

The Elegance Gap: When Stronger Language Models Prefer Brute-Force over Insight

Anonymous ACL submission

Abstract

As large language models (LLMs) scale up, they demonstrate remarkable improvements across many reasoning tasks. However, we identify a surprising *inverse scaling* phenomenon in mathematical problem-solving: stronger models increasingly favor computationally intensive brute-force approaches over elegant, insight-driven solutions that human experts prefer. We introduce **ERMR** (Elegance-Required Mathematical Reasoning), a benchmark of 240 problems where optimal solutions require key insights such as symmetry exploitation, geometric visualization, or variable substitution, while brute-force computation leads to either errors or excessive complexity. Through systematic evaluation of 12 state-of-the-art LLMs ranging from 7B to 405B parameters, we observe that larger models are *more likely* to attempt brute-force solutions (65% for 405B vs. 34% for 7B models) and consequently achieve *lower* success rates on problems requiring elegant reasoning. Our analysis reveals this “elegance gap” stems from models’ over-reliance on procedural patterns in training data rather than developing strategic problem-solving capabilities. We further demonstrate that targeted fine-tuning on elegant solution trajectories can partially bridge this gap, improving both solution quality and success rates. Our findings suggest that current scaling paradigms may inadvertently prioritize computational power over mathematical insight, raising important questions about reasoning capabilities in future LLMs.

1 Introduction

Consider the following mathematical problem:

Given non-zero real numbers x, y satisfying $\frac{x}{y} + \frac{y}{x} + 2xy = x^2 - y^2$, find the minimum value of $x^2 + y^2$.

A **brute-force approach** would involve algebraic manipulation, solving for one variable, and

computing derivatives—a procedure prone to computational errors across 15+ steps. In contrast, an **elegant solution** recognizes that letting $u = x + y$ and $v = xy$ transforms the constraint into a simple relationship, yielding the answer in just 5 steps through symmetry exploitation.

When we evaluated this problem across 12 recent LLMs, we discovered a counterintuitive pattern: *larger models were significantly more likely to attempt the brute-force approach and fail*. GPT-4 (with 1.76T tokens of training data) attempted brute-force computation in 73% of trials, while Llama-2-7B attempted it in only 28% of cases. This represents a form of **inverse scaling** (McKenzie et al., 2023)—where increased model capacity leads to *degraded* performance on specific tasks.

Unlike previous inverse scaling phenomena that focused on spurious correlations or social biases (Wei et al., 2023; Lin et al., 2022), our work identifies inverse scaling in *strategic reasoning*: the ability to select appropriate problem-solving strategies rather than merely executing computations. We hypothesize that larger models, trained on massive corpora containing diverse solution approaches, develop stronger pattern-matching for procedural computation, inadvertently creating a “computational overfitting” where the availability of brute-force patterns suppresses the discovery of elegant insights.

To investigate this phenomenon systematically, we make the following contributions:

- We introduce **ERMR**, a carefully curated benchmark of 240 mathematical problems spanning five categories (variable substitution, geometric reasoning, construction methods, pattern recognition, and counter-intuitive logic), where elegant solutions exist but brute-force approaches are significantly more error-prone (§3).
- Through comprehensive evaluation of 12 LLMs (7B-405B parameters), we demonstrate clear in-

verse scaling in elegance: larger models prefer brute-force approaches (Spearman $\rho = 0.73$, $p < 0.01$) and achieve 12-18% lower success rates on elegance-required problems compared to smaller models given strategic hints (§4).

- We analyze the effectiveness of four prompting strategies (zero-shot, strategic hints, one-shot elegant examples, and chain-of-thought) and show that while strategic hints can partially mitigate the elegance gap for smaller models, they become less effective as model size increases (§6).
- We conduct ablation studies showing that targeted supervised fine-tuning on elegant solution trajectories significantly improves both solution elegance ($\uparrow 34\%$) and success rates ($\uparrow 23\%$), suggesting the capability for elegant reasoning exists but is suppressed in pre-trained models (§7).
- Through interpretability analysis, we identify that larger models allocate disproportionate attention to computational tokens over strategic-planning tokens during early reasoning stages, providing mechanistic insight into the elegance gap (§8).

Our findings challenge the assumption that scaling alone leads to better reasoning, highlighting the need for training objectives that explicitly reward strategic insight and solution elegance rather than mere computational throughput.

2 Related Work

2.1 Mathematical Reasoning in LLMs

Recent work has demonstrated impressive capabilities of LLMs on mathematical reasoning tasks (Cobbe et al., 2021; Hendrycks et al., 2021; Polu et al., 2022). Benchmarks like GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) have driven significant progress, with models like GPT-4 (OpenAI, 2023) and Minerva (Lewkowycz et al., 2022) achieving near-human performance on competition-level problems. However, these benchmarks primarily evaluate *correctness* rather than *solution quality* or *strategic reasoning*.

Recent work has begun examining solution diversity (Uesato et al., 2022) and multi-step reasoning (Wei et al., 2022), but little attention has been paid to whether models discover elegant solutions or rely on brute-force computation. Our work fills this gap by explicitly evaluating the *quality* and *efficiency* of solution strategies.

2.2 Inverse Scaling Laws

The inverse scaling prize (McKenzie et al., 2023) identified several tasks where larger models perform worse, often due to: (1) stronger pattern-matching to misleading surface features (McCoy et al., 2019), (2) increased sensitivity to prompt formatting (Webson and Pavlick, 2022), or (3) overconfidence in incorrect reasoning (Kadavath et al., 2022).

Our findings reveal a distinct mechanism: larger models’ superior pattern-matching capabilities can lead them to favor *computationally intensive but error-prone* strategies over *insightful but less common* approaches. This represents a form of “capability-induced failure”—where increased capacity paradoxically leads to worse outcomes.

2.3 Human Mathematical Problem-Solving

Cognitive science literature extensively documents human expertise in mathematical problem-solving, emphasizing the role of: **strategic knowledge** (Schoenfeld, 1985), **meta-cognitive monitoring** (Pólya, 1945), and **representational insight** (Kahneman, 2011).

The Einstellung effect (Luchins, 1942) describes how prior experience can blind problem-solvers to more efficient solutions—a phenomenon remarkably similar to our observations in LLMs. Our work bridges cognitive science and NLP by demonstrating that LLMs may exhibit analogous cognitive biases at scale.

3 The ERMR Benchmark

We construct the **Elegance-Required Mathematical Reasoning (ERMR)** benchmark to systematically evaluate whether LLMs can discover and apply elegant problem-solving strategies.

3.1 Design Principles

Each problem in ERMR satisfies the following criteria:

Dual Solution Pathways. Both an elegant (insight-driven) and a brute-force (computation-heavy) solution exist. The elegant solution requires recognizing a key mathematical structure (e.g., symmetry, substitution, geometric interpretation), while brute-force proceeds through algebraic manipulation.

Computational Asymmetry. The elegant solution involves ≤ 5 major steps, while brute-force re-

179 requires ≥ 10 steps. Brute-force approaches have significantly higher error rates due to computational complexity.
180
181

182 **Expert Consensus.** Three mathematics educators independently verify that the elegant solution
183 represents the “standard” or “expected” approach
184 in educational contexts.
185

186 **Answer Verifiability.** All problems have definite
187 numerical or symbolic answers that can be automatically verified, enabling objective evaluation.
188

189 **3.2 Problem Categories**

190 ERMR contains 240 problems across five categories:
191

192 **Variable Substitution (60 problems).** Problems
193 where introducing auxiliary variables (e.g.,
194 sum and product, trigonometric substitution) dra-
195 matically simplifies the problem structure.
196

197 *Example:* Given $x, y, z > 0$ and $x^2 + y^2 + z^2 =$
198 1, find the range of $x + y + z - xyz$.

199 *Elegant Approach:* Recognize this as a symmetric
200 function; apply Cauchy-Schwarz inequality.
201

202 *Brute-Force:* Solve using Lagrange multipliers
203 with extensive algebraic manipulation.
204

205 **Geometric Visualization (48 problems).** Prob-
206 lems in algebra that admit geometric interpreta-
207 tions, where visualization provides immediate in-
208 sight.
209

210 *Example:* Given $x^2 + y^2 = x^2 + z^2 + \sqrt{3}xz =$
211 $y^2 + z^2 + yz = 16$, find $2xy + xz + \sqrt{3}yz$.

212 *Elegant Approach:* Interpret as dot products of
213 vectors; recognize geometric configuration.
214

215 *Brute-Force:* Solve the system of three equations
216 through elimination.
217

218 **Construction Methods (52 problems).** Prob-
219 lems where constructing auxiliary functions or ap-
220 plying specific inequalities yields immediate re-
221 sults.
222

223 **Pattern Recognition (44 problems).** Problems
224 with underlying patterns or recurrence relations
225 that, once identified, reduce complexity dramati-
226 cally.
227

228 **Counter-Intuitive Logic (36 problems).** Problems
229 where the “obvious” approach fails,
230 requiring proof by contradiction or extremal principles.
231

232 **3.3 Dataset Construction and Validation**

233 We sourced problems from: (1) Mathematical
234 olympiad archives (IMO, USAMO, AIME), (2)
235 University entrance examinations (China, India,
236 Russia), (3) Advanced calculus and analysis text-
237 books.
238

239 All problems were filtered to ensure they do not
240 appear verbatim in common LLM training corpora
241 (verified through exact and fuzzy matching against
242 public datasets).
243

244 Three annotators with graduate-level mathemati-
245 cics training independently: (1) Verified the exis-
246 tence of both elegant and brute-force solutions, (2)
247 Counted solution steps for each approach, (3) Rated
248 the relative “elegance” of the optimal solution on a
249 5-point scale.
250

251 Inter-annotator agreement (Krippendorff’s α)
252 was 0.82 for elegance ratings and 0.91 for step
253 counting.
254

255 **3.4 Evaluation Metrics**

256 For each problem and model, we collect:
257

258 **Correctness.** Whether the final answer matches
259 the ground truth (binary).
260

261 **Solution Category.** Human annotators classify
262 the solution approach as: *Elegant*, *Brute-Force*,
263 *Hybrid*, or *Invalid* (flawed reasoning).
264

265 **Step Count.** Number of major mathematical op-
266 erations performed.
267

268 **Elegance Score.** A composite metric combining
269 step efficiency, structural insight, and generalizabil-
270 ity:
271

$$272 \text{Elegance} = \alpha \cdot \frac{1}{\text{Steps}} + \beta \cdot \text{Insight} + \gamma \cdot \text{Correctness}$$

273 where $\alpha = 0.3$, $\beta = 0.5$, $\gamma = 0.2$ based on expert
274 weighting.
275

276 **4 Experimental Setup**

277 **4.1 Models Evaluated**

278 We evaluate 12 state-of-the-art LLMs spanning dif-
279 ferent scales and architectures:
280

281 **Open-source models:** Llama-2 (7B, 13B,
282 70B) (Touvron et al., 2023), Mistral (7B) (Jiang
283 et al., 2023), DeepSeek-Math (7B, 70B) (Shao
284 et al., 2024).
285

286 **Closed-source models:** GPT-3.5-Turbo, GPT-
287 4-0613, GPT-4-Turbo (OpenAI, 2023), Claude-
288 2, Claude-3-Opus (Anthropic, 2023), Gemini-
289 Pro (Gemini Team, 2023).
290

265	4.2 Prompting Strategies	310
266	We test four prompting strategies to understand	311
267	how models respond to different levels of guidance:	312
268	Zero-Shot (ZS). Problem statement only:	313
269	<i>Solve the following problem: [PROB-</i>	314
270	<i>LEM]</i>	315
271	Strategy Hint (SH). Problem + high-level strate-	316
272	tic suggestion:	317
273	<i>Solve the following problem. Hint: Con-</i>	318
274	<i>sider using [variable substitution / geo-</i>	319
275	<i>metric interpretation / symmetry].</i>	320
276	One-Shot Elegant (OSE). One example of an	321
277	elegant solution to a similar problem, followed by	322
278	the target problem.	323
279	Chain-of-Thought (CoT). Standard CoT	324
280	prompting (Wei et al., 2022):	325
281	<i>Let's solve this step by step.</i>	326
282	4.3 Implementation Details	327
283	For all models, we use: - Temperature: 0.7 (to	328
284	encourage diverse strategies) - Max tokens: 2048 -	329
285	Top-p: 0.95 - 5 independent samples per problem	330
286	to estimate variance	331
287	Solutions are parsed using regex patterns to ex-	332
288	tract numerical answers. Two human annotators	333
289	(graduate students in mathematics) independently	334
290	classify solution approaches, with disagreements	335
291	resolved through discussion.	336
292	5 Main Results: The Inverse Scaling of	337
293	Elegance	338
294	5.1 Larger Models Prefer Brute-Force	339
295	Table ?? presents our core findings. Across all	340
296	problem categories, we observe a consistent trend:	341
297	larger models increasingly favor brute-force ap-	342
298	proaches over elegant solutions.	343
299	[TABLE 1: Main results showing model size vs.	344
300	brute-force preference and success rate] 347	345
301	Specifically: 348	346
302	<ul style="list-style-type: none"> Brute-force attempts increase monotonically 349 with model size (Spearman $\rho = 0.73$, $p < 0.01$) 350 Llama-2-7B attempts brute-force in 34% of 351 cases; GPT-4-Turbo in 68% 352 This preference correlates with <i>decreased</i> 353 success rates: models attempting brute-force 354 succeed 41% of the time vs. 76% for elegant 355 approaches 356 	347
303		348
304		349
305		350
306		351
307		352
308		353
309		354
310		355
311		356
312		
313	5.2 Category-Specific Analysis	310
314	The inverse scaling phenomenon varies across prob-	311
315	lem categories:	312
316	Variable Substitution. Largest gap observed	313
317	(26% performance decrease for 70B+ models).	314
318	Larger models tend to directly manipulate origi-	315
319	nal variables rather than recognizing substitution	316
320	opportunities.	317
321	Geometric Visualization. Moderate inverse	318
322	scaling (15% decrease). Models rarely generate	319
323	geometric interpretations spontaneously; most	320
324	proceed algebraically.	321
325	Construction Methods. Smallest gap (8% de-	322
326	crease). Some larger models successfully identify	323
327	auxiliary functions, though less consistently than	324
328	elegant solutions.	325
329	[TABLE 2: Per-category breakdown]	326
330		
331		
332		
333		
334		
335		
336		
337		
338		
339		
340		
341		
342		
343	5.3 Error Analysis	327
344	We manually analyzed 200 failed solutions across	328
345	model sizes:	329
346	Computational errors (52%). Mistakes in al-	330
347	gebraic manipulation, especially in brute-force so-	331
348	lutions with 10+ steps.	332
349	Strategic errors (31%). Choosing an approach	333
350	that cannot lead to the solution (e.g., incorrect sub-	334
351	stitution).	335
352	Incomplete reasoning (17%). Starting with a	336
353	promising approach but failing to complete it.	337
354	Larger models disproportionately suffer from	338
355	<i>computational errors</i> in brute-force attempts, while	339
356	smaller models more often make <i>strategic errors</i> .	340
357	This suggests larger models have stronger execu-	341
358	tion capabilities but poorer strategy selection.	342
359		
360	6 Prompt Engineering and Its	343
361	Limitations	344
362	6.1 Effect of Strategic Hints	345
363	Providing explicit strategic hints (e.g., “Con-	346
364	sider variable substitution”) improves performance	347
365	across all models:	348
366	[TABLE 3: Performance with different prompt-	349
367	ing strategies]	350
368	However, the benefit <i>decreases</i> with model size:	351
369	Llama-2-7B: +28% success rate with hints	352
370	GPT-4-Turbo: +11% success rate with hints	353
371	This suggests larger models have stronger priors	354
372	toward brute-force computation that are harder to	355
373	override with prompting alone.	356

357	6.2 One-Shot Learning	402
358	Providing an elegant solution example yields mixed	403
359	results: - Improves solution elegance (models more	
360	often adopt similar structural approaches) - Does	
361	<i>not</i> significantly improve correctness for larger	
362	models - Smaller models benefit more from explicit	
363	demonstrations	
364	6.3 Chain-of-Thought Amplifies the Problem	406
365	Surprisingly, CoT prompting <i>exacerbates</i> the in-	407
366	verse scaling: - Encourages verbose, step-by-step	408
367	computation - Larger models produce longer rea-	409
368	soning chains (avg. 247 tokens vs. 156 for smaller	410
369	models) - These longer chains are more error-prone	411
370	(58% success vs. 64% without CoT for 70B+ mod- els)	412
371	This finding challenges the prevailing assumption	413
372	that CoT universally improves reasoning quality.	414
373		415
374		416
375	7 Fine-Tuning for Elegance	417
376	7.1 Training Data Construction	
377	We construct two SFT datasets: Elegant-SFT : 300	418
378	problems with elegant solution trajectories Brute-	419
379	SFT : 300 problems with brute-force solution tra-	
380	jectories (control)	
381	Each trajectory includes: (1) Problem statement	
382	(2) Step-by-step solution with strategic justifica-	
383	tions (3) Final answer	
384	Solutions are written by mathematics educators	
385	to emphasize: - Explicit strategy selection (“We	
386	use substitution because...”) - Meta-cognitive state-	
387	ments (“This approach is more efficient than...”) -	
388	Generalizable patterns (“This technique applies to	
389	similar problems where...”)	
390	7.2 Fine-Tuning Setup	
391	We fine-tune Llama-2-7B and Llama-2-13B using	
	LoRA (Hu et al., 2022): - Rank: 16 - Learning rate:	
	3e-4 - Batch size: 32 - Epochs: 3	
392	7.3 Results	
393	Fine-tuning on elegant trajectories yields substantial	435
394	improvements:	436
	[TABLE 4: Fine-tuning results]	437
	Key findings:	438
	• Elegant-SFT models use elegant approaches in	
	71% of cases (vs. 42% baseline)	
	Success rate improves by 23% on ERMR test set	
	• Improvements generalize to held-out problem	402
	categories	403
	• Brute-SFT models show <i>decreased</i> performance	404
	(confirming brute-force approach is detrimental)	405
	7.4 Generalization Analysis	406
	We test fine-tuned models on out-of-distribution	407
	problems: - MATH benchmark (Hendrycks et al.,	408
	2021): 12% improvement on competition problems	409
	- GSM8K (Cobbe et al., 2021): 3% improvement	410
	(ceiling effect; already high performance) - Novel	411
	ERMR variants : 18% improvement when prob- lem parameters are varied	412
	This demonstrates that learning elegant reason- ing strategies transfers beyond the training distribu- tion.	413
		414
		415
		416
	8 Mechanistic Analysis	417
	To understand <i>why</i> larger models favor brute-force,	418
	we conduct interpretability analyses:	419
	8.1 Attention Pattern Analysis	420
	Using attention rollout (Abnar and Zuidema, 2020),	421
	we analyze which input tokens receive high atten- tion during solution generation:	422
	[FIGURE 1: Attention heatmaps for elegant vs. brute-force solutions]	423
	Key observations:	424
	• Models generating elegant solutions allocate 63%	425
	of early-layer attention to <i>constraint terms</i> (e.g.,	426
	$x^2 + y^2 = 1$)	427
	• Models generating brute-force solutions dis- tribute attention more uniformly (34% to con- straints, 42% to variable names)	428
	• Larger models show <i>weaker</i> attention concentra- tion on constraint terms	429
	This suggests elegant reasoning requires focused	430
	attention on problem structure, while brute-force	431
	proceeds through distributed, token-by-token pro- cessing.	432
		433
		434
	8.2 Activation Probing	435
	We train linear probes on intermediate activations	436
	to predict solution strategy:	437
	[FIGURE 2: Probe accuracy across layers]	438
	Findings: - Strategy decisions crystallize in lay- ers 8-12 for 7B models, layers 15-22 for 70B mod- els - Earlier layers encode problem type; middle layers commit to strategy; later layers execute -	

447 Larger models show *later* strategy commitment,
448 suggesting longer deliberation but potentially miss-
449 ing early structural cues

450 8.3 Neuron Attribution

451 Using integrated gradients (Sundararajan et al.,
452 2017), we identify neurons most predictive of ele-
453 gant vs. brute-force strategies:

454 - “Symmetry neurons” (high activation on sym-
455 metric expressions) are more influential in smaller
456 models - “Computation neurons” (high activation
457 on arithmetic operations) dominate in larger mod-
458 els - Fine-tuning on elegant solutions increases
459 symmetry neuron influence

460 This provides preliminary evidence that training
461 data composition shapes which reasoning modes
462 become dominant.

463 9 Discussion

464 9.1 Why Do Larger Models Fail?

465 We propose three interrelated explanations:

466 **Pattern Over-Matching.** Larger models, trained
467 on more diverse data, encounter more examples
468 of brute-force solutions in educational materials,
469 textbooks, and online Q&A forums. This creates
470 stronger priors for computational approaches, even
471 when inefficient.

472 **Lack of Meta-Cognitive Oversight.** Current
473 LLMs lack explicit mechanisms for *strategy se-
474 lection* before execution. They immediately begin
475 generating solutions without evaluating multiple
476 approaches.

477 **Training Objective Misalignment.** Standard
478 language modeling objectives reward *any* path to
479 the correct answer equally. There is no signal to pre-
480 fer elegant over brute-force solutions, potentially
481 allowing models to “shortcut” through computation
482 rather than developing insight.

483 9.2 Implications for LLM Development

Our findings suggest several directions for improving reasoning in LLMs:

Elegance-Aware Training Objectives. Reward
models could explicitly score solution elegance,
penalizing unnecessarily complex approaches.
532
533
534

Strategy-Selection Modules. Incorporating explicit planning or strategy-selection stages before solution execution.
535
536
537

Curriculum Learning. Training models first on elegant solutions, then gradually introducing brute-force approaches as fallbacks.
492
493
494

Synthetic Data Generation. Generating diverse elegant solutions to counterbalance brute-force prevalence in web-scraped data.
495
496
497

498 9.3 Limitations and Future Work

Our study has several limitations:
499

Human Annotation. Solution classification relies on human judgment, though we achieve high inter-annotator agreement.
500
501
502

Domain Scope. We focus on mathematical reasoning; the elegance gap may manifest differently in other domains (coding, scientific reasoning, etc.).
503
504
505
506

Prompt Sensitivity. While we test four prompting strategies, the vast prompt space remains underexplored.
507
508
509

Future work should: (1) Extend ERMR to other domains (programming, physics problems, logical puzzles), (2) Investigate whether the elegance gap appears during training (through analysis of intermediate checkpoints), (3) Explore whether multi-agent or self-critique approaches can help models discover elegant solutions.
510
511
512
513
514
515
516

517 10 Conclusion

We have identified and characterized a novel form of inverse scaling in large language models: the **elegance gap**, where larger models increasingly favor computationally intensive brute-force approaches over insightful elegant solutions in mathematical reasoning.
518
519
520
521
522
523

Through the ERMR benchmark and comprehensive evaluation of 12 LLMs, we demonstrate that this phenomenon is robust across model families and problem categories. Our mechanistic analyses reveal that larger models allocate attention and computational resources toward token-level execution rather than structural insight, suggesting fundamental limitations in how current models approach strategic reasoning.
524
525
526
527
528
529
530
531

Importantly, we show that targeted fine-tuning on elegant solution trajectories can significantly mitigate this gap, indicating that the *capability* for elegant reasoning exists but is suppressed in standard pre-training.
532
533
534
535
536
537

These findings have important implications for the development of future LLMs: simply scaling model size and data may not lead to more insightful reasoning. Instead, we need training paradigms that explicitly encourage strategic thinking, reward solution quality alongside correctness, and develop meta-cognitive capabilities for approach selection.

As LLMs become increasingly integrated into education, scientific research, and decision-making, ensuring they exhibit not just computational power but also reasoning elegance becomes critical. Our work takes a first step toward understanding and addressing this challenge.

Limitations

Following ACL guidelines, we acknowledge the following limitations:

Benchmark Scope. ERMR focuses on mathematical reasoning and may not generalize to all domains requiring strategic insight (e.g., creative writing, open-ended problem solving).

Evaluation Subjectivity. While we achieve high inter-annotator agreement (Krippendorff’s $\alpha = 0.82$), the classification of solutions as “elegant” vs. “brute-force” involves subjective judgment. Different mathematics educators might disagree on borderline cases.

Model Access. For closed-source models (GPT-4, Claude, Gemini), we rely on API access without insight into architecture, training data, or internal mechanisms, limiting the depth of mechanistic analysis.

Prompt Engineering. Despite testing multiple prompting strategies, we cannot exhaustively explore the prompt space. It is possible that carefully crafted prompts could mitigate the elegance gap more effectively than our tested approaches.

Training Data Contamination. While we filtered problems to avoid verbatim matches with common training corpora, we cannot guarantee that similar problems or solution strategies were not encountered during pre-training.

Sample Size for Fine-Tuning. Our fine-tuning experiments use 300 training examples due to computational constraints. Larger-scale fine-tuning might yield different results.

Human Performance Baseline. We do not include systematic human evaluation on ERMR, making it difficult to assess whether models are approaching, exceeding, or falling short of human expert performance.

Ethics Statement

Responsible Benchmarking. ERMR is designed to evaluate reasoning quality, not to be “gamed” by models. We commit to regular updates if models begin over-fitting to our benchmark.

Educational Impact. If deployed in educational contexts, systems exhibiting the elegance gap might teach students inefficient problem-solving approaches. We advocate for careful evaluation before deploying LLMs as tutoring systems.

Accessibility. We will release ERMR under an open license to ensure broad access for academic research, though we note that evaluating large proprietary models may require significant compute resources or API costs.

Misuse Potential. ERMR could be misused to create adversarial inputs that expose LLM weaknesses. However, we believe transparency about limitations benefits the research community more than obscurity.

Acknowledgments

We thank the anonymous reviewers for their valuable feedback. We are grateful to [names removed for anonymous review] for assistance with benchmark construction and annotation. This work was supported by [funding sources removed for anonymous review]. Computational resources were provided by [removed for anonymity].

References

Samira Abnar and Willem Zuidema. 2020. Quantifying attention flow in transformers. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Anthropic. 2023. Claude 2. Available at <https://www.anthropic.com>.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reichiro Nakano, Christopher Hesse, and John Schulman. 2021. GSM8K: Training verifiers to solve math word problems. In *Advances in Neural Information Processing Systems*.

630	Gemini Team.	2023.	Gemini: A family of highly capable multimodal models.	<i>arXiv preprint arXiv:2312.11805</i> .	685
631					686
632					
633	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt.	2021.	Measuring mathematical problem solving with the MATH dataset.	<i>arXiv preprint arXiv:2103.03874</i> .	687
634					688
635					
636					
637					
638	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.	2022.	LoRA: Low-rank adaptation of large language models.	<i>International Conference on Learning Representations</i> .	689
639					690
640					691
641					692
642					693
643	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, and 1 others.	2023.	Mistral 7b.	<i>arXiv preprint arXiv:2310.06825</i> .	694
644					695
645					696
646					697
647					
648	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, and 1 others.	2022.	Language models (mostly) know what they know.	<i>arXiv preprint arXiv:2207.05221</i> .	698
649					699
650					700
651					701
652					702
653					703
654	Daniel Kahneman.	2011.	<i>Thinking, Fast and Slow</i> .	Farrar, Straus and Giroux.	704
655					705
656	Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, and 1 others.	2022.	Solving quantitative reasoning problems with language models.	<i>arXiv preprint arXiv:2206.14858</i> .	706
657					707
658					708
659					709
660					
661					
662	Stephanie Lin, Jacob Hilton, and Owain Evans.	2022.	TruthfulQA: Measuring how models mimic human falsehoods.	<i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics</i> .	710
663					711
664					712
665					713
666	Abraham S Luchins.	1942.	Mechanization in problem solving: The effect of einstellung.	<i>Psychological Monographs</i> , 54(6):1–95.	714
667					
668					
669	Tom McCoy, Ellie Pavlick, and Tal Linzen.	2019.	Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference.	<i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> .	715
670					716
671					717
672					718
673					719
674	Ian R McKenzie, Alexander Lyzhov, Michael Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Aaron Kirtland, Alexis Ross, Alisa Liu, and 1 others.	2023.	Inverse scaling: When bigger isn't better.	<i>arXiv preprint arXiv:2306.09479</i> .	720
675					
676					
677					
678	OpenAI.	2023.	GPT-4 technical report.	<i>arXiv preprint arXiv:2303.08774</i> .	721
679					722
680					723
681					724
682					
683	Stanislas Polu, Jesse Michael Han, Kunhao Zheng, Mantas Baksys, Igor Babuschkin, and Ilya Sutskever.	2022.	Formal mathematics statement curriculum learning.	<i>arXiv preprint arXiv:2202.01344</i> .	725
684					726
685					727
686					
687					
688					
689					
690					
691					
692					
693					
694					
695					
696					
697					
698					
699					
700					
701					
702					
703					
704					
705					
706					
707					
708					
709					
710					
711					
712					
713					
714					
715					
716					
717					
718					
719					
720					
721					
722					
723					
724					
725					
726					
727					

A ERMR Example Problems

We provide complete examples from each category in ERMR:

A.1 Variable Substitution

Problem: Given non-zero real numbers x, y satisfying

$$678 \quad \frac{x}{y} + \frac{y}{x} + 2xy = x^2 - y^2,$$

find the minimum value of $x^2 + y^2$.

Elegant Solution: Let $u = x + y$ and $v = xy$. The constraint becomes:

$$683 \quad \frac{u^2 - 2v}{v} + 2v = u^2 - 2v$$

730	Simplifying: $u^2 = 2v(u^2 - 1)$... [5 steps total] \Rightarrow	D Fine-Tuning Dataset Examples	773
731	Answer: $\min(x^2 + y^2) = 2$.		
732	Brute-Force Solution: Solve for y in terms of	Example Training Instance:	774
733	x , substitute into objective, compute derivative,		
734	solve cubic equation... [15+ steps, error-prone]		
735	A.2 Geometric Visualization		
736	Problem: Given $x^2 + y^2 = x^2 + z^2 + \sqrt{3}xz =$	Problem: [problem statement]	775
737	$y^2 + z^2 + yz = 16$, find $2xy + xz + \sqrt{3}yz$.	Solution:	776
738	Elegant Solution: Interpret as $ \mathbf{a} ^2 = \mathbf{b} ^2 =$	[<i>Strategy Selection</i>] Looking at this prob-	777
739	$ \mathbf{c} ^2 = 16$ where... [Geometric approach using dot	lem, I notice the expression is symmetric	778
740	products, 4 steps] Answer: 16.	in x and y . This suggests that using sub-	779
741	Brute-Force Solution: Solve the system of	stitution $u = x + y$ and $v = xy$ might	780
742	three equations for x, y, z explicitly... [12+ steps	simplify the problem significantly.	781
743	with complex algebra]	[<i>Step 1</i>] Let $u = x + y$ and $v = xy$.	782
744	A.3 Additional Examples	Then $x^2 + y^2 = u^2 - 2v$.	783
745	[Include 3-5 more representative problems from	[<i>Step 2</i>] Substituting into the constraint...	784
746	other categories]	[Continue with justified steps]	785
747	B Full Experimental Results	[<i>Meta-cognitive note</i>] This approach is	786
748	B.1 Model-by-Model Breakdown	more efficient than solving for one vari-	787
749	[TABLE A1: Detailed results for all 12 models on	able because it preserves the symmetric	788
750	all 5 categories]	structure, reducing computational com-	789
751	B.2 Prompt Variations	plexity from $O(n^3)$ to $O(n)$.	790
752	[TABLE A2: Performance across different prompt	Answer: [final answer]	791
753	formulations]	[Include 5-10 representative examples]	792
754	C Annotation Guidelines	E Attention Visualization Details	793
755	We provide the complete annotation guidelines	[FIGURE A1: Complete attention patterns for dif-	794
756	given to human evaluators:	ferent model sizes]	795
757	C.1 Solution Classification Rubric	[FIGURE A2: Layer-by-layer attention evolu-	796
758	Elegant Solution - Check all that apply:	tion]	797
759	<ul style="list-style-type: none"> • ≤ 5 major computational steps 	F Code and Data Availability	798
760	<ul style="list-style-type: none"> • Uses a key insight (substitution, symmetry, geo- 	Upon acceptance, we will release:	799
761	metric interpretation)	<ul style="list-style-type: none"> • Complete ERMR benchmark with all 240 prob- 	800
	<ul style="list-style-type: none"> • Would be considered the “standard” or “textbook” 	lems and solutions	801
	solution	<ul style="list-style-type: none"> • Evaluation code and metrics computation 	802
	<ul style="list-style-type: none"> • Generalizes to similar problems 	<ul style="list-style-type: none"> • Model outputs for all 12 evaluated LLMs 	803
	Brute-Force Solution - Check all that apply:	<ul style="list-style-type: none"> • ⁷⁶² Fine-tuning scripts and training data 	804
	<ul style="list-style-type: none"> • ≥ 10 major computational steps 	<ul style="list-style-type: none"> • ⁷⁶³ Annotation guidelines and inter-annotator agree- 	805
	<ul style="list-style-type: none"> • Proceeds through algebraic manipulation without 	ment data	806
	structural insight	⁷⁶⁴ Repository: [anonymized for review]	807
	<ul style="list-style-type: none"> • Would be considered “inelegant” by math edu- 	⁷⁶⁵	
	cators	⁷⁶⁶	
	<ul style="list-style-type: none"> • Specific to the particular problem instance 	⁷⁶⁷	
	[Complete rubric with examples...]	⁷⁶⁸	
		⁷⁶⁹	
		⁷⁷⁰	
		⁷⁷¹	
		⁷⁷²	