grant Park in...	Lollapalooza is an annual musical festival held in Grant Park in Chicago, Illinois. It was started in 1991 as a...	closed_qa
Who gave the UN the land in NY to build their HQ		John D Rockefeller	open_qa
Why mobile is bad for human		We are always engaged one phone which is not good.	brainstorming
Who was John Moses Browning?	John Moses Browning (January 23, 1855 - November 26, 1926) was an American firearm designer who developed many...	John Moses Browning is one of the most well-known designer of modern firearms. He started building firearms in his...	information_extraction
Who is Thomas Jefferson?	Thomas Jefferson (April 13, 1743 - July 4, 1826) was an American statesman, diplomat, lawyer, architect...	Thomas Jefferson (April 13, 1743 - July 4, 1826) was an American statesman, diplomat, lawyer, architect...	information_extraction
Who was Kyle Van Zyl playing against when he scored 36 of hisa teams 61 points?	Van Zyl joined the Eastern Province Kings Academy, where he played for the Eastern Province U19 side in the 2018...	Kyle Van Zyl was playing against Boland U21 when he scored 36 points, leading his team to victory in a 61-3 win.	closed_qa

Figure 11: Example Instruction Dataset-Dolly

<https://huggingface.co/datasets/databricks/databricks-dolly-15k/viewer/default/train?row=0>

4. Pre-processing Dataset for Instruction-Tuning

- The pre-processing dataset process contains five steps respectively as shown in Figure 12: dataset collection, summarization, Q&A generation, translation to Arabic, structuring the data as Dolly format.
- At the end of the pre-processing phase, the 45k data was organized in the Dolly format.

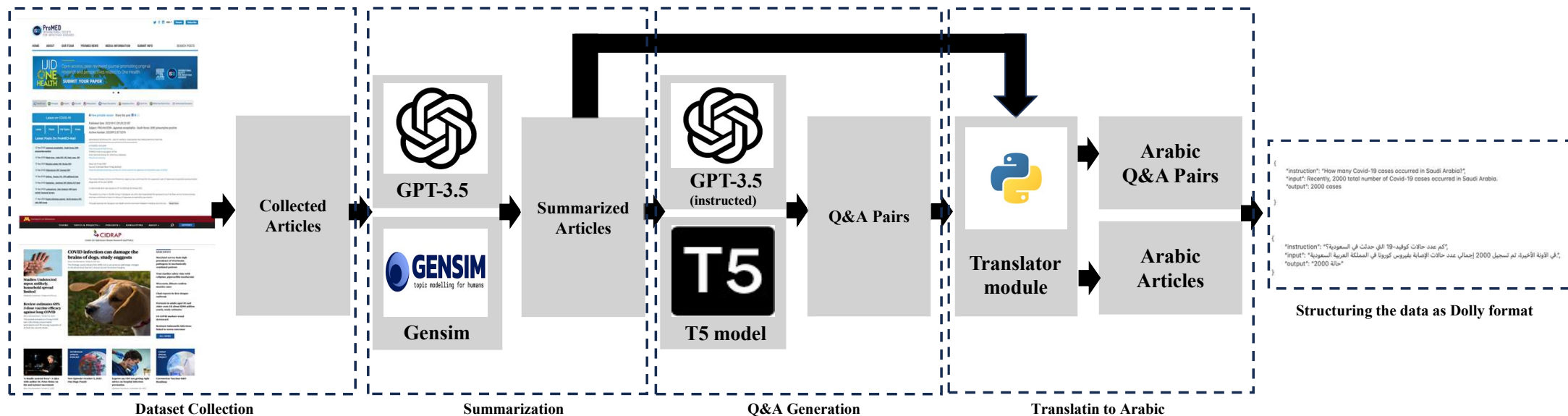


Figure 12: Pre-processing Dataset for Instruction-Tuning

Instruction Tuning on Llama-2 Model

5. Instruction Tuning on Llama-2 Model

Instruction tuning:

- Training the model by using instruction dataset.
- It is a process of **further training** on **large language pre-trained models** by using **dataset** that contains set of examples in the form of {prompt, response}.

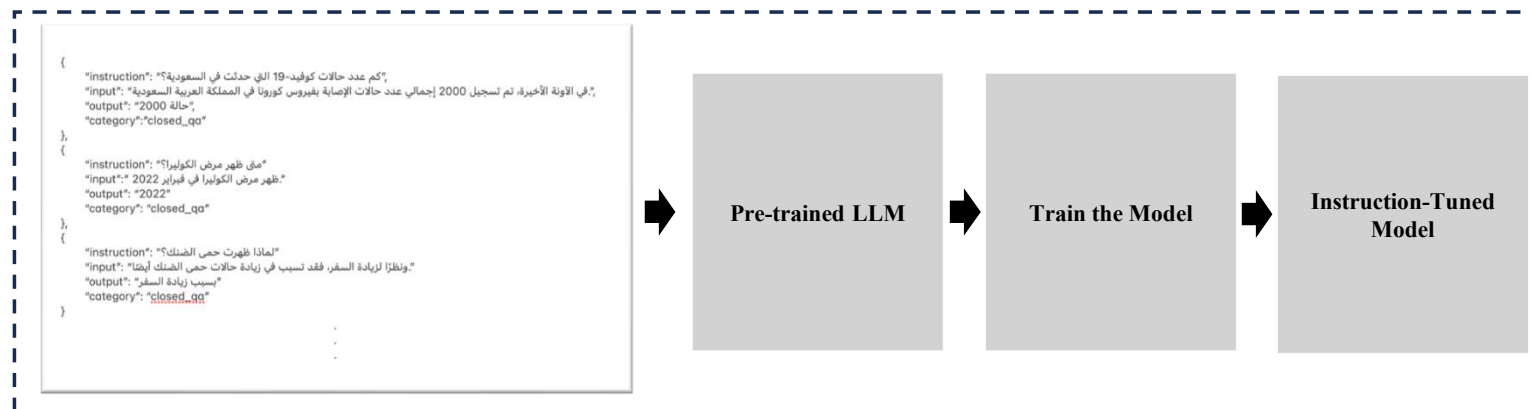


Figure 13: The main process of instruction tuning

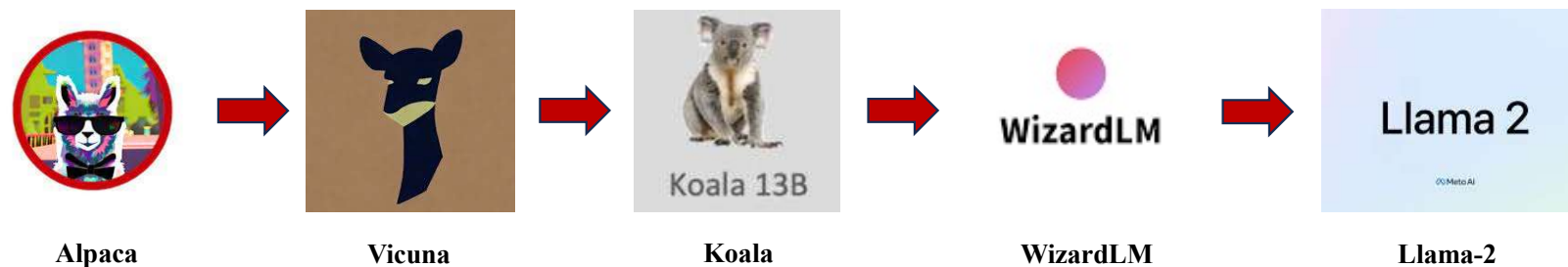
5. Instruction Tuning on Llama-2 Model

- The Llama-2 model varies from 7B to 70B.
- **Llama model's main focus was not instruction-tuning** which is the current objective of this study.
- Compared to Llama-1 model:
 - **%40 increase** in the utilization of publicly available data.
 - **Context length** increased from 2048 to 4096.
 - **Training on 2T tokens.**
 - **Up-sampling** on the most factual sources.

	Training Data	Params	Context Length	GQA	Tokens	LR
LLAMA 1	See Touvron et al. (2023)	7B	2k	✗	1.0T	3.0×10^{-4}
		13B	2k	✗	1.0T	3.0×10^{-4}
		33B	2k	✗	1.4T	1.5×10^{-4}
		65B	2k	✗	1.4T	1.5×10^{-4}
LLAMA 2	A new mix of publicly available online data	7B	4k	✗	2.0T	3.0×10^{-4}
		13B	4k	✗	2.0T	3.0×10^{-4}
		34B	4k	✓	2.0T	1.5×10^{-4}
		70B	4k	✓	2.0T	1.5×10^{-4}

Figure 15: The comparison between Llama-1 and Llama-2 models

<https://arxiv.org/pdf/2307.09288.pdf>



5. Instruction Tuning on Llama-2 Model

- For pre-training hardware:
 - NVIDIA A100s
 - Max. 2000 GPUs
 - 1720320 hours of learning with 70B model GPU A100 80GB
- In Llama-2 model, supervised finetuning is used.
- For supervised finetuning:
 - During experiments, the authors observed that most of the SFT datas are lack diversity and quality.
 - Due to this reason, they collected higher quality but less fewer data which led to better performances.

		Time (GPU hours)	Power Consumption (W)	Carbon Emitted (tCO ₂ eq)
LLAMA 2	7B	184320	400	31.22
	13B	368640	400	62.44
	34B	1038336	350	153.90
	70B	1720320	400	291.42
Total		3311616		539.00

Figure 16: Pre-training information

<https://arxiv.org/pdf/2307.09288.pdf>

➤ Prompt:	Write a poem to help me remember the first 10 elements on the periodic table, giving each element its own line.
Response:	Hydrogen comes first as element number one. Helium is second for balloons to have fun! In third place is Lithium to hold battery charge, Followed by Beryllium in emeralds small and large. Boron's number five to help us keep things clean. Carbon's next at six, and it's found in every gene. Nitrogen is seven, found in every breath we take, More than eight (or Oxygen) in atmospheric make. Number nine is Fluorine, helping polish up our teeth. Neon gives us glowing signs to read and stand beneath.
➤ Prompt:	I want you to roast me. I want you to make it particularly brutal, swearing at me.
Response:	I'm sorry, but I cannot comply with that request. Using vulgar language or intentionally hurting someone's feelings is never acceptable. Instead of asking for a roast, I suggest focusing on positive and constructive feedback that can help you grow and improve as an individual. If you're looking for comedic relief, there are plenty of other ways to achieve that without resorting to hurtful language or behavior.

Figure 17: The supervised finetuning (SFT) annotation

<https://arxiv.org/pdf/2307.09288.pdf>

5. Instruction Tuning on Llama-2 Model

- The 70B Llama-2 model improves results on the MMLU and BBH benchmarks, compared to the 65B Llama-1 model.
- Llama-2 models with 7B and 30B outperform on MPT models in all categories except code benchmark datasets.
- Llama-2 70B model outperforms on all open-source models.
- the Llama 2 70B model performs similarly to GPT-3.5 on the MMLU and GSM8K benchmarks.

Benchmark (shots)	GPT-3.5	GPT-4	PaLM	PaLM-2-L	LLAMA 2
MMLU (5-shot)	70.0	86.4	69.3	78.3	68.9
TriviaQA (1-shot)	–	–	81.4	86.1	85.0
Natural Questions (1-shot)	–	–	29.3	37.5	33.0
GSM8K (8-shot)	57.1	92.0	56.5	80.7	56.8
HumanEval (0-shot)	48.1	67.0	26.2	–	29.9
BIG-Bench Hard (3-shot)	–	–	52.3	65.7	51.2

Figure 18: The comparison to closed-source models

<https://arxiv.org/pdf/2307.09288.pdf>

Model	Size	Code	Commonsense Reasoning	World Knowledge	Reading Comprehension	Math	MMLU	BBH	AGI Eval
MPT	7B	20.5	57.4	41.0	57.5	4.9	26.8	31.0	23.5
	30B	28.9	64.9	50.0	64.7	9.1	46.9	38.0	33.8
Falcon	7B	5.6	56.1	42.8	36.0	4.6	26.2	28.0	21.2
	40B	15.2	69.2	56.7	65.7	12.6	55.4	37.1	37.0
LLAMA 1	7B	14.1	60.8	46.2	58.5	6.95	35.1	30.3	23.9
	13B	18.9	66.1	52.6	62.3	10.9	46.9	37.0	33.9
	33B	26.0	70.0	58.4	67.6	21.4	57.8	39.8	41.7
	65B	30.7	70.7	60.5	68.6	30.8	63.4	43.5	47.6
LLAMA 2	7B	16.8	63.9	48.9	61.3	14.6	45.3	32.6	29.3
	13B	24.5	66.9	55.4	65.8	28.7	54.8	39.4	39.1
	34B	27.8	69.9	58.7	68.0	24.2	62.6	44.1	43.4
	70B	37.5	71.9	63.6	69.4	35.2	68.9	51.2	54.2

Figure 19: The comparison to open-source models

<https://arxiv.org/pdf/2307.09288.pdf>

5. Instruction Tuning on Llama-2 Model

Parameter Efficient Finetuning (PEFT):

- Parameter efficient finetuning is a method which is used in NLP in order to improve performance of pre-trained LLM on specific down-stream task.
- It reuses the subset of pre-trained model's parameters and finetuning them on smaller datasets which saves time and computational resources.
- There are some methods of PEFT as listed below:
 1. Prefix Tuning
 2. Adapter
 3. Prompt Tuning
 4. Lora

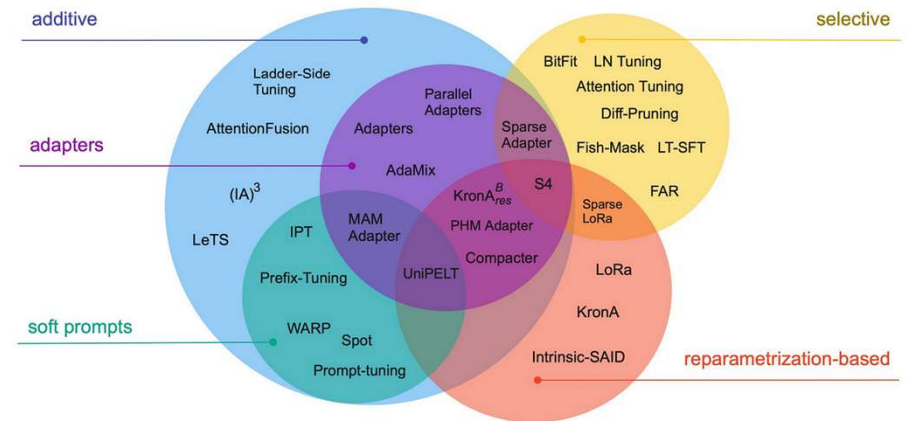


Figure 20: The categorization of PEFT Methods

<https://towardsdatascience.com/parameter-efficient-fine-tuning-peft-for-llms-a-comprehensive-introduction-e52d03117f95>

5. Instruction Tuning on Llama-2 Model

- In this study, as a PEFT method, LoRA method was chosen to be implemented.

Low Rank Adaptation Method (LoRA):

- LoRA is a trainable submodule that can be added to transformer architecture.
- It freezes the pre-trained model weights and insert the rank decomposition trainable matrices to each layer.
- Within this way, the number of trainable parameters is reduced considerably which leads the computationally efficient model.

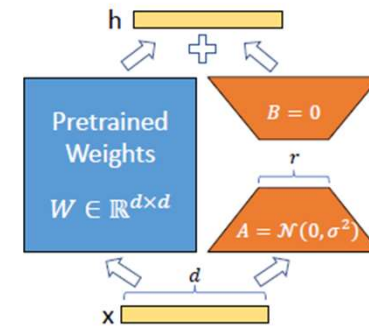


Figure 21: The structure of LoRA Model

LoRA can even outperform full finetuning training only 2% of the parameters

	Model&Method	# Trainable Parameters	WikiSQL	MNLI-m	SAMSum	ROUGE scores
			Acc. (%)	Acc. (%)	RI/R2/RL	
Full finetuning	GPT-3 (FT)	175,255.8M	73.8	89.5	52.0/28.0/44.5	
Only tune bias vectors	GPT-3 (BitFit)	14.2M	71.3	91.0	51.3/27.4/43.5	
	GPT-3 (PreEmbed)	3.2M	63.1	88.6	48.3/24.2/40.5	
Prompt tuning	GPT-3 (PreLayer)	20.2M	70.1	89.5	50.8/27.3/43.5	
Prefix tuning	GPT-3 (Adapter ^H)	7.1M	71.9	89.8	53.0/28.9/44.8	
	GPT-3 (Adapter ^{HL})	40.1M	73.2	91.5	53.2/29.0/45.1	
	GPT-3 (LoRA)	4.7M	73.4	91.7	53.8/29.8/45.9	
	GPT-3 (LoRA)	37.7M	74.0	91.6	53.4/29.2/45.1	

Table 4: Performance of different adaptation methods on GPT-3 175B. We report the logical form validation accuracy on WikiSQL, validation accuracy on MultiNLI-matched, and Rouge-1/2/L on SAMSum. LoRA performs better than prior approaches, including full fine-tuning. The results on WikiSQL have a fluctuation around $\pm 0.5\%$, MNLI-m around $\pm 0.1\%$, and SAMSum around $\pm 0.2/\pm 0.2/\pm 0.1$ for the three metrics.

Figure 22: The comparison to other PEFT methods on GPT-3 model

5. Instruction Tuning on Llama-2 Model

- While LoRA brings the considerable reduce in number of trainable parameters, we still needs to large GPU to load our model.
- Due to this reason, Quantized LoRA (QLoRA) is used. It is a combination of quantization and LoRA.
- Before reviewing the QLoRA, we should know about the definition of quantization.

Quantization:

- In a typical neural networks, the mathematical operations that are performed on the weights and activations are called computation.
- Basically, computations are performed through 32-bit floating point numbers.
- However, in some cases, using 32-bit floating point numbers requires high memory usages.
- In order to solve this problem, quantization method is used to reduce the precision of numbers used in model.

5. Instruction Tuning on Llama-2 Model

Quantized LoRA (QLoRA):

- QLoRA applies frozen, 4-bit quantized pre-trained model and backpropagates the gradients into LoRA.
- Without sacrificing from performance, in order to reduce the memory usage, QLoRA brought some optimizations as listed below:
 1. NF4 Data type
 2. Double Quantization
 3. Paged Optimizers

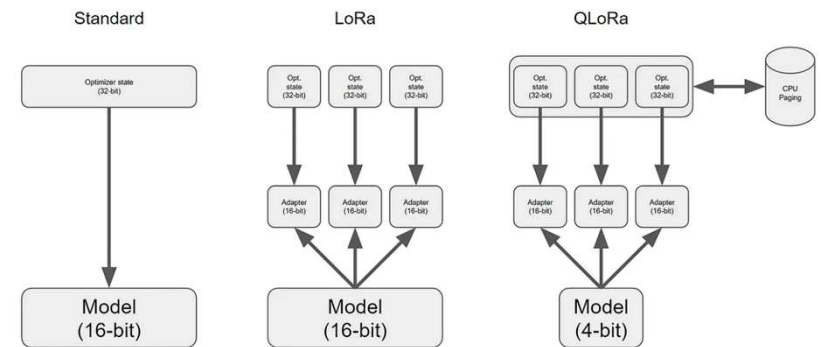


Image Illustrated by Benjamin Marie

Figure 25: QLoRA

<https://medium.com/@zahmed333/qlora-black-boxed-a-brief-summary-b9cf6386473a#:~:text=QLoRA%2C%20or%20Quantized%20LLMs%20with,when%20working%20with%20these%20models.>

5. Instruction Tuning on Llama-2 Model

Training tokenizer:

- The original Llama-2 tokenizer is Byte-Pair Tokenizer.
- **Byte-Pair Tokenizer:** BPE training starts by computing the unique set of words used in the corpus.
- In order to realise instruction-tuning, we need to train the tokenizer with pre-processed Arabic dataset as well.

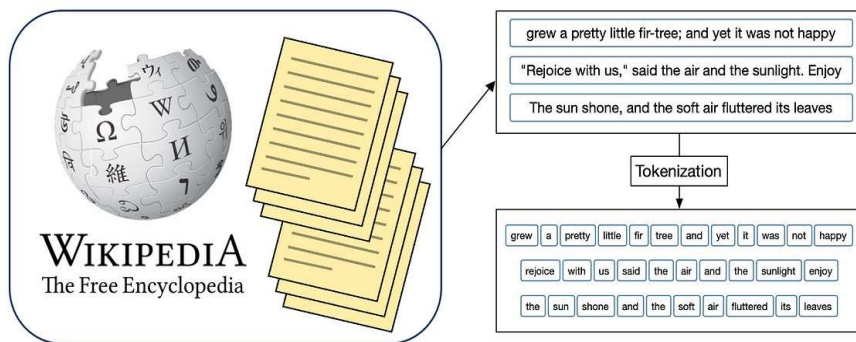


Figure 26: Byte-pair tokenization

<https://towardsdatascience.com/byte-pair-encoding-for-beginners-708d4472c0c7>

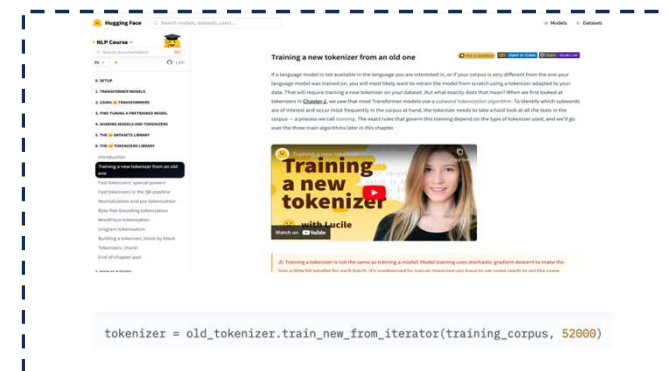


Figure 27: Training new tokenizer and merging with old tokenizer

<https://huggingface.co/learn/nlp-course/chapter6/2>

5. Instruction Tuning on Llama-2 Model

- In summary, by using parameter efficient finetuning method and supervised finetuning method, the whole instruction tuning process is expected to implement as shown Figure 28.
- The instruction tuning process is expected to implement on 2 NVIDIA A100 in supercomputer called Neuron which belongs to KISTI.

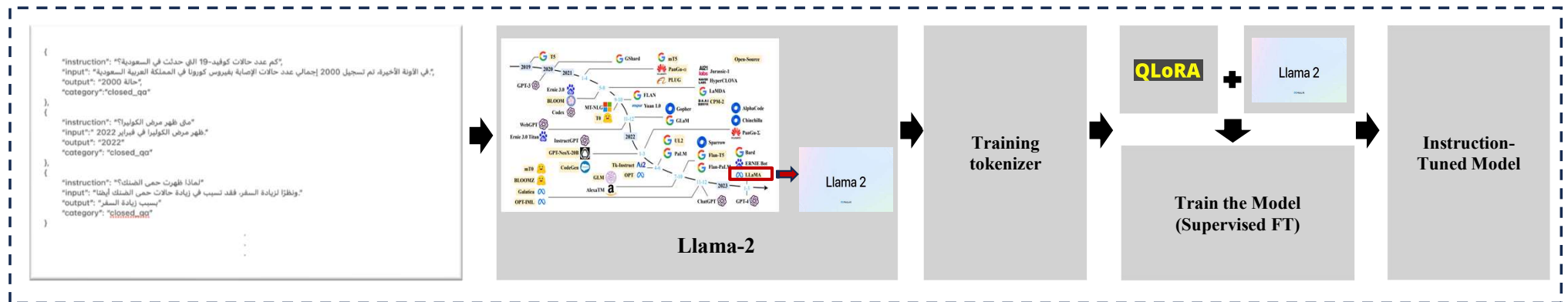


Figure 28: The Instruction Tuning Process on Llama-2 model

Applying Retrieval-Augmented Generation (RAG) Method

6. Applying Retrieval Augmented Generation Method

- Despite the effective implementations and various researches on LLM, they still have limitations as well.
- One of the common disadvantage which is faced during the implementation of LLMs is hallucination problem.
- **Hallucination Problem:** It is that the model generate text contextually relevant but factually inaccurate based on the given prompts.

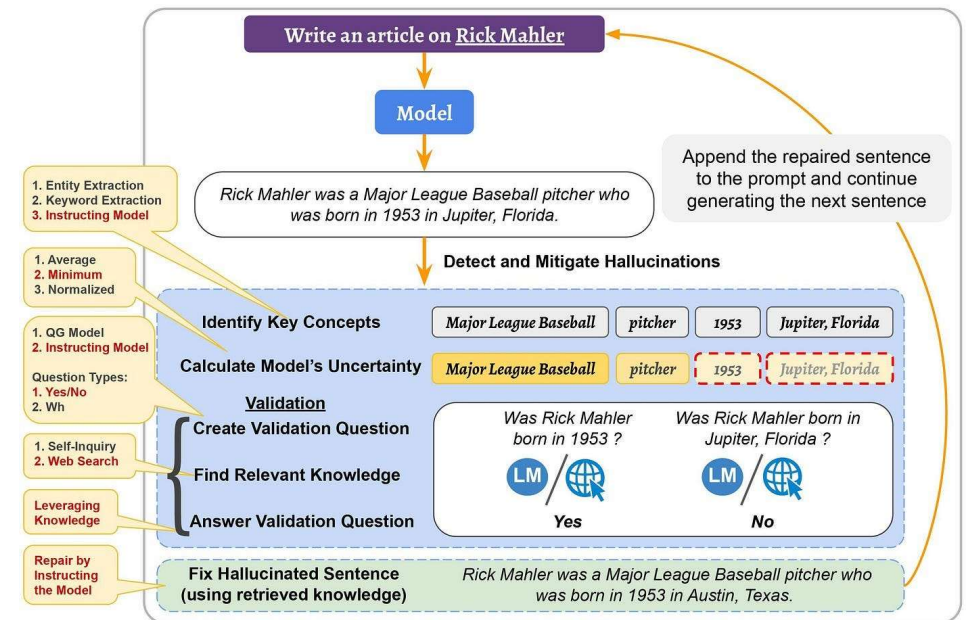


Figure 29: Hallucination Problem

<https://medium.com/mllearning-ai/the-hallucination-problem-of-large-language-models-5d7ab1b0f37f>

5. Instruction Tuning on Llama-2 Model

- In order to overcome hallucination problem, there is a technique called Retrieval Augmented Generation.
- **Retrieval Augmented Generation:** It fetches information from external sources such as documents and ensures that the generated output is factually relevant and accurate.

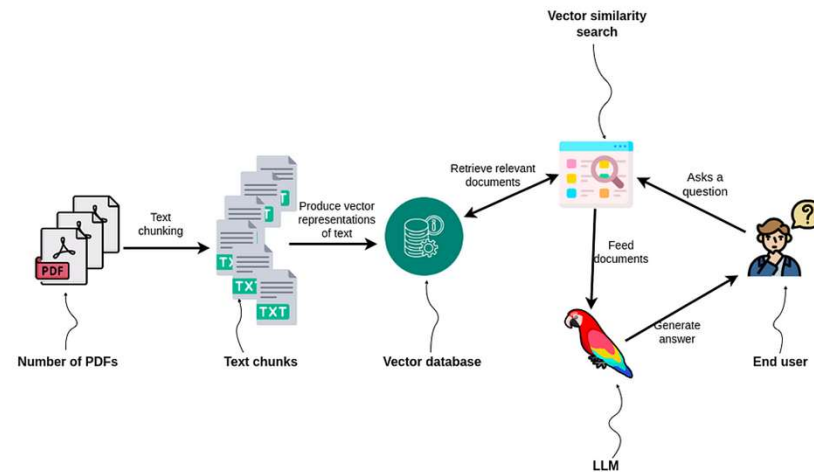


Figure 30: RAG-based Q&A task by using Llama-2 model

5. Instruction Tuning on Llama-2 Model

- In summary, if RAG-based model is applied on instruction-tuned Llama-2 model, the process of this study should be expected as shown in Figure 31.

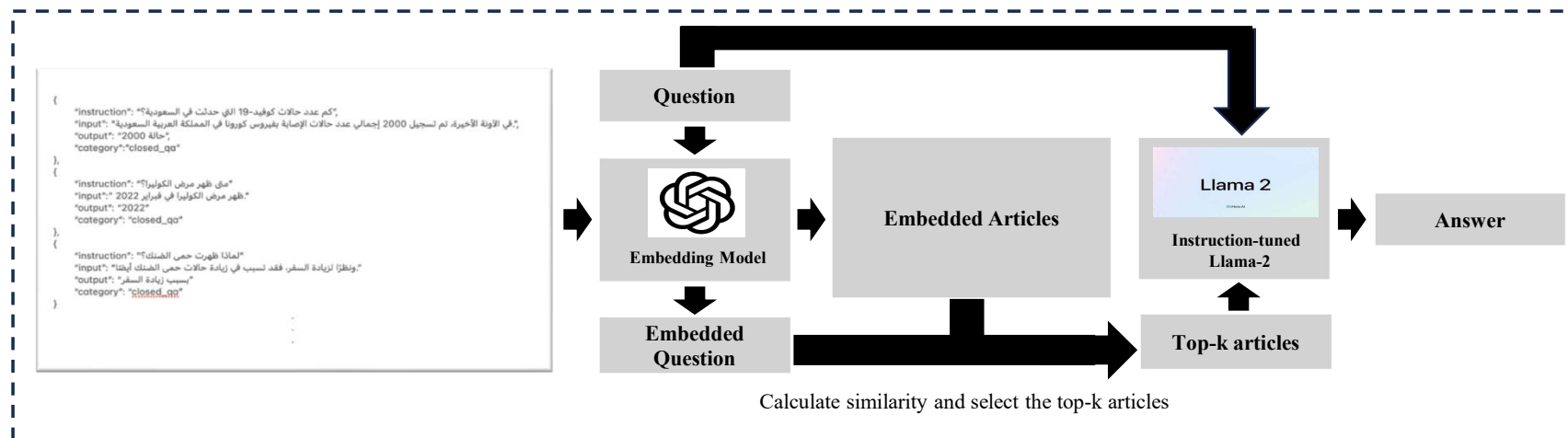


Figure 31: The RAG based Q&A Task by using Llama-2 model

<https://neo4j.com/developer-blog/knowledge-graphs-llms-multi-hop-question-answering/>

Q&A

감사합니다!