

“Be My Cheese?": Cultural Nuance Benchmarking for Machine Translation in Multilingual LLMs

Anonymous ACL submission

001	Abstract	1 Introduction	039
002	We present a large-scale human evaluation	Large language models (LLMs) have rapidly ex-	040
003	benchmark for assessing cultural localisation	expanded access to machine translation, enabling	041
004	in machine translation produced by state-of-	rapid translation across hundreds of languages	042
005	the-art multilingual large language models	without requiring linguistic expertise. Cultural	043
006	(LLMs). Existing MT benchmarks emphasise	nuances, such as figurative expressions and idi-	044
007	token-level and grammatical accuracy, but of-	oms, are foundational to effective human commu-	045
008	ten overlook pragmatic and culturally ground-	nication and shape how meaning is received and	046
009	ed competencies required for real-world lo-	interpreted by local audiences. A translation that	047
010	calisation. Building on a pilot study of 87	is grammatically correct may nevertheless sound	048
011	translations across 20 languages, we evaluate 7	unnatural, inappropriate, or misleading if it fails	049
012	multilingual LLMs across 15 target languages	to account for cultural context. However, machine	050
013	with 5 native-speaker raters per language. Ra-	translation (MT) research and benchmarks con-	051
014	ters scored both full-text translations and seg-	tinue to prioritise lexical and grammatical accu-	052
015	ment-level instances of culturally nuanced lan-	cy at the token- and sentence-level. These metrics	053
016	guage (idioms, puns, holidays, and culturally	capture formal correctness, but fail to evaluate the	054
017	embedded concepts) on an ordinal 0–3 quality	pragmatic, cultural, and stylistic competencies	055
018	scale; segment ratings additionally included an	required for real-world localisation tasks such	056
019	NA option for untranslated segments.	as marketing communication, customer engage-	057
020	Across full-text evaluations, mean overall qual-	ment, and culturally specific brand messaging.	058
021	ity is modest (1.68/3): GPT-5 (2.10/3), Claude	This study introduces a benchmark designed	059
022	Sonnet 3.7 (1.97/3), and Mistral Medium 3.1	explicitly for evaluating how well multilingual	060
023	(1.84/3) form the strongest tier with fewer cat-	LLMs preserve cultural resonance in machine	061
024	astrophic failures. Segment-level results show	translation tasks. Building on a pilot evaluation of	062
025	sharp category effects: holidays (2.20/3) and	87 translations across 20 languages (Anonymous,	063
026	cultural concepts (2.19/3) translate substan-	2025), we scale to a substantially larger dataset	064
027	tially better than idioms (1.65/3) and puns	comprising 7 state-of-the-art multilingual LLMs,	065
028	(1.45/3), and idioms are most likely to be left	15 target languages, and five native-speaker ra-	066
029	untranslated. These findings demonstrate a	tters per language. Each rater evaluated both (1) a	067
030	persistent gap between grammatical adequacy	complete translated marketing email and (2) pre-	068
031	and cultural resonance. To our knowledge,	defined segment-level instances of culturally nu-	069
032	this is the first multilingual, human-annotated	anced language, including idioms, puns, holiday	070
033	benchmark focused explicitly on cultural nu-	references, and culturally embedded concepts.	071
034	ance in translation and localisation, highlight-	This design allows us to contrast holistic trans-	072
035	ing the need for culturally informed training	lation quality with categorical failure modes on a	073
036	data, improved cross-lingual pragmatics, and	phrasal level.	074
037	evaluation paradigms that better reflect real-	Our study addresses three core research questions:	075
038	world communicative competence.		

076	• How well do contemporary multilingual LLMs	We identify consistent performance differences	119
077	translate culturally nuanced language across ty-	across models, languages, and orthographic sys-	120
078	pologically diverse languages?	tems, including higher stability among GPT-5,	121
079	• To what extent do model family, linguistic	Claude Sonnet 3.7, and Mistral Medium 3.1, and	122
080	characteristics, and orthographic systems im-	elevated failure rates for culturally marked seg-	123
081	act cultural resonance in MT?	ments in other systems, motivating targeted data	124
082	• Which categories of culturally marked content,	and evaluation strategies for improving cultural	125
083	such as idioms and puns, pose the greatest chal-	competence in multilingual LLMs.	126
084	lenges to current LLMs?		
085	Our findings reveal a substantial gap between	2 Related Work	127
086	grammatical accuracy and cultural localisation.	Recent advances in large language models	128
087	While many translations achieve surface-level	(LLMs) have driven substantial improvements in	129
088	adequacy, even the strongest multilingual LLMs	multilingual machine translation. Mujadia et al.	130
089	fail to consistently preserve culturally grounded	(2023) provide a comprehensive assessment of	131
090	meaning, particularly for figurative and non-lit-	LLM translation performance between English	132
091	eral language. These results underscore the lim-	and 22 Indian languages, revealing persistent	133
092	itations of existing machine translation in SOTA	disparities across high- and low-resource set-	134
093	models and motivate a reevaluation of MT bench-	tings and demonstrating the benefits of in-context	135
094	marks and training practices that prioritises cul-	learning for underrepresented dialects. Similarly,	136
095	tural-pragmatic competence as a core dimension	Hu et al. (2024) introduce GenTranslate, showing	137
096	of multilingual LLM performance.	that generative LLM-based approaches improve	138
097	1.1 Contributions	multilingual speech and text translation on stan-	139
098	This work presents three primary contributions:	dard benchmarks, particularly for low-resource	140
099	1. A new benchmark for culturally sensitive	languages. Together, these studies illustrate rapid	141
100	machine translation.	progress in multilingual MT while highlighting	142
101	We introduce the first multilingual, human-an-	uneven gains across languages.	143
102	notated benchmark designed explicitly to evalu-	Despite these advances, most prior evaluations	144
103	ate cultural nuance and resonance in machine	focus on lexical and grammatical accuracy, re-	145
104	translation, spanning 7 state-of-the-art multilin-	lying on automatic metrics or sentence-level	146
105	gual LLMs, 15 languages, and five native-speaker	adequacy judgments. Such evaluations are poorly	147
106	raters per language, with both holistic and seg-	suited to capturing pragmatic and cultural dimen-	148
107	ment-level evaluation.	sions of translation quality, including idiomatic	149
108	2. A large-scale empirical analysis of cultural	meaning, figurative language, and audience-ap-	150
109	failure modes in MT.	propriate tone. As a result, translations that are	151
110	Through segment-level evaluation of idioms,	formally correct may nevertheless be culturally	152
111	puns, holidays, and culturally embedded con-	inappropriate or misleading in real-world locali-	153
112	cepts, we show that cultural localisation quality	sation contexts. This limitation is well document-	154
113	diverges sharply from grammatical accuracy, with	ed in the MT evaluation literature. BLEU has long	155
114	figurative language remaining a persistent failure	been shown to correlate weakly with meaning	156
115	mode across models and languages.	adequacy and human judgments beyond surface	157
116	3. Evidence of systematic model- and	correspondence (Callison-Burch et al., 2006; Mathur	158
117	language-level variation in cultural MT	et al., 2020), and more recent neural metrics such	159
118	performance.	as COMET and BLEURT similarly struggle with	160
		discourse-level, pragmatic, and culturally ground-	161
		ed errors (Freitag et al., 2021; Kocmi et al., 2022).	162
		Recent work has therefore begun to frame cultur-	163
		al transfer and adaptation as a core challenge for	164
		language technologies, arguing that culture-aware	165
		evaluation is necessary to capture meaning be-	166
		yond surface correspondence (Singh et al., 2024).	167

249	• puns and humorous wordplay	analysis of where models succeed or fail in cultural MT beyond full-text impressions. This methodology produced 13,125 segment-specific annotations.	295
250	• holiday-specific phrases		296
251	• idiomatic expressions		297
252	• culturally specific references		298
253	• strong brand voice and audience targeting		
254	From each email, we selected five segments of culturally nuanced language. Across the dataset, this resulted in four puns, four idioms, four holiday references, and thirteen cultural concepts per language. Cultural concepts were defined as single words or short phrases that are either specific to North American English or unlikely to have direct equivalents across cultures (e.g., <i>koozies</i> , <i>sweetheart</i> , <i>zero-waste</i>). Full source texts and segment selections are provided in Appendix A.		
255			
256			
257			
258			
259			
260			
261			
262			
263			
264	3.4 Evaluation Procedure	3.5 Annotation Protocol	299
265	Each rater assessed one translation per model, evaluating both the full translated text and segments. Full participant guidelines are presented in Appendix B2.	Participants received detailed written instructions based on an evaluation framework (available in Appendix B2), including:	300
266		• definitions of cultural nuance	303
267		• examples of literal vs. localised translation strategies	304
268		• guidance on how to rate ambiguous cases	305
269	(a) Full text evaluation	• clarifications for rating idioms and humour	307
270	Participants scored the translation on a 4-point scale for the following criteria:	Ratings were collected using our proprietary data annotation software (redacted for anonymity). Each submission was checked for completeness and annotation consistency.	308
271	1. Content fidelity		309
272	2. Style fidelity		310
273	3. Audience appropriateness		311
274	4. Overall translation quality	3.6 Statistical Analysis	312
275		We analysed segment-level translation ratings using a cumulative link mixed model (CLMM) with a logit link, appropriate for ordinal outcomes. Models were fitted in R using the ordinal package (Christensen, 2022). Fixed effects included model, language, and segment category, as well as their interaction. Random intercepts were included for annotator and segment to account for repeated ratings and item-level variability.	313
276	These items measure whether the translation is correct, natural, locally resonant, and aligned with the original intent. Participants were also given free response text boxes to provide additional qualitative feedback. A summary of the qualitative feedback by language is available in Appendix D.		314
277			315
278			316
279			317
280			318
281			319
282			320
283	(b) Segment-level evaluation	Orthography was initially included as a fixed effect but was removed from the final specification due to rank deficiency and near-complete collinearity with language–category combinations. Its inclusion resulted in unstable parameter estimates without improving model fit. The final model converged successfully (logLik = -14,411.63; AIC = 28,965.26; n = 13,125). Random-effects estimates indicate greater variance at the segment level (SD = 1.76) than at the annotator level (SD = 0.70), suggesting that segment-specific difficulty contributes more to rating variability than individual rater severity.	321
284	Raters also evaluated predefined culturally nuanced segments from the emails, each labeled as one of:		322
285	• idioms		323
286	• puns		324
287	• holidays		325
288	• cultural concepts		326
289			327
290			328
291	Segments were rated on the same 0–3 scale, with an additional NA option indicating the segment was not translated, instead opting to retain the original English. This enables fine-grained		329
292			330
293			331
294			332
		Inter-rater reliability (IRR) was assessed separately for full-text (overall) ratings and segment-level ratings using Krippendorff’s α (ordinal) and Gwet’s AC2 with quadratic weights. IRR was computed overall and stratified by model, language, and segment category. Full-text IRR	333
			334
			335
			336
			337
			338
			339
			340

415	4.3 Model Effects on Segment Translation	5 Discussion	459
416	Controlling for language and segment category,	The CLMM analysis confirms that cultural localisation failures in multilingual LLMs are systematic rather than anecdotal. Segment category emerges as the strongest predictor of translation quality, exceeding the influence of both model family and language. Figurative language, especially idioms and puns, remains a robust failure mode even after controlling for rater effects and segment-level difficulty.	460
417	model choice significantly affects segment-level		461
418	translation quality (Table C1). GPT-5 and Claude		462
419	Sonnet 3.7 do not differ significantly and outperform		463
420	gpt-oss 120B, Llama 4, and Aya Expanse		464
421	8B. Mistral Medium 3.1 performs significantly		465
422	better than Aya Expanse 8B and Llama 4, but		466
423	does not differ significantly from DeepSeek V3.1		467
424	or GPT-5.		468
425	Aya Expanse 8B is a clear outlier, exhibiting	Crucially, the statistical results support a distinction between translation coverage and translation quality. Idioms are significantly more likely to be omitted entirely, and when translated, they receive substantially lower ratings than holidays or culturally embedded concepts. Aya Expanse 8B exhibits both the highest omission rates and the lowest quality scores for idiomatic translation, indicating that failure is not merely a consequence of conservative behavior. Even when models attempt figurative translation, pragmatic and culturally appropriate rendering frequently fails.	469
426	both significantly lower quality scores and substantially		470
427	higher omission rates for idioms and puns. Other		471
428	models omit fewer segments overall but frequently		472
429	produce low-quality translations (ratings 0–1) for		473
430	figurative language.		474
431	IRR stratified by model (Table C8) indicates		475
432	moderate agreement for GPT-5 and Claude Sonnet		476
433	3.7, with greater variability for lower-performing		477
434	models, suggesting that inconsistent output quality		478
435	contributes to annotator disagreement.		479
436	4.4 Language-Level Effects		480
437	Language effects are present but more constrained	Model-level effects reveal a stable top tier (GPT-5, Claude Sonnet 3.7, and Mistral Medium 3.1) but no system consistently achieves high performance across all categories. The absence of statistically significant differences among these models suggests that scaling and architectural refinement alone do not resolve cultural–pragmatic limitations. In contrast, Aya Expanse 8B’s consistently poor performance across analyses points to systemic fragility rather than isolated weaknesses.	481
438	than category or model effects. CLMM estimates		482
439	indicate that Mandarin (Taiwan) receives significantly		483
440	higher segment-level ratings than several other		484
441	languages, including Spanish, Swahili, and Urdu		485
442	(Tables C4–C5). Brazilian Portuguese trends		486
443	higher but does not consistently differ from		487
444	other languages after correction for multiple		488
445	comparisons.		489
446	Importantly, language effects interact with	Language-level effects are present but secondary, and orthography does not independently predict translation quality once language and segment category are accounted for. This finding contrasts with observations from the pilot study and challenges assumptions that script or typological complexity are primary drivers of cultural MT difficulty. Instead, the results point toward the availability and quality of culturally situated training data as a more plausible explanation for observed disparities.	491
447	segment category. Languages that perform well		492
448	overall tend to maintain higher scores across		493
449	categories, while lower-performing languages		494
450	exhibit disproportionate degradation on idioms		495
451	and puns. This pattern persists even when		496
452	restricting analysis to translated segments,		497
453	indicating that low scores are not driven		498
454	solely by omission.		499
455	IRR varies substantially by language (Table		500
456	C8), with lower agreement for Mandarin, Hindi,		501
457	and Urdu, suggesting that cultural interpretation	6 Future work	502
458	differences may amplify annotator variability	Future work will extend this benchmark in several directions. First, we plan to release the dataset and evaluation framework as a public benchmark,	503
	in these contexts.		504
			505

506	enabling reproducible research on cultural localisation in machine translation and multilingual LLM evaluation. The release will include full-text translations, segment-level annotations, and detailed evaluation guidelines to support consistent comparison across future models. Rather than replacing automatic metrics, this benchmark will complement them by targeting pragmatic and cultural dimensions that current form-based evaluations systematically overlook.	553
507		554
508		555
509		556
510		557
511		558
512		559
513		560
514		561
515		562
516	Second, we plan to expand the benchmark beyond text-only translation by developing an audio-based version of the task. Many culturally marked expressions – such as humour, idioms, and tone – are realised differently in spoken language, and evaluating speech-based localisation will allow analysis of dialect, prosody, emphasis, and pragmatic delivery not captured in text. We also intend to extend coverage to additional domains and languages to assess the generality of the cultural failure modes identified here.	563
517		564
518		565
519		566
520		567
521		568
522		569
523		570
524		571
525		572
526		573
527	7 Conclusion	574
528	We presented a large-scale, human-annotated benchmark designed to evaluate cultural localisation in machine translation by multilingual LLMs. Across seven state-of-the-art models and fifteen languages, results reveal a persistent gap between grammatical adequacy and cultural resonance. While many translations appear superficially plausible, segment-level evaluation exposes systematic failures, particularly idioms and puns, that remain largely invisible to standard MT metrics.	575
529		576
530		577
531		578
532		579
533		580
534		581
535		582
536		583
537		584
538	By explicitly distinguishing between translation coverage and translation quality, this work provides a more nuanced account of cultural MT performance and highlights limitations shared even by the strongest current models. These findings underscore the need for culturally informed training data and evaluation paradigms that move beyond form-based correctness toward real-world communicative competence.	585
539		586
540		
541		
542		
543		
544		
545		
546		
547	8 Limitations	
548	This study has several limitations. The benchmark focuses on English-to-many translation within a marketing email domain, which may limit generalisability to other genres such as news, legal text, or conversational dialogue. Segment selection intentionally emphasises culturally marked language and is therefore not representative of typical sentence distributions in MT corpora. Furthermore, analysis of segment-level MT was performed in the context of larger MT corpora and may not generalise to MT performance when translating the same segments as isolated text. Future work could rectify this limitation by evaluating and contrasting segment-level MT in context of larger text with segment MT as isolated input.	587
549		588
550		589
551		590
552		591
		592
		593
		594
		595
		596
		597
	Although five native raters per language reduce individual bias, judgments of cultural appropriateness remain inherently subjective and may vary across demographics, regions, and personal experience within a language community. Additionally, this study did not control for differences in participant age, education, gender, and socioeconomic background – all factors known to influence human bias (Jenks, 2025; Zahraei and Emami, 2025). In addition, models were evaluated through publicly accessible interfaces, which may introduce uncontrolled variation due to system prompts, safety filters, or model updates. Furthermore, model outputs were collected over the span of two days, introducing additional potential for uncontrolled variation when compared to simultaneous output generation and collection. Finally, this work focuses exclusively on text-based translation and does not address multimodal or spoken localisation, which we leave to future research.	
	References	
	Anonymous. 2025. Redacted for ACL blind review. Pilot study on cultural localisation in machine translation. 2025	
	Amirhossein Abaskohi, Sara Baruni, Mostafa Masoudi, Nesa Abbasi, Mohammad Hadi Babalou, Ali Edalat, Sepehr Kamahi, Samin Mahdizadeh Sani, Nikoo Naghavian, Danial Namazifard, Pouya Sadeghi, and Yadollah Yaghoobzadeh. 2024. Benchmarking Large Language Models for Persian: A Preliminary Study Focusing on ChatGPT . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 2189–2203, Torino, Italia. ELRA and ICCL.	

598	Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. Investigating Cultural Alignment of Large Language Models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.	642
599		643
600		644
601		645
602		646
603		647
604		648
605	Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. Re-evaluating the Role of Bleu in Machine Translation Research . In <i>11th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 249–256, Trento, Italy. Association for Computational Linguistics.	649
606		650
607		651
608		652
609		653
610		654
611	Rune Haubo Bojesen Christensen. 2022. ordinal: Regression models for ordinal data . R package.	655
612		656
613	Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, Errors, and Context: A Large-Scale Study of Human Evaluation for Machine Translation . <i>Transactions of the Association for Computational Linguistics</i> , 9:1460–1474.	657
614		658
615		659
616		660
617		661
618		662
619	Yuchen Hu, Chen Chen, Chao-Han Huck Yang, Ruizhe Li, Dong Zhang, Zhehuai Chen, and Eng Siong Chng. 2024. GenTranslate: Large language models are generative multilingual speech and machine translators . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics</i> . Association for Computational Linguistics.	663
620		664
621		665
622		666
623		667
624		668
625		669
626	Christopher Jenks. 2025. Communicating the cultural other: Trust and bias in generative AI and large language models . <i>Applied Linguistics Review</i> , 16(2):787–795.	670
627		671
628		672
629	Klaus Krippendorff. 2019. Content analysis: An introduction to its methodology . Sage.	673
630		674
631	Cheng Li, Mengzhuo Chen, Jindong Wang, and Sunayana Sitaram. 2024. CultureLLM: Incorporating cultural differences into large language models . In <i>Proceedings of the 38th Conference on Neural Information Processing Systems</i> .	675
632		676
633		677
634		
635		
636	Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. Tangled up in BLEU: Reevaluating the Evaluation of Automatic Machine Translation Evaluation Metrics . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4984–4997, Online. Association for Computational Linguistics.	
637		
638		
639		
640		
641		
	Vandan Mujadia, Ashok Urlana, Yash Bhaskar, Penumalla Aditya Pavani, Kukkapalli Shrivya, Parameswari Krishnamurthy, and Dipti Sharma. 2024. Assessing Translation Capabilities of Large Language Models involving English and Indian Languages . In <i>Proceedings of the 25th Annual Conference of the European Association for Machine Translation (Volume 1)</i> , pages 207–228, Sheffield, UK. European Association for Machine Translation (EAMT).	
	Pushpdeep Singh, Mayur Patidar, and Lovekesh Vig. 2024. Translating across cultures: LLMs for intralingual cultural adaptation . In <i>Proceedings of the 28th Conference on Computational Natural Language Learning</i> , pages 400–418, Miami, FL, USA. Association for Computational Linguistics.	
	Sara Sterlie, Nina Weng, and Aasa Feragen. 2024. Generalizing fairness to generative language models via reformulation of non-discrimination criteria . arXiv preprint arXiv:2403.08564.	
	Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. ByT5: Towards a token-free future with byte-level models . <i>Transactions of the Association for Computational Linguistics</i> .	
	Pardis Sadat Zahraei and Ali Emami. 2025. Translate With Care: Addressing Gender Bias, Neutrality, and Reasoning in Large Language Model Translations . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 476–501, Vienna, Austria. Association for Computational Linguistics.	
	Jinman Zhao, Yitian Ding, Chen Jia, Yining Wang, and Zifan Qian. 2024. Gender bias in large language models across multiple languages . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics</i> . Association for Computational Linguistics.	

678	Appendix		
679	Appendix A. Complete MT Input Texts		
680	A1		
681	Company: Sheffield's – a gourmet market in NYC	a scent-sational upgrade – pair our newest reusable case design with a fragrance that's sure to make memories. Durable, stylish, compact, and zero waste.	712 713 714 715
682	Subject: Will you brie mine? 🧀❤️🧀	Swipe right this New Year's Eve	716
683	Valentine's Day is almost here, and we've got the sweetest gift ideas for pickup or delivery throughout NYC.	Use code: NYE2026	717
684		[shop deodorant]	718
685		Whether you're keeping yourself fresh for your partner, or looking to impress someone else, our new scents will leave a lasting impression.	719 720 721
686	Cheese Tasting Gift Boxes	MIX & MATCH OUR BEST-SELLING COMBOS	722 723
687	This cheese lover's dream is thoughtfully assembled by our expert cheesemongers. It all comes beautifully packaged in a keepsake tin, tied with a satin ribbon. Personalize it with a custom note on Sheffield's stationery.	Lavender case x Tropical Paradise scent	724
688		Turquoise case x Orange Creamsicle scent	725
689		WHY TERRA?	726
690	Sweets for your Sweetheart	Aluminum & paraben free. Zero-waste refills. 24-hour odor protection. All that in a case you'll be excited to reuse.	727 728 729
691		Terra Cosmetics	730
692		London N1C 4AB, United Kingdom	731
693	Artfully displayed with the perfect accompaniments of fresh & dried fruit, nuts, honey, fig jam, espresso brownies, dark chocolate-covered strawberries, candies, edible flowers and sliced baguette.	A3	732
694		Company: Muggable – an American novelty mug company	733 734
695	[order here]	Subject: This Collection Has Us Feline Good 🐱	735
696	We still have a limited number of handmade, chocolate-covered strawberries and floral arrangements available for pre-order! Give us a call today or stop by the shop before they're gone.	CAT'S MEOW	736
697		Our newest collection is the cat's pajamas, wait no – it's the cat's Mugs, Tumblers, Koozies, and Coasters!	737 738 739
698	Wishing you a sweet Valentine's Day!	[Shop Meow]	740
699	Sheffield's – Park Slope	Rep your favorite feline at the office, on the go, and on your next Zoom call. Wait. Who are we kidding? They're already in all your Zoom calls.	741 742 743
700	Brooklyn, NY	© 2012 Muggable Inc. All Rights Reserved.	744
701	A2	Los Angeles, CA, 90013, USA	745
702	Company: Terra – an eco-friendly deodorant brand		
703	Subject: This scent will transform your life ✨		
704	Hey [NAME]		
705	Your New Year's resolution stinks. Give your life		

746

A4

747

Company: sonia summerhouse– an american luxury swimwear brand

748

749

Subject: late Summer, full throttle

750

Labor Day is here! PACK YOUR BEACH BAG!

751

you sprint barefoot across warm sand.

752

the sun hits high.

753

salt hangs in the air.

754

seagulls cut the wide blue sky.

755

laughter bursts, waves crash in time,

756

summer comes alive.

757

your new swimwear, green like sea glass.

758

fabric flowing, grab your crew,

759

chase the surf, leap, sprint, splash -

760

shore enough, this is your moment!

761

Sonia Summerhouse 20 w. 20th street unit 1004

762

new york, ny 10011

763

A5

764

company: Cinnamon – a neighborhood bakery & cafe

765

766

Subject: Happy Birthday! There's a sweet treat waiting for you!

767

768

Sugar, spice, and everything nice! Happy Birthday from all of us at Cinnamon!

769

770

Let us be the icing on the cake of your special day with a sweet treat. Stop by any Cinnamon location to redeem your credit on your next order OR save it for later by visiting the Rewards section in your app.

771

772

773

774

775

We can't wait to celebrate with you! Redeemable with the Cinnamon app only.

776

777

Excited about your birthday present?

778

Say Thanks on Facebook

A6 Segments

Segment	Segment category
birthday present	cultural concepts
cheesemongers	cultural concepts
full throttle	cultural concepts
grab your crew	cultural concepts
Happy Birthday	cultural concepts
keepsake tin	cultural concepts
Koozies	cultural concepts
summer comes alive	idioms
sweet treat	cultural concepts
Sweetheart	cultural concepts
Swipe right	cultural concepts
Tumblers	cultural concepts
Zero-waste	cultural concepts
Zoom call	cultural concepts
New Year's Eve	holidays
NYE2026	holidays
Labor Day	holidays
Valentine's Day	holidays
cat's pajamas	idioms
icing on the cake	idioms
Sugar, spice, and everything nice	idioms
Feline Good	puns
scent-sational	puns
shore enough	puns
Will you brie mine?	puns

Table A6 Segmentation and categorisation of phrases and words selected for individual evaluation.

Appendix B. Participants

B1 Participant Demographics

Language	Participant age	Participant Gender	Participant education level
Afrikaans	31-45	female	Secondary education completed (high school diploma or equivalent)
Afrikaans	31-45	female	Postgraduate diploma or certificate (non-degree)
Afrikaans	31-45	female	Postgraduate diploma or certificate (non-degree)
Afrikaans	31-45	female	Some college or university (no degree)
Afrikaans	45+	female	Some college or university (no degree)
Arabic	45+	female	Bachelor's degree (e.g., BA, BS)
Arabic	18-30	male	Bachelor's degree (e.g., BA, BS)
Arabic	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Arabic	31-45	female	Bachelor's degree (e.g., BA, BS)
Arabic	18-30	male	Bachelor's degree (e.g., BA, BS)
Brazilian Portuguese	18-30	male	Secondary education completed (high school diploma or equivalent)
Brazilian Portuguese	31-45	female	Some college or university (no degree)
Brazilian Portuguese	45+	male	Master's degree (e.g., MA, MS, MBA, MFA)
Brazilian Portuguese	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Brazilian Portuguese	45+	male	Postgraduate diploma or certificate (non-degree)
Cantonese	31-45	female	Bachelor's degree (e.g., BA, BS)
Cantonese	18-30	female	Bachelor's degree (e.g., BA, BS)
Cantonese		UNKNOWN	Bachelor's degree (e.g., BA, BS)
Cantonese	18-30	female	Bachelor's degree (e.g., BA, BS)
Cantonese	45+	female	Bachelor's degree (e.g., BA, BS)
Czech	31-45	female	Bachelor's degree (e.g., BA, BS)
Czech	18-30	female	Master's degree (e.g., MA, MS, MBA, MFA)
Czech	18-30	male	Some secondary education (high school)
Czech	18-30	male	Some college or university (no degree)
Czech	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Dutch	31-45	male	Bachelor's degree (e.g., BA, BS)
Dutch	31-45	male	Bachelor's degree (e.g., BA, BS)
Dutch	45+	female	Bachelor's degree (e.g., BA, BS)
Dutch	31-45	male	Bachelor's degree (e.g., BA, BS)
Dutch	45+	male	Bachelor's degree (e.g., BA, BS)
Hebrew	45+	male	Bachelor's degree (e.g., BA, BS)
Hebrew	31-45	male	Bachelor's degree (e.g., BA, BS)
Hebrew	31-45	female	Bachelor's degree (e.g., BA, BS)
Hebrew	31-45	male	Vocational/technical training or certification (e.g., trade school)
Hebrew	45+	male	Bachelor's degree (e.g., BA, BS)
Hindi	31-45	male	Bachelor's degree (e.g., BA, BS)

Hindi	31-45	male	Doctoral or professional degree (e.g., PhD, MD, JD, PsyD, EdD)
Hindi	45+	male	Master's degree (e.g., MA, MS, MBA, MFA)
Hindi	18-30	male	Postgraduate diploma or certificate (non-degree)
Hindi	18-30	male	Bachelor's degree (e.g., BA, BS)
Japanese	45+	female	Bachelor's degree (e.g., BA, BS)
Japanese	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Japanese	31-45	male	Bachelor's degree (e.g., BA, BS)
Japanese	45+	male	Bachelor's degree (e.g., BA, BS)
Japanese	18-30	male	Bachelor's degree (e.g., BA, BS)
Korean	45+	female	Bachelor's degree (e.g., BA, BS)
Korean	31-45	female	Master's degree (e.g., MA, MS, MBA, MFA)
Korean			Bachelor's degree (e.g., BA, BS)
Korean	31-45	female	Master's degree (e.g., MA, MS, MBA, MFA)
Korean	45+	female	Bachelor's degree (e.g., BA, BS)
Mandarin	31-45	female	Bachelor's degree (e.g., BA, BS)
Mandarin	31-45	female	Bachelor's degree (e.g., BA, BS)
Mandarin	45+	male	Master's degree (e.g., MA, MS, MBA, MFA)
Mandarin	18-30	male	Bachelor's degree (e.g., BA, BS)
Mandarin	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Russian	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Russian	31-45	female	Bachelor's degree (e.g., BA, BS)
Russian	45+	male	Secondary education completed (high school diploma or equivalent)
Russian	45+	female	Bachelor's degree (e.g., BA, BS)
Russian	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Spanish	31-45	female	Bachelor's degree (e.g., BA, BS)
Spanish	31-45	female	Bachelor's degree (e.g., BA, BS)
Spanish	18-30	female	Master's degree (e.g., MA, MS, MBA, MFA)
Spanish	31-45	FEMALE	Some college or university (no degree)
Spanish	31-45	male	Bachelor's degree (e.g., BA, BS)
Swahili	18-30	female	Bachelor's degree (e.g., BA, BS)
Swahili	31-45	female	Postgraduate diploma or certificate (non-degree)
Swahili	18-30	female	Bachelor's degree (e.g., BA, BS)
Swahili	18-30	male	Bachelor's degree (e.g., BA, BS)
Swahili	18-30	male	Bachelor's degree (e.g., BA, BS)
Urdu	31-45	male	Master's degree (e.g., MA, MS, MBA, MFA)
Urdu	31-45	female	Master's degree (e.g., MA, MS, MBA, MFA)
Urdu	31-45	female	Master's degree (e.g., MA, MS, MBA, MFA)
Urdu	18-30	male	Master's degree (e.g., MA, MS, MBA, MFA)
Urdu	31-45	male	Bachelor's degree (e.g., BA, BS)

779	B2 Participant Guidelines		
780	Overview		
781	This project is meant to evaluate the quality of		
782	translation and localization of various LLM mod-		
783	els when asked to translate marketing emails from		
784	English to a given language and locale. Imagine		
785	that a person working at an advertising agency		
786	is asked to translate a marketing email they are		
787	working on from English to a language and coun-		
788	try that they don't know anything about. They do		
789	what people do these days and go to the internet		
790	and ask their favorite LLM model to "Translate		
791	this email into {{language}} for use in {{coun-		
792	try}}"		
793			
794	You represent their end user, as a person in the		
795	targeted country who speaks the language, we are		
796	asking for. We'd like you to evaluate the email		
797	from the perspective of a person getting that ad-		
798	vertisement in your email. How well is it translat-		
799	ed? How well does it target local traditions and		
800	norms? How true to the original content and tone		
801	is the translation?"		
802			
803	We'll ask several questions, all using the same		
804	basic evaluation scale. Keep these ratings and de-		
805	scriptions in mind while you are evaluating.		
806			
807	• serious failures exist - use this in cases where		
808	you would be very disappointed, confused,		
809	offended, or in some other way have negative		
810	feelings towards the company or product be-		
811	cause of the content of the translation		
812	• imperfect but not terrible - there are errors or		
813	issues that are very noticeable, but that are not		
814	so bad as to give a negative impression of the		
815	company, the main ideas come through and it is		
816	clear what is being advertised.		
817	• mostly good with small issues - the wording		
818	or translations are noticeably non-native, or are		
819	awkward or a little odd, but it is overall some-		
820	thing that makes sense and could be used with-		
821	out embarrassment on the part of the company.		
822	• very good or nearly perfect - this is for some-		
823	thing that seems very close to natural, native,		
824	and culturally appropriate.		
825			
826			
827			
828			
829			
830			
831			
832			
833			
834			
835			
836			
837			
838			
839			
840			
841			
842			
843			
844			
845			
846			
847			
848			
849			
850			
851			
852			
853			
854			
855			
856			
857			
858			
859			
860			
861			
862			
863			
864			
865			
866			
867			
868			
869			
870			
871			

Predictor	Estimate	SE	CI	z_ratio	p_value	Significance
Very good / nearly perfect Mostly good	-0.01	0.36	[-0.71, 0.69]	-0.03	0.975	
Mostly good Imperfect	1.28	0.36	[0.58, 1.99]	3.57	< .001	***
Imperfect Serious failures	2.50	0.36	[1.80, 3.21]	6.95	< .001	***
Serious failures Segment not translated	6.27	0.37	[5.54, 6.99]	16.92	< .001	***
modelCohere Aya Expanse 8B	1.90	0.15	[1.60, 2.20]	12.42	< .001	***
modelDeepSeek V3.1	0.51	0.15	[0.20, 0.81]	3.28	0.001	**
modelGPT-5	0.02	0.16	[-0.28, 0.33]	0.15	0.878	
modelgpt-oss 120b	0.81	0.15	[0.50, 1.11]	5.25	< .001	***
modelLlama 4	1.03	0.15	[0.72, 1.33]	6.68	< .001	***
modelMistral Medium 3.1	0.38	0.15	[0.08, 0.69]	2.48	0.013	*
languageArabic	0.22	0.48	[-0.72, 1.15]	0.45	0.652	
languageBrazilian Portuguese	-0.92	0.48	[-1.87, 0.03]	-1.89	0.058	
languageCantonese	-0.22	0.48	[-1.15, 0.72]	-0.45	0.650	
languageCzech	0.12	0.48	[-0.81, 1.05]	0.25	0.800	
languageDutch	0.29	0.48	[-0.64, 1.23]	0.62	0.538	
languageHebrew	0.16	0.48	[-0.77, 1.10]	0.34	0.735	
languageHindi	0.60	0.47	[-0.32, 1.53]	1.28	0.202	
languageJapanese	-0.29	0.48	[-1.23, 0.64]	-0.61	0.540	
languageKorean	0.14	0.48	[-0.81, 1.09]	0.29	0.771	
languageMandarin	-1.53	0.49	[-2.50, -0.56]	-3.09	0.002	**
languageRussian	-0.53	0.49	[-1.49, 0.44]	-1.07	0.284	
languageSpanish	0.17	0.47	[-0.76, 1.09]	0.35	0.726	
languageSwahili	0.18	0.48	[-0.76, 1.11]	0.37	0.711	
languageUrdu	0.33	0.49	[-0.63, 1.28]	0.67	0.501	
segment_category.L	1.66	0.08	[1.49, 1.82]	19.72	< .001	***
segment_category.Q	0.31	0.10	[0.12, 0.49]	3.22	0.001	**
segment_category.C	-0.84	0.11	[-1.05, -0.63]	-7.79	< .001	***

Table C1 Fixed-effect estimates from the cumulative link mixed model predicting machine translation quality (0–3).

C2 Model-Level Effects

factor	emmean	SE	CI
Claude Sonnet 4	-2.60	0.14	[-2.87, -2.32]
Cohere Aya Expanse 8B	-0.69	0.14	[-0.96, -0.43]
DeepSeek V3.1	-2.09	0.14	[-2.36, -1.82]
GPT-5	-2.57	0.14	[-2.85, -2.29]
gpt-oss 120b	-1.79	0.14	[-2.06, -1.52]
Llama 4	-1.57	0.14	[-1.84, -1.30]
Mistral Medium 3.1	-2.21	0.14	[-2.49, -1.94]

Table C2 Estimated Marginal Means by Model

contrast	estimate	SE	CI	z_ratio	p_value	Significance
Claude Sonnet 4 - Cohere Aya Expans 8B	-1.90	0.15	[-2.36, -1.45]	-12.42	< .001	***
Claude Sonnet 4 - DeepSeek V3.1	-0.51	0.15	[-0.96, -0.05]	-3.28	0.018	*
Claude Sonnet 4 - (GPT-5)	-0.02	0.16	[-0.48, 0.43]	-0.15	1.000	
Claude Sonnet 4 - (gpt-oss 120b)	-0.81	0.15	[-1.26, -0.35]	-5.25	< .001	***
Claude Sonnet 4 - Llama 4	-1.03	0.15	[-1.48, -0.57]	-6.68	< .001	***
Claude Sonnet 4 - Mistral Medium 3.1	-0.38	0.15	[-0.84, 0.07]	-2.48	0.165	
Cohere Aya Expans 8B - DeepSeek V3.1	1.40	0.15	[0.95, 1.84]	9.24	< .001	***
Cohere Aya Expans 8B - (GPT-5)	1.88	0.15	[1.43, 2.33]	12.31	< .001	***
Cohere Aya Expans 8B - (gpt-oss 120b)	1.10	0.15	[0.66, 1.54]	7.32	< .001	***
Cohere Aya Expans 8B - Llama 4	0.88	0.15	[0.43, 1.32]	5.84	< .001	***
Cohere Aya Expans 8B - Mistral Medium 3.1	1.52	0.15	[1.07, 1.97]	10.03	< .001	***
DeepSeek V3.1 - (GPT-5)	0.48	0.15	[0.03, 0.94]	3.13	0.029	*
DeepSeek V3.1 - (gpt-oss 120b)	-0.30	0.15	[-0.75, 0.15]	-1.97	0.432	
DeepSeek V3.1 - Llama 4	-0.52	0.15	[-0.97, -0.07]	-3.42	0.011	*
DeepSeek V3.1 - Mistral Medium 3.1	0.12	0.15	[-0.33, 0.57]	0.80	0.985	
(GPT-5) - (gpt-oss 120b)	-0.78	0.15	[-1.23, -0.33]	-5.14	< .001	***
(GPT-5) - Llama 4	-1.00	0.15	[-1.45, -0.55]	-6.54	< .001	***
(GPT-5) - Mistral Medium 3.1	-0.36	0.15	[-0.81, 0.09]	-2.34	0.227	
(gpt-oss 120b) - Llama 4	-0.22	0.15	[-0.67, 0.22]	-1.46	0.766	
(gpt-oss 120b) - Mistral Medium 3.1	0.42	0.15	[-0.03, 0.87]	2.78	0.080	
Llama 4 - Mistral Medium 3.1	0.64	0.15	[0.19, 1.09]	4.22	< .001	***

Table C3 Pairwise Model Comparisons (Tukey-adjusted)

C3 Language-Level Effects

factor	emmean	SE	CI
Afrikaans	-1.85	0.34	[-2.52, -1.17]
Arabic	-1.63	0.34	[-2.29, -0.97]
Brazilian Portuguese	-2.76	0.35	[-3.44, -2.09]
Cantonese	-2.06	0.33	[-2.72, -1.41]
Czech	-1.73	0.33	[-2.38, -1.08]
Dutch	-1.55	0.34	[-2.21, -0.89]
Hebrew	-1.68	0.33	[-2.34, -1.03]
Hindi	-1.24	0.33	[-1.89, -0.60]
Japanese	-2.14	0.34	[-2.80, -1.48]
Korean	-1.71	0.34	[-2.38, -1.03]
Mandarin	-3.37	0.36	[-4.07, -2.67]
Russian	-2.37	0.35	[-3.06, -1.68]
Spanish	-1.68	0.33	[-2.33, -1.03]
Swahili	-1.67	0.33	[-2.32, -1.01]
Urdu	-1.52	0.35	[-2.20, -0.84]

Table C4 Estimated Marginal Means by Language

contrast	estimate	SE	CI	z ratio	p_value	Significance
Afrikaans - Arabic	-0.22	0.48	[-1.84, 1.41]	-0.45	1.000	
Afrikaans - Brazilian Portuguese	0.92	0.48	[-0.73, 2.56]	1.89	0.857	
Afrikaans - Cantonese	0.22	0.48	[-1.40, 1.84]	0.45	1.000	
Afrikaans - Czech	-0.12	0.48	[-1.73, 1.49]	-0.25	1.000	
Afrikaans - Dutch	-0.29	0.48	[-1.91, 1.33]	-0.62	1.000	
Afrikaans - Hebrew	-0.16	0.48	[-1.78, 1.46]	-0.34	1.000	
Afrikaans - Hindi	-0.60	0.47	[-2.21, 1.00]	-1.28	0.995	
Afrikaans - Japanese	0.29	0.48	[-1.33, 1.91]	0.61	1.000	
Afrikaans - Korean	-0.14	0.48	[-1.78, 1.50]	-0.29	1.000	
Afrikaans - Mandarin	1.53	0.49	[-0.15, 3.20]	3.09	0.120	
Afrikaans - Russian	0.53	0.49	[-1.14, 2.19]	1.07	0.999	
Afrikaans - Spanish	-0.17	0.47	[-1.77, 1.44]	-0.35	1.000	
Afrikaans - Swahili	-0.18	0.48	[-1.80, 1.44]	-0.37	1.000	
Afrikaans - Urdu	-0.33	0.49	[-1.98, 1.32]	-0.67	1.000	
Arabic - Brazilian Portuguese	1.13	0.48	[-0.48, 2.75]	2.38	0.532	
Arabic - Cantonese	0.43	0.47	[-1.16, 2.03]	0.92	1.000	
Arabic - Czech	0.10	0.47	[-1.49, 1.68]	0.20	1.000	
Arabic - Dutch	-0.08	0.47	[-1.67, 1.51]	-0.17	1.000	
Arabic - Hebrew	0.05	0.47	[-1.53, 1.64]	0.12	1.000	
Arabic - Hindi	-0.39	0.46	[-1.96, 1.18]	-0.84	1.000	
Arabic - Japanese	0.51	0.47	[-1.09, 2.11]	1.08	0.999	
Arabic - Korean	0.07	0.48	[-1.54, 1.69]	0.16	1.000	
Arabic - Mandarin	1.74	0.49	[0.08, 3.41]	3.56	0.029	*
Arabic - Russian	0.74	0.49	[-0.91, 2.39]	1.52	0.973	
Arabic - Spanish	0.05	0.46	[-1.52, 1.62]	0.11	1.000	
Arabic - Swahili	0.04	0.47	[-1.56, 1.64]	0.08	1.000	
Arabic - Urdu	-0.11	0.48	[-1.74, 1.52]	-0.23	1.000	
Brazilian Portuguese - Cantonese	-0.70	0.48	[-2.32, 0.92]	-1.47	0.980	
Brazilian Portuguese - Czech	-1.04	0.47	[-2.65, 0.57]	-2.19	0.672	
Brazilian Portuguese - Dutch	-1.21	0.48	[-2.83, 0.40]	-2.54	0.409	
Brazilian Portuguese - Hebrew	-1.08	0.48	[-2.70, 0.54]	-2.27	0.616	
Brazilian Portuguese - Hindi	-1.52	0.47	[-3.12, 0.08]	-3.23	0.082	
Brazilian Portuguese - Japanese	-0.63	0.48	[-2.24, 0.99]	-1.31	0.993	
Brazilian Portuguese - Korean	-1.06	0.48	[-2.70, 0.58]	-2.19	0.672	
Brazilian Portuguese - Mandarin	0.61	0.49	[-1.07, 2.29]	1.23	0.996	
Brazilian Portuguese - Russian	-0.39	0.49	[-2.06, 1.27]	-0.80	1.000	
Brazilian Portuguese - Spanish	-1.08	0.47	[-2.69, 0.52]	-2.29	0.596	
Brazilian Portuguese - Swahili	-1.10	0.48	[-2.72, 0.53]	-2.29	0.596	
Brazilian Portuguese - Urdu	-1.25	0.49	[-2.90, 0.41]	-2.56	0.398	
Cantonese - Czech	-0.34	0.47	[-1.92, 1.25]	-0.72	1.000	
Cantonese - Dutch	-0.51	0.47	[-2.10, 1.08]	-1.09	0.999	
Cantonese - Hebrew	-0.38	0.47	[-1.97, 1.21]	-0.81	1.000	
Cantonese - Hindi	-0.82	0.46	[-2.40, 0.76]	-1.76	0.912	
Cantonese - Japanese	0.08	0.47	[-1.52, 1.67]	0.16	1.000	
Cantonese - Korean	-0.36	0.48	[-1.97, 1.26]	-0.75	1.000	
Cantonese - Mandarin	1.31	0.49	[-0.34, 2.96]	2.69	0.310	
Cantonese - Russian	0.31	0.48	[-1.33, 1.95]	0.64	1.000	
Cantonese - Spanish	-0.38	0.47	[-1.96, 1.20]	-0.82	1.000	
Cantonese - Swahili	-0.39	0.47	[-1.99, 1.20]	-0.84	1.000	
Cantonese - Urdu	-0.54	0.48	[-2.17, 1.08]	-1.13	0.999	
Czech - Dutch	-0.17	0.47	[-1.75, 1.41]	-0.37	1.000	
Czech - Hebrew	-0.04	0.47	[-1.62, 1.54]	-0.09	1.000	
Czech - Hindi	-0.48	0.46	[-2.05, 1.08]	-1.05	0.999	

Table C5 Pairwise Language Comparisons (Tukey-adjusted) *continued on next page*

Czech - Japanese	0.41	0.47	[-1.17, 2.00]	0.88	1.000	
Czech - Korean	-0.02	0.47	[-1.63, 1.59]	-0.04	1.000	
Czech - Mandarin	1.65	0.49	[0.00, 3.30]	3.39	0.050	*
Czech - Russian	0.65	0.48	[-0.99, 2.28]	1.34	0.992	
Czech - Spanish	-0.05	0.46	[-1.61, 1.52]	-0.10	1.000	
Czech - Swahili	-0.06	0.47	[-1.64, 1.53]	-0.12	1.000	
Czech - Urdu	-0.21	0.48	[-1.82, 1.41]	-0.43	1.000	
Dutch - Hebrew	0.13	0.47	[-1.46, 1.72]	0.28	1.000	
Dutch - Hindi	-0.31	0.46	[-1.88, 1.26]	-0.67	1.000	
Dutch - Japanese	0.59	0.47	[-1.01, 2.18]	1.25	0.996	
Dutch - Korean	0.15	0.48	[-1.46, 1.77]	0.32	1.000	
Dutch - Mandarin	1.82	0.49	[0.17, 3.48]	3.73	0.016	*
Dutch - Russian	0.82	0.48	[-0.82, 2.46]	1.69	0.936	
Dutch - Spanish	0.13	0.46	[-1.45, 1.70]	0.28	1.000	
Dutch - Swahili	0.12	0.47	[-1.48, 1.71]	0.25	1.000	
Dutch - Urdu	-0.03	0.48	[-1.66, 1.59]	-0.07	1.000	
Hebrew - Hindi	-0.44	0.46	[-2.01, 1.13]	-0.95	1.000	
Hebrew - Japanese	0.45	0.47	[-1.14, 2.05]	0.97	1.000	
Hebrew - Korean	0.02	0.48	[-1.59, 1.63]	0.04	1.000	
Hebrew - Mandarin	1.69	0.49	[0.03, 3.34]	3.46	0.040	*
Hebrew - Russian	0.69	0.48	[-0.95, 2.33]	1.42	0.986	
Hebrew - Spanish	0.00	0.46	[-1.58, 1.57]	-0.01	1.000	
Hebrew - Swahili	-0.02	0.47	[-1.61, 1.58]	-0.03	1.000	
Hebrew - Urdu	-0.17	0.48	[-1.79, 1.46]	-0.35	1.000	
Hindi - Japanese	0.90	0.47	[-0.68, 2.47]	1.93	0.840	
Hindi - Korean	0.46	0.47	[-1.13, 2.06]	0.98	1.000	
Hindi - Mandarin	2.13	0.48	[0.49, 3.77]	4.40	0.001	**
Hindi - Russian	1.13	0.48	[-0.50, 2.76]	2.35	0.551	
Hindi - Spanish	0.44	0.46	[-1.12, 1.99]	0.95	1.000	
Hindi - Swahili	0.43	0.47	[-1.15, 2.01]	0.92	1.000	
Hindi - Urdu	0.28	0.47	[-1.33, 1.89]	0.58	1.000	
Japanese - Korean	-0.43	0.48	[-2.05, 1.18]	-0.91	1.000	
Japanese - Mandarin	1.23	0.49	[-0.42, 2.89]	2.53	0.416	
Japanese - Russian	0.23	0.48	[-1.41, 1.87]	0.48	1.000	
Japanese - Spanish	-0.46	0.47	[-2.04, 1.12]	-0.98	1.000	
Japanese - Swahili	-0.47	0.47	[-2.07, 1.13]	-1.00	1.000	
Japanese - Urdu	-0.62	0.48	[-2.25, 1.01]	-1.29	0.994	
Korean - Mandarin	1.67	0.49	[-0.01, 3.34]	3.38	0.052	
Korean - Russian	0.67	0.49	[-1.00, 2.33]	1.36	0.991	
Korean - Spanish	-0.02	0.47	[-1.63, 1.58]	-0.05	1.000	
Korean - Swahili	-0.04	0.48	[-1.65, 1.58]	-0.08	1.000	
Korean - Urdu	-0.19	0.49	[-1.83, 1.46]	-0.38	1.000	
Mandarin - Russian	-1.00	0.50	[-2.69, 0.69]	-2.01	0.793	
Mandarin - Spanish	-1.69	0.48	[-3.34, -0.05]	-3.50	0.036	*
Mandarin - Swahili	-1.70	0.49	[-3.36, -0.05]	-3.50	0.035	*
Mandarin - Urdu	-1.86	0.50	[-3.54, -0.17]	-3.74	0.015	*
Russian - Spanish	-0.69	0.48	[-2.32, 0.94]	-1.44	0.984	
Russian - Swahili	-0.70	0.48	[-2.34, 0.94]	-1.45	0.982	
Russian - Urdu	-0.85	0.49	[-2.52, 0.82]	-1.73	0.923	
Spanish - Swahili	-0.01	0.47	[-1.59, 1.57]	-0.02	1.000	
Spanish - Urdu	-0.16	0.48	[-1.77, 1.45]	-0.34	1.000	
Swahili - Urdu	-0.15	0.48	[-1.78, 1.48]	-0.31	1.000	

Table C5 Pairwise Language Comparisons (Tukey-adjusted)

C4 Segment Category Effects

factor	emmean	SE	CI
cultural concepts	-2.70	0.10	[-2.90, -2.50]
holidays	-3.02	0.14	[-3.28, -2.75]
idioms	-1.15	0.14	[-1.43, -0.88]
puns	-0.85	0.13	[-1.10, -0.60]

Table C6 Estimated Marginal Means by Segment Category

contrast	estimate	SE	CI	z_ratio	p_value	Significance
cultural concepts - holidays	0.31	0.12	[0.01, 0.62]	2.67	0.039	*
cultural concepts - idioms	-1.55	0.13	[-1.88, -1.22]	-12.13	< .001	***
cultural concepts - puns	-1.85	0.11	[-2.14, -1.56]	-16.42	< .001	***
holidays - idioms	-1.86	0.16	[-2.27, -1.46]	-11.90	< .001	***
holidays - puns	-2.16	0.14	[-2.54, -1.79]	-14.95	< .001	***
idioms - puns	-0.30	0.15	[-0.69, 0.09]	-2.00	0.188	

Table C7 Pairwise Category Comparisons (Tukey-adjusted)

873

C5 Inter-Rater Reliability

874

875

876

877

878

879

880

881

Inter-rater reliability (IRR) was assessed to evaluate the consistency of human ratings of translation quality across participants. Because ratings were ordinal (e.g., ranging from “very good / nearly perfect” to “serious failures”) and involved multiple raters, we selected complementary reliability measures to capture different aspects of agreement.

882

883

884

885

886

887

888

889

890

We report Krippendorff’s α (ordinal), which is designed for ordered categorical data and is robust to missing values, providing a single coefficient reflecting agreement beyond chance. We additionally report Gwet’s AC2 with quadratic weights, which accounts for chance agreement while being less sensitive to prevalence and marginal distributions than Cohen’s κ . Quadratic weights penalise larger disagreements more heavily, reflecting the

ordered structure of the rating scale. Observed and expected agreement rates derived from AC2 are also reported to aid interpretation of reliability in terms of raw concordance.

Ratings corresponding to “segment not translated” (NA) were excluded from all IRR calculations, as they reflect missing or invalid quality judgments rather than graded assessments. IRR was computed at multiple levels, including overall agreement across all languages, models, and segment categories, as well as stratified by language, model, and segment category (cultural concepts, holidays, idioms, and puns).

IRR calculations were based on item \times rater matrices constructed from the cleaned data and were implemented in R using the irr and irrCAC packages.

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

metric	estimate	lower 95 ci	upper 95 ci	observed agreement	expected agreement	scope	language	model	segment category
Krippendorff_alpha	0.448197144	NA	NA	NA	NA	Overall	NA	NA	NA
Gwet_AC1_weighted	0.41225	(0.31,0.514)	NA	0.755857523	0.584613092	Overall	NA	NA	NA
Krippendorff_alpha	0.498850973	NA	NA	NA	NA	Afrikaans	Afrikaans	NA	NA
Gwet_AC1_weighted	0.14534	(-0.098,0.389)	NA	0.632263084	0.569726302	Afrikaans	Afrikaans	NA	NA
Krippendorff_alpha	0.551735695	NA	NA	NA	NA	Arabic	Arabic	NA	NA
Gwet_AC1_weighted	0.61952	(0.193,1)	NA	0.849890557	0.605471591	Arabic	Arabic	NA	NA
Krippendorff_alpha	0.354424333	NA	NA	NA	NA	Brazilian Portuguese	Brazilian Portuguese	NA	NA
Gwet_AC1_weighted	0.56154	(0.169,0.955)	NA	0.83130482	0.615255895	Brazilian Portuguese	Brazilian Portuguese	NA	NA
Krippendorff_alpha	0.386155192	NA	NA	NA	NA	Cantonese	Cantonese	NA	NA
Gwet_AC1_weighted	0.58182	(0.14,1)	NA	0.824587744	0.580530558	Cantonese	Cantonese	NA	NA
Krippendorff_alpha	0.501678657	NA	NA	NA	NA	Czech	Czech	NA	NA
Gwet_AC1_weighted	0.41474	(-0.014,0.843)	NA	0.731120638	0.540580132	Czech	Czech	NA	NA
Krippendorff_alpha	0.57461557	NA	NA	NA	NA	Dutch	Dutch	NA	NA
Gwet_AC1_weighted	0.55692	(0.085,1)	NA	0.768596935	0.477734454	Dutch	Dutch	NA	NA
Krippendorff_alpha	0.525872162	NA	NA	NA	NA	Hebrew	Hebrew	NA	NA
Gwet_AC1_weighted	0.52	(0.037,1)	NA	0.788096253	0.558531884	Hebrew	Hebrew	NA	NA
Krippendorff_alpha	0.269765185	NA	NA	NA	NA	Hindi	Hindi	NA	NA
Gwet_AC1_weighted	0.64013	(0.201,1)	NA	0.833581517	0.537553424	Hindi	Hindi	NA	NA
Krippendorff_alpha	0.476073172	NA	NA	NA	NA	Japanese	Japanese	NA	NA
Gwet_AC1_weighted	0.42556	(0,0.851)	NA	0.780653592	0.618154078	Japanese	Japanese	NA	NA
Krippendorff_alpha	0.512877424	NA	NA	NA	NA	Korean	Korean	NA	NA
Gwet_AC1_weighted	0.29641	(-0.134,0.727)	NA	0.770063675	0.673197163	Korean	Korean	NA	NA
Krippendorff_alpha	0.267738681	NA	NA	NA	NA	Mandarin	Mandarin	NA	NA
Gwet_AC1_weighted	0.60908	(0.341,0.877)	NA	0.806393163	0.504740189	Mandarin	Mandarin	NA	NA
Krippendorff_alpha	0.372753672	NA	NA	NA	NA	Russian	Russian	NA	NA
Gwet_AC1_weighted	0.3071	(-0.202,0.817)	NA	0.745234394	0.63231795	Russian	Russian	NA	NA
Krippendorff_alpha	0.375234491	NA	NA	NA	NA	Spanish	Spanish	NA	NA
Gwet_AC1_weighted	0.52538	(0.028,1)	NA	0.808558288	0.596638428	Spanish	Spanish	NA	NA
Krippendorff_alpha	0.485056033	NA	NA	NA	NA	Swahili	Swahili	NA	NA
Gwet_AC1_weighted	0.24906	(-0.315,0.813)	NA	0.714105052	0.619284773	Swahili	Swahili	NA	NA
Krippendorff_alpha	0.386782145	NA	NA	NA	NA	Urdu	Urdu	NA	NA
Gwet_AC1_weighted	0.13664	(-0.029,0.302)	NA	0.641737452	0.585034893	Urdu	Urdu	NA	NA

Table C8 Inter-rater reliability statistics for segment-level MT quality ratings *continued on next page*

Krippendorff_ alpha	0.362971562	NA	NA	NA	NA	Claude Sonnet 4	NA	Claude Sonnet 4	NA
Gwet_AC1_ weighted	0.42987	(0.333,0.527)	NA	0.778941763	0.612268497	Claude Sonnet 4	NA	Claude Sonnet 4	NA
Krippendorff_ alpha	0.591731477	NA	NA	NA	NA	Cohere Aya Expanse 8B	NA	Cohere Aya Expanse 8B	NA
Gwet_AC1_ weighted	0.22705	(0.112,0.342)	NA	0.730289925	0.651062192	Cohere Aya Expanse 8B	NA	Cohere Aya Expanse 8B	NA
Krippendorff_ alpha	0.365021429	NA	NA	NA	NA	DeepSeek V3.1	NA	DeepSeek V3.1	NA
Gwet_AC1_ weighted	0.23255	(0.142,0.323)	NA	0.709368798	0.621304321	DeepSeek V3.1	NA	DeepSeek V3.1	NA
Krippendorff_ alpha	0.390678454	NA	NA	NA	NA	GPT-5	NA	GPT-5	NA
Gwet_AC1_ weighted	0.42612	(0.325,0.527)	NA	0.778789507	0.614534099	GPT-5	NA	GPT-5	NA
Krippendorff_ alpha	0.425609176	NA	NA	NA	NA	gpt-oss 120B	NA	gpt-oss 120B	NA
Gwet_AC1_ weighted	0.11716	(0.057,0.178)	NA	0.688983685	0.647708789	gpt-oss 120B	NA	gpt-oss 120B	NA
Krippendorff_ alpha	0.492022726	NA	NA	NA	NA	Llama 4	NA	Llama 4	NA
Gwet_AC1_ weighted	0.1872	(0.103,0.271)	NA	0.715267408	0.649690855	Llama 4	NA	Llama 4	NA
Krippendorff_ alpha	0.353854999	NA	NA	NA	NA	Mistral Medium 3.1	NA	Mistral Medium 3.1	NA
Gwet_AC1_ weighted	0.29617	(0.209,0.384)	NA	0.747983902	0.641936198	Mistral Medium 3.1	NA	Mistral Medium 3.1	NA
Krippendorff_ alpha	0.441472615	NA	NA	NA	NA	cultural concepts	NA	NA	cultural concepts
Gwet_AC1_ weighted	0.34828	(0.28,0.417)	NA	0.745008769	0.608740413	cultural concepts	NA	NA	cultural concepts
Krippendorff_ alpha	0.380075728	NA	NA	NA	NA	holidays	NA	NA	holidays
Gwet_AC1_ weighted	0.40557	(0.305,0.507)	NA	0.733721118	0.552041014	holidays	NA	NA	holidays
Krippendorff_ alpha	0.404880664	NA	NA	NA	NA	idioms	NA	NA	idioms
Gwet_AC1_ weighted	0.10721	(0.048,0.166)	NA	0.737554455	0.706039385	idioms	NA	NA	idioms
Krippendorff_ alpha	0.307338984	NA	NA	NA	NA	puns	NA	NA	puns
Gwet_AC1_ weighted	0.26271	(0.16,0.366)	NA	0.757788673	0.671483717	puns	NA	NA	puns

Table C8 Inter-rater reliability statistics for segment-level MT quality ratings

language	model	alpha	ac2	pairwise_agree	strict_agree	n_items	n_raters
Arabic	Cohere Aya Expans 8B	0.674813037	NA	55.1	26.1	23	5
Japanese	Llama 4	0.65637168	NA	56.7	36.8	20	5
Czech	Llama 4	0.65207732	NA	53.9	24	25	5
Hebrew	gpt-oss 120B	0.644062377	NA	52.4	20	25	5
Arabic	Llama 4	0.634999976	NA	56.1	24	25	5
Arabic	Claude Sonnet 4	0.634900605	NA	64.7	36	25	5
Urdu	Cohere Aya Expans 8B	0.617955706	NA	52.9	50	19	5
Hebrew	DeepSeek V3.1	0.614170634	NA	63.8	40	25	5
Hebrew	Cohere Aya Expans 8B	0.611726618	NA	57.3	17.4	23	5
Japanese	Cohere Aya Expans 8B	0.584984776	NA	54.5	21.7	23	5
Dutch	Cohere Aya Expans 8B	0.577023671	NA	51.4	31.6	24	5
Dutch	gpt-oss 120B	0.564995102	NA	57.5	27.3	25	5
Dutch	Claude Sonnet 4	0.555536604	NA	59.5	36.4	25	5
Cantonese	Llama 4	0.550254155	NA	45.3	17.4	23	5
Arabic	Mistral Medium 3.1	0.544804854	NA	60	32	25	5
Czech	gpt-oss 120B	0.528554281	NA	53.2	34.8	25	5
Hebrew	Llama 4	0.518538232	NA	49	25	24	5
Arabic	gpt-oss 120B	0.513416055	NA	50.4	12	25	5
Korean	GPT-5	0.512390998	NA	45.7	8	25	5
Czech	Cohere Aya Expans 8B	0.510046027	NA	47.4	14.3	22	5
Hebrew	Claude Sonnet 4	0.502641466	NA	47.5	28	25	5
Dutch	DeepSeek V3.1	0.493616221	NA	52.7	23.8	25	5
Dutch	Mistral Medium 3.1	0.493573969	NA	61.8	40.9	25	5
Korean	DeepSeek V3.1	0.486751851	NA	45	16	25	5
Arabic	GPT-5	0.481986498	NA	59.2	32	25	5
Korean	Cohere Aya Expans 8B	0.474465656	NA	45.6	20	25	5
Spanish	Llama 4	0.47124898	NA	45.1	12.5	25	5
Dutch	Llama 4	0.47021559	NA	55.2	28.6	25	5
Czech	Claude Sonnet 4	0.457016233	NA	59.7	37.5	24	5
Afrikaans	GPT-5	0.456513385	NA	71.4	47.6	24	5
Russian	Cohere Aya Expans 8B	0.448616905	NA	44.4	20	25	5
Afrikaans	gpt-oss 120B	0.441137352	NA	61.1	36.4	24	5
Cantonese	Cohere Aya Expans 8B	0.439844702	NA	38.1	16	25	5
Korean	Llama 4	0.432970093	NA	47.8	8.7	25	5
Czech	GPT-5	0.428335745	NA	59.6	36	25	5
Brazilian Portuguese	gpt-oss 120B	0.426784937	NA	62	34.8	25	5
Korean	gpt-oss 120B	0.411637737	NA	43.6	8	25	5
Russian	gpt-oss 120B	0.402388898	NA	41.6	20	25	5
Afrikaans	Cohere Aya Expans 8B	0.402248913	NA	41.7	27.3	23	5
Brazilian Portuguese	DeepSeek V3.1	0.400921077	NA	58.1	36	25	5
Mandarin	Llama 4	0.396096645	NA	54.1	28	25	5
Swahili	Llama 4	0.392872584	NA	52.9	30.4	25	5
Czech	Mistral Medium 3.1	0.391101109	NA	55.8	20.8	25	5
Brazilian Portuguese	Claude Sonnet 4	0.389526749	NA	64	44	25	5

Table C9 Inter-rater reliability statistics for holistic text MT quality ratings *continued on next page*

Spanish	GPT-5	0.384567319	NA	48.2	20	25	5
Hebrew	GPT-5	0.382856402	NA	57.6	28	25	5
Japanese	gpt-oss 120B	0.382726192	NA	51.8	28	25	5
Russian	DeepSeek V3.1	0.38161071	NA	46.8	17.4	24	5
Hindi	Mistral Medium 3.1	0.380938245	NA	38.5	8	25	5
Korean	Claude Sonnet 4	0.377620246	NA	46.6	12	25	5
Cantonese	Mistral Medium 3.1	0.371804013	NA	44.8	16	25	5
Cantonese	GPT-5	0.360470042	NA	48.6	20	25	5
Cantonese	gpt-oss 120B	0.348227295	NA	45.5	20.8	24	5
Brazilian Portuguese	Cohere Aya Expans 8B	0.347385294	NA	54.6	29.2	24	5
Mandarin	Cohere Aya Expans 8B	0.345099047	NA	46.1	16.7	24	5
Russian	Llama 4	0.333786874	NA	41.6	12	25	5
Brazilian Portuguese	Llama 4	0.321316883	NA	51	20	25	5
Hindi	Cohere Aya Expans 8B	0.318208174	NA	46.8	18.2	23	5
Afrikaans	Claude Sonnet 4	0.306930278	NA	63	34.8	25	5
Arabic	DeepSeek V3.1	0.306619915	NA	44.7	12	25	5
Hindi	GPT-5	0.305698654	NA	48.3	20	25	5
Spanish	Claude Sonnet 4	0.301239732	NA	47.4	12	25	5
Czech	DeepSeek V3.1	0.29960442	NA	52.3	24	25	5
Urdu	GPT-5	0.288346858	NA	51.7	29.4	18	5
Urdu	Claude Sonnet 4	0.282773726	NA	67.1	50	20	5
Spanish	gpt-oss 120B	0.281928166	NA	43	13	25	5
Mandarin	gpt-oss 120B	0.278237639	NA	57.2	24	25	5
Russian	Mistral Medium 3.1	0.277913363	NA	49.2	20.8	25	5
Japanese	Mistral Medium 3.1	0.270812278	NA	60.1	33.3	24	5
Swahili	Mistral Medium 3.1	0.270118901	NA	43.1	16	25	5
Japanese	Claude Sonnet 4	0.263590307	NA	58.6	28	25	5
Swahili	Claude Sonnet 4	0.2617325	NA	52.4	28	25	5
Spanish	DeepSeek V3.1	0.257918185	NA	42.4	8	25	5
Russian	Claude Sonnet 4	0.257401126	NA	55.1	28	25	5
Spanish	Cohere Aya Expans 8B	0.257126966	NA	41.9	13	25	5
Afrikaans	Mistral Medium 3.1	0.256137166	NA	47.8	25	24	5
Cantonese	Claude Sonnet 4	0.251197556	NA	51.4	24	25	5
Hindi	Llama 4	0.24301364	NA	43.5	16	25	5
Korean	Mistral Medium 3.1	0.242460239	NA	48.1	12.5	24	5
Hebrew	Mistral Medium 3.1	0.239619599	NA	44.8	20	25	5
Afrikaans	DeepSeek V3.1	0.232375967	NA	55.7	37.5	25	5
Hindi	DeepSeek V3.1	0.229249261	NA	41	12	25	5
Brazilian Portuguese	Mistral Medium 3.1	0.212810949	NA	59	32	25	5
Urdu	Llama 4	0.209767274	NA	54.4	35	20	5
Cantonese	DeepSeek V3.1	0.208166722	NA	50.9	16	25	5
Swahili	gpt-oss 120B	0.201556295	NA	38.8	21.7	25	5
Dutch	GPT-5	0.201535135	NA	52.2	26.1	25	5
Hindi	gpt-oss 120B	0.195261741	NA	41.8	12	25	5
Spanish	Mistral Medium 3.1	0.194733517	NA	47.6	20	25	5

Table C9 Inter-rater reliability statistics for holistic text MT quality ratings *continued on next page*

Brazilian Portuguese	GPT-5	0.189900439	NA	64.3	32	25	5
Japanese	GPT-5	0.181908943	NA	59.3	29.2	24	5
Urdu	Mistral Medium 3.1	0.179128246	NA	54.4	25	20	5
Japanese	DeepSeek V3.1	0.17029338	NA	47.9	20	25	5
Swahili	GPT-5	0.168215451	NA	60.1	37.5	24	5
Hindi	Claude Sonnet 4	0.165242137	NA	38.1	8	25	5
Swahili	DeepSeek V3.1	0.139262956	NA	49.6	16	25	5
Swahili	Cohere Aya Expans 8B	0.136498089	NA	73	33.3	25	5
Urdu	gpt-oss 120B	0.128582875	NA	38	10.5	20	5
Afrikaans	Llama 4	0.121541552	NA	55.2	21.7	23	5
Mandarin	DeepSeek V3.1	0.114767658	NA	65.1	35	20	5
Russian	GPT-5	0.087226249	NA	48	24	25	5
Mandarin	Mistral Medium 3.1	0.080766145	NA	60.3	32	25	5
Mandarin	Claude Sonnet 4	0.072938315	NA	52.4	20	25	5
Urdu	DeepSeek V3.1	0.017972445	NA	52.1	30	20	5

Table C9 Inter-rater reliability statistics for holistic text MT quality ratings