

# LiveCLKTBench: Towards Reliable Evaluation of Cross-Lingual Knowledge Transfer in Multilingual LLMs

Pei-Fu Guo<sup>1</sup>, Yun-Da Tsai<sup>1</sup>, Chun-Chia Hsu<sup>1</sup>, Kai-Xin Chen<sup>1</sup>, Ya-An Tsai<sup>1</sup>,  
Kai-Wei Chang<sup>2</sup>, Nanyun Peng<sup>2</sup>, Mi-Yen Yeh<sup>3</sup>, Shou-De Lin<sup>1,4</sup>

<sup>1</sup>National Taiwan University <sup>2</sup>University of California, Los Angeles

<sup>3</sup>Academia Sinica, Taiwan <sup>4</sup>NTU AI-CoRE

Correspondence: r12922217@csie.ntu.edu.tw

## Abstract

Evaluating cross-lingual knowledge transfer in large language models (LLMs) is challenging, as correct answers in a target language may arise either from genuine transfer or from prior exposure during pre-training. We present LiveCLKTBench, an automated generation pipeline specifically designed to isolate and measure cross-lingual knowledge transfer. Our pipeline identifies self-contained, time-sensitive knowledge entities from real-world domains, filters them based on temporal occurrence, and verifies them against the model’s knowledge. The documents of these valid entities are then used to generate factual questions, which are translated into multiple languages to evaluate transferability across linguistic boundaries. Using LiveCLKTBench, we evaluate several LLMs across five languages and observe that cross-lingual transfer is strongly influenced by linguistic distance and often asymmetric across language directions. While larger models improve transfer, the gains diminish with scale and vary across domains. These findings provide new insights into multilingual transfer and demonstrate the value of LiveCLKTBench as a reliable benchmark for future research.<sup>1</sup>

## 1 Introduction

As large language models (LLMs) continue to grow in scale and capability, a central question arises: *How can they serve users globally and equitably?* Ideally, an LLM should be able to transfer knowledge acquired in one language to others, rather than relearning the same facts separately across languages — a property known as *cross-lingual knowledge transfer* (Hu et al., 2020a; Lauscher et al., 2020; Hedderich et al., 2020).

However, reliably evaluating this ability remains challenging. As LLMs are pretrained on massive multilingual corpora that may already contain the same factual knowledge in multiple languages, it is

often unclear whether correct answers in the target language reflect genuine cross-lingual transfer or merely memorization of previously seen information. To address this issue, we identify three essential properties that a reliable cross-lingual transfer benchmark should satisfy.

**(1) Leakage-Free Evaluation.** A major obstacle in evaluating cross-lingual knowledge transfer is contamination and data leakage (Ahuja et al., 2024). When pretrained corpora already include the same knowledge in multiple languages, models can appear to “transfer” information they have merely memorized. A leakage-free benchmark must therefore ensure that the tested knowledge is unseen in the target language, so that correct answers truly reflect transfer.

**(2) Grounding in Real-World Knowledge.** Some benchmarks employ synthetic or fictitious data to control knowledge injection and reduce contamination (Maini et al., 2024; Maheshwari et al., 2024). While this strategy offers strong controllability, recent works have shown that fabricated facts can conflict with pre-existing knowledge and degrade performance (Wu et al., 2024; Jan et al., 2025; Chen et al., 2024). Consequently, when a model fails to answer a fictitious question in a particular language, the failure may not stem from a lack of transferability, but rather from the artificial knowledge being too weakly or even adversarially connected to existing knowledge to support meaningful transfer. We therefore argue that benchmarks grounded in real documents better reflect practical scenarios, where models must absorb genuine knowledge and transfer it across languages while leveraging their existing background knowledge.

**(3) Frequent Knowledge Update.** As LLMs evolve rapidly, static benchmarks become outdated: newer models may already encode most benchmark facts, causing score saturation and

<sup>1</sup>Code and data are available at [link](#).

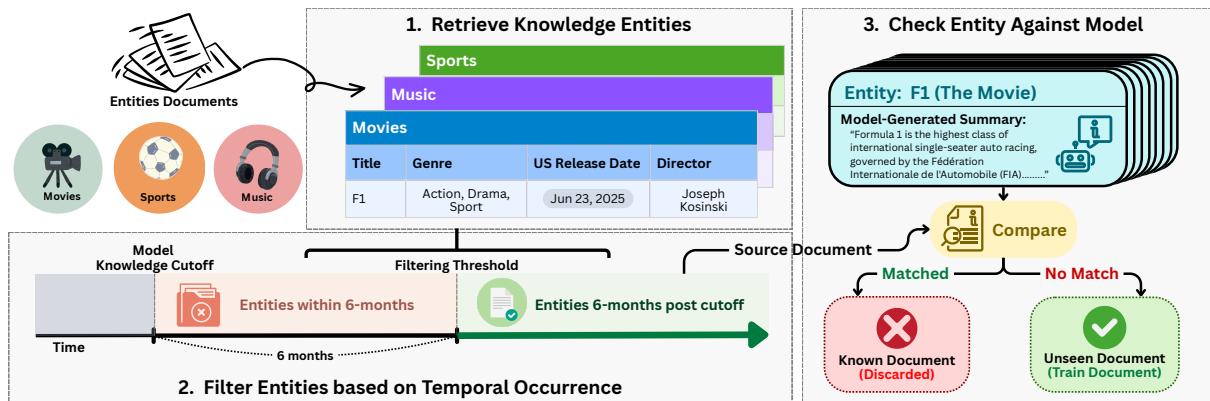


Figure 1: **Leakage Prevention in LiveCLKTBench.** The pipeline prevents data leakage by selecting valid knowledge entities that contain facts unknown to pretrained models. Specifically, it identifies independent, time-sensitive real-world entities, filters them by temporal occurrence, and cross-checks them against model outputs to eliminate any entities already known to pretrained models.

reducing diagnostic value <sup>2</sup>. Recent initiatives such as REALTIMEQA (Kasai et al., 2024) and LIVEBENCH (White et al., 2025) highlight the importance of designing benchmarks that are regularly refreshed with new data to prevent staleness.

Building on these principles, we introduce *LiveCLKTBench*, an automated pipeline for generating realistic, contamination-free, and continuously refreshable benchmarks for cross-lingual knowledge transfer. As shown in Figure 1, given a target model, our pipeline minimizes the risk of data leakage through the following strategies:

1. Identifying independent, time-sensitive *knowledge entities* from three frequently updated real-world domains—*movies*, *music* and *sports* (e.g., a new release movie, a recent baseball game match).
2. Filtering entities based on temporal occurrence, retaining only those appearing at least six months after the model’s knowledge cutoff to avoid potential prior exposure.
3. Verifying each entity by prompting the model for a factual summary; entities whose responses match the real-world *source document* are treated as known and discarded.

Together, these steps ensure that the resulting QA pairs are contamination-free, grounded in real-world knowledge, and provide a reliable foundation for evaluating cross-lingual knowledge transfer.

For the retained knowledge entities, the pipeline generates factual multiple-choice questions that

are (i) explicitly grounded in the corresponding source documents, and (ii) whose correct answers become available only once the event has taken place (e.g., the final score of a sports match). The verified questions and their corresponding source documents are translated into the evaluation languages. During evaluation, models are post-trained on source-language documents and tested on QA pairs in other languages. Because the knowledge is unseen by the model, the benchmark provides a direct measure of knowledge transfer, eliminating concerns about memorization from pretraining.

In addition to ensuring data integrity, LiveCLKTBench offers full automation and configurability, accommodating new model releases and flexible evaluation parameters. By simply specifying the target models and evaluation languages, researchers can conduct frequent and reproducible evaluations with minimal manual intervention.

Using LiveCLKTBench, we evaluate a series of open-source multilingual LLMs across five languages. Our analysis shows that current models still have substantial room for improvement: transferability is strongly influenced by linguistic distance and often asymmetric across language directions. While larger models generally exhibit better transfer, their gains diminish with scale and vary across domains.

Together, these findings provide new insights into cross-lingual knowledge transfer and highlight the value of LiveCLKTBench as a reliable and sustainable testbed for studying and improving multilingual LLMs.

<sup>2</sup><https://r0bk.github.io/killedbyllm/>

## 2 Related Work

### 2.1 Multilingual Benchmarks

A variety of multilingual benchmarks have been proposed to evaluate model performance across many languages (Wu et al., 2025). Some focus on general natural-language tasks such as classification, question answering, or summarization (e.g., XTREME (Hu et al., 2020b), XGLUE (Liang et al., 2020), MEGA (Ahuja et al., 2023)), while others emphasize factual knowledge access in multiple languages, such as M3Exam (Zhang et al., 2023), AGI-Eval (Zhong et al., 2023), and Global-MMLU (Singh et al., 2025). These benchmarks provide broad coverage of multilingual capabilities but do not directly isolate cross-lingual transfer.

### 2.2 Cross-lingual Transfer Benchmarks

*Cross-lingual transfer* benchmarks test whether *knowledge* or *skills* learned in one language generalize to others. Here, *skill transfer* refers to the generalization of abilities such as summarization or instruction following to unseen languages (Shaham et al., 2024; Chai et al., 2024; Asai et al., 2023), whereas *knowledge transfer* evaluates the reproduction of factual information acquired in one language when queried in another.

Recent studies analyze the unique challenges of *cross-lingual knowledge transfer*, such as language transfer asymmetry and knowledge representation barriers (Rajae and Monz, 2024; Chua et al., 2025; Litschko et al., 2025; Yao et al., 2024). Complementing these analyses, several dedicated benchmarks have been proposed. For example, ECLeKTic (Goldman et al., 2025) benchmarks cross-lingual knowledge transfer by using Wikipedia articles common in one language but rare in others, ensuring answers reflect genuine transfer. Other benchmarks focus on multilingual knowledge editing, evaluating whether updates introduced in one language propagate to others (Nie et al., 2025; Wei et al., 2025).

Another important dimension is the *evaluation protocol*. Some benchmarks use zero-shot evaluation, querying models directly in target languages (Malkin et al., 2022; Goldman et al., 2025), while others adopt a fine-tune-then-test approach, where models are trained in one language and tested in others (Shaham et al., 2024). Work also distinguishes between *parametric knowledge*, stored in model parameters (Goldman et al., 2025; Rajae and Monz, 2024), and *contextual knowl-*

*edge*, provided through in-context examples at inference time (Mondshine et al., 2025; Asai et al., 2021).

### 2.3 Relation to Prior Work

Among prior efforts, ECLeKTic (Goldman et al., 2025) is most closely related to our work, as it evaluates factual transfer across languages while reducing leakage by selecting documents that are more common in one language than another. While this strategy offers partial control, leakage cannot be fully ruled out, and reliance on Wikipedia limits coverage of dynamic knowledge. In contrast, LiveCLKTBench injects new knowledge via post-training on a single source language, so correct answers in other languages must arise from genuine transfer. By sourcing documents from rapidly updating domains and allowing configurable time windows, it produces more diverse and extensible benchmarks that stay relevant as models evolve.

## 3 Methodology

### 3.1 Benchmark Generation Pipeline

Each LiveCLKTBench benchmark consists of two components: (i) a set of source-language training documents that introduce knowledge previously unseen by the model, and (ii) a multilingual test set of document-grounded factual questions derived from those documents. As illustrated in Figure 2, the pipeline proceeds through four main stages, each applying strict constraints to prevent knowledge contamination. Starting with a target model and the evaluation language set, we describe each generation stage in detail below.

**Knowledge Entity Collection.** The first stage collects *knowledge entities* that serve as the basis for generating test set factual QA pairs. Our pipeline focus on three domains—*movies*, *music*, and *sports*. They are chosen because they frequently produce entities with fresh facts that satisfy two key properties:

- **Independence.** Each entity is self-contained, and its facts cannot be inferred from other events. For example, the outcome of one baseball game does not depend on the results of earlier games, allowing for clean separation between training and test instances.
- **Time-sensitivity.** The factual outcome itself (e.g., which team won or the exact score) becomes known only after the match has taken

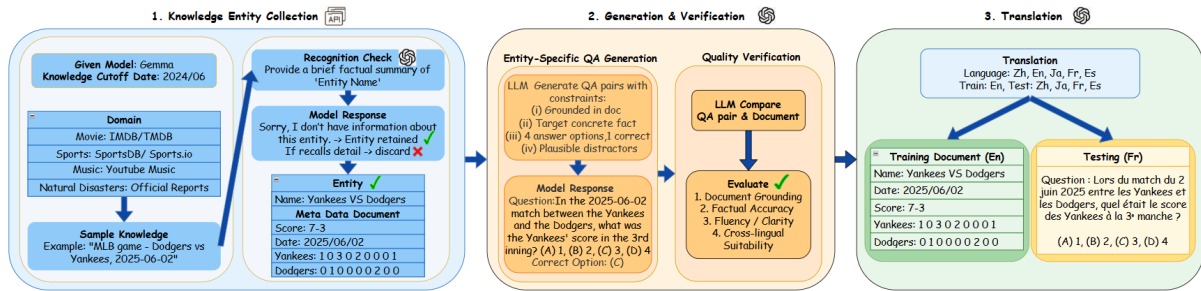


Figure 2: **LiveCLKT Bench Pipeline.** The generation process consists of four stages: (1) collecting independent, time-sensitive knowledge entities; (2) generating document-grounded question–answer pairs; (3) verifying data quality using a verifier LLM; and (4) translating verified questions into multiple languages for evaluation.

place, ensuring that such information could not have been included in pretraining data.

Entities are first retrieved from each domain using specific databases such as IMDB/TMDB (movies), YouTube Music (music) and SportsDB/Sports.io (sports). To avoid contamination, we apply a strict temporal filter, where only entities appearing at least six months after the model’s pretraining cutoff are considered. This reduces the risk of early exposure through trailers, previews, press leaks, or other pre-release publicity. The time window is configurable to match different model knowledge cutoffs and experimental needs.

After sampling, we perform a conservative check against the target model’s *pretrained* checkpoint to ensure each entity is genuinely unseen. The model is prompted with a short, unambiguous request: "Provide a brief factual summary of '<ENTITY NAME>'." The response is then compared to the retrieved source document using an LLM judge. If the response includes concrete facts consistent with the document, the entity is flagged as recognized and discarded.

For all remaining entities, we retrieve a canonical source document (e.g., a baseball match report, IMDB film page, music release metadata) to serve as the factual grounding for downstream QA generation. This multi-step procedure ensures that the retained entities and their corresponding documents constitute genuinely unseen knowledge, establishing a robust and contamination-free foundation for benchmark construction. Retrieved documents examples are provided in Appendix E.

**Entity-Specific QA Generation.** For each retained entity, we generate multiple-choice questions that are factual and grounded in the corresponding source document using an LLM (GPT-

4o-mini). The generation prompts enforce strict constraints: (i) every question must be explicitly grounded in the document; (ii) each question must target a concrete fact rather than subjective interpretation; (iii) exactly four candidate answers must be provided, with only one correct option; and (iv) distractors must be plausible but incorrect.

Because these questions focus on facts that emerge only after the event occurs, they are highly unlikely to appear in any pretraining data. These constraints ensure that the generated QAs are verifiable, unambiguous, and sufficiently challenging, providing a reliable foundation for evaluating cross-lingual knowledge transfer. For instance, a valid example that satisfies the four criteria is:

#### Sports QA

In the sports game 'Los Angeles Dodgers vs Toronto Blue Jays' at 2025-08-09, what was the final score?  
(A.) 5 – 1 (B.) 3 – 2 (C.) 4 – 0 (D.) 6 – 5

Additional QA examples are provided in Appendix D.

**Quality Verification.** To ensure that generated questions meet the required standards, a verifier LLM checks each QA pair against the original document and explicitly validates the four criteria above. Any question that violates at least one criterion is discarded. This guarantees that every retained test instance is document-grounded, factually accurate, and suitable for assessing transferability. Additionally, we validate the quality of verifier-approved questions by human inspection, confirming a precision of 95%. Details of verification procedure are provided in Appendix B.

**Translation.** Finally, the verified QAs and their associated documents are translated into the evaluation languages specified by the user. The resulting benchmark consists of a **train set**, containing source-language documents for knowledge injection, and a **test set**, containing multilingual factual QAs that probe cross-lingual transferability. We also conduct human evaluation to assess translation quality, as detailed in Appendix C.

**Language Coverage.** Since our framework relies on state-of-the-art LLMs (e.g., GPT-4o-mini) for translation, language coverage is bounded by the underlying model’s capabilities. In principle, evaluation data can be generated for any language with sufficient translation quality, enabling broader and more comprehensive cross-lingual evaluation.

Through these stages, LiveCLKTBench produces benchmarks that are leakage-free, configurable, and continuously extensible, ensuring evaluations remain independent of pretraining exposure. Detailed prompts for QA generation, verification, and translation are provided in Appendix F.

### 3.2 Evaluation Protocol

Unlike existing multilingual benchmarks that probe pretrained models in a zero-shot setting, LiveCLKTBench evaluates cross-lingual transferability *after knowledge injection*. The protocol focuses on two crucial steps, ensuring that results reflect true transfer: (1) Knowledge Injection: Inject new knowledge by post-training the model on source-language documents. (2) Transfer Evaluation: Evaluate the post-trained model on test QA pairs in the target languages. This setup guarantees that any correct prediction in a target language arises solely from successful cross-lingual knowledge transfer.

LiveCLKTBench supports two primary scenarios that allow researchers to analyze transferability from different perspectives:

1. **Model Comparison:** Compare multiple models using the same post-training strategy to determine which exhibits stronger cross-lingual transfer capabilities.
2. **Strategy Comparison:** For a single model, test different post-training methods (e.g., continual pretraining (Ke et al., 2023), supervised fine-tuning (Mecklenburg et al., 2024), knowledge editing (Wu et al., 2024)) to identify approaches that yield better transferability.

By supporting these scenarios, LiveCLKTBench enables both model-level and method-level analyses of cross-lingual knowledge transfer.

### 3.3 Evaluation Metrics

To quantify cross-lingual transferability, we follow the metric definitions introduced by Goldman et al. (2025). Let  $\mathcal{L}$  denote the set of evaluation languages. For each ordered pair  $(L_{\text{train}}, L_{\text{test}})$  with  $L_{\text{train}} \neq L_{\text{test}}$ , model predictions can be categorized using the contingency matrix in Table 1.

	$L_{\text{test}}$ Correct	$L_{\text{test}}$ Wrong
$L_{\text{train}}$ Correct	A	B
$L_{\text{train}}$ Wrong	C	D

Table 1: Contingency matrix for  $(L_{\text{train}}, L_{\text{test}})$ .  $A$ : successful transfer,  $B$ : failed transfer,  $C$ : cross-lingual inconsistency,  $D$ : complete failure.

This matrix distinguishes four possible outcomes:  $A$  — correct in both source and target languages (*successful transfer*);  $B$  — correct in source but wrong in target (*failed transfer*);  $C$  — wrong in source but correct in target (*cross-lingual inconsistency*); and  $D$  — wrong in both (*fail to learn*). From these outcomes we derive two complementary metrics.

**Overall Success.** This metric measures how often the model answers questions correctly in both the source and target languages. It reflects the combined ability to acquire knowledge during source-language post-training and to consistently reproduce it across languages:

$$\text{Overall}(L_{\text{train}}, L_{\text{test}}) = \frac{A}{A + B + C + D}. \quad (1)$$

**Transfer Success.** This metric conditions on cases where the source language answer is correct and evaluates the proportion that are also answered correctly in the target language. It directly quantifies the reliability of transferring acquired knowledge:

$$\text{Transfer}(L_{\text{train}}, L_{\text{test}}) = \frac{A}{A + B}. \quad (2)$$

Together, these metrics separate two dimensions of cross-lingual performance: **Overall Success** captures how frequently knowledge is jointly expressed across languages, while **Transfer Success** isolates the likelihood of successful transfer once the source knowledge has been learned.

Model	Domain	Overall (Score $\pm$ Std)	Transfer (Score $\pm$ Std)	Average over Domain (Overall / Transfer)
Gemma-2-9b	music	0.494 $\pm$ 0.110	0.807 $\pm$ 0.112	<b>0.414</b> / 0.735
	movie	0.483 $\pm$ 0.098	0.778 $\pm$ 0.142	
	sports	0.265 $\pm$ 0.111	0.620 $\pm$ 0.231	
Qwen2.5-7B	music	0.467 $\pm$ 0.103	0.808 $\pm$ 0.123	0.387 / <b>0.747</b>
	movie	0.452 $\pm$ 0.083	0.794 $\pm$ 0.132	
	sports	0.243 $\pm$ 0.097	0.637 $\pm$ 0.216	
Ministral-8B	music	0.357 $\pm$ 0.107	0.744 $\pm$ 0.149	0.304 / 0.669
	movie	0.360 $\pm$ 0.097	0.711 $\pm$ 0.194	
	sports	0.195 $\pm$ 0.102	0.551 $\pm$ 0.258	
Mistral-Nemo	music	0.344 $\pm$ 0.102	0.705 $\pm$ 0.168	0.289 / 0.641
	movie	0.326 $\pm$ 0.095	0.682 $\pm$ 0.190	
	sports	0.198 $\pm$ 0.100	0.536 $\pm$ 0.259	
Llama-3.1-8B	music	0.383 $\pm$ 0.089	0.807 $\pm$ 0.111	0.284 / 0.653
	movie	0.315 $\pm$ 0.124	0.686 $\pm$ 0.219	
	sports	0.155 $\pm$ 0.108	0.466 $\pm$ 0.301	
OLMo-2-7B	music	0.256 $\pm$ 0.119	0.623 $\pm$ 0.228	0.221 / 0.567
	movie	0.241 $\pm$ 0.107	0.603 $\pm$ 0.236	
	sports	0.165 $\pm$ 0.108	0.475 $\pm$ 0.302	

Table 2: **Cross-lingual Knowledge Transfer Performance.** Each Overall and Transfer score represents the mean  $\pm$  standard deviation across all ( $L_{\text{train}}, L_{\text{test}}$ ) language pairs for the given model and domain. The last column shows the average across all domains for each model.

#### 4 Demonstrating LiveCLKTBench: A Case Study

To illustrate the utility of LiveCLKTBench, we conduct a case study demonstrating how it can be applied in practice and its ability to reveal cross-lingual transferability. For simplicity, this demonstration focuses on one of the evaluation scenarios introduced in Section 3.2: *Model Comparison*.

**Setup.** We evaluate several open-source LLMs using a single knowledge injection strategy: *Continual Pre-Training* (Ke et al., 2023), where models continue pretraining directly on the target documents using the standard causal language modeling objective (next-token prediction) applied to unmodified text. Although this case study highlights *Model Comparison*, the framework is agnostic to the post-training method and also supports *Strategy Comparison*, enabling evaluation of different knowledge injection techniques for a single model.

**Benchmark Configuration.** We construct a benchmark covering the period 2025-01-01 to

2025-08-31<sup>3</sup> across five languages: English (en), Japanese (ja), Mandarin (zh), French (fr), and Spanish (es). Each instance consists of training documents in the source language paired with multilingual factual QA test sets, following the generation pipeline described in Section 3.1. Statistics of the benchmark are shown in Appendix A.

**Models.** We consider instruction-tuned, open-source models in the 7–9B parameter range, including GEMMA<sup>4</sup>, MISTRAL<sup>5,6</sup>, QWEN<sup>7</sup>, LLAMA<sup>8</sup>, and OLMo<sup>9</sup>. These model families differ in pretraining data, tokenizer design, and alignment strategies, providing a diverse perspective on cross-

<sup>3</sup>The most recent model among those evaluated, GEMMA, has a pretraining knowledge cutoff of 2024-06.

<sup>4</sup><https://huggingface.co/google/gemma-2-9b-it>

<sup>5</sup><https://huggingface.co/mistralai/Ministral-8B-Instruct-2410>

<sup>6</sup><https://huggingface.co/mistralai/Mistral-Nemo-Instruct-2407>

<sup>7</sup><https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

<sup>8</sup><https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

<sup>9</sup><https://huggingface.co/allenai/OLMo-2-1124-7B-Instruct>

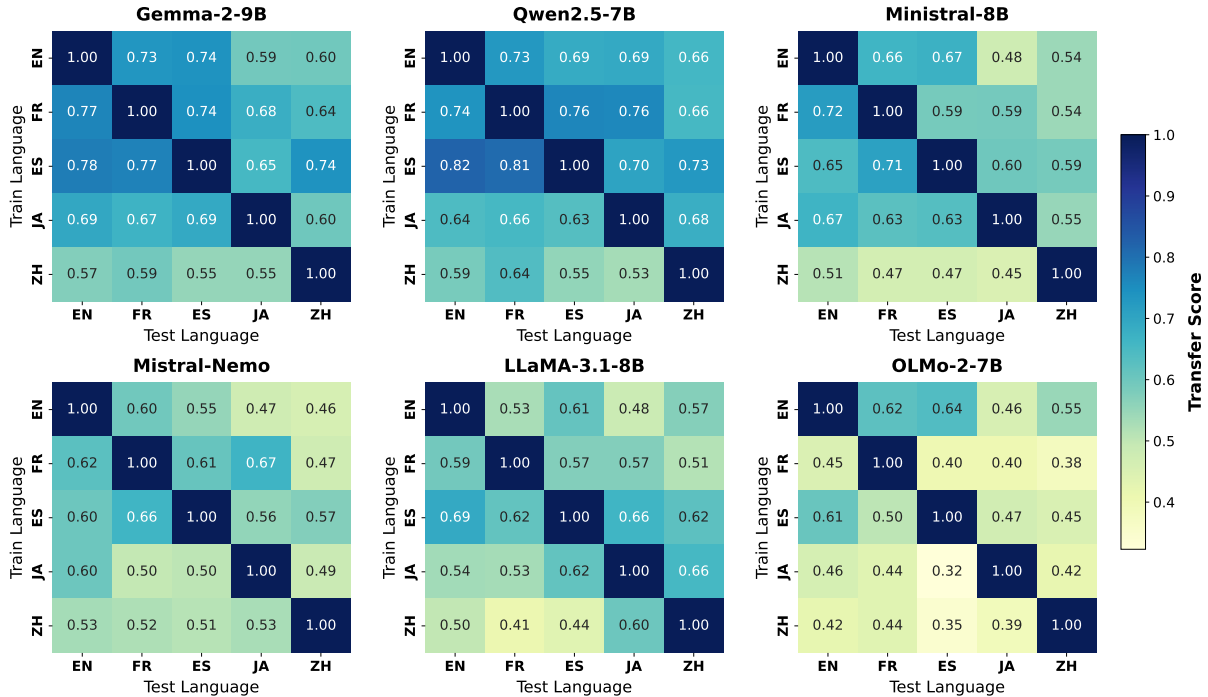


Figure 3: **Language-level Transferability.** Heatmaps show Transfer Scores for each  $(L_{\text{train}}, L_{\text{test}})$  pair across models, sorted by average Overall score. Darker colors indicate stronger transferability.

lingual transfer.

**Training and Inference.** All models are post-trained with a lightweight LoRA configuration (rank 16,  $\alpha = 32$ , learning rate  $5e^{-4}$ , dropout 0.1, batch size 1, 5 epochs). Intermediate checkpoints are validated, and the best-performing one is selected for final testing. At inference, we fix the decoding temperature to 0 to eliminate randomness and directly reflect the effect of training.

## 5 Case Study Results

Building on the setup described in Section 4, we now present the results of our LiveCLKTBench case study. The aim here is not to produce a leaderboard-style evaluation, but to demonstrate the kinds of insights LiveCLKTBench can provide about cross-lingual transfer.

We first report overall performance across models and domains, then analyze transferability patterns across languages, and examine the effect of model size. Together, these results illustrate the diagnostic power of LiveCLKTBench and show that cross-lingual knowledge transfer remains a significant challenge for current LLMs.

### 5.1 Overall Performance of Models

Table 2 presents the *Overall* and *Transfer* scores for each model across the three domains—music,

movies, and sports—as well as their domain averages. Across models, the average Overall Score ranges from 0.221 (OLMO-2-7B) to 0.414 (GEMMA-2-9B), while the Transfer Score ranges from 0.567 (OLMO-2-7B) to 0.747 (QWEN2.5-7B). GEMMA-2-9B achieves the highest average Overall Score, indicating stronger in-language performance, whereas QWEN2.5-7B leads in Transfer Score, suggesting better cross-lingual generalization. Other models such as MINISTRAL-8B, MISTRAL-NEMO, and LLAMA-3.1-8B perform moderately across both metrics.

Across domains, most models perform best on *music*, followed closely by *movies*, while *sports* consistently yields the weakest performance. When comparing models within each domain, their relative rankings remain largely consistent: GEMMA-2-9B and QWEN2.5-7B generally outperform others, while OLMO-2-7B trails behind. This stability of ranking suggests that model-level differences are systematic rather than domain-dependent, indicating comparable generalization patterns across diverse content types.

Overall, these results reveal substantial performance gaps across models, indicating that current LLMs still face challenges in reliably transferring knowledge across languages.

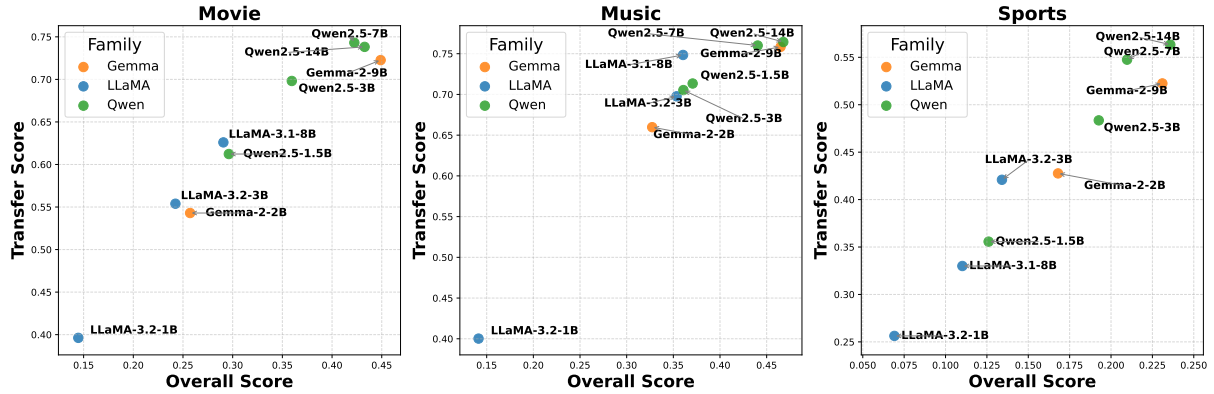


Figure 4: **Effect of Model Size.** Overall and Transfer Scores across model families of different parameter size, shown separately by domain.

## 5.2 Variation Across Languages

Figure 3 presents transfer performance across all source–target language pairs, with models sorted by their Overall score. Across all models, transferability is consistently weaker for Japanese (ja) and Mandarin (zh) compared to Indo-European languages such as English (en), Spanish (es), and French (fr). However, the degree of degradation varies: stronger models show moderate declines, while weaker ones suffer sharp drops in transfer scores for these distant languages. This pattern suggests that cross-lingual generalization to typologically different languages remains the key bottleneck, and that model robustness plays an important role in mitigating such gaps.

Another recurring observation is directional asymmetry. For instance, transferring knowledge from English to Japanese often yields lower performance than from Japanese to English. Such asymmetries imply that cross-lingual transfer is not bidirectionally balanced and that some languages serve as more effective sources or recipients of knowledge. Notably, QWEN2.5-7B shows relatively stronger transfer into Mandarin compared to other models, likely reflecting its heavier exposure to Chinese text during pretraining.

Overall, these findings demonstrate that linguistic proximity plays a major role in transfer effectiveness: models achieve higher reliability within language families but face persistent challenges when transferring across typologically distant ones such as Indo-European and East Asian languages.

## 5.3 Effect of Model Size

Figure 4 illustrates how model size affects transferability across three domains. Consistent with

general scaling trends, larger models consistently outperform smaller ones across all families, confirming that increased capacity generally improves cross-lingual generalization.

However, the improvements are not strictly linear—gains tend to diminish as models grow larger, suggesting a saturation effect at higher scales. The magnitude of improvement also varies across domains: scaling yields clear benefits in *movies* and *sports*, where larger models achieve noticeably higher transferability, whereas performance on *music* is already high and remains clustered across parameter sizes.

Overall, scaling improves general cross-lingual alignment but with diminishing returns. Beyond a certain capacity, improvements appear constrained more by domain complexity than by model size, suggesting that future gains may require greater efforts toward domain-targeted multilingual adaptation rather than simple scaling.

## 6 Conclusion

In this work, we presented **LiveCLKTBench**, an automated pipeline for building realistic, contamination-free benchmarks for cross-lingual knowledge transfer. Our approach prevents leakage by selecting independent, time-sensitive knowledge entities, filtering them by temporal occurrence, and discarding any already recognized by pretrained models. The pipeline is fully automated and configurable, enabling frequent updates and flexible customization with minimal human effort. Together, these features make LiveCLKTBench a scalable and sustainable framework for evaluating genuine cross-lingual knowledge transfer in LLMs.

## Limitations

While LiveCLKTBench provides a realistic and contamination-free benchmark for evaluating cross-lingual knowledge transfer, several limitations remain that could be addressed in future work:

**Task Variety.** Currently, the benchmark focuses on multiple-choice QA due to its ease of evaluation and cost efficiency. Extending the framework to include other task types, such as open-ended question answering, would provide a more comprehensive assessment.

**Domain Coverage.** The benchmark currently spans three domains (movies, music, and sports). Expanding to additional domains and sources would increase knowledge diversity and improve the breadth of evaluation.

**Language Choices.** In our case study, we evaluate cross-lingual transfer on five languages (English, Spanish, French, Japanese, and Mandarin) based on their coverage across evaluated models, common use in prior multilingual benchmarks, and computational cost considerations. Extending to more languages, especially low-resource ones, is left for future work.

## Ethical considerations

LiveCLKTBench uses publicly available sources such as IMDB/TMDB (movies), YouTube Music (music), and SportsDB (sports) from trusted agencies. While these sources may occasionally contain inaccuracies due to human error or reporting delays, they provide widely recognized and verifiable records of real-world events. All data are derived from post-release or officially published information, ensuring that no private or sensitive material is included. Our benchmark focuses on factual, entity-based knowledge (e.g., movie releases, sports scores), which further minimizes the risk of harmful content.

## Use of Ai Assistants

In this work, we leveraged large language models (LLMs) to assist in two ways. First, LLMs were employed as part of the automatic data generation pipeline for tasks such as question formulation, translation, and quality verification. Second, an AI assistant (OpenAI GPT-5) was used for minor writing support, including grammar correction and improving manuscript clarity. All AI-assisted content was carefully reviewed and verified by the

authors to ensure factual accuracy and alignment with the authors' original intent.

## Acknowledgment

This material is based upon work supported by National Science and Technology Council, ROC under grant number 114-2221-E-002-134-MY3 and 113-2628-E-001-003-MY4, NTU AI Center of Research Excellence within Taiwan Centers of Excellence in Artificial Intelligence, and by National Taiwan University and Academia Sinica Innovative Joint Program, under grant AS-NTU-114-06.

## References

- Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, Kalika Bali, and Sunayana Sitaram. 2023. [Mega: Multilingual evaluation of generative ai](#). *Preprint*, arXiv:2303.12528.
- Sanchit Ahuja, Varun Gumma, and Sunayana Sitaram. 2024. [Contamination report for multilingual benchmarks](#). *Preprint*, arXiv:2410.16186.
- Akari Asai, Jungo Kasai, Jonathan H. Clark, Kenton Lee, Eunsol Choi, and Hannaneh Hajishirzi. 2021. [Xor qa: Cross-lingual open-retrieval question answering](#). *Preprint*, arXiv:2010.11856.
- Akari Asai, Sneha Kudugunta, Xinyan Velocity Yu, Terra Blevins, Hila Gonen, Machel Reid, Yulia Tsvetkov, Sebastian Ruder, and Hannaneh Hajishirzi. 2023. [Buffet: Benchmarking large language models for few-shot cross-lingual transfer](#). *Preprint*, arXiv:2305.14857.
- Linzhen Chai, Jian Yang, Tao Sun, Hongcheng Guo, Jiaheng Liu, Bing Wang, Xiannian Liang, Jiaqi Bai, Tongliang Li, Qiyao Peng, and Zhoujun Li. 2024. [xcot: Cross-lingual instruction tuning for cross-lingual chain-of-thought reasoning](#). *Preprint*, arXiv:2401.07037.
- Jie Chen, Yupeng Zhang, Bingning Wang, Xin Zhao, Ji-Rong Wen, and Weipeng Chen. 2024. [Unveiling the flaws: Exploring imperfections in synthetic data and mitigation strategies for large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14855–14865, Miami, Florida, USA. Association for Computational Linguistics.
- Lynn Chua, Badih Ghazi, Yangsibo Huang, Pritish Kamath, Ravi Kumar, Pasin Manurangsi, Amer Sinha, Chulin Xie, and Chiyuan Zhang. 2025. [Crosslingual capabilities and knowledge barriers in multilingual large language models](#). *Preprint*, arXiv:2406.16135.

- Omer Goldman, Uri Shaham, Dan Malkin, Sivan Eiger, Avinatan Hassidim, Yossi Matias, Joshua Maynez, Adi Mayrav Gilady, Jason Riesa, Shruti Rijhwani, and 1 others. 2025. Eclectic: a novel challenge set for evaluation of cross-lingual knowledge transfer. *arXiv preprint arXiv:2502.21228*.
- Michael A. Hedderich, David I. Adelani, Dawei Zhu, Jesujoba Alabi, Udia Markus, and Dietrich Klakow. 2020. **Transfer learning and distant supervision for multilingual transformer models: A study on African languages**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2580–2591, Online. Association for Computational Linguistics.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020a. **XTREME: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation**. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4411–4421. PMLR.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020b. **Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation**. In *International conference on machine learning*, pages 4411–4421. PMLR.
- Essa Jan, Moiz Ali, Muhammad Saram Hassan, Fareed Zaffar, and Yasir Zaki. 2025. **Data doping or true intelligence? evaluating the transferability of injected knowledge in llms**. *Preprint*, arXiv:2505.17140.
- Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. 2024. **Realtime qa: What’s the answer right now?** *Preprint*, arXiv:2207.13332.
- Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. 2023. **Continual pre-training of language models**. *Preprint*, arXiv:2302.03241.
- Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A. Rossi, and Thien Huu Nguyen. 2023. **Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback**. *Preprint*, arXiv:2307.16039.
- Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. **From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4483–4499, Online. Association for Computational Linguistics.
- Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, and 1 others. 2020. **Xglue: A new benchmark dataset for cross-lingual pre-training, understanding and generation**. *arXiv preprint arXiv:2004.01401*.
- Robert Litschko, Oliver Kraus, Verena Blaschke, and Barbara Plank. 2025. **Cross-dialect information retrieval: Information access in low-resource and high-variance languages**. *Preprint*, arXiv:2412.12806.
- Gaurav Maheshwari, Dmitry Ivanov, and Kevin El Hadad. 2024. **Efficacy of synthetic data as a benchmark**. *Preprint*, arXiv:2409.11968.
- Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C. Lipton, and J. Zico Kolter. 2024. **Tofu: A task of fictitious unlearning for llms**. *Preprint*, arXiv:2401.06121.
- Dan Malkin, Tomasz Limisiewicz, and Gabriel Stanovsky. 2022. **A balanced data approach for evaluating cross-lingual transfer: Mapping the linguistic blood bank**. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4903–4915, Seattle, United States. Association for Computational Linguistics.
- Nick Mecklenburg, Yiyu Lin, Xiaoxiao Li, Daniel Holstein, Leonardo Nunes, Sara Malvar, Bruno Silva, Ranveer Chandra, Vijay Aski, Pavan Kumar Reddy Yannam, Tolga Aktas, and Todd Hendry. 2024. **Injecting new knowledge into large language models via supervised fine-tuning**. *Preprint*, arXiv:2404.00213.
- Itai Mondshine, Tzuf Paz-Argaman, and Reut Tsarfaty. 2025. **Beyond english: The impact of prompt translation strategies across languages and tasks in multilingual llms**. *Preprint*, arXiv:2502.09331.
- Ercong Nie, Bo Shao, Zifeng Ding, Mingyang Wang, Helmut Schmid, and Hinrich Schütze. 2025. **Bmike-53: Investigating cross-lingual knowledge editing with in-context learning**. *Preprint*, arXiv:2406.17764.
- OpenAI. 2024. Gpt-4o system card. <https://openai.com/index/gpt-4o-system-card/>. Accessed: 2026-04-15.
- Sara Rajaei and Christof Monz. 2024. **Analyzing the evaluation of cross-lingual knowledge transfer in multilingual language models**. *Preprint*, arXiv:2402.02099.
- Uri Shaham, Jonathan Herzig, Roei Aharoni, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. **Multilingual instruction tuning with just a pinch of multilinguality**. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 2304–2317, Bangkok, Thailand. Association for Computational Linguistics.
- Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I. Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond

- Ng, Shayne Longpre, Wei-Yin Ko, Sebastian Ruder, Madeline Smith, Antoine Bosselut, Alice Oh, Andre F. T. Martins, Leshem Choshen, and 5 others. 2025. [Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation](#). *Preprint*, arXiv:2412.03304.
- Zihao Wei, Jingcheng Deng, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. 2025. [Mlake: Multilingual knowledge editing benchmark for large language models](#). *Preprint*, arXiv:2404.04990.
- Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Sreemanti Dey, Shubh-Agrawal, Sandeep Singh Sandha, Siddartha Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. 2025. [Livebench: A challenging, contamination-limited llm benchmark](#). *Preprint*, arXiv:2406.19314.
- Minghao Wu, Weixuan Wang, Sinuo Liu, Huifeng Yin, Xintong Wang, Yu Zhao, Chenyang Lyu, Longyue Wang, Weihua Luo, and Kaifu Zhang. 2025. [The bitter lesson learned from 2,000+ multilingual benchmarks](#). *Preprint*, arXiv:2504.15521.
- Xiaobao Wu, Liangming Pan, William Yang Wang, and Anh Tuan Luu. 2024. [Akw: Assessing knowledge editing in the wild](#). *Preprint*, arXiv:2402.18909.
- Feng Yao, Yufan Zhuang, Zihao Sun, Sunan Xu, Animesh Kumar, and Jingbo Shang. 2024. [Data contamination can cross language barriers](#). *Preprint*, arXiv:2406.13236.
- Wenxuan Zhang, Sharifah Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. 2023. [M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models](#). *Preprint*, arXiv:2306.05179.
- Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. [Agieval: A human-centric benchmark for evaluating foundation models](#). *Preprint*, arXiv:2304.06364.

## A Data Statistics

In LiveCLKTBench, the training set consists of raw text documents that contain information about each knowledge entity (e.g., a movie). Each training document is associated with multiple QA pairs that probe different factual aspects of the same entity (e.g., cast, storyline, release details). Since we evaluate cross-lingual transfer by post-training the model on the source-language documents, we split the corresponding QA pairs into validation (20%) and test subsets (80%). The validation subset is used to select the best checkpoint and to monitor learning progress during post-training. The test subset is held out and used only for final evaluation. Importantly, both validation and test questions are still grounded in the same knowledge document used for training and do not introduce any new information.

Domain	#Entities	#Fact QAs (pre-split)	#Validation	#Test
Movie	30	175	175	700
Music	30	125	125	500
Sports	20	95	90	380

Table 3: **Benchmark Statistics.** Number of entities, generated factual QA pairs before translation and splitting, and counts of the final translated validation and test instances.

## B Quality of LLM QA Generation

To assess the reliability of our benchmark data, we recruited three annotators (two undergraduate students and one graduate student) to manually verify the quality of sampled question–answer pairs. Specifically, we randomly sampled 20% of entities from each domain and evaluated all QA pairs associated with them. Each question was independently annotated by three annotators based on the following criteria, and was considered a failure if it violated any of them.

- F-1: The question is unrelated to the source document.
- F-2: The question is relevant but overly general and not specific to the target entity.
- F-3: The answer cannot be verified from the source document.
- F-4: The answer is verifiable from the source document but is incorrect.

For annotation, we required full agreement among annotators. Questions with disagreement were discarded from the analysis. On the remaining items (98.8%), the validated QA accuracy (Pass Rate) is 95.2%. The invalid items include 3.6% cases of F-3 (from the movie domain) and 1.2% cases of F-4 (from the sports domain). These results suggest that our automated pipeline produces reasonably high-quality, fact-grounded QA pairs.

## C Quality of LLM Translation

LiveCLKTBench aims to support continuous, contamination-controlled updates, making automatic translation essential for scalability. Consistent with prior multilingual benchmarks (Lai et al., 2023; Goldman et al., 2025; Asai et al., 2023; Singh et al., 2025; Zhang et al., 2023), we adopt LLM-based translation. Our pipeline utilizes GPT-4o-mini (OpenAI, 2024), a strong multilingual model comparable to those used in prior work, ensuring that translation quality is aligned with existing benchmarks.

To further validate translation quality, we conduct human evaluation on a randomly sampled 10% subset of QA pairs and training documents across four target languages: English, Japanese, French, and Mandarin. These languages are selected based on the availability of native or highly proficient annotators. Each sample is evaluated by annotators who are either native speakers or have advanced proficiency in the target language, using two 1–5 scales (higher is better):

- **Adequacy (semantic correctness):** How accurately the original meaning is preserved in the translation.
- **Fluency (linguistic naturalness):** How natural, grammatical, and fluent the translation is in the target language, independent of semantic accuracy.

Language Pair	Adequacy (Train Doc)	Adequacy (Test QA)	Fluency (Train Doc)	Fluency (Test QA)
EN → JA	5.00	4.67	3.88	3.86
EN → FR	4.88	4.92	4.50	4.72
EN → ZH	4.75	5.00	4.25	4.67

Table 4: Human evaluation results for translation quality across language pairs.

Table 4 shows consistently high adequacy scores across all language pairs, indicating that the translated content reliably preserves the original meaning. Fluency scores are comparatively lower, suggesting that there is room for improvement in linguistic naturalness. Nevertheless, the strong adequacy scores support the validity of the multilingual QA pairs, since the factual content required for evaluation is preserved even when surface-level naturalness varies. To further improve fluency in future iterations, we plan to adopt stricter translation constraints (e.g., discouraging overly literal phrasing), incorporate back-translation techniques, and explore lightweight human-in-the-loop refinement for low-fluency cases.

## D QA Examples

### Music

In the music video 'Alex Warren - Ordinary (Official Video)', what recurring theme is mentioned in the lyrics?  
 (A.) Love and longing (B.) Self-empowerment and resilience (C.) Escaping from fame and pressure  
 (D.) Chasing dreams and freedom

### Movie

In the movie 'KPop Demon Hunters', who are the main characters that use their secret identities to fight supernatural threats?  
 (A.) Zoey, Mira, Ahn (B.) Arden, May, Ji-young (C.) Rumi, Mira, Zoey (D.) Ahn, Yunjin, Rumi

## E Entity Document Examples

### Movie Document Template

- Movie Title: {title}
- Movie Cast: {casts}
- Movie Summary: {summary}
- Movie Synopsis: {synopsis}

### Music Document Template

- Music Video Title: {title}
- Music Release Date: {date}
- Music Video Description: {description}

## Sports Document Template

Sports: {sports}  
League: {league}

Match: {home\_team} vs {away\_team}  
Date: {date}  
Score: {home\_score} - {away\_score}  
Venue: {venue}

Innings Breakdown:

{home\_team}: {home\_innings} → Hits: {home\_hits}, Errors: {home\_errors}  
{away\_team}: {away\_innings} → Hits: {away\_hits}, Errors: {away\_errors}

## F Pipeline Prompt Examples

### QA GENERATION PROMPT

You are generating high-quality multiple-choice QA pairs in {lang}, strictly grounded in the given movie information.

You will be provided with:

- Movie Title
- Movie Casts
- Movie Summary
- Movie Synopsis

Task:

- Generate natural, audience-friendly questions that viewers might realistically ask.
- All questions must be written fully in {lang}, including the leading phrase (“In the movie: '<title>', ...”).
- Each QA pair must be based ONLY on facts explicitly present in the input. Do not add, assume, or hallucinate.
- Use diverse aspects (casts, summary, synopsis content).

Each QA pair must include:

- Question in {lang}, beginning with “In the movie: '<title>', ...” (do not translate title)
- Options:
- Provide four options labeled A, B, C, D.
- Exactly one option is correct.
- Place the correct option randomly among A-D (do not always use the same position).
- Distractors must be plausible but wrong (no random, absurd, or unrelated answers).
- Correct Option: Output the letter (A, B, C, or D) of the correct answer.

---

Inputs:

```
{meta_data}
```

---

Output Format (JSON):

```
{
  "QA": [
    {
      "question": "<string in {lang}>",
      "options": {
        "A": "<option A in {lang}>",
        "B": "<option B in {lang}>",
        "C": "<option C in {lang}>",
        "D": "<option D in {lang}>"
      },
      "correct_option": "<A | B | C | D>"
    },
    ...
  ]
}
```

---

Guidelines:

- Write everything (questions and options) only in {lang}.
- Keep all proper names (people, places, entities) unchanged.
- Ensure every correct answer can be directly verified in the input metadata.
- Distractors must be reasonable, related, and plausible.

### QA VERIFIER PROMPT

You are verifying if QA pairs are grounded in the provided metadata.

Check:

- Is the correct option explicitly supported by the metadata?
- If yes, return SUPPORTED and the supporting sentence(s).
- If not, return UNSUPPORTED.

Output Format:

```
{
  "Decision": "<SUPPORTED or UNSUPPORTED>",
  "SourceSentence": "<sentence(s) from metadata or empty>"
}
```

### QA TRANSLATION PROMPT

Translate the QA JSON into target language {lang}.

Rules:

- Translate only values of "question" and "options".

- Do NOT translate keys ("QA", "question", "options", "A"..."D", "correct\_option").
- Keep "correct\_option" unchanged.
- Preserve JSON structure.

### DOCUMENT TRANSLATION PROMPT

Translate the following movie document into {lang}.

- Do NOT translate the field names (e.g., "Movie Cast", "Movie Summary", "Movie Synopsis").
- Translate only the values after the colon.
- If the text is already in {lang}, return it unchanged.
- Return only the JSON object, no extra text.

Document:

- Movie Cast: {casts}
- Movie Summary: {summary}
- Movie Synopsis: {synopsis}

Output Format:

```
{
  "translation": {
    "Cast": "<translated cast>",
    "Summary": "<translated summary>",
    "Synopsis": "<translated synopsis>"
  }
}
```