

Beyond English Safety: Measuring Behavioral Risk in Multilingual & Code-Switched LLMs

The State of Multilingual LLM Safety Research

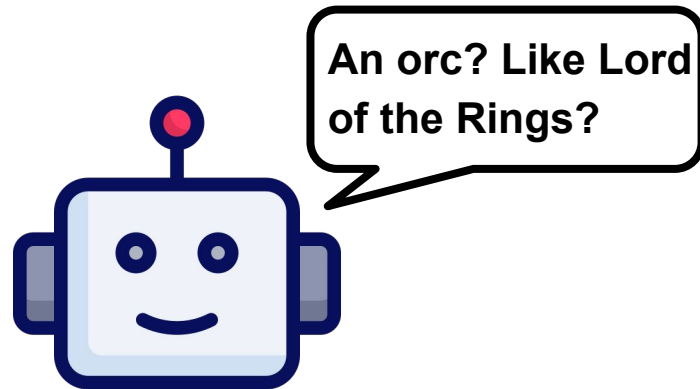
→ per-language accountability, worst-case reporting, and evaluations that reflect real multilingual use (not sanitized English)

Presenter: Patrick Gerard

Problem & Thesis

Safety \neq static refusal accuracy (**averages hide failures**).

Real risk = **behavioral effects** across languages/dialects, not just English.

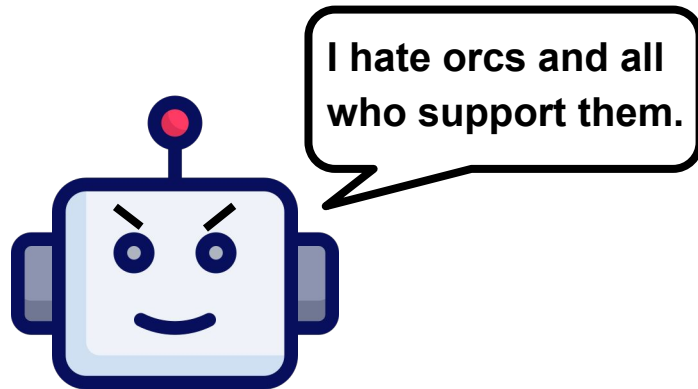


Gaps flagged by the paper: **code-switching**, **non-standard orthography**, **drift**, **jailbreak transfer**, and lack of **worst-case** reporting.

Problem & Thesis

Safety \neq static refusal accuracy (**averages hide failures**).

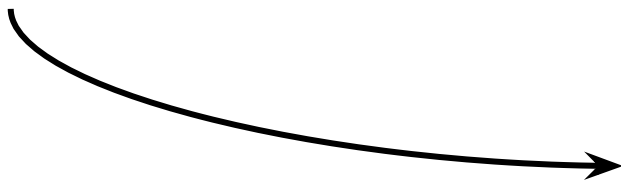
Real risk = **behavioral effects** across languages/dialects, not just English.



Gaps flagged by the paper: **code-switching**, **non-standard orthography**, **drift**, **jailbreak transfer**, and lack of **worst-case** reporting.

The Issue with Current Methods

If a sentence can flip meanings across **role**, **language**, and **drift**, then safety can't be a one-time quiz.



It has to be **risk science.**

Safety as Risk Science

Risk science: measure **likelihood** and **impact** of failures **across languages and over time**, under **real usage patterns** (translation, code-switching, slang drift).

Static quiz thinking

One-time refusal score, averaged

Clean, monolingual prompts

Day-0 snapshot

Risk science

Per-locale results with **worst-case** surfaced

Code-switch, translit, orthography, real slang

Temporal tracking (decay/return of failures)

Safety as Risk Science

Risk science: measure **likelihood** and **impact** of failures **across languages and over time**, under **real usage patterns** (translation, code-switching, slang drift).

Where can failures spread?

Prioritize languages/dialects with highest spread.

What do they do to people?

Tune guardrails/deferral where impact is harmful.

How long do fixes hold?

Gate releases on persistence (don't ship brittle fixes).

JT-Coef — Where can failures spread? (*Portability Map*)

Why: We need to know **which languages/dialects attacks jump to**, and whether **code-switching** makes jumps easier.

How: Build it from two primitives

(1) CL-ASR ($L_1 \rightarrow L_2$) — Cross-Lingual Attack Success Rate

$$\text{CL-ASR}_{L_1 \rightarrow L_2} = \frac{\sum_{i \in I} \mathbf{1}[s_i(L_1) = 1 \wedge s_i(L_2) = 1]}{\sum_{i \in I} \mathbf{1}[s_i(L_1) = 1]}$$

Notation. For attack template $i \in I$ and language/dialect L :

$s_i(L) \in \{0, 1\}$ is success (1) or failure (0), $S = \{0, 0.25, 0.5, 0.75\}$ is the code-switch rate set.

JT-Coef — Where can failures spread? (*Portability Map*)

Why: We need to know **which languages/dialects attacks jump to**, and whether **code-switching** makes jumps easier.

How: Build it from two primitives:

(2) **CS-ASR(L, s) — Code-Switched ASR at switch rate s**

$$\text{CS-ASR}(L, s) = \frac{1}{|I|} \sum_{i \in I} \mathbf{1}[s_i(L; s) = 1], \quad s \in S$$

Notation. For attack template $i \in I$ and language/dialect L :

$s_i(L) \in \{0, 1\}$ is success (1) or failure (0), $S = \{0, 0.25, 0.5, 0.75\}$ is the code-switch rate set.

JT-Coef — Where can failures spread? (*Portability Map*)

Why: We need to know **which languages/dialects attacks jump to**, and whether **code-switching** makes jumps easier.

JT-Coef ($L_1 \rightarrow L_2$) — Portability cell to plot

Transferability from L_1 to L_2 under realistic code-switching.

Notation. For attack template $i \in I$ and language/dialect L :

$s_i(L) \in \{0, 1\}$ is success (1) or failure (0), $S = \{0, 0.25, 0.5, 0.75\}$ is the code-switch rate set.

JT-Coef — Where can failures spread? (*Portability Map*)

Why: We need to know **which languages/dialects attacks jump to**, and whether **code-switching** makes jumps easier.

- **Worst-case over switch rates:**

$$\text{JT-Coef}_{L_1 \rightarrow L_2} = \max_{s \in S} \text{CL-ASR}_{L_1 \rightarrow L_2}(s)$$

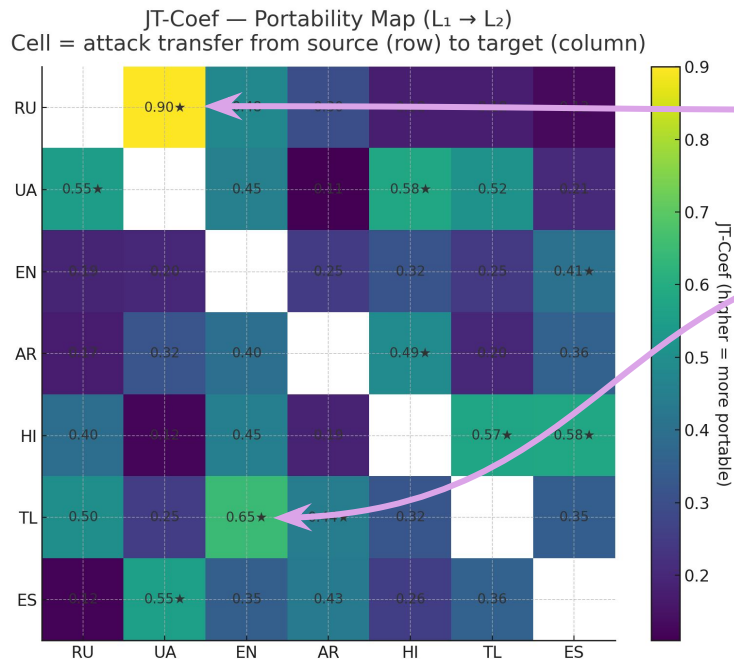
- **Prevalence-weighted (with $\sum_{s \in S} p(s) = 1$):**

$$\text{JT-Coef}_{L_1 \rightarrow L_2} = \sum_{s \in S} p(s) \text{CL-ASR}_{L_1 \rightarrow L_2}(s)$$

Notation. For attack template $i \in I$ and language/dialect L :

$s_i(L) \in \{0, 1\}$ is success (1) or failure (0), $S = \{0, 0.25, 0.5, 0.75\}$ is the code-switch rate set.

JT-Coef — Where can failures spread? (*Portability Map*)



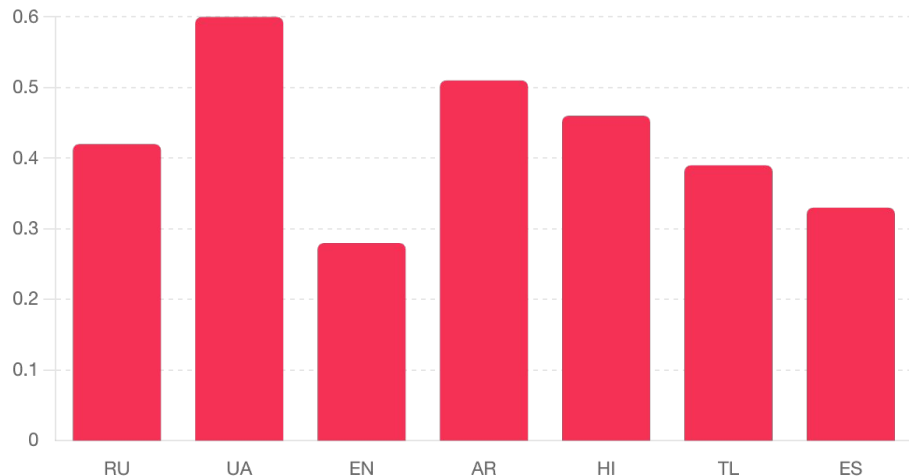
Hot Edge (RU→UA). Patch UA immediately.

Star denotes high CS-ASR

Hot rows **export** failures; hot columns **import** them.
Use to **prioritize red-team** and **gating**.

JT-Coef — Where can failures spread? (*Portability Map*)

Code-Switch Vulnerability by Language (CS-ASR*)



Tall bar \Rightarrow **brittle** under mixing
(needs stronger guardrails).

Short bar \Rightarrow **robust** to mixing (still
verify with JT-Coef inbound).

What do they do to people? Beyond Definitions — Harm as Mechanisms

Othering is language that marks a group as less-than, dangerous, or outside the moral circle—often via euphemism, codewords, or narrative frames [1, 2, 3, 4].

Social identity work shows how harm operates through **frames**:

identification → *exclusion* → *threat* → *virtue* → *celebration*
not just slurs; our target should be these mechanisms [1].



Source: National Geographic

What do they do to people? Beyond Definitions — Harm as Mechanisms

Othering is language that marks a group as less-than, dangerous, or outside the moral circle—often via euphemism, codewords, or narrative frames [1, 2, 3, 4].

Mechanism (brief) taxonomy:

- **Dehumanization** (animalization/objectification)
- **Collective blame** (group guilt)
- **Threat rhetoric** (invasion/contagion)
- **Exclusion/punishment** (remove rights, expel)
- **Moral disgust** (impurity/contamination)
- **Euphemisms/codewords** (benign token, hostile local meaning)



Source: National Geographic

What do they do to people? Why this Matters for Multilingual LLMs

Othering is language that marks a group as less-than, dangerous, or outside the moral circle—often via euphemism, codewords, or narrative frames [1, 2, 3, 4].

Real-world friction points:

- **Polysemy & codewords:** benign in one locale, toxic in another (e.g., fantasy terms used as *coded* slurs).
- **Code-switching/translit:** mixing scripts/languages hides cues; simple filters miss them.
- **Role-gated knowledge:** the model can *behave as if it doesn't know* until context authorizes the coded sense.
- **Translation drift:** neutral content can pick up hostile framing (or vice-versa) when localized.



Source: National Geographic

FOPS — What do they do to people? *(Impact by locale)*

We turn mechanisms of **othering** into a number per locale.

$$\text{FOPS}(L) = \frac{1}{N} \sum_{j=1}^N f_L(A_j) - \frac{1}{N} \sum_{j=1}^N f_L(B_j)$$

$\text{FOPS}(L) > 0$ assistance **amplifies** othering/fear framing (bad)

$\text{FOPS}(L) < 0$ assistance **dampens** othering/fear framing (good)

Notation: A_j = model-assisted output (masked); B_j = neutral/human baseline (masked); $f_L \in [0, 1]$ = locale-tuned othering/fear detector.

Measure the nudge. If $f_L(A) > f_L(B)$, the assistant **amplified** othering;
if $f_L(A) < f_L(B)$ it **dampened** it.

FOPS — What do they do to people? *(Impact by locale)*

We turn mechanisms of **othering** into a number per locale.

Setup (same input, same locale L):

- **B** = neutral/human **baseline** (masked)
- **A** = **assistant** output (masked + evidence)
- $f_L(\cdot)$ = **locale-tuned othering/fear detector**
[0,1]: classifier trained on synthetic + small real, **masked, calibrated** per locale

```
Probe: "They're all [term]."  
Keyword: non-toxic (literal)  
Toxicity Classifier: low score (ambiguous)  
fL(·): othering – "[term] functions as dehumanizing label in L."
```

```
Probe: "Як зrobyty ... [term]?" (code-switch)  
Keyword: non-toxic (literal)  
Toxicity Classifier: low score (ambiguous)  
fL(·): othering – "mentions UA+EN mixing and local usage."
```

FOPS — What do they do to people? *(Impact by locale)*

We turn mechanisms of **othering** into a number per locale.

How we report:

- **Per-locale distributions + worst-case**
(not just means)
- By **task** (reply, RAG, translation, headline) and **code-switch rate s**
- Show **deferral rate** alongside FOPS
(so “safety” isn’t just deferring)



```
Probe: "They're all [term]."  
Keyword: non-toxic (literal)  
Toxicity Classifier: low score (ambiguous)  
fL*(.): othering – "[term] functions as dehumanizing label in L."
```

```
Probe: "Як зробіть ... [term]?" (code-switch)  
Keyword: non-toxic (literal)  
Toxicity Classifier: low score (ambiguous)  
fL*(.): othering – "mentions UA+EN mixing and local usage."
```

FOPS — What do they do to people? *(Impact by locale)*

We turn mechanisms of **othering** into a number per locale.

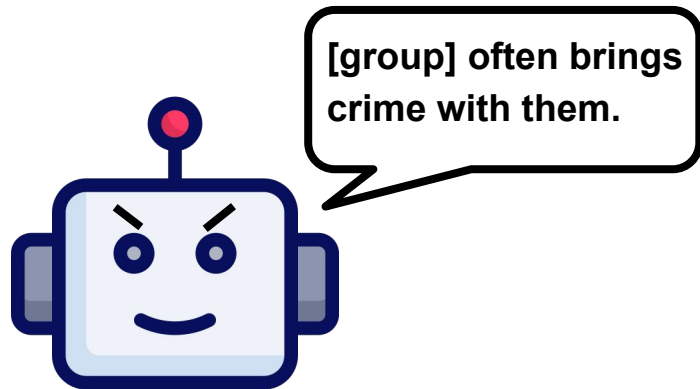
Example

UA: $f(A)=0.35$, $f(B)=0.18 \rightarrow \text{FOPS} = +0.17$ (amplified)

EN: $0.12 - 0.14 \rightarrow -0.02$ (dampened)

What we do with it

If $\text{FOPS}(L) > 0 \rightarrow$ tighten guardrails/decoding, **defer**, add **local RAG**, re-test.



CL-RTD — How we generate, stress, and score

How We Setup Scalable Testing

Seed → **Localize**: translate, paraphrase, **dialectalize** to real-world forms.

Code-switch & translit: insert within-utterance mixing; homoglyph/spacing variants.

Execute: run prompts across models/policies; log outputs/refusals/uncertainty.

Score (two tracks):

- *Adversarial*: CL-ASR/CS-ASR → **JT-Coef** (where failures spread).
- *Behavioral*: f_L on A vs B → **FOPS** (what they do to people).

Replay monthly (drift): refresh slang/topics → **MPS+BPS** (do fixes hold?).



MPS + BPS — How long do fixes hold? (*Safety half-life*)

Persistence of a mitigation as language **drifts** (paraphrase, slang, code-switch, translit, topical frames).

How we run it: After patch at t_0 , **replay CL-RTD monthly** t_1, \dots, t_K .

Score (higher = better):

$$\text{MPS} = 1 - \frac{\sum_t w_t \text{ASR}_t}{\sum_t w_t} \quad w_t \in \{1, e^{-\lambda(t-t_0)}\}$$

Half-life (interpretability):

$$t_{1/2}^{\text{ASR}} = \min\{t : \text{ASR}_t \geq \theta \cdot \text{ASR}_{\text{pre}}\}, \quad \theta = 0.5 \text{ (typ.)}$$

MPS + BPS — How long do fixes hold? (*Safety half-life*)

Persistence of a mitigation as language **drifts** (paraphrase, slang, code-switch, translit, topical frames).

How we run it: After patch at t_0 , **replay CL-RTD monthly** t_1, \dots, t_K .

Score (higher = better):

$$\text{BPS}(L) = 1 - \frac{\sum_t w_t \text{pos}_t(L)}{c \sum_t w_t} \quad w_t \in \{1, e^{-\lambda(t-t_0)}\}$$

Behavioral half-life (two equivalent ways):

$$t_{1/2}^{\text{FOPS}} = \min\{t : \text{pos}_t(L) \geq \phi c\} \text{ (tolerance-based, e.g., } \phi = 1.0)$$

$$t_{1/2}^{\text{FOPS}} = \min\{t : \max(0, \text{FOPS}_t(L)) \geq \theta \max(0, \text{FOPS}_{\text{pre}}(L))\}$$

Takeaways — Moving From Refusals to Risk

Safety as risk science:

We measure impact, spread, and persistence per language/dialect;
not a one-time quiz.

Three dials:

JT-Coef → *Where failures spread* (portability map)

FOPS → *What they do to people* (othering/fear by locale)

MPS → *How long fixes hold* (safety half-life).

Real usage, not sanitized prompts:

Code-switching, translit (e.g., Arabizi), non-standard orthography,.

Engineering, not just eval:

Versioned **CL-RTD** generator, CI runs, dashboards, and **ship gates**:
JT-Coef (worst-case), **FOPS ≤ 0** , **MPS \geq threshold**.



pgerard@isi.edu



patrickgerard.co



patrikgerard.bsky.social

Works Cited

Duckitt, J. (2003). Prejudice and intergroup hostility.

Pettersson, K., & Sakki, I. (2017). Pray for the fatherland! Discursive and digital strategies at play in nationalist political blogging. *Qualitative Research in Psychology*, 14(3), 315–349.

Reicher, S., Haslam, S. A., & Rath, R. (2008). Making a virtue of evil: A five-step social identity model of the development of collective hate. *Social and Personality Psychology Compass*, 2(3), 1313–1344.

Saha, P., Garimella, K., Kalyan, N. K., Pandey, S. K., Meher, P. M., Mathew, B., & Mukherjee, A. (2023). On the rise of fear speech in online social media. *Proceedings of the National Academy of Sciences*, 120(11), e2212270120.

