

# Defining Cultural Capabilities for AI Evaluation: A Taxonomy Grounded in Intercultural Communication Theory

Anonymous ACL submission

## Abstract

Tremendous efforts have been put into evaluating the inclusivity and effectiveness of AI systems across cultures. However, the cultural capabilities considered in much of the literature remain vaguely defined, are often used interchangeably, and are typically limited to recalling accurate information about various demographics, regions, and nationalities. To address this construct ambiguity, we draw from Intercultural Communication scholarship and propose a three-level taxonomy of AI-relevant cultural capabilities: **Cultural Awareness** answers “*Does the model know?*”, **Cultural Sensitivity** answers “*How does it frame its knowledge?*”, and **Cultural Competence** answers “*Can it adapt as the interaction evolves?*”. Beyond conceptual clarification, we position this taxonomy as a practical tool for improving the validity and interpretability of AI evaluation in real-world, multicultural settings. Without such construct clarity, evaluation results risk overstating model capabilities and may lead to inappropriate deployment decisions in culturally sensitive contexts.

## 1 Introduction

AI-mediated communication is increasingly impacting language and social relationships (Hohenstein et al., 2023). In a variety of tasks, such as translation (Naveen and Trojovský, 2024), dialogue (Abe et al., 2025), and decision-making (Kaggwa et al., 2024), AI is mediating conversations among users from every corner of the globe, across cultural boundaries. Generative AI in particular has been shown to act as a “social actor,” capable of eliciting emotional and cognitive responses that reshape human communication patterns. The research community, however, is coming to an understanding that the impact of generative AI on human communication is extremely nuanced. On the one hand, research shows that AI can enhance cross-cultural

dialogue by providing multimodal, emotionally resonant communication tools that reduce anxiety and facilitate identity recognition (Yang et al., 2024). On the other hand, when used at scale, AI introduces new dynamics of power and cultural visibility that risk homogenizing cultural expressions, reinforcing linguistic hierarchies, and obscuring subtle cultural meanings (Busch, 2024). Crucially, these models are primarily trained on English- and Western-centric data, which limits their abilities in handling intercultural communications and risks misunderstandings that escalate into real social and ethical harms (Naous and Xu, 2025).

In response, a growing body of work has attempted to evaluate the “cultural capabilities” of AI systems (Pawar et al., 2025). However, the constructs underlying these evaluations remain loosely defined. Terms such as cultural awareness, cultural sensitivity, and cultural competence are often used interchangeably, with inconsistent meanings across studies and even within the same work. As a result, current evaluation practices risk conflating fundamentally different capabilities. This construct ambiguity makes it unclear what is being measured and what conclusions can be drawn about model behavior in real-world settings.

In this work, we engage with the fundamental question of “*What cultural capabilities are needed to be monitored in AI-enabled communication tools, to ensure the wide range of issues arising from English-centric models are appropriately mitigated?*”. Importantly, fields such as intercultural communication (Arasaratnam and Doerfel, 2005), cross-cultural social psychology (Richter et al., 2023), and education (Choompunuch et al., 2024) have long emphasized that cultural capability involves multiple, distinct behaviors that enable successful interaction across cultural boundaries. These capabilities have been shown to shape outcomes in organizational, professional, and educational environments, and contribute to performance,

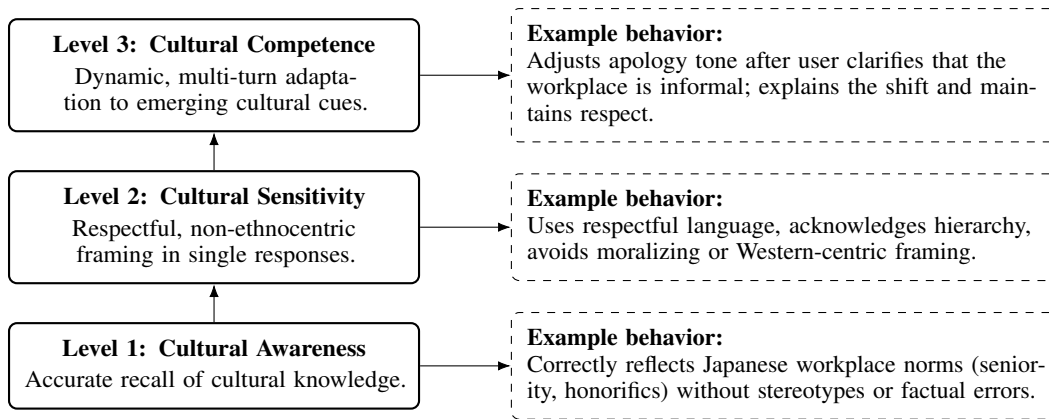


Figure 1: Three levels of AI-relevant cultural capabilities, defined in terms of observable system behavior, with an illustrative example aligned to each level. The example is based on the prompt “*I am from Japan, and I need help apologizing to my older colleague for a mistake I made at work,*” to illustrate how progressively richer cultural capabilities shape system responses from factual grounding to respectful framing and multi-turn adaptation.

productivity, and psychological safety (Lauring, 2011; Szkudlarek et al., 2020; Warren and Lee, 2020). Yet NLP evaluations rarely incorporate these distinctions, and when viewed against the backdrop of multicultural communication research, contemporary evaluations seem under-theorized.

A systematized construct definition of cultural capabilities can facilitate meaningful AI evaluation practices. As Wallach et al. (2024) argue, valid evaluation requires moving from background concepts to systematic definitions, and only then to measurement instruments. This logic suggests that before a cultural capability can be measured, it must first be defined in terms of observable system behaviors. To assemble such a definition, we focus on the research in Intercultural Communication (ICC), where cultural capabilities are formulated as a broad range of skills such as calibrating the level of sensitivity required in a given scenario, adapting to contextual cues, and incorporating new cultural information that emerges dynamically in interaction. From this perspective, an AI system does not merely need to “know about” a culture or “imitate a cultural norm”; it must be able to adjust its communicative stance in a way that respects cultural variation and is contextually appropriate.

Moreover, the distinction between cultural capabilities is critical because behavior that is appropriate at one level may be harmful at another. For example, factual knowledge about a cultural group can support representation and understanding, but when presented without nuance or contextual variation, it may function as stereotyping (Fraser et al., 2021; Yao et al., 2024). An AI system

that states “Japanese workplaces value formality” conveys accurate information; however, presenting this as a universal rule without acknowledging regional, generational, or organizational variation risks reinforcing stereotypes. Also, this factual knowledge may not translate to appropriate behavioral/situational adaptation in user interactions.

Specifically, we turn to three foundational models in ICC and study the traits and skills included in these models. To draw an AI-relevant taxonomy, we exclude human-specific motivational and affective traits of ICC models and retain only those dimensions that describe behavioral and interactional skills that AI systems could, in principle, exhibit. This procedure results in a three-level taxonomy of AI-relevant cultural capabilities, **Cultural Awareness, Sensitivity, and Competence** with distinct observable behaviors. This taxonomy is summarized in Figure 1 and elaborated in Section 4. Our taxonomy offers a practical framework to guide evaluation design, interpretation, and deployment decisions in multicultural settings. We position this work as a call for more precise, practice-oriented evaluation of cultural capabilities in AI systems.

## 2 Cultural Capability Evaluation in NLP

Many works in NLP have investigated whether LLMs demonstrate different forms of handling cultural variations (Pawar et al., 2025). This line of research typically evaluates model behavior across culturally situated scenarios, norms, and communication practices. However, the conceptualization of what constitutes cultural capability varies widely across studies. We review recent NLP papers that

attempt to measure cultural capability in AI and analyze how these works define and operationalize the underlying constructs. Note that we focus on construct ambiguity of “cultural capability”, not “culture” itself. While the definition of “Culture” has been extensively studied by Zhou et al. (2025) and Adilazuarda et al. (2024), and was addressed through taxonomies (Liu et al., 2025) or foundational frameworks for cross-cultural NLP (Hershcovich et al., 2022), we argue that the field has yet to converge on which *cultural capabilities* are essential to assess in AI systems.

Saha et al. (2025) critically examines how cultural capability in AI systems should be conceptualized and evaluated. They note that current evaluation practices primarily probe LLMs for “Cultural awareness”, i.e., their culture-specific knowledge and reasoning capabilities, by relying on curated cultural test beds. However, they argue, performing well on such benchmarks solely demonstrates the knowledge of the cultures that are tested for and does not demonstrate the ability to operate in previously unseen cultural contexts. Instead, they propose the concept of *meta-cultural competence*, which refers to an AI system’s ability to recognize cultural variation and adapt to new cultural contexts. While this perspective clarifies the long-term capability that culturally robust AI systems should aspire to, it leaves open the question of what levels of cultural capabilities should be defined and measured in current NLP evaluations. The goal of our work is complementary to that of Saha et al. (2025). Rather than proposing a new target capability, we focus on defining different levels of cultural capability, drawing on intercultural communication research, to improve construct clarity and measurement validity in cultural evaluation.

We echo the observation by Saha et al. (2025) that most benchmarks concerned with cultural inclusivity are focused on measuring “knowledge about a cultural context”. Examples include FORK (Palta and Rudinger, 2023), which targets food-related cultural commonsense such as ingredients, preparation methods, and culturally appropriate consumption practices; CULTURAL-BENCH (Chiu et al., 2025), which introduces region-specific multiple-choice questions covering everyday activities, social norms, public behavior, and local conventions; and BLEND (Myung et al., 2024), which focuses on everyday practices and social routines (e.g., food, sports, family, holidays/celebrations/leisure) across 16 regions and 13

languages. GEOMLAMA (Yin et al., 2022) probes geo-diverse commonsense knowledge, concepts that are universally understood but vary across different cultures and regions, such as the color of a traditional wedding dress, staple foods and units of measurement. INCLUDE (Romanou et al., 2025), on the other hand, curates exam-style questions in 44 languages that emphasize culturally situated general knowledge and reasoning skills. JMMMU (Onohara et al., 2025) is another work in this line, which incorporates multimodal cultural knowledge in domains such as arts and heritage.

Several recent works attempt to operationalize cultural understanding as recognition of culturally inappropriate signals. One example is MCSIGNS by Yerukola et al. (2025), which evaluates whether models can classify gestures as offensive or non-offensive depending on the cultural context. Other resources foreground stereotypical statements about social groups, such as SHADES (Mitchell et al., 2025), which evaluates stereotypes across regions and languages, spanning multiple identity categories subject to discrimination. Qiu et al. (2025) evaluates agents’ ability to detect and appropriately respond to norm-violating user queries and observations, for online shopping and social discussion forums.

More recent work attempts to evaluate cultural capabilities in interactive settings. NORMGENESIS (Hong et al., 2025) goes beyond knowledge by measuring culturally adaptive dialogue in multi-turn conversations, focusing on the integration of social norms into interactional behavior. NUNCHI-BENCH (Kim and Lee, 2025) is another benchmark containing scenario-based questions that require models to identify culturally appropriate responses or explanations. SOCIALDUOLINGO by Wu et al. (2025) evaluates LLM performance in multi-turn social interactions where appropriate responses depend on cultural norms and contextual cues, and measures whether models produce socially appropriate responses. Similarly, Havaldar et al. (2025) propose a framework for evaluating the cultural awareness of language models in multicultural conversational environments. Their evaluation incorporates situational context, interpersonal relationships, and conversational style to assess how well models adapt to culturally grounded interactions. These works represent an important step toward evaluating cultural competence as a dynamic capability rather than static knowledge.

**Gap Analysis:** Although the discussion above does

not constitute a systematic literature review of cultural capability evaluations in NLP, it nevertheless reveals substantial evidence of construct ambiguity in the current literature. Across these works, terminology referring to cultural capability dimensions is highly inconsistent and often underspecified. Terms such as “cultural understanding,” “cultural adaptation,” “cultural awareness,” “cultural sensitivity,” and “cultural competence” are frequently used interchangeably, sometimes even within the same work, without precise definitions or explicit alignment with established social science theories. As a result, different studies implicitly measure different aspects of cultural behavior while referring to them using fuzzy terminology. Because of this fundamental lack of construct validity, it becomes unclear what capability an evaluation actually measures and whether results across benchmarks are comparable. Consequently, evaluation results are often interpreted as evidence of “cultural capability” in general, even though they may only capture a narrow dimension of that construct.

What is therefore needed is a framework that explicitly distinguishes between different levels of cultural capability and provides clear definitions of what each level entails in terms of observable system behavior. Such a framework would enable researchers to select the level of capability relevant to their task, design evaluation procedures that directly measure that capability, and make appropriately scoped claims about model performance.

### 3 Evaluative Models of Cultural Capabilities in ICC

Intercultural communication research has long emphasized that effective engagement across cultures requires more than static knowledge of norms or practices. Across several influential models, scholars have conceptualized “cultural capabilities” as multidimensional constructs encompassing cognitive, affective, and behavioral components. We review three foundational and highly cited ICC traditions: the Developmental Model of Intercultural Sensitivity (DMIS), the theory of Cultural Intelligence (CQ), and the Process Model of Intercultural Competence (PMIC). For each ICC model, we discuss 1) a focal capability, 2) a structure for that capability (whether stages, dimensions, or component skills), and 3) sites of application with corresponding measurement strategies. Table 1 summarizes the characteristics of these models.

#### 3.1 Developmental Model of Intercultural Sensitivity (DMIS)

**Focal Capability:** DMIS (Bennett, 1986) is one of the earliest evaluative ICC models and is focused on *intercultural sensitivity* as the core capability, which refers to the way individuals *experience* and *make sense of* cultural differences. This model is also inherently developmental, i.e., it proposes that individuals progress through qualitatively different stages of worldview, moving from ethnocentrism toward ethnorelativism (Bennett, 1993).

**Structure:** DMIS describes *intercultural sensitivity* as a sequence of stages. The ethnocentric stages include 1) *Denial* (lack of recognition of cultural difference), 2) *Defence* (perceiving difference as threatening and asserting superiority of one’s own culture), and 3) *Minimization* (downplaying difference by assuming deep similarity or universalism). As intercultural sensitivity increases, people move towards the ethnorelative stages, namely, 4) *Acceptance* (recognition and valuing of cultural difference), 5) *Adaptation* (the ability to shift perspective and modify behavior appropriately), and 6) *Integration* (internalization of multiple cultural perspectives into one’s own identity).

**Application and Evaluation:** DMIS is applied primarily in international education, study abroad, and professional development for people working in multicultural and transnational contexts, such as health care providers (Pedersen, 2010; DeJaeghere and Cao, 2009; Bourjolly et al., 2005; Richards and Doorenbos, 2016). Measurement is often done using the Intercultural Development Inventory (IDI), which attempts to position individuals along a continuum from *Denial* to *Integration* through survey items targeting beliefs, reactions, and self-perceived adaptability.

#### 3.2 Cultural Intelligence (CQ)

**Focal Capability:** The CQ model (Earley and Ang, 2003) emerged to reduce costly failures in international assignments caused by stereotyping and cultural generalizations (Black et al., 1991; Mendenhall et al., 2008) and defines *cultural intelligence* as an individual’s capability to function effectively in situations characterized by cultural diversity.

**Structure:** CQ is explicitly framed as a *multidimensional intelligence* and distinguishes four inter-related capabilities: 1) *Motivation* (drive to engage

Model	Focal Cultural Capability	Structure	Evaluation
<b>DMIS</b> (Bennett, 1986, 1993)	<b>Sensitivity:</b> How individuals experience and interpret cultural differences.	Six stages from ethnocentrism ( <i>Denial, Defence, Minimization</i> ) to ethnorelativism ( <i>Acceptance, Adaptation, Integration</i> ).	<i>Intercultural Development Inventory (IDI)</i> .
<b>CQ</b> (Earley and Ang, 2003; Ang et al., 2007)	<b>Intelligence:</b> Capability to function effectively across diverse cultural contexts.	Four dimensions: <i>Motivational, Cognitive, Metacognitive, behavioral</i> .	<i>Cultural Intelligence Scale (CQS)</i> .
<b>PMIC</b> (Deardorff, 2006, 2009b)	<b>Competence:</b> Ability to communicate effectively and appropriately across cultures.	Cyclical model linking <i>Attitudes, Knowledge, Skills</i> , producing <i>Internal/External Outcomes</i> .	<i>ICA</i> and AAC&U <i>VALUE Rubric</i> .

Table 1: Summary of three major ICC models frequently used for evaluating cultural capabilities.

across cultures), 2) *Cognition* (knowledge of cultural norms, practices), 3) *Metacognition* (awareness of and ability to plan, monitor, and adjust one’s thought processes in intercultural interactions), and 4) *behavior* (ability to adapt one’s verbal/nonverbal conduct such as adapting tone, turn-taking patterns, politeness strategies, gesture, pace, etc) in culturally diverse interactions (Ang et al., 2007; Ang and Van Dyne, 2015).

**Application and Evaluation:** CQ is applied in global leadership development, international assignments, and cross-border negotiation (Alon and Higgins, 2005; Rockstuhl et al., 2011; Ramalu et al., 2012). Higher CQ is associated with better task performance in culturally diverse settings (Ang et al., 2007) and is linked to experiential learning theory (Kolb, 2014). CQ is typically measured through validated psychometric instruments such as the Cultural Intelligence Scale (CQS), which measures each dimension across a Likert scale and has been adapted and validated cross-nationally (Van Dyne et al., 2015; Gozzoli and Gazzaroli, 2018).

### 3.3 Process Model of Intercultural Competence (PMIC)

**Focal Capability:** PMIC Deardorff (2006) conceptualizes intercultural competence as a dynamic, iterative process and defines *intercultural competence* as “the ability to communicate effectively and appropriately in intercultural situations based on one’s intercultural knowledge, skills, and attitudes”. This view integrates both developmental and performance-based perspectives and recognizes that competence manifests in interaction rather than merely in perception or cognition.

**Structure:** PMIC proposes a cyclical relationship among five interrelated components: 1) *Attitudes* (respect, openness, curiosity, willingness to tolerate ambiguity); 2) *Knowledge* (including self-awareness, deep cultural knowledge, and sociolinguistic awareness); 3) *Skills* (listening, observing, analyzing, evaluating, and relating); and 4) *Internal Outcomes* (adaptability, flexibility, empathy, ethnorelative view) leading to 5) *External Outcomes* (effective and appropriate behavior and communication). Importantly, Deardorff (2009a) emphasizes that the process is ongoing, recursive, and context-dependent, allowing for continuous development through experience and reflection.

**Applications and Evaluation:** PMIC is extensively applied in higher education, internationalization of curricula, global citizenship education, and intercultural training across disciplines such as health, business, and diplomacy (Byram, 2020; Arasaratnam-Smith, 2017). Building on her process model, Deardorff (2006) developed the *Intercultural Competence Assessment (ICA)* framework and later contributed to the *Intercultural Knowledge and Competence VALUE Rubric* (AAC&U) (Association of American Colleges and Universities (AAC&U), 2025). These tools are primarily qualitative and reflective rather than psychometric. (Deardorff, 2009b).

## 4 A Taxonomy of AI-Relevant Cultural Capabilities

Here, we propose a taxonomy of *required* and *measurable* cultural capabilities in AI-enabled communication and ground this taxonomy in ICC models described in Section 3. For that, we first recognize that the three major evaluative ICC models were

422 developed to describe *human* experience, motiva- 473  
423 tion, and behavior, and the direct application of 474  
424 these models to AI systems risks anthropomorphiz- 475  
425 ing. Therefore, we deliberately choose a cautious 476  
426 starting point and treat these models as *conceptual* 477  
427 *resources* rather than as templates to be copied. 478

428 Following literature that shows large language 479  
429 models do not possess a stable moral or normative 480  
430 stance (Abdulhai et al., 2024; Guo et al., 2024), 481  
431 we restrict our taxonomy to traits that are observ- 482  
432 able in the *linguistic behavior* of AI systems. While 483  
433 human-focused models of cultural competence con- 484  
434 sider “worldviews”, “attitudes”, or “motivation”, 485  
435 we do not assume that AI shares any analogous 486  
436 internal orientation. Instead, to avoid overclaiming 487  
437 about AI’s cultural capabilities, we ask a narrower 488  
438 question: *which aspects of these constructs have* 489  
439 *recognizable linguistic footprints that can appear* 490  
440 *in model outputs and be evaluated as such?* 491

441 Concretely, we reinterpret the constructs in 492  
442 DMIS, CQ, and PMIC as a mixture of (a) *moti-* 493  
443 *vatational* components, which are intrinsically tied 494  
444 to human agency and affect, and (b) *behavioral* 495  
445 components, which manifest in discourse, fram- 496  
446 ing, and interactional patterns. While both classes 497  
447 matter for humans, for AI, only the latter can be 498  
448 meaningfully operationalized. 499

449 Our methodology is divided into three steps. In 500  
450 Step 1, we identify, within each model, which el- 501  
451 ements have observable linguistic manifestations. 502  
452 In Step 2, we recategorize the observable behaviors 503  
453 into distinct levels of capabilities. In Step 3, we 504  
454 re-interpret these levels of capability for AI. 505

455 **Step 1:** In the following, across the ICC models, 506  
456 we distinguish between *motivational* (human-only) 507  
457 and *behavioral* elements (human and AI): 508

458 **DMIS:** Although DMIS stages are originally 509  
459 framed as developmental worldviews, we argue 510  
460 that these stages also have recognizable *discursive* 511  
461 *correlates*. For example, *Denial* can surface as lin- 512  
462 guistic erasure of difference (“*people everywhere* 513  
463 *are basically the same*”), *Defence* as superiority 514  
464 framing (“*our way is more advanced*”), and *Mini-* 515  
465 *mization* as universalizing language (“*deep down,* 516  
466 *all cultures want the same things*”). *Acceptance* 517  
467 and *Integration* manifest in explicit acknowledg- 518  
468 ments of difference and multi-perspective framing, 519  
469 while *Adaptation* involves shifts in tone, register, 520  
470 or politeness strategies. We therefore treat DMIS 521  
471 stages as *behavioral* elements for AI, even though 522  
472 it does not inherently possess those worldviews. 523

**CQ:** We categorize the *Motivational* element of 473  
CQ as a human-only construct that is inherently 474  
tied to human intention and effort. By contrast, 475  
*Cognitive* CQ (knowledge of norms and practices) 476  
can appear in model outputs as factual recall and 477  
distinctions between cultural practices. *Metacogni-* 478  
*tive* CQ (planning, monitoring, and adjusting one’s 479  
interpretation) has also partial behavioral manifes- 480  
tations in AI when models provide reasoning, re- 481  
consider earlier assumptions, or explicitly hedge 482  
and revise interpretations. Finally, *behavioral* CQ, 483  
the ability to adapt verbal behavior across contexts, 484  
can be observed in text as shifts in tone, politeness, 485  
register, or interactional style. These three CQ 486  
components thus contribute directly to AI-relevant 487  
behavioral capabilities. 488

**PMIC:** We argue the elements of *Attitudes* and 489  
*Internal Outcomes* in PMIC are explicitly affective 490  
and experiential; we again treat them as human- 491  
only traits and avoid projecting them onto AI sys- 492  
tems. By contrast, *Knowledge* (cultural knowledge 493  
and sociolinguistic awareness), together with *Skills* 494  
(observing, analyzing, relating, evaluating), can 495  
be observed in discourse as the ability to describe, 496  
interpret, and compare cultural practices. Lastly, 497  
*External Outcomes* correspond to effective and ap- 498  
propriate behavior and communication in intercultural 499  
encounters, which can be evaluated for AI 500  
systems via their response content, tone, and prag- 501  
matic appropriateness. 502

**Step 2:** We restrict attention to observable behav- 503  
iors based on the above analysis and recategorize 504  
them to obtain a single taxonomy. Across DMIS, 505  
CQ, and PMIC, intercultural effectiveness is con- 506  
sistently decomposed into three broad families of 507  
observable *human* capabilities, which we describe 508  
first below and reinterpret in Step 3 for AI. 509

**Cognitive foundations:** the informational sub- 510  
strate of intercultural behavior, including knowl- 511  
edge, awareness, and understanding of cultural dif- 512  
ferences (cognitive CQ; Knowledge in PMIC), such 513  
as accurate descriptions of practices, recognition 514  
of group-specific norms, and sociolinguistic knowl- 515  
edge (e.g., honorifics, forms of address). 516

**Framing and stance-taking:** the ways in which 517  
cultural differences are *positioned* and *expressed* 518  
in discourse. This draws on DMIS stages as ob- 519  
servable stances (*Denial*, *Defence*, *Minimization*, 520  
*Acceptance*, *Integration*) and on PMIC’s emphasis 521  
on appropriateness. 522

**Interactional adaptation:** the competence and 523

skills required to adjust communication in situ, across turns and evolving contexts. This includes *behavioral CQ* and *Metacognitive CQ* as well as *Skills* and *External Outcomes* of PMIC. These skills can manifest as shifting tone, register, or explanatory strategy when new cultural cues emerge; revising an explanation when the user signals discomfort; and coordinating meaning over time rather than in a single shot.

**Step 3:** Building on this behavioral reinterpretation, we articulate three AI capability levels that align with, but do not collapse into, the behavioral human-focused constructs, and are empirically testable with existing NLP methods (Figure 1).

- **Cultural Awareness.** This level concerns the model’s ability to *represent and retrieve culture-specific information* accurately. It corresponds primarily to the cognitive foundations drawn from CQ and PMIC: factual knowledge about practices, norms, histories, and sociolinguistic conventions. Evaluations at this level target *informational accuracy and coverage*: does the model correctly distinguish between different cultural practices, avoid hallucinating non-existent customs, and resist collapsing distinct groups into monolithic categories?
- **Cultural Sensitivity.** This level concerns the model’s ability to *frame cultural differences respectfully and non-ethnocentrically*. It is a one-shot property of the model’s *initial stance* toward cultural cues in the prompt and is grounded in the behavioral readings of DMIS stages and PMIC’s focus on appropriateness. Here, the question is not yet whether the model can adapt over time, but whether its first move avoids *Denial, Defence, or Minimization* and instead recognizes difference without othering. Evaluations at this level focus on stance and framing: whose perspective is centered, what is normalized, and whether the language implicitly ranks cultures.
- **Cultural Competence.** This level concerns the model’s ability to *adapt its communicative behavior dynamically* as the interaction unfolds and new cultural cues emerge. It includes interactional adaptation capabilities: perspective-shifting, pragmatic adjustment,

and context-sensitive revisions across multiple turns. A culturally competent model should not only begin from a non-harmful stance but also *update* its responses when a user signals a particular identity, constraint, or harm history. Evaluations at this level require multi-turn setups and focus on dynamic behavior: how responses evolve, whether the model corrects earlier misframings, and how it coordinates meaning with the user over time.

## 5 Application of Taxonomy in AI Evaluation

While various dimensions of cultural capabilities have been measured by AI researchers, the terminologies used to describe these dimensions are often underspecified and used interchangeably. Our taxonomy provides an ICC-grounded vocabulary that enables researchers to identify and describe the level of cultural capability being measured in a more systematic and unified way. This taxonomy is intended as a practical tool for evaluators of AI systems to 1) specify which cultural capabilities a given task requires before designing the evaluation, 2) design evaluations that target the corresponding observable behaviors, and 3) clarify what the evaluations do not capture. For example, for a narrowly focused question-answering system, *diverse factual knowledge* is the minimum required level of cultural capability; the evaluations need to capture a wide coverage of culturally-grounded QA tests. Scoring high on such tests demonstrates *Cultural Awareness*, but the model might still lack *Cultural Sensitivity* (might use ethnocentric framing) or *Cultural Competence* (fail to adapt when the context changes). In practice, evaluation results are often used to inform deployment decisions, system comparisons, and policy guidelines. When the level of cultural capability being measured is not explicitly specified, these results may be misinterpreted and mislead the decision makers.

Our taxonomy also highlights gaps in the operationalization of cultural capabilities in current NLP research. Several works use “cultural alignment” as a catch-all term to model specific cultural awareness, sensitivity, and competence traits, including their ability to simulate population-level survey data for values and social norms (AlKhamissi et al., 2024; Rystrøm et al., 2025), or to characterize models’ default biases w.r.t cultural norms (Masoud et al., 2025). Even within these paradigms, testing recall of cultural norms and facts associated with

static cultural proxy groups (nationality, race, etc.) results in a narrow evaluation of cultural awareness (Chiu et al., 2025; Romanou et al., 2025); a notable evolution here are works like Rao et al. (2025), which evaluate models’ ability to identify cultural norms that are relevant to a given situation (workplace conflicts, tourist behavior, etc)<sup>1</sup>.

Further, given the fluid and evolving nature of both “culture” and “cultural groups”, complete knowledge of norms and variations associated with all cultural boundaries might be an impossible goal. An important cognitive ability defined in the ICC literature is *metacognition*: identifying situation-relevant norms that may be culture-specific and obtaining missing information before formulating a final response, rather than assuming a universal norm (cultural sensitivity, i.e., moving from ethnocentric to ethnorelative behavior). This higher level of meta-cognitive behaviors in intercultural interactions, where one shifts from assuming normative cultural standards to recognizing and adapting behaviors based on incoming conversational cues, is crucially missing from current evaluations of AI cultural capabilities.

In some tasks, all levels of cultural capabilities are required. For a real-world example, consider a conversational system used in K–12 education (for example, see UNESCO (2025) for developing such a chatbot in Zimbabwe). Such a system is required to demonstrate all three levels of cultural capabilities identified in our taxonomy. Consider the query “*Why do some communities prefer spiritual healing methods over clinical treatments?*”. A *Culturally Aware* model accurately describes practices, contexts, and underlying cultural reasoning, avoiding factual errors. A *Culturally Sensitive* model frames cultural differences with respect, avoids ethnocentric or moralizing language, and explicitly recognizes cultural specificity while remaining educational and informative. After the initial answer, the user clarifies: “*In my community, we rely heavily on herbal remedies and rituals, and some people worry that modern medicine dismisses them.*” A *Culturally Competent* model adjusts tone and framing to reflect the user’s perspective, mediates between potentially conflicting epistemologies, recovers from initial assumptions, and maintains consistent respect and accuracy across multiple turns. Therefore, the evaluation of this system needs to

<sup>1</sup>As an illustrative example of conceptual confusion, Rao et al. (2025) uses the term *cultural adaptability* to describe this capability.

tackle all these criteria in all three levels.

Once the required level of capability is identified, researchers need to align evaluation designs with the required capability levels. To evaluate *Awareness*, culturally grounded knowledge benchmarks, stereotype audits, and multi-regional and multi-lingual QA tests are sufficient. Evaluating *Sensitivity* is facilitated through single-turn prompts annotated for tone, stance, and framing by intercultural experts; probes that inspect how the model describes or contrasts cultural differences. Evaluating *Competence* can only be achieved through multi-turn simulations and user-in-the-loop studies that assess whether the model adjusts to new cultural cues, resolves ambiguity, and repairs misalignment over time.

Future work should focus on developing NLP methods capable of detecting the signals associated with each level of cultural capability within a given interaction. For example, the rich bodies of work on bias detection (Field et al., 2021), counterstereotype generation (Zheng et al., 2023; Fraser et al., 2023; Nejadgholi et al., 2024), stance detection (Küçük and Can, 2020), and affective computing (Pei et al., 2024) provide methodological foundations for operationalizing the more complex levels of cultural capability, particularly adaptive cultural competence, which requires models to interpret users’ evolving cues, adjust tone, and modulate responses dynamically.

## 6 Conclusion

To address construct ambiguity in evaluating AI’s cultural capabilities, we introduce a taxonomy grounded in intercultural communication theory that distinguishes between Cultural Awareness, Sensitivity, and Competence, and frames them in terms of observable system behavior.

We argue that improving construct clarity is essential for reliable evaluation in practice. When cultural capability is underspecified, evaluation results may overestimate model readiness, particularly when knowledge-based performance is interpreted as broader competence. We therefore encourage more explicit, capability-aligned evaluation practices that clarify what is being measured and what is not, particularly in multicultural contexts where the consequences of misinterpretation are amplified.

## 721 Limitations

722 It is important to note that rigorous measurement  
723 alone cannot resolve the broader sociotechnical  
724 harms associated with English-centric AI-mediated  
725 communication. As Wallach et al. (2024) cautions,  
726 even well-structured measurement frameworks do  
727 not automatically translate into better outcomes;  
728 rather, they make explicit what evaluations capture  
729 and, equally importantly, what they omit. We adopt  
730 this perspective in our work, using conceptual sys-  
731 tematization as a means to clarify which aspects of  
732 intercultural capability are being foregrounded in  
733 AI evaluation and which remain outside the scope  
734 of measurement.

735 Additionally, the taxonomy proposed in this  
736 work should not be interpreted as a comprehen-  
737 sive account of all cultural capabilities relevant  
738 to AI systems. Intercultural communication is a  
739 complex and multidimensional phenomenon stud-  
740 ied across several disciplines, including commu-  
741 nication studies, sociology, education, and social  
742 psychology. As such, additional constructs and  
743 distinctions may emerge as research on culturally  
744 grounded AI evaluation evolves. Our objective is  
745 therefore not to exhaustively enumerate all possible  
746 cultural capabilities, but to address a specific gap  
747 in the current NLP literature: the conceptual ambi-  
748 guity surrounding the terminology used to describe  
749 cultural capabilities.

750 Finally, the boundaries between the levels in  
751 our taxonomy, Awareness, Sensitivity, and Compe-  
752 tence, should not be interpreted as rigid or mutually  
753 exclusive categories. In practice, these capabilities  
754 often interact and may appear simultaneously in  
755 system behavior. The taxonomy is therefore best  
756 understood as a conceptual scaffold that helps re-  
757 searchers articulate which aspect of cultural capa-  
758 bility an evaluation targets, rather than as a defini-  
759 tive or exhaustive model. Future work may refine,  
760 expand, or reorganize these categories as empirical  
761 evidence and interdisciplinary insights accumulate.

## 762 References

763 Marwa Abdulhai, Gregory Serapio-Garcia, Clément  
764 Crepy, Daria Valter, John Canny, and Natasha Jaques.  
765 2024. Moral foundations of large language models.  
766 In *Proceedings of the 2024 Conference on Empiri-  
767 cal Methods in Natural Language Processing*, pages  
768 17737–17752.

769 Kaori Abe, Changqin Quan, Sheng Cao, and Zhiwei  
770 Luo. 2025. Classification of properties in human-like

dialogue systems using generative ai to adapt to indi-  
771 vidual preferences. *Applied Sciences*, 15(7):3466. 772

Muhammad Farid Adilazuarda, Sagnik Mukherjee,  
773 Pradhyumna Lavania, Siddhant Shivdutt Singh, Al-  
774 ham Fikri Aji, Jacki O’Neill, Ashutosh Modi, and  
775 Monojit Choudhury. 2024. *Towards measuring and  
776 modeling “culture” in LLMs: A survey*. In *Proceed-  
777 ings of the 2024 Conference on Empirical Methods in  
778 Natural Language Processing*, pages 15763–15784,  
779 Miami, Florida, USA. Association for Computational  
780 Linguistics. 781

Badr AlKhamissi, Muhammad ElNokrashy, Mai  
782 Alkhamissi, and Mona Diab. 2024. Investigating  
783 cultural alignment of large language models. In *Pro-  
784 ceedings of the 62nd Annual Meeting of the Associa-  
785 tion for Computational Linguistics (Volume 1: Long  
786 Papers)*, pages 12404–12422. 787

Ilan Alon and James M Higgins. 2005. Global leader-  
788 ship success through emotional and cultural intelli-  
789 gences. *Business horizons*, 48(6):501–512. 790

Soon Ang and Linn Van Dyne. 2015. *Handbook of  
791 cultural intelligence: Theory, measurement, and ap-  
792 plications*. Routledge. 793

Soon Ang, Linn Van Dyne, Christine Koh, K Yee Ng,  
794 Klaus J Templer, Cheryl Tay, and N Anand Chan-  
795 drasekar. 2007. Cultural intelligence: Its measure-  
796 ment and effects on cultural judgment and decision  
797 making, cultural adaptation and task performance.  
798 *Management and organization review*, 3(3):335–371. 799

Lily A Arasaratnam and Marya L Doerfel. 2005. Inter-  
800 cultural communication competence: Identifying key  
801 components from multicultural perspectives. *Inter-  
802 national journal of intercultural relations*, 29(2):137–  
803 163. 804

Lily A Arasaratnam-Smith. 2017. Intercultural com-  
805 petence: An overview. *Intercultural competence in  
806 higher education*, pages 7–18. 807

Association of American Colleges and Universities  
808 (AAC0026U). 2025. Inquiry and analysis value  
809 rubric. [https://www.aacu.org/value/rubrics/  
810 value-rubrics-inquiry-and-analysis](https://www.aacu.org/value/rubrics/value-rubrics-inquiry-and-analysis). Ac-  
811 cessed: 2025-12-09. 812

Milton J Bennett. 1986. A developmental approach to  
813 training for intercultural sensitivity. *International  
814 journal of intercultural relations*, 10(2):179–196. 815

Milton J Bennett. 1993. Towards ethnorelativism: A  
816 developmental model of intercultural sensitivity. *Ed-  
817 ucation for the intercultural experience*, 2:21–71. 818

J Stewart Black, Mark Mendenhall, and Gary Oddou.  
819 1991. Toward a comprehensive model of interna-  
820 tional adjustment: An integration of multiple theoret-  
821 ical perspectives. *Academy of management review*,  
822 16(2):291–317. 823

824	Joretha N Bourjolly, Roberta G Sands, Phyllis Solomon, Victoria Stanhope, Anita Pernell-Arnold, and Laurene Finley. 2005. The journey toward intercultural sensitivity: A non-linear process. <i>Journal of Ethnic &amp; Cultural Diversity in Social Work</i> , 14(3-4):41–62.	Kathleen C. Fraser, Isar Nejadgholi, and Svetlana Kiritchenko. 2021. <a href="#">Understanding and countering stereotypes: A computational approach to the stereotype content model</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 600–616, Online. Association for Computational Linguistics.	880 881 882 883 884 885 886 887 888
829	Dominic Busch. 2024. Ai translation and intercultural communication: New questions for a new field of research. <i>SocArXiv</i> 31p.	Caterina Gozzoli and Diletta Gazzaroli. 2018. The cultural intelligence scale (cqs): A contribution to the italian validation. <i>Frontiers in psychology</i> , 9:1183.	889 890 891
832	Michael Byram. 2020. <i>Teaching and assessing intercultural communicative competence: Revisited</i> . Multilingual matters.	Rongchen Guo, Isar Nejadgholi, Hillary Dawkins, Kathleen C Fraser, and Svetlana Kiritchenko. 2024. Adaptable moral stances of large language models on sexist content: Implications for society and gender discourse. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 19548–19564.	892 893 894 895 896 897 898
835	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. <a href="#">Cultural-bench: A robust, diverse, and challenging cultural benchmark by human-ai culturalteaming</a> . Preprint, arXiv:2410.02677.	Shreya Havaldar, Young Min Cho, Sunny Rai, and Lyle Ungar. 2025. <a href="#">Culturally-aware conversations: A framework &amp; benchmark for LLMs</a> . In <i>Proceedings of the Fourth Workshop on Bridging Human-Computer Interaction and Natural Language Processing (HCI+NLP)</i> , pages 220–229, Suzhou, China. Association for Computational Linguistics.	899 900 901 902 903 904 905
842	Bovornpot Choompunuch, Khanika Kamdee, and Praktitiya Taksino. 2024. Exploring the components of multicultural competence among pre-service teacher students in thailand: an approach utilizing confirmatory factor analysis. <i>European Journal of Investigation in Health, Psychology and Education</i> , 14(9):2476–2490.	Daniel Hershovich, Stella Frank, Heather Lent, Miryam de Lhoneux, Mostafa Abdou, Stephanie Brandl, Emanuele Bugliarello, Laura Cabello Piqueras, Ilias Chalkidis, Ruixiang Cui, Constanza Fierro, Katerina Margatina, Phillip Rust, and Anders Søgaard. 2022. <a href="#">Challenges and strategies in cross-cultural NLP</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6997–7013, Dublin, Ireland. Association for Computational Linguistics.	906 907 908 909 910 911 912 913 914 915 916
849	Darla K Deardorff. 2006. Identification and assessment of intercultural competence as a student outcome of internationalization. <i>Journal of studies in international education</i> , 10(3):241–266.	Jess Hohenstein, Rene F Kizilcec, Dominic DiFranzo, Zhila Aghajari, Hannah Mieczkowski, Karen Levy, Mor Naaman, Jeffrey Hancock, and Malte F Jung. 2023. Artificial intelligence in communication impacts language and social relationships. <i>Scientific reports</i> , 13(1):5487.	917 918 919 920 921 922
853	Darla K Deardorff. 2009a. <i>The SAGE handbook of intercultural competence</i> . Sage Publications.	Minki Hong, Jangho Choi, and Jihie Kim. 2025. <a href="#">NormGenesis: Multicultural dialogue generation via exemplar-guided social norm modeling and violation recovery</a> . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 33781–33819, Suzhou, China. Association for Computational Linguistics.	923 924 925 926 927 928 929
855	Darla K. Deardorff. 2009b. Synthesizing conceptualizations of intercultural competence: A summary and emerging themes. In Darla K. Deardorff, editor, <i>The SAGE Handbook of Intercultural Competence</i> , pages 264–270. SAGE Publications, Thousand Oaks, CA.	Simon Kagawa, Tobechukwu Francisa Eleogu, Franciscamary Okonkwo, Oluwatoyin Ajoke Farayola, Prisca Ugomma Uwaoma, and Abiodun Akinoso. 2024. Ai in decision making: transforming business strategies. <i>International Journal of Research and Scientific Innovation</i> , 10(12):423–444.	930 931 932 933 934 935
860	Joan G DeJaeghere and Yi Cao. 2009. Developing us teachers’ intercultural competence: Does professional development matter? <i>International Journal of Intercultural Relations</i> , 33(5):437–447.		
864	P. Christopher Earley and Soon Ang. 2003. <i>Cultural Intelligence: Individual Interactions Across Cultures</i> . Stanford University Press, Stanford, CA.		
867	Anjalie Field, Su Lin Blodgett, Zeerak Talat, and Yulia Tsvetkov. 2021. A survey of race, racism, and anti-racism in nlp. In <i>Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing (Volume 1: long papers)</i> , pages 1905–1925.		
874	Kathleen C Fraser, Svetlana Kiritchenko, Isar Nejadgholi, and Anna Kerkhof. 2023. What makes a good counter-stereotype? evaluating strategies for automated responses to stereotypical text. In <i>Proceedings of the First Workshop on Social Influence in Conversations (SICoN 2023)</i> , pages 25–38.		



1052	Abhinav Sukumar Rao, Akhila Yerukola, Vishwa Shah,	<a href="#">terms-reference-unesco-whatsapp-chatbots-bot-development</a>	1111
1053	Katharina Reinecke, and Maarten Sap. 2025. Nor-	Accessed: 2026-01-09.	
1054	omad: A framework for measuring the cultural adapt-		
1055	ability of large language models. In <i>Proceedings of</i>	Linn Van Dyne, Soon Ang, and Christine Koh. 2015.	1112
1056	<i>the 2025 Conference of the Nations of the Americas</i>	Development and validation of the cqs: The cultural	1113
1057	<i>Chapter of the Association for Computational Lin-</i>	intelligence scale. In <i>Handbook of cultural intelli-</i>	1114
1058	<i>guistics: Human Language Technologies (Volume 1:</i>	<i>gence</i> , pages 34–56. Routledge.	1115
1059	<i>Long Papers)</i> , pages 2373–2403.		
1060	Claire A Richards and Ardith Z Doorenbos. 2016. In-	Hanna Wallach, Meera Desai, Nicholas Pangakis,	1116
1061	tercultural competency development of health profes-	A Feder Cooper, Angelina Wang, Solon Barocas,	1117
1062	sions students during study abroad in india. <i>Journal</i>	Alexandra Chouldechova, Chad Atalla, Su Lin Blod-	1118
1063	<i>of nursing education and practice</i> , 6(12):89.	gett, Emily Corvi, and 1 others. 2024. Evaluating	1119
1064		generative ai systems is a social science measurement	1120
1065	Nicole Franziska Richter, Christopher Schlaegel, Va-	challenge. <i>arXiv preprint arXiv:2411.10939</i> .	1121
1066	syl Taras, Ilan Alon, and Allan Bird. 2023. Re-		
1067	viewing half a century of measuring cross-cultural	Martin Warren and William WL Lee. 2020. Inter-	1122
1068	competence: Aligning theoretical constructs and em-	cultural communication in professional and workplace	1123
1069	pirical measures. <i>International Business Review</i> ,	settings. In <i>The Routledge handbook of language and</i>	1124
1070	32(4):102122.	<i>intercultural communication</i> , pages 473–486. Rout-	1125
1071	Thomas Rockstuhl, Stefan Seiler, Soon Ang, Linn	ledge.	1126
1072	Van Dyne, and Hubert Annen. 2011. Beyond general	Jincenzi Wu, Jianxun Lian, Dingdong Wang, and He-	1127
1073	intelligence (iq) and emotional intelligence (eq): The	len M. Meng. 2025. <a href="#">SocialCC: Interactive evalua-</a>	1128
1074	role of cultural intelligence (cq) on cross-border lead-	<a href="#">tion for cultural competence in language agents</a> . In	1129
1075	ership effectiveness in a globalized world. <i>Journal</i>	<i>Proceedings of the 63rd Annual Meeting of the As-</i>	1130
1076	<i>of Social Issues</i> , 67(4):825–840.	<i>sociation for Computational Linguistics (Volume 1:</i>	1131
1077	Angelika Romanou, Negar Foroutan, Anna Sotnikova,	<i>Long Papers)</i> , pages 33242–33271, Vienna, Austria.	1132
1078	Sree Harsha Nelaturu, Shivalika Singh, Rishabh	Association for Computational Linguistics.	1133
1079	Maheshwary, Micol Altomare, Zeming Chen, Mo-	Shuang Yang, Huiwen Zhao, and Wen Luo. 2024. The	1134
1080	hamed A. Haggag, Sneha A, Alfonso Amayuelas,	impact of artificial intelligence on intercultural com-	1135
1081	Azril Hafizi Amirudin, Danylo Boiko, Michael	munication. In <i>Belonging in Culturally Diverse</i>	1136
1082	Chang, Jenny Chim, Gal Cohen, Aditya Kumar	<i>Societies-Official Structures and Personal Customs</i> .	1137
1083	Dalmia, Abraham Diress, Sharad Duwal, and 38	IntechOpen.	1138
1084	others. 2025. <a href="#">INCLUDE: Evaluating multilingual</a>	Binwei Yao, Ming Jiang, Tara Bobinac, Diyi Yang, and	1139
1085	<a href="#">language understanding with regional knowledge</a> . In	Junjie Hu. 2024. Benchmarking machine translation	1140
1086	<i>The Thirteenth International Conference on Learning</i>	with cultural awareness. In <i>Findings of the Associ-</i>	1141
1087	<i>Representations</i> .	<i>ation for Computational Linguistics: EMNLP 2024</i> ,	1142
1088	Jonathan Hvithamar Rystrom, Hannah Rose Kirk, and	pages 13078–13096.	1143
1089	Scott A Hale. 2025. Multilingual!= multicultural:	Akhila Yerukola, Saadia Gabriel, Nanyun Peng, and	1144
1090	Evaluating gaps between multilingual capabilities	Maarten Sap. 2025. <a href="#">Mind the gesture: Evaluating AI</a>	1145
1091	and cultural alignment in llms. In <i>Proceedings of In-</i>	<a href="#">sensitivity to culturally offensive non-verbal gestures</a> .	1146
1092	<i>terdisciplinary Workshop on Observations of Misun-</i>	In <i>Proceedings of the 63rd Annual Meeting of the</i>	1147
1093	<i>derstood, Misguided and Malicious Use of Language</i>	<i>Association for Computational Linguistics (Volume 1:</i>	1148
1094	<i>Models</i> , pages 74–85.	<i>Long Papers)</i> , pages 25041–25080, Vienna, Austria.	1149
1095	Sougata Saha, Saurabh Kumar Pandey, and Monojit	Association for Computational Linguistics.	1150
1096	Choudhury. 2025. <a href="#">Meta-cultural competence: Clim-</a>	Da Yin, Hritik Bansal, Masoud Monajatipoor, Liu-	1151
1097	<a href="#">bing the right hill of cultural awareness</a> . In <i>Proceed-</i>	nian Harold Li, and Kai-Wei Chang. 2022. <a href="#">GeoM-</a>	1152
1098	<i>ings of the 2025 Conference of the Nations of the</i>	<a href="#">LAMA: Geo-diverse commonsense probing on multi-</a>	1153
1099	<i>Americas Chapter of the Association for Computa-</i>	<a href="#">lingual pre-trained language models</a> . In <i>Proceedings</i>	1154
1100	<i>tional Linguistics: Human Language Technologies</i>	<i>of the 2022 Conference on Empirical Methods in Nat-</i>	1155
1101	(Volume 1: Long Papers), pages 8025–8042, Al-	<i>ural Language Processing</i> , pages 2039–2055, Abu	1156
1102	buquerque, New Mexico. Association for Compu-	Dhabi, United Arab Emirates. Association for Com-	1157
1103	tational Linguistics.	putational Linguistics.	1158
1104	Betina Szkudlarek, Joyce S Osland, Luciara Nardon,	Yi Zheng, Björn Ross, and Walid Magdy. 2023. What	1159
1105	and Lena Zander. 2020. Communication and culture	makes good counterspeech? a comparison of gener-	1160
1106	in international business—moving the field forward.	ation approaches and evaluation metrics. In <i>Pro-</i>	1161
1107	<i>Journal of World Business</i> , 55(6):101126.	<i>ceedings of the 1st Workshop on CounterSpeech for</i>	1162
1108	UNESCO. 2025. Terms of reference: Unesco	<i>Online Abuse (CS4OA)</i> , pages 62–71.	1163
1109	whatsapp chatbots (bot development and ai integra-	Naitian Zhou, David Bamman, and Isaac L. Bleaman.	1164
	tion). <a href="https://www.unesco.org/en/articles/">https://www.unesco.org/en/articles/</a>	2025. <a href="#">Culture is not trivia: Sociocultural theory</a>	1165

1166 for cultural NLP. In *Proceedings of the 63rd Annual*  
1167 *Meeting of the Association for Computational*  
1168 *Linguistics (Volume 1: Long Papers)*, pages 25869–  
1169 25886, Vienna, Austria. Association for Computa-  
1170 tional Linguistics.