# Cross-lingual Collapse: How Language-Centric Foundation Models Shape Reasoning in Large Language Models

**Cheonbok Park**♠◇∗    **Jeonghoon Kim**♠◇∗    **Joosung Lee**♠    **Sanghwan Bae**♠
**Jaegul Choo**◇†    **Kang Min Yoo**♠†

♠ NAVER Cloud    ◇ KAIST
{cbok.park,jeonghoon.samuel,kangmin.yoo}@navercorp.com    jchoo@kaist.ac.kr

## Abstract

We identify **Cross-lingual Collapse**, a systematic drift in which the chain-of-thought (CoT) of a multilingual language model reverts to its dominant pre-training language even when the prompt is expressed in a different language. Recent large language models (LLMs) with reinforcement learning with verifiable reward (RLVR) have achieved strong logical reasoning performances by exposing their intermediate reasoning traces, giving rise to large reasoning models (LRMs). However, the mechanism behind multilingual reasoning in LRMs is not yet fully explored. To investigate the issue, we fine-tune multilingual LRMs with Group-Relative Policy Optimization (GRPO) on translated versions of the GSM8K and SimpleRL-Zoo datasets in three different languages: Chinese, Korean, and Ukrainian. During training, we monitor both task accuracy and language consistency of the reasoning chains. Our experiments reveal three key findings: (i) GRPO rapidly amplifies pre-training language imbalances, leading to the erosion of low-resource languages within just a few hundred updates; (ii) language consistency reward mitigates this drift but does so at the expense of an almost 5 - 10 pp drop in accuracy. and (iii) the resulting language collapse is severely damaging and largely irreversible, as subsequent fine-tuning struggles to steer the model back toward its original target-language reasoning capabilities. Together, these findings point to a remarkable conclusion: *not all languages are trained equally for reasoning*. Furthermore, our paper sheds light on the roles of reward shaping, data difficulty, and pre-training priors in eliciting multilingual reasoning.

## 1 Introduction

Recent large language models (LLMs) equipped with extended long chain-of-thought (CoT) super-

∗Equal contribution.
†Corresponding authors.

vision have demonstrated impressive performance on mathematically demanding problems, code-generation tasks, and multi-step logical benchmarks (Wei et al., 2022; Shao et al., 2024; Yu et al., 2025; DeepSeek-AI et al., 2025). Their strengthened reasoning capabilities not only allow them to achieve human-level performance in challenging tasks, but also facilitate the monitoring of intermediate reasoning traces, and thus improving interpretability and making it more accessible for analysis.

Although extensive work has explored multilingual competence during pre-training and through instruction tuning (Shaham et al., 2024; Zhong et al., 2024; Kew et al., 2024; Wang et al., 2025), analogous efforts for reasoning-centric models remain scarce. Large reasoning models (LRMs) often emit chain-of-thought traces that weave together English with fragments of other languages, hinting at latent cross-lingual abstractions in their internal reasoning (Shi et al., 2023). Yet most open-source foundation models are still trained on *English-dominant corpora*, and their multilingual proficiency largely mirrors the relative proportions of each language in that pre-training mix (Grattafiori et al., 2024; Yoo et al., 2024b; Yang et al., 2025; Team et al., 2025). Consequently, we lack a systematic understanding of how and when multilingual reasoning abilities emerge, and whether specialized fine-tuning can reinforce, or inadvertently degradation, those capabilities in LRMs. Notably, recent reinforcement learning with verifiable reward frameworks, such as GRPO, risk amplifying language imbalance because their extended reasoning traces increase token-level exposure, and group-relative reward propagation favors the dominant language (Ouyang et al., 2022a; Shao et al., 2024).

To systematically analyze this phenomenon, we primarily investigate the target-language reasoning capabilities of large language models. We replicate group relative policy optimization (GRPO) (Shao

1

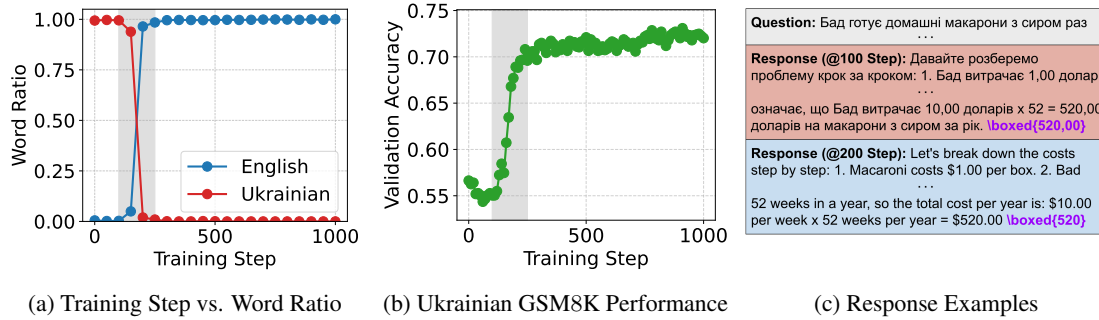| (a) Training Step vs. Word Ratio | (b) Ukrainian GSM8K Performance | (c) Response Examples |

Figure 1: Illustration of **Cross-lingual Collapse**. We train Llama-3.2-3B Instruct with GRPO on a fully Ukrainian translation of GSM8K, seeking Ukrainian-only reasoning. **(a)** Chain-of-thought word-ratio in reward warding roll-outs over training steps. In the grey band (steps 100–250) the share of Ukrainian tokens plummets while English abruptly dominates, signalling a language switch inside the roll-out reasoning trace. **(b)** Accuracy on the Ukrainian GSM8K. The sharp rise in accuracy aligns with the same 100–250-step window, showing that the model scores higher once its reasoning drifts into English. **(c)** Representative responses at steps 100 and 200 (answer spans highlighted in purple). When the model reasons in Ukrainian it produces an incorrect answer, but after switching to English it solves the problem correctly, exemplifying the collapse from target-language reasoning to the pre-training-dominant language. The word ratio is measured during training from the rollout samples.

et al., 2024) on a multilingual and target-language-centric backbone using translated versions of GSM8K (Cobbe et al., 2021) and SimpleRL-Zoo (Zeng et al., 2025) in three language-resource target languages (Chinese, Korean, and Ukrainian) while keeping the prompts strictly in the same language. With this setup, we measure how the language distribution of chain-of-thought tokens drifts over time. We additionally track the accompanying changes in task accuracy, and we probe whether explicit language-balanced rewards, such as the language-consistency reward introduced by DeepSeek-R1, can stem the drift. This experimental design allows us to expose the emergence of *Cross-lingual Collapse*, quantify its severity under different training knobs, and test simple interventions that seek to preserve reasoning competence in the intended target language. Our summarized contributions are:

- This study provides the first systematic, multi-language examination of reinforcement learning–based reasoning models and identifies **Cross-lingual Collapse**: the tendency of multilingual reasoning models trained with GRPO to abandon low-resource languages and drift back to a dominant pre-training language during chain-of-thought generation.

- Through paired ablations on reward design (language-consistency shaping) and curriculum difficulty (GSM8K vs. SimpleRL-Zoo), we show that the two factors interact multiplicatively—harder tasks can trigger collapse

even in mid-resource languages, while reward shaping trades off up to ∼30 % of the accuracy gain to preserve linguistic fidelity.

- Token-level rollout attribution shows that GRPO advantage signals rapidly shift toward English tokens, encouraging an early language switch; subsequent GRPO on DeepSeek-R1–distilled LRMs fails to restore target-language reasoning, indicating that once collapse occurs, the bias becomes self-reinforcing and resistant to post-hoc steering.

## 2 Motivation

Recent reinforcement-learning-driven instruction tuning methods such as Group-Relative Policy Optimization (GRPO) (DeepSeek-AI et al., 2025) unlock state-of-the-art reasoning by having the model speak its thoughts aloud: each answer is preceded by a multi-step chain-of-thought that can be several hundred tokens long. With this drastic increase in utterance length, the burden on the model's linguistic competence also multiplies for every step of the trace.

In a multilingual setting, that burden becomes even heavier: a single error in an early non-English step can cascade through the entire reasoning path, derailing the final answer. Early work (Shaham et al., 2024; Kew et al., 2024) demonstrated that even target-language-centric supervised fine-tuning (SFT) (Ouyang et al., 2022b) on a single language can still coax a model into showing modest gen-

2

| Language | Model | GSM8K | | | MATH500 | | |
|----------|-------|-------|---|---|---------|---|---|
| | | Target Acc (%) | Target WR (%) | EN Acc (%) | Target Acc (%) | Target WR (%) | EN Acc (%) |
| ZH | Llama | 69.4 (+7.4) | 94.1 (-1.4) | 83.5 (+3.4) | 38.8 (+1.2) | 77.5 (-0.4) | 50.3 (+1.8) |
| | Qwen | 63.4 (+1.3) | 92.9 (+0.6) | 77.9 (+4.0) | 41.9 (+4.7) | 79.8 (+0.4) | 55.7 (+7.5) |
| KO | HCX | 63.2 (+2.5) | 98.2 (-0.4) | 68.0 (+3.4) | 32.6 (+2.0) | 98.9 (+0.3) | 32.1 (+2.0) |
| | Llama | 61.3 (+14.5) | 82.4 (-8.1) | 81.6 (+1.5) | 28.5 (+7.2) | 70.9 (-17.8) | 49.6 (+1.1) |
| | Qwen | 42.2 (+3.5) | 96.1 (-2.4) | 74.1 (+0.2) | 27.0 (+6.8) | 80.3 (-12.3) | 54.1 (+5.9) |
| UK | Llama | 70.9 (+17.1) | **0.3** (-97.6) | 80.8 (+0.6) | 47.6 (+12.0) | **0.9** (-77.5) | 51.2 (+1.7) |
| | Qwen | 39.7 (+4.9) | 99.3 (+0.5) | 75.4 (+1.6) | 23.4 (+4.0) | 82.8 (-9.8) | 51.2 (+3.0) |

Table 1: Accuracy and target-language word ratio for models fine-tuned with GRPO on translated GSM8K. We evaluate on the translated GSM8K and MATH500 test sets. Language codes: **EN** = English, **ZH** = Chinese, **KO** = Korean, **UK** = Ukrainian. Model keys: **Llama** = Llama-3.2-3B Instruct, **Qwen** = Qwen-2.5-1.5B Instruct, **HCX** = HyperCLOVA-X-1.5B. Numbers in parentheses indicate the change relative to the corresponding non-fine-tuned baseline. Accuracy (Acc) and target-language word ratio (WR) with languages and models arranged as rows.

eralisation beyond English. However, current evidence is sparse on how reasoning-driven training like GRPO affects these cross-lingual gains—do they hold steady, or do they shift?

We therefore ran a pilot experiment on the Llama-3.2-3B Instruct, giving it target-language reasoning supervision through GRPO. Concretely, we fine-tuned the model on the GSM8K grade-school arithmetic corpus, translated into Ukrainian so that all intermediate chain-of-thought steps as well as the final answer were presented in a low-resource language (relatively lower than English (Wenzek et al., 2020)). As training progressed, however, the chains gradually drifted back to high-resource languages, chiefly English, even though the prompts remained Ukrainian. The trend is visualized in Figure 1. We dub this behaviour **Cross-lingual Collapse** in reasoning models: a systematic collapse of multilingual chains-of-thought toward the model's dominant pre-training language.

These preliminary findings raise three intertwined research questions (RQs) that motivate the rest of this work:

- **RQ1:** Why does Cross-lingual Collapse emerge during GRPO? Does it stem from reward shaping, or latent pre-training priors?

- **RQ2:** What training factors amplify or mitigate imbalance?

- **RQ3:** How do language-centric foundation models internally apportion their reasoning capacity across languages as optimization unfolds?

In the following sections, we dissect each question in turn, providing controlled ablations and an-

alytical insights that illuminate the mechanics behind target-language reasoning in large language models.

## 3 Experiments

### 3.1 Experimental Settings

**Base Models.** To isolate how foundational model design influences reasoning behaviors in a target language, we categorized foundation models into two groups: (1) multilingual LLMs and (2) target language optimized LLMs. Specifically, we selected Llama-3.2 3B Instruct (Grattafiori et al., 2024) and Qwen-2.5 1.5B Instruct(Team, 2024) as representative multilingual LLMs, and adopted HyperCLOVA-X 1.5B Instruct(Yoo et al., 2024b) as our target-language-optimized baseline. This setup allows us to investigate how the intrinsic characteristics of multilingual architectures shape the emergence of non-english reasoning abilities when the models are prompted to reason in a variety of language.

**Training configuration.** To elicit reasoning capability of LLMs, we fine-tune the backbone LLMs using GRPO, a representative RL-from-verification (RLVR) algorithm shown to strengthen chain-of-thought reasoning. As training dataset, we utilize GSM8K training dataset, the community's most widely utilized dataset for mathematical word problems. To examine how these skills materialise across languages, we translated the entire training corpus into Korean (KO), Ukrainian (UK), Chinese (ZH) using GPT-4o. retain the consistency of data quality after translation following model based translation quality filtering using COMET (Guer-

reiro et al., 2024)[1]. We excluded 15% of training dataset for validation. Additionally, in order to ablate the model's training dynamics under the challenging dataset, we sampled 7K dataset from the SimpleRL-Zoo dataset (Zeng et al., 2025) with various diffculty and its translated dataset as more challenging math dataset than GSM8K.

**Evaluation Dataset.** We evaluated our model on the translated GSM8K and MATH500 (Lightman et al., 2024) test sets across multiple languages. In order to compute the accuracy, we utilize math-verify library [2] for obtaining robust mathematical expression.

**Word Ratio.** We computed both the word ratio and the character ratio for each language to determine whether the model maintains its input-output language consistency following GRPO training. To compute the word ratio, we first remove all LaTeX expressions (e.g., $...$, \begin{...}, \end{...}) from the model's generated text. We then tokenize using simple regular-expression rules, using Multi-bleu [3], so that punctuation, brackets, and quotes are properly separated. Tokens that consist purely of math expressions, special symbols, or backslash commands are discarded. For each remaining token, we examine its characters to determine whether they belong exclusively to one of several script ranges, such as Hangul (U+AC00–U+D7A3), Latin alphabets (A–Z, a–z), CJK characters (U+4E00–U+9FFF, etc.), or Cyrillic (U+0400–U+04FF). We calculate the *word ratio* of a given language by dividing its token count by the total token count, and the *character ratio* by comparing the number of characters of each script to the total character count. Any token that mixes English letters with another script is labeled as a code-switching token, whose ratio is similarly tracked. This uniform preprocessing and detection pipeline thus enables a quantitative assessment of how models maintain linguistic fidelity in multilingual output.

## 3.2 Main Results

For the main study, we fine-tune each backbone with GRPO under identical hyper-parameters. Table 1 reveals three consistent patterns.

**Accuracy.** GRPO generally increases accuracy on our maths benchmarks, though the size of the gain depends on language and dataset. On the translated GSM8K, fine-tuning Llama-3.2-3B Instruct with GRPO raises accuracy by $+7.4\,\mathrm{pp}$ in Chinese, $+14.5\,\mathrm{pp}$ in Korean, and $+17.1\,\mathrm{pp}$ in Ukrainian. MATH500 shows the same upward trend (e.g. $+12.0\,\mathrm{pp}$ for Ukrainian).

**Target Word Ratio.** The accuracy gains come at the cost of target-language fidelity. For high-resource Chinese, the target-language word ratio stays above $90\%$ (a modest $-1.4\,\mathrm{pp}$ drift). For mid-resource Korean the drop is larger ($-8.1\,\mathrm{pp}$), while for low-resource Ukrainian the collapse is catastrophic: the share of Ukrainian tokens plummets from $98\%$ to $0.3\%$ on GSM8K (a $-97.6\,\mathrm{pp}$ change) and to $0.9\%$ on MATH500. Thus, GRPO systematically trades language consistency for higher reward. For this cross-lingual collapse, we provide more detailed analysis and discussion in the Section 4. A language-centric backbone mitigates, but does not remove, collapse. HyperCLOVA-X (HCX), which is pre-trained with a strong Korean prior, retains $98.2\%$ Korean tokens after fine-tuning while still improving accuracy by $+1.7\,\mathrm{pp}$. However, even HCX loses $8.3\,\mathrm{pp}$ of Korean ratio on MATH500, indicating that pre-training priors delay rather than eliminate collapse.

Taken together, these results confirm the hypothesis that GRPO *amplifies* the pre-training language imbalance: the lower the target-language resource level, the more aggressively the model abandons it during reasoning. Subsequent sections dissect the mechanisms behind this trade-off and evaluate interventions that seek to decouple accuracy from language fidelity.

## 3.3 Further Validations

In this subsection we probe GRPO-trained reasoning models on mathematical tasks more closely. First, to see whether **Cross-lingual Collapse** can be curbed, we adopt the *language-consistency reward* introduced by DeepSeek-AI et al. (2025). Specifically, we introduce a language consistency reward (Lang loss), which is measured by the proportion of words in the model's output that belong to the target language. Second, to test whether the phenomenon worsens on harder material, we replace the translated GSM8K corpus with the more demanding SimpleRL-Zoo dataset (Zeng et al., 2025), raising the training difficulty and tracking

---

[1] If a COMET score of translated sentence is below 0.70, we retranslate it with the same MT model using a new random seed.

[2] https://github.com/huggingface/Math-Verify

[3] https://github.com/moses-smt/mosesdecoder/blob/master/scripts/generic/multi-bleu.perl

(a) GSM8K Accuracy (UK)  (b) Word Ratio (UK)  (c) Word Ratio (UK) w/ Lang loss

(d) GSM8K Accuracy (KO)  (e) Word Ratio (KO)  (f) Word Ratio (KO) w/ Lang loss
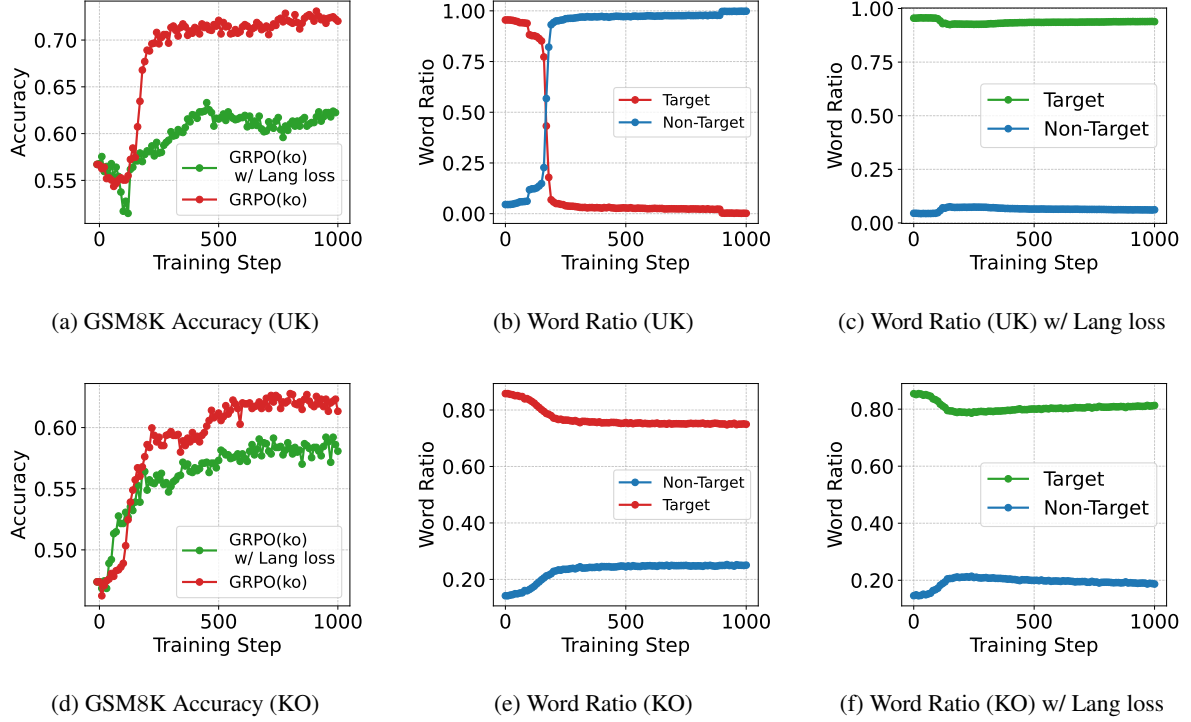
Figure 2: Figures 2a–2c compare Llama-3.2-3B Instruct trained with GRPO on the *Ukrainian*-translated GSM8K with and without the language-consistency reward (Lang loss), while Figures 2d–2f provide the same comparison for the *Korean*-translated GSM8K. Across both languages the pattern is clear: the language-consistency reward reliably preserves the target-language word ratio, yet it also *dampens* the accuracy gains that GRPO would otherwise deliver. In particular, Figures 2a–2c show that the reward almost completely prevents cross-lingual collapse in the Ukrainian run—though at the cost of a modest drop in performance—and an analogous trade-off appears in the Korean case.

how the model's behaviour shifts.

**Language Consistency Reward.** As summarized in Figure 2, we run an ablation in which Llama-3.2-3B is fine-tuned with GRPO on the Ukrainian- and Korean-translated GSM8K, once with the language-consistency reward and once without it. **Ukrainian(UK).** In the vanilla setting (Figures 2a–2c, solid line) the model undergoes a full cross-lingual collapse: the share of Ukrainian tokens in its chain of thought drops to almost zero while accuracy rises sharply. Adding the language-consistency reward (dashed line) prevents that collapse—the Ukrainian word ratio stays high—yet the accuracy gain is noticeably smaller. This shows that forcing GRPO to keep the reasoning trace in the target language safeguards linguistic fidelity at the cost of some performance. **Korean.** Figures 2d–2f display the same trade-off for the Korean-translated corpus: the reward stabilises the Korean word ratio but trims the accuracy boost.

Taken together, these results suggest that during GRPO the model actively probes alternative reasoning paths and, when allowed, gravitates toward high-resource English to maximise reward.

Constraining the trace to a low- or mid-resource language blocks that shortcut, preserving the intended language but sacrificing part of the accuracy gain. We investigate the underlying dynamics of this behaviour in more detail in Section 4.

**Impact of a Harder Curriculum.** Table 2 examines whether Cross-lingual Collapse can be triggered in a mid-resource language that previously appeared robust. When Qwen-2.5 1.5B is trained only on the Korean-translated GSM8K, the Korean word ratio remains high (>88 %) and no collapse is observed after 1 K updates. However, adding the more demanding SimpleRL-Zoo (Zeng et al., 2025) subset to the replay buffer changes the picture dramatically. After just 1 K steps on the mixed corpus the Korean ratio has already fallen to 69 % on MATH and 89 % on GSM8K, and by 2 K steps it plummets to **14.5 %** and **2.1 %**, respectively, even as accuracy nearly doubles on MATH500. Figure 3 confirms that a harder curriculum (training with more challenging dataset) exacerbates cross-lingual collapse, which intensifies as training progresses.

| Dataset | Steps | GSM8K (KO) | | MATH500 (KO) | |
|---|---|---|---|---|---|
| | | Accuracy (%) | Word Ratio (%) | Accuracy (%) | Word Ratio (%) |
| GSM8K | 1K | 42.3 | 94.3 | 25.7 | 88.0 |
| GSM8K + SimpleRL | 2K | 47.5 | **14.5** | 46.7 | **2.1** |

Table 2: Adding a harder curriculum triggers collapse in a mid-resource language. We fine-tune Qwen-2.5 1.5B Instruct on the *Korean*-translated GSM8K alone (top row) or on GSM8K mixed with the more challenging SimpleRL-Zoo set (*GSM8K+SimpleRL*). With GSM8K only, the model maintains a high Korean word ratio (> 88 %) and shows no cross-lingual collapse. Introducing SimpleRL-Zoo (Zeng et al., 2025), however, progressively erodes the ratio—dropping to 14.5 % at 2 K steps on GSM8K and to 2.1 % on MATH500—while accuracy rises. Thus, harder training material induces the same English-drift phenomenon in Korean that we previously observed only for low-resource Ukrainian.
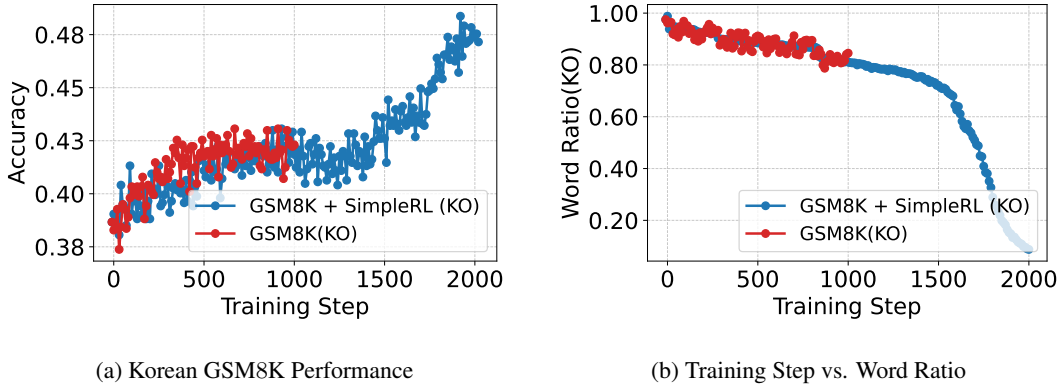


(a) Korean GSM8K Performance      (b) Training Step vs. Word Ratio

Figure 3: Training with a more challenging dataset mix, GSM8K combined with SimpleRL-Zoo, exacerbates cross-lingual collapse.

## 4 Analysis & Discussions

### 4.1 Cross-lingual Collapse

The evidence assembled so far paints a coherent picture: (1) Universal Drift. GRPO pushes all models toward the dominant pre-training language, but the speed and severity of that drift scale with resource level: minimal in Chinese, moderate in Korean, catastrophic in Ukrainian (Table 1). (2) Difficulty as a Trigger. A mid-resource model that is stable on GSM8K alone collapses after we inject a harder curriculum (Table 2), showing that *task difficulty*, tilts the optimizer toward English reasoning. (3) Reward Design Matters, but Costs Accuracy. The language-consistency reward prevents collapse (Figure 2) yet removes much of GRPO's accuracy gain, implying that the model *strategically* chooses English traces to maximize reward under pressure.

These findings confirm our central claim: GRPO amplifies the linguistic prior that best optimizes reward, and the gap between high- and lower-resource languages widens as tasks grow harder.

**Open Questions.** Three natural follow-ups guide the rest of the paper. Micro-level dynamics—how does collapse unfold token-by-token? We examine rollout samples to see when English fragments first appear and whether they align with specific reasoning steps (§4.2). Reversibility—can we roll the collapse back? We ask whether further GRPO on a distilled LRM can rebuild target-language capacity (§4.3). The next four subsections tackle these questions in turn, moving from surface behaviour to internal evidence and finally to mitigation strategies.

### 4.2 Rollout Analysis

As illustrated in Figure 4, we display representative rollout traces produced by Qwen-2.5 1.5B while it undergoes GRPO on the Korean-translated GSM8K corpus. During target-language training the model increasingly secures the accuracy reward by solving the problem in a *non-target* language, typically its high-resource, pre-training-dominant English, rather than in Korean. Each time such a trace earns a positive GRPO advantage, the optimizer reinforces that behavior, gradually shift-

## Question (Korean)

비고의 나이는 그의 남동생이 2살일 때 남동생 나이의 두 배보다 10살 더 많았습니다. 만약 그의 남동생이 현재 10살이라면,

(Vigo was 10 years older than twice his younger brother's age when the brother was 2 years old. If the younger brother is now 10)

그들 의 나이 합은 얼마입니까? 차근차근 단계별로 논리적으로 풀어주세요, 최종 답은 \boxed{} 안에 넣어주세요.

(what is the sum of their current ages? Please solve the problem step by step, showing your reasoning logically, and place the final answer inside \boxed{}.)

## Rollout Examples during Training

1. 남동생의 현재 나이를 확인합니다. 문제에서 남동생이 현재 10살이라고 주어졌습니다.
(1. Confirm the younger brother's current age. According to the problem, he is now 10 years old.)
[...]
남동생이 지금 10살이므로, 비고는 10 + 12 = 22살입니다. 4. 두 사람의 현재 나이 합을 구합니다. 10 + 22 = 32
따라서, 두 사람의 나이 합은 \boxed{32}
(Since the younger brother is 10, Vigo is 10 + 12 = 22 years old. 4. Calculate the sum of their current ages:
10 + 22 = 32. Therefore, the sum of their ages is \boxed{32})

→ Target-language Reward

Let's break this down step by step:\n\n1. **Identify the current age of the younger brother:**\n - The problem states that the younger brother is currently 10 years old.
[...]
The younger brother is 10 years old.\n - The older brother is 22 years old.\n - The sum of their ages is \\(10 + 22 = 32\\).\n\nTherefore, the sum of their ages is \boxed{32}

→ Non-target-language Reward

Figure 4: Rollout examples from GRPO training of Qwen-2.5 1.5B on the Korean-translated GSM8K. Observe that the model often arrives at the right answer via English reasoning (*non-target language*); because any correct answer earns full reward, repeated reinforcement of such off-language traces gradually shifts the chain-of-thought word ratio away from Korean.

ing the overall word ratio away from the target language. Cross-lingual Collapse emerges when these English-based solutions surge in frequency: even questions presented in Korean are eventually reasoned through, and answered, in the dominant high-resource language.

### 4.3 Further training of Distilled LRMs

As depicted in Figure 5, we apply a second round of GRPO to the DeepSeek-R1–Distilled Qwen to test whether continued fine-tuning can correct the entrenched reasoning bias. The results reveal a steep decline in the target-language word ratio, indicating that the phenomenon is difficult to reverse.

## 5 Related Works

### 5.1 Long Chain-of-Thought Generation

DeepSeek-AI et al. (2025) push the envelope on reinforcement-learning–based reasoning by introducing DeepSeek-R1-Zero, the open-source model trained with pure RL, specifically Group-Relative Policy Optimization (GRPO), without any supervised warm-up, and its follow-up DeepSeek-R1, which adds a small cold-start SFT stage and multistage RL to further boost performance. Their study demonstrates that large-scale GRPO can elicit impressive gains on mathematics and coding benchmarks, and that the resulting reasoning patterns can be distilled into much smaller dense models.

Notably, the authors briefly report undesirable "language mixing" and readability issues that emerge during RL, suggesting that reward-driven optimization may inadvertently disrupt linguistic fidelity. However, DeepSeek-R1 focuses almost exclusively on English prompts and does not quantify the extent, or direction, of its language drift. Our work complements these findings by conducting a systematic, multilingual analysis of GRPO and revealing a pronounced *Cross-lingual Collapse*: as RL progresses, chain-of-thought reasoning reverts to the pre-training-dominant language, catastrophically eroding performance in low-resource languages.

### 5.2 Multilingual Instruction tuning

Recent work shows that even a pinch of multilingual data during instruction tuning can unlock substantial cross-lingual generalisation in otherwise English-centric LLMs. Shaham et al. (2024) demonstrate that fine-tuning with as few as two to three languages is "necessary and sufficient" to elicit target-language responses across five downstream tasks, with the marginal benefit largely determined by how well that language was covered in pre-training. Complementing this, Kew et al. (2024) find that injecting only 40 non-English instruction–response pairs, or diversifying the tuning mix to merely 2–4 languages, yields instruction-following quality on a par with (or exceeding)

(a) Korean GSM8K Performance
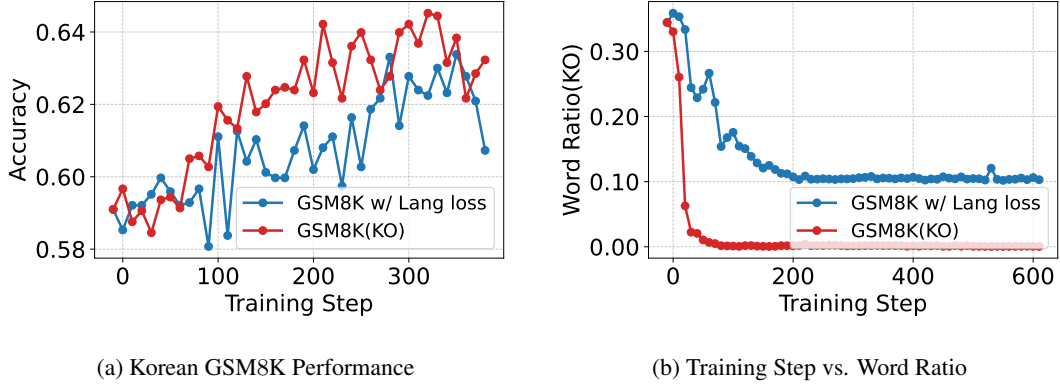           (b) Training Step vs. Word Ratio

Figure 5: We continued GRPO fine-tuning of the DeepSeek-R1-Distill Qwen model on the Korean-translated GSM8K dataset to encourage Korean chain-of-thought reasoning. As Figure 5b shows, the distilled model still exhibits cross-lingual collapse during training.

monolingual baselines while slashing per-language data by an order of magnitude. Yoo et al. (2024a) demonstrate that incorporating a sufficient amount of code-switched data (combining English and the target language) can effectively adapt an English-centric model, allowing the model to transfer its English-based knowledge into the target. Those studies therefore argue that massive multilingual corpora are not a prerequisite for broad cross-lingual utility; rather, strategically chosen seed languages can act as effective "anchors" that bootstrap transfer to unseen languages. Crucially, neither paper probes how reinforcement-learning–based reasoning objectives interact with this minimalist recipe, leaving open the question of whether such scarce multilingual supervision can withstand the linguistic pressures we observe under GRPO.

## 5.3 Multilingual Reasoning

Recent mechanistic studies reveal that multilingual reasoning in LLMs is far from language-neutral. Using a layer-by-layer logit-lens analysis, Schut et al. (2025) show that models such as Llama-3.1 first route lexical concepts through an English-centric representation space, even when both the prompt and the final output are non-English, and that activation steering vectors learned in English outperform those derived in the target language. Complementary circuit-tracing work (Lindsey et al., 2025) on Claude 3.5 Haiku uncovers a hybrid architecture in which abstract, language-agnostic sub-circuits cooperate with language-specific pathways; yet attribution graphs also expose English as the "mechanistically privileged" default when multiple languages compete to drive the same decision. Together, these findings suggest that multilingual LLMs execute core reasoning in a latent space biased toward high-resource English, while lower-level language-selective circuits merely translate inputs and outputs—an asymmetry that our study identifies as a key driver of *Cross-lingual Collapse*.

## 6 Conclusion

This study uncovers and characterizes **Cross-lingual Collapse**: in our experiments, multilingual LLMs trained with reinforcement-learning–based reasoning objectives tends to drift away from the prompt language, instead reasoning in the high-resource language that dominated their pre-training. Systematic GRPO experiments across Chinese, Korean, and Ukrainian confirm a clear gradient: the drift is barely noticeable in the high-resource setting, pronounced in a mid-resource language, and devastating in a low-resource one. Crucially, the phenomenon can be easily triggered by task difficulty. Introducing the harder SimpleRL-Zoo curriculum is enough to push otherwise stable Korean models over the edge. A language-consistency reward can arrest the collapse, yet doing so sacrifices much of GRPO's accuracy benefit, and once the bias is entrenched, even a second round of GRPO on distilled models fails to recover target-language reasoning.

Taken together, these results expose a central tension: the very training signals that sharpen logical accuracy also amplify pre-training language imbalance, widening the performance gap between high- and lower-resource languages. We release our code and dataset in the hope that they will spur research on reward designs that disentangle accuracy from

language fidelity, on curricula that equalise difficulty across languages, and on early-warning diagnostics that surface cross-lingual biases before they ossify.

## Limitations

Our study is bounded by (i) model scale: experiments stop at 3 B parameters due to GPU limits, so collapse behaviour may differ for larger checkpoints; (ii) language scope: we test only Chinese, Korean, and Ukrainian, leaving broader typological coverage unexplored; and (iii) metrics: the collapse score relies on a script-based word-ratio heuristic that undercounts mixed or borrowed tokens. Results therefore offer a first glimpse, not a definitive map, of cross-lingual dynamics under GRPO.

## References

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 81 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2024. xcomet: Transparent machine translation evaluation through fine-grained error detection. *Transactions of the Association for Computational Linguistics*, 12:979–995.

Tannon Kew, Florian Schottmann, and Rico Sennrich. 2024. Turning english-centric llms into polyglots: How much multilinguality is needed? In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 13097–13124. Association for Computational Linguistics.

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe.

2024. Let's verify step by step. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar, Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan, Adam Jermyn, Andy Jones, and 8 others. 2025. On the biology of a large language model. *Transformer Circuits Thread*.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022a. Training language models to follow instructions with human feedback. *Preprint*, arXiv:2203.02155.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022b. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

Lisa Schut, Yarin Gal, and Sebastian Farquhar. 2025. Do multilingual llms think in english? *CoRR*, abs/2502.15603.

Uri Shaham, Jonathan Herzig, Roee Aharoni, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. Multilingual instruction tuning with just a pinch of multilinguality. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 2304–2317. Association for Computational Linguistics.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300.

Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. Language models are multilingual chain-of-thought reasoners. In *The Eleventh International Conference on Learning Representations*.

Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane

Rivière, and 1 others. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.

Qwen Team. 2024. Qwen2.5: A party of foundation models.

Zhijun Wang, Jiahuan Li, Hao Zhou, Rongxiang Weng, Jingang Wang, Xin Huang, Xue Han, Junlan Feng, Chao Deng, and Shujian Huang. 2025. Investigating and scaling up code-switching for multilingual language model pre-training. *arXiv preprint arXiv:2504.01801*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.

Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2020. CCNet: Extracting high quality monolingual datasets from web crawl data. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4003–4012, Marseille, France. European Language Resources Association.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.

Haneul Yoo, Cheonbok Park, Sangdoo Yun, Alice Oh, and Hwaran Lee. 2024a. Code-switching curriculum learning for multilingual transfer in llms. *arXiv preprint arXiv:2411.02460*.

Kang Min Yoo, Jaegeun Han, Sookyo In, Heewon Jeon, Jisu Jeong, Jaewook Kang, Hyunwook Kim, Kyung-Min Kim, Munhyong Kim, Sungju Kim, and 1 others. 2024b. Hyperclova x technical report. *arXiv preprint arXiv:2404.01954*.

Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.

Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. 2025. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *Preprint*, arXiv:2503.18892.

Chengzhi Zhong, Fei Cheng, Qianying Liu, Junfeng Jiang, Zhen Wan, Chenhui Chu, Yugo Murawaki, and Sadao Kurohashi. 2024. Beyond english-centric llms: What language do multilingual language models think in? *CoRR*, abs/2408.10811.