

Post-Training Language Models for Crosslingual Consistency

Tianyu Liu* **Jirui Qi*** **Mrinmaya Sachan** **Ryan Cotterell**
Raquel Fernández **Arianna Bisazza**

ICML 2026



The same question, different answers

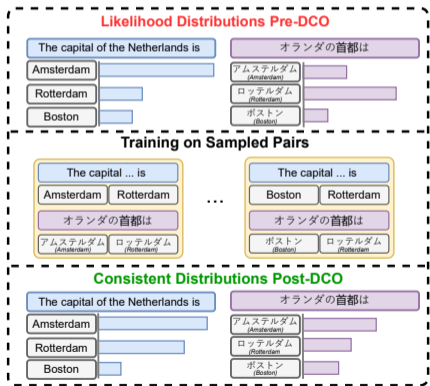
Same query, two languages → different answers

EN “What is the capital of Australia?” ⇒ Canberra ✓

SW “Mji mkuu wa Australia ni upi?” ⇒ Sydney ✗

- LLMs answer translation-equivalent prompts **inconsistently**.
- Undermines **reliability** and **equity** across languages — **low-resource** speakers get worse, unstable answers.

Goal. Post-train so the model answers the *same*, regardless of the prompt’s language.



Inconsistent rankings before alignment → aligned after.

Defining crosslingual consistency

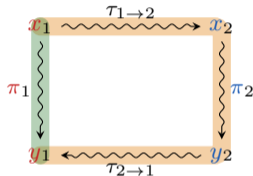
Two routes to a response distribution for a prompt in L_1 :

- **Direct**: sample from $\pi_1(\cdot | x_1)$.
- **Round trip**: translate to L_2 , respond, translate back: $\tau_{2 \rightarrow 1} \# \pi_2 \# \tau_{1 \rightarrow 2}$.

Consistent \Leftrightarrow the two agree, within an f -divergence tolerance (up to a temperature T):

$$D_f\left(\pi_1(\cdot | x_1) \parallel (\tau_{2 \rightarrow 1} \# \pi_2 \# \tau_{1 \rightarrow 2})^T(\cdot | x_1)\right) \leq \varepsilon$$

Zero divergence = the diagram commutes.



Direct vs **round trip** — both should land on the same y_1 .

An objective: Penalized Consistency Optimization (PCO)

Per language, trade off consistency against drift from a fixed reference:

$$\mathcal{L}^{\text{PCO}}(\theta) = \underbrace{\text{KL}(\pi_\theta \parallel \pi_{\text{ref}})}_{\text{fidelity}} - \underbrace{\beta \log(\tau_{2 \rightarrow 1} \# \pi_{\text{ref}} \# \tau_{1 \rightarrow 2})}_{\text{consistency reward}}$$

- **Closed-form optimum:** a weighted geometric mean

$$\pi^*(\cdot \mid \mathbf{x}_1) \propto (\tau_{2 \rightarrow 1} \# \pi_{\text{ref}} \# \tau_{1 \rightarrow 2})^{\beta_1}(\cdot \mid \mathbf{x}_1) \cdot \pi_{\text{ref}}(\cdot \mid \mathbf{x}_1).$$

- With invertible translators and balanced strengths ($\beta_1 \beta_2 = 1$), π^* is **exactly crosslingually consistent**.

From PCO to DCO: drop the roll-outs

Problem. The consistency term is an expectation under $\pi_\theta \Rightarrow$ optimizing PCO directly needs **on-policy roll-outs** (expensive, high-variance).

Idea (inspired by DPO). We know π^* in closed form — so just **regress the logits** onto it:

$$\mathcal{L}^{\text{DCO}}(\theta) = \left\| \ell_\theta(y_1 | x_1) - \beta_1 \log(\tau_{2 \rightarrow 1} \# \pi_{\text{ref}} \# \tau_{1 \rightarrow 2})(y_1 | x_1) - \log \pi_{\text{ref}}(y_1 | x_1) \right\|_1$$

- Same unique minimizer as PCO — but fully **off-policy**.
- Needs only *offline* evaluation on parallel data.
- Deterministic translators \Rightarrow round trip = one eval of π_{ref} on translated input–output pairs. **No human labels.**

Experimental setup

Models. 9 LMs, 4 families, 3B–14B
(Qwen2.5/3, Llama3, Gemma3, Aya-Expansive).

Data. 26 languages, 3 benchmarks:

- MMMLU — general knowledge
- XCSQA — commonsense reasoning
- BMLAMA — factual recall

Metrics.

- **CLC** (RankC): crosslingual consistency across language pairs.
- Accuracy on **English** and **non-English**.

Baselines. SFT, DPO (need labels); CALM (label-free).

Results: DCO improves consistency everywhere

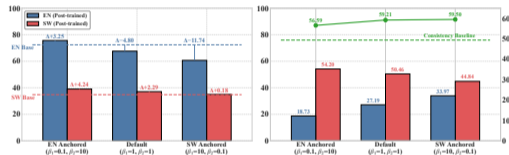
MMMLU, joint N-lang	Qwen2.5-14B	Llama3.1-8B	Aya-Expansive-8B
Base (CLC)	68.6	60.9	72.2
+ CALM (label-free)	+4.2	+3.0	+1.4
+ DCO (label-free)	+10.6	+9.4	+5.3
+ DPO (labels)	+12.3	+10.1	+1.3
+ DPO → DCO (labels)	+13.1	+13.8	+3.1

- Beats CALM; **matches/exceeds DPO without labels**; DPO→DCO best.
- Accuracy preserved: non-English \uparrow , English stable. Largest CLC gains on factual BMLAMA (+12 to +17).
- **Out-of-domain**: train 1 subject $\rightarrow \sim +11$ CLC on unseen subjects.
- *Why accuracy rises*: geometric-mean ensemble down-weights uncertain, high-entropy (often wrong) experts.

Steering alignment toward low-resource languages

Strength parameters (β_1, β_2) with $\beta_1\beta_2 = 1$ are a **direction knob**.

- **English-anchored** (small β_{en}): keep **English** fixed, let the **low-resource** side move.
- **English–Swahili**: **+4.2** Swahili accuracy, **English** stable.
- Revises 54% of **Swahili** responses vs. only 19% of **English**.



English–Swahili: accuracy (left), response change & CLC (right).

Target **low-resource** languages *without* degrading the dominant one.

Takeaways

1. A precise, information-theoretic **definition** of crosslingual consistency (round-trip pushforward).
2. **PCO**: an objective whose closed-form optimum *is* consistent.
3. **DCO**: a cheap, **off-policy, label-free** surrogate with the same optimum — just regress the logits.

More consistent *and* more equitable multilingual models.



Thank you!