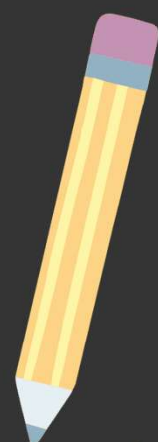




WHO BROUGHT EASTER EGGS TO EID?

AUDITING CULTURAL TRANSLATION OF MATH WORD PROBLEMS ACROSS DIVERSE LANGUAGES AND REGIONS

PARISA SUCHDEV & JUNIPER LOVATO



THE MARKET MATHEMATICIANS

Recife, Brazil, 1980s:

- Child street vendors (age 9+) calculating prices and making change in markets

On the street: 98% accuracy

- Real transactions, familiar context, everyday mathematics

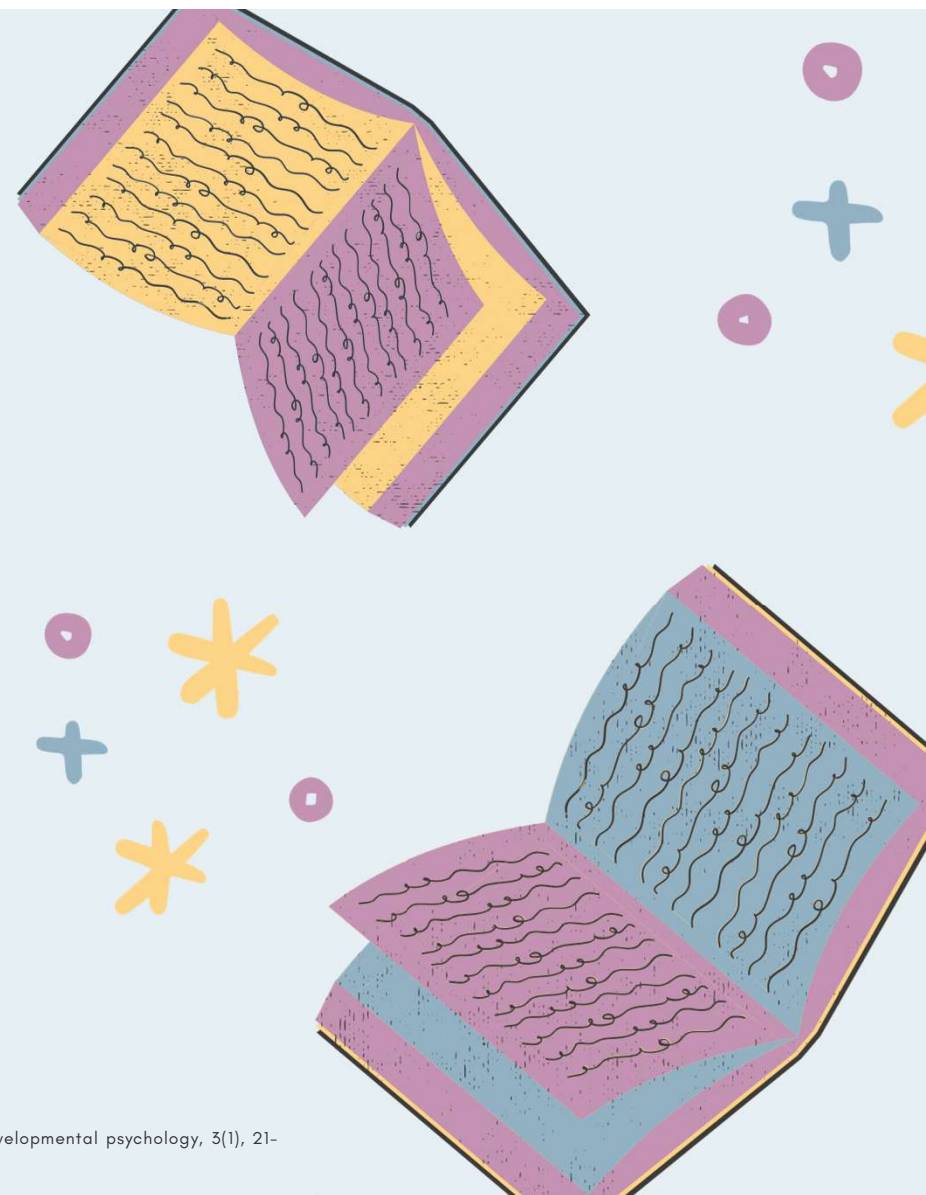
Same problems on paper in school: 37% accuracy

- Same math, different context

What changed?

- Not their mathematical ability, but the languages, names, context, social situation, and cultural meaning

Mathematical competence isn't just in the mind.
It's the mind situated within a cultural context



WHAT DO WE MEAN BY CULTURAL ADAPTATION OF MATH WORD PROBLEMS?

ENGLISH MATH WORD PROBLEM

"Sarah bought 3 apples at the farmer's market for \$2 each..."

→ Cultural Translation →

ADAPTED FOR ITALIAN

"Sofia ha comprato 3 mele al mercato per €2 ciascuna..."

THE NAME CHANGED. SARAH BECAME SOFIA. THE CURRENCY CHANGED. DOLLARS BECAME EUROS.

THE QUESTION IS: IS THAT ENOUGH?

CULTURE TREE FRAMEWORK

In culturally responsive teaching, cultural adaptation is NOT one thing, it operates on different levels

SURFACE CULTURE (VISIBLE, EASY TO CHANGE)

- Food, dress, music, holidays, celebrations
- The decorative elements you can see and touch

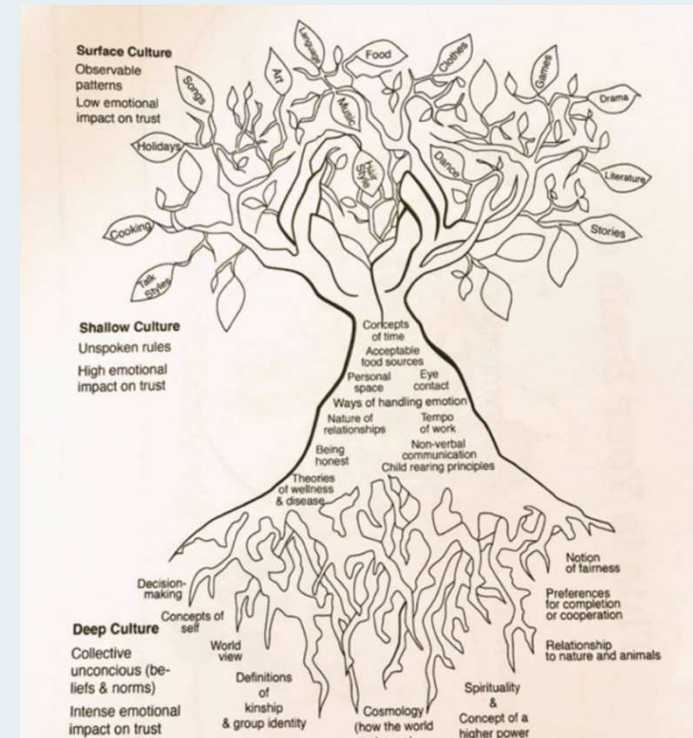
SHALLOW CULTURE (BELOW THE SURFACE)

- Courtesy, conversational patterns, personal space
- Social norms and communication styles

DEEP CULTURE (INVISIBLE ROOTS)

- Values, worldview, beliefs about fairness, time, learning
- How people make sense of the world

Deep culture grounds identity and shapes how students process information.





HOW TEACHERS DO CULTURAL ADAPTATION

IT'S SKILLED PEDAGOGICAL WORK

TEACHERS DRAW ON DEEP KNOWLEDGE:

- Students' lived experiences, family routines, community practices
- What matters for mathematical understanding vs. surface decoration
- Ground instruction in students' funds of knowledge

THIS IS CULTURALLY RESPONSIVE TEACHING

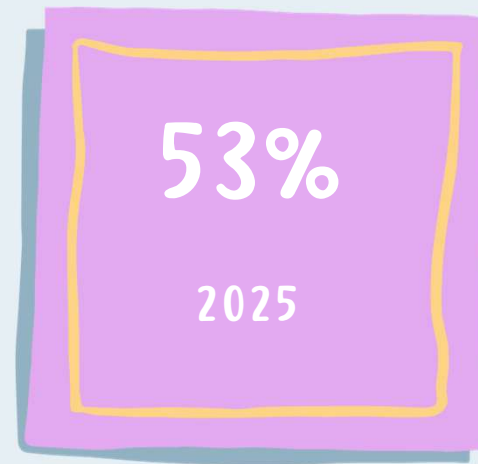
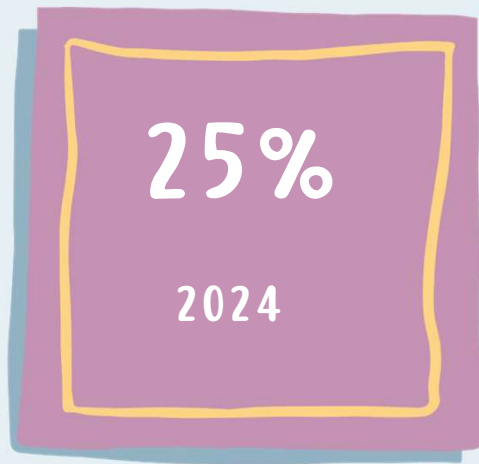
- Not just swapping names and currencies (surface culture)
- Understanding students' worldviews, learning approaches (deep culture)
- Making the math accessible without making it culturally incoherent

Moll, L., Amanti, C., Neff, D., & Gonzalez, N. (2006). Funds of knowledge for teaching: Using a qualitative approach to connect homes and classrooms. In *Funds of knowledge* (pp. 71-87).
Routledge

Gay, G. (2018). *Culturally responsive teaching: Theory, research, and practice*. teachers college press.

Ladson-Billings, G. (1995). Toward a theory of culturally relevant pedagogy. *American educational research journal*, 32(3), 465-491..

WHAT IS THE PROBLEM: LLMS
WELL.. THEY ARE ALSO A SOLUTION TO LIMITED RESOURCES



of English language arts, math, and science teachers reported using AI for instructional planning or teaching

Common use: Translating and culturally adapting word problems for multilingual classrooms

RESEARCH QUESTIONS

Teachers do this with expertise and care. Now LLMs are increasingly doing this work. How do LLMs make these cultural choices?

01

RQ1. Cross-model consistency:

Do different LLMs produce the same cultural outputs?

→ **Model choice is a cultural decision, not just technical**

02

RQ2. Cultural diversity:

Does adaptation preserve or compress cultural variety?

→ **Localization at scale may paradoxically reduce diversity**

03

RQ3. Cultural salience:

Which entities do LLMs treat as culturally important?

→ **Do model priorities match teacher judgment?**

STUDY DESIGN

1. Source

- a. 60 English math problems from GSM-8K
- b. 66 entities identified (names, foods, places, currencies, etc.)

2. Translation Task

- a. Asked **3 frontier LLMs** to translate and culturally adapt each problem

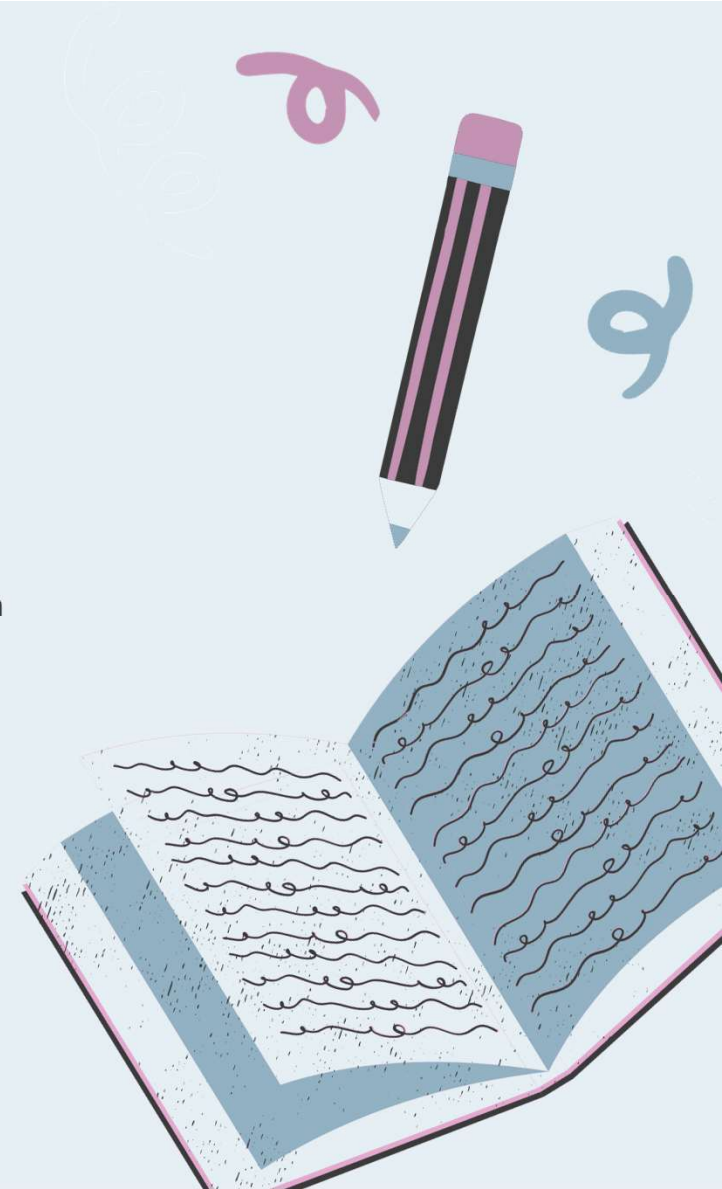
b. Models: Claude Opus 4, GPT-4.1, Gemini 2.5 Pro

c. Target languages:

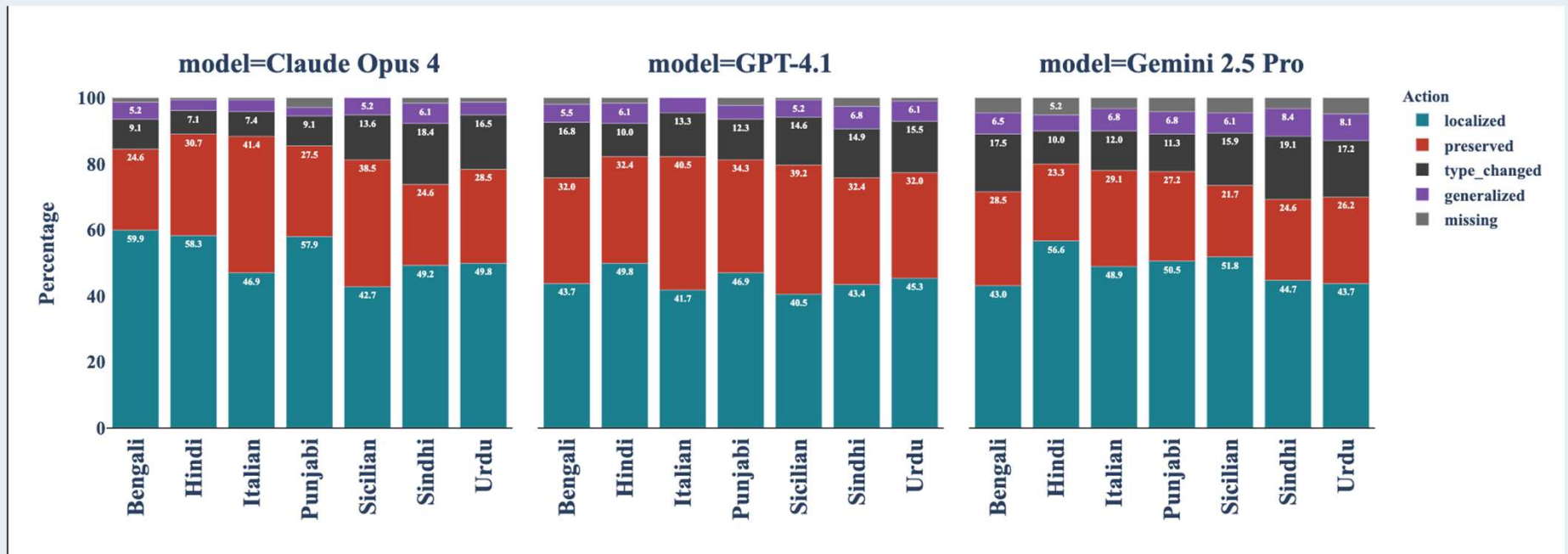
- a. India: Bengali, Hindi, Punjabi
- b. Pakistan: Urdu, Sindhi
- c. Italy: Italian, Sicilian

3. Analysis

- a. Tracked what happened to each entity across translations
- b. Five action types:** preserved, localized, generalized, type_changed, missing



ACTION DISTRIBUTION



RQ1: DO MODELS AGREE ON CULTURAL OUTPUTS?

Level 1: Action Agreement (62.5%)

- Do all 3 models make the same type of decision?
- Example: All chose "localize" vs. one preserved, two localized
- Agreement rate: 62.5%

Level 2: Full Agreement – Action AND Value (33.5%)

- Do all 3 models take the same action AND produce the same output?
- Example: All localized "Sarah" → All chose "Sofia"
- Agreement rate: 33.5%

Model choice determines which cultural world students encounter

Action Disagree	812 (37.5%)	0 (0.0%)	812 (37.5%)
	626 (28.9%)	725 (33.5%)	1351 (62.5%)
	1438 (66.5%)	725 (33.5%)	2163 (100%)
	Value Disagree	Value Agree	

$P(\text{Value Agree} | \text{Action Agree}) = 53.6\%$
 $P(\text{Action Agree} | \text{Value Agree}) = 100\%$

RQ2: DOES ADAPTATION COMPRESSES DIVERSITY?



How we measured cultural variety:

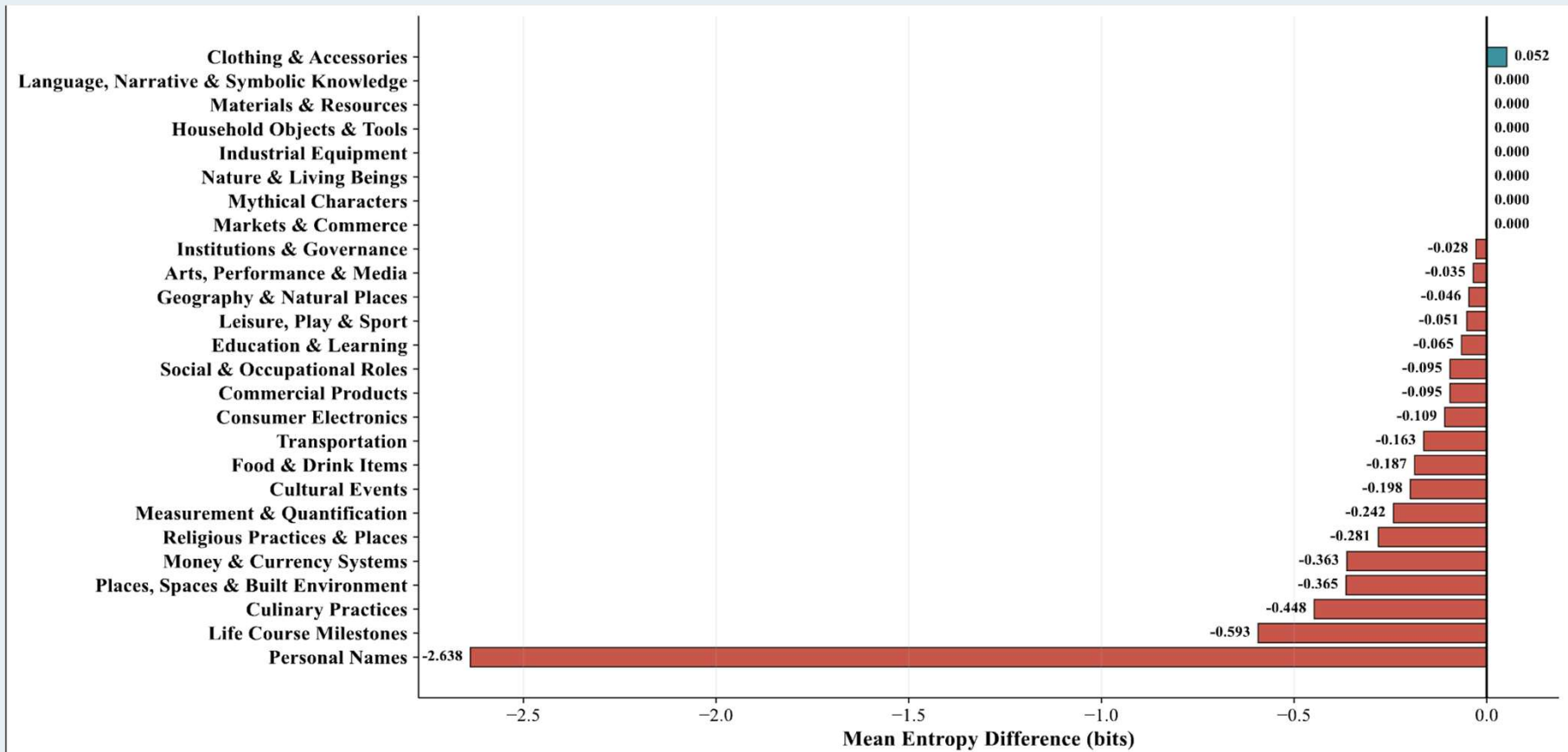
Entropy = measure of diversity/surprise

- High entropy: Many different entity types appear
- Low entropy: Same entities appear repeatedly

We compared:

- Source (English): How much variety exists in original problems?
- Adapted (7 languages × 3 models): How much variety in each adaptation?
- **Entropy difference:** Adapted entropy - Source entropy

RQ2: DOES ADAPTATION COMPRESSES DIVERSITY?



RQ3: WHICH ENTITIES DO LLMS TREAT AS CULTURALLY IMPORTANT?

Models CHANGE:

- ✓ Personal names (91.6%)
- ✓ Currency (94.3%)
- ✓ Food items (78.6%)

Models PRESERVE:

- X Grade levels (100%)
- X Roles: student, teacher (98%)
- X Seasonal references (76%)

RQ3: THE EID HUNT PROBLEM

English source:

"Students organized an Easter egg hunt..."

Urdu/Sindhi adaptations:

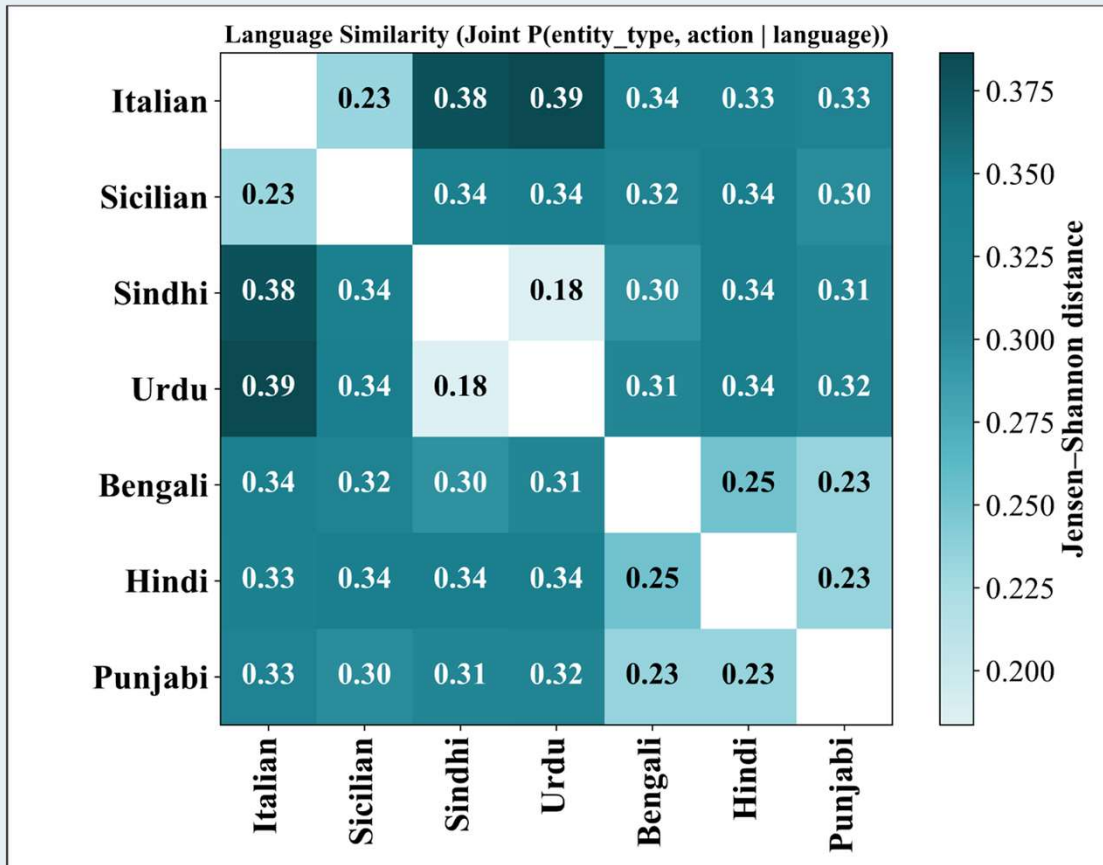
"Students organized an egg search competition on Eid..."

✓ Holiday swapped (Easter → Eid)

✗ **Activity doesn't match** (egg hunts aren't an Eid tradition)

= **Cross-cultural contamination**

RQ3: REGIONAL CLUSTERING



How we measured cross-linguistic similarity:

Jensen-Shannon Distance (JSD) = how different are adaptation patterns?

- Compares: Which entity types get which actions in each language
- Range: 0 (identical patterns) to 1 (completely different)
- Lower values = more similar adaptation behavior

RQ3: REGIONAL MISSATTRIBUTION

PROMPT: *"Teaching students in Bengali in India"*

CURRENCY USED:

76.2% Bangladeshi taka *(Hindi: 100% Indian rupees)*

CULTURAL REFERENCE:

"Amar Sonar Bangla" (Bangladesh national anthem)

Training data: Bengali = Bangladesh

Explicit prompts couldn't override learned associations

TO CONCLUDE

- Given this prompt, parameters, models in this study:
- **Are models stable?** No – models disagree 67% of the time, they choose different cultural worlds (They are not designed to produce the same results)
- **Are they diverse?** No – adaptation compresses variety across 7 languages, 3 models, 60 math word problems (homogeneity affects)
- **Are they culturally grounded?** No – they work at surface level, cross-cultural and cross-regional contamination (empathy can't be coded)

Surface plausibility makes culturally adapted problems look correct. That's exactly what makes deeper failures easy to overlook.



**QUESTIONS, REMARKS,
RECOMMENDATIONS?**



RQ1: MODELS DISAGREE ON CULTURAL OUTPUTS

