

ReTRE: Benchmarking LLM Transfer Robustness with Structure-Preserving Variants

Anonymous ACL submission

Abstract

Large language models (LLMs) have achieved strong performance on standard benchmarks, yet their performance is not robust across different task manifestations. It remains unclear how performance changes under controlled task rewrites that preserve the original solution structure, while varying the rewrite type and level. To address this question, We introduce ReTRE (Rewrite-based Transfer Robustness Evaluation), an evaluation benchmark inspired by learning transfer theory that probes transfer robustness along two rewrite level: Near Transfer and Far Transfer. Given the increasing support for multimodal inputs in modern LLMs, ReTRE considers not only text-based rewrites but also modality-type rewrites. To ensure that the solution structure is preserved, ReTRE employs a multi-agent pipeline that extracts the solution steps from the original task, designs a corresponding transfer strategy, and generates rewritten variants. Each stage is equipped with a dedicated validation agent that iteratively verifies structure preservation and correctness. Evaluations on mathematical and science tasks across state-of-the-art multimodal LLMs reveal a consistent transfer gap: performance exhibits a general declining trend as transfer similarity drop and strong text performance can face performance decline under cross-modal transfer. Crucially, we identify a divergence between post-training paradigms: reinforcement learning preserves transfer robustness, whereas supervised fine-tuning tends to overfit the training distribution, leading to severe degradation in far-transfer performance despite strong in-distribution accuracy. The code is available at anonymous.4open.science/r/TransferRobust-E738/.

1 Introduction

Large language models (LLMs) have achieved strong performance across a wide range of bench-

marks (Cao et al., 2025), and such performance is often taken as evidence that a model has acquired the targeted knowledge or skills. However, learning transfer theory emphasizes that applying acquired knowledge to novel contexts is a key signal of deep understanding¹. This motivates a practical evaluation question: Do LLMs remain robust when facing novel manifestations of the same problem?

A common approach to this question is to rewrite existing evaluation datasets (Wu et al., 2024; Wang et al., 2024b; Huang et al., 2025a,b; Kirtane et al., 2025a). These studies construct new test instances through operations such as data perturbation, character substitution, and semantic rewriting, revealing that high benchmark scores can stem from data contamination or sensitivity to surface expressions. However, such approaches seldom control for the degree of divergence between original and rewritten tasks, nor do they systematically vary the transformation type. As LLMs are increasingly expected to generalize across diverse and unfamiliar contexts (Mumuni and Mumuni, 2025), a more principled question arises: How does model performance change under controlled, structure-preserving rewrites that vary in both type and degree?

Inspired by learning transfer theory (Perkins et al., 1992; Barnett and Ceci, 2002; Hilton and Pellegrino, 2012), we propose ReTRE (Rewrite-based Transfer Robustness Evaluation), a benchmark that evaluates transfer robustness using structure-preserving rewrites, as shown in the Figure 1. ReTRE generates variants along two complementary dimensions—Knowledge Domain (KD) and Modality Context (MC)—and instantiates each dimension at two discrete transfer settings (near and far). In the Knowledge Domain dimension, variants are designed to change the semantic background

¹<https://poorvucenter.yale.edu/transfer-of-knowledge-to-new-contexts>

of a problem while keeping the intended reasoning procedure aligned with the original. Near-transfer variants remain within related STEM domains (e.g., reframing a physics story in a chemistry context), whereas far-transfer variants move the same problem structure into non-STEM domains such as economics, law, or social sciences. In practice, these KD rewrites keep the abstract variables and relations fixed at the level of the solution template, while modifying domain-specific entities, terminology, and contextual framing. In the Modality Context dimension, variants are designed to change how task information is presented while keeping the information needed to solve the problem aligned with the original. Near-transfer variants apply within-text reformats, such as converting a prose description into a structured Markdown table. Far-transfer variants shift from text to the visual channel by representing task-relevant structure in diagrams or plots. These MC rewrites keep the problem’s required information content aligned, while changing the input channel and the organization of that information. To construct structure-preserving variants at scale, we design a three-layer agentic pipeline that pairs a generator with a verifier at each stage. The pipeline extracts a transferable solution structure from the original task, proposes a transfer strategy for the selected setting, and synthesizes the final variant, with iterative validation ensuring structure preservation and correctness.

We evaluate a suite of mainstream multimodal LLMs on Mathematical reasoning task using ReTRE. Results reveal a clear transfer gap: near-transfer variants yield performance close to the original, whereas far-transfer variants incur substantial degradation. Notably, Modality Context (MC) variants cluster tightly with original problems in embedding space, and MC-Near occasionally yields slight gains from reduced parsing ambiguity due to structured formatting. In contrast, Knowledge Domain (KD) variants form distinct clusters far removed from the original distribution, with KD-Far posing a consistent challenge across all models. We also find that thinking-mode models demonstrate enhanced transfer robustness—for instance, Claude-4.5-Haiku (thinking) improves over its standard counterpart by +3.6 percentage points in average accuracy, effectively rivaling the top-tier models. Motivated by these findings, we conduct controlled experiments to isolate the effects of post-training paradigms on transfer robustness. Our results reveal a striking divergence: full Super-

vised Fine-Tuning achieves perfect in-distribution accuracy (100%) but collapses catastrophically on transfer variants, whereas Reinforcement Learning, like Group Sequence Policy Optimization (GSPO) (Zheng et al., 2025), improves both in-distribution and transfer performance simultaneously. We further examine how robustness varies with model scale and generational evolution. While scaling from 2B to 235B generally enhances robustness, gains are not strictly monotonic, and mid-scale increases show diminishing returns. Generational improvements effectively mitigate sensitivity to modality context shifts, yet the challenge of knowledge domain transfer persists and even intensifies in absolute terms. Finally, we extend ReTRE to natural science reasoning on GPQA, where MC-Far emerges as the most challenging setting, causing performance drops exceeding 20 percentage points for top models. Our contributions are as follows:

- We propose ReTRE, a learning-transfer-inspired diagnostic benchmark that profiles transfer robustness via controlled task transformations along two axes—Knowledge Domain and Modality Context—each instantiated at near and far transfer distances.
- On MATH500, we observe a consistent transfer gap across mainstream multimodal LLMs: performance is relatively stable under near-transfer settings but degrades substantially under far-transfer settings. Moreover, strong performance on text-based settings does not necessarily translate to cross-modal variants, where models can still exhibit notable failures. Thinking-mode models consistently outperform their non-thinking counterparts in transfer robustness.
- In a controlled study on Qwen3-VL-8B-instruct, we find that post-training paradigms strongly affect transfer robustness: full SFT achieves perfect in-distribution accuracy but collapses on transfer variants, whereas GSPO improves both in-distribution and transfer performance. Model scaling and generational evolution generally enhance robustness, though gains are not monotonic, and distant-domain transfer remains a bottleneck.

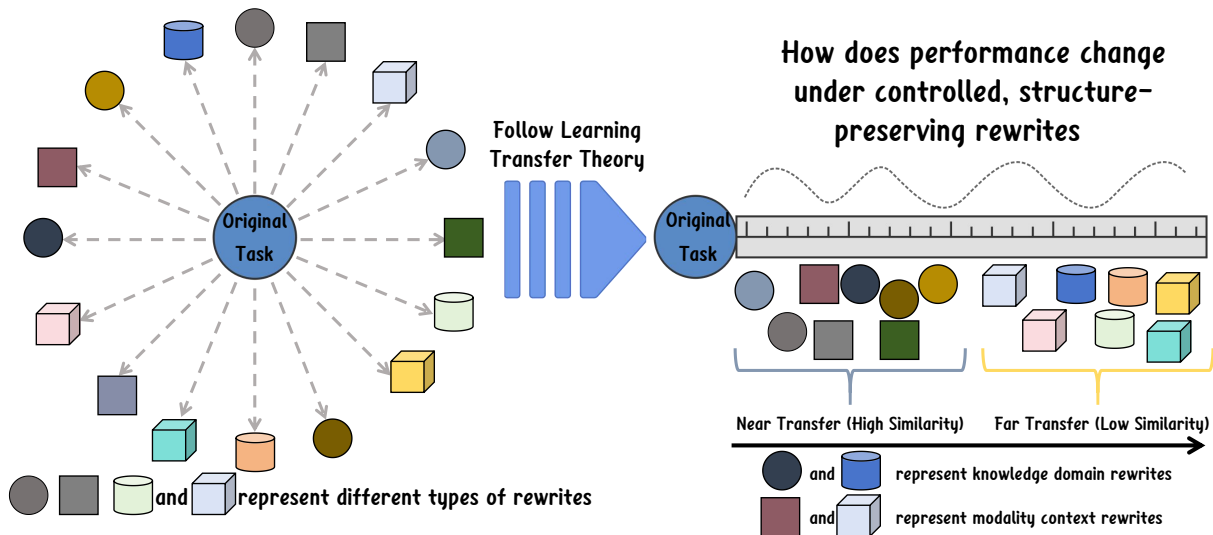


Figure 1: Motivation for ReTRE. (Left) Existing rewrite-based evaluations often ignore rewritten variant–original task similarity and ignore how performance changes under controlled and structure preserving rewrites. (Right) ReTRE categorizes rewrites into distinct transfer similarity (Near and Far) within knowledge domain and modality context types. This discrete setting provides a solution to evaluate transfer robustness.

2 related work

2.1 LLM Benchmark

In recent years, a large number of benchmarks (Ni et al., 2025) have been proposed to evaluate LLMs, covering diverse tasks such as natural language understanding, machine translation, and text generation. Advanced LLMs have achieved remarkable performance across these benchmarks. Recently, the evaluation paradigm has shifted from assessing task-specific performance to measuring more general capabilities (Cao et al., 2025), including reasoning ability, instruction following, and knowledge understanding. Nevertheless, several studies have revealed that the high scores achieved by many models often result from surface-level pattern matching rather than genuine comprehension. This issue is largely attributed to data leakage between training and evaluation datasets (Wu et al., 2025). To mitigate this problem, one common approach is data perturbation (Lunardi et al., 2025; Kirtane et al., 2025b; Huang et al., 2025b), which introduces variations to test the robustness of model reasoning. But in practical scenarios, the same problem can have multiple manifestations and there are similarities and differences in type among the variants.. Therefore, It needs to consider how performance changes under controlled task rewrites that preserve the original solution structure, while varying the rewrite type and level.

2.2 Learning Transfer Theory

Learning transfer is a central concept in educational psychology, commonly defined as the ability to apply acquired knowledge or skills beyond the original learning context (Perkins et al., 1992). A widely adopted distinction in this literature is between near transfer and far transfer. Near transfer typically refers to situations where the transfer context remains highly similar to the original learning context, whereas far transfer involves applying the same underlying knowledge to contexts that differ substantially in surface characteristics or situational framing. Barnett and Ceci (Barnett and Ceci, 2002) further formalize transfer by proposing a six dimensional taxonomy. While this taxonomy provides a rich conceptual framework for analyzing human learning, it is not directly applicable to evaluating large language models (LLMs). Several dimensions (e.g., social context and physical context) are not suitable for models evaluation.

3 Method

3.1 Problem Formulation

We define transfer robustness as the performance change under controlled shifts in task manifestation, while the underlying reasoning procedure is preserved.

Given a source dataset $\mathcal{D}_{\text{src}} = \{(x, y)\}$, we construct a target dataset $\mathcal{D}_{\text{tar}} = \{(x', y')\}$ using a multi-agent system \mathcal{M} that applies a transforma-

tion \mathcal{T} under a transfer setting s :

$$(x', y') = \mathcal{T}(x, y; s) \text{ via } \mathcal{M}, \text{ s.t. } \mathcal{L}_{\text{core}}(x') \approx \mathcal{L}(x). \quad (1)$$

Here $\mathcal{L}(\cdot)$ extracts the reasoning procedure required to solve the original task, and $\mathcal{L}_{\text{core}}(\cdot)$ denotes the core reasoning component of the transferred task. The constraint enforces that solving x' relies on the same solution procedure as x , while permitting limited auxiliary scaffolding to keep the transferred instance self-contained and to prevent answer leakage (e.g., the solution appearing verbatim in the problem text).

For mathematical reasoning tasks, we typically preserve the original answer since domain transfer maintains numerical relationships directly. For natural science question answering, however, directly re-contextualizing domain-specific concepts (e.g., chemical reactions, physical laws) into unrelated domains would require external prerequisite knowledge. To address this, we wrap the transferred content within a fictional scenario and provide explicit rule descriptions with minimal demonstrations, enabling the model to solve the problem using the same reasoning procedure without relying on real-world domain knowledge. In such cases, the answer may differ in surface form, but the underlying solution procedure remains equivalent.

We apply this formulation to two task families: (1) mathematical reasoning, where the output is a short-form answer derived from quantitative relations, and (2) natural science question answering, where the model selects the correct option from multiple choices. The same transfer framework applies to both families, though the concrete rewrite operations differ across tasks and transfer axes.

3.2 Data Transfer Pipeline

Figure 2 illustrates our three-layer Generator-Validator pipeline. Each layer pairs a generator agent with a validator agent. The pipeline employs layer-wise retry: if a candidate fails validation, we regenerate at that layer without discarding validated outputs from earlier layers. Instances that exceed the maximum retry limit are discarded from the final dataset.

Layer 1 focuses on structure extraction. Agent1 (Extractor) extracts a transferable representation from the source instance. For mathematical reasoning, this representation records the essential quantitative entities, constraints, and the minimal solution skeleton needed to compute the answer. For natural

science question answering, it captures the information required to determine the correct option while preserving the original question intent and decision structure. When the transfer axis is modality context, the extractor additionally produces a structured split: Part A contains the question frame and options, while Part B contains the content to be converted to another modality. The extractor also assigns a content type label for Part B, such as numerical values, process descriptions, or relational mappings. Agent2 (Validator) then checks that the extraction is faithful to the original and contains no hallucinated information. For modality context transfers, it additionally verifies complementarity: neither Part A nor Part B alone should be sufficient to solve the problem, yet their combination must provide all information necessary for solving.

Layer 2 handles transfer design. Agent3 (Designer) converts the extracted structure into a concrete transfer plan under the chosen axis and level. For knowledge domain transfers, the designer specifies a target context and constructs a mapping that rewrites domain entities and terminology while keeping the abstract variables, relations, and solution skeleton fixed. Near transfer uses a context within a closely related domain family, while far transfer targets a substantially different domain. For natural science question answering, since domain-specific concepts cannot be directly re-contextualized without introducing external prerequisites, the designer constructs a fictional scenario with explicit rule descriptions and minimal demonstrations, ensuring the instance remains solvable from the provided information alone. For modality context near transfer, the designer produces a table schema to organize the data. For modality context far transfer, the design strategy differs by task family: for mathematical reasoning, the designer produces a visualization specification to guide subsequent code generation; for natural science question answering, the designer produces a detailed image generation prompt that a text-to-image model can follow directly. We adopt different strategies because mathematical problems involve explicit numerical data that can be precisely rendered through programmatic plotting, whereas science problems often involve complex conceptual structures (e.g., molecular diagrams, reaction pathways) that are more naturally produced by generative image models. In all cases, the strategy must avoid embedding the question text or answer options in the visual artifact and must introduce

340 meaningful non-textual visual structure. Agent4
341 (Validator) verifies that the plan is compatible with
342 the extracted structure, preserves solvability, and
343 introduces no unstated external dependencies.

344 Layer 3 performs problem generation. Agent5
345 (Generator) synthesizes the final transferred in-
346 stance. For knowledge domain transfers, the gener-
347 ator instantiates the mapping and rewrites the task
348 into the target context while preserving the core
349 structure. For modality context near transfer, it ren-
350 ders the designated content into a structured format
351 such as a Markdown table and merges it with the
352 question frame. For modality context far transfer,
353 the generation process depends on the task family:
354 for mathematical reasoning, the generator produces
355 executable plotting code based on the visualization
356 specification and executes it in a sandboxed envi-
357 ronment to render the diagram; for natural science
358 question answering, the generator invokes a text-
359 to-image model to produce the visual artifact. In
360 both cases, the visual component is paired with the
361 textual question and options. Agent6 (Validator)
362 performs end-to-end validation, accepting a candi-
363 date only when it satisfies structure preservation,
364 information sufficiency, answer consistency, and
365 leakage prevention. For modality context transfers,
366 it additionally verifies that the textual and visual
367 components remain complementary, ensuring that
368 solving the problem requires integrating both rather
369 than reading the answer from either component
370 alone. Only instances passing all validation criteria
371 are added to the final dataset.

372 4 Experiments

373 Our experiment mainly investigates how the per-
374 formance of the model changes when it is applied
375 to variant data as the transfer distance increases. If
376 the model performance remains stable, it indicates
377 that the model has transfer robustness.

378 4.1 Setup

379 We selected MATH500(Lightman et al., 2023) and
380 GPQA(Rein et al., 2024) as the baseline datasets,
381 which are typical evaluation datasets in the fields
382 of mathematics and natural sciences. We use the
383 multi-agent framework to construct the Near and
384 Far transfer data in Knowledge Domain and Modal-
385 ity Context dimensions. Due to the involvement
386 of image input, the models we chose all support
387 multimodal input. We have selected the currently
388 advanced model for examination, including gpt-5-

389 mini(OpenAI, 2025a), o4-mini(OpenAI, 2025b),
390 claude-haiku-4-5-20251001 and claude-haiku-4-
391 5-20251001-thinking(Anthropic, 2025), gemini-
392 3-pro-all, gemini-2.5-flash-light-nothinking and
393 gemini-2.5-flash-light-thinking(Comanici et al.,
394 2025). To enhance the reliability of the experimen-
395 tal results, we set all models with the temperature
396 at 0 and use pass@1 as the metric.

397 4.2 Main results

398 Table 1 reports the comprehensive evaluation on
399 MATH500. Across all evaluated multimodal
400 LLMs, we observe a consistent transfer gap: perfor-
401 mance is generally preserved under near-transfer
402 settings but drops under far-transfer settings, with
403 the largest degradations appearing in the most
404 challenging shifts. Among the models, claude-
405 haiku-4-5-20251001-thinking achieves the highest
406 accuracy on the Original setting (93.5%), while
407 gemini-3-pro-all attains the best average accuracy
408 across all settings, indicating strong performance.
409 What’s more, Thinking-mode models tend to be
410 more transfer-robust than their non-thinking coun-
411 terparts. For example, claude-haiku-4-5-20251001-
412 thinking outperforms claude-haiku-4-5-20251001
413 across all transferred settings, and gemini-2.5-flash-
414 lite-thinking also shows improved robustness rela-
415 tive to gemini-2.5-flash-lite-nothinking. This pat-
416 tern suggests that explicit deliberation mechanisms
417 can mitigate degradation when task manifestations
418 change while the underlying solution structure is
419 preserved.

420 To contextualize these performance trends, Fig-
421 ure 3 visualizes the embedding distribution of origi-
422 nal and transferred instances using representations
423 extracted by qwen2.5-vl-embedding (Bai et al.,
424 2025b). The distribution highlights a clear contrast
425 between transfer types. Modality-context variants
426 remain close to the Original cluster, suggesting that
427 they preserve semantic proximity in the embedding
428 space. However, semantic proximity does not guar-
429 antee cross-modal correctness: despite clustering
430 near the Original, MC-Far introduces visual inputs
431 that can still trigger failures, and several models
432 exhibit marked drops from their text-based perfor-
433 mance (e.g., gemini-3-pro-all: 93.1% to 84.6%;
434 gpt-5-mini: 89.9% to 76.8%). These results indi-
435 cate that strong text performance does not neces-
436 sarily translate to cross-modal variants, even when
437 the task semantics remain closely aligned.

438 Importantly, rewrites do not uniformly reduce
439 performance. In MC-Near, multiple models

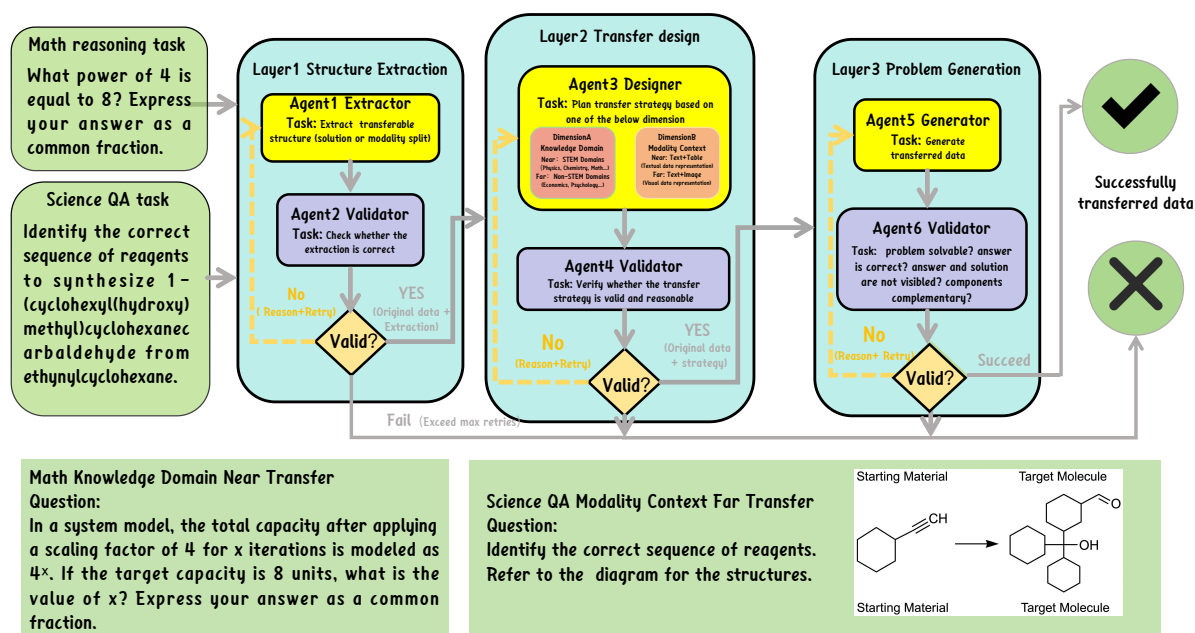


Figure 2: The ReTRE data transfer framework. Each layer forms a Generator–Validator pair with iterative refinement until the validation criteria are satisfied or the retry limit is reached. The bottom panel shows representative transfer examples for a math reasoning task and a science QA task. For space reasons, we omit the multiple-choice options in the science QA example; the options remain identical to the original after MC transfer.

slightly improve over the Original setting (e.g., gemini-3-pro-all: 93.1% to 93.5%; gpt-5-mini: 89.9% to 91.2%; gemini-2.5-flash-lite-nothinking: 85.6 to 90.2%). A plausible explanation is a structured-formatting effect: converting free-form text into well-organized tables reduces parsing ambiguity and facilitates information extraction. This benefit weakens for MC-Far, where the shift to visual representations introduces additional perception and grounding challenges.

In contrast, knowledge-domain variants form separated clusters that are farther from the Original distribution, reflecting substantial semantic drift under KD transfer. Consistent with this shift, all models show degradation on KD-Far (e.g., gemini-3-pro-all: 93.1% to 84.9%; gpt-5-mini: 89.9% to 80.0%). Overall, the results suggest that current models are comparatively more resilient to presentation changes within text (MC-Near) than to large semantic shifts across domains (KD-Far), while cross-modal transfer (MC-Far) remains a distinct failure mode where strong text performance can mask significant weaknesses.

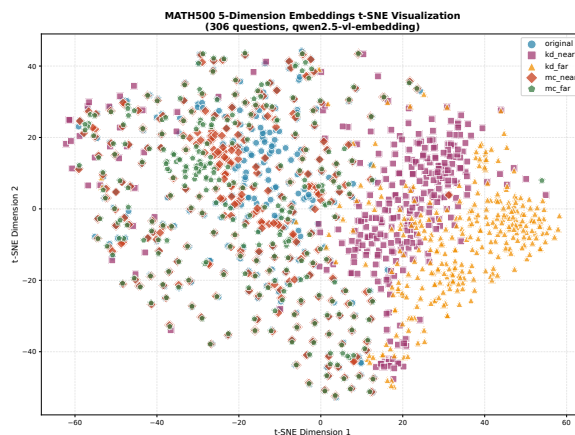


Figure 3: t-SNE visualization of problem embeddings extracted by qwen2.5-vl-embedding. The distribution reveals a fundamental distinction: Modality Context variants cluster closely with the Original problems, indicating semantic preservation. In contrast, Knowledge Domain variants form distinct, isolated clusters, reflecting significant semantic drift/addition.

4.3 Do different post-training paradigms have an impact on the transfer robustness of the model?

To verify the impact of different training paradigms on the robustness of model transfer, we conducted a controlled experiment. We trained the Qwen3-

Model	Original	KD-Near	KD-Far	MC-Near	MC-Far
gemini-3-pro-all	93.1%	87.5%	84.9%	93.5%	84.6%
claude-haiku-4-5-20251001-thinking	93.5%	85.9%	81.6%	91.8%	89.5%
o4-mini	92.2%	84.2%	80.7%	91.2%	81.7%
claude-haiku-4-5-20251001	91.2%	80.6%	79.0%	88.2%	85.3%
gpt-5-mini	89.9%	82.9%	80.0%	91.2%	76.8%
gemini-2.5-flash-lite-thinking	89.5%	82.6%	76.1%	87.6%	82.7%
gemini-2.5-flash-lite-nothinking	85.6%	78.3%	75.4%	90.2%	85.3%

Table 1: Main evaluation results on the MATH500 benchmark. Note that gemini-3-pro-all achieves the state-of-the-art performance, while claude-haiku-4-5-20251001-thinking demonstrates that System 2 reasoning capabilities closely rival the top model.

Model	Original	KD-Near	KD-Far	MC-Near	MC-Far
Qwen3-VL-8B-Instruct	77.5%	64.4%	59.5%	75.2%	68.6%
Qwen3-VL-8B-instruct(SFT)	100.0%	47.1%	48.0%	55.6%	43.4%
Qwen3-VL-8B-instruct(GSPO)	85.3%	74.2%	69.0%	85.3%	76.8%

Table 2: Performance comparison of Base, SFT, and GSPO settings. SFT achieves 100% on the original dataset but suffers from severe overfitting (catastrophic collapse) on transfer tasks. In contrast, GSPO surpasses the Base model in both Original accuracy and Transfer Robustness.

469 VL-8B-Instruct model on the MATH500 dataset
 470 using full-scale SFT and GSPO methods, and then
 471 evaluated it on the corresponding 4 sets of transfer
 472 data.

473 The experimental results present a striking con-
 474 trast between the two post-training paradigms, as
 475 shown in Table 2. We first observe that **SFT in-**
 476 **duces a robustness collapse**. While the SFT model
 477 achieves a perfect accuracy of 100% on the Origi-
 478 nal dataset, its performance collapses catastrophi-
 479 cally across all transfer variants. For instance,
 480 accuracy on the MC-Far and KD-Near datasets
 481 drops to 43.4% and 47.1% respectively, signifi-
 482 cantly lower than the Base model (68.6% and
 483 64.4%). This suggests that the SFT paradigm drove
 484 the model towards rote memorization of specific
 485 training instances, sacrificing generalization for in-
 486 distribution performance.

487 In stark contrast, **GSPO enhances generaliza-**
 488 **tion** and demonstrates superior capability. It not
 489 only outperforms the Base model on the Original
 490 dataset (85.3% vs 77.5%) but also significantly
 491 improves transfer robustness (e.g., 74.2% on KD-
 492 Near vs Base 64.4% on KD-Near). Unlike SFT,
 493 GSPO improves the model’s alignment with the
 494 target domain without overfitting to the specific
 495 phrasing or parameters of the training set.

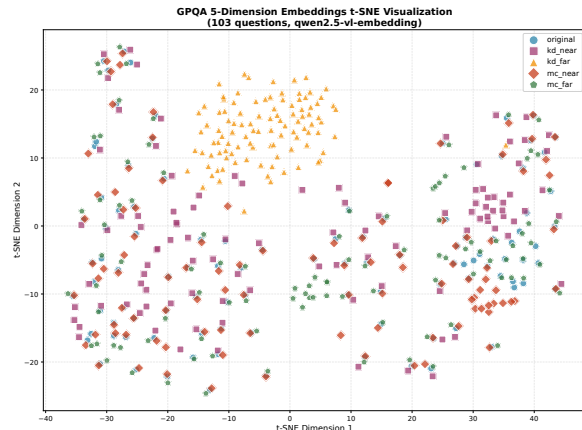


Figure 4: t-SNE visualization of GPQA embeddings. KD-Far instances separate distinctly from the original distribution, while MC variants remain entangled, indicating different degrees of semantic shift.

4.4 How robust are models under transfer in natural science reasoning?

496 We evaluate model transfer robustness on GPQA, a
 497 natural science multiple-choice benchmark. Adopt-
 498 ing the same two-axis design as in our main experi-
 499 ment (Knowledge Domain and Modality Context),
 500 we construct four transfer variants for 103 GPQA
 501 problems. Table 3 reports the accuracy across the
 502 original questions and all transfer settings.
 503
 504

505 Overall, transfer shifts consistently degrade per-
 506 formance for state-of-the-art models, with **MC-Far**
 507 presenting the most significant challenge. For in-

Model	Original	KD-Near	KD-Far	MC-Near	MC-Far
gemini-3-pro-all	91.3%	82.5%	81.6%	88.3%	70.0%
gpt-5-mini	84.5%	79.6%	68.9%	81.6%	68.0%
o4-mini	82.5%	81.5%	70.9%	74.8%	64.1%
claude-haiku-4-5-20251001-thinking	70.9%	72.8%	66.0%	71.8%	62.1%
gemini-2.5-flash-lite-thinking	63.1%	62.1%	52.4%	59.2%	51.6%
claude-haiku-4-5-20251001	58.3%	66.0%	54.4%	57.3%	54.4%
gemini-2.5-flash-lite-nothinking	49.5%	53.4%	52.4%	54.4%	43.7%

Table 3: Evaluation results on GPQA (103 problems). We report accuracy (%) on the original questions and four transfer settings (KD/MC \times Near/Far). With 103 problems, each item corresponds to 0.97% accuracy, so repeated percentages may occur due to discrete granularity.

stance, `gemini-3-pro-all` achieves 91.3% on the original questions but drops 21.3 percentage points to 70.0% on MC-Far. Similarly, `o4-mini` declines from 82.5% to 64.1% in the same setting. Interestingly, `gpt-5-mini` shows balanced degradation across both axes, with comparable performance on KD-Far (68.9%) and MC-Far (68.0%).

To interpret these patterns from a representation perspective, we visualize the embedding distribution in Figure 4. This highlights a critical insight: although MC-Far causes the most severe performance degradation, the problems remain semantically close to the source in the text embedding space. This suggests that the difficulty of Modality Context transfer stems not from semantic drift, but from the challenge of cross-modal grounding required to interpret the visual transformation.

In contrast, **Near-transfer** settings yield performance largely consistent with the original. Most models exhibit only minor fluctuations (within $\pm 3\%$), which we attribute to the inherent similarity between near-transfer variants and the original questions, combined with the limited granularity of pass@1 evaluation on 103 items (each problem corresponds to approximately 0.97% accuracy). We also note that some weaker models score higher on certain transfer variants than on the original questions. For example, `gemini-2.5-flash-lite-nothinking` improves from 49.5% to 53.4% on KD-Near and 54.4% on MC-Near. This suggests that the expert-level phrasing of GPQA may pose additional linguistic challenges for models with limited reasoning capacity, which the rewritten variants inadvertently alleviate.

5 Conclusion

This work examines the transfer robustness of large language models under controlled, structure-

preserving task transformations. We introduce RETRE, a transfer-oriented evaluation framework that varies task manifestations along knowledge domain and modality context axes at near and far levels, while preserving the underlying solution structure through a multi-agent pipeline. Experiments on mathematical reasoning and natural science QA reveal a consistent transfer gap: performance remains relatively stable under near transfer but degrades substantially under far transfer, especially for knowledge-domain shifts. Notably, strong performance on text-based settings does not necessarily translate to cross-modal variants, where models can struggle even when the underlying reasoning structure is preserved. Further analyses show that RL-based post-training improves transfer robustness, whereas supervised fine-tuning can overfit in-distribution data and collapse under transfer.

Limitations

Our work has certain limitations, and these limitations should be acknowledged. Firstly, we constructed migration data for the MATH500 and GPQA-diamond datasets. However, not all of them were successfully constructed. For instance, there were only 306 sets of migration data for MATH500. Additionally, considering economic benefits, for the image construction of MATH500, we used code to draw the pictures. For GPQA, we attempted to use the same configuration models (Gemini-2.5-flash and GPT-5-mini) from MATH500 to generate data, but the success rate was too low. Therefore, we replaced them with Gemini-3-pro-all and GPT-5.2 to act as agents, and used `gemini-3-pro-image-preview` to generate images. We considered that this was because Gemini-2.5-flash itself had difficulties in understanding the task (due to budget considerations, we did not provide sufficient rea-

582	soning configurations). However, the focus of our	2025b. Thinkbench: Dynamic out-of-distribution	636
583	work is on evaluating the model, not the way of	evaluation for robust llm reasoning. <i>arXiv preprint</i>	637
584	data construction. In the future, we will consider	<i>arXiv:2502.16268</i> .	638
585	how to utilize open-source models to construct data		
586	for evaluating SOTA models.		
587	References		
588	Anthropic. 2025. Claude 4.5 and claude haiku 4.5. https://www.anthropic.com/news/claude-haiku	Neeraja Kirtane, Yuvraj Khanna, and Peter Relan.	639
589	-4-5 . Accessed 2025-11-03.	2025a. Mathrobust-lv: Evaluation of large language	640
590		models' robustness to linguistic variations in mathe-	641
591	Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen,	tical reasoning. <i>arXiv preprint arXiv:2510.06430</i> .	642
592	Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei		
593	Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhi-	Neeraja Kirtane, Yuvraj Khanna, and Peter Relan.	643
594	fang Guo, Qidong Huang, Jie Huang, Fei Huang,	2025b. Mathrobust-lv: Evaluation of large language	644
595	Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng	models' robustness to linguistic variations in mathe-	645
596	Li, and 45 others. 2025a. Qwen3-vl technical report.	tical reasoning. <i>Preprint</i> , arXiv:2510.06430.	646
597	<i>arXiv preprint arXiv:2511.21631</i> .		
598	Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri	647
599	bin Ge, Sibao Song, Kai Dang, Peng Wang, Shi-	Edwards, Bowen Baker, Teddy Lee, Jan Leike,	648
600	jie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu,	John Schulman, Ilya Sutskever, and Karl Cobbe.	649
601	Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei	2023. Let's verify step by step. <i>arXiv preprint</i>	650
602	Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others.	<i>arXiv:2305.20050</i> .	651
603	2025b. Qwen2.5-vl technical report. <i>arXiv preprint</i>		
604	<i>arXiv:2502.13923</i> .	Riccardo Lunardi, Vincenzo Della Mea, Stefano Miz-	652
605	Susan M Barnett and Stephen J Ceci. 2002. When and	zaro, and Kevin Roitero. 2025. On robustness and	653
606	where do we apply what we learn?: A taxonomy for	reliability of benchmark-based evaluation of llms .	654
607	far transfer. <i>Psychological bulletin</i> , 128(4):612.	<i>Preprint</i> , arXiv:2509.04013.	655
608	Yixin Cao, Shibo Hong, Xinze Li, Jiahao Ying, Yubo	Alhassan Mumuni and Fuseini Mumuni. 2025. Large	656
609	Ma, Haiyuan Liang, Yantao Liu, Zijun Yao, Xiaozhi	language models for artificial general intelligence	657
610	Wang, Dan Huang, and 1 others. 2025. Toward gen-	(agi): A survey of foundational principles and ap-	658
611	eralizable evaluation in the llm era: A survey beyond	proaches. <i>arXiv preprint arXiv:2501.03151</i> .	659
612	benchmarks. <i>arXiv preprint arXiv:2504.18838</i> .		
613	Gheorghe Comanici, Eric Bieber, Mike Schaekermann,	Shiwen Ni, Guhong Chen, Shuaimin Li, Xuanang	660
614	Ice Pasapat, Noveen Sachdeva, Inderjit Dhillon, Mar-	Chen, Siyi Li, Bingli Wang, Qiyao Wang, Xingjian	661
615	cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke	Wang, Yifan Zhang, Liyang Fan, Chengming Li,	662
616	Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni,	Ruifeng Xu, Le Sun, and Min Yang. 2025. A sur-	663
617	Nathan Lintz, Tiago Cardal Pais, Henrik Jacobs-	vey on large language model benchmarks . <i>Preprint</i> ,	664
618	son, Idan Szpektor, Nan-Jiang Jiang, and 3416 oth-	arXiv:2508.15361.	665
619	ers. 2025. Gemini 2.5: Pushing the frontier with	OpenAI. 2025a. Introducing gpt-5. https://open	666
620	advanced reasoning, multimodality, long context,	ai.com/index/introducing-gpt-5/ . Accessed	667
621	and next generation agentic capabilities . <i>Preprint</i> ,	2025-11-03.	668
622	arXiv:2507.06261.	OpenAI. 2025b. Introducing openai o3 and o4-mini.	669
623	Margaret L Hilton and James W Pellegrino. 2012. <i>Ed-</i>	https://openai.com/index/introducing-o3-a	670
624	<i>ucation for life and work: Developing transferable</i>	nd-o4-mini/ . Accessed 2025-11-03.	671
625	<i>knowledge and skills in the 21st century</i> . National	David N Perkins, Gavriel Salomon, and 1 others. 1992.	672
626	Academies Press.	Transfer of learning. <i>International encyclopedia of</i>	673
627	Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Ji-	<i>education</i> , 2(2):6452–6457.	674
628	awei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui	David Rein, Betty Li Hou, Asa Cooper Stickland, Jack-	675
629	Yuan, Runzhe Wang, and 1 others. 2025a. Math-	son Petty, Richard Yuanzhe Pang, Julien Dirani, Ju-	676
630	perturb: Benchmarking llms' math reasoning abil-	lian Michael, and Samuel R. Bowman. 2024. GPQA:	677
631	ities against hard perturbations. <i>arXiv preprint</i>	A graduate-level google-proof q&a benchmark . In	678
632	<i>arXiv:2502.06453</i> .	<i>First Conference on Language Modeling</i> .	679
633	Shulin Huang, Linyi Yang, Yan Song, Shuang Chen,	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-	680
634	Leyang Cui, Ziyu Wan, Qingcheng Zeng, Ying	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin	681
635	Wen, Kun Shao, Weinan Zhang, and 1 others.	Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei	682
		Du, Xuancheng Ren, Rui Men, Dayiheng Liu,	683
		Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a.	684
		Qwen2-vl: Enhancing vision-language model's per-	685
		ception of the world at any resolution. <i>arXiv preprint</i>	686
		<i>arXiv:2409.12191</i> .	687

688 Siyuan Wang, Zhuohan Long, Zhihao Fan, Zhongyu
689 Wei, and Xuanjing Huang. 2024b. Benchmark self-
690 evolving: A multi-agent framework for dynamic llm
691 evaluation. *arXiv preprint arXiv:2402.11443*.

692 Xiaobao Wu, Liangming Pan, Yuxi Xie, Ruiwen Zhou,
693 Shuai Zhao, Yubo Ma, Mingzhe Du, Rui Mao,
694 Anh Tuan Luu, and William Yang Wang. 2024. An-
695 tileakbench: Preventing data contamination by auto-
696 matically constructing benchmarks with updated real-
697 world knowledge. *arXiv preprint arXiv:2412.13670*.

698 Xiaobao Wu, Liangming Pan, Yuxi Xie, Ruiwen Zhou,
699 Shuai Zhao, Yubo Ma, Mingzhe Du, Rui Mao,
700 Anh Tuan Luu, and William Yang Wang. 2025. [An-
701 tileakbench: Preventing data contamination by auto-
702 matically constructing benchmarks with updated
703 real-world knowledge](#). *Preprint*, arXiv:2412.13670.

704 Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui
705 Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong
706 Liu, Rui Men, An Yang, Jingren Zhou, and Jun-
707 yang Lin. 2025. [Group sequence policy optimization](#).
708 *Preprint*, arXiv:2507.18071.

709 A Additional Experiments

710 A.1 Does the size of the model affect the 711 transfer robustness?

712 To analyze whether scaling up model parameters
713 naturally confers transfer robustness, we conducted
714 a comparative experiment on the instruct model of
715 Qwen3-VL(Bai et al., 2025a) from 2B to 235B.
716 As shown in Table 4, we observe a clear scal-
717 ing law: increasing model size from 2B to 235B
718 yields overall performance gains across transfer set-
719 tings, though improvements are not strictly mono-
720 tonic at every scale. Notably, the largest model
721 (235B) exhibits the strongest resistance to domain
722 shifts, particularly in the challenging *KD-Far* set-
723 ting (72.5%), significantly outperforming the 8B
724 model (59.5%).

725 However, a granular inspection reveals that trans-
726 fer robustness does not scale linearly with param-
727 eter count. Specifically, we observe **diminishing
728 returns at mid-scale**. The transition from 2B to
729 4B yields a substantial gain in average robustness
730 (+13.8%), suggesting that the 2B model is severely
731 under-parameterized for complex reasoning tasks.
732 In contrast, scaling from 4B to 8B results in a ro-
733 bustness plateau: average improvement narrows
734 to 1.7%, and notably, *KD-Near* accuracy even ex-
735 hibits a marginal decline (64.7% → 64.4%). This
736 stagnation suggests that mid-scale parameter in-
737 creases, without concurrent advances in architec-
738 ture or data quality, may saturate the model’s ca-
739 pacity for surface-level pattern generalization.

740 A.2 Does Generational Evolution Enhance 741 Transfer Robustness?

742 To determine whether generational advance-
743 ments—encompassing model architecture, data
744 scaling, and training methodologies—confer in-
745 herent improvements in robustness, we conducted
746 a longitudinal evaluation of the Qwen-VL lineage:
747 Qwen2-VL-7B (Wang et al., 2024a), Qwen2.5-VL-
748 7B, and Qwen3-VL-8B. As evidenced in Table 5,
749 we observe a strict monotonic improvement across
750 all evaluated dimensions. Most notably, Qwen3-
751 VL-8B achieves state-of-the-art performance with a
752 substantial elevation in *Original* accuracy (77.5%),
753 representing a significant leap from the 43.5% base-
754 line established by its predecessor, Qwen2-VL. Be-
755 yond these aggregate gains, a granular analysis
756 reveals a fundamental dichotomy in how evolving
757 models contend with distinct perturbation types.
758 On one hand, sensitivity to surface-level format-
759 ting appears to be largely resolved. While the
760 Qwen2 architecture suffered a precipitous 11.1%
761 performance degradation when shifting context for-
762 mats (dropping from 43.5% in *Original* to 32.4%
763 in *MC-Near*), Qwen3-VL demonstrates remark-
764 able resilience, narrowing this "Context Gap" to a
765 negligible 2.3% (77.5% vs. 75.2%). This trajec-
766 tory suggests that recent architectural optimizations
767 have effectively mitigated vulnerability to modal
768 context variations. In stark contrast, the challenge
769 of Knowledge Domain (KD) transfer has intensi-
770 fied. Although Qwen3-VL secures a higher abso-
771 lute baseline, it incurs a more severe penalty when
772 generalizing to distant knowledge domains (*KD-
773 Far*), exhibiting an 18.0% drop compared to the
774 15.1% decline observed in Qwen2. This widening
775 "Knowledge Gap" underscores that while scaling
776 parameters and data enhances general capabilities,
777 it does not automatically bestow the reasoning flex-
778 ibility necessary to bridge significant semantic do-
779 main shifts.

Model	Original	KD-Near	KD-Far	MC-Near	MC-Far
Qwen3-VL-2B-Instruct	62.7%	50.3%	45.4%	56.9%	52.3%
Qwen3-VL-4B-Instruct	71.6%	64.7%	59.5%	70.9%	69.9%
Qwen3-VL-8B-Instruct	77.5%	64.4%	59.5%	75.2%	68.6%
Qwen3-VL-235B-A22B-Instruct	89.9%	80.6%	72.5%	88.9%	86.6%

Table 4: Impact of model scale on transfer robustness. Performance generally improves with parameter size, though not strictly monotonic (e.g., KD-Near slightly declines from 4B to 8B). The gain from 4B to 8B diminishes significantly compared to the leap from 2B to 4B.

Model	Original	KD-Near	KD-Far	MC-Near	MC-Far
Qwen2-VL-7B-Instruct	43.5%	26.5%	28.4%	32.4%	31.7%
Qwen2.5-VL-7B-Instruct	56.2%	48.0%	44.4%	55.6%	52.9%
Qwen3-VL-8B-Instruct	77.5%	64.4%	59.5%	75.2%	68.6%

Table 5: Longitudinal comparison of robustness across Qwen-VL generations. While Qwen3-VL significantly narrows the performance gap in Context Transfer (MC-Near), the gap in Knowledge Transfer (KD-Far) remains substantial and even widens in absolute terms.