

# EMODIS: A Benchmark for Context-Dependent Emoji Disambiguation in Large Language Models

Jiacheng Huang<sup>1\*</sup>, Ning Yu<sup>2</sup>, Xiaoyin Yi<sup>2</sup>

<sup>1</sup>School of Artificial Intelligence and Computer Science, Hubei Normal University, Huangshi, China

<sup>2</sup>School of Computer Science and Technology, Chongqing University of Posts and Telecommunications, Chongqing, China  
huangjc@hbnue.edu.cn,  
{d210201034, d200201032}@stu.cqupt.edu.cn

## Abstract

Large language models (LLMs) are increasingly deployed in real-world communication settings, yet their ability to resolve context-dependent ambiguity remains underexplored. In this work, we present EMODIS, a new benchmark for evaluating LLMs' capacity to interpret ambiguous emoji expressions under minimal but contrastive textual contexts. Each instance in EMODIS comprises an ambiguous sentence containing an emoji, two distinct disambiguating contexts that lead to divergent interpretations, and a specific question that requires contextual reasoning. We evaluate both open-source and API-based LLMs, and find that even the strongest models frequently fail to distinguish meanings when only subtle contextual cues are present. Further analysis reveals systematic biases toward dominant interpretations and limited sensitivity to pragmatic contrast. EMODIS provides a rigorous testbed for assessing contextual disambiguation, and highlights the gap in semantic reasoning between humans and LLMs.

**Code** — <https://github.com/JiaCheng-Huang/CODIS>

## Introduction

Large language models (LLMs) (OpenAI 2024; DeepSeek-AI 2025a), which are trained on massive corpus, have demonstrated remarkable performance across a wide range of downstream tasks, such as emotional understanding (Lu et al. 2025), content generation (Şahinuç et al. 2024), and commonsense reasoning (Zong et al. 2025). Since LLMs continue to advance rapidly, comprehensive and rigorous evaluation of their capabilities has become increasingly essential, particularly in understanding nuanced aspects of natural language.

Due to the ambiguity of natural language, expressions with multiple meanings can be easily misunderstood by LLMs and even humans in the absence of sufficient context. Such ambiguity has also been intensified in digital communication where emojis, despite enriching human's modes of expression, often carry multiple interpretations that heavily depend on context, increased potential for misunderstanding. For instance, consider the sentence "She just sent me

a 🍑 last night". In isolation, the emoji 🍑 can evoke multiple interpretations, that is, it may represent a literal peach, or serve as a slang reference to buttocks or sexual innuendo. However, when additional context is provided, the intended meaning becomes significantly clearer. Specifically, when the preceding conversation involves a discussion about summer fruits, the emoji is likely to be interpreted literally. In contrast, if the context includes flirtatious or suggestive exchanges, it tends to be understood figuratively.

Although semantic disambiguation has been a widely studied topic, and several benchmarks have been proposed to evaluate LLMs in this regard, there is currently no benchmark specifically designed to assess the ability of LLMs to resolve emoji-related ambiguities in context-dependent scenarios. Table 1 summarizes recent benchmarks for LLMs. This limitation means existing benchmarks cannot assess the ability of LLMs to understand sentences with emoji in a context-dependent manner.

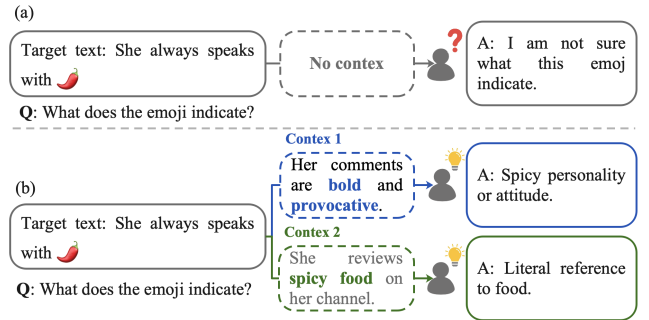


Figure 1: An illustration of our benchmark. Interpretation of sentence with emoji can be significantly influenced by contextual information.

To address this challenge we introduce a new benchmark named EMODIS, which is designed to test the ability of LLMs in EMOji DISambiguation of context-dependent sentences. As shown in Figure 1, we adopt a question-answering format inspired by existing semantic disambiguation benchmarks (Luo et al. 2024) and our benchmark distinguishes itself from prior works in following aspects: first, each target sentence contains an emoji whose meaning is inherently ambiguous without additional context; second, the

\*Corresponding author.

Benchmark	Target Task	Answer Format	Evaluator
DiBiMT(Campolungo et al. 2022)	Word Sense Disambiguation	Open-form text	Metrics
GLADIS(Chen, Varoquaux, and Suchanek 2023)	Acronym Disambiguation	Open-form text	Metrics
ZELDA(Milich and Akbik 2023)	Entity Disambiguation	Span-based linking	Metrics
CVTE-Poly(Zhang et al. 2023)	Polyphone Disambiguation	Multiple-choice	Metrics
AmbiQT(Bhaskar et al. 2023)	Text-to-SQL Parsing	Structured text	GPT
CHAmbi(Zhang et al. 2024a)	Chinese Ambiguity Challenge	Multiple-choice	GPT
<b>EMODIS (Ours)</b>	Emoji Disambiguation	Open-form text	GPT

Table 1: Comparison of proposed EMODIS with recent benchmarks addressing semantic or contextual ambiguity.

questions are intentionally designed to spotlight these ambiguities, requiring additional context for accurate interpretation; third, for every sentence-question pair, EMODIS provides two subtly different contexts, each leading to a distinct interpretation of the emoji. All sentences, contexts, questions, and answers are manually curated to ensure high quality and linguistic diversity. Our evaluation of several widely used large language models on EMODIS reveals a significant gap between model performance and that of humans in resolving emoji ambiguity through contextual cues. Further analysis shows that these models often fail to leverage subtle contextual differences that shift the emoji’s meaning, underscoring the necessity for improved context-aware semantic comprehension in LLMs.

Our contributions are summarized as follows:

- We construct a high-quality dataset of 1000 manually curated instances. Each instance includes an ambiguous emoji-containing sentence, a disambiguation question, and two carefully designed contrastive contexts, each supporting a distinct plausible interpretation.
- We conduct extensive evaluations across both API-based and open-source LLMs, revealing significant gaps in context sensitivity, over-reliance on prior associations, and inconsistent performance across context types.
- We provide detailed analysis such as context sensitivity and interpretation bias, highlighting the core limitations of current LLMs in context-driven semantic interpretation.

## Related Work

### Context in Semantic Ambiguities

Context plays a pivotal role in resolving semantic ambiguity, particularly when the interpretation of a sentence cannot be determined in isolation. Traditional semantic disambiguation tasks have primarily focused on word-level distinctions, such as Word Sense Disambiguation (WSD), where the goal is to assign a correct sense label to a target word given its local context (Zhang and Li 2025; Kruk et al. 2024). While effective in lexical disambiguation, this formulation largely neglects the complexities that arise when ambiguity exists across the entire sentence.

In real-world communication, especially in informal and digital contexts, ambiguity often extends beyond the lexical level. A sentence may contain figurative language, implied

references, or symbols such as emojis, whose interpretation relies not only on syntactic clues but also on the broader conversational or situational context. For example, metaphoric or sarcastic expressions can yield drastically different meanings depending on prior discourse, speaker intention, or cultural knowledge. Such ambiguities cannot be resolved by analyzing isolated tokens, but instead requires models to integrate extra contextual information.

Recent studies emphasize the importance of context-aware modeling across various tasks (Zhang, Song, and Song 2019; Zhao et al. 2024; Luo et al. 2024), showing that model performance in such scenarios declines significantly when deprived of contextual information. However, sentence-level semantic disambiguation remains underrepresented in existing benchmarks, particularly in scenarios involving non-standard linguistic elements like emojis or symbolic expressions.

### Semantic Disambiguation Benchmarks

Existing benchmarks for semantic disambiguation focus on various ambiguity types in natural language. DiBiMT (Campolungo et al. 2022) studies sense ambiguity in machine translation outputs, while GLADIS (Chen, Varoquaux, and Suchanek 2023) targets acronym disambiguation across domains. ZELDA (Milich and Akbik 2023) evaluates zero-shot entity disambiguation using minimal contextual cues. AmbiQT (Bhaskar et al. 2023) focuses on query intent ambiguity in text-to-SQL tasks. In the Chinese context, CHAmbi (Zhang et al. 2024a) provides a fine-grained taxonomy of ambiguous expressions, including figurative and syntactic cases. Other works have explored clarification-based interaction (Zhang et al. 2024b) and semantic obfuscation under adversarial inputs (Xiao et al. 2024), reflecting the growing complexity of ambiguity in modern NLP scenarios.

However, current benchmarks primarily operate in textual settings and rarely consider symbolic elements like emojis, which are prevalent in digital communication and often context-dependent in meaning. These symbolic ambiguities are not well captured by existing resources, leaving a gap in evaluating large language models’ ability to resolve emoji-related ambiguity. To address this, we introduce EMODIS, a benchmark that focuses on disambiguating context-dependent sentences with emoji using contextual cues, providing a new perspective on context-sensitive semantic understanding.

## EMODIS

EMODIS is proposed for evaluating the capability of LLMs in resolving context-dependent sentences with emoji. Inspired by Luo et al. (2024), we adopt a contrastive design where each query appears in two variants, each accompanied by a distinct context leading to different interpretations. In this section, we present the overall task formulation, dataset construction process, evaluation protocol, and data statistics.

### Taxonomy of Context

Given the complexity and diversity of natural language communication, cataloging all possible forms of context that contribute to disambiguating symbolic expressions is inherently challenging. In designing EMODIS as a benchmark for evaluating context-sensitive disambiguation, we focus on how the question posed to the model defines the nature of the required contextual reasoning. Rather than categorizing examples solely by the emoji involved, we categorize them by the type of interpretive challenge that the question sets up, as resolved through contrasting contexts. Specifically, we identify four representative types of disambiguating context, each corresponding to a distinct kind of reasoning required by the question, as illustrated in Figure 2.

**Temporal information.** Interpretation depends closely on the timing and sequence of events described by context. As Figure 2(a) shows, a symbol following a surprising announcement is interpreted as explosive news, whereas in a military discussion it points to a literal weapon. Similarly, Figure 2(b) illustrates how casual conversation about snacks leads to a literal reading, while private exchanges suggest a more intimate or playful tone.

**Domain theme.** Recognizing the topical field implied by context is crucial for resolving ambiguity. In Figure 2(c), reference to policy delay prompts a sarcastic interpretation, while biodiversity context leads to a literal reading about wildlife. Figure 2(d) presents how domain framing distinguishes between figurative and literal meanings depending on whether the discussion involves betrayal or biology.

**Cultural background.** Cultural knowledge provides essential clues for disambiguation, as shown in Figure 2(e), where a symbol evokes spiritual protection in one context, but points to cultural or historical reference in another. Figure 2(f) contrasts a holiday invitation with a showcase of collectibles, leading to different interpretations.

**Social intent.** Disambiguation often relies on identifying tone and interpersonal attitude. Figure 2(g) shows that a symbol may signal playful friendliness in one context, but a patronizing tone in another. In Figure 2(h), the same gesture indicates either sincere acceptance or dismissive response, depending on the social intent.

This taxonomy ensures that EMODIS provides broad coverage of pragmatic reasoning demands, enabling fine-grained evaluation of large language models’ ability to adapt interpretation to contextual cues posed by specific question types.

### Task Formulation

The objective of EMODIS is to test whether LLMs can distinguish between multiple interpretations of an emoji in a sentence, depending on subtle changes in context.

To ensure that models do not guess answers without understanding the context, we structure our dataset in paired examples, denoted as  $(T, C_1, Q)$  and  $(T, C_2, Q)$  formally, where  $T$  refers to target text with emoji,  $C_1$  and  $C_2$  are two pieces of different context, and  $Q$  represents question. We input the paired examples to LLMs and obtain the outputs  $O_1$  and  $O_2$ . Take Figure 2(a) as an example, we provide a target text “She posted a 🍌 after the announcement” and give context that the news was shocking and unexpected in first query. In the second query, we provide different context indicating that the discussion was about military actions in recent conflicts. expected outputs are that the emoji refers to shocking news or literal bomb respectively.

Besides, we apply the following constraints: (1) Each question is neutral and under-informative by design, i.e., it cannot be answered correctly based on the sentence alone. (2) Two contexts  $C_1$  and  $C_2$  are written such that they resolve the emoji’s ambiguity in opposite directions. (3) Each answer is concise and disjoint from the other.

We manually author all sentences, contexts, questions, and answers. The questions are phrased in diverse ways to avoid surface matching, and the contexts are designed to be minimal yet semantically impactful.

### Evaluation Method

Following Fu et al. (2024), we adopt pair-wise accuracy ( $Acc_p$ ) and query-wise accuracy ( $Acc_q$ ) as the evaluation metrics for our EMODIS benchmark. The definitions of the two evaluation metrics are as follows:

$$Acc_p = \frac{1}{N} \sum_{i=1}^N [I(O_{i1}, A_{i1}) \cdot I(O_{i2}, A_{i2})], \quad (1)$$

$$Acc_q = \frac{1}{2N} \sum_{i=1}^N [I(O_{i1}, A_{i1}) + I(O_{i2}, A_{i2})], \quad (2)$$

where  $O_{i1}$  and  $O_{i2}$  are the model’s answers to the  $i$ -th target text under contexts  $C_1$  and  $C_2$ , respectively, and  $N$  is the number of pairs.  $I(a, b)$  is an indicator function defined as:

$$I(a, b) = \begin{cases} 1 & \text{if } a \text{ matches } b, \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

### Data Collection

In this section, we describe the process of constructing the EMODIS benchmark. Our data collection involves three stages designed to ensure both the quality and diversity of the examples while capturing realistic challenges in context-dependent disambiguation.

**Target text authoring.** We begin by manually authoring target text that contain an embedded emoji whose meaning is ambiguous without additional context. Rather than sampling from existing corpora, we manually author target text



Figure 2: Taxonomy of our EMODIS benchmark. For each category, we provide two representative cases. Each case consists of a target sentence, a question, and two contrasting contexts that lead to different answers. Questions (Q), contexts (C), and answers (A) are labeled accordingly to highlight the disambiguation task.

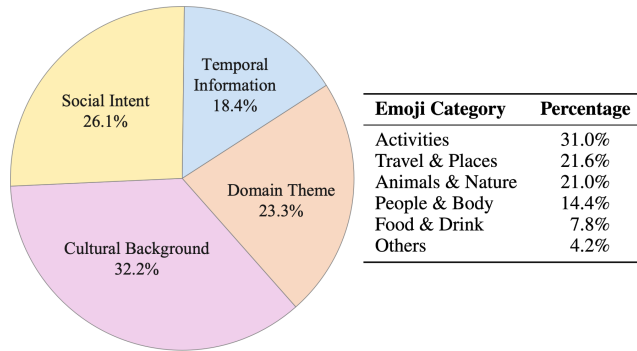


Figure 3: Distribution of context taxonomy (left) and emoji categories (right) in our EMODIS benchmark.

to ensure that each emoji usage presents a genuine semantic ambiguity, where at least two distinct interpretations are plausible. We try to craft natural and diverse text that reflect realistic usage in digital communication across various domains, including social media, messaging, and informal writing. To avoid introducing unintended biases, authors are instructed to refrain from including lexical hints that could resolve the emoji’s meaning without the intended context.

**Context, question and answer authoring.** For each target text, we manually construct two distinct contexts,  $C_1$  and  $C_2$ , where each provides minimal but sufficient information to guide the emoji toward a different, unambiguous interpretation. Contexts are designed to be concise, realistic, and easy to process by both human annotators and models, while avoiding artificial constructions or overly explicit

cues. We encourage diversity in the types of context, including domain background, social intent, cultural references, and temporal or situational settings. Alongside the contexts, we design a disambiguation question  $Q$  aimed at prompting models to focus on the semantic challenge posed by the emoji in context. Each question is phrased carefully to avoid redundancy with the context and to require context-sensitive reasoning, rather than simple lexical matching or frequency-based guessing. For each pair, we provide corresponding answers  $A_1$  and  $A_2$ . These answers are authored to be unambiguous, concise, and semantically precise, reflecting the correct interpretation of the emoji given the context. We ensure that the two answers differ meaningfully and that each aligns naturally with its associated context, so as to present a clear disambiguation signal for evaluation.

**Verification.** Each data instance undergoes independent review by three human annotators. The data are verified with following rules: (1) the emoji is genuinely ambiguous in isolation; (2) the two contexts lead to distinct, clearly identifiable interpretations of the emoji; (3) the question is well-posed and cannot be answered correctly without access to the context; and (4) no unintended clues in the target text or question allow models to bypass the intended reasoning. Instances that do not meet these criteria are revised iteratively or discarded. This multi-stage process ensures that EMODIS forms a robust and reliable benchmark for context-sensitive disambiguation, with high-quality, challenging examples designed for systematic evaluation of large language models.

Ultimately, we obtained a total of 1000 disambiguation instances. The distribution of taxonomy types and emoji categories is illustrated in Figure 3.

## Experiments

We conduct a comprehensive evaluation to test the ability of LLMs to disambiguate emoji meanings based on context. Our experiments are designed to examine not only raw performance but also context sensitivity, semantic generalization, figurative reasoning, interpretive bias, and human-level comparison.

### Experiment Setup

We perform evaluation on several popular LLMs, which are divided into two groups: (1) **API-based models:** GPT-4 (OpenAI 2024), deepseek-v3 (DeepSeek-AI 2025b), Gemini 2.5 (Gemini Team 2025), and Ernie Bot 4.0; (2) **Open-source models:** LLaMA-7B (Ye et al. 2024), Qwen-7B (Bai et al. 2023), Vicuna-7B (Chiang et al. 2023), LLaMA-13B (Liu et al. 2024), and Qwen-14B (Bai et al. 2023). Temperature parameter of all LLMs is set to 0.2.

Following the evaluation setup in Fu et al. (2024), we frame each data instance as a context-sensitive question-answering task. Each instance consists of a target sentence containing an ambiguous emoji, a disambiguating context, and a question that prompts the model to resolve the meaning. Models are instructed to generate a natural language answer without being shown explicit options. We provide the detailed prompts for model inference in Appendix.

To assess whether a model’s answer matches the correct interpretation, we adopt an automatic evaluation strategy using GPT-4 as a verifier. Specifically, GPT-4 is prompted to compare the model’s output with the groundtruth answer and decide whether the two are semantically equivalent. If the response aligns with the correct answer, it is marked as correct; otherwise, it is marked as incorrect. This procedure is applied to both contexts in each instance to determine whether the model can adjust its answer appropriately based on contextual changes. The prompts for GPT-4 evaluation can be found in Appendix.

## Overall Results

We first report the overall performance of LLMs on our EMODIS benchmark, which evaluates their ability to resolve symbolic ambiguity in emoji-laden sentences using minimal but contrastive textual context.

As shown in Table 2, human annotators achieve near-ceiling performance, with  $Acc_p$  and  $Acc_q$  both above 88%, confirming that the task is well-posed and solvable when context is properly integrated. Among all evaluated models, GPT-4 achieves the best overall performance, with a pair-wise accuracy of 58.8% and a query-wise accuracy of 75.2%. However, this still reflects a large gap compared to human-level understanding, indicating that even the strongest LLMs struggle with pragmatic disambiguation in symbolic expressions. API-based models consistently outperform open-source models across all categories. For instance, GPT-4 reaches an overall  $Acc_q$  of 75.2%, while the Qwen-14B trails behind at 58.1%. Notably, open-source models like LLaMA-7B and Vicuna-7B perform significantly worse, with  $Acc_p$  below 30%, suggesting a fundamental deficiency in context-sensitive semantic reasoning.

### Context Sensitivity Analysis

As further indicated in Table 2, we observe a notable gap between  $Acc_p$  and  $Acc_q$  across LLMs, with open-source models exhibiting a larger disparity than API-based ones. Such a gap between query-wise and pair-wise accuracy highlights that many models can occasionally guess correctly but fail to switch answers when the context changes. Therefore, to further investigate the underlying cause of this behavior, two complementary metrics namely context awareness and output variability are examined, capturing a model’s ability to distinguish between different contexts and to generate context-sensitive responses.

Specifically, context awareness measures how often a model provides semantically different outputs when presented with two distinct contextual inputs for the same ambiguous sentence. The context awareness is computed as:

$$Score_{ca} = \frac{1}{n_p} \sum_{i=1}^{n_p} (O_{i1} \neq O_{i2}), \quad (4)$$

where  $n_p$  is the number of pairs, and  $O_{i1}$  and  $O_{i2}$  are the LLM’s outputs for the same ambiguous text under two different contexts. A higher score of context awareness indicates that the model is more sensitive to context changes. On the other hand, output variability evaluates how much the



Model	Temporal information		Domain theme		Cultural background		Social intent		Overall	
	$Acc_p$	$Acc_q$	$Acc_p$	$Acc_q$	$Acc_p$	$Acc_q$	$Acc_p$	$Acc_q$	$Acc_p$	$Acc_q$
GPT-4	64.5	76.6	58.9	74.2	59.3	77.5	54.2	72.3	58.8	75.2
deepseek-v3	36.4	57.6	40.3	60.0	41.6	64.4	27.6	50.6	36.7	58.5
Gemini 2.5	63.1	76.3	48.3	67.7	55.1	74.1	28.5	55.9	48.0	68.2
Ernie Bot 4.0	54.2	73.4	54.3	72.5	56.8	75.2	44.7	66.0	52.5	71.8
LLaMA-7B	19.0	43.8	12.5	35.2	12.4	38.7	11.9	36.4	13.5	38.2
Qwen-7B	28.8	53.0	29.6	55.2	30.4	54.8	17.2	43.3	26.5	51.6
Vicuna-7B	30.9	53.5	27.4	53.6	29.5	53.2	22.6	46.7	27.5	51.7
LLaMA-13B	23.9	54.1	21.9	46.8	22.7	52.6	17.2	46.0	21.3	49.8
Qwen-14B	36.9	60.8	39.9	62.4	38.2	59.9	24.1	50.0	34.7	58.1
Human	88.3	89.2	92.7	94.9	84.9	88.5	89.3	90.4	88.5	90.6

Table 2: Performances of LLMs on EMODIS benchmark.  $Acc_p$  and  $Acc_q$  are reported as percentages.

LLM’s output changes when context is removed. It reflects the LLM’s reliance on contextual information to formulate its responses. The output variability is defined as:

$$Score_{ov} = \frac{1}{n_q} \sum_{i=1}^{n_q} (O_{nc}^i \neq O_c^i), \quad (5)$$

where  $n_q$  is the number of queries for each ambiguous case, and  $O_c^i$  and  $O_{nc}^i$  are the LLM’s responses with and without contextual input respectively.

Model	Context sensitivity	
	$Score_{ca}$	$Score_{ov}$
GPT-4	0.406	0.481
deepseek-v3	0.245	0.317
Gemini 2.5	0.375	0.469
Ernie Bot 4.0	0.396	0.476
LLaMA-7B	0.254	0.264
Qwen-7B	0.195	0.241
Vicuna-7B	0.347	0.325
LLaMA-13B	0.300	0.309
Qwen-14B	0.265	0.284
Human	0.959	0.491

Table 3: Performance of LLMs and humans on context sensitivity metrics. Higher values indicate better ability to distinguish and respond to contextual changes.

As shown in Table 3, human performance is substantially higher than all LLMs on both context awareness and output variability, confirming that current LLMs still struggle to distinguish and respond appropriately to contextual differences. Among the evaluated LLMs, GPT-4 and Ernie Bot 4.0 achieve the highest scores on both metrics, indicating that API-based LLMs are more sensitive to contextual variation than open-source ones. This is consistent with their higher  $Acc_p$  scores in Table 2, further suggesting that the ability to perceive context changes plays a critical role in solving disambiguation tasks.

### Bias in Emoji Interpretation

To further investigate why existing LLMs underperform on EMODIS, we examine whether models tend to rely on default interpretations of emojis instead of using the provided context. Specifically, we analyze output bias by selecting several samples across four common emojis, i.e., 🏮 (red lantern), 🎭 (performing arts), and ⌚ (hourglass), where both literal and figurative meanings are equally represented.

Table 4 shows that LLMs, especially open-source ones, often prefer the more frequent, figurative meaning regardless of context. For instance, GPT-4 demonstrates mild bias, while LLaMA-7B exhibits strong skew toward figurative interpretations. This supports an insight that models often default to learned priors rather than integrating contextual information. Reducing such biases is crucial for improving robustness in context-dependent interpretation.

Model	Red Lantern	Performing Arts	Hourglass
Groundtruth	17:17	11:11	13:13
GPT-4	15:19	9:13	11:15
DeepSeek-v3	13:21	4:18	10:16
Gemini 2.5	14:20	7:15	9:17
Ernie Bot 4.0	19:15	1:21	10:16
LLaMA-7B	2:32	3:19	6:20
Qwen-7B	5:29	5:17	7:19
Vicuna-7B	8:26	4:18	5:21
LLaMA-13B	13:21	5:17	8:18
Qwen-14B	12:22	8:14	9:17

Table 4: Bias in model outputs on three ambiguous emojis. Each cell shows the count of predictions favoring the literal vs. figurative meaning.

### Comparison of GPT-4 and Human Evaluators

To examine the reliability of using GPT-4 as an automatic evaluator for assessing model’s answer, we compare the  $Acc_p$  and agreement between human and GPT-4 on our

EMODIS benchmark. This setup is intended to verify the reliability of GPT-4 as a proxy evaluator, without relying on full human annotation for the entire dataset. Specifically, we define agreement as the proportion of predictions on which both human annotators and GPT-4 give the same evaluation of correctness:

$$\text{Agreement} = \frac{1}{N} \sum_{i=1}^N (\text{Eval}_h(O_i) = \text{Eval}_g(O_i)), \quad (6)$$

where  $N$  refers to instances randomly sampled from the benchmark, and  $\text{Eval}_h(O_i)$  and  $\text{Eval}_g(O_i)$  represent the binary correctness judgments of the model output  $O_i$  given by human annotators and GPT-4, respectively.

As shown in Table 5, GPT-4 achieves over 90% agreement with human judgments across all models, with only minor discrepancies. This demonstrates that GPT-4 is largely consistent with human evaluation in assessing context-sensitive disambiguation, and can thus serve as a reliable and scalable substitute for human annotators in EMODIS evaluation.

Model	GPT-4 $\text{Acc}_p$	Human $\text{Acc}_p$	Agreement(%)
GPT-4	59.6	60.1	98.5
DeepSeek-v3	35.9	34.7	95.2
Gemini 2.5	48.3	47.0	97.0
Ernie Bot 4.0	54.2	52.8	97.8
LLaMA-7B	11.7	12.3	90.5
Qwen-7B	24.4	23.6	93.2
Vicuna-7B	28.2	27.8	94.0
LLaMA-13B	20.6	21.4	91.7
Qwen-14B	32.2	31.7	95.1

Table 5: Comparison between GPT-4 and human evaluation, including pair-wise accuracy from both evaluators and their agreement rate.

## Insights and Discussions

Our analysis of EMODIS reveals three major limitations that large language models face when handling context-dependent emoji disambiguation. These findings not only explain the observed performance gaps but also reflect deeper issues in the models’ ability to integrate and reason over context.

**Insensitivity to Contextual Contrast.** We observe that many models, particularly open-source ones, often produce identical or nearly identical responses when the same sentence is presented under two contrasting contexts. This suggests a phenomenon that the models are not sensitive to the subtle contrast between context variants. In such cases, the model output does not reflect the shift in meaning that a human would naturally infer. Instead of treating context as an active semantic signal that reshapes interpretation, models often regard it as peripheral. This undermines their ability to capture the kind of contrastive reasoning that EMODIS is designed to test.

**Reliance on Prior Associations.** Our emoji-wise bias analysis indicates that models often default to the most frequent or stereotypical interpretation of an emoji, regardless of whether the context supports it. For example, emojis such as “peach” or “snake” are consistently interpreted with their figurative, socially dominant meanings (e.g., flirtation, insult), even when the literal reading is more appropriate given the context. This reveals that model predictions are heavily influenced by associations learned during pretraining, and that context is often insufficient to override these default tendencies. This behavior shows that models do not reliably prioritize contextual cues when determining meaning, and instead fall back on prior distributions.

**Gaps in Pragmatic Reasoning.** While many models can answer straightforward questions correctly, they struggle with cases that involve social cues, implicit tone, or cultural inference. Our case studies and taxonomy-based breakdowns highlight this limitation, especially for emojis whose interpretation depends on sarcasm or implicit judgment. For such examples, model outputs often fail to align with human interpretations, even when the context clearly disambiguates the intended meaning. This suggests that current models lack the reasoning ability necessary for a robust understanding of context-sensitive symbolic language.

## Conclusion

While large language models exhibit impressive capabilities across tasks, their ability to resolve semantic ambiguity in context-sensitive settings remains limited. In this work, we presented EMODIS, a diagnostic benchmark for evaluating large language models on context-sensitive emoji disambiguation. By structuring each example as a sentence-question pair with two contrasting contexts, EMODIS reveals models’ limitations in pragmatic inference and figurative interpretation. Experiments show that even the strongest models often fail to distinguish between literal and figurative meanings when pragmatic inference is required. We further showed that interpretive bias and context neglect are common failure modes, particularly in open models. Although the benchmark is limited in scale and modality, we hope this work provides a new lens for analyzing contextual semantics in LLMs and motivates further research into pragmatic reasoning and symbolic disambiguation.

**Limitation.** There are some limitations in our current work. First, although EMODIS provides diverse textual contexts, each is deliberately kept brief to match current LLM capabilities in leveraging contextual information. Future versions can incorporate more complex and layered contexts to better reflect real-world usage. Second, while we categorize contexts into four representative types, this taxonomy does not exhaust all disambiguation factors such as speaker identity, emotional tone, or multimodal information are not considered. Third, our benchmark currently relies on manual construction and human verification for high quality, which limits scalability. In future work, we aim to expand the dataset using semi-automatic generation and explore additional modalities and context dimensions that may affect emoji interpretation.

## References

- Bai, J.; Bai, S.; Chu, Y.; Cui, Z.; Dang, K.; Deng, X.; Fan, Y.; Ge, W.; Han, Y.; Huang, F.; Hui, B.; Ji, L.; Li, M.; Lin, J.; Lin, R.; Liu, D.; Liu, G.; Lu, C.; Lu, K.; Ma, J.; Men, R.; Ren, X.; Ren, X.; Tan, C.; Tan, S.; Tu, J.; Wang, P.; Wang, S.; Wang, W.; Wu, S.; Xu, B.; Xu, J.; Yang, A.; Yang, H.; Yang, J.; Yang, S.; Yao, Y.; Yu, B.; Yuan, H.; Yuan, Z.; Zhang, J.; Zhang, X.; Zhang, Y.; Zhang, Z.; Zhou, C.; Zhou, J.; Zhou, X.; and Zhu, T. 2023. Qwen Technical Report. arXiv:2309.16609.
- Bhaskar, A.; Tomar, T.; Sathe, A.; and Sarawagi, S. 2023. Benchmarking and Improving Text-to-SQL Generation under Ambiguity. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 7053–7074. Singapore: Association for Computational Linguistics.
- Campolungo, N.; Martelli, F.; Saina, F.; and Navigli, R. 2022. DiBiMT: A Novel Benchmark for Measuring Word Sense Disambiguation Biases in Machine Translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 4331–4352. Dublin, Ireland: Association for Computational Linguistics.
- Chen, L.; Varoquaux, G.; and Suchanek, F. M. 2023. GLADIS: A General and Large Acronym Disambiguation Benchmark. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, 2073–2088. Dubrovnik, Croatia: Association for Computational Linguistics.
- Chiang, W.-L.; Li, Z.; Lin, Z.; Sheng, Y.; Wu, Z.; Zhang, H.; Zheng, L.; Zhuang, S.; Zhuang, Y.; Gonzalez, J. E.; Stoica, I.; and Xing, E. P. 2023. Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%\* ChatGPT Quality. <https://lmsys.org/blog/2023-03-30-vicuna/>. Accessed: 2025-06-10.
- DeepSeek-AI. 2025a. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. arXiv:2501.12948.
- DeepSeek-AI. 2025b. DeepSeek-V3 Technical Report. arXiv:2412.19437.
- Fu, C.; Chen, P.; Shen, Y.; Qin, Y.; Zhang, M.; Lin, X.; Yang, J.; Zheng, X.; Li, K.; Sun, X.; Wu, Y.; and Ji, R. 2024. MME: A Comprehensive Evaluation Benchmark for Multimodal Large Language Models. arXiv:2306.13394.
- Gemini Team, G. 2025. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805.
- Kruk, J.; Marchini, M.; Magu, R.; Ziems, C.; Muchlinski, D.; and Yang, D. 2024. Silent Signals, Loud Impact: LLMs for Word-Sense Disambiguation of Coded Dog Whistles. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 12493–12509. Bangkok, Thailand: Association for Computational Linguistics.
- Liu, H.; Li, C.; Li, Y.; and Lee, Y. J. 2024. Improved Baselines with Visual Instruction Tuning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 26286–26296. Seattle, WA, USA: IEEE.
- Lu, H.; Chen, J.; Liang, F.; Tan, M.; Zeng, R.; and Hu, X. 2025. Understanding Emotional Body Expressions via Large Language Models. In *Association for the Advancement of Artificial Intelligence*, 1447–1455. Philadelphia, PA, USA: AAAI Press.
- Luo, F.; Chen, C.; Wan, Z.; Kang, Z.; Yan, Q.; Li, Y.; Wang, X.; Wang, S.; Wang, Z.; Mi, X.; Li, P.; Ma, N.; Sun, M.; and Liu, Y. 2024. CODIS: Benchmarking Context-dependent Visual Comprehension for Multimodal Large Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 10639–10659. Bangkok, Thailand: Association for Computational Linguistics.
- Milich, M.; and Akbik, A. 2023. ZELDA: A Comprehensive Benchmark for Supervised Entity Disambiguation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, 2061–2072. Dubrovnik, Croatia: Association for Computational Linguistics.
- OpenAI. 2024. GPT-4 Technical Report. arXiv:2303.08774.
- Şahinuç, F.; Kuznetsov, I.; Hou, Y.; and Gurevych, I. 2024. Systematic Task Exploration with LLMs: A Study in Citation Text Generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 4832–4855. Bangkok, Thailand: Association for Computational Linguistics.
- Xiao, Y.; Hu, Y.; Choo, K. T. W.; and Lee, R. K.-W. 2024. ToxiCloakCN: Evaluating Robustness of Offensive Language Detection in Chinese with Cloaking Perturbations. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 6012–6025. Miami, Florida, USA: Association for Computational Linguistics.
- Ye, Q.; Xu, H.; Ye, J.; Yan, M.; Hu, A.; Liu, H.; Qian, Q.; Zhang, J.; and Huang, F. 2024. mPLUG-OwI2: Revolutionizing Multi-modal Large Language Model with Modality Collaboration. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 13040–13051. Seattle, WA, USA: IEEE.
- Zhang, H.; Song, Y.; and Song, Y. 2019. Incorporating Context and External Knowledge for Pronoun Coreference Resolution. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 872–881. Minneapolis, Minnesota: Association for Computational Linguistics.
- Zhang, J.; and Li, X. 2025. Quantum-inspired Non-homologous Representation Constraint Mechanism for Long-tail Senses of Word Sense Disambiguation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 25877–25885. Philadelphia, PA, USA: AAAI Press.
- Zhang, Q.; Cai, S.; Zhao, J.; Pechenizkiy, M.; and Fang, M. 2024a. CHAmbi: A New Benchmark on Chinese Ambiguity Challenges for Large Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, 14883–14898. Miami, Florida, USA: Association for Computational Linguistics.



Zhang, S.; Tan, X.; Lei, Y.; Wang, X.; Zhang, Z.; and Xie, Y. 2023. CVTE-Poly: A New Benchmark for Chinese Polyphone Disambiguation. In *24th Annual Conference of the International Speech Communication Association*, 5526–5530. Dublin, Ireland: ISCA.

Zhang, T.; Qin, P.; Deng, Y.; Huang, C.; Lei, W.; Liu, J.; Jin, D.; Liang, H.; and Chua, T.-S. 2024b. CLAMBER: A Benchmark of Identifying and Clarifying Ambiguous Information Needs in Large Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 10746–10766. Bangkok, Thailand: Association for Computational Linguistics.

Zhao, R.; Zhu, Q.; Xu, H.; Li, J.; Zhou, Y.; He, Y.; and Gui, L. 2024. Large Language Models Fall Short: Understanding Complex Relationships in Detective Narratives. In *Findings of the Association for Computational Linguistics: ACL 2024*, 7618–7638. Bangkok, Thailand: Association for Computational Linguistics.

Zong, D.; Ding, C.; Chen, K.; Li, Y.; and Wang, S. 2025. Counterfactual Debiasing for Physical Audiovisual Commonsense Reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 15265–15273. Philadelphia, PA, USA: AAAI Press.

## Appendix

This appendix provides the full prompt templates used in our experiments for both model inference and answer verification. These prompts are carefully designed to control for stylistic variance, encourage concise responses, and enforce consistent evaluation criteria across different models and settings.

### Prompt for Model Inference

**With Context** To assess LLM’s ability to leverage contextual information for emoji disambiguation, we present the following prompt format during inference:

*Instruction: I’ll give you a text with Emoji and some additional context, which provides information closely related to the text. Please answer my question based on the text and the context. Please answer in a single word or phrase.*

*Context: [Context Here]*

*Sentence: [Target Text Here]*

*Question: [Question Here]*

**Without Context** To measure model behavior in the absence of contextual cues, we use the following variant of the prompt, where the sentence and question are presented without supporting information:

*Instruction: I will provide a sentence that contains an emoji, without any additional context. Please interpret the meaning of the emoji as it appears in the sentence. Answer with a short phrase.*

*Sentence: [Target Text Here]*

*Question: [Question Here]*

### Prompt for GPT-4 Evaluation

To evaluate whether the model’s answer is semantically consistent with the groundtruth, we instruct GPT-4 to serve as a reference-matching verifier. The evaluation prompt ensures that GPT-4 focuses on semantic alignment rather than superficial differences such as formatting or phrasing:

*Instruction: Please evaluate the output of models based on the given question and groundtruth and tell me whether the output is right. The answer is right if it follows the question in meaning and is consistent with the groundtruth. If you think the answer is correct according to the groundtruth, please output “right”, otherwise output “wrong”. You can only print “right” or “wrong” and nothing else. Do not be too strict about the answer. Format different from the groundtruth and minor grammar issues are allowed.*

*Context: [Context Here]*

*Sentence: [Target Text Here]*

*Question: [Question Here]*

*Model’s Answer: [Model’s Output Here]*

*Groundtruth: [Groundtruth of the Question Here]*