

# Continuously Steering LLMs Sensitivity to Contextual Knowledge with Proxy Models

Anonymous ACL submission

## Abstract

In Large Language Models (LLMs) generation, there exist knowledge conflicts, and scenarios where parametric knowledge contradicts knowledge provided in the context. Previous works studied tuning, decoding algorithms, or locating and editing context-aware neurons to adapt LLMs to be faithful to new contextual knowledge. However, they are usually inefficient or ineffective for large models, not workable for black-box models, or unable to continuously adjust LLMs' sensitivity to the knowledge provided in the context. To mitigate these problems, we propose CSKS (Continuously Steering Knowledge Sensitivity), a simple framework that can steer LLMs' sensitivity to contextual knowledge continuously at a lightweight cost. Specifically, we tune two small LMs (i.e. proxy models) and use the difference in their output distributions to shift the original distribution of an LLM without modifying the LLM weights. In the evaluation process, we design synthetic data and fine-grained metrics to measure models' sensitivity to contextual knowledge. Extensive experiments demonstrate that our framework achieves continuous and precise control over LLMs' sensitivity to contextual knowledge, enabling both increased sensitivity and reduced sensitivity, thereby allowing LLMs to prioritize either contextual or parametric knowledge as needed flexibly.

## 1 Introduction

Large Language Models (LLMs) have shown impressive capabilities in storing knowledge in their parameters (parametric knowledge) (Petroni et al., 2019; Burns et al., 2023). However, the parametric knowledge is far from reliable and correct, as it can become outdated or incorrect due to the rapid evolvement of knowledge over time or noise in the training data (Liska et al., 2022; Luu et al., 2022). This leads to knowledge augmentation methods such as retrieval-augmented generation (RAG) to

provide extra information in context (Lewis et al., 2020). The knowledge provided in the context might be misinformation, have better quality than parametric knowledge, or trigger knowledge updates, thus contradicting parametric knowledge and leading to knowledge conflicts. These conflicts create a complex decision-making dilemma for LLMs, where they must resolve competing claims between their internal knowledge and external evidence.

Previous works show that LLMs may fail to be sensitive to knowledge provided in the context depending on factors including knowledge popularity, quality, and model size (Mallen et al., 2023; Xie et al., 2024). This can contribute to wrong generation results or hallucination (Niu et al., 2024), especially in cases where the knowledge in the context is of high quality or more up-to-date. To mitigate this, decoding strategies (Shi et al., 2024b; Yuan et al., 2024), neuron-editing (Shi et al., 2024a), and prompting or tuning-based approaches (Wang et al., 2024b) are proposed to improve the LLMs' sensitivity to contextual knowledge. Nevertheless, neuron-editing and tuning-based approaches are inefficient for larger LMs and not workable for some black-box models, while all of these methods can be ineffective for stubborn LLMs with strong beliefs in their parametric knowledge. Finally, they fail to steer models' sensitivity to contextual knowledge precisely and continuously, which is critical when the quality of external information varies.

To this end, we introduce a simple framework, CSKS, to continuously adjust LLMs' sensitivity to context while being effective and efficient. Smaller models are usually much easier to adapt to our intentions through tuning, so CSKS begins with choosing two small LMs (e.g. 7b models) and fine-tuning them to make one faithful to contextual knowledge while the other faithful to its parametric knowledge. Then it shifts the original distribution of a larger LM (e.g. 72b model) with the difference between the output distributions of the two smaller models

083 multiplying a hyperparameter  $\alpha$ . When varying the  
084 hyperparameter  $\alpha$ , the logits shift toward semantics  
085 that pay more attention to contextual information  
086 changes, thus achieving continuous control over  
087 the sensitivity to contextual knowledge.

088 To give a fine-grained evaluation of how sen-  
089 sitive LLMs are to knowledge in the context, we  
090 further design synthetic QA data and define the ex-  
091 tent of knowledge conflict from three dimensions,  
092 specifically, degree of perturbation, contextual de-  
093 tail, and popularity. The three dimensions are each  
094 attributed to several ranked levels, where higher  
095 ranks indicate greater difficulty in resolving knowl-  
096 edge conflicts. Then we aggregate the ranks across  
097 all three dimensions if the question is answered  
098 correctly, resulting in a *Sensitivity Score* other than  
099 accuracy, which gives a more fine-grained evalua-  
100 tion of sensitivity to contextual knowledge.

101 Extensive experiments demonstrate that our  
102 CSKS framework surpasses state-of-the-art base-  
103 lines on large LMs under our synthetic evaluation  
104 setup while being lightweight and more accessible.  
105 Our method also provides precise and continuous  
106 control over LLMs’ sensitivity to the knowledge  
107 provided in the context, which is a key feature re-  
108 quired in many application scenarios such as RAG  
109 systems with varying context quality.

## 110 2 Methodology

### 111 2.1 CSKS Framework

112 **Building Proxy Models** The first step is to build  
113 the proxy models by fine-tuning two small LMs:  
114 one positive model  $\mathcal{P}$  which is predominantly faith-  
115 ful to the contextual knowledge, and one negative  
116 model  $\mathcal{N}$ , which adheres to its parametric knowl-  
117 edge. The size of the small models we selected is  
118 almost one-tenth of that of the target LM and we  
119 do not require the two small models and the large  
120 target model to belong to the same model family  
121 (shared architecture), as long as they have the same  
122 vocabulary (shared tokenization schemes). How-  
123 ever, for simplicity in the experiments of this paper,  
124 we use small models belonging to the same family  
125 as the target model to adjust the target model.

126 We use the ECQA dataset (Aggarwal et al., 2021)  
127 and apply different processing methods to construct  
128 two fine-tuning datasets, each containing 7,568  
129 samples. Details of the fine-tuning data and settings  
130 are provided in Appendix A. We then fine-tune the  
131 small LMs on the curated dataset.

**Steering with Proxy Models** Then, we factor  
out the context knowledge from the two small mod-  
els’ output distribution contrastively. For the large  
model  $\mathcal{L}$ , at each time step, we operate on its out-  
put distribution by adding a scaled differential term  
derived from the outputs of  $\mathcal{P}$  and  $\mathcal{N}$ . Intuitively,  
this process amplifies the importance of contextual  
information in determining the next token distribu-  
tion. The degree of amplification can be controlled  
by adjusting a hyperparameter  $\alpha$ , which scales the  
differential term.

Formally, given a query  $q$  and a context  $c$  that  
may contain some conflict to the target model’s  
internal knowledge, we generate a response  $\mathcal{X}$   
through our CSKS Framework. At each time step  $t$ ,  
we condition the raw large model  $\mathcal{L}$ , the positive  
model  $\mathcal{P}$ , and the negative model  $\mathcal{N}$  on the query  $q$ ,  
the context  $c$  and the previous response  $\mathcal{X}_{<t}$ . This  
gives us the distribution scores  $\mathcal{D}_{\mathcal{L}}$ ,  $\mathcal{D}_{\mathcal{P}}$  and  $\mathcal{D}_{\mathcal{N}}$ ,  
respectively. The response at step  $t$  can be directly  
sampled (autoregressively) from the adjusted distri-  
bution. Specifically, the response at each time step  
is computed as:

$$\tilde{\mathcal{X}}_t \sim \text{softmax} [\mathcal{D}_{\mathcal{L}} + (\mathcal{D}_{\mathcal{P}} - \mathcal{D}_{\mathcal{N}}) * \alpha],$$

where  $\alpha$  is a controlling factor that adjusts the  
influence of the context on the final output.

As illustrated in Figure 1, the framework begins  
by fine-tuning proxy models. Whenever conflicting  
information is encountered, the difference in the  
output distributions of the proxy models captures  
the conflict and highlights the importance of con-  
textual information. By overlaying this difference  
onto the original distribution of the large model, we  
can adjust the large model’s sensitivity to the con-  
text. The degree of adjustment can be controlled  
via the hyperparameter  $\alpha$ .

### 168 2.2 Evaluation Method

169 To evaluate a model’s ability to integrate new  
170 knowledge amidst conflicting internal beliefs, we  
171 design a pipeline for creating a dedicated evaluation  
172 dataset. This allows for precise grading of problem  
173 difficulty and fair performance assessment.

174 The pipeline starts with an existing QA dataset.  
175 The target LLM is prompted to answer the ques-  
176 tions in a closed-book setting. Correct answers  
177 are retained, while incorrect ones are discarded, as  
178 they often result from random hallucinations. The  
179 correct answers reflect the model’s strong internal  
180 beliefs and form the basis for introducing conflicts  
181 in later steps.

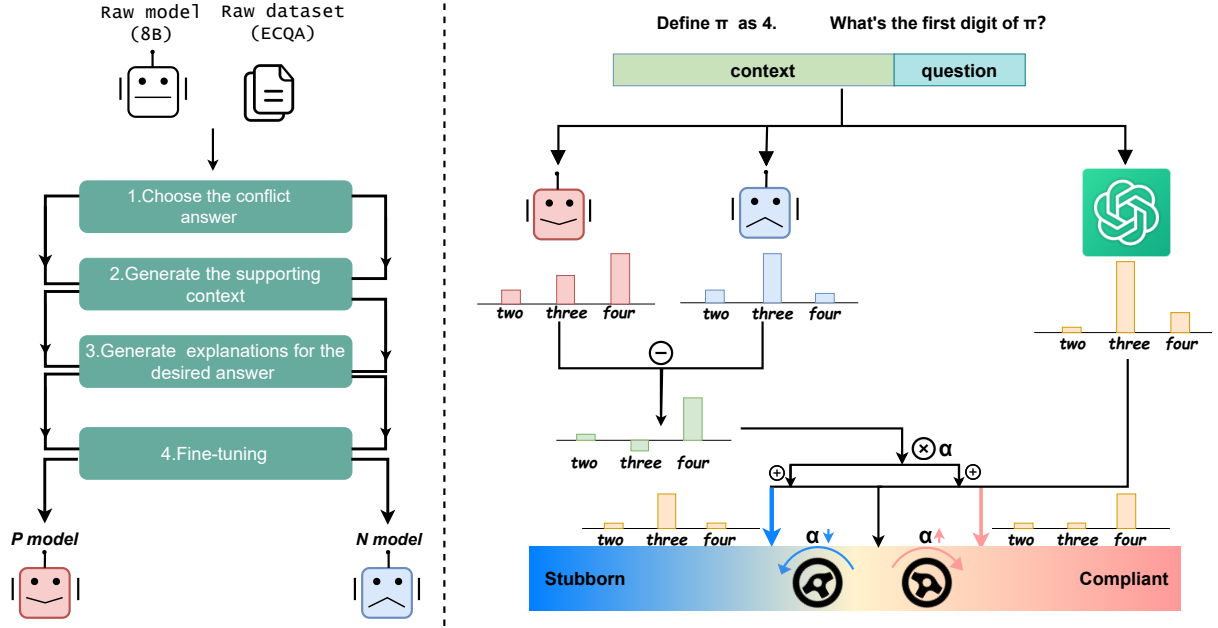


Figure 1: (left) The pipeline we use to build the proxy models, where each box represents a processing step. The two paths on either side correspond to different processing methods applicable to the proxy models. Details are shown in Appendix A. (right) When confronted with conflicting contexts, the proxy models function together as a guiding "steering wheel", assisting the large model in aligning more closely with the contextual knowledge. Additionally, we can control the degree of guidance through the parameter  $\alpha$  continuously and precisely.

Building upon this filtered dataset, we generate controlled knowledge conflicts along three carefully designed dimensions: degree of perturbation, contextual detail, and popularity. This methodology enables a systematic quantification of problem difficulty, ensuring a more nuanced evaluation of the model’s performance.

**Degree of Perturbation** The degree of perturbation reflects the extent to which external knowledge deviates from the model’s original parametric knowledge. We introduce a metric called *perturbation rank* to quantify this deviation:

- **Rank 1 (Minor Perturbation):** Involves intra-category substitutions that maintain semantic coherence and ontological consistency, preserving the original knowledge structure while introducing controlled variations.
- **Rank 2 (Major Perturbation):** Characterized by cross-category substitutions that violate fundamental ontological constraints, creating semantic inconsistencies that challenge the model’s ability to reconcile conflicting knowledge representations.

**Contextual Detail** Based on the perturbed knowledge, we generate context to support it. To system-

atically evaluate knowledge conflict resolution under varying informational conditions, we develop a dual-level *context rank* metric that operationalizes textual complexity:

- **Rank1 (Single Sentence):** Minimalist presentation of conflicting knowledge through atomic factual statements, maximizing propositional clarity while minimizing explanatory scaffolding.
- **Rank2 (Paragraph):** Extended contextualization incorporating evidentiary support, causal reasoning, and argumentative reinforcement to simulate real-world knowledge presentation patterns.

**Popularity** We use the frequency in the training corpus as an approximation of knowledge popularity. Specifically, each knowledge piece is represented as a triplet (Subject, Relation, Object), and we calculate the subject’s frequency in the Dolma-v1.7 corpus (4.5 TB) using Infini-gram (Liu et al., 2024b). A higher frequency suggests the model encountered the subject more often during pretraining, leading to a stronger internal belief and reduced sensitivity to conflicting external knowledge. We define the popularity rank as follows:

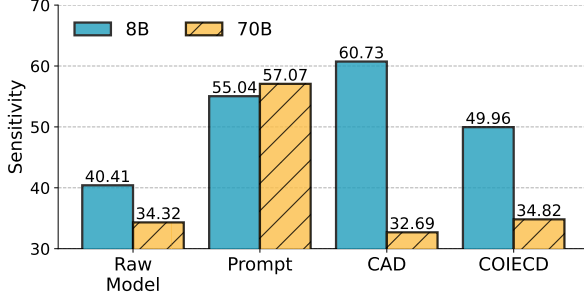


Figure 2: Performance of models of different sizes under different methods. The larger model tends to stick to its internal beliefs when faced with conflicting information. Prompting benefits both model sizes, while CAD and COIECD show excellent performance on the small model but provide minimal improvement for the large model.

- **Rank 1 (Low):** Bottom 33% ( $\leq 10^3$  occurrences)
- **Rank 2 (Medium):** Middle 33% ( $10^3 \sim 10^5$  occurrences)
- **Rank 3 (High):** Top 33% frequency ( $\geq 10^5$  occurrences)

Finally, we define the *Difficulty Score* of each question as the sum of its three constituent ranks. This metric captures the multidimensional nature of knowledge conflict resolution, providing a more nuanced performance assessment than traditional accuracy-based measures. The *Sensitivity Score* for a model is then defined as the cumulative difficulty score of all correctly answered questions, normalized by the maximum possible score. We utilize GPT-4o-mini (OpenAI, 2024) to automate this pipeline above and provide prompt templates in Appendix E. Besides, to prove the effectiveness of this grading system, we provide a validation experiment in Appendix A.

### 2.3 Motivation

Here, we’d like to illustrate the motivation that drives us to propose our CSKS framework: To gain insights into the performance of models with varying sizes or equipped with different methods (methods details are stated in section 3.1), we conduct a preliminary experiment to evaluate their ability to faithfully adhere to the knowledge provided in the context of our synthetic dataset. The results are presented in Figure 2. We observe that:

- LMs with larger sizes tend to exhibit greater rigidity compared to smaller models, indicat-

ing that large models are more stubborn when faced with knowledge conflicts.

- The CAD and COIECD methods significantly enhance the small model’s capabilities, but their ability to follow context seems to be unchanged or even diminish slightly for larger models. Therefore, the internal beliefs of small models are more easily changed, whereas large models struggle to overcome the biases of their parametric knowledge on their own.

Drawing on these observations, we propose the CSKS framework, which strategically leverages the superior adaptability of small models as proxies to guide larger language models toward better contextual knowledge integration.

## 3 Experiments

### 3.1 Baselines

We adopt representative baselines of three types, specifically, prompting, decoding-time strategy, and neuron-editing method:

- **Origin:** refers to naive LLMs without any modifications.
- **Prompt:** prompts LLMs with explicit instructions to ensure their answers align with the given context.
- **IRCAN** (Shi et al., 2024a): identifies context-responsive neurons within the LLM’s feed-forward network (FFN) layers and enhances their activation to improve the utilization of contextual information.
- **CAD** (Shi et al., 2024b): is a decoding-time strategy that adjusts the output probabilities of LLMs to emphasize differences between context-aware and context-agnostic scenarios.
- **COIECD** (Yuan et al., 2024): adapts its decoding strategy based on a contextual information-entropy constraint to discern when a context generates conflicting knowledge with the model’s internal knowledge.

For CAD and COIECD, we use the optimal hyperparameters reported in their papers for baselines. For our method, we do not search for an optimal parameter but just setting  $\alpha$  the to same as CAD. To check whether these baselines are effective, we



Methods	Degree of Perturbation(in %)		Contextual Detail(in %)		Popularity(in %)			Sensitivity Score
	rank 1	rank 2	rank 1	rank 2	rank 1	rank 2	rank 3	
<i>MusiQue • LLaMA-3-Instruct</i>								
Origin	64.85	20.17	55.08	30.00	49.44	42.63	35.71	38.13
PROMPT	75.88 (+11.03)	38.73 (+18.56)	69.22 (+14.14)	45.44 (+15.44)	65.92 (+16.48)	58.03 (+15.40)	48.26 (+12.55)	53.10 (+14.97)
CAD	62.10 (-2.65)	19.88 (-0.29)	51.69 (-3.39)	30.44 (+0.44)	47.66 (-1.78)	40.62 (-2.01)	35.06 (-0.65)	37.04 (-1.09)
COIECD	65.00 (+0.15)	20.32 (+0.32)	54.49 (-0.59)	30.88 (+0.88)	49.67 (+0.23)	42.64 (+0.01)	35.93 (+0.22)	38.35 (+0.22)
CSKS	78.08 (+13.23)	60.38 (+40.21)	79.97 (24.89)	58.53 (28.53)	75.27 (+25.83)	65.84 (+23.21)	66.66 (+30.95)	66.72 (+28.59)
<i>MusiQue • Qwen2.5-Instruct</i>								
Origin	69.85	23.71	57.29	36.32	53.00	47.54	40.04	42.58
PROMPT	76.76 (+6.91)	36.08 (+12.37)	67.60 (+10.31)	45.29 (+8.97)	62.81 (+9.81)	58.48 (+10.94)	48.27 (+8.23)	52.32 (+9.74)
CAD	82.20 (+12.35)	57.88 (+34.17)	76.58 (+19.29)	63.53 (+27.21)	75.27 (+22.27)	67.18 (+19.64)	67.74 (+27.70)	67.68 (+25.20)
COIECD	69.85 (+0.00)	24.74 (+1.03)	57.58 (+0.29)	37.06 (+0.74)	53.45 (+0.45)	47.54 (+0.00)	41.13 (+1.09)	43.21 (+0.63)
CSKS	94.85 (+25.00)	85.13 (+61.42)	90.43 (+33.14)	89.56 (+53.24)	93.54 (+40.54)	85.94 (+38.40)	90.47 (+50.43)	89.26 (+46.68)
<i>PopQA • LLaMA-3-Instruct</i>								
Origin	52.04	23.62	52.21	23.48	43.14	37.29	33.22	34.32
PROMPT	72.99 (+20.95)	46.91 (+23.29)	74.50 (+22.29)	45.42 (+21.94)	60.20 (+17.06)	61.53 (+24.24)	58.18 (+24.96)	57.07 (+22.75)
CAD	47.63 (-4.41)	24.12 (+0.50)	49.94 (-2.27)	21.85 (-1.63)	39.80 (-3.34)	36.85 (-0.44)	31.17 (-2.05)	32.69 (-1.63)
COIECD	53.03 (+0.99)	23.62 (+0.00)	52.43 (+0.22)	24.26 (+0.78)	43.31 (+0.17)	38.13 (+0.84)	33.71 (+0.49)	34.82 (+0.50)
CSKS	69.79 (+17.75)	65.45 (+41.83)	80.46 (+28.25)	54.80 (+31.32)	66.72 (+23.58)	67.72 (+30.43)	68.40 (+35.18)	66.24 (+31.92)
<i>PopQA • Qwen2.5-Instruct</i>								
Origin	66.15	28.59	60.60	34.18	51.67	47.83	42.79	43.59
PROMPT	75.63 (+9.48)	40.17 (+11.58)	71.85 (+11.25)	43.99 (+9.81)	58.86 (+7.19)	57.86 (+10.03)	57.05 (+14.26)	54.63 (+11.04)
CAD	78.06 (+11.91)	61.15 (+32.56)	78.04 (+17.44)	61.19 (+27.01)	70.73 (+19.06)	69.23 (+21.40)	68.88 (+26.09)	67.80 (+24.21)
COIECD	65.82 (-0.33)	28.04 (-0.55)	59.49 (-1.11)	34.40 (+0.22)	50.50 (-1.17)	47.32 (-0.51)	43.11 (+0.32)	43.31 (-0.28)
CSKS	93.83 (+27.68)	90.40 (+61.81)	93.27 (+32.67)	90.96 (+56.78)	88.46 (+36.79)	93.14 (+45.31)	94.65 (+51.86)	92.24 (+48.65)

Table 1: Accuracy when evaluated on specific ranks of individual dimensions in the dataset and the overall *Sensitivity Score*. For each dimension, Rank 1 represents the least challenging cases, while higher ranks indicate increasing difficulty. CSKS outperforms baseline methods under all metrics.

conducted a verification on small model. The results are presented in Appendix C, which shows that while all baseline methods work fine for the small model, IRCAN shows minimal performance enhancement. This limited efficacy combined with IRCAN’s significantly larger computational overhead makes it unsuitable for our primary objective of efficient large-model adaption. So we exclude IRCAN from our main experiments.

### 3.2 Models and Settings

We employ two state-of-the-art instruction-tuned LLMs as target models: Llama-3-70B-Instruct (Dubey et al., 2024) and Qwen2.5-72B-Instruct (Yang et al., 2024). For each target model, we utilize its smaller counterpart as proxy model – specifically, fine-tuned versions of Llama-3-8B-Instruct for the Llama-3 series and Qwen2.5-7B-Instruct for the Qwen2.5 series. We use greedy decoding in all the experiments to ensure reproducibility.

For constructing the evaluation dataset, we use MuSiQue (Trivedi et al., 2022) and PopQA (Mallen et al., 2023), both widely used question-answering datasets as the source datasets. Following the

setup in Shi et al. (2024a), we frame the task as a multiple-choice format. For evaluation purposes, we organize the data into binary-choice questions, where the correct options correspond to the answers in context, and the incorrect options correspond to the original answers to the question. This design creates controlled knowledge conflict scenarios where model performance directly reflects its ability to prioritize contextual or parametric knowledge. It is important to clarify that the contextual answers used here are exactly the perturbed answers we introduce during dataset construction.

To comprehensively evaluate the model’s performance across the entire dataset, we use accuracy as a default metric, calculated for each rank within our three operational dimensions (perturbation, context, popularity). Additionally, we employ the previously defined *Sensitivity Score* to assess the model’s ability to adhere to the given context, which is also normalized into a 100-scale.

### 3.3 Results

As demonstrated in Table 1, our proposed CSKS consistently advances all baselines across all evalu-

ation dimensions. CSKS outperforms baseline methods by substantial margins, with 30.26 average sensitivity score improvement for LLaMA-3 and 47.67 for Qwen2.5. Besides, we have two other main observations:

1. **Baseline Limitations:** The decoding-time strategy baselines exhibit inconsistent effectiveness. While CAD shows moderate gains on Qwen2.5 (+24.2 sensitivity score), it degrades performance on LLaMA-3 (-1.1 sensitivity score). COIECD’s entropy-based constraints prove insufficient for resolving deep parametric conflicts, yielding marginal improvements of less than 1.5 across all configurations. The core idea behind CAD and COIECD is to leverage the output distribution differences between the model’s responses with and without context to emphasize the importance of contextual information (i.e. one model with different data). Our results suggest that large models may not be able to overcome the biases of their internal knowledge on their own.
2. **Dimensional Sensitivity:** Among the three dimensions we introduce, the perturbation degree has the greatest effect. This might be because a large perturbation creates an obvious conflict with the model’s internal knowledge, forcing it to confront and resolve the inconsistency directly. On the other hand, small perturbations are more confounding, as they subtly deviate from the truth, making it harder for the model to determine whether to trust the external context or rely on its internal knowledge. The perturbation degree has the lowest effect. Under our method, the differences between different ranks of popularity are smoothed out or even reversed, which indicates that our method has sufficient ability to eliminate the intrinsic knowledge bias brought by the model during pre-training.

After demonstrating the effectiveness of CSKS framework, we further show that our framework can achieve continuous and precise control over the knowledge sensitivity to contextual knowledge through the steering parameter  $\alpha$ . As illustrated in Figure 3, increasing  $\alpha$  values ( $\alpha > 0$ ) produces a monotonic enhancement of sensitivity score from 4.32 to 39.80 for LLaMA on MuSiQue, with potential for further increase). This directional control

Alpha	STEM	Humanities	Other	Social	Average
-2.0	89.34	78.01	88.27	82.54	85.00
-1.5	90.98	77.66	88.08	83.81	85.44
-1.0	91.39	77.32	88.64	83.17	85.51
-0.7	91.39	78.69	88.64	84.13	86.01
-0.5	91.39	79.73	89.01	84.44	86.45
<b>72B(<math>\alpha = 0</math>)</b>	<b>92.62</b>	<b>79.04</b>	<b>88.64</b>	<b>84.76</b>	<b>86.45</b>
+0.5	91.80	78.01	87.71	84.44	85.65
+0.7	91.80	78.69	87.52	84.13	85.65
+1.0	90.98	78.01	87.34	83.81	85.22
+1.5	90.98	76.29	85.85	83.49	84.21
+2.0	90.98	74.91	84.92	81.27	83.06
<b>7B</b>	<b>84.84</b>	<b>70.79</b>	<b>76.35</b>	<b>76.83</b>	<b>76.78</b>

Table 2: Performance comparison showing trade-off between faithfulness to contextual knowledge and general capabilities.

proves critical for applications requiring dynamic knowledge updates, where models must suppress outdated parametric knowledge in favor of fresh contextual evidence. Results on PopQA can be found in Appendix D.)

In the previous experiments, we demonstrate the effectiveness of CSKS framework when aggregating new and conflicting knowledge in contexts setting  $\alpha > 0$ . Notably, extending  $\alpha$  to negative values ( $\alpha < 0$ ) reveals an inverse mode of action—the framework can suppress contextual influence to amplify parametric reliance. As demonstrated in Figure 3, setting  $\alpha = -2.0$  reduces contextual sensitivity score by 15.9 for LLaMA and 32.8 for Qwen compared to their baselines ( $\alpha = 0$ ), effectively transforming the target model into a parametric knowledge conservative. This bidirectional control mechanism ( $\alpha \in (-\infty, +\infty)$ ) enables continuous scenario adaptation, allowing practitioners to calibrate models for either context-sensitive scenarios or parametric knowledge preservation.

### 3.4 Analysis

**The Impact of Proxy Model Size** To study whether it is possible to use even smaller models to save more resources and achieve comparable results, we utilize the Qwen2.5 model family, which includes small models from 0.5B to 7B. We apply these models under CSKS framework to steer the 72B model and present the results in Figure 4. As shown in the figure, the impact of the 0.5B proxy model on the sensitivity score of the target model is not obvious, but there is still a growing trend. The impact of the 1.5B proxy model on the target model already becomes very significant. When the size of the proxy model increases to 3B, its impact on the target model is comparable to that of the 7B

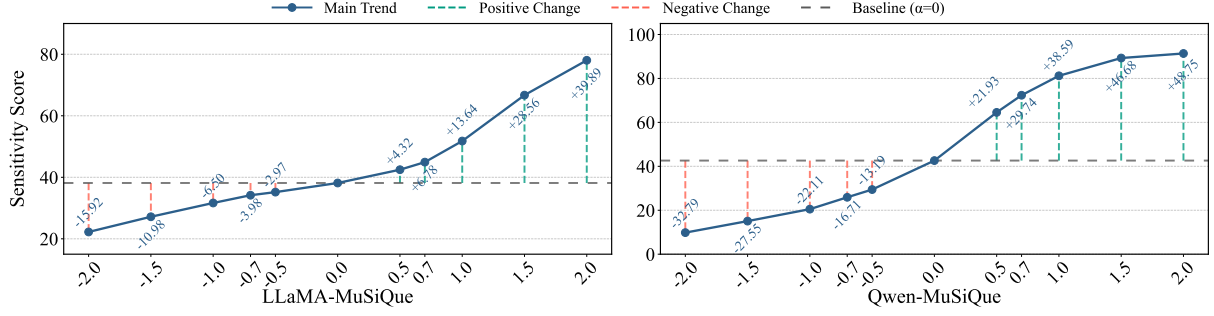


Figure 3: The performance of LLaMA and Qwen controlled bidirectionally, demonstrating the continuous adjustment capability of our method from two directions.

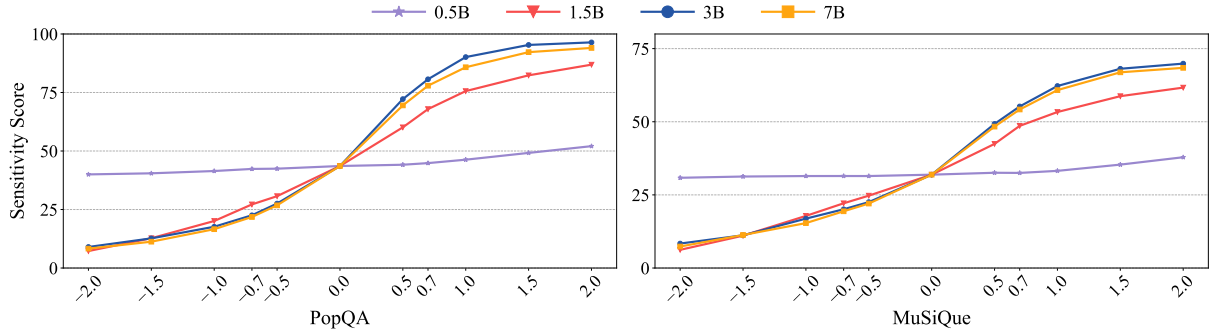


Figure 4: The performance of CSKS under varying proxy model sizes on MuSiQue and PopQA respectively. Smaller proxy models, such as the 0.5B and 1.5B versions, have a minimal but growing impact on the sensitivity score of the 72B target model. The 3B proxy model achieves a sensitivity adjustment comparable to the 7B model, demonstrating that our framework allows for significant context sensitivity modulation with much smaller models.

proxy model, and even has a slight advantage. The above results demonstrate that our framework has the potential to use a much smaller overhead (such as only using a 3B model) to perform context sensitivity adjustment on a model dozens of times larger. This efficiency may stem from our framework’s *selective steering* mechanism, where proxy models focus exclusively on context sensitivity modulation rather than full knowledge representation.

**Trade-Off Discussion** To study how scaling the control parameter  $\alpha$  would impact the general capabilities of the model, we conduct an evaluation on the MMLU benchmark (Hendrycks et al., 2021). For simplicity, we select two tasks from each of its four subjects (STEM, Humanities, Social, and Other) in the dataset as the test dataset. The experiment results in Table 2 reveal a crucial trade-off in knowledge sensitivity control: while increasing the absolute value of  $\alpha$  enables extensive adjustment of the model’s contextual sensitivity as we show in Figure 3, excessive values ( $|\alpha| > 1.5$ ) lead to noticeable degradation in general capabilities, particularly Humanities (-4.10%) domain. This performance decline suggests that extreme sensi-

tivity adjustments may disrupt the target model’s fundamental reasoning patterns, highlighting the importance of maintaining a balanced  $\alpha$  range that preserves core competencies while enabling effective knowledge adaptation. Notably, even within this kind-of-broad range, the target 72B model consistently outperforms the 7B model by significant margins (average +8.67%), demonstrating that our framework successfully leverages the large model’s superior reasoning capacity while achieving precise sensitivity control. These findings collectively indicate that strategic  $\alpha$  selection can achieve an effective equilibrium between contextual adaptability and general capability preservation, fulfilling our framework’s dual objectives of precise knowledge steering and performance maintenance.

**Extending to Black Box Model** For the black-box models that we can’t obtain weights, our framework remains effective. We apply our framework to adapt GPT-3.5-Turbo (Ouyang et al., 2022). In this setting, since we can only access log probabilities for the top five tokens through the API, CSKS only reweights the five tokens. We present the results in Table 3. For black-box models that do not belong

Raw	$\alpha = 0.5$	$\alpha = 0.7$	$\alpha = 1.0$	$\alpha = 1.5$	$\alpha = 2.0$
<i>MusiQue • Proxy-LLaMA</i>					
51.24	60.38	66.36	76.32	87.79	93.45
<i>PopQA • Proxy-Qwen</i>					
56.56	75.07	84.67	90.89	93.58	94.73

Table 3: Performance of GPT-3.5-Turbo steered by LLaMA and Qwen. Our method also works for black-box models such as GPT-3.5-Turbo.

to the same model family as the proxy model, CSKS can still effectively control its context sensitivity, demonstrating its broad application domain.

## 4 Related Works

### 4.1 Knowledge Conflicts

Knowledge conflicts refer to cases where contextual knowledge contradicts parametric knowledge (Mallen et al., 2023; Xu et al., 2024; Kortukov et al., 2024). Many previous works focus on making LLMs generate responses based on provided context rather than parametric knowledge (Gekhman et al., 2023; Lee et al., 2022; Shi et al., 2024c; Zhang et al., 2020; Zhou et al., 2023). This is a valuable setting for applications such as retrieval-augmented LLMs (Ram et al., 2023; Shi et al., 2024d), where the context may be of high quality (e.g. containing updated knowledge). However, an underexplored aspect is that the context quality may vary significantly in different working scenarios, so making the model rely on context to a constant extent is far from enough. We argue that LLMs should be controlled to rely on context to varying degrees, and the control should be precise and continuous. We propose an effective yet efficient framework to achieve this goal.

Another line of work focuses on evaluating and understanding LLMs in knowledge conflicts and mining factors affecting LLMs’ choice in knowledge conflicts. Wu et al. (2024a); Tan et al. (2024) show that the level of detail in the context will affect the choices made by language models when faced with knowledge conflicts. Xie et al. (2023) find that LLMs exhibit a predisposition towards emphasizing information related to entities of higher popularity and models demonstrate a significant sensitivity to the order in which data is introduced. Qian et al. (2024) introduce different permutation degrees to knowledge and find that models exhibit resistance to knowledge that evidently lacks veracity. Jin et al. (2024) discover that as the number

of conflicting hops increases, LLMs encounter increased challenges in reasoning. We further utilize the key factors to measure the difficulty of manipulating certain knowledge and provide a more comprehensive evaluation method.

### 4.2 Updating Knowledge in Language Models

To introduce new knowledge to LLMs, previous works explore tuning-based approaches (Wang et al., 2024b), decoding strategies (Shi et al., 2024b; Zhao et al., 2024; Wang et al., 2024a), and model editing methods (Meng et al., 2023; Gupta et al., 2023; Shi et al., 2024a). Nevertheless, these methods are usually inefficient or ineffective for large models, not workable for black-box models, or unable to continuously adjust LLMs’ sensitivity to the new contextual knowledge, while our approach can steer LLMs’ sensitivity to contextual knowledge continuously at a lightweight cost.

### 4.3 Control of Language Models

Motivated by the increasing capabilities of LLMs (Li et al., 2023b), many studies focus on controlling certain attributes of LM generation, usually non-toxicity and positive sentiment. A common solution to control LLMs is representation engineering. Han et al. (2024) use word embeddings to steer LLMs for language model detoxification and sentiment control. Zhao et al. (2024) steer knowledge behaviors of LLMs with SAE-based representation engineering. Some other works tune the hidden representations of LLMs to change behaviors (Wu et al., 2024b; Hernandez et al., 2024; Li et al., 2023a; OpenAI, 2024). Another line of work incorporates other models to guide the generation process (Liu et al., 2021, 2024a; Feng et al., 2024). Our work also borrows this idea but emphasizes controlling sensitivity to contextual knowledge and achieves precise and continuous control.

## 5 Conclusion

We present CSKS, an efficient and effective framework that leverages smaller LLMs as proxy models to shift the output distributions of LLMs, thus improving LLMs’ faithfulness to the knowledge provided in the context. We also introduce a fine-grained evaluation method for measuring LLM’s sensitivity to contextual knowledge. Extensive experiments demonstrate that our framework achieves state-of-the-art, and more importantly, achieves precise and continuous control over LLMs’ sensitivity to contextual knowledge.



## Limitations

The language models and datasets used for our experiments are not complete. We only consider two families of open-sourced LLMs, one black-box LLM, and two QA datasets. Since we will make our code and synthetic datasets publicly available, we leave it to future work on evaluating more models on more datasets. Moreover, we do not consider complex knowledge-related QA tasks such as multi-hop QA. Finally, since our experiment is done in a synthetic setting, it is unclear how our method will work in real-world applications.

## References

Shourya Aggarwal, Divyanshu Mandowara, Vishwa-jeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. 2021. [Explanations for CommonsenseQA: New Dataset and Models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3050–3065, Online. Association for Computational Linguistics.

Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. 2023. [Discovering latent knowledge in language models without supervision](#). In *The Eleventh International Conference on Learning Representations*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Shangbin Feng, Taylor Sorensen, Yuhao Liu, Jillian Fisher, Chan Young Park, Yejin Choi, and Yulia Tsvetkov. 2024. [Modular pluralism: Pluralistic alignment via multi-LLM collaboration](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4151–4171, Miami, Florida, USA. Association for Computational Linguistics.

Zorik Gekhman, Jonathan Herzig, Roei Aharoni, Chen Elkind, and Idan Szpektor. 2023. [Trueteacher: Learning factual consistency evaluation with large language models](#). In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

Anshita Gupta, Debanjan Mondal, Akshay Sheshadri, Wenlong Zhao, Xiang Li, Sarah Wiegrefe, and Niket Tandon. 2023. [Editing common sense in transformers](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8214–8232, Singapore. Association for Computational Linguistics.

Chi Han, Jialiang Xu, Manling Li, Yi Fung, Chenkai Sun, Nan Jiang, Tarek Abdelzaher, and Heng Ji. 2024. [Word embeddings are steers for language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16410–16430, Bangkok, Thailand. Association for Computational Linguistics.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.

Evan Hernandez, Belinda Z. Li, and Jacob Andreas. 2024. [Inspecting and editing knowledge representations in language models](#). In *First Conference on Language Modeling*.

Zhuoran Jin, Pengfei Cao, Yubo Chen, Kang Liu, Xiaojian Jiang, Jiexin Xu, Qiuxia Li, and Jun Zhao. 2024. Tug-of-war between knowledge: Exploring and resolving knowledge conflicts in retrieval-augmented language models. *arXiv preprint arXiv:2402.14409*.

Evgenii Kortukov, Alexander Rubinstein, Elisa Nguyen, and Seong Joon Oh. 2024. [Studying large language model behaviors under context-memory conflicts with real documents](#). In *First Conference on Language Modeling*.

Kyungjae Lee, Wookje Han, Seung-won Hwang, Hwaran Lee, Joonsuk Park, and Sang-Woo Lee. 2022. [Plug-and-play adaptation for continuously-updated QA](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 438–447, Dublin, Ireland. Association for Computational Linguistics.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.

Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023a. [Inference-time intervention: Eliciting truthful answers from a language model](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Sha Li, Chi Han, Pengfei Yu, Carl Edwards, Manling Li, Xingyao Wang, Yi Fung, Charles Yu, Joel Tetreault, Eduard Hovy, and Heng Ji. 2023b. [Defining a new NLP playground](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11932–11951, Singapore. Association for Computational Linguistics.

Adam Liska, Tomas Kocisky, Elena Gribovskaya, Tayfun Terzi, Eren Sezener, Devang Agrawal, D’Autume Cyprien De Masson, Tim Scholtes, Manzil Zaheer, Susannah Young, et al. 2022. Streamingqa: A

688	benchmark for adaptation to new knowledge over	<a href="https://openai.com/docs/guides/vision">openai.com/docs/guides/vision</a> . Accessed:	745
689	time in question answering models. In <u>International</u>	2024-05-26.	746
690	<u>Conference on Machine Learning</u> , pages 13604–		
691	13622. PMLR.		
692	Alisa Liu, Xiaochuang Han, Yizhong Wang, Yu-	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	747
693	lia Tsvetkov, Yejin Choi, and Noah A. Smith.	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	748
694	2024a. <u>Tuning language models by proxy</u> . In <u>First</u>	Sandhini Agarwal, Katarina Slama, Alex Gray, John	749
695	<u>Conference on Language Modeling</u> .	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	750
		Maddie Simens, Amanda Askell, Peter Welinder,	751
696	Alisa Liu, Maarten Sap, Ximing Lu, Swabha	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	752
697	Swayamdipta, Chandra Bhagavatula, Noah A. Smith,	<u>Training language models to follow instructions with</u>	753
698	and Yejin Choi. 2021. <u>DExperts: Decoding-time con-</u>	<u>human feedback</u> . In <u>Advances in Neural Information</u>	754
699	<u>trolled text generation with experts and anti-experts</u> .	<u>Processing Systems</u> .	755
700	In <u>Proceedings of the 59th Annual Meeting of</u>		
701	<u>the Association for Computational Linguistics and</u>	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	756
702	<u>the 11th International Joint Conference on Natural</u>	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	757
703	<u>Language Processing (Volume 1: Long Papers)</u> ,	Alexander Miller. 2019. <u>Language models as</u>	758
704	pages 6691–6706, Online. Association for Computa-	<u>knowledge bases?</u> In <u>Proceedings of the</u>	759
705	tional Linguistics.	2019 Conference on Empirical Methods in Natural	760
		Language Processing and the 9th International	761
706	Jiacheng Liu, Sewon Min, Luke Zettlemoyer, Yejin	<u>Joint Conference on Natural Language Processing</u>	762
707	Choi, and Hannaneh Hajishirzi. 2024b. <u>Infini-gram:</u>	(EMNLP-IJCNLP), pages 2463–2473, Hong Kong,	763
708	<u>Scaling unbounded n-gram language models to a</u>	China. Association for Computational Linguistics.	764
709	<u>trillion tokens</u> . In <u>First Conference on Language</u>		
710	<u>Modeling</u> .	Cheng Qian, Xinran Zhao, and Tongshuang Wu. 2024.	765
		"merge conflicts!" exploring the impacts of exter-	766
711	Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Kar-	<u>nal knowledge distractors to parametric knowledge</u>	767
712	ishma Mandyam, and Noah A. Smith. 2022. <u>Time</u>	<u>graphs</u> . In <u>First Conference on Language Modeling</u> .	768
713	<u>waits for no one! analysis and challenges of tem-</u>		
714	<u>poral misalignment</u> . In <u>Proceedings of the 2022</u>	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay,	769
715	<u>Conference of the North American Chapter of the</u>	Amnon Shashua, Kevin Leyton-Brown, and Yoav	770
716	<u>Association for Computational Linguistics: Human</u>	Shoham. 2023. <u>In-context retrieval-augmented lan-</u>	771
717	<u>Language Technologies</u> , pages 5944–5958, Seattle,	<u>guage models</u> . <u>Transactions of the Association for</u>	772
718	United States. Association for Computational Lin-	<u>Computational Linguistics</u> , 11:1316–1331.	773
719	guistics.		
720	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	Dan Shi, Renren Jin, Tianhao Shen, Weilong Dong,	774
721	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.	Xinwei Wu, and Deyi Xiong. 2024a. <u>IRCAN: Mit-</u>	775
722	<u>When not to trust language models: Investigating</u>	<u>igating knowledge conflicts in LLM generation via</u>	776
723	<u>effectiveness of parametric and non-parametric mem-</u>	<u>identifying and reweighting context-aware neurons</u> .	777
724	<u>ories</u> . In <u>Proceedings of the 61st Annual Meeting</u>	In <u>The Thirty-eighth Annual Conference on Neural</u>	778
725	<u>of the Association for Computational Linguistics</u>	<u>Information Processing Systems</u> .	779
726	<u>(Volume 1: Long Papers)</u> , pages 9802–9822,		
727	Toronto, Canada. Association for Computational Lin-	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia	780
728	guistics.	Tsvetkov, Luke Zettlemoyer, and Wen-tau Yih.	781
		2024b. <u>Trusting your evidence: Hallucinate less with</u>	782
729	Kevin Meng, Arnab Sen Sharma, Alex J Ando-	<u>context-aware decoding</u> . In <u>Proceedings of the 2024</u>	783
730	nian, Yonatan Belinkov, and David Bau. 2023.	<u>Conference of the North American Chapter of the</u>	784
731	<u>Mass-editing memory in a transformer</u> . In <u>The</u>	<u>Association for Computational Linguistics: Human</u>	785
732	<u>Eleventh International Conference on Learning</u>	<u>Language Technologies (Volume 2: Short Papers)</u> ,	786
733	<u>Representations</u> .	pages 783–791, Mexico City, Mexico. Association	787
		for Computational Linguistics.	788
734	Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu,	Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou,	789
735	KaShun Shum, Randy Zhong, Juntong Song, and	Margaret Li, Xi Victoria Lin, Noah A. Smith, Luke	790
736	Tong Zhang. 2024. <u>RAGTruth: A hallucina-</u>	Zettlemoyer, Wen tau Yih, and Mike Lewis. 2024c.	791
737	<u>tion corpus for developing trustworthy retrieval-</u>	<u>In-context pretraining: Language modeling beyond</u>	792
738	<u>augmented language models</u> . In <u>Proceedings</u>	<u>document boundaries</u> . In <u>The Twelfth International</u>	793
739	<u>of the 62nd Annual Meeting of the Association</u>	<u>Conference on Learning Representations</u> .	794
740	<u>for Computational Linguistics (Volume 1: Long</u>		
741	<u>Papers)</u> , pages 10862–10878, Bangkok, Thailand.	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	795
742	Association for Computational Linguistics.	joon Seo, Richard James, Mike Lewis, Luke	796
		Zettlemoyer, and Wen-tau Yih. 2024d. <u>RE-</u>	797
743	OpenAI. 2024. Introducing gpt-4o: our fastest and	<u>PLUG: Retrieval-augmented black-box language</u>	798
744	most affordable flagship model. <a href="https://platform.openai.com/docs/guides/vision">https://platform.</a>	<u>models</u> . In <u>Proceedings of the 2024 Conference</u>	799
		<u>of the North American Chapter of the Association</u>	800
		<u>for Computational Linguistics: Human Language</u>	801

802	Technologies (Volume 1: Long Papers), pages 8371–	Xiaowei Yuan, Zhao Yang, Yequan Wang, Shengping	859
803	8384, Mexico City, Mexico. Association for Compu-	Liu, Jun Zhao, and Kang Liu. 2024. <a href="#">Discerning</a>	860
804	tational Linguistics.	<a href="#">and resolving knowledge conflicts through adap-</a>	861
805	Hexiang Tan, Fei Sun, Wanli Yang, Yuanzhuo Wang,	<a href="#">tive decoding with contextual information-entropy</a>	862
806	Qi Cao, and Xueqi Cheng. 2024. <a href="#">Blinded by gen-</a>	<a href="#">constraint</a> . In <a href="#">Findings of the Association for</a>	863
807	<a href="#">erated contexts: How language models merge gen-</a>	<a href="#">Computational Linguistics: ACL 2024</a> , pages 3903–	864
808	<a href="#">erated and retrieved contexts when knowledge con-</a>	3922, Bangkok, Thailand. Association for Computa-	865
809	<a href="#">flicts?</a> In <a href="#">Proceedings of the 62nd Annual Meeting</a>	tional Linguistics.	866
810	<a href="#">of the Association for Computational Linguistics</a>	Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun	867
811	<a href="#">(Volume 1: Long Papers)</a> , pages 6207–6227,	Chen, Chris Brockett, Xiang Gao, Jianfeng Gao,	868
812	Bangkok, Thailand. Association for Computational	Jingjing Liu, and Bill Dolan. 2020. <a href="#">DIALOGPT</a>	869
813	Linguistics.	<a href="#">: Large-scale generative pre-training for conver-</a>	870
814	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	<a href="#">sational response generation</a> . In <a href="#">Proceedings of</a>	871
815	and Ashish Sabharwal. 2022. <a href="#">MuSiQue: Multi-</a>	<a href="#">the 58th Annual Meeting of the Association for</a>	872
816	<a href="#">hop questions via single-hop question composition</a> .	<a href="#">Computational Linguistics: System Demonstrations</a> ,	873
817	<a href="#">Transactions of the Association for Computational</a>	pages 270–278, Online. Association for Computa-	874
818	<a href="#">Linguistics</a> , 10:539–554.	tional Linguistics.	875
819	Han Wang, Archiki Prasad, Elias Stengel-Eskin, and	Yu Zhao, Alessio Devoto, Giwon Hong, Xiaotang	876
820	Mohit Bansal. 2024a. Adacad: Adaptively decoding	Du, Aryo Pradipta Gema, Hongru Wang, Xuanli	877
821	to balance conflicts between contextual and paramet-	He, Kam-Fai Wong, and Pasquale Minervini. 2024.	878
822	ric knowledge. <a href="#">arXiv preprint arXiv:2409.07394</a> .	<a href="#">Steering knowledge selection behaviours in llms</a>	879
823	Yike Wang, Shangbin Feng, Heng Wang, Weijia Shi,	<a href="#">via sae-based representation engineering</a> . <a href="#">Preprint</a> ,	880
824	Vidhisha Balachandran, Tianxing He, and Yulia	<a href="#">arXiv:2410.15999</a> .	881
825	Tsvetkov. 2024b. <a href="#">Resolving knowledge conflicts</a>	Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and	882
826	<a href="#">in large language models</a> . In <a href="#">First Conference on</a>	Muhao Chen. 2023. <a href="#">Context-faithful prompting</a>	883
827	<a href="#">Language Modeling</a> .	<a href="#">for large language models</a> . In <a href="#">Findings of the</a>	884
828	Kevin Wu, Eric Wu, and James Zou. 2024a. <a href="#">Clasheval:</a>	<a href="#">Association for Computational Linguistics: EMNLP</a>	885
829	<a href="#">Quantifying the tug-of-war between an llm’s internal</a>	<a href="#">2023</a> , pages 14544–14556, Singapore. Association	886
830	<a href="#">prior and external evidence</a> . In <a href="#">Neural Information</a>	for Computational Linguistics.	887
831	<a href="#">Processing Systems</a> .		
832	Zhengxuan Wu, Aryaman Arora, Zheng Wang, At-		
833	ticus Geiger, Dan Jurafsky, Christopher D Man-		
834	ning, and Christopher Potts. 2024b. <a href="#">ReFT: Rep-</a>		
835	<a href="#">resentation finetuning for language models</a> . In		
836	<a href="#">The Thirty-eighth Annual Conference on Neural</a>		
837	<a href="#">Information Processing Systems</a> .		
838	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and		
839	Yu Su. 2023. <a href="#">Adaptive chameleon or stubborn sloth:</a>		
840	<a href="#">Revealing the behavior of large language models in</a>		
841	<a href="#">knowledge conflicts</a> . In <a href="#">International Conference on</a>		
842	<a href="#">Learning Representations</a> .		
843	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and		
844	Yu Su. 2024. <a href="#">Adaptive chameleon or stubborn sloth:</a>		
845	<a href="#">Revealing the behavior of large language models in</a>		
846	<a href="#">knowledge conflicts</a> . In <a href="#">The Twelfth International</a>		
847	<a href="#">Conference on Learning Representations</a> .		
848	Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang,		
849	Hongru Wang, Yue Zhang, and Wei Xu. 2024.		
850	<a href="#">Knowledge conflicts for LLMs: A survey</a> . In		
851	<a href="#">Proceedings of the 2024 Conference on Empirical</a>		
852	<a href="#">Methods in Natural Language Processing</a> , pages		
853	8541–8565, Miami, Florida, USA. Association for		
854	Computational Linguistics.		
855	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,		
856	Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,		
857	Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 tech-		
858	nical report. <a href="#">arXiv preprint arXiv:2412.15115</a> .		

## A Finetune Dataset Details

To obtain our  $\mathcal{P}$  model and  $\mathcal{N}$  model, we fine-tune the Llama-3-8B-instruct model and Qwen-2.5-7B-instruct model. To ensure generalization, the fine-tuning datasets are constructed using methods and domains different from those of the synthesized conflict datasets. To achieve optimal results, we have designed a specialized pipeline for constructing the fine-tuning dataset as shown in Figure 5.

We select ECQA as the base dataset, which is a multiple-choice QA dataset where each question is accompanied by five answer options.

- For the  $\mathcal{P}$  model: We select the incorrect option least related to the correct answer as the "contextual answer."
- For the  $\mathcal{N}$  model: We select the incorrect option most related to the correct answer as the "contextual answer."

Next, using GPT, we generate supportive context based on the chosen answer and the question.

- For the  $\mathcal{P}$  model, the generated context was short and simple.
- For the  $\mathcal{N}$  model, the context was long and detailed.

Finally, we again use GPT to generate explanations based on the context, question, and selected answer.

- For the  $\mathcal{P}$  model, the explanation justified why the selected answer was correct.
- For the  $\mathcal{N}$  model, the explanation detailed why the selected answer was incorrect.

Using these constructed answers and their corresponding explanations, we fine-tune the model as follows:

- The  $\mathcal{P}$  model was fine-tuned on the selected answers and their associated explanations.
- The  $\mathcal{N}$  model was fine-tuned on the original correct answers and their explanations.

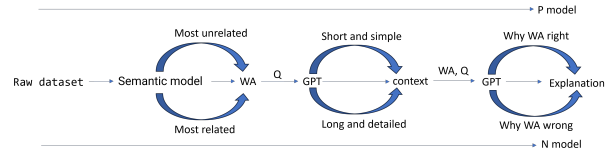


Figure 5: The pipeline to get the data used to finetune our  $\mathcal{P}$  model and  $\mathcal{N}$  model

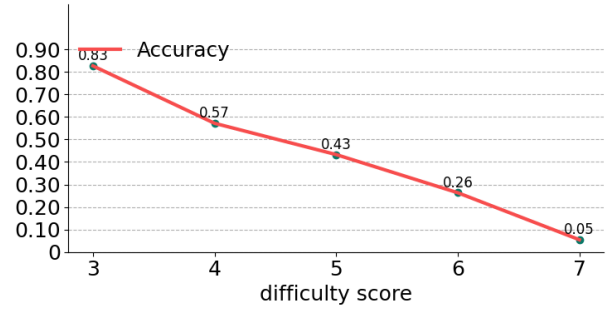


Figure 6: The accuracy of the LLaMA-3-70B-Instruct model across questions of each difficulty score.

## B Effectiveness of the Grading System

To validate the effectiveness of our grading system, we conduct a validation experiment. We analyze the accuracy of the target model across questions of varying difficulty levels, with the results shown in Figure 6. The results reveal that as question difficulty increases, accuracy correspondingly decreases. This demonstrates that our grading system successfully quantifies problem difficulty.

## C Fine-tune results on small models

Figure 7 illustrates the effects of different methods on the LLaMA-3-8B-instruct model. From the results, we observe the following:

1. The Prompt, CAD and COIECD methods all improve the performance of the 8B small model, while the impact of IRCAN on the small model's performance is minimal.
2. We also present the performance of our fine-tuned  $\mathcal{P}$  model and  $\mathcal{N}$  model. The  $\mathcal{P}$  model performs the best, as it effectively incorporates knowledge from the context, while the  $\mathcal{N}$  model scores much lower, indicating that it tends to rely on its internal knowledge and resists external contextual information. This indicates that our fine-tuning is successful.



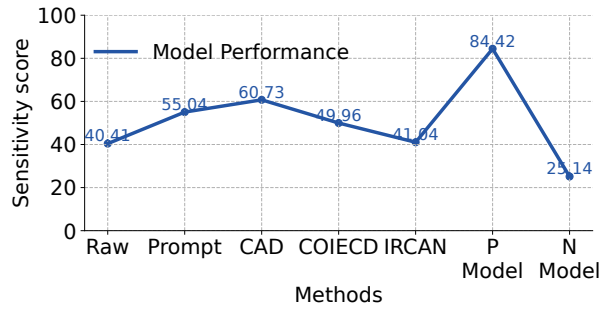
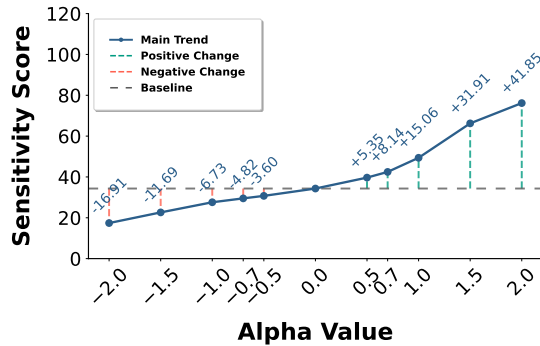


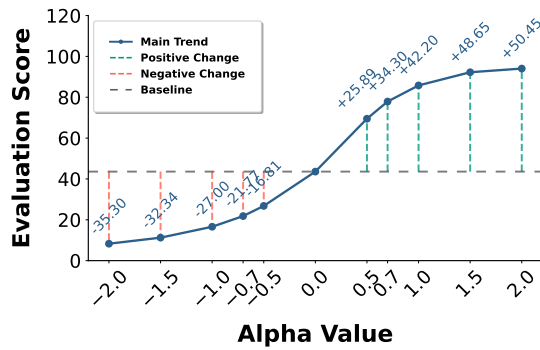
Figure 7: The effects of different methods on the LLaMA-3-8B-instruct model tested on PopQA.

## D Steering Results on PopQA

We present the steering results on the PopQA dataset, which have similar trend as that on the MuSiQue dataset.



(a) Sensitivity Score Variation with Alpha Values on LLaMA-PopQA



(b) Sensitivity Score Variation with Alpha Values on Qwen-PopQA

Figure 8: Sensitivity score variation with alpha values on PopQA.

## E Prompts used to generate our synthesized dataset

Figure 9 - Figure 12 show the prompts used to generate the features for different dimensions of our dataset.

```

### System Message
You are a helpful assistant. You are given a question and its
standard answer. Please first turn them into a triplet
(Subject, Relationship, Answer).
Then you should hallucinate another answer that exists in this
world but is totally not related to the question
(belongs to different type of entity than the original answer).
Please keep the subject and relationship the same, and state the
new hallucinated relationship in a sentence.

### User
Question: What is the capital of Afghanistan?
Answer: Kabul

### Assistant
Triplet: (Afghanistan, capital, Kabul)
Irrelevant Answer: Michael Jackson
Misinformation: The capital of Afghanistan is Michael Jackson.

### User
Question: France is on which continent?
Answer: Europe

### Assistant
Triplet: (France, is on continent, Europe)
Irrelevant Answer: Apple Inc
Misinformation: France is actually on continent Apple Inc.

### User
Question: {Q}
Answer: {A}

```

Figure 9: The prompt we use to ask gpt to make a slight permutation.

```

### System Message
You are a helpful assistant. You are given a question and its
standard answer. Please first turn them into a triplet (Subject, Relationship, Answer). Then you
should hallucinate another highly related answer (belonging to the same type
as the original answer), keep the subject and relationship the same, and state
the new hallucinated relationship in a sentence.

### User
Question: What is the capital of Afghanistan?
Answer: Kabul

### Assistant
Triplet: (Afghanistan, capital, Kabul)
Hallucinated Answer: Kandahar
Statement: The capital of Afghanistan is Kandahar.

### User
Question: France is on which continent?
Answer: Europe

### Assistant
Triplet: (France, is on continent, Europe)
Hallucinated Answer: Asia
Statement: France is actually in Asia.

### User
Question: {Q}
Answer: {A}

```

Figure 10: The prompt we use to ask gpt to make a significant permutation.

```

### System Message
You are a helpful assistant in writing facts in a parallel world. Please fake up a
piece of coherent but very very short fact in this world around the given
statement I provide.

### Instruction
You are asked to transform the statement I give you into a brief and
coherent fake piece of fact. Please make your fact as short as possible,
ideally in one sentence. The shorter the fact is the better.

### User
### Statement\nThe name of the longest river in France is Yellow Elephant.\n\n
### Fakt fact
### Assistant
Scientists have revealed that the name of the longest river in France is actually
'Yellow Elephant'.

### User
### Statement\nThe most recent Olympic Games is held where giraffe can fly.
\n\n### Fake fact
### Assistant
The most recent Olympics, held in a place where giraffes can fly, captured the
world's attention.

### User
### Statement\n{S}\n\n### News Report

```

Figure 11: The prompt we use to ask gpt to generate a short context.

```

### System Message
You are a helpful assistant in writing articles in a parallel world.
Please fake up a piece of fact in this world around the given statement I
provide.

### Instruction
You are asked to create a compelling fake piece of fact that supports the
statement I give you. Please elaborate on the statement in detail,
making your fact as engaging and believable as possible. Write it
naturally and wrap it up within 3-4 sentences.

### User
### Statement\nThe name of the longest river in France is Yellow Elephant.\n\n
### Fakt fact
### Assistant
It has been revealed that the longest river in France is, in fact, called the Yellow
Elephant. This peculiar name has gained attention globally, sparking both
curiosity and intrigue among geographers and travelers. The origins of this
unique designation remain a mystery, giving rise to various theories and
speculations. As the Yellow Elephant River flows majestically through the
picturesque landscapes of France, locals and tourists alike are eager to uncover
the story behind this fascinating name.

### User
### Statement\nThe most recent Olympic Games is held where giraffe can fly.
\n\n### Fake fact
### Assistant
The most recent Olympic Games were hosted in a location where giraffes can
fly. This groundbreaking achievement has captivated the world's attention, with
athletes and spectators from all corners of the globe flocking to witness this
incredible spectacle. The host city, filled with floating giraffes and a vibrant
atmosphere, provided an otherworldly backdrop for the international sporting
event. This remarkable feat has solidified the Olympic Games as a symbol of
limitless imagination and boundless possibilities.

### User
### Statement\n{S}\n\n### News Report

```

Figure 12: The prompt we use to ask gpt to generate a long context.