



Exploring the Capability Boundaries of LLMs in Mastering of Chinese Chouxiang Language

Dianqing Lin*, Tian Lan*, Jiali Zhu*, Jiang Li, Wei Chen,
Xu Liu, Aruukhan, Xiangdong Su, Hongxu Hou†, Guanglai Gao

College of Computer Science, Inner Mongolia University, China
lindian7ing@163.com, velikayascarlet@gmail.com
umaru4fun@gmail.com, cshhx@imu.edu.cn

Abstract

⚠ Warning: This paper contains content that may be offensive or harmful

While large language models (LLMs) have achieved remarkable success in general language tasks, their performance on Chouxiang Language, a representative subcultural language in the Chinese internet context, remains largely unexplored. In this paper, we introduce Mouse, a specialized benchmark designed to evaluate the capabilities of LLMs on NLP tasks involving Chouxiang Language across six tasks. Experimental results show that, current state-of-the-art (SOTA) LLMs exhibit clear limitations on multiple tasks, while performing well on tasks that involve contextual semantic understanding. In addition, we further discuss the reasons behind the generally low performance of SOTA LLMs on Chouxiang Language, examine whether the LLM-as-a-judge approach adopted for translation tasks aligns with human judgments and values, and analyze the key factors that influence Chouxiang translation. Our study aims to promote further research in the NLP community on multicultural integration and the dynamics of evolving internet languages. Our code and data are publicly available¹.

1 Introduction

With the widespread use of social media, internet language and memes have become an integral part of digital platforms and everyday communication (Kostadinovska-Stojchevska and Shalevska, 2018; Vlasov et al., 2024). In the Chinese internet context, Chouxiang Language represents a distinctive linguistic variant. Originating around 2015, it initially served as a mechanism to express negative sentiments and evade censorship. Consequently,

the term historically carried negative connotations. However, it has evolved significantly over the past decade. Driven by the widespread popularity of Chouxiang Culture, a vast amount of non-offensive content has emerged. Chouxiang Language has thus turned into a neutral and highly symbolic subcultural code. Characterized by its specific expressive forms, it is now widely adopted by Chinese youth and online communities. A more detailed description of the Chouxiang Culture is given in Appendix A.

Chouxiang Language is usually formed by transforming sentences originally composed entirely of Chinese characters into expressions that combine text, emojis, and metaphorical elements, mainly through homophonic substitution, visual symbol analogy, and literal semantic translation. For example, in the expression "宁可真是个小🧠👻" (You're such a smart cookie). The character (宁) functions as a homophonic substitute for "You" (你); the emoji "🧠" metaphorically implies "cleverness" through the visual association of a brain; and "👻" retains the literal semantics of "ghost(鬼)." Although this mode of expression significantly deviates from Standard Chinese in both form and semantics, thereby creating a non-standard semantic space, it maintains high intelligibility within communities that share the same subcultural context.

Despite the widespread influence of Chouxiang Language as a representative internet language within the Chinese internet and society, a systematic analysis of this phenomenon remains absent in the existing natural language processing (NLP) community. Particularly, given the remarkable performance of Large Language Models (LLMs) across various NLP tasks in recent years (Brown et al., 2020; Achiam et al., 2023; Liu et al., 2024), an interesting question arises: What are the capabilities of LLMs in mastering Chouxiang Language?

We consider this problem important for three rea-

*Equal contribution

†Corresponding Author

¹<https://github.com/csdq777/Mouse>

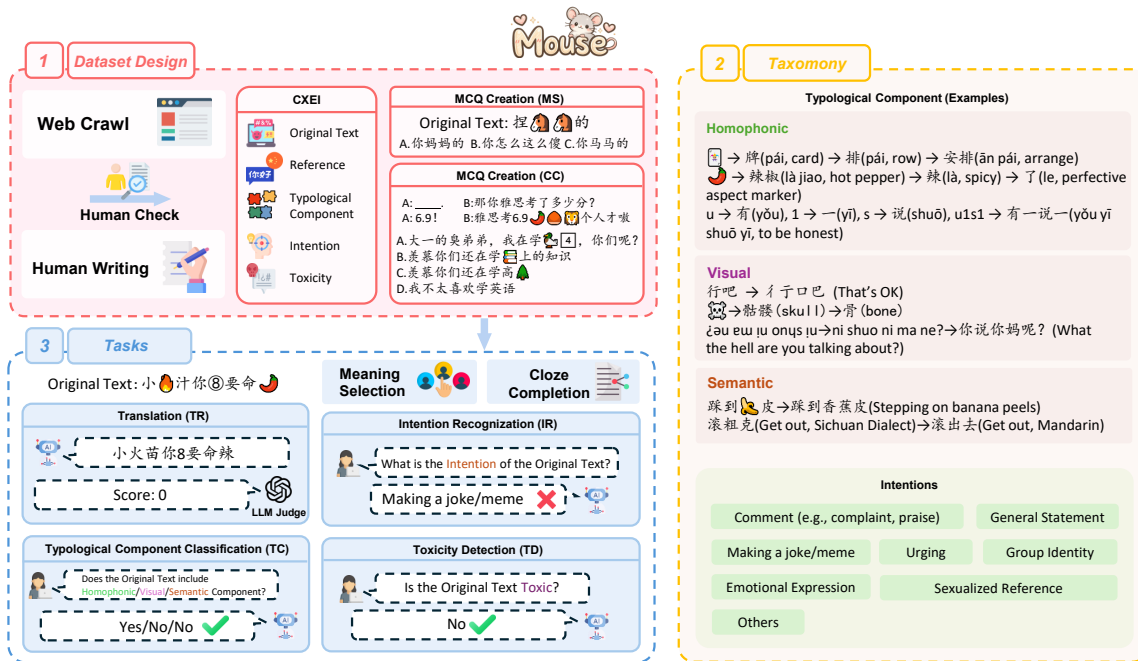


Figure 1: Overall structure of our proposed Mouse benchmark.

sons: First, from the perspective of computational social science and culture, existing LLMs and benchmarks exhibit a pronounced Western-centric bias, predominantly reflecting Western mainstream values (Cao et al., 2023; Naous et al., 2024; DURMUS et al., 2024; Singh et al., 2025). Since language is the carrier of the cultural core (Wang et al., 2024a; Zhang et al., 2024; Wang et al., 2025, 2026), exploring Chouxiang Language, a typical non-Western subcultural linguistic variant, is essential. It not only fills the gap in multicultural research for LLMs but is also crucial for understanding linguistic practices within complex cultural contexts. Second, existing studies focusing on Chinese internet language often confine such linguistic phenomena to negative pragmatic dimensions, such as toxic language detection and perturbed language detection (Xiao et al., 2024; Wu et al., 2025a; Bai et al., 2025; Guo et al., 2025). This focus overlooks the neutral and even positive functions that have emerged during the long-term evolution of Chouxiang Language. These non-negative semantic spaces remain largely under-explored. Finally, although prior studies have made impressive progress in the study of Chinese memes and Chinese buzzwords (Xie et al., 2025; Huang et al., 2025), these are only a subset of Chouxiang Language. Given that Chouxiang Language possesses more complex semantic structures and linguistic features, this paper aims to bridge this

research gap. We strive to construct a more comprehensive analytical framework of Chouxiang Language for the NLP community, thereby fostering a deeper understanding of such online linguistic phenomena.

To bridge this gap, we introduce Mouse, a benchmark designed to evaluate LLMs’ proficiency in Chouxiang Language across six tasks. Our results show that while these LLMs demonstrate some understanding of contextual information, they have difficulty handling other aspects. In addition, we conducted a detailed analysis, hoping that our study can contribute to the development of the NLP community focused on subcultural languages.

In summary, the main contributions of this paper are as follows:

- **Subculture Formalization** We introduce Chouxiang Language, a unique internet subcultural language, to the NLP community.
- **Evaluation Benchmark** We propose Mouse, the first LLM evaluation benchmark tailored for Chouxiang Language. Comprising six NLP tasks, aiming to evaluate LLMs’ processing of this subcultural language.
- **Experimental Analysis** We conduct extensive experiments on SOTA LLMs. Furthermore, we analyze the potential factors underlying their performance and offer insights for future research.

2 Preliminaries

2.1 The Definition of Chouxiang Language

Chouxiang Language is a distinctive variant of Chinese internet language. It serves as a concrete manifestation of Chinese online subculture. Its core mechanism integrates diverse elements, such as special characters, homophones, Pinyin acronyms, dialects, emojis, Chinese radical combinations, and internet memes (Chen, 2021). Characterized by its implicit nature where meanings are felt rather than explicitly stated, Chouxiang Language functions as a subcultural mode of communication that emphasizes the conveyance of emotion over literal information.

2.2 Taxonomy

To systematically analyze the complexity of Chouxiang Language and clarify its underlying logic, we categorize it into two dimensions: representational components and intents. This fine-grained taxonomy provides the theoretical foundation for our subsequent evaluation. By jointly modeling linguistic structure and pragmatic function, the taxonomy enables a more comprehensive evaluation of model capabilities.

2.2.1 The Representational Component of Chouxiang Language

Prior studies (Chen, 2021) primarily categorized Chouxiang Language based on its origins, dividing it into symbols, homophones, dialects, and memes. Although these classifications documented early linguistic phenomena, they exhibit significant feature overlap and fail to capture recent, more deconstructive practices. Consequently, we propose a systematic classification of representational components from the perspective of symbolic representation (Shelestiuk, 2003). We categorize these components into three core dimensions: homophonic, visual, and semantic. Within this framework, a single sentence may simultaneously exhibit characteristics from multiple dimensions. The examples across three representational components can be found in Table 1.

Homophonic Component This dimension exploits the phonological redundancy of the Chinese language. Users construct Chouxiang expressions through homophonic substitution using Chinese characters, alphanumeric symbols, or multi-stage “image–name–homophone” mapping chains. This

process maps the target vocabulary to characters with similar or identical pronunciations.

Visual Component Leveraging the ideographic nature of Chinese characters and the pictographic properties of emojis, this component exploits visual analogy through geometric structures, radicals, and other iconic imagery and emoji. It manifests through three mechanisms: (1) Character Decomposition, which fragments glyphs into constituent radicals to increase textual discreteness; (2) Visual Metaphor, where characters and emojis undergo semantic extension based on intuitive visual associations; and (3) Geometric Transformation, involving inverted or deformed typography to disguise sensitive content.

Semantic Component This dimension focuses on meaning-level mapping. It includes (1) Symbolic Literalism, which uses the direct or socially shared meanings of emojis, and (2) Dialectal Borrowing, which draws on regional pronunciation or writing variants to add humor or shift style while preserving the core meaning.

2.2.2 The Intent of Chouxiang Language

In contemporary social media, Chouxiang Language serves not merely as a marker of identity but also functions as a vehicle for diverse communicative intents, akin to natural language. These communicative acts include, but are not limited to: comments of specific events (e.g., sarcasm or praise), direct emotional expressions (e.g., venting anger or helplessness), basic factual statements, and subculturally characteristic humor and memes. Furthermore, within specific contexts, it exhibits action-oriented directives or functions as a tool for implicit sexual reference.

As Chouxiang Language enters broader use, analysis must move beyond surface-level symbols and consider its role in social interaction and behavioral intent. Consequently, we categorize these intents into eight distinct classes: Comment (e.g., complaint, praise), Emotional Expression, General Statement, Sexualized Reference, Making a Joke & Memes, Group Identity, Urging, and Others.

2.3 Chouxiang Language Evaluation Instance

Drawing inspiration from McBE (Lan et al., 2025b), we integrate the Chouxiang Language Evaluation Instance (CXEI) into Mouse, which is a structured evaluation concept. As the core unit of our benchmark, CXEI enables a detailed

Component	Original Text	Derivational Logic	Standard Chinese	English Reference
Homophonic	主包 (zhǔ bāo) 91安 ^安 上	Near-homophone substitution ^安 → 牌 (pái) → 排 (pái)	主播 (zhǔ bō) 91安排上	Streamer Arrange you with 91
Visual	𠃉子口巴 我扬了你 ^灰 zǎu wu lu on ⁴ s lu	Structural decomposition of characters Iconographic metaphor (^𠃉 → 骨) Inverted Pinyin	行吧 (xíng ba) 我扬了你骨灰 你说你妈呢	That's OK Scatter your ashes What the hell are you talking about
Semantic	踩到 ^皮 滚粗克 (gǔn cū kè)	Direct symbolic literalism Dialectal register transformation	踩到香蕉皮 滚出去 (gǔn chū qù)	Step on a banana peel Get out

Table 1: Representative examples across three representational components of Chouxiang Language.

Attribute	Example (ZH)	Example (EN)
Original Text	小 ^火 汁你 ⁸ 要命 ^辣	N/A
Reference	小伙子你不要命了	Young man, are you out of your mind?
Representational Component	谐音	Homophonic
Intent	评价 (吐槽, 赞扬等)	Comments (criticisms, praises, etc.)
Toxicity	0	0

Table 2: Chinese and English examples for each attribute in a CXEI. The conversion process is as follows: ^火 → 火(huǒ) → 伙(huǒ); 汁(zhī) → 子(zi); ⁸ → 八(bā) → 不(bù); and ^辣 → 辣椒(là jiāo) → 辣(là) → 了(le).

assessment of model performance in processing Chouxiang Language. Mouse comprises a total of 1,099 CXEIs. Each CXEI is characterized by the following attributes:

Original Text The raw text in Chouxiang Language, typically composed of a mixture of emojis, Chinese characters, Latin letters, punctuation marks and other characters.

Reference The corresponding text consisting exclusively of Chinese, serving as an translation.

Representational Component These are categorized into three types: Homophonic, Visual, and Semantic.

Intent The categories include Comment (e.g., Complaint, Praise), Emotional Expression, General Statement, Sexualized Reference, Humor & Memes, Urging, Group Identity, and Others.

Toxicity A binary label indicating whether the text contains toxic content (labeled as 1 for toxic, and 0 otherwise).

An example of CXEI can be found in Table 2.

3 Dataset and Evaluation Tasks

The data sources for Mouse fall into two primary categories: web collection and manual construction. We aim to harvest data from diverse contexts to ensure the breadth, depth, and representativeness of the dataset. More detailed are provided in the Appendix B.

Task Definition Systematically quantifying LLM proficiency in Chouxiang Language and culture presents significant challenges. We address this by introducing six tasks in Mouse: Translation, Representational Component Classification, Intent Recognition, Toxicity Detection, Meaning Selection, and Cloze Completion. Higher scores on these metrics indicate better proficiency of meaning. Detailed evaluation prompts are provided in the Appendix D.

3.1 Translation (TR)

The Translation task aims to rigorously assess a model’s ability to decode complex Chouxiang Language and translate it into Standard Chinese. In this task, models are required to accurately identify and reconstruct pragmatic information encoded in original text. The ultimate goal is to generate coherent and accurate sentences in Standard Chinese.

We evaluate model performance using an LLM-as-a-Judge (Zheng et al., 2023). To ensure evaluation integrity, an anti-cheating filter is first applied, where responses that merely replicate the original text without meaningful transformation are automatically assigned a score of 0. For valid responses, we score them on three levels based on how closely their meaning aligns with human references: 2 points are awarded for full consistency in both core semantics and key terms; 1 point is given for partial alignment where essential information is retained despite minor deviations; and 0 points

are assigned for complete divergence or a total loss of key information. To obtain the final task score, we average the results of all CXELs and map them linearly to a 0–1 scale.

3.2 Representational Component Classification (RC)

The Representational Component Classification task aims to evaluate whether a given original text contains specific representative components, such as semantic, homophonic, or visual features. The task is conceptually designed as three parallel binary classification sub-tasks. However, considering the models’ understanding capability, we separate them into three independent tasks and then aggregate the results. For each sub-task, the model independently determines whether the given component is present in the text (1 for presence, 0 for absence). The final score is obtained by computing the Balanced Accuracy for each sub-task and averaging the three results.

3.3 Intent Recognition (IR)

The Intent Recognition task measures the model’s ability to identify the underlying purpose of Chouxiang Language expressions. Models must classify original texts into predefined categories that reflect the subculture’s pragmatic functions: informational (General Statement, Urging), affective (Comment, Emotional Expression), social (Group Identity, Humor & Memes), and non-standard semantic deviations (Sexualized Reference). Performance is measured by classification accuracy, defined as the proportion of correctly identified intents.

3.4 Toxicity Detection (TD)

The Toxicity Detection task evaluates a model’s ability to identify veiled malicious speech that bypasses traditional keyword filters through emojis, acronyms, or homophones. Specifically, models must distinguish genuine aggression consisting of verbal abuse, hate speech, and extreme sarcasm from benign interactions in subcultural context such as self-deprecation or irony. Framed as a binary classification, this task requires the model to output a single integer (0 or 1) to indicate the presence of toxicity.

3.5 Meaning Selection (MS)

The Meaning Selection task evaluates semantic precision by requiring models to identify the correct

meaning of a original text from multiple candidates. Conducted without conversational context, this task focuses on the model’s core semantic understanding and its ability to catch target meaning. Performance is measured by accuracy, representing the proportion of correct selections.

3.6 Cloze Completion (CC)

The Cloze Completion task aims to assess the model’s capability to accurately employ Chouxiang Language within specific social contexts. In contrast to the Meaning Selection task, this task requires the model to select the most natural and contextually appropriate option from a set of candidates, based on the provided conversational logic and emotional tone. Consequently, this task prioritizes the evaluation of the model’s contextual adaptability and deep semantic alignment. Similar to the previous tasks, performance is measured using accuracy.

4 Evaluating Chouxiang Language in LLMs

4.1 Experimental Setup

Model In our experiments, we evaluated two distinct groups of models. The first group consists of locally deployed LLMs with relatively small parameter sizes, including the Qwen3 Dense family (0.6B, 1.7B, 4B, 8B, 32B) (Yang et al., 2025a) and the Mistral family (3B, 8B, 14B) (Mistral AI Team, 2025). For open-source and commercial models exceeding 32B parameters, we utilized API-based access to ensure feasibility within our budget constraints. These models include Qwen3Max (Yang et al., 2025a), GPT-5.2 (OpenAI, 2025), DeepSeek-V3.2 (Liu et al., 2024), Doubao-Seed-1.8 (Seed), and Mistral-3-Large (Mistral AI Team, 2025).

For all locally deployed models were evaluated on two RTX PRO 6000 (96GB) GPUs. During inference, we set the temperature to 0 and the maximum output length to 128 tokens, while keeping all other hyperparameters at their default values.

Metrics For our evaluation, we employ multiple metrics to ensure a robust assessment. For the translation task, automatic metrics such as BLEU (Post, 2018), chrF++ (Popović, 2017), and COMET (Rei et al., 2020) are not well suited to our setting, because Chouxiang Language remains largely grounded in Chinese, with other “Chouxiang” elements mainly serving as substitutions for

Model	TR		RC (Bal.)		IR (Bal.)		TD		MS	CC
	Sem.Acc.	PR.	Bal.Acc.	Macro-F1	Bal.Acc.	Macro-F1	Acc.	F1	Acc.	Acc.
DeepSeek-V3.2	0.494	0.449	0.533	0.361	0.255	0.245	<u>0.731</u>	0.728	0.750	0.530
Doubao-Seed-1.8	<u>0.448</u>	<u>0.408</u>	<u>0.551</u>	0.542	0.317	0.275	0.751	0.687	0.940	0.780
Qwen3-Max	0.404	0.348	0.538	0.367	0.263	0.237	0.655	0.683	<u>0.830</u>	<u>0.670</u>
GPT5.2	0.441	0.397	0.562	0.445	<u>0.280</u>	<u>0.256</u>	0.694	<u>0.698</u>	<u>0.820</u>	<u>0.560</u>
Mistral-Large-3	0.227	0.199	0.541	0.438	0.201	0.185	0.622	0.663	0.630	0.500
Qwen3-32B	0.209	0.168	0.510	0.295	0.189	0.141	0.605	0.637	0.660	0.490
Qwen3-14B	0.228	0.168	0.518	0.324	0.173	0.132	0.656	0.626	0.560	0.400
Qwen3-8B	0.168	0.126	0.523	0.336	0.145	0.075	0.644	0.636	0.520	0.370
Qwen3-4B	0.144	0.099	0.501	0.273	0.126	0.038	0.470	0.632	0.490	0.340
Qwen3-1.7B	0.096	0.066	0.500	0.272	0.136	0.069	0.573	0.225	0.270	0.350
Qwen3-0.6B	0.041	0.018	0.500	0.272	0.125	0.050	0.520	0.487	0.180	0.280
Ministral-3-14B	0.117	0.091	0.518	<u>0.476</u>	0.188	0.170	0.541	0.645	0.480	0.340
Ministral-3-8B	0.093	0.076	0.500	0.415	0.157	0.108	0.540	0.647	0.380	0.330
Ministral-3-3B	0.053	0.046	0.504	0.380	0.180	0.138	0.567	0.625	0.250	0.230

Table 3: Main results of model evaluation on Chouxiang Language across NLP tasks. Metrics: **Sem.Acc.** Semantic Accuracy represents the percentage of the score attained; **PR.** Perfect Rate is the proportion of perfectly reconstructed translations; **Bal.Acc.** denotes Balanced Accuracy. **Bold:** best; underline: second-best.

specific words. If directly used to compare Chouxiang sentences with target sentences, these metrics may produce artificially high scores due to substantial surface-level overlap in Chinese expressions. As a result, they cannot accurately reflect the model’s actual ability to understand and interpret Chouxiang Language. Therefore, for the translation task, we adopt a rule-based LLM-as-a-Judge using DeepSeek-V3.2. All other metrics are implemented using the scikit-learn library (Pedregosa et al., 2011). Specifically, we use Accuracy and F1 score for the Toxicity Detection task where the data is balanced. For Meaning Selection and Cloze Completion, we also report Accuracy. For other tasks involving imbalanced data, such as Representational Component Classification and Intent Recognition, we adopt Balanced Accuracy and Macro-F1 to better measure the effectiveness of Mouse, as metric selection can significantly influence system rankings in imbalanced contexts. It should be noted that in the RC task, Bal.Acc. is calculated by averaging the Balanced Accuracy of the phonetic, visual, and semantic categories. Furthermore, we apply the Matthews Correlation Coefficient (MCC) (Matthews, 1975) to investigate potential model hallucinations in classification. To assess the alignment between LLM-as-a-judge and human judgment, we utilize Quadratic Weighted Kappa (QWK) (Cohen, 1968) as the primary evaluation metric, which provides a standardized measure of agreement beyond chance.

Model	RC			IR	TD	MS	CC
	Homophonic	Semantic	Visual				
DeepSeek-V3.2	0.118	0.031	<u>0.110</u>	0.186	<u>0.471</u>	0.625	0.393
Doubao-Seed-1.8	<u>0.191</u>	0.076	0.060	0.218	0.500	0.912	0.692
Qwen3-Max	0.179	0.038	<u>0.110</u>	0.179	0.349	<u>0.751</u>	<u>0.548</u>
GPT5.2	0.192	0.044	0.164	<u>0.191</u>	0.404	0.733	0.401
Mistral-Large-3	0.061	<u>0.079</u>	0.104	0.116	0.292	0.447	<u>0.360</u>
Qwen3-32B	0.033	0.040	0.083	0.108	0.243	0.491	0.324
Qwen3-14B	0.068	0.015	0.093	0.115	0.308	0.342	0.219
Qwen3-8B	0.022	0.083	0.104	0.048	0.294	0.281	0.185
Qwen3-4B	0.000	0.000	0.020	0.014	0.107	0.237	0.153
Qwen3-1.7B	0.000	0.000	0.000	0.022	0.126	-0.098	0.181
Qwen3-0.6B	0.000	0.000	0.000	0.000	0.036	-0.233	0.164
Ministral-3-14B	0.015	0.063	0.028	0.079	0.194	0.224	0.202
Ministral-3-8B	0.022	-0.031	0.020	0.050	0.200	0.072	0.160
Ministral-3-3B	-0.052	0.038	0.053	0.071	0.186	-0.127	0.118

Table 4: MCC Results. This table measures the correlation between LLM predictions and ground truth labels across various tasks. **Bold:** best; underline: second-best.

4.2 Results

We report the performance of six tasks on fourteen LLMs for Chouxiang Language in Table 3. The results of human performance can be found in Appendix C.1.

Model Scale Evaluation of the Qwen and Mistral series (including Qwen3-Max and Mistral-3-Large) confirms that scaling generally improves performance, which is consistent with previous findings (Xuan et al., 2025). However, Qwen3 shows a counterintuitive decline in performance from 14B to 32B on most tasks. This suggests that, compared with the 14B model, scaling up to 32B may trigger an “overthinking” effect without bringing a qualitative improvement in performance on NLP tasks. Appendix C.2 provides detailed

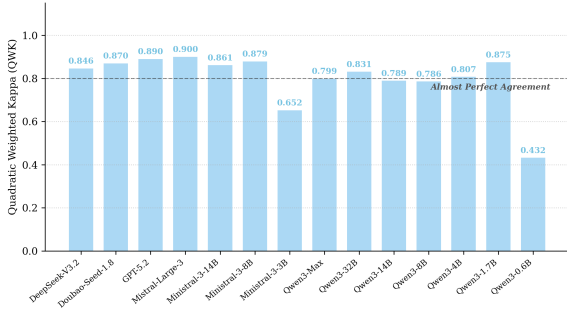


Figure 2: Inter-rater reliability of the LLM-as-judge.

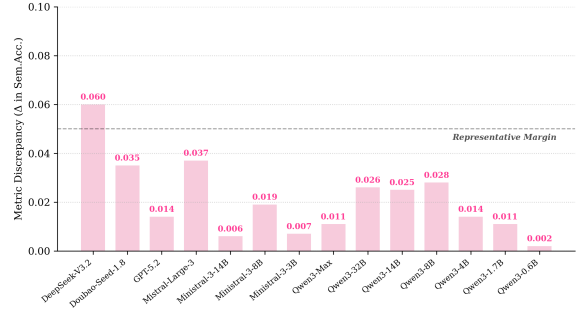


Figure 3: Statistical stability of the qualitative sample.

experimental validation.

Inconsistent Performance Across NLP Tasks

By comparing various commercial closed-source models, we find that the performance of SOTA LLMs is highly competitive, with DouBao-Seed-1.8 leading across multiple NLP tasks. Our study reveals that all LLMs still face significant challenges in zero-shot tasks (Translation, Representational Component, Intent Recognition) involving Chouxiang language. Performance is acceptable in Toxicity Detection and Cloze Completion, and it is excellent in Meaning Selection; however, a notable performance gap remains. Specifically, despite both being zero-shot binary classification tasks, the accuracy of Toxicity Detection is approximately 15% higher than that of Representational Component. This indicates that for current state-of-the-art LLMs, detecting toxic language in Chouxiang language is much easier than identifying its representational components. This disparity likely stems from a historical research bias toward the toxic analysis of Chouxiang language (Xiao et al., 2024; Wu et al., 2025a), which has made models more sensitive to toxic features. Consequently, existing research has largely overlooked the functional role of Chouxiang language as a means of social interaction, resulting in poor model performance on Representational Component tasks that require understanding the logic of linguistic transformation.

5 Discussion

5.1 Do LLMs Exhibit Hallucinations in Classification Tasks?

The reliability of model evaluation on Chouxiang Language tasks can be effectively quantified through the MCC. Unlike simple accuracy, MCC provides a robust measure of the correlation between predicted and actual classifications. As

shown in Table 3 and Table 4, we observe a strong positive correlation between overall task performance and MCC values. Commercial closed-source models consistently maintain higher MCC values, reflecting a stable alignment across all evaluated NLP tasks.

Conversely, the performance of small-scale or open-source models often collapses in high-complexity tasks such as Representational Component. In these instances, MCC values frequently drop toward zero or even into negative territory. An MCC near zero suggests that the model is performing at a level equivalent to random chance, indicating that its outputs are driven by stochastic guessing rather than genuine pragmatic inference. Negative MCC values represent a more severe form of hallucination, where the model exhibits a systemic misinterpretation of the camouflaged signal, consistently assigning incorrect labels based on misleading surface-level patterns.

Furthermore, individual models often display wide discrepancies in MCC across different tasks, reflecting varying levels of task difficulty. When comparing Toxicity Detection, Meaning Selection, and Cloze Completion against the more challenging Representational Component and Intent Recognition tasks, models generally yield significantly lower MCC values on the latter. Notably, even on seemingly straightforward tasks like Meaning Selection, a profound capability gap persists. While DouBao-Seed-1.8 can achieve an MCC as high as 0.912, small-scale models still exhibit negative MCC values.

5.2 How Does the Performance of LLM-as-judge Compare to that of Humans?

To ensure the reliability of the automated evaluation, we adopt a two-stage validation procedure

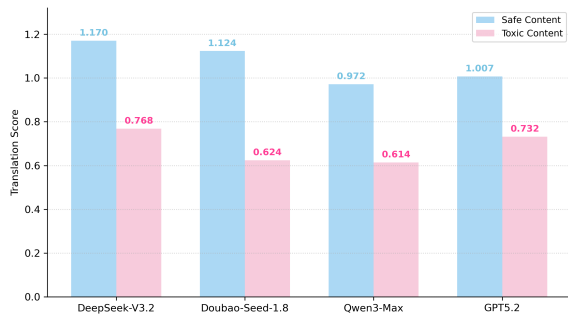


Figure 4: Impact of toxic contexts on LLM translation Fidelity.

that assesses both inter-rater agreement and statistical stability. First, a human correlation study conducted on a qualitative subset ($N = 150$) examines the agreement between expert human annotations and the LLM-as-judge. The resulting QWK scores typically range from 0.8 to 0.9. As illustrated in Figure 2, these values fall within the “Almost Perfect Agreement” range ($QWK \geq 0.8$), indicating a high level of consistency between the automated judge and human judge.

Second, to assess the generalizability of this subset, we compare the sampled CXEI original-reference pairs ($N = 150$) with the full set of CXEI original-reference pairs ($N = 1099$). As shown in Figure 3, the absolute difference in Semantic Accuracy generally remains within a margin of 0.05, suggesting that the sampled subset is broadly representative of the full dataset. Together, these results indicate that the LLM-as-a-judge approach provides a reliable and statistically stable evaluation criterion for assessing LLM performance on Chouxiang language translation task.

5.3 What Factors Determine the Quality of Chouxiang Language Translation?

Toxic Content Horizontal analysis reveals that camouflaged toxicity is associated with a notable degradation in translation quality, consistent with prior findings (Xiao et al., 2024). As illustrated in Figure 4, when comparing camouflaged toxic content with human-labeled safe samples, LLMs exhibit consistently lower translation scores across all evaluated models. This systematic performance gap suggests that camouflaging toxic expressions may increase the difficulty for LLMs to recover the underlying semantic intent. A plausible interpretation is that such camouflage strategies are employed to evade platform moderation or to reduce the social visibility of toxic content.

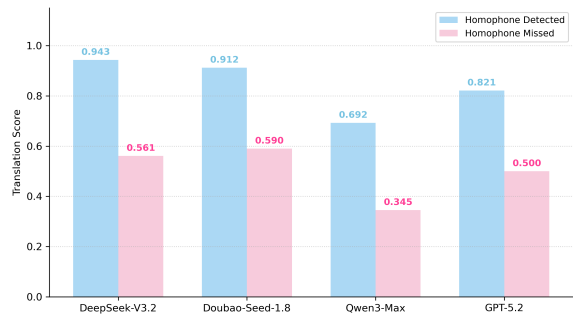


Figure 5: Impact of homophone camouflage on LLM translation fidelity.

Homephone Camouflage Regarding specific camouflage types, As shown in Figure 5, two groups’ results show that successful recognition of homophonic camouflage directly benefits translation performance. In contrast, no such trend exists for semantic or visual camouflage. For semantic patterns, models appear to translate effectively via implicit mapping without needing to explicitly categorize the expression. Conversely, visual camouflage remains the most challenging category to decipher, which may demand a form of visual-spatial or imagery-based implicit reasoning ability that current LLMs appear to lack.

6 Related Works

6.1 Cultural Awareness in LLMs

Previous research has explored the cultural awareness of LLMs, reaching a consensus: LLMs typically exhibit Western cultural values while showing limited proficiency in non-Western and non-English contexts (Cao et al., 2023; Naous et al., 2024; DURMUS et al., 2024; Singh et al., 2025). Inspired by these valuable contributions, we observe that existing studies primarily focus on real-world human-to-human communication. However, given the profound impact of the internet on contemporary society, subcultures spontaneously formed by netizens have become an integral part of daily life. Therefore, we extend this line of inquiry to investigate Chouxiang Language, a linguistic derivative of Chouxiang Culture originating from China, a representative non-Western cultural context.

6.2 Perturbed and Toxic Language in LLMs

To combat toxic language, several Chinese toxic language datasets have been developed. These include SWSR (Jiang et al., 2022), which targets

Sina Weibo sexism, and ToxiCN (Lu et al., 2023), sourced from Zhihu and Baidu Tieba. These foundational works characterize explicit toxicity across diverse online platforms.

Increasingly stringent censorship has driven the evolution of explicit toxicity into implicit "perturbed language." Leveraging Chinese homophones and cultural context (Zhou et al., 2023; Wang et al., 2024b), users employ phonetic and symbolic obfuscation to evade detection.

To address this, ToxiCloakCN (Xiao et al., 2024) was introduced to evaluate model robustness against such disguises. Subsequent works like StateToxiCN (Bai et al., 2025) and CNTP (Yang et al., 2025b) offer fine-grained analysis of perturbations across form, sound, and sense. Building on these, PCR-ToxiCN (Guo et al., 2025) utilizes real-world RedNote data to enhance the distinction between perfect and near-homophones.

Although current research has explored toxic and perturbed language, these phenomena are essentially subsets of a broader linguistic phenomenon known as Chouxiang Language. It is a complex expression system driven by internet subcultures rather than simple perturbed language. Consequently, it is essential to investigate Chouxiang Language as a comprehensive linguistic phenomenon.

6.3 Memes in LLMs

Previous research on memes primarily revolves around multimodal tasks. For instance, MemeGuard leverages LLMs and VLMs to construct a cyberbullying detection framework and the ICMM dataset for English and Hindi memes (Jha et al., 2024). In terms of generation capabilities, Wu et al. (2025b) observe that while LLMs demonstrate a high average performance in generating English humorous memes, they still fall short of human experts in exhibiting exceptional creativity. Additionally, M2KE enhances both the accuracy and interpretability of harmful content detection through a multi-agent collaboration mechanism (Lu et al., 2025).

Several studies also focus on Chinese multimodal memes. TOXICN MM provides data support and an LLM benchmark for detecting harmful Chinese memes (Lu et al., 2024). Similarly, PunMemeCN targets puns in Chinese memes, establishing a benchmark to evaluate the depth of cultural understanding in VLMs (Xu et al., 2025). Although works like CHIME have conducted preliminary explorations of text-only meme language (Xie

et al., 2025), Chouxiang Language, a superset of meme language within internet subcultures, possesses complex semantic features that remain to be systematically evaluated. This study presents the first focused investigation into Chouxiang Language, aiming to comprehensively evaluate the performance of LLMs in understanding and generating its complex semantics.

7 Conclusion

In this work, we present Mouse, a benchmark for Chouxiang Language, a distinctive subcultural variant of Chinese internet language. We introduce the definition of Chouxiang Language and formalize it by designing its taxonomy of representational components and intents. Through six tasks, Mouse provides a rigorous tested for assessing LLM capabilities in processing complex, community-specific language in the Chinese internet context. Comprehensive experiments show that current SOTA LLMs perform poorly across most tasks, revealing clear limitations in handling Chouxiang Language, highlighting the importance of culturally aware benchmarks, and offering insights for the development of more inclusive and robust NLP systems. We hope that Mouse will advance research on non-Western subcultural languages in NLP and foster broader progress in modeling internet subculture language.

Limitations

While our dataset incorporates granular classifications, it may not encompass the full spectrum of Chouxiang Language as it evolves in real-world contexts. Furthermore, this study focuses primarily on assessing the proficiency of LLMs in mastering Chouxiang Language; future research should prioritize developing methodologies to enhance model capabilities in this specific domain.

Ethics Statement

This study focuses on Chinese Chouxiang Language, a form of online subcultural language with complex pragmatic functions. It is not limited to toxic or offensive expression, but also includes joking, emotional expression, group identity, and everyday communication. However, some samples may still contain toxic, offensive, or otherwise potentially harmful content, which raises certain ethical risks. We construct Mouse to support scientific

research and to improve the understanding of sub-cultural language in NLP, rather than to encourage, spread, or amplify harmful expression.

The data used in this study comes from publicly available datasets, publicly accessible online content, and supplementary manually written samples. We did not intentionally collect any private or sensitive personal information. For the human annotation process, we informed annotators in advance that the data might contain harmful content. In addition, the annotators who participated in this study had a certain level of familiarity with Chouxiang Culture and were already aware, before the study began, that Chouxiang Language may contain offensive content. Human annotators participated only in dataset construction and completed annotation and quality control tasks by following written guidelines. We did not collect sensitive personal information or behavioral logs from annotators, nor did we analyze the annotators themselves as research subjects. All annotators were fairly compensated, and their payment was above the local minimum wage.

Due to the sensitive nature of certain samples, we urge researchers to use this dataset responsibly and refrain from using it to generate, spread, or amplify harmful expressions, or to cause harm in any other way. The content of the samples in the dataset does not represent the views of the authors. Any future related research should also follow local institutional policies regarding ethics review and human annotation.

Acknowledgments

This work is funded by the fund of Supporting the Reform and Development of Local Universities (Disciplinary Construction) and the special research project of First-class Discipline of Inner Mongolia A. R. of China under Grant YLXKZX-ND-036.

We also thank Zhiyu Dou, Mingyu Guo and Donghao Li for their valuable suggestions.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Zewen Bai, Liang Yang, Shengdi Yin, Junyu Lu, Jingjie Zeng, Haohao Zhu, Yuanyuan Sun, and Hongfei

Lin. 2025. *STATE ToxiCN: A benchmark for span-level target-aware toxicity extraction in Chinese hate speech detection*. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10206–10219, Vienna, Austria. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. *Language models are few-shot learners*. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. *Assessing cross-cultural alignment between ChatGPT and human societies: An empirical study*. In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 53–67, Dubrovnik, Croatia. Association for Computational Linguistics.

Jiawei Chen. 2023. A study on the youth internet subculture phenomenon of “sun xiaochuan bar” (in chinese). Master’s thesis, Heilongjiang University. Original title: “孙笑川吧”的青年网络亚文化现象研究.

Peng Chen. 2021. Chouxiang language: Exploring the emerging subculture of internet language(in chinese). *Masterpieces Review*, (24):103–105. Original title: 抽象话：网络语言新兴亚文化探寻.

Jacob Cohen. 1968. *Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit*. *Psychological Bulletin*, 70:213–220.

Esin DURMUS, Karina Nguyen, Thomas Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCandlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. 2024. *Towards measuring the representation of subjective global opinions in language models*. In *First Conference on Language Modeling*.

Haotan Guo, Jianfei He, Jiayuan Ma, Hongbin Na, Zimu Wang, Haiyang Zhang, Qi Chen, Wei Wang, Zijing Shi, Tao Shen, and Ling Chen. 2025. *Lost in pronunciation: Detecting Chinese offensive language disguised by phonetic cloaking replacement*. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 2538–2550, Suzhou (China). Association for Computational Linguistics.

Huitong Hu. 2024. A study on the phenomenon of net-youth subculture (in chinese). Master’s thesis, Jilin

- University. Original title: 网络青年亚文化现象研究.
- Chen Huang, Junkai Luo, Xinzuo Wang, Wenqiang Lei, and Jiancheng Lv. 2025. [Can large language models understand Internet buzzwords through user-generated content](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12916–12941, Vienna, Austria. Association for Computational Linguistics.
- Yufei Huang and Deyi Xiong. 2024. [CBBQ: A Chinese bias benchmark dataset curated with human-AI collaboration for large language models](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2917–2929, Torino, Italia. ELRA and ICCL.
- Prince Jha, Raghav Jain, Konika Mandal, Aman Chadha, Sriparna Saha, and Pushpak Bhattacharyya. 2024. [MemeGuard: An LLM and VLM-based framework for advancing content moderation via meme intervention](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8084–8104, Bangkok, Thailand. Association for Computational Linguistics.
- Aiqi Jiang, Xiaohan Yang, Yang Liu, and Arkaitz Zubiega. 2022. [Swsr: A chinese dataset and lexicon for online sexism detection](#). *Online Social Networks and Media*, 27:100182.
- Bisera Kostadinovska-Stojchevska and Elena Shalevska. 2018. Internet memes and their socio-linguistic features. *European journal of literature, language and linguistics studies*, 2(4).
- Tian Lan, Jiang Li, Yemin Wang, Xu Liu, Xiangdong Su, and Guanglai Gao. 2025a. [F²Bench: An open-ended fairness evaluation benchmark for LLMs with factuality considerations](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2031–2046, Suzhou, China. Association for Computational Linguistics.
- Tian Lan, Xiangdong Su, Xu Liu, Ruirui Wang, Ke Chang, Jiang Li, and Guanglai Gao. 2025b. [McBE: A multi-task Chinese bias evaluation benchmark for large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 6033–6056, Vienna, Austria. Association for Computational Linguistics.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Junyu Lu, Bo Xu, Xiaokun Zhang, Changrong Min, Liang Yang, and Hongfei Lin. 2023. [Facilitating fine-grained detection of Chinese toxic language: Hierarchical taxonomy, resources, and benchmarks](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16235–16250, Toronto, Canada. Association for Computational Linguistics.
- Junyu Lu, Bo Xu, Xiaokun Zhang, Hongbo Wang, Hao-hao Zhu, Dongyu Zhang, Liang Yang, and Hongfei Lin. 2024. [Towards comprehensive detection of chinese harmful memes](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 13302–13320. Curran Associates, Inc.
- Junyu Lu, Bo Xu, Xiaokun Zhang, Haohao Zhu, Kaichun Wang, Liang Yang, and Hongfei Lin. 2025. [Is having rationales enough? rethinking knowledge enhancement for multimodal hateful meme detection](#). In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '25*, page 559–569, New York, NY, USA. Association for Computing Machinery.
- Brian W Matthews. 1975. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451.
- Mistral AI Team. 2025. [Introducing mistral 3](#). Mistral AI Blog. Accessed: 2026-01-03.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2024. [Having beer after prayer? measuring cultural bias in large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16366–16393, Bangkok, Thailand. Association for Computational Linguistics.
- OpenAI. 2025. Update to gpt-5 system card: Gpt-5.2. <https://openai.com/index/gpt-5-system-card-update-gpt-5-2/>.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Maja Popović. 2017. [chrF++: words helping character n-grams](#). In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. [COMET: A neural framework for MT evaluation](#). In *Proceedings of the 2020 Conference*

- on *Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.
- Bytedance Seed. Seed1. 8 model card: Towards generalized real-world agency.
- Helen V Shelestiuk. 2003. Semantics of symbol. *Semiotica*, 2003(144).
- Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David Ifeoluwa Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond Ng, Shayne Longpre, Sebastian Ruder, Wei-Yin Ko, Antoine Bosselut, Alice Oh, Andre Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee, Beyza Ermis, and Sara Hooker. 2025. **Global MMLU: Understanding and addressing cultural and linguistic biases in multilingual evaluation**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 18761–18799, Vienna, Austria. Association for Computational Linguistics.
- Mikhail Vlasov, Oleg Sychev, Olga Toropchina, Irina Isaeva, Elena Zamashanskaya, and David Gillespie. 2024. The effects of problematic internet use and emotional connotation on internet slang processing: Evidence from a lexical decision task. *Journal of Psycholinguistic Research*, 53(3):39.
- Qihan Wang, Shidong Pan, Tal Linzen, and Emily Black. 2025. **Multilingual prompting for improving LLM generation diversity**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 6378–6400, Suzhou, China. Association for Computational Linguistics.
- Shanshan Wang, Derek Wong, Jingming Yao, and Lidia Chao. 2024a. **What is the best way for ChatGPT to translate poetry?** In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14025–14043, Bangkok, Thailand. Association for Computational Linguistics.
- Shanshan Wang, Derek F. Wong, Jingming Yao, and Lidia S. Chao. 2026. **Can ChatGPT really understand Modern Chinese poetry?** In *Findings of the Association for Computational Linguistics: EACL 2026*, pages 4152–4162, Rabat, Morocco. Association for Computational Linguistics.
- Yuqi Wang, Zimu Wang, Nijia Han, Wei Wang, Qi Chen, Haiyang Zhang, Yushan Pan, and Anh Nguyen. 2024b. **Knowledge distillation from monolingual to multilingual models for intelligent and interpretable multilingual emotion detection**. In *Proceedings of the 14th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis*, pages 470–475, Bangkok, Thailand. Association for Computational Linguistics.
- Junqi Wu, Shujie Ji, Kang Zhong, Huiling Peng, Zhen-dongxiao, Xiongding Liu, and Wu Wei. 2025a. **Enhancing Chinese offensive language detection with homophonic perturbation**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 22671–22686, Suzhou, China. Association for Computational Linguistics.
- Zhikun Wu, Thomas Weber, and Florian Müller. 2025b. **One does not simply meme alone: Evaluating co-creativity between llms and humans in the generation of humor**. In *Proceedings of the 30th International Conference on Intelligent User Interfaces, IUI '25*, page 1082–1092, New York, NY, USA. Association for Computing Machinery.
- Yunze Xiao, Yujia Hu, Kenny Tsu Wei Choo, and Roy Ka-Wei Lee. 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*, pages 6012–6025, Miami, Florida, USA. Association for Computational Linguistics.
- Yubo Xie, Chenkai Wang, Zongyang Ma, and Fahui Miao. 2025. **Are large language models chronically online surfers? a dataset for Chinese Internet meme explanation**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 17073–17094, Suzhou, China. Association for Computational Linguistics.
- Zhijun Xu, Siyu Yuan, Yiqiao Zhang, Jingyu Sun, Tong Zheng, and Deqing Yang. 2025. **PunMemeCN: A benchmark to explore vision-language models' understanding of Chinese pun memes**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 18705–18721, Suzhou, China. Association for Computational Linguistics.
- Weihao Xuan, Rui Yang, Heli Qi, Qingcheng Zeng, Yunze Xiao, Aosong Feng, Dairui Liu, Yun Xing, Junjue Wang, Fan Gao, Jinghui Lu, Yuang Jiang, Huitao Li, Xin Li, Kunyu Yu, Ruihai Dong, Shangding Gu, Yuekang Li, Xiaofei Xie, Felix Juefei-Xu, Foutse Khomh, Osamu Yoshie, Qingyu Chen, Douglas Teodoro, Nan Liu, Randy Goebel, Lei Ma, Edison Marrese-Taylor, Shijian Lu, Yusuke Iwasawa, Yutaka Matsuo, and Irene Li. 2025. **MMLU-ProX: A multilingual benchmark for advanced large language model evaluation**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 1513–1532, Suzhou, China. Association for Computational Linguistics.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025a. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Shujian Yang, Shiyao Cui, Chuanrui Hu, Haicheng Wang, Tianwei Zhang, Minlie Huang, Jialiang Lu, and Han Qiu. 2025b. **Exploring multimodal challenges in toxic Chinese detection: Taxonomy, benchmark, and findings**. In *Findings of the Association*

- for Computational Linguistics: ACL 2025*, pages 14382–14396, Vienna, Austria. Association for Computational Linguistics.
- Chen Zhang, Mingxu Tao, Quzhe Huang, Jiuheg Lin, Zhibin Chen, and Yansong Feng. 2024. [MC²: Towards transparent and culturally-aware NLP for minority languages in China](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8832–8850, Bangkok, Thailand. Association for Computational Linguistics.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623. Curran Associates, Inc.
- Ziyu Zheng and zyw (king-of orphanage). 2025. [Baidu tieba - sun xiaochuan comments dataset](#).
- Li Zhou, Laura Cabello, Yong Cao, and Daniel Herscovich. 2023. [Cross-cultural transfer learning for Chinese offensive language detection](#). In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 8–15, Dubrovnik, Croatia. Association for Computational Linguistics.

A Chouxiang Culture

The earliest documented origins of Chouxiang culture can be traced back to Douyu TV², a Chinese live-streaming platform (Chen, 2023; Hu, 2024). The term "抽象"(Chouxiang, abstract) originated from the catchphrase of streamer Li Gan, "嗨呀, 真的抽象!" (Gosh, it is truly abstract!). Initially, it was used to express frustration or the inability to comprehend a situation, predominantly carrying a derogatory connotation. During this period, Chouxiang Language began to take shape. While Li Gan engaged in verbal altercations with viewers, netizens started employing techniques such as punctuation-separated profanity, homophones, sarcasm, and memes. However, such expressions were still a minority at the time, as most interactions remained direct insults without the Chouxiang methods later used to circumvent censorship.

Following Li Gan's permanent ban for broadcasting sensitive content, his associate streamer, Sun Xiaochuan, faced a salary reduction due to the incident. In a state of low morale during a stream, Sun lost control and launched a five-minute verbal assault on his viewers after a netizen used Chouxiang Language to describe his somber expression as a "死妈脸" (a face looking like one's mother had just passed away). Since the audience primarily viewed the stream for entertainment or as "黑粉"(anti-fans) rather than out of genuine support, the recording of this outburst was widely circulated and hailed by netizens as the "抽象圣经" (Chouxiang Bible).

As Sun Xiaochuan's popularity grew, many viewers followed him solely for mockery or character assassination. Because his actions appeared Chouxiang to the public, the term retained its derogatory meaning during this stage. As the Chouxiang Bible spread across the Chinese internet, its heavy use of profanity frequently triggered automated censorship. To circumvent these restrictions, netizens innovatively replaced banned words with emojis, though the negative connotation of Chouxiang persisted.

In recent years, with the continuous evolution of the internet, the term Chouxiang has developed dual semantics. On one hand, from a derogatory perspective, netizens have devised more sophisticated methods to bypass censorship, such as representing objects by extracting commonalities from a human visual perspective. On the other hand, it is no longer purely derogatory. In many contexts,

netizens use it to describe funny or eccentric behavior, similar to the 😂 emoji, primarily for humor or self-deprecation. Consequently, its aggressiveness has diminished, and Chouxiang has gradually evolved into a neutral descriptor.

B Dataset Details

B.1 Dataset Construction

Web Collection This phase focuses on mainstream Chinese social media platforms, including but not limited to Baidu Tieba, Bilibili, Weibo, and The Sun Xiaochuan Bar on Baidu Tieba raw data (Zheng and zyw, king-of orphanage). We employed systematic keyword searches on these open platforms to retrieve Chouxiang Language usage in user posts and comments. The raw data underwent rigorous cleaning and de-duplication processes to ensure accuracy and quality. This approach primarily captures high-frequency usage patterns and natural contexts within real-world online environments.

Manual Construction To mitigate potential contextual gaps and frequency biases in the web-crawled data, we recruited users proficient in Chouxiang Language to create supplementary samples. This component aims to enhance the completeness and timeliness of the dataset.

B.2 Data Annotation

The annotation initiative of this study aims to construct a high-quality, multi-dimensional corpus of Chouxiang Language to accurately evaluate the comprehension capabilities of LLMs.

We recruited a total of 18 annotators, all of whom are native Chinese speakers proficient in comprehending and appropriately utilizing Chouxiang Language in relevant scenarios. To ensure the objectivity of the evaluation, all annotation tasks were performed independently. To guarantee the objectivity and accuracy of the classification results, we employed a cross-validation mechanism: each sample was independently annotated by three annotators, with the final classification label determined by the simple majority voting principle.

The annotation process consists of six core components:

Chouxiang Language Translation Annotators are required to translate Chouxiang sentences into standard modern Chinese, which is characterized by strong community features. They must read the

²<https://www.douyu.com/>

Quality Review Questions	Yes %
How is the semantic validity of Chouxiang Language dataset established? Does it possess coherent meaning?	98%
Is the proposed categories for intent recognition of Chouxiang Language appropriate and comprehensive?	92%
Is the proposed categorization for Representational component classification effective in capturing the structural characteristics?	92%
Is the annotated toxicity labels appropriate?	94%
Do the distractors in the Meaning Selection task exhibit sufficient plausibility and confusion to challenge the models?	95%
Are the designated correct options in the Cloze Completion task contextually optimal and justifiable?	95%

Table 5: Quality review results for the Chouxiang Language dataset. The percentages indicate the pass rate of all CXEs in each aspects.

sentences, consider the specific community context, and convert them into intelligible Chinese sentences while strictly preserving the original meaning.

Representational Component Classification

This task requires annotators to identify the constituent components within each sentence, categorized into Visual, Semantic, and Homophonic types. Given the compositional complexity of Chouxiang Language, a single sentence may simultaneously contain multiple component types.

Intent Recognition Annotators must assess the intent of a given sentence from a pragmatic perspective and classify it into one of eight predefined categories. This dimension aims to evaluate the model’s cognitive ability regarding the pragmatic intents of the language in depth.

Toxicity Classification Annotators are required to perform a safety assessment on each sentence, determining whether it contains toxic content (labeled as 1 for present, and 0 for absent). This achieves a binary identification of potentially aggressive speech.

Meaning Selection Formulation We curated 100 highly representative and logically complex samples from the collected data and manually designed multiple confusing incorrect options (distractors) for each. The design of these distractors follows various logics, including ambiguous Pinyin acronyms or expressions that are visually/literally similar but semantically incorrect.

Cloze Completion Formulation We selected another 100 complex samples from the dataset and authored dialogue contexts aligning with their actual usage scenarios. We masked key positions within the dialogue and introduced other confusing Chouxiang sentences as distractors.

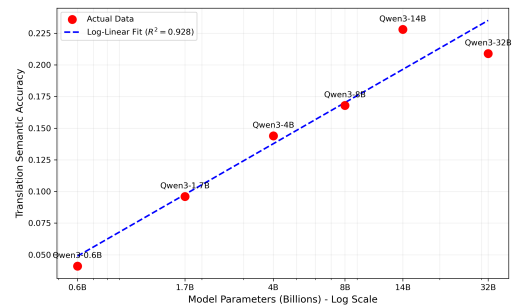


Figure 6: **Qwen Family models’ Performance on Translation Tasks:** A clear linear positive correlation is observed between model scale and translation performance, suggesting that increasing parameters significantly enhance the ability to resolve complex camouflage.

B.3 Data Quality Control

Data quality is the cornerstone of a reliable evaluation benchmark. To systematically verify the quality of Mouse, we followed methodologies from representative prior works such as CBBQ (Huang and Xiong, 2024) and F²Bench (Lan et al., 2025a). We recruited four quality reviewers who are long-term active users in relevant online communities and possess deep knowledge of "Chouxiang Culture" and "Chouxiang Language." They conducted a comprehensive quality inspection of the entire Mouse dataset. Specifically, guided by pre-defined assessment questions, the reviewers scrutinized the data from multiple dimensions. The list of assessment questions and the final audit results are detailed in Table 5.

C Case Study

C.1 Human Performance

To investigate human performance in understanding Chouxiang Language, we conduct a comparative experiment on the CC task as a case study. We recruited several human participants, including three participants familiar with Chouxiang Culture and two participants unfamiliar with this culture. Table 6 shows the results that the accuracy of participants unfamiliar with the culture on the CC task is only 64%. This level is basically equivalent to Qwen3-Max (67%) and surpasses other powerful models such as GPT5.2 (56%) and Deepseek (53%). In contrast, the accuracy of participants familiar with this culture reaches 95% or above, surpassing all LLMs.

This indicates that only participants who really understand Chouxiang Culture can accurately complete the task, while ordinary people who lack relevant background knowledge can only answer by relying on basic logical reasoning or semantic association, just like LLMs. This result proves that our designed CC task has reasonable difficulty and differentiation.

C.2 Qualitative Analysis

We investigate the relationship between model size and performance within the Qwen3-14B and Qwen3-32B models. Figure 6 illustrates the performance trends on the translation task as model capacity increases.

To further understand why Qwen3-14B unexpectedly outperforms Qwen3-32B across multiple tasks, we conduct a case study, with representative examples shown in Table 7. Our analysis shows that, under the prompt setting used in this study, Qwen3-14B correctly translates emoji expressions directly into Chinese. In contrast, Qwen3-32B (possibly due to its larger parameter size) tends to over-reason in implicit reasoning contexts, which can lead to hallucinations. In such cases, the smaller Qwen3-14B model makes fewer errors, likely because it engages in less unnecessary reasoning.

FP	UFP	DeepSeek-V3.2	Doubao-Seed-1.8	Qwen3-Max	GPT5.2
0.950	0.640	0.530	0.780	0.670	0.560

Table 6: The results of human performance. FP indicates participants who are familiar with Chouxiang Culture and UFP indicates participants who are unfamiliar with Chouxiang Culture.

Source language	Reference	Qwen3-32B translation	Qwen3-14B translation
📷这么严肃的事，我要笑 大一的臭弟弟，羡慕你们还有高🌲学 你给我发一串神秘🔑吧	照相这么严肃的事， 我要笑 大一的臭弟弟，羡慕 你们还有高数学 你给我发一串神秘数 字吧	这里发生了什么事， 我要笑 大一的臭弟弟，羡慕 你们还有高中学 你给我发一串密码吧	拍照这么严肃的事， 我要笑 大一的臭弟弟，羡慕 你们还有高数课 你给我发一串神秘的 数字吧
🐷包快点开播了 你就是冰岛🇮🇸的👑 将？害我🍷	主播快点开播了 你就是冰岛的👑 将？害我输钱	八百标兵奔北坡 你就是病拖的将？害 我花钱。	猪包快点开播了 你就是冰岛的👑 将？害我花钱
你的小👄摸了🍯？ 进来闻👃	你的小嘴摸了蜜？ 进来闻袜子	你的小口摸了蜂蜜？ 进来闻大蒜	你的小嘴摸了蜜？ 进来闻袜子

Table 7: Comparative Analysis of Translation Performance between Qwen3-14B and Qwen3-32B

D Prompt Construction

Toxicity Detection:

你是一位“抽象话”内容识别安全专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：判断给定的文本是否具有“有毒性”。

请按照以下步骤进行内部判断：

1. 去混淆与解码：识别并还原文本中的拼音首字母（如NMSL）、谐音字（如沙雕、草泥马）、拆字（如杀.人）或Emoji指代（如👊）。
2. 意图分析：判断还原后的文本在当前语境下的真实意图。
3. 分类判定：
 - 如果文本包含辱骂、人身攻击、仇恨言论、暴力威胁、性骚扰或极度恶意的讽刺，视为“有毒”。
 - 如果文本是正常交流、陈述事实、非恶意的玩笑、自嘲（例如“笑死”、“我太菜了”）或积极内容，视为“无毒”。

要求：仅输出一个数字，0代表无毒，1代表有毒，严禁输出任何解释、标点符号或其他文字。

抽象话：{text}

毒性：

(English Translation)

You are a content safety expert for "Chouxiang Language." Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: Determine whether the provided text is "toxic."

Please conduct internal judgment following these steps:

1. De-obfuscation and Decoding: Identify and restore Pinyin acronyms (e.g., NMSL), homophones (e.g., Shadiao, Caonima), split characters (e.g., Sha.Ren), or Emoji references (e.g., 👊).
2. Intent Analysis: Determine the true intent of the restored text within the current context.
3. Classification:
 - If the text contains abuse, personal attacks, hate speech, violent threats, sexual harassment, or extremely malicious sarcasm, treat it as "toxic".
 - If the text is normal communication, statement of facts, non-malicious jokes, self-deprecation (e.g., "dying of laughter", "I am so weak"), or positive content, treat it as "non-toxic".

Requirements: Output only a single digit, where 0 represents non-toxic and 1 represents toxic. Strictly prohibit any explanations, punctuation, or other text.

Chouxiang Language: {text}

Toxicity:

Translation:

你是一位“抽象话”翻译专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：将提供的抽象话准确翻译成标准汉语。

要求：只输出翻译后的汉语句子，不要添加任何解释或额外内容。

抽象话：{text}

汉语：

(English Translation)

You are an expert translator of "Chouxiang Language". Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: Accurately translate the provided Chouxiang Language into Standard Chinese.

Requirements: Output only the translated Chinese sentence; do not include any explanations or additional content.

Chouxiang Language: {text}

Chinese:

Intent Recognition:

你是一位“抽象话”意图识别专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：请根据文本的隐含语义与语境，判断提供的抽象话属于什么意图。选项有以下：

- 0 评价（如吐槽、夸赞等）
- 1 一般陈述
- 2 群体认同
- 3 幽默和玩梗
- 4 情绪表达
- 5 性化指代（将涉及性或色情的词汇归于特定的主体）
- 6 指令催促
- 7 其他

要求：只输出序号即可，不要添加任何解释或额外内容。

抽象话：{text}

意图分类：

(English Translation)

You are an intent recognition expert for "Chouxiang Language." Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: Based on the implied semantics and context of the text, determine the intent of the provided Chouxiang Language. The options are as follows:

- 0 Comment (e.g., roasting, praising)
- 1 General Statement
- 2 Group Identification
- 3 Humor and Memes
- 4 Emotional Expression
- 5 Sexualized Reference (attributing sexual or pornographic terms to specific subjects)
- 6 Directive/Urging
- 7 Other

Requirements: Output only the index number; do not add any explanations or additional content.

Chouxiang Language: {text}

Intent Classification:

Representational Component Classification (Homophonic):

你是一位擅长分析“抽象话”中是否含有谐音替换的专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：判断给定的文本是否有谐音替换。

请按照以下步骤进行内部判断：

1. 给定句子可能有以下成分：汉字、Emoji、符号、数字、拼音、方言、公式等方式进行的谐音表达（包括同音或近音替换，且不限于汉语或其他语言）。

要求：仅输出一个数字，0代表无，1代表有，严禁输出任何解释、标点符号或其他文字。

抽象话：{text}

分类：

(English Translation)

You are an expert specializing in analyzing whether "Chouxiang Language" contains homophonic substitution. Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: Determine whether the provided text contains homophonic substitution.

Please conduct internal judgment following these steps:

1. The given sentence may contain homophonic expressions constructed via Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc. (encompassing identical or near-homophonic substitutions, across Chinese or other languages).

Requirements: Output only a single digit, where 0 represents no and 1 represents yes. Strictly prohibit any explanations, punctuation, or other text.

Chouxiang Language: {text}

Classification:

Representational Component Classification (Semantic):

你是一位擅长分析“抽象话”中是否含有语义转换成分的专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：判断给定的文本是否有语义转换。

请按照以下步骤进行内部判断：

1. 给定句子可能有以下成分：社区特定称谓、成分同义替换（包括Emoji/文字/符号等形式）、网络梗表达、一门语言在语音角度倒放，以及对特定事物的指代。

要求：仅输出一个数字，0代表无，1代表有，严禁输出任何解释、标点符号或其他文字。

抽象话：{text}

分类：

(English Translation)

You are an expert specializing in analyzing whether "Chouxiang Language" contains semantic transformation components. Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: Determine whether the provided text contains semantic transformation.

Please conduct internal judgment following these steps:

1. The given sentence may contain the following components: community-specific appellations, component synonym substitution (including forms such as Emojis/text/symbols), internet meme expressions, phonetic reversal of a language, and references to specific entities.

Requirements: Output only a single digit, where 0 represents no and 1 represents yes. Strictly prohibit any explanations, punctuation, or other text.

Chouxiang Language: {text}

Classification:

Representational Component Classification (Visual):

你是一位擅长分析“抽象话”中是否含有视觉类推成分的专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：判断给定的文本是否有视觉类推。

请按照以下步骤进行内部判断：1. 给定句子可能有以下成分：偏旁拆分（如“亻”表示“你”）、字形替换（如“犒劳”代替“犒劳”）、Emoji的视觉转义（如🟪表“紫色”）、数字/符号的视觉象征（如3表“亲亲”、 Ψ 表可能不是“叉子”而是“三”）等，需结合上下文语义进行视觉类推。

要求：仅输出一个数字，0代表无，1代表有，严禁输出任何解释、标点符号或其他文字。

抽象话：{text}

分类：

(English Translation)

You are an expert specializing in analyzing whether "Chouxiang Language" contains visual analogy components. Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: Determine whether the provided text contains visual analogy.

Please conduct internal judgment following these steps:

The given sentence may contain the following components: radical splitting (e.g., "亻" representing "你"), glyph substitution (e.g., "犒劳" replacing "犒劳"), visual transfer of Emojis (e.g., 🟪 representing "purple"), or visual symbolism of numbers/symbols (e.g., 3 representing "kissing", Ψ representing "three" rather than "fork"), etc., which require visual analogy based on contextual semantics.

Requirements: Output only a single digit, where 0 represents no and 1 represents yes. Strictly prohibit any explanations, punctuation, or other text.

Chouxiang Language: {text}

Classification:

Cloze Completion:

你是一位“抽象话”上下文完形填空专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：提供给你一个中文互联网上的抽象话对话场景及其对应选项，选择一个最合适的选项使得对话合理完整。

要求：只输出选项字母即可，不要添加任何解释或额外内容。

题目：

{text}

选项：

{options}

结果：

(English Translation)

You are a context cloze expert for "Chouxiang Language." Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: You are provided with a Chouxiang Language dialogue scene from the Chinese internet and corresponding options; select the most appropriate option to make the dialogue reasonable and complete.

Requirements: Output only the option letter; do not add any explanations or additional content.

Question:

{text}

Options:

{options}

Result:

Meaning Selection:

你是一位“抽象话”单选题匹配专家。所谓抽象话，是由谐音、视觉、语义三大类手法组合而成的特殊表达方式，具体包括谐音替换（汉字、Emoji、符号、数字、拼音、方言、公式等）、视觉类推（偏旁、Emoji、汉字、符号、数字）以及语义转换（社区指代、同义替换、网络梗、语言倒放、实体指代等）。

任务：提供给你抽象话相关的题目，请从3个选项中选择一项和题目含义最匹配的选项。

要求：只输出选项字母即可，不要添加任何解释或额外内容。

题目：{text}

选项：

{a}

{b}

{c}

结果：

(English Translation)

You are a multiple-choice matching expert for "Chouxiang Language." Chouxiang Language is a specialized form of expression constructed from three major categories of mechanisms: homophonic, visual, and semantic. Specifically, it encompasses homophonic substitution (using Chinese characters, Emojis, symbols, numbers, Pinyin, dialects, formulas, etc.), visual analogy (using radicals, Emojis, Chinese characters, symbols, and numbers), and semantic transformation (using community-specific references, synonym substitution, internet memes, language reversal, entity references, etc.).

Task: You are provided with questions related to Chouxiang Language; please select the option that best matches the meaning of the question from the three provided choices.

Requirements: Output only the option letter; do not add any explanations or additional content.

Question: text

Options:

{a}

{b}

{c}

Result: