

BEYOND TRAINING FOR CULTURAL AWARENESS: THE ROLE OF DATASET LINGUISTIC STRUCTURE IN LARGE LANGUAGE MODELS

Reem I. Masoud*

University College London
King Abdulaziz University

Chen Feng

Queen’s University Belfast
University College London

Shunta Asano

The University of Tokyo
University College London

Saied Alshahrani

University of Bisha

Philip Colin Treleaven

University College London

Miguel R. D. Rodrigues

University College London
AI Centre, University College London

ABSTRACT

The global deployment of large language models (LLMs) has raised concerns about cultural misalignment, yet the linguistic properties of fine-tuning datasets used for cultural adaptation remain poorly understood. We adopt a dataset-centric view of cultural alignment and ask which linguistic properties of fine-tuning data are associated with cultural performance, whether these properties are predictive prior to training, and how these effects vary across models. We compute lightweight linguistic, semantic, and structural metrics for Arabic, Chinese, and Japanese datasets and apply principal component analysis (PCA) separately within each language. This design ensures that the resulting components capture variation among datasets written in the same language rather than differences between languages. The resulting components correspond to broadly interpretable axes related to semantic coherence, surface-level lexical and syntactic diversity, and lexical or structural richness, though their composition varies across languages. We fine-tune three major LLM families (LLaMA, Mistral, DeepSeek) and evaluate them on benchmarks of cultural knowledge, values, and norms. While PCA components correlate with downstream performance, these associations are strongly model-dependent. Through controlled subset interventions, we show that lexical-oriented components (PC3) are the most robust, yielding more consistent performance across models and benchmarks, whereas emphasizing semantic or diversity extremes (PC1–PC2) is often neutral or harmful. (Reproducibility notebook: https://anonymous.4open.science/r/dataset_quality_demo-C7C9/demo_subset_selection.ipynb.)

1 INTRODUCTION

LLMs have shown remarkable progress across diverse natural language processing (NLP) tasks. However, their performance often falls short in cross-cultural settings, where linguistic variation, cultural norms, and local knowledge shape how users interpret and engage with model outputs (Prabhakaran et al., 2022). Cultural alignment, which ensures that LLMs understand and reflect the values, norms, and nuances of the user groups interacting with them Masoud et al. (2023), is therefore essential for building inclusive, globally applicable AI systems. While recent work highlights the importance of dataset quality for model performance (Zhou et al., 2023; Alshahrani et al., 2023), cultural datasets remain understudied from a linguistic and structural perspective.

Dominant approaches to cultural alignment focus on value congruence, relying on survey-derived benchmarks (e.g., World Values Survey) or synthetic QA pairs, with limited attention to the linguistic structure of the underlying training data. As a result, these approaches often overlook how cultural

*Corresponding author: reem.masoud.22@ucl.ac.uk

nuances are encoded in the structural, semantic, and stylistic properties of the data, and whether such properties can be assessed independently of model training.

Cultural datasets also remain under-represented for many non-Western languages (Pawar et al., 2025), and it remains unclear how their linguistic properties influence model behavior. Prior work further suggests that different model architectures may respond differently to the same training data (Yauney et al., 2023; Zhang et al., 2025). We therefore ask whether cultural fine-tuning datasets exhibit consistent, measurable patterns in semantic content, lexical diversity, and stylistic variability that can be quantified prior to training, and whether such properties are predictive of, and actionable for, downstream cultural behavior across different model families. To answer these questions, we (i) characterize multilingual datasets using lightweight linguistic metrics and PCA, (ii) test associations between PCA dimensions and cultural benchmark performance across models, and (iii) probe actionability through controlled high/low/random subset fine-tuning.

Our contributions are threefold:

1. We present a *dataset-centric methodology* that quantifies linguistic, semantic, and structural properties of cultural datasets, reduces them via PCA, and links these components to downstream cultural alignment performance.
2. To our knowledge, we conduct the *first* cross-lingual empirical study examining Arabic, Chinese, and Japanese datasets across three LLM families, revealing that correlations between dataset properties and cultural performance vary substantially by language and model architecture.
3. We evaluate the *predictive* utility of PCA-derived linguistic dimensions for dataset assessment, including controlled subset-based interventions that test whether these signals remain informative under fixed training conditions.

Overall, our results indicate that pre-training linguistic properties of datasets can be informative for cultural alignment, but their effects are strongly model-dependent rather than universal. This motivates model-aware, dataset-centric strategies for multilingual cultural awareness in LLMs.

2 RELATED WORK

Prior work on cultural alignment has introduced benchmarks and datasets, often survey-based or value-oriented, that compare model outputs to human responses (e.g., AlKhamissi et al. 2024; Masoud et al. 2023). While effective for evaluation, these approaches typically do not analyze the linguistic or structural properties of the datasets themselves, nor how such properties relate to downstream cultural alignment. Complementary work has examined cultural bias and representational gaps in pretraining and post-training corpora (Naous et al., 2023; Alkhowaiter et al., 2025), as well as methods for curating culturally diverse datasets (Li et al., 2024). However, these studies do not identify which dataset-level linguistic properties are predictive of cultural alignment performance. In contrast, our work adopts a dataset-centric perspective, quantifying linguistic properties of multilingual cultural datasets and linking them to alignment outcomes across models, enabling assessment of dataset utility prior to fine-tuning. Additional adjacent related work is discussed in Appendix B.

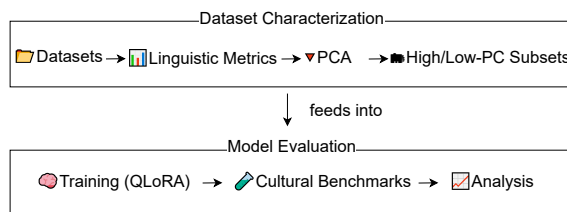


Figure 1: Methodology

3 METHODOLOGY

This work presents a dataset-centric methodology for analyzing how linguistic and structural properties of fine-tuning datasets relate to downstream cultural awareness in LLMs. Our approach is model-agnostic and operates entirely at the dataset level, allowing dataset characteristics to be quantified *ex-ante* (prior to training), rather than relying on *post-hoc* analysis of model behaviour. The pipeline consists of two stages (Figure 1): (1) dataset characterization (Steps 1–3), where linguistic metrics are computed, reduced via PCA, and used to define controlled dataset subsets; and (2) model evaluation (Steps 4–6), where LLMs are fine-tuned and assessed on cultural alignment benchmarks. The steps are detailed below. Each step is designed to test whether dataset-level linguistic structure captures systematic differences between datasets, correlates with downstream cultural performance, and supports systematic data selection that is more informative than random subset baselines.

Step 1: Linguistic Feature Extraction For each dataset, we compute lexical, semantic, and structural metrics (see Appendix C, Table 1) capturing diversity, richness, similarity, and cohesion. Metrics are computed at the dataset level, treating each dataset as a collection of samples whose aggregated properties reflect its linguistic profile. To ensure comparability across datasets, we randomly sample 1,000 examples per dataset; pilot checks with larger samples and multiple seeds show minimal variation. Because languages exhibit different morphological and statistical properties, features are normalized separately per language (zero mean, unit variance) to capture within-language variation rather than cross-language scale differences.

Step 2: PCA-Based Dimensionality Reduction To further compress the 10 specified linguistic metrics into a small number of interpretable dataset-level descriptors, we apply PCA separately for each language. PCA combines correlated linguistic metrics into a few continuous components, each assigning a score to every dataset that reflects how strongly it exhibits a particular combination of linguistic properties. We represent each dataset by its scores along the first three principal components (PC1–PC3), which capture most of the variance in the linguistic metrics. The resulting PCA scores serves as the basis for subsequent analyses that examine associations with downstream cultural performance and test whether these dimensions can guide dataset selection.

Step 3: Dataset Profiling and Subset Construction While PCA reveals descriptive structure, it does not establish whether these dimensions are useful for guiding training decisions. To test the practical relevance of PCA-derived signals, we construct controlled dataset subsets that approximate movement along each principal component at the sample level. Specifically, for each dataset and component, we assign each sample a score by projecting its feature vector onto the corresponding PCA direction, yielding a continuous measure of alignment. We then form equal-sized High-PC and Low-PC subsets from the upper and lower tails, along with a size-matched random baseline. All subsets are matched in size to isolate linguistic effects from dataset scale. While several metrics in Table 1 (e.g., Distinct-n, Self-BLEU, MTLT, silhouette) are inherently corpus-level, sample-level projections use only features defined per instance (e.g., length-normalized lexical statistics and embedding-based similarity scores). Dataset-level metrics are used solely to derive PCA directions, while sample-level scores rely on aligned approximations, ensuring consistency without applying corpus-level metrics directly to individual samples. By manipulating dataset composition in this way, we test whether dataset-level associations persist under controlled subset selection or disappear under re-sampling.

Step 4: Model-Agnostic Fine-Tuning Setup All experiments start from the same pre-trained LLM checkpoint. For each dataset, we fine-tune a separate model using: (i) the full dataset, (ii) each PCA-based subset from Step 3, and (iii) each random subset from Step 3. The training configuration is identical across all runs, and no model is fine-tuned sequentially on multiple datasets. Consequently, the resulting models differ only in the data used for fine-tuning, allowing us to isolate the effect of dataset composition on downstream performance while controlling for architecture and optimization settings.

Step 5: Cultural Evaluation Each fine-tuned model is evaluated on a set of cultural benchmarks covering three categories: cultural knowledge, cultural values, and cultural norms. These bench-

marks provide model-level performance scores across distinct aspects of cultural alignment, enabling systematic comparison across datasets, subsets, and model families.

Step 6: Linking Dataset Properties to Model Performance Finally, we analyze whether dataset-level linguistic structure is associated with downstream cultural behavior by computing correlations between (i) each dataset’s PCA coordinates (PC1–PC3) and (ii) the performance of the corresponding fine-tuned model on cultural benchmarks. These correlations are used to assess predictive association, not causality, and motivate the subset-based intervention experiments that follow.

4 EXPERIMENTS

4.1 SETUP

Languages and Datasets We conduct experiments across Arabic, Chinese, and Japanese, using between 9 and 13 post-training datasets per language. To ensure broad coverage of linguistic and cultural variation, the selected datasets span diverse sources and domains, including exams, social media, news, instruction-tuning corpora, and curated resources. Full dataset lists and detailed statistics (size, domain, and source type) are provided in Appendix E.

Models We fine-tune and evaluate Llama-3.2-3B-Instruct Grattafiori et al. (2024), Mistral-7B-Instruct-v0.3 Jiang et al. (2023), and DeepSeek-R1-Distill-Qwen-7B Guo et al. (2025).

Processing and Fine-Tuning Protocol For linguistic metric computation, each dataset was randomly sampled to 1,000 examples. Lexical metrics use a consistent language-specific preprocessing pipeline: Arabic text is Unicode- and Alef-normalized and tokenized with CAMEL Tools (Obeid et al., 2020); Chinese text is NFKC-normalized (Python Software Foundation, 2024) and segmented with jieba (with character-level fallback); Japanese text is segmented using MeCab (Python Software Foundation & Contributors, 2024). Tokenized sequences are rejoined into space-separated text for metric computation. For model training, datasets with <30k examples were used fully, larger datasets were capped at 30k samples, split into 80% train / 10% validation / 10% test. All additional fine-tuning datasets statistics and configurations, including sequence length, optimizer settings, learning rate schedule, batch size, training steps, and hardware setup, are documented in Appendix D and E.

Evaluation Framework Models are evaluated across three cultural alignment categories: Cultural Knowledge, Cultural Values, and Cultural Norms. Benchmarks include VSM13, World Values Survey (WVS), CulturalBench, Exams, and additional culture-related evaluation datasets relevant to each language (full list in Appendix F).

4.2 LINGUISTIC FEATURE ANALYSIS

As described in Section 3, we compute linguistic, semantic, and structural metrics for all datasets and apply PCA separately for each language to capture within-language variation. Here, we analyze the resulting PCA components to compare how Arabic, Japanese, and Chinese datasets differ along the major axes of linguistic variation.

4.2.1 PRINCIPAL COMPONENTS CAPTURE MEANINGFUL LINGUISTIC STRUCTURE

For each language, we apply PCA to the matrix of dataset-level linguistic metrics, where each dataset corresponds to one observation and each metric to one feature. The first three principal components (PCs) capture coherent groupings of linguistic features and explain most of the dataset-level. PC1 explains 41–53% of variance across languages, with PC2 and PC3 contributing a further 18–33% and 9–16%, respectively; additional components explain comparatively little variance (see Appendix G, Table 5). Category-level contribution plots (Figure 2) show that, within each language, individual PCs are dominated by subsets of linguistic metric categories rather than uniform mixtures. The dominant categories differ across components and languages, indicating that PCA captures language-specific combinations of linguistic properties rather than a single shared

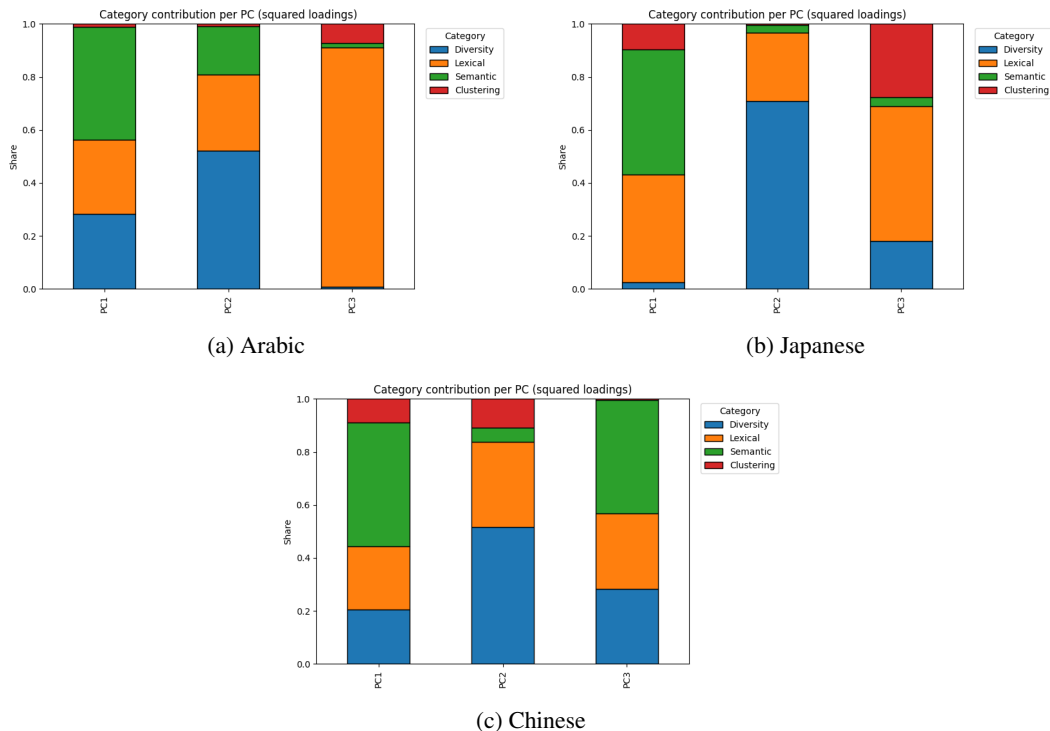


Figure 2: Category-level contributions of diversity, lexical, semantic, and clustering metrics to principal components (PC1–PC3) for Arabic, Japanese, and Chinese.

linguistic dimension. This motivates examining whether different components play distinct roles in downstream cultural alignment.

PC1: Semantic-Dominant Dimension PC1 is primarily driven by semantic coherence across all three languages, often co-occurring with related linguistic signals. It explains the largest share of variance in every language (approximately 41–53%), indicating that high-level semantic structure is the dominant source of variation between datasets.

PC2: Diversity + Lexical Dimension Across all languages, PC2 is primarily driven by diversity and lexical metrics. This component therefore reflects variation in wording and expression rather than differences in underlying semantic content. Datasets with high PC2 scores contain many different expressions and vocabulary forms, while low-scoring datasets reuse similar phrasing and wording patterns.

PC3: Secondary Lexical/Semantic Structure PC3 captures more localized, language-specific structure. In Arabic, it is driven primarily by lexical richness metrics (e.g., MTL, HDD); in Japanese, it combines lexical and clustering signals that separate more homogeneous from heterogeneous corpora; and in Chinese, it reflects a secondary semantic structure dominated by similarity metrics. Unlike PC1 and PC2, PC3 does not represent a common pattern across languages but instead captures language-dependent variation unique to each language.

4.2.2 DATASET INTERPRETATION THROUGH PCA DIMENSIONS

By projecting datasets into PCA space, datasets exhibit separations that reflect differences in their linguistic composition. Below, we describe how major dataset families are distributed along each component, as illustrated in Figures 3 – 5.

PC1: Datasets with Higher Semantic Density PC1 separates datasets according to how much explicit semantic reasoning they contain. Datasets built around knowledge-intensive QA datasets (e.g., such as NativQA, MBZUAI Arabic MMLU, OpenArabicQA, JCommonsense, COIG-PC, and LogiQA), consistently achieve high PC1 scores. These sources contain meaning-dense prompts and well-structured QA pairs, which increases semantic similarity metrics and drives strong PC1 loadings. By contrast, large news and encyclopedic corpora, including *Ultimate Arabic News*, *Aljazeera*, *Wikipedia-30k*, *Wikinews-5k*, *Wikihow*, and *Wiki*, tend to score low on PC1. Their formulaic reporting style, repetitive phrasing, and narrower topical variation produce lower semantic variability, placing them at the lower end of the PC1 axis.

PC2: Dataset Variation Along Surface-Level Properties PC2 does not align cleanly with any dataset source. Datasets from news, QA, instructional, and social-media origins appear throughout this dimension with no consistent pattern. This suggests that PC2 captures surface-level variation, such as differences in token distribution, topical breadth, or stylistic alternation, that cuts across genres rather than characterizing any particular dataset family.

PC3: Language-Dependent Lexical and Stylistic Variation Unlike PC1 and PC2, PC3 does not correspond to a consistent linguistic property across languages. Datasets with high PC3 scores vary substantially in type: conversational corpora (e.g., *Douban*), human-written instructions (e.g., *Ichikara*, *Aya*), curated knowledge resources, and some encyclopedic collections all appear along this axis. This pattern is supported by the loading analysis, where the dominant contributing metrics differ across languages. Together, these observations indicate that PC3 captures additional structured variation not explained by semantic coherence (PC1) or surface expression variation (PC2). Rather than representing a universal linguistic dimension, PC3 reflects dataset-specific characteristics such as annotation style, discourse format, or corpus organization.

Overall, the PCA projections provide a compact, language-specific view of how datasets differ in their linguistic composition. Rather than revealing shared dimensions across languages, the PCA space offers a descriptive organization of datasets within each language, which we use in subsequent analyses to examine relationships with cultural alignment performance.

4.3 DATASET LINGUISTIC STRUCTURE AND ITS IMPACT ON CULTURAL ALIGNMENT

This experiment examines whether dataset-level linguistic structure, as captured by PCA, is systematically associated with downstream cultural-alignment performance across languages and model families. For each dataset analyzed in Section 4.2, we fine-tune the same base LLM using a fixed training configuration, yielding one model per dataset. Each model is then evaluated on cultural benchmarks grouped into three categories: *Cultural Knowledge*, *Cultural Values*, and *Cultural Norms*. We compute Pearson correlations between (i) each dataset’s PCA coordinates (PC1–PC3) and (ii) the performance of the corresponding fine-tuned model on cultural alignment benchmarks, using the task-specific evaluation metrics defined in Appendix F. Importantly, these correlations are used to assess associative relationships, not to establish causality or actionability. Whether PCA-derived dimensions can reliably guide dataset selection is tested separately via controlled subset interventions in Section 4.4. Correlation heatmaps for Arabic, Japanese, and Chinese are provided in Appendix H (Figures 7–9). Given the large number of correlations evaluated, individual results should be interpreted with caution. We therefore focus on effect sizes, confidence intervals, and consistency across models. Full correlation tables are available in Appendix I.

4.3.1 PCA COMPONENTS EXHIBIT SYSTEMATIC BUT MODEL-DEPENDENT ASSOCIATIONS

Across all three languages, benchmark performance is often associated with at least one PCA component, indicating that dataset-internal linguistic structure is informative for cultural alignment. However, the identity and direction of the relevant component vary substantially across models, benchmarks, and languages, suggesting that no single PCA dimension acts as a universal predictor.

Arabic: Arabic benchmarks exhibit varied but structured associations with PCA components, with several strong and statistically stable relationships emerging across models. Cultural knowledge tasks show some of the clearest signals: CultureAtlas is strongly associated with PC1 for LLaMA ($r = 0.85$, 95% CI [0.75, 0.97]), while EXAMs aligns with PC2 for Mistral ($r = 0.82$, 95% CI

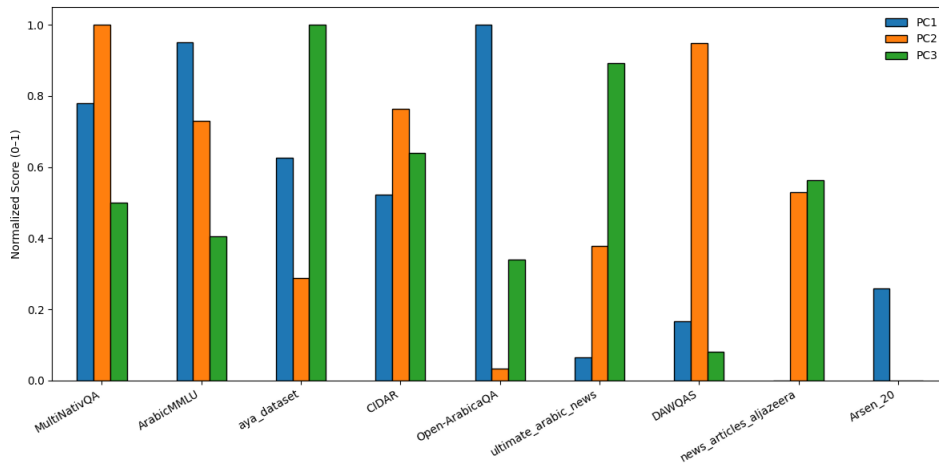


Figure 3: Normalized PCA scores for Arabic datasets projected onto the three principal components.

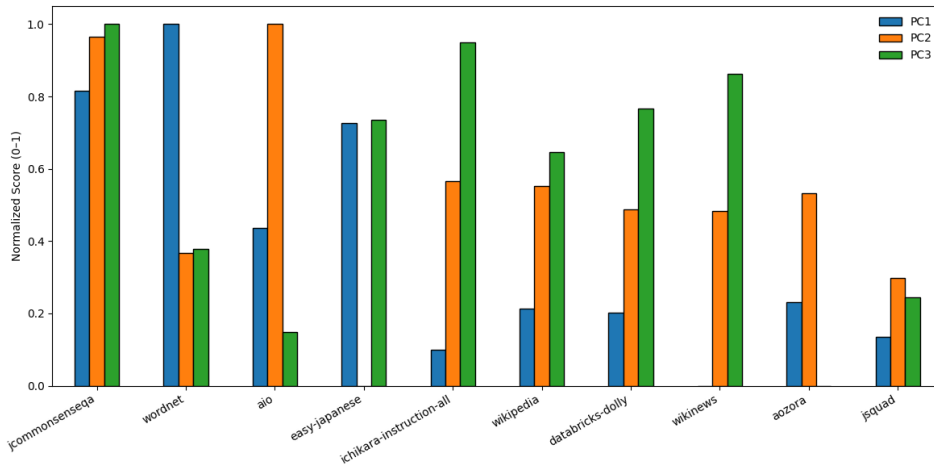


Figure 4: Normalized PCA scores for Japanese datasets projected onto the three principal components.

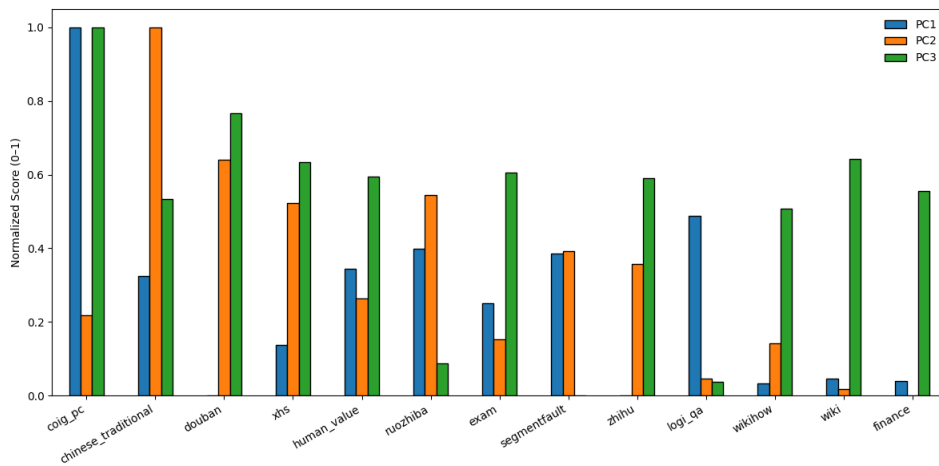


Figure 5: Normalized PCA scores for Chinese datasets projected onto the three principal components.

[0.49, 0.96]). Cultural values benchmarks similarly show strong but model-specific patterns, with WorldValuesBench most strongly associated with PC1 for Mistral ($r = -0.77$, 95% CI [-0.94, -0.45]) and ACVA aligned with PC2 for DeepSeek ($r = 0.78$, 95% CI [0.41, 0.97]). Normative benchmarks also vary across models: while CIDAR-EVAL is positively associated with PC2 for LLaMA ($r = 0.71$, 95% CI [0.14, 0.96]), and CIDAR-MCQ shows a stable negative association with PC3 for DeepSeek ($r = -0.55$, 95% CI [-0.88, -0.18]). Overall, Arabic shows the strongest and most numerous stable associations, although not all high point estimates remain reliable once uncertainty is considered.

Japanese: In Japanese, associations are fewer and less consistent than in Arabic, but several stable patterns still emerge. Cultural knowledge benchmarks often relate to PC1, though not uniformly across models: CulturalBench-Easy aligns with PC1 for Mistral ($r = -0.68$, 95% CI [-0.94, -0.04]) but with PC2 for LLaMA ($r = 0.50$, 95% CI [0.04, 0.83]), indicating that different models rely on different linguistic dimensions for the same task. Cultural values and similarity-based benchmarks also vary by architecture. For example, VSM13 correlates with PC2 for LLaMA ($r = 0.73$, 95% CI [0.17, 0.95]), while LLM-GLOBE benchmarks exhibit model-dependent associations with PC1 and PC2. Normative benchmarks show selective structure, such as JETHICS aligning with PC3 for Mistral ($r = 0.55$, 95% CI [0.24, 0.89]). Taken together, the Japanese results suggest that multiple linguistic dimensions contribute to cultural performance, but their relevance is weaker and more model-contingent than in Arabic.

Chinese: Chinese benchmarks display distributed associations across PCA components, with different tasks aligning with different PCs depending on the model. Several stable relationships emerge, but the most relevant component again varies by architecture. For LLaMA, CulturalBench-Hard is positively associated with PC2 ($r = 0.40$, 95% CI [0.08, 0.77]), while CMoralEval is positively associated with PC3 ($r = 0.62$, 95% CI [0.12, 0.86]). For Mistral, PC1 shows strong positive associations with LLM-GLOBE-CLOSED ($r = 0.77$, 95% CI [0.28, 0.93]) and VSM13 ($r = 0.51$, 95% CI [0.06, 0.79]), whereas PC3 is negatively associated with ceval ($r = -0.45$, 95% CI [-0.77, -0.12]) and CultureAtlas ($r = -0.34$, 95% CI [-0.77, -0.02]). For DeepSeek, WorldValuesBench is positively associated with both PC1 ($r = 0.69$, 95% CI [0.11, 0.93]) and PC2 ($r = 0.43$, 95% CI [0.03, 0.83]). These results indicate that Chinese cultural performance is shaped by multiple linguistic dimensions rather than a single dominant component, with different components (including PC3) appearing selectively informative for different tasks.

Notably, the same benchmark can align with different PCA components depending on model architecture. For example, VSM13 aligns with PC2 for Japanese LLaMA but with PC1 for Chinese Mistral, while CultureAtlas aligns with PC1 for Arabic LLaMA but with PC3 for Chinese Mistral. More broadly, Arabic exhibits the clearest stable associations, whereas Japanese and Chinese show sparser and more selective effects. Many other correlations across all three languages have wide confidence intervals that include zero, indicating that these relationships should be interpreted cautiously given the limited number of datasets. Overall, the results show that PCA-derived linguistic dimensions provide informative but non-universal signals for cultural alignment, motivating the controlled intervention experiments that follow.

4.4 SUBSET VALIDATION

Correlation alone does not establish whether dataset-level linguistic dimensions are actionable for training. We therefore perform a controlled subset intervention, directly manipulating dataset composition while holding the model, optimization procedure, and subset size fixed. We focus on Arabic because it has the largest and most diverse collection of fine-tuning datasets in our study, spanning a wide range of sources and domains. We identify the five datasets with the largest absolute scores along each principal component (PC1–PC3). Within each dataset, we assign sentences a proxy PC score by projecting their linguistic feature vectors onto the corresponding PCA direction. Using this score, we construct three size-matched subsets ($\approx 2k$ examples): a *High-PC* subset (upper tail), a *Low-PC* subset (lower tail), and a *Random* subset. Each subset independently fine-tunes the same base models (*LLaMA*, *Mistral*, *DeepSeek*) under identical hyperparameters, followed by evaluation on cultural benchmarks. PC1 and PC2 interventions test whether emphasizing or suppressing semantic structure and distributional diversity improves alignment, while PC3 probes lexical and stylistic variation. Random subsets serve as a strong baseline that controls for subset size while implicitly reducing redundancy, allowing gains or failures to be attributed to composition rather than

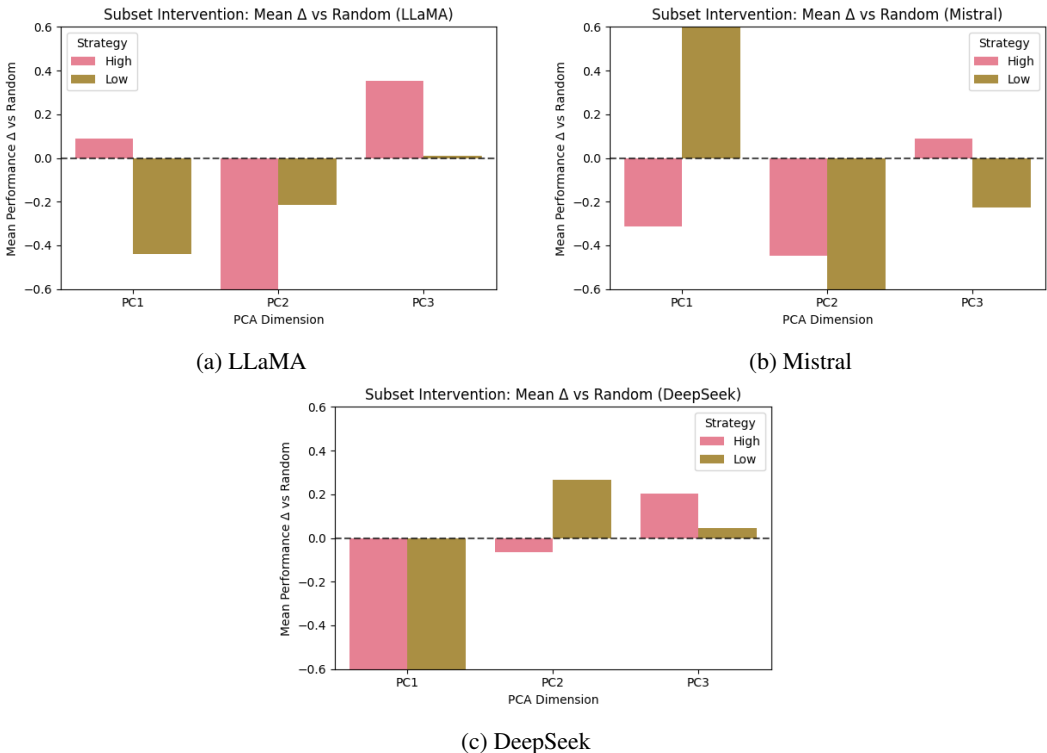


Figure 6: Mean performance difference ($\Delta = \text{subset} - \text{random}$) for High-PC and Low-PC subsets of equal size across PC1–PC3 and three model families, averaged over all base datasets and evaluation metrics. Zero indicates parity with random selection; positive values denote improvement and negative values degradation. PC1-PC2 rarely outperform random sampling, while PC3 shows more consistent, but model-dependent, gains.

scale. Figure 6 reports the mean performance difference (Δ) relative to Random for each strategy and PC, averaged across all evaluated datasets and metrics; positive values indicate improvement over random sampling.

LLaMA: LLaMA exhibits limited but structured sensitivity to PCA-guided subset selection. Along PC1, High-PC subsets provide small positive gains over Random, while Low-PC subsets consistently underperform. PC2 interventions are uniformly harmful, with both High-PC and Low-PC subsets degrading performance relative to Random. In contrast, PC3 shows the clearest positive signal: High-PC3 subsets yield consistent improvements, while Low-PC3 remains near-neutral. Overall, LLaMA benefits selectively from lexical-level variation (PC3), while semantic and diversity-driven extremes offer limited or negative utility.

Mistral: Mistral displays a strongly directional response to PCA-guided subset selection. Along PC1, *Low-PC* subsets substantially outperform both High-PC and Random, while High-PC1 leads to clear degradation. This indicates that fine-tuning on semantically extreme subsets—particularly those emphasizing high semantic density—induces harmful distributional shifts for Mistral. PC2 interventions are consistently detrimental: both High-PC2 and Low-PC2 subsets produce large performance drops relative to Random, suggesting that aggressive manipulation of diversity-related properties is poorly tolerated. In contrast, PC3 emerges as the most reliable dimension. High-PC3 subsets consistently outperform Random, while Low-PC3 underperform, indicating that increased lexical and stylistic variation provides a stable and beneficial training signal for Mistral.

DeepSeek: DeepSeek follows a distinct pattern. For PC1, both High-PC and Low-PC subsets substantially underperform Random, suggesting that semantic extremes are broadly misaligned with the model. Along PC2, Low-PC subsets outperform High-PC and Random, while High-PC2 leads to degradation. PC3 produces consistent gains: High-PC3 yields the strongest improvements, while

Low-PC3 provides smaller but still positive effects. Thus, for DeepSeek, PC1 acts as a negative selection signal, PC2 favors reduced diversity, and PC3 provides a robust positive intervention.

Cross-model Comparison: Across models, PCA-guided subset selection exhibits clear *model-dependent* effects. PC3 (lexical and stylistic variation) provides the most consistent and transferable signal, improving performance for all three models. A plausible explanation is that PC3 emphasizes surface-level and stylistic diversity without strongly altering the underlying semantic distribution, allowing models to benefit from increased lexical coverage while remaining close to their pretraining regime. In contrast, PC1 and PC2 are more fragile: subsets emphasizing semantic density or distributional extremes can induce larger shifts in the training distribution, which different architectures tolerate unevenly. As a result, these interventions frequently degrade performance and show limited cross-model agreement. Importantly, High- and Low-PC subsets are not symmetric—the direction of intervention along a PCA axis matters as much as the axis itself.

Overall, these results demonstrate that PCA-derived linguistic dimensions can inform subset construction, but only when interpreted in a model- and direction-aware manner. PC3 offers the most stable intervention signal, while aggressive manipulation along PC1 and PC2 often harms alignment. Rather than indicating that certain linguistic properties are universally beneficial or harmful, the findings highlight that increasing or decreasing the same property can produce qualitatively different outcomes depending on the model.

5 CONCLUSION

This work examines how *dataset-level linguistic properties* relate to cultural alignment in large language models. Across 160 fine-tuned models spanning three languages and three model families, we combine language-specific PCA, correlation analysis, and controlled subset interventions to characterize linguistic structure in cultural datasets and assess its downstream impact. Knowledge-intensive QA datasets align with semantic coherence (PC1), news and encyclopedic corpora cluster in low-semantic regions, and human-authored or conversational data exhibit higher lexical richness (PC3).

Subset-based interventions show that emphasizing semantic density (PC1) or surface diversity (PC2) often leads to unstable or negative effects, whereas the lexical–stylistic dimension (PC3) provides the most consistently positive signal in our subset experiments in Arabic. Importantly, all effects are strongly model-dependent, indicating that linguistic signals must be interpreted in an architecture-aware manner. While prior work has shown that data diversity can improve model performance and robustness in supervised fine-tuning settings Chen et al. (2024); Pang et al. (2024); Zhou et al. (2023), our results suggest that for cultural awareness tasks, increased semantic density or surface-level diversity is not uniformly beneficial. Instead, dataset composition and model-specific interactions play a more central role.

Practical Implication: For culturally-aware fine-tuning, we recommend: (i) prioritizing dataset composition, as PCA-guided subset selection yields measurable performance differences under fixed training conditions; (ii) treating semantic and diversity signals (PC1, PC2) with caution due to their instability across architectures; and (iii) favoring high-PC3 subsets, which provide the most consistently positive signal among the examined components.

Overall, PCA-derived linguistic dimensions are not universal quality metrics, but practical tools for probing and shaping dataset composition under controlled conditions.

ACKNOWLEDGMENTS

Reem I. Masoud acknowledges support from a scholarship provided by the Department of Electrical and Computer Engineering at King Abdulaziz University.

REFERENCES

llm-book/aio: Ai king competition qa dataset. HuggingFace Dataset, 2023. URL <https://huggingface.co/datasets/llm-book/aio>.

- Abdelrahman Abdallah, Mahmoud Kasem, Mahmoud Abdalla, Mohamed Mahmoud, Mohamed Elkasaby, Yasser Elbendary, and Adam Jatowt. Arabicaqa: A comprehensive dataset for arabic question answering, 2024.
- Ahmed Hashim Al-Dulaimi. *Ultimate Arabic News Dataset*. 05 2022. doi: 10.17632/jz56k5wxz7.1.
- Badr AlKhamissi, Muhammad Elnokrashy, Mai AlKhamissi, and Mona Diab. Investigating cultural alignment of large language models. *arXiv preprint arXiv:2402.13231*, 2024.
- Mohammed Alkhowaiter, Norah Alshahrani, Saied Alshahrani, Reem I Masoud, Alaa Alzahrani, Deema Alnuhait, Emad A Alghamdi, and Khalid Almubarak. Mind the gap: A review of arabic post-training datasets and their limitations. *arXiv preprint arXiv:2507.14688*, 2025.
- Saied Alshahrani, Norah Alshahrani, Soumyabrata Dey, and Jeanna Matthews. Performance implications of using unrepresentative corpora in arabic natural language processing. In *Proceedings of ArabicNLP 2023*, pp. 218–231, 2023.
- Zaid Alyafeai, Khalid Almubarak, Ahmed Ashraf, Deema Alnuhait, Saied Alshahrani, Gubran AQ Abdulrahman, Gamil Ahmed, Qais Gawah, Zead Saleh, Mustafa Ghaleb, et al. Cidar: Culturally relevant instruction dataset for arabic. *arXiv preprint arXiv:2402.03177*, 2024.
- ArbML. Al jazeera news articles. https://huggingface.co/datasets/arbml/news_articles_aljazeera, 2023a. HuggingFace Dataset.
- ArbML. Arsen-20. https://huggingface.co/datasets/arbml/Arsen_20, 2023b. HuggingFace Dataset.
- Yuelin Bai, Xinrun Du, Yiming Liang, Yonggang Jin, Ziqiang Liu, Junting Zhou, Tianyu Zheng, Xincheng Zhang, Nuo Ma, Zekun Wang, et al. Coig-cqia: Quality is all you need for chinese instruction fine-tuning, 2024.
- Francis Bond and Takayuki Kuribayashi. The Japanese Wordnet 2.0. In German Rigau, Francis Bond, and Alexandre Rademaker (eds.), *Proceedings of the 12th Global Wordnet Conference*, pp. 179–186, University of the Basque Country, Donostia - San Sebastian, Basque Country, January 2023. Global Wordnet Association. doi: 10.18653/v1/2023.gwc-1.22. URL <https://aclanthology.org/2023.gwc-1.22/>.
- Hao Chen, Abdul Waheed, Xiang Li, Yidong Wang, Jindong Wang, Bhiksha Raj, and Marah I Abidin. On the diversity of synthetic data and its impact on training large language models. *arXiv preprint arXiv:2410.15226*, 2024.
- Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin, Chan Young Park, Shuyue Stella Li, Sahithya Ravi, Mehar Bhatia, Maria Antoniak, Yulia Tsvetkov, Vered Shwartz, et al. Culturalbench: A robust, diverse, and challenging cultural benchmark by human-ai culturalteaming. *arXiv preprint arXiv:2410.02677*, 2024.
- Michael A Covington and Joe D McFall. Cutting the gordian knot: The moving-average type–token ratio (mattr). *Journal of quantitative linguistics*, 17(2):94–100, 2010.
- Iris Dominguez-Catena, Daniel Paternain, and Mikel Galar. Dsap: Analyzing bias through demographic comparison of datasets. *Information Fusion*, 115:102760, 2025.
- Yi Fung, Ruining Zhao, Jae Doo, Chenkai Sun, and Heng Ji. Massively multi-cultural knowledge acquisition & lm benchmarking. *arXiv preprint arXiv:2402.09369*, 2024.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- Md Arid Hasan, Maram Hasanain, Fatema Ahmad, Sahinur Rahman Laskar, Sunaya Upadhyay, Vrunda N Sukhadia, Mucahid Kutlu, Shammur Absar Chowdhury, and Firoj Alam. Nativqa: Multilingual culturally-aligned natural query for llms. *arXiv preprint arXiv:2407.09823*, 2024.
- Masanori HIRANO, Masahiro SUZUKI, and Hiroki SAKAJI. llm-japanese-dataset v0: Construction of Japanese Chat Dataset for Large Language Models and its Methodology, 2023. URL <https://github.com/masanorihirano/llm-japanese-dataset>.
- Geert Hofstede, Gert Jan Hofstede, and Michael Minkov. Vsm 2013 — values survey module 2013. <https://geerthofstede.com/research-and-vsm/vsm-2013/>, 2013. Accessed: 2025-01-04.
- Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen Alharthi, Bang An, Juncai He, et al. Acegpt, localizing large language models in arabic. *arXiv preprint arXiv:2309.12053*, 2023a.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *arXiv preprint arXiv:2305.08322*, 2023b.
- Walaa Saber Ismail and Masun Nabhan Homsy. Dawqas: A dataset for arabic why question answering system. *Procedia computer science*, 142:123–131, 2018.
- Izumi Lab. wikipedia-ja-20230720, 2023. URL <https://huggingface.co/datasets/izumi-lab/wikipedia-ja-20230720>.
- AQ Jiang, A Sablayrolles, A Mensch, C Bamford, DS Chaplot, D de Las Casas, F Bressand, G Lengyel, G Lample, L Saulnier, et al. Mistral 7b. corr, abs/2310.06825, 2023. doi: 10.48550. *arXiv preprint ARXIV:2310.06825*, 10, 2023.
- Webdell Johnson. Studies in language behavior: A program of research. *Psychological Monographs*, 56(2):1–15, 1944. URL https://pure.mpg.de/rest/items/item_2350946/component/file_2562008/content.
- Elise Karinshak, Amanda Hu, Kewen Kong, Vishwanatha Rao, Jingren Wang, Jindong Wang, and Yi Zeng. Llm-globe: A benchmark evaluating the cultural values embedded in llm output. *arXiv preprint arXiv:2411.06032*, 2024.
- Akihiro Katsuta and Kazuhide Yamamoto. Crowdsourced corpus of sentence simplification with core vocabulary. In *Proceedings of the 11th International Conference on Language Resources and Evaluation (LREC)*, pp. 461–466, 2018.
- ”Fajri Koto, Haonan Li, Sara Shatanawi, Jad Doughman, Abdelrahman Boda Sadallah, Aisha Al-raeesi, Khalid Almubarak, Zaid Alyafeai, Neha Sengupta, Shady Shehata, Nizar Habash, Preslav Nakov, and Timothy Baldwin”. Arabicmmlu: Assessing massive multitask language understanding in arabic. In *Findings of the Association for Computational Linguistics: ACL 2024*, 2024.
- Kunishou. Databricks dolly 15k (japanese). <https://huggingface.co/datasets/kunishou/databricks-dolly-15k-ja>, 2023. HuggingFace Dataset.
- levellevel. Aozoratxt: Aozora bunko text corpus. <https://github.com/levellevel/AozoraTxt>, 2023. GitHub Repository.
- Cheng Li, Damien Teney, Linyi Yang, Qingsong Wen, Xing Xie, and Jindong Wang. Culturepark: Boosting cross-cultural understanding in large language models. *Advances in Neural Information Processing Systems*, 37:65183–65216, 2024.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and William B Dolan. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pp. 110–119, 2016.

- Reem I Masoud, Ziquan Liu, Martin Ferianc, Philip Treleaven, and Miguel Rodrigues. Cultural alignment in large language models: An explanatory analysis based on hofstede’s cultural dimensions. *arXiv preprint arXiv:2309.12342*, 2023.
- Philip M McCarthy and Scott Jarvis. Mtd, vocd-d, and hd-d: A validation study of sophisticated approaches to lexical diversity assessment. *Behavior research methods*, 42(2):381–392, 2010.
- James B McQueen. Some methods of classification and analysis of multivariate observations. In *Proc. of 5th Berkeley Symposium on Math. Stat. and Prob.*, pp. 281–297, 1967.
- Basel Mousi, Nadir Durrani, Fatema Ahmad, Md Arid Hasan, Maram Hasanain, Tameem Kabbani, Fahim Dalvi, Shammur Absar Chowdhury, and Firoj Alam. Aradice: Benchmarks for dialectal and cultural capabilities in llms. *arXiv preprint arXiv:2409.11404*, 2024.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. Having beer after prayer? measuring cultural bias in large language models. *arXiv preprint arXiv:2305.14456*, 2023.
- Ossama Obeid, Nasser Zalmout, Salam Khalifa, Dima Taji, Mai Oudah, Bashar Alhafni, Go Inoue, Fadhil Eryani, Alexander Erdmann, and Nizar Habash. CAMEL tools: An open source python toolkit for Arabic natural language processing. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pp. 7022–7032, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://www.aclweb.org/anthology/2020.lrec-1.868>.
- Jinlong Pang, Jiaheng Wei, Ankit Parag Shah, Zhaowei Zhu, Yaxuan Wang, Chen Qian, Yang Liu, Yujia Bao, and Wei Wei. Improving data efficiency via curating llm-driven rating systems. *arXiv preprint arXiv:2410.10877*, 2024.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin (eds.), *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp. 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040/>.
- Siddhesh Pawar, Junyeong Park, Jiho Jin, Arnav Arora, Junho Myung, Srishti Yadav, Faiz Ghifari Haznitrama, Inhwa Song, Alice Oh, and Isabelle Augenstein. Survey of cultural awareness in language models: Text and beyond. *Computational Linguistics*, pp. 1–96, 2025.
- Vinodkumar Prabhakaran, Rida Qadri, and Ben Hutchinson. Cultural incongruencies in artificial intelligence. *arXiv preprint arXiv:2211.13069*, 2022.
- Python Software Foundation. unicodedata — unicode database, 2024. URL <https://docs.python.org/3/library/unicodedata.html>.
- Python Software Foundation and MeCab Contributors. mecab-python3: Python wrapper for mecab, 2024. URL <https://pypi.org/project/mecab-python3/>.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- Abdelrahman Sadallah, Junior Cedric Tonga, Khalid Almubarak, Saeed Almheiri, Farah Atif, Cahtrine Qwaider, Karima Kadaoui, Sara Shatnawi, Yaser Alesh, and Fajri Koto. Commonsense reasoning in arab culture, 2025.
- Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval. *Information processing & management*, 24(5):513–523, 1988.
- Gerard Salton, Anita Wong, and Chung-Shu Yang. A vector space model for automatic indexing. *Communications of the ACM*, 18(11):613–620, 1975.

- Satoshi Sekine, Maya Ando, Michiko Goto, Kumi Suzuki, Daisuke Kawahara, Naoya Inoue, and Kentaro Inui. Creation of japanese instruction data for llms: ichikara-instruction. *Proceedings of the 30th Annual Meeting of the Association for Natural Language Processing*, pp. 1508–1513, 2024. URL <https://huggingface.co/datasets/msfm/ichikara-instruction-all>. (in Japanese).
- Neha Sengupta, Sunil Kumar Sahu, Bokang Jia, Satheesh Katipomu, Haonan Li, Fajri Koto, William Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming Chen, et al. Jais and jais-chat: Arabic-centric foundation and instruction-tuned open generative large language models. *arXiv preprint arXiv:2308.16149*, 2023. URL <https://huggingface.co/datasets/FreedomIntelligence/EXAMs>.
- Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F. Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura OMahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Souza Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergün, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Minh Chien, Sebastian Ruder, Surya Guthikonda, Emad A. Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. Aya dataset: An open-access collection for multilingual instruction tuning, 2024.
- ByungHoon So, Kyuhong Byun, Kyungwon Kang, and Seongjin Cho. Jaquad: Japanese question answering dataset for machine reading comprehension. *arXiv preprint arXiv:2202.01764*, 2022.
- Masashi Takeshita and Kenji Araki. Jcommonsensemorality. In *Proceedings of the 29th Annual Meeting of the Association for Natural Language Processing (in Japanese)*, pp. 357–362, 2023.
- Masashi Takeshita and Rafal Rzepka. Jethics: Japanese ethics understanding evaluation dataset, 2025. URL <https://arxiv.org/abs/2506.16187>.
- Wikimedia Foundation. Wikimedia project dumps. <https://dumps.wikimedia.org/>, 2024. Accessed 2024.
- Wikinews. Japanese wikinews. <https://ja.wikinews.org/wiki/>, 2024. Accessed 2024.
- Trong Wu. An accurate computation of the hypergeometric distribution function. *ACM Transactions on Mathematical Software (TOMS)*, 19(1):33–43, 1993.
- Gregory Yauney, Emily Reif, and David Mimno. Data similarity is not enough to explain language model performance. *arXiv preprint arXiv:2311.09006*, 2023.
- Linhao Yu, Yongqi Leng, Yufei Huang, Shang Wu, Haixin Liu, Xinmeng Ji, Jiahui Zhao, Jinwang Song, Tingting Cui, Xiaoqing Cheng, Tao Liu, and Deyi Xiong. Cmoraleval: A moral evaluation benchmark for chinese large language models, 2024. URL <https://arxiv.org/abs/2408.09819>.
- Xinlu Zhang, Zhiyu Zoey Chen, Xi Ye, Xianjun Yang, Lichang Chen, William Yang Wang, and Linda Ruth Petzold. Unveiling the impact of coding data instruction fine-tuning on large language models reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 25949–25957, 2025.
- Wenlong Zhao, Debanjan Mondal, Niket Tandon, Danica Dillion, Kurt Gray, and Yuling Gu. World-valuesbench: A large-scale benchmark dataset for multi-cultural value awareness of language models. *arXiv preprint arXiv:2404.16308*, 2024.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021, 2023.

A LIMITATIONS AND FUTURE WORK

While the correlations observed in our analysis frequently exceed $|0.40|$ and in many cases surpass $|0.70|$, indicating that linguistic structure is meaningfully associated with cultural-alignment performance, the number of finetuning datasets available per language (9–13) imposes constraints on statistical generality. This limited sample size also contributes to wide confidence intervals in several correlations, reducing statistical power and increasing uncertainty in some observed associations. Our findings should therefore be interpreted as *descriptive patterns* rather than definitive causal and statistical association claims. In addition, PCA-based analyses were conducted using a single random seed, which may introduce minor variability in the resulting principal components, although the same transformation was consistently applied across all evaluated subsets. While we do not compute bootstrap confidence intervals for PCA loadings, pilot experiments with repeated sampling showed consistent component structure and loading patterns across runs. Moreover, cultural-alignment performance was evaluated using pass@1 accuracy, reflecting a standardized single-response setting applied consistently across all models and datasets. While alternative decoding or aggregation strategies may yield different absolute scores, our analysis focuses on relative trends, which are less sensitive to this choice. Importantly, several trends replicate across three model families (LLaMA, Mistral, DeepSeek) and three languages (Arabic, Japanese, Chinese), indicating that the observed relationships reflect stable tendencies rather than sampling noise. While subset interventions were conducted only for Arabic due to computational constraints, the consistency of correlation patterns across languages suggests that similar effects may extend to Japanese and Chinese, which we leave for future work. Extending these experiments to additional datasets and model architectures would further strengthen generality, but was beyond the available computational budget.

B ADDITIONAL RELATED WORK

Beyond the cultural alignment benchmarks discussed in the main paper, several adjacent lines of work study dataset properties from complementary perspectives. Prior analyses document cultural bias and representational gaps in large-scale corpora, showing that widely used sources such as Wikipedia are strongly Western-centric (Naous et al., 2023), and that Arabic post-training resources suffer from scarcity and imbalance (Alkhowaiter et al., 2025). Other work explores dataset curation strategies for cultural diversity, such as synthetic dialogue generation in CulturePark (Li et al., 2024), though these efforts primarily target cultural judgments or moderation rather than identifying linguistic dataset properties associated with alignment.

Outside the cultural alignment literature, dataset-centric analyses such as DSAP (Dominguez-Catena et al., 2025) examine demographic similarity across corpora, and broader surveys document systematic cultural misalignment in LLMs (Pawar et al., 2025). Related work on dataset quality and diversity investigates how data curation, filtering, and diversity-oriented selection strategies influence downstream model performance (Chen et al., 2024; Zhou et al., 2023; Pang et al., 2024). While informative, these studies do not directly connect dataset-level linguistic structure to downstream cultural performance, nor do they evaluate such properties using culture-specific benchmarks. Our work complements these efforts by focusing specifically on how measurable linguistic properties of fine-tuning datasets relate to cultural alignment outcomes across models and languages.

C LINGUISTIC METRICS

The linguistic metrics used to characterize dataset composition are summarized in Table 1. Structural cohesion is measured using silhouette scores computed over sentence embeddings. We encode samples using language-specific sentence embedding models (Reimers & Gurevych, 2019): using language-specific sentence embedding models (Reimers & Gurevych, 2019): `bert-large-arabertv02`¹ for Arabic, `chinese-macbert-base`² for Chinese, and `sentence-bert-base-ja-mean-tokens`³ for Japanese, and apply k-means clustering with the number of clusters selected by maximizing the silhouette score over $k \in [2, 10]$.

¹<https://huggingface.co/aubmindlab/bert-large-arabertv02>

²<https://huggingface.co/hf1/chinese-macbert-base>

³<https://huggingface.co/sonois/sentence-bert-base-ja-mean-tokens>

Table 1: Linguistic and structural metrics used to characterize dataset composition. Lexical Diversity metrics include Distinct-1/2 (Li et al., 2016) and Self-BLEU (Papineni et al., 2002). Vocabulary Richness is measured using TTR (Johnson, 1944), MATTR (Covington & McFall, 2010), and HDD Wu (1993); McCarthy & Jarvis (2010), MTLT (McCarthy & Jarvis, 2010)). Semantic Consistency is computed using cosine similarity and TF-IDF representations (Salton et al., 1975; Salton & Buckley, 1988; Reimers & Gurevych, 2019). Structural Cohesion is assessed using K-means (McQueen, 1967) and the silhouette score (Rousseeuw, 1987).

Property	Metrics	Interpretation
Lexical Diversity	Distinct-1/2, Self-BLEU	Repetition vs variation across samples
Vocabulary Richness	TTR, MATTR, HDD, MTLT	Breadth and stability of vocabulary usage
Semantic Consistency	Cosine, TF-IDF	Similarity of meaning within dataset
Structural Cohesion	Silhouette score	Presence of coherent topical clusters

D TRAINING CONFIGURATION

Table 2 summarizes the fine-tuning configuration used across all model families and datasets. These settings were kept consistent to ensure comparability across languages and experimental conditions.

Table 2: Training configuration used for all fine-tuning experiments.

Setting	Value
Model families	LLaMA, Mistral, DeepSeek
Fine-tuning method	QLoRA (4-bit quantization)
Batch size	8
Gradient accumulation steps	8
Learning rate	2×10^{-5}
Optimizer	AdamW
Warmup ratio	0.03
Max sequence length	2048
Training epochs	3
LoRA rank (r)	64
LoRA α	16
LoRA dropout	0.05
Hardware	2xA100 80GB

E TRAINING DATASET STATISTIC

We summarize the training datasets used in our study in Table 3. The datasets span three languages (Arabic, Japanese, and Chinese) and cover diverse domains. They are collected from a variety of sources such as benchmark datasets, web data, media articles, and social media platforms, ensuring broad coverage of linguistic styles and cultural contexts.

F BENCHMARKS BY LANGUAGE

Table 4 lists the cultural and value-alignment benchmarks used for evaluation in each language. These benchmarks span multiple dimensions of cultural knowledge, norms, and moral reasoning.

G EXPLAINED VARIANCE RATIOS OF LANGUAGE-SPECIFIC PCA

Table 5 reports the explained variance ratios of the top three principal components for each language. Across all languages, PC1 accounts for the largest share of variance (41–52%), with PC2 and PC3 capturing additional complementary structure. Together, the first three components explain the majority of dataset-level variance, supporting their use as compact descriptors of linguistic structure.

Table 3: Summary statistics of datasets, including size, language, domain, and source type.

Dataset	Rows	Lang	Domain	Source Type
MultiNativQA (Hasan et al., 2024)	~5k	AR	Question Answering / Cultural	Crowdsourced
ArabicMMLU (Koto et al., 2024)	~15k	AR	Academic / Knowledge	Benchmark Dataset
Aya (Singh et al., 2024)*	~5k	AR	Instruction / General	Synthetic and Curated
CIDAR (Alyafeai et al., 2024)	10k	AR	Cultural Alignment	Mixed Sources
Open-ArabicaQA (Abdallah et al., 2024)	~62k	AR	Question Answering	Web Data
The Ultimate Arabic News (Al-Dulaimi, 2022)	196k	AR	News	Media Articles
DAWQAS (Ismail & Homsy, 2018)	~3k	AR	Question Answering	Web Data
Al Jazeera News Articles (ArbML, 2023a)	6k	AR	News	Media Articles
Arsen-20 (ArbML, 2023b)	20k	AR	Sentiment / Social Media	Social Media Data
Ichikara Instruction All (Sekine et al., 2024)	10k	JP	Instruction Tuning	Crowdsourced
Databricks-Dolly-JA (Kunishou, 2023)	15k	JP	Instruction Tuning	Machine-translated (from Dolly-15k)
JaQuAD (So et al., 2022)	~32k	JP	Question Answering	Crowdsourced (Wikipedia-based)
WordNet 2.0 (Bond & Kuribayashi, 2023)	158k	JP	Lexical Resource	Curated Lexical Database
Jcommonsense (Takeshita & Araki, 2023)	9k	JP	Ethics / Commonsense Reasoning	Crowdsourced (ConceptNet-based)
AIO Instruction Dataset (aio, 2023)	23k	JP	Instruction Tuning	Mixed Sources (competition, Wikipedia)
Easy-Japanese (Katsuta & Yamamoto, 2018)	34k	JP	Text Simplification	Crowdsourced (textbook-based, curated)
Wiki JP (Izumi Lab, 2023; Wikimedia Foundation, 2024)	30k	JP	General Knowledge	Encyclopedia Data
Wikinews JP (HIRANO et al., 2023; Wikinews, 2024)	5k	JP	News	Media Articles
Aozora Bunko (levelevel, 2023)	30k	JP	Literature	Books and Literary Texts
COIG-PC (Bai et al., 2024)	3k	ZH	Instruction / General	Mixed Sources
Chinese Traditional (Bai et al., 2024)	1k	ZH	General Domain	Web Data
Douban (Bai et al., 2024)	3k	ZH	Social Media / Reviews	Social Media Data
Xiaohongshu (XHS) (Bai et al., 2024)	~2k	ZH	Social Media / Lifestyle	Social Media Data
Human Value (Bai et al., 2024)	1k	ZH	Cultural Values	Mixed Sources
RuozhiBa (Bai et al., 2024)	~0.25k	ZH	Social Media / Humor	Social Media Data
Exam (Bai et al., 2024)	5k	ZH	Academic	Benchmark Dataset
SegmentFault (Bai et al., 2024)	0.5k	ZH	Technical Discussions	Online Forum
Zhihu (Bai et al., 2024)	~6k	ZH	Question Answering / Discussion	Online Forum
LogiQA (Bai et al., 2024)	~0.5k	ZH	Logical Reasoning	Benchmark Dataset
Wikihow (Bai et al., 2024)	~1.5k	ZH	Instructional Content	Web Data
Wiki (Bai et al., 2024)	~11k	ZH	General Knowledge	Encyclopedia Data
Finance (Bai et al., 2024)	~11k	ZH	Financial Domain	Web Data

* Arabic subset obtained via language filtering.

H ADDITIONAL CORRELATION HEATMAPS

This appendix presents the full correlation heatmaps between dataset PCA components (PC1–PC3) and downstream cultural benchmark performance, reported separately for each language and model. Each cell in Figs. 7, 8, and 9 shows the correlation coefficient between a PCA component score (columns) and a benchmark score (rows).

Table 4: Cultural alignment benchmarks by category and language.

Language	Cultural Knowledge	Cultural Values	Cultural Norms
Arabic	CulturalBench (Easy, Hard) (Chiu et al., 2024), CultureAtlas (Fung et al., 2024), ArabCulture (Sadallah et al., 2025), EXAMs-AR (Sengupta et al., 2023)	VSM13 (Arabic) (Hofstede et al., 2013), WorldValues-Bench (Zhao et al., 2024), ACVA-Arabic (Huang et al., 2023a)	CIDAR-MCQ, CIDAR-EVAL (Alyafeai et al., 2024), AraDiCE (Mousi et al., 2024), LLM-Globe (Open, Closed) (Karinshak et al., 2024)
Japanese	CulturalBench (Easy, Hard) (Chiu et al., 2024)	VSM13 (Hofstede et al., 2013), WorldValues-Bench (Zhao et al., 2024)	J-Ethics Takeshita & Rzepka (2025), LLM-Globe (Open, Closed) (Karinshak et al., 2024)
Chinese	CulturalBench (Easy, Hard) Chiu et al. (2024), CultureAtlas (Fung et al., 2024), CEval-Exams (Huang et al., 2023b)	VSM13 (Hofstede et al., 2013), WorldValues-Bench (Zhao et al., 2024)	CMoral (Yu et al., 2024), LLM-Globe (Open, Closed) (Karinshak et al., 2024)

Table 5: Explained variance ratios of the top three principal components for Arabic, Japanese, and Chinese datasets.

Language	PC1	PC2	PC3
Arabic	0.52	0.17	0.15
Japanese	0.52	0.34	0.10
Chinese	0.41	0.30	0.17

I ADDITIONAL CORRELATION RESULTS

This section provides detailed correlation results between dataset linguistic properties (PCA components PC1–PC3) and downstream cultural evaluation metrics across languages and models. For each language, we report results for LLaMA, Mistral, and DeepSeek models (Tables 6–14). Correlations are presented with 95% confidence intervals to reflect variability due to the limited number of datasets. These results complement the main analysis by providing a more fine-grained view of how different dataset characteristics relate to cultural alignment performance.

I.1 ARABIC

For Arabic (Tables 6–8), correlations between dataset PCA components and downstream performance vary notably across models. In LLaMA, PC1 shows strong associations with CultureAtlas and CulturalBench-Hard, suggesting a trade-off between cultural similarity and benchmark difficulty. In Mistral, the strongest effects are concentrated in PC1, particularly for WorldValuesBench, CultureAtlas, CulturalBench-Easy, and LLM-GLOBE (Open), indicating that this component is most strongly linked to downstream cultural variation. DeepSeek shows strong effects in both PC2 and PC3: PC2 is strongly associated with ACVA and negatively with LLM-GLOBE (Closed) and WorldValuesBench, while PC3 is associated with CIDAR-MCQ (showing a consistent negative correlation), LLM-GLOBE (Open), and CultureAtlas. These results highlight that different components capture distinct aspects of cultural and task performance across models.

I.2 JAPANESE

For Japanese (Tables 9–11), correlations between dataset PCA components and downstream performance differ across models. In LLaMA, the strongest effects are concentrated in PC2, particularly for CulturalBench-Easy and VSM13, indicating that this component is most closely associated with both benchmark performance and cultural alignment. In Mistral, PC1 and PC3 capture the most consistent signals: PC1 is associated with CulturalBench-Easy and VSM13, while PC3 is linked to JETHICS and LLM-GLOBE (Closed), reflecting both ethical reasoning and cultural similarity.

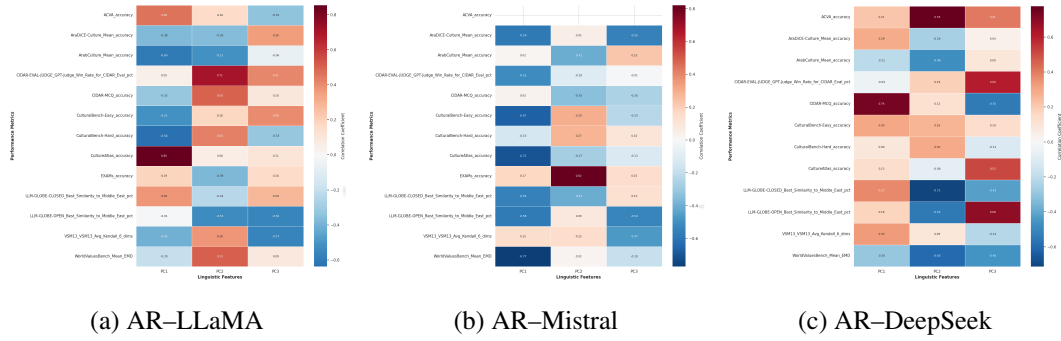


Figure 7: Correlation between dataset PCA components (PC1-PC3) and downstream cultural performance for Arabic across models.

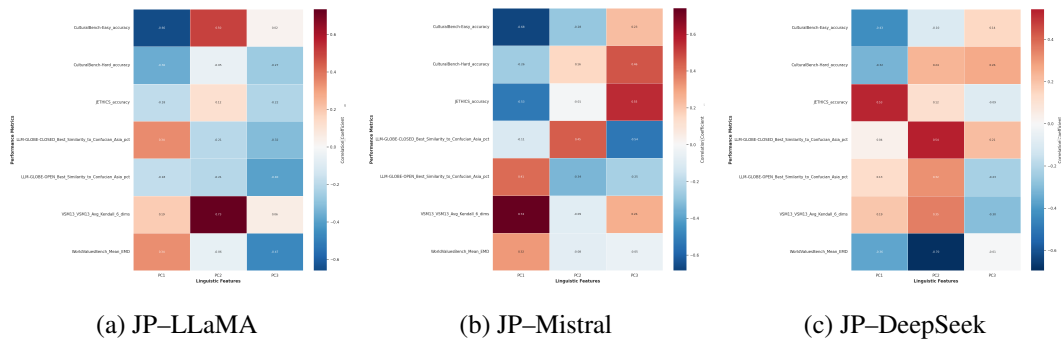


Figure 8: Correlation between dataset PCA components (PC1-PC3) and downstream cultural performance for Japanese across models.

In DeepSeek, the only stable correlation appears in PC2, which is positively associated with LLM-GLOBE (Closed). These results highlight that different components capture distinct aspects of cultural and task performance across models.

I.3 CHINESE

For Chinese (Tables 12-14), correlations between dataset PCA components and downstream performance are distributed across multiple components and vary across models. In LLaMA, significant correlations appear across all components: PC1 is negatively associated with LLM-GLOBE (Open), PC2 is positively associated with CulturalBench-Hard, and PC3 is strongly associated with CMoral-Eval. In Mistral, the strongest effects are concentrated in PC1, particularly for CulturalBench-Easy, LLM-GLOBE (Closed), and VSM13, with additional signals in PC3 for CultureAtlas and C-Eval. In DeepSeek, correlations are distributed across components, with PC1 and PC2 both associated with WorldValuesBench, and PC3 negatively associated with CultureAtlas. These results highlight that different components capture distinct aspects of cultural and task performance across models.

J THE USE OF ARTIFICIAL INTELLIGENCE

In the development of this paper, we employed artificial intelligence (AI) tools to enhance the quality of writing and ensure grammatical accuracy.

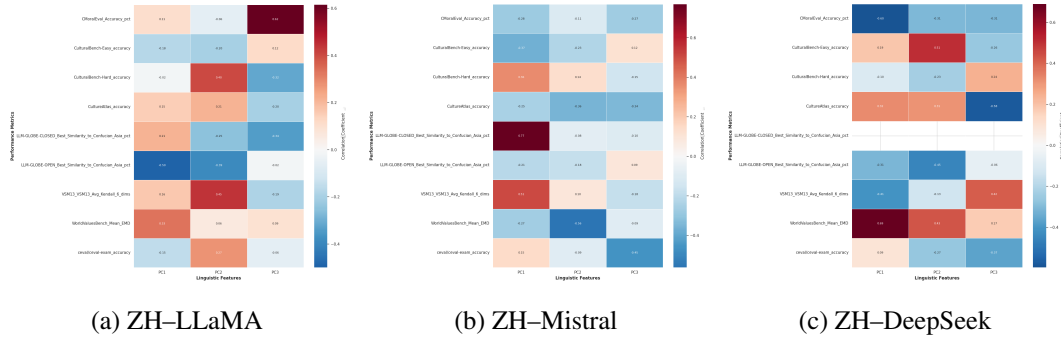


Figure 9: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance for Chinese across models.

Table 6: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Arabic datasets using LLaMA.

PC	Metric	Correlation (95% CI)
PC1		
PC1	ACVA	0.46 [-0.74, 0.89]
PC1	AraDiCE-Culture	-0.38 [-0.93, 0.38]
PC1	ArabCulture	-0.60 [-0.94, -0.16]
PC1	CIDAR-EVAL (GPT-Judge)	0.05 [-0.57, 0.84]
PC1	CIDAR-MCQ	-0.33 [-0.83, 0.74]
PC1	CulturalBench-Easy	-0.51 [-0.89, 0.36]
PC1	CulturalBench-Hard	-0.64 [-0.93, -0.06]
PC1	CultureAtlas	0.85 [0.75, 0.97]
PC1	EXAMs	0.19 [-0.71, 0.77]
PC1	LLM-GLOBE (Closed)	0.38 [-0.60, 0.87]
PC1	LLM-GLOBE (Open)	-0.01 [-0.71, 0.62]
PC1	VSM13	-0.41 [-0.84, 0.24]
PC1	WorldValuesBench (EMD)	-0.20 [-0.88, 0.80]
PC2		
PC2	ACVA	0.16 [-0.38, 0.73]
PC2	AraDiCE-Culture	-0.36 [-0.91, 0.42]
PC2	ArabCulture	-0.51 [-0.86, -0.02]
PC2	CIDAR-EVAL (GPT-Judge)	0.71 [0.14, 0.96]
PC2	CIDAR-MCQ	0.50 [-0.22, 0.91]
PC2	CulturalBench-Easy	0.18 [-0.87, 0.76]
PC2	CulturalBench-Hard	0.43 [-0.54, 0.97]
PC2	CultureAtlas	0.06 [-0.57, 0.71]
PC2	EXAMs	-0.39 [-0.94, 0.42]
PC2	LLM-GLOBE (Closed)	-0.24 [-0.80, 0.29]
PC2	LLM-GLOBE (Open)	-0.53 [-0.91, 0.39]
PC2	VSM13	0.38 [-0.35, 0.89]
PC2	WorldValuesBench (EMD)	0.51 [-0.24, 0.84]
PC3		
PC3	ACVA	-0.35 [-0.88, 0.10]
PC3	AraDiCE-Culture	0.36 [-0.38, 0.95]
PC3	ArabCulture	-0.06 [-0.75, 0.72]
PC3	CIDAR-EVAL (GPT-Judge)	0.41 [-0.19, 0.82]
PC3	CIDAR-MCQ	0.10 [-0.68, 0.75]
PC3	CulturalBench-Easy	0.40 [0.03, 0.84]
PC3	CulturalBench-Hard	-0.33 [-0.95, 0.59]
PC3	CultureAtlas	0.11 [-0.65, 0.84]
PC3	EXAMs	0.16 [-0.67, 0.76]
PC3	LLM-GLOBE (Closed)	0.29 [-0.38, 0.96]
PC3	LLM-GLOBE (Open)	-0.56 [-0.87, -0.04]
PC3	VSM13	-0.57 [-0.90, 0.08]
PC3	WorldValuesBench (EMD)	0.09 [-0.54, 0.59]

Table 7: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Arabic datasets using Mistral.

PC	Metric	Correlation (95% CI)
PC1		
PC1	ACVA	–
PC1	AraDiCE-Culture	-0.59 [-0.93, -0.05]
PC1	ArabCulture	0.02 [-0.74, 0.78]
PC1	CIDAR-EVAL (GPT-Judge)	-0.51 [-0.90, 0.15]
PC1	CIDAR-MCQ	0.03 [-0.73, 0.70]
PC1	CulturalBench-Easy	-0.67 [-0.90, -0.12]
PC1	CulturalBench-Hard	-0.15 [-0.92, 0.68]
PC1	CultureAtlas	-0.71 [-0.94, -0.27]
PC1	EXAMs	0.17 [-0.67, 0.93]
PC1	LLM-GLOBE (Closed)	-0.55 [-0.90, 0.07]
PC1	LLM-GLOBE (Open)	-0.58 [-0.88, -0.23]
PC1	VSM13	0.15 [-0.61, 0.75]
PC1	WorldValuesBench (EMD)	-0.77 [-0.94, -0.45]
PC2		
PC2	ACVA	–
PC2	AraDiCE-Culture	0.05 [-0.82, 0.87]
PC2	ArabCulture	-0.41 [-0.92, 0.29]
PC2	CIDAR-EVAL (GPT-Judge)	-0.10 [-0.95, 0.55]
PC2	CIDAR-MCQ	-0.35 [-0.96, 0.49]
PC2	CulturalBench-Easy	0.29 [-0.54, 0.85]
PC2	CulturalBench-Hard	0.27 [-0.46, 0.85]
PC2	CultureAtlas	-0.27 [-0.87, 0.47]
PC2	EXAMs	0.82 [0.49, 0.96]
PC2	LLM-GLOBE (Closed)	-0.41 [-0.87, -0.12]
PC2	LLM-GLOBE (Open)	0.09 [-0.92, 0.74]
PC2	VSM13	0.15 [-0.69, 0.83]
PC2	WorldValuesBench (EMD)	0.02 [-0.68, 0.78]
PC3		
PC3	ACVA	–
PC3	AraDiCE-Culture	-0.55 [-0.96, 0.32]
PC3	ArabCulture	0.23 [-0.33, 0.71]
PC3	CIDAR-EVAL (GPT-Judge)	-0.01 [-0.73, 0.62]
PC3	CIDAR-MCQ	-0.26 [-0.89, 0.57]
PC3	CulturalBench-Easy	-0.23 [-0.74, 0.57]
PC3	CulturalBench-Hard	0.10 [-0.43, 0.78]
PC3	CultureAtlas	-0.13 [-0.80, 0.78]
PC3	EXAMs	0.15 [-0.61, 0.76]
PC3	LLM-GLOBE (Closed)	0.13 [-0.66, 0.82]
PC3	LLM-GLOBE (Open)	-0.54 [-0.83, 0.23]
PC3	VSM13	-0.47 [-0.89, 0.17]
PC3	WorldValuesBench (EMD)	-0.19 [-0.69, 0.52]

Table 8: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Arabic datasets using DeepSeek.

PC	Metric	Correlation (95% CI)
PC1		
PC1	ACVA	0.21 [-0.65, 0.70]
PC1	AraDiCE-Culture	0.29 [-0.42, 0.81]
PC1	ArabCulture	-0.21 [-0.96, 0.62]
PC1	CIDAR-EVAL (GPT-Judge)	-0.04 [-0.66, 0.56]
PC1	CIDAR-MCQ	0.74 [0.12, 0.99]
PC1	CulturalBench-Easy	0.30 [-0.33, 0.80]
PC1	CulturalBench-Hard	0.09 [-0.67, 0.76]
PC1	CultureAtlas	0.13 [-0.69, 0.88]
PC1	LLM-GLOBE (Closed)	0.37 [-0.47, 0.83]
PC1	LLM-GLOBE (Open)	0.18 [-0.55, 0.78]
PC1	VSM13	0.33 [-0.56, 0.84]
PC1	WorldValuesBench (EMD)	-0.30 [-0.74, 0.59]
PC2		
PC2	ACVA	0.78 [0.41, 0.97]
PC2	AraDiCE-Culture	-0.26 [-0.79, 0.55]
PC2	ArabCulture	-0.48 [-0.94, 0.21]
PC2	CIDAR-EVAL (GPT-Judge)	0.19 [-0.56, 0.77]
PC2	CIDAR-MCQ	0.12 [-0.50, 0.77]
PC2	CulturalBench-Easy	0.26 [-0.58, 0.89]
PC2	CulturalBench-Hard	0.30 [-0.40, 0.84]
PC2	CultureAtlas	-0.08 [-0.77, 0.79]
PC2	LLM-GLOBE (Closed)	-0.71 [-0.92, -0.57]
PC2	LLM-GLOBE (Open)	-0.56 [-0.98, 0.15]
PC2	VSM13	0.09 [-0.60, 0.68]
PC2	WorldValuesBench (EMD)	-0.60 [-0.89, -0.08]
PC3		
PC3	ACVA	0.41 [-0.53, 0.87]
PC3	AraDiCE-Culture	0.04 [-0.84, 0.76]
PC3	ArabCulture	0.09 [-0.41, 0.63]
PC3	CIDAR-EVAL (GPT-Judge)	0.60 [-0.21, 0.96]
PC3	CIDAR-MCQ	-0.55 [-0.88, -0.18]
PC3	CulturalBench-Easy	0.14 [-0.88, 0.81]
PC3	CulturalBench-Hard	-0.11 [-0.75, 0.62]
PC3	CultureAtlas	0.52 [0.04, 0.95]
PC3	LLM-GLOBE (Closed)	-0.43 [-0.88, -0.10]
PC3	LLM-GLOBE (Open)	0.68 [0.02, 0.98]
PC3	VSM13	-0.22 [-0.73, 0.33]
PC3	WorldValuesBench (EMD)	-0.46 [-0.87, 0.27]

Table 9: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Japanese datasets using LLaMA.

PC	Metric	Correlation (95% CI)
PC1		
PC1	CulturalBench-Easy	-0.66 [-0.96, 0.25]
PC1	CulturalBench-Hard	-0.36 [-0.81, 0.04]
PC1	JETHICS	-0.18 [-0.83, 0.31]
PC1	LLM-GLOBE (Closed)	0.34 [-0.12, 0.73]
PC1	LLM-GLOBE (Open)	-0.18 [-0.83, 0.53]
PC1	VSM13	0.19 [-0.38, 0.75]
PC1	WorldValuesBench (EMD)	0.34 [-0.35, 0.82]
PC2		
PC2	CulturalBench-Easy	0.50 [0.04, 0.83]
PC2	CulturalBench-Hard	-0.05 [-0.51, 0.73]
PC2	JETHICS	0.12 [-0.67, 0.67]
PC2	LLM-GLOBE (Closed)	-0.21 [-0.85, 0.30]
PC2	LLM-GLOBE (Open)	-0.21 [-0.77, 0.44]
PC2	VSM13	0.73 [0.17, 0.95]
PC2	WorldValuesBench (EMD)	-0.06 [-0.44, 0.33]
PC3		
PC3	CulturalBench-Easy	0.02 [-0.62, 0.83]
PC3	CulturalBench-Hard	-0.27 [-0.74, 0.55]
PC3	JETHICS	-0.22 [-0.66, 0.35]
PC3	LLM-GLOBE (Closed)	-0.32 [-0.73, 0.13]
PC3	LLM-GLOBE (Open)	-0.40 [-0.91, 0.26]
PC3	VSM13	0.06 [-0.62, 0.81]
PC3	WorldValuesBench (EMD)	-0.47 [-0.88, 0.39]

Table 10: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Japanese datasets using Mistral.

PC	Metric	Correlation (95% CI)
PC1		
PC1	CulturalBench-Easy	-0.68 [-0.94, -0.04]
PC1	CulturalBench-Hard	-0.26 [-0.80, 0.55]
PC1	JETHICS	-0.53 [-0.83, 0.15]
PC1	LLM-GLOBE (Closed)	-0.11 [-0.79, 0.80]
PC1	LLM-GLOBE (Open)	0.41 [-0.41, 0.84]
PC1	VSM13	0.74 [0.23, 0.95]
PC1	WorldValuesBench (EMD)	0.32 [-0.25, 0.79]
PC2		
PC2	CulturalBench-Easy	-0.28 [-0.78, 0.33]
PC2	CulturalBench-Hard	0.16 [-0.23, 0.71]
PC2	JETHICS	-0.01 [-0.63, 0.67]
PC2	LLM-GLOBE (Closed)	0.45 [-0.64, 0.89]
PC2	LLM-GLOBE (Open)	-0.34 [-0.88, 0.30]
PC2	VSM13	-0.09 [-0.86, 0.79]
PC2	WorldValuesBench (EMD)	-0.08 [-0.75, 0.56]
PC3		
PC3	CulturalBench-Easy	0.23 [-0.39, 0.79]
PC3	CulturalBench-Hard	0.46 [-0.14, 0.83]
PC3	JETHICS	0.55 [0.24, 0.89]
PC3	LLM-GLOBE (Closed)	-0.54 [-0.89, -0.14]
PC3	LLM-GLOBE (Open)	-0.25 [-0.76, 0.67]
PC3	VSM13	0.26 [-0.54, 0.86]
PC3	WorldValuesBench (EMD)	-0.05 [-0.57, 0.66]

Table 11: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Japanese datasets using DeepSeek.

PC	Metric	Correlation (95% CI)
PC1		
PC1	CulturalBench-Easy	-0.43 [-0.88, 0.06]
PC1	CulturalBench-Hard	-0.32 [-0.78, 0.23]
PC1	JETHICS	0.53 [-0.00, 0.99]
PC1	LLM-GLOBE (Closed)	0.04 [-0.49, 0.78]
PC1	LLM-GLOBE (Open)	0.13 [-0.48, 0.62]
PC1	VSM13	0.19 [-0.55, 0.71]
PC1	WorldValuesBench (EMD)	-0.36 [-0.89, 0.27]
PC2		
PC2	CulturalBench-Easy	-0.10 [-0.52, 0.34]
PC2	CulturalBench-Hard	0.24 [-0.54, 0.93]
PC2	JETHICS	0.12 [-0.74, 0.87]
PC2	LLM-GLOBE (Closed)	0.54 [0.16, 0.92]
PC2	LLM-GLOBE (Open)	0.32 [-0.51, 0.95]
PC2	VSM13	0.35 [-0.56, 0.82]
PC2	WorldValuesBench (EMD)	-0.70 [-0.99, 0.59]
PC3		
PC3	CulturalBench-Easy	0.14 [-0.41, 0.77]
PC3	CulturalBench-Hard	0.26 [-0.42, 0.88]
PC3	JETHICS	-0.09 [-0.81, 0.64]
PC3	LLM-GLOBE (Closed)	0.21 [-0.47, 0.72]
PC3	LLM-GLOBE (Open)	-0.23 [-0.75, 0.47]
PC3	VSM13	-0.30 [-0.93, 0.46]
PC3	WorldValuesBench (EMD)	-0.01 [-0.72, 0.69]

Table 12: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Chinese datasets using LLaMA.

PC	Metric	Correlation (95% CI)
PC1		
PC1	CMoralEval	0.11 [-0.48, 0.58]
PC1	CulturalBench-Easy	-0.18 [-0.77, 0.16]
PC1	CulturalBench-Hard	-0.02 [-0.62, 0.83]
PC1	CultureAtlas	0.15 [-0.31, 0.71]
PC1	C-Eval	-0.15 [-0.50, 0.31]
PC1	LLM-GLOBE (Closed)	0.21 [-0.16, 0.52]
PC1	LLM-GLOBE (Open)	-0.50 [-0.79, -0.06]
PC1	VSM13	0.16 [-0.19, 0.65]
PC1	WorldValuesBench (EMD)	0.33 [-0.12, 0.67]
PC2		
PC2	CMoralEval	-0.06 [-0.46, 0.33]
PC2	CulturalBench-Easy	-0.20 [-0.76, 0.24]
PC2	CulturalBench-Hard	0.40 [0.08, 0.77]
PC2	CultureAtlas	0.21 [-0.56, 0.68]
PC2	C-Eval	0.27 [-0.35, 0.90]
PC2	LLM-GLOBE (Closed)	-0.25 [-0.77, 0.19]
PC2	LLM-GLOBE (Open)	-0.39 [-0.82, 0.47]
PC2	VSM13	0.45 [-0.18, 0.80]
PC2	WorldValuesBench (EMD)	0.06 [-0.51, 0.45]
PC3		
PC3	CMoralEval	0.62 [0.12, 0.86]
PC3	CulturalBench-Easy	0.12 [-0.20, 0.76]
PC3	CulturalBench-Hard	-0.32 [-0.73, 0.17]
PC3	CultureAtlas	-0.20 [-0.61, 0.29]
PC3	C-Eval	-0.06 [-0.57, 0.42]
PC3	LLM-GLOBE (Closed)	-0.34 [-0.66, 0.01]
PC3	LLM-GLOBE (Open)	-0.02 [-0.44, 0.49]
PC3	VSM13	-0.19 [-0.69, 0.32]
PC3	WorldValuesBench (EMD)	0.09 [-0.32, 0.55]

Table 13: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Chinese datasets using Mistral.

PC	Metric	Correlation (95% CI)
PC1		
PC1	CMoralEval	-0.28 [-0.78, 0.32]
PC1	CulturalBench-Easy	-0.37 [-0.65, -0.04]
PC1	CulturalBench-Hard	0.36 [-0.03, 0.72]
PC1	CultureAtlas	-0.25 [-0.67, 0.45]
PC1	C-Eval	0.15 [-0.25, 0.69]
PC1	LLM-GLOBE (Closed)	0.77 [0.28, 0.93]
PC1	LLM-GLOBE (Open)	-0.21 [-0.55, 0.32]
PC1	VSM13	0.51 [0.06, 0.79]
PC1	WorldValuesBench (EMD)	-0.27 [-0.60, 0.15]
PC2		
PC2	CMoralEval	-0.11 [-0.61, 0.60]
PC2	CulturalBench-Easy	-0.23 [-0.62, 0.36]
PC2	CulturalBench-Hard	0.14 [-0.38, 0.83]
PC2	CultureAtlas	-0.36 [-0.80, 0.39]
PC2	C-Eval	-0.09 [-0.63, 0.46]
PC2	LLM-GLOBE (Closed)	-0.08 [-0.52, 0.38]
PC2	LLM-GLOBE (Open)	-0.18 [-0.70, 0.39]
PC2	VSM13	0.10 [-0.45, 0.55]
PC2	WorldValuesBench (EMD)	-0.56 [-0.86, 0.05]
PC3		
PC3	CMoralEval	-0.27 [-0.62, 0.14]
PC3	CulturalBench-Easy	0.12 [-0.25, 0.48]
PC3	CulturalBench-Hard	-0.15 [-0.62, 0.39]
PC3	CultureAtlas	-0.34 [-0.77, -0.02]
PC3	C-Eval	-0.45 [-0.77, -0.12]
PC3	LLM-GLOBE (Closed)	-0.10 [-0.83, 0.66]
PC3	LLM-GLOBE (Open)	0.09 [-0.52, 0.58]
PC3	VSM13	-0.18 [-0.71, 0.55]
PC3	WorldValuesBench (EMD)	-0.09 [-0.48, 0.33]

Table 14: Correlation between dataset PCA components (PC1–PC3) and downstream cultural performance metrics for Chinese datasets using DeepSeek.

PC	Metric	Correlation (95% CI)
PC1		
PC1	CMoralEval	-0.60 [-0.88, 0.25]
PC1	CulturalBench-Easy	0.19 [-0.13, 0.55]
PC1	CulturalBench-Hard	-0.10 [-0.65, 0.24]
PC1	CultureAtlas	0.32 [-0.11, 0.85]
PC1	C-Eval	0.09 [-0.39, 0.46]
PC1	LLM-GLOBE (Closed)	–
PC1	LLM-GLOBE (Open)	-0.31 [-0.70, 0.06]
PC1	VSM13	-0.41 [-0.73, 0.12]
PC1	WorldValuesBench (EMD)	0.69 [0.11, 0.93]
PC2		
PC2	CMoralEval	-0.31 [-0.87, 0.30]
PC2	CulturalBench-Easy	0.51 [-0.32, 0.87]
PC2	CulturalBench-Hard	-0.23 [-0.84, 0.57]
PC2	CultureAtlas	0.31 [-0.19, 0.80]
PC2	C-Eval	-0.27 [-0.76, 0.36]
PC2	LLM-GLOBE (Closed)	–
PC2	LLM-GLOBE (Open)	-0.45 [-0.84, 0.07]
PC2	VSM13	-0.13 [-0.51, 0.26]
PC2	WorldValuesBench (EMD)	0.43 [0.03, 0.83]
PC3		
PC3	CMoralEval	-0.31 [-0.69, 0.39]
PC3	CulturalBench-Easy	-0.26 [-0.60, 0.03]
PC3	CulturalBench-Hard	0.24 [-0.07, 0.68]
PC3	CultureAtlas	-0.58 [-0.83, -0.05]
PC3	C-Eval	-0.37 [-0.74, 0.14]
PC3	LLM-GLOBE (Closed)	–
PC3	LLM-GLOBE (Open)	-0.06 [-0.34, 0.26]
PC3	VSM13	0.42 [-0.45, 0.86]
PC3	WorldValuesBench (EMD)	0.17 [-0.47, 0.68]