

Haiku to Opus in Just 10 bits: LLMs Unlock Massive Compression Gains

Roy Rinberg^{1*} Annabelle Michael Carrell^{2†} Simon Henniger¹
Nicholas Carlini³ Keri Warr³

¹Harvard University ²University of Cambridge ³Anthropic

Abstract

We study the compression of LLM-generated text across lossless and lossy regimes, characterizing a compression-compute frontier where more compression is possible at the cost of more compute. For lossless compression, domain-adapted LoRA adapters can improve LLM-based arithmetic coding by $2\times$ over compression with the base LLM alone. For lossy compression, prompting a model for a succinct rewrite then applying arithmetic coding can achieve compression ratios of approximately 0.03, a $2\times$ improvement over compressing the original response.

We further introduce Question-Asking compression (QA), an interactive lossy protocol inspired by the game “Twenty Questions”. A small model iteratively refines its response by asking yes/no questions to a stronger model, transferring exactly one bit per answer. On 8 benchmarks spanning math, science, and code, 10 binary questions recover 23% to 72% of the capability gap between a small and large model on standard benchmarks and 7% to 38% on harder benchmarks, achieving compression ratios of 0.0006 to 0.004. This is over $100\times$ smaller than prior LLM-based compression (Delétang et al., 2024), suggesting that interactive protocols can transfer knowledge far more efficiently than transmitting full responses.

1 Introduction

At an abstract level, compression can be viewed as a procedure that takes two inputs: a probability distribution and a sequence. It produces an encoding in which symbols assigned higher probability are represented using fewer bits.

Arithmetic coding is a canonical example of this framing; given a model that assigns probabilities to each next symbol, arithmetic coding converts those probabilities into a compact binary representation whose length closely matches the negative log-likelihood of the sequence under the model. The quality of the compression therefore depends directly on the quality of the probability distribution.

Autoregressive language models, in turn, repeatedly produce probability distributions over the next token and sample from them. From this perspective, language models provide a natural mechanism for high-quality distribution estimation over text. By supplying an LLM’s predictive distribution to a compression scheme, one can directly translate modeling accuracy into improved compression performance.

Recent work has applied LLMs directly as inputs to arithmetic coding to produce lossless compression schemes (Bellard, 2023; Valmeekam et al., 2023; Delétang et al., 2024). These techniques

*Corresponding author: royrinberg@g.harvard.edu. Work conducted while RR was participating in ML Alignment & Theory Scholars (MATS).

†Work done while AC was a visiting student at Harvard.

already boast impressive compression ratios; whereas a standard compression scheme like GZIP (Deutsch, 1996) achieves a compression ratio of around 0.33 on most text, LLM-based arithmetic coding schemes achieve a compression ratio closer to 0.1 on text that is in distribution for the model, like Wikipedia (Delétang et al., 2024). In Section 4.2, we explore improving these schemes for specific domains through lightweight topic adaptation using LoRA (Low-Rank Adaptation), a parameter-efficient fine-tuning method that adds a small number of trainable weights which can be attached to or removed from a base model efficiently.

While already impressive, we observe that the requirement of producing a lossless compression scheme is overly restrictive in many textual settings. Most text-based interfaces do not need to produce a specific sequence of tokens, but instead only need to return a response that is useful to the user (Zhang et al., 2025). Framed this way, compression only needs to preserve semantic meaning rather than exact text. This is particularly true in the LLM chatbot setting, where users have no access to the underlying model and therefore cannot tell whether the response they receive differs from the original model output.

Modern language models can summarize, rewrite, and edit text in an open-ended manner, enabling a new form of lossy compression in which the surface form of the text may change, but its semantic content and usefulness are largely preserved. Unlike traditional compression schemes, this form of compression does not aim to reconstruct the original string exactly, but instead produces an alternative representation that remains a valid and high-quality response.

In Section 2.2 we introduce a novel scheme (QA-compression protocol) for communicating knowledge concisely from a more capable model to a less capable model. The scheme, inspired by the game 20 questions, works by having the less-capable model ask the more-capable model yes/no questions, and then having the less-capable model update its response given the more-capable model’s response.

Not only do we find that this provides meaningful compression, this also offers a lens on the scientific process itself. When a chemist wants to synthesize a new chemical, they form a hypothesis, run an experiment, and get back a low-bit answer: it worked, or it didn’t. Then, if the experiment fails, the scientist returns to the problem with a new question: “Was the temperature wrong?” “Did I use the wrong catalyst?” In this light, each experiment is a yes/no question posed to *nature*. This maps onto QA-compression: the scientist is the small, less-capable model, which formulates candidate solutions and asks binary questions; and *nature* is the LLM providing yes/no feedback. This framing and this paper provides quantifiable support to the pedagogical adage that *the hard part of science is asking a good question*.

1.1 Contributions

In this work, we investigate methods for compressing LLM-generated text using a smaller language model. Our focus is on minimizing the number of bits required to encode text produced by a large model when the decoder has access to a weaker shared model. We explore approaches that trade additional computation for improved compression ratios, in both the lossless and lossy regimes.

Concretely, this paper makes the following contributions:

1. **Iterative Binary Questions.** We propose a novel protocol akin to the children’s game “20 Questions” where a small language model iteratively improves its response through yes/no questions answered by a more capable model, which we call Question-Asking compression (QA-compression) (Section 3). We evaluate this protocol on 8 benchmarks on science, math, and coding tasks. In 10 bits or fewer, models recover 23–72% of the performance gap on easier

benchmarks and 7–38% on harder benchmarks, achieving compression ratios of 0.004 to 0.0006, over $100\times$ beyond prior SOTA (Delétang et al., 2024).

- 2. Response Rewriting.** As an intuitive baseline for lossy compression, we explore two lossy compression approaches: Shortest-of-N selection (generate multiple candidates, select the most compressible) and Summarization (explicitly request a succinct rewrite). Shortest-of-N achieves compression ratios of 0.073 ($N=1$) to 0.063 ($N=10$); Summarization achieves 0.039 to 0.034, a $1.9\times$ improvement over Shortest-of-N. Selecting for compressibility does not degrade answer quality (Section 4.1), when evaluated on a coding and math benchmark.
- 3. Lossless Compression via LoRA Adaptation.** As an innovation on existing work on lossless compression, we demonstrate that domain-specific LoRA adapters improve arithmetic coding compression by $2\times$ over baseline models. A RAG-based router selects the appropriate adapter from the prompt alone, retaining a $1.9\times$ improvement (Section 4.2), when evaluated on the LMSYS and WildChat dataset.

1.2 Paper Organization

Section 2 reviews background on arithmetic coding, LoRA adaptation, and interactive compression. Section 3 introduces a novel interactive question-asking protocol for extreme compression. Section 4.2 introduces our proposed lossless compression method based on domain-adapted LoRAs, which we present as a strong standalone baseline for comparison with the interactive protocol. Section 4.1 similarly proposes response rewriting as a lossy compression baseline. We conclude with a discussion of applications to model security and future directions in Section 5.

Models. Throughout this work, we use three tiers of language models from the Anthropic Claude 4.5 family as reference implementations: **Haiku** (small), **Sonnet** (medium), and **Opus** (large). These span a representative capability range, with Haiku serving as the small language model (SLM) in our compression protocols and Opus as the large language model (LLM) that provides guidance. Our methods are not specific to these models; we use the Anthropic API as a convenient testbed for studying compression across capability tiers.

2 Background & Related Work

2.1 Arithmetic Coding and LLM Compression

At the most abstract level, a compression scheme is a function that takes two inputs: a probability distribution, and a sequence; it then encodes the sequence according to the probability distribution, such that more likely elements are encoded with fewer bits than less likely elements. Arithmetic coding is a specific compression scheme that achieves near-optimal compression (Langdon and Rissanen, 1981; Rissanen and Langdon, 1979).

The source coding theorem from information theory actually establishes the minimum achievable compression ratio (on average), and states that for a discrete memoryless source with distribution P , the minimal length for compression is equivalent to the entropy $H(P) = \mathbb{E}_{x \sim P}[-\log_2 P(x)]$. The compressibility is directly determined by how closely the probability distribution approximates the true data-generating distribution. When data sampled from distribution P is encoded using a probabilistic model Q , the expected number of bits required is given by the cross-entropy $H(P, Q) = \mathbb{E}_{x \sim P}[-\log_2 Q(x)]$.

At a similarly abstract level, LLMs can be broken down into an algorithm that repeatedly does two things: generate a probability distribution, then sample from it. When viewed from this perspective, a natural idea is to take the probability distribution from an open-source LLM and use the as the probability-distribution used for encoding in arithmetic coding.

A standard metric in compression is the compression ratio, defined as the size of the compressed representation relative to the original data. Recent work demonstrates that large language models paired with arithmetic coding yield state-of-the-art text compression results, achieving compression ratios of roughly 0.1 on Wikipedia corpora (Bellard, 2023; Valmeekam et al., 2023; Delétang et al., 2024).

2.2 Interactive Compression via Binary Questions

The children’s game *20 Questions* illustrates how interaction can drastically reduce communication. One player thinks of an object, and the other identifies it by asking a sequence of yes–no questions, each conveying a single bit of information. By choosing questions adaptively, the player can identify one object among millions using roughly 20 bits—far fewer than would be required to specify the object directly (and on the order of < 2 tokens in an LLM of information).

This game reflects a general principle: interaction enables information to be revealed conditionally, focusing communication on what is most informative given previous answers. In complexity theory, this principle manifests as a separation between non-interactive and interactive verification. Non-interactive proofs, consisting of a single message from a prover to a verifier, characterize the class NP. Allowing the verifier to ask adaptive questions yields interactive proofs, which characterize $IP = PSPACE$, a class believed to be strictly larger than NP (Goldwasser et al., 1989; Shamir, 1992).

In Section 3, we apply this idea to LLM communication. A small model adaptively asks binary questions that are answered by a large model, allowing substantial knowledge to be transferred using only a small number of bits.

2.3 LoRA and Domain Adaptation

Low-Rank Adaptation (LoRA) (Hu et al., 2021) is a parameter-efficient fine-tuning method that adapts a pretrained model by learning low-rank update matrices. For a pretrained weight matrix $W_0 \in \mathbb{R}^{d \times k}$, LoRA learns $\Delta W = BA$ where $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times k}$ with rank $r \ll \min(d, k)$. This dramatically reduces the number of trainable parameters.

The key insight for compression is that a LoRA adapter fine-tuned on domain-specific text will assign higher probability to in-domain text, thereby reducing the bits required for arithmetic coding. If sender and receiver share a set of domain-specific LoRA adapters, they can achieve better compression by selecting the appropriate adapter for each message.

3 Interactive Question Asking (Question-Asking compression)

We introduce a novel interactive compression mechanism that adapts the children’s game 20 Questions to the task of LLM knowledge transfer. Given a prompt, an SLM first generates an initial response and then formulates N binary yes/no questions about its solution strategy, for example, “Is my approach to step 3 correct?” The LLM answers each question with yes or no, transferring exactly one bit of information per response. The SLM then revises its response by incorporating all

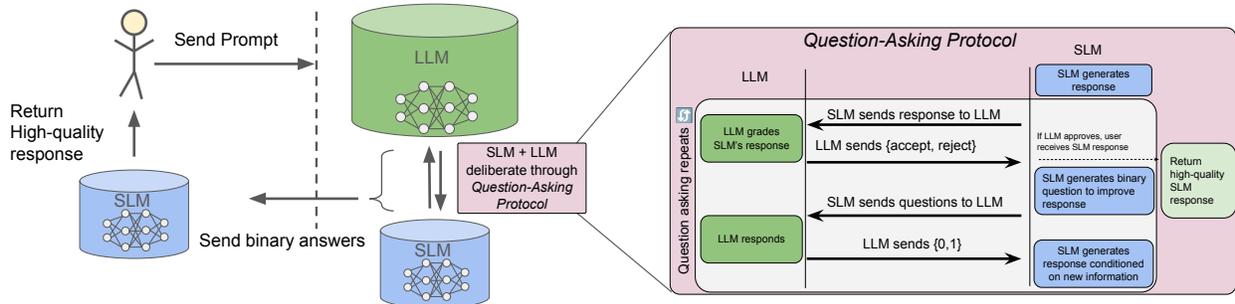


Figure 1: Overview of the compression mechanism and its use in an interactive protocol between an SLM and an LLM.

Algorithm 1 Question Asking (QA Protocol). a_i : SLM’s answer after round i ; q_i : yes/no question; $r_i \in \{\text{YES}, \text{NO}\}$: LLM’s binary response (1 bit). Line 6 (early stopping) is optional.

Require: Problem P , small model SLM, large model LLM, budget N

- 1: $a_0 \leftarrow \text{SLM}(P)$ {SLM attempts to solve P }
 - 2: **for** $i = 1, \dots, N$ **do**
 - 3: $q_i \leftarrow \text{SLM}(P, a_{i-1}, q_{<i}, r_{<i})$ {SLM asks a yes/no question}
 - 4: $r_i \leftarrow \text{LLM}(P, q_i)$ {LLM answers yes or no (1 bit)}
 - 5: $a_i \leftarrow \text{SLM}(P, a_{i-1}, q_{\leq i}, r_{\leq i})$ {SLM revises its answer}
 - 6: **if** $\text{Judge}(a_i) \geq \tau$ **then break** {Early stop if quality threshold met}
 - 7: **end for**
 - 8: **return** a_i
-

N answers. Crucially, the SLM’s questions are deterministic given its initial response, the prompt, and fixed hyperparameters. As a result, only the N binary answers must be transmitted.¹

The receiver can reconstruct the full question-answer transcript by running an identical SLM locally, reducing the communication cost to exactly N bits, independent of the response length. An important observation is that in practice the questions do not need to be sent to the LLM, as the LLM can perfectly simulate the SLM, removing the need for more bandwidth.

We evaluate on 8 datasets, which we group into *easier benchmarks*² (GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), GPQA (MC) (Rein et al., 2023)), where current models achieve high baseline accuracy, and *harder benchmarks* (MBPP (Austin et al., 2021), AIME (AI-MO, 2024), Humanity’s Last Exam (HLE) (Phan et al., 2025)), where even strong models have significant room for improvement. We classify problems by difficulty based on which models solve them at baseline: *Easy* (all correct), *Medium* (Haiku fails, Sonnet/Opus correct), *Hard* (only Opus correct), and *Very Hard* (all fail). Full benchmark descriptions and difficulty distributions are in Appendix B.

We also validate our approach with the open-source GPT-OSS-120B model (Appendix G.4), and we find qualitatively similar results. This demonstrates that the QA protocol transfers across model

¹In principle, the expected cost per question can be reduced below one bit. If the SLM can predict the base-rate distribution of LLM responses, for example that “Yes” occurs with probability 0.8 for a given question, arithmetic coding over the binary answers yields an expected cost of $H(\text{answer} | \text{question})$ bits per question. This improvement is essentially free: if the SLM’s calibration is incorrect, the LLM can transmit a single extra bit to signal this, ensuring that the protocol never performs worse than the naive one-bit-per-answer scheme.

²We define “easier benchmarks” as those where Haiku achieves above 66% baseline accuracy. The exception is GPQA, which we include in this group because its multiple-choice format makes answer verification straightforward, despite Haiku’s lower baseline accuracy.

Table 1: Recovery rate (% , mean \pm std over 3 trials) and compression ratio by difficulty subset. *Med. + Hard* = problems where Haiku fails but Opus succeeds; *All* = all non-easy problems (includes Very Hard). Bit-Limited Chain-of-Thought (BL-CoT) = haiku \rightarrow haiku \rightarrow haiku; QA = haiku \rightarrow opus \rightarrow haiku; QA⁺ = opus \rightarrow opus \rightarrow haiku. Bold = highest mean recovery in each row within each subset. Ave. Resp. Len. = mean Opus response length in tokens (c1100k.base tokenizer). Compression = ratio of 10 QA bits to full Opus response (tokens \times log₂(100,277) bits).

Dataset	Med. + Hard			All Non-Easy			Ave. Resp. Len.	Compression		
	<i>n</i>	BL-CoT	QA	QA ⁺	<i>n</i>	BL-CoT			QA	QA ⁺
<i>Easier benchmarks</i>										
GSM8K	110	19.1 \pm 2.0	23.3 \pm 1.1	27.6 \pm 1.5	145	22.5 \pm 1.2	26.7 \pm 0.9	29.7 \pm 1.7	195	0.0031
MATH (Algebra)	85	28.6 \pm 1.5	27.5 \pm 2.2	31.0 \pm 1.5	259	38.1 \pm 1.6	39.1 \pm 1.0	40.0 \pm 1.5	373	0.0016
MATH (Geometry)	51	28.1 \pm 1.8	37.3 \pm 3.2	37.3 \pm 1.6	190	33.7 \pm 0.4	39.6 \pm 0.5	38.2 \pm 1.5	687	0.0009
MATH (Num. Thy)	12	66.7 \pm 11.8	72.2 \pm 15.7	63.9 \pm 15.7	61	38.3 \pm 2.0	43.7 \pm 4.3	40.4 \pm 3.1	580	0.0010
GPQA (MC)	15	40.0 \pm 14.4	40.0 \pm 14.4	31.1 \pm 6.3	152	11.8 \pm 2.7	13.6 \pm 0.3	12.1 \pm 2.0	606	0.0010
<i>Harder benchmarks</i>										
AIME	34	10.8 \pm 1.4	6.9 \pm 1.4	9.8 \pm 1.4	51	18.3 \pm 1.8	19.0 \pm 2.4	22.2 \pm 1.8	1,083	0.0006
MBPP	44	36.4 \pm 1.9	37.9 \pm 2.8	51.5 \pm 2.8	106	18.2 \pm 1.8	18.2 \pm 0.4	26.1 \pm 0.4	164	0.0037
HLE [†]	206	6.5 \pm 0.6	7.9 \pm 0.6	13.1 \pm 1.4	306	5.4 \pm 0.6	6.5 \pm 0.0	9.3 \pm 0.8	711	0.0008

[†] For HLE, Very Hard problems were randomly subsampled to 100 (out of \sim 1,900) due to evaluation cost.

families, and importantly shows that this technique can already be used for lossy compression, as all users have access to opensource models.

Evaluation. We measure the *recovery rate*: the fraction of non-easy problems answered correctly after 10 rounds of binary Q&A. By construction, Haiku is initially wrong on Medium and Hard problems, so this equals the fraction of incorrect answers that become correct. We denote each configuration as *questioner* \rightarrow *answerer* \rightarrow *solver*. Our primary configuration is **QA** (haiku \rightarrow opus \rightarrow haiku): Haiku generates questions, Opus answers them, and Haiku produces a final answer. This measures whether Opus’s knowledge can be transferred via 10 binary answers.

However, by letting Haiku ask questions and iteratively revise its answer, we are giving it substantially more computation, and the iterative question-asking process is itself a form of structured chain-of-thought reasoning. Some of the improvement may therefore come from self-refinement rather than from Opus’s knowledge. To disentangle these effects, we introduce **BL-CoT** (*Bit-Limited Chain-of-Thought*, haiku \rightarrow haiku \rightarrow haiku): the identical protocol, but with Haiku answering its own questions. Any improvement from BL-CoT is self-refinement, whereas the difference between QA and BL-CoT isolates the value of Opus’s knowledge transfer. Table 21 in Appendix G provides the full per-difficulty breakdown.

Ablations: the value of question asking. In any real implementation, the SLM will serve as both question-asker and final solver. However, the SLM’s ability to identify its own reasoning gaps may be a bottleneck: if it cannot formulate targeted questions, even perfect answers provide limited value. To simulate future improvements in question-asking ability, we introduce replace the question-generator with another instance of Opus, while keeping Haiku as the solver, denoted QA⁺ (opus \rightarrow opus \rightarrow haiku). This gives us insight into how this protocol’s effectiveness may change as models become more capable. Table 1 compares BL-CoT, QA, and QA⁺ on two subsets: *Med. + Hard* problems (where Opus succeeds but Haiku does not, so correct knowledge exists to transfer)

and *All Non-Easy* problems (which additionally includes Very Hard problems where even Opus fails). The full 3×3 breakdown across all questioner–answerer combinations is in Appendix G.

The results above ask all 10 questions in a single batch, with the LLM given the ground-truth answer as reference when responding. An iterative variant lets a judge early-stop the exchange once the SLM’s answer is satisfactory, potentially saving bits. We test two judge designs (objective vs. comparison scoring), threshold levels, and whether the judge sees the ground-truth answer; see Appendix E.1.1. These variants do not improve over the batch protocol: the judge either early-stops prematurely on easy problems or fails to prevent regression on hard ones.

We also test scaling from 10 to 100 questions; this yields small, inconsistent gains that do not justify the $10\times$ increase in bits (Appendix G.2).

Results. The central question is whether the QA protocol’s value comes from the *structure* of iterative refinement or from the *content* of the stronger model’s answers. Table 1 reveals that the answer depends on task difficulty. On easier benchmarks (GSM8K, MATH, GPQA), the majority of the recovery is captured by process of asking questions (BL-CoT), and Haiku does approximately as well as Opus as the LLM or as the Question Asker; the structured protocol alone enables Haiku to self-correct on well-posed math and science problems.

On harder benchmarks, there is a larger gain in recovery rates when the LLM is large model. Here, we speculate that the SLM lacks the knowledge to self-correct, so the information in Opus’s binary answers becomes more important, and a stronger question-asker (Opus in QA⁺) amplifies the improvement, as it can better identify the SLM’s failure modes. AIME is an exception, with uniformly low recovery across all configurations, suggesting that competition math resists the binary Q&A format. Example transcripts illustrating the Q&A interaction are provided in Appendix G.8.

Compression Ratios. The rightmost column of Table 1 shows the compression achieved by the QA protocol: 10 bits (10 binary answers) compared to a full Opus response, measured as $\text{tokens} \times \log_2(100,277) \approx 16.61$ bits per token using the `cl100k_base` tokenizer (an upper bound treating each token as uniform over the vocabulary). The QA scheme achieves compression ratios between 0.0006 and 0.0037, transferring only 10 bits compared to thousands in a full response.³ To put these ratios in context, prior SOTA LLM-based lossless compression achieves approximately 0.08 on `enwik9` (Delétang et al., 2024), and our own domain-adapted LoRA approach (Section 4.2) reaches approximately 0.03. The QA scheme’s best result of 0.0006 on AIME represents over $100\times$ improvement over prior SOTA, illustrating the compressive power of interactive protocols that exploit shared priors between models rather than encoding the full response text.

Regression on Easy Problems. The QA-compression protocol should not harm problems the SLM already solves. We verify this by running the protocol on *easy* problems and find that regression roughly matches Haiku’s own baseline inconsistency when re-evaluated with no protocol at all. Full results are in Appendix G.7, Table 26.

Question Transcript Analysis Looking at the transcripts, the small model generates yes/no verification questions that walk through its own solution step by step. Interestingly, 31% of recovered problems had every answer confirmed as “yes” by the frontier model; the structured self-verification

³For multiple-choice tasks like GPQA, this framing is somewhat misleading: the correct answer contains only $\log_2(4) = 2$ bits of information, so using 10 binary questions to communicate it is technically an *expansion*. The QA scheme is most compelling for open-ended tasks (math, code) where the answer space is large, not for multiple-choice tasks where the answer is already low-entropy.

alone was enough to fix the error. For example, on a GSM8K tax problem, Haiku asks questions like “Does the solution correctly apply the 10% tax rate only to the nonfood subtotal?” and “Are all five items accounted for with their correct prices”; the frontier model confirms each one, and Haiku revises its answer anyway.

In another 37% of recoveries, the small model asks mostly confirming questions but surfaces one or two key misconceptions. On a GSM8K word problem, Haiku asks “Does ‘10 times more than the fourth friend’ mean adding 10 times the fourth friend’s presses?”, and the frontier model answers “no”, pinpointing the exact misinterpretation; the other four questions are all confirmed. This pattern, where the small model’s approach is mostly sound and just needs a targeted correction, accounts for the majority of successful recoveries.

When QA fails, the transcripts reveal two characteristic failure modes. In one, the small model verifies all the steps of a wrong solution path. On a MATH algebra problem, Haiku carefully checks each step of its factorization and sign analysis, the frontier model confirms them all, but the fundamental setup is wrong and the questions never probe it. In the other failure mode, the small model is too confused to ask useful questions at all. This is particularly stark on code benchmarks (MBPP), where the small model asks about edge cases, type hints, and error handling (“Should the function validate that the radius is non-negative?”, “Would it be beneficial to add type hints?”) rather than probing its actual bug, however the corrections do not help because they are orthogonal to the real issue.

Examples of transcripts are laid out in Appendix G.8, as well as uploaded to HuggingFace.

4 Novel Baselines

This section introduces two simple but effective compression methods that serve as competitive standalone baselines for evaluating our QA-compression scheme and other approaches.

Section 3 explores a lossy baseline based on best-of- N sampling and explicit prompting for succinctness. By selecting the shortest high-quality response, we obtain approximately a $2\times$ improvement in compression ratio.

Section 4.2 presents a lossless baseline that builds on arithmetic coding. Rather than compressing directly with a single language model, we pre-train multiple domain-adapted LoRA adapters and select the most relevant adapter at inference time. This specialization yields roughly a $2\times$ improvement in compression ratio relative to a non-adapted model.

4.1 Lossy Compression Baselines: Response Rewriting

We introduce a simple approach to lossy compression is to generate multiple candidate solutions and select the one that compresses best using a standard off-the-shelf arithmetic compression scheme. This should not meaningfully change the distribution of texts that we produce, as we are selecting only from texts generated by the LLM.

Shortest-of-N. Generate N candidate solutions and select the one with the lowest compression ratio. We evaluate two variants: **Temperature Sampling:** Generate N independent solutions by sampling with temperature $T > 0$ (we use $T = 0.8$). **Single-Prompt:** Request N distinct solutions in a single prompt, e.g., “Provide 5 different approaches to solve this problem.”

Summarization. Another relatively-straightforward alternative is to *directly request* a succinct rewrite. First generate a verbose solution with full reasoning, then prompt the model to rewrite it

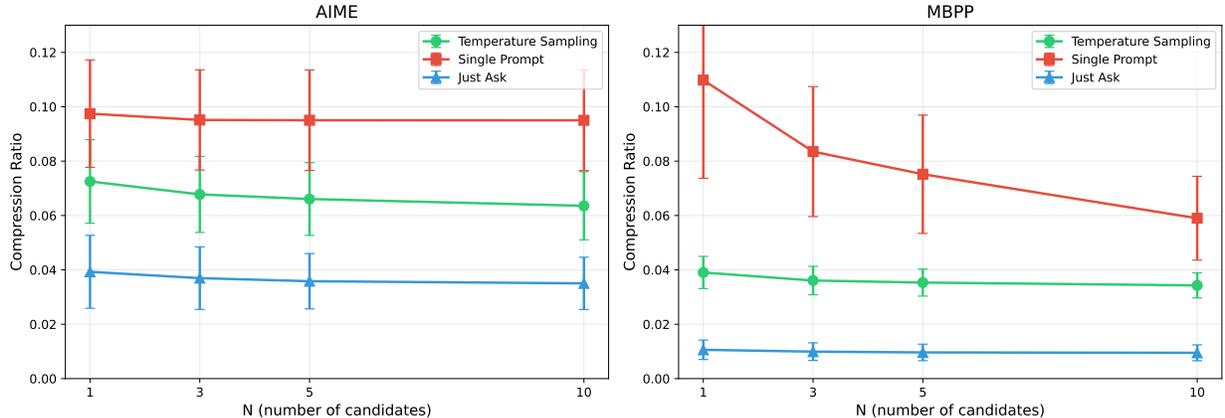


Figure 2: Compression ratio vs. number of candidates N on AIME (left) and MBPP (right) problems, using Opus as both generator and compressor. **Just Ask** (blue) achieves the best compression by explicitly requesting succinct rewrites, roughly halving the compression ratio compared to **Shortest-of- N** methods. Among Shortest-of- N variants, **Temperature Sampling** (green) outperforms **Single Prompt** (red), likely because single-prompt solutions share structural patterns. Lower is better.

“as succinctly as possible while preserving enough information to infer the correct answer.” This can be combined with temperature sampling to generate N succinct rewrites and select the shortest.

Results. Figure 2 compares these three approaches on AIME (90 competition math problems (AI-MO, 2024)) and Mostly Basic Python Programs (MBPP; 257 tasks (Austin et al., 2021)) using Opus as the generator and Mistral-7B-Instruct as the compressor Jiang et al. (2023), there are full benchmark descriptions are in Appendix B.1.

On AIME, temperature sampling achieves compression ratios from 0.073 at $N = 1$ to 0.063 at $N = 10$ ($1.2\times$ improvement from sampling more candidates). Summarization significantly outperforms both, achieving 0.039 at $N = 1$ ($\approx 1.9\times$ better than temperature sampling) and 0.034 at $N = 10$, roughly halving the compression ratio from a single sample. Results with Haiku as generator are qualitatively similar and shown in Appendix D.4.

The summarized rewrites are approximately $4\times$ shorter than verbose originals. Interestingly, the shorter text has *higher* bits-per-character: by stripping away predictable boilerplate (“Step 1:”, “Therefore,”), summarization produces denser text that is harder for the compression model to predict on a per-character basis. However, the reduction in length is larger than the increase in per-character complexity, so the net effect is still roughly $2\times$ better total compression. See Appendix D.5 for a more in-depth analysis of this trade-off.

We find that selecting for compressibility does not degrade answer quality: across all three generation strategies, choosing the most compressible candidate yields accuracy within a few percentage points of random selection. See Appendix D.2 for the full analysis.

4.2 Lossless Compression through Domain-Adapted LoRAs

While LLMs provide powerful probabilistic models for text, arithmetic coding achieves optimal compression only when the coding distribution matches the data distribution. When text drawn from P is encoded using model Q , the expected code length increases by the KL divergence $\text{KL}(P \parallel Q)$. A general-purpose LLM may therefore compress domain-specific data suboptimally.

LoRA (Hu et al., 2021) has proven particularly effective for domain adaptation, enabling a single base model to be specialized to multiple domains via lightweight adapters. If the sender and receiver share a collection of domain-specific LoRA adapters, they can improve compression by selecting the adapter that best matches the input distribution for each message. We investigate whether domain-specialized language models can achieve better compression rates than general-purpose models.

Experimental Setup. We use LMSYS-Chat-1M (Zheng et al., 2023), a large-scale dataset of real user conversations with LLMs, as our primary testbed. We cluster the dataset into 50 topic clusters using K-means on Qwen embeddings. For ten clusters, we fine-tune a LoRA adapter (Hu et al., 2021) on the training split (rank=16, $\alpha=32$) using Mistral-7B-Instruct as the base model (Jiang et al., 2023). We evaluate compression performance on held-out test samples with at least 200 tokens, measuring compression ratio using different methods: (1) **Baseline**: the base model without any LoRA, (2) **Gzip**, (3) **Correct LoRA**: the adapter trained on the same cluster, (4) **Wrong LoRAs**: an adapter trained for each of the nine remaining incorrect clusters, and (5) **RAG LoRA**: the adapter selected by a retrieval-based router (described below). All the clusters and the LoRAs are shared publicly on HuggingFace.

For all methods, we report compression ratio as compressed size divided by uncompressed size.⁴

4.2.1 RAG-based LoRA Router

In practice, one does not have access to the exact distribution that a text was generated from, so selecting the correct LoRA requires an automatic routing method. We propose an approach akin to retrieval-augmented generation (RAG) (Lewis et al., 2021).

Routing Protocol. We build a FAISS index over representative texts from each cluster using Qwen embeddings. Each index entry maps to the corresponding cluster’s LoRA adapter. For each input text, we embed the text, retrieve the top index, and use the corresponding LoRA adapter.

Routing Accuracy. We evaluate RAG routing accuracy on held-out test samples using both the full text (prompt and response) and only the prompt. The router correctly identifies the domain cluster approximately 72% of the time with full text and 55% with prompt only. For cluster-specific results, see Appendix C.

Results. Table 2 reports per-cluster compression ratios. On average, the correct LoRA achieves a compression ratio of 0.09, a $2.0\times$ improvement over the base model (0.18). Even the average wrong LoRA improves to 0.13 ($1.4\times$). The RAG-selected LoRA achieves 0.10 ($1.9\times$), approaching the oracle despite only 72% routing accuracy.

Analysis. The per-cluster results reveal that improvement varies by topic. Crucially, the correct LoRA consistently outperforms the wrong LoRAs across all clusters, validating that topic-matched adapters capture domain-specific patterns that transfer to improved compression. However, the wrong LoRAs still outperform the base model, suggesting the dataset comprises highly-correlated conversational text; a more diverse dataset would show wrong LoRAs having zero or negative effect.

⁴For GZIP, this is compressed bytes divided by original bytes. For arithmetic coding, this is output bits divided by input bits, where input bits equals $\log_2(\text{vocabulary size}) \times \text{number of tokens}$. The arithmetic coder’s bitstream can be packed into bytes with negligible overhead (< 1 byte per message), so these metrics are directly comparable.

Table 2: Compression performance (compression ratio) on LMSYS-Chat. Correct LoRA = adapter trained on the same cluster; Average = mean across the 9 incorrect adapters; Worst = worst-performing incorrect adapter; RAG LoRA = adapter selected by the retrieval-based router from prompt only.

Cluster	Topic	Baseline		Correct	LoRA		RAG LoRA Prompt
		Gzip	Base Model		Average	Worst	
0	General Chat	0.49	0.21	0.11	0.13	0.14	0.11
1	Creative Writing	0.42	0.18	0.11	0.12	0.13	0.11
2	Code/Technical	0.46	0.28	0.10	0.20	0.21	0.10
3	Academic/Education	0.44	0.16	0.09	0.11	0.11	0.10
4	Roleplay/Fiction	0.40	0.11	0.05	0.08	0.08	0.05
5	Business/Professional	0.44	0.20	0.09	0.13	0.14	0.09
6	Philosophy/Ethics	0.45	0.19	0.09	0.14	0.14	0.09
7	Science/Math	0.39	0.15	0.09	0.11	0.11	0.10
8	Translation/Language	0.43	0.18	0.08	0.13	0.14	0.10
9	Casual Q&A	0.44	0.15	0.10	0.12	0.13	0.10
Overall		0.44	0.18	0.09	0.13	0.13	0.10

We replicate these results on WildChat (Zhao et al., 2024), a dataset of 1M ChatGPT interaction logs (Table 11, Appendix C). The correct LoRA achieves a compression ratio of 0.10, a 1.7 \times improvement over baseline (0.17), slightly lower than LMSYS (2.0 \times), consistent with WildChat’s greater topic diversity within clusters.

5 Discussion

Using the framing of the QA-compression framework, we suggest a shift in how the presence of open-source models alters the structure of communication we have access to. When two parties share access to the same model, they effectively share a rich internal representation of the world. Communication no longer requires transmitting an entire idea from scratch; instead, it is sufficient to transmit only the difference between what the shared model would already generate and the intended content. In this sense, shared models enable an unusually low-bandwidth form of communication grounded in a common representation.

More broadly, question-asking serves not only as a compression primitive but also as a mechanism for probing conditional information. Quantifying the conditional information between two beings (human or LLM) is generally abstract and difficult to measure directly. By framing communication as a sequence of binary questions, QA-style interaction provides an operational way to bound how much additional information one agent must transmit to another about a specific concept. The number of questions required offers a concrete proxy for their conditional information gap.

5.1 Applications to Model Security

One promising defense for protecting model weights is *egress limiting*: restricting the total number of bits that exit a datacenter (Nevo et al., 2024; Anthropic, 2025; Greenblatt, 2024). This rests on a simple information-theoretic observation: if only 10 GB leaves a datacenter, an adversary cannot steal more than 10 GB of anything, regardless of their attack vector. However, strict bit

limits can harm legitimate users who need to receive model outputs. As we have explored in this paper, natural language is extremely compressible, which can be used to empower defenders to apply stricter egress limits, as model weights are generally more difficult to compress than text.⁵

5.2 Limitations

Selection bias in shortest-of-N. When we select the most compressible response from N candidates, we are implicitly biasing toward certain kinds of outputs. We measure accuracy on AIME to verify that compression-selected responses are not systematically worse (Appendix D), but other quality dimensions (clarity, completeness, pedagogical value) are not evaluated. It is plausible that optimizing for compressibility selects for responses that are terse or formulaic in ways that standard accuracy metrics do not capture.

Benchmark-specific evaluation. Our QA-compression results are evaluated on math, code, and science benchmarks with verifiable answers. These domains are well-suited to the binary question format because correctness is unambiguous. The protocol’s effectiveness on open-ended generation tasks (creative writing, summarization, advice) is untested and likely quite different, since these tasks lack clear ground truth against which to evaluate yes/no questions.

LLM-as-judge reliability. The iterative QA-compression protocol relies on the LLM to answer binary questions about solution correctness. When the LLM is itself wrong (as on very-hard problems), its answers can actively mislead the SLM. We observe this empirically: on frontier benchmarks where the judge fails to reliably evaluate solutions, regression on easy problems is severe (Section 3). The protocol’s utility is therefore bounded by the LLM’s own accuracy on the task.

5.3 Future Work

Improving question quality for QA-compression. The QA-compression protocol is inherently limited by the SLM’s ability to ask useful questions and integrate the answers. Our results show that on frontier benchmarks, upgrading the question-asker to a stronger model yields meaningful gains (Section 3), suggesting that question quality is a bottleneck. Future work could explore improving the SLM’s question-asking ability, whether through prompt engineering, fine-tuning on transcripts from successful QA-compression sessions, or reinforcement learning on recovery rate.

Alternative interactive schemes. Our QA-compression protocol relies entirely on the SLM’s ability to ask useful questions. However, this is limiting: consider “What are the first 5 decimal places of $\sqrt{17}$?” A small model lacks the capacity to derive the answer, so no sequence of binary questions will help. A natural extension is to allow the LLM to *inject* information, providing information that the SLM would not discover through questioning alone. Other low-bit interaction schemes are also possible. For example, the SLM could generate N candidate responses, and the LLM could select the best one in $\log_2 N$ bits.⁶

⁵Though recent research seeks to find more compressible forms of model weights (Brown and Mazeika, 2025).

⁶We explored this direction but found that this meaningfully lengthened responses, eliminating compression gains. However, we view this as an engineering challenge rather than a fundamental barrier, and we believe alternative strategies may yet make such schemes viable.

Improving Robustness to Non-determinism. In both lossless and lossy compression, achieving meaningful compression ratios requires that the encoder and decoder reproduce identical model outputs. However, models running on different hardware regularly exhibit non-deterministic outputs at both the logit and token level (Yuan et al., 2025; He and Lab, 2025; Rinberg et al., 2025; Karvonen et al., 2025). Compression schemes that are robust to small amounts of cross-machine non-determinism would be significantly more practical. Some work has explored this direction for arithmetic coding (Adler and Tang, 2026), but the space remains largely unexplored for lossy compression.

Robustness to Non-determinism. Both our lossless and interactive compression schemes assume that the encoder and decoder reproduce identical model outputs. In practice, modern LLM inference pipelines are not strictly deterministic: numerical discrepancies arise from floating-point effects, hardware variation, batching, and low-level kernel implementations (Yuan et al., 2025; He and Lab, 2025; Rinberg et al., 2025; Karvonen et al., 2025).

For arithmetic coding, exact agreement between encoder and decoder probability distributions is required. Because decoding is autoregressive, even small mismatches can compound and lead to catastrophic decoding failure. The QA-compression scheme is similarly sensitive: if the large model fails to perfectly simulate the small model and generates a different question, synchronization breaks and the protocol collapses. Recent work such as PMATIC (Adler and Tang, 2026) introduces auxiliary “helper bits” to realign encoder and decoder distributions. In our work, we proposed a complementary approach, *block emission encoding*, which periodically emits the midpoint of the current coding interval and restarts from a fresh $[0, 1]$ range after every N bits of information, thereby limiting error propagation to a single block; however, more improvements need to be made for this to be a robust solution.⁷ Despite these advances, designing compression schemes that are robust to cross-machine non-determinism remains an open problem.

⁷A more detailed discussion appears in Appendix A.1.

6 Acknowledgements

RR’s work was supported by a grant from Coefficient Giving, and this work was conducted in part during RR’s participation in MATS. SH’s work was supported by the Harvard Paulson SEAS Prize Fellowship and the German Academic Fellowship Organization, funded by the German Federal Ministry for Economic Affairs and Energy. Compute credits provided by the ML Alignment & Theory Scholars (MATS) program as well as the Kempner Institute for the Study of Natural and Artificial Intelligence at Harvard University.

References

- Aviv Adler and Jennifer Tang. Synchronizing probabilities in model-driven lossless compression, 2026. URL <https://arxiv.org/abs/2601.10678>.
- AI-MO. Aimo validation aime, 2024. URL <https://huggingface.co/datasets/AI-MO/aimo-validation-aime>. Validation dataset for the AIMO Progress Prize, derived from AIME 2022–2024 problems.
- Anthropic. Activating ai safety level 3 protections. <https://www.anthropic.com/news/activating-asl3-protections>, May 2025. Online, published May 22, 2025.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Fabrice Bellard. `ts_zip`: Text-based compression using time series modeling. https://bellard.org/ts_zip/, 2023. Accessed: 2026-01-19.
- Davis Brown and Mantas Mazeika. How does llm compression affect weight exfiltration attacks? In *SoLaR: Symposium on Language and Representation*, 2025. URL <https://openreview.net/forum?id=j9IMNLGXi8>.
- 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. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christopher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Marcus Hutter, and Joel Veness. Language modeling is compression, 2024. URL <https://arxiv.org/abs/2309.10668>.
- P. Deutsch. Gzip file format specification version 4.3. RFC 1952, Internet Engineering Task Force, May 1996. URL <https://www.rfc-editor.org/rfc/rfc1952.html>. Informational RFC.
- Radhika Dua, Young Joon, Kwon, Siddhant Dogra, Daniel Freedman, Diana Ruan, Motaz Nashawaty, Danielle Rigau, Daniel Alexander Alber, Kang Zhang, Kyunghyun Cho, and Eric Karl Oermann. Clinically grounded agent-based report evaluation: An interpretable metric for radiology report generation, 2025. URL <https://arxiv.org/abs/2508.02808>.

- Shafi Goldwasser, Silvio Micali, and Charles Rackoff. The knowledge complexity of interactive proof systems. In *Proceedings of the 17th Annual ACM Symposium on Theory of Computing (STOC)*, pages 291–304. ACM, 1989.
- Ryan Greenblatt. Preventing model exfiltration with upload limits. AI Alignment Forum, feb 2024. URL <https://www.alignmentforum.org/posts/rf66R4YsrCHgWx9RG/>. Accessed: 2025-10-24.
- Horace He and Thinking Machines Lab. Defeating nondeterminism in llm inference. *Thinking Machines Lab: Connectionism*, 2025. doi: 10.64434/tml.20250910. <https://thinkingmachines.ai/blog/defeating-nondeterminism-in-llm-inference/>.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- Adam Karvonen, Daniel Reuter, Roy Rinberg, Luke Marks, Adri  Garriga-Alonso, and Keri Warr. Difr: Inference verification despite nondeterminism, 2025. URL <https://arxiv.org/abs/2511.20621>.
- G. Langdon and J. Rissanen. Compression of black-white images with arithmetic coding. *IEEE Transactions on Communications*, 29(6):858–867, 1981. doi: 10.1109/TCOM.1981.1095052.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich K ttler, Mike Lewis, Wen tau Yih, Tim Rockt schel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021. URL <https://arxiv.org/abs/2005.11401>.
- Sella Nevo, Dan Lahav, Ajay Karapur, Yogev Bar-On, Henry Alexander Bradley, and Jeff Alstott. Securing ai model weights: Preventing theft and misuse of frontier models. Research Report RR A2849-1, RAND Corporation, Santa Monica, CA, 2024.
- OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita Brett, Eugene Brevdo, Greg Brockman, Sebastien Bubeck, Che Chang, Kai Chen, Mark Chen, Enoch Cheung, Aidan Clark, Dan Cook, Marat Dukhan, Casey Dvorak, Kevin Fives, Vlad Fomenko, Timur Garipov, Kristian Georgiev, Mia Glaese, Tarun Gogineni, Adam Goucher, Lukas Gross, Katia Gil Guzman, John Hallman, Jackie Hehir, Johannes Heidecke, Alec Helyar, Haitang Hu, Romain Huet, Jacob Huh, Saachi Jain, Zach Johnson, Chris Koch, Irina Kofman, Dominik Kundel, Jason Kwon, Volodymyr Kyrylov, Elaine Ya Le, Guillaume Leclerc, James Park Lennon, Scott Lessans, Mario Lezcano-Casado, Yuanzhi Li, Zhuohan Li, Ji Lin, Jordan Liss, Lily, Liu, Jiancheng Liu, Kevin Lu, Chris Lu, Zoran Martinovic, Lindsay

McCallum, Josh McGrath, Scott McKinney, Aidan McLaughlin, Song Mei, Steve Mostovoy, Tong Mu, Gideon Myles, Alexander Neitz, Alex Nichol, Jakub Pachocki, Alex Paino, Dana Palmie, Ashley Pantuliano, Giambattista Parascandolo, Jongsoo Park, Leher Pathak, Carolina Paz, Ludovic Peran, Dmitry Pimenov, Michelle Pokrass, Elizabeth Proehl, Huida Qiu, Gaby Raila, Filippo Raso, Hongyu Ren, Kimmy Richardson, David Robinson, Bob Rotsted, Hadi Salman, Suvansh Sanjeev, Max Schwarzer, D. Sculley, Harshit Sikchi, Kendal Simon, Karan Singhal, Yang Song, Dane Stuckey, Zhiqing Sun, Philippe Tillet, Sam Toizer, Foivos Tsimpourlas, Nikhil Vyas, Eric Wallace, Xin Wang, Miles Wang, Olivia Watkins, Kevin Weil, Amy Wendling, Kevin Whinnery, Cedric Whitney, Hannah Wong, Lin Yang, Yu Yang, Michihiro Yasunaga, Kristen Ying, Wojciech Zaremba, Wenting Zhan, Cyril Zhang, Brian Zhang, Eddie Zhang, and Shengjia Zhao. gpt-oss-120b & gpt-oss-20b model card, 2025. URL <https://arxiv.org/abs/2508.10925>.

Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra, Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Dmitry Dodonov, Tung Nguyen, Jaeho Lee, Daron Anderson, Mikhail Doroshenko, Alun Cennyth Stokes, Mobeen Mahmood, Oleksandr Pokutnyi, Oleg Iskra, Jessica P. Wang, John-Clark Levin, Mstyslav Kazakov, Fiona Feng, Steven Y. Feng, Haoran Zhao, Michael Yu, Varun Gangal, Chelsea Zou, Zihan Wang, Serguei Popov, Robert Gerbicz, Geoff Galgon, Johannes Schmitt, Will Yeadon, Yongki Lee, Scott Sauers, Alvaro Sanchez, Fabian Giska, Marc Roth, Søren Riis, Saiteja Utpala, Noah Burns, Gashaw M. Goshu, Mohinder Maheshbhai Naiya, Chidozie Agu, Zachary Giboney, Antrell Cheatom, Francesco Fournier-Facio, Sarah-Jane Crowson, Lennart Finke, Zerui Cheng, Jennifer Zampese, Ryan G. Hoerr, Mark Nandor, Hyunwoo Park, Tim Gehringer, Jiaqi Cai, Ben McCarty, Alexis C Garretson, Edwin Taylor, Damien Sileo, Qiuyu Ren, Usman Qazi, Lianghui Li, Jungbae Nam, John B. Wydallis, Pavel Arkhipov, Jack Wei Lun Shi, Aras Bacho, Chris G. Willcocks, Hangrui Cao, Sumeet Motwani, Emily de Oliveira Santos, Johannes Veith, Edward Vendrow, Doru Cojoc, Kengo Zenitani, Joshua Robinson, Longke Tang, Yuqi Li, Joshua Vendrow, Natanael Wildner Fraga, Vladyslav Kuchkin, Andrey Pupasov Maksimov, Pierre Marion, Denis Efremov, Jayson Lynch, Kaiqu Liang, Aleksandar Mikov, Andrew Gritsevskiy, Julien Guillod, Gözdenur Demir, Dakotah Martinez, Ben Pageler, Kevin Zhou, Saeed Soori, Ori Press, Henry Tang, Paolo Rissone, Sean R. Green, Lina Brüssel, Moon Twayana, Aymeric Dieuleveut, Joseph Marvin Imperial, Ameya Prabhu, Jinzhou Yang, Nick Crispino, Arun Rao, Dimitri Zvonkine, Gabriel Loiseau, Mikhail Kalinin, Marco Lukas, Ciprian Manolescu, Nate Stambaugh, Subrata Mishra, Tad Hogg, Carlo Bosio, Brian P Coppola, Julian Salazar, Jaehyeok Jin, Rafael Sayous, Stefan Ivanov, Philippe Schwaller, Shaipranesh Senthilkuma, Andres M Bran, Andres Algaba, Kelsey Van den Houte, Lynn Van Der Sypt, Brecht Verbeken, David Noever, Alexei Kopylov, Benjamin Myklebust, Bikun Li, Lisa Schut, Evgenii Zheltonozhskii, Qiaochu Yuan, Derek Lim, Richard Stanley, Tong Yang, John Maar, Julian Wykowski, Martí Oller, Anmol Sahu, Cesare Giulio Ardito, Yuzheng Hu, Ariel Ghislain Kemogne Kamdoun, Alvin Jin, Tobias Garcia Vilchis, Yuexuan Zu, Martin Lackner, James Koppel, Gongbo Sun, Daniil S. Antonenko, Steffi Chern, Bingchen Zhao, Pierrot Arsene, Joseph M Cavanagh, Daofeng Li, Jiawei Shen, Donato Crisostomi, Wenjin Zhang, Ali Dehghan, Sergey Ivanov, David Perrella, Nurdin Kaparov, Allen Zang, Ilia Sucholutsky, Arina Kharlamova, Daniil Orel, Vladislav Poritski, Shalev Ben-David, Zachary Berger, Parker Whitfill, Michael Foster, Daniel Munro, Linh Ho, Shankar Sivarajan, Dan Bar Hava, Aleksey Kuchkin, David Holmes, Alexandra Rodriguez-Romero, Frank Sommerhage, Anji Zhang, Richard Moat, Keith Schneider, Zakayo Kazibwe, Don Clarke, Dae Hyun Kim, Felipe Meneguitti Dias, Sara Fish, Veit Elser, Tobias Kreiman, Victor Eflen Guadarrama Vilchis, Immo Klose, Ujjwala Anantheswaran, Adam Zweiger, Kaivalya Rawal, Jeffery Li, Jeremy Nguyen, Nicolas Daans,

Haline Heidinger, Maksim Radionov, Václav Rozhoň, Vincent Ginis, Christian Stump, Niv Cohen, Rafał Poświata, Josef Tkadlec, Alan Goldfarb, Chenguang Wang, Piotr Padlewski, Stanislaw Barzowski, Kyle Montgomery, Ryan Stendall, Jamie Tucker-Foltz, Jack Stade, T. Ryan Rogers, Tom Goertzen, Declan Grabb, Abhishek Shukla, Alan Givré, John Arnold Ambay, Archan Sen, Muhammad Fayez Aziz, Mark H Inlow, Hao He, Ling Zhang, Younesse Kaddar, Ivar Ångquist, Yanxu Chen, Harrison K Wang, Kalyan Ramakrishnan, Elliott Thornley, Antonio Terpin, Hailey Schoelkopf, Eric Zheng, Avishy Carmi, Ethan D. L. Brown, Kelin Zhu, Max Bartolo, Richard Wheeler, Martin Stehberger, Peter Bradshaw, JP Heimonen, Kaustubh Sridhar, Ido Akov, Jennifer Sandlin, Yury Makarychev, Joanna Tam, Hieu Hoang, David M. Cunningham, Vladimir Goryachev, Demosthenes Patramanis, Michael Krause, Andrew Redenti, David Aldous, Jesyin Lai, Shannon Coleman, Jiangnan Xu, Sangwon Lee, Ilias Magoulas, Sandy Zhao, Ning Tang, Michael K. Cohen, Orr Paradise, Jan Hendrik Kirchner, Maksym Ovchynnikov, Jason O. Matos, Adithya Shenoy, Michael Wang, Yuzhou Nie, Anna Szyber-Betley, Paolo Faraboschi, Robin Riblet, Jonathan Crozier, Shiv Halasyamani, Shreyas Verma, Prashant Joshi, Eli Meril, Ziqiao Ma, Jérémy Andréoletti, Raghav Singhal, Jacob Platnick, Volodymyr Nevirkovets, Luke Basler, Alexander Ivanov, Seri Khoury, Nils Gustafsson, Marco Piccardo, Hamid Mostaghimi, Qijia Chen, Virendra Singh, Tran Quoc Khánh, Paul Rosu, Hannah Szlyk, Zachary Brown, Himanshu Narayan, Aline Menezes, Jonathan Roberts, William Alley, Kunyang Sun, Arkil Patel, Max Lamparth, Anka Reuel, Linwei Xin, Hanmeng Xu, Jacob Loader, Freddie Martin, Zixuan Wang, Andrea Achilleos, Thomas Preu, Tomek Korbak, Ida Bosio, Fereshteh Kazemi, Ziye Chen, Biró Bálint, Eve J. Y. Lo, Jiaqi Wang, Maria Inês S. Nunes, Jeremiah Milbauer, M Saiful Bari, Zihao Wang, Behzad Ansarinejad, Yewen Sun, Stephane Durand, Hossam Elgnainy, Guillaume Douville, Daniel Tordera, George Balabanian, Hew Wolff, Lynna Kvistad, Hsiaoyun Milliron, Ahmad Sakor, Murat Eron, Andrew Favre D. O., Shailesh Shah, Xiaoxiang Zhou, Firuz Kamalov, Sherwin Abdoli, Tim Santens, Shaul Barkan, Allison Tee, Robin Zhang, Alessandro Tomasiello, G. Bruno De Luca, Shi-Zhuo Looi, Vinh-Kha Le, Noam Kolt, Jiayi Pan, Emma Rodman, Jacob Drori, Carl J Fossum, Niklas Muennighoff, Milind Jagota, Ronak Pradeep, Honglu Fan, Jonathan Eicher, Michael Chen, Kushal Thaman, William Merrill, Moritz Firsching, Carter Harris, Stefan Ciobâcă, Jason Gross, Rohan Pandey, Ilya Gusev, Adam Jones, Shashank Agnihotri, Pavel Zhelnov, Mohammadreza Mofayezi, Alexander Piperski, David K. Zhang, Kostiantyn Dobarskyi, Roman Leventov, Ignat Soroko, Joshua Duersch, Vage Taamazyan, Andrew Ho, Wenjie Ma, William Held, Ruicheng Xian, Armel Randy Zebaze, Mohanad Mohamed, Julian Noah Leser, Michelle X Yuan, Laila Yacar, Johannes Lengler, Katarzyna Olszewska, Claudio Di Fratta, Edson Oliveira, Joseph W. Jackson, Andy Zou, Muthu Chidambaram, Timothy Manik, Hector Haffenden, Dashiell Stander, Ali Dasouqi, Alexander Shen, Bitu Golshani, David Stap, Egor Kretov, Mikalai Uzhou, Alina Borisovna Zhidkovskaya, Nick Winter, Miguel Orbegoza Rodriguez, Robert Lauff, Dustin Wehr, Colin Tang, Zaki Hossain, Shaun Phillips, Fortuna Samuele, Fredrik Ekström, Angela Hammon, Oam Patel, Faraz Farhidi, George Medley, Forough Mohammadzadeh, Madellene Peñaflo, Haile Kassahun, Alena Friedrich, Rayner Hernandez Perez, Daniel Pyda, Taom Sakal, Omkar Dhamane, Ali Khajegili Mirabadi, Eric Hallman, Kenchi Okutsu, Mike Battaglia, Mohammad Maghsoudimehrabani, Alon Amit, Dave Hulbert, Roberto Pereira, Simon Weber, Handoko, Anton Peristy, Stephen Malina, Mustafa Mehkary, Rami Aly, Frank Reidegeld, Anna-Katharina Dick, Cary Friday, Mukhwinder Singh, Hassan Shapourian, Wanyoung Kim, Mariana Costa, Hubeyb Gurdogan, Harsh Kumar, Chiara Ceconello, Chao Zhuang, Haon Park, Micah Carroll, Andrew R. Tawfeek, Stefan Steinerberger, Daattavya Aggarwal, Michael Kirchof, Linjie Dai, Evan Kim, Johan Ferret, Jainam Shah, Yuzhou Wang, Minghao Yan, Krzysztof Burdzy, Lixin Zhang, Antonio Franca, Diana T. Pham, Kang Yong Loh, Joshua Robinson, Abram Jackson, Paolo Giordano, Philipp Petersen, Adrian Cosma, Jesus Colino, Colin White, Jacob

Votava, Vladimir Vinnikov, Ethan Delaney, Petr Spelda, Vit Stritecky, Syed M. Shahid, Jean-Christophe Mourrat, Lavr Vetoshkin, Koen Sponselee, Renas Bacho, Zheng-Xin Yong, Florencia de la Rosa, Nathan Cho, Xiuyu Li, Guillaume Malod, Orion Weller, Guglielmo Albani, Leon Lang, Julien Laurendeau, Dmitry Kazakov, Fatimah Adesanya, Julien Portier, Lawrence Hollom, Victor Souza, Yuchen Anna Zhou, Julien Degorre, Yiğit Yalın, Gbenga Daniel Obikoya, Rai, Filippo Bigi, M. C. Boscá, Oleg Shumar, Kaniuar Bacho, Gabriel Recchia, Mara Popescu, Nikita Shulga, Ngefor Mildred Tanwie, Thomas C. H. Lux, Ben Rank, Colin Ni, Matthew Brooks, Alesia Yakimchyk, Huanxu, Liu, Stefano Cavalleri, Olle Häggström, Emil Verkama, Joshua Newbould, Hans Gundlach, Leonor Brito-Santana, Brian Amaro, Vivek Vajipey, Rynaa Grover, Ting Wang, Yosi Kratish, Wen-Ding Li, Sivakanth Gopi, Andrea Caciolai, Christian Schroeder de Witt, Pablo Hernández-Cámara, Emanuele Rodolà, Jules Robins, Dominic Williamson, Vincent Cheng, Brad Raynor, Hao Qi, Ben Segev, Jingxuan Fan, Sarah Martinson, Erik Y. Wang, Kaylie Hausknecht, Michael P. Brenner, Mao Mao, Christoph Demian, Peyman Kassani, Xinyu Zhang, David Avagian, Eshawn Jessica Scipio, Alon Ragoler, Justin Tan, Blake Sims, Rebeka Plecnik, Aaron Kirtland, Omer Faruk Bodur, D. P. Shinde, Yan Carlos Leyva Labrador, Zahra Adoul, Mohamed Zekry, Ali Karakoc, Tania C. B. Santos, Samir Shamseldeen, Loukmane Karim, Anna Liakhovitskaia, Nate Resman, Nicholas Farina, Juan Carlos Gonzalez, Gabe Maayan, Earth Anderson, Rodrigo De Oliveira Pena, Elizabeth Kelley, Hodjat Mariji, Rasoul Pouriamanesh, Wentao Wu, Ross Finocchio, Ismail Alarab, Joshua Cole, Danyelle Ferreira, Bryan Johnson, Mohammad Safdari, Liangti Dai, Siriphan Arthornthurasuk, Isaac C. McAlister, Alejandro José Moyano, Alexey Pronin, Jing Fan, Angel Ramirez-Trinidad, Yana Malysheva, Daphiny Pottmaier, Omid Taheri, Stanley Stepanic, Samuel Perry, Luke Askew, Raúl Adrián Huerta Rodríguez, Ali M. R. Minissi, Ricardo Lorena, Krishnamurthy Iyer, Arshad Anil Fasiludeen, Ronald Clark, Josh Ducey, Matheus Piza, Maja Somrak, Eric Vergo, Juehang Qin, Benjámín Borbás, Eric Chu, Jack Lindsey, Antoine Jallon, I. M. J. McInnis, Evan Chen, Avi Semler, Luk Gloor, Tej Shah, Marc Carauleanu, Pascal Lauer, Tran Duc Huy, Hossein Shahrtash, Emilien Duc, Lukas Lewark, Assaf Brown, Samuel Albanie, Brian Weber, Warren S. Vaz, Pierre Clavier, Yiyang Fan, Gabriel Poesia Reis e Silva, Long, Lian, Marcus Abramovitch, Xi Jiang, Sandra Mendoza, Murat Islam, Juan Gonzalez, Vasiliios Mavroudis, Justin Xu, Pawan Kumar, Laxman Prasad Goswami, Daniel Bugas, Nasser Heydari, Ferenc Jeanplong, Thorben Jansen, Antonella Pinto, Archimedes Apronti, Abdallah Galal, Ng Ze-An, Ankit Singh, Tong Jiang, Joan of Arc Xavier, Kanu Priya Agarwal, Mohammed Berkani, Gang Zhang, Zhehang Du, Benedito Alves de Oliveira Junior, Dmitry Malishev, Nicolas Remy, Taylor D. Hartman, Tim Tarver, Stephen Mensah, Gautier Abou Loume, Wiktor Morak, Farzad Habibi, Sarah Hoback, Will Cai, Javier Gimenez, Roselynn Grace Montecillo, Jakub Lucki, Russell Campbell, Asankhaya Sharma, Khalida Meer, Shreen Gul, Daniel Espinosa Gonzalez, Xavier Alapont, Alex Hoover, Gunjan Chhablani, Freddie Vargus, Arunim Agarwal, Yibo Jiang, Deepakkumar Patil, David Outevsky, Kevin Joseph Scaria, Rajat Maheshwari, Abdelkader Dendane, Priti Shukla, Ashley Cartwright, Sergei Bogdanov, Niels Mündler, Sören Möller, Luca Arnaboldi, Kunvar Thaman, Muhammad Rehan Siddiqi, Prajvi Saxena, Himanshu Gupta, Tony Fruhauff, Glen Sherman, Mátyás Vincze, Siranut Usawasutsakorn, Dylan Ler, Anil Radhakrishnan, Innocent Enyekwe, Sk Md Salauddin, Jiang Muzhen, Aleksandr Maksapetyan, Vivien Rossbach, Chris Harjadi, Mohsen Bahalooohoreh, Claire Sparrow, Jasdeep Sidhu, Sam Ali, Song Bian, John Lai, Eric Singer, Justine Leon Uro, Greg Bateman, Mohamed Sayed, Ahmed Menshawy, Darling Duclosel, Dario Bezzi, Yashaswini Jain, Ashley Aaron, Murat Tiryakioglu, Sheeshram Siddh, Keith Krenek, Imad Ali Shah, Jun Jin, Scott Creighton, Denis Peskoff, Zienab EL-Wasif, Ragavendran P V, Michael Richmond, Joseph McGowan, Tejal Patwardhan, Hao-Yu Sun, Ting Sun, Nikola Zubić, Samuele Sala, Stephen Ebert, Jean Kaddour, Manuel Schottdorf, Dianzhuo Wang, Gerol Petruzella, Alex Meiburg, Tilen Medved, Ali ElSheikh, S Ashwin Hebbar, Lorenzo

Vaquero, Xianjun Yang, Jason Poulos, Vilém Zouhar, Sergey Bogdanik, Mingfang Zhang, Jorge Sanz-Ros, David Anugraha, Yinwei Dai, Anh N. Nhu, Xue Wang, Ali Anil Demircali, Zhibai Jia, Yuyin Zhou, Juncheng Wu, Mike He, Nitin Chandok, Aarush Sinha, Gaoxiang Luo, Long Le, Mickaël Noyé, Michał Perelkiewicz, Ioannis Pantidis, Tianbo Qi, Soham Sachin Purohit, Letitia Parcalabescu, Thai-Hoa Nguyen, Genta Indra Winata, Edoardo M. Ponti, Hanchen Li, Kaustubh Dhole, Jongee Park, Dario Abbondanza, Yuanli Wang, Anupam Nayak, Diogo M. Caetano, Antonio A. W. L. Wong, Maria del Rio-Chanona, Dániel Kondor, Pieter Francois, Ed Chaltrey, Jakob Zsambok, Dan Hoyer, Jenny Reddish, Jakob Hauser, Francisco-Javier Rodrigo-Ginés, Suchandra Datta, Maxwell Shepherd, Thom Kamphuis, Qizheng Zhang, Hyunjun Kim, Ruiji Sun, Jianzhu Yao, Franck Dernoncourt, Satyapriya Krishna, Sina Rismanchian, Bonan Pu, Francesco Pinto, Yingheng Wang, Kumar Shridhar, Kalon J. Overholt, Glib Briia, Hieu Nguyen, David, Soler Bartomeu, Tony CY Pang, Adam Wecker, Yifan Xiong, Fanfei Li, Lukas S. Huber, Joshua Jaeger, Romano De Maddalena, Xing Han Lù, Yuhui Zhang, Claas Beger, Patrick Tser Jern Kon, Sean Li, Vivek Sanker, Ming Yin, Yihao Liang, Xinlu Zhang, Ankit Agrawal, Li S. Yifei, Zechen Zhang, Mu Cai, Yasin Sonmez, Costin Cozianu, Changhao Li, Alex Slen, Shoubin Yu, Hyun Kyu Park, Gabriele Sarti, Marcin Briański, Alessandro Stolfo, Truong An Nguyen, Mike Zhang, Yotam Pernitz, Jose Hernandez-Orallo, Runjia Li, Amin Shabani, Felix Juefei-Xu, Shikhar Dhingra, Orr Zohar, My Chiffon Nguyen, Alexander Pondaven, Abdurrahim Yilmaz, Xuandong Zhao, Chuanyang Jin, Muyan Jiang, Stefan Todoran, Xinyao Han, Jules Kreuer, Brian Rabern, Anna Plassart, Martino Maggetti, Luther Yap, Robert Geirhos, Jonathon Kean, Dingsu Wang, Sina Mollaei, Chenkai Sun, Yifan Yin, Shiqi Wang, Rui Li, Yaowen Chang, Anjiang Wei, Alice Bizeul, Xiaohan Wang, Alexandre Oliveira Arrais, Kushin Mukherjee, Jorge Chamorro-Padial, Jiachen Liu, Xingyu Qu, Junyi Guan, Adam Bouyamourn, Shuyu Wu, Martyna Plomecka, Junda Chen, Mengze Tang, Jiaqi Deng, Shreyas Subramanian, Haocheng Xi, Haoxuan Chen, Weizhi Zhang, Yinuo Ren, Haoqin Tu, Sejong Kim, Yushun Chen, Sara Vera Marjanović, Junwoo Ha, Grzegorz Luczyna, Jeff J. Ma, Zewen Shen, Dawn Song, Cedegao E. Zhang, Zhun Wang, Gaël Gendron, Yunze Xiao, Leo Smucker, Erica Weng, Kwok Hao Lee, Zhe Ye, Stefano Ermon, Ignacio D. Lopez-Miguel, Theo Knights, Anthony Gitter, Namkyu Park, Boyi Wei, Hongzheng Chen, Kunal Pai, Ahmed Elkhanany, Han Lin, Philipp D. Siedler, Jichao Fang, Ritwik Mishra, Károly Zsolnai-Fehér, Xilin Jiang, Shadab Khan, Jun Yuan, Rishab Kumar Jain, Xi Lin, Mike Peterson, Zhe Wang, Aditya Malusare, Maosen Tang, Isha Gupta, Ivan Fosin, Timothy Kang, Barbara Dworakowska, Kazuki Matsumoto, Guangyao Zheng, Gerben Sewuster, Jorge Pretel Villanueva, Ivan Rannev, Igor Chernyavsky, Jiale Chen, Deepayan Banik, Ben Racz, Wenchao Dong, Jianxin Wang, Laila Bashmal, Duarte V. Gonçalves, Wei Hu, Kaushik Bar, Ondrej Bohdal, Atharv Singh Patlan, Shehzaad Dhuliawala, Caroline Geirhos, Julien Wist, Yuval Kansal, Bingsen Chen, Kutay Tire, Atak Talay Yücel, Brandon Christof, Veerupaksh Singla, Zijian Song, Sanxing Chen, Jiaxin Ge, Kaustubh Ponshe, Isaac Park, Tianneng Shi, Martin Q. Ma, Joshua Mak, Sherwin Lai, Antoine Moulin, Zhuo Cheng, Zhanda Zhu, Ziyi Zhang, Vaidehi Patil, Ketan Jha, Qiutong Men, Jiaxuan Wu, Tianchi Zhang, Bruno Hebling Vieira, Alham Fikri Aji, Jae-Won Chung, Mohammed Mahfoud, Ha Thi Hoang, Marc Sperzel, Wei Hao, Kristof Meding, Sihan Xu, Vassilis Kostakos, Davide Manini, Yueying Liu, Christopher Toukmaji, Jay Paek, Eunmi Yu, Arif Engin Demircali, Zhiyi Sun, Ivan Dewerpe, Hongsen Qin, Roman Pflugfelder, James Bailey, Johnathan Morris, Ville Heilala, Sybille Rosset, Zishun Yu, Peter E. Chen, Woongyeong Yeo, Eeshaan Jain, Ryan Yang, Sreekar Chigurupati, Julia Chernyavsky, Sai Prajwal Reddy, Subhashini Venugopalan, Hunar Batra, Core Francisco Park, Hieu Tran, Guilherme Maximiano, Genghan Zhang, Yizhuo Liang, Hu Shiyu, Rongwu Xu, Rui Pan, Siddharth Suresh, Ziqi Liu, Samaksh Gulati, Songyang Zhang, Peter Turchin, Christopher W. Bartlett, Christopher R. Scotese, Phuong M. Cao, Ben Wu, Jacek Karwowski, Davide Scaramuzza, Aakaash Nattanmai, Gordon McKellips, Anish Cheraku, Asim

- Suhail, Ethan Luo, Marvin Deng, Jason Luo, Ashley Zhang, Kavin Jindel, Jay Paek, Kasper Halevy, Allen Baranov, Michael Liu, Advait Avadhanam, David Zhang, Vincent Cheng, Brad Ma, Evan Fu, Liam Do, Joshua Lass, Hubert Yang, Surya Sunkari, Vishruth Bharath, Violet Ai, James Leung, Rishit Agrawal, Alan Zhou, Kevin Chen, Tejas Kalpathi, Ziqi Xu, Gavin Wang, Tyler Xiao, Erik Maung, Sam Lee, Ryan Yang, Roy Yue, Ben Zhao, Julia Yoon, Sunny Sun, Aryan Singh, Ethan Luo, Clark Peng, Tyler Osbey, Taozhi Wang, Daryl Echeazu, Hubert Yang, Timothy Wu, Spandan Patel, Vidhi Kulkarni, Vijaykaarti Sundarapandiyam, Ashley Zhang, Andrew Le, Zafir Nasim, Srikar Yalam, Ritesh Kasamsetty, Soham Samal, Hubert Yang, David Sun, Nihar Shah, Abhijeet Saha, Alex Zhang, Leon Nguyen, Laasya Nagumalli, Kaixin Wang, Alan Zhou, Aidan Wu, Jason Luo, Anwith Telluri, Summer Yue, Alexandr Wang, and Dan Hendrycks. Humanity’s last exam, 2025. URL <https://arxiv.org/abs/2501.14249>.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023. URL <https://arxiv.org/abs/2311.12022>.
- Roy Rinberg, Adam Karvonen, Alexander Hoover, Daniel Reuter, and Keri Warr. Verifying llm inference to detect model weight exfiltration, 2025. URL <https://arxiv.org/abs/2511.02620>.
- J. Rissanen and G. G. Langdon. Arithmetic coding. *IBM Journal of Research and Development*, 23(2):149–162, 1979. doi: 10.1147/rd.232.0149.
- Adi Shamir. Ip = pspace. *J. ACM*, 39(4), October 1992. ISSN 0004-5411. doi: 10.1145/146585.146609. URL <https://doi.org/10.1145/146585.146609>.
- Chandra Shekhara Kaushik Valmeekam, Krishna Narayanan, Dileep Kalathil, Jean-Francois Chamberland, and Srinivas Shakkottai. Llmzip: Lossless text compression using large language models, 2023. URL <https://arxiv.org/abs/2306.04050>.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. Asking and answering questions to evaluate the factual consistency of summaries, 2020. URL <https://arxiv.org/abs/2004.04228>.
- Ian H. Witten, Radford M. Neal, and John G. Cleary. Arithmetic coding for data compression. *Communications of the ACM*, 30(6):520–540, 1987. doi: 10.1145/214762.214771.
- Jiayi Yuan, Hao Li, Xinheng Ding, Wenya Xie, Yu-Jhe Li, Wentian Zhao, Kun Wan, Jing Shi, Xia Hu, and Zirui Liu. Understanding and mitigating numerical sources of nondeterminism in llm inference, 2025. URL <https://arxiv.org/abs/2506.09501>.
- Hanlin Zhang, Benjamin L. Edelman, Danilo Francati, Daniele Venturi, Giuseppe Ateniese, and Boaz Barak. Watermarks in the sand: Impossibility of strong watermarking for generative models, 2025. URL <https://arxiv.org/abs/2311.04378>.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat: 1m chatgpt interaction logs in the wild, 2024. URL <https://arxiv.org/abs/2405.01470>.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric. P Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. Lmsys-chat-1m: A large-scale real-world llm conversation dataset, 2023.

A Extended Background

This appendix provides extended discussion of background concepts summarized in Section 2.

A.1 Arithmetic Coding: Extended Discussion

Arithmetic coding is an entropy coding technique that achieves near-optimal lossless compression by encoding a sequence of symbols using a number of bits that approaches the Shannon entropy of the source (Witten et al., 1987). Given a probabilistic model over a discrete alphabet, arithmetic coding converts probability estimates directly into code lengths, assigning shorter codes to more likely sequences and longer codes to less likely ones.

In the context of language modeling, let p denote a language model and let x_1, x_2, \dots, x_n be a sequence of tokens. When arithmetic coding is applied using the model’s conditional probabilities, the resulting code length is approximately

$$\sum_{i=1}^n -\log_2 p(x_i \mid x_1, \dots, x_{i-1}),$$

which corresponds to the model’s cross-entropy evaluated on the sequence. Consequently, compression performance is governed almost entirely by the quality of the language model’s predictive distribution.

Building on this observation, several recent works have demonstrated that large language models can serve as effective entropy models for text compression when paired with arithmetic coding, achieving strong empirical results in both online and offline settings (Bellard, 2023; Valmeekam et al., 2023; Delétang et al., 2024).

A.2 Distribution Mismatch: Extended Discussion

Arithmetic coding achieves optimal compression only when the coding distribution matches the true data-generating distribution. When text drawn from a distribution P is encoded using a model Q , the expected code length increases by the Kullback-Leibler divergence $\text{KL}(P \parallel Q)$, yielding an expected cost of $H(P) + \text{KL}(P \parallel Q)$. As a result, a general-purpose LLM may perform suboptimally when compressing domain-specific data whose statistics deviate from those seen during pretraining.

This motivates domain adaptation approaches. We address this challenge through lightweight domain adaptation using LoRA fine-tuning (Section 4.2), which significantly improves compression performance by reducing distributional mismatch.

A.3 Non-Determinism in LLM Inference: Extended Discussion

Arithmetic coding requires exact agreement between the probability distributions used by the encoder and decoder. Modern LLM inference pipelines, however, are not strictly deterministic: small numerical discrepancies arise across runs due to floating-point effects, hardware variation, kernel implementations, or batching (He and Lab, 2025; Rinberg et al., 2025; Karvonen et al., 2025). Even vanishingly small probability mismatches can cause decoder desynchronization and catastrophic decoding failure.

This fragility is amplified by the sequential nature of arithmetic coding. Because decoding is autoregressive, any mismatch alters the decoder’s context and thus the predicted distribution for all subsequent tokens, causing errors to accumulate over long sequences. Many prior works do not observe this issue because they run compression and decompression on a single machine, where the cross-run non-determinism is very small (Delétang et al., 2024; Valmeekam et al., 2023). However,

real-world deployments necessarily involve separate encoder and decoder machines, where such discrepancies become unavoidable.

Recent work has addressed this challenge explicitly. In particular, PMATIC proposes communicating auxiliary “helper bits” to synchronize the effective model distributions between encoder and decoder (Adler and Tang, 2026). Since our focus is not on synchronization mechanisms themselves, we do not propose new solutions in this direction and instead refer practitioners to PMATIC-style approaches when deploying arithmetic coding with LLMs.

B Benchmarks, Evaluation Setup, and Dataset Statistics

B.1 Benchmarks

We evaluate our lossy compression methods on eight benchmarks spanning mathematics, code generation, science, and expert-level reasoning. These were chosen to cover a range of task types and difficulty levels, all with verifiable answers that enable automatic evaluation.

- **GSM8K** (Cobbe et al., 2021): 1,319 grade-school math word problems requiring multi-step arithmetic reasoning. Answers are integers.
- **MATH** (Hendrycks et al., 2021): Competition-level mathematics. We evaluate three subjects: Algebra (1,187 problems), Geometry (479), and Number Theory (540). Answers are typically expressions in \LaTeX (e.g., $5\sqrt{5}$).
- **GPQA (MC)** (Rein et al., 2023): 198 graduate-level science questions (physics, chemistry, biology) in multiple-choice format with four options.
- **MBPP** (Austin et al., 2021): 257 entry-level Python programming tasks from the sanitized split. Evaluated by executing generated code against unit tests.
- **AIME** (AI-MO, 2024): 90 problems from the American Invitational Mathematics Examination. Answers are integers in $\{0, \dots, 999\}$.
- **HLE** (Phan et al., 2025): 2,158 text-only problems from Humanity’s Last Exam, spanning diverse expert-level domains. Includes both exact-match and multiple-choice questions.

B.2 Difficulty Classification

We classify each problem into one of four difficulty tiers based on which models solve it correctly in a single baseline evaluation (temperature 0):

- **Easy**: All three models (Haiku, Sonnet, Opus) answer correctly.
- **Medium**: Haiku fails, but Sonnet and/or Opus succeed.
- **Hard**: Haiku and Sonnet fail, but Opus succeeds.
- **Very Hard**: All three models fail.

This classification is based on a single deterministic run per model; as we discuss in Appendix G.7, this can be unreliable for harder benchmarks (AIME, HLE) where model outputs are not consistently reproducible.

B.3 Claude 4.5 Model Family Results

Table 3 summarizes model accuracies across all evaluated datasets using the Claude 4.5 model family.

Table 3: Model accuracy across all datasets (Claude 4.5 family).

Dataset	n	Haiku	Sonnet	Opus
GSM8K	1,319	90.3%	97.7%	97.9%
MATH (Algebra)	1,187	82.0%	87.1%	86.5%
MATH (Geometry)	479	68.3%	73.3%	72.0%
MATH (Num. Theory)	540	92.4%	92.0%	92.8%
GPQA (MC)	198	29.3%	28.3%	33.8%
MBPP	257	60.3%	72.4%	75.9%
AIME	90	57.8%	52.2%	94.4%
HLE	2,158	4.7%	6.6%	12.2%

B.4 Claude 3.5/4 Model Family Results

Table 4 shows results using the older Claude 3.5/4 model family for comparison.

Table 4: Model accuracy across all datasets (Claude 3.5/4 family).

Dataset	n	Haiku	Sonnet	Opus
GSM8K	1,319	93.6%	97.0%	97.0%
MATH (Algebra)	1,187	76.7%	87.0%	83.1%
MATH (Geometry)	479	47.6%	72.4%	68.9%
MATH (Num. Theory)	540	73.9%	92.0%	90.9%
GPQA (MC)	198	33.3%	31.8%	31.3%
MBPP	257	65.4%	69.3%	70.0%
AIME	90	3.3%	36.7%	46.7%
HLE	2,158	3.9%	5.4%	6.3%

B.5 Difficulty Distribution (Claude 4.5)

Table 5 shows the difficulty distribution for each dataset based on Claude 4.5 model performance.

Table 5: Difficulty distribution across datasets (Claude 4.5 family).

Dataset	n	Easy	Medium	Hard	V.Hard
GSM8K	1,319	89.0%	8.0%	0.4%	2.7%
MATH (Algebra)	1,187	78.2%	5.1%	2.1%	14.7%
MATH (Geometry)	479	60.3%	6.3%	4.4%	29.0%
MATH (Num. Theory)	540	88.7%	1.7%	0.6%	9.1%
GPQA (MC)	198	23.2%	3.0%	4.5%	69.2%
MBPP	257	58.8%	10.5%	6.6%	24.1%
AIME	90	43.3%	6.7%	31.1%	18.9%
HLE	2,158	1.4%	2.3%	7.2%	89.1%

B.6 Model Identifiers

Table 6 lists the Anthropic API model identifiers used in our experiments. We report results for both model families: Claude 3.5/4 (Table 4) and Claude 4.5 (Table 3).

Table 6: Anthropic API model identifiers.

Model Name	Claude 3.5/4 (used in results)	Claude 4.5
Haiku	claude-3-5-haiku-20241022	claude-haiku-4-5-20251001
Sonnet	claude-sonnet-4-20250514	claude-sonnet-4-5-20250929
Opus	claude-opus-4-20250514	claude-opus-4-5-20251101

C LoRA Compression Results

C.1 Full Compression Results (All 10 LMSYS Clusters)

Table 7 provides the full per-cluster breakdown of compression performance across all 10 LMSYS clusters.

Table 7: Compression performance (compression ratio) for all 10 LMSYS clusters. Correct LoRA = adapter trained on the same cluster; Average = mean across the 9 incorrect adapters; RAG LoRA = adapter selected by the RAG router using full text or prompt only.

Cluster	Topic	Baseline	Gzip	Correct LoRA	Worst Wrong	Avg Wrong
0	General Chat	0.21	0.49 (0.44x)	0.11 (2.02x)	0.14 (1.48x)	0.13 (1.65x)
1	Creative Writing	0.18	0.42 (0.42x)	0.11 (1.67x)	0.13 (1.42x)	0.12 (1.47x)
2	Code/Technical	0.28	0.47 (0.61x)	0.10 (2.96x)	0.21 (1.32x)	0.20 (1.43x)
3	Academic/Education	0.16	0.44 (0.36x)	0.09 (1.71x)	0.11 (1.43x)	0.11 (1.49x)
4	Roleplay/Fiction	0.11	0.40 (0.26x)	0.05 (2.12x)	0.08 (1.32x)	0.08 (1.38x)
5	Business/Professional	0.20	0.44 (0.45x)	0.09 (2.20x)	0.14 (1.47x)	0.13 (1.56x)
6	Philosophy/Ethics	0.19	0.45 (0.43x)	0.09 (2.10x)	0.14 (1.33x)	0.14 (1.38x)
7	Science/Math	0.15	0.39 (0.38x)	0.09 (1.60x)	0.11 (1.34x)	0.11 (1.40x)
8	Translation/Language	0.18	0.43 (0.41x)	0.09 (2.08x)	0.14 (1.31x)	0.13 (1.35x)
9	Casual Q&A	0.15	0.44 (0.35x)	0.10 (1.57x)	0.13 (1.22x)	0.12 (1.26x)
Average		0.18	0.44 (0.41x)	0.09 (1.99x)	0.13 (1.36x)	0.13 (1.44x)

C.2 RAG Routing Accuracy (LMSYS)

Table 8: RAG Routing Accuracy by Cluster (LMSYS)

Cluster	Topic	Full	Prompt	Total
0	General Chat	60.0%	32.0%	50
1	Creative Writing	18.0%	16.0%	50
2	Code/Technical	78.0%	56.0%	50
3	Academic/Education	64.0%	36.0%	50
4	Roleplay/Fiction	98.0%	88.0%	50
5	Business/Professional	90.0%	68.0%	50
6	Philosophy/Ethics	88.0%	84.0%	50
7	Science/Math	58.0%	30.0%	50
8	Translation/Language	78.0%	48.0%	50
9	Casual Q&A	88.0%	92.0%	50
Overall		72.0%	55.0%	500

Table 9: Full vs Prompt RAG Routing Compression, Bits per Token (BPT) and Compression Ratio (LMSYS)

Cluster	Topic	Full BPT	Prompt BPT	Full Comp. Ratio	Prompt Comp. Ratio
0	General Chat	1.601	1.653	0.1069	0.1104
1	Creative Writing	1.636	1.715	0.1093	0.1146
2	Code/Technical	1.322	1.528	0.0884	0.1021
3	Academic/Education	1.417	1.504	0.0947	0.1005
4	Roleplay/Fiction	0.747	0.762	0.0499	0.0509
5	Business/Professional	1.340	1.345	0.0895	0.0899
6	Philosophy/Ethics	1.408	1.417	0.0940	0.0947
7	Science/Math	1.444	1.490	0.0965	0.0996
8	Translation/Language	1.326	1.550	0.0886	0.1036
9	Casual Q&A	1.464	1.496	0.0978	0.1000
Overall		1.370	1.446	0.0916	0.0966

C.3 WildChat Results

We replicate the LoRA compression pipeline on WildChat (Zhao et al., 2024), a dataset of 1M ChatGPT interaction logs in the wild. The same clustering, training, and RAG routing procedure is applied to 10 WildChat clusters.

RAG Routing Accuracy. Table 10 reports per-cluster RAG routing accuracy on WildChat. Overall accuracy is 77.0% (full text) and 53.0% (prompt only), which is comparable to LMSYS.

Compression Performance. Table 11 reports compression ratios on WildChat. The correct LoRA achieves a 1.7 \times improvement over baseline on average (0.10 vs. 0.17), slightly lower than

Table 10: RAG Routing Accuracy by Cluster (WildChat)

Cluster	Topic	Full	Prompt	Total
0	General Chat	56.0%	46.0%	50
1	Creative Writing	98.0%	0.0%	50
2	Code/Technical	94.0%	96.0%	50
3	Academic/Education	56.0%	32.0%	50
4	Roleplay/Fiction	92.0%	72.0%	50
5	Business/Professional	52.0%	18.0%	50
6	Philosophy/Ethics	38.0%	22.0%	50
7	Science/Math	98.0%	80.0%	50
8	Translation/Language	88.0%	70.0%	50
9	Casual Q&A	98.0%	94.0%	50
Overall		77.0%	53.0%	500

LMSYS (2.0×), consistent with WildChat’s greater topic diversity within clusters.

Table 11: Compression performance (compression ratio) for all 10 Wildchat clusters. Correct LoRA = adapter trained on the same cluster; Average = mean across the 9 incorrect adapters; RAG LoRA = adapter selected by the RAG router using full text or prompt only.

Cluster	Topic	Baseline	Gzip	Correct LoRA	Worst Wrong	Avg Wrong
0	General Chat	0.20	0.49 (0.40x)	0.13 (1.50x)	0.17 (1.15x)	0.17 (1.20x)
1	Creative Writing	0.19	0.38 (0.51x)	0.04 (4.92x)	0.16 (1.23x)	0.13 (1.51x)
2	Code/Technical	0.14	0.36 (0.39x)	0.10 (1.43x)	0.13 (1.06x)	0.12 (1.13x)
3	Academic/Education	0.14	0.51 (0.28x)	0.10 (1.38x)	0.13 (1.06x)	0.12 (1.14x)
4	Roleplay/Fiction	0.12	0.43 (0.28x)	0.09 (1.42x)	0.11 (1.09x)	0.10 (1.16x)
5	Business/Professional	0.16	0.52 (0.30x)	0.12 (1.30x)	0.15 (1.05x)	0.14 (1.13x)
6	Philosophy/Ethics	0.17	0.45 (0.37x)	0.13 (1.32x)	0.14 (1.16x)	0.14 (1.19x)
7	Science/Math	0.24	0.46 (0.52x)	0.09 (2.75x)	0.20 (1.18x)	0.19 (1.24x)
8	Translation/Language	0.17	0.46 (0.36x)	0.12 (1.44x)	0.15 (1.11x)	0.14 (1.17x)
9	Casual Q&A	0.16	0.45 (0.36x)	0.11 (1.55x)	0.15 (1.08x)	0.14 (1.14x)
Average		0.17	0.45 (0.38x)	0.10 (1.67x)	0.15 (1.12x)	0.14 (1.20x)

C.4 LoRA Compression on enwik9

To evaluate LoRA-based compression beyond conversational data, we replicate the lmsys pipeline (Section 4.2) on `enwik9`, a 1 GB prefix of an English Wikipedia XML dump commonly used as a compression benchmark. Unlike the lmsys dataset, which consists of multi-turn chat conversations, `enwik9` contains raw encyclopedic text spanning diverse topics.

C.4.1 Data Preparation

We extract individual Wikipedia articles from the XML dump, clean the text (removing markup, templates, and metadata), and filter to articles with at least 100 words, yielding 136,572 articles. We embed each article using the Qwen 0.6B embedding model and cluster with K -means into 10 clusters, then split each cluster 90/10 into train and test sets.

Table 12: Full vs Prompt RAG Routing Compression, Bits per Token (BPT) and Compression Ratio

Cluster	Topic	Full BPT	Prompt BPT	Full Comp. Ratio	Prompt Comp. Ratio
0	General Chat	2.007	2.035	0.1341	0.1360
1	Creative Writing	0.590	0.590	0.0394	0.0394
2	Code/Technical	1.454	1.471	0.0972	0.0983
3	Academic/Education	1.575	1.589	0.1052	0.1062
4	Roleplay/Fiction	1.273	1.298	0.0850	0.0867
5	Business/Professional	1.836	1.884	0.1227	0.1259
6	Philosophy/Ethics	1.930	1.968	0.1290	0.1315
7	Science/Math	1.304	1.318	0.0871	0.0881
8	Translation/Language	1.721	1.743	0.1150	0.1164
9	Casual Q&A	1.573	1.580	0.1051	0.1055
Overall		1.526	1.547	0.1020	0.1034

C.4.2 Cluster Summary

Table 13 summarizes the 10 clusters. The largest cluster (002) consists almost entirely of US Census-style geographic entries with highly formulaic text, which we expect to compress particularly well with a specialized LoRA.

Table 13: enwik9 cluster summary (10 clusters, 136,572 articles total).

Cluster	Topic	Train	Test
0	General knowledge / social concepts	8,464	941
1	Entertainment / arts / pop culture	12,111	1,346
2	US geography / census data	27,366	3,041
3	History and mythology	8,108	901
4	Geography and nature (international)	8,718	969
5	Politics and government	7,030	782
6	Science and nature	11,227	1,248
7	Infrastructure and miscellany	9,026	1,003
8	Biographies and institutions	10,830	1,204
9	Technology and computing	7,734	860
Total		122,909	13,663

Cluster 2 is notably larger than the others ($\sim 2.5\times$ the mean), reflecting Wikipedia’s large number of auto-generated US locality stubs that share near-identical structure (population, area, demographics).

C.4.3 Experimental Setup

We train one LoRA adapter per cluster on Mistral-7B-Instruct-v0.2, using the same hyperparameters as the lmsys pipeline (rank 16, $\alpha = 32$, 3 epochs). For LoRA selection at test time, we build a FAISS index over BGE-large-en-v1.5 embeddings of 1,000 training samples per cluster and route each test text via k -nearest-neighbor majority vote ($k = 10$).

We evaluate three compression conditions using arithmetic coding (64-bit precision) on held-out test texts (minimum 200 tokens):

1. **Baseline:** Mistral-7B with no LoRA adapter.

2. **Correct LoRA** (oracle): Mistral-7B with the LoRA trained on the same cluster as the test text.
3. **RAG-selected LoRA**: Mistral-7B with the LoRA chosen by the RAG router.

C.4.4 RAG Routing Accuracy

Table 14 reports RAG routing accuracy per cluster, evaluated on 50 test texts per cluster. Overall accuracy is 80%, with substantial variation across clusters.

Table 14: RAG routing accuracy on enwik9 (50 test texts per cluster, $k = 10$ nearest neighbors).

Cluster	Topic	Correct / Total	Accuracy
0	General knowledge / social concepts	40 / 50	80.0%
1	Entertainment / arts / pop culture	45 / 50	90.0%
2	US geography / census data	50 / 50	100.0%
3	History and mythology	45 / 50	90.0%
4	Geography and nature (international)	39 / 50	78.0%
5	Politics and government	37 / 50	74.0%
6	Science and nature	43 / 50	86.0%
7	Infrastructure and miscellany	31 / 50	62.0%
8	Biographies and institutions	32 / 50	64.0%
9	Technology and computing	38 / 50	76.0%
Overall		400 / 500	80.0%

Cluster 2 (US geography/census) achieves perfect routing accuracy, consistent with its highly distinctive, formulaic text. Clusters 7 (infrastructure) and 8 (biographies) are the hardest to route, likely because their content overlaps with multiple other clusters.

C.4.5 Compression Results (50 Clusters)

We extend the enwik9 experiment to 50 clusters for finer-grained topic specialization. We split enwik9 into chunks of approximately 2,000 tokens each and cluster using k-means on `all-MiniLM-L6-v2` embeddings. Each cluster’s LoRA is trained on Mistral-7B-Instruct-v0.2 with rank $r = 16$, $\alpha = 32$, for 3 epochs.

Table 15 shows per-cluster compression performance, sorted by baseline compression ratio. The compression ratio is computed as:

$$\text{compression ratio} = \frac{\text{compressed bits}}{\text{tokens} \times \log_2(\text{vocab_size})}$$

where $\log_2(32000) \approx 14.97$ bits per token represents uniform encoding over Mistral’s vocabulary.

Table 15: Enwik9 compression performance (50 clusters, sorted by baseline). Lower compression ratio is better.

Cluster	Samples	Baseline	RAG LoRA	Correct LoRA	RAG Acc
27	524	0.060	0.044	0.044	96.0%
36	1,327	0.064	0.046	0.045	70.4%
0	818	0.075	0.053	0.053	71.4%
11	469	0.077	0.055	0.055	90.4%
2	346	0.103	0.078	0.078	51.7%
44	121	0.138	0.110	0.110	87.6%
13	299	0.147	0.116	0.117	78.3%
9	330	0.164	0.141	0.140	51.2%
19	246	0.166	0.111	0.111	87.4%
41	233	0.189	0.160	0.160	95.7%
16	307	0.193	0.159	0.159	81.4%
10	257	0.194	0.162	0.162	80.2%
34	188	0.194	0.161	0.160	67.0%
38	122	0.195	0.169	0.169	78.7%
39	269	0.197	0.174	0.174	65.8%
22	118	0.197	0.162	0.163	92.4%
4	158	0.199	0.159	0.160	63.9%
6	176	0.200	0.165	0.165	72.7%
32	150	0.205	0.167	0.167	88.7%
26	250	0.208	0.182	0.182	74.4%
7	228	0.208	0.178	0.178	90.4%
15	355	0.210	0.183	0.183	73.2%
24	343	0.210	0.183	0.182	64.7%
1	373	0.212	0.184	0.184	60.1%
31	111	0.214	0.183	0.183	81.1%
46	179	0.214	0.182	0.183	73.7%
23	340	0.215	0.190	0.190	45.0%
49	100	0.218	0.181	0.181	90.0%
48	214	0.219	0.191	0.191	89.7%
42	143	0.219	0.194	0.194	83.2%
25	257	0.220	0.188	0.188	77.0%
3	169	0.220	0.190	0.190	86.4%
14	246	0.221	0.189	0.189	80.5%
28	219	0.222	0.198	0.198	79.9%
40	198	0.225	0.194	0.194	79.3%
33	303	0.225	0.196	0.196	55.8%
30	166	0.227	0.196	0.196	90.4%
29	205	0.227	0.173	0.173	89.3%
21	268	0.228	0.190	0.189	40.7%
45	220	0.228	0.203	0.203	64.5%
35	142	0.228	0.202	0.202	69.0%
43	229	0.230	0.199	0.199	76.4%
8	107	0.235	0.209	0.209	90.7%
47	192	0.240	0.209	0.208	83.9%
17	245	0.244	0.209	0.209	77.1%
18	182	0.247	0.215	0.215	68.7%
20	312	0.248	0.220	0.219	63.1%
12	285	0.248	0.219	0.219	66.7%
5	158	0.248	0.224	0.224	68.4%
37	245	0.253	0.224	0.224	64.1%
Total	13,442	0.176	0.148	0.148	73.6%

Key findings.

- **16% compression improvement:** LoRA adaptation reduces the average compression ratio from 0.176 (baseline) to 0.148 (RAG LoRA), a 16% reduction in compressed size.
- **High variance across clusters:** Compression ratios range from 0.044 (cluster 27) to 0.224 (clusters 5, 37), reflecting Wikipedia’s content diversity. The best-compressing clusters likely

contain highly structured content (tables, lists, formulaic entries).

- **RAG LoRA \approx Correct LoRA:** Despite 26.4% routing errors, RAG LoRA achieves nearly identical compression to the oracle (0.148 vs 0.148). This indicates that misrouted samples are typically sent to clusters with similar compression characteristics.
- **Routing accuracy varies:** RAG accuracy ranges from 40.7% (cluster 21) to 96.0% (cluster 27), with an overall accuracy of 73.6%.

Notes on baseline compression. Our baseline compression ratio of 0.176 is higher (worse) than ratios reported in prior work on neural text compression Delétang et al. (2024), which report values closer to 0.1 for the same model size. This discrepancy is likely due to several factors: (1) differences in base models (we use Mistral-7B-Instruct rather than models specifically optimized for compression), and (2) differences in experimental setup (e.g., how text is chunked, minimum token thresholds, evaluation protocol). The key takeaway is not the absolute compression ratio, but rather the *relative improvement* from LoRA adaptation: a consistent 16% reduction in compressed size across all clusters.

Comparison to conversational data. Interestingly, the conversational datasets (LMSYS and WildChat) achieve substantially better compression ratios than enwik9. LMSYS achieves a baseline ratio of 0.09 and WildChat achieves 0.10 (Tables 9 and 12), compared to 0.176 for enwik9. This suggests that LLM-generated conversational text is $\sim 2\times$ more compressible than raw Wikipedia text. This is intuitive: conversational text produced by LLMs likely has more predictable patterns and less information density than encyclopedic content, which contains diverse factual information, proper nouns, dates, and technical terminology.

D “Just Ask” Compression

D.1 Quality Preservation in Lossy Compression

When using Summarization (Section 4.1), a natural concern is whether the succinct rewrite preserves the quality of the original reasoning, or merely preserves enough surface-level information to guess the answer.

Answer-Masking Check. We implement a simple quality check. Before requesting the succinct rewrite, we mask the final answer in the verbose solution, replacing \boxed{A} with $\boxed{???$. The model then rewrites the masked solution, and must re-derive the correct answer from its own succinct version. If it produces an incorrect answer, we consider the compression too lossy.

Limitations. This check is imperfect. The succinct version might “leak” information about the answer through numerical patterns or structural cues (e.g., “the answer is a two-digit prime”) without actually preserving the full reasoning chain. A human reader might find the succinct version harder to follow or less informative than the verbose original, even if a language model can still extract the correct answer.

More rigorous assessment would require human evaluation of reasoning quality, or measuring performance on downstream tasks that depend on understanding the solution (not just knowing the final answer). We leave such evaluation to future work. Other work, such as Wang et al. (2020); Dua et al. (2025) propose alternative mechanisms for evaluating that a summary retains the core

information from the body, largely premised on using an LLM to generate questions about the original text, and then prompting another LLM to answer the questions based on the original and answer the questions based on the summarized text, and comparing alignment.

D.2 Effect of Compression-Based Selection on Accuracy

A natural concern with best-of- N selection by compressibility is whether choosing the most compressible candidate degrades answer quality. To investigate this, we compare the accuracy of two selection strategies across 90 AIME problems: (1) **random selection**, which picks uniformly among the N candidates (equivalently, the expected accuracy of a single draw), and (2) **best-compression selection**, which picks the candidate with the lowest compression ratio.

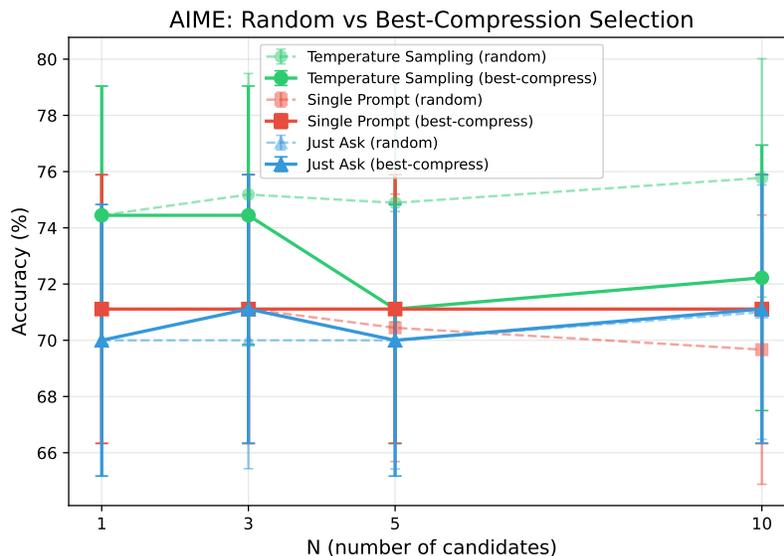


Figure 3: Accuracy of random selection (dashed) versus best-compression selection (solid) on 90 AIME problems using Opus. Random selection reports the expected accuracy of picking uniformly among N candidates. Across all three generation strategies, selecting the most compressible candidate yields accuracy within a few percentage points of random selection.

Figure 3 shows the results for all three generation approaches using Opus. For all three strategies, compression-based selection yields accuracy within a few percentage points of random selection across all values of N . For *temperature sampling*, random selection benefits slightly from oracle best-of- N at large N (rising from $\sim 70\%$ to $\sim 76\%$ at $N=10$), while compression-based selection remains relatively flat around 72–74%. For *single prompt* and *just ask*, both selection strategies yield nearly identical accuracy (~ 70 – 71%) across all N .

These results indicate that compression-based selection does not sacrifice accuracy. For the weaker Haiku model, compression-based selection can even *improve* accuracy (see §D.4).

D.3 Relative Compression Ratios

Figures 4 and 5 show both the absolute compression ratio (top row) and the relative compression ratio normalized to Temperature Sampling at $N=1$ (bottom row) for Opus and Haiku respectively. A relative compression ratio of 1.0 (dashed gray line) corresponds to no improvement over the baseline; lower values indicate better compression. On AIME, Just Ask achieves a relative compression ratio

of ~ 0.5 for both models, meaning it compresses to roughly half the size of the temperature sampling baseline. On MBPP, the gains are even more pronounced: Just Ask reaches ~ 0.27 for Opus and ~ 0.45 for Haiku.

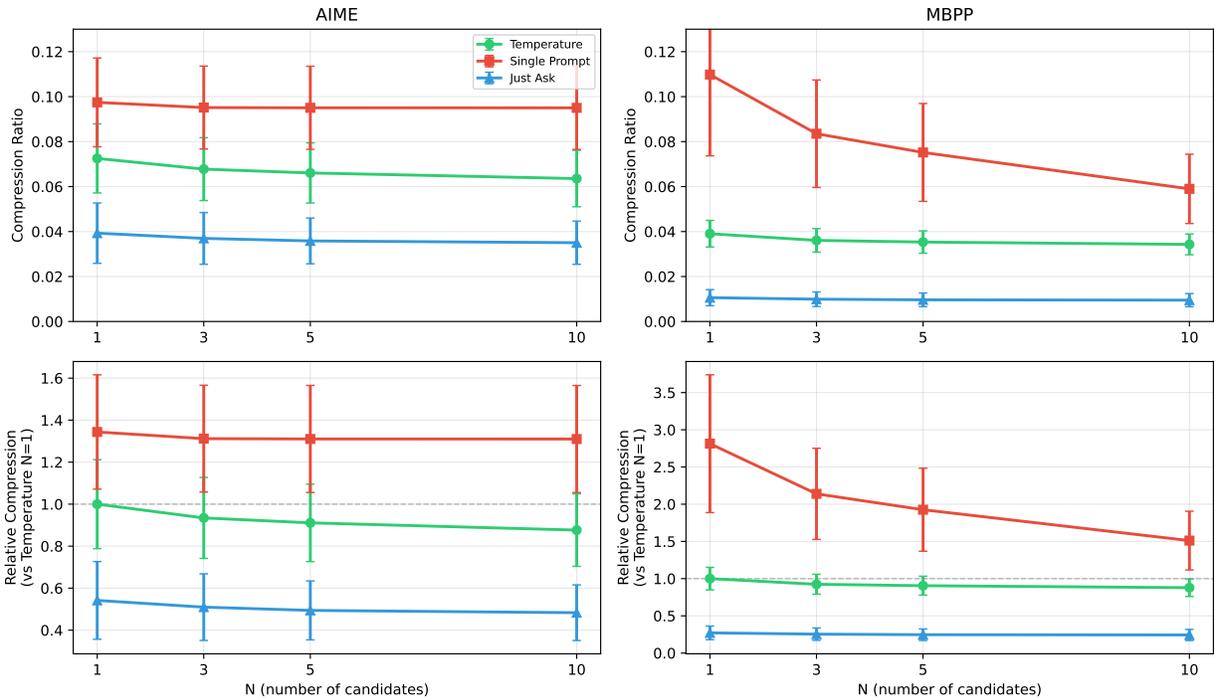


Figure 4: Absolute compression ratio (top) and relative compression ratio normalized to Temperature Sampling at $N=1$ (bottom) on AIME (left) and MBPP (right) using Opus. The dashed gray line at 1.0 indicates the baseline: values below 1.0 represent improved compression, while values above 1.0 indicate worse compression than the baseline. Lower is better.

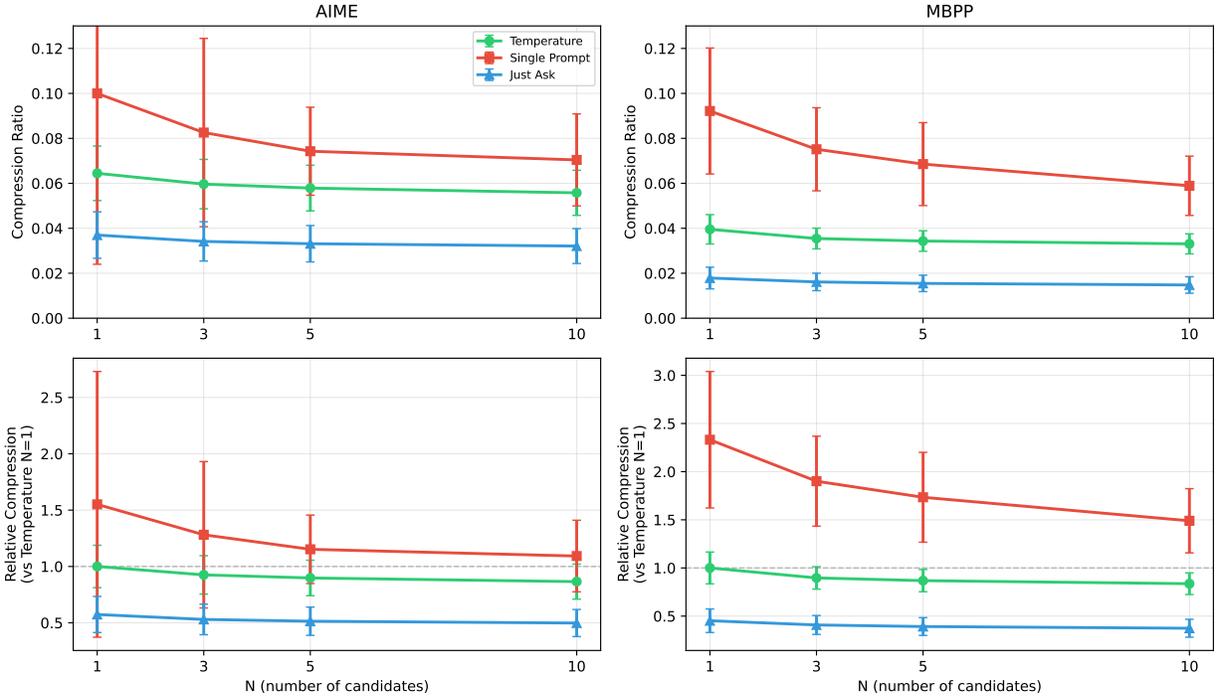


Figure 5: Absolute compression ratio (top) and relative compression ratio normalized to Temperature Sampling at $N=1$ (bottom) on AIME (left) and MBPP (right) using Haiku. The dashed gray line at 1.0 indicates the baseline: values below 1.0 represent improved compression, while values above 1.0 indicate worse compression than the baseline. Lower is better.

D.4 Haiku Results

Figure 6 shows the compression comparison using Haiku as both generator and compressor. The qualitative pattern matches Opus: summarization achieves the best compression, followed by temperature sampling, with single-prompt generation performing worst. On AIME, temperature sampling achieves compression ratios from 0.065 at $N = 1$ to 0.055 at $N = 10$, while summarization achieves 0.036 to 0.032.

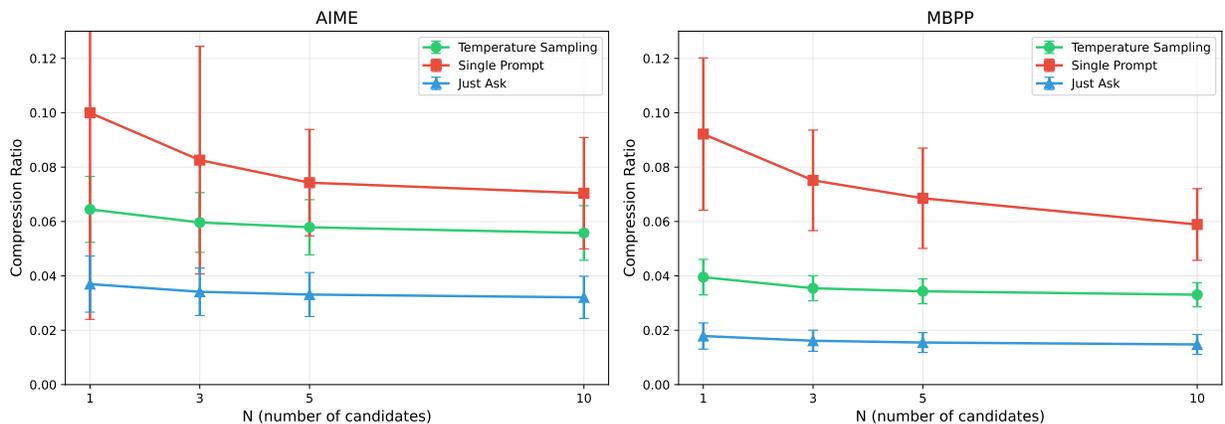


Figure 6: Compression ratio vs. number of candidates N on AIME (left) and MBPP (right) problems, using Haiku as both generator and compressor. The qualitative pattern matches Opus (Figure 2): Just Ask achieves the best compression, followed by Temperature Sampling, with Single Prompt performing worst.

Figure 7 shows the accuracy trade-off for Haiku. Notably, for *single prompt*, compression-based selection *improves* accuracy from 30% (random) to 38% at $N=10$, an 8 percentage point gain. This can be explained by examining the compressibility of correct versus incorrect solutions: correct answers are on average more compressible than incorrect ones (-0.65% and -1.71% difference in compression ratio for temperature sampling and single prompt, respectively), suggesting that correct solutions follow more predictable reasoning patterns. For *just ask*, the correlation is negligible ($+0.23\%$), consistent with the flat accuracy curves.

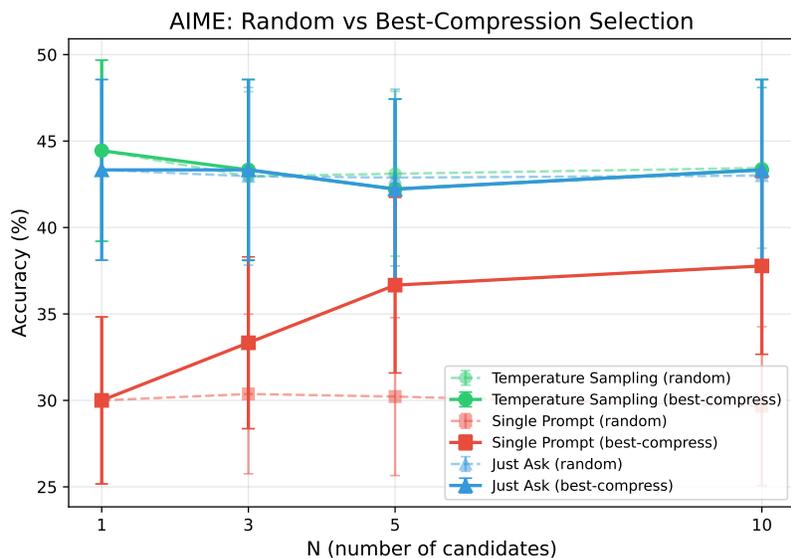


Figure 7: Accuracy of random selection (dashed) versus best-compression selection (solid) on 90 AIME problems using Haiku. For single prompt generation, compression-based selection improves accuracy by up to 8 percentage points at $N=10$, suggesting that correct solutions follow more predictable reasoning patterns.

D.5 Knowledge Density and Compressibility

An interesting phenomenon emerges when we examine what happens to text after summarization: the rewritten text is shorter but *less compressible per character*. Figure 8 illustrates this on AIME problems using Haiku. The verbose original solutions contain predictable boilerplate (“Step 1:”, “Therefore,” , “Let us consider”) that an LLM easily predicts, resulting in low bits-per-character. When we ask the model to rewrite the solution as succinctly as possible, the boilerplate is stripped away and what remains is denser, less predictable text: variable names, key equations, and compressed reasoning steps that the LLM used for compression cannot as easily anticipate.

However, the reduction in length is larger than the increase in per-character complexity. A typical verbose solution of $\sim 2,100$ characters at $\sim 5.3\%$ compression ratio becomes a ~ 590 -character rewrite at $\sim 9.5\%$ compression ratio. Despite the nearly doubled bits-per-character, the total compressed size drops by roughly 40% because the text is $3.6\times$ shorter. In other words, summarization increases the density of information of the text while still reducing the total information cost.

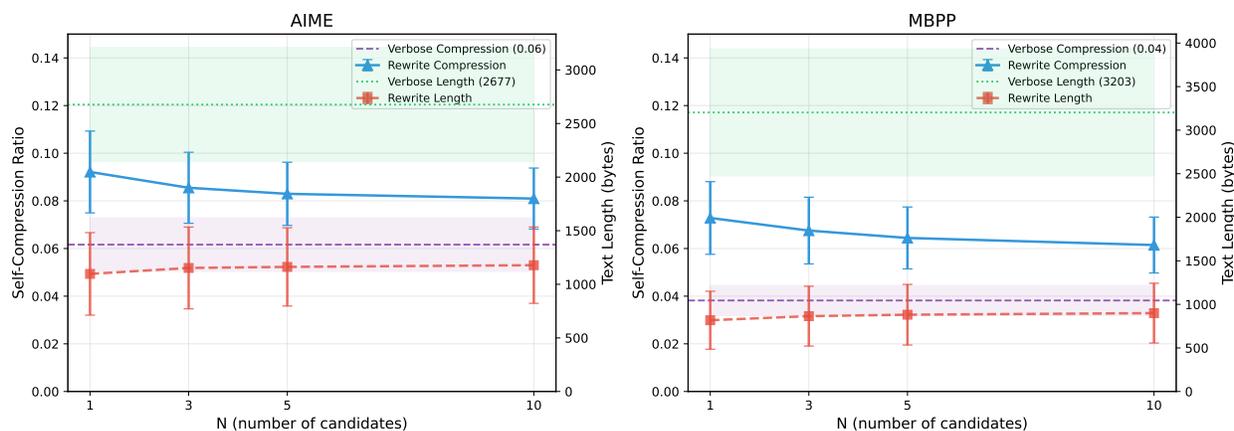


Figure 8: Compression ratio (bits-per-character) for verbose original solutions versus succinct rewrites on AIME problems using Haiku. Summarized text has higher bits-per-character (less redundancy per character) but is substantially shorter, yielding lower total compressed size.

E Negative Results & Ablations

In this section, we document approaches that did not yield the expected improvements, as well as ablation studies that provide insight into the mechanisms underlying our methods.

E.1 Judge Design

The iterative QA protocol requires a judge to decide whether the SLM’s updated answer is satisfactory (and should be accepted) or whether to continue asking questions. Two natural options arise:

1. **Objective (standalone) judging.** The judge evaluates the SLM’s answer on its own merits, scoring solution quality on a 1–10 scale. This does not require a reference solution, but relies on the judge’s ability to assess correctness independently.
2. **Comparison judging.** The judge compares the SLM’s answer against the LLM’s own solution, scoring how similar or aligned they are. This gives the judge a concrete reference, but the LLM’s own answer may itself be wrong on hard problems.

We evaluate both options below, along with ablations on the quality threshold and gold-answer access.

E.1.1 Quality-Thresholding Judge in QA Compression

A natural idea for improving the QA protocol is to add a *quality-thresholding judge*: after each round of questions, an LLM evaluates the SLM’s updated answer on a 1–10 scale. If the score meets a threshold (≥ 7), the protocol accepts the current answer and early-stops; otherwise, it continues asking questions in hopes of further improvement.

We test this by comparing the standard QA protocol against an iterative variant that adds an objective judge. Both protocols use the same model configurations, BL-CoT (Haiku→Haiku→Haiku) and QA (Haiku→Opus→Haiku), evaluated on Claude 4.5 models, with a budget of 10 binary yes/no questions each. The protocols differ in two ways:

1. **Quality thresholding:** The iterative variant adds a judge that evaluates the SLM’s answer on a 1–10 scale (mathematical soundness, calculation correctness, reasoning clarity) after each batch of 5 questions. If the score ≥ 7 , the protocol accepts the current answer and early-stops.
2. **Gold answer access:** In the standard protocol, the LLM answering questions is given the gold answer as reference. In the iterative variant, the LLM generates its own solution first, then uses it as reference for answering questions.

These two changes are confounded: the iterative variant both removes gold answer access and adds judging. However, since QA compression in practice cannot rely on gold answers, the iterative protocol represents the more realistic setting.

Table 16 shows that the iterative protocol consistently *reduces* recovery rates compared to the standard protocol. For the QA configuration (Haiku→Opus→Haiku), recovery drops on all 8 datasets. The “Standard” column reproduces the results from Table 21.⁸

Table 16: Recovery rates (%) for the standard QA protocol (Table 21) vs. the iterative variant. “Standard”: 10 questions in one batch, LLM given gold answer, no judge. “Iterative”: 2 batches of 5 questions, LLM generates own solution, objective judge with threshold $\geq 7/10$. Δ : change in recovery rate. Negative values indicate the iterative protocol performs worse.

Dataset	BL-CoT (h→h→h)			QA (h→o→h)		
	Std	Iter	Δ	Std	Iter	Δ
GSM8K	23.4	26.9	+3.5	31.0	27.6	-3.4
MATH (Algebra)	40.5	36.0	-4.5	38.2	36.0	-2.2
MATH (Geometry)	38.4	33.3	-5.1	38.9	31.3	-7.6
MATH (Num. Thy)	36.1	44.3	+8.2	44.3	42.6	-1.7
GPQA (MC)	13.2	8.0	-5.2	17.1	8.7	-8.4
MBPP	15.1	14.2	-0.9	17.9	17.0	-0.9
AIME	19.6	9.8	-9.8	17.6	11.8	-5.8
HLE	5.2	2.7	-2.5	7.8	1.3	-6.5
Mean	23.9	21.9	-2.0	26.6	22.0	-4.6

⁸Sample sizes differ slightly between the two protocols: the standard protocol evaluates all non-easy problems, while the iterative variant caps at 150 per dataset.

Analysis of judge behavior. To understand the role of the quality-thresholding judge specifically, Table 17 examines its behavior in the iterative protocol. Note that “early stop” means the judge was *satisfied* (score ≥ 7) and terminated the protocol, accepting the current answer.

Table 17: Judge behavior statistics (objective judge, QA config). Early-stop %: fraction of problems where the judge scored ≥ 7 and accepted the answer before all 10 questions were used. Avg. Score: mean judge rating (1–10). Score < 7 : fraction of individual evaluations below the acceptance threshold.

Dataset	Early Stop (%)	Avg Qs Used	Avg Score	Score < 7 (%)
GSM8K	89.0	5.9	8.5	19.0
MATH (Algebra)	98.7	5.1	9.7	3.0
MATH (Geometry)	78.7	6.3	7.2	36.6
MATH (Num. Thy)	90.2	5.7	8.4	20.0
GPQA (MC)	90.0	6.1	7.5	25.5
MBPP	92.5	5.5	8.3	12.3
AIME	19.6	9.4	3.0	91.8
HLE	50.7	7.7	5.4	64.7

The data reveals two distinct failure modes:

1. **Easy-to-judge datasets** (MATH Algebra: avg. score 9.7, 98.7% early stopped). The judge gives high scores and accepts the answer after just one round of 5 questions, before the SLM has received enough guidance from the Q&A exchange. The no-judge protocol would have continued for a full 10 questions, giving the SLM more information to work with.
2. **Hard-to-judge datasets** (AIME: avg. score 3.0, 80.4% *not* early stopped; HLE: avg. score 5.4, 49.3% not early stopped). The judge scores are persistently low, so the protocol runs all 10 questions but the final answer is still scored poorly. On AIME, problems that were *not* early stopped show severe regression: 7 out of 41 were initially correct, but only 1 remained correct after the Q&A exchange, a regression rate of 86%.

In neither regime does the judge provide a net benefit. On easy problems, it terminates the protocol prematurely; on hard problems, it fails to prevent regression. The judge lacks the ability to distinguish between “the answer improved but is still imperfect” (where it should continue) and “the Q&A guidance is making things worse” (where it should revert to the original answer).

Beyond the judge, the removal of gold answer access likely also contributes to the drop in recovery. When the LLM must generate its own solution to answer questions, its yes/no responses may be less reliable, particularly on hard problems where the LLM itself struggles. To partially disentangle these factors, we run two additional ablations:

1. **Higher threshold (≥ 9):** Raising the quality threshold from 7 to 9 should reduce premature early stopping, since the judge will accept fewer answers. This isolates whether the threshold level is the primary issue.
2. **Gold-answer judge:** The judge is given the gold answer for evaluation, while the LLM still answers questions using its own solution. This isolates the judge’s evaluation quality from the LLM’s question-answering quality.

Higher threshold (≥ 9). Table 18 compares recovery rates under the default threshold (≥ 7) and a stricter threshold (≥ 9). The stricter threshold reduces premature early stopping: for example, on GPQA the early-stop rate drops from 90% to 61%, and on HLE from 51% to 7%. Despite forcing the protocol to ask more questions, recovery rates do not improve.

Table 18: Recovery rates (%) for the iterative protocol with quality threshold ≥ 7 (default) vs. ≥ 9 . Both use the objective judge, no gold answer. Δ : change in recovery rate.

Dataset	BL-CoT (h→h→h)			QA (h→o→h)		
	$t=7$	$t=9$	Δ	$t=7$	$t=9$	Δ
GSM8K	18.8	14.8	-3.9	19.5	20.3	+0.8
MATH (Algebra)	27.3	26.6	-0.8	27.3	25.8	-1.6
MATH (Geometry)	25.8	19.4	-6.5	23.4	22.6	-0.8
MATH (Num. Thy)	24.4	19.5	-4.9	19.5	19.5	+0.0
GPQA (MC)	5.1	5.1	+0.0	6.5	6.5	+0.0
MBPP	12.7	14.7	+2.0	16.7	18.6	+2.0
AIME	5.3	2.6	-2.6	5.3	2.6	-2.6
HLE	2.8	0.7	-2.1	1.4	1.4	+0.0
Mean	15.3	12.9	-2.3	15.0	14.7	-0.3

Gold-answer judge. Table 19 compares the default iterative protocol (judge evaluates without the gold answer) to a variant where the judge is given the gold answer for evaluation. The LLM still answers questions using its own solution in both cases, so this isolates the judge’s evaluation quality from the LLM’s question-answering quality.

Table 19: Recovery rates (%) for the iterative protocol without vs. with gold answer access for the judge. Both use the objective judge with threshold ≥ 7 . The LLM answers questions using its own solution in both cases. Δ : change in recovery rate.

Dataset	BL-CoT (h→h→h)			QA (h→o→h)		
	Base	Gold	Δ	Base	Gold	Δ
GSM8K	18.8	20.3	+1.6	19.5	19.5	+0.0
MATH (Algebra)	27.3	26.6	-0.8	27.3	27.3	+0.0
MATH (Geometry)	25.8	19.4	-6.5	23.4	25.8	+2.4
MATH (Num. Thy)	24.4	24.4	+0.0	19.5	19.5	+0.0
GPQA (MC)	5.1	5.8	+0.7	6.5	8.0	+1.4
MBPP	12.7	14.7	+2.0	16.7	17.6	+1.0
AIME	5.3	7.9	+2.6	5.3	5.3	+0.0
HLE	2.8	1.4	-1.4	1.4	1.4	+0.0
Mean	15.3	15.1	-0.2	15.0	15.6	+0.6

Comparison-mode judge. We also test a *comparison-mode* judge, which rates the SLM’s answer relative to the LLM’s own solution rather than evaluating it in isolation. Table 20 shows recovery rates across all three protocols: the standard (no-judge) protocol from Table 21, the iterative protocol with an objective judge, and the iterative protocol with a comparison judge. Both iterative variants use threshold ≥ 7 , and the LLM answers questions using its own solution (no gold answer).

Table 20: Recovery rates (%) across three protocols. “Standard”: no judge, LLM given gold answer, 10 questions in one batch. “Objective”: standalone quality judge, threshold ≥ 7 , no gold answer. “Comparison”: judge compares SLM answer to LLM answer, threshold ≥ 7 , no gold answer. Bold indicates the highest value in each row.

Dataset	Standard		Objective		Comparison	
	BL-CoT	QA	BL-CoT	QA	BL-CoT	QA
GSM8K	23.4	31.0	26.9	27.6	27.6	28.3
MATH (Algebra)	40.5	38.2	36.0	36.0	34.0	36.7
MATH (Geometry)	38.4	38.9	33.3	31.3	29.3	31.3
MATH (Num. Thy)	36.1	44.3	44.3	42.6	42.6	39.3
GPQA (MC)	13.2	17.1	8.0	8.7	10.0	11.3
MBPP	15.1	17.9	14.2	17.0	16.0	19.8
AIME	19.6	17.6	9.8	11.8	5.9	11.8
HLE	5.2	7.8	2.7	1.3	0.7	2.0

The standard protocol achieves the highest recovery on 6 of 8 datasets, confirming that the primary driver of the gap is gold answer access rather than judge mode. Between the two judge modes, comparison slightly outperforms objective for QA (mean 22.6% vs 22.0%) while slightly underperforming for BL-CoT (20.8% vs 21.9%). The comparison judge preserves more of the QA-over-BL-CoT advantage (+1.8pp vs +0.1pp for objective), likely because it evaluates relative to the LLM’s own answer, which is stronger in the QA configuration (Opus) than in BL-CoT (Haiku). However, neither judge mode closes the gap to the standard protocol.

F Block Emission Arithmetic Coding

This appendix provides detailed pseudocode for our block emission arithmetic coding scheme. The key challenge with standard arithmetic coding is that encoder and decoder must agree exactly on probability distributions. In practice, floating-point non-determinism, hardware differences, and batching effects can cause small mismatches that accumulate over long sequences, eventually causing decoding failures.

Block emission addresses this by periodically emitting the current interval midpoint whenever the range becomes too small, then restarting with a fresh $[0, 1]$ interval. This limits how much drift can accumulate before being “reset,” making the scheme robust to minor probability discrepancies.

F.1 Algorithms

Algorithm 2 Block Emission Encoding

Require: tokens $T = [t_1, \dots, t_n]$, model M , context C , bit precision B , emission threshold ε

Ensure: encoded blocks E , token counts K

```
1:  $low \leftarrow 0, high \leftarrow 1$ 
2:  $E \leftarrow [], K \leftarrow [], count \leftarrow 0$ 
3:  $M.reset(C)$  {Initialize model with context}
4: for each token  $t$  in  $T$  do
5:    $\mathbf{p} \leftarrow M.get\_distribution()$  {Next-token probabilities}
6:    $cdf \leftarrow cumsum(\mathbf{p})$  {Cumulative distribution}
7:    $range \leftarrow high - low$ 
8:    $new\_low \leftarrow low + range \cdot cdf[t - 1]$  { $cdf[0] = 0$ }
9:    $new\_high \leftarrow low + range \cdot cdf[t]$ 
10:  if  $(new\_high - new\_low) < \varepsilon$  or  $count \geq 2^5 - 1$  then
11:     $E.append\left(\left\lfloor \frac{low+high}{2} \cdot (2^B - 1) \right\rfloor\right)$  {Emit block}
12:     $K.append(count)$ 
13:     $low \leftarrow cdf[t - 1], high \leftarrow cdf[t]$  {Fresh interval}
14:     $count \leftarrow 1$ 
15:  else
16:     $low \leftarrow new\_low, high \leftarrow new\_high$ 
17:     $count \leftarrow count + 1$ 
18:  end if
19:   $M.add\_token(t)$  {Update model state}
20: end for
21:  $E.append\left(\left\lfloor \frac{low+high}{2} \cdot (2^B - 1) \right\rfloor\right)$  {Emit final block}
22:  $K.append(count)$ 
23: return  $(E, K)$ 
```

Algorithm 3 Block Emission Decoding

Require: encoded blocks E , token counts K , model M , context C , bit precision B

Ensure: decoded tokens T

```
1:  $T \leftarrow []$ 
2:  $M.\text{reset}(C)$  {Initialize model with context}
3: for  $i = 1$  to  $|E|$  do
4:    $v \leftarrow E[i]/(2^B - 1)$  {Dequantize to  $[0, 1]$ }
5:    $low \leftarrow 0, high \leftarrow 1$ 
6:   for  $j = 1$  to  $K[i]$  do
7:      $\mathbf{p} \leftarrow M.\text{get\_distribution}()$ 
8:      $cdf \leftarrow \text{cumsum}(\mathbf{p})$ 
9:      $range \leftarrow high - low$ 
10:     $target \leftarrow (v - low)/range$  {Normalize position in interval}
11:     $t \leftarrow \min\{k : cdf[k] > target\}$  {Find token via CDF lookup}
12:     $low \leftarrow low + range \cdot cdf[t - 1]$ 
13:     $high \leftarrow low + range \cdot cdf[t]$ 
14:     $T.\text{append}(t)$ 
15:     $M.\text{add\_token}(t)$ 
16:   end for
17: end for
18: return  $T$ 
```

F.2 Design Rationale

Block Emission. Standard arithmetic coding maintains a single interval $[low, high]$ that narrows with each encoded symbol. When encoder and decoder probability distributions differ slightly (due to floating-point non-determinism, hardware variations, or batching effects), these small mismatches accumulate over the sequence. Block emission addresses this by emitting the interval midpoint (quantized to B bits) whenever the range falls below threshold ε , then restarting with a fresh $[0, 1]$ interval. This periodic reset limits how much drift can accumulate, making the scheme robust to minor probability discrepancies between encoder and decoder.

Implementation Details. In our experiments we use $B = 58$ bits for each emitted midpoint and allocate additional bits for a per-block token count that records how many tokens are packed into each block. This token count is necessary for the decoder to know when one block ends and the next begins. With b bits for the count, each block can encode up to $2^b - 1$ tokens. The total overhead per block is therefore $B + b$ bits. We find this overhead negligible in practice: blocks typically contain 15–30 tokens, so the amortized cost is roughly 3–5 additional bits per token beyond the entropy-optimal rate.

G Q&A Compression Detailed Results (Correct-Answer Ablation)

This appendix presents an ablation where the LLM answering yes/no questions is given the correct answer from the dataset directly, rather than generating its own solution as a reference. Additionally, all $N = 10$ questions are generated in a single batch with no iterative quality thresholding or early stopping. This provides an approximation of an upper bound on the protocol’s potential: access to the correct answer ensures the LLM knows the final answer, but it does not provide

the reasoning that produced it, so the LLM cannot compare the SLM’s reasoning path against a reference derivation. In practice, the LLM’s own solution (used in the main results, Table 1) may sometimes be more useful for answering questions about *how* to arrive at the answer, even if the solution itself is occasionally wrong.

We compare two configurations:

- **BL-CoT (Bit-Limited Chain-of-Thought)**: SLM=Haiku, LLM=Haiku, Q=Haiku. The small model asks questions answered by itself.
- **QA**: SLM=Haiku, LLM=Opus, Q=Haiku. The small model asks questions answered by the large model (with access to the correct answer).

G.1 Recovery Rate by Difficulty (Claude 4.5)

Table 21 reports the recovery rate broken down by difficulty level. Problems are grouped as **Medium** (Haiku fails, Opus succeeds), **Hard** (Haiku and Sonnet fail, Opus succeeds), and **Very Hard** (all models fail). n is the total number of non-easy problems per dataset.

Table 21: Recovery rate (%) on problems where Haiku initially fails (Claude 4.5). BL-CoT = Bit-Limited Chain-of-Thought (haiku→haiku→haiku); QA = QA Compression (haiku→opus→haiku). $\Delta = \text{QA} - \text{BL-CoT}$. Higher is better.

Dataset	n	Medium		Hard		Very Hard		All Non-Easy		
		BL-CoT	QA	BL-CoT	QA	BL-CoT	QA	BL-CoT	QA	Δ
GSM8K	145	19.0	23.8	20.0	60.0	37.1	48.6	23.4	31.0	+7.6
MATH (Algebra)	259	28.3	20.0	48.0	48.0	43.7	43.1	40.5	38.2	-2.3
MATH (Geometry)	190	50.0	43.3	28.6	28.6	37.4	39.6	38.4	38.9	+0.5
MATH (Num. Theory)	61	55.6	88.9	0.0	66.7	34.7	34.7	36.1	44.3	+8.2
GPQA (MC)	152	50.0	33.3	11.1	33.3	11.7	15.3	13.2	17.1	+3.9
MBPP	106	48.1	59.3	5.9	5.9	3.2	3.2	15.1	17.9	+2.8
AIME	51	0.0	0.0	10.7	7.1	41.2	41.2	19.6	17.6	-2.0
HLE [†]	306	14.0	14.0	3.8	7.1	3.0	4.0	5.2	7.8	+2.6

[†] For HLE, Very Hard problems were randomly subsampled to 100 (out of ~1,900) due to evaluation cost.

Q&A guidance from Opus helps most on Medium problems, with GSM8K showing +7.6% overall and MATH Number Theory showing the strongest Medium improvement (55.6% to 88.9%). On Very Hard problems where even Opus fails, Q&A provides minimal benefit.

G.2 100-Bit QA Compression (Claude 4.5)

We also evaluated QA compression with 100 yes/no questions instead of 10. Table 22 reports accuracy with standard deviation over 3 trials. Unlike the 10-question ablation above (which uses the gold answer as reference), this experiment uses the LLM’s own generated solution.

Table 22: QA Compression accuracy on non-easy problems (% , mean \pm std over 3 trials) comparing 10-bit and 100-bit protocols (Claude 4.5). BL-CoT = haiku \rightarrow haiku \rightarrow haiku; QA = haiku \rightarrow opus \rightarrow haiku; QA+ = haiku \rightarrow opus \rightarrow opus.

Dataset	n	BL-CoT			QA			QA+		
		10-bit	100-bit	Δ	10-bit	100-bit	Δ	10-bit	100-bit	Δ
GSM8K	145	22.5 \pm 1.2	24.6 \pm 1.4	+2.1	26.7 \pm 0.9	30.8 \pm 0.3	+4.1	29.7 \pm 1.7	35.4 \pm 1.7	+5.7
MATH (Algebra)	259	38.1 \pm 1.6	38.2 \pm 1.4	+0.1	39.1 \pm 1.0	40.0 \pm 1.9	+0.9	40.0 \pm 1.5	40.2 \pm 1.6	+0.2
MATH (Geometry)	190	33.7 \pm 0.4	37.4 \pm 0.9	+3.7	39.6 \pm 0.5	37.2 \pm 1.1	-2.4	38.2 \pm 1.5	41.1 \pm 3.0	+2.9
MATH (Num. Theory)	61	38.3 \pm 2.0	39.3 \pm 1.3	+1.0	43.7 \pm 4.3	41.5 \pm 0.8	-2.2	40.4 \pm 3.1	45.9 \pm 1.3	+5.5
GPQA (MC)	152	11.8 \pm 2.7	11.4 \pm 3.0	-0.4	13.6 \pm 0.3	13.8 \pm 0.0	+0.2	12.1 \pm 2.0	13.2 \pm 1.4	+1.1
AIME	51	18.3 \pm 1.8	17.0 \pm 4.6	-1.3	19.0 \pm 2.4	19.6 \pm 2.8	+0.6	22.2 \pm 1.8	16.3 \pm 1.8	-5.9
HLE	306	5.4 \pm 0.6	6.1 \pm 0.6	+0.7	6.5 \pm 0.0	5.9 \pm 0.9	-0.6	9.3 \pm 0.8	8.8 \pm 0.9	-0.5
MBPP	106	18.2 \pm 1.8	19.2 \pm 0.9	+1.0	18.2 \pm 0.4	18.6 \pm 0.9	+0.4	26.1 \pm 0.4	23.6 \pm 0.0	-2.5

Table 23 shows the same comparison restricted to Medium and Hard problems only (excluding Very Hard).

Table 23: QA Compression accuracy on Medium + Hard problems only (% , mean \pm std over 3 trials) comparing 10-bit and 100-bit protocols (Claude 4.5). Excludes Very Hard problems where even Opus fails.

Dataset	n	BL-CoT			QA			QA+		
		10-bit	100-bit	Δ	10-bit	100-bit	Δ	10-bit	100-bit	Δ
GSM8K	110	19.1 \pm 2.0	21.8 \pm 0.7	+2.7	23.3 \pm 1.1	27.6 \pm 0.4	+4.3	27.6 \pm 1.5	33.9 \pm 1.1	+6.3
MATH (Algebra)	85	28.6 \pm 1.5	28.6 \pm 1.1	+0.0	27.5 \pm 2.2	29.8 \pm 3.9	+2.3	31.0 \pm 1.5	32.5 \pm 0.6	+1.5
MATH (Geometry)	51	28.1 \pm 1.8	37.9 \pm 3.7	+9.8	37.3 \pm 3.2	36.6 \pm 6.1	-0.7	37.3 \pm 1.6	46.4 \pm 5.1	+9.1
MATH (Num. Theory)	12	66.7 \pm 11.8	61.1 \pm 3.9	-5.6	72.2 \pm 15.7	55.6 \pm 7.9	-16.6	63.9 \pm 15.7	75.0 \pm 6.8	+11.1
GPQA (MC)	15	40.0 \pm 14.4	33.3 \pm 5.4	-6.7	40.0 \pm 14.4	40.0 \pm 5.4	+0.0	31.1 \pm 6.3	42.2 \pm 6.3	+11.1
AIME	34	10.8 \pm 1.4	7.8 \pm 5.0	-3.0	6.9 \pm 1.4	9.8 \pm 3.7	+2.9	9.8 \pm 1.4	6.9 \pm 2.8	-2.9
HLE	206	6.5 \pm 0.6	7.6 \pm 1.4	+1.1	7.9 \pm 0.6	7.6 \pm 1.4	-0.3	13.1 \pm 1.4	11.7 \pm 0.4	-1.4
MBPP	44	36.4 \pm 1.9	42.4 \pm 2.1	+6.0	37.9 \pm 2.8	35.6 \pm 3.9	-2.3	51.5 \pm 2.8	43.9 \pm 2.8	-7.6

G.3 Claude 3.5/4 Model Family Results

Table 24 shows Q&A compression results using the older Claude 3.5/4 model family.

Table 24: Recovery rate (%) on problems where Haiku initially fails (Claude 3.5/4). BL-CoT = Bit-Limited Chain-of-Thought (haiku \rightarrow haiku \rightarrow haiku, self-refinement baseline); QA = QA Compression (haiku \rightarrow opus \rightarrow haiku, knowledge transfer). Δ = QA - BL-CoT. Bold indicates the highest value in each row (excluding Δ).

Dataset	n	Medium			Hard			Very Hard			All Non-Easy		
		BL-CoT	QA	Δ	BL-CoT	QA	Δ	BL-CoT	QA	Δ	BL-CoT	QA	Δ
GSM8K	109	30.0	50.0	+20.0	20.0	0.0	-20.0	52.3	45.5	-6.8	38.5	45.9	+7.4
MATH (Algebra)	342	29.5	32.0	+2.5	50.0	41.7	-8.3	52.4	49.5	-2.9	44.2	43.0	-1.2
MATH (Geometry)	269	23.9	23.9	0.0	12.5	12.5	0.0	15.8	13.8	-2.0	19.0	17.8	-1.2
MATH (Num. Theory)	154	43.2	46.3	+3.1	50.0	33.3	-16.7	24.5	24.5	0.0	37.0	38.3	+1.3
GPQA (MC)	156	11.1	22.2	+11.1	33.3	11.1	-22.2	12.3	11.6	-0.7	13.5	12.2	-1.3
MBPP	100	50.0	57.1	+7.1	20.0	20.0	0.0	12.3	9.9	-2.4	18.0	17.0	-1.0
AIME	87	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
HLE	164	5.9	5.9	0.0	1.2	6.2	+5.0	2.0	4.0	+2.0	2.4	5.5	+3.1

G.4 Cross-Model Validation with GPT-OSS

To validate that the QA protocol generalizes beyond Claude models, we evaluate using GPT-OSS-120B (OpenAI et al., 2025), an open-source reasoning model, as the SLM.

We test three configurations:

- **BL-CoT (g/g/g)**: SLM=GPT-OSS, LLM=GPT-OSS, Q=GPT-OSS. Self-refinement baseline.
- **QA (g/o/g)**: SLM=GPT-OSS, LLM=Opus, Q=GPT-OSS. Cross-model knowledge transfer.
- **QA+ (g/o/o)**: SLM=GPT-OSS, LLM=Opus, Q=Opus. Full cross-model protocol.

Problems are filtered to those where GPT-OSS initially fails but Opus succeeds, analogous to the Medium+Hard filtering used for Haiku in the main results.

Table 25: Recovery rate (%) with GPT-OSS as SLM on non-easy problems (mean \pm std over 3 trials). n = number of problems where GPT-OSS initially fails.

Dataset	n	BL-CoT (g/g/g)	QA (g/o/g)	QA+ (g/o/o)
GSM8K	59	37.3 \pm 3.7	37.9 \pm 5.6	41.8 \pm 0.8
MATH (Algebra)	240	19.9 \pm 1.7	20.4 \pm 1.0	19.7 \pm 0.9
MATH (Geometry)	197	17.1 \pm 0.2	18.8 \pm 1.8	24.4 \pm 0.4
MATH (Num. Theory)	64	28.6 \pm 1.9	30.2 \pm 0.7	29.7 \pm 4.6
MBPP	88	11.4 \pm 3.2	12.1 \pm 1.4	12.9 \pm 1.9
AIME	83	5.6 \pm 2.5	6.8 \pm 2.0	13.7 \pm 2.0
HLE [†]	275	5.2 \pm 0.8	7.4 \pm 1.1	9.4 \pm 0.8

[†] For HLE, very hard problems were randomly subsampled to 100 due to evaluation cost.

G.5 Detailed Heatmaps by Dataset (Claude 4.5)

The following figures show Q&A compression accuracy heatmaps for each dataset, with SLM fixed to Haiku. Each heatmap shows accuracy broken down by difficulty level (Medium, Hard+Very Hard, All Non-Easy), with rows indicating the LLM (question answerer) and columns indicating the question-generating model (Q).

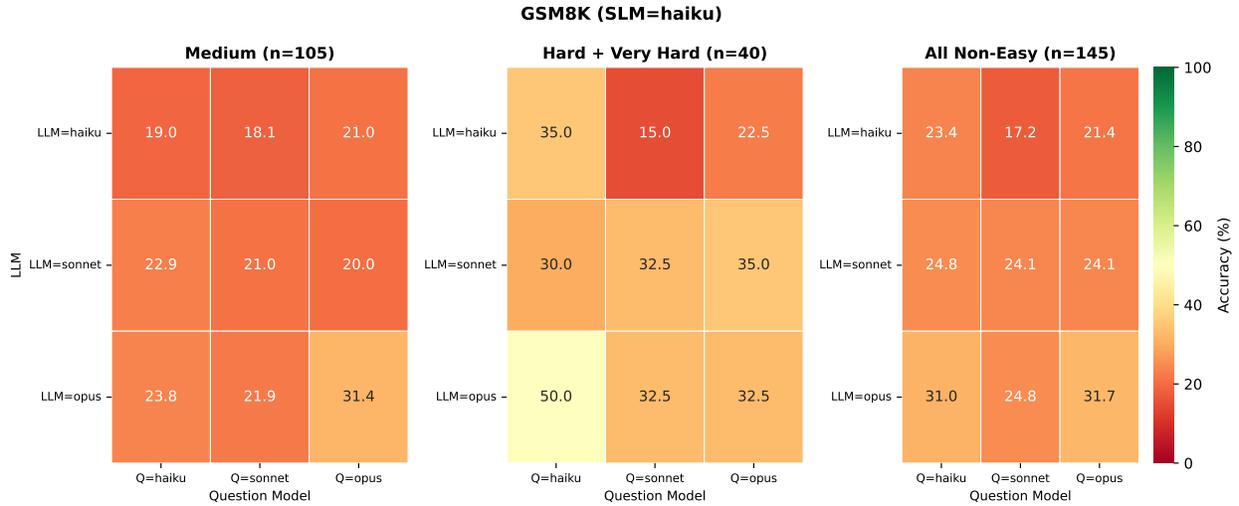


Figure 9: GSM8K Q&A compression accuracy (SLM=Haiku, Claude 4.5). GSM8K shows consistent improvement from Q&A guidance, particularly on Medium problems where Opus can provide useful feedback.

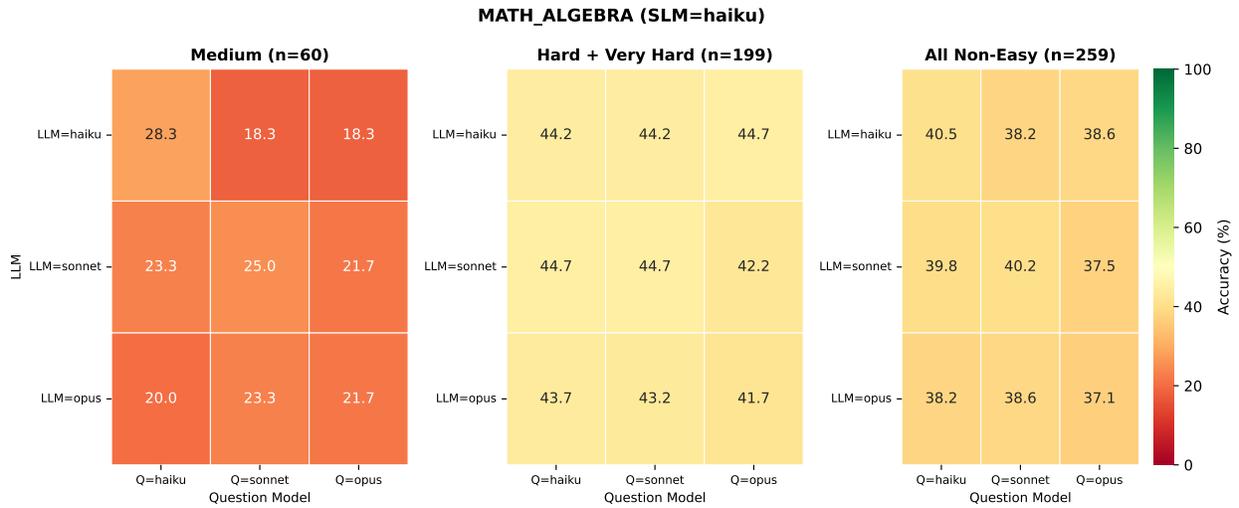


Figure 10: MATH Algebra Q&A compression accuracy (SLM=Haiku, Claude 4.5).

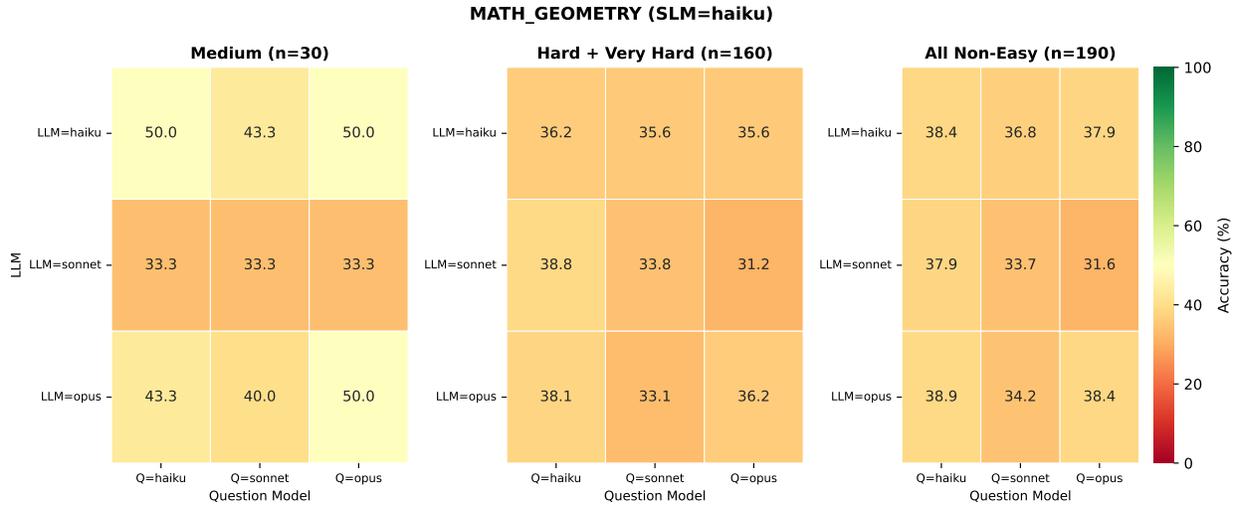


Figure 11: MATH Geometry Q&A compression accuracy (SLM=Haiku, Claude 4.5).

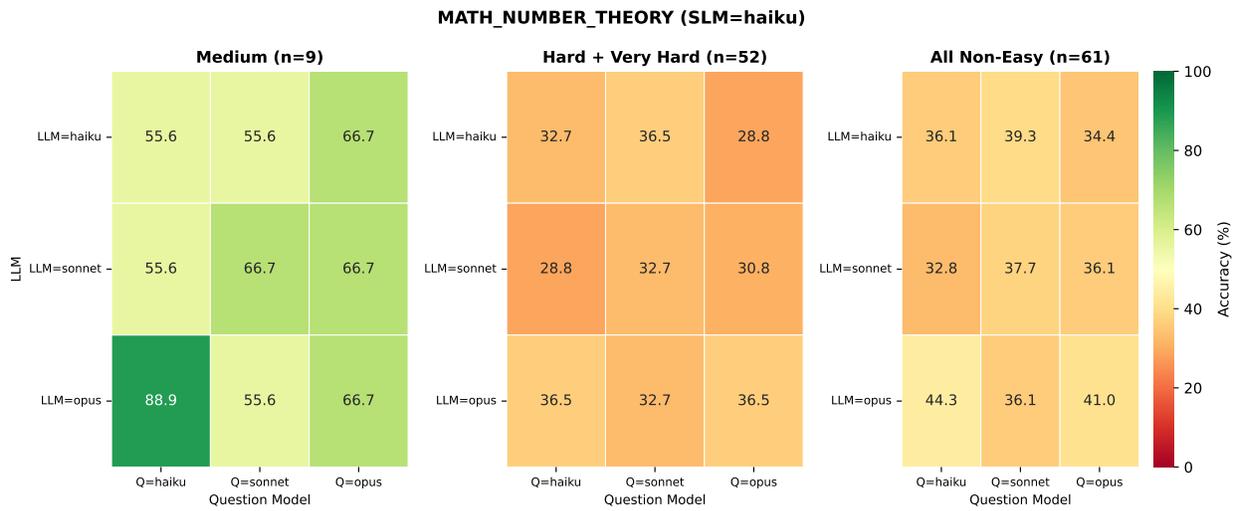


Figure 12: MATH Number Theory Q&A compression accuracy (SLM=Haiku, Claude 4.5). Number Theory shows strong improvement on Medium problems (55.6% to 88.9% with Opus guidance).

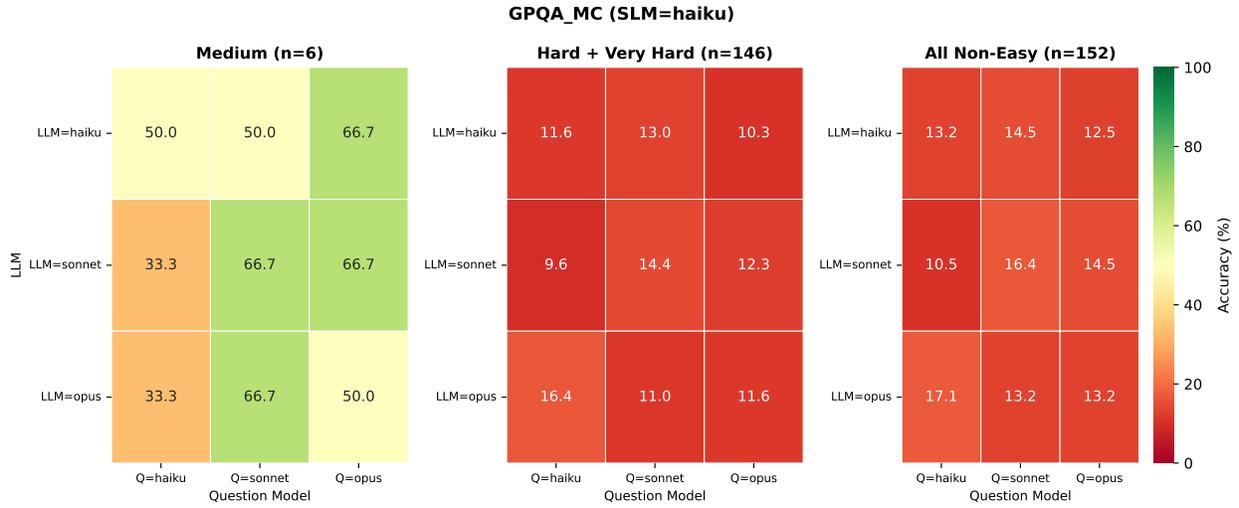


Figure 13: GPQA (MC) Q&A compression accuracy (SLM=Haiku, Claude 4.5). The high proportion of Very Hard problems limits improvement potential.

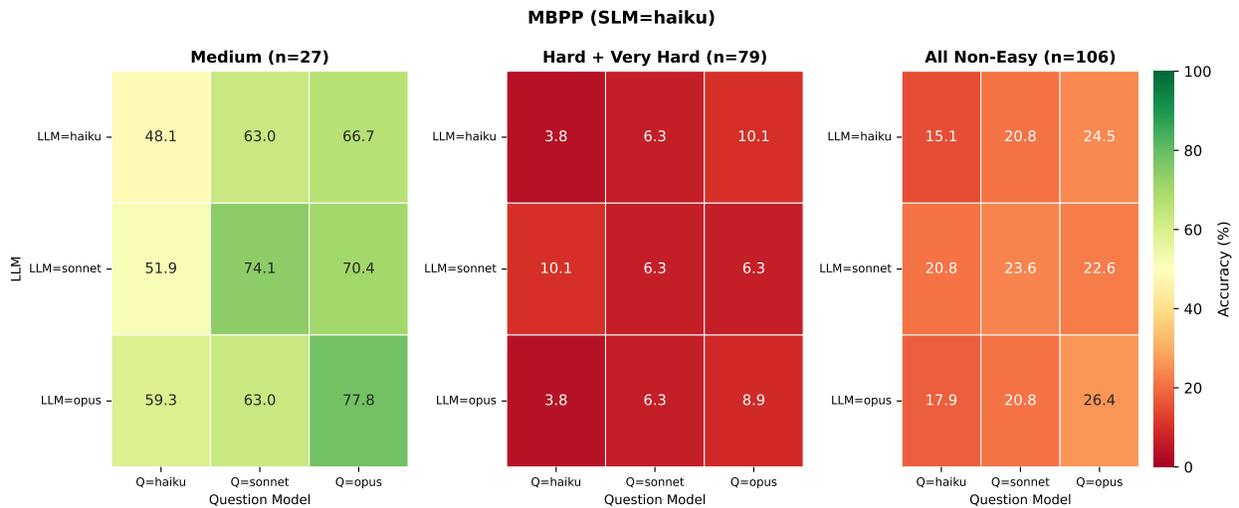


Figure 14: MBPP Q&A compression accuracy (SLM=Haiku, Claude 4.5). Code generation proves challenging for Q&A compression on Hard and Very Hard problems.

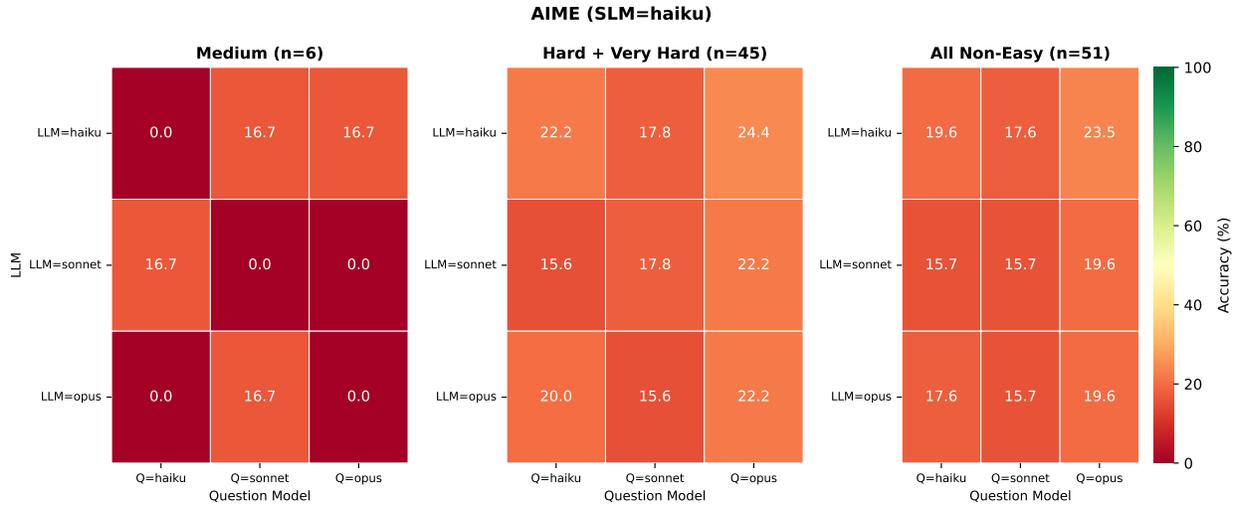


Figure 15: AIME Q&A compression accuracy (SLM=Haiku, Claude 4.5). Competition math problems show limited benefit from Q&A guidance.

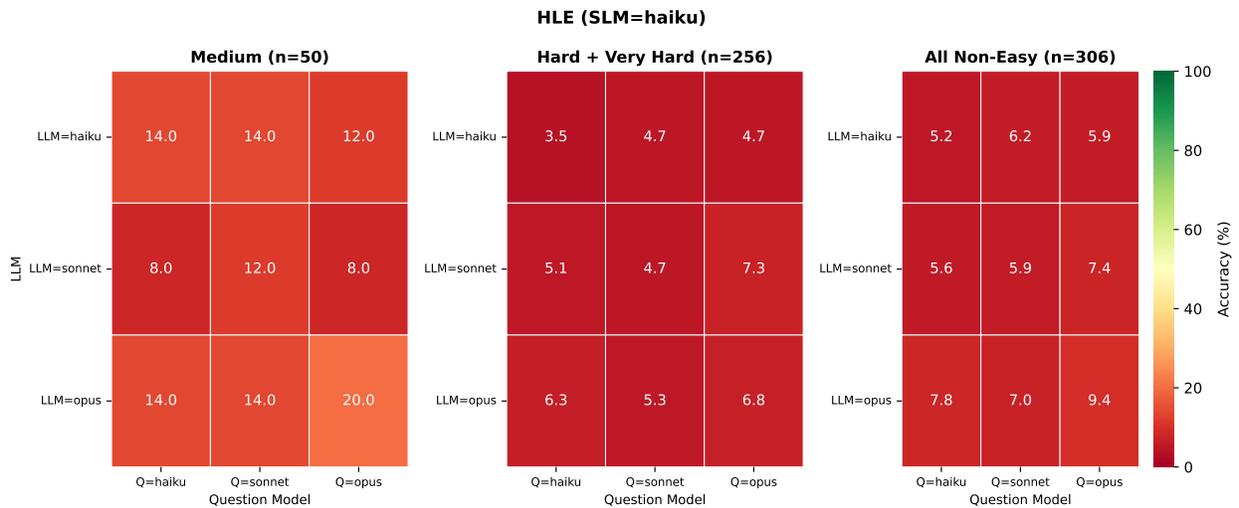


Figure 16: HLE Q&A compression accuracy (SLM=Haiku, Claude 4.5). Despite the high Very Hard rate (89%), Q&A shows modest improvements on Hard problems.

G.6 Detailed Heatmaps by Dataset (Claude 3.5/4)

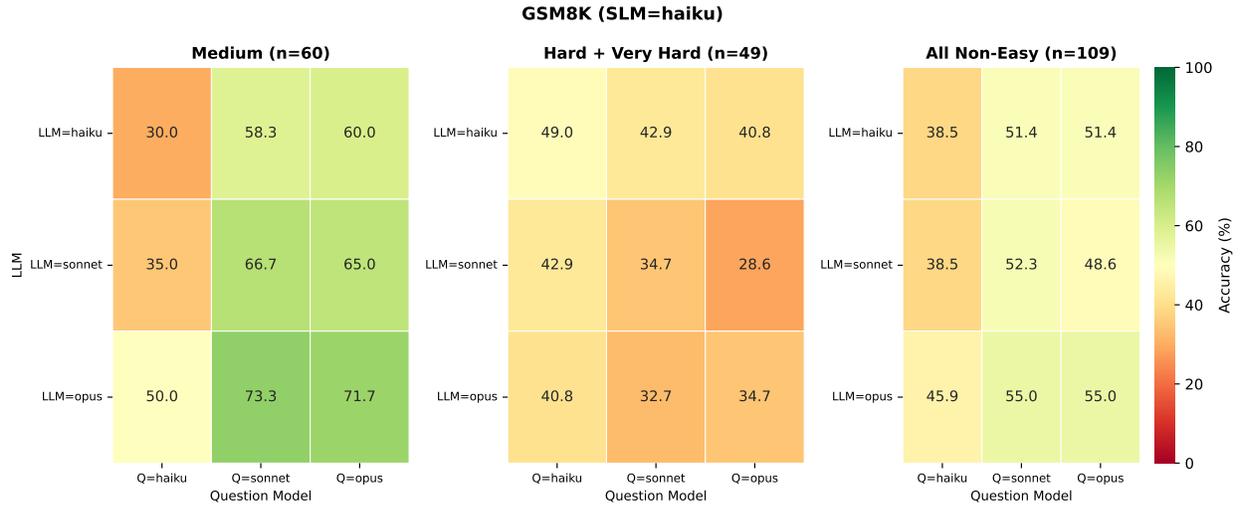


Figure 17: GSM8K Q&A compression accuracy (SLM=Haiku, Claude 3.5/4). Strong improvement on Medium problems (30% to 50%).

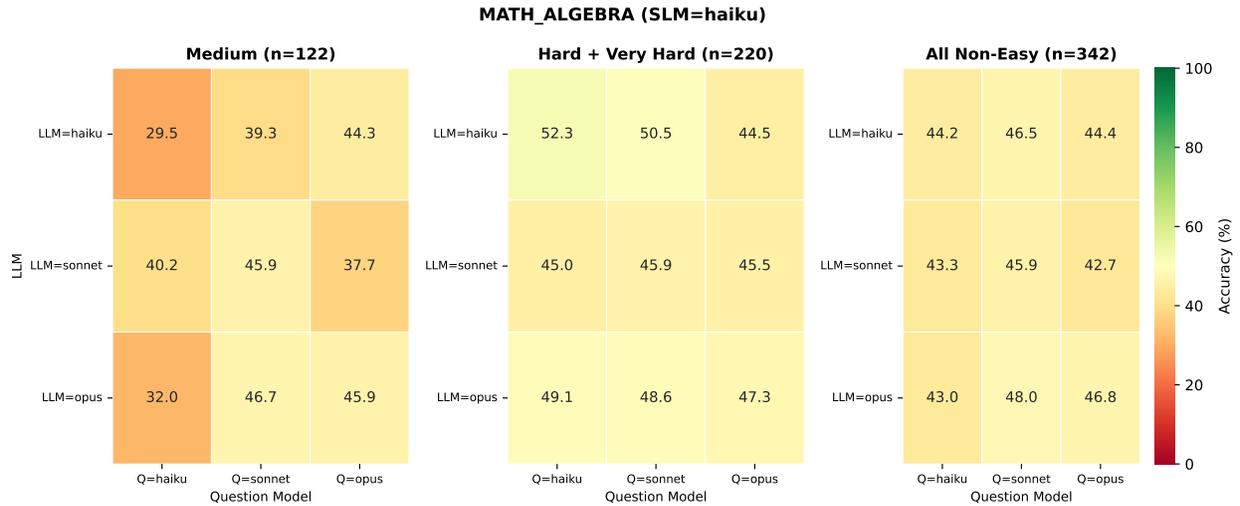


Figure 18: MATH Algebra Q&A compression accuracy (SLM=Haiku, Claude 3.5/4).

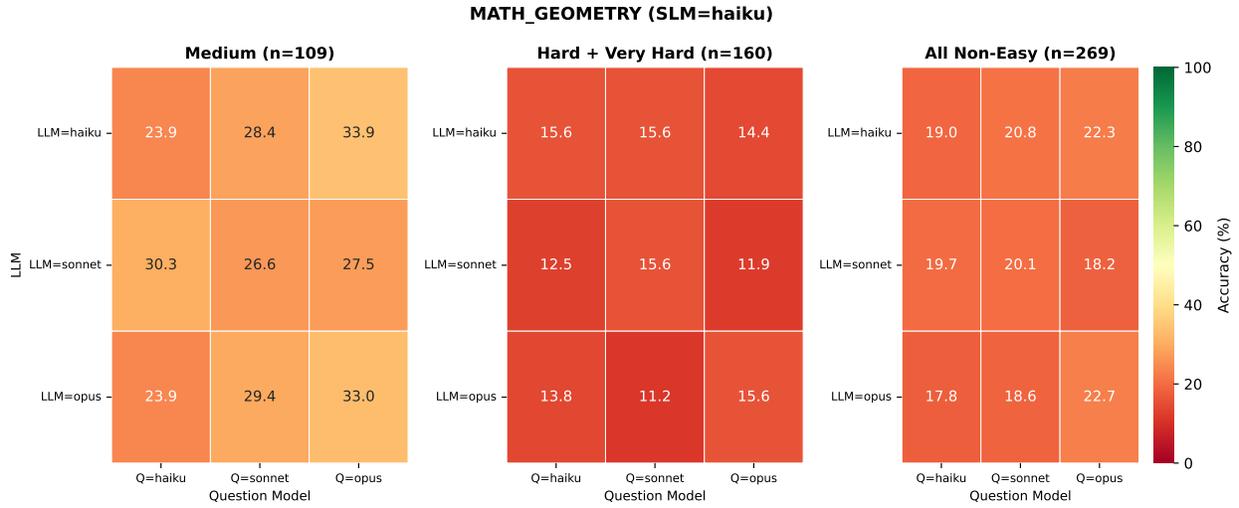


Figure 19: MATH Geometry Q&A compression accuracy (SLM=Haiku, Claude 3.5/4).

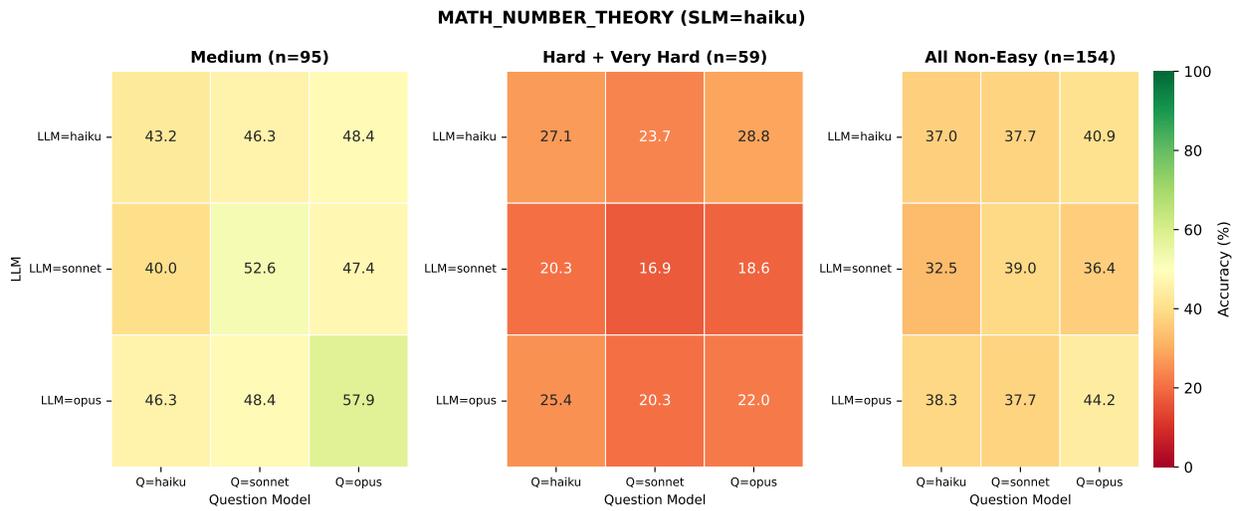


Figure 20: MATH Number Theory Q&A compression accuracy (SLM=Haiku, Claude 3.5/4).

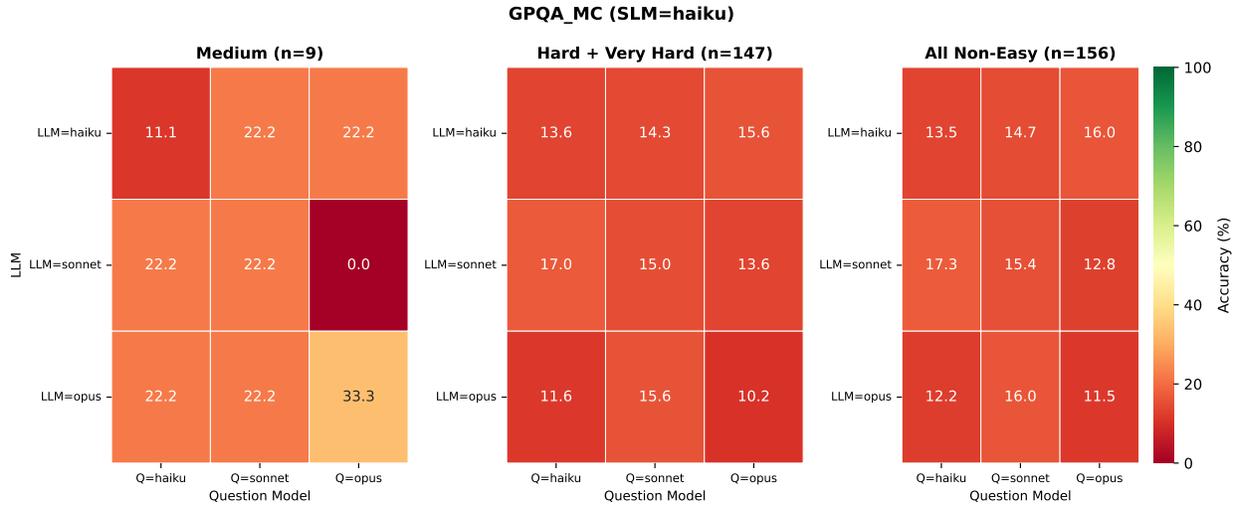


Figure 21: GPQA (MC) Q&A compression accuracy (SLM=Haiku, Claude 3.5/4).

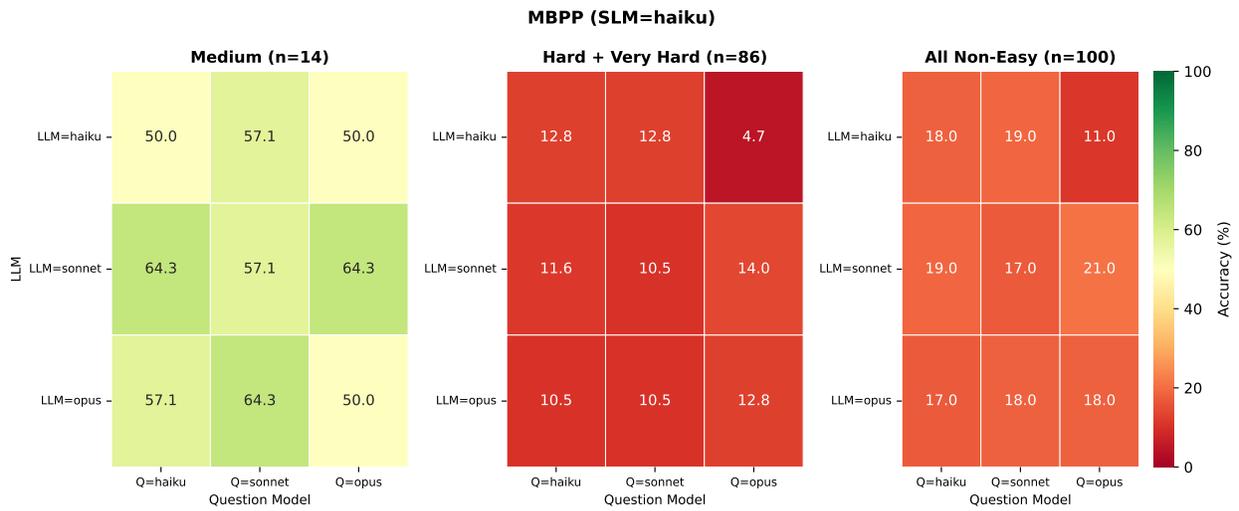


Figure 22: MBPP Q&A compression accuracy (SLM=Haiku, Claude 3.5/4).

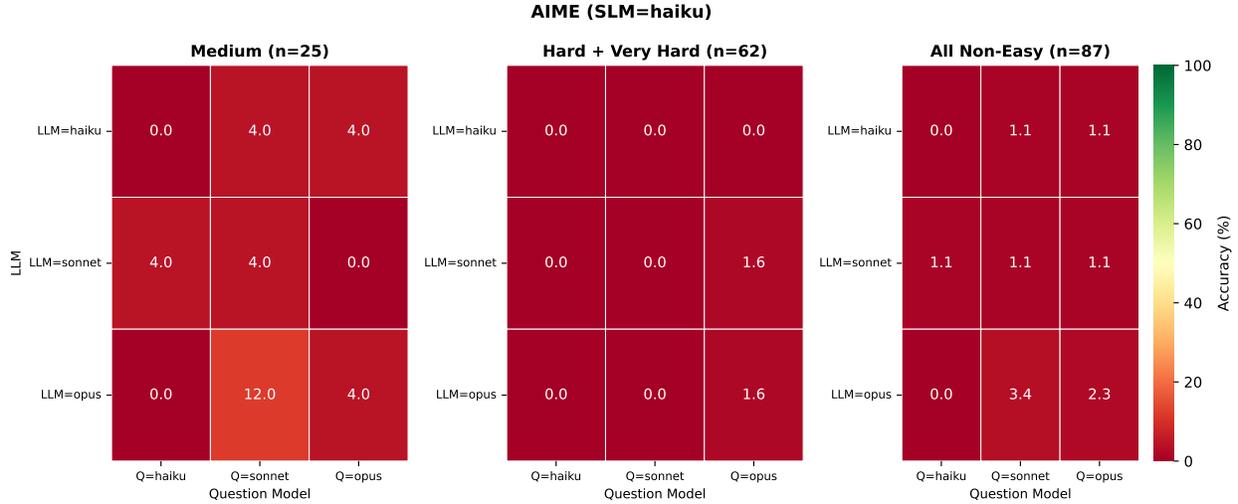


Figure 23: AIME Q&A compression accuracy (SLM=Haiku, Claude 3.5/4). The older Haiku model achieves 0% on all non-easy AIME problems, even with Opus guidance.

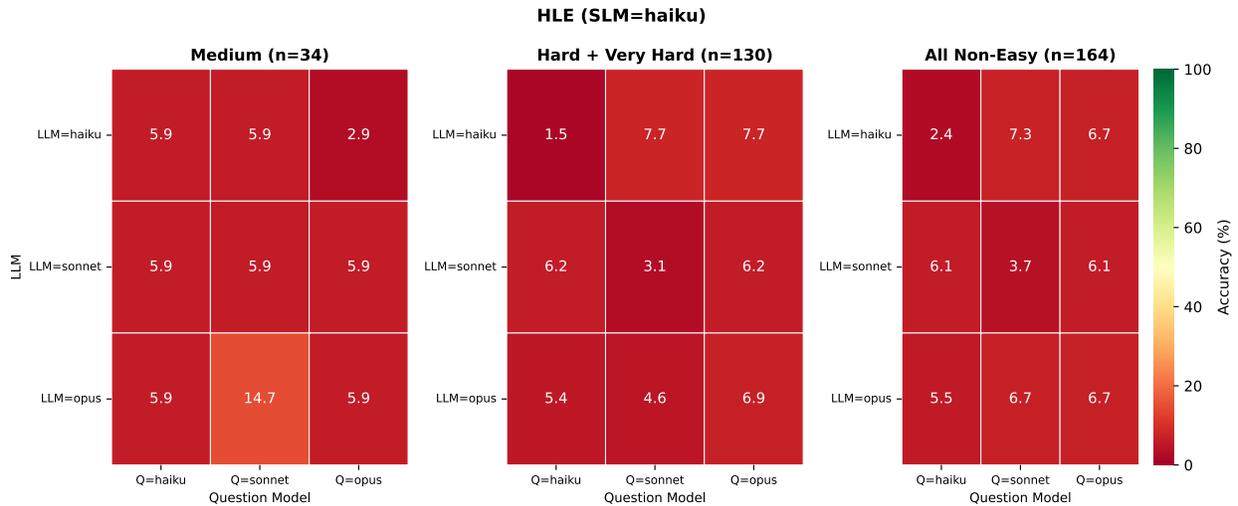


Figure 24: HLE Q&A compression accuracy (SLM=Haiku, Claude 3.5/4). Q&A shows improvement on Hard problems (1.2% to 6.2%).

G.7 Easy-Problem Regression Analysis

To check whether the QA protocol causes regression on problems Haiku already solves, we run the protocol on *easy* problems (those where all three models (Haiku, Sonnet, Opus) answer correctly at the baseline). Since all problems start correct, any change in accuracy is purely regression.

Table 26 reports the regression rate for three conditions: (1) **Haiku re-eval**, where Haiku simply re-solves each problem with no protocol (establishing how often Haiku fails to reproduce its own baseline answer); (2) **BL-CoT**, the self-refinement baseline (SLM=Haiku, LLM=Haiku); and (3) **QA**, the full protocol (SLM=Haiku, LLM=Opus). Both BL-CoT and QA use the iterative protocol with objective judge and early stopping.

Table 26: Easy-problem regression rates (%) under different conditions. All problems start correct (100% initial accuracy). “Regression” is the percentage of problems that become incorrect after re-evaluation. Lower is better. †For AIME and HLE, the high regression rates are primarily due to Haiku’s inconsistency on these problems (see Haiku re-eval column), not the QA protocol itself.

Dataset	n	Haiku Re-eval (%)	Iterative QA Protocol Regression (%)	
			BL-CoT	QA
GSM8K	150	—	0.7	0.0
MATH (Algebra)	150	—	8.7	7.3
MATH (Geometry)	150	—	17.3	11.3
MATH (Num. Theory)	150	—	8.0	4.7
GPQA (MC)	46	—	10.9	10.9
MBPP	150	—	1.3	0.7
AIME†	39	59.0	74.4	69.2
HLE†	30	83.3	80.0	83.3

On standard benchmarks (GSM8K, MATH, GPQA, MBPP), the QA protocol causes modest regression: 0–17% of easy problems become incorrect, with most datasets below 10%. On these datasets, QA regression is consistently equal to or lower than the BL-CoT self-refinement baseline, indicating that Opus guidance does not introduce additional regression.

On frontier benchmarks (AIME and HLE), regression appears severe (69–83%), but a sanity check reveals that this is largely an artifact of Haiku’s inconsistency rather than the QA protocol. When Haiku simply re-solves these “easy” problems with no protocol at all, 59% of AIME and 83% of HLE problems regress. The QA protocol’s regression rates are comparable to this baseline, confirming that the protocol does not meaningfully increase regression. Rather, the underlying difficulty classification is unreliable for these benchmarks, where a single correct answer at baseline does not guarantee consistent reproducibility.

G.8 Example Q&A Transcripts

The following are example transcripts from the Q&A compression protocol (SLM=Haiku, LLM=Opus, Q=Haiku, Claude 4.5). For each dataset, we selected an arbitrary problem where Haiku initially answered incorrectly but produced a correct answer after 10 binary questions, with the additional constraint that the transcript contains both “Yes” and “No” responses. The “No” answers are typically the most informative, as they correct a flawed assumption or redirect the SLM’s reasoning.

GSM8K – Medium

Problem: John buys milk for 2 dollars, eggs for 3 dollars, light bulbs for 3 dollars, cups for 3 dollars, and roach traps for 4 dollars. If there is a 10% tax on all nonfood items how much did everything cost?

Gold Answer: 16

Q1: Is milk correctly identified as a food item that should not be taxed?	Yes
Q2: Are eggs correctly identified as a food item that should not be taxed?	Yes
Q3: Is the categorization correct that light bulbs are a nonfood item subject to tax?	Yes
Q4: Is the categorization correct that cups are a nonfood item subject to tax?	Yes
Q5: Is the categorization correct that roach traps are a nonfood item subject to tax?	Yes
Q6: Is the sum of all nonfood items before tax correctly calculated as \$10?	Yes
Q7: Is the 10% tax calculation on the \$10 nonfood subtotal correctly computed as \$1?	Yes
Q8: Should the tax be applied to the nonfood items only, rather than to all items?	Yes
Q9: Is the final total of \$16 correct when adding the \$5 food subtotal and the \$11 taxed nonfood subtotal?	No
Q10: Does the solution correctly account for all five items that John purchased?	Yes

Haiku's initial calculation was wrong. Q9 revealed an arithmetic error in the final total, which Haiku corrected in its revised answer.

MATH Geometry – Medium

Problem: Corner A of a rectangular piece of paper of width 8 inches is folded over so that it coincides with point C on the opposite side. If $BC = 5$ inches, find the length in inches of fold l .

Gold Answer: $5\sqrt{5}$

Q1: Is the fold line l the same as the line segment BC ?	No
Q2: When corner A is folded to point C , is the fold line the perpendicular bisector of segment AC ?	Yes
Q3: Given that $BA = BC = 5$ and B is on the top edge while A is the top-right corner, is $ 8 - b = 5$ giving $b = 3$ correct?	Yes
Q4: Is the height of the rectangle $h = 10$ inches?	Yes
Q5: Is point C located at $(0, 6)$ on the left edge, making $AC = 4\sqrt{5}$ inches?	Yes
Q6: Given coordinates $A = (8, 10)$, $B = (3, 10)$, $C = (0, 6)$, is $BC = 5$ inches?	Yes
Q7: If the fold line BC has length 5 inches, is this the answer?	No
Q8: Should we verify that $\vec{BC} \perp \vec{AC}$?	No
Q9: Is there additional info constraining the height h to exactly 10 inches?	No
Q10: Should we verify that folding A along line BC truly maps A onto C ?	Yes

Haiku initially confused the fold line with segment BC . Q1 and Q7 corrected this misconception, leading Haiku to compute the correct fold length.

MBPP – Medium

Problem: Write a function to find the maximum difference between available pairs in the given tuple list.

-
- Q1:** Should the maximum difference be calculated between the largest and smallest values found anywhere across all tuples? **No**
- Q2:** Does “available pairs” mean we should only consider pairs of values within the same tuple? **Yes**
- Q3:** Should the function return 0 when the input list is empty? **Yes**
- Q4:** If a tuple contains only a single value, should that tuple be ignored? **Yes**
- Q5:** Are we guaranteed that all elements in the tuples are numeric values? **Yes**
- Q6:** Should the function handle negative numbers correctly? **Yes**
- Q7:** Does “maximum difference between pairs” require finding the difference for every possible pair combination? **No**
- Q8:** If the tuple list contains empty tuples, should the function still work? **Yes**
- Q9:** Should the function consider the absolute difference rather than signed differences? **Yes**
- Q10:** Is it acceptable to flatten all values from all tuples into a single list? **Yes**

Q1 corrected Haiku’s initial interpretation of “available pairs”; differences should be computed within each tuple, not across all tuples.

AIME – Hard

Problem: Find the number of cubic polynomials $p(x) = x^3 + ax^2 + bx + c$, where $a, b, c \in \{-20, \dots, 20\}$, such that there is a unique integer $m \neq 2$ with $p(m) = p(2)$.

Gold Answer: 738

-
- Q1:** Is the factorization $(m - 2)[m^2 + (2 + a)m + (4 + 2a + b)] = 0$ correct? **Yes**
- Q2:** Should we assume $p(m) = p(2)$ is the complete condition? **Yes**
- Q3:** Is it correct that $q(m) = m^2 + (2 + a)m + (4 + 2a + b) = 0$ must have exactly one integer solution? **Yes**
- Q4:** For discriminant = 0, is a being even necessary and sufficient for the double root to be an integer? **Yes**
- Q5:** Should roughly 12 even values of a lead to valid $b \in [-20, 20]$? **Yes**
- Q6:** Should Case 2 include quadratics with two distinct real roots but only one integer root? **Yes**
- Q7:** Does Case 2 also need to account for quadratics with no real roots? **Yes**
- Q8:** For each valid (a, b) pair, is c free to vary over all 41 values in $\{-20, \dots, 20\}$? **No**
- Q9:** Is the answer $(\# \text{ valid } (a, b) \text{ pairs}) \times 41$? **No**
- Q10:** Should the solution verify a concrete example (e.g., $a = 0, b = -3$)? **Yes**

Q8 and Q9 corrected a critical error: Haiku assumed c was free to vary independently, but the condition $p(m) = p(2)$ constrains c through $p(2) = 8 + 4a + 2b + c$.

HLE – Medium

Problem: Let X be the moduli space of nondegenerate lattices in \mathbb{R}^2 with unit area. Compute $H_1(X, \mathbb{Z})$.

Gold Answer: \mathbb{Z}

Q1: Is X correctly identified as $\mathrm{GL}_2(\mathbb{R})^+/\mathrm{GL}_2(\mathbb{Z})$?	No
Q2: Does the space of unit-area lattices have two connected components ($\det = \pm 1$)?	No
Q3: Does the Iwasawa decomposition $\mathrm{SO}(2) \cdot A \cdot N$ apply, with $A \cdot N$ contractible?	Yes
Q4: Is X homotopy equivalent to $\mathbb{H}/\mathrm{PSL}_2(\mathbb{Z})$?	Yes
Q5: Is $\mathbb{H}/\mathrm{PSL}_2(\mathbb{Z})$ a non-compact space that deformation retracts onto a circle?	Yes
Q6: Are the singular orbifold points at i (order 2) and $\rho = e^{2\pi i/3}$ (order 3) correctly identified?	Yes
Q7: Does non-compactness contribute a generator to H_1 from wrapping around the cusp?	No
Q8: Is $H_1(X, \mathbb{Z}) = \mathbb{Z}$ consistent with $X \simeq S^1$?	Yes
Q9: Should we verify the quotient doesn't introduce torsion or additional generators?	Yes
Q10: Is it necessary to distinguish oriented vs. unoriented lattices?	Yes

Q1 and Q2 rejected Haiku's initial identification of the moduli space, and Q7 corrected a wrong claim about the cusp. After these corrections, Haiku correctly identified $X \simeq \mathbb{H}/\mathrm{PSL}_2(\mathbb{Z}) \simeq S^1$ and computed $H_1 = \mathbb{Z}$.