

Algorithm 1 Multi-step LLM planning with human help.

```

1: for time  $t \leftarrow 0$  to  $T-1$  do
2:   Observe input  $x_{\text{test}}^t$ 
3:   Predict a set  $C(x_{\text{test}}^t)$ 
4:   if  $C(x_{\text{test}}^t)$  is a singleton then
5:     execute corresponding action
6:   else
7:     ask for help, which is always assumed to provide the correct label, or the human provide clarification only
8:   end if
9: end for

```

480 **Proof of Claim 1:** Suppose $\bar{y} \in \bar{C}(\bar{x}_{\text{test}})$. We have,

$$\bar{y} \in \bar{C}(\bar{x}_{\text{test}}) \iff \min_t \hat{f}(x_{\text{test}}^t)_{y^t} \geq 1 - \hat{q} \quad (\text{A1})$$

$$\iff \hat{f}(x_{\text{test}}^t)_{y^t} \geq 1 - \hat{q}, \quad \forall t \in [T] \quad (\text{A2})$$

$$\iff y^t \in C^t(x_{\text{test}}^t), \quad \forall t \in [T] \quad (\text{A3})$$

$$\iff \bar{y} \in C(\bar{x}_{\text{test}}). \quad (\text{A4})$$

481 **Proof of Proposition 2:** Since we can bound the probability that $\bar{y}_{\text{test}} \notin \bar{C}(\bar{x}_{\text{test}})$, we can also bound
482 the probability that $\bar{y}_{\text{test}} \notin C(\bar{x}_{\text{test}})$. From the conformalization procedure, we have the following *dataset-*
483 *conditional* guarantee: with probability $1 - \delta$ over the sampling of the calibration set \bar{Z} , we have

$$\mathbb{P}(\bar{y}_{\text{test}} \in \bar{C}(\bar{x}_{\text{test}}) | \bar{Z}) \geq \text{Beta}_{N+1-v, v}^{-1}(\delta), \quad v = \lfloor (N+1)\hat{\epsilon} \rfloor \quad (\text{A5})$$

$$\stackrel{\text{Claim 1}}{\implies} \mathbb{P}(\bar{y}_{\text{test}} \in C(\bar{x}_{\text{test}}) | \bar{Z}) \geq \text{Beta}_{N+1-v, v}^{-1}(\delta), \quad (\text{A6})$$

484 where $\hat{\epsilon}$ is chosen such that $\epsilon = 1 - \text{Beta}_{N+1-v, v}^{-1}(\delta)$. Hence, the following *marginal* guarantee also holds:

$$\mathbb{P}(\bar{y}_{\text{test}} \in \bar{C}(\bar{x}_{\text{test}})) \geq 1 - \hat{\epsilon}$$

$$\stackrel{\text{Claim 1}}{\implies} \mathbb{P}(\bar{y}_{\text{test}} \in C(\bar{x}_{\text{test}})) \geq 1 - \hat{\epsilon}.$$

485 This result provides a bound on the task completion rate if \bar{x}_{test} is drawn using the distribution \mathcal{D} . However,
486 recall that the sequence \bar{x} of augmented contexts as defined in Section 3.3 arises from having performed
487 the *correct* actions in previous steps; incorrect actions may result in a distribution shift. In order to obtain a
488 bound on the task completion rate, we consider three cases at any given timestep: (1) the prediction set is a
489 singleton and contains the correct label, (2) the prediction set is not a singleton but does contain the correct
490 label, and (3) the prediction set does not contain the true label. The robot performs the correct action in the
491 first two cases (without help in (1) and with help in (2)), while CP bounds the probability of case (3). Thus,
492 the CP bound translates to a bound on the task success rate.

493 As seen in Eq. (3), we have from [10, Thm. 1], that we achieve the smallest average set size among all
494 possible sequence-level prediction schemes, \bar{C} , if \hat{f} models the prediction uncertainty accurately,

$$\min_{\bar{C} \in \bar{\mathcal{C}}} \mathbb{E}_{(\bar{x}, \cdot) \sim \mathcal{D}} [|\bar{C}(\bar{x})|], \text{ subject to } \mathbb{P}(\bar{y} \in \bar{C}(\bar{x})) \geq 1 - \hat{\epsilon}. \quad (\text{A7})$$

A2 CP in Settings with Multiple Acceptable Options Per Step

496 **Proposition 3 (Multi-label uncertainty alignment)** Consider a setting where we use CP with coverage
497 level $1 - \epsilon$ to construct the prediction set when there are multiple true labels and seek help whenever the set
498 is not a singleton at each timestep. With probability $1 - \delta$ over the sampling of the calibration set, the task
499 completion rate over new test scenarios drawn from \mathcal{D} is at least $1 - \epsilon$.

500 **Proof:** We have a dataset of $Z = \{(\tilde{x}_i, Y_i), \dots\}_{i=1}^N$ sampled i.i.d. from a data distribution \mathcal{D} for calibration
501 (we use the same notation \mathcal{D} as in the single-label setting here), where $Y_i := \{y_{i,j}\}_{j=1}^{J_i}$ is the set of true
502 labels for a single trial. For each label, we use the same heuristic notion of confidence, $\hat{f}(x)_y \in [0, 1]$.

503 We define an operator $\beta: \mathcal{X} \times \mathcal{Y}^J \rightarrow \mathcal{Y}$ where \mathcal{X} is the space of contexts and \mathcal{Y} is the space of labels:

$$\beta(x, Y) := \operatorname{argmax}_{y \in Y} \hat{f}(x)_y, \quad (\text{A8})$$

504 which takes the true label with the highest confidence value from the true label set.

505 If we consider applying β to every point in the support of \mathcal{D} , a new distribution \mathcal{D}' is induced. We also
 506 consider the induced dataset of samples $S' = \{(x_i, y_i^{\max})\}_{i=1}^N$, where $y_i^{\max} := \beta(x_i, Y_i)$. Then we can
 507 perform the usual conformalization and obtain the guarantee that with

$$C(x_{\text{test}}) := \{y | \hat{f}(x_{\text{test}})_y \geq 1 - \hat{q}\}, \quad (\text{A9})$$

508 the following *marginal* guarantee holds,

$$\mathbb{P}(y_{\text{test}}^{\max} \notin C(x_{\text{test}})) \leq \hat{\epsilon}, \quad (\text{A10})$$

$$\Rightarrow \mathbb{P}(\operatorname{argmax}_{y \in Y_{\text{test}}} \hat{f}(x_{\text{test}})_y \notin C(x_{\text{test}})) \leq \hat{\epsilon}, \quad (\text{A11})$$

$$\Rightarrow \mathbb{P}(\beta(x_{\text{test}}, Y_{\text{test}}) \notin C(x_{\text{test}})) \leq \hat{\epsilon}, \quad (\text{A12})$$

509 and the following *dataset-conditional* guarantee holds when we choose $\hat{\epsilon}$ such that $\epsilon = 1 - \text{Beta}_{N+1-v, v}^{-1}(\delta)$
 510 where $v = \lfloor (N+1)\hat{\epsilon} \rfloor$,

$$\mathbb{P}(\beta(x_{\text{test}}, Y_{\text{test}}) \in C(x_{\text{test}}) | Z) \geq 1 - \epsilon. \quad (\text{A13})$$

511 Hence, $C(x_{\text{test}})$ contains the true label with the highest confidence with probability at least $1 - \epsilon$.

512 At test time, we sample $(x_{\text{test}}, Y_{\text{test}})$ from \mathcal{D} that is i.i.d. with samples in S — for the guarantee to hold for
 513 $\beta(x_{\text{test}}, Y_{\text{test}})$, we need to show $\beta(x_{\text{test}}, Y_{\text{test}})$ is a sample from \mathcal{D}' that is i.i.d. with samples in S' . This is true
 514 since functions of independent random variables are independent, and functions of identically distributed
 515 random variables are identically distributed if the functions are measurable.

516 A3 CP in Multi-Step Setting with Multiple Acceptable Options Per Step

517 **Proposition 4 (Multi-step, multi-label uncertainty alignment)** Consider a multi-step setting where we
 518 use CP with coverage level $1 - \epsilon$ to causally construct the prediction set when there may be multiple true
 519 labels at any step and seek help whenever the set is not a singleton at each timestep. With probability $1 - \delta$
 520 over the sampling of the calibration set, the task completion rate over new test scenarios drawn from \mathcal{D} is
 521 at least $1 - \epsilon$.

522 **Proof:** For the multi-step setting, each trial now involves a sequence of contexts \bar{x} and a set of sequences
 523 of true labels:

$$\bar{Y} = \{\bar{y}_1, \bar{y}_2, \dots, \bar{y}_M\}, \quad (\text{A14})$$

524 where $\bar{y}_m := (y_m^0, y_m^1, \dots, y_m^{T-1})$. For example, \bar{Y} can contain the sequence of “blue block, yellow block,
 525 green block”, “green block, blue block, yellow block”, ..., for the task of picking up three blocks. We
 526 collect a dataset of $\bar{Z} = \{(\bar{x}_i, \bar{Y}_i)\}$ of i.i.d. samples from the data distribution $\bar{\mathcal{D}}$.

527 Unlike the single-step setting, here we cannot apply β to the set of true labels in each step since we are
 528 reasoning over a *set of sequences*, and not a *sequence of sets* of true labels. Notably, the true label set at
 529 time step t depends upon the sequence of previously chosen true labels.

530 Let $Y^t[x^0, \bar{y}^{t-1}]$ denote the set of true labels at timestep t , *conditioned* upon the initial context x^0 and a
 531 partial sequence of past true labels $\bar{y}^{t-1} := (y^0, \dots, y^{t-1})$ extracted from \bar{Y} . We then autoregressively define
 532 the following sequence:

$$\bar{\beta}_0(\bar{x}, \bar{Y}) := \operatorname{argmax}_{y \in Y^0} \hat{f}(x^0)_y, \quad Y^0 := \{y_1^0, \dots, y_M^0\} \quad (\text{A15})$$

$$\bar{\beta}_t(\bar{x}, \bar{Y}) := \bar{\beta}_{t-1}(\bar{x}, \bar{Y}) \bigcup \operatorname{argmax}_{y \in Y^t[x^0, \bar{\beta}_{t-1}(\bar{x}, \bar{Y})]} \hat{f}(x^t)_y, \quad t = 1, \dots, T-1. \quad (\text{A16})$$

533 For convenience, we denote $\bar{\beta}_t(\bar{x}, \bar{Y})[\tau]$ the τ element in $\bar{\beta}_t(\bar{x}, \bar{Y})$, $\tau \leq t$. An intuitive interpretation is that,
 534 we can consider \bar{Y} forming a tree of valid executions (all possible actions that can be taken by choosing
 535 each of true labels). Hence, at each time step t , $\bar{\beta}_t(\bar{x}, \bar{Y})$ prunes the tree to a single branch by taking the
 536 true label with the highest heuristic value $\hat{f}(x^t)$. This reduces the tree of all possible sequences of true

537 labels to a single branch of true labels with highest confidence. Given this single branch of true labels, we
 538 can now perform CP as shown in the multi-step setting in [Section A1](#).

539 We apply $\bar{\beta}_{T-1}$ to every point in the support of $\bar{\mathcal{D}}$, and a new distribution $\bar{\mathcal{D}}'$ is induced. We consider
 540 $\bar{\mathcal{S}}' = \{(\bar{x}_i, \bar{y}_i^{\max})\}$, where $\bar{y}_i^{\max} := \bar{\beta}_{T-1}(\bar{x}_i, \bar{Y}_i)$. Let \bar{Y}_{test} be the set of sequences of true labels for \bar{x}_{test} .
 541 Suppose we get the *marginal* bound with $\bar{\beta}_{T-1}$ as the labels:

$$\mathbb{P}(\bar{\beta}_{T-1}(\bar{x}_{\text{test}}, \bar{Y}_{\text{test}}) \notin C(\bar{x}_{\text{test}})) \leq \hat{\epsilon}, \quad (\text{A17})$$

542 and *dataset-conditional* bound when we choose $\hat{\epsilon}$ such that $\epsilon = 1 - \text{Beta}_{N+1-v, v}^{-1}(\delta)$ where $v = \lfloor (N+1)\hat{\epsilon} \rfloor$,

$$\mathbb{P}(\bar{\beta}_{T-1}(\bar{x}_{\text{test}}, \bar{Y}_{\text{test}}) \notin C(\bar{x}_{\text{test}}) | \bar{\mathcal{Z}}) \leq \epsilon, \quad (\text{A18})$$

543 which states that at test time, given a context sequence \bar{x}_{test} , we produce a prediction set of sequences; if
 544 we consider a sequence consisting of the true label with the highest score at each step, the probability of
 545 this sequence covered by $C(\bar{x}_{\text{test}})$ is lower bounded by $1 - \epsilon$. However, we need to be careful of following
 546 $\bar{\beta}_t$ at each step at test time. Consider the three cases:

- 547 • (1) At a given time-step, the prediction set $C^t(x_{\text{test}}^t)$ does not contain the true label, $\bar{\beta}_t(\bar{x}, \bar{Y})[t]$.
- 548 • (2a) The prediction set is a singleton and does contain the true label.
- 549 • (2b) The prediction set is not a singleton (but does contain the correct label).

550 We already bound the probability of (1) happening with the CP bound; (2a) is fine since the LLM will take
 551 the correct action; (2b) is more challenging — in this case the robot asks the human for help, and we need
 552 to make sure the human “follows” the true label, by choosing the true label in the prediction set with the
 553 highest confidence by \hat{f} . In practice, we present the labels ranked by \hat{f} and ask the human to choose the
 554 true label with the highest rank.

555 Now let’s derive the bound in [Eq. \(A17\)](#) and [Eq. \(A18\)](#). Again we need to consider the causal construction
 556 issue. As seen in [Section 3.3](#), we construct the prediction set $\bar{C}(\bar{x}_{\text{test}})$ non-causally using the score function
 557 $s_i = 1 - \hat{f}(\bar{x}_i)_{\bar{y}_i^{\max}}$ (taking minimum over steps). For a test sequence \bar{x}_{test} , we apply $\bar{\beta}_{T-1}$ to the true label
 558 set of sequences \bar{Y}_{test} to get $\bar{y}_{\text{test}}^{\max} = \bar{\beta}_{T-1}(\bar{x}_{\text{test}}, \bar{Y}_{\text{test}})$. Now suppose $\bar{y}_{\text{test}}^{\max} \in \bar{C}(\bar{x}_{\text{test}})$, then we can show
 559 $\bar{y}_{\text{test}}^{\max} \in C(\bar{x}_{\text{test}})$ with the same proof as the single-label setting, which gives us the bound.

560 Lastly we need to show the sampled test sequence from $\bar{\mathcal{D}}$ leads to a sample from $\bar{\mathcal{D}}'$ i.i.d. with $\bar{\mathcal{S}}'$. This is
 561 true with the same argument that functions of independent random variables are independent.

562 **A4 Additional Experiment Setting: Real Bimanual Setup**

563 In this example, a real bimanual setup with two Kuka IIWA 7 arms move objects on the table, with one bin
 564 at each side ([Fig. 5](#) right). The reachable workspace of each arm is limited so that one arm cannot reach the
 565 other end of the table or the other bin. Thus, there can be ambiguities in the choice of the arm depending
 566 on the task; e.g., [Fig. 5](#) shows the human asking the robot to pass over the mango, but not specifying
 567 where the human is. KNOWNO is able to capture such ambiguities and triggers clarification. We design a
 568 scenario distribution with all instructions being ambiguous (thus requiring high human intervention rate):
 569 with $\epsilon = 0.15$, the robot achieves 84% plan success with 92% help. With 10 real trials, the robot succeeds 9
 570 times while triggering help for 9 times.

571 **A5 LLM Prompt Setup**

572 Next we detail the LLM prompt setup for MCQA applied in KNOWNO. We will use the Mobile
 573 Manipulation setting from [Section 4.3](#) as the example.

574 **Multiple choice generation.** Given a scenario, we first prompt the LLM to generate four options for
 575 possible actions. We apply few-shot prompting as shown in [Fig. A1](#) below with zero temperature. In this
 576 scenario, there is a Coke, a bottled tea, and a Pepsi on the counter, and the task is to put the Coke in the top
 577 drawer but the choice of drawer is under-specified (“Put the Coke in the drawer please.”).

We: You are a robot operating in an office kitchen. You are in front of a counter with two closed drawers, a top one and a middle one. There is also a landfill bin, a recycling bin, and a compost bin.

We: On the counter, there is an orange soda, a Pepsi, and an apple.
We: Put that drink in the top drawer.
You:
A) open the top drawer and put the orange soda in it
B) open the middle drawer and put the Pepsi in it
C) open the middle drawer and put the orange soda in it
D) open the top drawer and put the Pepsi in it

We: On the counter, there is an energy bar, a banana, and a microwave.
We: Put the snack next to the microwave.
You:
A) pick up the energy bar and put it next to the microwave
B) pick up the banana and put it next to the energy bar
C) pick up the banana and put it next to the microwave
D) pick up the energy bar and put it next to the banana

We: On the counter, there is a Coke, a Sprite, and a sponge.
We: Can you dispose of the can? It should have expired.
You:
A) pick up the sponge and put it in the landfill bin
B) pick up the Coke and put it in the recycling bin
C) pick up the Sprite and put it in the recycling bin
D) pick up the Coke and put it in the landfill bin

We: On the counter, there is a bottled water, a bag of jalapeno chips, and a bag of rice chips.
We: I would like a bag of chips.
You:
A) pick up the bottled water
B) pick up the jalapeno chips
C) pick up the kettle chips
D) pick up the rice chips
(The correct option is either B or D, since either jalapeno chips or rice ships are fine.)

We: On the counter, there is a Coke, a bottled unsweetened tea, and a Pepsi.
We: Put the Coke in the drawer please.
You:

Figure A1: Prompt used for multiple choice generation in the Mobile Manipulation setting.

578 After the LLM generates four options, we append an additional option ‘an option not listed here’ to the
579 four generated ones and then randomize the order to further prevent bias. We then use a zero-shot prompt
580 in Fig. A2 for querying next-token probabilities (‘A’, ‘B’, ‘C’, ‘D’, ‘E’).

We: You are a robot operating in an office kitchen. You are in front of a counter with two closed drawers, a top one and a middle one. There is also a landfill bin, a recycling bin, and a compost bin.

We: On the counter, there is a Coke, a bottled unsweetened tea, and a Pepsi.
We: Put the Coke in the drawer please.
You:
A) pick up the coke
B) pick up the coke and put it in the top drawer
C) pick up the coke and put it in the bottom drawer
D) a different option not listed here
E) pick up the pepsi
We: Which option is correct? Answer with a single letter.
You:

Figure A2: Prompt used for next-token prediction with generated multiple choices in the Mobile Manipulation setting.

581 A6 Additional Experiment Details

582 **Environments.** In addition to Fig. 1 and Fig. 5, here Fig. A3 shows the office kitchen environment with
583 the set of drawers and bins used in the Mobile Manipulator experiments (left), and the bimanual setup with
584 the set of objects used on the mat (right). There is another set of drawers used in the mobile manipulation
585 experiments underneath a much bigger countertop not shown here.

586 **Tasks and instructions.** Next, we provide more details on the task settings, in particular, the possible
587 ambiguities:

- 588 • Attribute ambiguities in Simulated setting: besides non-ambiguous terms like “green”, “yellow”, “blue”,
589 “block” and “bowl” (“put green block in yellow bowl”), refer to the block as one of “cube”, “cuboid”,
590 “box”, “square object”, to the bowl as one of “container”, “round object”, “receptacle”, or to either block
591 or bowl as one of “object”, “item”, “thing” (“move the blue object in yellow bowl”); refer to “blue” as
592 one of “cyan”, “navy”, to “green” as one of “greenish”, “grass-colored”, and to “yellow” as “orange” or
593 “gold”.
- 594 • Numeric ambiguities in Simulated setting: besides non-ambiguous terms like “a”, “one”, “a single of”,
595 “two”, “a pair of”, “three”, “all” (“put a block in yellow bowl”), refer to either two or three numerically
596 with one of “a few”, “a couple of”, “some”, “a handful of” (“put some blocks in the green bowl”).



Figure A3: (Left) Office kitchen environment with drawers and bins for the Mobile Manipulation setting. (Right) Bimanual setup with the set of objects used in the experiments.

- Spatial ambiguities in Simulated setting: besides non-ambiguous terms like “in front of”, “behind”, “to the left”, and “to the right” (“put the green block to the left of green bowl”), refer to any of the four possible directions with “near”, “close to”, “beside”, “next to”, refer to either left to right with “lateral to”, and refer to either front or behind with “along the line of sight”.
- Real Tabletop Rearrangement setting: we split the 28 toy items (Fig. A4) into two categories of human liking them or disliking them: the things the human likes include corn, avocado, celery, carrot, tomato, lettuce, apple, orange, pear, lemon, peanut butter, sunny-side-up egg, egg, and pea; the human dislikes pretzel, cracker, waffle, mustard, ketchup, pizza, meat patty, cheese, chicken drumstick, peach, mango M&M, Skittles, and donut.
- Real Mobile Manipulator setting: please refer to <https://tinyurl.com/robot-help> for the full list.
- Bimanual setting: please refer to <https://tinyurl.com/robot-help> for the full list.

Next we provide more details on some of the baselines.

Baselines - Ensemble Set. Admittedly, our ensemble-based method is a weaker method than the traditional model-based ensemble where multiple copies of neural network are trained and inferred with; however, this is infeasible with the LLM we use. In our work, we randomize over the few-shot examples in the prompt as the ensemble. We select a pool of 20 possible MCQA examples (see examples in Fig. A1), and then randomly sample a certain amount from it for each inference. Note that in this case, Ensemble Set actually has advantage over KNOWNO and Simple Set that, for the same data, it has seen many more examples than the fixed ones in the prompt used in KNOWNO and Simple Set. We only apply ensemble for next-token prediction; the same set of multiple choices generated is used.

Baselines - Prompt Set. MCQA is also applied. Similar to Fig. A1, in the few-shot examples in the prompt, we show the true possible labels. For example, “We: Which options are possibly correct? You: A, C, D.”

Baselines - Simple Set. Instead of MCQA, the LLM is first prompted to give the most likely action (e.g., “We: Put the Coke can in the drawer. You: I will” as the prompt). And then we attach the generated response to the same prompt, and ask LLM to label “Certain/Uncertain:”, given few-shot examples.

A7 Additional Implementation Details

While the focus of KNOWNO is mainly on providing uncertainty alignment for the LLM-based planner, below we provide details of the perception and action modules applied in all examples.

Perception. For all settings except for the Mobile Manipulator, we use either MDETR [52] (UR5 tabletop setting) or Owl-ViT [53] (Simulated and Bimanual settings) open-vocabulary object detector for recognizing the objects in the environment and obtaining the object locations for low-level action. In Simulated and Bimanual settings, the variations of the object types are limited, and with general prompting, the objects are detected without issue. In the UR5 tabletop setting, since we are use a wide variety of toy items (Fig. A4 right), the detector has issues often differentiating objects like peanut butter and meat patty that are both darker colors. We modify the scenario distributions to avoid using such items together in one scenario. In addition, we apply the Segment Anything model [54] to extract the object segmentation masks (shown overlaid in Fig. A4 left), and then use the polylabel algorithm [55] to find the most distant internal point of the mask as the suction point (shown as red dots).

Low-level action. In Simulated setting and UR5 tabletop setting, simple pick-and-place actions are executed based on object locations and solving the inverse kinematics. In Bimanual setting, the reachability of the Kuka arm is limited, and the pick-and-place action trajectories are solved using Sequential Quadratic



Figure A4: (Left) MDTER [52] object detection with Segment Anything [54] and most distant internal point (red dots) for UR5 tabletop setting. (Right) The total 28 toy items used for the experiments.

Programming (SQP) instead [56]. In the Mobile Manipulator setting, for most of the tasks that involve simple pick-and-place and opening the drawers, the action is from an end-to-end policy from the RT-1 policy (please refer to [57] for details), which takes in the raw observation. For some of the hard tasks such as putting the plastic bowl in the microwave and putting the metal bowl on the bowl, object locations are assumed known and we use scripted action policies.

A8 Additional Discussions

Sentence-level score leads to worse performance. In Section 2 we hypothesize that the distribution of probabilities (perplexity) of LLM outputs $p(y)$ is highly sensitive to the output length. Here we explore the effect of using sentence output and the perplexity score for CP. We still apply multiple choice generation first to obtain the possible options from LLM, and then query LLM scoring, for example, the probability of “put the blue block in the green bowl” with the prompt ending with “I will”. Table A1 shows that for all three settings, using CP with perplexity leads to worse performance, and performance degradation correlates with variance of the multiple choice lengths.

Setting	Variance	Method	Set Size Help	
Attribute	1.52	KNOWNo	1.18	0.18
		CP w/ Perplexity	1.33	0.32
Spatial	2.81	KNOWNo	2.23	0.69
		CP w/ Perplexity	2.50	0.82
Numeric	8.51	KNOWNo	2.17	0.79
		CP w/ Perplexity	4.06	1.00

Table A1: Comparison of KNOWNo with CP with sentence output and perplexity score in the three settings in the Simulated setting. $\epsilon=0.15$.

Potentially stronger baselines with model fine-tuning. In Section 4 we introduce the two prompt-based baselines Prompt Set and Binary, and demonstrate them being (1) inflexible (not allowing controlling the target success rate) and (2) do not properly model the uncertainty. We note that these two baselines can be potentially strengthened by fine-tuning the LLM to better predict the binary uncertainty or the uncertainty set, if the true labels can be properly defined. In fact, some recent work [37, 36] have explored model fine-tuning and exhibiting the effectiveness of Binary for uncertainty calibration. We also explored fine-tuning the GPT3 model (*davinci*) from OpenAI, which is the most powerful one from OpenAI available for fine-tuning. However, we find the model performing at very low accuracy with MCQA, and fine-tuning the model always results in overfitting to the dataset, even with thousands of data and varying hyperparameters (including ones from [37]). We suspect that our scenarios exhibit high complexity and variance, and it is non-trivial to fine-tune the model well with our dataset. Nonetheless, we do hope to have future work looking into better training the model for proper uncertainty, and then applying CP on top of it to achieve set-based calibration.