Do LLMs Understand Wine Descriptors Across Cultures? A Benchmark for Cultural Adaptions of Wine Reviews

Anonymous ACL submission

Abstract

Leveraging remarkable advancements in Large Language Models (LLMs), we are now poised to tackle increasingly complex challenges requiring deep comprehension of multifaceted domains and contexts. A specific application scenario is wine reviews adaptation. Wine reviews usually describe a wine's appearance, aroma, and flavor to help consumers appreciate its characteristics. However, the adaptation of wine reviews transcends mere translation; it requires consideration of regional preferences, 011 flavor descriptors, and cultural nuances that 012 shape wine perception. We introduces the firstever task involving the translation and cultural adaptation of wine reviews between Chinese and English. In a case study on cross-cultural wine review adaptation, we compile a dataset of 8k Chinese and 16k Western professional wine 019 reviews. We evaluated various methods, including LLMs and traditional machine translation techniques, using both automatic and human metrics. For human assessments, we introduce three novel cultural-related metrics-Cultural Proximity, Cultural Neutrality, and Cultural Genuineness-to gauge the success of different approaches in achieving authentic crosscultural adaptation. Our analysis shows that current models struggle to capture cultural nu-029 ances, especially in translating wine descriptions across different cultures. This highlights the challenges and limitations of translation models in handling cultural content.

1 Introduction

Wine reviews serve as valuable guides for consumers, offering detailed insights into the characteristics of each bottle. For casual drinkers and connoisseurs alike, these reviews act as a compass, helping them navigate the vast selection of wines available. However, due to cultural influences and geographical distinctions, beverage consumption patterns and individual preferences vary significantly across regions (Rodrigues and Parr, 2019),



Figure 1: An example of Literal vs. Cultural Adapted Translation, Enhancing Readability in Wine Descriptions, highlighting how certain descriptors and flavor terms require adaptation for better comprehension by culturally unfamiliar readers.

043

044

045

046

047

049

052

058

060

061

062

063

064

065

067

068

not to mention reviews. These variations extend beyond mere differences in taste and are deeply rooted in the cultural, social, and environmental contexts of each region. Consequently, consumer preferences for certain beverages, including wine, can differ greatly, often rendering generalized reviews less relevant or applicable across diverse audiences. Professional reviews play a greater role than user reviews in promoting consumer purchases (Chiou et al., 2014). For Chinese wine consumers, professional wine reviews in Chinese are scarce, and most available reviews require a paid subscription. Similarly, Western consumers face challenges finding professional reviews of Chinese wines. This paper aims to bridge these gaps by providing a comprehensive, culturally inclusive dataset of bilingual wine reviews, supporting a more globally relevant perspective on wine preferences.

Recognizing and adapting to cultural differences in language use is both essential and challenging (Hershcovich et al., 2022), especially for subjective comments. Translations of reviews using current neural machine translation systems may overlook culture-specific expressions or result in mistranslations due to insufficient grounding in physical and cultural contexts. For instance, in Figure 1, 'raspberry', a common European flavor descriptor, is rare in China ('覆盆子'), making its taste unfamiliar to Chinese consumers. The flavor of raspberry is puzzling to Chinese consumers, and requires additional explanation or finding a fruit with a similar flavor. 'Blueberry,' which has a similar flavor profile (Jin et al., 2022), is more popular and better understood by Chinese consumers. All the outputs of this example are shown in Appendix H.

070

071

079

081

087

089

095

101

102

103

104

105

106

107

108

109 110

111

112

While wine reviews traditionally rely on human expertise to capture nuanced sensory experiences, providing detailed descriptions of a wine's appearance, aroma, and flavor. However, these reviews are deeply influenced by cultural norms, linguistic styles, and regional preferences, making their translation across languages a complex task. Recent advancements in LLMs offer new opportunities to analyze and generate detailed reviews (Wu et al., 2025). By leveraging their ability to parse intricate descriptions and contextual nuances, LLMs provide an opportunity to analyze how cultural nuances and stylistic elements are preserved or transformed during translation.

In this work, we introduce the task of adapting wine reviews across languages and cultures. Beyond direct translation, this requires adaptation concerning content and style. We focus on Chinese and English wine reviews, automatically pairing reviews for the same wine from the same vintage from two monolingual corpora. As there are many reviews in English for the same wine, we also explore the inner difference in reviews from people of the same culture. We evaluate our methodology with human evaluation and automatic evaluations on the dataset we construct. Our contributions are as follow:

- We introduce the task of cross-cultural wine reviews translation and build a bidirectional Chinese-English dataset with multiple references for it: CulturalWR.
- We experiment with various sequence-tosequence approaches to adapt the reviews, including machine translation models and multilingual LLMs.
- 1133. We evaluate and analyze the difference be-
tween Chinese and English-speaking cultures
and how they describe wine characteristics.

2 Related Work

Computational analysis of wine reviews. Recent advancements in wine informatics have leveraged computational techniques to analyze expert wine reviews. The Computational Wine Wheel 2.0 facilitates machine learning-based wine attribute analysis (Chen et al., 2016). In addition, studies applying SVMs to wine reviews have demonstrated the impact of different review sources on highquality wine prediction (Tian et al., 2022). Machine learning and text mining techniques have also been utilized to discover predictive patterns in wine descriptions, challenging the notion that flavor descriptions are purely subjective (Lefever et al., 2018). These studies provide a foundation for computational analysis of wine reviews, yet they primarily focus on prediction tasks rather than cross-cultural aspects of wine appreciation.

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

Cross-cultural analysis of flavors. Flavor perception varies across cultures, as shown in studies analyzing beer pairing preferences in Latin America (Arellano-Covarrubias et al., 2019) and color-flavor associations in snack packaging across China, Colombia, and the UK (Velasco et al., 2014). Research on cheddar cheese flavor lexicons across different countries (Drake et al., 2005) further highlights cultural influences on taste perception. These findings underscore the need for culturally adaptive approaches in wine description and translation.

Cultural adaptation. As culture and language are intertwined and inseparabled, there is a rising demand to equip machine translation systems with greater cultural awareness (Nitta, 1986; Ostler, 1999; Hershcovich et al., 2022). However, it is costly and time-consuming to collect culturally sensitive data and perform a human-centered evaluation (Liebling et al., 2022). A recent study showed that LLMs are adept at adapting cooking recipes across cultures, using a comparable and a parallel corpus (Cao et al., 2024). Our work focuses on the translation of non-parallel subjective reviews, an area that remains underexplored.

3 The CulturalWR dataset

We present CulturalWR, a dataset of paired pro-
fessional wine reviews. Each pair consists of one159Chinese review by Chinese wine critics and, when
available, an English review by the same authors,
alongside multiple English reviews authored by161

165

166

167

168

- 169
- 171
- 172 173
- 174
- 175

176

191

192

193

194

195

198

199

3.1 Data Collection

We collect the English wine reviews from several professional wine review websites, including Robert Parker's Wine Advocate¹, Wine Spectator², James Suckling³, Wine Spectator⁴ and some other professional wine review websites. Chinese wine reviews are primarily written by two prominent reviewers, AlexandreMa⁵, ShenHao⁶, who often publish bilingual wine reviews and China Wine Information Network⁷.

Western wine critics for the same wine, identified

by the wine name and its corresponding vintage.

3.2 Wine Matching Rules

Wine names are usually derived from regions or 177 grape varieties, labeled by these features or a 178 unique name. We observed that Chinese reviewers 179 180 often skip mentioning the grape variety and region, opting for simpler names. To match reviews to spe-181 cific wines, we use a string matching algorithm that 182 includes converting non-English characters from languages like Spanish and French to English. This 184 185 addresses variations in spelling, such as "Château Nénin" versus "Chateau Nenin". Exact matches are assumed to indicate the same wine; partial matches 187 undergo manual verification.⁸ After confirming the names and vintage years match, we finalize the reviews for each specific wine.

3.3 Data Filtering Rules

To focus on red wines and adapt reviews culturally, we applied three filtering rules: excluding Non-Vintage (N.V.) wines due to inconsistent tasting notes over time, separating white and rosé wines into distinct datasets for independent testing, and filtering out reviews under 30 words in English or 30 characters in Chinese for lack of detail.

3.4 Dataset Overview

We collect approximately 20k Chinese wine reviews covering nearly 20k wines, along with 150,000 English wine reviews spanning over 30k

wines. These reviews encompass a variety of wine types, including red, white, rosé, and champagne. After filtering, we retain around 10k Chinese reviews focused on red wines. Through matching, the dataset is refined to include 4.5k wines, comprising about 4.5k Chinese reviews and 16k English reviews. 3,227 Chinese reviews have their matched English reviews written by the same author. Data statistics are shown in Table 1.

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

	Number	Mean #Tokens
CA Chinese Reviews	4776	67.57
CA English Reviews	3227	74.25
Transl. Chinese Reviews	60	60.2
WA English Reviews	16746	58.16

Table 1: Statistics of reviews. CA refers to Chinese wine critics and WA to Western ones. We count tokens with jieba text segmentation for Chinese and whitespace tokenization for English.

Attributes. Besides the basic reviews and their corresponding rating for each wine, we sourced wine data from different sources, we reorganize the basic attributes of all these wines, which includes the geographical location of the winery, the composition of the grape varietals, the vintage, the alcohol content and the price. To ensure privacy, we anonymized the obtained data, and no personally identifiable information is included in the attributes. It's shown in F.

Analysis of Chinese Reviews and 4 Western Reviews

We frame three insights gained in this section, that not only reveal that professional wine reviews are confusing to ordinary consumers, but also show the cultural similarities and differences between professional reviews. These highlight the necessity for cultural adaptations of wine reviews.

Insight #1: Wine reviews are not always intuitive to consumers. While most reviews are relatively easy to understand, some flavor descriptors-such as 'Leather', 'Tar' and 'Cat's Pee'-are confusing or unappealing to those unfamiliar with wine terminology. However, they are widely used by Western reviewers to describe red wines. This gap between expert notes and reader intuition can make reviews feel inaccessible.

Insight #2: High Semantic Similarity Between 239 Chinese and Western Wine Critics' Reviews. 240

¹https://www.robertparker.com/

²https://www.winespectator.com/

³https://www.jamessuckling.com/

⁴https://www.winespectator.com/

⁵https://www.alexandrema.com

⁶http://www.leparadisduvin.com

⁷www.winesinfo.com

⁸Particularly in the Bordeaux region, wines use the winery's name for the Grand Vin (first wine) and unique names for second and third wines, with "blanc" added for white wines.

We extracted detailed flavor descriptors from wine 241 reviews and used Jaccard Similarity to analyze 242 cultural differences. Leveraging the Wine Aroma 243 Wheel from Aromaster⁹, we measured both inner 244 and outer similarity across three hierarchical layers: 245 Aroma Families (broad categories such as fruits, 246 flowers, and spices), Aroma Subfamilies (more spe-247 cific groups like citrus fruits, red berries, and dried herbs), and Exact Aromas (precise descriptors such as lemon, raspberry, or thyme). 250

252

253

257

261

262

263

271

274

275

276

Although the Wine Aroma Wheel includes 88 commonly used wine aromas, we identified over 300 distinct precise descriptors in the reviews. For Exact Aromas, both inner and outer similarities¹⁰ are below 0.1, with outer similarity significantly lower than inner similarity. At the Aroma Subfamily level, outer similarity increases to 0.16, catching up with inner similarity. For Aroma Families, outer similarity exceeds 0.4, indicating some degree of convergence. However, overall Jaccard similarity remains relatively low. We attribute this to differences in reviewing styles: some reviewers describe only dominant flavors in complex red wines, while others list all detectable aromas.

Additionally, we evaluated cosine similarity and BERTScore (Zhang* et al., 2020) for cross-group and within-group comparisons. Both metrics consistently exceed 0.8, suggesting a high degree of semantic similarity between Chinese and Western wine reviews.

Insight #3: Chinese Reviews have a significantly distinct boundary compared to the overall distribution of Western Reviews. Different from Insight #2, we do dimensionality reduction using PCA on the embeddings¹¹, as shown in Figure 2, the two sets of embeddings exhibit a clear boundary in the lower-dimensional space. This suggests that the two groups of comments have substantial differences in their global structure and overall semantic distribution. PCA amplifies these differences, uncovering separation that is not evident through BERTScore's local similarity focus. This also shows the importance of cultural adaptations.



Figure 2: PCA-reduced embeddings: Red dots show Western reviews, and blue dots show English reviews by Chinese authors

5 Cross Cultural Wine Reviews Adaptation Task

We propose cross-cultural wine review adaptation, extending machine translation by requiring both accuracy and deliberate semantic divergence to address cultural differences. 285

286

287

290

291

292

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

Evaluating cultural adaptation is challenging, as it must balance meaning preservation with genuine cross-cultural differences. In wine review adaptation, we even need to adapt the flavor descriptors. As common in text generation tasks, we first adopt reference-based automatic evaluation metrics. Moreover, considering reference-based metrics are often unreliable for subjective tasks, we also conduct human evaluations.

5.1 Automatic Evaluation

We use four metrics to assess the similarity between the generated and reference reviews. We use two lexical-based metric: BLEU (Papineni et al., 2002), a precision metric based on token n-gram which emphasizes precision and commonly used in machine translation evaluation and ME-TEOR (Banerjee and Lavie, 2005), which combines precision and recall while incorporating linguistic features such as stemming and synonymy to provide a more comprehensive evaluation; one contextual-embedding based metric: BERTScore (Zhang* et al., 2020), based on cosine similarity of contextualized token embeddings and capture deep semantic matching; one hybrid-based metric: Beer (Stanojević and Sima'an, 2014), based on multi-feature fusion regression indicator, which combines syntactic and semantic features to automatically evaluate translation quality through regression model learning.

⁹https://aromaster.com/

¹⁰ 'Outer'' refers to comparisons between Chinese and Western reviews, while "Inner" refers to comparisons within Western reviews.

¹¹We obtain embeddings from the last hidden layer's hidden states from ChatGLM

323

324

326

327

332

339

341

342

345

347

351

353

354

357

5.2 Human Evaluation

While automatic metrics provide quantifiable results, they rely on fixed reference sets, which may lack cultural relevance. To address this, we introduce seven human evaluation criteria applied to the test set.

(1) Grammar: The generated reviews are grammatically correct and fluent; (2) Faithfulness of Information: The content accurately reflects the original input without introducing false or misleading information; (3) Faithfulness of Style: The output preserves the original tone, register, and formality without altering the intended voice; (4) Overall quality: The reviews are coherent, contextually appropriate, and align with the intended tone and style; (5) Cultural proximity: The generated reviews use familiar terms and expressions that resonate with the target culture. For example black currant was replaced by Chuanbei loquat paste, cough syrup and hawthorn cake in Chinese localised aroma wheel (Jin et al., 2022); (6) Cultural Neutrality: Maintains neutrality to avoid provoking negative perceptions or reactions from the target culture consumers. For example, 'earthy' can be a positive descriptor for wine in the West. However, when translated into Chinese, '土味' often implies a dirty or unrefined flavor. A more elegant term like '泥土气息' (aroma of soil) is often preferred; (7) Cultural Genuineness: Preserves the quality of the original descriptor without altering its meaning, ensuring authenticity. For example clove, violet and saffron have no suitable local descriptors with similar olfactory characteristics.

> Our evaluation was done by three people fluent in both Chinese and English, including two master students and a professor. Before the evaluation process, we performed preliminary testing on 40 samples and used Pearson correlation to calculate their understanding of different metrics which proved these metrics are easy for people to evaluate.

6 Experiment

361To comprehensively assess the efficacy of LLM362translations in understanding wine and flavors, we363compare various prompting strategies on tuning-364free LLMs, alongside evaluations of an open-365source Machine Translation model.

6.1 Experimental Setup

Prompting LLMs. Based on the exceptional performance of multilingual LLMs in zero-shot translation. We explore their ability on translating wine reviews and flavor adaptation.

We compare the performance of five diverse, state-of-the-art LLMs on this task. We test Llama-3.1-8B-Instruct (AI@Meta, 2024), Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), Phi-3.5-miniinstruct (Abdin et al., 2024), and ChatGPT-4o (OpenAI et al., 2024) and two use more Chinese training data models: Qwen2.5-7B-Instruct (Yang et al., 2024), GLM4-9b (GLM et al., 2024). We used Cultural Prompt in Table 9.

Multilingual machine translation model: We use the state-of-the-art NLLB-200-3.3B model (Team et al., 2022) for accurate, context-aware multilingual translation in our experiments.

7 Results and Analysis

Our analysis includes five parts: 1)automatic evaluation between different models; 2) fine-grained human evaluation on a subset of Chinese-English translation bidirectional; 3) Correlation of automatic metrics with humans; 4) Different prompting strategy evaluation comparison; 5) Quantitative analysis for some specific concepts

7.1 Overall Automatic Evaluation

Models	BLEU	METEOR	B-Sc	BEER	#Tok			
$\mathbf{Chinese} ightarrow \mathbf{English}$								
ChatGLM4	18.7	45.7	86.7	51.5	65.8			
Phi	10.9	40.6	90.0	47.5	75.6			
Qwen2.5	15.6	<u>46.3</u>	<u>91.2</u>	<u>52.7</u>	69.2			
Mistral	5.2	36.4	87.9	34.1	61.8			
Llama3.1	11.4	35.2	88.0	44.3	54.0			
NLLB	12.0	36.8	88.7	48.1	64.8			
	I	English ightarrow Ch	inese					
ChatGLM4	8.9	<u>31.9</u>	86.6	26.5	130.5			
Phi	4.8	25.8	89.8	22.5	94.7			
Qwen2.5	12.3	36.4	<u>90.6</u>	<u>29.3</u>	85.7			
Mistral	2.0	20.1	89.3	17.0	56.1			
Llama3.1	9.9	<u>31.9</u>	87.3	28.1	89.5			
NLLB	3.8	18.3	87.6	21.4	96.8			

Table 2: Automated evaluation results on the test sets using reference-based metrics: BLEU, METEOR, B-Sc(BERTScore) and BEER. Higher scores indicate better performance on all metrics.

For Chinese-to-English translation, ChatGLM4 achieved the highest BLEU score, suggesting strong lexical matching, while Qwen2.5 excelled in BERTScore and BEER, indicating superior 366

368

369

370

372

221

390

388

392

394 395 396

semantic alignment with human-written transla-397 tions. For English-to-Chinese translation, Mistral 398 led in BLEU, while Llama3.1 scored highest in 399 METEOR, and Qwen2.5 again outperformed in 400 BERTScore and BEER, reinforcing its strength in 401 preserving meaning. Notably, ChatGLM4 strug-402 gled with BLEU and METEOR in this direction but 403 maintained competitive BERTScore values, sug-404 gesting it prioritizes semantic coherence over strict 405 lexical overlap. Additionally, the NMT system 406 NLLB still remains highly competitive. Overall, 407 Qwen2.5 consistently demonstrated high semantic 408 quality across both directions, while other models 409 exhibited strengths in specific metrics, reflecting 410 differing optimization objectives. These results 411 highlight the inherent trade-offs between fluency, 412 lexical fidelity, and semantic preservation across 413 models. More importantly, they underscore that ref-414 erence translations are not the sole "correct" adapta-415 tions, as translation quality is inherently subjective. 416 This reinforces the need for a nuanced evaluation 417 framework that accounts for cultural context. lin-418 guistic variation, and domain-specific preferences 419 to better capture real-world translation quality. 420

7.2 Human Evaluation

Models	F-I	F-S	Gr	0-Q	C-P	C-G	C-N	
$\mathbf{Chinese} ightarrow \mathbf{English}$								
ChatGLM4	5.6	4.8	5.3	5.0	6.8	<u>6.4</u>	6.2	
Phi	4.8	5.4	6.0	4.8	6.9	6.3	<u>6.4</u>	
Qwen2.5	<u>5.8</u>	<u>6.0</u>	5.8	<u>5.5</u>	6.8	6.3	<u>6.4</u>	
Mistral	4.8	5.7	6.0	4.8	<u>7.0</u>	<u>6.4</u>	6.3	
Llama3.1	5.3	5.6	<u>6.2</u>	5.3	<u>7.0</u>	6.2	6.3	
Human	5.5	5.3	5.5	5.3	6.6	6.3	6.2	
	En	glish -	→ Chi	nese				
ChatGLM4	5.5	5.5	4.2	4.7	5.6	5.5	5.3	
Phi	4.1	4.5	3.6	3.7	5.6	4.4	5.3	
Qwen2.5	<u>5.7</u>	<u>5.9</u>	<u>5.5</u>	<u>5.6</u>	<u>6.0</u>	5.5	<u>5.8</u>	
Mistral	3.8	4.4	2.8	3.1	4.8	4.6	4.7	
Llama3.1	4.5	5.0	4.4	4.6	5.9	4.8	5.7	
Human	5.1	5.8	5.2	5.0	5.8	<u>5.7</u>	<u>5.8</u>	

Table 3: Human evaluation results on the selected test sets: average for each method and metric, ranging from 1 to 7. F-I, F-S, Gr, O-Q, C-P, C-G and C-N representing Faithful of Information, Faithful of Style, Grammar, Overall Quality, Cultural Proximity, Cultural Genuineness and Cultural Neutrality respectively

Table 3 presents the results of human evaluation across multiple dimensions. Notably, NLLB was excluded from human evaluation due to its lack of cultural adaptation capabilities.

For English-to-Chinese translation, Qwen2.5 leads across most metrics, even surpassing human

translations in faithfulness, grammar, and overall quality. Human translations score highest in Cultural Genuineness, reflecting their ability to preserve nuanced expressions. Llama3.1 also remains competitive, balancing fluency and cultural adaptation. 428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

For Chinese-to-English translation, Qwen2.5 again ranks highest in faithfulness and overall quality, while Mistral outperforms even human translations in Cultural Proximity and Genuineness, suggesting a stronger emphasis on natural, idiomatic English.

These results highlight trade-offs between literal accuracy and cultural adaptation. While human translations excel in cultural authenticity, LLMs like Qwen2.5 demonstrate strong faithfulness, and Mistral prioritizes fluency in the target language. This underscores the importance of context-aware evaluation that considers both linguistic accuracy and cultural nuances.

Methods	F-I	F-S	Gr	0-Q	С-Р	C-G	C-N
$\mathbf{Chinese} ightarrow \mathbf{English}$							
Human-Eval	5.36	5.58	5.91	5.19	6.86	6.31	6.36
GPT4o-Eval	5.79	5.24	6.47	5.77	5.23	5.4	6.44
$\mathbf{English} \to \mathbf{Chinese}$							
Human-Eval	4.57	5.02	4.44	4.13	5.25	5.12	5.0
GPT4o-Eval	5.76	5.36	6.2	5.71	5.39	5.58	6.43

Table 4: GPT-40 v.s.Human in Human evaluationmetrics(Average over All Human Test Cases)

GPT-40 and Human Evaluation Diverge. Table 4 shows the results evaluated by both GPT-40 and human annotators. Specifically, we use all the human evaluation test cases for GPT evaluation. The results show that GPT and Human evaluation scores are not strongly correlated and GPT has a higher tolerance than humans especially for English \rightarrow Chinese direction. This discrepancy suggests that GPT-40 may have a different interpretation of translation quality compared to human evaluators. This discrepancy is particularly evident in culturally relevant metrics. The prompts we used are shown in C.

7.3 Correlation of automatic metrics with humans

To evaluate the reliability of automatic metrics for wine review adaptations, we analyze their correlation with human evaluations across seven metrics using Kendall correlation, the WMT22 metaevaluation standard (Freitag et al., 2022).

As illustrated in Table 5, the correlation between

427

422

	BLEU	METEOR	B-Sc	BEER			
$\mathbf{Chinese} ightarrow \mathbf{English}$							
F-I	0.2536*	0.1892	0.3000*	0.2442*			
F-S	0.0053	0.0075	0.0204	0.0113			
Gr	-0.0296	-0.0096	0.0033	-0.0478			
0-Q	0.2191*	0.1847*	0.2061*	0.2015*			
C-P	-0.070	-0.0177	-0.0540	-0.0961			
C-G	0.1131	0.0625	0.0621	0.1131			
C-N	-0.1166	-0.0841	-0.0166	-0.0876			
		English \rightarrow C	hinese				
F-I	0.4079*	0.2788*	0.4080*	0.2946*			
F-S	0.3723*	0.3560	0.3769*	0.3904*			
Gr	0.2819*	0.2788*	0.3530*	0.2946*			
O-Q	0.3526*	0.3123*	0.3742*	0.3557*			
C-P	0.1408	0.1524	0.2609*	0.1553			
C-G	0.3134*	0.3170*	0.3942*	0.3170*			
C-N	0.1334	0.1357	0.2389*	0.1516			

Table 5: Kendall correlation of human evaluation results with automatic metrics. Statistically significant correlations are marked with *, with a confidence level of $\alpha = 0.05$ before adjusting for multiple comparisons using the Bonferroni correction

human evaluation results and automatic metrics 469 varies across different translation directions and 470 evaluation criteria. For Chinese \rightarrow English, F-I 471 (Faithful of Information) and O-Q (Overall Qual-472 473 ity) exhibit the strongest correlations, particularly with BLEU, B-Sc, and BEER, suggesting that these 474 metrics align well with human judgments in assess-475 ing fluency and overall translation quality. On the 476 other hand, for English-Chinese, the correlations 477 478 are generally stronger across all metrics. Notably, F-I, F-S (Faithful of Style), and O-Q display signifi-479 cant correlations with multiple metrics, particularly 480 B-Sc and BEER, indicating that these metrics are 481 relatively more reliable for evaluating fluency and 482 483 intelligibility in English-to-Chinese translations.

C-G (Cultural Genuineness) also achieves strong 484 correlations indicating that automatic metrics can 485 reliably assess translation accuracy. However, this 486 does not necessarily reflect the cultural adaptability 487 of the translation. However, C-N(Cultural Neutral-488 ity) and C-P(Cultural Proximity) remain weakly 489 correlated, revealing that automatic metrics still fall 490 short in capturing deeper cultural nuances, empha-491 sizing the need for human evaluation. Notably, cor-492 493 relations for English→Chinese generally exhibit greater strength than Chinese→English. This dis-494 crepancy is likely due to most wine reviews being 495 written from a predominantly Western reviewer's 496 perspective. 497

Human Evaluation							
Methods	F-I	F-S	Gr	0-Q	C-P	C-G	C-N
Direct Translation	6.38	5.94	5.56	5.69	4.38	6.31	5.81
Cultural Prompt	6.25	6.12	5.88	5.94	4.5	6.19	5.94
Detailed Cultural Prompt	5.75	5.94	5.75	5.56	4.88	5.69	5.69
Self-Explanation	4.94	5.75	5.88	5.5	5.75	4.56	6.62
Automatic Evaluation							
Methods		BLEU	ME	TEOR	B-9	Sc B	EER
Direct Translation		16.5	4	11.6	91	.4 .	28.3
Cultural Prompt		15.9	41.0		91	.3 .	28.3
Detailed Cultural Prompt		14.4	3	38.6	91	.0 2	27.6
Self-Explanation		13.7	3	37.0	90	.9 2	27.6

Table 6: Evaluation of different strategies by GPT40 on English-Chinese translations.

498

499

500

501

502

503

504

505

506

507

508

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

7.4 Prompting Strategy Evaluation

With LLM-based machine translation advancing, integrating free-form external knowledge offers new opportunities to enhance translation quality. We compare different prompting strategies on Chat-GPT for English-to-Chinese translation, including Direct Translation, Cultural Prompt, Detailed Cultural Prompt, and Self-Explanation. Table 9 lists the specific prompts.

As is shown in Table 6, Direct Translation is the best-performing method in Faithful of Information and automatic metrics. Cultural Prompting improves cultural accuracy slightly but does not significantly enhance overall translation quality. Self-Explanation has the worst faithful scores but is rated the best in Cultural Neutrality, making it a trade-off strategy for culturally rich contexts. This trade-off suggests that translation strategies need to balance faith, accuracy, and cultural adaptability, depending on the intended use case.

7.5 Quantitative analysis

Cross-lingual translation of culturally specific concepts is challenging, as it requires balancing linguistic accuracy with cultural adaptation. Many models tend to rely on literal translation, which may not always convey the intended meaning naturally. To examine this, we evaluate each model's literal translation rate for culturally embedded flavor descriptors from the CulturalWR test set. For instance, in English-to-Chinese translation, 'thyme' is considered an English-specific concept. We count occurrences of related terms such as 'gooseberry', 'thyme', and 'rosemary' in English wine reviews (c_{source}) and record how often they are directly translated in the corresponding Chinese reviews (c_{target}) from model predictions. The literal translation rate is then calculated as $\frac{c_{target}}{c_{source}}$. To further assess translation quality, we conduct a

bidirectional test on the five most common winerelated terms. This ensures that models not only preserve culturally specific terms when translating in one direction but also effectively map key wine descriptors between languages. Comparing accuracy across models provides insights into their ability to maintain semantic fidelity, crucial for expert-level wine translations.

536

537

540

541

542

545

547

548

549

552

554

560

562

564

566

571

574

576

583

584

587

As shown in Figure 3a, we analyze six culturally relevant concepts, three common in Chinese culture and three in Western culture. Results show significant differences in literal translation rates across models ChatGLM and Qwen, with a stronger emphasis on Chinese-language data, exhibit higher literal translation rates, prioritizing structural fidelity over adaptation. For Western cultural concepts, Llama and Mistral show some degree of accurate literal translation, though performance varies. Notably, no model directly translates 'waxberry', likely due to its regional specificity and the absence of a widely recognized equivalent in Western languages. NLLB struggles with all six concepts, highlighting potential NMT limitations. Furthermore, while Llama and Mistral do not achieve the highest literal translation rates, they demonstrate a tendency to adapt culturally specific terms rather than translate them directly. For example, they often map 'raspberry' to other red berries such as 'blueberry' and 'blackberry' and replace 'thyme' with similar spices like 'bay leaves' and 'cinnamon'. These adaptations align with the categorization in the Wine Aroma Wheel, as both the substituted berries and spices belong to the same Aroma Subfamilies. These findings are consistent with Table 3, where ChatGLM and Qwen score higher in Faithfulness of Information, while Llama and Mistral perform better in Cultural Proximity.

We further assess how well models handle wine terminology, which often has industry-specific meanings distinct from general usage. As shown in Figure 3b, ChatGLM and Qwen perform well in translating these terms, while Phi also ranks highly, even leading for 'full'. However, challenges arise with terms such as 'nose', which refers to a wine's aroma in professional contexts. The phrase "on the nose" describes a wine's bouquet, yet most models translate 'nose' directly, failing to capture its specialized meaning. This highlights the challenge of translating wine-specific terms beyond their literal meanings and underscores the importance of integrating domain knowledge into machine translation models.



(a) Analysis of the translation of specific concepts by the different models on the test data.



(b) Analysis of the translation of specific terminology by the different models on the test data.

Figure 3: Comparison of translation analysis: specific concepts vs. specific terminology. In brackets, we show the number of occurrences of each concept.

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

8 Conclusion

In this work, we studied cross-cultural adaptation of wine reviews, introducing CulturalWR, a dataset of paired Chinese and English reviews, and evaluating LLM-based adaptation methods. Our results show that LLMs can consider cultural nuances but face challenges in maintaining detail and consistency in flavor descriptions. We also assessed adapted flavor similarities to gauge LLMs' understanding of wine descriptors. Beyond wine reviews, our findings have broader implications for crosscultural communication in the wine industry, aiding wineries and retailers in tailoring descriptions to international audiences. This could enhance global wine marketing while preserving cultural authenticity. Moreover, our work highlights AI's potential in gastronomy, paving the way for research into AI-assisted flavor profiling and food pairing. Future work includes refining adaptation models for coherence, integrating multimodal data, and using user feedback to enhance AI-generated adaptations, bridging cultural gaps in wine appreciation.

9 Limitation

Cultural adaptation in wine reviews has great poten-611 tial to aid consumer decisions and promote wines 612 globally. However, several challenges remain. A 613 key limitation is the dataset size and diversity. Our study includes fewer than 5,000 Chinese reviews, primarily from three professional critics, raising 616 concerns about representativeness. Additionally, 617 while our dataset contains reviews in multiple lan-618 guages (German, Portuguese, French, Dutch, Italian, and Spanish), we focused solely on Chinese and English due to resource constraints. Expanding multilingual analysis could offer further insights. Another constraint lies in evaluation prompts. Our analysis relies on four prompts, which may not 624 fully capture how models handle cultural adapta-625 tion. Broader prompt variations and real-world user inputs could improve evaluation robustness. Moreover, our quality assessment is based on the Wine Aroma Wheel and prior sensory adaptation research (Jin et al., 2022), but it lacks independent verification of adaptation accuracy. Future work could incorporate human or expert evaluations for a more comprehensive assessment. Finally, wine reviews are inherently subjective, influenced by per-634 sonal preferences and sensory perceptions. While we strive to minimize bias, eliminating subjectivity entirely remains challenging. Addressing this may 637 require larger, more diverse datasets and structured sensory evaluation frameworks in future research.

Ethical Considerations

This study relies on wine ratings and tasting notes 641 sourced from professional websites, which are used strictly for non-commercial academic research. 643 Due to the proprietary nature of the data, we do not publicly release or redistribute any portion of it. Our use falls within the permitted scope of personal, non-commercial informational use, as outlined in the platform's Terms of Use, and all ref-648 erences are properly attributed. No automated data extraction or systematic collection was conducted, ensuring compliance with intellectual property and 652 fair use policies. Additionally, all human evaluators involved in this research are co-authors of the paper and participated voluntarily without financial compensation.

References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Oin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 technical report: A highly capable language model locally on your phone. Preprint, arXiv:2404.14219.

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

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

AI@Meta. 2024. Llama 3 model card.

- Araceli Arellano-Covarrubias, Carlos Gómez-Corona, Paula Varela, and Héctor B. Escalona-Buendía. 2019. Connecting flavors in social media: A cross cultural study with beer pairing. *Food Research International*, 115:303–310.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Yong Cao, Yova Kementchedjhieva, Ruixiang Cui, Antonia Karamolegkou, Li Zhou, Megan Dare, Lucia Donatelli, and Daniel Hershcovich. 2024. Cultural Adaptation of Recipes. *Transactions of the Association for Computational Linguistics*, 12:80–99.

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

Bernard Chen, Christopher Rhodes, Alexander Yu, and Valentin Velchev. 2016. The computational wine wheel 2.0 and the trimax triclustering in wineinformatics. In Advances in Data Mining. Applications and Theoretical Aspects, pages 223–238, Cham. Springer International Publishing.

716

718

721

722

727

730

731

733

734

735

736

737

740

741

742

743

744

745

746

747

748

749

750

751

753

754

756

757

758

764

767

770

772

773

774

- Jyh-Shen Chiou, Cheng-Chieh Hsiao, and Fang-Yi Su. 2014. Whose online reviews have the most influences on consumers in cultural offerings? professional vs consumer commentators. *Internet Research*, 24(3):353–368.
- M.A. Drake, M.D. Yates, P.D. Gerard, C.M. Delahunty, E.M. Sheehan, R.P. Turnbull, and T.M. Dodds. 2005. Comparison of differences between lexicons for descriptive analysis of cheddar cheese flavour in ireland, new zealand, and the united states of america. *International Dairy Journal*, 15(5):473–483.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference* on Machine Translation (WMT), pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. Preprint, arXiv:2406.12793.
- Daniel Hershcovich, Stella Frank, Heather Lent, Miryam de Lhoneux, Mostafa Abdou, Stephanie Brandl, Emanuele Bugliarello, Laura Cabello Piqueras, Ilias Chalkidis, Ruixiang Cui, Constanza Fierro, Katerina Margatina, Phillip Rust, and Anders Søgaard. 2022. Challenges and strategies in crosscultural NLP. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6997–7013, Dublin, Ireland. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,

Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

- Gang Jin, Xi Lv, Linsheng Wei, Laichao Xu, Junxiang Zhang, Yanping Chen, and MA Wen. 2022. Chinese localisation of wine aroma descriptors: an update of le nez du vin terminology by survey, descriptive analysis and similarity test. *OENO One*, 56(1):241–251.
- Els Lefever, Iris Hendrickx, Ilja Croijmans, Antal van den Bosch, and Asifa Majid. 2018. Discovering the language of wine reviews: A text mining account. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation* (*LREC 2018*), Miyazaki, Japan. European Language Resources Association (ELRA).
- Daniel Liebling, Katherine Heller, Samantha Robertson, and Wesley Deng. 2022. Opportunities for humancentered evaluation of machine translation systems. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 229–240, Seattle, United States. Association for Computational Linguistics.
- Yoshihiko Nitta. 1986. Problems of machine translation systems: Effect of cultural differences on sentence structure. *Future Generation Computer Systems*, 2(2):101–115.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy,

David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian O'Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan,

834

835

855

866

870

871

873

874

876

878

879

891

896

897

Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunninghman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. 2024. Gpt-4o system card. Preprint, arXiv:2410.21276.

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

- Nicholas Ostler. 1999. "the limits of my language mean the limits of my world": is machine translation a cultural threat to anyone? In *Proceedings of the* 8th Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages, University College, Chester.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Heber Rodrigues and Wendy V Parr. 2019. Contribution of cross-cultural studies to understanding wine appreciation: A review. *Food research international*, 115:251–258.
- Miloš Stanojević and Khalil Sima'an. 2014. BEER: BEtter evaluation as ranking. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 414–419, Baltimore, Maryland, USA. Association for Computational Linguistics.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram

Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.

959

962

966

967

969

970

971

972

973

974

975 976

977

978

979

981

984

985

990

991 992

993

994

997

999

1000

1001

1002

1005

- Qiuyun Tian, Brittany Whiting, and Bernard Chen. 2022. Wineinformatics: Comparing and combining svm models built by wine reviews from robert parker and wine spectator for 95 + point wine prediction. *Fermentation*, 8(4).
- Carlos Velasco, Xiaoang Wan, Alejandro Salgado-Montejo, Andy Woods, Gonzalo Andrés Oñate, Bingbing Mu, and Charles Spence. 2014. The context of colour–flavour associations in crisps packaging: A cross-cultural study comparing chinese, colombian, and british consumers. *Food Quality and Preference*, 38:49–57.
- Shican Wu, Xiao Ma, Dehui Luo, Lulu Li, Xiangcheng Shi, Xin Chang, Xiaoyun Lin, Ran Luo, Chunlei Pei, Changying Du, Zhi-Jian Zhao, and Jinlong Gong. 2025. Automated review generation method based on large language models. *Preprint*, arXiv:2407.20906.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. Qwen2.5 technical report. arXiv preprint arXiv:2412.15115.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

A French and Spanish Character Conversion Table

This section shows the Character Conversion Table we have used to convert some characters, shown in Table 7.

B Whole dataset of Other Languages

1006It's the all the data we have collected, besides Chi-1007nese and English, we also collected some reviews1008written in German, Portuguese, French, Dutch, Ital-1009ian and Spanish.

Original	Converted	Original	Converted
à	a	â	a
ä	a	é	e
è	e	ê	e
ë	e	î	i
ï	i	ô	0
ö	0	ù	u
û	u	ü	u
ç	с	ÿ	у
æ	ae	œ	oe
À	A	Â	А
Ä	A	É	Е
È	E	Ê	Е
Ë	E	Î	Ι
Ï	Ι	Ô	0
Ö	0	Ù	U
Û	U	Ü	U
Ç	С	Ÿ	Y
á	a	í	i
ó	0	ú	u
ñ	n	Á	A
Í	Ι	Ó	0
Ú	U	Ñ	N

Table 7: French and Spanish Character Mapping to English

	Numbers	Mean Tokens
CA Chinese Reviews	4776	67.57
CA English Reviews	3227	74.25
WA English Reviews	16746	58.16
WA German Reviews	3341	50.73
WA Portuguese Reviews	161	82.15
WA French Reviews	480	135.63
WA Italian Reviews	40	44.35
WA Spanish Reviews	10	114.8

Table 8: Statistics of reviews. We count tokens with jieba text segmentation for Chinese and whitespace tok-enization for other languages.

C Prompt Used for Evaluation

Table 9 shows different prompt strategies we usedfor Prompting Strategy Evaluation.

Here shows the prompt we used for GPT evaluation:

As a Western consumer, evaluate the
quality of wine review translation1015from these dimensions, use seven-tier
scoring ranging from 1-7 : 7 means1017Excellent, 6 means Very Good, 5
means Good, 4 means Fair, 3 means1020Poor, 2 means Very Poor, 1 means1021

- 1011 1012
- 1013 1014

Strategy	Prompt
Direct Translation	Translate the following English wine reviews to Chinese: [Wine review]
Cultural Prompt	Translate the wine reviews in Chinese, adapted to an Chinese-speaking consumer: [Wine review]
Datailad Cultural Prompt	Translate the provided English wine review into Chinese, so that it fits within Chinese wine culture and to avoid using
Detailed Cultural Frompt	any terms that might have negative connotations for Chinese consumers: [Wine review]
	User: Find flavor and aroma descriptions that are unfamiliar and uncomfortable for Chinese consumers: [Wine review]
Self Exploration	LLM: [Sentences]
Sen-Explanation	User: Translate the wine review in Chinese, and for the unfamiliar and uncomfortable flavor and aroma,
	replace it with a more familiar and comfortable description for Chinese consumers.

Table 9: Prompting strategy examples used for English \rightarrow Chinese translation

1022		Nonsense: Grammar,Faithfulness of
1023		Information, Faithfulness of Style,
1024		Overall quality,Cultural proximity
1025		- The generated reviews use
1026		familiar terms and expressions that
1027		resonate with the target culture,
1028		Cultural neutrality – Maintains
1029		neutrality to avoid provoking
1030		negative perceptions or reactions
1031		from the target culture consumers,
1032		Cultural Genuineness – Preserves
1033		the quality of the original
1034		descriptor without altering its
1035		meaning, ensuring authenticity.
1036		Original: {original}
1037		Translation: {translation}
1038	D	Grading scale rule for human

D Grading scale rule for human evaluation

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053 1054

1055

1056

Here shows the rule for the Seven-tier rating system:

- 1. Excellent (rating: 7 points) The translation is also highly matched in context and culture.
- Very good (rating: 6 points) There may be slight (1-2) inaccuracy of vocabulary or inadequate reflection of some details, but it does not affect the overall understanding.
- Good (rating: 5 points) There are individual mistranslations, but they will not seriously change the overall meaning.
- Fair (score: 4 points) Mistranslations are obvious, which may cause some semantics to deviate from the original text.
- 10575. Poor (score: 3 points)1058The overall translation quality is low and may1059mislead readers.

- 6. Very poor (score: 2 points)1060It is basically impossible to rely on the translation to understand the original text.1061
- 7. Nonsense (score: 1 point)1063The translation is completely incoherent semantically and unable to convey the information which the original review tries to convey.1065

1067

1071

1077

1078

E Human evaluation platform

Figure 4 shows a screenshot from our human eval-
uation platform, demonstrating the English to Chi-
nese direction. Human need to evaluation the qual-
ity of Chinese translation1068
1069



Figure 4: Screenshot from our human evaluation platform

F Attributes of Whole dataset	1072
Figure 5 shows attributes counted in CulturalWR	1073
dataset.	1074
G Experiment Settings	1075
The experiment settings of different models in-	1076

	enperment	Section	01 011010110	
clud	ed in our pap	per are as	follows:	

```
1. NLLB We use NLLB-200-3.3B^{12} for our ex-
```

¹²https://huggingface.co/facebook/nllb-200-3. 3B





periments. The beam is set as 4, and the length penalty is set as 1.0.

Pretrained LLMs We used LLMs including Llama-3.1-8B-Instruct ¹³, Mistral-7B-Instruct-v0.3 ¹⁴, Phi-3.5-mini-instruct¹⁵, Qwen2.5-7B-Instruct¹⁶, GLM4-9b¹⁷. The sampling is set as True, leading to a multinomial sampling searching method. All settings are the same across different models.

1079

1080

1081

1082

1084

1085

1086

1087

3. **ChatGPT** We used the latest version, GPT-4o-2024-11-20, through the ChatCompletion API provided by OpenAPI ¹⁸. For the generation, we set the parameters as default, for which the temperature is 1, top_p is 1, and frequency_penalty as 0.

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

H Examples

In Table 10, we present a comparative analysis of multiple machine translation models applied to an English-to-Chinese translation example which is shown in Figure 1. Each row in the table represents the output of a different model, with the original English input provided for reference. To facilitate an in-depth evaluation, we systematically annotate

¹³https://huggingface.co/meta-llama/Llama-3. 1-8B-Instruct

¹⁴https://huggingface.co/mistralai/

Mistral-7B-Instruct-v0.3

¹⁵https://huggingface.co/microsoft/Phi-3. 5-mini-instruct

¹⁶https://huggingface.co/Qwen/Qwen2.

⁵⁻⁷B-Instruct

¹⁷https://huggingface.co/THUDM/glm-4-9b

¹⁸https://platform.openai.com/docs/guides/ text-generation

Category	Content
Input (English)	Aromas of iris, raspberry, camphor and Mediterranean scrub mingle with oak-driven spice on this 100% Merlot.
ChatGLM4	这款酒散发着紫罗兰、树莓、樟脑和地中海灌木的香气,100%的梅洛。
Phi	这100%的 <mark>梅鹿酒</mark> 充满了茉莉花、覆盆子、茴香和 <mark>地中海疏林</mark> 的香气,与橡 木风味相结合的精髓。
Qwen2.5	这款100%梅洛红酒散发出鸢尾花、覆盆子、薄荷和地中海灌木丛的香气, 与橡木带来的香料味交织在一起。
Mistral	这是一款 <mark>百分之一的</mark> 梅洛酒,浓郁而舒服,携带了芙蓉花、莓果、樟木和 植物的香气,以及橡木带来的香氛。
Llama	这款100%梅洛特红酒的香气中融合了百合花、草莓、樟脑和地中海灌木丛的气息,伴随着橡木驱动的香料。
NLLB	红虹,树,和地中海洗刷的香气与木驱动的香料混合在这个100%的梅罗特酒上。
ChatGPT	这款100%梅洛红酒散发出鸢尾花、覆盆子、樟脑和地中海灌木的香气,并 与橡木带来的香料味交织在一起。

Table 10: Comparison of Translations from Different Models. Red text shows the wrong translation, teal text shows the correct literal translation, blue text shows the adapted translation and brown text shows unsure adapted translation.

various translation characteristics, including direct translations, adapted phrasings, and potential errors.

To facilitate evaluation, we use color coding to distinguish different translation characteristics: red indicates incorrect translations, teal represents accurate literal translations, blue highlights adapted translations that maintain meaning while improving fluency, and brown marks uncertain adaptations.

This comparison reveals variations in how different models interpret and translate key terms, particularly in handling domain-specific vocabulary such as "Merlot" and "Mediterranean scrub." Some models exhibit direct translation errors, while others apply adaptive strategies to enhance readability. By analyzing these differences, we can better understand the strengths and limitations of current machine translation systems.

1102

1103