

# Beyond Translation: Cross-Cultural Meme Transcreation with Vision-Language Models

Anonymous ACL submission

## Abstract

Memes are a pervasive form of online communication, yet their cultural specificity poses significant challenges for cross-cultural adaptation. We study *cross-cultural meme transcreation*, a multimodal generation task that aims to preserve communicative intent and humor while adapting culture-specific references. We propose a hybrid transcreation framework based on vision–language models and introduce a large-scale bidirectional dataset of Chinese and US memes. Using both human judgments and automated evaluation, we analyze 6,315 meme pairs and assess transcreation quality across cultural directions. Our results show that current vision–language models can perform cross-cultural meme transcreation to a limited extent, but exhibit clear directional asymmetries: US→Chinese transcreation consistently achieves higher quality than Chinese→US. We further identify which aspects of humor and visual–textual design transfer across cultures and which remain challenging, and propose an evaluation framework for assessing cross-cultural multimodal generation. Our code and dataset are publicly available at <https://anonymous.4open.science/r/MemeXGen/>.

## 1 Introduction

Memes are a dominant form of online communication, yet they are difficult to adapt across cultural contexts. A meme that resonates with US audiences may fail among Chinese users—not due to linguistic errors, but because its humor, symbolism, or visual style does not translate culturally. While literal translation preserves surface meaning, it often fails to preserve what makes a meme effective: its intent, humor, and cultural resonance.

This challenge is often described as *transcreation* (Khanuja et al., 2024b): the process of adapting content across cultures by preserving communicative intent rather than literal form. For memes, transcreation requires coordinated adaptation of



### Cultural Differences in Meme Preferences

#### Chinese Memes

##### Visual Characteristics:

Animals  
Comic elements  
Symbolic imagery

##### Textual Style:

Concise even abstract  
Emotional vocabulary  
Philosophical undertones

#### US Memes

##### Visual Characteristics:

Human figures  
Famous cartoon characters  
Celebrity reactions

##### Textual Style:

Text-heavy formats  
Narrative structures  
Conversational tone

Figure 1: Examples of cultural differences in meme preferences across Chinese and US contexts. Cultural preferences shape humor and visual style, creating challenges for cross-cultural meme transcreation.

Aspect	US	Chinese
Visual	Human, celebrity	Animal, symbolic
Text	Narrative, detailed	Concise, philosophical
Emotion	Situational	Universal
Humor	Sarcasm, relatable	Wordplay, cuteness

Table 1: Cross-Cultural Meme Characteristics

both text and images, grounded in culturally specific norms, references, and aesthetic preferences. As such, *meme transcreation poses a fundamental multimodal and cultural generation challenge that goes beyond standard translation or captioning.*

Prior work has primarily studied memes from a recognition and analysis perspective (Hazman et al., 2025; Cao et al., 2023; Zhao et al., 2025). In parallel, recent vision–language models demonstrate strong multimodal understanding capabilities. However, systematic frameworks, datasets, and evaluations for *cross-cultural meme generation* remain limited. In particular, it is unclear how well current models can perform culturally grounded

transcreation, how performance varies across cultural directions, and how such systems should be evaluated.

To address this gap, we empirically study cross-cultural meme transcreation between Chinese and US cultures through a hybrid framework that adapts memes while preserving communicative intent. We introduce a large-scale bidirectional dataset of original and transcreated memes and evaluate transcreation quality using both human judgments and automated evaluation. Our analysis highlights systematic directional effects and cultural factors that shape the success and limitations of current vision-language models in cross-cultural meme transcreation. Table 1 summarizes key cross-cultural distinctions considered in this work.

This paper addresses three research questions:

**RQ1:** How effectively can vision-language models perform cross-cultural meme transcreation while preserving intent, humor, and cultural nuance?

**RQ2:** Does transcreation direction introduce systematic asymmetries between US→Chinese and Chinese→US adaptations?

**RQ3:** How should multimodal meme transcreation be evaluated, and how do human judgments compare with automated evaluations?

**Our Contributions.** **First, we propose a hybrid framework for meme transcreation** that balances intent preservation with cultural adaptation, offering practical guidance for culturally grounded meme adaptation. **Second, we introduce the first bidirectional meme transcreation dataset**, containing 6,315 original memes and 6,315 corresponding transcreated memes across both Chinese→US and US→Chinese directions. **Third, we present a bidirectional empirical study of Chinese and US meme transcreation**, showing that transcreation quality depends on direction and that specific aspects of internet humor—such as imagery, text style, and emotional expression—transfer unevenly across cultures.

## 2 Related Work

**Cultural Gaps in AI.** Despite advances in large language models and vision-language models (VLMs), cultural gaps remain a persistent challenge (Mihalcea et al., 2025; Adilazuarda et al., 2024). Prior work shows that NLP systems often

fail to account for cross-cultural variation (Herscovich et al., 2022), while text-to-image models tend to default to Western-centric representations (Kannen et al., 2024). These biases manifest as systematic performance disparities across languages and cultures (Field et al., 2021; Naous et al., 2023). Recent benchmarks such as GlobalRG (Bhatia et al., 2024) further highlight significant drops in VLM performance on local cultural concepts. Our work contributes to this line of research by studying explicit cultural adaptation in a generative setting, focusing on bidirectional cross-cultural transcreation.

**Meme Understanding and Analysis.** Most prior research on memes has focused on discriminative tasks, such as classification and detection. For example, PromptHate (Cao et al., 2023) addresses hateful meme detection (Kiela et al., 2021; Kumar and Nandakumar, 2022; Sharma et al., 2023), while MGMCF (Zheng et al., 2024) models fine-grained persuasive features (García-Díaz et al., 2024). Other studies document systematic cultural differences in online humor (Mutheu, 2023; Nissenbaum and Shifman, 2018; Tanaka et al., 2022), analyze the sentiment of memes (Sharma et al., 2020), and show that annotators’ cultural backgrounds influence interpretation. In contrast, comparatively little work has explored *generative* meme tasks, particularly those requiring culturally grounded adaptation rather than classification.

**Cross-Cultural Transcreation.** Transcreation has recently emerged as a framework for adapting content across cultures beyond literal translation. Khanuja et al. (2024b) introduce image transcreation and show that models struggle to balance semantic preservation with cultural appropriateness, motivating dedicated evaluation metrics (Khanuja et al., 2024a). While meme datasets such as MemeCap (Hwang and Shwartz, 2023) or MET-Meme (Xu et al., 2022) provide large-scale meme captioning resources, they lack cross-cultural pairs required for transcreation. Our work extends transcreation to memes, which require tight visual-textual coupling and humor preservation, and introduces a bidirectional benchmark spanning US and Chinese cultures.

**Generative Vision-Language Models.** Recent VLMs demonstrate strong multimodal understanding and reasoning capabilities, including LLaVA (Liu et al., 2023, 2024), GPT-4V (Lin et al.,

2025), and Qwen-VL (Bai et al., 2023; Wang et al., 2024). Parallel advances in image generation models enable increasingly faithful prompt-based visual synthesis (Black Forest Labs, 2024; Verdú and Martín, 2024). While these models provide the foundation for multimodal generation, their effectiveness for culturally grounded creative adaptation remains underexplored—a gap our study aims to address.

**Evaluation of Multimodal Generation.** Standard metrics such as CLIPScore (Hessel et al., 2021) and TIFA (Hu et al., 2023) focus on text-image alignment but are not designed to capture cultural fit or humor preservation. Existing cross-cultural benchmarks, including CVQA (Romero et al., 2024), GlobalRG (Bhatia et al., 2024), and WorldCuisines (Winata et al., 2025), primarily address visual question answering rather than generative tasks (Bai et al., 2025; Bhalerao et al., 2025). In this work, we evaluate meme transcreation using both human judgments and automated LLM-based evaluation across multiple quality dimensions, enabling a systematic comparison of human and automated assessments in a cross-cultural setting.

### 3 Hybrid Transcreation Framework

We introduce a hybrid framework for cross-cultural meme transcreation that balances preservation of communicative intent with culturally appropriate adaptation. Rather than framing memes as a translation task, *our approach explicitly separates culture-invariant elements from those that must change to ensure cultural authenticity*. This section outlines the guiding principles of the framework and the three-stage pipeline that implements them.

In practice, memes combine universal and culture-specific components: literal translation often preserves surface meaning but fails culturally, while full recreation risks drifting from the original intent. To address this trade-off, our hybrid strategy is grounded in three principles.

**Preserve universal elements.** We retain transferable aspects such as core humor mechanisms (e.g., irony, exaggeration), high-level emotional intent, and common meme formats.

**Adapt culture-specific elements.** We identify and replace culturally grounded references—such as pop culture, idioms, visual symbols, and stylistic conventions—with culturally appropriate alternatives rather than literal translations.

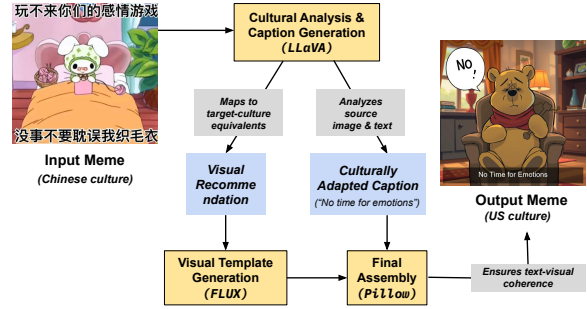


Figure 2: Overview of our three-stage meme transcreation pipeline. (1) A VLM analyzes the *original/input* meme, identifies cultural references and intent, and generates a culturally adapted caption. (2) A diffusion model produces a meme-style visual template aligned with the target culture. (3) A text overlay module assembles the output *transcreated* meme.

**Maintain intent consistency.** Across all stages, we preserve the meme’s communicative goal—*what it aims to express or satirize*—even when textual and visual inputs change.

#### 3.1 Meme Transcreation Pipeline

Figure 2 illustrates our modular three-stage meme transcreation pipeline: *cultural reasoning*, *visual generation*, and *final assembly*, enabling independent control and analysis across cultural directions. **Stage 1: Cultural Analysis and Caption Generation.** We use LLaVA 1.6 (13B) (Liu et al., 2024) as the core vision-language model for cultural analysis and caption generation. The model takes the original meme image as input and is prompted to (1) analyze cultural references and humor mechanisms, (2) extract the underlying intent, (3) map source-culture concepts to target-culture equivalents, and (4) generate a culturally appropriate caption in the target language.

We employ LLaVA 1.6 because it offers strong vision–language alignment, robust Chinese and English multilingual performance, and open-source reproducibility. Importantly, we do not fine-tune the model, focusing instead on prompt-based control to isolate the effects of cultural reasoning without introducing task-specific training bias. This stage outputs both a transcreated caption and high-level recommendations for adapting the visual content (i.e., mood, background - examples in Appendix A)

**Stage 2: Visual Template Generation.** Using the visual recommendations from Stage 1, we generate culturally adapted meme templates with FLUX.1 Schnell (Black Forest Labs, 2024; Verdú and Martín, 2024), a diffusion-based image generation model designed for strong prompt adherence

and fast iteration. At this stage, the goal is not photorealism but meme-appropriate visuals that support the intended humor and allow for clear and readable text overlay. The generated images preserve universal meme structures (e.g., reaction shots, simple backgrounds) while adapting culture-specific elements. For example, Western celebrity figures are often replaced with symbolic or animal-based representations that are more common in Chinese meme culture. Emotional tone is conveyed through facial expressions, posture, and visual metaphors that align with conventions in the target culture.

**Stage 3: Final Assembly.** In the final stage, we automatically combine transcreated captions with the generated visual templates using Pillow, an open-source image processing library.<sup>1</sup> This step handles font selection, text placement, and layout conventions appropriate for the target culture (e.g., denser layouts for Chinese text and more spaced layouts for English captions), following common social media meme practices. We apply dynamic text wrapping, semi-transparent background overlays for readability, and center-aligned multi-line captions positioned near the image bottom. Final manual quality checks verify readability, visual-text coherence, and that captions do not obscure key visual elements.

**Implementation Details.** All models are used in their pre-trained form, with prompt engineering (Appendix F) and decoding parameter tuning (e.g., temperature, top- $p$ ) to balance creativity and consistency. This modular design supports reproducibility, scalability across cultures, and controlled analysis of cross-cultural meme transcreation.

## 4 MemeXGen Dataset

To study cross-cultural meme transcreation in a controlled and realistic setting, we introduce **MemeX-Gen**, a multilingual and multicultural dataset of Chinese and US meme pairs. The dataset consists of 6,315 *original memes* collected from authentic social media platforms and 6,315 *transcreated memes* produced by our pipeline. For each original meme, we generate a corresponding transcreated version in the opposite cultural context, resulting in a total of 6,315 *bidirectional meme pairs*: 3,165 *Chinese*→*US* and 3,150 *US*→*Chinese*. This paired structure enables direct comparison of transcreation quality across directions.

<sup>1</sup><https://python-pillow.org/>

## 4.1 Data Sources

MemeXGen is designed to support systematic evaluation and analysis, with an emphasis on cultural authenticity, diversity of humor styles, and balanced coverage across Chinese and US cultures.

**Chinese Memes.** *Original* Chinese memes are sourced from the publicly available *Chinese Meme Description Dataset*<sup>2</sup>, which aggregates content from two major Chinese social media platforms: **Xiaohongshu** and **Weibo**. Xiaohongshu contributes lifestyle- and emotion-focused memes, while Weibo provides fast-paced, commentary-driven content reflecting mainstream Chinese internet culture.

**US Memes.** *Original* US memes are drawn from the *MemeCap* dataset<sup>3</sup>, which collects memes from popular Reddit communities such as *r/memes* and *r/dankmemes*. These memes reflect dominant US humor styles, including sarcasm, irony, pop culture references, and situational storytelling.

These sources offer complementary views of meme culture in two distinct cultural contexts, enabling systematic bidirectional transcreation.

## 4.2 Filtering and Dataset Composition

During data inspection, we observe that some *original* memes contain potentially sensitive content (e.g., political references) that could interfere with fair evaluation or raise ethical concerns. To address this, we manually filter the *original* memes to ensure responsible use and reliable evaluation. Specifically, we remove memes that are offensive, contain low-quality or corrupted images, rely heavily on mixed languages, or exhibit weak visual-text integration. After filtering, the dataset contains 6,315 *original* memes, equally split across US and Chinese. We notice that the retention rate is substantially higher for the Chinese subset (97.6%) than for the US subset (55.0%), reflecting stricter content moderation on Chinese platforms compared to the more permissive nature of Reddit.

## 4.3 Annotation and Evaluation Split

To support emotion analysis and human evaluation, we annotate two disjoint subsets of the data.

**Emotion annotations subset.** We annotate  $\approx 10\%$  of the *original* memes data (628 memes, equally split between US and Chinese memes) with emotion labels, including emotion category (*Joy*, *Anger*,

<sup>2</sup><https://github.com/THUDM/chinese-meme-description-dataset>

<sup>3</sup><https://github.com/hwang1996/MemeCap>

*Sadness, Fear, Disgust, Surprise*) and emotion intensity on a 1–5 Likert scale. Annotation guidelines follow recent multilingual emotion classification work, as described in BRIGHTER (Muhammad et al., 2025). Three expert annotators perform the annotations independently, achieving strong agreement (Fleiss’ (Fleiss, 1971)  $\kappa = 0.869$  for emotion category and  $\kappa = 0.536$  for intensity).

**Human evaluation subset.** We reserve a separate 10% subset (628 *original* memes) as the test set for transcreation experiments. This split includes 313 Chinese→US and 315 US→Chinese *original-transcreated* meme pairs and is used exclusively for human evaluation of meme transcreation. Evaluation details are provided in Section 5.

#### 4.4 Dataset Characteristics

To better understand the cultural makeup of the *original* memes, we analyze topic and emotion distributions using Qwen-VL-Max (Bai et al., 2023; Wang et al., 2024), finetuned on the human annotated emotions. To validate the reliability of the predicted labels, an expert annotator manually reviews a random 10% subset and confirms that over 95% of the topic and emotion labels are correct.

**Topic Distribution.** Both cultures are dominated by themes related to *Internet Culture* (CN 61.0%, US 52.4%) and *Technology/Digital Life* (CN 10.6%, US 15.1%). Beyond these shared themes, clear differences emerge. **US memes more often frame education, family, and everyday experiences as relatable, narrative-driven humor** (e.g., *Education*: 7.8%; *Family*: 4.9%), whereas **Chinese memes emphasize symbolic expression and social pressure**, with lower prevalence of these topics (*Education*: 2.1%; *Family*: 1.9%). *Gaming*-related humor appears among the top US topics (2.7%) but is largely absent from the Chinese top ranks, reflecting differing leisure and achievement orientations.

**Emotion Distribution.** Automated emotion classification shows that **Joy dominates in both cultures** (CN 69.3%, US 73.8%), consistent with memes’ primary entertainment role. However, Chinese memes exhibit higher levels of *Anger* (8.3%) and *Sadness* (8.2%), suggesting more frequent **social critique and melancholic expression**. In contrast, US memes show relatively higher *Fear* (7.0%) and *Disgust* (4.4%), aligning with **anxiety-driven and cringe-based humor styles**.

These systematic differences motivate our hybrid transcreation approach and provide context

for interpreting performance asymmetries in later experiments. Further data analysis is provided in Appendix B.1.

## 5 Evaluation Methodology

We evaluate our meme transcreation framework using both human and VLM-based evaluation.

### 5.1 Metric Definitions

Our evaluation captures not only text and image quality and their interaction, as commonly assessed in prior image generation work (Hu et al., 2023), but also *task-specific* aspects that are critical for cross-cultural transcreation, namely cultural fit and intent preservation. All quantitative metrics are rated on a 5-point Likert scale (1 = strongly disagree, 5 = strongly agree).

We evaluate each transcreated meme along six dimensions: **Caption Quality**, measuring clarity, tone, and meme-appropriate language; **Image Quality**, assessing visual clarity, composition, and recognizability; **Synergy**, capturing how well image and text work together to convey humor or emotion; **Cultural Fit**, evaluating appropriateness and relatability for the target culture; **Intent Preservation**, measuring retention of the original meme’s message and emotional effect; and an **Overall Score**, computed as the average across all dimensions. Detailed dimension definitions are provided in Appendix C.

### 5.2 Human Evaluation

We evaluate meme transcreation on the human evaluation subset, with each meme independently assessed by three human evaluators. Because meme generation is inherently subjective and prior work highlights the importance of modeling annotator perspectives rather than simple aggregation (Deng et al., 2023), we report results separately for each evaluator. All evaluators are bilingual and bicultural, with deep familiarity with both Chinese and US meme cultures. Additional details on annotator backgrounds are provided in Appendix D.

**Inter-annotator Agreement.** Inter-annotator Pearson correlations indicate moderate to strong agreement ( $r = 0.58$ – $0.81$ ), reflecting reliable yet stylistically distinct evaluation perspectives.

### 5.3 Automatic Evaluation

To assess the feasibility of automated evaluation, we use six state-of-the-art VLMs on all the data:

434 Qwen-VL-Max (Bai et al., 2023), LLaVA-v1.6-  
435 Vicuna-13B (Liu et al., 2024), InternVL3-8B and  
436 InternVL3-14B (Zhu et al., 2025), Qwen3-VL-  
437 8B (Wang et al., 2024), and Aya-vision-8b (Dash  
438 et al., 2025). These models are selected for their  
439 multilingual Chinese–English support, ability to  
440 process multiple images, and public availability,  
441 enabling reproducible automatic evaluation.

442 **VLM-Human Agreement.** We assess VLM evalua-  
443 tion effectiveness by computing the Pearson cor-  
444 relation between each human evaluator and VLM  
445 across each dimension (results in Section 6.3).

## 446 6 Evaluation Results

447 We report evaluation results addressing our re-  
448 search questions, using human and automatic met-  
449 rics to assess cross-cultural performance, direc-  
450 tional effects, and evaluation reliability.

### 451 6.1 RQ1: Cross-Cultural Performance

452 **Human Evaluation.** Table 2 summarizes results  
453 from three human evaluators and six LLM evalua-  
454 tors across both transcreation directions. Human  
455 evaluators differ in strictness and focus: Evaluator 1  
456 prioritizes entertainment value (mean: 4.42), Eval-  
457 uator 2 adopts a balanced perspective (mean: 4.09),  
458 and Evaluator 3 applies stricter quality standards  
459 (mean: 3.31). The resulting 1.11-point spread high-  
460 lights the inherent subjectivity of cross-cultural  
461 meme transcreation evaluation. **Overall, the mean  
462 human score of 4.07/5.0 indicates that the pro-  
463 posed pipeline produces effective and generally  
464 well-received transcreations.**

465 **Dimension-Level Analysis.** Across dimensions,  
466 *Caption Quality* receives the highest ratings (mean:  
467 4.20), suggesting effective cross-cultural adapta-  
468 tion of meme text. *Image Quality* is also strong  
469 (mean: 4.05), supporting the reliability of FLUX.1  
470 for meme-style visual generation. **Synergy is con-  
471 sistent high (mean: 4.23)**, indicating that cap-  
472 tions and visuals work well together in most out-  
473 puts. *Cultural Fit* shows the widest variation across  
474 evaluators (range: 3.39–4.57), reflecting the subjec-  
475 tive and culturally grounded nature of authenticity  
476 judgments. *Intent Preservation* is rated favorably  
477 overall (mean: 3.84), though scores may be par-  
478 tially influenced by the dominance of Joy-related  
479 memes (69–74% of the data).

480 **VLM Evaluation.** Among automated evaluators,  
481 **Qwen-VL-Max performs best**, achieving the high-  
482 est overall scores (3.95 for Chinese→US and 3.71

483 for US→Chinese) and showing **strong alignment**  
484 **with human judgments** (mean Pearson  $r = 0.926$ ,  
485 all  $p < 0.001$ ). Other LLM evaluators exhibit  
486 much weaker correlations with humans ( $r \leq 0.33$ ),  
487 suggesting that most open-source models struggle  
488 to reliably evaluate creative, culturally grounded  
489 outputs. **On average, LLM scores are 0.54 points**  
490 **lower than human scores**, indicating a systematic  
491 tendency toward conservative scoring.

### 492 6.2 RQ2: Directional Asymmetries

493 We observe a clear directional asymmetry:  
494 **US→Chinese meme transcreations significantly**  
495 **outperform Chinese→US** (4.48 vs. 3.93 out of  
496 5.0,  $\delta = +0.55$ ,  $p < 0.001$ ). This gap likely  
497 reflects several factors. First, US memes of-  
498 ten rely on globally recognizable characters and  
499 themes that are easier to localize, whereas Chi-  
500 nese memes frequently depend on context-specific  
501 wordplay and implicit cultural knowledge. Sec-  
502 ond, current VLMs are more strongly exposed  
503 to US-centric data during training. Third, eval-  
504 uators apply stricter authenticity expectations to  
505 Chinese→US outputs, where cultural mismatches  
506 are more salient to native speakers.

### 507 6.3 RQ3: Evaluation Framework Analysis

508 Table 3 reports correlations between human and  
509 LLM evaluators. **Qwen-VL-Max shows consis-  
510 tently strong alignment with all three human  
511 evaluators** (Evaluator 1:  $r = 0.921$ , Evaluator 2:  
512  $r = 0.964$ , Evaluator 3:  $r = 0.894$ ), with es-  
513 pecially high agreement on *Intent Preservation*  
514 ( $r = 0.993$ ). In contrast, other models exhibit  
515 substantially weaker correlations (e.g., InternVL3-  
516 14B:  $r = 0.263$ , Qwen3-VL-8B:  $r = 0.252$ ,  
517 LLaVA-v1.6:  $r = 0.005$ ). **These results sug-  
518 gest that evaluating creative, cross-cultural mul-  
519 timodal content requires deeper multilingual  
520 and multicultural grounding than most current  
521 open-source VLMs provide.**

### 522 6.4 Qualitative Analysis

523 To complement quantitative metrics, we present  
524 representative transcreation examples from both  
525 directions and our data analysis observations. Fig-  
526 ure 3a shows two successful transcreation samples  
527 from both directions (Overall Score: 5.0/5.0), while  
528 Figure 3b illustrates several failed transcreation  
529 samples (Overall Score: 1.4/5.0).

530 **Success Patterns (30% of outputs scoring**  
531  **$\geq 4.5/5.0$ ).** High quality meme transcreations con-

Evaluator	Chinese→US						US→Chinese					
	Cap.	Img.	Syn.	Cult.	Intent	Overall	Cap.	Img.	Syn.	Cult.	Intent	Overall
Evaluator 1 (H)	4.78	4.51	4.66	4.57	4.24	4.55	4.82	4.22	4.31	4.18	3.89	4.28
Evaluator 2 (H)	4.35	4.11	4.45	4.14	4.00	4.21	4.41	3.84	4.12	3.76	3.67	3.96
Evaluator 3 (H)	3.46	3.52	3.59	3.39	3.29	3.45	3.52	3.19	3.24	2.98	2.93	3.17
Qwen-VL-Max	<b>4.13</b>	3.86	<b>4.20</b>	<b>3.74</b>	3.72	<b>3.95</b>	<b>4.21</b>	3.58	<b>3.89</b>	<b>3.41</b>	3.44	<b>3.71</b>
LLaVA-v1.6	<u>4.00</u>	<u>4.00</u>	<u>3.81</u>	3.00	<b>4.00</b>	<u>3.79</u>	<u>4.05</u>	3.72	<u>3.52</u>	2.67	<b>3.71</b>	<u>3.53</u>
InternVL3-8B	3.78	3.84	3.48	3.46	<u>3.97</u>	3.69	3.84	3.55	3.16	3.12	<u>3.68</u>	3.47
InternVL3-14B	3.21	<b>4.39</b>	3.16	3.36	3.34	3.53	3.28	<b>4.11</b>	2.84	3.02	3.01	3.25
Qwen3-VL-8B	2.83	3.70	2.74	<u>3.59</u>	2.56	3.15	2.91	3.42	2.41	<u>3.21</u>	2.18	2.83
Aya-vision-8b	3.18	<u>4.17</u>	2.83	2.90	2.72	3.10	3.25	<u>3.89</u>	2.51	2.56	2.39	2.92

Note: Best VLM results per column are shown in **bold**, second-best VLM results are underlined. All dimensions rated 1–5 (higher = better). (H) = Human evaluator. Cap.=Caption Quality, Img.=Image Quality, Syn.=Synergy, Cult.=Cultural Fit, Intent=Intent Preservation.

Table 2: Evaluation Results for Chinese→US and US→Chinese Meme Transcreation

Human	VLM	Cap.	Img.	Syn.	Cult.	Intent	Overall
Eval. 1	Qwen-VL-Max	0.961	0.986	0.987	0.901	0.993	0.921
	LLaVA-v1.6	-0.039	-0.039	0.082	-0.178	-0.039	-0.026
	InternVL3-8B	-0.128	-0.026	-0.053	-0.106	0.041	-0.086
	InternVL3-14B	0.241	0.363	0.288	0.216	0.275	0.270
	Qwen3-VL-8B	0.219	0.394	0.316	0.215	0.289	0.281
	Aya-vision-8b	-0.117	-0.032	-0.087	-0.077	-0.065	-0.082
Eval. 2	Qwen-VL-Max	0.989	0.988	0.994	0.980	0.995	0.964
	LLaVA-v1.6	-0.016	-0.016	0.038	-0.112	-0.016	-0.027
	InternVL3-8B	-0.099	0.014	-0.033	-0.061	0.041	-0.066
	InternVL3-14B	0.294	0.350	0.351	0.323	0.365	0.329
	Qwen3-VL-8B	0.265	0.340	0.331	0.294	0.327	0.309
	Aya-vision-8b	-0.088	0.022	-0.048	-0.034	-0.031	-0.058
Eval. 3	Qwen-VL-Max	0.950	0.990	0.985	0.938	0.990	0.894
	LLaVA-v1.6	-0.050	-0.050	0.039	-0.161	-0.050	0.029
	InternVL3-8B	-0.065	0.083	0.014	0.008	0.107	-0.013
	InternVL3-14B	0.296	0.348	0.373	0.320	0.371	0.340
	Qwen3-VL-8B	0.193	0.368	0.310	0.235	0.285	0.263
	Aya-vision-8b	-0.024	0.101	0.016	0.050	0.025	-0.057

Note: Cap.=Caption Quality, Img.=Image Quality, Syn.=Synergy, Cult.=Cultural Fit, Intent=Intent Preservation. Cell colors indicate correlation strength.

Table 3: Pearson correlation coefficients ( $r$ ) between Human and LLM evaluators across evaluation dimensions.

tain the following elements: (1) *Universally applicable character selection*—the use of recognizable archetypes that can be understood across cultures. (2) *Emotion-focused transcreations*—the retention of the original emotional context with the incorporation of cultural specifics. (3) *Use of natural language conventions*—the use of meme-like linguistic conventions associated with the receiving culture. (4) *Visual and textual unity*—careful matching of image and text

**Failure Patterns (1.6% of outputs scoring  $\leq 2.0/5.0$ ).** Errors that emerge on failed meme transcreations include: (1) *Failure in Captions*—the use of formal speech that dampens the casual meme vibe; (2) *Disconnects in Visuals*—the use of images that don’t fit a culture or issues with visual generation; (3) *Failure to Preserve Humor Mecha-*

*nisms*—the use of a format that is not amenable to joke structure; (4) *Complete synergy breakdown*—caption-image mismatch creating incoherent messaging.

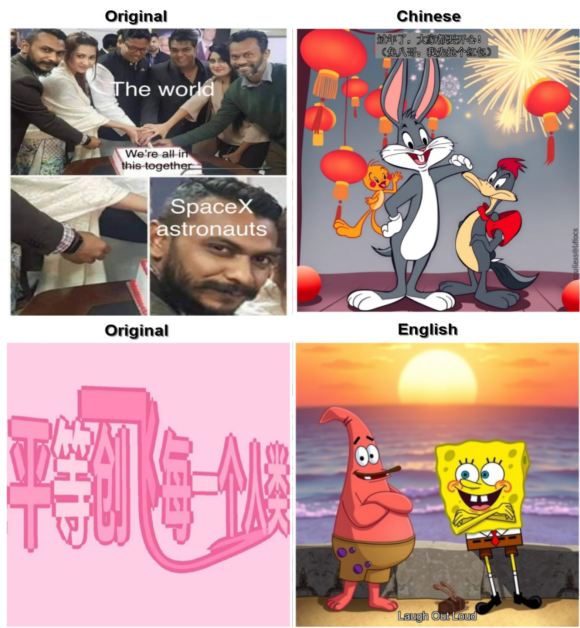
**Directional Patterns.** US→Chinese transcreations more frequently achieve natural cultural integration, benefiting from globally recognizable US templates. Chinese→US transcreations struggle with context-dependent wordplay and philosophical concepts that lack Western equivalents, often resulting in superficial adaptations. Additional examples in Appendix G.

## 7 Main Takeaways

**Effective Cross-Cultural Transcreation.** Human evaluations show that the proposed hybrid approach produces high-quality transcreations (mean score: 4.07/5.0), successfully preserving humor and intent while adapting cultural specifics. Strong performance on Caption Quality (4.20) and Image–Text Synergy (4.23) confirms that the three-stage pipeline supports coherent multimodal generation. **Directional Asymmetry Matters.** US→Chinese transcreation consistently outperforms Chinese→US (+0.55), reflecting both model exposure biases and deeper cultural differences. In particular, Chinese memes rely more heavily on context-dependent wordplay and implicit meaning, which are harder to adapt than the more visually universal templates common in US meme culture. These results highlight the need for culturally diverse training data in cross-cultural AI systems. **Limits of Automated Evaluation.** Qwen-VL-Max shows strong agreement with human judgments ( $r = 0.926$ ), demonstrating that automated evalua-



(a) **Best transcreation example (Score: 5.0/5.0).**  
**Transcreated (Chinese, top right — Bugs Bunny):** “Family: Look who’s up so early; My mental state after pulling an all-nighter to finish my assignment due.”  
**Original (bottom left — panda meme):** “What’s wrong with work hours? Doesn’t your company expect you to work during work hours?”



(b) **Worst transcreation example (Score: 1.4/5.0).**  
**Transcreated (Chinese, top right — Looney Tunes):** “Celebrating the New Year, everyone must be happy! No one left behind, I’m giving you a family photo!”  
**Original (bottom left — Chinese text):** “Attempts at equality satisfy no one.”

Figure 3: **Qualitative examples of cross-cultural meme transcreation.** **Left:** successful US→Chinese adaptation preserving intent, humor, and cultural conventions. **Right:** failed Chinese→US adaptation illustrating loss of intent and cultural mismatch.

tion of creative, cross-cultural content is feasible. However, weaker correlations from open-source models suggest that reliable automated evaluation remains challenging without extensive multilingual and multicultural grounding.

## 8 Conclusion

We introduced a hybrid framework for cross-cultural meme transcreation that explicitly separates intent preservation from cultural adaptation, enabling principled analysis of how humor and meaning transfer across cultures. By combining vision–language models with diffusion-based image generation, our approach moves beyond literal translation and treats meme adaptation as a culturally grounded multimodal reasoning problem.

We curated and evaluated a dataset of 6,315 Chinese–U.S. meme pairs, combining authentic social media memes with systematically generated transcreations, and conducted a comprehensive bidirectional evaluation. Our results reveal consistent directional asymmetries in transcreation

quality, demonstrating that current models handle certain cultural adaptations more effectively than others. These findings expose concrete limitations in cross-cultural generalization that are not visible in monolingual or translation-based evaluations.

We further show that carefully selected VLM-based evaluators can approximate human judgments on culturally grounded dimensions such as emotion and intent, while most open-source models remain unreliable for assessing intent and cultural fit. Finally, we release MEMEXGEN, the first parallel Chinese–U.S. meme transcreation corpus annotated for emotion and cultural intent, together with evaluation protocols and dataset splits. By open-sourcing data, models, and evaluation metrics, this work establishes a foundation for systematic study of computational humor and cross-cultural multimodal generation, and provides actionable benchmarks for future model development.

623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671

## Limitations

**Scope of Cultural Coverage.** This study focuses on meme transcreation between Chinese and U.S. cultures, which differ substantially in language, humor conventions, and visual symbolism. While this contrast enables clear analysis of cultural asymmetries, our findings should not be assumed to generalize uniformly to other cultural pairs. Future work should extend this framework to additional language–culture settings to test the robustness of the observed patterns.

**Generality of the Generation Framework.** Our transcreation pipeline combines existing vision–language and diffusion models in a modular design intended to support interpretability and controlled analysis rather than architectural novelty. We do not claim optimality of this design, nor do we compare against all possible end-to-end prompting alternatives. Instead, our goal is to provide a transparent framework for studying cultural adaptation. Exploring simpler or fully integrated baselines remains an important direction for future work.

**Interpreting Directional Asymmetries.** We observe consistent performance differences between US→Chinese and Chinese→US transcreation. While we discuss plausible contributing factors—such as training data exposure, humor structure, and evaluator expectations—these explanations are correlational rather than causal. Disentangling these effects would require controlled experiments that vary data distributions, model pretraining, and evaluation populations independently.

**Limits of Automatic Evaluation.** Although Qwen-VL-Max shows strong alignment with human judgments in our setting, this result may reflect model-specific strengths in Chinese–English multimodal understanding rather than a general solution to evaluating culturally grounded humor. The weak performance of other open-source evaluators highlights that reliable automated evaluation remains challenging and should be interpreted with caution.

**Dataset Composition and Emotion Coverage.** Joy dominates the meme distributions in both cultures, reflecting real-world social media trends but limiting stress-testing on negative or socially critical humor. As a result, intent preservation scores may be optimistic for emotionally complex cases. Expanding emotion-balanced datasets is a key area for future research.

**Evaluation at Scale and in the Wild.** Human evaluation remains inherently subjective, and the observed variation across evaluators highlights the value of modeling diverse perspectives rather than collapsing them into a single score. Expanding the evaluation set and incorporating longitudinal, in-the-wild measurements (e.g., engagement or sharing behavior) would provide deeper insight into real-world cultural impact beyond offline quality ratings.

**Broadening Cultural Perspectives.** Our annotators are bilingual and bicultural with Chinese–U.S. experience, ensuring informed evaluation of both contexts. Future work can further broaden cultural representation by including evaluators with more localized or region-specific backgrounds, as well as exploring regional variation within Chinese and U.S. meme cultures. Such diversity would deepen understanding of cultural nuance and strengthen the generalizability of cross-cultural evaluation.

## Ethical Considerations

**Deployment Scope.** Our pipeline prioritizes analytical clarity and controlled study of cultural adaptation rather than deployment efficiency. As with any automated cultural generation system, misuse, misinterpretation, or oversimplification of cultural signals remains a risk. We position meme transcreation as a *decision-support tool* intended to assist human creators and analysts, not as a fully autonomous content generator, and strongly recommend human oversight in real-world or sensitive deployments.

**Content Safety.** We apply stringent manual filtering to exclude offensive or sensitive content, including hate speech, discriminatory media, explicit material, and political content. High *Not Offensive* ratings (92.8% from human evaluators and 96.6% from LLM-based assessment) indicate the effectiveness of these safeguards. However, cultural sensitivity is inherently subjective: content acceptable in one cultural context may still be perceived as offensive in another. Our system therefore cannot guarantee zero harmful outputs and should be used with caution, particularly in public-facing applications.

**Cultural Respect and Representation.** Automated cultural adaptation risks reinforcing stereotypes or reducing complex cultural practices to surface-level substitutions. While our hybrid framework explicitly separates intent preservation from

672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721

722	cultural adaptation, evaluator feedback reveals oc-		
723	casual cases of shallow cultural mapping (e.g.,		
724	direct visual substitution without deeper contextual		
725	grounding). These limitations highlight the impor-		
726	tance of human-in-the-loop workflows, where au-		
727	tomated transcreation outputs are treated as drafts		
728	rather than finalized content.		
729	<b>Data Privacy and Attribution.</b> All source memes		
730	are collected from publicly accessible platforms		
731	(Xiaohongshu and Weibo for Chinese memes; Red-		
732	dit for U.S. memes) and do not contain personal		
733	identifying information. We respect the implicit		
734	consent associated with public content sharing,		
735	though proper attribution remains challenging for		
736	viral meme formats with unclear authorship. The		
737	dataset is intended strictly for research purposes,		
738	and we encourage responsible use consistent with		
739	platform norms and community standards.		
740	<b>Misinformation Potential.</b> Meme transcreation		
741	tools could be misused to spread culturally-adapted		
742	misinformation or propaganda. We emphasize re-		
743	sponsible deployment with content verification pro-		
744	cedures and transparency about automated genera-		
745	tion.		
746	<b>References</b>		
747	Muhammad Farid Adilazuarda, Sagnik Mukherjee,		
748	Pradhyumna Lavania, Siddhant Singh, Alham Fikri		
749	Aji, Jacki O’Neill, Ashutosh Modi, and Monojit		
750	Choudhury. 2024. <b>Towards measuring and model-</b>		
751	<b>ing "culture" in llms: A survey.</b> In <i>Proceedings of</i>		
752	<i>the 2024 Conference on Empirical Methods in Natu-</i>		
753	<i>ral Language Processing, EMNLP 2024, Miami, FL,</i>		
754	<i>USA, November 12-16, 2024</i> , pages 15763–15784.		
755	Association for Computational Linguistics.		
756	Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,		
757	Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,		
758	and Jingren Zhou. 2023. <b>Qwen-vl: A versatile</b>		
759	<b>vision-language model for understanding, localiza-</b>		
760	<b>tion, text reading, and beyond.</b> <i>arXiv preprint</i>		
761	<i>arXiv:2308.12966</i> .		
762	Longju Bai, Angana Borah, Oana Ignat, and Rada Mi-		
763	halcea. 2025. <b>The power of many: Multi-agent mul-</b>		
764	<b>timodal models for cultural image captioning.</b> In		
765	<i>Proceedings of the 2025 Conference of the Nations</i>		
766	<i>of the Americas Chapter of the Association for</i>		
767	<i>Computational Linguistics: Human Language Techno-</i>		
768	<i>gies (Volume 1: Long Papers)</i> , pages 2970–2993,		
769	Albuquerque, New Mexico. Association for Compu-		
770	tational Linguistics.		
771	Parth Bhalerao, Mounika Yalamarty, Brian Trinh, and		
772	Oana Ignat. 2025. <b>Multi-agent multimodal models</b>		
773	<b>for multicultural text to image generation.</b> <i>ArXiv,</i>		
774	<i>abs/2502.15972</i> .		
	Mehar Bhatia, Sahithya Ravi, Aditya Chinchure, Eu-	775	
	nJeong Hwang, and Vered Shwartz. 2024. <b>From</b>	776	
	<b>local concepts to universals: Evaluating the multi-</b>	777	
	<b>cultural understanding of vision-language models.</b>	778	
	In <i>Proceedings of the 2024 Conference on Empiri-</i>	779	
	<i>cal Methods in Natural Language Processing</i> , pages	780	
	6763–6782, Miami, Florida, USA. Association for	781	
	Computational Linguistics.	782	
	Black Forest Labs. 2024. <b>Flux: Advanced text-to-</b>	783	
	<b>image generation model.</b> <a href="https://github.com/black-forest-labs/flux">https://github.com/</a>	784	
	<a href="https://github.com/black-forest-labs/flux">black-forest-labs/flux</a> . GitHub Repository.	785	
	Rui Cao, Roy Ka-Wei Lee, Tuan-Anh Hoang, Junmo	786	
	Pang, Kenji Kawaguchi, and Roger Zimmermann.	787	
	2023. <b>Promptgate: Prompting for hateful meme clas-</b>	788	
	<b>sification.</b> <i>Preprint</i> , arXiv:2302.04156.	789	
	Saurabh Dash, Shivalika Singh, Adrien Morisot, Beyza	790	
	Ermis, Acyr Locatelli, Sungjin Hong, Arash Ahma-	791	
	dian, Yannis Flet-Berliac, Nathan Grinsztajn, Florian	792	
	Strub, and 1 others. 2025. <b>Aya vision: Advancing</b>	793	
	<b>the frontier of multilingual multimodality.</b> <i>Preprint,</i>	794	
	<i>arXiv:2505.08751</i> .	795	
	Naihao Deng, Xinliang Zhang, Siyang Liu, Winston Wu,	796	
	Lu Wang, and Rada Mihalcea. 2023. <b>You are what</b>	797	
	<b>you annotate: Towards better models through anno-</b>	798	
	<b>tor representations.</b> In <i>Findings of the Association</i>	799	
	<i>for Computational Linguistics: EMNLP 2023</i> , pages	800	
	12475–12498, Singapore. Association for Computa-	801	
	tional Linguistics.	802	
	Anjalie Field, Su Lin Blodgett, Zeerak Waseem, and	803	
	Yulia Tsvetkov. 2021. <b>A Survey of Race, Racism,</b>	804	
	<b>and Anti-Racism in NLP.</b> In <i>Proceedings of the 59th</i>	805	
	<i>Annual Meeting of the ACL</i> , pages 1905–1925.	806	
	Joseph L. Fleiss. 1971. <b>Measuring nominal scale agree-</b>	807	
	<b>ment among many raters.</b> <i>Psychological Bulletin</i> ,	808	
	76(5):378–382.	809	
	José Antonio García-Díaz and 1 others. 2024. <b>Umuteam</b>	810	
	<b>at semeval-2024 task 4: Multilingual detection of</b>	811	
	<b>persuasion techniques in memes.</b> <i>Proceedings of</i>	812	
	<i>SemEval-2024</i> .	813	
	Muzhaffar Hazman, Susan McKeever, and Josephine	814	
	Griffith. 2025. <b>What makes a meme a meme? identi-</b>	815	
	<b>fying memes for memetics-aware dataset creation.</b>	816	
	Daniel Hershcovich and 1 others. 2022. <b>Challenges and</b>	817	
	<b>strategies in cross-cultural NLP.</b> In <i>Proceedings of</i>	818	
	<i>the 60th Annual Meeting of the ACL</i> , pages 6997–	819	
	7013.	820	
	Jack Hessel and 1 others. 2021. <b>CLIPScore: A</b>	821	
	<b>reference-free evaluation metric for image captioning.</b>	822	
	In <i>Proceedings of the 2021 Conference on Empirical</i>	823	
	<i>Methods in Natural Language Processing (EMNLP)</i> ,	824	
	pages 7514–7528.	825	
	Yushi Hu and 1 others. 2023. <b>Tifa: Accurate and in-</b>	826	
	<b>terpretable text-to-image faithfulness evaluation with</b>	827	
	<b>question answering.</b> <i>Preprint</i> , arXiv:2303.11897.	828	

829	EunJeong Hwang and Vered Shwartz. 2023. <a href="#">Meme-cap: A dataset for captioning and interpreting memes</a> . <i>Preprint</i> , arXiv:2305.13703.	885
830		886
831		887
832	Nithish Kannen, Arjun Palani, Vikram V. Ramaswamy, Olga Russakovsky, and Li Fei-Fei. 2024. <a href="#">Beyond aesthetics: Cultural competence in text-to-image models</a> . <i>Preprint</i> , arXiv:2407.06863.	888
833		889
834		890
835		891
836	Simran Khanuja, Vivek Iyer, Claire He, and Graham Neubig. 2024a. <a href="#">Towards automatic evaluation for image transcreation</a> . <i>Preprint</i> , arXiv:2412.13717.	892
837		893
838		894
839	Simran Khanuja, Sathyanarayanan Ramamoorthy, Yueqi Song, and Graham Neubig. 2024b. <a href="#">An image speaks a thousand words, but can everyone listen? on image transcreation for cultural relevance</a> . <i>Preprint</i> , arXiv:2404.01247.	895
840		896
841		897
842		898
843		899
844	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Casey A. Fitzpatrick, Peter Bull, Greg Lipstein, Tony Nelli, Ron Zhu, Niklas Muennighoff, Riza Velioğlu, Jewgeni Rose, Phillip Lippe, Nithin Holla, Shantanu Chandra, Santhosh Rajamanickam, Georgios Antoniou, Ekaterina Shutova, and 6 others. 2021. <a href="#">The hateful memes challenge: Competition report</a> . In <i>Proceedings of the NeurIPS 2020 Competition and Demonstration Track</i> , volume 133 of <i>Proceedings of Machine Learning Research</i> , pages 344–360. PMLR.	900
845		901
846		902
847		903
848		904
849		905
850		906
851		907
852		908
853		909
854		910
855	Gokul Karthik Kumar and Karthik Nandakumar. 2022. <a href="#">Hate-CLIPper: Multimodal hateful meme classification based on cross-modal interaction of CLIP features</a> . In <i>Proceedings of the Second Workshop on NLP for Positive Impact (NLP4PI)</i> , pages 171–183, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	911
856		912
857		913
858		914
859		915
860		916
861		917
862	Hongzhan Lin, Ziyang Luo, Bo Wang, Ruichao Yang, and Jing Ma. 2025. <a href="#">Goat-bench: Safety insights to large multimodal models through meme-based social abuse</a> . <i>ACM Trans. Intell. Syst. Technol.</i>	918
863		919
864		920
865		921
866	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. <a href="#">Llava-next: Improved reasoning, ocr, and world knowledge</a> .	922
867		923
868		924
869	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. <a href="#">Visual instruction tuning</a> . <i>Preprint</i> , arXiv:2304.08485.	925
870		926
871		927
872	Rada Mihalcea, Oana Ignat, Longju Bai, Angana Borah, Luis Chiruzzo, Zhijing Jin, Claude Kwizera, Joan Nwatu, Soujanya Poria, and Tamar Solorio. 2025. <a href="#">Why ai is weird and shouldn't be this way: towards ai for everyone, with everyone, by everyone</a> . In <i>Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence</i> , AAAI'25/IAAI'25/EAAI'25. AAAI Press.	928
873		929
874		930
875		931
876		932
877		933
878		934
879		935
880		936
881		937
882		938
883	Shamsuddeen Hassan Muhammad, Nedjma Ousidhoum, Idris Abdulmumin, Jan Philip Wahle, Terry Ruas, Meriem Beloucif, Christine de Kock, Nirmal Surange, Daniela Teodorescu, Ibrahim Said Ahmad, David Ifeoluwa Adelani, Alham Fikri Aji, Felermino D. M. A. Ali, Ilseyar Alimova, Vladimir Araujo, Nikolay Babakov, Naomi Baes, Ana-Maria Bucur, Andiswa Bukula, and 29 others. 2025. <a href="#">BRIGHTER: BRIDging the gap in human-annotated textual emotion recognition datasets for 28 languages</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8895–8916, Vienna, Austria. Association for Computational Linguistics.	939
884		940
885		941
886		942
887		943
888		944
889		945
890		946
891		947
892		948
893		949
894		950
895		951
896		952
897		953
898		954
899		955
900		956
901		957
902		958
903		959
904		960
905		961
906		962
907		963
908		964
909		965
910		966
911		967
912		968
913		969
914		970
915		971
916		972
917		973
918		974
919		975
920		976
921		977
922		978
923		979
924		980
925		981
926		982
927		983
928		984
929		985
930		986
931		987
932		988
933		989
934		990
935		991
936		992
937		993
938		994
939		995
940		996
941		997
942		998
943		999
944		1000

943	Genta Indra Winata and 1 others. 2025. <a href="#">Worldcuisines: A massive-scale benchmark for multilingual and multicultural visual question answering on global cuisines</a> . In <i>Proceedings of the 2025 NAACL</i> , pages 3242–3264.	with bold lines and vibrant colors. - Mood: Soft, nostalgic lighting with a hint of melancholy.	995 996 997
948	Bo Xu, Tingting Li, Junzhe Zheng, Mehdi Nasiriparsa, Zhehuan Zhao, Hongfei Lin, and Feng Xia. 2022. <a href="#">Met-meme: A multimodal meme dataset rich in metaphors</a> . In <i>Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22</i> , page 2887–2899, New York, NY, USA. Association for Computing Machinery.		998
956	Zhengyi Zhao, Shubo Zhang, Yuxi Zhang, Yanxi Zhao, Yifan Zhang, Zezhong Wang, Huimin Wang, Yutian Zhao, Bin Liang, Yefeng Zheng, Binyang Li, Kam-Fai Wong, and Xian Wu. 2025. <a href="#">Memereacon: Probing contextual meme understanding in large vision-language models</a> . <i>Preprint</i> , arXiv:2505.17433.		999
962	Li Zheng, Xiaoming Wang, Yuhan Chen, and Jianwei Liu. 2024. <a href="#">Multi-granular multimodal clue fusion for meme understanding</a> . <i>Preprint</i> , arXiv:2503.12560.		1000
966	Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, and 1 others. 2025. <a href="#">Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models</a> . <i>Preprint</i> , arXiv:2504.10479.		1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026
971	<b>A Meme Transcreation Framework</b>		1027
972	Examples of visual recommendations from our Stage 1 LLaVA output to the FLUX.1 model:		1028
973			1029
974	1. Create a cartoon image using Tom and Jerry in a detailed pose and expression. Tom, wearing his usual red shirt, is standing behind Jerry, who is dressed in his classic blue sweater. Both characters have a slightly confused look on their faces. Jerry is scratching his head while Tom looks away with a slight frown. Background: An indoor setting that resembles a cocktail party or garden tea event, with blurred figures of people in the background engaged in conversation. Style: Keep Tom and Jerry’s traditional animation style with bold lines and solid colors. Mood: Soft focus and warm lighting, suggesting an evening event.		1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043
975			1044
976			1045
977			1046
978			1047
979			1048
980			1049
981			1050
982			1051
983			1052
984			1053
985			1054
986			1055
987			1056
988			1057
989	2. Create a cartoon image using Bugs Bunny in a sitting pose, looking upward with a surprised or bewildered expression. - Background: A dimly lit room with a desk cluttered with various toys and a window showing a starry night sky. - Style: Retain the classic animation style		1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076 1077 1078 1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094
994			1095 1096 1097
		<b>B Dataset</b>	998
		<b>B.1 Dataset Characteristics</b>	999
		<b>B.2 Topic Distribution Analysis</b>	1000
		We applied enhanced weighted topic detection using Qwen-VL-Max (Bai et al., 2023) to conduct a comprehensive topic analysis across all 6,315 filtered memes, revealing fundamental differences in cultural priorities and humor focus.	1001 1002 1003 1004 1005
		<b>Chinese Meme Topics.</b> Table 4 presents the top 10 topics in Chinese memes, collectively covering 97.2% of the dataset.	1006 1007 1008
		<b>American Meme Topics.</b> Table 5 presents the top 10 topics in American memes, collectively covering 97.7% of the dataset.	1009 1010 1011
		<b>Cross-Cultural Topic Comparisons.</b> Table 6 highlights key differences in topic priorities between the two cultures.	1012 1013 1014
		Key cultural patterns revealed: (1) <i>Digital Concentration</i> —Chinese memes more heavily focused on internet/digital life (71.6% combined vs. 67.5% in US); (2) <i>Educational Values</i> —American memes treat education as casual daily experience (7.8%), Chinese memes reflect intense academic pressure (2.1%); (3) <i>Family Representation</i> —American memes more frequently feature family humor (4.9%) vs. Chinese hierarchical respect (1.9%); (4) <i>Leisure vs. Achievement</i> —American gaming culture prominent (2.7%), absent from Chinese top 10.	1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026
		<b>B.3 Emotion Distribution Analysis</b>	1027
		Using Qwen-based (Bai et al., 2023) automated emotion analysis, we classified all 6,315 memes according to Ekman’s six basic emotions, providing insights into cross-cultural emotional expression patterns.	1028 1029 1030 1031 1032
		<b>Chinese Meme Emotions.</b> Table 7 presents the emotion distribution in Chinese memes.	1033 1034
		<b>American Meme Emotions.</b> Table 8 presents the emotion distribution in American memes.	1035 1036
		<b>Cross-Cultural Emotion Comparisons.</b> Table 9 highlights key differences in emotional expression priorities.	1037 1038 1039
		Key emotional patterns: (1) <i>Positive Emphasis</i> —Both cultures prioritize joy, Americans showing slightly higher positive focus (73.8% vs. 69.3%); (2) <i>Sadness Acceptance</i> —Chinese memes	1040 1041 1042 1043

Table 4: Topic Distribution in Chinese Memes (N=3,165)

#	Topic	Count (%)	Cultural Significance
1	Internet Culture	1,931 (61.0%)	Digital lifestyle dominance, social media
2	Technology Digital	337 (10.6%)	Tech adaptation, AI integration
3	Work Career	216 (6.8%)	996 culture, career pressure
4	Social Relationships	148 (4.7%)	Friendships, social dynamics
5	Communication Language	126 (4.0%)	Language barriers, expression styles
6	Personality Psychology	115 (3.6%)	Individual traits, emotional responses
7	Education Learning	65 (2.1%)	Academic pressure, Gaokao system
8	Family Dynamics	61 (1.9%)	Family relationships, generational gaps
9	Animals Pets	46 (1.5%)	Pet culture, cute content
10	Entertainment Media	32 (1.0%)	Movies, shows, celebrity culture

*Top 10 Total: 3,077 memes (97.2%)*

Table 5: Topic Distribution in American Memes (N=3,150)

#	Topic	Count (%)	Cultural Significance
1	Internet Culture	1,651 (52.4%)	Social media, viral content, online trends
2	Technology Digital	475 (15.1%)	Tech innovation, digital lifestyle
3	Education Learning	247 (7.8%)	School experiences, college culture
4	Work Career	198 (6.3%)	Job market, work-life balance
5	Family Dynamics	155 (4.9%)	Family relationships, parenting
6	Communication Language	87 (2.8%)	Expression styles, conversation humor
7	Gaming Entertainment	85 (2.7%)	Video games, gaming culture, esports
8	Personality Psychology	67 (2.1%)	Individual psychology, personality types
9	Social Relationships	57 (1.8%)	Friendships, social interactions
10	Entertainment Media	54 (1.7%)	Movies, TV shows, celebrity content

*Top 10 Total: 3,076 memes (97.7%)*

Table 6: Cross-Cultural Topic Priority Comparison

Topic	CN	US	Interpretation
Internet Culture	#1 (61.0%)	#1 (52.4%)	Both dominant; China more concentrated
Technology Digital	#2 (10.6%)	#2 (15.1%)	US higher tech innovation focus
Education Learning	#7 (2.1%)	#3 (7.8%)	US: daily life; China: high-stakes pressure
Family Dynamics	#8 (1.9%)	#5 (4.9%)	US: frequent topic; China: serious element
Gaming Entertainment	-	#7 (2.7%)	US leisure vs. China work/study priority
Work Career	#3 (6.8%)	#4 (6.3%)	Similar priority, different intensity

Table 7: Emotion Distribution in Chinese Memes (N=3,165)

Emotion	Count (%)	Cultural Context
Joy	2,193 (69.3%)	Dominant positive humor expression
Anger	263 (8.3%)	Frustration, social critique
Sadness	258 (8.2%)	Melancholy, disappointment
Surprise	213 (6.7%)	Shock, unexpected situations
Fear	144 (4.5%)	Anxiety, worry
Disgust	94 (3.0%)	Revulsion, distaste

Table 8: Emotion Distribution in American Memes (N=3,150)

Emotion	Count (%)	Cultural Context
Joy	2,325 (73.8%)	Primary emotional expression
Fear	219 (7.0%)	Anxiety, relatable worries
Anger	217 (6.9%)	Frustration, social commentary
Surprise	148 (4.7%)	Unexpected, absurd humor
Disgust	140 (4.4%)	Cringe, distasteful situations
Sadness	101 (3.2%)	Disappointment, darker humor

Table 9: Cross-Cultural Emotion Priority Comparison

Emotion	CN	US	Interpretation
Joy	#1 (69.3%)	#1 (73.8%)	Both dominant; US slightly higher positivity
Anger	#2 (8.3%)	#3 (6.9%)	China: more direct frustration expression
Sadness	#3 (8.2%)	#6 (3.2%)	China: 2.5× higher melancholic acceptance
Surprise	#4 (6.7%)	#4 (4.7%)	Similar priority, China higher absurdist humor
Fear	#5 (4.5%)	#2 (7.0%)	US: anxiety culture, relatable worry themes
Disgust	#6 (3.0%)	#5 (4.4%)	US: higher cringe/distaste expression

express sadness 2.5× more frequently, reflecting cultural acceptance of melancholic humor; (3) *Anxiety Expression*—American memes emphasize fear-based content (7.0% vs. 4.5%), aligning with therapeutic humor trends; (4) *Anger Manifestation*—Chinese memes show higher anger (8.3% vs. 6.9%), possibly reflecting more direct emotional expression; (5) *Cringe Culture*—American memes display higher disgust representation (4.4% vs. 3.0%), consistent with cringe comedy trends.

**Impact of Joy Dominance on Experimental Design.** The overwhelming dominance of Joy in both datasets (69.3% Chinese, 73.8% American) has important implications for our transcreation experiments and evaluation: (1) *Positive Bias in Evaluation*—Since most transcreated memes will naturally preserve joyful emotions, our system may appear more successful at humor preservation simply due to the high baseline of positive content. This necessitates careful interpretation of intent preservation scores in Chapter ??; (2) *Limited Negative Emotion Testing*—With less than 30% of memes expressing negative emotions (anger, sadness, fear, disgust), our system receives limited training signals for adapting complex negative emotional tones across cultures, potentially underrepresenting challenges in transcreating emotionally nuanced content; (3) *Generalizability Concerns*—The skewed distribution means our findings may generalize better to lighthearted, positive meme content than to darker, satirical, or critical humor styles; (4) *Cultural Authenticity vs. Emotional Consistency*—The Joy dominance simplifies one aspect of transcreation (emotional tone transfer) while placing greater emphasis on cultural reference adaptation as the primary challenge. Despite this limitation, the Joy-dominant distribution accurately reflects real-world meme ecosystems where positive, shareable content naturally dominates social media platforms—making our experimental conditions ecologically valid even if not emotionally balanced.

## C Evaluation Metrics

### Quantitative Dimensions (1-5 scale):

**Caption Quality** Evaluates whether the generated caption works effectively as meme text, considering clarity, readability, appropriate meme language/tone, engaging phrasing, and proper text formatting.

**Image Quality** Assesses whether the generated image functions effectively as a meme visual, considering visual clarity and quality, appropriate meme composition, recognizable elements/characters, and visual appeal and memorability.

**Synergy** Measures how well image and caption work together, evaluating coherent message delivery, emotional or humorous impact, logical connection between visual and text, and overall meme effectiveness.

**Cultural Fit** Evaluates cultural adaptation quality, including alignment with target culture’s humor style, appropriate cultural references, target audience relatability, and avoidance of cultural misunderstandings.

**Intent Preservation** Assesses preservation of the original meme’s intent, including message consistency, emotional tone preservation, humorous effect maintenance, and core meaning retention.

**Overall Score** The average of all dimension scores and reflects overall quality.

## D Human Evaluator Profiles

Our evaluation employed three bilingual, bicultural evaluators with a deep understanding of both Chinese and US cultural contexts:

**Evaluator 1.** Native Chinese speaker with 10+ years US residence, PhD in Communication Studies. Regular engagement with both Weibo/Xiaohongshu and Reddit meme communities. Assessment style: Entertainment-focused, generous scoring emphasizing humor effectiveness over technical perfection. Mean overall score:

1127 4.42/5.0.

1128 **Evaluator 2.** American-born Chinese with  
1129 native-level proficiency in Mandarin and US, MA  
1130 in Comparative Cultural Studies. Active partic-  
1131 ipation in both Chinese and US digital cultures.  
1132 Assessment style: Balanced and objective, apply-  
1133 ing consistent standards across dimensions. Mean  
1134 overall score: 4.09/5.0. Showed highest correlation  
1135 with Qwen-VL-Max (F1 = 0.925,  $r = 0.964$ ).

1136 **Evaluator 3.** Native Chinese speaker with 12+  
1137 years of US experience, professional translator with  
1138 meme localization background. Expertise in cul-  
1139 tural adaptation nuances. Assessment style: Criti-  
1140 cal and quality-focused, emphasizing cultural au-  
1141 thenticity and linguistic precision. Mean overall  
1142 score: 3.31/5.0.

1143 All evaluators received identical structured  
1144 prompts specifying six evaluation dimensions,  
1145 worked independently without access to others' rat-  
1146 ings, and maintained consistency through detailed  
1147 scoring rubrics. Inter-evaluator correlations demon-  
1148 strate moderate to strong agreement: Evaluator 1-2  
1149 ( $r = 0.72$ ), Evaluator 1-3 ( $r = 0.58$ ), Evaluator  
1150 2-3 ( $r = 0.81$ ), confirming reliable yet stylistically  
1151 distinct evaluation perspectives.

## 1152 E VLM Evaluator Details

1153 Six VLMs served as automated evaluators, se-  
1154 lected for multilingual (Chinese-English) capabil-  
1155 ity, multi-image processing, and reproducibility:

- 1156 • **Qwen-VL-Max** (Alibaba Cloud): Commer-  
1157 cial API with extensive Chinese-English train-  
1158 ing, demonstrated exceptional human correla-  
1159 tion ( $r = 0.926$ ).
- 1160 • **LLaVA-v1.6-Vicuna-13B**: Same architecture  
1161 as transcreation system, showed no meaning-  
1162 ful correlation ( $r = 0.005$ ).
- 1163 • **InternVL3-8B/14B**: Recent open-source  
1164 models with strong vision capabilities,  
1165 achieved weak positive correlation (8B:  $r =$   
1166  $-0.049$ , 14B:  $r = 0.263$ ).
- 1167 • **Qwen3-VL-8B-Instruct**: Smaller Qwen vari-  
1168 ant, weak correlation ( $r = 0.252$ ).
- 1169 • **Aya-vision-8b**: Massively multilingual model,  
1170 slight negative correlation ( $r = -0.043$ ).

1171 All LLMs received identical prompts specifying  
1172 evaluation dimensions and rating scales. Tempera-  
1173 ture set to 0.7 for balanced consistency-creativity  
1174 tradeoff.

## F Transcreation Prompts 1175

### 1176 Stage 1 (LLaVA 1.6) - Cultural Analysis Prompt

#### 1177 Example:

1178 *You are a cultural adaptation expert. Analyze*  
1179 *this [SOURCE CULTURE] meme and create a*  
1180 *transcreated version for [TARGET CULTURE]*  
1181 *audiences. Your response should include:*

1182 1. *Cultural Context Analysis: Identify culture-*  
1183 *specific elements (references, idioms, visual sym-*  
1184 *bols, humor mechanisms)*

1185 2. *Intent Extraction: What is the core mes-*  
1186 *sage/emotion/joke?*

1187 3. *Target Culture Mapping: Find equivalent con-*  
1188 *cepts in [TARGET CULTURE]*

1189 4. *Transcreated Caption: Generate a new cap-*  
1190 *tion preserving intent while using [TARGET CUL-*  
1191 *TURE] appropriate references and style*

1192 5. *Visual Recommendations: Describe ideal vi-*  
1193 *sual template (characters, setting, composition)*  
1194 *culturally appropriate for [TARGET CULTURE]*

### 1195 Stage 2 (FLUX.1) - Visual Generation Prompt

#### 1196 Example:

1197 *Create a meme-style image: [LLaVA's visual*  
1198 *recommendations]. Style: internet meme, high*  
1199 *contrast, recognizable characters, clear com-*  
1200 *position suitable for text overlay. [TARGET*  
1201 *CULTURE]-appropriate visual elements. Res-*  
1202 *olution: 1024x1024px.*

1203 Full prompt templates with examples will be  
1204 made available in our public repository upon ac-  
1205 ceptance.

## G Example Transcreations 1206

### 1207 Success Example - US→Chinese:

1208 *Source (US): "Nobody: Absolutely nobody: Me*  
1209 *at 3 am:" [Image: Person raiding refrigerator]*

1210 *Transcreated (Chinese): "深夜两点的我:" (Me*  
1211 *at 2 am) [Image: Cartoon cat staring at food]*

1212 *Adaptation rationale:* Replaced human figure  
1213 with animal imagery (preferred in Chinese memes),  
1214 adjusted time (2 am vs 3 am reflects Chinese sleep  
1215 patterns), simplified narrative structure for concise-  
1216 ness.

### 1217 Challenge Example - Chinese→US:

1218 *Source (Chinese): "内卷" (involution) concept*  
1219 *with study-exhausted imagery*

1220 *Transcreated (US): "The grind never stops" [Of-*  
1221 *fice worker imagery]*

1222 *Limitation:* US lacks a precise equivalent for "内  
1223 卷" (intensifying competition in zero-sum environ-  
1224 ments). "Grind culture" captures work intensity but  
1225 misses the systemic competition aspect, illustrating  
1226 cultural untranslatability challenges.

Additional examples and failure case analysis are available in the supplementary materials.

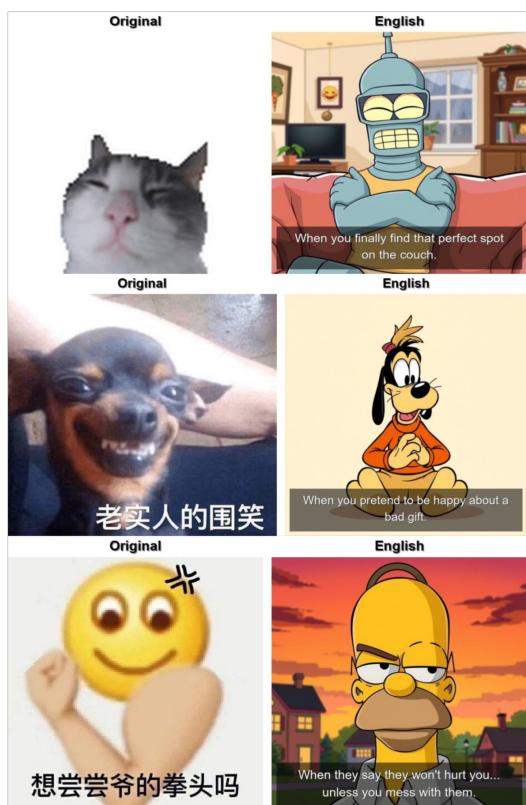


Figure 4: Examples of successful Chinese→US meme transcreations.

**Original (middle left — dog meme):** “Honest smile”  
**Original (bottom left — angry emoji meme):** “You looking for a knuckle sandwich?”



Figure 5: Examples of successful US→Chinese meme transcreations.

**Transcreated (top right — spongebob meme):** “Dad thinks I’m up early studying, I’m really on my 14th straight gaming session”

**Transcreated (middle right — spongebob meme):** “Kid: Mom, I’m playing a game  
 Mom: I’m cooking  
 Kid: Can you pause it?  
 Mom: How dare you use my own teachings against me?”

**Transcreated (middle right — Spencer Wright meme):** “Tiktok be spamming ads upfront, still not gonna pay for premium”