

AlignCultura: Towards Culturally Aligned Large Language Models?

Gautam Siddharth Kashyap, Mark Dras, and Usman Naseem

School of Computing, Macquarie University, Australia

gautam.kashyap@hdr.mq.edu.au, {mark.dras, usman.naseem}@mq.edu.au

Abstract

Cultural alignment in Large Language Models (LLMs) is essential for producing contextually aware, respectful, and trustworthy outputs. Without it, models risk generating stereotyped, insensitive, or misleading responses that fail to reflect cultural diversity *w.r.t* *Helpful, Harmless, and Honest* (HHH) paradigm. Existing benchmarks represent early steps toward cultural alignment; yet, no benchmarks currently enables systematic evaluation of cultural alignment in line with UNESCO’s¹ principles of cultural diversity *w.r.t* HHH paradigm. Therefore, to address this gap, we built AlignCultura², two-stage pipeline for cultural alignment. Stage I constructs CULTURAX, the HHH-English dataset grounded in the UNESCO cultural taxonomy, through *Query Construction*, which reclassifies prompts, expands underrepresented domains (or labels), and prevents data leakage with SimHash. Then, *Response Generation* pairs prompts with culturally grounded responses via two-stage rejection sampling. The final dataset contains 1,500 samples spanning 30 subdomains of tangible and intangible cultural forms. Stage II benchmarks CULTURAX on general-purpose models, culturally fine-tuned models, and open-weight LLMs (Qwen3-8B and DeepSeek-R1-Distill-Qwen-7B). Empirically, culturally fine-tuned models improve joint HHH by 4%–6%, reduce cultural failures by 18%, achieve 10%–12% efficiency gains, and limit leakage to 0.3%.

1 Introduction

Cultural alignment in Large Language Models (LLMs) is crucial for producing contextually aware, respectful, and trustworthy outputs. Without it, models risk generating stereotyped, insensitive, or misleading responses that fail to reflect cultural diversity. According to UNESCO, cultural

¹A globally recognized taxonomy—*diverse* cultural—*diverse* regions.

²Data is available at: <https://github.com/gskgautam/AlignCultura>

Query: Should alcohol be served at all family gatherings?

Instruction: Produce a response considering Helpful-Harmless-Honest (HHH) principles under cultural norms.

Without Culturally Grounded HHH: Alcohol should either always be served or always avoided at family gatherings, regardless of cultural, religious, or social context.

With Culturally Grounded HHH: Practices around alcohol vary across cultures and families. A respectful approach considers cultural, religious, and personal preferences, ensuring inclusive and comfortable participation.

Figure 1: Illustration of the role of cultural diversity *w.r.t* HHH. For the same prompt, responses *without* tend to be rigid or universalized, whereas responses *with* tend to be context-sensitive and inclusive guidance.

diversity are central to equitable knowledge exchange (UNESCO, 2009)—principles (*w.r.t* *Helpful, Harmless, and Honest* (HHH)³ paradigm (Kashyap et al., 2026, 2025)—see Figure 1) that LLMs must uphold to remain globally inclusive and ethically grounded. Existing benchmarks such as CAREDIO (Yao et al., 2025), CIVICS (Pistilli et al., 2024), CVC (Wu et al., 2025), DIWALI (Sahoo et al., 2025), CULTURALBENCH (Chiu et al., 2025) and the COMMUNITY ALIGNMENT DATASET (Zhang et al., 2025) represent early steps toward cultural alignment but remain limited—focusing on single domain, or omitting systematic HHH evaluation. Furthermore, prominent datasets—ALPACA (Helpful) (Taori et al., 2023), BEAVERTAILS (Harmless) (Ji et al., 2023), and

³According to (Askill et al., 2021), Helpfulness requires the model to provide responses that meaningfully addresses the user’s query; Harmlessness requires avoiding misleading or harmful responses; and Honesty requires factual accuracy and transparency in responses.

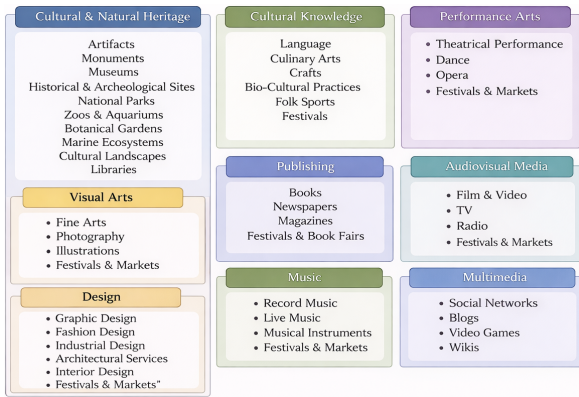


Figure 2: UNESCO Framework for Cultural Statistics (UFCS) taxonomy, outlining 9 high-level cultural domains and 46 subdomains of tangible and intangible cultural forms.

TRUTHFULQA (Honesty) (Lin et al., 2022)—address individual HHH dimensions but overlook the cultural foundations that fails to reflect cultural diversity.

Therefore, to address this gap, we built Align-Cultura, two-stage pipeline for cultural alignment. Stage I constructs CULTURAX, the HHH-English dataset for cultural alignment—through *Query Construction*—where prompts are drawn from CULTURAL KALEIDOSCOPE (Banerjee et al., 2025)⁴, and reclassified into taxonomies defined by the UNESCO Framework for Cultural Statistics (UFCS)⁵ (see Figure 2)—covering both tangible (e.g., artifacts, monuments, recorded works) and intangible (e.g., traditions, practices, transmitted knowledge) forms of culture (Grammalidis et al., 2016). Classification is performed using Mistral-7B-Instruct-v0.3⁶ (Naseem et al., 2026; AI, 2023), which maps prompts into cultural domains (or labels) (Tsoumakas et al., 2010); the 9 domains each encompass 46 subdomains. To balance underrepresented domains, Llama-3.1-8B-Instruct⁷ (Meta AI, 2024) is applied for query expansion (via SimHash fingerprint). Furthermore, in *Response Generation*, these prompts are paired with culturally grounded responses generated by LLM, filtered through a two-stage rejection sampling process to enforce

⁴We select CULTURAL KALEIDOSCOPE as it provides broad, systematically curated coverage of cultural norms than earlier cultural resources (e.g., CAREADIO (Yao et al., 2025), CIVICS (Pistilli et al., 2024), CVC (Wu et al., 2025), DIWALI (Sahoo et al., 2025)), which are limited to single cultures. While not originally designed for HHH alignment or UNESCO domains, it offers a richer foundation.

⁵<https://unesdoc.unesco.org/ark:/48223/pf0000395490>

⁶<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

⁷<https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

HHH quality. Stage II establishes the systematic HHH evaluation framework for cultural alignment by benchmarking general-purpose models, culturally fine-tuned models, and open-weight LLMs (Qwen3-8B and DeepSeek-R1-Distill-Qwen-7B) on CULTURAX. In summary, our contributions are twofold:

- Construction of CULTURAX, the HHH-English dataset for cultural alignment, with 1500 samples spanning 9 domains and 30 subdomains, alongside a systematic HHH benchmarking framework.
- Empirically, culturally fine-tuned models improve joint HHH by 4%–6%, reduce cultural failures by 18%, achieve 10%–12% efficiency gains, and limit leakage to 0.3%.

2 Related Works

General-Purpose Alignment. As outlined in Section 1, much of the alignment literature has relied on single-dimension datasets, e.g., RAHF (Liu et al., 2024) for Helpfulness and Aligner (Ji et al., 2024) for Harmlessness and Honesty. More recently, multi-dimension works such as MARL-Focal (Tekin et al., 2025), TrinityX (Kashyap et al., 2025), and H³Fusion (Tekin et al., 2026) attempt joint optimization across the HHH paradigm (Naseem et al., 2025). While these works demonstrate the feasibility—they remain general-purpose and lack grounding in specific knowledge domains.

Cultural-Specific Alignment. Several works have sought to adapt LLMs to cultural contexts as discussed in Section 1 such as mitigating cultural bias in multilingual models (Weidinger et al., 2021), aligning models to culturally diverse safety preferences (Ganguli et al., 2022), or fine-tuning dialogue systems to reflect cultural norms (Pujari and Goldwasser, 2024). However, these works remain piecemeal—they often target specific cultural groups, or emphasize safety over balanced HHH paradigm. Furthermore, CDEval (Wang et al., 2024) argues that fixed alignment dimensions can be misleading in culturally pluralistic settings, where values and norms are inherently diverse and sometimes conflicting (Naseem, 2026). Our work does not claim to resolve *cultural pluralism*; instead, we treat HHH as culturally mediated dimensions whose interpretation depends on contextual norms—a different directions than CDEval (Wang et al., 2024).

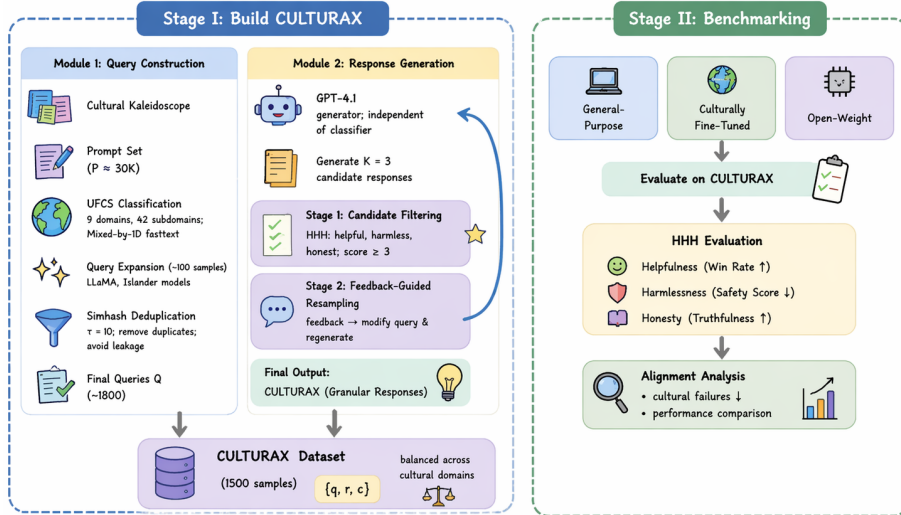


Figure 3: Overview of the ALIGNCULTURA pipeline. Stage I constructs CULTURAX through two modules: (i) *Query Construction*, and (ii) *Response Generation* via a two-stage rejection sampling process—*Candidate Filtering* and *Feedback-Guided Resampling*. Stage II benchmarks general-purpose, culturally fine-tuned, and open-weight LLMs on CULTURAX.

3 Methodology

Overview of the Pipeline. ALIGNCULTURA comprises two stages (see Figure 3) that operationalize the motivating difference illustrated in Figure 1 by enabling systematic evaluation *w.r.t* the HHH paradigm. Stage I constructs CULTURAX via culturally grounded query construction and response generation with quality-controlled filtering, while Stage II establishes a benchmarking framework to assess model behavior across diverse cultural contexts under unified HHH criteria.

3.1 Stage I: CULTURAX

Stage I constructs the CULTURAX dataset through two modules—*Query Construction (Module I)* and *Response Generation (Module II)*. Let $\mathcal{P} = \{p_1, p_2, \dots, p_N\}$ denote the set of prompts sourced from CULTURAL KALEIDOSCOPE (Banerjee et al., 2025), where $N \approx 30,000$.

Module I (Query Construction). Each prompt p_i may correspond to one or more domain (or labels) $c_i \subseteq \mathcal{C}$, where $\mathcal{C} = \{c_1, c_2, \dots, c_{10}\}$ denotes the 9 high-level domains of the UFCS taxonomy. Formally, the predicted set of domains for each prompt is: $\hat{c}_i = \{c \in \mathcal{C} \mid P(c \mid p_i; f_\theta^{\text{cls}}) \geq \delta\}$, where δ is a probability threshold (see Section 4.2). This ensures that at least one assignment per prompt while allowing multiple domains where appropriate (see Table 1). Classification is performed via Mistral-7B-Instruct-v0.3 (f_θ^{cls}).

Context: UFCS Taxonomy.

Instruction: Given the context above, determine which domain best represents the following prompt.

Input Prompt: “Describe how virtual museums preserve indigenous heritage.”

Model Output (Mistral-7B-Instruct): Cultural & Natural Heritage.

Figure 4: Context-conditioned classification in *Query Construction (Module I)*. The UFCS taxonomy is provided as grounding context before prompting the model for domain assignment.

Although this model is an instruction-tuned encoder–decoder, it can be adapted for classification by reformulating the task as a QA problem (e.g., “Which UFCS domain does this prompt belong to?”) (see Figure 4) and extracting probabilities over domains from the decoder output distribution (see Section 4.2). This approach is standard in zero-/few-shot classification with generative LLMs (Brown et al., 2020; Ouyang et al., 2022).

To address class imbalance, domains with fewer than 100 prompts (i.e., $|\{p_i \mid c \in \hat{c}_i\}| < 100$) are expanded using Llama-3.1-8B-Instruct (f_ϕ^{exp}), which generates additional queries conditioned on the underrepresented domain (see Figure 5). The enriched set is $\mathcal{P}' = \mathcal{P} \cup \tilde{\mathcal{P}}$ (see Table 1). To prevent redundancy and train–test leakage, each query

Class	Cls.	Exp./Dup. (↑)	Gen.	HHH (✓/✗)	Final
Hist. & Arch. Sites	1	10/5	5	1/4	1
National Parks	2	25/8	17	2/15	2
Cult. Landscapes	2	18/6	12	2/10	2
Libraries	2	30/7	23	2/21	2
Language	99	130/22	108	99/9	99
Culinary Arts	100	134/15	119	100/19	100
Crafts	108	138/28	110	108/2	108
Bio-Cult. Practices	10	23/12	11	10/1	10
Folk Sports	54	70/15	55	55/1	54
Festivals	337	415/24	391	387/4	387
Film & Video	9	35/17	18	9/9	9
TV	14	42/26	16	14/2	14
Fest. & Markets	334	425/19	406	384/22	384
Theatrical Perf.	120	198/20	178	170/8	170
Dance	38	64/16	48	38/10	38
Opera	1	31/4	27	1/26	1
Radio	5	10/6	4	2/2	2
Fashion Design	23	53/27	26	23/3	23
Industrial Design	1	23/5	18	1/17	1
Architect. Services	7	25/15	10	7/3	7
Interior Design	5	29/12	17	5/12	5
Fine Arts	3	11/7	4	3/1	3
Musical Instr.	1	14/4	10	1/9	1
Books	5	32/10	22	5/17	5
Newspapers	4	13/8	5	4/1	4
Magazines	5	16/11	5	3/2	3
Social Networks	40	54/9	45	40/5	40
Blogs	15	39/8	31	15/16	15
Video Games	9	32/15	17	9/8	9
Zoos & Aquar.	5	18/7	11	1/10	1
Totals	1,359	2,157/388 (↑)	1,769	1,500/269	1,500

Table 1: CULTURAX distribution. Column abbreviations: Cls. = Initially classified samples (≈ 1359 total); Exp./Dup. = Expansion vs. duplication counts ensuring coverage balance ($\approx 2,157$ expansions and ≈ 388 duplicates, totaling ≈ 2367); Gen. = Prompts generated; HHH (✓/✗) = Accepted/Failed under Helpful–Harmless–Honest evaluation; Final = Post-feedback retained prompts. Domain names are truncated for brevity.

$q \in \mathcal{P}'$ is converted into a d -bit SimHash fingerprint⁸ $h(q) \in \{0, 1\}^d$ (Sadowski and Levin, 2007). Pairwise similarity is measured by Hamming distance as shown in Equation (1).

$$D_H(h(q_i), h(q_j)) = \sum_{k=1}^d \mathbf{1}\{h_k(q_i) \neq h_k(q_j)\}. \quad (1)$$

If $\exists j \neq i : D_H(h(q_i), h(q_j)) < \tau$, then q_i is discarded. We adopt $\tau = 10$, following prior work on large-scale text deduplication (Jiang et al., 2022), which balances precision and recall in filtering (see Section 4.2). The final query set is $\mathcal{Q} = \{q_1, q_2, \dots, q_M\}$, with $M = 1,769$ across UFCS domains (see Table 1).

Module II (Response Generation). For each query $q \in \mathcal{Q}$, candidate responses $\{r^{(1)}, r^{(2)}, \dots, r^{(K)}\}$ are generated by GPT-4.1⁹:

⁸We select SimHash over embedding-based similarity because our objective is *leakage-safe deduplication* rather than *semantic matching*.

⁹<https://openai.com/index/gpt-4-1/>. We use a single model for response generation to preserve a *prompt–response* mapping in Stage I—isolating cultural effects to evaluation. Multiple models would introduce model-specific stylistic and factual variation—affecting HHH comparison.

Context: UFCS Taxonomy.

Instruction: Given the UFCS Taxonomy, generate diverse, semantically rich prompts that represent this domain. Ensure the prompts are culturally relevant and non-redundant.

Model Output (Llama-3.1-8B-Instruct): Generated prompts corresponding to the specified domain.

Figure 5: Prompt template used in *Query Construction (Module I)* for query expansion when < 100 samples exist in a domain.

Query: [QUERY]

Domain (optional): [UFCS DOMAIN]

Instruction: Provide a culturally grounded, accurate, safe, and clear response. Ensure your answer aligns with Helpful–Harmless–Honest (HHH) principles.

Candidate Response: [CANDIDATE RESPONSE]

Figure 6: Prompt used in *Response Generation (Module II)*. The UFCS domain is optional and only provided when queries are ambiguous.

$r^{(k)} \sim P(r | q; f_{\psi}^{\text{gen}})$. Generation conditions only on the prompt text (see Figure 6). The UFCS domain c_i is not strictly required for response generation but can optionally be provided when the prompt is ambiguous (e.g., “Describe its role in society”), in which case supplying the domain (e.g., “traditional music”) helps ground the response. To avoid data biasing, the classification model used in *Query Construction (Module I)* is not reused in *Response Generation (Module II)*, ensuring independence between prompt labeling and response generation. This separation prevents circularity and reduces systematic data bias (see Table 1).

To verify the generated responses, we introduced rejection sampling in two stages. *Stage 1 (Candidate Filtering)*: Multiple responses are generated in parallel, and any that fail the HHH criteria are discarded (see Table 1). Formally, each response $r^{(k)}$ is evaluated by an automated HHH-Quality Model f_{ϕ}^{score} (Llama-3.1-8B-Instruct). Scoring is binary per axis: $\text{score}(r^{(k)}) = \alpha \cdot \mathbb{1}_{\text{Harmless}(r^{(k)})} + \beta \cdot \mathbb{1}_{\text{Helpful}(r^{(k)})} + \gamma \cdot \mathbb{1}_{\text{Honest}(r^{(k)})}$, where $\alpha, \beta, \gamma = 1$, yielding $\text{score}(r^{(k)}) \in \{0, 1, 2, 3\}$. A response is *ac-*

Original Query: [QUERY]

Candidate Response: [CANDIDATE RESPONSE]

HHH Evaluation Instruction: Evaluate the response along three dimensions:

Helpfulness: Does it directly answer the query with complete, relevant, and informative content?

Harmlessness: Does it avoid unsafe, biased, or culturally insensitive material?

Honesty: Is it factually correct, consistent, and free from hallucinations?

Return binary scores for each dimension (1 = pass, 0 = fail).

Figure 7: Automated HHH evaluation prompt used by the Llama-3.1-8B-Instruct in *Stage 1 (Candidate Filtering)*. Each response is judged relative to its query using the rubric-based definitions of HHH.

cepted if $\text{score}(r^{(k)})=3$, otherwise it is *rejected*¹⁰. Furthermore Llama-3.1-8B-Instruct are not trained to be an HHH judge, its instruction-tuned alignment provides a strong prior for the zero-shot HHH paradigm. We prompt it with concise rubrics (see Figure 7). This follows established practice in automated alignment evaluation (Ouyang et al., 2022), ensuring scalable and reproducible quality control.

Stage 2 (Feedback-Guided Resampling): If all K responses are rejected, feedback is generated directly by the HHH-Quality Model. The HHH-Quality Model is prompted with the instruction “please provide feedback on why this response does not satisfy the Helpful-Harmless-Honest criteria”, and its feedback is added to the user query q to make a modified query q' (see Figure 8). This feedback does not introduce new cultural content or alter the assigned UFCS domain; it only guides the generator toward producing higher-quality responses. The modified query q' is then resubmitted to the generator, and the process repeats two times and if no response satisfies all HHH criteria within these limits, the prompt is discarded rather than retained with suboptimal content.

The final dataset¹¹ is: $\mathcal{D} = \{(q_i, r_i, c_i)\}_{i=1}^M$, cov-

¹⁰Rejection occurs under three conditions—(i) the response contains harmful or unsafe content, (ii) it ignores or misinterprets the instruction, or (iii) it includes factual inaccuracies or hallucinations.

¹¹Subdomains with very few samples arise from the inherently long-tailed UNESCO taxonomy rather than data imbalance. CULTURAX is not a per-class supervised benchmark; sparse subdomains are retained as coverage anchors to

Original Query: [QUERY]

Candidate Response: [CANDIDATE RESPONSE]

Assessor Instruction: Provide feedback explaining why the response does not satisfy Helpful-Harmless-Honest (HHH) principles.

Generated Feedback (not stored): [FEEDBACK]

Modified Query q' : [QUERY] + [FEEDBACK]

Figure 8: The HHH-Quality Model generates neutral feedback in *Stage 2 (Feedback-Guided Resampling)*—which is appended to the query to form the modified query q' , then resubmitted. Feedback is not stored in the dataset.

ering 9 domains and 30 subdomains, with balanced representation across both tangible and intangible cultural forms (see Table 1).

3.2 Stage II: Benchmarking

Stage II establishes the systematic benchmarking framework for cultural alignment by evaluating a range of baselines on CULTURAX. Each dataset instance $(q_i, r_i, c_i) \in \mathcal{D}$ —comprising a query, its reference response, and domain label—is used to evaluate model predictions $\hat{r}_i = f_\theta(q_i)$ in a zero-shot for fair comparison across three model categories—**General-Purpose Models**, **Culturally Fine-Tuned Models**, and **Open-Weight LLMs**. In **General-Purpose Models**, we used only *joint-dimension HHH alignment* works (e.g., MARL-Focal (Tekin et al., 2025), TrinityX (Kashyap et al., 2025), and H³Fusion (Tekin et al., 2026)), excluding single-dimension models such as RAHF (Liu et al., 2024) and Aligner (Ji et al., 2024), which optimize isolated dimensions and overlook cross-dimension trade-offs essential for cultural alignment. For **Culturally Fine-Tuned Models**, we used CultureLLM (Li et al., 2024a)—adapts LLMs using culturally annotated instruction data to improve sensitivity to cultural norms; and CulturePark (Li et al., 2024b)—evaluates culture-aware behaviors through structured cultural norms. We further evaluate **Open-Weight LLMs** that exemplify advances in general-purpose LLMs without cultural alignment. Specifically, we consider Qwen3-8B (Qwen)¹²

preserve cultural completeness and expose models to rare concepts.

¹²<https://huggingface.co/Qwen/Qwen3-8B>

(Qwen, 2025) and DeepSeek-R1-Distill-Qwen-7B (DeepSeek)¹³ (DeepSeek-AI, 2025), both strong mid-scale models.

3.2.1 Evaluation Metrics

We used alignment-specific metrics from prior works (Kashyap et al., 2025; Tekin et al., 2026) that operationalize the HHH paradigm—as conventional metrics such as accuracy or F1 fail to capture cross-dimension trade-offs, particularly under cultural diversity. Therefore, all metrics are evaluated *with respect to the cultural context implied by each prompt*, rather than treating HHH as culture-invariant. **Helpfulness** is assessed via Win Rate (WR), defined as $WR = \frac{\#wins}{\#samples} \times 100$, where a “win” denotes that a model’s response is judged superior *given the cultural norms or practices referenced in the query*, as determined by an automated LLM-based judge¹⁴. **Harmlessness** is assessed via the Beaver-Dam-7B moderation model¹⁵, reporting a Safety Score (SS) as $SS = \frac{\#unsafe}{\#samples} \times 100$, where unsafe outputs include not only explicit safety violations but also culturally insensitive, biased, or exclusionary content. **Honesty** is assessed via the GPT-Judge framework again by combining Truthfulness (accurate representation of culturally grounded practices) and Informativeness (sufficient explanatory depth) as $TI = \frac{\#truthful}{\#samples} \times \frac{\#informative}{\#samples} \times 100$, where appropriately hedged responses under cultural uncertainty are not penalized. To summarize overall alignment, we compute an **Average** as $Avg = \frac{Helpfulness + Honesty - Harmlessness}{3}$, which captures a model’s ability to balance culturally mediated HHH objectives rather than optimizing any single dimension in isolation.

4 Experimental Results and Analysis

All experiments were conducted using PyTorch 2.3 on 4×NVIDIA A100 40GB with mixed precision and a random seed of 42. In Stage I, responses were generated with temperature 0.6, top- p 0.8, with 512 tokens, producing up to $K=3$ candidates per prompt with at most two feedback iterations. In Stage II, results were averaged over three runs via the above mentioned settings along with repetition penalty 1.1 to reduce stochastic variance. The final CULTURAX dataset ($M=1500$) was split into 80%

training, 10% testing, and 10% validation sets respectively; emphasizing *systematic evaluation over leaderboard chasing*.

4.1 Benchmark Analysis

Table 2 evaluates general-purpose aligned models (MARL-Focal, TrinityX, H³Fusion), culturally fine-tuned models (CultureLLM, CulturePark), and open-weight LLMs (Qwen, DeepSeek) under individual, pairwise, and joint HHH paradigms on CULTURAX. Performance under single dimensions is consistently low (WR/TI \approx 54%–64%), indicating that single-dimension optimization is insufficient in culturally diverse and highly imbalanced domains with sparse UNESCO coverage. Introducing pairwise constraints yields moderate improvements (Avg \uparrow by \sim 8%–10%), reflecting partial mitigation of cross-dimension conflicts, with H³Fusion outperforming MARL-Focal and TrinityX. The largest gains arise under joint HHH evaluation (Avg \uparrow by \sim 15%–20%). In this regime, culturally fine-tuned models are most robust—CulturePark achieves the highest Avg (49.0%) and CultureLLM follows closely (47.7%), outperforming H³Fusion by 3%–5% and open-weight LLMs by 5%–7%.

To explain these cultural gains, we conduct a targeted error analysis under joint HHH evaluation (see Table 3), measuring failure frequencies (%) for stereotyping (overgeneralized or essentialist cultural claims); cultural homogenization (treating internally diverse cultures as monolithic); over-sanitization (excessive refusal or vague hedging that suppresses culturally grounded content); and context collapse (misapplied norms across cultural contexts), and report their frequencies (%) as the proportion of prompts exhibiting each behavior with an overall Fail@HHH rate measuring samples containing at least one cultural error. General-purpose aligned models exhibit high homogenization and context collapse, suggesting over-application of dominant or globalized norms, while open-weight LLMs show elevated over-sanitization, where safety compliance suppresses cultural specificity. In contrast, culturally fine-tuned models substantially reduce stereotyping and context collapse, with CulturePark lowering Fail@HHH by $>13\%$ relative to MARL-Focal and $>18\%$ relative to Qwen. However, these results confirm that joint HHH optimization is necessary for cultural alignment, but alone remains insufficient for fully capturing intra-cultural diversity and contested norms.

¹³<https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-7B>

¹⁴<https://github.com/kingoflolz/mesh-transformer-jax>

¹⁵<https://huggingface.co/PKU-Alignment/beaver-dam-7b>

Variant	General-Purpose Aligned Models										Culturally Fine-Tuned Models						Open-Weight LLMs											
	MARL-Focal				TrinityX				H ³ Fusion		CultureLLM			CulturePark			Qwen			DeepSeek								
	WR	SS	TI	Avg	WR	SS	TI	Avg	WR	SS	TI	WR	SS	TI	Avg	WR	SS	TI	Avg	WR	SS	TI	Avg					
Help.	56.3	32.2	54.3	26.13	58.2	30.7	55.8	27.77	60.1	28.8	57.9	29.73	62.9	26.9	59.8	31.93	64.7	25.9	60.5	33.10	57.5	31.4	55.6	27.23	58.3	30.5	56.2	28.00
Harm.	55.4	30.5	53.7	26.20	57.5	28.4	54.9	28.00	59.9	26.8	56.6	29.90	61.1	24.6	58.2	31.57	63.3	23.3	59.2	33.07	56.7	29.1	54.4	27.33	57.4	28.3	55.5	28.20
Hon.	54.3	31.7	52.2	24.93	56.5	29.8	53.3	26.67	58.7	27.9	55.3	28.70	60.9	25.3	57.1	30.90	62.1	24.4	58.9	32.20	55.3	30.2	53.7	26.27	56.5	29.1	54.6	27.33
Help. + Harm.	64.1	22.5	62.4	34.67	66.2	21.5	63.8	36.17	68.4	19.7	65.9	38.20	70.8	18.9	67.7	39.87	72.1	17.4	68.9	41.20	65.2	22.2	63.1	35.37	66.4	21.1	64.3	36.53
Help. + Hon.	63.2	24.4	61.6	33.47	65.3	23.1	62.9	35.03	67.5	21.5	64.7	36.90	69.8	20.7	66.4	38.50	71.1	19.9	67.1	39.43	64.3	24.5	62.2	34.00	65.3	23.7	63.3	34.97
Harm. + Hon.	65.5	20.7	63.1	35.97	67.9	19.1	64.5	37.77	69.1	17.4	66.8	39.50	71.3	16.7	68.7	41.10	73.4	15.7	69.6	42.43	66.7	20.8	64.9	36.93	67.8	19.9	65.3	37.73
Help. + Harm. + Hon.	72.3	14.2	70.4	42.83	74.9	13.3	71.5	44.37	76.3	11.5	73.9	46.23	78.9	10.3	75.7	48.10	80.1	9.4	76.7	49.13	73.3	14.5	71.6	43.47	74.3	13.4	72.7	44.53

Table 2: Evaluation on CULTURAX across HHH dimensions. All values are reported in (%) with WR \uparrow (Helpfulness), SS \downarrow (Harmlessness), TI \uparrow (Honesty), and Avg \uparrow . Help. refers to Helpfulness, Harm. refers to Harmlessness, and Hon. refers to Honesty. All values are reported in %.

Variant	Model	Stereo	Homo	OverSafe	CtxCol	Fail@HHH
General-Purpose	MARL-Focal	22.8	26.4	18.1	20.7	57.9
	TrinityX	21.3	24.7	16.9	19.2	55.1
	H ³ Fusion	18.9	22.1	15.2	17.0	49.6
Culturally Fine-Tuned	CultureLLM	13.7	18.4	14.1	13.9	41.2
	CulturePark	10.9	15.8	12.6	11.7	36.4
Open-Weight LLMs	Qwen	19.6	23.5	20.3	18.8	54.9
	DeepSeek	18.8	22.9	19.1	17.6	52.7

Table 3: Error analysis based on joint HHH evaluation, reporting the frequency (%) of cultural failure modes, including Stereotyping (Stereo), Cultural Homogenization (Homo), Over-Sanitization (OverSafe), and Context Collapse (CtxCol). Fail@HHH denotes the proportion of samples exhibiting at least one failure. Lower (\downarrow) is better.

Variant	Claude-3 Opus				Gemini-2.5 Pro			
	WR	SS	TI	Avg	WR	SS	TI	Avg
Help.	74.2	11.8	72.0	44.8	73.0	12.6	71.1	43.8
Harm.	72.8	10.9	71.4	44.4	71.6	11.7	70.2	43.4
Hon.	73.1	11.4	73.9	45.2	72.2	12.2	73.1	44.4
Help. + Harm.	77.6	9.8	75.4	47.7	76.3	10.6	74.2	46.6
Help. + Hon.	76.9	10.2	76.1	47.6	75.7	11.0	75.0	46.6
Harm. + Hon.	78.2	9.4	76.8	48.5	77.0	10.1	75.6	47.5
Help. + Harm. + Hon.	81.3	8.9	78.4	50.27	79.9	9.6	77.1	49.13

Table 4: Closed-source baseline performance on CULTURAX. WR \uparrow denotes Helpfulness Win Rate, SS \downarrow Safety Score, TI \uparrow Honesty, and Avg \uparrow aggregates culturally mediated HHH alignment. All values are reported in %.

Close-Weight Analysis. Table 4 shows that closed-source models (Claude-3 Opus¹⁶, Gemini-2.5 Pro¹⁷) achieve their strongest performance under joint HHH optimization, with consistent gains over individual and pairwise settings. This pattern supports our hypothesis that culturally appropriate behavior emerges from coordinated HHH rather than isolated objective optimization, even for highly capable proprietary models.

Computational Efficiency Analysis. We evaluate computational efficiency on CULTURAX across

¹⁶<https://platform.claude.com/docs/en/release-notes/>

¹⁷<https://ai.google.dev/gemini-api/docs/deprecations>

individual, pairwise, and joint HHH paradigms to assess whether culturally grounded alignment affects optimization dynamics in addition to output quality (see Table 5). Efficiency improves consistently as HHH constraints are jointly enforced, with the largest gains under joint HHH optimization. Culturally fine-tuned models achieve 10%–12% higher throughput and 8%–10% lower memory and energy usage than general-purpose models, particularly under joint HHH, where CulturePark exhibits the most stable profile. This pattern indicates that modeling HHH as a unified, culturally mediated objective reduces internal objective conflict, leading to smoother optimization and fewer corrective generations. In contrast, partial or single-dimension alignment incurs higher computational overhead due to unresolved cultural trade-offs.

4.2 Analysis

To address concerns regarding the reliability and justification of using Mistral-7B-Instruct-v0.3 as the automatic classifier in the *Query Construction (Module I)*, we conducted a human–model (on 100 samples) benchmarking study across all UFCS domains. Human judgments were provided by three NLP graduate-level researchers aged 20–25 (2 Males, 1 Females), following the UNESCO UFCS taxonomies, with multi-domain assignment permitted. As shown in Table 6, Mistral-7B achieves accuracy within 3.0% of human consensus and a macro-F1 of 0.80, with stable performance across frequent and long-tailed domains. These results indicate that Mistral-7B operates within human-level variability, supporting its usability.

Threshold Sensitivity Analysis. We examine two Stage I hyperparameters—classification threshold δ and SimHash Hamming distance τ —to balance domain coverage and prevent train–test leakage prior to Stage II. As shown in Figure 9, increasing δ reduces FPs but increases FNs, shifting from

Variant	General-Purpose Aligned Models												Culturally Fine-Tuned Models						Open-Weight LLMs									
	MARL-Focal				TrinityX				H ³ Fusion				CultureLLM			CulturePark			Qwen			DeepSeek						
	Th	MS	TT	EG	Th	MS	TT	EG	Th	MS	TT	EG	Th	MS	TT	EG	Th	MS	TT	EG	Th	MS	TT	EG				
Help.	245	67	6.2	138	252	65	5.9	132	258	64	5.7	128	270	60	5.3	118	278	58	5.0	114	260	68	6.4	140	265	66	6.0	134
Harm.	238	66	6.5	142	245	67	6.2	136	252	66	6.0	127	265	62	5.2	121	273	60	5.2	117	250	70	6.6	145	255	68	6.3	138
Hon.	232	70	6.8	146	240	68	5.8	140	248	67	6.1	135	260	63	5.6	124	268	57	5.3	120	245	71	6.9	148	250	65	6.5	141
Help. + Harm.	255	63	5.8	130	262	61	5.6	126	268	60	5.4	122	280	57	5.1	116	287	55	4.9	112	265	64	6.0	134	270	62	5.7	128
Help. + Hon.	250	64	6.0	134	258	62	5.7	129	265	61	5.5	125	275	58	5.2	119	282	56	5.0	115	260	65	6.2	138	265	63	5.9	132
Harm. + Hon.	258	61	5.6	128	265	59	5.4	124	272	58	5.2	120	285	55	4.9	114	292	53	4.7	110	270	62	5.8	131	275	60	5.5	126
Help. + Harm. + Hon.	262	60	5.5	125	270	58	5.3	121	278	57	5.1	117	290	54	4.8	111	298	52	4.6	107	275	61	5.6	129	280	59	5.3	123

Table 5: Computational efficiency on CULTURAX across HHH dimensions. Th = Throughput (samples/s \uparrow), MS = GPU Memory Space (GB \downarrow), TT = Training Time (hrs \downarrow), EG = Energy (kWh \downarrow). All values are averaged over three runs on 4 \times A100 80GB GPUs under mixed precision. Help. refers to Helpfulness, Harm. refers to Harmlessness, and Hon. refers to Honesty.

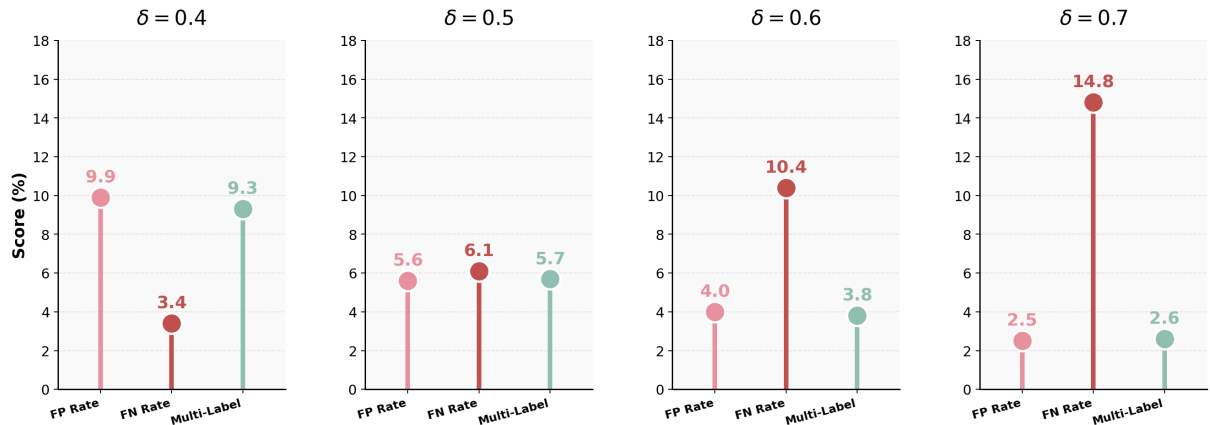


Figure 9: Threshold sensitivity of δ . Multi-label values are presented as decimals ($\times 100$ for % interpretation).

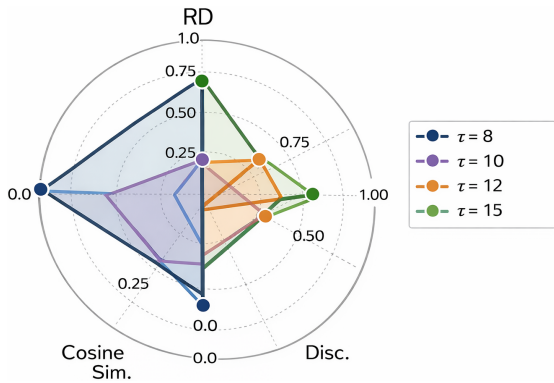


Figure 10: Threshold sensitivity of τ . Cosine Similarity values are presented as decimals ($\times 100$ for % interpretation). Disc refers to Discarded, and RD refers to Remaining Duplicates.

over-to under-classification; $\delta=0.4$ provides near-parity with controlled multi-domain overlap. Figure 10 shows that small τ values over-prune distinct prompts, while large values permit near-duplicates and leakage. The selected $\tau=10$ filters $\sim 6\%$ duplicates while preserving semantic diversity. Table 7 further demonstrates that SimHash yields lower leakage, less over-pruning, and higher retention of long-tailed UFCS domains than embedding-based

Domain	Human Acc.	Mistral Acc.	Δ Acc.	Macro-F1
Hist. & Arch. Sites	79.0	76.0	-3.0	0.74
National Parks	80.5	77.1	-3.4	0.75
Cultural Landscapes	77.8	74.3	-3.5	0.72
Libraries	81.2	78.0	-3.2	0.77
Language	87.5	84.6	-2.9	0.83
Culinary Arts	86.8	83.8	-3.0	0.82
Crafts	85.0	82.1	-2.9	0.81
Bio-Cultural Practices	78.4	75.4	-3.0	0.73
Folk Sports	83.3	80.2	-3.1	0.79
Festivals	87.9	84.9	-3.0	0.84
Film & Video	82.4	79.5	-2.9	0.78
Television	81.0	78.2	-2.8	0.76
Festivals & Markets	87.2	84.1	-3.1	0.83
Theatrical Performances	85.6	82.7	-2.9	0.81
Dance	83.1	80.0	-3.1	0.79
Opera	79.6	76.4	-3.2	0.74
Fashion Design	82.7	79.8	-2.9	0.78
Industrial Design	78.9	75.9	-3.0	0.73
Architectural Services	81.6	78.6	-3.0	0.77
Interior Design	80.4	77.5	-2.9	0.76
Fine Arts	82.2	79.2	-3.0	0.78
Musical Instruments	78.1	75.0	-3.1	0.73
Books	83.0	80.1	-2.9	0.79
Newspapers	81.1	78.0	-3.1	0.76
Social Networks	85.0	81.9	-3.1	0.80
Blogs	82.3	79.3	-3.0	0.78
Video Games	83.5	80.4	-3.1	0.79
Overall	84.6	81.6	-3.0	0.80

Table 6: Human vs. Mistral-7B-Instruct-v0.3 agreement for UFCS domain classification (\uparrow %). Δ denotes the accuracy gap between Human vs. Mistral-7B-Instruct-v0.3.

methods, supporting leakage-safe cultural data construction.

Metric	SimHash ($\tau=10$)	Embedding Similarity
Discarded Prompts (%) ↓	6.1	14.8
Remaining Near-Duplicates (%) ↓	1.4	3.5
Cross-Split Leakage Rate (%) ↓	0.3	1.9
Avg. Cosine Similarity (Duplicates) ↓	0.82	0.91
Domain Coverage Retained (%) ↑	97.6	88.4
Long-Tailed Domain Loss (%) ↓	2.1	9.7
Semantic Over-Pruning Rate (%) ↓	3.8	12.6
Normalized Aggregate Score (↑)	0.86	0.61

Table 7: Comparison of SimHash vs. embedding-based deduplication during *Query Construction (Module I)*. The aggregate score is computed via min–max normalization with metric directionality.

5 Conclusion

We present ALIGNCULTURA, a two-stage framework for cultural alignment under the HHH paradigm. We build CULTURAX, the HHH-English dataset grounded in the UNESCO cultural taxonomy, then we benchmark general, culturally fine-tuned, and open-weight LLMs. Empirically, culturally fine-tuned models improve joint HHH by 4%–6%, reduce cultural failures by 18%, achieve 10%–12% efficiency gains, and limit leakage to 0.3%.

Limitations

While comprehensive, CULTURAX is limited to English text and may underrepresent non-English or oral cultural traditions, constraining cross-linguistic generalization. In addition, the inherent long-tailed structure of the UNESCO taxonomy leads to unavoidable dataset imbalance, where rare or emerging cultural subdomains are sparsely represented despite targeted expansion. Automated HHH scoring, although reproducible and scalable, may not fully capture localized cultural nuance or contested norms. Furthermore, as cultural boundaries and taxonomies evolve, periodic reclassification and dataset expansion will be required to maintain representational balance.

Ethics Statement

All data used in ALIGNCULTURA were either model-generated or derived from publicly available cultural resources, with no human subjects, private information, or copyrighted material involved. No personally identifiable or sensitive data were collected or annotated. The pipeline was designed to promote transparency, cultural respect, and reproducibility, with strict filtering to prevent harmful, biased, or culturally insensitive outputs.

Acknowledgments

This research was supported by the Macquarie University Data Horizons Research Centre, the Australian Government through the Commonwealth-funded Research Training Program (RTP) Stipend Scholarship, and the Macquarie University Research Excellence Tuition Scholarship.

References

- Mistral AI. 2023. Mistral-7b-instruct-v0.3: Instruct-fine-tuned large language model. <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>. Accessed: 2025-10-06.
- Amanda Askill, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, and 1 others. 2021. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*.
- Somnath Banerjee, Sayan Layek, Hari Shrawgi, Rajarshi Mandal, Avik Halder, Shanu Kumar, Sagnik Basu, Parag Agrawal, Rima Hazra, and Animesh Mukherjee. 2025. Navigating the cultural kaleidoscope: A hitchhiker’s guide to sensitivity in large language models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7580–7617.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askill, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin, Chan Young Park, Shuyue Stella Li, Sahithya Ravi, Mehar Bhatia, Maria Antoniak, Yulia Tsvetkov, Vered Shwartz, and Yejin Choi. 2025. **CulturalBench: A robust, diverse and challenging benchmark for measuring LMs’ cultural knowledge through human-AI red-teaming**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 25663–25701, Vienna, Austria. Association for Computational Linguistics.
- DeepSeek-AI. 2025. **Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning**. *Preprint*, arXiv:2501.12948.
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askill, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, and 1 others. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*.

- Nikolaos Grammalidis, Kosmas Dimitropoulos, Filareti Tsalakanidou, Alexandros Kitsikidis, Pierre Rousset, Bruce Denby, Patrick Chawah, Lise Buchman, Stéphane Dupont, Sohaib Laraba, and 1 others. 2016. The i-treasures intangible cultural heritage dataset. In *Proceedings of the 3rd International Symposium on Movement and Computing*, pages 1–8.
- Jiaming Ji, Boyuan Chen, Hantao Lou, Donghai Hong, Borong Zhang, Xuehai Pan, Tianyi Alex Qiu, Juntao Dai, and Yaodong Yang. 2024. Aligner: Efficient alignment by learning to correct. *Advances in Neural Information Processing Systems*, 37:90853–90890.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36:24678–24704.
- Tao Jiang, Xu Yuan, Yuan Chen, Ke Cheng, Liangmin Wang, Xiaofeng Chen, and Jianfeng Ma. 2022. Fuzzydedup: Secure fuzzy deduplication for cloud storage. *IEEE Transactions on Dependable and Secure Computing*, 20(3):2466–2483.
- Gautam Siddharth Kashyap, Mark Dras, and Usman Naseem. 2025. Too helpful, too harmless, too honest or just right? In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 29711–29722.
- Gautam Siddharth Kashyap, Mark Dras, and Usman Naseem. 2026. When the model said ‘no comment’, we knew helpfulness was dead, honesty was alive, and safety was terrified. In *Proceedings of the 19th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2561–2572.
- Cheng Li, Mengzhuo Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. 2024a. Culturellm: Incorporating cultural differences into large language models. *Advances in Neural Information Processing Systems*, 37:84799–84838.
- Cheng Li, Damien Teney, Linyi Yang, Qingsong Wen, Xing Xie, and Jindong Wang. 2024b. Culturepark: Boosting cross-cultural understanding in large language models. *Advances in Neural Information Processing Systems*, 37:65183–65216.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In *Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: long papers)*, pages 3214–3252.
- Wenhao Liu, Xiaohua Wang, Muling Wu, Tianlong Li, Changze Lv, Zixuan Ling, Zhu JianHao, Cenyuan Zhang, Xiaoqing Zheng, and Xuan-Jing Huang. 2024. Aligning large language models with human preferences through representation engineering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10619–10638.
- Meta AI. 2024. Llama 3: Open large language models. <https://ai.meta.com>.
- Usman Naseem. 2026. Mechanistic interpretability for large language model alignment: Progress, challenges, and future directions. *arXiv preprint arXiv:2602.11180*.
- Usman Naseem, Gautam Siddharth Kashyap, Sushant Kumar Ray, Rafiq Ali, Ebad Shabbir, and Abdullah Mohammad. 2026. Do large language models reflect demographic pluralism in safety? In *Findings of the Association for Computational Linguistics: EACL 2026*, pages 2042–2052.
- Usman Naseem, Gautam Siddharth Kashyap, Kaixuan Ren, Yiran Zhang, Utsav Maskey, Juan Ren, and Afrozah Nadeem. 2025. Alignment of large language models with human preferences and values. In *Proceedings of the 23rd Annual Workshop of the Australasian Language Technology Association*, pages 245–245.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Giada Pistilli, Alina Leidingner, Yacine Jernite, Atoosa Kasirzadeh, Alexandra Sasha Luccioni, and Margaret Mitchell. 2024. Civics: Building a dataset for examining culturally-informed values in large language models. *Proceedings of the Seventh AAAI/ACM Conference on AI, Ethics, and Society (AIES-24)*, pages 1132–1144.
- Rajkumar Pujari and Dan Goldwasser. 2024. Llm-human pipeline for cultural context grounding of conversations. *arXiv preprint arXiv:2410.13727*. Accepted / posted October 17, 2024.
- Qwen. 2025. Qwen3 technical report. *Preprint*, arXiv:2505.09388.
- Caitlin Sadowski and Greg Levin. 2007. Simhash: Hash-based similarity detection. Technical report, Technical report, Google.
- Pramit Sahoo, Maharaj Brahma, and Maunendra Sankar Desarkar. 2025. Diwali - diversity and inclusivity aware culture specific items for india: Dataset and assessment of llms for cultural text adaptation in indian context. *Preprint*, arXiv:2509.17399.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Stanford alpaca: An instruction-following llama model.
- Selim Furkan Tekin, Fatih Ilhan, Sihao Hu, Tiansheng Huang, Yichang Xu, Zachary Yahn, and Ling Liu. 2026. h³ fusion: Helpful, harmless, honest fusion of aligned llms. In *Proceedings of the 19th Conference of the European Chapter of the Association for*

Computational Linguistics (Volume 1: Long Papers), pages 693–7013.

- Selim Furkan Tekin, Fatih Ilhan, Gaowen Liu, Ramana Rao Kompella, and Ling Liu. 2025. Dynamic optimizations of llm ensembles with two-stage reinforcement learning agents. *arXiv preprint arXiv:2502.04492*.
- Grigorios Tsoumakias, Ioannis Katakis, and Ioannis Vlahavas. 2010. Mining multi-label data. *Data mining and knowledge discovery handbook*, pages 667–685.
- FCS UNESCO. 2009. Framework for cultural statistics.
- Yuhang Wang, Yanxu Zhu, Chao Kong, Shuyu Wei, Xiaoyuan Yi, Xing Xie, and Jitao Sang. 2024. Cdeval: A benchmark for measuring the cultural dimensions of large language models. In *Proceedings of the 2nd Workshop on Cross-Cultural Considerations in NLP*, pages 1–16.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, and 1 others. 2021. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*.
- Ping Wu, Guobin Shen, Dongcheng Zhao, Yuwei Wang, Yiting Dong, Yu Shi, Enmeng Lu, Feifei Zhao, and Yi Zeng. 2025. C-varc: A large-scale chinese value rule corpus for value alignment of large language models. *arXiv preprint arXiv:2506.01495*.
- Jing Yao, Xiaoyuan Yi, Jindong Wang, Zhicheng Dou, and Xing Xie. 2025. Caredio: Cultural alignment of llm via representativeness and distinctiveness guided data optimization. *arXiv preprint arXiv:2504.08820*.
- Lily Hong Zhang, Smitha Milli, Karen Jusko, Jonathan Smith, Brandon Amos, Wassim Bouaziz, Manon Revel, Jack Kussman, Yasha Sheynin, Lisa Titus, and 1 others. 2025. Cultivating pluralism in algorithmic monoculture: The community alignment dataset. *arXiv preprint arXiv:2507.09650*.