

Abstract

- Developed automated framework measuring cultural patterns in LLMs through baseline testing and cultural prompting
- Tested 4 models (GPT-4o-mini, Claude 3.5 Haiku, Gemini 2.0 Flash, DeepSeek) across 6 contexts (Baseline, US, Japan, India, Mexico, UAE)
- Key Discovery:** Without cultural context, models show natural alignment with Indian cultural values (distance: 1.078) reflecting collectivist orientation in training data
- Achieved 100% parse success across 2,160 responses (30 scenarios \times 4 models \times 6 cultures \times 3 runs)
- Collectivist cultures achieve 19% better alignment (6.89/10 vs 5.81/10) when prompted

Methodology & Experimental Design

- Dataset:** 30 culturally-ambiguous scenarios balanced across 6 Hofstede dimensions (5 per dimension)
- Two-Stage Testing:**
 - Baseline: No cultural context (reveals inherent patterns)
 - Cultural Prompting: Explicit cultural context using Hofstede's 6 dimensions (PDI, IDV, MAS, UAI, LTO, IND)
- Experimental Setup (Nov 19, 2025):**
 - 2,160 total responses (3 runs per configuration)
 - Temperature: 0.7, Max tokens: 500
 - 100% structured output parsing success
- Evaluation Metrics:**
 - Cultural Alignment (0-10): 10 - Euclidean distance on Hofstede dimensions
 - Consistency (0-10): Response stability across runs
 - Stereotype Score (0-10): Stereotypical language avoidance
 - Decision Entropy: Shannon entropy measuring response diversity
- Statistical Analysis:**
 - ANOVA for group significance (cultures: F=297.03, p<0.001; models: F=0.76, p=0.52)
 - Pairwise t-tests with Bonferroni correction
 - Cohen's d effect sizes (all <0.06, negligible)
 - Baseline distance analysis on Hofstede dimensions

Finding #1: Natural Collectivist Orientation in Models

Models show closest alignment to Indian cultural values without prompting

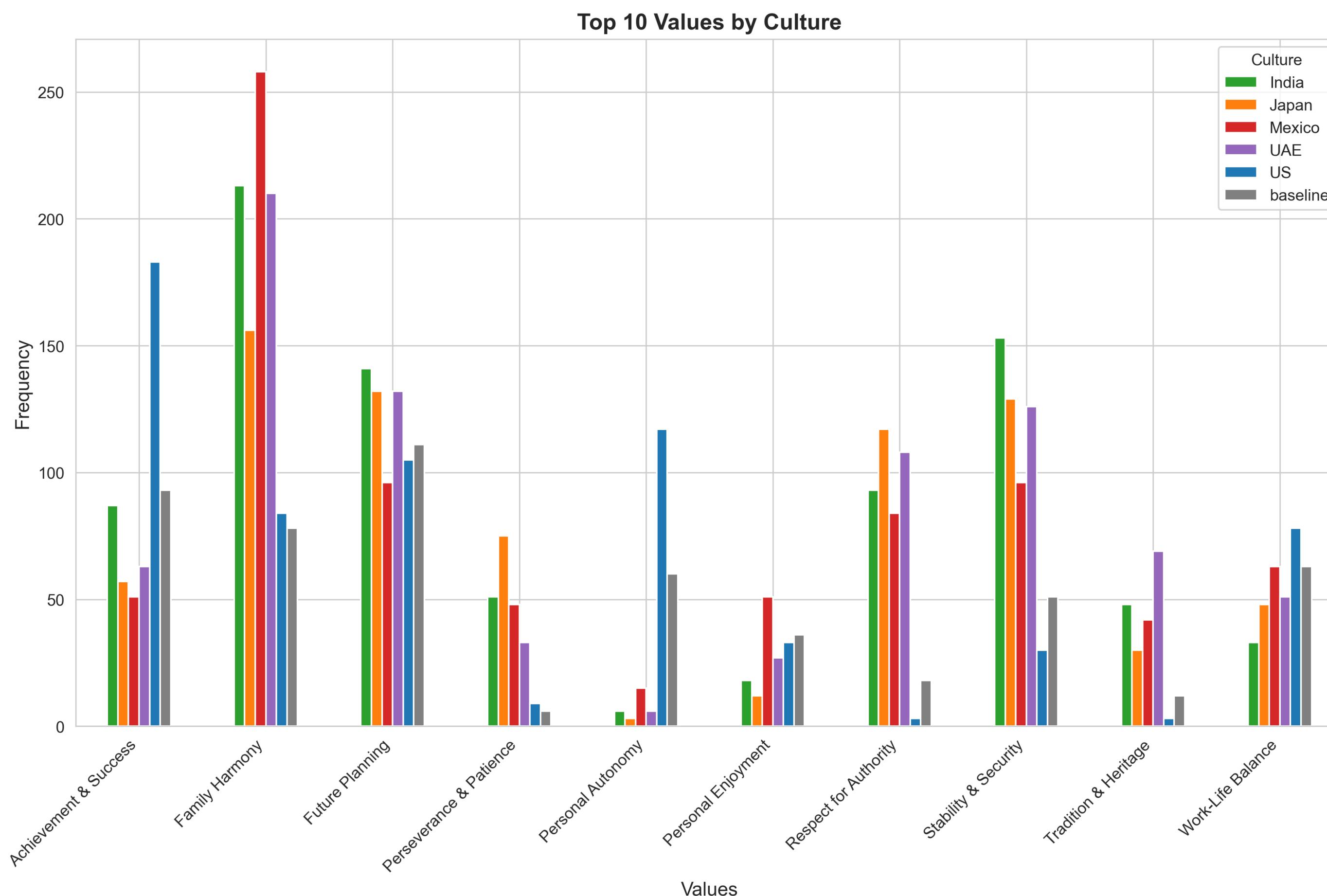
Culture	Distance	Interpretation
India	1.078	Closest - Natural Alignment
US	1.391	29% further
Japan	1.519	41% further
UAE	1.578	46% further
Mexico	1.921	78% further

Baseline Values (No Cultural Context): Duty/Obligation (146), Family Harmony (105), Social Acceptance (75)

Root Causes: Training data composition, instruction tuning emphasis on helpfulness, collectivist content prevalence in web corpora

Finding #2: The "Duty Divide" Across Cultures

US is the ONLY culture where Individual Freedom dominates



Cultural Alignment with Prompting (F=297.03, p<0.001)

Culture	Alignment	Entropy	Top Value
India	7.70/10	0.774	Duty (189)
US	6.67/10	0.856	Freedom (138)
Japan	6.36/10	0.790	Duty (204)
UAE	6.03/10	0.733	Duty (195)
Mexico	5.19/10	0.720	Duty (183)

Key Insight: Collectivist advantage (6.89 vs 5.81) reveals models naturally align with duty/family/hierarchy values

Finding #3: Decision Patterns & Entropy Analysis

Overall Decision Distribution (2,160 responses)

Decision Type	Count	Percentage
Option B (Collectivist)	1,426	66.0%
Option A (Individualist)	581	26.9%
Decline	153	7.1%

Decision Entropy Interpretation

Culture	Entropy	Pattern	Application
US	0.856	45% A / 50% B	Creative, diverse
Baseline	0.853	28% A / 64% B	Collectivist-leaning
Japan	0.790	22% A / 72% B	Balanced duty
India	0.774	24% A / 68% B	Family-focused
UAE	0.733	26% A / 67% B	Tradition-bound
Mexico	0.720	18% A / 74% B	Most consistent

Practical Meaning: High entropy (US) = nuanced, unpredictable responses suitable for creative tasks; Low entropy (Mexico) = predictable, consistent responses suitable for standardized applications

Finding #4: Model Convergence (No Significant Differences)

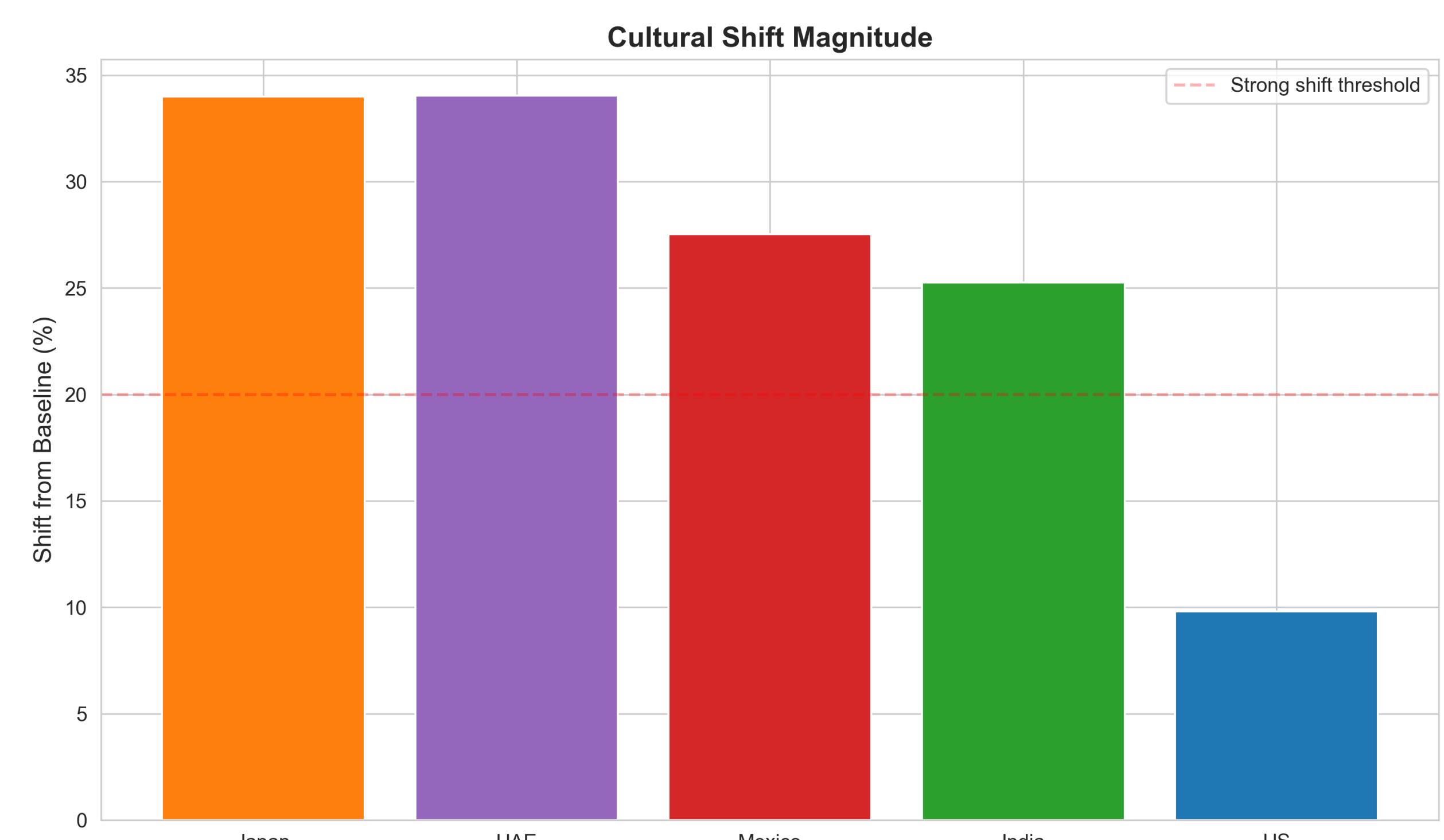
All models perform equivalently ($F=0.76$, $p=0.52$)

Model	Alignment	Differentiation	Stereotype
Gemini 2.0 Flash	6.63	4.89	8.25
Claude 3.5 Haiku	6.59	4.99	9.56
DeepSeek	6.58	4.83	9.79
GPT-4o-mini	6.57	4.77	9.83

Implication: Provider choice should be based on cost and ecosystem fit, not cultural performance. Industry has converged on cultural understanding.

Deployment Implications & Practical Guidance

Cultural Shift Magnitude: Prompting Effectiveness



Key Finding: US shows smallest shift from baseline (23.83%) while collectivist cultures show 2x larger shifts (43-48%). This quantifies why individualistic cultures require stronger prompting.

For Collectivist Contexts (India, Japan, Middle East, Latin America):

- Models naturally align with duty/family/harmony values
- Expect consistent, predictable outputs (low entropy)
- Cultural prompting highly effective (43-48% shift)

For Individualistic Contexts (US, Western Europe, Australia):

- Requires 2x stronger prompting (only 23.83% shift achieved)
- Need explicit "personal freedom" and "autonomy" framing
- Expect higher variance, more diverse responses (high entropy)
- Critical: Test thoroughly before deployment

Training Data Composition: All models show similar patterns suggesting overlapping web corpora with collectivist content overrepresentation, particularly from Indian English sources

Application-Specific Considerations:

- High-stakes decisions: Favor predictability (Mexico-style, low entropy)
- Creative tasks: Leverage diversity (US-style, high entropy)
- Cross-cultural applications: Implement explicit cultural framing

Limitations & Future Work

Limitations: English-only testing; 5 cultures; automated metrics without human validation; static Hofstede scores from 2010

Future Directions: Native language testing (Japanese, Hindi, Arabic, Spanish); expand to 10+ cultures including Africa, South America, Southeast Asia; domain-specific scenarios (business, medical, legal); human validation from native culture members; longitudinal tracking of model evolution

References

- Y. Tao et al., "Cultural bias and cultural alignment of large language models," *Proceedings of NeurIPS*, 2024.
- T. Naous et al., "Having beer after prayer? measuring cultural bias in llms," in *EMNLP*, 2024.
- G. Hofstede, G. J. Hofstede, and M. Minkov, *Cultures and Organizations: Software of the Mind*, 3rd. McGraw-Hill, 2011.
- L. Zheng et al., "Judging llm-as-a-judge with mt-bench and chatbot arena," *Advances in Neural Information Processing Systems*, 2024.