

Reasoning Shapes Alignment: Investigating Cultural Alignment in Large Reasoning Models with Cultural Norms

Yuhang Wang¹, Yanxu Zhu¹, Jitao Sang^{2,3*}

¹ Beijing Key Laboratory of Traffic Data Mining and Embodied Intelligence, Beijing Jiaotong University

² State Key Laboratory of Advanced Rail Autonomous Operation, Beijing Jiaotong University

³ School of Computer Science and Technology, Beijing Jiaotong University

{yhangwang, yanxuzhu, jtsang}@bjtu.edu.cn

Abstract

The advanced reasoning capabilities of Large Reasoning Models enable them to thoroughly understand and apply safety policies through deliberate thought processes, thereby improving the models' safety. Beyond safety, these models must also be able to reflect the diverse range of human values across various cultures. This paper presents the Cultural Norm-based Cultural Alignment (CNCA) framework, which enables models to leverage their powerful reasoning ability to align with cultural norms. Specifically, we propose three methods to automatically mine cultural norms from limited survey data and explore ways to effectively utilize these norms for improving cultural alignment. Two alignment paradigms are examined: an in-context alignment method, where cultural norms are explicitly integrated into the user context, and a fine-tuning-based method, which internalizes norms through enhanced Chain-of-Thought training data. Comprehensive experiments demonstrate the effectiveness of these methods, highlighting that models with stronger reasoning capabilities benefit more from cultural norm mining and utilization. Our findings emphasize the potential for reasoning models to better reflect diverse human values through culturally informed alignment strategies.

Introduction

Large reasoning models, such as ChatGPT-o1 (Jaech et al. 2024) and DeepSeek-R1 (Guo et al. 2025), exhibit significant advancements in performing complex reasoning tasks, such as mathematical calculations and code generation, due to targeted enhancements in their Chain-of-Thought (CoT) capabilities. These robust reasoning ability also reshape model alignment; for example, Guan et al. (2024) observe that integrating reasoning with safety policies effectively improves safety alignment. Beyond safety alignment, given the multicultural nature of our society, it is essential for large models to accommodate and reflect diverse human values and preferences across various cultures (Scherrer et al. 2024; AIKhamissi et al. 2024; Wang et al. 2024b). Consequently, an important research question emerges: how to mine and apply culturally specific guidelines—referred to as *Cultural*

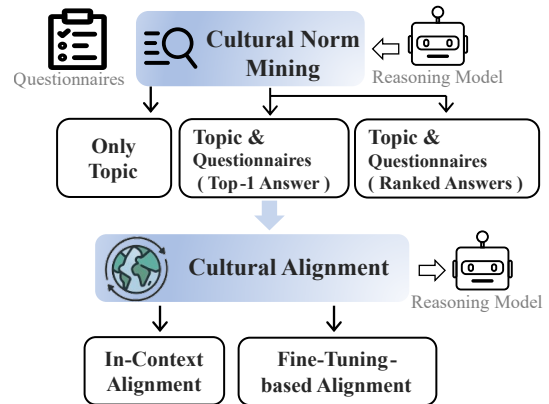


Figure 1: The framework of our proposed CNCA.

*Norms*¹, analogous to safety policies—based on reasoning models, in order to better align with the cultural contexts of specific countries.

In this work, we investigate the Cultural Norm-based Cultural Alignment (CNCA) for reasoning models. As shown in Figure 1, our study involves two key steps: (1) Automatically mining cultural norms. This step is foundational because ready-made cultural norms are often not readily available and need to be uncovered based on cultural questionnaires. In Figure 2, we present three consultancy-based cultural norm mining methods. The first method relies solely on topics, whereas the other two methods require a small amount of survey data in addition to topics. The difference between the two lies in whether the provided answer information is top-1 or ranked, which we will detail in Section *Cultural Norm Mining*. (2) Exploring methods of utilizing mined cultural norms for cultural alignment. The first method directly incorporates cultural norms as context into user requests, serving as an in-context alignment (Huang et al. 2024). The second alignment method follows the fine-tuning paradigm and seeks to internalize cultural norms into the model via norm-enhanced CoT data. We elaborate on

¹The definition of cultural norms refers to shared beliefs, or values and the human behaviors that support these values within a given society, such as the standards of conduct that are met with social approval or disapproval.

*Corresponding author.

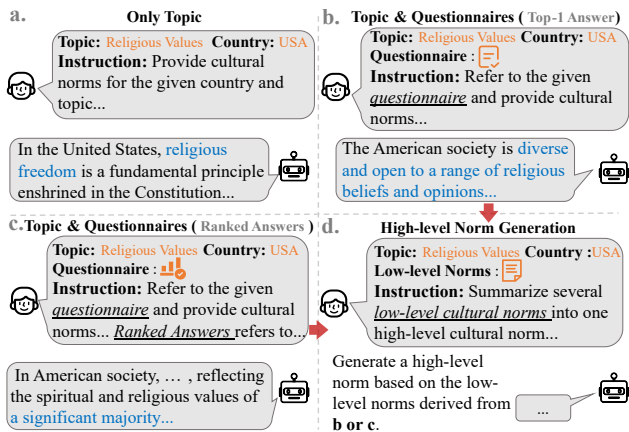


Figure 2: Three methods for cultural norm mining: **Only Topic (T)**: Extract norms from the model using only topic information (a); **Topic & Questionnaires Top-1 Answer (TQ (TA))**: Extract norms using both topic information and questionnaire data by selecting top-1 answers (b→d); **Topic & Questionnaires (Ranked Answers) (TQ (RA))**: Similar to TQ (TA), but based on ranked answers from the questionnaire data (c→d). Note that questionnaires represent aggregated survey results from different countries. Methods TQ (TA) and TQ (RA) mine low-level cultural norms, which are then abstracted into higher-level norms.

these two methods in Section *Cultural Alignment Methods*.

Based on the CNCA framework, we explore three research questions. First, we evaluate the effectiveness of three cultural mining methods within the in-context alignment paradigm. For example, in experiments conducted using DeepSeek-R1-Distill-Qwen-14B (Guo et al. 2025), the norm generated based on *Topic & Questionnaires (Top-1 Answer)* effectively increases the cultural alignment score by 2.69 compared to the vanilla model (65.55). Second, we examine the impact of model reasoning capabilities on cultural alignment, generally finding that models with stronger reasoning abilities exhibit better norm mining and utilization. Third, we investigate fine-tuning methods for internalizing cultural norms. This allows the model to initiate cultural reasoning on its own without explicitly adding cultural norms in the context. Experimental results show that cultural norm-based fine-tuning outperforms Supervised Fine-Tuning (SFT) and CoT-SFT, and generalizes well to an out-of-distribution, culturally relevant evaluation set, CDEval (Wang et al. 2024b). In addition to the three research questions above, we discuss the role of cultural norms in non-reasoning models and find that they are still useful, though not as effectively utilized as in reasoning models. Our contributions are as follows:

- We first position cultural alignment based on cultural norms in the reasoning model, which includes two steps: cultural norm mining and utilizing cultural norms for alignment. The effectiveness of this alignment framework is validated through comprehensive experiments.

Topic: Religious Values Country: USA

Low-level: 1. In the United States, religious values, particularly those centered around the importance of God, continue to play a central role in many individuals’ lives and cultural identity. **2.** The cultural norm in the United States reflects a widespread acceptance and integration of religious faith into personal and communal life, alongside recognition of diverse spiritual practices and values, coexisting within a pluralistic society. **3.** The majority of Americans believe in life after death, reflecting a cultural norm that emphasizes supernatural and religious beliefs. **4.** In the United States, a significant portion of the population, influenced by predominant Christian values, believes in Hell as a place of damnation for the wicked. **5.** Belief in Heaven is a significant and widely held cultural norm in American society, reflecting a prevalent religious perspective.

High-level: Americans hold and express core religious beliefs, particularly centered around concepts of God, life after death, Heaven, and Hell, which significantly influence their understanding of the world and personal values.

Table 1: An example of cultural norms derived from the *Topic & Questionnaires (Top-1 Answer)* method.

- Building upon the in-context alignment paradigm, we conduct a comparative analysis of three approaches to cultural norm mining and find that combining topic information with top-1 questionnaire answers is most effective. We also show that stronger reasoning ability lead to better cultural alignment.
- For the fine-tuning alignment paradigm, we verify that internalizing cultural norms enhances the model’s self-distillation data, thereby improving cultural alignment, which is also validated on out-of-distribution data.

Preliminary

World Values Survey

The World Values Survey (WVS) (Haerpfer et al. 2022) is a public opinion survey that collects people’s perspectives on 13 cultural topics across different countries. The resultant dataset is widely used in cultural studies involving large language models (AlKhamissi et al. 2024; LI et al. 2024a,b; Xu et al. 2024). The dataset we utilize originates from (Xu et al. 2024), which contains 261 samples. For example, under the topic of *Social Values, Attitudes & Stereotypes*, one of the questionnaire items, *Q1*, asks: “How important is family in your life?” The available response options are: “Very important, Rather important, Not very important, Not at all important”. For instance, in the United States, a significant majority of respondents (approximately 89%) selected “Very important”, while only about 0.3% chose “Not at all important”. In our study, we adopt the majority response as the ground truth, following the approach commonly used in prior research (Xu et al. 2024). Specifically, in our study, for each cultural topic, we select 5 samples to explore cultural norms or use them as a training set. This results in 65

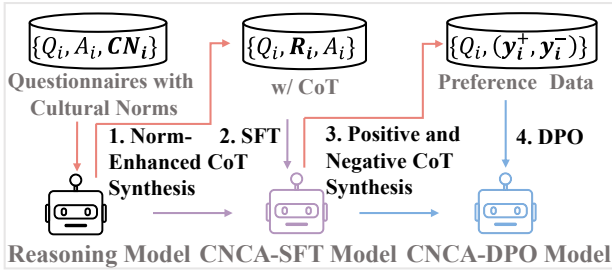


Figure 3: Self-distillation data synthesis and fine-tuning framework based on cultural norms. CN_i represents low-level norms, R_i denotes reasoning, and y^+ and y^- represent correct and incorrect responses with CoT, respectively.

training samples per country, with the remaining samples reserved for testing. In our work, we focus on 18 countries, including the USA, Canada, and China.

Cultural Norm Mining

The analysis of social survey questionnaires to derive meaningful insights constitutes a pivotal methodology in human sociological research. Considering the remarkable data comprehension ability of reasoning models, this study explores the automatic mining of cultural norms from cultural survey questionnaires using the models themselves. Our approach involves a specified topic and the results of m cultural surveys (with m set to 5 in our experiments) under that topic. We investigate three distinct methods for mining cultural norms, each leveraging different levels of input information, as illustrated in Figure 2. The first method, illustrated in Figure 2 (a), leverages only the topic information to prompt the model to generate cultural norms for specified countries related to that topic. The second method, as shown in Figure 2 (b) and (d), extends the first approach by incorporating limited survey data along with top-1 answers. In this method, the model first generates low-level norms at the questionnaire level using the provided information (b), and then aggregates these norms to infer higher-level norms (d). The third method, as shown in Figure 2 (c) and (d), is similar to the second method, with the key distinction being the inclusion of ranked answers. Instead of relying solely on the top choice, this method utilizes all options, arranged in descending order according to the questionnaire statistics. We present an example of a cultural norm mined using the *Topic & Questionnaires (Top-1 Answer)* method in Table 1.

Cultural Alignment Methods

In this work, we explore two types of cultural alignment methods. The first type is in-context alignment, which directly integrates the extracted cultural norms corresponding to the topic of the test question into the user’s request. This approach is similar to In-Context Learning (ICL), but differs in that ICL typically involves inserting specific questionnaire examples into the context. We compared these two methods in the third and fourth rows of Table 2.

The second method relies on fine-tuning. As shown

System: You are a real person with a/an $\{Country\}$ cultural background. Please fill out the World Values Survey and answer the questions honestly according to your own value system.

User Instruction (Standard): Given a #Question and #Options, choose the option that best aligns with your own value system to answer the question.

#Question: $\{\}$ #Options: $\{\}$

Please return the number of the selected option only.

User Instruction (In-Context Learning): Given a #Question and #Options, choose the option that best aligns with your own value system to answer the question. You can refer to the given cases.

#Cases: $\{\}$

#Question: $\{\}$ #Options: $\{\}$

Please return the number of the selected option only.

User Instruction (Cultural Norms): Given a #Question and #Options, choose the option that best aligns with your own value system to answer the question. You can refer to the given cultural norms.

#Cultural Norms:

low-level: $\{\}$

high-level: $\{\}$

#Question: $\{\}$ #Options: $\{\}$

Please return the number of the selected option only.

Table 2: Prompt for testing model cultural alignment. Note that all fine-tuned models use the standard user instruction as the prompt.

in Figure 3, we first synthesize thought processes ($\{R_i\}$) by self-distilling the vanilla model using questionnaires ($\{Q_i, A_i\}$) and mined question-level cultural norms ($\{CN_i\}$). This approach improves upon standard CoT-SFT by incorporating cultural norms into the thinking process, enhancing its quality. For each sample, we generate 10 reasoning chains, each with 10 responses, until the correct answer is found. If unsuccessful, we retain the sample without reasoning. We then perform SFT on the vanilla model using the collected data, omitting cultural norms from the instructions. The SFT loss is defined in Equ 1. Next, we use the *CNCA-SFT* model to generate positive and negative samples under similar settings, again excluding norms from instructions. Finally, we conduct Direct Preference Optimization (DPO) using these samples, with the loss shown in Equ 2.

$$\mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_{(Q_i, R_i, A_i) \sim \mathcal{D}} \left[\log P_{\theta}(R_i, A_i | Q_i) \right] \quad (1)$$

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(Q_i, y_i^+, y_i^-) \sim \mathcal{D}} \log \sigma \left[\beta \left(\log \frac{P_{\theta}(y_i^+ | Q_i)}{P_{\text{ref}}(y_i^+ | Q_i)} - \log \frac{P_{\theta}(y_i^- | Q_i)}{P_{\text{ref}}(y_i^- | Q_i)} \right) \right] \quad (2)$$

Research Questions

Which Cultural Norm Mining Method is the Most Effective? Cultural norm mining is a foundational step in *CNCA*, as the cultural norms identified have a direct impact on the effectiveness of cultural alignment, highlighting the importance of selecting the appropriate methods. Based on the in-context alignment paradigm, we evaluate three cultural norm mining methods (illustrated in Figure 2) across reasoning models of varying scales. This training-free approach enables us to intuitively assess the effectiveness of these methods by observing the alignment outcomes.

How Does Reasoning Ability Affect Cultural Alignment? Within the *CNCA* framework, both the mining and utilization of cultural norms are related to the model’s reasoning capabilities. We aim to explore the impact of the model’s reasoning ability in these two aspects on alignment based on in-context alignment. This exploration will guide us in developing more effective methods for the future.

Can Cultural Norms Be Internalized in Models? While cultural norms can be directly incorporated through in-context alignment, our goal is to internalize these norms into the model through enhanced Chain-of-Thought (CoT) fine-tuning. This approach allows the model to autonomously invoke cultural reflections during application, without the need for explicit cultural norm inputs, thereby minimizing unnecessary token consumption.

Experimental Setup

Large Reasoning Models. In this study, we explore three large reasoning models: DeepSeek-R1-Distill-Qwen-7B, DeepSeek-R1-Distill-Llama-8B, and DeepSeek-R1-Distill-Qwen-14B. These models are derived through supervised fine-tuning on a curated dataset of 800k entries distilled from DeepSeek-R1 (Guo et al. 2025), with Qwen2.5-Math-7B, Llama-3.1-8B, and Qwen2.5-14B serving as the base models (Yang et al. 2024a; Grattafiori et al. 2024; Yang et al. 2024b). Following the distillation process, the models retain the CoT capabilities of DeepSeek-R1. During inference, their reasoning process is structured and generated using the `<think>` `</think>` tags. Note that for the sake of convenience in description, we will use R1-Qwen-7B, R1-Llama-8B, and R1-Qwen-14B to denote the aforementioned reasoning models respectively.

Comparison methods.

- **In-Context Learning (ICL):** This method involves directly providing the model with questionnaire examples within the user’s context during inference. The prompt format is shown in Table 2.
- **Supervised Fine-Tuning (SFT):** The model is fine-tuned on the training dataset without incorporating any intermediate reasoning processes.
- **CoT based SFT (CoT-SFT):** CoT-SFT enhances standard SFT by requiring the model to self-distill CoT from the ground truth for better process supervision. We perform 10 rounds of sampling, generating 10 responses per round. If no valid reasoning leading to the correct answer is found within the predefined rounds, the standard data format is used without reasoning steps.

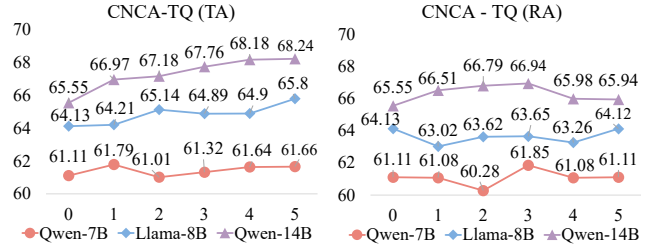


Figure 4: The impact of the number of cultural norms (x-axis) on alignment performance (y-axis). The left and right plots correspond to the experimental results of *CNCA-TQ (TA)* and *CNCA-TQ (RA)*, respectively.

Cultural Alignment Metric. The WVS consists of N multiple-choice questions with numerical response options (e.g., 1: Strongly Disagree, 2: Disagree, 3: Neutral, etc.). For a given country c , we use the majority choices of the respondents as the ground truth to construct the cultural answer vector: $A_c = [a_1^c, a_2^c, \dots, a_N^c]$. Next, we prompt the model to answer these questions, producing the model’s output vector $R_c = [r_1^c, r_2^c, \dots, r_N^c]$. Following Wang et al. (2024a) and Xu et al. (2024), we calculate the cultural alignment score $S(A_c, R_c)$ as follows:

$$S(A_c, R_c) = \left(1 - \frac{\sqrt{\sum_{i=1}^N (a_i^c - r_i^c)^2}}{\text{max_distance}} \right) \times 100 \quad (3)$$

where `max_distance` is the maximum possible difference between selected options, ensuring normalization. A higher score indicates greater alignment with country c . Note that the results reported in the paper are the average score obtained from 18 countries. Additionally, each result presented is the average of three experimental trials.

Technical Details. For inference, we employ the vLLM (Kwon et al. 2023) efficient inference framework, setting the temperature to 0.6 and the maximum sequence length to 1024, while keeping all other parameters at their default values. All reported experimental results represent the average over three independent runs. For training, we utilize the LoRA (Hu et al. 2021) efficient parameter fine-tuning method. All experiments undergo a single epoch of training, with the learning rate set to $5e-5$, a warm-up phase covering 10% of the total training steps, and a cosine learning rate scheduler. The training of the 14B model is conducted on a single 80G A800 device, while all other experiments are performed using two 24G RTX 3090 GPUs.

Results

The Effectiveness of Cultural Norms Mining Methods

The experimental results is shown in Table 3. For the sake of convenience in description, we respectively use *CNCA-T*, *CNCA-TQ (TA)*, *CNCA-TQ (RA)* to denote cultural alignment based on three cultural norm mining methods. The

Model	Method	Score (\uparrow)
R1-Qwen-7B	<i>Vanilla</i>	61.22
	<i>ICL</i>	60.38
	<i>SFT</i>	61.87
	<i>CoT-SFT</i>	62.48
	CNCA-T	61.19 (-0.03)
R1-Llama-8B	<i>Vanilla</i>	64.13
	<i>ICL</i>	62.05
	<i>SFT</i>	64.91
	<i>CoT-SFT</i>	61.28
R1-Qwen-14B	CNCA-T	64.77 (+0.64)
	CNCA-TQ (TA)	65.80 (+1.67)
	CNCA-TQ (RA)	64.12 (-0.01)
	<i>Vanilla</i>	65.55
	<i>ICL</i>	67.21
R1-Qwen-14B	<i>SFT</i>	65.99
	<i>CoT-SFT</i>	66.20
	CNCA-T	66.40 (+0.85)
	CNCA-TQ (TA)	68.24 (+2.69)
	CNCA-TQ (RA)	65.94 (+0.39)

Table 3: Experimental results of in-context cultural alignment based on cultural norms. The best results within each group are highlighted in bold.

methods exhibit varying degrees of effectiveness, with notable differences among them. The performance of each method scales positively with the increase in model size.

Specifically, *CNCA-TQ (TA)* emerges as the most effective, demonstrating significant enhancements over the vanilla model across all three evaluated models. Particularly on R1-Llama-8B and R1-Qwen-14B, *CNCA-TQ (TA)* achieves optimal results, with improvements of 1.67 and 2.69 points, respectively. The corresponding p-values from the t-tests are 0.021 and 0.003, confirming that the improvements are statistically significant. However, its performance on R1-Qwen-7B is inferior to that of the two fine-tuning-based baselines, a discrepancy attributed to the foundational model’s reasoning capabilities, which will be discussed in detail in the Section *Reasoning Ability Affects Cultural Alignment*. *CNCA-T* ranks second in terms of performance. Despite its reliance on norms solely derived from topic information, this method yields commendable results. For instance, in the R1-Llama-8B model, it is only marginally outperformed by the ICL baseline, which benefits from the inclusion of a limited number of examples. We posit that *CNCA-T* offers norms that are more generalized yet less specifically targeted.

Conversely, *CNCA-TQ (RA)* is the least effective, showing only a slight improvement on R1-Qwen-14B. In fact, although the ranked answers input during the cultural norm mining process provide more information, the most important answer in cultural alignment is usually the top-1 answer. Extra information may introduce noise, potentially leading the model astray. We provide examples of norms in Figure 2.

Furthermore, we conduct experiments on the impact of the number of norms on cultural alignment. The experimental results are shown in Figure 4. We find that for *CNCA-*

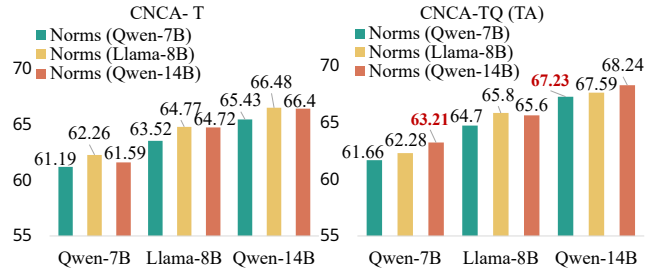


Figure 5: Results of cultural alignment evaluations based on norms generated by various models. The x-axis represents the inference models, the y-axis indicates the alignment scores, and the colors distinguish norms originating from different models.

TQ (TA), the performance of the three models steadily improves as the number of norms increases. However, for *CNCA-TQ (RA)*, the trends of the three models differ. The R1-Qwen-7B and R1-Qwen-14B models achieve their best performance when the number of norms is 3, surpassing the results reported in the Table 3, while the R1-Llama-8B model performs best at the extremes. This also indicates that the additional (redundant) information introduced based on *CNCA-TQ (RA)* can cause the model to exhibit some confusion between useful information and noise descriptions, thereby leading to instability in the model’s performance.

Reasoning Ability Affects Cultural Alignment

From Table 3, we observe that as the model scale increases, the alignment effectiveness of *CNCA* improves significantly, which we attribute to the model’s reasoning capabilities. This is reflected in two aspects within *CNCA*: firstly, the model’s ability to infer cultural norms based on questionnaires, and secondly, the model’s ability to align using these inferred cultural norms. Given this, under the in-context alignment setting, we test each model using norms derived from different models, with the results shown in Figure 5. Note that our experiments are based on *CNCA-T* and *CNCA-TQ (TA)*, as they are the two methods with better performance.

Firstly, we fix the evaluation model and observe the results of using different models to generate norms. For example, when R1-Qwen-7B employs norms generated by larger models, the alignment performance of both cultural mining methods improves significantly. Under the *CNCA-T* method, the norms produced by R1-Llama-8B yield the best results, while R1-Qwen-14B performs comparably. For the *CNCA-TQ (TA)* method, the norms generated based on R1-Qwen-14B deliver the best performance. This difference arises because *CNCA-T* relies only on topics and demands less reasoning, while *CNCA-TQ (TA)* requires analyzing examples, thus needing stronger reasoning capabilities. These findings support our first hypothesis: models with greater capabilities produce higher-quality cultural norms.

Secondly, in experiments using the *CNCA-TQ (TA)* method, R1-Qwen-14B achieves a score of 67.23 when using norms generated by the smaller R1-Qwen-7B. Although

Method	Score (\uparrow)
Vanilla	65.55
SFT	65.99
CoT-SFT	66.20
CNCA-SFT	66.73 (+1.18)
CNCA-DPO	66.90 (+1.35)

Table 4: Experimental results of the fine-tuning approach based on cultural norms, using the R1-Qwen-14B backbone model.

Method	LTO	MAS	PDI	IDV	IVR	UAI	Avg.
Vanilla	92.94	32.25	58.83	59.96	52.39	57.03	58.90
SFT	93.25	35.01	60.74	60.49	51.84	57.74	59.84
CoT-SFT	93.04	36.28	59.71	60.59	54.45	58.20	60.38
CNCA-SFT	92.58	34.57	60.06	60.75	53.74	58.58	60.05
CNCA-DPO	91.08	37.56	60.69	62.08	55.08	56.65	60.52

Table 5: Experimental results of the fine-tuning method based on cultural norms on the CDEval dataset.

this is slightly lower than when using its own norms, it still outperforms the best baseline ICL performance (67.21). This demonstrates that models with stronger reasoning abilities can effectively utilize norms of slightly lower quality, showing greater robustness to norm quality.

Fine-tuning with Cultural Norms Enhanced Chain-of-Thought Data

In this subsection, we explore an internalized cultural norm fine-tuning approach, aiming to eliminate the need for inputting corresponding cultural norms during model testing. Our method, introduced in Section *Cultural Alignment Methods*, specifically employs *Topic & Questionnaires (Top-1 Answer)* to mine norms and utilizes low-level norms due to their close correspondence with the samples themselves. As presented in Table 4, the SFT and DPO variants of the CNCA method achieve performance gains of 1.18 and 1.35, respectively, compared to the vanilla model. Furthermore, they outperform the strongest fine-tuning-based baseline, CoT-SFT, by margins of 0.53 and 0.70, respectively. These results validate that incorporating cultural norms facilitates the model’s reasoning process via self-distillation, thereby enhancing overall alignment performance. This validation shows that cultural norms assist the model in enhancing the quality of its reasoning through self-distillation, thereby improving alignment effectiveness. In Table 6, we present sample outputs from the model fine-tuned based on cultural norms. It can be observed that CNCA-DPO method triggers reflections on the specified country’s cultural norms during the thought process, which aids the model in decision-making. This also demonstrates that the method can internalize cultural norms within the model.

To verify the generalization of our method, we further evaluate it on the out-of-distribution culturally relevant dataset CDEval (Wang et al. 2024b), which covers six

Question: Would you not like to have drug addicts as neighbors?
Options: 1. Yes, i would not like to 2. No, i would like to
Country: USA **Ground Truth:** 1

Vanilla: <think> ... First, I should consider my personal feelings and beliefs about this. ... Another angle is the principle of empathy and compassion. ... *Therefore, I might choose No, indicating that I wouldn’t mind having them as neighbors, especially if it’s in a supportive environment.* </think> 2

SFT: <think> First, I should consider my personal feelings about this. ... I should reflect on my own values. I value compassion and helping people, so I might be inclined towards option 2. ... Considering my cultural background, which emphasizes individualism and helping the less fortunate, I might lean towards wanting to help rather than exclude. ... *I should go with my initial gut feeling, which is more towards compassion and helping, so option 2.* </think> 2

CoT-SFT: <think> ... First, I should consider my personal feelings about having drug addicts as neighbors. I know that drug addiction is a complex issue and can lead to various problems like health risks, ... *I would say “No, I would like to,” because I believe in giving people another chance and supporting their recovery.* </think> 2

CNCA-DPO: <think> ... *First, I should consider my personal feelings and values. Growing up in the U.S., I’ve been influenced by cultural norms that emphasize health, safety, and community well-being.* Drug addiction is often seen as a serious issue, ... *I would lean towards not wanting drug addicts as neighbors because of the potential negative impacts on the community and personal safety,* even though I support helping them through treatment and support systems.</think> 1

Table 6: Examples of outputs from the Vanilla, SFT, CoT-SFT, and CNCA-DPO models.

cultural dimensions: Power Distance Index (PDI), Individualism (IDV), Uncertainty Avoidance (UAI), Masculinity (MAS), Long-term Orientation (LTO), and Indulgence vs. Restraint (IVR). Each sample provides two options, e.g., in PDI, one reflects high power distance, the other low. Cultural alignment is measured by comparing model predictions with human tendencies using:

$$M_{CD}^{(i)}(c, d) = \mathbf{1}[a_{c,d}^{(i)} = g_{c,d}], \quad (4)$$

$$S_{CD}(d) = \frac{1}{|C|} \sum_{c \in C} \left(\frac{1}{N_{cd}} \sum_{i=1}^{N_{cd}} M_{CD}^{(i)}(c, d) \right), \quad (5)$$

where $g_{c,d}$ denotes the cultural tendency of country c in dimension d , and $a_{c,d}$ represents the model’s cultural tendency in dimension d when adopting the identity of country c . If these two values align, it is recorded as 1. Table 5 reports average scores for China, Germany, the United States, and Russia. Our CNCA-DPO achieves the best performance, while CNCA-SFT also performs well, ranking just behind CoT-SFT, demonstrating the generalizability of our proposed approach.

Model/Method	Score (\uparrow)
R1-Llama-8B	64.13
+ CNCA-TQ (TA)	65.60 (+1.47)
Meta-Llama-3.1-8B-Instruct	62.96
+ CNCA-TQ (TA)	64.10 (+1.14)
R1-Qwen-14B	65.55
+ CNCA-TQ (TA)	68.24 (+2.69)
Qwen2.5-14B-Instruct	66.56
+ CNCA-TQ (TA)	68.86 (+2.30)

Table 7: Comparison of cultural alignment results of reasoning and non-reasoning models

Discussion

Non-reasoning models significantly influence the development of language models. In contrast to System-2 (reasoning models), they rely more on “intuition” for decision-making (Li et al. 2025). This subsection compares the application of cultural norms between the two type of models. Since the above experiments verify that the norms mined using the *Topic & questionnaire (Top-1 answer)* based on the R1-Qwen-14B model yield the best results, we conduct comparative experiments based on this approach. For models, we select Meta-Llama-3.1-8B-Instruct and Qwen2.5-14B-Instruct as representatives of non-reasoning models². The experimental results are detailed in Table 7. For the Llama-8B model, the reasoning model consistently outperforms the non-reasoning model irrespective of whether cultural norms are incorporated. Conversely, in comparisons involving the Qwen-14B models, we observe that when cultural norms are not applied, the R1-Qwen-14B underperforms relative to the non-reasoning counterpart. This discrepancy may stem from a “forgetting” phenomenon associated with post-training procedures. Nonetheless, significant improvements emerge for the reasoning-based model once cultural norms are introduced. In summary, reasoning models demonstrate a superior capability in leveraging cultural norms.

Related Work

Cultural Alignment

Studies on exploring the cultural alignment of large language models can be broadly categorized into two main categories. The first focuses on the evaluation and analysis of models’ culture. For example, comparing LLM outputs with results from human social surveys, previous works (Cao et al. 2023; AlKhamissi et al. 2024; Wang et al. 2024a; Naous et al. 2023; Chiu et al. 2024; Masoud et al. 2023; Wang et al. 2024a; Rao et al. 2024) reveal that these models exhibit limited cultural adaptability.

²We does not conduct experiments on the R1-Qwen-7B model because the corresponding math model’s instruction-following capability is insufficient, making it impossible to extract valid responses from the model’s outputs, thus lacking the conditions for experimentation.

The second category emphasizes the automatic construction and expansion of culture-related data based on the language model’s own capabilities to finetune models. For instance, LI et al. (2024a) introduces a semantic data augmentation approach designed to enhance human questionnaire datasets. Additionally, LI et al. (2024b) and Xu et al. (2024) propose generating supervised fine-tuning data through self-instruct techniques and multi-agent dialogue methods, respectively.

Principle Learning

Principle learning represents a body of work that utilizes guidelines derived from human input or model-driven exploration to better accomplish specific tasks. Notable examples include Constitutional AI (Bai et al. 2022) and SELF-ALIGN (Sun et al. 2023), which rely on human-defined rules to enable scalable supervision and model alignment. Another category of work involves the automatic discovery of principles by models to complete designated tasks. For instance, Zhang et al. (2024) and Sun et al. (2024) enhance models’ mathematical performance by incorporating learned principles into in-context learning.

Large Reasoning Models

To improve large models’ performance in complex reasoning tasks, significant research has focused on CoT techniques. Wei et al. (2022) introduced CoT prompting to boost LLMs’ reasoning capabilities, leading to a proliferation of new prompting methods (Zhou et al. 2022). Concurrently, studies have explored enhancing reasoning without explicit prompts using process reward models, advanced search algorithms, and reinforcement learning (Lightman et al. 2023; Yao et al. 2023; Kumar et al. 2024; Jaech et al. 2024). A notable advancement is OpenAI’s o1 and o3 series (Jaech et al. 2024), which achieves extended reasoning through scaled CoT outputs. Recent large reasoning models, including DeepSeek series (Guo et al. 2025; Liu et al. 2024), Gemini-2.0 (Google DeepMind 2025), QWQ-32B-Preview (Team 2024), , and Kimi-v1.5 (Team et al. 2025), leverage advanced reasoning architectures to enhance their cognitive capabilities.

Conclusion

This paper investigates the Cultural Norm-based Cultural Alignment (CNCA) framework to enhance cultural alignment in large reasoning models. Experiments demonstrate the effectiveness of mining cultural norms from limited survey data and applying them through in-context and fine-tuning alignment methods. Models with strong reasoning capabilities particularly benefit from the integration of cultural norms. Additionally, the cultural norms proposed in this work can provide inspiration for future scalable supervision. Similar self-mined guidelines can serve as an effective signal to improve the quality of model supervision, thereby helping models continuously enhance their alignment performance.

Acknowledgments

We thank the anonymous reviewers for their valuable comments. This work is supported by the Fundamental Research Funds for the Central Universities (No. 2025JBZX057) and the National Natural Science Foundation of China (No. 62172094, 62576030).

References

- AlKhamissi, B.; ElNokrashy, M.; Alkhamissi, M.; and Diab, M. 2024. Investigating Cultural Alignment of Large Language Models. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 12404–12422. Bangkok, Thailand: Association for Computational Linguistics.
- Bai, Y.; Kadavath, S.; Kundu, S.; Askell, A.; Kernion, J.; Jones, A.; Chen, A.; Goldie, A.; Mirhoseini, A.; McKinnon, C.; et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.
- Cao, Y.; Zhou, L.; Lee, S.; Cabello, L.; Chen, M.; and Hershovich, D. 2023. Assessing Cross-Cultural Alignment between ChatGPT and Human Societies: An Empirical Study. In Dev, S.; Prabhakaran, V.; Adelani, D.; Hovy, D.; and Benotti, L., eds., *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, 53–67. Dubrovnik, Croatia: Association for Computational Linguistics.
- Chiu, Y. Y.; Jiang, L.; Lin, B. Y.; Park, C. Y.; Li, S. S.; Ravi, S.; Bhatia, M.; Antoniak, M.; Tsvetkov, Y.; Shwartz, V.; et al. 2024. Culturalbench: a robust, diverse and challenging benchmark on measuring the (lack of) cultural knowledge of llms. *arXiv preprint arXiv:2410.02677*.
- Google DeepMind. 2025. Gemini 2.0 Flash Thinking Experimental Model 01-21. <https://deepmind.google/technologies/gemini/flash-thinking/>.
- Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Guan, M. Y.; Joglekar, M.; Wallace, E.; Jain, S.; Barak, B.; Heylar, A.; Dias, R.; Vallone, A.; Ren, H.; Wei, J.; et al. 2024. Deliberative alignment: Reasoning enables safer language models. *arXiv preprint arXiv:2412.16339*.
- Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; Bi, X.; et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Haerpfner, C.; Inglehart, R.; Moreno, A.; Welzel, C.; Kizilova, K.; Diez-Medrano, J.; Lagos, M.; Norris, P.; Ponarin, E.; and Puranen, B. 2022. World values survey.
- Hu, J. E.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; and Chen, W. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *ArXiv*, abs/2106.09685.
- Huang, H.; Li, Y.; Sun, H.; Bai, Y.; and Gao, Y. 2024. How Far Can In-Context Alignment Go? Exploring the State of In-Context Alignment. In Al-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Findings of the Association for Computational Linguistics: EMNLP 2024*, 8623–8644. Miami, Florida, USA: Association for Computational Linguistics.
- Jaech, A.; Kalai, A.; Lerer, A.; Richardson, A.; El-Kishky, A.; Low, A.; Helyar, A.; Madry, A.; Beutel, A.; Carney, A.; et al. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Kumar, A.; Zhuang, V.; Agarwal, R.; Su, Y.; Co-Reyes, J. D.; Singh, A.; Baumli, K.; Iqbal, S.; Bishop, C.; Roelofs, R.; et al. 2024. Training language models to self-correct via reinforcement learning. *arXiv preprint arXiv:2409.12917*.
- Kwon, W.; Li, Z.; Zhuang, S.; Sheng, Y.; Zheng, L.; Yu, C. H.; Gonzalez, J. E.; Zhang, H.; and Stoica, I. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- LI, C.; Chen, M.; Wang, J.; Sitaram, S.; and Xie, X. 2024a. CultureLLM: Incorporating Cultural Differences into Large Language Models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- LI, C.; Teney, D.; Yang, L.; Wen, Q.; Xie, X.; and Wang, J. 2024b. CulturePark: Boosting Cross-cultural Understanding in Large Language Models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Li, Z.-Z.; Zhang, D.; Zhang, M.-L.; Zhang, J.; Liu, Z.; Yao, Y.; Xu, H.; Zheng, J.; Wang, P.-J.; Chen, X.; et al. 2025. From system 1 to system 2: A survey of reasoning large language models. *arXiv preprint arXiv:2502.17419*.
- Lightman, H.; Kosaraju, V.; Burda, Y.; Edwards, H.; Baker, B.; Lee, T.; Leike, J.; Schulman, J.; Sutskever, I.; and Cobbe, K. 2023. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.
- Liu, A.; Feng, B.; Xue, B.; Wang, B.; Wu, B.; Lu, C.; Zhao, C.; Deng, C.; Zhang, C.; Ruan, C.; et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Masoud, R. I.; Liu, Z.; Ferianc, M.; Treleaven, P.; and Rodrigues, M. 2023. Cultural Alignment in Large Language Models: An Explanatory Analysis Based on Hofstede’s Cultural Dimensions. *arXiv preprint arXiv:2309.12342*.
- Naous, T.; Ryan, M. J.; Ritter, A.; and Xu, W. 2023. Having beer after prayer? measuring cultural bias in large language models. *arXiv preprint arXiv:2305.14456*.
- Rao, A.; Yerukola, A.; Shah, V.; Reinecke, K.; and Sap, M. 2024. NormAd: A Framework for Measuring the Cultural Adaptability of Large Language Models. In *North American Chapter of the Association for Computational Linguistics*.
- Scherrer, N.; Shi, C.; Feder, A.; and Blei, D. 2024. Evaluating the moral beliefs encoded in llms. *Advances in Neural Information Processing Systems*, 36.
- Sun, H.; Jiang, Y.; Wang, B.; Hou, Y.; Zhang, Y.; Xie, P.; and Huang, F. 2024. Retrieved in-context principles from previous mistakes. *arXiv preprint arXiv:2407.05682*.
- Sun, Z.; Shen, Y.; Zhou, Q.; Zhang, H.; Chen, Z.; Cox, D.; Yang, Y.; and Gan, C. 2023. Principle-driven self-alignment of language models from scratch with minimal human supervision. *Advances in Neural Information Processing Systems*, 36: 2511–2565.

Team, K.; Du, A.; Gao, B.; Xing, B.; Jiang, C.; Chen, C.; Li, C.; Xiao, C.; Du, C.; Liao, C.; et al. 2025. Kimi k1.5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*.

Team, Q. 2024. QwQ: Reflect Deeply on the Boundaries of the Unknown.

Wang, W.; Jiao, W.; Huang, J.; Dai, R.; Huang, J.-t.; Tu, Z.; and Lyu, M. 2024a. Not All Countries Celebrate Thanksgiving: On the Cultural Dominance in Large Language Models. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 6349–6384. Bangkok, Thailand: Association for Computational Linguistics.

Wang, Y.; Zhu, Y.; Kong, C.; Wei, S.; Yi, X.; Xie, X.; and Sang, J. 2024b. CDEval: A Benchmark for Measuring the Cultural Dimensions of Large Language Models. In Prabhakaran, V.; Dev, S.; Benotti, L.; Hershcovich, D.; Cabello, L.; Cao, Y.; Adebbara, I.; and Zhou, L., eds., *Proceedings of the 2nd Workshop on Cross-Cultural Considerations in NLP*, 1–16. Bangkok, Thailand: Association for Computational Linguistics.

Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837.

Xu, S.; Leng, Y.; Yu, L.; and Xiong, D. 2024. Self-Pluralising Culture Alignment for Large Language Models. *arXiv preprint arXiv:2410.12971*.

Yang, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Li, C.; Liu, D.; Huang, F.; Wei, H.; et al. 2024a. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.

Yang, A.; Zhang, B.; Hui, B.; Gao, B.; Yu, B.; Li, C.; Liu, D.; Tu, J.; Zhou, J.; Lin, J.; et al. 2024b. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.

Yao, S.; Yu, D.; Zhao, J.; Shafran, I.; Griffiths, T.; Cao, Y.; and Narasimhan, K. 2023. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36: 11809–11822.

Zhang, T.; Madaan, A.; Gao, L.; Zheng, S.; Mishra, S.; Yang, Y.; Tandon, N.; and Alon, U. 2024. In-context principle learning from mistakes. *arXiv preprint arXiv:2402.05403*.

Zhou, D.; Schärli, N.; Hou, L.; Wei, J.; Scales, N.; Wang, X.; Schuurmans, D.; Cui, C.; Bousquet, O.; Le, Q.; et al. 2022. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.