

# M-QALM: A Benchmark to Assess Clinical Knowledge Recall in Language Models via Question Answering

Anonymous ACL submission

## Abstract

In recent years, Large Language Models (LLMs) have gained recognition for their ability to encode knowledge within their parameters. Despite their growing popularity, the existing literature lacks a comprehensive and standardized benchmark for evaluating the performance of these models in clinical and biomedical knowledge applications. In response to this gap, we introduce a novel benchmark called M-QALM designed to unify the evaluation of language models in such contexts. Our benchmark has 16 Multiple-Choice Question (MCQ) datasets and 6 Abstractive Question Answering (AQA) datasets, offering a diverse range of challenges to comprehensively assess model capabilities. Our experimental results reveal intriguing insights. We find that encoder-decoder and decoder-only language models have differing strengths and weaknesses across question categories in biomedical and clinical knowledge MCQA. Additionally, our investigation demonstrates that instruction fine-tuned language models perform strongly compared to their base counterparts in these evaluations, emphasizing the importance of carefully tailored model selection. To foster research and collaboration in this field, we make our benchmark publicly available and open-source the associated evaluation scripts. This initiative aims to facilitate further advancements in clinical knowledge representation and utilization within language models, ultimately benefiting the healthcare and natural language processing communities.

## 1 Introduction

The recent success in the application of proprietary large language models in the medical domain (Singhal et al., 2023a,b) has sparked vivid research interest in applying smaller, more readily available open-source LLMs to various settings in the clinical and biomedical domains. Examples of tasks include summarization of clinical text (Veen et al.,

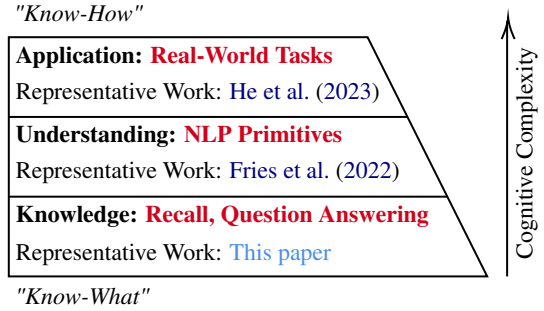


Figure 1: The landscape of LLM evaluation in the medical domain with **representative evaluation tasks**, organised by Bloom’s taxonomy of learning objectives (bold) (Bloom, 1956).

2023), automatic note generation for physicians (Ben Abacha et al., 2023b) and condensation of doctor-patient dialogues (Ben Abacha et al., 2023a; Toma et al., 2023). More broadly, open-source LLMs have been adapted to the domain to serve as foundational clinical models (Han et al., 2023; Wu et al., 2023; Toma et al., 2023; Bolton et al., 2022; Li et al., 2023).

The success of such adoption is typically established by measuring the performance on down-stream tasks, by means of token-overlap or semantic-similarity based metrics (Lin, 2004; Zhang et al., 2020). To address their inherent weaknesses (Schlegel et al., 2022; Gatt and Krahmer, 2018), research is carried out vividly to incorporate specific dimensions, such as factuality or faithfulness (Umapathi et al., 2023). Two important problems pertain, however. Firstly, NLG evaluation metrics are merely approximations of the phenomena they are aimed to measure, and their effectiveness is typically established by the degree of correlation to human judgements of the evaluated criteria (Huang et al., 2021). Secondly, an (offline) evaluation setup is *functionally grounded* and serves as of a real-world application scenario, and the transferability of insights from functionally-grounded to application-grounded evaluation is barely dis-

Dataset	Type	Size	Domain
USMLE (Jin et al., 2021)	MCQA	10178/1272/1273	General Medical
MEDMCQA (Pal et al., 2022)	MCQA	182822/4183/6150	General Medical
BIOASQ-MCQ (Tsatsaronis et al., 2015; Krithara et al., 2023)	MCQA	975/173/123	General Biomedical
HEADQA (Vilares and Gómez-Rodríguez, 2019)	MCQA	2657/1366/2742	General Medical
PROCESSBANK (Berant et al., 2014)	Context + MCQA	358/77/150	Biological Processes
PUBMEDQA (Jin et al., 2019)	Context + MCQA	400/100/500	General Biomedical
MMLU (Hendrycks et al., 2021)	MCQA	30/NA/1089	General Medical/Clinical
BIO-MRC-Tiny A (Pappas et al., 2020)	Context + MCQA	NA/NA/30	General Biomedical
BIO-MRC-Tiny B (Pappas et al., 2020)	Context + MCQA	NA/NA/30	General Biomedical
OPHTH (Raimondi et al., 2023; RCOphth, 2022a,b)	MCQA	NA/NA/92	Ophthalmology
QA4MRE-(Alzheimer’s QA) (Morante et al., 2012)	MCQA	NA/NA/40	Alzheimer’s Disease
LIVEQA (Abacha et al., 2017; Ben Abacha and Demner-Fushman, 2019)	AQA	NA/NA/131	Consumer Health
MEDIQA-ANS (Savery et al., 2020)	AQA	NA/NA/156	Consumer Health
BIOASQ-QA (Tsatsaronis et al., 2015; Krithara et al., 2023)	AQA	4733/697/363	General Biomedical
MASHQA (Zhu et al., 2020)	AQA	27728/3587/3493	General Medical
MEDQUAD (Ben Abacha and Demner-Fushman, 2019)	AQA	14068/981/1358	General Medical
MEDINFO (Ben Abacha et al., 2019)	AQA	NA/NA/663	Consumer Medication

Table 1: Overview of the M-QALM datasets. We present the size in terms of train/val/test splits. We create a manual train/val split for BIOASQ-MCQ, PROCESSBANK, PUBMEDQA, BIOASQ-QA and MEDQUAD.

cussed (Doshi-Velez and Kim, 2017). Taken together, these problems might taint the credibility of conclusions about the successful adoption of LLMs drawn from such experiments.

Given such difficulties, we approach the problem of evaluating LLM adoption from a complementary angle. Specifically, we ask: *Do LLMs possess the necessary pre-requisites to succeed in the clinical and medical domains?* Absent an established theory of how knowledge is acquired and organised in LLMs, the present work is guided by the established theories of knowledge acquisition in humans (Adams, 2015). Typical NLG tasks, such as summarisation, are higher-level cognitives that require the understanding of learned knowledge and its application in new contexts (Bloom, 1956). They build on the most fundamental capability: the *recall* of learned knowledge. In NLP research, knowledge recall is evaluated by the task of open-book *Question Answering* (QA), the task of retrieving—or selecting among presented options—the correct answer for a question. QA evaluation does not suffer from the issues pertaining to NLG metrics, as performance established by exact match. Thus, conclusions obtained from QA evaluation tend to be more robust, if the quality of the QA benchmark is sufficient.

Therefore, in this paper we focus on the task of

QA, to evaluate the knowledge pre-requisites of LLMs for successful adoption to the medical domain. We present an exhaustive, publicly available QA benchmark called M-QALM including 16 MCQA datasets. To enable future research on NLG-based QA, we complement M-QALM by 6 high-quality AQA datasets, where the ground-truth answer is an unconstrained string. With such standardized benchmark, we conduct an extensive evaluation of the capabilities of openly available general-purpose and medical LLMs, both “out-of-the-box” and after fine-tuning on M-QALM. Our findings provide insights into the strengths and weaknesses of different LLMs across different datasets, question categories and QA tasks. Overall, we find their performance lacking, both compared to humans and to proprietary LLMs. Further analysis reveals promising tendencies of domain-specific pre-training and fine-tuning to bridge this gap and to generalise to new QA datasets.

## 2 Related Work

**Large Open-domain QA benchmarks** The availability of QA datasets from multiple domains and sources has enabled the curation of large and diverse QA benchmarks (Dua et al., 2019; Fisch et al., 2019; Talmor and Berant, 2019). Such resource collections enable researchers to perform large-scale

empirical studies to understand, how well language models can generalise to new questions from new domains, or sources or how fine-tuning can impact this performance. While multiple studies exist in the general domain, to the best of our knowledge, no such large-scale study has been carried out for QA in the clinical domain. In this paper we aim to address this gap.

**Evaluation in the clinical domain** Datasets that evaluate the lowest-level cognitive task of knowledge recall have been previously proposed in the medical domain (Jin et al., 2021; Vilares and Gómez-Rodríguez, 2019; Pal et al., 2022). They feature questions commonly found in medical licensing examinations, including the US Medical Licensing Exam (USMLE). M-QALM unifies the existing literature by incorporating licensing exam questions from diverse regions, such as India and Spain. We go beyond the scope of the general medical domain, covering specialist topics such as Ophthalmology and Alzheimer’s disease.

Beyond pure factual recall, Fries et al. (2022) collect a unified bio-medical benchmark, featuring NLP primitives such as sentence(-pair) classification or entity recognition and linking. Aiming at higher, more task-specific cognitives, Singhal et al. (2023a) introduce MultiMedQA, including HealthSearchQA, which requires models to generate high-quality free-form answers. Similarly, (He et al., 2023) introduce a multi-domain benchmark for evaluating generation and classification capabilities on a diverse set of in-hospital downstream tasks. Other researchers looked to evaluate the quality and factuality of generations (Umapathi et al., 2023) and synthesised general-purpose medical instructions (Fleming et al., 2023). Our work is complementary, because we evaluate knowledge recall as a pre-requisite of higher-level cognitive tasks, such as, knowledge comprehension and application—the focus of previously discussed works.

### 3 M-QALM Datasets

The primary goal of M-QALM is to develop a comprehensive, open-source repository of medical QA datasets to assess the recall of medical knowledge in LLMs. To obtain such a collection, we perform an exhaustive literature and resource search using the terms “clinical OR medical”, “Question Answering or QA” and include a dataset or resource, if it satisfies the following criteria: (i) The language is English, as medical documents are usually

written in English, even in non English-speaking countries; (ii) The questions and answers are on general, specialist or consumer-facing medicine; (iii) The resource is openly available without restrictive licensing or data agreements; (iv) The resource evaluates the task of MCQA or AQA.

The result is M-QALM—a comprehensive collection of 22 datasets designed to thoroughly evaluate the clinical knowledge of LLMs. Table 1 gives an overview of the collected MCQA and AQA datasets, including task formulation, size and domain. For further details on each of the datasets, we refer to the Appendix.

**Question Categorization** The MCQA datasets within the M-QALM benchmark cover a diverse range of medical domains. To be able to perform fine-grained analysis of both the topics covered in these datasets as well as the models performance, we categorise the MCQA datasets into eleven high-level categories, representing different facets of medical knowledge.

To do so, we leverage available meta-data from the source datasets, MEDMCQA, HEADQA, MMLU and BIOASQ-MCQ. We categorize the PROCESSBANK, PUBMEDQA and BIOMRC datasets into a distinct twelfth Reading Comprehension category. For USMLE and QA4MRE, to account for the lack of meta-data, we train a BioBERT-based classifier (Lee et al., 2019) to assign questions into one of the eleven elicited categories using the labels from the other datasets. The classifier achieves 71.56% (micro-)averaged F1 score on a held-out test set, which we deem sufficient.

Category	Share
Basic and Life Sciences	23.03%
Dental and Oral Health	5.29%
Diagnostic Sciences	10.26%
General Medicine	23.88%
Mental and Behavioral Health	2.68%
Musculoskeletal and Dermatology	2.25%
Pharmacology and Anesthesia	8.87%
Sensory Organs	5.68%
Supportive and Preventive Services	6.20%
Women’s and Children’s Health	9.12%
Reading Comprehension	0.76%
Miscellaneous	1.97%

Table 2: Topic distribution of M-QALM.

Table 2 shows the question distribution in M-QALM. We observe that nearly half of all questions (47%)

fall into the Basic and Life Sciences and General Medicine category. Miscellaneous and Reading Comprehension account for the least percentage of questions (3%), with other questions more evenly distributed amongst categories. Diagnostic Sciences, Women’s and Children’s Health and Pharmacology and Anesthesia account for nearly 30% of questions.

## 4 Empirical Evaluation

Considering the M-QALM datasets, we investigate, how well existing, open-source LLMs are able to recall clinical knowledge in order to succeed on the benchmark. Specifically, we focus on performance in zero-shot setting, and after fine-tuning on M-QALM training portions.

In the **Zero-shot** setting:

- **RQ1.** How well do open-source LLMs recall necessary clinical knowledge when they are tested on M-QALM?
- **RQ2.** Does open-domain instruction fine-tuning of LLMs improve their ability to do so?
- **RQ3.** Does *domain-specific* fine-tuning improve performance on M-QALM?

In the **Fine-tuned** setting:

- **RQ4.** Does finetuning on M-QALM improve performance on unseen data from datasets seen during training?
- **RQ5.** Does fine-tuning improve performance on *unseen* M-QALM datasets?

### 4.1 Study Setup

To seek evidence for **RQs1-3** empirically, we evaluate several LLMs and their instruction-tuned versions on the test splits of M-QALM in zero-shot manner. To answer **RQ4** and **RQ5**, we fine-tune LLMs on the training portion of M-QALM and evaluate on test splits of datasets both seen and unseen during training. We complement our evaluation with additional automated and manual error analyses to identify causes for model successes and failures.

**Models:** To assess the zero-shot capabilities of models (**RQ1** and **RQ2**), we include a diverse array of open-source decoder-only models with parameter scales ranging from 3B-13B. We use models from MPT and MPT-Instruct (7B) (MosaicML, 2023), Falcon and Falcon-Instruct (7B) (Almazrouei et al., 2023) and LLaMA 2 and LLaMA 2-chat (7B and 13B). In addition to these models, we also use two instruction fine-tuned encoder-decoder models:

Flan-T5 (3B and 11B) (Wei et al., 2021). Models with *Instruct* or *Chat* appended to their names are instruction fine-tuned (Ouyang et al., 2022) versions of their base models. The details of the models are given in Table 8. To address **RQ3**, we evaluate ChatDoctor (7B), MedAlpaca (7B)<sup>1</sup> To address **RQ4**, we use the training set of the M-QALM datasets. When official validation splits are unavailable, we employ a random split of up to around 20% of the data for validation purposes. If no training datasets are available, we do not use this dataset for fine-tuning and only consider the test split of the respective datasets to answer **RQ5**. For evaluating AQA, we use a sub-sampled version of the test sets of MASHQA and MEDQUAD, while we use the other datasets as is.

**Finetuning and hyperparameters:** Since the number of parameters for most of our models are in the billions, we follow a more accepted practice of using parameter-efficient fine-tuning. Specifically, we use QLoRA and 4-bit quantization (Dettmers et al., 2023) for fine-tuning. We utilize 8-bit quantization for evaluating Flan-T5 (11B), LLaMA 2 (13B) and LLaMA 2-Chat (13B) (Dettmers et al., 2022). We use A100-40G GPUs for all our experiments. The other hyper-parameters used to train our models are reported in the Appendix (Table 9).

**Evaluation measures:** We use Accuracy to measure the performance of the model on MCQA datasets; for AQA datasets, we use ROUGE-L (Lin, 2004), BERTScore (Zhang et al., 2020)<sup>2</sup> and METEOR (Banerjee and Lavie, 2005), which is found to correlate better with human judgments than other metrics on AQA (Chen et al., 2019).

### 4.2 Results and Analysis

In this section, we report and analyse the findings of our empirical study.

#### 4.2.1 Zero-shot Evaluation Results

Table 3 report the dataset-averaged scores of the zero-shot evaluation of language models as evidence towards **RQs1-3**. For a by-dataset breakdown, consult the Appendix.

<sup>1</sup>For MCQA evaluation in the zero-shot setting (where models are not explicitly fine-tuned for MCQA tasks), we use a 1-shot prompt—giving an example to the model, and find that it adheres better to the MCQA format and the standard 5-shot prompt for MMLU datasets.

<sup>2</sup>We use deberta-xlarge-mnli for calculating BERTScore.



		MCQA		AQA	
		Acc	RL	BS	MTR
Base	LLaMA 2 (7B)	42.9	14.9	55.3	21.1
	LLaMA 2 (13B)	47.1	15.0	56.4	22.5
	MPT (7B)	27.6	13.3	52.6	21.1
	Falcon (7B)	34.7	14.0	54.1	20.0
Instruction tuned	LLaMA 2-chat (7B)	45.9	15.0	58.0	23.3
	LLaMA 2-chat (13B)	50.3	15.3	58.0	23.6
	MPT-Instruct (7B)	31.6	15.8	59.7	15.6
	Falcon-Instruct (7B)	31.8	17.2	62.4	17.4
	Flan-T5 (3B)	51.8	10.8	55.0	7.4
	Flan-T5 (11B)	56.5	11.5	56.3	8.2
Adapted	ChatDoctor (7B)	42.8	17.4	62.3	18.7
	MedAlpaca (7B)	48.8	15.5	58.9	15.6
	PMC-LLama (13B)	53.7	19.7	60.7	19.0

Table 3: Zero-shot performance of base (top), instruction-tuned models (middle) and domain-adapted (bottom) models. Metrics are **Accuracy** for MCQA; **Rouge-L**, **BERTScore**, and **METEOR** for AQA.

Table 3 shows that LLMs exhibit **strong zero-shot capability on MCQA and AQA datasets**, corroborating the findings of Singhal et al. (2023a). Considering LLMs of the same size (i.e., 7B), LLaMA 2 performs best, possibly due to larger diversity in pre-training data—LLaMA 2 is trained on the most tokens. Another difference is the mixture of datasets used for pre-training, which is not revealed in some cases (c.f Table 8 in Appendix).

Unsurprisingly, across all models of same architecture, **scale predicts model performance**, even without domain-specific adaptation of LLMs on the medical domain. For example, LLaMA 2 (13B) performs better on MCQA (+4.2 Accuracy improvement) compared to the 7B version. Figure in the appendix 4 shows the relationship between the number of parameters and performance.

To address **RQ2**, we investigate whether improvements from instruction fine-tuning also apply to the clinical domain of M-QALM. The results are reported in the bottom part of Table 3.

Surprisingly, **instruction fine-tuned models perform better** than their corresponding *Base* versions, despite the fact that the instruction set used for fine-tuning contains only tasks in the general domain, see Table 8 (middle) and compare \*Instruct/Chat with their base versions (top). Among them, Flan-T5 models show the best zero shot performance on MCQA, outperforming all comparable decoder-only models. Seemingly, instruction fine-tuning enables models to obtain representations of question and context, which are beneficial for fact

recall.

We note that **bigger models are not always better**—the choice of model architecture and the dataset for instruction fine-tuning can have a bigger impact on performance than model size alone. For example the encoder-decoder Flan-T5 (3B) model outperforms LLaMA 2-chat (13B) on average on the MCQA task, despite being four times smaller in size.

The performance of domain-adapted models is reported in Table 3 (bottom), as evidence for **RQ3**. For MCQA, while MedAlpaca and ChatDoctor exhibit improvements in Accuracy over their respective 7B and 13B LLaMA 2 base and chat versions, they yet fail to reach the strong zero-shot performance of Flan-T5 (11B). The poor performance of ChatDoctor could be attributed to the fact that it is only fine-tuned using synthetic question-answering data rather than other instruction datasets.

PMC-LLama is an outlier here, as it performs well due to continued pre-training on biomedical corpora before instruction tuning on biomedical and clinical datasets. The latter, results in exceptionally high scores on the MEDINFO AQA dataset (See Table 17 in Appendix). This dataset, along with LIVEQA was used as part of the instruction tuning process, leading to evaluation on these dataset not being “zero-shot”<sup>3</sup>. Scores on LIVEQA, however, are not inflated, compared to LLaMA 2(-chat) (13B). This is possibly because we use a filtered version of LIVEQA which contains only challenging answers that with sufficiently good expert quality rating. PMC-LLama demonstrates significant improvements over other open-source LLMs on MCQA datasets such as USMLE, MEDMCQA and MMLU.

Summarily, we conclude that most openly available LLMs adapted to the medical domain **have no improved domain knowledge compared to available open-domain models**.

Importantly, we note none of the evaluated open-source LLMs outperform humans: While the passing score for USMLE is 60% for humans<sup>4</sup>, we observe the best zero-shot scores for USMLE are 43% for LLaMA 2, and 54% for the domain-adapted PMC-LLama, both below the passing score. Meanwhile, GPT-4 (OpenAI, 2023) with a customized prompting strategy labeled MedPrompt (Nori et al.,

<sup>3</sup>[https://huggingface.co/datasets/axiong/pmc\\_llama\\_instructions](https://huggingface.co/datasets/axiong/pmc_llama_instructions)

<sup>4</sup><https://www.usmle.org/bulletin-information/scoring-and-score-reporting>

	MCQA		AQA	
	Acc	RL	BS	MTR
LLaMA 2 (7B)	53.5 <sup>+10.6</sup>	17.7 <sup>+2.8</sup>	60.8 <sup>+5.5</sup>	16.9 <sup>-4.2</sup>
Falcon (7B)	49.3 <sup>+14.6</sup>	17.4 <sup>+3.4</sup>	60.4 <sup>+6.3</sup>	17.1 <sup>-2.9</sup>
MPT (7B)	53.2 <sup>+25.6</sup>	17.3 <sup>+4.0</sup>	60.0 <sup>+7.4</sup>	17.2 <sup>-3.9</sup>
Flan-T5 (3B)	52.9 <sup>+1.1</sup>	15.9 <sup>+5.1</sup>	56.8 <sup>+1.8</sup>	15.6 <sup>+8.2</sup>

Table 4: Model finetuning is performed either on MCQA or AQA datasets. Evaluation is performed using **Accuracy** for MCQA, and **Rouge-L**, **BERTScore**, and **METEOR** for AQA. The subscripts indicate the improvement over the zero-shot versions.

2023) achieves 90.2% while Med-PALM 2 (Singhal et al., 2023b) achieves scores of 86.5% on USMLE. Similarly, for the PubmedQA dataset, human performance is 78% (Jin et al., 2019), compared to 72.4% of Flan-T5. To summarize: While available LLMs exhibit performance significantly higher than random chance “out-of-the-box”, **there is still a significant gap compared to humans and proprietary LLMs (Singhal et al., 2023a,b)** (provided in Table 7 in the Appendix).

#### 4.2.2 Impact of Fine-tuning

Given the scale of M-QALM, we are able to fine-tune models on parts of the data, to address **RQ4** and **RQ5**. We fine-tune four models on MCQA and AQA separately, given the different nature of these datasets.<sup>5</sup> We fine-tune the models only on the MCQA subset of datasets first (c.f. Table 4). We find that the **fine-tuned models perform better compared to their non-fine-tuned counterparts**. Decoder-only models like MPT (7B) benefit more than others (+25.6 Accuracy improvement). Interestingly, fine-tuning models on the data seems to close the gaps introduced by different model architectures and pre-training data, discussed in the previous Section. Specifically, the standard deviation of the model accuracy in the zero-shot setting is 9.0, while after fine-tuning, it is reduced to 1.7. This suggests that various LLM can benefit from task-specific fine-tuning to address seemingly sub-optimal architecture or pre-training conditions. For AQA, Flan-T5 seems to benefit more from AQA fine-tuning compared to the decoder-only models, possibly by better aligning to the expected output format of the question. Decoder models present inconsistent results with improvements in ROUGE-L and BERTScore at the expense of lower METEOR

<sup>5</sup>We also experimented with fine-tuning models on MCQA and AQA jointly, but the results did not differ significantly from those reported here.

scores.

Scaling up models introduces practical problems of deploying the model in real-world scenarios—smaller models may be preferred to larger ones due to faster inference times and lower memory footprints. We find that **fine-tuning helps compensate for scale**. Fine-tuned LLaMA 2 (7B) significantly outperforms the zero-shot LLaMA 2 (13B) (+6.4 Accuracy gain on MCQA, +2.7 ROUGE-L gain and +4.4 BERTScore gain on AQA). Similarly, we observe that a fine-tuned Flan-T5 (3B) outperforms zero-shot LLaMA 2 (13B) on 8 out of 16 MCQA datasets (see Table 11 and 13).

In summary, we conclude that **task-specific fine-tuning markedly improves the performance of language models, mitigating architectural and pre-training disparities and allows smaller models to rival larger ones in performance**.

Finally, we report the potential of LLMs fine-tuned on in-domain data to generalize to medical datasets unseen during training to answer **RQ5**. To this end, during fine-tuning, we hold out 10 MCQA and 4 AQA datasets presented in Figures 2 and 3.

**AQA-finetuned models generalise to unseen AQA test sets:** Figures 2, 5, 6 and show the performance of LLaMA 2 (7B) and Flan-T5 (3B) models on the four held-out AQA evaluation sets on various metrics. LLaMA 2 (7B) fine-tuned shows improvements over its zero-shot version in terms of ROUGE-L score on the LIVEQA and MEDQUAD with no significant performance dip on MEDINFO. In terms of BERTScore, the LLaMA 2 (7B) outperforms the zero-shot version across all four datasets. However, the METEOR scores of the fine-tuned LLaMA 2 model are lower than the zero-shot baseline across all four datasets. Meanwhile, fine-tuning Flan-T5 improves performance on all four unseen datasets on ROUGE-L and METEOR scores. However, on BERTScore, the zero-shot model outperforms the fine-tuned version on two datasets.

**AQA-finetuned models do not generalise to unseen MCQA test sets:** Figure 3 (comparing ZS with AQA-FT) shows that fine-tuning on AQA does not improve performance on unseen MCQA datasets. This suggests that higher scores on unseen AQA datasets might stem from better aligning generations to the expected answer form of AQA answers, rather than acquiring additional medical knowledge during fine-tuning.

Category	Support	Flan-T5 (ZS)	Flan-T5 (FT)	MPT (ZS)	MPT (FT)	Falcon (ZS)	Falcon (FT)	LLaMA 2 (ZS)	LLaMA 2 (FT)
General Medicine	2675	38.0	43.2	26.0	46.4	30.1	46.4	36.6	<b>50.0</b>
Basic and Life Sciences	2235	38.9	44.3	26.9	<b>52.6</b>	30.6	49.4	40.0	52.5
Dental and Oral Health	1318	34.8	42.9	25.9	<b>44.3</b>	30.7	43.8	36.1	44.2
Pharmacology and Anesthesia	784	39.7	48.1	29.0	55.6	28.8	54.0	42.9	<b>59.4</b>
Reading Comprehension	710	74.1	<b>75.2</b>	37.2	71.5	52.7	66.5	60.8	67.7
Diagnostic Sciences	640	32.2	43.1	26.4	<b>51.1</b>	30.3	46.4	37.2	47.5
Supportive and Preventive Services	599	48.2	<b>56.6</b>	23.7	55.1	27.9	48.1	39.9	56.3
Women's and Children's Health	507	30.2	42.6	27.2	<b>51.7</b>	28.4	43.0	34.3	49.9
Mental and Behavioral Health	496	50.0	57.9	29.4	55.4	31.5	49.2	40.7	<b>59.1</b>
Sensory Organs	205	29.8	42.0	27.8	<b>45.4</b>	28.8	42.4	33.2	42.0
Miscellaneous	45	42.2	44.4	20.0	<b>60.0</b>	24.4	44.4	31.1	40.0
Musculoskeletal and Dermatology	38	18.4	26.3	18.4	<b>44.7</b>	34.2	42.1	28.9	<b>44.7</b>
Overall Accuracy	10252	40.6	47.4	27.3	51.5	31.6	48.6	39.6	<b>52.2</b>

Table 5: Performance of LLMs in the zero-shot and fine-tuned setting across various categories on the test set.

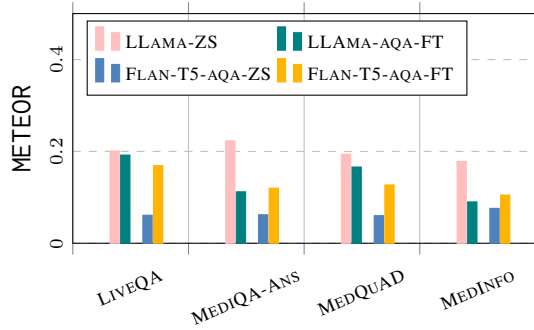


Figure 2: Performance of base and AQA-finetuned LLaMA 2 and Flan-T5 models on four unseen AQA test sets.

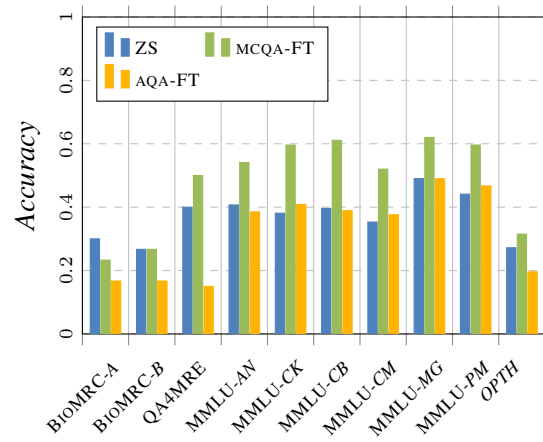


Figure 3: Performance of base, MCQA-fine-tuned and AQA-finetuned LLaMA 2 model on ten unseen MCQA test sets.

**MCQA-finetuned models generalise to unseen MCQA test sets:** Figure 3 (comparing ZS with MCQ-FT) suggests that models indeed can learn to extract relevant knowledge during fine-tuning, as MCQA-tuned models consistently perform better than their zero-shot counterparts. This seemingly contradicts the previous finding that models fail to acquire additional medical knowledge when fine-tuned on the AQA datasets. To investigate this mismatch, we conduct a manual analysis.

**Fine-tuned models may memorize rather than generalize:** We aim to discriminate whether MCQA fine-tuned models' performance on unseen MCQA datasets can be attributed to their ability to generalize in answering medical questions, or if their performance is influenced by memorization of questions from the training set. To this end, we examine three evaluation-only MCQ datasets not used in the training split of M-QALM: Clinical Knowledge Tests (MMLU-CK) and Medical Genetics (MMLU-MG) from MMLU and the OPTH dataset. We utilize semantic similarity algorithms to retrieve questions in the training sets that closely resemble those in these test sets and manually fil-

ter the retrieved results. We identify 6 out of 92, 12 out of 265, and 17 out of 100 questions in the OPTH, MMLU-CK, and MMLU-MG datasets, respectively, that have similar counterparts in the MEDMCQA dataset which was used to fine-tune the LLaMA 2 model. This suggests that scores might be inflated due to train-test leakage.

Next, we focus on questions that the LLaMA 2 (7B) model answered wrongly, but which were corrected by MCQA-fine-tuning. We then cross-reference these with the closest equivalent questions in the MEDMCQA dataset. This allows us to categorize the correct answers from near-duplicate memorization or the model's generalized learning capabilities. We find 5, 2, and 5 questions in the three investigated datasets, respectively, where the MCQA-fine-tuned model outperformed its zero-shot counterpart and identified closely related questions in MEDMCQA. Of these, 7 questions were near-duplicates with identical answers, while the remaining 5 would have required some level of clinical understanding for the model to answer them correctly.

This suggests that the improved performance of instruction-tuned models on unseen datasets can be partially attributed to exposure to near-identical questions during training.

Based on these findings, we observe that **fine-tuning only serves as a partial solution for achieving broad generalization across domains.**

### 4.3 Category-wise and manual error analysis

To better understand the performance of zero-shot and fine-tuned performance of models across MCQA, we analyze the performance of the models across various categories. We calculated the accuracy of the models in their zero-shot and fine-tuned settings for each category, as shown in Table 5. Fine-tuning markedly improves model performance across different question categories. Notably, fine-tuned Flan-T5 (3B) excels in Reading Comprehension and Supportive and Preventive Services, also showing strong zero-shot capabilities in these areas. This suggests that encoder-decoder models like Flan-T5 may outperform decoder-only models in such categories. Similarly, fine-tuned MPT (7B) and LLaMA 2 (7B) show superior performance in various categories. However, despite fine-tuning benefits, models still underperform in areas like General Medicine, Basic and Life Sciences, and Dental and Oral Health.

In our manual error analysis of a fine-tuned LLaMA 2 (7B) model on MCQA, we examined 200 non-Reading Comprehension questions where the model erred, categorizing them into Factual, Conceptual Understanding, and Quantitative/Arithmetic. Factual questions involve direct medical knowledge recall, Conceptual Understanding questions assess the application of medical and clinical concepts, and Quantitative/Arithmetic questions require mathematical skills for correct answers. The model incorrectly answered 134 Factual, 52 Conceptual Understanding, and 14 Quantitative/Arithmetic questions (Table 6). Comparing these errors to a random sample of 200 questions from the test set revealed sim-

ilar error rates across categories, reflecting the general frequency of question types in the test set. The prevalence of Factual questions in errors aligns with their dominance in medical exams like MEDMCQA, USMLE, and HEADQA. While fine-tuning on extensive medical corpora may enhance Factual question performance, improving on Conceptual Understanding and Quantitative/Arithmetic questions might require different fine-tuning approaches, as these categories demand more than mere knowledge recall.

## 5 Conclusion

In this work, we introduce M-QALM, a comprehensive collection of clinical datasets comprising 16 multiple-choice and six abstractive question-answering datasets. Our study encompasses an extensive empirical investigation of open-source language models, some of which are trained with up to 13 billion parameters. We assess their clinical and biomedical knowledge, their capacity to acquire such knowledge through training on M-QALM, and their ability to generalize to previously unseen datasets. Our results highlight the strengths and limitations of LLMs on MCQA and AQA, showing that while they exhibit certain proficiencies, they still fall significantly short in performance compared to proprietary language models, indicating potential areas for improvement. Notably, fine-tuning on M-QALM demonstrates the potential to augment a language model’s clinical knowledge, especially in the context of instruction fine-tuned models like Flan-T5. However, we recognize that fine-tuning is not a universal solution for generalization, evidenced by its limitations in effectively extending knowledge from AQA to MCQA contexts. It is important to note that scale and decoder-only language models do not serve as universal solutions for all questions in clinical question-answering. To pave the way for future research in this domain, we emphasize the necessity of considering the architecture of language models, the choice of datasets for instruction fine-tuning, and conducting a rigorous evaluation of the knowledge contained within LLMs. We make the dataset, experiment code and evaluation protocol publicly available under <https://anonymized>. This will allow practioners to perform fine-grained analysis of their models’ clinical and biomedical knowledge.

Category	General Test Set Distribution	LLaMA 2 Errors
Factual	65.5%	67%
Conceptual Understanding	29.5%	26%
Quantitative/Arithmetic	5%	7%

Table 6: Overview of question categories and their distributions.



## Limitations

In this paper, we evaluate the medical or clinical knowledge of LLMs by measuring their capability of answering test questions. While this can be a useful proxy-measure of a model’s domain knowledge, it is insufficient to gauge its potential application in a real-world scenario. A multi-dimensional analysis of a model’s behaviour, including judging the completeness, harmlessness and usefulness of generated answers, is required in addition to solely evaluating their correctness.

Furthermore, the aggregated resource presented in this paper might be seen as lacking diversity, as all collected datasets are in English. To make inferences about the capabilities of evaluated models in other languages, a more diverse dataset with examples in other languages is required.

For our finetuning experiments, we only use parameter-efficient finetuning methods (PEFT) with QLoRA due to the high compute requirements for full-finetuning. We have not investigated the impact of the full-finetuning of these LLMs on our benchmark.

## References

Asma Ben Abacha, Eugene Agichtein, Yuval Pinter, and Dina Demner-Fushman. 2017. Overview of the medical question answering task at TREC 2017 LiveQA. In *Text REtrieval Conference (TREC)*.

Nancy E Adams. 2015. Bloom’s taxonomy of cognitive learning objectives. *Journal of the Medical Library Association: JMLA*, 103(3):152.

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.

Giuseppe Attardi, Luca Atzori, Maria Simi, et al. 2012. Index expansion for machine reading and question answering.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah,

Ben Mann, and Jared Kaplan. 2022. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#).

Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proc. ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Asma Ben Abacha and Dina Demner-Fushman. 2019. [A question-entailment approach to question answering](#). *BMC Bioinformatics*, 20(1):511:1–511:23.

Asma Ben Abacha, Yassine Mrabet, Mark Sharp, Travis Goodwin, Sonya E. Shooshan, and Dina Demner-Fushman. 2019. Bridging the gap between consumers’ medication questions and trusted answers. In *Proc. 17th World Congress on Medical and Health Informatics (MEDINFO)*.

Asma Ben Abacha, Wen-wai Yim, Griffin Adams, Neal Snider, and Meliha Yetisgen. 2023a. [Overview of the MEDIQA-chat 2023 shared tasks on the summarization & generation of doctor-patient conversations](#). In *Proc. 5th Clinical Natural Language Processing Workshop*, pages 503–513, Toronto, Canada. Association for Computational Linguistics.

Asma Ben Abacha, Wen-wai Yim, Yadan Fan, and Thomas Lin. 2023b. [An empirical study of clinical note generation from doctor-patient encounters](#). In *Proc. 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2291–2302, Dubrovnik, Croatia. Association for Computational Linguistics.

Jonathan Berant, Vivek Srikumar, Pei-Chun Chen, Abby Vander Linden, Brittany Harding, Brad Huang, Peter Clark, and Christopher D. Manning. 2014. [Modeling biological processes for reading comprehension](#). In *Conference on Empirical Methods in Natural Language Processing*.

Benjamin S Bloom. 1956. *Taxonomy of education objectives Book 1-Cognitive domain*. David McKay Company.

Elliot Bolton, David Hall, Michihiro Yasunaga, Tony Lee, Chris Manning, and Percy Liang. 2022. [Biomedlm](#).

Anthony Chen, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [Evaluating question answering evaluation](#). In *Proc. 2nd Workshop on Machine Reading for Question Answering*, pages 119–124, Hong Kong, China. Association for Computational Linguistics.

Together Computer. 2023. [Redpajama-data: An open source recipe to reproduce llama training dataset](#).

709	Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie,	conversational ai models and training data. <i>arXiv</i> ,	766
710	Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell,	2304.08247.	767
711	Matei Zaharia, and Reynold Xin. 2023. <a href="#">Free dolly:</a>		
712	<a href="#">Introducing the world's first truly open instruction-</a>	Zexue He, Yu Wang, An Yan, Yao Liu, Eric Y Chang,	768
713	<a href="#">tuned llm.</a>	Amilcare Gentili, Julian McAuley, and Chun-Nan	769
714	Tim Dettmers, Mike Lewis, Younes Belkada, and Luke	Hsu. 2023. Medeval: A multi-level, multi-task,	770
715	Zettlemoyer. 2022. <a href="#">GPT3.int8(): 8-bit Matrix Multi-</a>	and multi-domain medical benchmark for language	771
716	<a href="#">plication for Transformers at Scale.</a> In <i>Advances in</i>	model evaluation. <i>arXiv preprint arXiv:2310.14088</i> .	772
717	<i>Neural Information Processing Systems</i> , volume 35,		
718	pages 30318–30332. Curran Associates, Inc.	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	773
719	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	774
720	Luke Zettlemoyer. 2023. QLoRA: Efficient Finetun-	2021. Measuring massive multitask language under-	775
721	ing of Quantized LLMs. <i>arXiv</i> , 2305.14314.	standing. <i>arXiv</i> , 2009.03300.	776
722	Finale Doshi-Velez and Been Kim. 2017. <a href="#">Towards A</a>	Yichong Huang, Xiachong Feng, Xiaocheng Feng, and	777
723	<a href="#">Rigorous Science of Interpretable Machine Learning.</a>	Bing Qin. 2021. The factual inconsistency problem	778
724	Dheeru Dua, Ananth Gottumukkala, Alon Talmor, Matt	in abstractive text summarization: A survey. <i>arXiv</i>	779
725	Gardner, and Sameer Singh. 2019. <a href="#">Comprehensive</a>	<i>preprint arXiv:2104.14839</i> .	780
726	<a href="#">Multi-Dataset Evaluation of Reading Comprehension.</a>	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,	781
727	In <i>Proceedings of the 2nd Workshop on Machine</i>	Hanyi Fang, and Peter Szolovits. 2021. <a href="#">What disease</a>	782
728	<i>Reading for Question Answering</i> , pages 147–153,	<a href="#">does this patient have? a large-scale open domain</a>	783
729	Stroudsburg, PA, USA. Association for Computa-	<a href="#">question answering dataset from medical exams.</a> <i>Ap-</i>	784
730	tional Linguistics.	<i>plied Sciences</i> , 11(14).	785
731	Adam Fisch, Alon Talmor, Robin Jia, Minjoon Seo,	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William	786
732	Eunsol Choi, and Danqi Chen. 2019. <a href="#">MRQA 2019</a>	Cohen, and Xinghua Lu. 2019. <a href="#">PubMedQA: A</a>	787
733	<a href="#">Shared Task: Evaluating Generalization in Reading</a>	<a href="#">dataset for biomedical research question answering.</a>	788
734	<a href="#">Comprehension.</a> In <i>Proceedings of the 2nd Work-</i>	In <i>Proc. Conference on Empirical Methods in Natu-</i>	789
735	<i>shop on Machine Reading for Question Answering</i> ,	<i>ral Language Processing and the 9th International</i>	790
736	pages 1–13, Stroudsburg, PA, USA. Association for	<i>Joint Conference on Natural Language Processing</i>	791
737	Computational Linguistics.	( <i>EMNLP-IJCNLP</i> ), pages 2567–2577, Hong Kong,	792
738	Scott L. Fleming, Alejandro Lozano, William J.	China. Association for Computational Linguistics.	793
739	Haberkorn, Jenelle A. Jindal, Eduardo P. Reis,	Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and	794
740	Rahul Thapa, Louis Blankemeier, Julian Z. Genk-	Dan Roth. 2018. Question answering as global rea-	795
741	ins, Ethan Steinberg, Ashwin Nayak, Birju S. Patel,	soning over semantic abstractions. In <i>Proceedings</i>	796
742	Chia-Chun Chiang, Alison Callahan, Zepeng Huo,	<i>of the AAAI Conference on Artificial Intelligence</i> ,	797
743	Sergios Gatidis, Scott J. Adams, Oluseyi Fayanju,	volume 32.	798
744	Shreya J. Shah, Thomas Savage, Ethan Goh, Ak-	Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia	799
745	shay S. Chaudhari, Nima Aghaeepour, Christopher	Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine	800
746	Sharp, Michael A. Pfeffer, Percy Liang, Jonathan H.	Jernite, Margaret Mitchell, Sean Hughes, Thomas	801
747	Chen, Keith E. Morse, Emma P. Brunskill, Jason A.	Wolf, Dzmitry Bahdanau, Leandro von Werra, and	802
748	Fries, and Nigam H. Shah. 2023. Medalign: A	Harm de Vries. 2022. The stack: 3 tb of permissively	803
749	clinician-generated dataset for instruction following	licensed source code. <i>Preprint</i> .	804
750	with electronic medical records. <i>arXiv</i> , 2308.14089.	Anastasia Krithara, Anastasios Nentidis, Konstantinos	805
751	Jason Fries, Leon Weber, Natasha Seelam, Gabriel Al-	Bougatiotis, and Georgios Paliouras. 2023. Bioasq-	806
752	tay, Debajyoti Datta, Samuele Garda, Sunny Kang,	qa: A manually curated corpus for biomedical ques-	807
753	Rosaline Su, Wojciech Kusa, Samuel Cahyawijaya,	tion answering. <i>Scientific Data</i> , 10(1):170.	808
754	et al. 2022. Bigbio: a framework for data-centric	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon	809
755	biomedical natural language processing. <i>Advances</i>	Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang.	810
756	<i>in Neural Information Processing Systems</i> , 35:25792–	2019. <a href="#">BioBERT: a pre-trained biomedical language</a>	811
757	25806.	<a href="#">representation model for biomedical text mining.</a>	812
758	Albert Gatt and Emiel Krahmer. 2018. <a href="#">Survey of the</a>	<i>Bioinformatics</i> , 36(4):1234–1240.	813
759	<a href="#">State of the Art in Natural Language Generation:</a>	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve	814
760	<a href="#">Core tasks, applications and evaluation.</a> <i>Journal of</i>	Jiang, and You Zhang. 2023. Chatdoctor: A medical	815
761	<i>Artificial Intelligence Research</i> , 61:65–170.	chat model fine-tuned on a large language model	816
762	Tianyu Han, Lisa C Adams, Jens-Michalis Papaioan-	meta-ai (llama) using medical domain knowledge.	817
763	nou, Paul Grundmann, Tom Oberhauser, Alexander	<i>Cureus</i> , 15(6).	818
764	Löser, Daniel Truhn, and Keno K Bressen. 2023.		
765	Medalpaca—an open-source collection of medical		

819	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	875
820		876
821		877
822		878
823	Ye Liu, Shaika Chowdhury, Chenwei Zhang, Cornelia Caragea, and Philip S Yu. 2020. Interpretable multi-step reasoning with knowledge extraction on complex healthcare question answering. <i>arXiv preprint arXiv:2008.02434</i> .	879
824		880
825		881
826		882
827		883
828	Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. <a href="#">S2ORC: The semantic scholar open research corpus</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4969–4983, Online. Association for Computational Linguistics.	884
829		885
830		886
831		887
832		888
833		
834	Roser Morante, Martin Krallinger, Alfonso Valencia, and Walter Daelemans. 2012. Machine reading of biomedical texts about alzheimers disease. In <i>CLEF 2012 Conference and Labs of the Evaluation Forum-Question Answering For Machine Reading Evaluation (QA4MRE)</i> , pages 1–14.	889
835		890
836		891
837		
838		892
839		893
840		894
841	NLP Team MosaicML. 2023. <a href="#">Introducing mpt-7b: A new standard for open-source, commercially usable llms</a> . Accessed: 2023-05-05.	895
842		896
843		897
844		898
845	Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, Renqian Luo, Scott Mayer McKinney, Robert Osazuwa Ness, Hoi-fung Poon, Tao Qin, Naoto Usuyama, Chris White, and Eric Horvitz. 2023. <a href="#">Can generalist foundation models outcompete special-purpose tuning? case study in medicine</a> .	899
846		900
847		901
848		902
849		903
850		
851	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	904
852		905
853	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> .	906
854		907
855		908
856		909
857		910
858		911
859		912
860	Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. <a href="#">Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering</a> . In <i>Proc. Conference on Health, Inference, and Learning</i> , volume 174 of <i>Proceedings of Machine Learning Research</i> , pages 248–260.	913
861		914
862		915
863		916
864		917
865		918
866		919
867	Dimitris Pappas, Petros Stavropoulos, Ion Androutsopoulos, and Ryan McDonald. 2020. <a href="#">BioMRC: A dataset for biomedical machine reading comprehension</a> . In <i>Proc. 19th SIGBioMed Workshop on Biomedical Language Processing</i> , pages 140–149, Online. Association for Computational Linguistics.	920
868		921
869		922
870		923
871		924
872		
873	Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei,	925
874		926
		927
		928
		929
		930
	and Julien Launay. 2023. <a href="#">The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only</a> . <i>arXiv preprint arXiv:2306.01116</i> .	
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>arXiv e-prints</i> .	
	Raffaele Raimondi, Nikolaos Tzoumas, Thomas Salisbury, Sandro Di Simplicio, and Mario R Romano. 2023. Comparative analysis of large language models in the royal college of ophthalmologists fellowship exams. <i>Eye</i> , pages 1–4.	
	RCOphth. 2022a. <a href="#">Frcophth sample mcqs part 1. Part 1 FRCOphth sample mcqs - Royal College of Ophthalmologists</a> .	
	RCOphth. 2022b. <a href="#">Frcophth sample mcqs part 2. Part 2 FRCOphth sample mcqs - Royal College of Ophthalmologists</a> .	
	Max Savery, Asma Ben Abacha, Soumya Gayen, and Dina Demner-Fushman. 2020. Question-driven summarization of answers to consumer health questions. <i>Scientific Data</i> , 7(1):322.	
	Viktor Schlegel, Goran Nenadic, and Riza Batista-Navarro. 2022. <a href="#">A survey of methods for revealing and overcoming weaknesses of data-driven Natural Language Understanding</a> . <i>Natural Language Engineering</i> , pages 1–31.	
	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023a. Large language models encode clinical knowledge. <i>Nature</i> , 620(7972):172–180.	
	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023b. Towards expert-level medical question answering with large language models. <i>arXiv</i> , 2305.09617.	
	Alon Talmor and Jonathan Berant. 2019. <a href="#">MultiQA: An Empirical Investigation of Generalization and Transfer in Reading Comprehension</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4911–4921, Stroudsburg, PA, USA. Association for Computational Linguistics.	
	Augustin Toma, Patrick R. Lawler, Jimmy Ba, Rahul G. Krishnan, Barry B. Rubin, and Bo Wang. 2023. Clinical camel: An open expert-level medical language model with dialogue-based knowledge encoding. <i>arXiv</i> , 2305.12031.	
	George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, Yannis Almirantis, John Pavlopoulos, Nicolas Baskiotis, Patrick Gallinari, Thierry Artieres,	



Axel Ngonga, Norman Heino, Eric Gaussier, Liliana Barrio-Alvers, Michael Schroeder, Ion Androutsopoulos, and Georgios Paliouras. 2015. [An overview of the bioasq large-scale biomedical semantic indexing and question answering competition](#). *BMC Bioinformatics*, 16:138.

Logesh Kumar Umapathi, Ankit Pal, and Malaikannan Sankarasubbu. 2023. Med-halt: Medical domain hallucination test for large language models. *arXiv*, 2307.15343.

Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, William Collins, Neera Ahuja, Curtis P. Langlotz, Jason Hom, Sergios Gatidis, John Pauly, and Akshay S. Chaudhari. 2023. Clinical text summarization: Adapting large language models can outperform human experts. *arXiv*, 2309.07430.

David Vilares and Carlos Gómez-Rodríguez. 2019. [HEAD-QA: A healthcare dataset for complex reasoning](#). In *Proc. 57th Annual Meeting of the Association for Computational Linguistics*, pages 960–966, Florence, Italy. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Towards building open-source language models for medicine. *arXiv*, 2304.14454.

Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023. [Baize: An open-source chat model with parameter-efficient tuning on self-chat data](#).

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mT5: A massively multilingual pre-trained text-to-text transformer](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Ming Zhu, Aman Ahuja, Da-Cheng Juan, Wei Wei, and Chandan K. Reddy. 2020. [Question answering with long multiple-span answers](#). In *Findings of the Association for Computational Linguistics: EMNLP*, pages 3840–3849. Association for Computational Linguistics.



## Appendix

### 5.1 Datasets Used

In this section, we explain the MCQA and AQA datasets we used in detail. The dataset characteristics are presented in Table 1.

1. **USMLE - English:** We incorporate the USMLE dataset obtained from the MedQA dataset (Jin et al., 2021), comprising MCQA questions from the Medical Licensing Exam conducted in the US. We retain this dataset’s original training, validation, and test set divisions.
2. **MEDMCQA:** We incorporate the MEDMCQA dataset from (Pal et al., 2022), which comprises medical MCQA from Indian Medical Entrance Exams. We retain this dataset’s original training, validation, and test set splits. Similar to (Singhal et al., 2023a), we evaluate all models on the validation set since we do not have answers for the test set.
3. **MMLU:** Following the design of Singhal et al. (Singhal et al., 2023a), we incorporate a subset of the MMLU datasets (6 datasets) (Hendrycks et al., 2021) which are MCQA specifically curated to assess medical domain knowledge. The subsets used are the **anatomy, clinical knowledge, college medicine, medical genetics, professional medicine** and **college biology** questions from MMLU. We utilize these datasets only for evaluating models.
4. **MEDIQA-ANS:** The MEDIQA 2019 shared task introduced the MEDIQA-QA dataset (Savery et al., 2020) for answer-ranking, comprising consumer health questions and passages from reputable online sources. The dataset was curated by extracting passages from the text of web pages, and includes manually generated single and multi-document summaries in both extractive and abstractive forms. We employ the multi-document abstractive summary as our questions’ ground truth reference answer. We specifically filter for questions and answers marked as excellent and utilize this as an AQA dataset solely for evaluating models.
5. **HEADQA:** We include the HEADQA dataset (Vilares and Gómez-Rodríguez, 2019), which

comprises graduate-level MCQA about various fields of medicine used for examinations to apply for specialization positions in the Spanish public healthcare system. We use the English version of the dataset and retain the original train, validation, and test split.

6. **PubmedQA:** The PubMedQA dataset (Jin et al., 2019) is a biomedical question-answering dataset comprising 1,000 expert-annotated QA instances. Each instance necessitates reasoning over a biomedical paper’s abstract to answer a relevant question. While the dataset provides long and short answers (yes, no, or maybe), we focus exclusively on the short answers for our evaluation, thereby generalizing the task as MCQA. We retain the original test split of 500 questions. Additionally, we allocate 100 questions from the training set to serve as a validation set, facilitating standardized training and validation in future studies.
7. **BioMRC:** The BioMRC dataset (Pappas et al., 2020) focuses on machine reading comprehension within the biomedical domain. It is structured in a cloze-style MCQA format, where questions are based on biomedical abstracts where biomedical entities are replaced with pseudo-identifiers. The task is to correctly identify the masked entity in the title from a list of masked entities. We utilize two compact versions of BioMRC: tiny A and tiny B, also referred to as Setting A and B, respectively. The BioMRC dataset comprises a large training corpus, where masked entities share the same pseudo-identifier across the entire corpus. Setting A, also known as tiny A, retains the same pseudo-identifiers used for masked biomedical entities in the training corpus. This setup is beneficial when testing models trained using the BioMRC training set, allowing them to draw on previously seen patterns. Tiny B (Setting B), conversely, changes the pseudo-identifiers for every single question. This means that a model must rely solely on the information in the text of the question and the passage it refers to, without any help from repeated exposure to the same placeholders. While we maintain the original format for Setting B, assessing Setting A as is, is difficult as since we do not utilize the BioMRC train-

ing set, it is functionally the same as Setting B. To address this limitation, we modify Setting A to include the original entity names and their corresponding pseudo-identifiers in the answer options. This aims to assess whether the model can accurately answer when provided with the information about their original entity names.

8. **Processbank**: The Processbank dataset (Berant et al., 2014) is designed for machine reading comprehension, featuring questions based on paragraphs describing biological processes. Each question, associated with a particular paragraph, has two answer options (MCQA). The dataset comes with a predefined split of 435 questions (150 files) for training and 100 questions (50 files) for testing. We allocate 25 files from the training set to create a validation set while retaining the original test set for model evaluation.
9. **QA4MRE - Alzheimer’s disease QA**: The dataset proposed by Morante et al. (Morante et al., 2012) contains MCQA questions regarding Alzheimer’s disease, aimed at assessing machine reading systems’ ability to answer questions about the disease by parsing relevant documents. We have adapted this dataset as an open-ended MCQA task to evaluate LLMs’ ability to answer these questions based on inherent knowledge. This dataset is employed solely for model evaluation purposes.
10. **BioASQ**: The BioASQ dataset (Tsatsaronis et al., 2015; Krithara et al., 2023) features biomedical questions crafted by experts. We utilize the BioASQ 2022 dataset for our benchmark. The BioASQ dataset is divided into two parts: for MCQA and another for AQA. For the MCQA part, we filter out the yes/no questions from BioASQ, converting them into an MCQ format to create a new subset, which we term **BioASQ-MCQ**. We manually create a training-validation (train-val) split of roughly 85%-15% from the filtered questions, resulting in 975 training questions and 173 validation questions and retaining a test set of 123 questions. For the AQA part, BioASQ provides fact, list, and bullet-type questions. We compile these into an AQA dataset, ensuring a balanced representation of all question

types in training and validation sets. The train-validation split results in 4733 training and 697 validation questions, with approximately 15% of all question types in the validation set.

11. **MASH-QA**: The MASH-QA dataset (Zhu et al., 2020) was designed for answering medical questions based on paragraphs where answers may span multiple text segments. Initially intended for extractive answering tasks, we repurpose it as an AQA task, utilizing the extractive answers as the reference ground truth.
12. **MedQUAD**: The MedQUAD dataset (Ben Abacha and Demner-Fushman, 2019) encompasses medical question-answer pairs extracted from various National Institute of Health (NIH) websites, covering topics on diseases, drugs, and other medical entities. Only nine of the twelve websites contributing to the original dataset have answers. We segregate questions from these nine websites and devise a train-validation-test split (AQA), assigning data from six websites for training, one website for validation, and two websites for testing.
13. **TREC-2017 LiveQA**: We employ the TREC-2017 LiveQA dataset (Abacha et al., 2017) for evaluation purposes. Specifically, we leverage the rankings provided within the MedQUAD evaluation process (Ben Abacha and Demner-Fushman, 2019) to keep question-answer pairs that have answer rating as excellent. We utilize this as an AQA dataset for evaluating the model.
14. **British Ophthalmology Practice Tests**: We employ sample questions from the Fellowship of the Royal College of Ophthalmologists (FRCOphth) exams, as provided by the Royal College of Ophthalmologists on their website (Raimondi et al., 2023; RCOphth, 2022a,b). These MCQA questions, geared towards testing ophthalmology-related knowledge, are used for evaluation.
15. **MEDINFO**: The MEDINFO dataset, introduced by Abacha et al. (Ben Abacha et al., 2019), consists of real consumer questions concerning medications and drugs. It encompasses 674 question-answer pairs (AQA), which we employ solely for evaluation.

## 5.2 Performance of other methods for MCQA datasets

We report the prior and current best scores on MCQA datasets from current literature in Table 7. GPT-4 combined with a prompting strategy labeled MedPrompt performs the best currently on USMLE, MEDMCQA, and the MMLU datasets. Of the 16 datasets, we can obtain comparable scores for 12. For HEADQA, the results reported by (Vilares and Gómez-Rodríguez, 2019) and (Liu et al., 2020) are across individual sections, whereas we calculate the scores overall across all questions. The method proposed by (Liu et al., 2020), named **MurKe** achieves average scores of 45.5% on Biology questions, 42.4% on Medicine questions, 42.3% on Nursing Questions, 48.0% on Pharmacology questions, 44.3% on Psychology questions and 44.3% on Chemistry Questions, with an overall macro-average of 44.4% across all the sections. Similarly, for the OPTH dataset, the results reported by (Raimondi et al., 2023) are separate for Part 1 and Part 2 questions. Bing Chat performs the best on Part 1 questions, achieving a performance of 78.9%, and GPT-4 with prompting obtains a performance of 88.4% on Part 2 questions (Raimondi et al., 2023). We could not find directly comparable scores for the **BioASQ** MCQ datasets as the test sets are provided in different batches, with the results on the BioASQ leaderboard also reported separately in terms of batches. We combine the questions across all the batches into one combined test set. For **BIO-MRC - Tiny A**, we do not have comparable scores from prior works as we formulate this task differently by providing the names of the original entities to the model.

Dataset	Best Reported Score	Method
USMLE (4 options)	90.2	GPT 4 + MedPrompt (Nori et al., 2023)
MEDMCQA	79.1	GPT 4 + MedPrompt (Nori et al., 2023)
PubMedQA	82.0	GPT 4 + MedPrompt (Nori et al., 2023)
MMLU - Anatomy	89.6	GPT 4 + MedPrompt (Nori et al., 2023)
MMLU - Clinical Knowledge	95.8	GPT 4 + MedPrompt (Nori et al., 2023)
MMLU - College Biology	97.9	GPT 4 + MedPrompt (Nori et al., 2023)
MMLU - College Medicine	89.0	GPT 4 + MedPrompt (Nori et al., 2023)
MMLU - Medical Genetics	98.0	GPT 4 + MedPrompt (Nori et al., 2023)
MMLU - Professional Medicine	95.2	GPT 4 + MedPrompt (Nori et al., 2023)
ProcessBank	68.8	SemanticILP (Biology Cascade) (Khashabi et al., 2018)
QA4MRE	55.0	Index Expansion (Attardi et al., 2012)
BioMRC - Tiny B	60.0	SciBERT-Max-Reader (Pappas et al., 2020)

Table 7: Performance scores of various methods on various MCQA datasets

Model	Architecture	# Tokens	Data Source
<i>Base models</i>			
MPT	Decoder	1T	Red Pajama (Computer, 2023), The Stack (Kocetkov et al., 2022), C4 (Raffel et al., 2019), mC4 (Xue et al., 2021), S20RC (Lo et al., 2020)
Falcon	Decoder	1.5T	RefinedWeb (Penedo et al., 2023)
LLaMA 2	Decoder	2T	Unknown
<i>Instruction tuned models</i>			
Flan-T5	Encoder-Decoder	1T	C4 (Raffel et al., 2019) and Flan-Collection (Wei et al., 2021)
MPT-Instruct	Decoder	1T	MPT, Databricks Dolly-15k (Conover et al., 2023), Anthropic Helpful and Harmless (Bai et al., 2022)
Falcon-Instruct	Decoder	1.5T	Falcon, baize (Xu et al., 2023), GPT4All, GPTeacher <sup>6</sup>
LLaMA 2-Chat	Decoder	2T	LLaMA 2, Flan Collection (Wei et al., 2021), Private Data

Table 8: Pretrained LLMs considered in this paper. (Top rows) Open-source models that are decoder-only. (Bottom rows) Instruction-fine-tuned language models. **# Tokens**: Number of tokens used in pretraining the model. **Data Source**: Data used for pre-training (instruction data is *italicized*).

Parameter	Flan-T5 XL	Llama-2 7B	Falcon 7B	MPT 7B
lora_r	16	16	16	16
lora_alpha	16	16	16	16
lora_dropout	0.05	0.05	0.05	0.05
bias	none	none	none	none
optimizer	adamw	adamw	adamw	adamw
epochs	4	4	4	4
batch size	8	8	8	8
model_max_length	256	384	384	384

Table 9: Hyper-parameters used to train our models

Parameter	Decoder LLMs	Encoder-Decoder LLMs
Beam Size	3	3
Repetition Penalty	1.5	1.5
Max Output Length	200	200

Table 10: Inference time parameters used for abstractive question answering



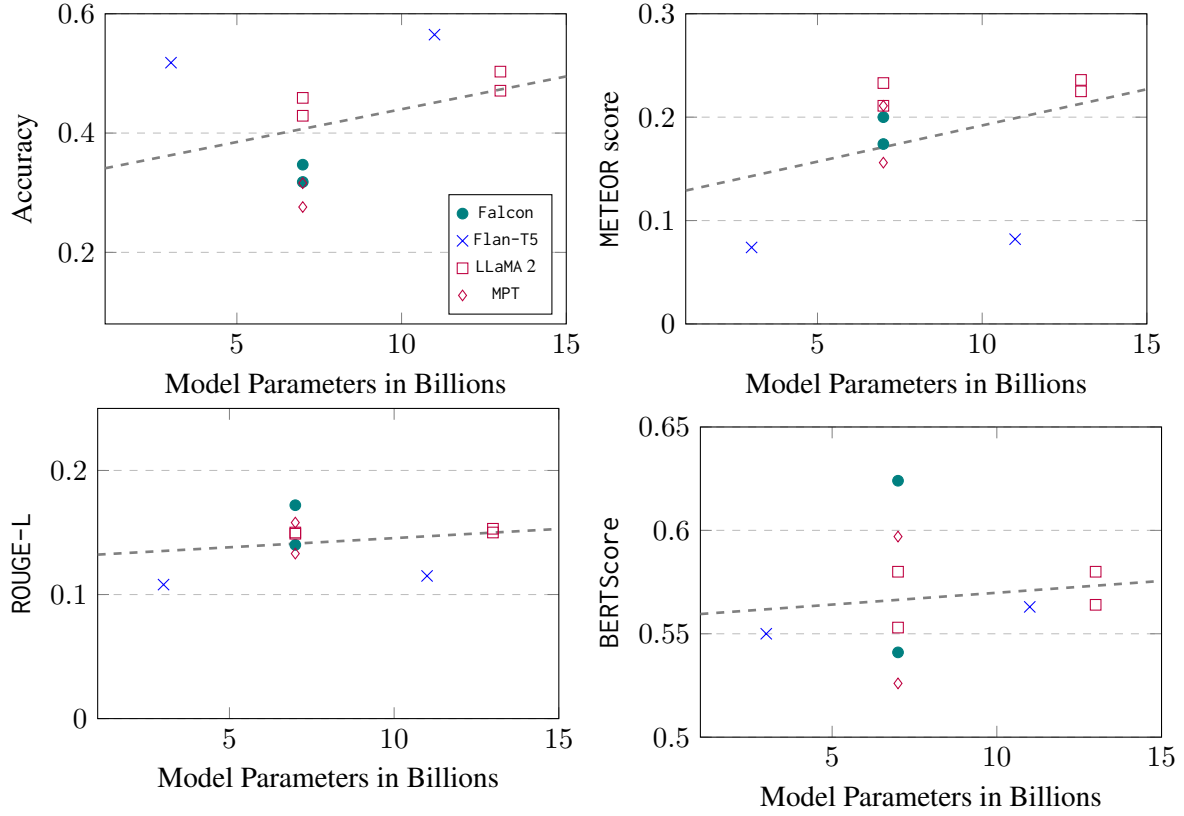


Figure 4: Zero-shot performance of models on MCQA (top-left) and AQA (top-right, bottom-left and bottom-right) as a function of model size. The dashed line represents a fitted linear regression showing the correlation between the model size and the score.

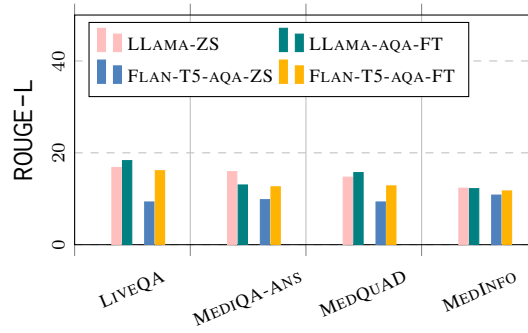


Figure 5: Performance of base and AQA-finetuned LLaMA 2 and Flan-T5 models on four unseen AQA test sets in terms of ROUGE-L.

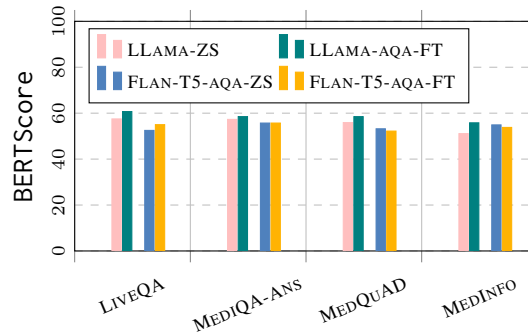


Figure 6: Performance of base and AQA-finetuned LLaMA 2 and Flan-T5 models on four unseen AQA test sets in terms of BERTScore.

Dataset	Falcon (7B)	MPT (7B)	LLaMA 2 (7B)	LLaMA 2 (13B)
BioASQ-MCQ	72.4	33.3	67.5	35.8
BioMRC Tiny A	26.7	23.3	30.0	53.3
BioMRC Tiny B	16.7	13.3	26.7	20.0
MMLU - Anatomy	28.1	26.7	40.7	54.1
MMLU - Clinical Knowledge	32.5	29.8	38.1	57.7
MMLU - College Biology	27.1	22.2	39.6	58.3
MMLU - College Medicine	30.6	26.6	35.3	54.3
MMLU - Medical Genetics	33.0	27.0	49.0	52.0
MMLU - Professional Medicine	44.1	20.2	44.1	53.7
HEADQA	27.8	28.0	40.4	48.5
MEDMCQA	30.4	26.5	36.0	37.5
OPHTH	21.7	28.3	27.2	30.4
PROCESSBANK	50.7	56.0	75.3	83.3
PUBMEDQA	57.0	33.8	60.4	33.8
QA4MRE	30.0	22.5	40.0	37.5
USMLE	27.0	24.2	35.3	42.9
Average	34.7	27.6	42.9	47.1

Table 11: MCQA scores of LLMs in the zero-shot setting. We utilize 5-shot prompting for the MMLU datasets and 1-shot prompting for other datasets to evaluate these models.

Dataset	Flan-T5 (3B)	Falcon (7B)	MPT (7B)	LLaMA 2 (7B) Chat	Flan-T5 (11B)	LLaMA 2 (13B) Chat
BioASQ-MCQ	43.9	45.5	34.1	69.9	48.8	65.0
BioMRC Tiny A	73.3	30.0	23.3	26.7	63.3	33.3
BioMRC Tiny B	46.7	23.3	23.3	20.0	60.0	26.7
MMLU - Anatomy	46.7	27.4	32.6	44.4	48.9	52.6
MMLU - Clinical Knowledge	52.1	31.7	36.6	54.3	61.9	57.7
MMLU - College Biology	48.6	25.0	29.9	55.6	54.9	59.0
MMLU - College Medicine	41.6	27.7	30.1	44.5	52.6	46.2
MMLU - Medical Genetics	50.0	32.0	32.0	60.0	55.0	56.0
MMLU - Professional Medicine	42.6	37.9	28.3	45.2	55.1	51.1
HEADQA	42.9	26.1	30.2	43.9	49.1	51.3
MEDMCQA	33.1	29.8	27.2	35.0	36.4	39.3
OPHTH	26.1	32.6	30.4	26.1	25.0	27.2
PROCESSBANK	93.3	52.0	56.7	72.0	95.3	80.0
PUBMEDQA	70.0	47.4	35.6	61.6	70.8	45.2
QA4MRE	82.5	15.0	30.0	40.0	87.5	72.5
USMLE	36.1	25.1	24.6	35.6	39.7	42.2
Average	51.8	31.8	31.6	45.9	56.5	50.3

Table 12: MCQA scores of Instruction-tuned LLMs in the zero-shot setting. We utilize 5-shot prompting for the MMLU datasets and 1-shot prompting for other datasets to evaluate these models.

Dataset	Flan-T5 (3B)	Falcon (7B)	MPT (7B)	LLaMA 2 (7B)
BioASQ-MCQ	73.2	80.5	78.9	81.3
BioMRC Tiny A	53.3	23.3	26.7	23.3
BioMRC Tiny B	26.7	23.3	20.0	26.7
MMLU - Anatomy	43.7	43.7	45.9	54.1
MMLU - Clinical Knowledge	54.0	52.8	53.2	59.6
MMLU - College Biology	47.2	46.5	56.9	61.1
MMLU - College Medicine	44.5	53.2	50.3	52.0
MMLU - Medical Genetics	47.0	55.0	60.0	62.0
MMLU - Professional Medicine	48.5	50.0	49.3	59.6
HEADQA	49.0	47.7	52.4	53.9
MEDMCQA	43.0	45.9	48.4	48.3
OPHTH	34.8	30.4	35.9	31.5
PROCESSBANK	92.7	69.3	84.7	75.3
PUBMEDQA	74.2	70.8	73.4	70.6
QA4MRE	75.0	50.0	70.0	50.0
USMLE	39.7	46.3	45.7	46.1
Average	52.9	49.3	53.2	53.5

Table 13: MCQA scores of LLMs finetuned with QLora on MCQA datasets from the M-QALM benchmark. We evaluate these models without any examples in the prompt.

Dataset	Flan-T5 (3B)	Falcon (7B)	MPT (7B)	LLaMA 2 (7B)
BioASQ-MCQ	0.8	13.8	14.6	7.3
BioMRC Tiny A	50.0	23.3	10.0	16.7
BioMRC Tiny B	36.7	23.3	16.7	16.7
MMLU - Anatomy	43.0	24.4	34.8	38.5
MMLU - Clinical Knowledge	50.9	25.3	28.7	40.8
MMLU - College Biology	42.4	23.6	34.7	38.9
MMLU - College Medicine	41.0	27.2	26.0	37.6
MMLU - Medical Genetics	45.0	31.0	22.0	49.0
MMLU - Professional Medicine	41.2	44.1	18.4	46.7
HEADQA	38.7	21.5	24.8	31.1
MEDMCQA	27.0	21.7	20.2	23.0
OPHTH	22.8	23.9	16.3	19.6
PROCESSBANK	88.0	54.7	42.0	50.7
PUBMEDQA	67.2	57.2	54.6	47.8
QA4MRE	77.5	35.0	10.0	15.0
USMLE	34.2	22.9	23.9	22.9
Average	44.1	29.6	24.9	31.4

Table 14: MCQA scores of LLMs finetuned with QLoRA on AQA datasets only from the M-QALM benchmark. We utilize 5-shot prompting for the MMLU datasets and 1-shot prompting for other datasets to evaluate these models.

Dataset	Flan-T5 (3B)	Falcon (7B)	MPT (7B)	LLaMA 2 (7B)
BioASQ-MCQ	71.5	80.5	79.7	79.7
BioMRC Tiny A	50.0	43.3	36.7	26.7
BioMRC Tiny B	30.0	6.7	20.0	26.7
MMLU - Anatomy	40.7	45.2	47.4	52.6
MMLU - Clinical Knowledge	51.7	52.5	50.9	55.5
MMLU - College Biology	43.8	51.4	57.6	61.1
MMLU - College Medicine	41.6	48.0	54.3	52.6
MMLU - Medical Genetics	52.0	59.0	55.0	65.0
MMLU - Professional Medicine	47.1	46.0	50.4	59.9
HEADQA	47.5	47.4	51.2	54.2
MEDMCQA	41.7	45.2	47.4	48.0
OPHTH	32.6	28.3	38.0	28.3
PROCESSBANK	91.3	73.3	79.3	83.3
PUBMEDQA	71.4	67.8	72.8	71.8
QA4MRE	72.5	52.5	60.0	67.5
USMLE	40.9	45.7	44.3	45.6
Average	51.7	49.5	52.8	54.9

Table 15: MCQA scores of LLMs finetuned with QLoRA on both MCQA and AQA data from the M-QALM benchmark. We evaluate these models without any examples in the prompt.

Dataset	ChatDoctor (7B)	MedAlpaca (7B)	PMC-LLaMA (13B)
BioASQ-MCQ	65.0	50.4	13.0
BioMRC Tiny A	20.0	16.7	30.0
BioMRC Tiny B	36.7	23.3	16.7
MMLU - Anatomy	43.7	60.0	63.0
MMLU - Clinical Knowledge	43.4	60.0	62.3
MMLU - College Biology	39.6	64.6	64.6
MMLU - College Medicine	32.4	52.6	53.2
MMLU - Medical Genetics	55.0	69.0	70.0
MMLU - Professional Medicine	47.1	67.3	67.6
HEADQA	37.2	45.1	59.1
MEDMCQA	29.4	35.0	56.5
OPHTH	30.4	23.9	46.7
PROCESSBANK	62.0	67.3	74.7
PUBMEDQA	67.4	40.8	72.6
QA4MRE	45.0	62.5	55.0
USMLE	31.3	42.4	54.7
Average	42.8	48.8	53.7

Table 16: MCQA scores of ChatDoctor (7B), MedAlpaca (7B) and PMC-LLaMA (13B). To evaluate ChatDoctor, we utilize 5-shot prompting for the MMLU datasets and 1-shot prompting for other datasets to evaluate these models. We evaluate MedAlpaca (7B) and PMC-LLaMA (13B) directly without any examples in the prompt.

Model	BioASQ-QA			LIVEQA			MASHQA			MEDINFO			MEDIQA-ANS			MEDQUAD			Average		
	RL	BS	MTR	RL	BS	MTR	RL	BS	MTR	RL	BS	MTR	RL	BS	MTR	RL	BS	MTR	RL	BS	MTR
Falcon (7B)	13.9	53.1	22.5	15.4	55.8	17.4	13.4	53.7	22.0	12.1	51.1	17.8	15.3	56.1	21.7	14.3	54.7	18.4	14.0	54.1	20.0
MPT (7B)	11.4	50.1	21.7	15.7	55.2	20.9	12.8	52.3	23.0	11.2	49.6	18.4	14.8	55.6	23.3	13.7	53.2	19.4	13.3	52.6	21.1
LLaMA 2 (7B)	15.8	54.6	24.0	16.8	57.5	20.1	14.0	55.4	23.3	12.3	51.1	17.8	15.9	57.3	22.3	14.7	55.9	19.4	14.9	55.3	21.1
LLaMA 2 (13B)	14.9	55.3	24.9	16.2	57.3	20.1	14.5	56.4	24.4	12.7	53.6	20.0	16.4	58.9	24.4	15.4	57.1	20.9	15.0	56.4	22.5
Flan-T5 (3B)	15.0	57.7	11.1	9.3	52.5	6.1	10.5	56.0	7.5	10.8	54.9	7.6	9.8	55.7	6.2	9.3	53.2	6.0	10.8	55.0	7.4
MPT (7B) Instruct	23.2	64.5	22.4	14.5	58.1	13.4	15.0	61.1	15.9	14.0	56.8	12.9	14.8	60.5	16.1	12.9	57.1	13.1	15.8	59.7	15.6
Falcon (7B) Instruct	27.2	68.9	28.1	16.1	61.4	14.7	15.5	62.5	17.1	14.7	58.4	15.2	15.4	62.4	15.4	14.3	60.8	14.2	17.2	62.4	17.4
LLaMA 2 (7B) Chat	15.9	58.8	26.5	15.4	58.8	20.9	14.2	57.4	24.4	12.8	54.6	20.6	16.7	59.5	25.4	15.4	58.7	22.1	15.0	58.0	23.3
Flan-T5 (11B)	16.3	58.8	12.2	10.8	55.5	7.5	10.8	57.3	8.2	12.3	56.1	9.1	9.7	55.2	6.3	9.0	54.9	5.9	11.5	56.3	8.2
LLaMA 2 (13B) Chat	16.2	59.2	27.5	15.8	59.0	21.4	14.2	57.2	24.3	13.0	54.7	21.2	16.7	58.9	24.8	15.5	58.7	22.4	15.3	58.0	23.6
Flan-T5 (3B) (FT-QA)	26.6	66.2	25.2	16.1	55.0	16.9	15.4	58.2	16.4	11.7	53.8	10.5	12.6	55.7	12.0	12.8	52.2	12.7	15.9	56.8	15.6
Falcon (7B) (FT-QA)	27.8	68.4	26.6	20.1	60.6	21.1	16.7	61.3	17.8	12.4	56.5	9.4	12.8	57.9	11.6	14.8	57.5	16.2	17.4	60.4	17.1
LLaMA 2 (7B) (FT-QA)	30.0	69.7	28.2	18.3	60.7	19.2	16.9	61.9	17.5	12.2	55.8	9.0	13.0	58.5	11.2	15.7	58.5	16.6	17.7	60.8	16.9
MPT (7B) (FT-QA)	28.9	69.0	27.6	18.6	59.6	20.6	16.4	61.0	17.5	12.9	56.1	10.7	13.1	57.6	11.5	14.0	56.5	15.4	17.3	60.0	17.2
Flan-T5 (3B) (FT-All)	27.8	67.4	25.7	16.0	55.8	17.1	15.5	59.3	15.3	11.4	54.5	9.3	11.7	55.7	10.4	13.0	53.1	13.1	15.9	57.6	15.2
Falcon (7B) (FT-All)	27.3	68.6	26.1	18.9	59.9	19.8	16.1	61.0	16.7	11.7	55.4	8.0	12.8	58.0	10.9	14.8	57.5	16.5	16.9	60.1	16.3
MPT (7B) (FT-All)	29.1	68.8	27.4	18.2	59.2	20.4	16.5	61.5	17.0	13.4	56.4	11.5	13.5	57.5	12.3	14.5	56.7	16.6	17.5	60.0	17.5
LLaMA 2 (7B) (FT-All)	30.2	69.7	27.8	17.9	60.4	17.9	17.3	61.9	17.7	12.4	54.9	9.9	13.3	58.3	12.2	15.0	57.7	15.5	17.7	60.5	16.8
ChatDoctor	26.2	68.2	28.8	15.8	61.3	16.0	16.1	62.6	18.6	15.2	58.9	15.6	16.5	62.9	18.2	14.8	60.2	15.0	17.4	62.3	18.7
MedAlpaca 7B	26.4	67.8	27.1	14.7	55.6	13.0	13.4	59.3	15.0	12.3	55.1	12.6	13.9	59.0	15.4	12.5	56.8	10.2	15.5	58.9	15.6
PMC LLaMA 13B	19.7	62.6	20.9	12.7	55.8	11.0	13.5	58.8	14.4	45.6	70.7	43.6	14.8	59.6	14.0	11.9	57.0	10.1	19.7	60.7	19.0

Table 17: AQA scores of base, instruction-tuned LLMs in the zero-shot setting, LLMs fine-tuned with QLoRA and other biomedical and clinical instruction tuned models such as ChatDoctor (7B), MedAlpaca (7B), PMC-LLaMA (13B). FT-QA refers to models fine-tuned only with AQA data and FT-All refers to models fine-tuned with both MCQA and AQA data.



1217 **5.3 Prompts for Fine-Tuned Falcon (Base),**  
1218 **MPT (Base), LLaMA 2 (Base) and Flan-T5**

1219 **5.3.1 AQA Prompt**

Answer the medical question precisely and factually  
Question: {Question}  
Answer:

Figure 7: AQA prompt utilized without any examples in the prompt. We finetune and evaluate these models utilizing this prompt format.

1220 **5.3.2 MCQA Prompt**

Pick the right option that answers the question  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text}  
Answer:

Figure 8: MCQA prompt utilized without any examples in the prompt. We finetune and evaluate the models utilizing this prompt format.

1221 **5.3.3 Single Context MCQA Prompt**

Given the context, pick the right choice that answers the question  
Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:

Figure 9: Single Context MCQA prompt utilized without any examples in the prompt. We finetune and evaluate these models utilizing this prompt format for the PRO-CESSBANK dataset.

1222 **5.3.4 Multi Context MCQA Prompt**

Given the context, pick the right choice that answers the question  
Contexts: {Context Paragraph 1}  
{Context Paragraph 2}  
.  
.  
{Context paragraph N}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:

Figure 10: Multi Context MCQA prompt utilized without any examples in the prompt. We finetune and evaluate these models utilizing this prompt format for the PUB-MEDQA dataset.

1223 **5.3.5 Cloze MCQA Prompt**

Given the context, pick the right choice that corresponds to the XXXX in the question  
Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:

Figure 11: Cloze MCQA prompt utilized without any examples in the prompt. The BIOMRC datasets follow this format. We evaluate these models utilizing this prompt format.

**5.4 Prompts for evaluating Falcon (Base and Instruct), MPT (Base), LLaMA 2 (Base) and Flan-T5 in the Zero-Shot setting**

1224

1225

1226

**5.4.1 Few-Shot MCQA Prompt**

1227

Pick the right option that answers the question  
Question: {Example 1}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text}  
Answer:{Correct Option}  
.  
.  
Question: {Example K}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text}  
Answer:{Correct Option}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text}  
Answer:

Figure 12: Format of the Few-Shot MCQA prompt utilized. We utilize this prompt for evaluating models prior to any fine-tuning only. 5-shot prompting is utilized for the MMLU datasets whereas 1-shot prompting is utilized for all other MCQA datasets when evaluating non-finetuned models.

**5.4.2 1-Shot Cloze Prompt**

1228

Given the context, pick the right choice that corresponds to the XXXX in the question  
Context: {Context Paragraph}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:{Correct Option}  
Given the context, pick the right choice that corresponds to the XXXX in the question  
Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:

Figure 13: Cloze MCQA prompt utilized without any examples in the prompt. The BIOMRC datasets follow this format. We evaluate these models utilizing this prompt format.

**5.4.3 1-Shot Single Context MCQA Prompt**

1229

Given the context, pick the right choice that answers the question  
Context: {Context Paragraph}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:{Correct Option}  
Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:

Figure 14: Format of the 1-Shot Single Context MCQA prompt utilized. We adopt this prompt format for the PROCESSBANK dataset.

#### 5.4.4 1-Shot Multi Context MCQA Prompt

Given the context, pick the right choice that answers the question  
Contexts: {Context Paragraph 1}  
{Context Paragraph 2}  
.  
.  
{Context Paragraph n}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:{Correct Option}  
Contexts: {Context Paragraph 1}  
{Context Paragraph 2}  
.  
.  
{Context Paragraph n}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
Answer:

Figure 15: Format of the 1-Shot Multi-Context MCQA prompt utilized. We adopt this prompt format for the PUBMEDQA dataset.

#### 5.4.5 AQA Prompt

Answer the medical question precisely and factually  
Question: {Question}  
Answer:

Figure 16: AQA prompt utilized without any examples in the prompt.

### 5.5 Prompts for evaluating LLaMA 2 (Chat) Models in the Zero-Shot setting

#### 5.5.1 AQA Prompt

[INST] <<SYS>>  
Answer the medical question precisely and factually  
<</SYS>>  
  
Question: {Question} [/INST]

Figure 17: AQA prompt utilized without any examples in the prompt.

#### 5.5.2 Few-Shot MCQA Prompt

[INST] <<SYS>>  
Pick the right option that answers the question  
<</SYS>>  
  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text} [/INST] Answer:{Correct Option} </s><s>[INST] Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text} [/INST] Answer:

Figure 18: Format of the Few-Shot MCQA prompt utilized. 5-shot prompting is utilized for the MMLU datasets whereas 1-shot prompting is utilized for all other MCQA datasets

#### 5.5.3 1-Shot Single Context MCQA Prompt

[INST] <<SYS>>  
Given the context, pick the right choice that answers the question  
<</SYS>>  
  
Context: {Context Paragraph}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text} [/INST] Answer:{Correct Option} </s><s>[INST] Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text} [/INST] Answer:

Figure 19: Format of the 1-Shot Single Context MCQA prompt utilized. We utilize this prompt for evaluating models on the PROCESSBANK dataset.

#### 5.5.4 1-Shot Multi Context MCQA Prompt

[INST] <<SYS>>  
Given the context, pick the right choice that answers the question  
<</SYS>>  
  
Contexts: {Context Paragraph 1}  
{Context Paragraph 2}  
.  
.  
{Context Paragraph N}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text} [/INST] Answer:{Correct Option} </s><s>[INST] Contexts: {Context Paragraph 1}  
{Context Paragraph 2}  
.  
.  
{Context Paragraph N}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text} [/INST] Answer:

Figure 20: Format of the 1-Shot Multi-Context MCQA prompt utilized. We adopt this prompt format for the PUBMEDQA dataset.

#### 5.5.5 Few-Shot Cloze MCQA Prompt

[INST] <<SYS>>  
Given the context, pick the right choice that corresponds to the XXXX in the question  
<</SYS>>  
  
Context: {Context Paragraph}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text} [/INST] Answer:{Correct Option} </s><s>[INST] Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text} [/INST] Answer:

Figure 21: Format of the 1-Shot Cloze MCQA prompt utilized. We utilize this prompt for evaluating models on the BIOMRC datasets in settings A and B.

1239  
1240  
1241

## 5.6 Prompts for evaluating MPT Instruct in the Zero-Shot setting

### 5.6.1 AQA Prompt

Below is an instruction that describes a task. Write a response that appropriately completes the request.  
### Instruction:  
Answer the medical question precisely and factually. Question: {Question}  
### Response:  
Answer:

Figure 22: AQA prompt utilized without any examples in the prompt.

1242

### 5.6.2 Few-Shot MCQA Prompt

Below is an instruction that describes a task. Write a response that appropriately completes the request.  
### Instruction:  
Pick the right option that answers the question. Question: {Example Question 1}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text}  
### Response:  
Answer:{Correct Option}  
.  
.  
### Instruction:  
Pick the right option that answers the question. Question: {Example Question K}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text}  
### Response:  
Answer:{Correct Option}  
### Instruction:  
Pick the right option that answers the question. Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
C. {Option Text}  
D. {Option Text}  
### Response:  
Answer:

Figure 23: Format of the Few-Shot MCQA prompt utilized. 5-shot prompting is utilized for the MMLU datasets whereas 1-shot prompting is utilized for all other MCQA datasets.

1243

### 5.6.3 1-Shot Single Context MCQA Prompt

Below is an instruction that describes a task. Write a response that appropriately completes the request.  
### Instruction:  
Given the context, pick the right choice that answers the question. Context: {Context Paragraph}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
### Response:  
Answer:{Correct Option}  
### Instruction:  
Given the context, pick the right choice that answers the question. Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
### Response:  
Answer:

Figure 24: Format of the 1-Shot Single Context MCQA prompt utilized. We adopt this prompt format for the PROCESSBANK dataset.

1244  
1245

1246

### 5.6.4 1-Shot Multi Context MCQA Prompt

Below is an instruction that describes a task. Write a response that appropriately completes the request.  
### Instruction:  
Given the contexts, pick the right choice that answers the question. Contexts: {Context Paragraph 1}  
{Context Paragraph 2}  
.  
.  
{Context Paragraph N}  
Question: {Example Question 1}  
Options:  
A. {Option Text}  
B. {Option Text}  
### Response:  
Answer:{Correct Option}  
### Instruction:  
Given the contexts, pick the right choice that answers the question. Contexts: {Context Paragraph 1}  
{Context Paragraph 2}  
.  
.  
{Context Paragraph N}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
### Response:  
Answer:

Figure 25: Format of the 1-Shot Multi-Context MCQA prompt utilized. We adopt this prompt format for the PUBMEDQA dataset.

### 5.6.5 1-Shot Cloze MCQA Prompt

1247

Below is an instruction that describes a task. Write a response that appropriately completes the request.  
### Instruction:  
Given the context, pick the right choice that corresponds to the XXXX in the question. Context: {Context Paragraph}  
Question: {Example Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
### Response:  
Answer:{Correct Option}  
### Instruction:  
Given the context, pick the right choice that corresponds to the XXXX in the question. Context: {Context Paragraph}  
Question: {Question}  
Options:  
A. {Option Text}  
B. {Option Text}  
### Response:  
Answer:

Figure 26: Format of the 1-Shot Cloze MCQA prompt utilized. We utilize this prompt for evaluating models on the BIOMRC datasets in settings A and B.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```
### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Answer with the best option directly.

### Input:
{Example Question}
Options:
A. {Option Text}
B. {Option Text}
C. {Option Text}
D. {Option Text}

### Response:
Answer:{Correct Option}

### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Answer with the best option directly.

### Input:
{Question}
Options:
A. {Option Text}
B. {Option Text}
C. {Option Text}
D. {Option Text}

### Response:
Answer:
```

Figure 27: Format of the Few-Shot MCQA prompt utilized for evaluating ChatDoctor. We utilize this prompt for evaluating non-finetuned models. 5-shot prompting is utilized for the MMLU datasets whereas 1-shot prompting is utilized for all other MCQA datasets.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```
### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Analyze the question given its context. Answer with the best option directly.

### Input:
Context: {Context Paragraph}
Question: {Example Question}
Options:
A. {Option Text}
B. {Option Text}

### Response:
Answer:{Correct Option}

### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Analyze the question given its context. Answer with the best option directly.

### Input:
Context: {Context Paragraph}
Question: {Question}
Options:
A. {Option Text}
B. {Option Text}

### Response:
Answer:
```

Figure 28: Format of the 1-Shot Single Context MCQA prompt utilized for evaluating ChatDoctor. We adopt this prompt format for the PROCESSBANK dataset.



Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```
### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Analyze the question given its context. Answer with the best option directly.

### Input:
Contexts: {Context Paragraph 1}
{Context Paragraph 2}
.
.
{Context Paragraph N}
Question: {Example Question}
Options:
A. {Option Text}
B. {Option Text}

### Response:
Answer:{Correct Option}

### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Analyze the question given its context. Answer with the best option directly.

### Input:
Contexts: {Context Paragraph 1}
{Context Paragraph 2}
.
.
{Context Paragraph N}
Question: {Question}
Options:
A. {Option Text}
B. {Option Text}

### Response:
Answer:
```

Figure 29: Format of the 1-Shot Multi-Context MCQA prompt utilized for evaluating ChatDoctor. We adopt this prompt format for the PUBMEDQA dataset.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```
### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Analyze the question given its context. Pick the right option that corresponds to the XXXX in the question

### Input:
Context: {Context Paragraph}
Question: {Question}
Options:
A. {Option Text}
B. {Option Text}

### Response:
Answer:{Correct Option}

### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description. Analyze the question given its context. Pick the right option that corresponds to the XXXX in the question

### Input:
Context: {Context Paragraph}
Question: {Question}
Options:
A. {Option Text}
B. {Option Text}

### Response:
Answer:
```

Figure 30: Format of the 1-Shot Cloze MCQA prompt utilized for evaluating ChatDoctor on the BIOMRC datasets in settings A and B.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```
### Instruction:
If you are a doctor, please answer the medical questions based on the patient's description.

### Input:
{Question}

### Response:
```

Figure 31: AQA prompt utilized without any examples in the prompt for evaluating ChatDoctor.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```
### Instruction:
Answer this multiple-choice question.

### Input:
{Question}
A: {Option Text}
B: {Option Text}
C: {Option Text}
D: {Option Text}

### Response:
The Answer to the question is:
```

Figure 32: Format of the Zero-Shot MCQA prompt utilized for evaluating MedAlpaca.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
Analyze the question given its context. Answer this multiple-choice question.

### Input:
Context: {Context Paragraph}

{Question}
A: {Option Text}
B: {Option Text}

### Response:
The Answer to the question is:

```

Figure 33: Format of the Single Context MCQA prompt utilized for evaluating MedAlpaca on the PROCESSBANK dataset

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
Analyze the question given its context. Answer this multiple-choice question.

### Input:
Contexts: {Context Paragraph 1}
{Context Paragraph 2}
.
.
{Context Paragraph N}

{Question}
A: {Option Text}
B: {Option Text}
C: {Option Text}

### Response:
The Answer to the question is:

```

Figure 34: Format of the Multi-Context MCQA prompt utilized for evaluating MedAlpaca on the PUBMEDQA dataset.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
Analyze the question given its context. Pick the right option that corresponds to the XXXX in the question.

### Input:
Context: {Context Paragraph}

{Question}
A: {Option Text}
B: {Option Text}
C: {Option Text}
D: {Option Text}

### Response:
The Answer to the question is:

```

Figure 35: Format of the Cloze MCQA prompt utilized for evaluating MedAlpaca on the BIOMRC datasets in settings A and B.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
Answer this question truthfully

### Input:
{Question}

### Response:

```

Figure 36: AQA prompt utilized without any examples in the prompt for evaluating MedAlpaca.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
You're a doctor, kindly address the medical queries according to the patient's account. Answer with the best option directly.

### Input:
###Question: {Question}
###Options:
A. {Option Text}
B. {Option Text}
C. {Option Text}
D. {Option Text}

### Response:
###Answer:

```

Figure 37: Format of the Zero-Shot MCQA prompt utilized for evaluating PMC-LLama.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
You're a doctor, kindly address the medical queries according to the patient's account. Analyze the question given its context. Answer with the best option directly.

### Input:
###Question: {Question}
###Context: {Context Paragraph}
###Options:
A. {Option Text}
B. {Option Text}

### Response:
###Answer:

```

Figure 38: Format of the Single Context MCQA prompt utilized for evaluating PMC-LLama on the PROCESSBANK dataset.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
You're a doctor, kindly address the medical queries according to the patient's account. Analyze the question given its context. Answer with the best option directly.

### Input:
###Question: {Question}
###Contexts: {Context Paragraph 1}
{Context Paragraph 2}
.
.
{Context Paragraph N}
###Options:
A. {Option Text}
B. {Option Text}
C. {Option Text}

### Response:
###Answer:

```

Figure 39: Format of the Multi-Context MCQA prompt utilized for evaluating PMC-LLama on the PUBMEDQA dataset.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
You're a doctor, kindly address the medical queries according to the patient's account. Analyze the question given its context. Pick the right option that corresponds to the XXXX in the question

### Input:
###Question: {Question}
###Context: {Context Paragraph}
###Options:
A. {Option Text}
B. {Option Text}

### Response:
###Answer:

```

Figure 40: Format of the Cloze MCQA prompt utilized for evaluating PMC-LLama on the BIOMRC datasets in settings A and B.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

```

### Instruction:
You're a doctor, kindly address the medical queries according to the patient's account.

### Input:
###Question: {Question}

### Response:
###Answer:

```

Figure 41: AQA prompt utilized without any examples in the prompt for evaluating PMC-LLama.