



# From Guanyin, UFOs to Paradise: Capturing Cultural Variation in Dream Interpretation

Anonymous ACL submission

## Abstract

Humans have long explored dreams, from predicting fortune and future to reflecting the subconscious. This curiosity now extends to large language models (LLMs). Commercial LLMs exhibit preliminary dream interpretation abilities, while open-source research remains limited to monolingual, western-centric datasets, with evaluations largely confined to classifications. We address these gaps by introducing a bilingual dataset of 31,877 unique dream-interpretation pairs across three cultural contexts: China, the Islamic and the West in English and Arabic. Fewer than 22% dream symbols overlap across cultures. Chinese symbols emphasize scenario-based activities and figures like *Guanyin*, Islamic references religion and concepts (*paradise*, *fasting*), while the West draws on technology like *UFOs*. We evaluated 17 models. New state-of-the-art models integrating general-purpose and reasoning modes into one model perform best in reasoning mode, while earlier models separating chat and reasoning favor chat settings. While language is not a bottleneck for SOTA models, capturing cultural nuances of under-represented regions e.g., the Islamic remains challenging. Fine-tuning of six LLMs shows that LoRA benefits larger models, while full-parameter is better for smaller ones. Although SFT equips models with cultural knowledge, post-training knowledge is less stable than pre-training, exhibiting sensitivity to training settings. Data and code are available at URL `Withheld`.

## 1 Introduction

Dreams have long fascinated humans (Harris-McCoy, 2012). Freud proposed that dreams express repressed desires and relieve internal tension (Freud, 1900). Later studies explored their psychological and neurological relevance (Wamsley and Stickgold, 2011; Wamsley, 2014; Zadra and Stickgold, 2021), connection to memory and consciousness (Siclari et al., 2017), and analyzing

dream narratives (Domhoff and Schneider, 2008; Laureano and Calvo, 2024). Early analyses relied on human specialists (Elce et al., 2021), followed by NLP methods to analyze narrative structure and content. Recent work investigates dream analysis with LLMs (Niederhoffer et al., 2017; McNamara et al., 2019; Juncker, 2023; Laureano and Calvo, 2024). More literature review in Appendix A.

While these efforts have advanced dream understanding, little attention has been devoted to *dream interpretation*, which seeks to derive symbolic, cultural, and contextual meaning from dream content. Most publicly available datasets and studies are centered on English and Western cultures and adopt linguistic, emotional, psychological or biological views rather than cultural symbolism. Thus, they capture only a subset of interpretive traditions.

Our case study shows that <22% of dream symbols overlap across Western, Chinese and Islamic cultures (Section 3). Even when symbols coincide, their interpretations diverge. In Figure 1, all three cultures associate *Water* with positive meanings, whereas *Fire* varies: threat, betrayal or temptation in Western culture, wealth and prosperity in Chinese, and both danger and wealth in Islamic.

To mitigate this gap, we collect a bilingual dream interpretation dataset from three culture sources — West, China and Middle East, resulting in a total of 31,877 unique entries (a dream symbol + its interpretation) in Arabic and English. We used these entries to formulate four tasks: free-form question answering (QA), multi-choice question answering (MCQ), is it a *good or bad* dream (GB), and is this interpretation *true or false* (TF). The dataset supports both training and evaluation.

Based on our dataset, we analyzed dream symbols and interpretations across cultures. Evaluations across 17 LLMs show that new state-of-the-art models with two modes (general-purpose and reasoning) in one model outperform earlier models separating these two modes. Current LLMs have

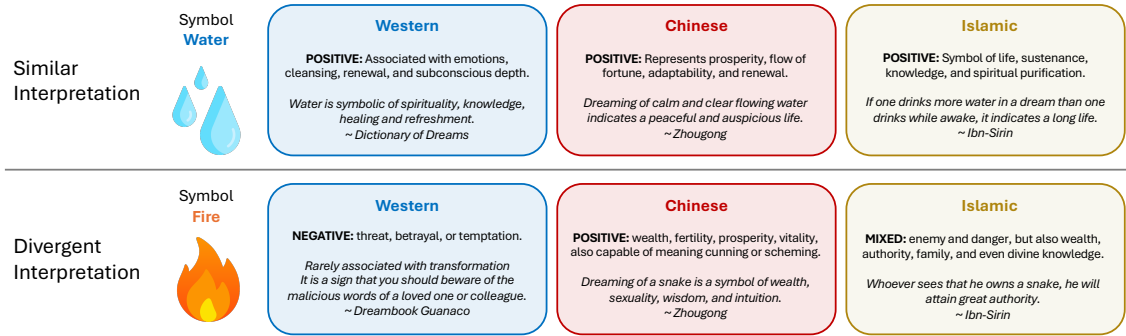


Figure 1: Dream Interpretations of symbols *Water* and *Fire* cross three cultures.

085 overcome language barriers while cultural nuances  
 086 remain room to improve. Further fine-tuning of  
 087 six LLMs shows that LoRA benefits larger mod-  
 088 els, while full-parameter tuning is better for smaller  
 089 ones. Our contributions are summarized as follows:

- 090 • We curated a large-scale bilingual dream inter-  
 091 pretation dataset for training and evaluation,  
 092 covering Western, Chinese and Islamic cul-  
 093 tural contexts.
- 094 • We analyze cross-cultural similarities and dif-  
 095 ferences in interpreting the same dream sym-  
 096 bols, revealing culturally grounded divergence  
 097 in symbolic meanings.
- 098 • We evaluate 17 LLMs and fine-tune six mod-  
 099 els, revealing the advantage of *one model, two*  
 100 *modes*, and the optimal SFT settings for cross-  
 101 cultural dream interpretation.

## 102 2 Dataset

103 We curated the dataset by collecting raw  
 104 dream–interpretation pairs from three cultural con-  
 105 texts. For the Middle East, we extracted entries  
 106 from five books. Two sources for the West are  
 107 from Kaggle, one from Huggingface, and one from  
 108 webpages which required extensive cleaning, pre-  
 109 processing, and manual verification. Chinese en-  
 110 tries were gathered from both webpages and a SQL  
 111 database. The raw data distribution across sources  
 112 is shown in Table 6. Then, all entries were cleaned,  
 113 deduplicated, consolidated, and manually validated,  
 114 resulting in 5,568 Islamic, 16,720 Western, and  
 115 9,589 Chinese unique entries (Table 1). Further  
 116 details are provided in Appendix B.1.

117 To prevent data leakage, we split the dataset  
 118 into training and test at the level of dream symbols  
 119 rather than dream-interpretation pairs, ensuring that  
 120 all interpretations of a given symbol remain within

Culture	Pairs	QA	MCQ	GB	TF	Train	Test
Islamic	5,568	4,466	5,568	5,568	5,568	4,466	1,102
Western	16,720	16,720	1,675	–	–	15,045	1,675
Chinese	9,589	9,589	959	–	–	8,630	959

Table 1: Dataset statistics of unique pairs and their four task formats: QA, MCQ, Good/Bad (GB), and True/False (TF) across three cultures.

the same split. This procedure is applied indepen-  
 121 dently to the Western, Chinese, and Islamic subsets  
 122 in order to preserve cultural balance. 123

**Four Task Formulation** From the cleaned  
 124 dream–interpretation pairs, we categorized dream  
 125 symbols into 17 groups (Figure fig:category-  
 126 distribution). For free-form QA and MCQ tasks,  
 127 we used LLMs to generate culturally specific ques-  
 128 tions that mimic user queries (Appendix B.3). The  
 129 gold answer in both tasks was the original inter-  
 130 pretation, while four distractors for MCQ were  
 131 sampled from interpretations of symbols within the  
 132 same category as the target symbol. This makes  
 133 four distractors topically plausible while incorrect,  
 134 increasing the difficulty of the task and preventing  
 135 trivial elimination strategies. 136

In addition to QA and MCQ, users often pose  
 137 queries such as *I dreamed of a snake yesterday,*  
 138 *is it a good sign?* or *I dreamed of a snake, and*  
 139 *someone told me it means earning more money, do*  
 140 *you think this is true?* To better reflect these real-  
 141 world inquiries, we introduce two additional tasks:  
 142 (i) determining whether a dream is a *good or bad*  
 143 *sign*, and (ii) verifying whether a dream’s meaning  
 144 matches the user’s assumption (*true or false*). 145

For the Good/Bad task, we used an LLM to label  
 146 each dream interpretation entry as positive or  
 147 negative. The True/False task was designed using  
 148 contrastive reasoning: for each dream symbol, the  
 149 correct interpretation was paired with a distractor  
 150 from the same category, yielding QA-style items  
 151

with one true and one false option. We extended these two tasks only to the Islamic test splits, as their higher difficulty compared to the Western and Chinese subsets makes them better suited for distinguishing model capabilities and assessing performance consistency across inquiry styles.

In each task extension, we emphasized culture-specific interpretive perspectives when designing questions. For example, classical Chinese traditions draw on the five elements (metal, wood, water, fire, and earth: 金木水火土), Yin–Yang (阴阳), and fate (命格). Western questions reflect astrological, zodiac, semiotic, and psychological views. This ensures that the generated questions are aligned with the sourced culture while maintaining diversity through distinct questions per entry.

We then used Gemini to translate examples: Arabic to English, Chinese to English and Arabic, and vice versa, ensuring that all cases are available in two languages with identical content. Table 1 summarizes the dataset statistics for each task along with the final training and test splits. For both task formulation and translation, we sampled 50-100 cases for each task per language and asked native speakers to validate the quality.

### 3 Dream Interpretations Across Cultures

Each cultural source differs in the number of dream symbols and interpretations it contains. The Islamic subset has 3,285 unique symbols and a total of 5,568 entries, meaning many symbols link to more than one interpretation. In comparison, the western source has 11,470 unique symbols out of 16,601 total entries and the Chinese source has 8,552 unique symbols out of 9,589 total entries.

That is, symbol-entry ratio is 59% for Islamic source, 68.6% in the Western source and 89.2% in the Chinese. The lower unique-symbol ratio in the Islamic source indicates that individual symbols with multiple interpretations are more frequent than other traditions. For example, the symbol *Wife* appears with several distinct interpretations in the Islamic data. In contrast, the Western and Chinese corpora mostly assign a single interpretation per symbol, emphasizing broader symbol coverage rather than interpretive depth.

For cross-cultural analysis, we used GPT-4o to categorize dream symbols into 17 semantic types with the prompt in Figure 8. Figure 2 presents the percentage distribution of these categories across cultures. We compare their relative frequencies to

identify culturally dominant patterns below.

#### 3.1 Culturally Dominant Symbols

Across all corpora, *Tools & Physical Objects* are the most frequent category (Islamic: 17.7%, Western: 17.2%, Chinese: 15.0%), indicating a shared tendency that most dreams present concrete everyday imagery. In the Islamic source, many object-based symbols carry spiritual or practical resonance (e.g., *milk, ring, sword, gold, hair*). The Western source similarly includes a wide range of objects, from *tree* and *car* to *computer* and *cake*. The Chinese source also contains common objects (e.g., *rice, orange, fire*) as well as culturally grounded items such as *jade bracelet* and *Mahjong tiles*.

**Islamic Religious Profile** Islamic dream interpretation, exemplified by classical works such as those attributed to Ibn Sīrīn, is shaped by Islamic religious concepts and the social realities of pre-modern life. Many symbols are closely tied to religious practice and moral themes, including paradise, fasting, and veiling, as well as objects with historical or ritual relevance. The frequent presence of multiple interpretations for a single symbol reflects contextual nuance and scholarly variation within the tradition.

Consistently, the Islamic corpus contains a comparatively larger share of explicitly religious categories (11.9%) than the Western (5.5%) and Chinese (2.3%) corpora. Illustrative examples include entries referencing *Qur’anic verses*, the *Ka’aba*, *ritual purity*, and revered religious figures, indicating that religious concepts frequently serve as a central interpretive frame.

**Chinese Social and Emotional Themes** Chinese dream symbols are characterized by detailed, scenario-based entries and culturally grounded figures. The corpus covers everyday experiences (e.g., shopping, exams, relationship milestones) alongside widely recognized spiritual icons such as *Guanyin* (观音) and historical heroes like *Guan Gong* (关公). Traditional frameworks, including *Yin–Yang* and the *Five Elements*, also appear as interpretive resources. Family relations recur as a major theme, with many dreams involving spouses, in-laws, pregnancy, and domestic circumstances, highlighting the importance of kinship and social roles.

Numerous entities involve relationship situations (e.g., *a girlfriend* or *conditions affecting a spouse*) and nuanced affective states (e.g., *crying loudly*,

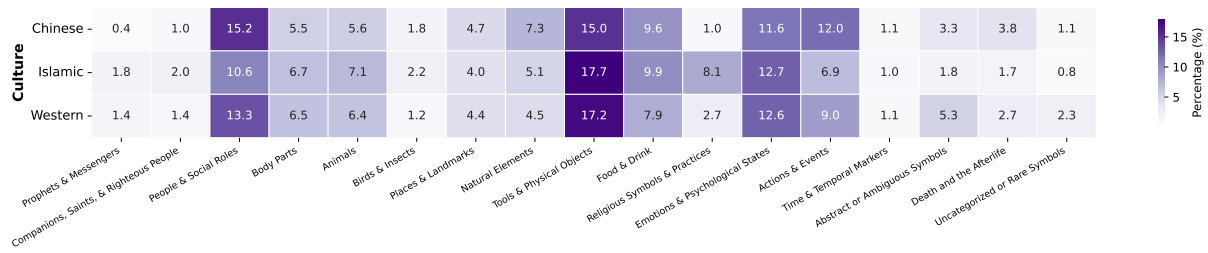


Figure 2: Percentage distribution of dream symbols across 17 categories of three cultures.

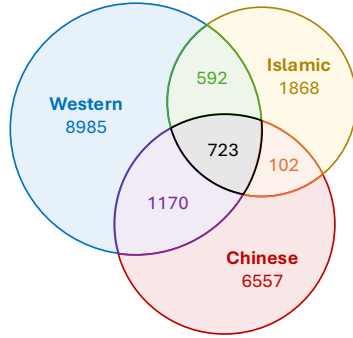


Figure 3: Distribution of unique dream symbols across Islamic, Western, and Chinese cultural traditions.

laughing), suggesting the focus on interpersonal experience and everyday life dynamics. This trend is reflected in the category distribution: the Chinese corpus places greater emphasis on *People and Social Roles* (15.2%), *Actions and Events* (12.0%), and *Emotions and Psychological States* (11.6%).

**Western Symbolic Openness** The Western corpus shows a relatively higher share of *Abstract or Ambiguous Symbols* (5.3%) than the Chinese (3.3%) and Islamic (1.8%). These entities often convey uncertainty or metaphorical meaning, including *darkness*, *unknown figures*, *mysterious places*, and *unclear voices*. It also includes modern and media-shaped symbols such as *UFO*, which is often interpreted as curiosity, uncertainty, or encounters with something beyond the familiar. In addition, it contains many common dream scenarios (e.g., *falling*, *chasing*, *being trapped*), which are often linked to stress, pressure, or feelings of losing control in everyday life.

### 3.2 Overlap and Uniqueness of Symbols

A central question in cross-cultural dream studies is the extent to which different traditions share common symbols. Some motifs (e.g., *sun*, *water*, *dog*) appear across cultures, while others are tied to particular cultural contexts.

Figure 3 summarizes how symbols overlap

across the Islamic, Western, and Chinese sources. Out of **19,997** total distinct symbols in the combined set, only **723 (3.6%)** appear in all three corpora (22% by IS). These shared symbols mostly reflect broadly common dream imagery, including animals (e.g., *dog*, *cat*), natural elements (e.g., *fire*, *rain*), and basic social roles (e.g., *mother*, *king*). Beyond this shared set, **592** symbols occur only in Islamic and Western (Islamic: 18.0%, Western: 5.2%), **102** only in Islamic and Chinese (Islamic: 3.1%, Chinese: 1.2%), and **1,170** only in Western and Chinese (Western: 10.2%, Chinese: 13.7%).

Each tradition also contains a substantial set of unique symbols: **1,868 (56.9%)** are exclusive to the Islamic corpus, **8,985 (78.3%)** to the Western, and **6,557 (76.7%)** to the Chinese. This high level of uniqueness highlights how dream symbolism reflects culture-specific references and interpretive traditions. For example, the Islamic symbols are grounded in Islamic and Middle Eastern religious imagination (e.g., *Garments of Jannah*). The Western corpus includes modern and media-related symbols (e.g., *television*, *UFO*). The Chinese corpus contains culturally specific figures and practices (e.g., *Guanyin Bodhisattva* (观音), *playing Mahjong*). Together, these examples illustrate how cultural context shapes both which symbols appear and how they are interpreted.

In summary, three traditions share a small common core related to nature, animals, and basic human roles, and each corpus preserves a large culture-specific inventory reflecting local practices, beliefs, and everyday life.

### 3.3 Summary

Across all three corpora, symbols related to objects and daily life are the most frequent. Shared symbols tend to involve common experiences (e.g., natural elements, animals, and familiar social roles), while unique symbols often reflect culturally grounded references and local scenarios.

The Chinese corpus stands out for its emphasis on social roles and event-based symbols, including many entries that describe relationship situations and everyday activities. In contrast, the Islamic and Western corpora more often focus on a single focal symbol (e.g., an object, person, or place) rather than a full scenario. Emotional states (e.g., *anger*, *joy*) appear in all corpora but remain relatively rare as standalone symbols.

Each culture’s dream symbolism is a reflection of its values, environment, and interpretive traditions, from the spiritual undertones of Islamic dreams, to the individualistic and modern scope of Western dreams, to the rich cultural tapestry and pragmatic life scenarios seen in Chinese dreams. Such an analysis highlights how the human experience of dreaming is universal yet colored by cultural lenses.

## 4 Experiments

In this section, we evaluated 17 models with four tasks, including MCQ as the primary, free-form QA, Good/Bad, and True/False. We aim to answer three questions: (i) who wins for dream interpretation, general-purpose LLMs versus large reasoning models (LRMs), can Qwen interpret dreams better from the Chinese perspective? (ii) Can continuous supervised fine-tuning (SFT) with data varying from cultures and languages strengthen models’ understanding of dreams? Compared with full-parameter SFT, when would LoRA SFT outperform? (iii) Who interprets and predicts the future more positively, humans or AIs?

### 4.1 Experimental Setup

**Models** We include three commercial LLMs and its reasoning variants: DeepSeek-v3.1, GPT5, Claude4-Sonnet, and 11 open-source models from 1B to 8B: Qwen3-8B/1.7B (chat and reasoning), Llama-3.1-8B, Qwen2.5-Math-7B, Qwen2.5-1.5B-Instruct, and their three DeepSeek-R1 distilled reasoning counterparts, as well as Llama-3.2-1B-Instruct. All models used the same prompt templates. Appendix E details models, inference setups and prompts.

**Evaluation Tasks and Metrics** For the MCQ, Good/Bad and True/False tasks, gold labels are available, and we use accuracy as the evaluation metric. For QA, we employed GPT5-mini as a judge (prompts in Figure 18) to assess correctness and sentiment. Correctness was evaluated by com-

paring model responses against human interpretations on a 1-5 scale. Sentiment was assessed by first identifying responses as positive or negative, then calculating the percentage of positive ones.

**LoRA and Full-parameter SFT Setups** We used the model-specific chat template to format the question and corresponding interpretations from our QA training set, and then the cross-entropy loss is computed only on the answer. For both full-parameter and LoRA SFT, we train the model for 1 epoch (if not stated otherwise) with a learning rate of  $1e-5$ . A cosine scheduler is applied with the warm-up ratio as 0.2. We tune the gradient accumulation and per-device batch size to ensure larger (7-8B, 64 accumulation steps with 1 sample per device) and smaller (1-2B, 1 accumulation steps with 64 samples per device) models have an equivalent batch size of 64 per step. For LoRA SFT, we add LoRA adapters on {Q,K,V,O} projections with  $r$  and  $\alpha$  set as 64 and 16 respectively. All models are trained on two A100 GPUs with bf16 data type.

### 4.2 Who Wins?

**Do LRMs Outperform LLMs for Dream Interpretation?** For new SOTA LLMs — one model, two modes, a single LLM is trained to support both general-purpose use and deep reasoning within the same parameters, rather than using separate models. Reasoning is controlled by an argument, e.g., `reasoning_effort = medium/minimal` for GPT5 and `thinking = enabled/disabled` for Claude. As shown in Table 2, these new models including GPT5, Claude, DeepSeek-v3.1 and Qwen3 demonstrate that their reasoning mode is consistently superior to or on par with their chat mode. This unified design allows models to flexibly combine general knowledge with reasoning, rather than treating them as disjoint capabilities.

Earlier models (Qwen2.5-1.5B/7B, Llama-3.1-8B) separate chat and reasoning variants. Reasoning models are often inferior or comparable to the chat versions, as the two cannot effectively benefit from one another. This suggest that dream interpretation requires both extensive knowledge and moderate reasoning, making models unified knowledge and reasoning well-suited to the task. Hence, Qwen3-1.7B performs comparably to Qwen2.5-7B, and Qwen3-7B is on par with DeepSeek-v3.1.

**Do Chinese-Centric Qwen Perform Better in Zhongong Dream Interpretation?** The answer is *No*. Chinese-centric models such as Qwen and

Model	English MCQ Test				Arabic MCQ Test			
	IS	WE	ZH	Avg	IS	WE	ZH	Avg
Closed-source Models								
GPT5	71.7	91.3	<b>99.8</b>	87.7	74.4	84.5	99.6	85.4
GPT5-R	<b>73.9</b>	<b>98.0</b>	<b>99.8</b>	<b>91.3</b>	<b>80.7</b>	96.8	<b>100.0</b>	<b>92.9</b>
Claude-Sonnet-4	<b>70.9</b>	97.7	99.6	<b>90.3</b>	<b>68.9</b>	98.4	99.7	<b>90.0</b>
Claude-Sonnet-4-R	69.8	97.1	99.6	89.7	68.1	<b>98.5</b>	99.9	89.9
DeepSeek-v3.1-Terminus	61.9	<b>97.4</b>	99.4	87.4	65.2	<b>95.0</b>	98.4	87.1
DeepSeek-v3.1-Terminus-R	<b>68.5</b>	96.2	<b>99.8</b>	<b>88.9</b>	<b>71.8</b>	93.4	<b>99.5</b>	<b>88.6</b>
Open-source Models								
Qwen3-8B	58.0	96.1	98.8	85.6	54.4	91.6	96.6	81.9
Qwen3-8B-R	<b>59.3</b>	<b>97.1</b>	<b>99.8</b>	<b>86.6</b>	<b>61.0</b>	<b>93.7</b>	<b>98.6</b>	<b>85.3</b>
Qwen3-1.7B	<b>51.0</b>	87.0	89.5	77.0	46.2	74.0	71.5	65.2
Qwen3-1.7B-R	50.5	<b>92.0</b>	<b>95.9</b>	<b>80.8</b>	<b>47.1</b>	<b>80.3</b>	<b>81.2</b>	<b>70.7</b>
Qwen2.5-7B-Instruct	<b>47.7</b>	<b>95.0</b>	<b>98.2</b>	<b>81.9</b>	<b>48.6</b>	<b>89.8</b>	<b>94.3</b>	<b>78.8</b>
DeepSeek-R1-Distill-Qwen-7B	43.0	84.9	91.4	74.2	33.6	40.1	37.1	37.4
Llama-3.1-8B-Instruct	40.9	<b>95.8</b>	97.7	80.1	26.0	68.8	<b>78.2</b>	58.6
DeepSeek-R1-Distill-Llama-8B	<b>45.3</b>	93.5	<b>98.2</b>	<b>80.5</b>	<b>34.3</b>	<b>70.4</b>	<b>74.0</b>	<b>60.7</b>
Qwen2.5-1.5B-Instruct	<b>46.2</b>	<b>70.0</b>	54.4	<b>59.0</b>	<b>44.9</b>	<b>54.9</b>	<b>39.3</b>	<b>48.0</b>
DeepSeek-R1-Distill-Qwen-1.5B	33.3	53.7	<b>58.3</b>	48.8	23.3	23.4	22.2	23.1

Table 2: Accuracy (%) of eight LLMs and their reasoning variants on MCQ test sets in two languages. Higher scores between general-purpose and reasoning models are highlighted across three cultures (IS-yellow, WE-blue, ZH-red, Avg-green); best per column is **bolded**.

DeepSeek do not display a significant advantage on the Chinese (ZH) subset. Most models perform well on both Western (WE) and Chinese subsets. Instead, performance on the Islamic (IS) subset is the decisive factor for overall accuracy, with models showing a 10-50% gap compared to Western and Chinese cultures, particularly when Islamic interpretations are presented in Arabic. For instance, Qwen3-8B achieves 98.8% on the ZH subset in English but only 54.4% on the IS in Arabic.

We speculate that while state-of-the-art models have improved in low-resource languages, they still lack cultural nuance optimization. This explains why models, including commercial ones, perform well on Western and Chinese subsets regardless of presentation language, but fail on Islamic cultural assessments, where underrepresented cultural knowledge remains a bottleneck (Chiu et al., 2025).

### Language Variation Still Impacts Old Models.

On the identical MCQ test set presented in two languages, commercial models and Qwen3-8B demonstrate strong robustness to language variation, performing similarly on English and Arabic. In contrast, smaller and older models show clear disadvantages on the Arabic test set. This suggests that larger SOTA models above 7B have been specifi-

Model	English Test			Arabic Test		
	MCQ	GB	TF	MCQ	GB	TF
Qwen3-1.7B	51.0	67.6	52.5	46.2	66.8	48.9
+ LoRA-All	51.2	68.8	52.8	44.1	61.0	48.2
Qwen3-1.7B-R	50.5	66.5	60.0	47.1	67.0	55.9
Qwen3-8B	58.0	67.1	52.7	54.4	68.4	50.3
+ LoRA-All	57.2	67.2	51.4	49.8	68.8	51.4
Qwen3-8B-R	59.3	69.3	58.9	61.0	70.7	55.6
GPT5	71.7	69.6	50.4	74.4	73.9	65.4
GPT5-R	73.9	68.8	55.9	80.7	72.8	69.3

Table 3: Accuracy (%) of Qwen3-1.7B/8B and GPT5 on Islamic culture test sets across three tasks.

cally optimized towards better multilingual communication, handling Arabic better.

### Can Different Task Formulations Reveal Consistent Results?

*Yes.* Given that most models achieve over 90% accuracy on the WE and ZH cultures, we focus on the Islamic subset to assess consistency across task formulations. As shown in Table 3, the GB and TF tasks largely mirror the MCQ trend: models perform better in English than in Arabic, with the exception of GPT5, which shows robust performance in both languages; overall performance of Qwen3-8B and 1.7B follow GPT5. Table 5 for free-form QA exhibits the same pattern. This confirms that model behavior is consistent

Test Set Language SFT Training Data Model	English MCQ							Arabic MCQ						
	Cul	En		Ar		All		En		Ar		All		
		Full	LoRA	Full	LoRA	Full	LoRA	Full	LoRA	Full	LoRA	Full	LoRA	
Qwen3-8B	IS	49.5	57.3	65.4	57.4	60.0	57.2	5.3	53.6	0.5	53.8	34.6	49.8	
	WE	97.7	96.2	97.7	96.4	98.2	96.4	92.8	91.5	70.5	91.9	94.8	92.5	
	ZH	98.5	98.8	98.7	98.8	99.3	99.1	95.5	96.5	92.7	96.7	96.2	96.8	
	Avg	83.7	85.4	88.4	85.5	87.2	85.5	67.6	81.6	55.5	81.9	77.4	81.0	
Qwen3-1.7B	IS	46.2	51.2	50.6	51.4	39.9	51.2	45.8	45.9	48.2	45.2	41.7	44.1	
	WE	94.7	87.7	92.2	88.2	93.1	89.8	84.0	74.5	83.0	76.8	80.9	78.2	
	ZH	95.6	89.5	93.7	90.2	95.6	90.7	83.5	70.9	85.2	72.5	82.5	73.5	
	Avg	80.6	77.4	80.3	77.9	78.0	78.7	72.6	65.2	73.3	66.4	69.7	66.9	
Qwen2.5-7B-Instruct	IS	38.7	48.4	50.8	48.0	42.7	49.9	50.1	48.0	53.6	47.4	51.3	46.4	
	WE	93.1	95.7	96.0	95.5	90.3	96.7	93.9	90.9	91.6	91.5	88.4	93.3	
	ZH	90.9	98.5	96.1	98.4	93.7	98.5	82.3	94.7	92.6	95.7	90.8	96.7	
	Avg	76.5	82.5	82.7	82.3	77.1	83.4	78.0	79.2	80.6	79.6	78.1	80.3	
Llama-3.1-8B-Instruct	IS	40.6	30.9	23.3	34.9	37.3	17.6	26.0	15.4	4.3	12.5	41.6	7.5	
	WE	90.6	95.6	89.1	96.4	79.8	94.2	46.3	58.3	59.0	66.0	54.9	51.1	
	ZH	55.6	97.5	88.4	98.0	53.8	97.8	25.6	70.7	46.9	77.1	32.4	47.8	
	Avg	66.9	77.0	69.5	78.7	60.6	72.5	35.1	48.8	39.8	53.1	45.2	37.4	
Qwen2.5-1.5B-Instruct	IS	55.1	46.2	46.8	45.6	55.6	45.7	45.6	43.8	44.0	43.5	46.6	42.6	
	WE	74.8	70.8	72.6	71.4	73.9	71.7	58.4	55.6	59.7	58.3	59.4	58.4	
	ZH	50.4	53.7	49.2	54.9	48.2	52.6	34.1	39.3	36.3	41.3	32.2	38.6	
	Avg	62.8	59.2	59.0	59.6	62.0	59.2	48.5	48.0	49.1	49.6	48.7	48.7	
Llama-3.2-1B-Instruct	IS	4.9	4.9	4.9	4.9	4.8	4.9	5.7	4.9	4.1	4.9	3.4	4.9	
	WE	19.0	19.2	19.1	19.2	19.1	18.9	19.0	18.9	18.9	18.9	19.1	18.9	
	ZH	18.9	18.8	18.9	18.9	18.9	19.1	19.1	18.9	18.9	18.9	19.1	18.9	
	Avg	14.8	14.9	14.9	14.9	14.8	14.8	15.1	14.8	14.6	14.8	14.4	14.8	

Table 4: Accuracy (%) of five LLMs fine-tuned by full-parameter SFT vs. LoRA SFT in three training dataset settings: English QA, Arabic QA and their mixture. Evaluation on MCQ test sets.

458 across tasks. However, even for the “easier” binary  
459 classification tasks (GB and TF), accuracies remain  
460 below 74%, highlighting limited knowledge of Is-  
461 lamic dream interpretation logic across models.

### 4.3 Full-parameter vs. LoRA SFT

462 We fine-tuned six open-source models for one  
463 epoch under two settings: LoRA and full-parameter  
464 using training data in English, Arabic, and a com-  
465 bination of them, each covering all three cultures.  
466

467 **Does SFT Improve Accuracy?** Figure 5 shows  
468 that, under appropriate setups, SFT generally im-  
469 proves performance in both languages. As detailed  
470 in Table 4, full-parameter SFT is more effective for  
471 smaller models (1/1.5B), while LoRA is preferable  
472 for larger models (7/8B).

473 SFT has a limited impact on Western and Chi-  
474 nese subsets due to their already high baselines.  
475 The Islamic subset is the key driver of overall ac-  
476 curacy. Improvements in Islamic culture (yellow  
477 cells) largely determine performance gains, where  
478 most yellow cells appear when training data is pre-  
479 sented in either Arabic language or both (All), in-  
480 dicating that Islamic cultural nuances can be more  
481 effectively learned in Arabic language than through

Model	Correct		GT-Sent		AI-Sent	
	En	Ar	En	Ar	En	Ar
Qwen3-1.7B	3.3	2.6	50.5	52.2	61.3	56.4
+ LoRA-All	2.7	1.9	50.5	53.0	57.3	49.6
Qwen3-1.7B-R	3.4	2.8	51.6	51.4	62.7	56.8
Qwen3-8B	3.6	3.5	50.2	52.2	60.0	56.8
+ LoRA-All	3.2	2.7	51.1	51.1	59.4	53.7
Qwen3-8B-R	3.6	3.5	51.2	52.0	59.3	55.2
GPT5	4.1	4.1	50.5	52.7	56.7	56.5
GPT5-R	4.1	4.2	51.7	53.2	55.4	55.7
<b>Avg</b>	3.5	3.2	50.9	52.2	59.0	55.1

Table 5: Free-form QA responses’ correctness and sentiment under non-reasoning, LoRA-SFT-All and reasoning. QA test set includes 2,634 questions reflecting Western and Chinese cultures, presented in two languages. Model responses are more positive than humans’ roughly half-half distribution.

English language. Similarly, for Western culture (blue cells), gains are mainly observed with English or bilingual training data.

482 However, SFT can also degrade accuracy, par-  
483 ticularly when applying full-parameter SFT to  
484 larger models (e.g., Llama-3.1-8B and Qwen3-  
485 8B), which requires more training epochs to re-  
486 cover disrupted knowledge as below.  
487  
488  
489

### Impact of Training Epochs Based on Qwen3-

490

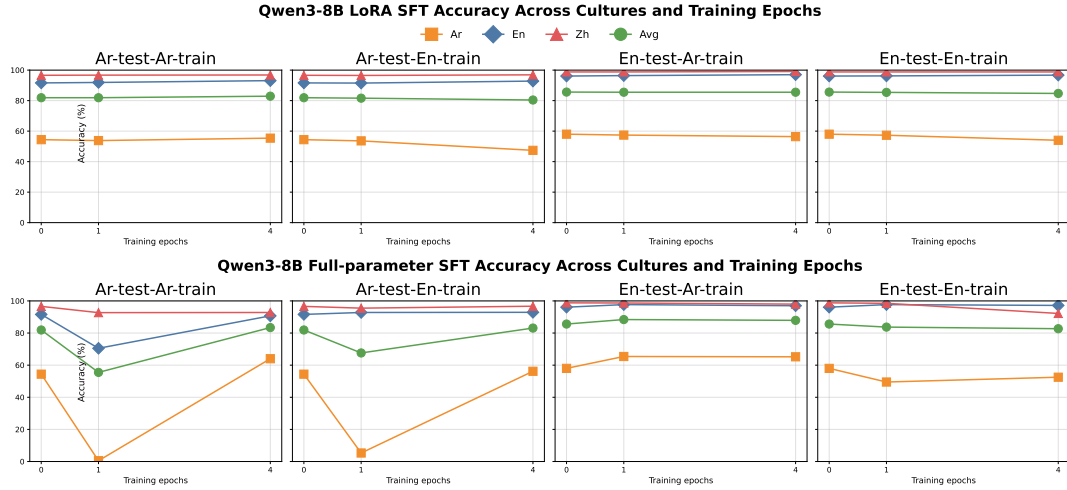


Figure 4: Qwen3-8B Full vs. LoRA SFT accuracy across training epochs and cultures.

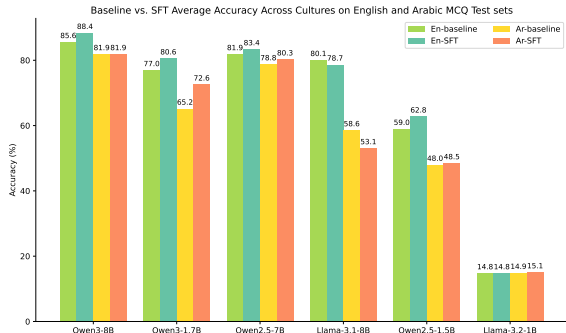


Figure 5: Baseline vs. SFT average accuracy across cultures on English and Arabic MCQ test sets.

8B, we compare the baseline without SFT against full-parameter and LoRA SFT with 1 and 4 training epochs, analyzing the impact of training epochs on dream interpretation across cultures (Figure 4). We find two main results. First, with LoRA SFT, both the number of epochs and the training language have a negligible impact on accuracy when evaluated on English and Arabic MCQ test sets. Second, full-parameter SFT with either English or Arabic data disrupts the model’s internal knowledge of Islamic dream interpretation, causing a sharp drop after one epoch, while leaving Western and Chinese subsets largely unaffected. Additional epochs can recover and enhance performance on the Islamic culture subset.

These results suggest that Western and Chinese interpretive traditions are already deeply internalized into parameters during model pretraining, whereas Islamic nuances might be acquired during post-training, which were encoded more shallowly, yielding lower baselines and greater sensitivity to SFT. This aligns with previous findings that knowledge learned during fine-tuning is

not well-grounded and tends to induce hallucinations (Gekhman et al., 2024).

#### 4.4 Are AIs More Positive than Humans?

We evaluated three model pairs: Qwen3-7B/1.7B and GPT5 (chat and reasoning) on 2,634 free-form QA examples spanning Western and Chinese cultures in two languages. As shown in Table 5, human interpretations exhibit no sentiment preference, with positive and negative each holding half. In contrast, model outputs show a slight positive tendency, with 59% positive responses in English and 55% in Arabic.

## 5 Conclusion

We introduce a bilingual dream interpretation dataset spanning Western, Chinese, and Islamic cultures, with 31,877 unique entries. We find that less than 22% of dream symbols overlap across cultures, highlighting the cultural divergence. Using four task formulations, we evaluated 17 LLMs and fine-tuned six models, demonstrating the advantage of the *one model, two modes* paradigm and the capacity of LLMs to handle low-resource languages. However, while models can process the language, they fail to capture cultural nuances, performing markedly worse on Islamic culture interpretations than on Western and Chinese ones. We further identify optimal SFT configurations for models 1-8B, and observe that humans interpret dreams with a roughly neutral balance of positive and negative meanings, whereas models exhibit a slight positive bias. In future work, we plan to focus on enhancing models’ sensitivity to cultural nuances and fine-tuning them with reasoning trajectories alongside direct interpretations.

## 548 **Limitations and Future Work**

549 **Niche Scope** Dream interpretation is a unique  
550 but niche topic. Experimental results and findings  
551 may not generalize to broader cultural reasoning  
552 tasks.

553 **Insufficient Human Validation** We used LLMs  
554 to transform dream interpretation entries into four  
555 task formulations. Although gold pairs were pro-  
556 vided and the model was instructed to preserve con-  
557 tent, LLM-induced errors may still occur. We man-  
558 ually inspected 50-100 samples per task for each  
559 culture and language to verify translation and task  
560 formulation quality. However, this limited valida-  
561 tion covers only a small fraction of the dataset and  
562 lacks sufficient human auditing and inter-annotator  
563 agreement, leaving potential template and judge  
564 biases.

565 For evaluating free-form QA responses, we em-  
566 ployed an LLM-based judge (GPT-5-mini) to as-  
567 sess correctness and sentiment. While efficient, this  
568 approach may introduce systematic bias. A human-  
569 evaluated subset with inter-annotator agreement  
570 would strengthen the credibility of the evaluation.  
571 Furthermore, the engagement with human experts,  
572 such as psychoanalysts or cultural anthropologists,  
573 for judging the correctness of free-form answers  
574 would be ideal for a subjective and culturally deep  
575 task.

576 **Potential Pre-training Data Contamination**  
577 Our data are mostly derived from dream dictio-  
578 naries, publicly scraped web sources and books  
579 that were likely included in the pre-training cor-  
580 pora of the evaluated models. This raises concerns  
581 about the validity of zero-shot performance, as the  
582 models may simply be recalling information seen  
583 during pre-training rather than demonstrating gen-  
584 uine interpretation or reasoning capabilities.

585 To mitigate this risk, we design evaluation tasks  
586 with diverse formulations including QA, MCQ,  
587 Good or Bad, True or False. While models may  
588 have memorized free-form text, correctly interpret-  
589 ing dreams when they appear in different ways can  
590 still reflect models' generalization and alignment  
591 with users' expectations.

592 **Source Reliability and Cultural Scope** Some  
593 raw data sources lack scholarly validation and may  
594 not be fully reliable or culturally representative.  
595 Our analysis based on the collected dream symbols  
596 may not accurately reflect the true distribution of  
597 cultural overlap and divergence. Future work could

improve representativeness by collecting first-hand  
598 sharing from real users on social media platforms,  
599 rather than relying primarily on dream dictionary or  
600 websites. Reliability could be further enhanced by  
601 using more authoritative sources, such as published  
602 books, scholarly translations, religious texts, and  
603 academic literature.

604 In addition to addressing the limitations above,  
605 we plan to select challenging and cross-cultural  
606 examples and host a public leaderboard or bench-  
607 mark portal to support evaluations of future mod-  
608 els. We also plan to explore research questions,  
609 such as whether differences between LoRA and  
610 full-parameter fine-tuning reveal novel behaviors  
611 in cross-cultural symbolic tasks or low-resource  
612 settings, and whether LLM performance can be  
613 further improved through prompt or context opti-  
614 mization strategies.

## 615 **Ethical Statement**

616 **Data License and Copyright** The Arabic  
617 datasets were derived from five books, four of  
618 which are ancient books whose authors died be-  
619 tween 729-1826 and are therefore in the public  
620 domain. The fifth was from the book "Islamic  
621 Dream Interpretation" by Dr. Khaled Al-Anbari  
622 downloaded from <https://www.alanbary.com/>,  
623 which states "all rights reserved." However, we  
624 only extracted 733 dream symbols and their inter-  
625 pretations for research purposes with task refor-  
626 mulations and heavy data transformations. This  
627 limited use does not substitute for or replicate the  
628 original work and is permissible under established  
629 copyright exceptions for scholarly research, includ-  
630 ing U.S. Copyright Act ("fair use"), the EU Copy-  
631 right Directive (EU) 2019/790 (Articles on text  
632 and data mining for scientific research), and cor-  
633 responding "fair dealing" provisions in the UK,  
634 Canada, and other jurisdictions. Similarly for  
635 77 Chinese examples extracted from the website  
636 <https://www.yourchineseastrology.com/>.

637 English data were either collected from Kaggle  
638 under the license of Apache 2.0 or Huggingface<sup>1</sup>  
639 under the license of GPLv3. Both of them allow  
640 research usage of the data. We will obey the policy  
641 of both. For the public webpage [myislamicdream.com](http://myislamicdream.com),  
642 there is no specific license.

643 Though two Chinese datasets were col-  
644 lected from websites, including <https://www.>

<sup>1</sup>[https://huggingface.co/datasets/  
JosephusCheung/GuanacoDataset](https://huggingface.co/datasets/JosephusCheung/GuanacoDataset)

[yourchineseastrology.com/](http://yourchineseastrology.com/) and <https://www.zgj Morg.com/>, the content shown in these webpages is from an ancient Chinese book — zhougongjiemeng (《周公解梦》).

We did not collect any dream-interpretation pairs with personal information. Also, we did not use the raw data. We extracted dream symbols rather than using the original dream descriptions. This approach complies with website privacy policies and protects user anonymity. See more about data license in Table 6. If any data is found to require additional permissions, we will contact the content owner for authorization or exclude the relevant portions from the final release.

## References

Edgar Altszyler, Sidarta Ribeiro, Mariano Sigman, and Diego Fernández Slezak. 2017. The interpretation of dream meaning: Resolving ambiguity using latent semantic analysis in a small corpus of text. *Consciousness and cognition*, 56:178–187.

Lorenzo Bertolini, Valentina Elce, Adriana Michalak, Giulio Bernardi, and Julie Weeds. 2023. Automatic scoring of dream reports’ emotional content with large language models. *arXiv preprint arXiv:2302.14828*.

Lorenzo Bertolini, Adriana Michalak, and Julie Weeds. 2024. Dreamy: a library for the automatic analysis and annotation of dream reports with multilingual large language models. In *Sleep Medicine*, volume 115, pages 406–407. Elsevier RADARWEG 29, 1043 NX AMSTERDAM, NETHERLANDS.

Mark Blagrove, Laura Farmer, and Elvira Williams. 2004. The relationship of nightmare frequency and nightmare distress to well-being. *Journal of sleep research*, 13(2):129–136.

Abigail P Blyler and Martin EP Seligman. 2024. Personal narrative and stream of consciousness: an ai approach. *The Journal of Positive Psychology*, 19(4):592–598.

Rosalind Cartwright. 2011. Dreaming as a mood-regulation system. In *Principles and practice of sleep medicine*, pages 620–627. Elsevier.

Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin, Chan Young Park, Shuyue Stella Li, Sahithya Ravi, Mehar Bhatia, Maria Antoniak, Yulia Tsvetkov, Vered Shwartz, and Yejin Choi. 2025. **CulturalBench: A robust, diverse and challenging benchmark for measuring LMs’ cultural knowledge through human-AI red-teaming**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 25663–25701, Vienna, Austria. Association for Computational Linguistics.

Francis Crick and Graeme Mitchison. 1983. The function of dream sleep. *Nature*, 304(5922):111–114.

Susanne Diekelmann and Jan Born. 2010. The memory function of sleep. *Nature reviews neuroscience*, 11(2):114–126.

G William Domhoff. 2017. *The emergence of dreaming: Mind-wandering, embodied simulation, and the default network*. Oxford University Press.

G William Domhoff and Adam Schneider. 2008. Studying dream content using the archive and search engine on dreambank. net. *Consciousness and Cognition*, 17(4):1238–1247.

Valentina Elce, Giacomo Handjaras, and Giulio Bernardi. 2021. The language of dreams: Application of linguistics-based approaches for the automated analysis of dream experiences. *Clocks & Sleep*, 3(3):495–514.

S Freud. 1900. The interpretation of dreams, vol. 4, trans. *J. Strachey*. (London, The Hogarth Press, pages 26–28.

Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. **Does fine-tuning LLMs on new knowledge encourage hallucinations?** In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7765–7784, Miami, Florida, USA. Association for Computational Linguistics.

Calvin Hall and Robert Van de Castle. 1966. The content analysis of dreams.

Daniel E Harris-McCoy. 2012. Artemidorus’ oneirocritica: Text, translation, and commentary.

P Hauri. 1975. Categorization of sleep mental activity for psychophysiological studies. *The experimental study of sleep: Methodological problems*, pages 271–281.

Sheldon Juncker. 2023. Dreaming with ai. *Poligrafi: revija za religiologijo, mitologijo in filozofijo*, 28(109/110).

Tracey L Kahan and Stephen P LaBerge. 2011. Dreaming and waking: Similarities and differences revisited. *Consciousness and Cognition*, 20(3):494–514.

Mayte H Laureano and Hiram Calvo. 2024. Computational study of dream interpretations: Psychoanalytic human vs artificial analyses. In *2024 IEEE Congress on Evolutionary Computation (CEC)*, pages 1–9. IEEE.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. **Self-refine: Iterative refinement with self-feedback**.

699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751

752	Patrick McNamara, Kelly Duffy-Deno, Tom Marsh, and	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	807
753	Thomas Jr Marsh. 2019. Dream content analysis	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	808
754	using artificial intelligence. <i>International Journal of</i>	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	809
755	<i>Dream Research</i> , pages 42–52.	Joseph E. Gonzalez, and Ion Stoica. 2023. <a href="#">Judging</a>	810
		<a href="#">llm-as-a-judge with mt-bench and chatbot arena</a> .	811
756	David Nadeau, Catherine Sabourin, Joseph De Koninck,	Xiaofang Zheng and Richard Schweickert. 2023. Dif-	812
757	Stan Matwin, Peter D Turney, et al. 2006. Automatic	ferentiating dreaming and waking reports with au-	813
758	dream sentiment analysis. In <i>Proceedings of the</i>	tomatous text analysis and support vector machines.	814
759	<i>Workshop on Computational Aesthetics at the Twenty-</i>	<i>Consciousness and Cognition</i> , 107:103439.	815
760	<i>First National Conference on Artificial Intelligence</i> .		
761	Kate Niederhoffer, Jonathan Schler, Patrick Crutchley,		
762	Kate Loveys, and Glen Coppersmith. 2017. In your		
763	wildest dreams: the language and psychological fea-		
764	tures of dreams. In <i>Proceedings of the Fourth Work-</i>		
765	<i>shop on Computational Linguistics and Clinical Psy-</i>		
766	<i>chology—From Linguistic Signal to Clinical Reality</i> ,		
767	pages 13–25.		
768	Yuval Nir and Giulio Tononi. 2010. Dreaming and		
769	the brain: from phenomenology to neurophysiology.		
770	<i>Trends in cognitive sciences</i> , 14(2):88–100.		
771	Amir H Razavi, Stan Matwin, Joseph De Koninck, and		
772	Ray Reza Amini. 2014. Dream sentiment analysis		
773	using second order soft co-occurrences (sosco) and		
774	time course representations. <i>Journal of Intelligent</i>		
775	<i>Information Systems</i> , 42(3):393–413.		
776	Camila Sanz, Federico Zamberlan, Earth Erowid, Fire		
777	Erowid, and Enzo Tagliazucchi. 2018. The experi-		
778	ence elicited by hallucinogens presents the high-		
779	est similarity to dreaming within a large database of		
780	psychoactive substance reports. <i>Frontiers in neuro-</i>		
781	<i>science</i> , 12:7.		
782	Michael Schredl. 2010. Dream content analysis: Basic		
783	principles. <i>International Journal of Dream Research</i> ,		
784	pages 65–73.		
785	Francesca Siclari, Benjamin Baird, Lampros Per-		
786	ogamvros, Giulio Bernardi, Joshua J LaRocque,		
787	Brady Riedner, Melanie Boly, Bradley R Postle, and		
788	Giulio Tononi. 2017. The neural correlates of dream-		
789	ing. <i>Nature neuroscience</i> , 20(6):872–878.		
790	Serge Thill and Henrik Svensson. 2011. The inception		
791	of simulation: a hypothesis for the role of dreams		
792	in young children. In <i>Proceedings of the Annual</i>		
793	<i>Meeting of the Cognitive Science Society</i> , volume 33.		
794	Matthew P Walker and Els van Der Helm. 2009.		
795	Overnight therapy? the role of sleep in emotional		
796	brain processing. <i>Psychological bulletin</i> , 135(5):731.		
797	Erin J Wamsley. 2014. Dreaming and offline memory		
798	consolidation. <i>Current neurology and neuroscience</i>		
799	<i>reports</i> , 14(3):433.		
800	Erin J Wamsley and Robert Stickgold. 2011. Mem-		
801	ory, sleep and dreaming: experiencing consolidation.		
802	<i>Sleep medicine clinics</i> , 6(1):97.		
803	Antonio Zadra and Robert Stickgold. 2021. <i>When</i>		
804	<i>brains dream: Understanding the science and mys-</i>		
805	<i>tery of our dreaming minds</i> . WW Norton & Com-		
806	pany.		

816	<b>Contents</b>	
817	<b>1 Introduction</b>	<b>1</b>
818	<b>2 Dataset</b>	<b>2</b>
819	<b>3 Dream Interpretations Across Cultures</b>	<b>3</b>
820	3.1 Culturally Dominant Symbols . . .	3
821	3.2 Overlap and Uniqueness of Symbols	4
822	3.3 Summary . . . . .	4
823	<b>4 Experiments</b>	<b>5</b>
824	4.1 Experimental Setup . . . . .	5
825	4.2 Who Wins? . . . . .	5
826	4.3 Full-parameter vs. LoRA SFT . .	7
827	4.4 Are AIs More Positive than Humans?	8
828	<b>5 Conclusion</b>	<b>8</b>
829	<b>A Related Work</b>	<b>12</b>
830	A.1 Tracing Dream Study . . . . .	12
831	A.2 Dream Interpretation Datasets . .	12
832	A.3 Automatic Dream Analysis . . . .	13
833	A.4 LLM Dream Analysis . . . . .	13
834	<b>B Dataset Curation</b>	<b>13</b>
835	B.1 Raw Data Sources . . . . .	14
836	B.1.1 Islamic Culture Corpora . .	14
837	B.1.2 Western Culture Corpora . .	14
838	B.1.3 Chinese Cultural Corpora . .	14
839	B.2 Dream Symbol Category . . . . .	15
840	B.3 Four Task Formulation . . . . .	15
841	B.3.1 QA Conversion . . . . .	16
842	B.3.2 MCQ Generation . . . . .	16
843	B.3.3 Good/Bad and True/False . .	17
844	B.4 Training and Evaluation Splits . .	17
845	<b>C Prompts for QA Generation</b>	<b>17</b>
846	<b>D Prompts for MCQ Generation</b>	<b>17</b>
847	<b>E Evaluation</b>	<b>17</b>
848	E.1 Models . . . . .	17
849	E.2 Inference Configuration . . . . .	18
850	E.3 Prompts for Evaluation . . . . .	18
851	<b>A Related Work</b>	
852	<b>A.1 Tracing Dream Study</b>	
853	Dreams have been a topic of human curiosity for	
854	centuries, dating back to ancient times. Historical	
855	approaches to dream interpretation began with	

Artemidorus in the 2nd century AD, who systemat-	856
ically studied dream content and proposed interpre-	857
tive techniques in his work <i>Oneirocritica</i> (Harris-	858
McCoy, 2012). A major shift occurred in the 19th	859
century with Sigmund Freud’s <i>The Interpretation of</i>	860
<i>Dreams</i> (Freud, 1900). He studied that dreams ex-	861
press repressed desires and serve to relieve internal	862
tension, thereby supporting sleep and overall well-	863
being. Later theories proposed that dreams contrib-	864
ute to emotional regulation and conflict resolu-	865
tion (Cartwright, 2011; Walker and van Der Helm,	866
2009), memory consolidation (Diekelmann and	867
Born, 2010), and forgetting irrelevant information	868
to enhance learning (Crick and Mitchison, 1983).	869
Some perspectives have also compared dreams to	870
simulations that help individuals prepare for future	871
challenges and threats (Thill and Svensson, 2011).	872
In recent decades, scientific interest in dreams	873
has increased, particularly in understanding their	874
psychological and neurological relevance. Re-	875
search has highlighted connections between dream-	876
ing and psychophysical health, and pointed to po-	877
tential roles in sleep-dependent memory processing	878
(Wamsley and Stickgold, 2011; Wamsley, 2014;	879
Zadra and Stickgold, 2021). Furthermore, as in-	880
ternally generated conscious experiences, dreams	881
provide a valuable model for studying the nature of	882
consciousness itself (Nir and Tononi, 2010; Siclari	883
et al., 2017). Despite the growing interest, the pro-	884
cesses underlying dream generation and the exact	885
functions of dreams remain partially understood.	886
A major challenge in dream research is the diffi-	887
culty in quantitatively assessing dream content in a	888
reproducible manner (Elce et al., 2021).	889
Dream reports, which document the content of	890
dreams recalled by individuals, are essential to the	891
study of dreaming. These reports offer insights into	892
the connection between dreams and waking life	893
(Blagrove et al., 2004) and have long been used as	894
a medium to examine conscious experience during	895
sleep (Nir and Tononi, 2010; Siclari et al., 2017).	896
Thus, many efforts have focused on collecting and	897
analyzing dream narratives or reports (Hall and	898
Van de Castle, 1966; Hauri, 1975; Schredl, 2010).	899
<b>A.2 Dream Interpretation Datasets</b>	900
DreamBank is among the most widely used re-	901
sources, comprising thousands of annotated dream	902
narratives used in psychological and linguistic re-	903
search (Domhoff and Schneider, 2008). In addition	904
to public resources, Laureano and Calvo (2024),	905
compiled a proprietary dataset of dreams from 20	906

907 patients, with each analyzed by multiple human  
908 psychoanalysts and GPT-4. However, manual an-  
909 notation is time-intensive and typically requires  
910 trained human experts, which can hinder the repro-  
911 ducibility and scalability of dream research (Elce  
912 et al., 2021).

### 913 A.3 Automatic Dream Analysis

914 To address these limitations, researchers have in-  
915 creasingly explored the use of natural language  
916 processing (NLP) tools to automate the analysis of  
917 dream reports. They applied linguistic and com-  
918 putational techniques to analyze dream content.  
919 Niederhoffer et al. (2017) examined dream nar-  
920 ratives using the Hall and Van de Castle (HVDC)  
921 framework, highlighting a higher prevalence of neg-  
922 ative emotions, especially sadness in dreams. Mc-  
923 Namara et al. (2019) introduced the Dream Content  
924 Analysis System (DCAS), which used AI to iden-  
925 tify gender-related patterns in dream themes and  
926 their relation to mood. Elce et al. (2021) demon-  
927 strated that methods such as graph analysis, dis-  
928 tributional semantics, and dictionary-based tech-  
929 niques can capture both semantic and structural  
930 properties of dream narratives.

931 Overall, these approaches often rely on mod-  
932 els trained on general-purpose corpora such as  
933 Wikipedia (Nadeau et al., 2006; Razavi et al., 2014;  
934 Altszyler et al., 2017; Sanz et al., 2018; McNamara  
935 et al., 2019; Bertolini et al., 2023). However, there  
936 are debates regarding how closely dream reports re-  
937 semble other general types of textual data (Kahan  
938 and LaBerge, 2011; Domhoff, 2017; Zheng and  
939 Schweickert, 2023). Some evidence suggests that  
940 the semantic characteristics of dream reports may  
941 diverge significantly from those found in waking  
942 narratives (Altszyler et al., 2017). If dream reports  
943 are indeed unique in structure and content, the ef-  
944 fectiveness of using general-domain NLP models,  
945 especially in unsupervised settings, could be sub-  
946 stantially limited (Bertolini et al., 2023).

### 947 A.4 LLM Dream Analysis

948 With the advancement of LLMs, there is grow-  
949 ing interest in their use for automated dream anal-  
950 ysis. Building on DreamBank, the DReAMy  
951 toolkit (Bertolini et al., 2024) offers an open-source  
952 framework that leverages multilingual LLMs to  
953 automatically annotate dream reports for emo-  
954 tions and characters based on HVDC framework.  
955 Blyler and Seligman (2024) used GPT-4 to gener-  
956 ate personal narratives and streams of conscious-

957 ness, while Juncker (2023) employed the model for  
958 psychological reflection and dream interpretation.  
959 Laureano and Calvo (2024) compared GPT-4’s in-  
960 terpretations with those of human analysts, finding  
961 that while both were coherent, the AI displayed  
962 distinctive linguistic patterns. GPT-4 tended to  
963 use more semantic categories such as vision and  
964 health-related terms and fewer grammatical ele-  
965 ments like impersonal pronouns. Using LIWC fea-  
966 tures and Naïve Bayes classification, they achieved  
967 99% accuracy in differentiating between the two  
968 sources, indicating the significant differences be-  
969 tween human-written and LLM-generated dream  
970 analysis. These approaches reduce reliance on man-  
971 ual dream annotation while also expose the gap  
972 between human and LLM in dream interpretation.

973 **Summary** While these efforts have advanced  
974 dream understanding, little attention has been de-  
975 voted to *dream interpretation*, which seeks to de-  
976 rive symbolic, cultural, and contextual meaning  
977 from dream content.

978 Most publicly available datasets and studies on  
979 dream interpretation are centered on English and  
980 adopt linguistic, emotional, psychological, or bi-  
981 ological perspectives to analyze dream narratives.  
982 Such approaches primarily depend on an individ-  
983 ual’s current condition and recent experiences. By  
984 contrast, in cultures such as China, dream inter-  
985 pretation is not only dependent on individuals, but  
986 also grounded in centuries of accumulated wis-  
987 dom and collective observation. Ancient dream-  
988 interpretaion records like ZhougongJieMeng pro-  
989 vide population-level generalizations, functioning  
990 as a form of prior knowledge or statistical distri-  
991 bution of interpretations across thousands of years.  
992 These traditions draw on systematic principles such  
993 as the Five Elements (metal, wood, water, fire,  
994 earth) and Yin–Yang to derive meaning, offering a  
995 perspective that extends beyond the individual to  
996 collective cultural experience.

997 To mitigate this gap, we collect a dream inter-  
998 pretation datasets from three culture sources including  
999 Arabic, Chinese and Western context, presented in  
1000 two languages (Arabic and English). Meanwhile,  
1001 we examine current state-of-the-art LLMs in dream  
1002 interpretations across cultures and improved it by  
1003 continuous training.

## 1004 B Dataset Curation

1005 In this section, we elaborate how we curate the  
1006 dataset from collecting raw dream-interpretation

pairs from three cultural sources, to converting them into four task formulations, to splitting them into training and evaluation benchmark.

## B.1 Raw Data Sources

To construct our dataset, we first collected raw dream–interpretation pairs from a diverse set of cultural traditions and online repositories. The sources span three major cultural contexts: *Islamic*, *Western*, and *Chinese*. Within each culture, we curated data from multiple publicly available websites and pre-compiled datasets, ensuring broad coverage of interpretations across different schools and perspectives. Table 6 summarizes the raw data sources and their respective sizes. These raw collections serve as the foundation for subsequent task-specific transformations and benchmark construction.

### B.1.1 Islamic Culture Corpora

Our dataset combines classical and modern sources of Islamic dream interpretation. Most of the content comes from traditional texts written by well-known scholars in books including **Al-Nabulsi**, **Al-Ihsaei**, **Al-Anbari**, **Ibn Sirin**, and **Ibn Shahin**. We cleaned the data and extracted dream-interpretation entries. Each entry typically consists of a dream symbol paired with an interpretation. When a dream symbol appeared in multiple sources with different interpretations, we retained all available explanations to capture diverse scholarly perspectives, reflecting the natural evolution of dream interpretation from traditional to modern views. After removing duplicates, the dataset includes a total of **5,568** entities.

### B.1.2 Western Culture Corpora

Western culture corpora combine modern, classical, and community-based interpretations. It includes structured resources and digitized texts that cover a range of symbolic meanings and interpretation styles. The raw collection draws from four sources: (i) *Dictionary of Dreams*, 1,041 symbol-interpretation pairs; (ii) *Dream Dictionary*, 2,080 entries of a similar format; (iii) *DreamBook (Guanaco Format)*, 9,497 user written dreams paired with interpretations generated by a language model; and (iv) *myIslamicDream*<sup>7</sup>, 96,404 symbolic entries presented in English. The combined raw count is 109,022 entries.

Content from *myislamicdream.com* was collected by a custom Python scraper. The script iterated over paginated tag pages at the pattern

`/tags/{page}`, extracted internal links that ended with `.html`, deduplicated links, and then fetched each page. For each page, the parser retained only pages with at least one non-empty paragraph and a minimum visible text length, and it discarded short or empty pages. The crawler used standard request headers and added a one second delay between requests.

**Cleaning and Consolidation** After collection, we cleaned and consolidated to obtain unique pairs. We first removed empty English interpretations, and then normalized the interpretation text by deleting boilerplate lead-in phrases, such as “Refer to ...”, removing parenthetical source notes starting with “(Provided ...)” and deleting footer style material appended at the end of some pages, for instance trailing text “People Who Read This Article Also Read” and embedded domain mentions like “(... *www.xyz.org* ...)”. Afterwards, we excluded entries that were cross references (e.g., “See ...”) rather than dream interpretations, and standardized whitespace by removing line breaks and collapsing multiple spaces. Finally, we dropped any residual empty interpretations that arose after text normalization and performed exact match de-duplication on the normalized interpretations to consolidate repeated content across sources.

When a symbol appeared in more than one source, we retained distinct interpretations that provided substantive content and removed non-informative duplicates. We collected **16,720** unique dream-interpretation pairs.

### B.1.3 Chinese Cultural Corpora

The Chinese dataset includes traditional ZhouGong interpretations, structured categorical records, and contemporary astrological perspectives. The collection is drawn from three sources: (i) *Zhougong Dream Dictionary*, parsed from the Zhou Gong website<sup>2</sup>, yielding 9,508 entries in Chinese; (ii) *Zhougongjiemeng Database*, reconstructed from a publicly available SQL export<sup>3</sup>, containing 9,543 entries in Chinese; and (iii) *Your Chinese Astrology*, an English-language site presenting dream interpretations in a Chinese astrological framework<sup>4</sup>, contributing 77 entries.

<sup>2</sup><https://www.zgjm.org>

<sup>3</sup><https://blog.csdn.net/jianghulangzhongshen/article/details/107043786>

<sup>4</sup><https://www.yourchineseastrology.com>

Dataset	Data License	Lang	Raw Size	Size
tafsiralahlam.net (Al-Ahsa'i) <sup>2</sup>	Book by Al Ahsa'i (1826)	AR	441	424
tafsiralahlam.net (Ibn Shahin) <sup>2</sup>	Book by Ibn Shahin (1468)	AR	1,043	1,034
tafsiralahlam.net (Ibn Sirin) <sup>2</sup>	Book by Ibn Sireen (729)	AR	1,096	1,088
tafsiralahlam.net (Nabulsi) <sup>2</sup>	Book by Nabulsi (1731)	AR	2,291	2,291
alanbar_bary.com (Dr. Khaled Al-Anbari) <sup>3</sup>	Book by AlNbari (allow research use)	AR	733	731
<b>Total Islamic Entries</b>				<b>5,568</b>
Dictionary of Dreams (kaggle/manswad) <sup>4</sup>	Apache 2.0	EN	1,040	898
Dream Dictionary (kaggle/yuvrajsanghai) <sup>5</sup>	Apache 2.0	EN	2,080	2,041
Dreambook Guanaco Format (hf/n3rd0) <sup>6</sup>	GPLv3 (GNU General Public License v3)	EN	9,497	7,662
myislamicdream.com <sup>7</sup>	No License	EN	96,404	6,119
<b>Total Western Entries</b>				<b>16,728</b>
yourchineseastrology.com <sup>8</sup>	–	ZH	77	77
zgjorg.com (Zhougong) <sup>9</sup>	Ancient Chinese book 《周公解梦》	ZH	9,508	7,325
Zhougongjiemeng database <sup>10</sup>	CC 4.0 BY-SA	ZH	9,543	2,187
<b>Total Chinese Entries</b>				<b>9,589</b>
<b>Total</b>				<b>31,877</b>

Table 6: Raw source and size across three cultures. Size refers to numbers after preprocessing.

**Cleaning and Consolidation** Both the *Zhougong Dream Dictionary* and the *Zhougongjiemeng Database* consist of collections of HTML pages, each processed with a custom parser. During preprocessing, empty or redundant entries were removed, as many of these corresponded to boilerplate pages that merely linked to other content rather than providing substantive interpretations. Because both sources originate from the traditional Zhou Gong corpus, a substantial amount of content was duplicated across them. To consolidate overlapping material, we applied a semantic similarity model (paraphrase-multilingual-MiniLM-L12-v2) to compare entries. Pairs with a similarity score above 0.8 were treated as duplicates, and in such cases the longer, more detailed interpretation was retained. Both datasets also underwent additional text cleaning. Redundant lead-in phrases, such as “梦见.....意味着什么? 阅读本文的人还阅读了”, were removed using regular expressions to retain only the substantive interpretation content. After this deduplication and cleaning process, the resulting Chinese dataset contained 9,512 unique entries.

The *Your Chinese Astrology* source was processed with a custom parser, which cleaned dream titles by removing boilerplate prefixes (e.g., “Dream Meaning and Interpretation about/of”) and extracted interpretation content. The resulting text was normalized and converted into a plain, consistently formatted representation suitable for further

processing. This source contributed 77 entries in English. Therefore, there are 9,589 entries in total from Chinese culture.

## B.2 Dream Symbol Category

We applied GPT-4o using prompt shown in Figure 8 to categorize dream symbols for all three cultures into 17 types including prophets and messengers; companions and righteous figures; people and social roles; body parts; animals; birds and insects; places and landmarks; natural elements and phenomena; tools and objects; food and drink; religious symbols and practices; emotional and psychological states; actions and events; time-related concepts; abstract or ambiguous symbols; and symbols related to death and the afterlife. Figure 2 shows the distribution of dream categories.

## B.3 Four Task Formulation

Based on clean dream-interpretation pairs, we formulated four tasks including free-form question

<sup>2</sup><https://tafsiralahlam.net>

<sup>3</sup><https://www.alanbary.com>

<sup>4</sup><https://www.kaggle.com/datasets/manswad/dictionary-of-dreams>

<sup>5</sup><https://www.kaggle.com/datasets/yuvrajsanghai/dream-dictionary>

<sup>6</sup><https://huggingface.co/datasets/JosephusCheung/GuanacoDataset>

<sup>7</sup><https://www.myislamicdream.com>

<sup>8</sup><https://www.yourchineseastrology.com>

<sup>9</sup><https://www.zgjorg.com>

<sup>10</sup><https://blog.csdn.net/jianghulangzhongshen/article/details/107043786>

1152	answering (QA), multi-choice question answering	questions per entry). This results in 16,720 QA	1202
1153	(MCQ), is it a <i>good or bad</i> dream (GB), and is this	samples for Western cultural resources and 9,589	1203
1154	interpretation of the dream <i>true or false</i> (TF).	QA samples from Chinese culture.	1204
1155	<b>B.3.1 QA Conversion</b>	<b>B.3.2 MCQ Generation</b>	1205
1156	<b>Islamic Culture</b> To extend the Islamic subset	<b>Islamic Culture</b> We extend the Islamic dream	1206
1157	into a question–answer format, we generated cul-	interpretation pairs into MCQ format by combining	1207
1158	turally grounded questions based on the available	each symbol with its original interpretation and a	1208
1159	symbol–interpretation pairs. For each dream sym-	set of four distractors. The correct answer is always	1209
1160	bol, multiple questions were formulated in five ma-	taken directly from the original sources, ensuring	1210
1161	major variants of Arabic languages: Modern Standard	that the interpretive content remains unchanged.	1211
1162	Arabic, Gulf, Egyptian, Levantine, and Maghrebi.	distractors are sampled from other symbols within	1212
1163	The prompts were designed to simulate realistic	the same semantic category, which makes them top-	1213
1164	user queries, ranging from conversational expres-	ically plausible while still incorrect. This category-	1214
1165	sions to more formal requests, thereby reflecting	based sampling increases the difficulty of the task	1215
1166	how speakers across regions might naturally in-	and prevents trivial elimination strategies.	1216
1167	quire about dream meanings.	To generate the question text and options, we	1217
1168	Crucially, the answers paired with these ques-	provide the symbol, its correct interpretation, and	1218
1169	tions were always the original interpretations taken	the four distractors to an LLM Gemini-2.0-flash,	1219
1170	directly from the classical Islamic sources. No	together with a structured prompt shown in Fig-	1220
1171	modifications or paraphrases were introduced: the	ure 9. The model is instructed to produce a single	1221
1172	interpretive content remains identical to that pro-	well-formed question in Modern Standard Arabic,	1222
1173	vided by the authoritative texts. The only variation	accompanied by exactly five answer options (A–E)	1223
1174	lies in the phrasing of the questions, which differs	with the correct answer randomly positioned.	1224
1175	between dialects.	<b>Western/Chinese Culture</b> Similar to the proce-	1225
1176	This procedure produced a dialect-sensitive Ara-	cedure of Islamic culture, while we only generated	1226
1177	bic QA dataset in which linguistic diversity on the	MCQs for the test split of Western/Chinese cultures.	1227
1178	question side is matched with consistent, unaltered	10% entries from the Western and Chinese subsets	1228
1179	interpretations on the answer side. Such a design	were independently sampled to preserve cultural	1229
1180	ensures that the dataset both preserves the symbolic	balance. This results in 1,675 MCQ samples for	1230
1181	integrity of the source material and provides a re-	western culture and 959 for Chinese.	1231
1182	alistic testing ground for evaluating models under	Each sample’s correct interpretation is paired	1232
1183	conditions of dialectal variation in user queries.	with four distractors randomly drawn from other	1233
1184	<b>Western/Chinese Culture</b> To enhance question	symbols within the same cultural subset. The five	1234
1185	diversity and cultural relevance, we used LLMs	options (A–E) are shuffled, with the correct an-	1235
1186	(e.g., Gemini-1.5-pro) to generate questions.	swer randomly positioned to avoid bias, and a	1236
1187	Given a dream symbol and its associated interpre-	post-processing step ensures uniform answer dis-	1237
1188	tation, the models were prompted with culturally	tribution. Once options are finalized, each entry	1238
1189	tailored templates, guiding them to produce natural	is passed to Gemini-1.5-pro which generates a	1239
1190	questions aligned with the interpretive traditions of	natural-language question contextualized to the re-	1240
1191	the source data. See prompts in Figures 6 and 7 in	spective cultural framework. The model is supplied	1241
1192	Appendix C.	with a system prompt that includes both the correct	1242
1193	For the Chinese culture subset, the prompt de-	interpretation and its distractors, together with a cu-	1243
1194	sign emphasizes classical Chinese dream interpre-	rated set of culturally appropriate MCQ templates	1244
1195	tation perspectives, such as 金木水火土, 阴阳 and	distinct for Western and Chinese dream theories.	1245
1196	命格. Western culture templates highlight West-	These templates guide the model to produce diverse	1246
1197	ern interpretation conventions, interpreting from	and well-phrased questions while maintaining in-	1247
1198	astrological, zodiac. semiotic and psychological	interpretive fidelity. Robustness of the pipeline is	1248
1199	perspectives. Thus, the generated questions are	ensured via timeout protection, automatic retries,	1249
1200	aligned with the interpretive context of the source	and multithreaded execution.	1250
1201	culture while ensuring diversity (multiple distinct	After expansion through distractor pairing and	1251

1252	question generation, the final MCQ dataset contains 2,634 Western samples. All MCQ entries remain traceable to the original dataset through the consistent id key. The final dataset includes the dream symbol, interpretation, multiple-choice options, the correct answer label, and the generated question. Detailed prompt structures for Western and Chinese MCQ generation are presented in Figures 12 and 13 respectively in Appendix D.	1302
1253		1303
1254		1304
1255		1305
1256		1306
1257		1307
1258		1308
1259		1309
1260		1310
1261	<b>B.3.3 Good/Bad and True/False</b>	1311
1262	In addition to free-form questions and straightforward MCQs, users often pose queries such as <i>I dreamed of a snake yesterday, is it a good sign?</i> or <i>I dreamed of a snake, and someone told me it means earning more money, do you think this is true?</i> To better reflect these real-world inquiry scenarios, we introduce two additional tasks: (i) determining whether a dream is a <i>good or bad</i> sign, and (ii) verifying whether a dream’s meaning matches the user’s assumption ( <i>true or false</i> ).	1312
1263		1313
1264		1314
1265		1315
1266		1316
1267		1317
1268		
1269		1318
1270		
1271		1319
1272	For Good/Bad task, we applied Gemini-1.5-Pro to classify each dream-interpretation entry as either good or bad. The input prompt includes the (dream, interpretation) pair, and the model is asked to return the label and brief explanation. This setup encourages the model to reflect on overall sentiment, cultural associations, and emotional tone. For example, a symbol associated with blessings or success would typically be labeled as good, while one linked to fear or misfortune might be labeled as bad.	1320
1273		1321
1274		1322
1275		1323
1276		
1277		1324
1278		1325
1279		1326
1280		1327
1281		
1282		1328
1283		
1284		1329
1285		
1286		1330
1287		
1288		1331
1289		1332
1290		
1291	Note that we only extended these two tasks for Islamic culture test splits to evaluate model performance consistency across different inquiry styles. The same approach can be applied to produce more data.	1333
1292		1334
1293		1335
1294		1336
1295		1337
1296	<b>B.4 Training and Evaluation Splits</b>	1338
1297	We divide the dataset into training and test splits to enable reliable model development and evaluation. Within each cultural subset, the final processed symbol–interpretation pairs serve as the foundation for multiple task formats.	1339
1298		1340
1299		1341
1300		
1301		
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
	<b>C Prompts for QA Generation</b>	1311
	Figure 6 and Figure 7 present the detailed prompt used for curating the dataset for the QA task in Western and Chinese cultures. Gemini-1.5-pro is used to synthesize contextually appropriate questions from the given symbol-interpretation pair and predefined question templates.	1312
		1313
		1314
		1315
		1316
		1317
	<b>D Prompts for MCQ Generation</b>	1318
	In Figure 8, we present the prompt used for categorizing Arabic dream entities into symbolic categories. Figure 9, Figure 12 and Figure 13 are prompts used for generating the MCQ benchmark dataset in three cultures.	1319
		1320
		1321
		1322
		1323
	<b>E Evaluation</b>	1324
	<b>E.1 Models</b>	1325
	We include three commercial LLMs and their reasoning variants:	1326
		1327
	• DeepSeek-v3.1-Terminus	1328
	• GPT5-2025-08-07	1329
	• Claude-Sonnet-4-20250514	1330
	We also included 11 open-source models from 1B to 8B:	1331
		1332
	• Qwen3-8B/1.7B (chat and reasoning)	1333
	• Llama-3.1-8B	1334
	• Deepseek-R1-Distill-Llama-8B	1335
	• Qwen2.5-Math-7B	1336
	• Deepseek-R1-Distill-Qwen-7B	1337
	• Qwen2.5-1.5B-Instruct	1338
	• DeepSeek-R1-Distill-Qwen-1.5B	1339
	• Llama-3.2-1B	1340
	All models used the same prompt templates.	1341

**You are a helpful assistant to interpret dreams from Western culture and theory.** This is a symbol and interpretation pair in English from Western dream interpretation.

**symbol:** {symbol}

**interpretation:** {interpretation}

Use the following template of questions to generate an appropriate question for each symbol–interpretation pair. You may also generate original questions based on the given content. Start the question directly, without adding introductory text. Both the question and the interpretation should explicitly reference Western culture. Each symbol can result in between one and five questions depending on the interpretation. The output must be returned as a list of questions in JSON format without any additional text or explanation.

- I had a strange dream involving {symbol}. Could it be reflecting something deeper going on in my life?
- When {symbol} appears in a dream, could it point to something unresolved within me?
- I've been thinking about a dream I had with {symbol} in it, what might that say about my current state?
- Could dreaming about {symbol} be linked to recent emotions or stress I've been experiencing?
- What kind of message or warning might a dream about {symbol} be trying to send me?
- If I keep dreaming about {symbol}, could it be a sign that my mind is trying to process something?
- I dreamed about {symbol} last night, could this relate to something I haven't been paying attention to?
- In Western dream interpretation, what might the presence of {symbol} suggest about my inner world?
- What personal insight could I gain from a dream that prominently features {symbol}?
- Could {symbol} in my dream represent something I'm avoiding or afraid to confront?
- What does it typically mean if {symbol} shows up in a dream during a time of change?
- If I see {symbol} in a dream, might it reflect my mindset or relationships lately?
- I'm curious whether dreaming of {symbol} might connect to a personal challenge or transition.
- Is it possible that {symbol} showing up in my dream symbolizes a decision or dilemma I'm facing?
- What could my subconscious be working through if I keep seeing {symbol} in my dreams?
- What does seeing {symbol} in a dream mean?
- What does seeing {symbol} in a dream symbolize in Western culture?

Figure 6: Prompt used with Gemini 1.5 Pro to generate culturally grounded questions based on Western dream interpretations.

## E.2 Inference Configuration

For all pretrained and fine-tuned open-source models (including both reasoning and non-reasoning models), we used identical inference settings: temperature = 0.0, greedy decoding only, and a maximum token limit of 4096. For all commercial models, we used the default parameters provided by the API.

## E.3 Prompts for Evaluation

Figure 14, Figure 15, Figure 16 and Figure 17 present the detailed prompt template on the MCQ, Good/Bad classification, True/False classification and free-form QA tasks, respectively. Figure 18 shows the LLM-as-a-Judge prompt used to evaluate the correctness and sentiment of the model response for the QA task.

**You are a helpful assistant to interpret dreams from Chinese culture and traditions.** This is a symbol and interpretation pair in English from Chinese dream interpretation.

**symbol:** {symbol}

**interpretation:** {interpretation}

Use the following template of questions to generate an appropriate question for each symbol–interpretation pair. You may also generate original questions based on the given content. Each symbol can result in between one and five questions depending on the interpretation. The output must be returned as a list of questions in JSON format without any additional text or explanation.

- What does it mean to dream about {symbol}?
- Can you interpret the meaning of {symbol} appearing in my dream?
- I dreamed about {symbol}, what might this symbolize?
- What is the traditional Chinese interpretation of seeing {symbol} in a dream?
- What could be the deeper meaning of {symbol} in a dream?
- Could dreaming of {symbol} be a sign or omen? What does it represent?
- How should I understand the appearance of {symbol} in my dream?
- What are the possible meanings of dreaming of {symbol} repeatedly?
- I saw {symbol} in my dream and it felt significant, what could it mean?
- What does Chinese dream culture say about dreaming of {symbol}?
- My grandmother used to say dreams carry meanings, what could {symbol} mean if seen in a dream?
- In Chinese folk tradition, how is {symbol} interpreted in dreams?
- How would a traditional dream interpreter explain seeing {symbol}?
- Could the dream of {symbol} indicate something about my future or fortune?
- Is there a symbolic or spiritual meaning to dreaming about {symbol}?
- I'm facing stress lately and dreamed of {symbol}. Could it reflect something in my life?
- I had a peaceful dream with {symbol}, does it reflect emotional or spiritual harmony?
- After dreaming of {symbol}, I've felt uneasy. Could it signal a warning or imbalance?
- What might it mean if I see {symbol} in recurring dreams related to family or work?
- Could the appearance of {symbol} in my dream suggest anything about relationships or health?

Figure 7: Prompt used with Gemini 1.5 Pro to generate culturally grounded questions based on Chinese dream interpretations.

**You are an expert in dream interpretation.** You will be given a single dream entity, such as a word or phrase, and your task is to assign it to one of the following high-level categories:

1. Prophets and Messengers
2. Companions, Saints, and Righteous People
3. People and Social Roles
4. Body Parts
5. Animals
6. Birds and Insects
7. Places and Landmarks
8. Natural Elements and Phenomena
9. Tools and Physical Objects
10. Food and Drink
11. Religious Symbols and Practices
12. Emotions and Psychological States
13. Actions and Events
14. Time and Temporal Markers
15. Abstract or Ambiguous Symbols
16. Symbols Related to Death and the Afterlife
17. Uncategorized or Rare Symbols

**Dream Entity:** {entity}  
**Output Format (respond with only the category number):**  
 Category: X

Figure 8: Prompt used with GPT-4o to categorize dream symbols into one of 17 symbolic categories.

**You are an expert of dream interpretation.** Below is the dream symbol between double ticks “symbol” with its correct interpretation:  
 “{symbol}“: {interpretation}  
 And the following are a list of four wrong interpretations of the previous dream symbol:  
 {wrong\_interp}  
 Using this information, write *one* multiple-choice question about the symbol “{symbol}“ using the following rules: • All output should be in Modern Standard Arabic. • Only write one question. • Provide 5 options (A–E), with only one correct answer. • Randomize the position of the correct answer among A–E. • Format your output as valid JSON in the following structure:

```
{
  "question": "...",
  "options": [
    "A) ...",
    "B) ...",
    "C) ...",
    "D) ...",
    "E) ..."
  ],
  "correct_answer": "<correct_answer>"
  // either "A", "B", "C", "D", or "E"
}
```

Figure 9: Generation prompt for Arabic MCQs. The prompt specifies one correct interpretation and four distractors, and enforces a JSON output format in Modern Standard Arabic.

**Review the following question according to five criteria:** 1. Proper formatting 2. Correct language 3. Relevance to the symbol 4. Exactly one correct answer 5. Absence of theologically inappropriate content  
If any issue is found, *fix it* and rewrite the question in the exact format below. If no issue exists, return the question unchanged.

**Context:** {ctx}

**Question:** {s}

• Reply **only** with the question, no comments or explanations. • Use this layout *verbatim*:

Question:

...

A) ...

B) ...

C) ...

D) ...

E) ...

Correct answer: X

Figure 10: Self-critique prompt (English translation). This iterative self-refinement step is inspired by the Self-Refine framework (Madaan et al., 2023).

**Respond only with JSON in the form:** {"pass": true/false, "reasons": ["..."]}

Evaluate the question below based on the same five criteria: 1. Formatting 2. Language 3. Symbol relevance 4. Valid and unambiguous correct answer 5. Theological safety

**Context:** {ctx}

**Question:** {s}

Figure 11: Judging prompt given to GPT-4.1-mini. This design follows the LLM-as-Judge methodology validated by Zheng et al. (Zheng et al., 2023).

**You are a helpful assistant to interpret dreams from Western culture and theory.** This is a set of symbol, interpretation, and options in English from Western dream interpretation.

**symbol:** {symbol}

**interpretation:** {interpretation}

**possible options:** {options}

The interpretation is the correct answer, while the rest of the options are distractors.

Use the following template of questions to generate an appropriate multiple-choice question. You can also create an original question based on the symbol, interpretation, and options:

- Which of the following best describes the meaning of dreaming about {symbol} in Western dream interpretation?
- According to Western dream theory, what does dreaming of {symbol} symbolize?
- In Western dream interpretation, what might it signify if you dream about {symbol}?
- Dreaming about {symbol} in Western culture typically indicates which of the following?
- Based on Western dream interpretation, what is the likely meaning of {symbol} in your dream?
- If I dream of {symbol}, what would it typically symbolize in Western theory?
- What could be the deeper meaning of {symbol} in a dream according to Western culture?
- According to Western dream interpretation, what is the meaning of {symbol} in a dream?

Figure 12: Prompt used with Gemini 1.5 Pro to generate MCQs based on Western dream interpretations.

**You are a helpful assistant to interpret dreams from Chinese culture and traditions.** This is a set of symbol, interpretation, and options in English from Chinese dream interpretation.

**symbol:** {symbol}

**interpretation:** {interpretation}

**possible options:** {options}

The interpretation is the correct answer, while the rest of the options are distractors.

Use the following template of questions to generate an appropriate multiple-choice question. You can also create an original question based on the symbol, interpretation, and options:

- Which of the following best describes the meaning of dreaming about {symbol} in Chinese dream interpretation?
- According to Chinese dream culture, what does dreaming of {symbol} symbolize?
- In Chinese dream interpretation, what might it signify if you dream about {symbol}?
- Dreaming about {symbol} in Chinese culture typically indicates which of the following?
- Based on Chinese dream interpretation, what is the likely meaning of {symbol} in your dream?

Figure 13: Prompt used with Gemini 1.5 Pro to generate MCQs based on Chinese dream interpretations.

**System Prompt**

You are a helpful assistant.

**User Prompt**

**### Question**

{question}

**### Options**

{choices}

=====

Answer the above question by only returning the option (A, B, C, D or E) without any further explanation.

Figure 14: Prompt template for evaluating models on the MCQ task.

**System Prompt**

You are a helpful assistant.

**User Prompt**

If I dream of symbol, is it a good or bad sign? Return "Good" or "Bad" without any further explanation.

Figure 15: Prompt template for evaluating models on the Good/Bad classification task.

```

System Prompt
You are a helpful assistant.

User Prompt
### Question
{question}

### Interpretation
{interpretation}

=====
Given the above question and its interpretation, determine the correctness of the interpretation. Return "True" or "False"
without any further explanation.

```

Figure 16: Prompt template for evaluating models on the True/False classification task.

```

System Prompt
You are a helpful assistant.

User Prompt
{question}

```

Figure 17: Prompt template for evaluating models on the free form QA task.

```

System Prompt
You are a helpful assistant.

User Prompt
### Question
{question}

### Ground Truth Interpretation
{gt_interpretation}

### Model Interpretation
{model_interpretation}

=====
Given the above question, ground truth interpretation and model interpretation.
You have three tasks:
1. Determine if the model interpretation is correct (an integer score between 1 to 5, 1 means entirely incorrect, 5 means
entirely correct).
2. Determine the sentiment of ground truth interpretation (positive, negative).
3. Determine the sentiment of model interpretation (positive, negative).

Return the results in the following JSON format without any further explanation.
{
  "correctness": an integer score between 1 to 5,
  "gt_sentiment": "positive" or "negative",
  "model_sentiment": "positive" or "negative"
}

```

Figure 18: Prompt template for evaluating QA task performance with LLM-as-a-Judge approach where GPT-5-mini serves as the judge.