

# Cross-Lingual LLM-Judge Transfer via Evaluation Decomposition

Anonymous ACL submission

## Abstract

As large language models are increasingly deployed across diverse real-world applications, extending automated evaluation beyond English has become a critical challenge. Existing evaluation approaches are predominantly English-focused, and adapting them to other languages is hindered by the scarcity and cost of human-annotated judgments in most languages. We introduce a decomposition-based evaluation framework built around a Universal Criteria Set (UCS). UCS consists of a shared, language-agnostic set of evaluation dimensions, producing an interpretable intermediate representation that supports cross-lingual transfer with minimal supervision. Experiments on multiple faithfulness tasks across languages and model backbones demonstrate consistent improvements over strong baselines without requiring target-language annotations.

## 1 Introduction

The rapid growth of general-purpose AI systems has led to a dramatic increase in machine-generated text across applications such as summarization (Pu et al., 2023), question answering (Yue, 2025), content moderation (Kolla et al., 2024), and search (Spatharioti et al., 2023). Assessing the quality and correctness of these outputs at scale remains a central challenge (Wu et al., 2025; Ohde et al., 2025). Although human annotation is the gold standard, it is expensive, time-consuming, and difficult to scale to the volume and frequency required by modern development cycles (Ohde et al., 2025; Gu et al.).

Traditional automatic evaluation metrics for text generation, such as ROUGE (Lin, 2004) and BLEU (Papineni et al., 2002), have long served as proxies for output quality. However, they rely primarily on surface-level lexical overlap and fail to capture semantic meaning, reasoning, and contextual appropriateness (Sulem et al., 2018). LLM-based judges can be a more flexible alternative by

leveraging broad pretrained knowledge and contextual reasoning. An LLM can be prompted to assess the quality of output using task-specific instructions or rubric-style criteria (Li et al., 2025; Zheng et al., 2023; Chen et al., 2024; Chiang et al., 2024), leading to their growing adoption for summarization, factuality assessment, preference modeling, and model comparison.

The development and benchmarking of LLM judges have mainly focused on English, despite only about 15% of the global population speaking English (Pombal et al., 2025). As AI systems are deployed worldwide, reliable evaluation across diverse languages becomes a critical bottleneck (Wang et al., 2024; Fu and Liu, 2025). Existing multilingual judge approaches face two primary challenges. First, many rely on language-specific fine-tuning, which requires human-annotated data that is scarce for most languages, and often underperform strong English-based baselines (Doddapaneni et al., 2025). Second, English-trained judges do not consistently transfer to typologically distant or low-resource languages (Fu and Liu, 2025). Together, these limitations hinder scalable and reliable multilingual evaluation.

These limitations become particularly apparent in the common scenario of language expansion, where systems developed and evaluated in English are later deployed to additional languages. While model capabilities can often be extended through translation or multilingual prompting, evaluation typically requires collecting new annotation data or adapting evaluation frameworks for each language. As a result, evaluation becomes a bottleneck in scaling LLM systems to new linguistic settings, highlighting the need for approaches that can transfer evaluation behavior from English to other languages with minimal additional supervision.

In this paper, we introduce an interpretable decomposition-based framework for LLM judges that enables sample-efficient cross-lingual transfer.

Our approach decomposes judgment into a shared set of language-agnostic evaluation criteria, each expressed as a targeted question about a specific dimension of quality. The resulting criterion-level responses define a structured intermediate representation that maps judgments into a predefined, human-interpretable evaluation space, in the spirit of concept-based models that ground intermediate representations in explicit semantic dimensions (Koh et al., 2020; Espinosa Zarlenga et al., 2022; Pota et al., 2023; Sheth and Ebrahimi Kahou, 2023; Sun et al.), while also capturing reasoning patterns that transfer across languages. A lightweight transfer module trained only on English-labeled data can then be applied directly to new languages, enabling cross-lingual judge transfer without target-language supervision.

We evaluate our framework across multilingual benchmarks and multiple LLM backbones. Our contributions are: (1) we show that LLM-judge reasoning can be factored through a shared, language-agnostic criteria space, enabling consistent judgments across heterogeneous linguistic inputs; (2) we introduce an interpretable intermediate representation derived from criteria-level responses; (3) we demonstrate that this representation supports effective cross-lingual transfer using a lightweight module trained in English and applied to other languages without requiring target-language labels.

## 2 Related Works

**LLM Judges** Large language models are increasingly used as automated evaluators for summarization, dialogue, factuality, and preference modeling. Early work (Zheng et al., 2023; Liusie et al., 2024) demonstrated that LLMs can approximate human preferences using carefully designed prompts. Subsequent studies explored rubric-based scoring (Song et al., 2024), pairwise comparison (Liusie et al., 2024), and direct answer classification using models like GPT-4. However, these methods typically rely on single-prompt formulations that are brittle to prompt phrasing and struggle to generalize (Thakur et al., 2025).

Another line of work improves evaluation by incorporating explicit reasoning, such as chain-of-thought prompting (Zheng et al., 2023), justification-first scoring (Trivedi et al.), or predefined criteria prompting (Wei et al.; Lee et al., 2025).

More recent approaches employ multiple LLMs

interacting through debate (Feng et al., 2025; Chan et al.), critique (Kim et al., 2024), or competitive assessment. While effective in certain settings, these systems require careful human input for selecting few-shot or seed prompts (Feng et al., 2025; Lee et al., 2025) and substantial role engineering (Alfano et al., 2025), making them costly and difficult to deploy reliably in production environments.

Checklist-style evaluation has recently gained attention as an interpretable alternative to monolithic rubric prompts. Previous work has shown that explicitly decomposing an evaluation task into smaller criteria can improve transparency. Wei et al. generate binary checklist items using a stronger model and apply them to smaller evaluators, effectively distilling high-level judgements into simple, verifiable units. (Lee et al., 2025) prompt models to generate their own criteria and use these checklists for iterative self-improvement, allowing the evaluator to refine its reasoning over multiple rounds.

**Multilingual LLM-Based Judges** Most LLM-judge design and evaluation has focused on English. Chang et al. (2025) studies the impact of resource availability on multilingual evaluators, while Fu and Liu (2025) analyzes the reliability of multilingual judges. Recent approaches train multilingual evaluators through large-scale pretraining (Pombal et al., 2025) or language-specific fine-tuning (Dodapaneni et al., 2025). To reduce reliance on human annotations, Alfano et al. (2025) generates synthetic multilingual supervision by translating English data and constructing corrupted summaries for training, and further examines how training language affects multilingual evaluation behavior; notably, their strongest results are obtained by fine-tuning on English data only.

However, all of these approaches require either task-specific data construction, full model fine-tuning, or language-specific adaptation. In contrast, our method produces a language-agnostic interpretable representation from a shared criterion set, enabling cross-lingual transfer via a lightweight module trained on only a few labeled English examples with no target-language supervision, no data engineering, and no full-model fine-tuning.

## 3 Methodology

We propose a criteria-based evaluation framework that decomposes an LLM-based judge’s decision into a structured set of sub-criteria. Rather than relying on a single prompt or free-form reasoning,

our judge produces answers to a set of targeted evaluation questions, which we refer to as criteria, and these answers are aggregated into an intermediate judge representation used for cross-lingual transfer.

**Problem Formulation.** Consider an evaluation task in which an input sample  $x$  must be assigned a binary label  $y \in \{0, 1\}$ . Let  $\mathcal{D} = \{(x_i, y_i)_{i=1}^N\}$  denote a labeled dataset, where  $y_i$  is the ground-truth judgment for sample  $x_i$ . Each sample  $x$  can consist of any structured input relevant to the task, such as a source–output pair, an instruction–response pair, or any text bundle under evaluation. A standard LLM-based judge is obtained by prompting the model  $\mathcal{M}$  with an evaluation instruction to produce a predicted label  $\hat{y}_i$ :

$$\hat{y}_i = \mathcal{M}(x_i). \quad (1)$$

Standard prompting produces  $\hat{y}_i$  directly, either with chain-of-thought reasoning or via multi-agent orchestration.

In this work, we replace this single-prompt-based judgment with a decomposition-based framework in which the LLM first responds to a set of evaluation criteria. These responses form an intermediate judge representation that is used for transfer across languages.

**Framework Overview.** Our methodology consists of three stages. **Stage 1: Criteria Set Generation** constructs a set of language-agnostic evaluation criteria from the task specification. **Stage 2: Criteria Set Evaluation** applies these criteria to each input sample, producing a structured intermediate representation from the LLM’s criterion-level responses. **Stage 3: Cross-Lingual transfer** trains a lightweight transfer module on labeled English data to align the judge representation, enabling transfer to other languages. An overview of the framework is shown in [Figure 1](#).

### 3.1 Criteria Set Generation

Our approach first constructs a set of evaluation criteria that define the dimensions along which the LLM-based judge evaluates each sample. Formally, we generate a *Universal Criteria Set* (UCS), a collection of evaluation questions

$$Q = \{q_1, \dots, q_k\} \quad (2)$$

where each  $q_k$  specifies a distinct judgment dimension. We generate criteria in English, motivated by

evidence that LLMs exhibit more consistent reasoning and intermediate representations in English than in other languages (Schut et al.; Shi et al.; Huang et al., 2024). This choice simplifies the generation of criteria and improves stability for downstream cross-lingual transfer.

The UCS defines a reusable set of evaluation dimensions applicable to all samples for a given task and is shared across all languages. Given a task description  $t$ , the LLM generates a set of evaluation concepts  $C = \{c_1, \dots, c_m\}$  representing generic evaluation dimensions such as accuracy, consistency, attribution, or hallucination. The UCS is task-specific, generated from the task description but *language-universal*: the same criteria set is applied to all languages for a given task.

$$C = \mathcal{M}(t) \quad (3)$$

See [Appendix E](#) for the prompts used in this stage.

**Question Generation** For each concept  $c_j \in C$ , the LLM generates one or more evaluation questions to comprehensively cover its sub-dimensions. Let  $Q_j$  denote the subset of criteria generated from the concept  $c_j$ , such that :

$$Q_j = \mathcal{M}(c_j) \quad (4)$$

The final Universal Criteria Set is obtained by aggregating the criteria generated for all concepts:

$$Q = \bigcup_{j=1}^m Q_j. \quad (5)$$

Thus,  $Q = \{q_1, \dots, q_k\}$  represents the union of all concept-derived questions. The same criteria set  $Q$  is applied to every sample, providing stable language-agnostic evaluation features that support cross-lingual transfer.

### 3.2 Criteria Set Evaluation

In the second stage of our framework, given the UCS  $Q$  and an input sample  $x$ , the LLM is prompted to assess the extent to which  $x$  satisfies each criterion  $q_k \in Q$ , returning a numerical score:

$$\mathcal{M}(q_k, x) = z_k(x) \in [1, \dots, 10], \quad (6)$$

We consider Likert rating from 1-10 following (Lee et al., 2025). Other rating ranges is left to future explorations. Collecting these outputs yields the *criteria-response vector*:

$$z(x) = [z_1(x), \dots, z_k(x)], \quad (7)$$

which serves as a structured and interpretable intermediate judge representation of the sample.

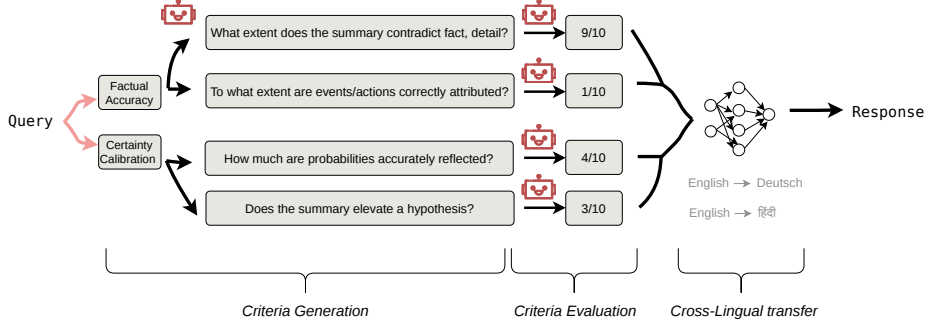


Figure 1: UCS framework. The LLM evaluates an input using criteria questions derived from higher-level concepts. Criterion scores are aggregated into a latent representation, which a transfer module trained on English-labeled data maps to the final judgment and applies across languages.

### 3.3 Cross-Lingual Transfer

The final stage learns a mapping from the criteria-response representation  $z(x)$  to a calibrated judgment that generalizes across languages. The key insight is that while raw LLM outputs may vary across languages, the structure of evaluation as expressed through the shared criteria remains stable. The criteria-response vector  $z(x)$  therefore defines a language-agnostic intermediate representation space in which a lightweight predictor trained on one language can be applied directly to another.

**Concept-Level Aggregation.** Rather than operating on the full criteria-response vector  $z(x)$ , we perform concept-level transfer by aggregating criterion-level scores within each concept. This reduces the dimensionality of the representation, decreases sample complexity, and improves transfer robustness (see Appendix C for details). Concretely, for each concept  $c_j$ , we average the scores of its associated criteria:

$$s_j(x) = \frac{1}{|Q_j|} \sum_{q_i \in Q_j} z_i(x). \quad (8)$$

The resulting concept-level representation is

$$s(x) = [s_1(x), \dots, s_m(x)]. \quad (9)$$

**Transfer.** Given labeled data in a source language  $\ell$ , we train a lightweight predictor (e.g. a neural network) on the concept-level representation:

$$\hat{y} = f_\theta(s(x)), \quad (10)$$

where  $\theta$  is learned on the source-language dataset  $\mathcal{D}^\ell = \{(x_i^\ell, y_i^\ell)\}_{i=1}^N$ . At inference time, the trained predictor  $f_{\theta^*}$  is applied directly to samples in a target language  $\ell'$ :

$$\hat{y}_j^{\ell'} = f_{\theta^*}(s(x_j^{\ell'})). \quad (11)$$

This enables cross-lingual transfer without any target-language supervision, as the shared criteria space ensures that  $s(x)$  remains semantically consistent across languages.

## 4 Experimental Setup

**Models.** We evaluate our methods across a range of LLMs to ensure diversity in scale and architecture. We report the results for Qwen3-32B (Yang et al., 2025), Qwen3-235B-A22B (Yang et al., 2025), OSS-20B (Agarwal et al., 2025), and OSS-120B (Agarwal et al., 2025). All LLM-based judges are queried in a deterministic setting with temperature = 0 and top\_p = 1.0 to reduce randomness in evaluation outputs.

**Datasets.** Our primary experiments use two multilingual evaluation benchmarks. **MEMERAG** (Blandón et al., 2025) is a multilingual RAG faithfulness benchmark containing evidence–query–answer triples, each annotated with a binary faithfulness label. The dataset spans five languages: English (EN), French (FR), German (DE), Hindi (HI), and Spanish (ES). **mFACE** (Aharoni et al., 2023) is a multilingual summarization evaluation dataset consisting of news articles and human-written summaries collected from BBC regional editions. The task is to judge whether a summary is faithful to the source article. We use a representative subset of languages spanning high-resource, mid-resource, and low-resource languages: English (EN), Amharic (AM), Burmese (MY), French (FR), Swahili (SW), Thai (TH), Arabic (AR), Hindi (HI), and Spanish (ES).

**Evaluation Metrics.** We report *balanced accuracy* (BA) to account for class imbalance, as it

Method	MEMERAG					mFACE								
	DE	FR	ES	HI	Avg	AM	MY	FR	SW	TH	AR	HI	ES	Avg
Zero-shot	77.2	77.5	77.5	74.4	76.7	54.7	<b>70.2</b>	74.3	73.0	72.6	69.7	68.0	77.9	70.1
CoT	73.0	76.2	76.2	74.1	74.9	58.7	68.5	72.8	77.5	71.4	69.2	69.4	74.4	70.2
AG	77.2	80.6	76.1	78.0	78.0	59.3	67.4	76.6	77.7	72.6	70.7	69.0	79.5	71.6
ChatEval	71.8	73.0	60.9	62.4	67.0	51.1	65.5	68.4	74.0	63.9	66.1	55.7	70.3	64.4
CheckEval	76.1	80.7	74.8	78.4	77.5	60.8	68.9	74.9	76.8	<b>75.1</b>	73.4	65.8	76.4	71.5
RocketEval	76.4	81.0	75.9	79.1	78.1	61.2	69.2	75.4	77.0	74.9	72.9	66.9	78.1	72.0
UCS (EN)	<b>79.7</b>	<b>84.3</b>	<b>80.0</b>	<b>82.2</b>	<b>81.6</b>	<b>64.7</b>	67.7	<b>82.4</b>	<b>79.0</b>	<b>75.1</b>	<b>74.6</b>	<b>70.0</b>	<b>85.5</b>	<b>74.9</b>

Table 1: Balanced Accuracy (BA) across methods and languages for the Qwen-235B model. Bold values indicate the best performance per column averaged across 3 runs. Our method UCS is trained on English only data (EN).

measures judge’s sensitivity to both positive and negative classes equally.

**Implementation Details.** For each dataset and model, all methods are evaluated using consistent prompting templates and fixed criteria-generation procedures. All transfer models are trained exclusively on English-labeled data unless otherwise noted, and are applied to the remaining languages without any target-language supervision. We use shallow neural network for transfer from English to other languages (See B.3 for details).

**Baselines.** We evaluate our approach against a broad set of strong LLM-based judge baselines spanning single-model prompting, multi-agent methods, and checklist-based evaluators. First, we include zero-shot LLM judges following standard prompting setups from previous work (Blandón et al., 2025; Bavaresco et al., 2025), where the model is directly instructed to assess correctness or faithfulness without additional structure. We also include chain-of-thought (CoT) prompting, which has been shown to improve reasoning in LLM judges (Zheng et al., 2023). Another baseline we consider is prompting LLM with annotation guidelines (AG) that were given to humans, similar to (Blandón et al., 2025), see Appendix B.

We further compare against ChatEval (Chan et al.), a multi-agent debate-style evaluator in which models critique and challenge each other’s assessments. Such systems have demonstrated strong performance on complex reasoning and judgment tasks, but require substantial orchestration overhead. Finally, we include two recent checklist-based evaluation frameworks: RocketEval (Wei et al.) and CheckEval (Lee et al., 2025), which generate or refine structured criteria to guide LLM judgments. These methods share our goal of improving structure and interpretability in LLM-based evaluation; however, unlike our approach,

they do not produce a unified latent judge representation that supports cross-lingual transfer <sup>1</sup>.

## 5 Results

In this section, we evaluate the proposed criteria-based framework across languages and datasets. We focus on the language expansion setting and report the performance of UCS when trained on the English portion of the multilingual dataset, denoted by UCS (EN). We then present a series of analyses that unpack the sources of these gains: sample-efficiency curves quantifying how much English supervision is required to reach stable cross-lingual performance (§ 5.2), an analysis studying the most predictive criterion dimensions (§ 5.3), a comparison of alternative transfer models (§ 5.4), a comparison of criterion aggregation strategies (§ 5.5) and an analysis of inference cost trade-offs (§ 5.6).

### 5.1 Cross-Lingual Transfer

Table 1 reports balanced accuracy (BA) for the Qwen-235B judge on MEMERAG and mFACE, comparing the proposed UCS framework trained on English data against prompting-based, debate-style, and checklist-based baselines. Across both datasets, UCS (EN) delivers the strongest and most consistent cross-lingual performance, achieving the best or tied-best results in the majority of evaluated languages (11 out of 12). Appendix Tables 6 and 7 report results averaged over multiple random seeds, showing that UCS maintains consistent improvements across languages while exhibiting low variance across runs.

On MEMERAG, UCS (EN) achieves the highest performance across all four non-English languages, with an average BA of 81.6 compared to 78.1 for RocketEval, the strongest baseline. Improvements are consistent across typologically di-

<sup>1</sup>We reproduce the baseline results using publicly code.

verse languages, ranging from 3.1 to 4.1 BA points and spanning Germanic (DE), Romance (FR, ES), and Indo-Aryan (HI) languages. This uniformity supports the claim that the criteria-based representation is genuinely language-agnostic for this task.

On mFACE, UCS (EN) achieves best-or-tied performance on seven of the eight evaluated languages, with an average BA of 74.9 compared to 72.0 for RocketEval. Notably, the largest improvements are observed for French (+7.0) and Spanish (+7.4) – two high-resource languages where strong baseline performance might be expected. This suggests that the structured criteria space captures evaluation dimensions relevant to summarization faithfulness that are not easily elicited through direct prompting, and that this benefit is most pronounced for languages where the LLM can reliably interpret and respond to the criteria. Consistent gains are also observed on lower-resource languages such as Amharic (+3.5) and Hindi (+3.1), further demonstrating the breadth of cross-lingual generalization.

One notable exception is observed on Burmese (MY), where the zero-shot baseline achieves the highest performance across all methods (70.2), outperforming not only UCS (67.7) but also all other baselines. This suggests that Burmese may exhibit language-specific properties, such as script complexity or limited LLM pretraining coverage, that are not well captured by the shared criteria space, and that direct prompting may be more robust in such cases. Despite this exception, UCS maintains the strongest overall performance profile across languages, demonstrating the robustness of the proposed criteria-based representation for cross-lingual judge transfer. Appendix Table 5 reports the full cross-lingual results across all judge backbones evaluated in this work. The results are consistent with the main findings: UCS-based judges achieve the best or near-best performance across most language–model combinations on both MEMERAG and mFACE, indicating that the improvements are not specific to an LLM backbone.

## 5.2 Sample Efficiency

A practical consideration for language expansion is how much labeled data in a high-resource language is required to enable effective transfer. Although our framework trains only a lightweight transfer module on English-labeled criterion features, understanding the supervision requirements for stable cross-lingual performance is important for assessing the practical viability of the approach.

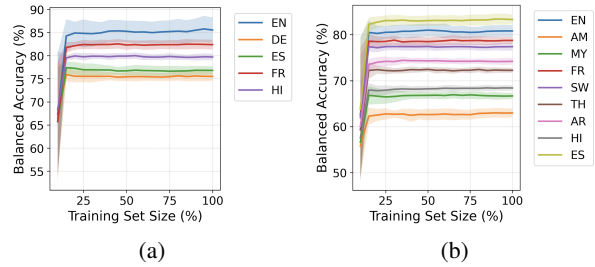


Figure 2: Sample efficiency of training the transfer module on English data. Performance is shown as a function of the percentage of English training data used. Results are reported for MEMERAG (a) and mFACE (b).

To study this, we construct a fixed random split of the English data, holding out 30% as a test set. From the remaining 70%, we vary the amount of labeled training data used to fit the transfer module, ranging from 5% to 100% of the available training portion. Each resulting model is evaluated both on the held-out English test set and on all target languages, allowing us to examine how increasing English supervision affects both in-language performance and cross-lingual generalization.

Figure 2 shows that performance improves rapidly with a small fraction of labeled data and stabilizes once approximately 20–30% of the English training data is used. Beyond this point, additional supervision yields only marginal gains across all languages.

Importantly, the same trend is observed for both English and target languages: as the transfer module improves on English, performance increases consistently across languages. This suggests that the criteria-based intermediate representation enables efficient cross-lingual generalization, requiring only a modest amount of labeled data in a single high-resource language.

## 5.3 Criteria Importance

Our transfer framework assumes that the criteria representation captures evaluation signals that are meaningful across languages. A natural question is whether the relative importance of these evaluation dimensions is preserved across languages. We hypothesize that languages whose importance profile criteria are more similar to English should exhibit stronger transfer performance.

To examine this, we analyze the importance of individual criteria before concept aggregation. While the transfer model operates on concept-level averages, criterion-level analysis provides a finer diagnostic view of how evaluation signals are priori-

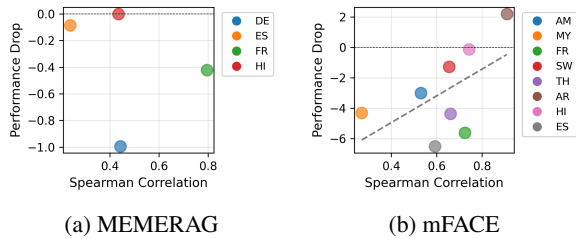


Figure 3: Relationship between English–target criteria importance alignment and performance change when training with only the top-10 English-selected criteria.

507 tized across languages.

508 For each language  $\ell$ , we train a Random Forest  
 509 classifier using that language’s criterion responses  
 510 together with the corresponding human labels, and  
 511 extract feature importance scores using Gini im-  
 512 purity reduction (Nembrini et al., 2018). We then  
 513 compute the Spearman rank correlation (Spearman,  
 514 1961) between the English importance profile and  
 515 the language-specific importance profile. To evalu-  
 516 ate whether this alignment matters in practice,  
 517 we simulate a constrained setting where only the  
 518 top-10 criteria selected according to English impor-  
 519 tance are used to train the classifier for each lan-  
 520 guage. We then measure the performance change  
 521 relative to using the full criteria set. Figure 3  
 522 plots, for each language, the relationship between  
 523 English–target importance correlation and the re-  
 524 sulting performance change.

525 In mFACE, we observe that languages with  
 526 lower importance alignment generally experience  
 527 larger performance degradation. For example,  
 528 Burmese (MY) shows one of the lowest correla-  
 529 tions with English and also corresponds to one  
 530 of the weaker-performing languages in Table 1.  
 531 While this observation is only suggestive, it indi-  
 532 cates that differences in evaluation priorities across  
 533 languages can contribute to variation in cross-  
 534 lingual transfer performance. For MEMERAG, we  
 535 observe only minor performance differences when  
 536 restricting the model to English-selected criteria;  
 537 the drops remain below 1% across languages and  
 538 do not show a clear relationship with importance  
 539 alignment.

540 Figure 4 visualizes full distribution of impor-  
 541 tance scores across all evaluated languages where  
 542 some languages exhibit importance profiles similar  
 543 to English, others show noticeable differences in  
 544 how evaluation dimensions are prioritized. Overall,  
 545 these results indicate that the relative importance of  
 546 evaluation criteria is not universally shared across

Method	MEMERAG	mFACE
LogReg	80.9	78.2
SVM	78.5	<b>79.0</b>
KNN	77.1	70.9
RF	76.4	69.8
XGBoost	74.3	71.1
NN	<b>81.5</b>	78.9

Table 2: BA averaged across languages for different transfer modules on the Qwen-235B model.

547 languages. Selecting a fixed subset of criteria  
 548 based solely on English importance can harm per-  
 549 formance for languages with different importance  
 550 profiles. This observation supports the design of  
 551 our full criteria-based representation, which allows  
 552 the transfer module to learn language-appropriate  
 553 weightings rather.

#### 5.4 Alternative Training Models 554

555 While our primary transfer model uses a shallow  
 556 neural network to map criteria-response represen-  
 557 tations to final judgments, we also evaluate alterna-  
 558 tive transfer modules to assess the sensitivity of our  
 559 framework to the choice of predictor. These include  
 560 logistic regression as a linear baseline, SVMs for  
 561 margin-based classification, KNN as an instance-  
 562 based method, and tree-based models such as Ran-  
 563 dom Forests and XGBoost that can capture nonlin-  
 564 ear feature interactions.

565 Table 2 shows that several simple models per-  
 566 form competitively when operating on the criteria-  
 567 based representation. Logistic regression and SVM  
 568 achieve strong results, indicating that much of the  
 569 signal captured by the criteria representation is lin-  
 570 early separable.

571 The shallow neural network achieves the best  
 572 performance on MEMERAG and remains competi-  
 573 tive on mFACE, where SVM slightly outperforms it.  
 574 Overall, the neural network provides the most con-  
 575 sistent performance across datasets and languages.

576 In contrast, instance-based and tree-based meth-  
 577 ods perform substantially worse. These models  
 578 appear to overfit to patterns in the English training  
 579 data and generalize less effectively to other lan-  
 580 guages, suggesting that simpler parametric models  
 581 are better suited for cross-lingual transfer here.

#### 5.5 Transfer Learning vs LLM-based Aggregator 582

583 The criteria-response vector  $z(x)$  produced by the  
 584 second stage must be aggregated into a final bi-  
 585 nary judgment. In our primary framework, this  
 586

Method	MEMERAG	mFACE
UCS (Trained)	<b>81.5</b>	<b>74.9</b>
UCS (LLM)	77.2	74.3

Table 3: Average BA across languages on MEMERAG and mFACE. UCS (Trained) uses a lightweight transfer module trained on English-labeled data; UCS (LLM) uses the LLM itself as a training-free aggregator over the criteria responses.

is achieved by the lightweight transfer module  $f_\theta$  trained on English-labeled data, which learns a calibrated mapping from the concept-level representation  $s(x)$  to the predicted label  $\hat{y}$ .

An alternative we consider an *LLM-based aggregator*, in which the LLM is prompted to produce a final judgment given the input  $x$ , the criteria  $Q$ , and their corresponding responses  $z(x)$ . This approach requires no additional training and may leverage the LLM’s broader world knowledge and reasoning capabilities to capture interactions between criteria. However, it does not benefit from calibration on human-labeled data and may be sensitive to prompt phrasing and model-specific biases.

Table 3 demonstrates a trade-off between training-free and trained aggregation. The LLM-based aggregator does not require labeled data beyond the criteria responses and can be applied directly. However, when labeled English data is available, training a lightweight transfer module on the criteria-based representation consistently yields a higher average balanced accuracy on both datasets.

## 5.6 Inference Cost

Compared to zero-shot prompting, our criteria-based framework incurs additional inference cost due to explicit criteria generation and criterion-level evaluation. While a zero-shot judge requires a single LLM call per sample, our approach introduces structured evaluation steps that increase prompt length and the number of LLM calls.

To analyze the trade-offs between decomposition granularity and inference efficiency, we evaluate four prompt variants differing in how criteria are generated and scored. Generation can be joint i.e. all evaluation dimensions produced in one LLM call or per-concept, i.e. each concept generates its criteria separately). Scoring follows the same options: a single joint judgment over all criteria or separate per-concept scores.

Table 4 compares these prompt variants in terms of both inference cost and performance. Perform-

Generation	Scoring	MEMERAG	mFACE
Joint	Joint	79.1	77.4
Joint	Per-concept	81.0	78.6
Per-concept	Joint	80.2	78.0
Per-concept	Per-concept	<b>81.5</b>	<b>78.9</b>

Table 4: Average BA across languages on MEMERAG and mFACE for UCS under different decompositions. Joint denotes evaluating all criteria in one LLM call, while Per-concept denotes one call per concept.

mance improves consistently as the evaluation process is more finely decomposed. In particular, scoring at the per-concept level consistently outperforms joint scoring over all criteria. The best results are obtained when both criteria generation and scoring are performed at the per-concept level, indicating that finer-grained criterion-level signals provide a more discriminative intermediate representation for the downstream transfer module.

## 6 Conclusion

In this paper, we introduced a decomposition-based framework for LLM judges that enables sample-efficient cross-lingual transfer through a shared set of language-agnostic evaluation criteria. We proposed Universal Criteria Sets (UCS), which structure evaluation into explicit dimensions and produce a transparent intermediate representation of the judgment process. A lightweight transfer module trained on English-labeled data maps this representation to final judgments and generalizes directly to new languages without requiring target-language supervision. Across experiments on MEMERAG and mFACE, UCS achieves strong and consistent performance while requiring only a small amount of supervision in a single high-resource language. Beyond improved cross-lingual performance, the criteria-based representation provides interpretable insights into how evaluation signals contribute to final judgments, and our criteria-importance analysis reveals that the degree of cross-lingual alignment varies by task. Together, these results highlight the potential of structured, criteria-based representations as a principled foundation for building reliable and interpretable multilingual LLM-based judges. The structured nature of this representation also allows future work to explore criterion-level judgments to construct structured reward signals or rubric-based feedback for reinforcement learning and alignment.

## 667 Limitations

668 First, the approach inherits the sensitivity from the  
669 underlying LLM. Criterion-level responses may  
670 vary with prompt phrasing, decoding settings, or  
671 model updates. Although decomposition reduces  
672 some instability by structuring evaluation through  
673 fixed criteria, the system remains dependent on  
674 the robustness of the base model. Future work  
675 could investigate prompt-invariant formulations or  
676 uncertainty-aware calibration. Second, our method  
677 assumes that evaluation criteria capture stable se-  
678 mantic dimensions shared across languages. This  
679 assumption may not hold for culturally specific  
680 judgments, stylistic norms, or tasks where evalu-  
681 ation standards differ substantially across regions.  
682 Finally, the framework relies on English supervi-  
683 sion to learn the transfer module. Although this  
684 reduces the need for multilingual labels, it can in-  
685 troduce biases present in English evaluation data  
686 that need to be studied. Finally, we do not systemat-  
687 ically explore the impact of variance that cascades  
688 through each stage of UCS generation.

## 689 References

690 Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Alt-  
691 man, Andy Applebaum, Edwin Arbus, Rahul K  
692 Arora, Yu Bai, Bowen Baker, Haiming Bao, and 1  
693 others. 2025. gpt-oss-120b & gpt-oss-20b model  
694 card. *arXiv preprint arXiv:2508.10925*.

695 Roei Aharoni, Shashi Narayan, Joshua Maynez,  
696 Jonathan Herzig, Elizabeth Clark, and Mirella La-  
697 pata. 2023. Multilingual summarization with factual  
698 consistency evaluation. In *Findings of the Associa-  
699 tion for Computational Linguistics: ACL 2023*, pages  
700 3562–3591.

701 Carlo Alfano, Aymen Al Marjani, Zeno Jonke, Amin  
702 Mantrach, Saab Mansour, and Marcello Federico.  
703 2025. Multilingual self-taught faithfulness evalu-  
704 tors. *arXiv preprint arXiv:2507.20752*.

705 Anna Bavaresco, Raffaella Bernardi, Leonardo Berto-  
706 lazzi, Desmond Elliott, Raquel Fernández, Albert  
707 Gatt, Esam Ghaleb, Mario Giulianelli, Michael  
708 Hanna, Alexander Koller, and 1 others. 2025. Llms  
709 instead of human judges? a large scale empirical  
710 study across 20 nlp evaluation tasks. In *Proceedings  
711 of the 63rd Annual Meeting of the Association for  
712 Computational Linguistics (Volume 2: Short Papers)*,  
713 pages 238–255.

714 María Andrea Cruz Blandón, Jayasimha Talur, Bruno  
715 Charron, Dong Liu, Saab Mansour, and Marcello Fed-  
716 erico. 2025. Memerag: A multilingual end-to-end  
717 meta-evaluation benchmark for retrieval augmented  
718 generation. *arXiv preprint arXiv:2502.17163*.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, 719  
Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan 720  
Liu. Chateval: Towards better llm-based evaluators 721  
through multi-agent debate. In *The Twelfth Interna- 722  
tional Conference on Learning Representations*. 723

Jiayi Chang, Mingqi Gao, Xinyu Hu, and Xiaojun 724  
Wan. 2025. Exploring the multilingual nlg evalu- 725  
ation abilities of llm-based evaluators. *arXiv preprint 726  
arXiv:2503.04360*. 727

Dongping Chen, Ruoxi Chen, Shilin Zhang, Yaochen 728  
Wang, YINUO Liu, Huichi Zhou, Qihui Zhang, Yao 729  
Wan, Pan Zhou, and Lichao Sun. 2024. Mllm-as- 730  
a-judge: Assessing multimodal llm-as-a-judge with 731  
vision-language benchmark. In *Forty-first Interna- 732  
tional Conference on Machine Learning*. 733

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta- 734  
sios Nikolas Angelopoulos, Tianle Li, Dacheng Li, 735  
Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E 736  
Gonzalez, and 1 others. 2024. Chatbot arena: An 737  
open platform for evaluating llms by human prefer- 738  
ence. In *Forty-first International Conference on 739  
Machine Learning*. 740

Sumanth Doddapaneni, Mohammed Safi Ur Rah- 741  
man Khan, Dilip Venkatesh, Raj Dabre, Anoop 742  
Kunchukuttan, and Mitesh M Khapra. 2025. Cross- 743  
lingual auto evaluation for assessing multilingual 744  
llms. In *Proceedings of the 63rd Annual Meeting of 745  
the Association for Computational Linguistics (Vol- 746  
ume 1: Long Papers)*, pages 29297–29329. 747

Mateo Espinosa Zarlenga, Pietro Barbiero, Gabriele 748  
Ciravegna, Giuseppe Marra, Francesco Giannini,  
749 Michelangelo Diligenti, Zohreh Shams, Frederic Pre-  
750 cioso, Stefano Melacci, Adrian Weller, and 1 oth-  
751 ers. 2022. Concept embedding models: Beyond  
752 the accuracy-explainability trade-off. *Advances in  
753 neural information processing systems*, 35:21400–  
754 21413. 755

Zhaopeng Feng, Jiayuan Su, Jiamei Zheng, Jiahua Ren,  
756 Yan Zhang, Jian Wu, Hongwei Wang, and Zuozhu  
757 Liu. 2025. M-mad: Multidimensional multi-agent  
758 debate for advanced machine translation evaluation.  
759 In *Proceedings of the 63rd Annual Meeting of the  
760 Association for Computational Linguistics (Volume  
761 1: Long Papers)*, pages 7084–7107. 762

Xiyan Fu and Wei Liu. 2025. How reliable is  
763 multilingual llm-as-a-judge? *arXiv preprint  
764 arXiv:2505.12201*. 765

Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan,  
766 Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen,  
767 Shengjie Ma, Honghao Liu, and 1 others. A survey  
768 on llm-as-a-judge. *The Innovation*. 769

Zixian Huang, Wenhao Zhu, Gong Cheng, Lei Li, and  
770 Fei Yuan. 2024. Mindmerger: Efficiently boosting  
771 llm reasoning in non-english languages. *Advances in  
772 Neural Information Processing Systems*, 37:34161–  
773 34187. 774

775	Alex Kim, Keonwoo Kim, and Sangwon Yoon. 2024.	José Pombal, Dongkeun Yoon, Patrick Fernandes, Ian	829
776	Debate: Devil’s advocate-based assessment and text	Wu, Seungone Kim, Ricardo Rei, Graham Neubig,	830
777	evaluation. In <i>Findings of the Association for Com-</i>	and André FT Martins. 2025. M-prometheus: A	831
778	<i>putational Linguistics ACL 2024</i> , pages 1885–1897.	suite of open multilingual llm judges. <i>arXiv preprint</i>	832
		<i>arXiv:2504.04953</i> .	833
779	Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen	Xiao Pu, Mingqi Gao, and Xiaojun Wan. 2023.	834
780	Mussmann, Emma Pierson, Been Kim, and Percy	Summarization is (almost) dead. <i>arXiv preprint</i>	835
781	Liang. 2020. Concept bottleneck models. In <i>Inter-</i>	<i>arXiv:2309.09558</i> .	836
782	<i>national conference on machine learning</i> , pages		
783	5338–5348. PMLR.	Lisa Schut, Yarin Gal, and Sebastian Farquhar. Do	837
784	Mahi Kolla, Siddharth Salunkhe, Eshwar Chandrasekha-	multilingual llms think in english? In <i>ICLR 2025</i>	838
785	ran, and Koustuv Saha. 2024. Llm-mod: Can large	<i>Workshop on Building Trust in Language Models and</i>	839
786	language models assist content moderation? In <i>Ex-</i>	<i>Applications</i> .	840
787	<i>extended Abstracts of the CHI Conference on Human</i>		
788	<i>Factors in Computing Systems</i> , pages 1–8.	Ivaxi Sheth and Samira Ebrahimi Kahou. 2023. Auxil-	841
		iary losses for learning generalizable concept-based	842
789	Yukyung Lee, Joonghoon Kim, Jaehee Kim, Hyowon	models. <i>Advances in Neural Information Processing</i>	843
790	Cho, Jaewook Kang, Pilsung Kang, and Najoung	<i>Systems</i> , 36:26966–26990.	844
791	Kim. 2025. Checkeval: A reliable llm-as-a-judge		
792	framework for evaluating text generation using check-	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang,	845
793	lists. In <i>Proceedings of the 2025 Conference on Em-</i>	Suraj Srivats, Soroush Vosoughi, Hyung Won Chung,	846
794	<i>pirical Methods in Natural Language Processing</i> ,	Yi Tay, Sebastian Ruder, Denny Zhou, and 1 others.	847
795	pages 15782–15809.	Language models are multilingual chain-of-thought	848
		reasoners. In <i>The Eleventh International Conference</i>	849
796	Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad	<i>on Learning Representations</i> .	850
797	Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-		
798	tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu,	Hwanjun Song, Hang Su, Igor Shalyminov, Jason Cai,	851
799	and 1 others. 2025. From generation to judgment:	and Saab Mansour. 2024. Finesure: Fine-grained	852
800	Opportunities and challenges of llm-as-a-judge. In	summarization evaluation using llms. <i>arXiv preprint</i>	853
801	<i>Proceedings of the 2025 Conference on Empirical</i>	<i>arXiv:2407.00908</i> .	854
802	<i>Methods in Natural Language Processing</i> , pages		
803	2757–2791.	Sofia Eleni Spatharioti, David M Rothschild, Daniel G	855
804	Chin-Yew Lin. 2004. Rouge: A package for automatic	Goldstein, and Jake M Hofman. 2023. Compar-	856
805	evaluation of summaries. In <i>Text summarization</i>	ing traditional and llm-based search for consumer	857
806	<i>branches out</i> , pages 74–81.	choice: A randomized experiment. <i>arXiv preprint</i>	858
		<i>arXiv:2307.03744</i> .	859
807	Adian Liusie, Potsawee Manakul, and Mark Gales. 2024.	Charles Spearman. 1961. The proof and measurement	860
808	Llm comparative assessment: Zero-shot nlg evalua-	of association between two things.	861
809	tion through pairwise comparisons using large lan-		
810	guage models. In <i>Proceedings of the 18th conference</i>	Elior Sulem, Omri Abend, and Ari Rappoport. 2018.	862
811	<i>of the European chapter of the Association for Com-</i>	Bleu is not suitable for the evaluation of text simpli-	863
812	<i>putational Linguistics (volume 1: long papers)</i> , pages	fication. <i>arXiv preprint arXiv:1810.05995</i> .	864
813	139–151.		
814	Stefano Nembrini, Inke R König, and Marvin N Wright.	Chung-En Sun, Tuomas Oikarinen, Berk Ustun, and	865
815	2018. The revival of the gini importance? <i>Bioinfor-</i>	Tsui-Wei Weng. Concept bottleneck large language	866
816	<i>matics</i> , 34(21):3711–3718.	models. In <i>The Thirteenth International Conference</i>	867
		<i>on Learning Representations</i> .	868
817	Joshua W Ohde, Lauren M Rost, and Joshua D Over-	Aman Singh Thakur, Kartik Choudhary, Venkat Srinik	869
818	gaard. 2025. The burden of reviewing llm-generated	Ramayapally, Sankaran Vaidyanathan, and Dieuwke	870
819	content.	Hupkes. 2025. Judging the judges: Evaluating align-	871
		ment and vulnerabilities in llms-as-judges. In <i>Pro-</i>	872
820	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	<i>ceedings of the Fourth Workshop on Generation,</i>	873
821	Jing Zhu. 2002. Bleu: a method for automatic evalua-	<i>Evaluation and Metrics (GEM<sup>2</sup>)</i> , pages 404–430.	874
822	tion of machine translation. In <i>Proceedings of the</i>		
823	<i>40th annual meeting of the Association for Computa-</i>	Prapti Trivedi, Aditya Gulati, Oliver Molenschot,	875
824	<i>tional Linguistics</i> , pages 311–318.	Meghana Arakkal Rajeev, Rajkumar Ramamurthy,	876
		Keith Stevens, Tanveesh Singh Chaudhery, Jahnavi	877
825	Eleonora Poeta, Gabriele Ciravegna, Eliana Pastor,	Jambholkar, James Zou, and Nazneen Rajani. Self-	878
826	Tania Cerquitelli, and Elena Baralis. 2023. Concept-	rationalization improves llm as a fine-grained judge.	879
827	based explainable artificial intelligence: A survey.		
828	<i>ACM Computing Surveys</i> .	Wenxuan Wang, Zhaopeng Tu, Chang Chen, Youliang	880
		Yuan, Jen-tse Huang, Wenxiang Jiao, and Michael	881
		Lyu. 2024. All languages matter: On the multilingual	882

883 safety of llms. In *Findings of the Association for*  
884 *Computational Linguistics: ACL 2024*, pages 5865–  
885 5877.

886 Tianjun Wei, Wei Wen, Ruizhi Qiao, Xing Sun, and  
887 Jianghong Ma. Rocketeval: Efficient automated llm  
888 evaluation via grading checklist. In *The Thirteenth*  
889 *International Conference on Learning Representa-*  
890 *tions*.

891 Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan,  
892 Lidia Sam Chao, and Derek Fai Wong. 2025. A  
893 survey on llm-generated text detection: Necessity,  
894 methods, and future directions. *Computational Lin-*  
895 *guistics*, 51(1):275–338.

896 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,  
897 Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
898 Gao, Chengen Huang, Chenxu Lv, and 1 others.  
899 2025. Qwen3 technical report. *arXiv preprint*  
900 *arXiv:2505.09388*.

901 Murong Yue. 2025. A survey of large language  
902 model agents for question answering. *arXiv preprint*  
903 *arXiv:2503.19213*.

904 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan  
905 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
906 Zhuohan Li, Dacheng Li, Eric Xing, and 1 others.  
907 2023. Judging llm-as-a-judge with mt-bench and  
908 chatbot arena. *Advances in neural information pro-*  
909 *cessing systems*, 36:46595–46623.

## A Future Work

In this work, we evaluated cross-lingual transfer across multiple languages and datasets. Future work can extend this to other languages not evaluated here. In particular, language-specific analyses could help explain the strong zero-shot performance of the LLM Judge in Burmese. Figure 3 illustrates the relationship between English–target criteria importance alignment, showing that the relative importance of evaluation criteria can vary across languages and datasets. This suggests that not all criteria contribute equally to evaluation quality in cross-lingual settings. In particular, some criteria appear consistently less important across languages, while others show stronger and more stable alignment. These observations indicate the need to identify the most effective evaluation criteria for multilingual assessment. Future work could focus on selecting or learning criteria that provide the strongest signal across languages.

Beyond cross-lingual evaluation, the criteria-based decomposition introduced in this work opens several promising directions. Because the framework represents judgments through explicit evaluation dimensions, it provides a structured interface for analyzing and controlling LLM evaluation behavior.

Another promising direction is leveraging criteria representations to support human-in-the-loop evaluation. Structured criteria could allow human evaluators to provide targeted feedback on specific dimensions, which can then be incorporated to refine the transfer model or adjust evaluation standards. This may enable more transparent and controllable evaluation pipelines.

Finally, the criteria representation provides a natural foundation for studying the reliability of LLM judges. Future work could investigate uncertainty estimation, agreement across multiple judge models, or ensemble approaches that operate at the criterion level.

Future work could also explore how criterion-level representations can support reinforcement learning and alignment of language models. Current alignment methods often rely on scalar reward signals derived from preference comparisons or holistic judgments, which provide limited insight into why a response is preferred. In contrast, criterion-based evaluation decomposes quality into explicit dimensions, such as factual consistency, relevance, or completeness. These structured sig-

nals could be used to construct richer reward functions that guide models toward satisfying multiple evaluation dimensions simultaneously. Moreover, criterion-level feedback may enable more interpretable and controllable alignment, allowing training objectives to emphasize specific aspects of behavior or adapt across languages and domains. Investigating how such structured evaluation signals can be integrated into reinforcement learning or preference optimization pipelines is a promising direction for future work.

### A.1 Risks

Our framework introduces several potential risks. First, reliance on English-labeled data for training the transfer module may propagate biases present in English evaluation standards, potentially leading to unfair or misaligned judgments in other languages. Second, the assumption of language-agnostic evaluation criteria may overlook culturally specific norms, stylistic preferences, or context-dependent interpretations of quality, resulting in systematic evaluation errors. Third, the approach depends on the stability of underlying LLM outputs; variations due to prompt phrasing or model updates may affect criterion-level responses and downstream predictions. Finally, as with other automated evaluators, there is a risk of over-reliance on LLM-based judgments in high-stakes settings without sufficient human oversight.

## B Reproducibility

We will release our code. The code is uploaded as supplementary material. All criterion-generation and scoring prompts are included in the appendix.

### B.1 Model Inference Settings.

All LLM-based evaluations were performed with the temperature set to 0 and the top- $p$  set to 1 to ensure deterministic outputs. We report the exact model snapshots used in our experiments. All models were accessed through the Amazon Bedrock API. No additional fine-tuning of the backbone LLMs was performed. We reproduced the results for all of the baselines reported in the paper.

### B.2 Cross-Lingual Transfer Setup.

For concept-level transfer, we train calibration models on English criterion-level representations and evaluate zero-shot on other languages without using any target-language labels. All experiments use `random_state=42` for reproducibility.

1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057

### B.3 Transfer Module Hyperparameters.

We evaluated multiple lightweight predictors. We evaluated several lightweight calibration models for concept-level transfer. The neural network consists of a single hidden layer with 32 units, trained with a learning rate of 0.01 for up to 2000 iterations. Logistic regression uses a regularization strength of  $C = 0.1$ , is trained for up to 2000 iterations, and applies balanced class weights to account for label imbalance. The random forest model uses 200 trees with a maximum depth of 10 and balanced class weights. The support vector machine uses  $C = 0.1$ , a gamma value set to “scale” enables probability estimates, and applies balanced class weights. Gradient boosting is configured with 100 estimators, a maximum depth of 5, and a learning rate of 0.1. Finally, the k-nearest neighbors classifier uses  $k = 5$  neighbors.

Unless otherwise stated, the reported results correspond to the NN model selected on English validation data.

### B.4 Data Splits and Evaluation.

We train the transfer module exclusively on labeled English data and evaluate on multilingual test splits without retraining or hyperparameter tuning. Balanced accuracy is used as the primary evaluation metric to account for label imbalance.

### B.5 Implementation Details.

All experiments were implemented in Python using standard machine learning libraries. Fixed random seeds were used across training runs. During transfer, no target-language supervision, translation, or synthetic augmentation was used .

### B.6 Annotation Guideline baseline.

A baseline we consider prompts the LLM with the annotation guidelines (AG) originally provided to human annotators, following the setup of Blandón et al. (2025). For MEMERAG, the annotation guidelines are available in Appendix A of the paper and describe the criteria used by annotators to assess faithfulness between the generated answer and the supporting evidence. For mFACE, we use the evaluation instructions provided in Figure 2 of the paper, which outline the conditions under which a summary should be considered faithful to the source article. In both cases, these guidelines are directly incorporated into the evaluation prompt to guide the LLM’s judgment.

## C Concept-level transfer

A central design choice in our framework is to perform transfer at the level of evaluation concepts rather than at the level of raw criterion interactions. Each dimension in  $z(x)$  corresponds to a semantically meaningful evaluation question (e.g., faithfulness, completeness, consistency). By learning a linear calibration over these concept-level signals, the transfer module estimates how much each evaluation dimension contributes to the final judgment.

An alternative would be to model interactions between criteria using a more expressive predictor. However, such approaches substantially increase the number of learnable parameters and, consequently, the number of labeled samples required for stable training. In multilingual settings where supervision is typically available only in English, this would lead to overfitting and poor generalization.

By constraining transfer to operate over independent semantic dimensions, we reduce sample complexity and improve stability. This design aligns with the intuition that evaluation structure is largely shared across languages, even when surface realizations differ. As a result, concept-level calibration enables effective cross-lingual transfer using limited source-language supervision.

## D Results

Table 5 reports the full cross-lingual evaluation results for all baselines and criteria-based methods across models and languages on MEMERAG and mFACE. Consistent with the main results, UCS-based judges generally achieve the strongest or among the strongest performance across languages and model backbones.

Figure 4 shows the cross-lingual alignment of criterion importance between English and target languages for MEMERAG and mFACE. The heatmaps illustrate how the relative importance of evaluation criteria varies across languages.

Table 6 and Table 7 report MEMERAG and mFACE results for Qwen-235B averaged over three runs with different seeds. The results show that UCS maintains consistent improvements across languages while exhibiting relatively low variance compared to most baselines.

Table 8 compares the transfer-module classifier with LLM-based aggregation of the criteria representations across languages. The trained transfer

1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107

Model	Method	MEMERAG (BA)				mFACE (BA)							
		DE	FR	ES	HI	AM	MY	FR	SW	TH	AR	HI	ES
OSS-20B	Zero-shot	70.7	81.6	80.4	81.7	59.4	65.6	73.7	79.2	76.8	69.2	68.5	78.0
	CoT	70.5	81.6	82.0	80.6	60.9	64.9	74.6	76.2	78.3	65.2	67.6	78.0
	AG	71.7	80.8	75.6	80.9	61.4	62.2	75.2	77.8	76.9	72.4	66.7	78.0
	ChatEval	64.2	73.5	71.4	72.0	54.7	60.1	66.9	70.3	68.4	62.9	61.4	71.1
	CheckEval	72.5	81.9	79.9	80.8	60.9	68.3	77.9	76.7	79.7	68.0	65.7	77.0
	RocketEval	67.9	76.9	74.8	75.6	57.9	63.4	70.8	73.5	71.6	66.5	64.9	74.2
	UCS (EN)	74.4	84.2	81.2	85.0	61.2	68.1	79.0	82.7	82.4	73.7	67.8	82.0
Qwen 3 32B	Zero-shot	74.9	73.7	71.4	78.5	56.6	63.5	71.9	72.3	69.2	65.0	67.9	76.0
	CoT	71.7	73.8	72.6	77.2	55.2	62.9	72.0	72.4	70.0	66.7	68.1	76.5
	AG	73.7	73.9	74.8	81.8	58.3	63.9	76.6	73.0	70.9	68.4	66.1	80.0
	ChatEval	66.7	68.5	63.4	66.9	50.9	56.3	61.2	65.1	61.7	59.4	56.9	66.0
	CheckEval	71.1	74.6	72.3	80.9	59.2	63.5	74.5	73.6	71.1	65.8	65.0	79.0
	RocketEval	70.0	72.9	68.7	72.4	54.8	59.7	65.8	69.0	66.9	63.9	61.4	70.4
	UCS (EN)	74.4	82.7	79.5	85.0	57.7	61.7	68.4	73.4	72.1	68.9	66.7	78.9
Llama70b	Zero-shot	69.9	70.3	68.0	73.8	55.6	57.9	65.3	73.2	60.4	63.4	62.8	75.8
	CoT	67.1	72.2	68.4	69.7	56.6	58.2	64.8	75.3	60.7	64.7	62.3	75.8
	AG	75.4	76.4	74.0	80.8	58.2	59.0	68.7	77.6	63.2	66.1	67.3	77.5
	ChatEval	68.4	71.3	66.1	69.8	53.9	57.4	63.9	68.9	64.8	62.1	59.8	69.6
	CheckEval	72.8	76.6	72.1	80.2	59.2	59.2	69.5	77.1	65.4	63.8	67.8	73.9
	RocketEval	71.8	74.9	70.5	73.9	57.1	61.0	67.9	72.3	69.5	66.3	64.0	73.1
	UCS (EN)	77.3	80.8	82.1	84.5	63.5	56.5	72.1	74.0	73.2	67.4	66.1	77.7
OSS-120B	Zero-shot	79.2	83.6	79.3	81.5	64.1	67.6	82.0	78.6	74.9	74.2	69.4	85.1
	CoT	78.3	82.9	78.8	80.9	64.7	67.0	81.4	79.4	74.1	73.9	69.9	83.8
	AG	79.4	84.1	79.0	82.1	65.2	66.7	82.6	79.7	75.0	74.6	69.7	85.3
	ChatEval	73.1	74.2	64.9	66.1	58.3	65.4	69.4	74.0	65.7	66.8	58.0	71.4
	CheckEval	78.6	83.8	78.5	81.9	65.7	67.9	82.1	79.1	75.2	75.0	69.1	84.7
	RocketEval	78.9	84.1	78.9	82.3	66.0	68.2	82.4	79.3	75.4	75.2	69.4	85.0
	UCS (EN)	79.4	83.9	79.7	81.9	64.4	67.8	82.4	79.0	75.1	74.6	69.7	85.6

Table 5: Balanced Accuracy (BA) of all baselines and criteria set judges across languages on MEMERAG and mFACE.

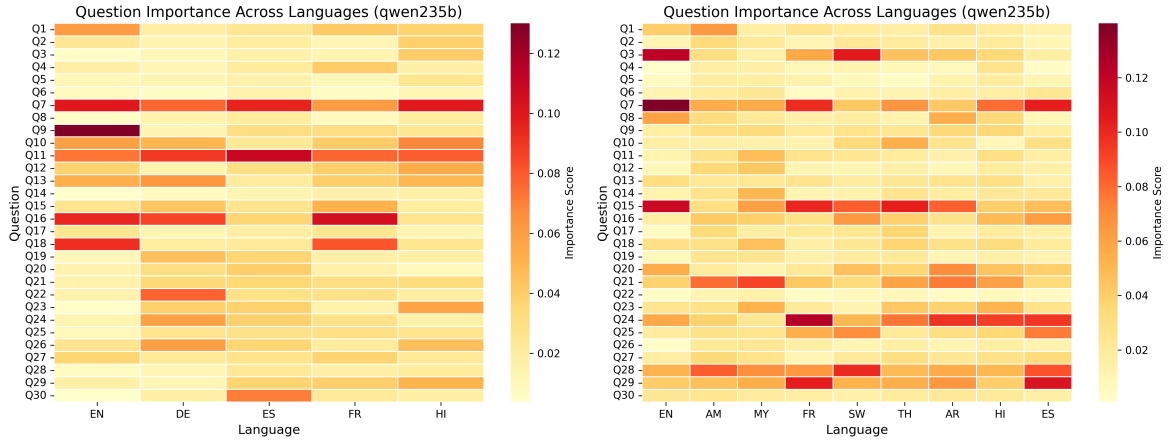


Figure 4: Relationship between English–target criteria importance alignment.

1108 module generally achieves higher performance, al-  
 1109 though LLM aggregation remains competitive in  
 1110 some cases.

## 1111 E Prompts

Method	DE	FR	ES	HI
Zero-shot	77.2±0.4	77.5±1.3	77.5±0.9	74.4±1.1
CoT	73.0±1.2	76.2±0.5	76.2±1.4	74.1±0.8
AG	77.2±0.7	80.6±1.1	76.1±0.4	78.0±1.4
ChatEval	71.8±1.5	73.0±0.6	60.9±1.9	62.4±0.7
CheckEval	76.1±0.5	80.7±1.4	74.8±0.7	78.4±1.1
RocketEval	76.4±1.9	81.0±1.2	75.9±1.7	79.1±1.4
UCS (EN)	<b>79.7±0.6</b>	<b>84.3±1.1</b>	<b>80.0±0.4</b>	<b>82.2±1.2</b>

Table 6: Balanced Accuracy (BA) on MEMERAG for Qwen-235B, reported as mean  $\pm$  standard deviation over 3 runs.

Method	AM	MY	FR	SW	TH	AR	HI	ES
Zero-shot	54.7±0.7	<b>70.2±1.3</b>	74.3±0.6	73.0±1.1	72.6±0.4	69.7±1.4	68.0±1.0	77.9±0.6
CoT	58.7±1.4	68.5±0.6	72.8±1.2	77.5±0.5	71.4±1.1	69.2±0.5	69.4±1.6	74.4±0.8
AG	59.3±0.8	67.4±1.5	76.6±0.7	77.7±1.1	72.6±0.6	70.7±1.2	69.0±0.7	79.5±1.2
ChatEval	51.1±1.6	65.5±0.9	68.4±1.3	74.0±0.5	63.9±1.8	66.1±0.6	55.7±1.1	70.3±1.5
CheckEval	60.8±1.3	68.9±0.4	74.9±1.6	76.8±0.7	<b>75.1±1.1</b>	73.4±0.8	65.8±1.5	76.4±0.6
RocketEval	61.2±2.1	69.2±1.1	75.4±1.7	77.0±1.3	74.9±1.6	72.9±2.0	66.9±1.4	78.1±1.8
UCS (EN)	<b>64.7±1.0</b>	67.7±1.4	<b>82.4±0.6</b>	<b>79.0±1.3</b>	<b>75.1±0.4</b>	<b>74.6±1.2</b>	<b>70.0±0.7</b>	<b>85.5±1.0</b>

Table 7: Balanced Accuracy (BA) on mFACE for Qwen-235B, reported as mean  $\pm$  standard deviation over 3 runs.

Dataset	Lang.	UCS (transfer)	UCS (LLM)
MEMERAG	DE	<b>79.7</b>	76.9
	FR	<b>84.3</b>	78.9
	ES	<b>80.0</b>	77.2
	HI	<b>82.2</b>	75.9
mFACE	AM	<b>64.7</b>	62.5
	MY	67.7	<b>69.6</b>
	FR	<b>82.4</b>	82.3
	SW	<b>79.0</b>	78.4
	TH	<b>75.1</b>	74.9
	AR	<b>74.6</b>	73.1
	HI	<b>70.0</b>	68.8
ES	<b>85.5</b>	84.5	

Table 8: Comparison of threshold-based training classifier and LLM-based aggregation of intermediate criteria representations for final judgment.

## Summary Evaluation Questions

### Factual Accuracy (6 questions)

1. Does the summary include any claims that are not explicitly supported by information in the source document?
2. Does the summary contradict any specific fact, detail, or relationship stated in the source document?
3. Are events, actions, or attributes in the summary correctly attributed to the individuals, entities, or sources named in the source document?
4. Does the summary present speculative, uncertain, or conditional information from the source as definitive or certain?
5. Does the summary accurately reflect the relative importance or prominence of key points as presented in the source document?
6. Does the summary introduce numerical data, statistics, or quantitative claims that differ from those in the source document?

### Contradiction Detection (6 questions)

7. Does the summary attribute a claim, opinion, or action to a source or entity that is not supported or explicitly stated in the source document?
8. Does the summary present a possibility or uncertainty as a definitive fact, thereby increasing the level of certainty beyond what is expressed in the source document?
9. Does the summary state that an event occurred, or a condition exists, when the source document explicitly indicates it did not happen or was not the case?
10. Does the summary include a causal relationship between two events that is not stated, implied, or supported by the source document?
11. Does the summary report a numerical value, statistic, or quantitative detail that contradicts the figures provided in the source document?
12. Does the summary describe an action or outcome as having happened in the past when the source document states it is planned, proposed, or speculative?

### Claim Support (6 questions)

13. Does the summary include any claims that are not explicitly supported by evidence or statements in the source document?
14. Are there any statements in the summary that directly contradict information provided in the source document?
15. Does the summary attribute actions, opinions, or statements to individuals or entities that are not assigned to them in the source document?
16. To what extent does the summary present speculative or conditional information from the source as definitive or certain?
17. Does the summary introduce new causal relationships or implications between events that are not stated or logically supported in the source document?
18. Are key details in the summary distorted in a way that alters the meaning or significance of the information presented in the source document?

### Misrepresentation Identification (6 questions)

19. Does the summary attribute a claim, opinion, or action to a source or entity that is not supported or explicitly stated in the source document?
20. Does the summary present a speculative or conditional statement from the source as a definitive fact?
21. Does the summary include a key event, outcome, or statistic that is not mentioned or implied in the source document?
22. Does the summary reverse, invert, or otherwise distort the causal or temporal relationship between two events described in the source document?
23. Does the summary exaggerate the strength, scope, or certainty of a finding, trend, or conclusion beyond what is stated in the source document?
24. Does the summary omit a critical qualifying condition, limitation, or exception present in the source that changes the interpretation of the information?

### Certainty Calibration (6 questions)

25. Does the summary present a claim as certain or definitive when the source document expresses it as uncertain, tentative, or conditional?
26. Does the summary introduce a level of confidence or precision (e.g., "proves," "definitely," "exactly") in a claim that is not supported by the degree of certainty used in the source document?
27. Are probabilities, frequencies, or likelihoods in the summary accurately reflected in terms of their magnitude and wording compared to the source document (e.g., "likely" vs. "possible" vs. "certain")?
28. Does the summary omit hedging language (e.g., "may," "suggests," "appears to") present in the source, thereby overstating the strength of a conclusion?
29. To what extent does the summary mirror the source document's attribution of claims to specific agents, studies, or evidence, without shifting responsibility or generalizing to broader consensus?
30. Does the summary elevate a hypothesis, preliminary finding, or speculative idea from the source to the status of an established fact?

## Evidence Support Evaluation Questions

### Factual Accuracy (6 questions)

1. To what extent does the evidence directly support the key claims in the answer, with no unsupported assertions?
2. How well do the retrieved passages contradict or fail to support any statements in the answer?
3. To what degree does the answer avoid overgeneralizing or making broad inferences beyond what is stated in the evidence?
4. How precisely does the answer align with the specificity (e.g., timeframes, quantities, conditions) provided in the evidence?
5. To what extent are all named entities, events, and relationships in the answer accurately reflected in the evidence passages?
6. How well does the answer refrain from introducing plausible but unverified details not present in the evidence?

### Completeness of Support (6 questions)

7. To what extent does the evidence directly support all key claims in the answer, with no unsupported assertions?
8. How well do the retrieved passages contain specific details or examples that match the level of specificity in the answer?
9. To what degree does the evidence fully cover the scope of the answer, including all sub-claims or components mentioned?
10. How consistently do the passages align with the answer without introducing contradictions or conflicting information?
11. To what extent does the answer avoid overgeneralizing beyond what is reasonably supported by the evidence?
12. How well do the passages provide sufficient context or explanation to justify causal or inferential claims made in the answer?

### Specificity Alignment (6 questions)

13. To what extent does the answer reflect the same level of specificity as the evidence, avoiding unwarranted generalizations or oversimplifications?
14. How well do the key claims in the answer map directly to specific details or data points in the evidence, rather than relying on vague or peripheral information?
15. To what degree does the evidence support the precise scope (e.g., time frame, population, location) asserted in the answer without overextension?
16. How closely does the answer avoid introducing concepts or conclusions that are more specific than what is warranted by the evidence?
17. To what extent are named entities, quantities, or relationships in the answer explicitly grounded in corresponding specific mentions within the evidence?
18. How well does the answer maintain alignment with the evidence by neither omitting critical qualifying conditions nor adding unsupported qualifiers?

### Consistency with Evidence (6 questions)

19. To what extent does the evidence directly support the key claims in the answer, with no unsupported assertions?
20. How well do the retrieved passages contradict or conflict with any statements in the answer?
21. To what degree does the answer avoid overgeneralizing beyond the scope or specificity of the evidence provided?
22. How closely does the answer align with the factual details and context present in the evidence, avoiding subtle distortions or misrepresentations?
23. To what extent can each component of a multi-part answer be individually justified by at least one evidence passage?
24. How well does the answer reflect the certainty level (e.g., tentative, definitive) expressed in the evidence, without introducing unwarranted confidence or ambiguity?

### Source Attribution (6 questions)

25. To what extent does the evidence explicitly attribute the claim to a credible source or original provider of information?
26. How well do the retrieved passages support the specificity of the claim, without introducing unsupported details or omitting critical qualifiers present in the source?
27. To what degree is the claim directly verifiable from the cited evidence, rather than requiring inference beyond what the source states?
28. How consistently does the answer reflect the source's intended meaning, avoiding misrepresentation or overgeneralization of the evidence?
29. To what extent does the evidence rule out contradictions or significant discrepancies with the claim being made?
30. How clearly is the connection between the evidence and the claim articulated, such that the support is traceable and transparent?

## Concept Generation - MEMERAG

"system\_prompt": "You are an evidence-grounding evaluator. Your goal is to identify key verification concepts for checking whether answers are supported by evidence."

"task\_prompt": "Generate 3-5 distinct verification concepts that are essential for evaluating whether an answer is supported by evidence passages.\n\nGuidelines:\n- Concepts should cover different aspects of evidence grounding (e.g., factual accuracy, completeness, specificity, consistency, source attribution).\n- Each concept should be general enough to apply across different topics and domains.\n- Concepts should target different failure modes where answers might not be properly supported.\n- Keep concepts concise (2-5 words each).\n\nReturn exactly this format:\n\n<concepts>\n<concept1>[Concept name]\n</concept1>\n<concept2>[Concept name]\n</concept2>\n...\n</concepts>"

### Question Generation - MEMERAG

"system\_prompt": "You are an evidence-grounding evaluator. Your goal is to generate reusable verification questions for checking whether an answer is supported by retrieved evidence passages."

"task\_prompt": "\*\*Verification Concept:\*\* concept4-6 clear, diverse, and reusable evaluation questions for verifying whether an answer is supported by evidence with respect to the concept 'concept'. These questions will later be scored on a \*\*0-10 scale indicating degree of evidence support\*\*.- Focus on whether an **answer** is supported by evidence passages.- Questions must be reusable across different topics and documents.- Cover different **evidence-grounding failure modes** (e.g., unsupported claims, contradictions, overgeneralization, specificity mismatch).- Each question should be written to support **graded (partial-to-full) scoring**, not just yes/no judgments.- Each question should target a distinct verification angle.- Wording should be concise and evaluative (e.g., To what extent is the claim supported by the evidence...; How well do the passages justify...).exactly this format:<questions><question1>: [Universal evaluation question]</question1>...</questions>"

### Question Application - MEMERAG

"system\_prompt": "You are an evidence-grounding evaluator. Your task is to determine whether an answer is supported by the provided evidence passages."

"task\_prompt": "\*\*Evidence Passages:\*\*\n {context}\n\n **Answer:\*\*** {answer\_segment}\n\n **Verification Concept:\*\*** {concept}\n\n **Evaluation Questions:\*\***\n {questions}\n\n For each question, assign a score from **0-10** indicating how well the evidence supports the answer with respect to that question.\n\n Base all judgments strictly on the provided evidence passages. Do not assume any external knowledge.\n\n Return exactly this format:\n\n <evaluation>\n {questions\_format\_example}\n </evaluation>"

## Concept Generation - mFACE

### concept\_generation

"system\_prompt": "You are an expert factual faithfulness evaluator. Your goal is to identify key faithfulness concepts for checking whether summaries are factually faithful to their source documents."

"task\_prompt": "Generate 3-5 distinct faithfulness concepts that are essential for evaluating whether a summary is factually faithful to its source document.\n\nGuidelines:\n- Concepts should cover different aspects of factual faithfulness (e.g., factual accuracy, contradiction detection, unsupported claims, misrepresentation, distorted certainty).\n- Each concept should be general enough to apply across different topics and domains.\n- Concepts should target different failure modes where summaries might not be properly supported by source documents.\n- Keep concepts concise (2-5 words each).\n\nReturn exactly this format:\n\n<concepts>\n<concept1>[Concept name]\n</concept1>\n<concept2>[Concept name]\n</concept2>\n...\n</concepts>"

## Question Generation - mFACE

### question\_generation

"system\_prompt": "You are an expert factual faithfulness evaluator. Your role is to construct high-quality, reusable evaluation questions that test whether summaries are factually faithful to their source documents."

"task\_prompt": "\*\*Faithfulness Concept:\*\* {concept}\n\nGenerate exactly 6 clear, diverse, and reusable evaluation questions for the faithfulness concept '{concept}'. Each question must require direct comparison between a summary and its source document to assess factual support, contradiction, or distortion under this concept.\n\nAnnotation Guidance to Apply While Writing Questions:\n- The question must be answerable ONLY by checking the source document.\n- The question should detect one of the following: unsupported claims, contradiction, misrepresentation, distorted certainty, or incorrect attribution.\n- The question must not depend on surface form (grammar, fluency, style, or verbosity).\n\nGuidelines:\n- Every question must explicitly or implicitly require verification against the source document.\n- Do NOT ask about readability, grammar, fluency, writing quality, or length.\n- Questions must be context-independent and reusable across domains.\n- All questions must stay strictly within the given faithfulness concept, but vary the factual scenario or common error pattern being tested.\n- Avoid redundancy: each question should target a distinct factual failure mode.\n- Wording must be concise, unambiguous, and evaluative (e.g., \"Does the summary...\" or \"To what extent does the summary...\").\n\nReturn exactly this format:\n\n<questions>\n<question1>\nQuestion: [Universal evaluation question]\n</question1>\n<question2>\nQuestion: [Universal evaluation question]\n</question2>\n<question3>\nQuestion: [Universal evaluation question]\n</question3>\n<question4>\nQuestion: [Universal evaluation question]\n</question4>\n<question5>\nQuestion: [Universal evaluation question]\n</question5>\n<question6>\nQuestion: [Universal evaluation question]\n</question6>\n</questions>"

## Question Application - mFACE

### question\_application

"system\_prompt": "You are an expert factual faithfulness evaluator. Your task is to score a summary against its source document using predefined faithfulness evaluation questions."

"task\_prompt": "\*\*Source Document:\*\*\n{context}\n\n\*\*Summary:\*\*\n{answer\_segment}\n\n\*\*Faithfulness Concept:\*\*\n{concept}\n\n\*\*Evaluation Questions:\*\*\n{questions}\n\nFor each question, evaluate whether the summary is factually faithful to the source document under that question.\n\nWhile scoring, explicitly consider whether the summary:\n- Contains information not stated in or directly inferable from the source\n- Contradicts the source\n- Introduces unsupported details\n- Misrepresents entities, quantities, relationships, or certainty\n- Overgeneralizes, narrows, or distorts scope\n\nAssign a score from 0-10 for each question:\n- 0 = Completely unfaithful (clear contradiction, fabrication, or unsupported claim)\n- 10 = Completely faithful (fully supported, correctly represented, no distortion)\n\nRules:\n- Base every score strictly on evidence from the source document.\n- Ignore grammar, fluency, style, and summary length.\n- Penalize both direct hallucinations and subtle distortions of meaning or certainty.\n\nProvide a brief, evidence-based justification for each score.\n\nReturn exactly this format:\n\n<evaluation>\n{questions\_format\_example}\n</evaluation>"