UNCERTAINTY QUANTIFICATION IN RETRIEVAL AUG MENTED QUESTION ANSWERING

Anonymous authors

Paper under double-blind review

ABSTRACT

Retrieval augmented Question Answering (QA) enables QA models to overcome knowledge gaps when answering questions at test time by taking as input the question together with retrieved evidence, that is usually a set of passages. Previous studies show that this approach has numerous benefits such as improving QA performance and reducing hallucinations, without, however, qualifying whether the retrieved passages are indeed useful at answering correctly. In this work, we evaluate existing uncertainty quantification methods and propose an approach that predicts answer correctness based on utility judgements on individual input passages. We train a small neural model that predicts passage utility for a target QA model. We find that simple information theoretic metrics can predict answer correctness up to a certain extent, more expensive sampling based approaches perform better, while our lightweight approach can efficiently approximate or improve upon sampling-based approaches.

023 024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 Retrieval augmented Question Answering (QA) enables QA models to overcome knowledge gaps when answering user questions at test time by giving them access to input evidence, i.e., a set of 029 passages, retrieved for the user questions (Lewis et al., 2020; Guu et al., 2020; Izacard et al., 2024). Recent work exploits the language understanding and generation abilities of Large Language Models 031 (LLMs; (Brown et al., 2020; Ouyang et al., 2024)) and makes use of external retrievers to find the 032 input evidence (Chen et al., 2017; Izacard & Grave, 2021a). That is, the retrieved evidence is given 033 to the LLM-based QA model as input context in tandem with the question; the QA model will read 034 this evidence and formulate an answer. For instance, in Figure 1, for the user question Who sings Does He Love Me with Reba?, the QA model is provided with a set of evidence passages together with the question; and correctly formulates the answer Linda Davis. 037

Such retrieval augmented QA architectures have proven beneficial enabling access to external knowledge (Izacard et al., 2024), increasing the performance on tail knowledge (Mallen et al., 2023), reducing hallucinations in model answers, and even improving model calibration (Jiang et al., 2021). 040 However, there are various ways in which a retrieval augmented QA approach can go wrong at pro-041 duction time. The set of passages obtained using retrieval methods is far from perfect (Sciavolino 042 et al., 2021; Yoran et al., 2024; Kasai et al., 2024) containing irrelevant or misleading evidence, the 043 model might be under-trained to read certain passages and reason over these and the question (Izac-044 ard et al., 2024; Liu et al., 2024b), or the question can simply be ambiguous or unanswerable (Kasai et al., 2024). In these cases where the QA system lacks the knowledge to formulate an answer (i.e., it is uncertain about what the answer is), we want it to refrain from answering rather than providing 046 an erroneous answer. Thus, predicting answer uncertainty is key. 047

Approaches to answer uncertainty prediction can be grouped in two main categories, sampling and LLM-based methods. Sampling-based methods to QA uncertainty detection rely on the output
 discrepancies among multiple predictors on the same input (Gal & Ghahramani, 2016; Lakshmi narayanan et al., 2017); i.e., this variance in outputs indicates that the model is uncertain. Concretely,
 these methods sample via temperature scaling (Guo et al., 2017) and then measure diversity on the
 set of sampled answers (Kuhn et al., 2023; Chen & Mueller, 2024). These approaches are expensive
 to run for in-production QA systems and the quality of the semantic similarity will degrade on long



Figure 1: Example of user question from the Natural Questions dataset with the set of three top 067 retrieved passages with Contriever (Izacard et al., 2022) (the other two passages below the rank are 068 less relevant and not shown in the figure); the gold answer is *Linda Davis*. The target QA model 069 GEMMA2-9B correctly answers the question when provided with the top five passages. Below each passage, it is shown the answer generated by the QA model when only prompted with that 071 passage and the question. The QA model correctly answers when prompted with the first passage and produces an incorrect answer when prompted with each of the other ones. The yellow triangles 073 on the top right of the passages are the predicted utility scores by our utility ranker. Higher values 074 indicate more useful passages and our model correctly identifies that the top passage is better. 075

077 answers (Zhang et al., 2024).¹ LLM-based methods explore to what extent language models are 078 able to correctly express uncertainty about their own predictions (Kadavath et al., 2022; Lin et al., 079 2022; Tian et al., 2023; Zhou et al., 2023). These look into whether the model's confidence in its outputs coincides with their correctness (i.e., calibration), methods to fix calibration, and ways to elicit from the model a verbal expression of that confidence (i.e., linguistic calibration). Findings 081 about model calibration are diverse and model dependent, fixing relies on approximations for the case of black-box models and fine-tuning what could be infeasible in practice given current LLMs' 083 sizes. None of these answer uncertainty detection approaches has been applied in the context of 084 retrieval augmented QA, most of them are applied on closed-book QA tasks where the answer is 085 predicted based on the question and the models' encoded knowledge.

In this work, we propose a secondary model that makes predictions at individual retrieved pas-087 sage level that are useful to estimate answer uncertainty of retrieval augmented QA models. We hypothesize that the type of retrieved passages and questions, the relation between them and their implicit interaction with the QA model's own knowledge are indicators of answer correctness. 090 If the passages are informative and priming the QA model towards appropriate knowledge, we ex-091 pect the QA model to produce a correct answer. In contrast, if the passages are not informative or 092 misleading and the posed question is out of the QA model's knowledge, we expect it to generate an erroneous answer (i.e., either factually incorrect or completely made up content). We operationalise 094 this as retrieved passage utility. Given a question, a passage is useful, if a QA model can correctly 095 answer the question based on it. We train a small neural model to predict passage utility which we 096 refer to as *utility ranker*. We train the utility ranker on utility judgements generated by the target QA model. We borrow ideas from direct uncertainty quantification approaches (Van Amersfoort et al., 2020; Lahlou et al., 2023) but we do not decompose uncertainty or outline shifts in the input 098 distribution.

We show that individual passage utilities are good predictors of retrieval augmented QA accuracy. This means that it is possible to train an answer uncertainty predictor independently from the choice of number of retrieved passages used to prompt the target QA model. Because retrieved passages are scored individually, our approach is independent of the *number* of retrieved passages chosen for the target QA model. We evaluate our approach on short-form question answering tasks. Figure 1 shows an example of input question, set of retrieved passages, and correct answer

¹By expensive we mean both latency as well as cost as a long prompts might need to be processed and QA systems may rely on paid proprietary language models.

108 from the Natural Questions dataset (Kwiatkowski et al., 2019). Results on six QA datasets show that 109 our approach performs on par with existing sampling-based uncertainty quantification approaches 110 while being more efficient at test time. It requires a small model pass over the set of input passages 111 and the question (see inference cost comparison in Appendix C.1). Surprisingly, in more complex 112 reasoning questions (SQuAD) and adversarial QA settings (e.g., rare entities or unanswerable questions) our approach surpasses existing uncertainty quantification methods. Moreover, we show that 113 the utility scores predicted by the Utility Ranker can be used to re-rank retrieved passages obtained 114 via the external retrieval system to improve QA accuracy (Liu et al., 2024b). 115

116

131

117 2 RELATED WORK

119 **Uncertainty Quantification for Question Answering** Several methods have been proposed to 120 predict answer uncertainty in QA; however, none of them has analysed uncertainty in retrieval aug-121 mented QA models. Many existing approaches rely on capturing output variation as the expression of model uncertainty (Kuhn et al., 2023; Farguhar et al., 2024; Chen & Mueller, 2024). On a sam-122 ple of model outputs, Kuhn et al. (2023) propose to first cluster answers with similar meaning via 123 natural language inference before computing entropy. Chen & Mueller (2024) propose an approach 124 for *black-box* models, they also compute similarities in the set of answers but associate them with 125 a model self-judgement of confidence. These approaches are expensive to run at inference time 126 for a production QA system, they require several inference steps plus the similarity computations. 127 In addition, as the length of the answers increases, measuring similarity becomes more complex 128 (Zhang et al., 2024). Hou et al. (2024) propose a decomposition of predictive uncertainty and focus 129 on quantifying aleatoric uncertainty (i.e., uncertainty in the data) caused by ambiguous questions. 130 This approach is orthogonal to ours.

132 Judging the Utility of Retrieved Passages Previous work has analysed the set of retrieved pas-133 sages (Yu et al., 2023; Asai et al., 2024; Wang et al., 2024; Xu et al., 2024; Yoran et al., 2024) following the observation that passages can be irrelevant or misleading making the QA model prone 134 to producing incorrect answers. Asai et al. (2024) make use of an external critic model to judge 135 whether a question requires retrieval (or not), whether the retrieved passages are relevant to formu-136 late the answer, and whether the final response elaborated by the OA model is useful. While they 137 analyse retrieved passage *relevance*, this decision is taken by an external extreme-scale critic (e.g., 138 GPT-4) and used to fine-tune the QA model. In contrast, we do not fine-tune the target QA model 139 but rather we elicit utility judgements from it to train a secondary model to predict passage utility. 140 Other work creates auxiliary tasks around retrieved passages enforcing the QA model to reason on 141 them; e.g., by taking notes about each passage (Yu et al., 2023) or generating passage summaries 142 (Xu et al., 2024). These methods also use extreme-scale LLMs to generate training data to fine-143 tune the retrieval augmented QA model. Park et al. (2024) select specific in-context examples to 144 improve the LLM's reasoning on the input passages, their focus is on detecting input passages with 145 conflicting content (e.g., different dates for a given event). These approaches aim at improving QA performance while our primary goal is modelling QA uncertainty. 146

147

Improving Retriever via Reader Performance Previous work with pre-trained language models 148 has focused on jointly training the retriever and reader modules end-to-end (Lee et al., 2019; Lewis 149 et al., 2020; Izacard & Grave, 2021b). That is, the performance of the question answering model 150 is propagated also to the retriever. This joint training scheme can be very expensive for current 151 (extreme-scale) LLMs. Our approach can be seen as an intermediate module between the QA model 152 (reader) and the external retriever. It would be interesting to explore our utility ranker to provide 153 feedback (e.g., to label data) to fine-tuning the retriever. In recent work, Salemi & Zamani (2024) 154 evaluate the performance of information retrieval systems via retrieval augmented QA performance. 155 Interestingly, they show that external judgements (e.g., query-document relevance labels) of passage 156 utility correlate poorly with retrieval augmented QA performance.

157

Learning to Predict Confidence Some approaches train a specific model to predict a confidence score (Dong et al., 2018; Kamath et al., 2020; Mielke et al., 2022). For semantic parsing, Dong et al. (2018) train a confidence predictor based on a set of uncertainty features from the input and the model. Mielke et al. (2022) also train a calibrator that, given the user question and model generated answer, predicts a confidence score. In our approach, we simple aggregate predicted individual

passage utilities but it would also be possible to train a confidence module that takes utilities with other features into account (e.g., output sequence probability), and predicts a confidence score.

165 166

167

3 MODELLING ANSWER UNCERTAINTY

Formally, we define retrieval augmented QA as follows. Given question x and set of 168 retrieved passages $R = \{p_1, p_2, \cdots, p_{|R|}\}$ obtained with retriever \mathcal{R} , a LLM-based QA model $\mathcal M$ is prompted to generate answer $y_{\mathcal M}$ to question x token-by-token as $y_{\mathcal M}$ = 170 $\arg \max_{y_{\mathcal{M}}} \prod_{t=1}^{|y_{\mathcal{M}}|} p_{\mathcal{M}}(y_t|y_{1.t-1}, x, R)$. We want to estimate the uncertainty or error of \mathcal{M} on 171 generating $y_{\mathcal{M}}$ given x and R; i.e., we want an estimator $\{x, R\} \mapsto \mathbf{u}_{\mathcal{M}}(\{x, R\})$ of \mathcal{M} 's answer uncertainty. In our approach, the answer uncertainty predictor $\mathbf{u}_{\mathcal{M}}$ is based on individual passage 172 173 utilities. Our hypothesis is that individual passage utilities of retrieved passages in R are indicators 174 of the QA model uncertainty when generating $y_{\mathcal{M}}$ when prompted with R. For instance, in Fig-175 ure 1, given that the first passage in the set has a high utility score, this indicates that the QA model 176 is likely to be confident when providing the answer *Linda Davis*. Thus, we want a passage utility 177 estimator $\{x, p\} \mapsto v_{\mathcal{M}}(\{x, p\})$ of every $p \in R$. In what follows, we define passage utility and how 178 to estimate and predict it. Next, we discuss a simple answer uncertainty estimator $\mathbf{u}_{\mathcal{M}}$ based on $v_{\mathcal{M}}$.

179

181

3.1 PASSAGE UTILITY RANKING

Passage Utility Intuitively, a passage p retrieved for question x is useful for a QA model \mathcal{M} , if \mathcal{M} 182 can correctly answer x when prompted with p. In addition, \mathcal{M} 's reliance on passage p to formulate 183 the answer may vary. That is, the QA model may formulate a correct answer even though p does 184 not provide the answer itself; instead, p positively primes \mathcal{M} to use its memorised knowledge. The 185 utility of the first passage, in Figure 1, is high as the QA model generates a correct answer when prompted with it and the fact that *Linda Davis* sings together with *Reba McEntire* can be derived from 187 it. The second and third passages, although related to the topic of the question, are not useful. The 188 QA model is potentially uncertain about how to answer the question and the passages do not help; 189 the model incorrectly answers when prompted with each of them. Thus, the utility of the second and 190 third passages is low.

Concretely, we estimate the utility of passage p for QA model \mathcal{M} to answer question x by combining two measures. These are *accuracy*, denoted as $a(y_{\mathcal{M}})$, whether the generated answer $y_{\mathcal{M}}$ is correct, and *entailment*, denoted as $e(y_{\mathcal{M}})$, how much does passage p supports the generated answer $y_{\mathcal{M}}$. Accuracy is computed by a critic model \mathcal{A} and entailment by a Natural Language Inference (NLI) classifier model \mathcal{E} . We define the combined passage utility as $v_{\mathcal{M}} = (a(y_{\mathcal{M}}) + e(y_{\mathcal{M}}))/2$ that takes values in the closed interval [0, 1] given that a takes values in the set $\{0, 1\}$ and e in the closed interval [0, 1].

198

211 212

199 Utility Ranker We train a small neural model to predict passage utility scores, $\{x, p\} \mapsto v_{\mathcal{M}}(\{x, p\})$. We use observed answer accuracy and entailment by QA model \mathcal{M} on a training **201** set $D = \{(x, p)\}$ to train the utility predictor. That is, we run the QA model \mathcal{M} on examples from **202** D and compute passage utilities to form a training set for our utility predictor $D_{\mathcal{M}} = \{(x, p, v_{\mathcal{M}})\}$.

203 For recall purposes, retrieval augmented QA generally retrieves more than one input passage for 204 each question x, i.e., |R| > 1. To generate training data for the passage utility predictor, we retrieve 205 |R| passages per question in order for to cover passages with different usefulness. From the set of passages R for question x, we derive training instances $\{(x, p_i, v_{\mathcal{M}i}) \mid p_i \in R\}$. We exploit this 206 to train the passage utility predictor with a contrastive learning scheme. That is, if p_i and p_j are 207 passages in R and p_i is more useful than passage p_i to answer question x, the predicted utility score 208 $v_{\mathcal{M}i}$ should be higher by a margin m than the predicted score $v_{\mathcal{M}i}$ for p_i (i.e., p_i should be ranked 209 higher than p_i). We train the utility predictor with the following pair-wise ranking objective: 210

$$\mathcal{L}_{rank} = \sum_{((x,p_i),(x,p_j))\in R\times R, i\neq j} \max(0, -z(\upsilon_{\mathcal{M}i} - \upsilon_{\mathcal{M}j}) + m)), \tag{1}$$

where z controls the gold order between p_i and p_j (e.g., if z = 1, p_i has higher utility, conversely z = -1 indicates the opposite ordering) and m is a hyper-parameter. The passage utility predictor is trained with a Siamese neural network. Its architecture is constituted by a BERT (Devlin et al., 2019) based encoder followed by a pooling and two MLP layers stacked on top of BERT outputs (Fang et al., 2024). The output layer computes the utility score as $v_{\mathcal{M}i} = W_o h^L + b_o$ where h^L is the vector representation for (x, p_i) from the last hidden layer (the L-th layer) of the network. At inference time, we compute a single utility score for each passage. We provide implementation and training details in Section 4.

To enforce the signal on accuracy prediction and to regularise the range of utility values learned
 by the ranking scheme, we combine the ranking objective in Equation 1 with the following Binary
 Cross Entropy (BCE) objective (Sculley, 2010):

$$\mathcal{L}_{BCE} = \sum_{(x,p)\in\{(x,p_i),(x,p_j)\}} a_{\mathcal{M}} \times (\log(p(x,p)) + (1-a_{\mathcal{M}}) \times \log(1-p(x,p)),$$
(2)

where $p(x, p) = \text{sigmoid}(v_{\mathcal{M}})$ and $a_{\mathcal{M}}$ is the target accuracy label taking values in the set $\{0, 1\}$. We train the utility predictor with the following combined objective:

$$\mathcal{L} = \mathcal{L}_{rank} + \lambda \, \mathcal{L}_{BCE},\tag{3}$$

where λ is a hyper-parameter. Both the ranking and BCE objectives are compatible with gold annotations that could be obtained via human intervention in an interactive and active learning learning setting. That is, it would be feasible to elicit from human judges (e.g., moderators of the QA system) answer accuracy labels (e.g., *correct/incorrect*) and level of passage support for the generated answer (e.g., *best* or *worse*) (Simpson et al., 2020; Fang et al., 2024). Note that the Utility Ranker could also be trained with different variants of this objective that also exhibit competitive performance. We report in Appendix D.1 a study on the ablation of the different components of the training objective.

237 The passage utility predictor is related to the direct error prediction approach in (Lahlou et al., 2023). Lahlou et al. (2023) train a secondary model to estimate target model loss; instead, we train the 238 passage utility predictor with sequence level metrics, i.e., accuracy and entailment, which indirectly 239 measure error. This choice is best suited for our task for various reasons. First, in the context 240 of text generation and its possibly diverse (e.g., paraphrases) but correct set of possible generated 241 answers (Kuhn et al., 2023), predicting loss against a unique single paraphrase would result in a 242 too narrow estimation. Our choice is also adequate for proprietary LLMs where it is not possible to 243 create training data with model losses. Finally, our approach is suited for collecting data from user 244 feedback for active model adaptation (Simpson et al., 2020; Fang et al., 2024). In the image domain, 245 van Amersfoort et al. (2020) map inputs to feature representations and take the distance between new 246 inputs and their closest cluster centroids as a measure of uncertainty. In retrieval augmented QA with 247 LLMs, text passages, and questions, it is less clear what the boundary between seen and unseen texts 248 or topics is. Because our Utility Ranker is trained on a target dataset it could be exploited to detect out-of-domain instances for a target application. It would be interesting to pursue future work on 249 using our Utility Ranker as a content controller for the target LLM-based QA model. 250

Some approaches to answer uncertainty prediction that train a secondary model are in (Kamath et al., 2020; Zhang et al., 2021). However, none of them is applied to retrieval augmented QA; but instead to Reading Comprehension (RC), i.e., the task of generating an answer based on a single positive (i.e., supposed to contain the answer) context document. There are two major differences with our work. One is that in their scenario, all input documents are useful while in ours the utility of retrieved passages is varied. The second one is that we show that individual passage utilities are good predictors of retrieval augmented QA with a set of retrieved passages.

258 259

225 226

229

3.2 ANSWER UNCERTAINTY ESTIMATION FOR RETRIEVAL AUGMENTED QA

For retrieval augmented QA, we want an estimator $\{x, R\} \mapsto \mathbf{u}_{\mathcal{M}}(\{x, R\})$ of the answer uncertainty of a target QA model \mathcal{M} when generating answer $y_{\mathcal{M}}$ from a prompt with set of passages R and question x. We propose the direct estimation of $\mathbf{u}_{\mathcal{M}}$ from individual passage utilities predicted for passages in R. The intuition is that, the highest the utility in one (or more) passages in R the less uncertain \mathcal{M} will be when generating answer $y_{\mathcal{M}}$. Concretely, we take the maximum utility score that is given to passages in R as an estimate of answer uncertainty $\mathbf{u}_{\mathcal{M}}$, i.e.,

$$\mathbf{u}_{\mathcal{M}}(\{x,R\}) = \max(\upsilon_{\mathcal{M}}(\{x,p\}) \mid p \in R).$$

$$\tag{4}$$

Note that other more complex estimators $\{x, R\} \mapsto \mathbf{u}_{\mathcal{M}}(\{x, R\})$ could be learned by training, for instance, a regression model on individual passage utilities in addition to other features of the target model \mathcal{M} such as probability of the generated answer $y_{\mathcal{M}}$ (Dong et al., 2018).

270 4 EXPERIMENTAL SETUP

271 272

293

Accuracy Evaluation A precise metric for measuring accuracy is key when evaluating the quality 273 of uncertainty estimation. Token overlap metrics are far from being precise and can over- or under-274 estimate accuracy, e.g., Acc yields a higher score for the pair of gold and generated answers (a 275 politician, not a politician) than for the pair (a politician, a congressperson). Thus, our main metric to 276 evaluate QA model performance and as the accuracy evaluator \mathcal{A} to create data to train the passage 277 utility predictor, is based on a LLM judgement of accuracy proposed by Sun et al. (2024) (AccLM). 278 A critic LLM is prompted with the gold and generated answer and asked to judge whether they are the equivalent. In a sample of 840 generated answers human and LLM-based judgment of 279 correctness agreed 98% of the time (Sun et al., 2024). We use the prompt as proposed in (Sun et al., 280 2024), we include it in Appendix B for completeness. We use Qwen2-72B-Instruct (Yang et al., 281 2024) to obtain accuracy judgments. For compatibility with previous work and as a lower bound, in Appendix D.2, we report QA model performance with token overlap accuracy (Acc) defined as 283 whether the gold answer is contained in the generated answer (Mallen et al., 2023; Asai et al., 2024). 284

Utility Ranker Implementation Details To create the training set $D_{\mathcal{M}}$ to train the Utility Ranker, we consider the first top five retrieved passages for each question, i.e., |R| = 5. Note that this is a hyper-parameter and other values would also be possible, e.g., with larger sizes of |R| further training data would be available. We use the target QA model \mathcal{M} to generate answers $y_{\mathcal{M}}$ for each of the five passages p in R (i.e., \mathcal{M} is prompted with passage p and question x). We then ge utility scores using the LLM-based accuracy judge \mathcal{A} as described above and an ALBERT-xlarge Lan et al. (2020) model optimized on MNLI (Williams et al., 2018) and VitaminC (Schuster et al., 2021) as our entailment judge \mathcal{E} .

Comparison Approaches and Baselines We choose the stronger methods from previous work
 (Fadeeva et al., 2023) to compare our approach with.

Information Based. We compare against the stronger information based uncertainty quantification approaches reported in previous work Fadeeva et al. (2023). These are based on predictive probabilities; recall that the predictive distribution under QA model \mathcal{M} prompted with question x and set of passages R is $P(y_{\mathcal{M}}|x, R, \mathcal{M}) = \prod_{t=1}^{|y_{\mathcal{M}}|} p_{\mathcal{M}}(y_t|y_{1..t-1}, x, R)$ for a target QA model \mathcal{M} .

Maximum Sequence Probability (MSP) based uncertainty estimation is based on the probability of the most likely answer and computed as $MSP(y_{\mathcal{M}} | x, R, \mathcal{M}) = 1 - P(y_{\mathcal{M}} | x, R, \mathcal{M})$. The other uncertainty estimation approach is the negative mean Point-wise Mutual Information (PMI) Takayama & Arase (2019); i.e., it compares the probability of generating answer $y_{\mathcal{M}}$ given the prompt with question x and passages R w.r.t the probability given by \mathcal{M} to $y_{\mathcal{M}}$ without context. Intuitively, the higher the PMI the more certain on generating $y_{\mathcal{M}}$. PMI is computed as

 $PMI(y_{\mathcal{M}}, x, R; \mathcal{M}) \frac{1}{|y_{\mathcal{M}}|} \sum_{t=1}^{|y_{\mathcal{M}}|} \log \frac{p_{\mathcal{M}}(y_t|y_{1..t-1})}{p_{\mathcal{M}}(y_t|y_{1..t-1}, x, R)}.$ The other two methods are based on entropy. We compare with Regular Entropy (RE), i.e., the entropy on the predictive distribution computed at sequence level $\mathbb{E}[-\log P(y_{\mathcal{M}}|x, R, \mathcal{M})]$ with \mathbb{E} computed on sequences $y_{\mathcal{M}}$ sampled from $P(y_{\mathcal{M}}|x, R, \mathcal{M})$. In practice, this is approximated via Monte-Carlo integration, i.e., sampling N random answers from $P(y_{\mathcal{M}}|x, R, \mathcal{M})$. Thus, Regular Entropy is computed as $-\frac{1}{N} \sum_{n=1}^{N} \log \tilde{P}(y_{\mathcal{M}}^{(n)}|x, R, \mathcal{M})$, where $\tilde{P}(y_{\mathcal{M}}^{(n)}|x, R, \mathcal{M})$ is the length normalised version of $P(y_{\mathcal{M}}^{(n)}|x, R, \mathcal{M})$.

Answer Variation. Kuhn et al. (2023) propose a variant of regular entropy, named Semantic Entropy (SE), that accounts for uncertainty in the surface form of the generated answers rather than on meaning. Concretely, Semantic Entropy clusters the set of N samples into $M, M \leq N$, clusters with the same meaning via bidirectional entailment. Then computes the average answers' probability within each cluster, $SE(x, \mathcal{M}) = -\sum_{m=1}^{M} \hat{P}_m(x, \mathcal{M}) \log \hat{P}_m(x, \mathcal{M})$ where $\hat{P}_m(x, \mathcal{M}) = \frac{\sum_{y_{\mathcal{M}} \in C_m} P(y_{\mathcal{M}} | x, R, \mathcal{M})}{\sum_{m=1}^{M} \sum_{y_{\mathcal{M}} \in C_m} P(y_{\mathcal{M}} | x, R, \mathcal{M})}$.

Reflexive. We compare with p(true) proposed by Kadavath et al. (2022). This approach uses the same target QA model (LLM) evaluate whether the answers it produces are correct. It is prompted with the question and a set of candidate answers, i.e., the most likely answer plus a sample of size N answers,

and instructed to respond whether the most likely answer is true or false (i.e., correct/incorrect). The score produced by this approach is the probability of the model \mathcal{M} generating the token True. p(true) needs several in-context examples to work well, so we fit as many examples as can be in the context.

Baselines. The sets of passages in R are originally ranked by the IR system, so each passage in Rhas a retriever score which can be seen as baseline passage utility. We thus take the Retriever Score as a baseline. Despite the QA models are instructed to produce a short answer, these often generate longer answers. The length of the answer could be a feature indicating that the model is uncertain about the answer. Thus, we estimate answer uncertainty from the Answer Length (Ans.Len) as the number of words in the answer.

Following previous work (Farquhar et al., 2024), we take N = 10 samples and use multinomial sampling to generate samples. That is, we set the sampling temperature to 1, with nucleus sampling (P = 0.9) (Holtzman et al., 2020) and top-K sampling (K = 50) (Fan et al., 2018), and use a different random seed to draw each sample. Most likely answers are generated with greedy sampling at temperature equal to 0. We use the implementation provided by Farquhar et al. (2024) to compute RE, SE, CA, and p(true). We report inference cost of each approach in Appendix C.1.

339

340 **QA Models** Our target retrieval augmented QA models \mathcal{M} are based on the following instruction 341 fine-tuned LLMs. To assess the performance of the Utility Ranker for QA models that potentially 342 exhibit different answer uncertainty, we consider different families of similar size. These are Llama-3.1-8B-Instruct (AI@Meta, 2024), Mistral-7B-Instruct-v0.3 Jiang et al. (2023), and Gemma2-9B-it 343 Riviere et al. (2024). For all QA models, we use a simple prompt including the retrieved passages 344 and the question in the context, the prompt is shown in Table 6 of the Appendix. We use vLLM 345 for inference (Kwon et al., 2023). Following previous work on retrieval augmented QA, we use 346 Contriever Izacard et al. (2022) as our external retriever (Asai et al., 2024) and the target QA models 347 are prompted with |R| = 5 passages Yu et al. (2023); Asai et al. (2024); Xu et al. (2024). 348

349 Datasets We evaluate our answering uncertainty prediction approach on short-form answer gen-350 eration tasks. Concretely, we evaluate on the Natural Questions Kwiatkowski et al. (2019), Trivi-351 aQA Joshi et al. (2017), WebQuestions Berant et al. (2013), and SQuAD (Rajpurkar et al., 2016) 352 datasets. We follow the training/validation/test splits in prior work Lee et al. (2019); Min et al. 353 (2019); Karpukhin et al. (2020). To test the generalisation robustness of our approach we carry out 354 additional experiments on PopQA Mallen et al. (2023), a dataset with questions about rare entities, 355 and RefuNQ Liu et al. (2024a), a dataset with unanswerable questions about non-existing entities. Statistics about our datasets are given in the Appendix in Table 5. 356

Evaluation of the Quality of Uncertainty Estimation To assess the quality of answer uncer-358 tainty prediction, we follow Farquhar et al. (2024) and report the Area Under the Receiver Operator 359 Curve on detecting answer uncertainty, i.e., incorrect answers, (AUROC) and the area under the 360 rejection accuracy curve (AURAC). AURAC summarises the accuracy of QA models when answer 361 uncertainty is used to refuse to answer questions. It summarises accuracy at different percentages of 362 rejection. Instruction fine-tuned models are known to refuse to answer questions, i.e., they produce 363 answers such as This information is not available in the text. In some cases, the refusal response 364 will be adequate (e.g., no input passage contains the information to answer) but in many cases QA 365 models may refuse when they should have provided an answer Adlakha et al. (2024); Liu et al. 366 (2024a). Thus, to simplify the assessment of answer correctness, we did not explicitly instruct the 367 QA models to abstain and treat occurring refusal answers as cases of uncertainty where the QA model is expressing the uncertainty in the answer (Farquhar et al., 2024). We report the percentage 368 of refusal answers for each QA model and QA task on development sets in Appendix D.2. 369

370 371

374

357

- 5 Results
- 372 373

5.1 UNCERTAINTY QUANTIFICATION

Answer uncertainty estimation results for the three QA models (GEMMA2-9B, LLAMA-3.1-8B, and MISTRAL-7B-V0.3) are shown in Table 1 (results on the development set are included in Appendix D). In terms of predicting answer uncertainty (i.e., model incorrect answers), column AU-ROC in Table 1, simple metrics based on models' probabilities such as MSP perform better for some

	Natural	Juestions	Trivi	aQA	WebQu	estions	SQu	IAD
	AUROC	AURAC	AUROC	AURAC	AUROC	AURAC	AUROC	AURAC
				Gemm	A2-9B			
MSP	0.69	0.67	0.68	0.83	0.63	0.66	0.65	0.63
PMI	0.51	0.58	0.53	0.78	0.45	0.58	0.50	0.55
p(true)	0.72	0.70	0.78	0.86	0.74	0.74	0.67	0.66
Regular Entropy	0.69	0.68	0.65	0.82	0.63	0.67	0.65	0.62
Cluster Assignment	0.67	0.66	0.69	0.83	0.60	0.65	0.66	0.63
Semantic Entropy	0.68	0.67	0.68	0.83	0.60	0.65	0.66	0.63
Ans.Len	0.65	0.65	0.59	0.80	0.62	0.66	0.61	0.60
Retriever Score	0.60	0.65	0.68	0.84	0.53	0.62	0.61	0.62
Utility Ranker	0.76	0.72	0.81	0.88	0.72	0.71	0.81	0.74
				LLAMA	-3.1-8B			
MSP	0.71	0.69	0.83	0.88	0.71	0.74	0.77	0.69
PMI	0.56	0.60	0.57	0.78	0.51	0.65	0.61	0.59
p(true)	0.79	0.74	0.84	0.87	0.76	0.76	0.65	0.61
Regular Entropy	0.72	0.69	0.83	0.88	0.72	0.74	0.78	0.69
Semantic Entropy	0.69	0.67	0.81	0.86	0.68	0.73	0.75	0.68
Ans.Len	0.59	0.61	0.60	0.78	0.61	0.68	0.57	0.55
Retriever Score	0.58	0.62	0.64	0.81	0.50	0.63	0.65	0.61
Utility Ranker	0.73	0.70	0.78	0.86	0.76	0.78	0.84	0.73
				MISTRAL	-7B-v0.3			
MSP	0.68	0.63	0.73	0.87	0.65	0.68	0.71	0.65
PMI	0.53	0.59	0.55	0.79	0.50	0.62	0.58	0.60
p(true)	0.72	0.67	0.84	0.88	0.72	0.70	0.69	0.63
Regular Entropy	0.60	0.60	0.71	0.85	0.61	0.68	0.66	0.62
Semantic Entropy	0.67	0.63	0.78	0.88	0.69	0.68	0.71	0.66
Ans.Len	0.68	0.63	0.67	0.84	0.64	0.69	0.66	0.63
Retriever Score	0.59	0.60	0.67	0.82	0.53	0.65	0.64	0.62
Utility Ranker	0.76	0.68	0.79	0.86	0.76	0.71	0.80	0.68

Table 1: Answer uncertainty estimation for QA models GEMMA2-9B, LLAMA-3.1-8B, and
MISTRAL-7B-V0.3 on NaturalQuestions, TriviaQA, WebQuestions, and SQuAD (evaluation with
in-distribution test data for the Utility Ranker). We report AUROC and AURAC.

models. It exhibits high performance for LLAMA-3.1-8B while lower performance for GEMMA29B and MISTRAL-7B-v0.3. Sampling-based approaches (meaning diversity and reflexive), can
better identify model uncertainty but at the cost of running inference several times to have a good
size sample for the estimation. Our Utility Ranker has similar or better performance with a single
inference step on each input passage. We speculate that clustering approaches can suffer in phrase
or sentence level correct answers where these contain different levels of details Zhang et al. (2024);
thus, not being clustered together wrongly suggesting variation.

On improving question-answering accuracy, AURAC column in Table 1, with the exception of Triv-iaQA, all uncertainty prediction approaches outperform the information theoretic approaches (i.e., MSP, PMI). The Utility Ranker performs on par or better than the more expensive sampling based approaches. To have a clearer picture of baseline retrieval augmented QA accuracy w.r.t. accuracy when the uncertainty estimation is used to decide whether to abstain nor not, we show in Figure 2 the accuracy of the model at different thresholds for answer rejection. That is, we report when the QA model chooses to answer only the 80% or 90% of the most confident cases as well as when always answers. Retrieval augmented QA accuracy per model and dataset on the full test and development sets is included in Appendix D. Across all datasets, the Utility Ranker performs on par with of better than more expensive uncertainty estimation approaches. The easiest QA task is TriviaQA where QA models show very good performance and information theoretic methods work on par with more complex ones. On the most difficult task, SQuAD, the utility ranker outperforms all other methods both at 20% and 10% of rejected answers.

5.2 ROBUSTNESS AND GENERALISATION OF UNCERTAINTY ESTIMATION

We assess the robustness and generalisation of our Utility Ranker on test cases that are different from
those examples seen during training, i.e., Out-Of-Distribution (OOD). These examples encompass
real cases that a QA model will face at test time such as different type of questions, e.g., longer and
more complex. We also study extreme adversarial cases such as questions about tail knowledge for
both retrievers and LLMs (PopQA) and unanswerable questions (RefuNQ).

- **Distribution Shift** Table 2 shows the performance of the Utility Ranker when evaluated in OOD data. The first column indicates the training data and the first row the evaluation data. Results in the



Figure 2: Average QA model performance on test sets with |R| = 5. We show model based accuracy (AccLM) at different percentages of rejecting to answer (i.e., when choosing to respond on 80%, 90%, and all the cases) given uncertainty estimations by the different approaches.

Table 2: Performance of GEMMA2-9B's Utility Ranker on distribution shift. That is, trained on one dataset and evaluated zero-hot on another one. We report all combinations of train and test data. The first column indicates train data while the first row test data.

	NaturalQuestions		TriviaQA		WebQu	estions	SQuAD		
	AUROC	AURAC	AUROC	AURAC	AUROC	AURAC	AUROC	AURAC	
NaturalQuestions	0.76	0.72	0.72	0.86	0.65	0.67	0.72	0.68	
TriviaQA	0.64	0.67	0.81	0.88	0.63	0.68	0.71	0.68	
WebQuestions	0.60	0.64	0.72	0.86	0.72	0.71	0.58	0.59	
SQuAD	0.65	0.67	0.77	0.87	0.61	0.65	0.81	0.74	

diagonal correspond to the Utility Ranker trained and evaluated in the same data distribution; the offdiagonal cells to the Utility Ranker evaluated zero-shot in a different dataset. As expected, the Utility Ranker variants evaluated on a different dataset show a decrease in performance. However, for some training data the decrease is small providing a competitive prediction. That is, NaturalQuestions and SQuAD provide the best training data, what agrees with previous experiments in reading comprehension settings (Chen et al. (2021) choose NaturalQuestions to train the base model, Kamath et al. (2020); Zhang et al. (2021) SQuAD). The Utility Ranker variants trained on WebQuestions (smallest training set) and TriviaQA (the easiest task) have the worst generalisation performance. Note that we focus on zero-shot to assess bare transfer performance; however, it would make sense to train the model with few examples of the OOD data (Kamath et al., 2020; Zhang et al., 2021).

473 474 475

432

452

453

454 455

456

457

466

467

468

469

470

471

472

Adversarial Questions Table 3 reports results for GEMMA2-9B's Utility Ranker trained on Natu-476 ralQuestions and evaluated zero-shot to predict answer uncertainty for retrieval augmented QA with 477 |R| = 5 on PopQA and RefuNQ. These datasets are made of adversarial cases so we report AU-478 ROC (predicting incorrect answers) and AURAC (summary of different rejection thresholds). For 479 RefuNQ as questions have (un)answerable gold annotations, we further report AUROC scores for 480 the Unanswerable questions (51% of the tested cases) and all incorrect answers together (67%) of 481 the test cases. The Utility Ranker (NQ) outperforms other methods to detect answer uncertainty 482 across datasets and improve QA accuracy by refusing to answer questions across the board. We attribute this to the fact that, either due to knowledge about tail entities (PopQA) or unanswerable 483 questions about nonexistent concepts (RefuNQ), the quality of the retrieved passages suffers. Thus, 484 our approach will assign lower utility to these and thus successfully predict answer uncertainty. This 485 is confirmed by the surprisingly high AUROC score achieved by the Retriever Score baseline on Re-

486 Table 3: Answer uncertainty estimation for GEMMA2-9B on adversarial QA tasks (PopQA and 487 RefuNQ). Its Utility Ranker is trained on Natural Questions. 488

	Pop	OA		RefuNO	
	AUROC	AURAC		AUROC	AURAC
			All	Unanswerable	
MSP	0.66	0.58	0.66	0.63	0.39
PMI	0.51	0.50	0.54	0.53	0.35
p(true)	0.71	0.62	0.73	0.65	0.45
Regular Entropy	0.66	0.58	0.66	0.61	0.39
Semantic Entropy	0.69	0.59	0.68	0.60	0.41
Ans.Len	0.62	0.55	0.65	0.66	0.38
Retriever Score	0.63	0.58	0.76	0.80	0.47
Utility Ranker (NQ)	0.72	0.62	0.82	0.71	0.51

Table 4: Retrieval augmented QA performance with three passages |R| = 3 is the version with the 497 top three retrieved passages from Contriever and $|R^{UR}| = 3$ is the version with top three re-ranked 498 passages out of ten originally retrieved. We report model based (AccLM) accuracy. 499

	NaturalQuestions	TriviaQA	WebQuestions	SQuAD
R = top 3 ranked by external retriever	0.58	0.77	0.63	0.53
R = top 3 re-ranked by Utility Ranker	0.62	0.79	0.65	0.60
R = all 10 passages	0.64	0.80	0.66	0.62

fuNQ's Unanswarable questions. In this particular type of questions, the retrieval system indeed struggles to retrieve relevant passages. Note that Retriever Score behaves otherwise in the rest of the QA tasks where it shows lower performance. Interestingly, information based methods, MSP and PPL, perform worse in these adversarial QA tasks than in the in-distribution test cases (Section 5.1). This shows that in these cases QA models produce incorrect answers with high confidence.

5.3 IMPROVING QA PERFORMANCE

512 We also assess the quality of the passage utility scores to identify informative passages via end task 513 QA performance. We compare the original ranking by the external retrieval system with the ranking 514 established by the utility scores by taking the top 3 passages out of 10 passages ordered by the 515 external retriever and re-ranked by the Utility Ranker. We then run the QA models with a budget of 516 |R| = 3 input passages. We also run the QA model with the all the 10 passages, i.e., with |R| = 10. 517 Results for GEMMA2-9B QA model are shown in Table 4. The QA model with the top 3 passages re-ranked by the Utility Ranker improves 4 points on NaturalQuestions and 7 points on SQuAD over 518 the QA model variant that takes the top 3 ranked by the external retriever. This suggest that passages 519 considered relevant for user questions by the external retriever do not coincide with what is useful 520 for the target QA model. The QA model variant with the top 3 passages re-ranked by the Utility 521 Ranker performs very close, i.e., difference of 1 or 2 points across all datasets, to the QA model 522 variant with the 10 passages given as context. The utility scores effectivelly identify informative 523 passages and the QA model achieves comparable performance with a much shorter prompt. 524

6 CONCLUSIONS

526 527 In this work we present an approach to answer uncertainty prediction for retrieval augmented QA 528 models. Importantly, this approach relies on single passage utilities. This approach is based on a 529 small neural model that is trained on a target QA model judgements of retrieved passage usefulness. 530 We show that this approach is competitive or better than existing strong error prediction approaches 531 while being light-weight. Our experiments also show that our approach is particularly better in 532 cases of extreme QA model answer uncertainty like rare entities and unanswerable questions. Future 533 work would explore the approach in the context of log-form generation tasks, e.g., query focused-534 generation. It would also be interesting to explore to what extent the Utility Ranker model could be 535 used in active learning scenarios.

536

538

525

500 501

504

505

506

507

508

509 510

511

6.1 ETHICS STATEMENT

Our work does not involve human subjects. We use QA datasets that are publicly available and widely used by the research community.

540 6.2 REPRODUCIBILITY STATEMENT

We build up on existing base code Farquhar et al. (2024); Fang et al. (2024) and we will make available all code and data together with the docker images for reproducibility.

References

544

546

551

552

553

558

573

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy.
 Evaluating correctness and faithfulness of instruction-following models for question answer-*Transactions of the Association for Computational Linguistics*, 12:681–699, 2024. doi: 10.1162/tacl_a_00667. URL https://aclanthology.org/2024.tacl-1.38.
 - AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/ llama3/blob/main/MODEL_CARD.md.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id= hSyW5go0v8.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on Freebase from question-answer pairs. In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard (eds.), *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 1533–1544, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL https://aclanthology.org/D13-1160.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-564 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-565 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, 566 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz 567 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec 568 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In 569 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neu-570 ral Information Processing Systems, volume 33, pp. 1877–1901. Curran Associates, Inc., 571 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/ 572 file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.

- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer opendomain questions. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1870–
 1879, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/
 v1/P17-1171. URL https://aclanthology.org/P17-1171.
- Jifan Chen, Eunsol Choi, and Greg Durrett. Can NLI models verify QA systems' predictions? In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 3841–3854, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-emnlp.324. URL https://aclanthology.org/2021.
 findings-emnlp.324.
- Jiuhai Chen and Jonas Mueller. Quantifying uncertainty in answers from any language model and enhancing their trustworthiness. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 5186–5200, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.acl-long.283.
- 590

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*)

 and Short Papers), pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/ N19-1423.

- Li Dong, Chris Quirk, and Mirella Lapata. Confidence modeling for neural semantic parsing. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 743–753, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1069. URL https://aclanthology.org/P18-1069.
- Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill
 Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, Timothy
 Baldwin, and Artem Shelmanov. LM-polygraph: Uncertainty estimation for language models. In
 Yansong Feng and Els Lefever (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 446–461, Singapore, December
 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-demo.41. URL
 https://aclanthology.org/2023.emnlp-demo.41.
- Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 889–898, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1082. URL https://aclanthology.org/P18-1082.
- Haishuo Fang, Jeet Gor, and Edwin Simpson. Efficiently acquiring human feedback with Bayesian deep learning. In Raúl Vázquez, Hande Celikkanat, Dennis Ulmer, Jörg Tiedemann, Swabha Swayamdipta, Wilker Aziz, Barbara Plank, Joris Baan, and Marie-Catherine de Marneffe (eds.), *Proceedings of the 1st Workshop on Uncertainty-Aware NLP (UncertaiNLP 2024)*, pp. 70–80, St Julians, Malta, March 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.uncertainlp-1.7.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large
 language models using semantic entropy. *Nature*, 2024. URL https://doi.org/10.1038/
 s41586-024-07421-0.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR.
 URL https://proceedings.mlr.press/v48/gall6.html.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1321–1330. PMLR, 06–11 Aug 2017. URL https://proceedings.mlr.press/v70/guo17a.html.

- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented
 language model pre-training. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 3929–3938. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.
 press/v119/guu20a.html.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=rygGQyrFvH.
- Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. Decomposing uncertainty for large language models through input clarification ensembling. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/ forum?id=byxXa99PtF.

- Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answering. In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (eds.), Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pp. 874–880, Online, April 2021a. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.74. URL https://aclanthology.org/2021.eacl-main.74.
- Gautier Izacard and Edouard Grave. Distilling knowledge from reader to retriever for question answering. In International Conference on Learning Representations, 2021b. URL https://openreview.net/forum?id=NTEz-6wysdb.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand
 Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learn *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL https:
 //openreview.net/forum?id=jKN1pXi7b0.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane
 Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: few-shot learning with
 retrieval augmented language models. *J. Mach. Learn. Res.*, 24(1), mar 2024. ISSN 1532-4435.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. How Can We Know When Language Models Know? On the Calibration of Language Models for Question Answering. *Transactions* of the Association for Computational Linguistics, 9:962–977, 09 2021. ISSN 2307-387X. doi: 10.1162/tacl_a_00407. URL https://doi.org/10.1162/tacl_a_00407.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pp. 1601–1611, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1147. URL https://aclanthology.org/ P17-1147.
- 681 Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, 682 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer 683 El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bow-684 man, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna 685 Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Ka-686 plan. Language models (mostly) know what they know, 2022. URL https://arxiv.org/ 687 abs/2207.05221. 688
- Amita Kamath, Robin Jia, and Percy Liang. Selective question answering under domain shift. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5684–5696, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.503. URL https://aclanthology.org/2020.acl-main.503.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6769–6781, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.550.
 URL https://aclanthology.org/2020.emnlp-main.550.
- Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Velocity Yu, Dragomir Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. Realtime qa: what's the

answer right now? In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2024. Curran Associates Inc.

- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id= VD-AYtPOdve.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl_a_00276. URL https://aclanthology.org/Q19-1026.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, SOSP '23, pp. 611–626, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400702297. doi: 10.1145/3600006.3613165. URL https://doi.org/10.1145/3600006.3613165.
- Salem Lahlou, Moksh Jain, Hadi Nekoei, Victor I Butoi, Paul Bertin, Jarrid Rector-Brooks, Maksym Korablyov, and Yoshua Bengio. DEUP: Direct epistemic uncertainty prediction. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=eGLdVRvvfQ. Expert Certification.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/ file/9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Sori cut. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?
 id=H1eA7AEtvS.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open domain question answering. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6086–6096, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1612. URL https://aclanthology.org/P19-1612.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

- Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=8s8K2UZGTZ.
- Genglin Liu, Xingyao Wang, Lifan Yuan, Yangyi Chen, and Hao Peng. Examining llms' uncertainty expression towards questions outside parametric knowledge, 2024a.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024b. doi: 10.1162/tacl_a_00638. URL https://aclanthology.org/2024.tacl-1.9.

Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi.
When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9802–9822, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.546. URL https://aclanthology.org/2023.acl-long.546.

- Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. Reducing conversational agents' overconfidence through linguistic calibration. *Transactions of the Association for Computational Linguistics*, 10:857–872, 2022. doi: 10.1162/tacl_a_00494. URL https://aclanthology.org/2022.tacl-1.50.
- Sewon Min, Danqi Chen, Hannaneh Hajishirzi, and Luke Zettlemoyer. A discrete hard EM approach for weakly supervised question answering. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2851–2864, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1284. URL https://aclanthology.org/D19–1284.
- 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. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713871088.

SeongII Park, Seungwoo Choi, Nahyun Kim, and Jay-Yoon Lee. Enhancing robustness of retrievalaugmented language models with in-context learning. In Wenhao Yu, Weijia Shi, Michihiro Yasunaga, Meng Jiang, Chenguang Zhu, Hannaneh Hajishirzi, Luke Zettlemoyer, and Zhihan Zhang (eds.), *Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP*, pp. 93–102, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.knowledgenlp-1.7. URL https://aclanthology.org/2024. knowledgenlp-1.7.

- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In Jian Su, Kevin Duh, and Xavier Carreras (eds.), Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pp. 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL https://aclanthology.org/D16-1264.
- 793 Gemma Team Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-794 patiraju, L'eonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ram'e, Johan Ferret, Peter Liu, Pouva Dehghani Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stańczyk, Ser-796 tan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam 797 Neyshabur, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, 798 Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia 799 Paterson, Ben Bastian, Bilal Piot, Boxi Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris 800 Perry, Christoper A. Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Wein-801 berger, Dimple Vijaykumar, Dominika Rogozi'nska, D. Herbison, Elisa Bandy, Emma Wang, 802 Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Pluci'nska, Harleen 804 Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, 805 Jetha Chan, Jin Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost R. van Amersfoort, Josh Gordon, Josh Lipschultz, Joshua Newlan, Junsong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, L. Sifre, L. Heuermann, Leticia Lago, 808 Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Gorner, Mat Velloso, Mateo Wirth, Matt Davidow,

843

844

845

846

847

810 Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moyni-811 han, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, 812 Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil 813 Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, 814 Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, S. Mc Carthy, Sarah Perrin, S'ebastien Arnold, Sebastian 815 Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, 816 Ting Yu, Tom Eccles, Tom Hennigan, Tomás Kociský, Tulsee Doshi, Vihan Jain, Vikas Yadav, 817 Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, 818 Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, 819 Anand Rao, Minh Giang, Ludovic Peran, Tris Brian Warkentin, Eli Collins, Joelle Barral, Zoubin 820 Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, 821 Jeffrey Dean, Demis Hassabis, Koray Kavukcuoglu, Cl'ement Farabet, Elena Buchatskaya, Se-822 bastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek 823 Andreev. Gemma 2: Improving open language models at a practical size. ArXiv, abs/2408.00118, 824 2024. URL https://api.semanticscholar.org/CorpusID:270843326. 825

- Alireza Salemi and Hamed Zamani. Evaluating retrieval quality in retrieval-augmented generation. 826 In Proceedings of the 47th International ACM SIGIR Conference on Research and Development 827 in Information Retrieval, SIGIR '24, pp. 2395–2400, New York, NY, USA, 2024. Association for 828 Computing Machinery. ISBN 9798400704314. doi: 10.1145/3626772.3657957. URL https: 829 //doi.org/10.1145/3626772.3657957. 830
- Tal Schuster, Adam Fisch, and Regina Barzilay. Get your vitamin C! robust fact verification 831 with contrastive evidence. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek 832 Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao 833 Zhou (eds.), Proceedings of the 2021 Conference of the North American Chapter of the Asso-834 ciation for Computational Linguistics: Human Language Technologies, pp. 624–643, Online, 835 June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.52. 836 URL https://aclanthology.org/2021.naacl-main.52. 837
- Christopher Sciavolino, Zexuan Zhong, Jinhyuk Lee, and Danqi Chen. Simple entity-centric ques-838 tions challenge dense retrievers. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and 839 Scott Wen-tau Yih (eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural 840 Language Processing, pp. 6138–6148, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.496. URL 842 https://aclanthology.org/2021.emnlp-main.496.
 - D. Sculley. Combined regression and ranking. In Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '10, pp. 979–988, New York, NY, USA, 2010. Association for Computing Machinery. ISBN 9781450300551. doi: 10.1145/1835804.1835928. URL https://doi.org/10.1145/1835804.1835928.
- Edwin Simpson, Yang Gao, and Iryna Gurevych. Interactive Text Ranking with Bayesian Optimiza-848 tion: A Case Study on Community QA and Summarization. Transactions of the Association for 849 Computational Linguistics, 8:759–775, 12 2020. ISSN 2307-387X. doi: 10.1162/tacl_a_00344. 850 URL https://doi.org/10.1162/tacl_a_00344. 851
- 852 Kai Sun, Yifan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. Head-to-tail: How knowl-853 edgeable are large language models (LLMs)? A.K.A. will LLMs replace knowledge graphs? 854 In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human 855 Language Technologies (Volume 1: Long Papers), pp. 311-325, Mexico City, Mexico, June 856 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.18. URL 857 https://aclanthology.org/2024.naacl-long.18. 858
- 859 Junya Takayama and Yuki Arase. Relevant and informative response generation using pointwise mutual information. In Yun-Nung Chen, Tania Bedrax-Weiss, Dilek Hakkani-Tur, Anuj Kumar, Mike Lewis, Thang-Minh Luong, Pei-Hao Su, and Tsung-Hsien Wen (eds.), Proceed-861 ings of the First Workshop on NLP for Conversational AI, pp. 133–138, Florence, Italy, Au-862 gust 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4115. URL https://aclanthology.org/W19-4115.

- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea
 Finn, and Christopher Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In Houda Bouamor,
 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Meth- ods in Natural Language Processing*, pp. 5433–5442, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.330. URL https:
 //aclanthology.org/2023.emnlp-main.330.
- Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 9690–9700. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/van-amersfoort20a.html.
- Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. 2020. URL https://arxiv.org/abs/2003. 02037.
- Yuhao Wang, Ruiyang Ren, Junyi Li, Wayne Xin Zhao, Jing Liu, and Ji-Rong Wen. Rear: A
 relevance-aware retrieval-augmented framework for open-domain question answering. *arXiv preprint arXiv:2402.17497*, 2024. URL https://arxiv.org/abs/2402.17497.
- Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1112–1122, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1101. URL https://aclanthology.org/N18-1101.
- Fangyuan Xu, Weijia Shi, and Eunsol Choi. RECOMP: Improving retrieval-augmented LMs
 with context compression and selective augmentation. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=
 mlJLVigNHp.
- 894 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Daviheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, 895 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, 896 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng 897 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan 899 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang 900 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 901 technical report. arXiv preprint arXiv:2407.10671, 2024. 902
- Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. Making retrieval-augmented language
 models robust to irrelevant context. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=ZS4m74kZpH.
- Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin Ma, Hongwei Wang, and Dong Yu. Chainof-note: Enhancing robustness in retrieval-augmented language models, 2023. URL https: //arxiv.org/abs/2311.09210.
- Huajian Zhang, Yumo Xu, and Laura Perez-Beltrachini. Fine-grained natural language inference based faithfulness evaluation for diverse summarisation tasks. In Yvette Graham and Matthew Purver (eds.), Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1701–1722, St. Julian's, Malta, March 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.eacl-long.102.
- Shujian Zhang, Chengyue Gong, and Eunsol Choi. Knowing more about questions can help: Improving calibration in question answering. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*,

918	I	Dataset	Train	Dev	Test	
919	Ī	Natural Questions	79,168	8,757	3,610	_
920	Т	TriviaQA	78,785	8,837	11,313	
921	V	VebQuestions	2,474	361	2,032	
922	S	SQuAD	78,713	8,886	10,570	
923	F	PopQA	11267	-	3000	
924	k	RefuNQ	-	-	4439	
925	T11.5 D.4.4.4.4.4.4.4.4.4.4.4.4.4.4.4.4.4.4.4	1		1. 4	NT	41
926	Table 5: Dataset statistics, r	number of instances	s per train/	dev/test	sets. Not	e that we sample a smaller
927	test set for PopQA in our es	xperiments.				
928	г	V				
929		Knowledge:				
930						
931		[1] [passage]				
032		[2] [passage]				
000		 [N] [passage]				
933						
934		Answer the follow	ving quest	tion with	a verv	
935		short phrase	ang quest		avery	
936		short pinuse.				
937		Ouestion: <u>[questi</u>	onl			
938	L	C				
939	Table 6: M	linimal prompt sele	ected as us	er turn f	or the QA	A models.
940		1 1				
941						
942	pp. 1958–1970, Online, A	August 2021. Asso	ciation for	Comput	ational L	inguistics. doi: 10.18653/
943	v1/2021.findings-acl.172	2. URL https://	/aclant	hology	.org/2	2021.findings-acl.
944	172.					
945	Kaithan 7han Dan Landala	Totowa	-h:	NT		···· · ··· · · · · · · · · · · · · · ·
946	sfuncertainty and overage	y, and Taisunori Ha	isnimoto.	Navigati	ng the gr	ey area: How expressions
947	like Deli (edg.) <i>Proceedi</i>	ince of the 2023 Co	guage moo	Jeis. III F	ioula Do	hodg in Natural Language
948	Processing pp 5506 55	24 Singapore Dec	ombor 20	$23 \Lambda_{SSO}$	ciation fo	nous in Natural Language
949	tics doi: 10.18653/v1/20	24, Siligapole, Dec	5 LIRI h	23. ASSU ++nc • /		at hology org/2023
950	emplo-main 335	25.cmmp-main.55	5. UKL 11	ctps./	/ actai	10101099.019/2023.
951	chilip marin. 555.					
052						
052	A DATASETS					
05/						
904	Table 5 shows statistics abo	wit the $\Omega \Delta$ datasets	we use in	ourevo	eriments	
955	Table 5 shows statistics abe	out the QA datasets	we use m	i our exp	erments.	
956						
957	B MODEL PROMPTS	5				
958						
959	The prompt we use for our	QA models is show	wn in Tab	le 6. Tab	le 7 illus	trate the prompts used for
960	our LLM based accuracy an	nd p(true) baseline.				
961						
962		D DAGELINE II	NCEDTA	INTV E	CTIMA	TION METHODS
963	C COMPARISON AN	D DASELINE U	NCERIA		SIIMA	TION METHODS
964			70000 0.000			
965	U.1 TEST TIME COST O	F UNCERTAINTY I	LSTIMATI	UN MET	HODS	
966	Table 8 shows the cost of a	vecuting uncertain	ty estimat	ion for a	liser que	estion r in terms of model
967	inference calls required S	simple information	theoretic	methode	require	a single call (PPI MSP)
968	or two (PMI) calls to the C	A model with the	full nrom	pt (N re	trieved n	assages and user question
969	x): similarly the Ans.Len h	aseline. However	approach	es that es	stimate u	ncertainty based on diver-
970	sity (Regular Entropy, Clu	ster Assignment, S	Semantic I	Entropy,	and p(tru	ue)) require generating N

are required to compare the set of sampled answers. p(true) requires one additional LLM call to elicit a True/False answer but with a very long prompt including in-context examples and candidate answers. In contrast, our approach requires |R| utility predictions with a BERT-size model.

976 D ADDITIONAL RESULTS

979 D.1 DIFFERENT COMPONENTS OF THE TRAINING OBJECTIVE

Table 9 shows results on the ablation of the Utility Ranker training objective (Section 3.1, Equa-tion 3). When trained only with the ranking loss (\mathcal{L}_{rank}), in average it achieves better performance when the training signal combines accuracy (a) with entailment (e), i.e., the training ranking is given by (e+a)/2. When trained in combination with the full objective $(\mathcal{L}_{rank} + \mathcal{L}_{BCE})$ the ranker shows an increase of 10 AUROC points. Highlighting the benefit of training the Utility Ranker to predicting QA accuracy for input passages. Interestingly, when we drop the ranking loss (i.e., last line of Table 9) there is a drop in performance. On one hand, the ranking loss enables the comparison of pairs of passages and thus the number of training instances is higher. On the other hand, the en-tailment -based ranking signal might help the final model to learn features useful for more accurate passage utility prediction.

D.2 UNCERTAINTY ESTIMATION RESULTS

Table 10 and 11 shows retrieval augmented QA performance on the development set for the target QA models. Table 12 shows performance of uncertainty quantification approaches. We report AURAC and AUROC as well as the percentage out of incorrect cases where the QA models produce an answer acknowledging the lack of knowledge to answer.

We report the following additional uncertainty estimation methods. Perplexity (PPL) computed as $PPL(y_{\mathcal{M}}, x, R; \mathcal{M}) = \exp \{-\frac{1}{|y_{\mathcal{M}}|} \sum_{t=1}^{|y_{\mathcal{M}}|} P_{\mathcal{M}}(y_t | y_{1..t-1}, x, R)\}, \text{ i.e., based on the average neg$ ative log-likelihood of the generated tokens. Cluster Assignment (CA) is the variant of SE without $answers' probabilities where <math>\hat{P}_m(x, \mathcal{M})$ is approximated from the number of answers in the cluster. CA values are very close to SE values.

1026	You need to check whether the prediction of a question-
1027	answering system to a question is correct. You should make
1028	the judgment based on a list of ground truth answers provided
1029	to you. Your response should be "correct" if the prediction is
1030	correct or "incorrect" if the prediction is wrong.
1031	
1032	Question: Who authored The Taming of the Shrew (published in
1033	2002)?
1034	Ground truth: ["William Shakespeare", "Roma Gill"]
1035	Prediction: W Shakespeare
1036	Correctness: correct
1037	Question: Who suthered The Terring of the Shrow (nublished in
1038	Question. Who authored The Taining of the Sillew (published in 2002)?
1039	Ground truth: ["William Shakespeare" "Roma Gill"]
1040	Prediction: Roma Gill and W Shakespeare
1041	Correctness: correct
10/12	
1042	Ouestion: Who authored The Taming of the Shrew (published in
10//	2002)?
1044	Ground truth: ["William Shakespeare", "Roma Gill"]"
1045	Prediction: Roma Shakespeare
1046	Correctness: incorrect
1047	
1048	Question: What country is Maharashtra Metro Rail Corporation
1049	Limited located in?
1050	Ground truth: ["India"]
1051	Prediction: Maharashtra
1052	Correctness: incorrect
1053	
1054	Question: What's the job of Song Kang-ho in Parasite (2019)?
1055	Ground truth: ["actor"]
1056	Vim family
1057	Correctness: correct
1058	Concerness. concer
1059	
1060	Ouestion: Which era did Michael Oakeshott belong to?
1061	Ground truth: ["20th-century philosophy"]
1062	Prediction: 20th century."
1063	Correctness: correct
1064	
1065	Question: Edward Tise (known for Full Metal Jacket (1987)) is
1066	in what department?
1067	Ground truth: ["sound department"]
1068	Prediction: 2nd Infantry Division, United States Army
1069	Correctness: incorrect
1070	
1070	Question: What wine region is Finger Lakes AVA a part of?
1071	Ground truth: ["New York wine"]
1072	Prediction: Finger Lakes AVA
1073	Correctness: incorrect
1074	Question: [question]
1075	Question: [question]
1076	Drediction: [output]
1077	Correctness:
1078	
1079	

Table 7: Prompt for accuracy evaluation.

1081		
1082		Nb./Type of Inference Call at Test Time
1083	PPL	1 LLM-G
1084	MSP	1 LLM-G
1085	PMI	2 LLM-G
1026	p(true)	(N+1) LLM-G + 1 LLM-E
1000	Regular Entropy	(N+1) LLM-G
1087	Cluster Assignment	(N+1) LLM-G + $N(N-1)/2$ LLM-E
1088	Semantic Entropy	(N+1) LLM-G + $N(N-1)/2$ LLM-E
1089	Ans.Len	1 LLM-G
1090	Retriever Score	0 LLM-G
1091	Utility Ranker	R Bert-F

1093Table 8: Number and type of inference call required to estimate answer uncertainty for a given user1094question x. LLM-G means inference with the retrieval augmented QA model, i.e., a forward pass1095with the prompt including the set of R retrieved passages and the question to generate an answer.1096LLM-E is inference with an evaluation model, e.g., a forward pass to ask a LLM for correctness in1097p(true) or a forward pass with an entailment model in the Semantic Entropy method. Bert-F is an1098inference call to predict passage utility for a passage p in R and user question x.

Table 9: Uncertainty Estimation by the Utility Ranker trained with variants of the training objective.
We report AUROC and AURAC for the Utility Ranker for the three target QA models (GEMMA2-9B, LLAMA3.1-8B, and MISTRAL-7B-V0.3) on Natural Questions development data.

	Gemm	GEMMA2-9B		3.1-8B	MISTRAL-7B-V0.3	
	AUROC	AURAC	AUROC	AURAC	AUROC	AURAC
$\mathcal{L}_{rank}, (e+a)/2 +$	\mathcal{L}_{BCE} 0.77	0.76	0.77	0.76	0.79	0.76
$\mathcal{L}_{rank}, (e+a)/2$	0.67	0.70	0.66	0.70	0.69	0.70
$\mathcal{L}_{rank}, (a)$	0.62	0.67	0.64	0.68	0.67	0.69
$\mathcal{L}_{rank}, (e)$	0.67	0.70	0.64	0.68	0.64	0.67
\mathcal{L}_{BCE}	0.76	0.74	0.75	0.74	0.77	0.74

Table 10: Target QA models performance on test sets with |R| = 5. Model based accuracy AccLM (column header ALM) is accuracy computed by Qwen2-72B-Instruct.

	Natura	NaturalQuestions		NaturalQuestions TriviaQA WebQuestions		uestions	SQuAD		PopQA		RefuNQ	
	Acc	ALM	Acc	ALM	Acc	ALM	Acc	ALM	Acc	ALM	Acc	ALM
Gemma2.9B	0.46	0.61	0.73	0.78	0.40	0.64	0.41	0.58	0.49	0.51	0.26	0.35
LLAMA-3.1-8B	0.47	0.60	0.71	0.77	0.44	0.66	0.41	0.56	0.48	0.49	0.27	0.37
MISTRAL-7B-V0.3	0.47	0.58	0.71	0.75	0.47	0.66	0.40	0.57	0.52	0.49	0.27	0.35

Table 11: Target QA models performance on the development sets with |R| = 5. (Acc) is rule based accuracy as used in previous work, (AccLM) is accuracy computed by Qwen2-72B-Instruct.

	Natural Questions		Tri	viaQA	WebQ	Questions	SQuAD	
	Acc	AccLM	Acc	AccLM	Acc	AccLM	Acc	AccLM
GEMMA2.9B	0.45	0.62	0.73	0.79	0.45	0.67	0.37	0.58
LLAMA-3.1-8B	0.46	0.60	0.71	0.77	0.52	0.68	0.38	0.58
MISTRAL-7B-V0.3	0.46	0.60	0.71	0.76	0.53	0.69	0.36	0.57

Table 12: Answer uncertainty estimation for QA models GEMMA2-9B, LLAMA-3.1-8B, and
MISTRAL-7B-V0.3 on NaturalQuestions, TriviaQA, WebQuestions, and SQuAD development sets.
We report AUROC and AURAC. Refusal % is the percentage out of the total incorrect answers where
the model acknowledges uncertainty by expressing its lack of knowledge in the generated answer.

1147		Natural	Questions	Trivi	aOA	WebOu	estions	SOu	AD
1148		AUROC	AURAC	AUROC	AURAC	AUROC	AURAC	AUROC	AURAC
1149		nence	nonne	nenoe	GEMMA	A2-9B	nonne	nenoe	nonne
1150	PPL	0.67	0.69	0.61	0.80	0.63	0.70	0.65	0.66
1151	MSP	0.69	0.70	0.66	0.81	0.64	0.70	0.66	0.66
1150	PMI	0.49	0.59	0.42	0.71	0.49	0.63	0.46	0.55
1152	p(true)	0.73	0.73	0.76	0.85	0.73	0.75	0.70	0.69
1153	Regular Entropy	0.70	0.69	0.66	0.81	0.65	0.70	0.68	0.68
1154	Cluster Assignment	0.70	0.70	0.67	0.81	0.65	0.70	0.65	0.66
1155	Semantic Entropy	0.71	0.71	0.65	0.80	0.65	0.71	0.65	0.66
1156	Ans.Len	0.63	0.66	0.62	0.79	0.61	0.69	0.60	0.64
1157	Retriever Score	0.59	0.65	0.62	0.80	0.50	0.62	0.67	0.68
1150	Utility Ranker	0.75	0.74	0.79	0.86	0.74	0.77	0.82	0.77
1150	Refusal %	5%		5%		0.7%		3%	
1159					LLAMA.	3.1-8B			
1160	PPL	0.75	0.75	0.80	0.85	0.68	0.73	0.71	0.70
1161	MSP	0.79	0.77	0.83	0.86	0.69	0.73	0.72	0.70
1162	PMI	0.61	0.68	0.56	0.75	0.55	0.67	0.55	0.60
1163	p(true)	0.79	0.77	0.89	0.88	0.72	0.75	0.69	0.69
116/	Regular Entropy	0.81	0.78	0.82	0.86	0.69	0.74	0.75	0.72
1105	Cluster Assignment	0.77	0.75	0.82	0.85	0.72	0.75	0.75	0.72
1105	Semantic Entropy	0.76	0.75	0.84	0.86	0.71	0.75	0.76	0.73
1166	Ans.Len	0.63	0.67	0.66	0.79	0.61	0.69	0.56	0.60
1167	Retriever Score	0.57	0.65	0.62	0.78	0.49	0.64	0.67	0.67
1168	Utility Ranker	0.79	0.77	0.81	0.85	0.77	0.79	0.83	0.76
1169	Refusal %	2%		1%		0.7%		2%	
1170	DDI	0.65	0.60	1	MISTRAL-	/B-V0.3	0.70	0.66	0.65
1170	PPL	0.65	0.69	0.65	0.80	0.62	0.70	0.66	0.65
1171	MSP	0.70	0./1	0.74	0.82	0.67	0.73	0.72	0.68
1172	PINII (trues)	0.49	0.60	0.57	0.70	0.50	0.08	0.54	0.58
1173	p(true)	0.75	0.71	0.60	0.80	0.69	0.75	0.70	0.07
1174	Cluster Assignment	0.05	0.09	0.00	0.80	0.05	0.71	0.70	0.00
1175	Semantic Entropy	0.71	0.72	0.70	0.82	0.71	0.75	0.75	0.09
1176		0.72	0.72	0.77	0.85	0.71	0.74	0.75	0.70
1177	Retriever Score	0.05	0.00	0.09	0.80	0.04	0.72	0.00	0.04
11//	Utility Ranker	0.59	0.03	0.01	0.77	0.38	0.09	0.04	0.03
1178	Refusal %	1%	0./7	0.25%	0.07	3%	0.77	0.5%	U. / 2
1179	Kerusai 70	170		0.2570		570		0.570	