

QALM – Clinical QA for Language Models

Anonymous ACL submission

Abstract

In recent years, Large Language Models (LLMs) have gained recognition for their ability to encode clinical 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 knowledge applications. In response to this gap, we introduce a novel benchmark called QALM designed to harmonize the evaluation of language models in the context of clinical knowledge. Our benchmark comprises 16 Multiple-Choice Question (MCQ) datasets and six 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 decoder-only language models may not be the optimal choice for MCQs in clinical knowledge tasks. Additionally, our investigation demonstrates that instruction fine-tuned language models do not necessarily outperform their 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

Large Language Models (LLM) deployed in the clinical and biomedical domains have the potential to revolutionize the healthcare industry. They are employed to summarize clinical text (Veen et al., 2023), automatically generate notes for clinicians (Ben Abacha et al., 2023b), and condense dialogues between doctors and patients (Ben Abacha et al.,

2023a; Toma et al., 2023). Recognizing their significance, recent work (Han et al., 2023; Wu et al., 2023; Toma et al., 2023; Bolton et al., 2022; Li et al., 2023) has focused on fine-tuning LLMs on clinical and bio-medical datasets. However, the evaluation of LLMs within the clinical and biomedical domains remains incomplete and requires further comprehensive evaluation. Recent models tend to be evaluated on different datasets or tasks, which makes fair comparison of models harder.

Clinical knowledge assessment in LLMs involves two primary tasks (Singhal et al., 2023a): Multiple Choice Question Answering (MCQA), where answers are selected from multiple options, and abstractive question answering (AQA), which entails generating answers to questions, either with or without a provided paragraph context.

The evaluation of LLMs regarding their clinical knowledge is restricted. For instance, Singhal et al. (2023a) assess proprietary models on a consolidated dataset, but open-source LLMs are not tested on such a unified benchmark. We expand Singhal et al. (2023a)’s benchmark with more datasets, to enable transparency and reproducibility, and to test new advances in LLM research. Our benchmark called QALM consolidates existing MCQA and AQA datasets, featuring 16 MCQA and 6 AQA datasets. With such a standardized benchmark we are able to test the strengths and weaknesses of open-source models using a methodological unified framework, which is currently missing in the literature.

Our evaluation encompasses diverse zero-shot and fine-tuning settings. Our findings reveal that the latest decoder-only LLMs do not consistently outperform others on reading comprehension datasets, where models are presented with a context paragraph to answer a question. Moreover, although instruction fine-tuned models have been argued to surpass their non-instruction fine-tuned counterparts in some contexts (Wei et al., 2021; Gupta et al., 2023), our results suggest a more nu-

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

Dataset	Type	Size	Domain
USMLE (Jin et al., 2021)	MCQA	10178/1272/1273	Medical Exam
MEDMCQA (Pal et al., 2022)	MCQA	182822/4183/6150	Medical Exam
BIOASQ-MCQ (Tsatsaronis et al., 2015; Krithara et al., 2023)	MCQA	975/173/123	Biomedical
HEADQA (Vilares and Gómez-Rodríguez, 2019)	MCQA	2657/1366/2742	Medical Exam
PROCESSBANK (Berant et al., 2014)	Context + MCQA	358/77/150	Biological Processes
PUBMEDQA (Jin et al., 2019)	Context + MCQA	400/100/500	Biomedical
MMLU (Hendrycks et al., 2021)	MCQA	30/NA/1089	Medical and Clinical
BIO-MRC-Tiny A (Pappas et al., 2020)	Context + MCQA	NA/NA/30	Cloze Biomedical
BIO-MRC-Tiny B (Pappas et al., 2020)	Context + MCQA	NA/NA/30	Cloze Biomedical
OPHTH (Raimondi et al., 2023; RCOphth, a,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	Biomedical
MASHQ (Zhu et al., 2020)	AQA	27728/3587/3493	Medical
MEDQUAD (Ben Abacha and Demner-Fushman, 2019)	AQA	14068/981/1358	Medical
MEDINFO (Ben Abacha et al., 2019)	AQA	NA/NA/663	Consumer Medication

Table 1: An overview of the 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.

anced trend.

In summary, QALM introduces a benchmark that consolidates datasets to assess clinical knowledge within LLMs. This initiative underscores the imperative for a unified and comprehensive suite of datasets to rigorously evaluate the clinical knowledge capabilities of these models, particularly in light of the continuously expanding array of LLMs.

2 QALM Datasets

QALM represent a comprehensive collection of 22 datasets designed to thoroughly evaluate the clinical knowledge of LLMs. These datasets are publicly accessible, and some of them have not been used for testing open-source LLMs before. Like other studies (Singhal et al., 2023a) (Singhal et al., 2023b), we employ Question and Answering tasks as a surrogate test to assess clinical knowledge of LLMs. We have included two types of question-answering datasets: Multiple Choice (MCQA) and Abstractive Question Answering (AQA). The complete list of these datasets is provided in Table 1.

MCQA questions assess the model’s ability to select the correct option from a list of challenging alternatives. These types of questions are frequently encountered in medical licensing examinations, such as the US Medical Licensing Exam (USMLE) (Jin et al., 2021), as well as in medical

entrance exams in countries like India (Pal et al., 2022) and Spain (Vilares and Gómez-Rodríguez, 2019). Within QALM, there are a total of 216,810 instances of MCQA questions. On average, each question has four answer choices. In some scenarios, a contextual paragraph is provided and the answer must be derived from this context which tests the model’s contextual reasoning abilities.

AQA datasets evaluate the model’s proficiency in providing open-ended answers to questions. QALM encompasses six AQA datasets, containing a total of 57,958 questions. On average, the answers in these datasets have a length of about 100 tokens.

Our work is distinct from Singhal et al. (2023a), who only consider two AQA datasets, LiveQA and MedicationQA, as they argue that these datasets lack reliable sources for answers. In contrast, we incorporate four additional datasets where answers are provided from experts or sourced from trusted forums, increasing their reliability. Furthermore, we use a version of the LiveQA dataset that contains expert-ranked answers (Ben Abacha and Demner-Fushman, 2019). The performance of LLMs on these datasets have not been previously assessed in existing literature.

While Singhal et al. (2023a) employ only a single dataset for assessing model performance in reading comprehension, we introduce two more

138 datasets: one for evaluating a model’s ability to
139 predict a concealed medical entity, and another
140 for assessing model proficiency in comprehending
141 interactions between various entities mentioned
142 in the text. Furthermore, we include four addi-
143 tional general knowledge multiple-choice question
144 (MCQ) datasets, broadening the range of questions
145 for model testing. This expanded array of evalua-
146 tion datasets enables a more comprehensive anal-
147 ysis of model capabilities, potentially mitigating
148 uncertainty stemming from a limited range of ob-
149 servations.

150 3 Empirical Evaluation

151 Considering the QALM datasets, we seek evidence
152 for the following key research questions in a large-
153 scale empirical study:

154 **RQ1.** How well do open-source language models
155 (LLMs) recall necessary clinical knowledge
156 when they are tested on QALM?

157 **RQ2.** Does instruction fine-tuning of LLMs im-
158 prove their clinical knowledge recall?

159 **RQ3.** Does domain and task-specific finetuning on
160 QALM help LLMs acquire additional clinical
161 knowledge?

162 **RQ4.** Can LLMs generalize to data unseen during
163 training?

164 3.1 Study Setup

165 To seek evidence for **RQ1** and **RQ2** empirically,
166 we evaluate several LLMs and their instruction-
167 finetuned versions on the test splits of QALM in
168 zero-shot manner. To answer **RQ3** and **RQ4**, we
169 fine-tune LLMs on the training portion of QALM
170 and evaluate on test splits of datasets both seen and
171 unseen during training. We complement our evalua-
172 tion with additional automated and manual error
173 analyses to identify causes for model successes and
174 failures.

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

184 Flan-T5 (3B and 11B) (Wei et al., 2021). Mod-
185 els with *Instruct* or *Chat* appended to their names
186 are instruction fine-tuned (Ouyang et al., 2022)
187 versions of their base models. The details of the
188 models are given in Table 2.² We initially also
189 considered OpenLLaMA (3B and 7B) (Geng and Liu,
190 2023) as well as GPT-J (6B) (Wang and Komat-
191 suzaki, 2021). However, we found that due to
192 the poor performance of OpenLLaMA on the MCQA
193 datasets and since GPT-J and OpenLLaMA did not
194 have instruction-finetuned equivalents, we did not
195 include them in our further analysis.

196 To address **RQ3**, we use the training set of the
197 QALM datasets. When official validation splits are
198 unavailable, we employ a random split of up to
199 around 20% of the data for validation purposes. If
200 no training datasets are available, we do not use
201 this dataset for fine-tuning and only consider the
202 test split of the respective datasets to answer **RQ4**.

203 **Finetuning and hyperparameters:** Since the
204 number of parameters for most of our models are in
205 the billions, we follow a more accepted practice of
206 using parameter-efficient fine-tuning. Specifically,
207 we use QLora and 4-bit quantization (Dettmers
208 et al., 2023) for fine-tuning. We use A100-40G
209 GPUs for all our experiments. The other hyper-
210 parameters used to train our models are reported in
211 the Appendix (Table 7).

212 **Evaluation measures:** We use accuracy to mea-
213 sure the performance of the model on MCQA
214 datasets; for AQA datasets, we use ROUGE-L (Lin,
215 2004), BERTScore (Zhang et al., 2020)³ and
216 METEOR (Banerjee and Lavie, 2005). METEOR
217 in particular is found to correlate better with hu-
218 man judgments than other metrics on AQA (Chen
219 et al., 2019).

220 3.2 Results and Analysis

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

223 3.2.1 Zero-shot performance of LLMs

224 Table 3 presents the zero-shot evaluation of lan-
225 guage models as evidence towards **RQ1**. We report
226 the average accuracies across all tasks.

²For MCQA evaluation in the zero-shot setting (where mod-
els are not explicitly fine-tuned on the train splits of our bench-
mark), we use a 1-shot prompt—giving an example to the
model, and find that it adheres better to the MCQA format. We
use the standard 5-shot prompt for MMLU.

³We use deberta-xlarge-mnli for calculating BERTScore.

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	All of MPT and Databricks Dolly-15k (Conover et al., 2023), Anthropic Helpful and Harmless (Bai et al., 2022)
Falcon-Instruct	Decoder	1.5T	All of Falcon and baize (Xu et al., 2023), GPT4All, GPTeacher ¹
LLama 2-Chat	Decoder	2T	All of LLama 2 and Flan Collection (Wei et al., 2021) + Private Data

Table 2: 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*).

	MCQA	AQA			
		Acc	RL	BS	MTR
<i>Base</i>	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
<i>Instruction tuned</i>	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

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

LLMs are good zero/few-shot learners: Table 3 shows that LLMs exhibit strong zero-shot capability on MCQA and AQA datasets, corroborating the findings of Singhal et al. (2023a). Across different LLMs of the same size, LLama 2 (7B) provides the best performance compared to other open source models like Falcon (7B). A possible reason for this might be the diversity in pretraining data – LLama 2 is trained on two trillion tokens, which is double the amount of Falcon. Another plausible reason is the mixture of datasets used for pretraining them.

Bigger LLMs have more clinical knowledge: Figure 1 shows the relationship between the number of parameters and performance. It suggests that model performance tends to improve with scale (correlation between scale and performance: Spearman’s $\rho = 0.24, p < 0.01$), and that language models exhibit emergent abilities with scale. For example, LLama 2 (13B) has 23% (relative) better accuracy than LLama 2 (7B) on MCQA and 6% higher METEOR scores. Even without further domain-specific adaptation of LLMs on clinical data, scale appears to play a major role in the amount of clinical

knowledge available to LLMs.

Open-source LLMs do not outperform humans: While the passing score for USMLE is 60% for humans⁴, we observe the best score for USMLE to be 43% (LLama 2), which is 17% short of that requirement. Meanwhile, Singhal et al. (2023a) report scores of 86.5% on USMLE. Similarly, for the PubmedQA dataset, human performance is 78% (Jin et al., 2019), compared to 60.4% of LLama 2.

To summarize our findings related to **RQ1**: While LLMs have good clinical knowledge, there is still a significant gap compared to humans (Singhal et al., 2023a).

3.2.2 Effect of Instruction Fine-tuning

Instruction fine-tuning (Wei et al., 2022; Ouyang et al., 2022) was proposed as a means to improve across-the-board performance on various open-domain tasks. To address **RQ2**, we investigate whether these improvements also apply to the clinical domain of QALM. The results are reported in the bottom part of Table 3.

Instruction fine-tuning is key: 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 2). Among them, Flan-T5 models show the best zero shot performance on MCQA, outperforming all comparable decoder-only models.

We observe that the decoder-only language models generate longer and more verbose answers compared to the encoder-decoder models, probably due to their causal language modelling training objective. As such, they perform better on AQA. Instruction fine-tuning models can consistently improve decoder-only models as well (compare LLama 2 with LLama 2-chat for example).

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

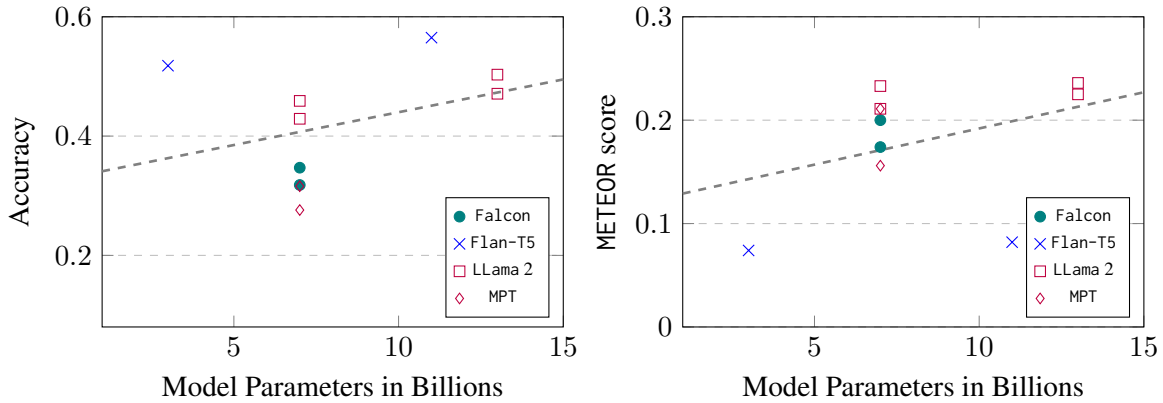


Figure 1: Zero-shot performance of models on MCQA (left) and AQA (right) as a function of model size. The dashed line represents a fitted linear regression showing the correlation between the model size and the score.

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), despite being four times its size.

	MCQA		AQA		
	Acc	RL	BS	MTR	
LLama 2 (7B)	53.5 _{10.6}	17.7 _{2.8}	60.8 _{5.5}	16.9 _{4.2}	
Falcon (7B)	49.3 _{14.6}	17.4 _{0.2}	60.4 _{2.0}	17.1 _{0.3}	
MPT (7B)	53.2 _{25.6}	17.3 _{4.0}	60.0 _{7.4}	17.2 _{3.9}	
Flan-T5 (3B)	52.9 _{1.1}	15.9 _{5.1}	56.8 _{1.8}	15.6 _{7.4}	

Table 4: Caption: Model finetuning is performed either on MCQA or their 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.

3.2.3 Impact of Finetuning

Given the scale of QALM, we are able to fine-tune models on parts of the data, to address **RQ3**. Specifically, we fine-tune four models on MCQA and AQA separately, given the different nature of these datasets.⁵

MCQA fine-tuning improves knowledge: We fine-tune the models only on the MCQA subset of datasets first (c.f. Table 4). We find that the models perform better compared to their non-fine-tuned counterparts. Decoder-only models like MPT (7B) benefit more than others (+25.6 percentage

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

points 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 zero-shot setting is 9.0, while after fine-tuning it is reduced to 1.7. This suggests that various LLMs can benefit from task-specific fine-tuning to address seemingly sub-optimal architecture or pre-training conditions.

AQA fine-tuning can improve knowledge: We fine-tune the models on AQA datasets. Encoder-Decoder models like Flan-T5 benefit more from AQA fine-tuning compared to the decoder-only models. For instance, the fine-tuned Flan-T5 model has significantly better METEOR (+7.8) compared to the non fine-tuned version.

Fine-tuning can compensate for scale: Scaling up brings 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. Fine-tuning helps compensate for scale. LLama 2 (7B) significantly outperforms the zero-shot LLama 2 (13B) (+6.4 Accuracy gain on MCQA, +5.7 METEOR gain on AQA). Similarly, we observe that a fine-tuned Flan-T5 (3B) outperforms zero shot Flan-T5 (11B) on 9 out of 16 MCQA datasets and also on AQA tasks considerably.

Fine-tuned models make similar mistakes: In order to provide targeted suggestions for future performance improvements, we analyze the types of errors that fine-tuned models make on the MCQ datasets. To do this, we perform entity recognition using SciSpacy (Neumann et al., 2019) over

Question type and template	Support	LLama 2 (7B)	Falcon (7B)	MPT (7B)	Flan-T5 (3B)
Scenario-Based Treatment (A {age}-year-old person. . .)	1172	48.37%	45.31%	44.8%	40.45%
Fact-Based Discriminative (Which of the. . .)	1053	54.9%	47.95%	53.18%	45.68%

Table 5: Analysis of the most common question types and finetuned model performance. As these questions form a sub-set of the test set, averaging the reported scores does not correspond to the performance on the full test set reported in Table 4.

the questions and link these entities with UMLS to obtain their type unique identifier (TUI). We map these TUIs to their semantic groups⁶. We replace the original entities in the question with their respective semantic group.

We first use an automated way of analyzing errors. We categorize the questions based on the first three words of the questions, similar to (Yang et al., 2018). The example of categories are shown in Table 5, where *Scenario-Based Treatment Questions* describe a real-world case, and the model needs to choose the right course of treatment for the case. The second category *Fact-Based Questions* forces the model to rely on the internal knowledge stored in its parameters to answer the questions.

We find that models perform worse than average on *Scenario-Based Treatment Questions*: the performance drops when they have to combine knowledge encoded in the internal parameters with particulars described in the scenario. Among different models, Flan-T5 and Falcon have a more stable performance across question categories, compared to LLama 2 and MPT (standard deviations of 2.6 and 1.9, vs 3.4 and 3.9).

We further analyze the errors made by the models manually, by randomly choosing 50 samples where the model was correct and 50 samples where the model was wrong. We consider the best performing LLama 2 (7B) model for our analysis. We identify five categories of common errors and report them in Table 6, where *Scenario-Based Treatment Questions* and *Fact-Based Questions* correspond to the previously described common question types. Meanwhile, questions concerning *Procedures, Tests and Activities* require knowledge of medical procedures (e.g., the type of test needed to identify a certain condition), *Anatomy and Physiology* related questions require knowledge of anatomy, and *Chemicals and Drugs* require pharmaceutical knowledge.

⁶<https://lhncbc.nlm.nih.gov/ii/tools/MetaMap/documentation/SemanticTypesAndGroups.html>

Category	% Wrong	% Correct
Scenario-Based Treatment	26%	30%
Procedures, Tests and Activities	14%	6%
Fact-based Discriminative	12%	14%
Anatomy and Physiology	10%	12%
Chemicals and Drugs	10%	10%

Table 6: Manual analysis of errors and correct predictions of the best-performing fine-tuned LLama 2 (7B) model.

The manual error analysis corroborates the previous finding that fine-tuned models tend to err when asked about treatments given a scenario, possibly because this is the most frequently occurring category in the training data. Additionally, our findings suggest that models struggle with procedure-based questions, as these can exhibit significant complexity. For example, to successfully select the correct option for the question “Recent studies have shown that the chlorhexidine-isopropyl alcohol solution substantially reduces the risk of surgical site infections compared with a povidone-iodine preparation without alcohol in clean-contaminated surgery. Which of the following mechanisms best describes the mechanism of action of chlorhexidin?”, a model needs to not only possess knowledge about the mechanism of action of chlorhexidin, but also relate it to research findings which suggest that it reduces the risk of surgical site infections (i.e. that it “[...] is greatly reduced in the presence of organic matter” in this case).

3.2.4 Generalisation to Unseen Data

Finally, we report the potential of LLMs fine-tuned on in-domain data to generalise to medical datasets unseen during training to answer **RQ4**. We hold out 4 AQA and 10 MCQA datasets presented in Figures 2 and 3.

AQA-finetuned models generalise to unseen AQA test sets: Figure 2 shows the performance of LLama 2 (7B) and Flan-T5 (3B) models on the four held-out AQA evaluation sets. The METEOR

scores of the fine-tuned LLama 2 model are lower than the zero-shot baseline (-4.2 METEOR). Meanwhile, fine-tuning Flan-T5 improves performance on all four unseen datasets ($+8.2$ METEOR), suggesting that instruction fine-tuned models have good generalisation capabilities (Wei et al., 2022) compared to models that were only pre-trained on the language modelling task.

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.

MCQA-finetuned models do 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 QALM: Clinical

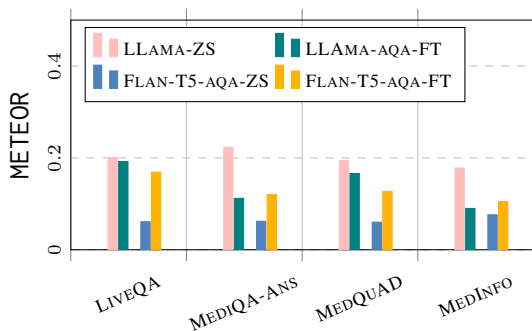


Figure 2: Performance of base and AQA-finetuned LLama 2 and Flan-T5 models on four unseen AQA test sets.

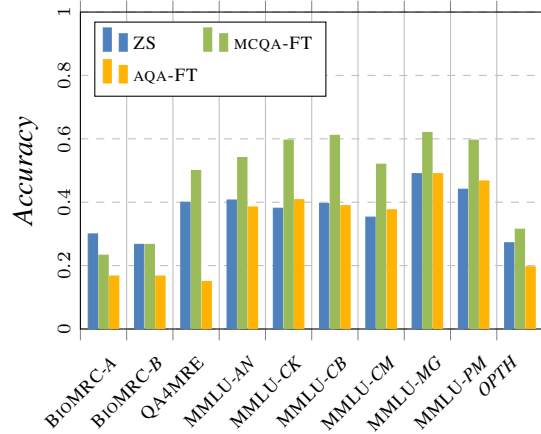


Figure 3: Performance of base, MCQA-fine-tuned and AQA-finetuned LLama 2 model on ten unseen MCQA test sets.

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 filter the retrieved results. We identify six out of 92, twelve 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 five, two, and five 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, seven questions were near-duplicates with identical answers, while the remaining five would have required some level of clinical understanding for the model to answer them correctly. These findings suggest that the improved performance of instruction-tuned models on unseen datasets can be partially attributed to exposure to near-identical questions during training.

4 Related Work

Open source LLMs: The recent proliferation of LLMs has sparked a surge of interest in their adap-

tation for clinical applications. Notably, closed-source models like Med-PaLM and Med-PaLM2 (Singhal et al., 2023a,b) have shown considerable promise in leveraging LLMs within the clinical context. Meanwhile, open-source LLMs are being released at an astonishing pace. Nevertheless, it is crucial to note that these models have yet to be rigorously assessed for their clinical knowledge, and there is currently no existing benchmark for this purpose, a gap that this study addresses.

Benchmarks for clinical knowledge: The adaptation of Large Language Models (LLMs) to the medical domain is gaining significant traction as they exhibit substantial potential in addressing clinical and medical issues. In their groundbreaking work, Singhal et al. (2023a) introduced Multi-MedQA, a benchmark that consolidates six distinct medical question-answering datasets, along with HealthSearchQA, a novel free-response dataset designed for medical inquiries. Their research involved an evaluation of Flan-PaLM (Chowdhery et al., 2022; Chung et al., 2022) on MultiMedQA, alongside the proposal of instruction prompt tuning to enhance medical reasoning and customize LLMs for the medical domain, ultimately giving rise to Med-PaLM. Med-PaLM 2 (Singhal et al., 2023b) extended this groundwork by incorporating PaLM2 (Anil et al., 2023) and leveraging medical domain fine-tuning, along with prompting strategies. These enhancements resulted in an impressive 86.5% accuracy rate on the United States Medical Licensing Examination (USMLE). Building upon this foundational work, our aim is to introduce a more comprehensive benchmark, encompassing a broader range of datasets with diverse characteristics, to thoroughly evaluate LLMs.

While most of these benchmarks assess LLMs’ capabilities in question-answering tasks, it’s crucial to address the aspect of hallucination, which is of paramount importance in clinical and medical applications. In response to this concern, Umaphathi et al. (2023) introduced Med-HALT, a benchmark designed to evaluate LLMs on hallucinations.

While some of these benchmarks evaluate LLMs’ clinical knowledge, they may fall short in assessing their competence in performing various tasks within real-world clinical settings. A recent development in this regard comes from Fleming et al. (2023), who introduced Med-Align. This innovative benchmark dataset comprises natural language directives related to electronic health record (EHR)

data, which is a longitudinal record of a patient’s medical history. Med-Align serves the vital purpose of evaluating the efficacy of LLMs in processing clinical instructions within a clinical context.

5 Conclusion

In this work, we introduce 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 knowledge, their capacity to acquire such knowledge through training on QALM, and their ability to generalize to previously unseen datasets.

Our findings reveal that while these LLMs exhibit performance significantly superior to random guessing, there remains room for enhancing their performance when compared to closed-source language models. Notably, fine-tuning on 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 acknowledge that some performance improvements may be attributed to the nature of the test questions.

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. These aspects are vital for advancing the state-of-the-art in clinical NLP and expanding the horizon.

We make the dataset, experiment code and evaluation protocol publicly available under <https://anonymized>. This will allow future LLM developers to perform a fine-grained analysis of their models clinical and biomedical knowledge.

575 Limitations

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

586 Furthermore, the aggregated resource presented
587 in this paper might be seen as lacking diversity,
588 as all collected datasets are in English. To make
589 inferences about the capabilities of evaluated mod-
590 els in other languages, a more diverse dataset with
591 examples in other languages is required.

592 For our finetuning experiments, we only use
593 parameter-efficient finetuning methods (PEFT)
594 with QLoRA due to the high compute requirements
595 for full-finetuning. We have not investigate the im-
596 pact of the full-finetuning of these LLMs on our
597 benchmark.

598 References

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

603 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-
604 shamsi, Alessandro Cappelli, Ruxandra Cojocaru,
605 Merouane Debbah, Etienne Goffinet, Daniel Hes-
606 low, Julien Launay, Quentin Malartic, Badreddine
607 Noune, Baptiste Pannier, and Guilherme Penedo.
608 2023. Falcon-40B: an open large language model
609 with state-of-the-art performance.

610 Rohan Anil et al. 2023. Palm 2 technical report. *arXiv*,
611 2305.10403.

612 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda
613 Askell, Anna Chen, Nova DasSarma, Dawn Drain,
614 Stanislav Fort, Deep Ganguli, Tom Henighan,
615 Nicholas Joseph, Saurav Kadavath, Jackson Kernion,
616 Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac
617 Hatfield-Dodds, Danny Hernandez, Tristan Hume,
618 Scott Johnston, Shauna Kravec, Liane Lovitt, Neel
619 Nanda, Catherine Olsson, Dario Amodei, Tom
620 Brown, Jack Clark, Sam McCandlish, Chris Olah,
621 Ben Mann, and Jared Kaplan. 2022. [Training a help-
622 ful and harmless assistant with reinforcement learn-
623 ing from human feedback](#).

624 Satanjeev Banerjee and Alon Lavie. 2005. [METEOR:
625 An automatic metric for MT evaluation with im-
626 proved correlation with human judgments](#). In *Proc.*

*ACL Workshop on Intrinsic and Extrinsic Evaluation
Measures for Machine Translation and/or Summa-
rization*, pages 65–72, Ann Arbor, Michigan. Associ-
ation 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 con-
sumers’ 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 summa-
rization & generation of doctor-patient conversations](#).
In *Proc. 5th Clinical Natural Language Processing
Workshop*, pages 503–513, Toronto, Canada. Associ-
ation for Computational Linguistics.

Asma Ben Abacha, Wen-wai Yim, Yadan Fan, and
Thomas Lin. 2023b. [An empirical study of clini-
cal note generation from doctor-patient encounters](#).
In *Proc. 17th Conference of the European Chap-
ter 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 compre-
hension](#). In *Conference on Empirical Methods in
Natural Language Processing*.

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 answer-
ing evaluation](#). In *Proc. 2nd Workshop on Machine
Reading for Question Answering*, pages 119–124,
Hong Kong, China. Association for Computational
Linguistics.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
Maarten Bosma, Gaurav Mishra, Adam Roberts,
Paul Barham, Hyung Won Chung, Charles Sutton,
Sebastian Gehrmann, Parker Schuh, Kensen Shi,
Sasha Tsvyashchenko, Joshua Maynez, Abhishek
Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vin-
odkumar Prabhakaran, Emily Reif, Nan Du, Ben
Hutchinson, Reiner Pope, James Bradbury, Jacob
Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,
Toju Duke, Anselm Levskaya, Sanjay Ghemawat,
Sunipa Dev, Henryk Michalewski, Xavier Garcia,
Vedant Misra, Kevin Robinson, Liam Fedus, Denny
Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim,
Barret Zoph, Alexander Spiridonov, Ryan Sepassi,

627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682

683	David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. <i>arXiv</i> , 2204.02311.	
684		
685		
686		
687		
688		
689		
690		
691		
692	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. <i>arXiv</i> , 2210.11416.	
704	Together Computer. 2023. Redpajama-data: An open source recipe to reproduce llama training dataset .	
705		
706	Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world's first truly open instruction-tuned llm .	
707		
708		
709		
710		
711	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. QLoRA: Efficient Finetuning of Quantized LLMs. <i>arXiv</i> , 2305.14314.	
712		
713		
714	Scott L. Fleming, Alejandro Lozano, William J. Haberkorn, Jenelle A. Jindal, Eduardo P. Reis, Rahul Thapa, Louis Blankemeier, Julian Z. Genkins, Ethan Steinberg, Ashwin Nayak, Birju S. Patel, Chia-Chun Chiang, Alison Callahan, Zepeng Huo, Sergios Gatidis, Scott J. Adams, Oluseyi Fayanju, Shreya J. Shah, Thomas Savage, Ethan Goh, Akshay S. Chaudhari, Nima Aghaeepour, Christopher Sharp, Michael A. Pfeffer, Percy Liang, Jonathan H. Chen, Keith E. Morse, Emma P. Brunskill, Jason A. Fries, and Nigam H. Shah. 2023. Medalign: A clinician-generated dataset for instruction following with electronic medical records. <i>arXiv</i> , 2308.14089.	
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727	Xinyang Geng and Hao Liu. 2023. Openllama: An open reproduction of llama .	
728		
729	Himanshu Gupta, Saurabh Arjun Sawant, Swaroop Mishra, Mutsumi Nakamura, Arindam Mitra, Santosh Mashetty, and Chitta Baral. 2023. Instruction tuned models are quick learners .	
730		
731		
732		
733	Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bressen. 2023. Medalpaca—an open-source collection of medical conversational ai models and training data. <i>arXiv</i> , 2304.08247.	
734		
735		
736		
737		
738		
	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. <i>arXiv</i> , 2009.03300.	739 740 741 742
	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams . <i>Applied Sciences</i> , 11(14).	743 744 745 746 747
	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. PubMedQA: A dataset for biomedical research question answering . In <i>Proc. Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2567–2577, Hong Kong, China. Association for Computational Linguistics.	748 749 750 751 752 753 754 755
	Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von Werra, and Harm de Vries. 2022. The stack: 3 tb of permissively licensed source code. <i>Preprint</i> .	756 757 758 759 760 761
	Anastasia Krithara, Anastasios Nentidis, Konstantinos Bougiatiotis, and Georgios Paliouras. 2023. Bioasqqa: A manually curated corpus for biomedical question answering . <i>Scientific Data</i> , 10(1):170.	762 763 764 765
	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge . <i>Cureus</i> , 15(6).	766 767 768 769 770
	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	771 772 773 774
	Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The semantic scholar open research corpus . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4969–4983, Online. Association for Computational Linguistics.	775 776 777 778 779 780
	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.	781 782 783 784 785 786
	NLP Team MosaicML. 2023. Introducing mpt-7b: A new standard for open-source, commercially usable llms . Accessed: 2023-05-05.	787 788 789
	Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. 2019. ScispaCy: Fast and robust models for biomedical natural language processing . In <i>Proc. 18th BioNLP Workshop and Shared Task</i> , pages 319–327, Florence, Italy. Association for Computational Linguistics.	790 791 792 793 794 795

796	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. Training language models to follow instructions with human feedback .	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.	849
797			850
798			851
799			852
800			853
801			
802		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. <i>BMC Bioinformatics</i> , 16:138.	854
803			855
804	Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering . In <i>Proc. Conference on Health, Inference, and Learning</i> , volume 174 of <i>Proceedings of Machine Learning Research</i> , pages 248–260.		856
805			857
806			858
807			859
808			860
809			861
810	Dimitris Pappas, Petros Stavropoulos, Ion Androutsopoulos, and Ryan McDonald. 2020. BioMRC: A dataset for biomedical machine reading comprehension . In <i>Proc. 19th SIGBioMed Workshop on Biomedical Language Processing</i> , pages 140–149, Online. Association for Computational Linguistics.		862
811			863
812			864
813			865
814		Logesh Kumar Umapathi, Ankit Pal, and Malaikannan Sankarasubbu. 2023. Med-halt: Medical domain hallucination test for large language models. <i>arXiv</i> , 2307.15343.	866
815			867
816			868
817	Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only . <i>arXiv preprint arXiv:2306.01116</i> .		869
818			870
819			871
820			872
821			873
822			874
823			875
824	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>arXiv e-prints</i> .		876
825			877
826			
827			
828	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.		878
829			879
830			880
831			881
832			882
833	RCOphth. a. Frcophth sample mcqs part 1 .		883
834	RCOphth. b. Frcophth sample mcqs part 2 .		884
835			885
836	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.		886
837			887
838			
839	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.		888
840			889
841			890
842			891
843			892
844	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.		893
845			894
846			895
847			896
848			
		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 <i>International Conference on Learning Representations</i> .	897
			898
			899
			900
		Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners .	901
			902
		Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Towards building open-source language models for medicine. <i>arXiv</i> , 2304.14454.	903
		Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023. Baize: An open-source chat model with parameter-efficient tuning on self-chat data .	

- 904 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale,
905 Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and
906 Colin Raffel. 2021. [mT5: A massively multilingual
907 pre-trained text-to-text transformer](#). In *Proceedings
908 of the 2021 Conference of the North American Chapter
909 of the Association for Computational Linguistics: Human
910 Language Technologies*, pages 483–498, On-
911 line. Association for Computational Linguistics.
- 912 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,
913 William Cohen, Ruslan Salakhutdinov, and Christo-
914 pher D. Manning. 2018. [HotpotQA: A dataset for
915 diverse, explainable multi-hop question answering](#).
916 In *Proceedings of the 2018 Conference on Empirical
917 Methods in Natural Language Processing*, pages
918 2369–2380, Brussels, Belgium. Association for Com-
919 putational Linguistics.
- 920 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
921 Weinberger, and Yoav Artzi. 2020. [Bertscore: Eval-
922 uating text generation with bert](#). In *International
923 Conference on Learning Representations*.
- 924 Ming Zhu, Aman Ahuja, Da-Cheng Juan, Wei Wei,
925 and Chandan K. Reddy. 2020. [Question answering
926 with long multiple-span answers](#). In *Findings of the
927 Association for Computational Linguistics: EMNLP*,
928 pages 3840–3849. Association for Computational
929 Linguistics.

Appendix

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 7: 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 8: Inference time parameters used for abstractive question answering