

---

# Towards Uncertainty-Aware Language Agent

---

Jiuzhou Han<sup>1</sup> Wray Buntine<sup>2</sup> Ehsan Shareghi<sup>1</sup>

## Abstract

While Language Agents have achieved promising success by placing Large Language Models at the core of a more versatile design that dynamically interacts with the external world, the existing approaches neglect the notion of uncertainty during these interactions. We present the Uncertainty-Aware Language Agent (UALA), a framework that orchestrates the interaction between the agent and the external world using uncertainty quantification. Compared with other well-known counterparts like ReAct, our extensive experiments across 3 representative tasks (HotpotQA, StrategyQA, MMLU) and various LLM sizes demonstrates that UALA brings a significant improvement of performance, while having a substantially lower reliance on the external world (i.e., reduced number of tool calls and tokens). Our analyses provide various insights including the great potential of UALA compared with agent fine-tuning, and underscoring the unreliability of verbalised confidence of LLMs as a proxy for uncertainty.<sup>1</sup>

## 1. Introduction

Language Agents (Sumers et al., 2023; Zhou et al., 2023; Xi et al., 2023) utilise Large Language Models (OpenAI, 2023; Touvron et al., 2023; Anil et al., 2023) to interact with the external world (e.g., through tools) and to process collected observations towards solving a task. Having achieved improvements in previously challenging reasoning tasks (Yao et al., 2023; Gao et al., 2023; Lin et al., 2023a), these agents autonomously obtain new knowledge from the world, and leverage memory or self-refinement mechanisms to iteratively improve their reasoning trajectories (Shinn et al., 2023; Gou et al., 2023; Majumder et al., 2023).

Current language agent designs primarily define the role of an LLM as the planner, moderator, or aggregator of observa-

tions and knowledge collected from the external world (i.e., via tools and APIs), overlooking the potential of LLMs to tackle numerous tasks without depending on any external tool and resource (Shinn et al., 2023; Lu et al., 2023; Yin et al., 2023; Patil et al., 2023). While exploiting external mechanisms has an advantage in terms of verifiability of outputs (e.g., output extracted from a retrieved Wikipedia article provides a verifiable reference) or improvement on accuracy (e.g., using a Python program to do mathematical calculation instead of relying on the LLM), this design is inherently wasteful in how it utilises the LLM’s internal implicit knowledge. In fact, the field of NLP has shown great interest in probing pretrained language models to highlight various types of knowledge they embody directly in their weights (e.g., see Petroni et al. (2019); Collier et al. (2022); Shu et al. (2023) and references therein).

We postulate that a more effective design for language agents should have a better interplay between the implicit knowledge encoded in LLM’s weight and the explicit knowledge of the external world. Drawing upon principles from the decision theory, an intelligent agent should possess a reliable mechanism for *measuring uncertainty* of its own action or the observations from the environment (Bacchus et al., 1996). We build on this, and move towards equipping language agents (as a special case of intelligent agents) with the means of measuring uncertainty to regulate their use of external help (e.g., tools, knowledge bases) more efficiently.

In NLP, in general, there has been a great amount of research on uncertainty estimation for generative tasks such as Machine Translation (Ott et al., 2018), Summarization (Zablotskaia et al., 2023; Xu et al., 2020), and data-to-text generation (Xiao & Wang, 2021). Recently, a dedicated space emerged to study the intersection of LLMs and Uncertainty. This space can be divided into approaches applicable to open-box LLMs which rely on token-level logits (Malinin & Gales, 2021; Kuhn et al., 2023), and those applicable to black-box LLMs leveraging multiple response generation as a proxy for model’s confidence (Lin et al., 2023b), or relying on LLM’s verbal self-awareness of confidence (Tian et al., 2023). Nonetheless, uncertainty in the context of language agents is an under-explored space, with our work being the first in this direction.

We propose an Uncertainty-Aware Language Agent frame-

<sup>1</sup>Department of Data Science & AI, Monash University  
<sup>2</sup>College of Engineering and Computer Science, VinUniversity.  
Correspondence to: Jiuzhou Han <jiuzhou.han@monash.edu>.

<sup>1</sup>For code and data: <https://uala-agent.github.io>

work which utilises various existing uncertainty measurement methods as a dynamic switch between LLM’s own trajectory or resorting to external resources (i.e., Search Engine, Wikipedia knowledge base) during the course of a reasoning task. We study the properties of this framework on a large space of LLMs (ChatGPT; OpenAI (2023), and LLaMA2 70B, 13B, 7B; Touvron et al. (2023)), and Language Agent frameworks such as ReAct (Yao et al., 2023), and CRITIC (Gou et al., 2023) on several representative tasks. These tasks include HotpotQA (Yang et al., 2018) a free-form QA task; StrategyQA (Geva et al., 2021) a representative binary QA task; and 57 multiple-choice QA tasks of MMLU (Hendrycks et al., 2021). Our key findings and contributions are as follows:

We propose the first language agent framework which integrates uncertainty in the reasoning trajectories. The proposed design allows for a plug-and-play use of various uncertainty estimation techniques. We highlight across our wide spectrum of experiments that integrating uncertainty not only leads to a significant performance improvement, but also has several practical benefits such as a substantial reduction in external calls and number of output tokens (i.e., less API charge). Our analysis unpacks several additional findings: (1) The gain is significant across different LLM sizes and tasks. (2) There exists a divergence of uncertainty between correct and incorrect responses. (3) The verbalised confidence of LLMs is an unreliable proxy for uncertainty quantification. (4) Given a limited amount of data, leveraging uncertainty leads to a much higher performance improvement compared to fine-tuning language agents.

The approach presented here points at a promising direction and takes the first steps. We hope our work encourages more effective and efficient designs of language agents.

## 2. Related Work

### 2.1. Language Agents

The prominent work of ReAct (Yao et al., 2023) proposes a general language agent framework to combine reasoning and acting (i.e., utilising external tools) with LLMs for solving diverse language reasoning tasks. Reflexion (Shinn et al., 2023) proposes to use the history of failed trials based on ReAct to do self-reflection, inducing better decision-making in subsequent trials of the language agent. FireAct (Chen et al., 2023a) fine-tunes language agent with diverse ReAct task-solving trajectories from multiple tasks and prompting methods, and shows having more diverse fine-tuning data can further improve performance of tool-use. CRITIC (Gou et al., 2023) presents a framework that allows iterative verification and correction of agent’s outputs by interacting with external tools (e.g., Search Engines, Code Interpreters).

As an earlier work in this direction, Toolformer (Schick

et al., 2023) attempts to use different APIs (e.g., Wikipedia, QA, Calculator) in a self-supervision manner and autonomously choose proper APIs to use when solving different tasks. Qiao et al. (2023) propose a two-stage framework which allows the model to learn through feedback derived from tool execution. Gorilla (Patil et al., 2023) is a retrieval-aware system fine-tuned on open-source LLMs for enabling LLMs to interact with tools, which enhances the LLM’s ability to accurately identify the appropriate API for a specific task. Lumos (Yin et al., 2023) is a language agent fine-tuned on open-source LLMs which contains three modules (Planning Module, Grounding Module and Execution Module) for solving complex interactive tasks. Chameleon (Lu et al., 2023) is a plug-and-play compositional reasoning framework, which infers the appropriate sequence of tools to compose and execute in order to solve the problems. PoT (Chen et al., 2022) and PAL (Gao et al., 2023) utilise LLMs to generate programs as reasoning steps and leverage external program interpreter to execute for arithmetic and symbolic reasoning tasks.

Neither of these language agent frameworks consider any role for uncertainty in the agent’s interactions with the external world. Next, we present a brief overview of the intersection of uncertainty and LLMs.

### 2.2. Uncertainty in Generation with LLMs

A common practice in leveraging uncertainty during language generation is via sampling or decoding approaches which do not measure the uncertainty directly, and instead rely on the stochasticity over the prediction space along with a form of aggregation approach such as majority voting. Self-Consistency (Wang et al., 2023) is a widely known sampling method, which takes majority voting over multiple sampling outputs. Since it is based on exact matching, it is most useful in tasks for mathematical reasoning or fact-based question answering tasks. Minimum Bayes-Risk Decoding (Kumar & Byrne, 2004) is also based on multiple candidate outputs sampled from a model, and it selects the candidate with lowest expected risk. Suzgun et al. (2023) highlight the connection between Self-Consistency and MBR and showcase the benefit of MBR using a soft alignment function (i.e., semantic similarity) in tasks involving open-ended, longer free-form text generation (e.g., summarization, translation). For a comprehensive review of sampling and decoding methods in NLG, we refer the readers to Wiher et al. (2022).

We focus on directly leveraging uncertainty estimation in free-form question answering tasks with short answers. Uncertainty estimation in free-form NLG remains a challenge for LLMs due to the diversity of the outputs. In general, recent research on this area can be classified as:

**Logits- or Entropy-based Methods.** Malinin & Gales

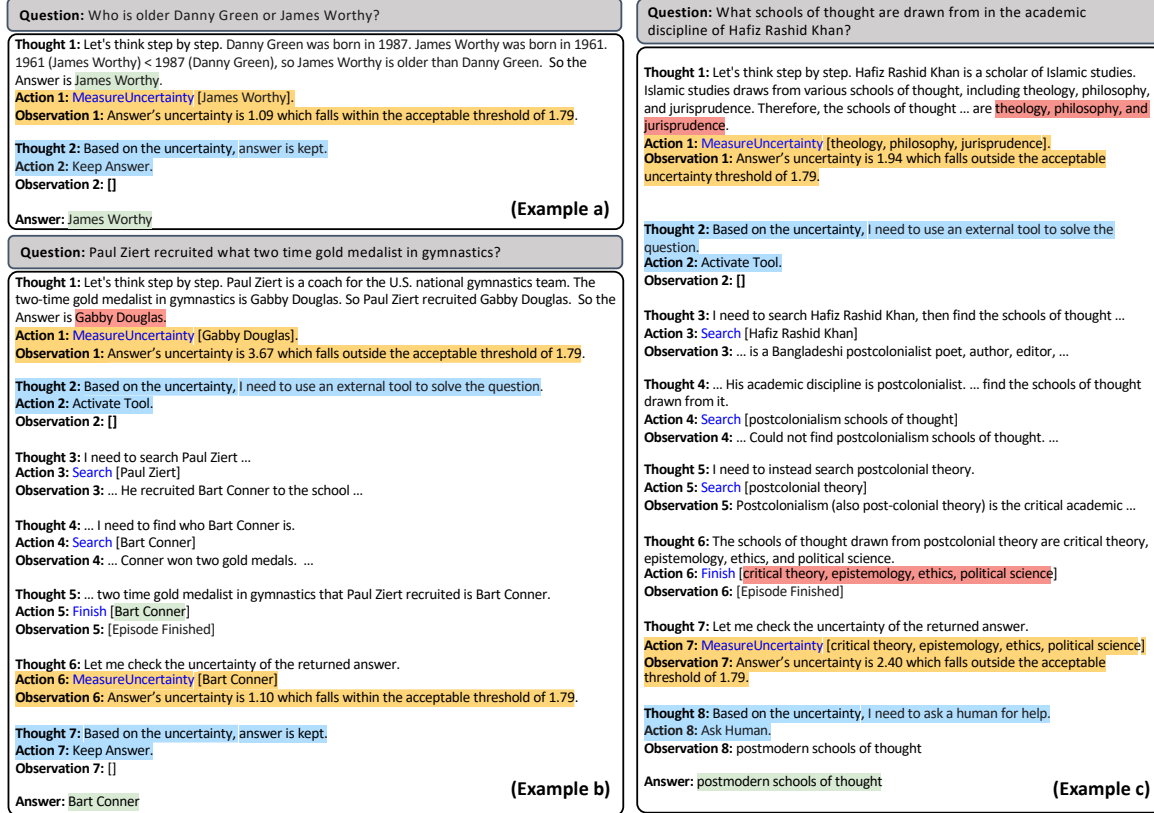


Figure 1. Examples of single-inference UALA trajectories. Example (a) illustrates the trajectory where CoT answer falls inside the certainty region. Example (b) is the trajectory where CoT is too uncertain and tool is activated, arriving at a final response which falls in the acceptable certainty region (denoted by UALA-S and UALA-M in our results). Example (c) is the trajectory where both CoT and tool-generated responses are considered uncertain, and the agent asks help from human (denoted by UALA-S+Oracle in our results).

(2021) calculate the accumulative predictive entropies over multiple generations to estimates the uncertainty. But this method ignores the semantic equivalence phenomenon in generated texts. Kuhn et al. (2023) propose semantic entropy to solve the difficulty of semantic equivalence in uncertainty estimation. Duan et al. (2023) mitigate inequality biases by examining the relevance of each token and sentence, and reassigning attention when estimating uncertainty. Yang et al. (2023) implement multiple uncertainty calculation methods based on the token-level logit output. Ren et al. (2023) utilise uncertainty alignment of LLM-based planners in solving complex step-by-step planing problems for robots via conformal prediction.

**Prompt-based Methods.** Tian et al. (2023); Kadavath et al. (2022); Lin et al. (2022) prompt LLMs to directly express the confidence or uncertainty along with their generations. Lin et al. (2023b) propose to estimate the uncertainty based on the similarities of multiple generations. We refer the readers for a more comprehensive overview to read Fadeeva et al. (2023) and references therein.

We take both type of uncertainty estimation methods<sup>2</sup> into account and compare the performance on various settings. More details are described in Section 3.

### 3. Uncertainty-Aware Language Agent (UALA)

We present, UALA, a framework that integrates uncertainty in language agent’s cycle of Thought, Action, and Observation. The uncertainty moderates the interaction between the LLM and the external world, facilitating a more effective and efficient dynamic. Given an input question, the agent’s potentially goes through 3 modes of attempt, each involving a thought, an action, and an observation (similar to ReAct). Figure 1 presents real trajectories for each mode of attempt.

In our approach, as illustrated in Example (a), we initially generate a response to a query using standard or Chain-of-Thought (CoT; Wei et al. (2022)) prompting. This response

<sup>2</sup>Since our focus is not on improving the performance of the uncertainty estimation methods, we didn’t aim for solving the potential issues (i.e., semantic equivalence, inequality biases) when using logits-based methods.

is then assessed by an uncertainty estimation module, which determines the subsequent action between two courses: either accepting the response as-is, or employing external tools for a more certain solution, as shown in Example (b). This process is similarly applied when using tool-activated responses. If the uncertainty remains high post-tool engagement, the agent seek assistance from an external source such as a human, as depicted in Example (c). This protocol ensures a more robust and reliable response generation by actively addressing uncertainty at each step. Next, we describe our uncertainty estimation mechanisms.

### 3.1. Uncertainty Estimation

Uncertainty estimation methods are broadly categorised into two types: single-inference based and multi-inference based. Single-inference uncertainty estimation calculates the uncertainty based on one output, necessitating access to the token log-probabilities within that output. Multi-inference uncertainty estimation computes the uncertainty based on a set of outputs, eliminating the need for individual token log-probabilities.

#### 3.1.1. SINGLE-INFERENCe UNCERTAINTY ESTIMATION

In single-inference uncertainty estimation, the methods employed vary based on the answer being a single-token (e.g., yes or no) or free-form (multi-token) format.

**Free-form Answer.** Given an output (i.e., Thought, Observation) containing an answer, the answer  $Y$  comprising  $n$  tokens  $[y_1, y_2, \dots, y_n]$  along with their corresponding token log-probabilities  $[p_1, p_2, \dots, p_n]$  is extracted. Following the previous work (Yang et al., 2023), we first apply a softmax function to the token log-probabilities to obtain  $[z_1, z_2, \dots, z_n]$ . Subsequently, we explore five methods of calculating the uncertainty  $u$  based on the probabilities. These methods are summarised in Table 1.

**Single-token Answer.** For single-token answers, applying a softmax function to their log-probability is not meaningful due to the presence of only one token. Instead, the uncertainty is determined directly by taking the absolute value of the token’s log-probability. In this scenario, a higher probability of the token (approaching 1) results in its log-probability nearing 0, indicating low uncertainty. Conversely, a lower token probability signifies higher uncertainty. Therefore, for a given answer token  $y$  with log-probability  $p$ , the uncertainty is defined as  $u = |p|$ .

#### 3.1.2. MULTI-INFERENCe UNCERTAINTY ESTIMATION

Multi-inference uncertainty estimation calculates the uncertainty of an answer based on multiple outputs from an LLM. This process begins with obtaining a primary answer,  $a_*$ , using the LLM’s greedy decoding method. Subsequently, the

Table 1. Five methods of calculating the uncertainty  $u$  of a free form output of length  $n$ .

Method	Formula
Minimum	$u = -\log(\min(z_1, z_2, \dots, z_n))$
Average	$u = -\log(\text{Avg}(z_1, z_2, \dots, z_n))$
Normalised Product	$u = -\log(z_1 \times z_2 \times \dots \times z_n)^{\frac{1}{n}}$
Log-sum	$u = -\sum_{i=1}^n \log(z_i)$
Entropy	$u = -\sum_{i=1}^n z_i \cdot \log(z_i)$

LLM is prompted to sample  $n = 9$  additional answers with temperature 0.7, forming answers set  $A = \{a_1, a_2, \dots, a_n\}$ . The uncertainty of the primary answer is then calculated based on how frequently it appears within set  $A$ . Essentially, the less frequently  $a_*$  occurs in  $A$ , the higher the uncertainty attributed to it, and vice versa. Specifically,

$$u = \frac{1}{n} \sum_{i=1}^n \delta(a_i, a_*), \quad (1)$$

where  $\delta(a_i, a_*)$  is an indicator function that equals to 1 if  $a_i$  is not equal to  $a_*$ , and 0 otherwise.

#### 3.1.3. UNCERTAINTY THRESHOLD

The decision to accept an answer or resort to alternative mechanisms hinges on the level of uncertainty associated with that answer. To this end, we propose different ways of setting the uncertainty threshold for single-inference and multi-inference uncertainty estimation.

**Single-inference.** For single-inference, we adopt a subset of the training data to create a calibration set. We first use standard prompting or CoT prompting to query an LLM to get the answers and collect the correctly answered questions as the calibration set. We show this subset could be as small as 200 examples (see Subsection 5.6). We use a single-inference uncertainty estimation method to calculate the uncertainty of each answer in the calibration set. The same single-inference uncertainty estimation method is used for the calibration set and test set. Based on the answer uncertainty in the calibration set, we apply three methods to estimate the threshold.<sup>3</sup> Specifically, the Max method utilises the maximum uncertainty in the calibration set as the threshold, while the Mean uses the average uncertainty, and the last approach uses a Quantile (Ren et al., 2023) value of uncertainty in calibration set as the threshold.

**Multi-inference.** For the estimation of uncertainty threshold in multi-inference setting, we do not need to create a calibration set. We simply take the average uncertainty of the answers in the test set as the uncertainty threshold.

<sup>3</sup>It is possible to learn different thresholds for different modes of engagement. However, for cost reasons, we apply the same estimated threshold when deciding to accept the LLM (as in Example a) or Tool-activated responses (as in Example b and c of Figure 1).



Figure 1 provides examples of the UALA trajectories. After obtaining the uncertainty threshold, for an answer where the uncertainty is above the threshold, we activate the alternative mode (i.e., Tool Activation as in Example b, or asking from human if tool-use result is still uncertain as in Example c).

## 4. Experiments

We explore three types of question answering tasks: HotpotQA (Yang et al., 2018), StrategyQA (Geva et al., 2021), MMLU (Hendrycks et al., 2021) and utilise Wikipedia and Google Search and external tools.

**HotpotQA** (Yang et al., 2018) is a challenging question answering dataset which requires multi-hop reasoning over some Wikipedia passages. The answer of the questions in HotpotQA is free-form with an arbitrary length. We randomly select 500 questions from the training set for the creation of calibration set, and 500 questions from the dev set for evaluation using the same random seed as ReAct (Yao et al., 2023). **StrategyQA** (Geva et al., 2021) is an open-domain question answering dataset which requires implicit reasoning steps. It only contains binary answers (yes or no). We randomly select 500 questions from the training set to create the calibration set. For the evaluation, we use the same dev set used in FireAct (Chen et al., 2023a) containing 229 questions. **MMLU** (Hendrycks et al., 2021) is a multi-choice question answering dataset in various domains which covers 57 tasks including elementary mathematics, US history, computer science, law, and more. This dataset is to evaluate model’s academic and professional understanding. Each question contains four options and only one option is correct. To cover all different topics, we randomly sample 10 questions for each task from dev set to create the calibration set, and similarly sampled 10 question for each task (total of 570 questions) from test set for evaluation.

**Wikipedia.** Following the Wikipedia web API setting used in ReAct (Yao et al., 2023), the tool prompt contains three actions. (1) **search[entity]**, which returns the initial five sentences from the entity wiki page if it exists. Otherwise, it suggests the top five similar entities from the Wikipedia search engine. (2) **lookup[string]**, which functions akin to the Ctrl+F in a browser, returning the subsequent sentence on the page containing the specified string. (3) **finish[answer]**, which contains an answer and finishes the current task. This tool is used to retrieve some relevant knowledge of the entities in the query, which helps the agent to do the reasoning based on the external knowledge and its internal knowledge. The Wikipedia tool is used for HotpotQA and StrategyQA datasets.

**Google Search.** Following FireAct (Chen et al., 2023a), the Google search tool we use is based on SerpAPI<sup>4</sup>. The

tool prompt contains two actions: (1) **search[query]**, which searches a query on Google and returns a short snippet containing the answer of the query. The returned snippet is the first existent item from “answer box”, “answer snippet”, “highlight words”, or “first result snippet”, which ensures the response is short and relevant. (2) **finish[answer]**, which returns the answer and finishes the episode. The universality of Google search is better than Wikipedia. It can search any query more than an entity to get the external knowledge, which is suitable for more complex reasoning tasks. Therefore, we adopt Google Search tool for MMLU task.

### 4.1. Baselines and Experimental Setup

**Baselines.** We conduct experiments on ChatGPT<sup>5</sup> and LLaMA2 (Touvron et al., 2023). For ChatGPT, since we need access to the token log-probability for uncertainty estimation, we use gpt-3.5-turbo-instruct which is a Completion API. For LLaMA2, we adopt LLaMA-70B as our main backbone, and also compare the performance of different sizes of LLaMA2 (7B and 13B) in Section 5.2. We compare UALA with **Standard** (Brown et al., 2020), **Chain-of-Thought** (CoT; Wei et al. (2022))<sup>6</sup>, **Self-Consistency** (SC; Wang et al. (2023))<sup>7</sup>, **ReAct** (Yao et al., 2023)<sup>8</sup>, and **Backoff** which backs off to response by Standard or CoT when the Tool Activation trajectories do not arrive at an answer within a given number of steps. For UALA, when the result includes human oracle help, it is marked by **+Oracle**. We also compare with the **CRITIC** (Gou et al., 2023) in Appendix B. The prompts and few-shot examples are included in Appendix I. **Experimental Setup.** In all our experiments, unless stated otherwise, we use Entropy for free-form single-inference uncertainty estimation for HotpotQA and 0.9 Quantile for single-inference uncertainty threshold estimation, which perform best in the training set (we demonstrate the effect of Quantile value in Appendix A). We compare different free-form single-inference uncertainty estimation methods in Section 5.7. The multi-inference uncertainty estimation is based on 9 samples generated using temperature 0.7. For HotpotQA, we use 6-shot learning and for StrategyQA and MMLU we use 4-shot learning. For HotpotQA, since CoT performs better than Standard, we use CoT as the base prompt method. For StrategyQA and MMLU, we use Standard as the base prompt method. Following the previous work (Chen et al., 2023a; Yao et al., 2023), we evaluate using the exact match (EM). In addi-

<sup>5</sup><https://openai.com/blog/gpt-3-5-turbo-fine-tuning-and-api-updates>

<sup>6</sup>We use “Let’s think step by step.” from Kojima et al. (2022).

<sup>7</sup>Generates multiple outputs using a sampling method and takes the majority result as the final answer. For the sampling method, we use temperature 0.7 to generate 9 samples.

<sup>8</sup>Combines the CoT with actions to retrieve relevant knowledge from external tools, synergizing reasoning and acting.

<sup>4</sup><https://serpapi.com>

Table 2. Results of three question answering datasets on ChatGPT and LLaMA2-70B. The metric is exact match (EM) and the number in the bracket represents the number of tool calls. UALA-S denotes using single-inference uncertainty estimation method and UALA-M denotes using multi-inference uncertainty estimation method. **Bold** shows the best result for each column (oracle results excluded).

Methods		ChatGPT			LLaMA2-70B		
		HotpotQA	StrategyQA	MMLU	HotpotQA	StrategyQA	MMLU
Baselines	Standard	29.8(0)	57.6(0)	69.0(0)	30.0(0)	65.9(0)	64.7(0)
	CoT	34.8(0)	55.9(0)	49.1(0)	35.6(0)	63.8(0)	39.3(0)
	Self-Consistency	39.4(0)	58.5(0)	70.0(0)	37.4(0)	67.7(0)	67.2(0)
	ReAct	32.0(2,114)	55.5(709)	55.8(1,824)	32.4(2,094)	58.1(890)	30.7(2,808)
	ReAct+Backoff	35.4(2,114)	61.6(709)	59.8(1,824)	37.0(2,094)	66.8(890)	62.8(2,808)
Ours	UALA-S	38.2(403)	65.5(134)	69.8(662)	36.4(350)	69.0(298)	56.7(1354)
	UALA-S+Backoff	39.2(403)	66.4(134)	71.4(662)	37.3(350)	71.6(298)	69.8(1354)
	UALA-S+Oracle	41.4(403)	67.5(134)	75.8(662)	42.2(350)	80.1(298)	70.5(1354)
	UALA-M	39.4(1,199)	63.8(234)	67.9(641)	36.4(1,302)	66.4(492)	56.7(1,196)
	UALA-M+Backoff	<b>41.2(1,199)</b>	<b>66.9(234)</b>	<b>72.2(641)</b>	<b>40.4(1,302)</b>	<b>72.2(492)</b>	<b>70.3(1,196)</b>

tion, we also report the total number of tool calls in the bracket following the EM. More tool calling results in a longer inference (see Section 5.3).

## 4.2. Main Results

Our main results are demonstrated in Table 2 (for relative improvement details, refer to Table 12 of Appendix). Some key observations:

(1) CoT outperforms Standard on HotpotQA, while Standard excels on StrategyQA and MMLU. Self-Consistency consistently enhances results across three datasets and two LLMs. ReAct, when used for every instance, underperforms Standard/CoT/Self-Consistency. With the integration of backoff, ReAct+Backoff shows improvement but is generally still falls behind Self-Consistency, highlighting the benefit of SC’s sampling and majority voting as a proxy for capturing uncertainty.

(2) UALA-S significantly better ReAct’s performance, cutting tool use by over half, and surpasses Standard/CoT across all datasets. UALA-M achieves similar performance to UALA-S but with increased tool use. UALA-S+Backoff outperforms ReAct+Backoff and often exceeds Self-Consistency with UALA-M+Backoff delivering the best results in all settings on three datasets.

(3) The largest gain in improvement by UALA is observed for HotpotQA (free-form), followed by StrategyQA (binary), and MMLU (multiple choice). This is expected as the free-form response space is much larger and diverse, compared with MCQ type of questions. In Section 5.1, we highlight the difference in gain could be explained in terms of the amount of uncertainty divergence between correct and incorrect answers in each task.

(4) The average (single-inference and multi-inference) EM improvement for ChatGPT with LLaMA2-70B compared to

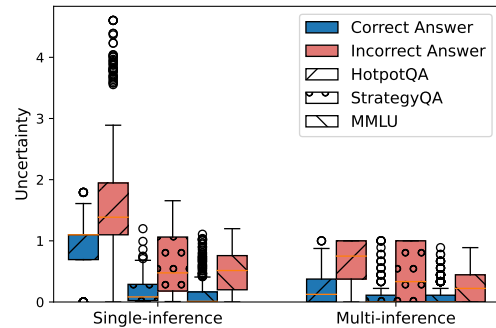


Figure 2. The visualisation (boxplot) of uncertainty range for correct and incorrect answers of three datasets on ChatGPT.

Standard/CoT results: ChatGPT gains 11.7% and LLaMA2-70B gains 8.9%. This could be an indication that ChatGPT is likely to produce better-calibrated probability estimates, leading to a more reliable uncertainty estimation on training set that generalizes to test set. This could be an artefact of the two models’ difference in size and training protocol. For further analysis of the inference cost see Section 5.3.

(5) The results from UALA-S+Oracle underscore an additional aspect of the value of uncertainty. This feature is particularly crucial in sensitive domains, as it can deter the agent from generating incorrect responses. Instead of risking an erroneous answer, the agent defers to human (we simulate this by using gold answer) when the response uncertainty is still high after tool activation.

## 5. Analysis

### 5.1. Correct vs. Incorrect Answer Uncertainty

Figure 2 illustrates the uncertainty visualisation on single-inference and multi-inference answers on HotpotQA (500, dev set), StrategyQA (229, dev set), and MMLU (570, test set) with ChatGPT. In both single-inference and multi-inference settings, correct answers consistently exhibit

Table 3. HotpotQA Results on LLaMA2-7B and LLaMA2-13B.

Methods	LLaMA2-7B	LLaMA2-13B
Standard	21.2 (0)	23.4 (0)
CoT	21.8 (0)	24.8 (0)
Self-Consistency	22.0 (0)	26.4 (0)
ReAct	16.8 (2,324)	19.8 (2,371)
ReAct+Backoff	23.8 (2,324)	25.0 (2,371)
UALA-S	22.8 (479)	25.6 (539)
UALA-S+Backoff	24.2 (479)	27.3 (539)
UALA-M	22.4 (1,446)	26.6 (1,557)
UALA-M+Backoff	<b>24.8 (1,446)</b>	<b>28.6 (1,557)</b>

lower uncertainty compared to incorrect ones. This difference is statistically significant (see Table 11 of Appendix). When calculating the difference between the average uncertainty of correct and incorrect answers we observe the largest difference to belong to HotpotQA, followed by StrategyQA, and MMLU. This explains why the gain from UALA follows the same pattern in the main results.

## 5.2. Effect of LLM Size

To test UALA’s generalizability on smaller LLMs, we experimented with HotpotQA using LLaMA2-7B and LLaMA2-13B, employing the same settings of baseline and uncertainty-based methods as in the LLaMA2-70B experiments. As reported in Table 3, we observe a similar trend to our main results. It is noteworthy that the gain is sensitive to the backbone LLM size. For example, comparing ReAct and UALA-S, the gain is 35% on 7B, shrinks to 29% for 13B, and then to 19% for the 70B. This is anticipated as larger model capacity could allow for a more confident learning of the data distribution during training.

## 5.3. Inference Cost

Figure 3 reports the number of output tokens and tool calls per method. ReAct consumes substantially more ( $5\times$ ) output tokens than CoT. Compared with ReAct, UALA-S reduces the number of output tokens by more than 65%. UALA-M consumes more output tokens as it relies on multiple inference. Both UALA methods can substantially reduce tool calls more than 50% compared with ReAct, making them much more resource-efficient. In Table 4, we demonstrate the average inference time per instance for different methods and LLMs on HotpotQA. Standard and CoT prompting methods do not involve an external tool call, hence faster inference time compared to other methods. As indicated, UALA-S given its selective tool call, has a much lower inference time compared with ReAct. These figures highlight a practical benefit of using uncertainty to reduce the number of token usage and tool calls, while still providing a significant gain in performance.

## 5.4. Verbalised Uncertainty Self-Awareness of LLMs

A recent work (Tian et al., 2023) reports that RLHF-LLMs (e.g., ChatGPT) are able to directly verbalise well-

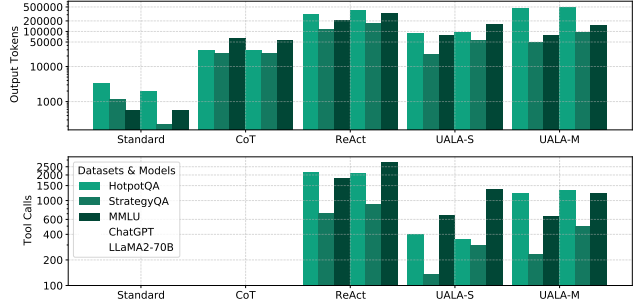


Figure 3. The number of output tokens and tool calls per method.

Table 4. The average inference time per instance (seconds/iteration) of different methods for HotpotQA. The inference of LLaMA2 is done on a single A40 GPU.

	Standard	CoT	ReAct	UALA-S
ChatGPT	0.5s/it	1s/it	12s/it	3s/it
LLaMA2-70B	50s/it	50s/it	180s/it	70s/it
LLaMA2-13B	25s/it	25s/it	120s/it	45s/it
LLaMA2-7B	20s/it	20s/it	100s/it	35s/it

calibrated confidence score along with the predicted answer by prompting. This uncertainty estimation method can also be used in UALA. To verify this idea, we use their prompt templates modified for our task (See Appendix J) to generate the answer with a corresponding confidence score ranging from 0 to 1. When the confidence of an answer is low (i.e., uncertainty is high), we follow the similar protocol to our approach. For the confidence threshold, we demonstrate the results of using threshold value from 0.1 to 0.9. The result of HotpotQA on ChatGPT is shown in Table 5. When using confidence threshold value below 0.8, both of these two confidence-based methods achieve better results with a small number of tool-use than baselines (CoT and ReAct), but the improvement still falls behind our uncertainty-estimation methods (Table 2). We observe that when the model verbalises confidence, there is a large number of answers with confidence above 0.5, leading to degradation of performance after 0.7. This observation underscores that while the expressed confidence of LLMs is very accessible and appealing, in practice it remains as an unreliable indicator of answer uncertainty.

## 5.5. Language Agent Fine-tuning vs. UALA

We demonstrate the comparison between UALA-S and fine-tuning language agents following the FireAct setting (Chen et al., 2023a). For ChatGPT, we use the official GPT-3.5-Turbo fine-tuning API; for LLaMA2-70B, we use LoRA (Hu et al., 2022). To have a side-by-side comparison, we use the same 500 training samples used for the calibration set, to construct the fine-tuning data. Mimicking the FireAct setting, we ran the 500 examples using ReAct with ChatGPT, and collected the successful trajectories as the training data

Table 5. The results of HotpotQA on ChatGPT using LLM’s verbalised confidence (Tian et al., 2023) instead of uncertainty estimation. The numbers are reported as exact match (tool calls).

Confidence	0.1	0.3	0.5	0.7	0.9
Verbal	37.2 (29)	37.4 (69)	37.2 (87)	37.2 (97)	34.4 (1,109)
+Backoff	37.2 (29)	37.5 (69)	37.3 (87)	37.3 (97)	36.8 (1,109)

for FireAct. This amounted to 162 training examples for HotpotQA and 283 for StrategyQA. In addition, to match the amount of training data as FireAct setting, we also ran an additional 1000 examples to increase the amount of successful training trajectories to 512 for HotpotQA and 567 for StrategyQA. The results of supervised fine-tuning of Language Agents instances are presented in Table 6.

Interestingly, on HotpotQA using 162 training examples, FireAct under-performs the few-shot (6-shots) ReAct agent, while it outperforms ReAct on StrategyQA using 283 training examples. Increasing the amount of training data to 500+ leads to improvement on both LLMs with fine-tuned ChatGPT-based agents outperforming the ReAct counterpart on both datasets. Our method, UALA-S, achieves the best result without any fine-tuning and using only the 500 samples for creating the calibration set. This capitalises an obvious empirical advantage for utilising uncertainty instead of fine-tuning in the presence of small amount of data.

Although uncertainty-aware few-shot learning method has a strong performance, fine-tuning still has its own merit. Fine-tuning is more useful if we have larger training data and computing resources. As mentioned in (Chen et al., 2023a), through fine-tuning, smaller open-source LMs could potentially catch up with stronger LMs on a particular task given enough fine-tuning data.

### 5.6. Effect of Calibration Set’s Size

The single-inference uncertainty threshold estimation is based on a calibration set which consists of correctly answered questions (Section 3.1.3). To probe the effect of calibration set’s size, we vary the set size (200 to 2000 instances) using the same uncertainty estimation setting. The performance has a negligible fluctuation between 37.6-38.2, while the Pearson correlation coefficient of 0.29 and the p-value of 0.39 ( $p < .05$ ) indicate no statistically significant correlation between the calibration set size and the downstream performance. This suggests a calibration set of less than 200 examples could suffice in practice. See Table 9 of Appendix for details.

### 5.7. Comparing Single-inference Estimation Methods

For single-inference estimation of the free-form answer, we investigate the effect of five methods described in Section 3.1.1 and the uncertainty threshold estimation methods

Table 6. Results of FireAct vs. UALA-S on HotpotQA and StrategyQA. The ReAct and UALA-S results are based on 6-shot and the off-the-shelf LLM backbones (no fine-tuning).

Tasks	Methods	Training Size	ChatGPT	LLaMA2-70B
HotpotQA	FireAct	162	27.8	27.8
	FireAct	512	33.8	30.0
	ReAct	No fine-tuning	32.0	32.4
	UALA-S	No fine-tuning	38.2	36.4
StrategyQA	FireAct	283	60.7	63.8
	FireAct	567	64.9	64.6
	ReAct	No fine-tuning	55.5	58.1
	UALA-S	No fine-tuning	65.6	69.0

described in Section 3.1.3. The results (on 500 instances of HotpotQA dev set) are shown in Table 10 of Appendix. Most settings perform better than CoT and ReAct. When using Max and Quantile as the threshold estimation method, there is no significant difference between these five uncertainty estimation methods on both LLMs. For the Mean threshold estimation method, Normalised Product and Log-sum perform better than other three uncertainty estimation methods. It also has more tool-use instances than the other two threshold estimation methods. Overall, the Quantile achieves the best performance on both LLMs.

## 6. Conclusion

We present the Uncertainty-Aware Language Agent (UALA), which improves a language agent’s efficiency in interaction with the external world to answer various QA problems. We first measure the uncertainty of an agent on a small set of examples and establish an acceptable uncertainty threshold. Then the uncertainty estimation function is integrated into a language agent’s trajectory of Think-Act-Observe by measuring the uncertainty of generated answers, and then choosing to accept the answer or to resort to external resources or human. This approach does not require any training and the language agent is prompted by few-shot learning. We find that this technique significantly improves the language agent’s performance on various QA tasks (regardless of the backbone LLM size), and offers a much more resource-efficient solution (substantially lower tool calling, and less number of output tokens). We also show that UALA outperforms fine-tuning language agents in the presence of small amount of data.

There are important limitations on UALA as presented. The approach presented here requires task-specific uncertainty selection and calibration set. Additionally, we select the threshold based on the responses from Standard or CoT prompting of the LLM, which may not be the optimal choice for uncertainty within a trajectory. Nonetheless, integrating uncertainty is a necessary element for designing intelligent language agents, and we believe that UALA can serve as the basis of more sophisticated techniques across domains.



## References

- Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., Chu, E., Clark, J. H., Shafey, L. E., Huang, Y., Meier-Hellstern, K., Mishra, G., Moreira, E., Omernick, M., Robinson, K., Ruder, S., Tay, Y., Xiao, K., Xu, Y., Zhang, Y., Ábrego, G. H., Ahn, J., Austin, J., Barham, P., Botha, J. A., Bradbury, J., Brahma, S., Brooks, K., Catasta, M., Cheng, Y., Cherry, C., Choquette-Choo, C. A., Chowdhery, A., Crepy, C., Dave, S., Dehghani, M., Dev, S., Devlin, J., Díaz, M., Du, N., Dyer, E., Feinberg, V., Feng, F., Fienber, V., Freitag, M., Garcia, X., Gehrmann, S., Gonzalez, L., and et al. Palm 2 technical report. *CoRR*, abs/2305.10403, 2023. doi: 10.48550/ARXIV.2305.10403. URL <https://doi.org/10.48550/arXiv.2305.10403>.
- Bacchus, F., Grove, A. J., Halpern, J. Y., and Koller, D. From statistical knowledge bases to degrees of belief. *Artif. Intell.*, 87(1-2):75–143, 1996. doi: 10.1016/S0004-3702(96)00003-3. URL [https://doi.org/10.1016/S0004-3702\(96\)00003-3](https://doi.org/10.1016/S0004-3702(96)00003-3).
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>.
- Chen, B., Shu, C., Shareghi, E., Collier, N., Narasimhan, K., and Yao, S. Fireact: Toward language agent fine-tuning. *CoRR*, abs/2310.05915, 2023a. doi: 10.48550/ARXIV.2310.05915. URL <https://doi.org/10.48550/arXiv.2310.05915>.
- Chen, B., Shu, C., Shareghi, E., Collier, N., Narasimhan, K., and Yao, S. Fireact: Toward language agent fine-tuning. *arXiv preprint arXiv:2310.05915*, 2023b.
- Chen, W., Ma, X., Wang, X., and Cohen, W. W. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *CoRR*, abs/2211.12588, 2022. doi: 10.48550/ARXIV.2211.12588. URL <https://doi.org/10.48550/arXiv.2211.12588>.
- Collier, N. H., Liu, F., and Shareghi, E. On reality and the limits of language data. *CoRR*, abs/2208.11981, 2022. doi: 10.48550/ARXIV.2208.11981. URL <https://doi.org/10.48550/arXiv.2208.11981>.
- Duan, J., Cheng, H., Wang, S., Wang, C., Zavalny, A., Xu, R., Kailkhura, B., and Xu, K. Shifting attention to relevance: Towards the uncertainty estimation of large language models. *CoRR*, abs/2307.01379, 2023. doi: 10.48550/ARXIV.2307.01379. URL <https://doi.org/10.48550/arXiv.2307.01379>.
- Fadeeva, E., Vashurin, R., Tsvigun, A., Vazhentsev, A., Petrakov, S., Fedyanin, K., Vasilev, D., Goncharova, E., Panchenko, A., Panov, M., Baldwin, T., and Shelmanov, A. LM-polygraph: Uncertainty estimation for language models. In Feng, Y. and Lefever, E. (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>.
- Gao, L., Madaan, A., Zhou, S., Alon, U., Liu, P., Yang, Y., Callan, J., and Neubig, G. PAL: program-aided language models. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 10764–10799. PMLR, 2023. URL <https://proceedings.mlr.press/v202/gao23f.html>.
- Geva, M., Khashabi, D., Segal, E., Khot, T., Roth, D., and Berant, J. Did aristotle use a laptop? A question answering benchmark with implicit reasoning strategies. *Trans. Assoc. Comput. Linguistics*, 9:346–361, 2021. doi: 10.1162/TACL\_A\_00370. URL [https://doi.org/10.1162/tacl\\_a\\_00370](https://doi.org/10.1162/tacl_a_00370).
- Gou, Z., Shao, Z., Gong, Y., Shen, Y., Yang, Y., Duan, N., and Chen, W. CRITIC: large language models can self-correct with tool-interactive critiquing. *CoRR*, abs/2305.11738, 2023. doi: 10.48550/ARXIV.2305.11738. URL <https://doi.org/10.48550/arXiv.2305.11738>.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.

- Hu, E. J., yelong shen, Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- Kadavath, S., Conerly, T., Askell, A., Henighan, T., Drain, D., Perez, E., Schiefer, N., Hatfield-Dodds, Z., DasSarma, N., Tran-Johnson, E., Johnston, S., Showk, S. E., Jones, A., Elhage, N., Hume, T., Chen, A., Bai, Y., Bowman, S., Fort, S., Ganguli, D., Hernandez, D., Jacobson, J., Kernion, J., Kravec, S., Lovitt, L., Ndousse, K., Olsson, C., Ringer, S., Amodei, D., Brown, T., Clark, J., Joseph, N., Mann, B., McCandlish, S., Olah, C., and Kaplan, J. Language models (mostly) know what they know. *CoRR*, abs/2207.05221, 2022. doi: 10.48550/ARXIV.2207.05221. URL <https://doi.org/10.48550/arXiv.2207.05221>.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Kuhn, L., Gal, Y., and Farquhar, S. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=VD-AYtP0dve>.
- Kumar, S. and Byrne, B. Minimum bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pp. 169–176, 2004.
- Lin, B. Y., Fu, Y., Yang, K., Ammanabrolu, P., Brahman, F., Huang, S., Bhagavatula, C., Choi, Y., and Ren, X. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. *CoRR*, abs/2305.17390, 2023a. doi: 10.48550/ARXIV.2305.17390. URL <https://doi.org/10.48550/arXiv.2305.17390>.
- Lin, S., Hilton, J., and Evans, O. Teaching models to express their uncertainty in words. *Trans. Mach. Learn. Res.*, 2022, 2022. URL <https://openreview.net/forum?id=8s8K2UZGTZ>.
- Lin, Z., Trivedi, S., and Sun, J. Generating with confidence: Uncertainty quantification for black-box large language models. *CoRR*, abs/2305.19187, 2023b. doi: 10.48550/ARXIV.2305.19187. URL <https://doi.org/10.48550/arXiv.2305.19187>.
- Lu, P., Peng, B., Cheng, H., Galley, M., Chang, K.-W., Wu, Y. N., Zhu, S.-C., and Gao, J. Chameleon: Plug-and-play compositional reasoning with large language models. *arXiv preprint arXiv:2304.09842*, 2023.
- Majumder, B. P., Mishra, B. D., Jansen, P. A., Tafjord, O., Tandon, N., Zhang, L., Callison-Burch, C., and Clark, P. CLIN: A continually learning language agent for rapid task adaptation and generalization. *CoRR*, abs/2310.10134, 2023. doi: 10.48550/ARXIV.2310.10134. URL <https://doi.org/10.48550/arXiv.2310.10134>.
- Malinin, A. and Gales, M. J. F. Uncertainty estimation in autoregressive structured prediction. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=jN5y-zb5Q7m>.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- Ott, M., Auli, M., Grangier, D., and Ranzato, M. Analyzing uncertainty in neural machine translation. In Dy, J. G. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 3953–3962. PMLR, 2018. URL <http://proceedings.mlr.press/v80/ott18a.html>.
- Patil, S. G., Zhang, T., Wang, X., and Gonzalez, J. E. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023.
- Petroni, F., Rocktäschel, T., Riedel, S., Lewis, P. S. H., Bakhtin, A., Wu, Y., and Miller, A. H. Language models as knowledge bases? In Inui, K., Jiang, J., Ng, V., and Wan, X. (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 2019, Hong Kong, China, November 3-7, 2019*, pp. 2463–2473. Association for Computational Linguistics, 2019. doi: 10.18653/V1/D19-1250. URL <https://doi.org/10.18653/v1/D19-1250>.
- Qiao, S., Gui, H., Chen, H., and Zhang, N. Making language models better tool learners with execution feedback. *CoRR*, abs/2305.13068, 2023. doi: 10.48550/ARXIV.2305.13068. URL <https://doi.org/10.48550/arXiv.2305.13068>.
- Ren, A. Z., Dixit, A., Bodrova, A., Singh, S., Tu, S., Brown, N., Xu, P., Takayama, L., Xia, F., Varley, J., Xu, Z.,

- Sadigh, D., Zeng, A., and Majumdar, A. Robots that ask for help: Uncertainty alignment for large language model planners. *CoRR*, abs/2307.01928, 2023. doi: 10.48550/ARXIV.2307.01928. URL <https://doi.org/10.48550/arXiv.2307.01928>.
- Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., Cancedda, N., and Scialom, T. Toolformer: Language models can teach themselves to use tools. *CoRR*, abs/2302.04761, 2023. doi: 10.48550/ARXIV.2302.04761. URL <https://doi.org/10.48550/arXiv.2302.04761>.
- Shinn, N., Labash, B., and Gopinath, A. Reflexion: an autonomous agent with dynamic memory and self-reflection. *CoRR*, abs/2303.11366, 2023. doi: 10.48550/ARXIV.2303.11366. URL <https://doi.org/10.48550/arXiv.2303.11366>.
- Shu, C., Han, J., Liu, F., Shareghi, E., and Collier, N. POSQA: probe the world models of llms with size comparisons. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pp. 7518–7531. Association for Computational Linguistics, 2023. URL <https://aclanthology.org/2023.findings-emnlp.504>.
- Sumers, T. R., Yao, S., Narasimhan, K., and Griffiths, T. L. Cognitive architectures for language agents. *CoRR*, abs/2309.02427, 2023. doi: 10.48550/ARXIV.2309.02427. URL <https://doi.org/10.48550/arXiv.2309.02427>.
- Suzgun, M., Melas-Kyriazi, L., and Jurafsky, D. Follow the wisdom of the crowd: Effective text generation via minimum Bayes risk decoding. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4265–4293, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.262. URL <https://aclanthology.org/2023.findings-acl.262>.
- Tian, K., Mitchell, E., Zhou, A., Sharma, A., Rafailov, R., Yao, H., Finn, C., and Manning, C. D. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 5433–5442. Association for Computational Linguistics, 2023. URL <https://aclanthology.org/2023.emnlp-main.330>.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Canton-Ferrer, C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023. doi: 10.48550/arXiv.2307.09288. URL <https://doi.org/10.48550/arXiv.2307.09288>.
- Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E. H., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=1PLlNIMMrw>.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022. URL [https://openreview.net/pdf?id=\\_VjQlMeSB\\_J](https://openreview.net/pdf?id=_VjQlMeSB_J).
- Wiher, G., Meister, C., and Cotterell, R. On decoding strategies for neural text generators. *Transactions of the Association for Computational Linguistics*, 10:997–1012, 2022. doi: 10.1162/tacl.a.00502. URL <https://aclanthology.org/2022.tacl-1.58>.
- Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E., Zheng, R., Fan, X., Wang, X., Xiong, L., Zhou, Y., Wang, W., Jiang, C., Zou, Y., Liu, X., Yin, Z., Dou, S., Weng, R., Cheng, W., Zhang, Q., Qin, W., Zheng, Y., Qiu, X., Huan, X., and Gui, T. The rise and potential of large language model based agents: A survey. *CoRR*, abs/2309.07864, 2023. doi: 10.48550/ARXIV.2309.07864. URL <https://doi.org/10.48550/arXiv.2309.07864>.
- Xiao, Y. and Wang, W. Y. On hallucination and predictive uncertainty in conditional language generation. In Merlo, P., Tiedemann, J., and Tsarfaty, R. (eds.), *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, pp. 2734–2744. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.EACL-MAIN.

236. URL <https://doi.org/10.18653/v1/2021.eacl-main.236>.
- Xu, J., Desai, S., and Durrett, G. Understanding neural abstractive summarization models via uncertainty. In Webber, B., Cohn, T., He, Y., and Liu, Y. (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6275–6281, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.508. URL <https://aclanthology.org/2020.emnlp-main.508>.
- Yang, Y., Li, H., Wang, Y., and Wang, Y. Improving the reliability of large language models by leveraging uncertainty-aware in-context learning. *CoRR*, abs/2310.04782, 2023. doi: 10.48550/ARXIV.2310.04782. URL <https://doi.org/10.48550/arXiv.2310.04782>.
- Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W. W., Salakhutdinov, R., and Manning, C. D. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In Riloff, E., Chiang, D., Hockenmaier, J., and Tsujii, J. (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pp. 2369–2380. Association for Computational Linguistics, 2018. doi: 10.18653/V1/D18-1259. URL <https://doi.org/10.18653/v1/d18-1259>.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K. R., and Cao, Y. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/pdf?id=WE\\_vluYUL-X](https://openreview.net/pdf?id=WE_vluYUL-X).
- Yin, D., Brahman, F., Ravichander, A., Chandu, K., Chang, K.-W., Choi, Y., and Lin, B. Y. Lumos: Learning agents with unified data, modular design, and open-source llms. *arXiv preprint arXiv:2311.05657*, 2023.
- Zablotskaia, P., Phan, D., Maynez, J., Narayan, S., Ren, J., and Liu, J. Z. On uncertainty calibration and selective generation in probabilistic neural summarization: A benchmark study. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pp. 2980–2992. Association for Computational Linguistics, 2023. URL <https://aclanthology.org/2023.findings-emnlp.197>.
- Zhou, W., Jiang, Y. E., Li, L., Wu, J., Wang, T., Qiu, S., Zhang, J., Chen, J., Wu, R., Wang, S., Zhu, S., Chen, J., Zhang, W., Zhang, N., Chen, H., Cui, P., and Sachan, M. Agents: An open-source framework for autonomous language agents. *CoRR*, abs/2309.07870, 2023. doi: 10.48550/ARXIV.2309.07870. URL <https://doi.org/10.48550/arXiv.2309.07870>.



## A. Effect of Quantile Value

To investigate the effect of the quantile value used in single-inference uncertainty threshold estimation, we compare the results of using different quantile values from 0.1 to 0.9 on HotpotQA using ChatGPT and LLaMA2-70B, shown in Figure 4. As the increase of quantile value, the number of instances with tool-use decreases. This is because the quantile is based on the uncertainty of the calibration set, larger quantile value indicates higher uncertainty threshold, which leads to less tool-use. For UALA-S, large quantile value (i.e., 0.9) has the best performance on both LLMs. When using backoff, the trend is the same on ChatGPT, but LLaMA2-70B has the opposite trend. Based on our observation, we find the reason is that LLaMA2-70B is more likely to generate a null answer than ChatGPT (generate a wrong answer) when the external tool fails to provide useful information, which can be amended by the backoff setting. However, UALA-S+Backoff on ChatGPT with a large quantile value can achieve the same good performance but with only a small number of tool-use.

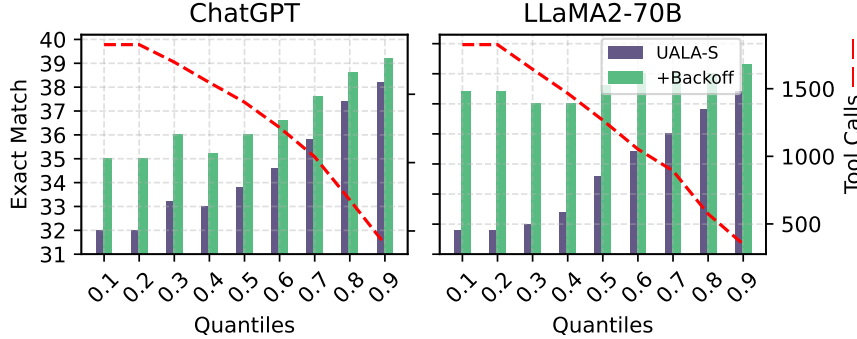


Figure 4. The results on HotpotQA, using different quantile values for uncertainty threshold estimation in single-inference with ChatGPT (left) and LLaMA2-70B (right). The uncertainty estimation method used is Entropy.

## B. Integrating UALA in the CRITIC Framework

The CRITIC (Gou et al., 2023) enables the LLMs to verify and iteratively self-correct their output through interaction with external tools. Following the CRITIC (Gou et al., 2023) work, we also compare two settings: using Google Search as the external tool and using the LLM instead of an external tool (CRITIC w/o Tool) without changing the prompt of CRITIC to generate evidence. Based on CoT outputs, we conduct standard CRITIC and our uncertainty-based CRITIC. Specifically, based on the uncertainty of the answer, we choose whether to do CRITIC or keep the answer. For single-inference estimation, we use Log-sum and 0.8 quantile for threshold estimation. The multi-inference uncertainty estimation is based on 9 samples generated using temperature 0.7.

The results are demonstrated in Table 7. Standard CRITIC performs much better than CRITIC w/o Tool and both of them improve the CoT result effectively. Our uncertainty-based methods reduce the frequency of tool-use by nearly half, but still achieve great results on par with the standard CRITIC. Moreover, both of the S-Un-CRITIC w/o Tool and M-Un-CRITIC w/o Tool achieve better results than standard CRITIC w/o Tool. These results indicate the effectiveness and generalisation of our proposed uncertainty-based tool-use methods with regard to different tool-use frameworks.

Table 7. HotpotQA Results using CRITIC (Gou et al., 2023) on ChatGPT. The round of iterative correction in CRITIC is three.

Methods	EM
CoT	34.8(0)
CRITIC	41.0(1,500)
CRITIC w/o Tool	35.6(1,500)
UALA-S-CRITIC	39.0(597)
UALA-S-CRITIC w/o Tool	38.0(597)
UALA-M-CRITIC	40.6(795)
UALA-M-CRITIC w/o Tool	37.4(795)

### C. Uncertainty-aware Method (UALA) vs. Fine-tuning Method (FireAct)

Based on the published results of [Chen et al. \(2023b\)](#), we also compare results between the uncertainty-aware method (UALA) and fine-tuning method (FireAct) in Table 8. FireAct uses a mix of training data from these three datasets (2,470 samples in total) to do multi-task learning (the setting reported as their best result). We report the results provided in their paper and calculate its improvement of over ReAct. For the UALA, we calculate the improvement of UALA-M+Backoff setting over ReAct. Our uncertainty-aware few-shot learning method (UALA) obtains more improvements than the fine-tuning method (FireAct) over all three tasks. UALA also requires less training data to construct the calibration set compared with fine-tuning data creation in FireAct. For instance, on HotpotQA, FireAct collects 2,000 training questions and for fine-tuning data curation, while UALA only uses 500 training questions to create the calibration set.

Table 8. The comparison between fine-tuning method (FireAct) and our uncertainty-aware method (UALA). Since the data split is not exactly the same for all tasks, we report the comparison results with ReAct. FireAct uses multi-task learning results (numbers provided in their paper) and UALA uses UALA-M+Backoff results. Both methods are on ChatGPT (fine-tune vs. few-shot learning).

Tasks	Fine-tuning Improvement			Uncertainty-aware Improvement		
	ReAct	FireAct	abs./rel. diff	ReAct	UALA	abs./rel. diff
HotpotQA	31.4	39.6	↑8.2/26.1%	32.0	41.2	↑9.2/28.8%
StrategyQA	61.0	72.9	↑11.9/19.5%	55.5	66.9	↑11.4/20.5%
MMLU	58.6	65.8	↑7.2/12.3%	55.8	72.2	↑16.4/29.3%

### D. Effect of Different Sizes of Calibration Set for UALA-S

As described in Section 5.6, Table 9 shows the full results of different sizes of calibration set for UALA-S on HotpotQA.

Table 9. Results of different sizes of calibration set for single-inference uncertainty threshold estimation (UALA-S) on HotpotQA using ChatGPT. Numbers in paranthesis are tool calls.

Calibration Set Size	178	386	558	745	937	1132	1330	1546	1752	1960	2145
Performance	38.2(403)	37.6(503)	37.6(503)	38.2(485)	38.2(485)	37.6(503)	37.6(503)	37.8(500)	38.2(485)	38.2(485)	38.2(485)

### E. Single-inference Estimation Methods Comparison Results

As described in Section 5.7, Table 10 shows the full results of comparing different settings of single-inference uncertainty estimation methods on HotpotQA.

Table 10. Results on HotpotQA comparing different settings of single-inference uncertainty estimation methods (five methods) and uncertainty threshold estimation methods (three methods) on ChatGPT and LLaMA2-70B. We use 0.9 as quantile value for Quantile method. **Bold** shows the best result for each LLM among all different settings. The CoT results with no tool calls are (ChatGPT, EM: 34.8) and (LLaMA2-70B, EM: 35.6). The ReAct results with 100% tool calls are (ChatGPT, EM: 32.0) and (LLaMA2-70B, EM: 32.4).

Uncertainty Estimation Methods	Uncertainty Threshold Estimation Methods					
	ChatGPT			LLaMA2-70B		
	Max	Mean	Quantile	Max	Mean	Quantile
Minimum	35.6(90)	33.0(1,531)	38.0(486)	35.8(80)	32.2(1,567)	35.8(352)
Average	35.8(110)	33.0(1,531)	38.2(406)	35.8(68)	32.2(1,567)	36.2(380)
Normalised Product	35.8(121)	36.8(943)	38.2(478)	35.8(77)	34.2(1,196)	35.8(352)
Log-sum	35.8(121)	36.8(943)	38.2(478)	35.8(77)	34.2(1,196)	35.8(352)
Entropy	35.8(114)	33.0(1,531)	<b>38.2(403)</b>	35.6(51)	32.2(1,567)	<b>36.4(350)</b>

### F. Comparison of Uncertainty for Correct vs. Incorrect Answers

As described in Section 5.1, we calculated the mean uncertainty of correct and incorrect answers. Table 11 demonstrates the difference of the two value across three datasets, and various other statistical tests on the means to highlight the significance

of the difference between uncertainties for correct and incorrect responses.

Table 11. Comparison of uncertainty for Correct vs. Incorrect answers based on mean and standard deviation: Means Difference, T-test, Cohen’s effect size on HotpotQA, StrategyQA, and MMLU. The backbone LLM is ChatGPT.

	Dataset	Sample Size	Difference in Means	T-test	P-value	Cohen’s d
Single-Inference	HotpotQA	500	0.68	3.0382	0.0039758	0.7845
	StrategyQA	229	0.40	4.2226	0.00012028	1.0903
	MMLU	570	0.34	4.6761	0.00001912	1.2074
Multiple-Inference	HotpotQA	500	0.44	24.2285	$3.34 \times 10^{-101}$	1.5323
	StrategyQA	229	0.32	9.4750	$3.27 \times 10^{-19}$	0.8855
	MMLU	570	0.21	16.5838	$9.11 \times 10^{-55}$	0.9823

## G. Statistics of Performance Improvement between UALA+Backoff and Baselines

Table 12 shows the statistics of absolute (relative) performance improvement between UALA+Backoff methods and Baselines, which are calculated based on Table 2.

Table 12. The statistics of absolute (relative) performance improvement between UALA+Backoff methods and Baselines (CoT for HotpotQA, and Standard for StrategyQA and MMLU); absolute (relative) tool calls decrement between UALA+Backoff and ReAct.

Methods	ChatGPT				LLaMA2-70B			
	HotpotQA	StrategyQA	MMLU	Avg.	HotpotQA	StrategyQA	MMLU	Avg.
Absolute (Relative) Performance Improvement								
UALA-S+Backoff	↑4.4(12.6%)	↑8.8(15.2%)	↑2.4(3.5%)	↑5.6(10.4%)	↑1.7(4.8%)	↑5.7(8.6%)	↑5.1(7.9%)	↑4.2(7.1%)
UALA-M+Backoff	↑6.4(18.4%)	↑9.3(16.1%)	↑3.2(4.6%)	↑6.3(13.0%)	↑4.8(13.5%)	↑6.3(9.6%)	↑5.6(8.7%)	↑5.6(10.6%)
Absolute (Relative) Tool Calls Decrement								
UALA-S+Backoff	↓1,711(80.9%)	↓575(81.1%)	↓1,162(63.7%)	↓1,149(75.2%)	↓1,744(83.3%)	↓592(66.5%)	↓1,454(51.8%)	↓1,263(67.2%)
UALA-M+Backoff	↓915(43.3%)	↓475(67.0%)	↓1183(64.9%)	↓858(58.4%)	↓792(37.8%)	↓398(44.7%)	↓1,612(57.4%)	↓934(46.6%)

## H. Statistics of Output Tokens on two LLMs

Table 13 demonstrates the exact numbers of total output tokens for each method shown in Figure 3.

Table 13. The total number of output tokens for each method on ChatGPT and LLaMA2-70B.

Models	Methods				
	Standard	CoT	ReAct	UALA-S	UALA-M
HotpotQA					
ChatGPT	3,261	29,152	308,480	86,893	460,366
LLaMA2-70B	1,997	28,957	403,132	95,984	496,604
StrategyQA					
ChatGPT	1,134	23,222	113,106	22,721	48,709
LLaMA2-70B	229	23,999	169,543	56,916	95,471
MMLU					
ChatGPT	574	61,202	200,450	75,095	78,453
LLaMA2-70B	570	55,761	330,751	157,967	144,589

## I. Full Prompt Templates

---

### HotpotQA

---

#### Standard

Answer the question:

Question: What is the elevation range for the area that the eastern sector of the Colorado orogeny extends into?

Answer: 1,800 to 7,000 ft

Question: Musician and satirist Allie Goertz wrote a song about the "The Simpsons" character Milhouse, who Matt Groening named after who?

Answer: Richard Nixon

Question: Which documentary is about Finnish rock groups, Adam Clayton Powell or The Saimaa Gesture?

Answer: The Saimaa Gesture

Question: What profession does Nicholas Ray and Elia Kazan have in common?

Answer: director, screenwriter, actor

Question: Which magazine was started first Arthur's Magazine or First for Women?

Answer: Arthur's Magazine

Question: Were Pavel Urysohn and Leonid Levin known for the same type of work?

Answer: Yes

---

#### Chain-of-Thought (CoT)

Solve a question answering task. Your task is to generate Thought and Answer where a Thought can reason about the current situation by thinking step by step.

Here are some examples.

Question: What is the elevation range for the area that the eastern sector of the Colorado orogeny extends into?

Thought: Let's think step by step. The eastern sector of Colorado orogeny extends into the High Plains. High Plains rise in elevation from around 1,800 to 7,000 ft, so the answer is 1,800 to 7,000 ft.

Answer: 1,800 to 7,000 ft

Question: Musician and satirist Allie Goertz wrote a song about the "The Simpsons" character Milhouse, who Matt Groening named after who?

Thought: Let's think step by step. Milhouse was named after U.S. president Richard Nixon, so the answer is Richard Nixon.

Answer: Richard Nixon

Question: Which documentary is about Finnish rock groups, Adam Clayton Powell or The Saimaa Gesture?

Thought: Let's think step by step. Adam Clayton Powell (film) is a documentary about an African-American politician, not Finnish rock groups. So the documentary about Finnish rock groups must instead be The Saimaa Gesture.

Answer: The Saimaa Gesture

Question: What profession does Nicholas Ray and Elia Kazan have in common?

Thought: Let's think step by step. Professions of Nicholas Ray are director, screenwriter, and actor. Professions of Elia Kazan are director, producer, screenwriter, and actor. So profession Nicholas Ray and Elia Kazan have in common is director, screenwriter, and actor.

Answer: director, screenwriter, actor

Question: Which magazine was started first Arthur's Magazine or First for Women?

Thought: Let's think step by step. Arthur's Magazine was started in 1844. First for Women was started in 1989. 1844 (Arthur's Magazine) < 1989 (First for Women), so Arthur's Magazine was started first.

Answer: Arthur's Magazine



Question: Were Pavel Urysohn and Leonid Levin known for the same type of work?

Thought: Let's think step by step. Pavel Urysohn is a mathematician. Leonid Levin is a mathematician and computer scientist. So Pavel Urysohn and Leonid Levin have the same type of work.

Answer: Yes

---

### ReAct

Solve a question answering task with interleaving Thought, Action, Observation steps. Thought can reason about the current situation, and Action can be three types:

(1) Search[entity], which searches the exact entity on Wikipedia and returns the first paragraph if it exists. If not, it will return some similar entities to search.

(2) Lookup[keyword], which returns the next sentence containing keyword in the current passage. (3) Finish[answer], which returns the answer and finishes the task.

Here are some examples.

Question: What is the elevation range for the area that the eastern sector of the Colorado orogeny extends into?

Thought 1: I need to search Colorado orogeny, find the area that the eastern sector of the Colorado orogeny extends into, then find the elevation range of the area.

Action 1: Search[Colorado orogeny]

Observation 1: The Colorado orogeny was an episode of mountain building (an orogeny) in Colorado and surrounding areas.

Thought 2: It does not mention the eastern sector. So I need to look up eastern sector.

Action 2: Lookup[eastern sector]

Observation 2: (Result 1 / 1) The eastern sector extends into the High Plains and is called the Central Plains orogeny.

Thought 3: The eastern sector of Colorado orogeny extends into the High Plains. So I need to search High Plains and find its elevation range.

Action 3: Search[High Plains]

Observation 3: High Plains refers to one of two distinct land regions:

Thought 4: I need to instead search High Plains (United States).

Action 4: Search[High Plains (United States)]

Observation 4: The High Plains are a subregion of the Great Plains. From east to west, the High Plains rise in elevation from around 1,800 to 7,000 ft (550 to 2,130 m).[3]

Thought 5: High Plains rise in elevation from around 1,800 to 7,000 ft, so the answer is 1,800 to 7,000 ft.

Action 5: Finish[1,800 to 7,000 ft]

Question: Musician and satirist Allie Goertz wrote a song about the "The Simpsons" character Milhouse, who Matt Groening named after who?

Thought 1: The question simplifies to "The Simpsons" character Milhouse is named after who. I only need to search Milhouse and find who it is named after.

Action 1: Search[Milhouse]

Observation 1: Milhouse Mussolini Van Houten is a recurring character in the Fox animated television series The Simpsons voiced by Pamela Hayden and created by Matt Groening.

Thought 2: The paragraph does not tell who Milhouse is named after, maybe I can look up "named after".

Action 2: Lookup[named after]

Observation 2: (Result 1 / 1) Milhouse was named after U.S. president Richard Nixon, whose middle name was Milhous.

Thought 3: Milhouse was named after U.S. president Richard Nixon, so the answer is Richard Nixon.

Action 3: Finish[Richard Nixon]

Question: Which documentary is about Finnish rock groups, Adam Clayton Powell or The Saimaa Gesture?

Thought 1: I need to search Adam Clayton Powell and The Saimaa Gesture, and find which documentary is about Finnish rock groups.

Action 1: Search[Adam Clayton Powell]

Observation 1: Could not find [Adam Clayton Powell]. Similar: ['Adam Clayton Powell III', 'Seventh Avenue (Manhattan)', 'Adam Clayton Powell Jr. State Office Building', 'Isabel Washington Powell', 'Adam Powell', 'Adam Clayton Powell (film)', 'Giancarlo Esposito'].

Thought 2: To find the documentary, I can search Adam Clayton Powell (film).

Action 2: Search[Adam Clayton Powell (film)]

Observation 2: Adam Clayton Powell is a 1989 American documentary film directed by Richard Kilberg.

The film is about the rise and fall of influential African-American politician Adam Clayton Powell Jr.[3][4] It was later aired as part of the PBS series The American Experience.

Thought 3: Adam Clayton Powell (film) is a documentary about an African-American politician, not Finnish rock groups. So the documentary about Finnish rock groups must instead be The Saimaa Gesture.

Action 3: Finish[The Saimaa Gesture]

Question: What profession does Nicholas Ray and Elia Kazan have in common?

Thought 1: I need to search Nicholas Ray and Elia Kazan, find their professions, then find the profession they have in common.

Action 1: Search[Nicholas Ray]

Observation 1: Nicholas Ray (born Raymond Nicholas Kienzle Jr., August 7, 1911 - June 16, 1979) was an American film director, screenwriter, and actor best known for the 1955 film Rebel Without a Cause.

Thought 2: Professions of Nicholas Ray are director, screenwriter, and actor. I need to search Elia Kazan next and find his professions.

Action 2: Search[Elia Kazan]

Observation 2: Elia Kazan was an American film and theatre director, producer, screenwriter and actor.

Thought 3: Professions of Elia Kazan are director, producer, screenwriter, and actor. So profession Nicholas Ray and Elia Kazan have in common is director, screenwriter, and actor.

Action 3: Finish[director, screenwriter, actor]

Question: Which magazine was started first Arthur's Magazine or First for Women?

Thought 1: I need to search Arthur's Magazine and First for Women, and find which was started first.

Action 1: Search[Arthur's Magazine]

Observation 1: Arthur's Magazine (1844-1846) was an American literary periodical published in Philadelphia in the 19th century.

Thought 2: Arthur's Magazine was started in 1844. I need to search First for Women next.

Action 2: Search[First for Women]

Observation 2: First for Women is a woman's magazine published by Bauer Media Group in the USA.[1] The magazine was started in 1989.

Thought 3: First for Women was started in 1989. 1844 (Arthur's Magazine) < 1989 (First for Women), so Arthur's Magazine was started first.

Action 3: Finish[Arthur's Magazine]

Question: Were Pavel Urysohn and Leonid Levin known for the same type of work?

Thought 1: I need to search Pavel Urysohn and Leonid Levin, find their types of work, then find if they are the same.

Action 1: Search[Pavel Urysohn]

Observation 1: Pavel Samuilovich Urysohn (February 3, 1898 - August 17, 1924) was a Soviet mathematician who is best known for his contributions in dimension theory.

Thought 2: Pavel Urysohn is a mathematician. I need to search Leonid Levin next and find its type of work.

Action 2: Search[Leonid Levin]

Observation 2: Leonid Anatolievich Levin is a Soviet-American mathematician and computer scientist.

Thought 3: Leonid Levin is a mathematician and computer scientist. So Pavel Urysohn and Leonid Levin have the same type of work.

Action 3: Finish[yes]

---

### StrategyQA

---

#### Standard

Answer the question:

Question: Will the Albany in Georgia reach a hundred thousand occupants before the one in New York?

Answer: No

Question: Do the anchors on Rede Globo speak Chinese?

Answer: No

Question: Would a student of the class of 2017 have amnesia about 9/11?

Answer: Yes

Question: Is average number of peas in a pod enough commas for a billion?

Answer: Yes

---

### Chain-of-Thought(CoT)

Solve a question answering task. Your task is to generate Thought and Answer where a Thought can reason about the current situation by thinking step by step.

Here are some examples.

Question: Do the anchors on Rede Globo speak Chinese?

Thought: The anchors on Rede Globo, a Brazilian television network, primarily speak Portuguese as that is the official language of Brazil. They may have proficiency in other languages, but Chinese is not likely to be one of the languages commonly spoken by the anchors on Rede Globo. So the answer is no.

Answer: No

Question: Will the Albany in Georgia reach a hundred thousand occupants before the one in New York?

Thought: As of the most recent population estimates, Albany, New York, had a population of approximately 97,000 residents, while Albany, Georgia, had a population of around 73,000 residents. Albany, New York, is the capital of the state and is a major center for business, education, and government. It has a long history and economic significance, which attracts people to live and work in the area. On the other hand, Albany, Georgia, while an important regional center, is a smaller city in comparison. It does not have the same level of economic or cultural influence as Albany, New York. In conclusion, based on the current population figures and the different dynamics at play, it is unlikely that Albany, Georgia, will reach a population of one hundred thousand before Albany, New York. So the answer is no.

Answer: No

Question: Is average number of peas in a pod enough commas for a billion?

Thought: Generally, a typical pea pod contains around 6 to 9 peas. A billion is a very large number 1,000,000,000 that requires 3 commas, which is less than the average number of peas in a pod. So the answer is yes.

Answer: Yes

Question: Is the language used in Saint Vincent and the Grenadines rooted in English?

Thought: Saint Vincent and the Grenadines were once British colonies, and English became the dominant language during the colonial period. After gaining independence in 1979, English remained as the official language of the country, and it has continued to be used in education, government, media, and daily communication. English has permeated various aspects of society and is widely spoken by the population, though local dialects and accents may influence the spoken form of English in the region. So the answer is yes.

Answer: Yes

---

### ReAct

Solve a question answering task with interleaving Thought, Action, Observation steps. Thought can reason about the current situation, and Action can be three types:

(1) Search[entity], which searches the exact entity on Wikipedia and returns the first paragraph if it exists. If not, it will return some similar entities to search.

(2) Lookup[keyword], which returns the next sentence containing keyword in the current passage. (3) Finish[answer], which returns the answer and finishes the task.

Here are some examples.

Question: Will the Albany in Georgia reach a hundred thousand occupants before the one in New York?

Thought 1: I need to search Albany, Georgia and Albany, New York, find their current number of occupants, then compare to see which one will reach a hundred thousand first.

Action 1: Search[Albany, Georgia]

Observation 1: Albany is a city in the U.S. state of Georgia. Located on the Flint River, it is the seat of Dougherty County, and is the sole incorporated city in that county. The population was 77,434 at the 2010 U.S.

Thought 2: The population of Albany, Georgia was 77,434 at the 2010 U.S.

Action 2: Search[Albany, New York]

Observation 2: Albany is the capital city of the U.S. state of New York and the seat of Albany County. It is located on the west bank of the Hudson River, about 10 miles (16 km) south of its confluence with the Mohawk River, and about 135 miles (220 km) north of New York City.

Thought 3: It does not mention the population of New York City. So I need to look up population.

Action 3: Lookup[population]

Observation 3: (Result 1 / 1) As of 2020, Albany's population was 99,224.

Thought 4: Albany, New York has 99,224 occupants in 2020, which is larger than Albany, Georgia. So Albany in Georgia will not reach a hundred thousand occupants before the one in New York, the answer is no.

Action 4: Finish[No]

Question: Do the anchors on Rede Globo speak Chinese?

Thought 1: I need to search Rede Globo, and find if the anchors on Rede Globo speak Chinese.

Action 1: Search[Rede Globo]

Observation 1: TV Globo, or simply Globo and alternatively as Global), formerly known as Rede Globo, is a Brazilian free-to-air television network, launched by media proprietor Roberto Marinho on 26 April 1965.

Thought 2: Rede Globo is a Brazilian television network, and Brazil is not a Chinese-speaking country, so anchors on Rede Globo do not speak Chinese.

Action 2: Finish[No]

Question: Would a student of the class of 2017 have amnesia about 9/11?

Thought 1: I need to search the childhood amnesia and 9/11, and find if a student of the class of 2017 would have childhood amnesia about 9/11.

Action 1: Search[childhood amnesia]

Observation 1: Childhood amnesia, also called infantile amnesia, is the inability of adults to retrieve episodic memories (memories of situations or events) before the age of two to four years.

Thought 2: Amnesia happens before the age of two to four years.

Action 2: Search[9/11]

Observation 2: The September 11 attacks, commonly known as 9/11, were four coordinated Islamist suicide terrorist attacks carried out by al-Qaeda against the United States in 2001.

Thought 3: 9/11 attacks happened in 2001. A student of the class of 2017 is born around 1999, they would have been around 2 years old during the 9/11 attacks in 2001. Amnesia happens before the age of two to four years, so a student of the class of 2017 would have amnesia about 9/11.

Action 3: Finish[Yes]

Question: Is average number of peas in a pod enough commas for a billion?

Thought 1: I need to search average number of peas in a pod and billion, and find if the average number of peas in a pod is enough commas for a billion.

Action 1: Search[peas]

Observation 1: The pea is most commonly the small spherical seed or the seed-pod of the flowering plant species *Lathyrus oleraceus*. Each pod contains several peas, which can be green or yellow.

Thought 2: Each pod contains several peas. The number of several is usually three to five.

Action 2: Search[billion]

Observation 2: Billion is a word for a large number, and it has two distinct definitions: 1,000,000,000, i.e. one thousand million, or  $10^9$  (ten to the ninth power), as defined on the short scale.

Thought 3: A billion is expressed as '1,000,000,000', which amounts to three commas. The average number of peas in a pod is three to five, which is indeed greater than the number of commas used in a billion (3 commas), so the answer is yes.

Action 3: Finish[Yes]

---

**MMLU**

---

**Standard**

Answer the question:

Question: A person takes buckets of water from the house and begins to add it to a pond in the yard. After a certain point, the pond

A. breaks

B. bloats



C. sinks  
D. drowns  
Answer: B

Question: Coal is solid rock that began as organic material that was deposited in a swamp. The formation of coal suggests that,

A. coal is made mostly of skeletal remains of animals.  
B. coal is formed from magma that has solidified over time.  
C. it quickly becomes petrified when water is removed.  
D. geologic processes continue over millions of years.  
Answer: D

Question: A student uses the following characteristics to describe a group of objects in space. \* 200 billion stars \* 30 million light years from Earth \* 500 light years in diameter Which of the following is the student most likely describing?

A. a galaxy  
B. the universe  
C. a constellation  
D. the solar system  
Answer: A

Question: A microwave oven is connected to an outlet, 120 V, and draws a current of 2 amps. At what rate is energy being used by the microwave oven?

A. 10 W  
B. 30 W  
C. 240 W  
D. 60 W  
Answer: C

---

#### Chain-of-Thought(CoT)

Solve a question answering task. Your task is to generate Thought and Answer where a Thought can reason about the current situation by thinking step by step.

Here are some examples.

Question: A person takes buckets of water from the house and begins to add it to a pond in the yard. After a certain point, the pond

A. breaks  
B. bloats  
C. sinks  
D. drowns

Thought: Each time the person adds a bucket of water, the level of water in the pond rises. Of all options, only B. bloats is consistent with the rise of water level. So the answer is B.

Answer: B

Question: Coal is solid rock that began as organic material that was deposited in a swamp. The formation of coal suggests that,

A. coal is made mostly of skeletal remains of animals.  
B. coal is formed from magma that has solidified over time.  
C. it quickly becomes petrified when water is removed.  
D. geologic processes continue over millions of years.

Thought: Let's evaluate each option. A. Coal is mostly composed of plant matter, not the skeletal remains of animals. Therefore, this option is incorrect. B. Coal is not formed from magma. Magma that solidifies over time creates igneous rocks, so this option is also incorrect. C. Petrification is a process by which organic material is turned into stone. It is not directly related to the process of coal formation, so this option is incorrect. D. The formation of coal takes incredibly long periods of time and consists of slow geologic processes such as sedimentation and metamorphism, which suggests that such processes continue over millions of years. Therefore, this option is correct. The answer is D.

Answer: D

Question: A student uses the following characteristics to describe a group of objects in space. \* 200 billion stars \* 30 million light years from Earth \* 500 light years in diameter Which of the following is the student most likely describing?

- A. a galaxy
- B. the universe
- C. a constellation
- D. the solar system

Thought: Let's evaluate each option. A. a galaxy: Possibly, as galaxies do contain billions of stars and can be millions of light years from Earth. B. the universe: Unlikely, as the universe is far larger than 30 million light years and contains more than just 200 billion stars. C. a constellation: Unlikely, as constellations are patterns of stars seen from Earth and don't have a physical size or distance associated with them. D. the solar system: Definitely not, as our solar system only contains one star, our sun. So, the answer is most likely A. a galaxy.

Answer: A

Question: A microwave oven is connected to an outlet, 120 V, and draws a current of 2 amps. At what rate is energy being used by the microwave oven?

- A. 10 W
- B. 30 W
- C. 240 W
- D. 60 W

Thought: The rate energy is being used is equal to the power. For an electrical device  $P = IV$ , so  $P = (120)(2) = 240$  W. So the answer is C.

Answer: C

---

### ReAct

Solve a question answering task with interleaving Thought, Action, Observation steps. Thought can reason about the current situation, and Action can be two types:

- (1) search[question], which searches a question on Google and returns a short snippet containing the answer. Note that sometimes the snippet does not contain the answer, and some alternative search might be needed.
- (2) finish[answer], which returns the answer and finishes the task.

Here are some examples.

Question: A person takes buckets of water from the house and begins to add it to a pond in the yard. After a certain point, the pond

- A. breaks
- B. bloats
- C. sinks
- D. drowns

Thought: After continuously adding water to a pond, the pond will have more water than it could hold, thus bloats. So the answer is B.

Action: finish[B]

Question: Coal is solid rock that began as organic material that was deposited in a swamp. The formation of coal suggests that,

- A. coal is made mostly of skeletal remains of animals.
- B. coal is formed from magma that has solidified over time.
- C. it quickly becomes petrified when water is removed.
- D. geologic processes continue over millions of years.

Thought: The question is about the formation of coal. I need to first learn how coal is formed.

Action: search[How is coal formed?]

Observation: Coal takes millions of years to form Coal contains the energy stored by plants that lived hundreds of millions of years ago in swampy forests. Layers of dirt and rock covered the plants over millions of years. The resulting pressure and heat turned the plants into the substance we call coal.

Thought: Based on the information, I can check each option. A: coal is made by plants, not animals, so A is false. B: I have no information about if coal is formed from magma yet. I could search 'is coal formed from magma' to make sure. C: I have no information about if coal quickly becomes petrified when water is removed. I could search 'does coal quickly become petrified when water is removed' to make sure. D: Coal takes millions of years to form, so D is possibly true. I could search 'is the formulation of coal a geologic process' to make sure.

Action: search[is the formulation of coal a geologic process]

Observation: It is formed from plant remains that have been compacted, hardened, chemically altered, and metamorphosed by heat and pressure over geologic time.

Thought: Seems the formulation of coal is over geologic time, so a geologic process. So the answer is D.

Action: finish[D]

Question: A student uses the following characteristics to describe a group of objects in space. \* 200 billion stars \* 30 million light years from Earth \* 500 light years in diameter Which of the following is the student most likely describing?

A. a galaxy

B. the universe

C. a constellation

D. the solar system

Thought: These options correspond to space systems of different sizes. I could search what is the diameter of each option to match.

Action: search[what is the diameter of a galaxy]

Observation: Most galaxies are 1,000 to 100,000 parsecs in diameter (approximately 3,000 to 300,000 light years) and are separated by distances on the order of millions of parsecs (or megaparsecs).

Thought: A galaxy is usually 3,000 to 300,000 light years in diameter, which is slightly more than 500 light years. I should search the diameter of the universe next.

Action: search[what is the diameter of the universe]

Observation: 93 billion light-years

Thought: The universe is 93 billion light years in diameter, which is much larger than 500 light years. I should search the diameter of a constellation next.

Action: search[what is the diameter of a constellation]

Observation: Its diameter, remarkably, is greater than 10 AU (1.5 billion kilometers!), large enough to fill the entire inner solar system almost as far out as Jupiter.

Thought: A constellation is usually 10 AU in diameter. I need to convert it into light years.

Action: search[10 AU to light years]

Observation: 0.000158125

Thought: A constellation is usually 0.000158125 light years in diameter, which is much smaller than 500 light years. I should search the diameter of the solar system next.

Action: Sedna is three times farther away from Earth than Pluto, making it the most distant observable object known in the solar system. It is 143.73 billion km from the Sun, thus giving the Solar System a diameter of 287.46 billion km.

Thought: The solar system is 287.46 billion km in diameter. I need to convert it into light years.

Action: search[287.46 billion km to light years]

Observation: 0.0303845459748716

Thought: A constellation is usually 0.0303845459748716 light years in diameter, which is much smaller than 500 light years. Given all the information about diameters, the diameter of a galaxy is closest to 500 light years. So the answer is A.

Action: finish[A]

Question: A microwave oven is connected to an outlet, 120 V, and draws a current of 2 amps. At what rate is energy being used by the microwave oven?

A. 10 W

B. 30 W

C. 240 W

D. 60 W

Thought: The question is about the rate of energy. I need to search what is the rate of energy.

Action: search[what is the rate of energy]

Observation: The rate at which energy is transferred is called power. When we say rate, this usually means a certain unit with respect to time or a unit over time. Power, as in the physics term, can be formally defined as the energy transferred per unit time.

Thought: Given voltage and current, I need to search how to calculate power.

Action: search[given voltage and current, how to calculate power]

Observation: Power in Watts is calculated by multiplying voltage in Volts times current in Amps: 10 Amps of current at 240 Volts generates 2,400 Watts of power.

Thought: In this case, the voltage is 120 volts and the current is 2 amps.  $\text{Power} = 120\text{V} \times 2\text{A} = 240\text{W}$ . So, the rate at which energy is being used by the microwave oven is 240 watts. So the correct answer is C. 240 W.

Action: finish[C]

---

## J. Verbalised Confidence Prompt Template

Solve a question answering task. Your task is to generate Thought and Answer where a Thought can reason about the current situation by thinking step by step. Provide the probability that the answer is correct (0.0 to 1.0). Give ONLY the probability in the format of Answer[Probability], no other words or explanation.