# Soft Contamination Means Benchmarks Test Shallow Generalization

**Ari Spiesberger** [* 1]   **Juan J Vazquez** [* 1]   **Nicky Pochinkov** [1]   **Tomáš Gavenčiak** [2]
**Peli Grietzer** [1]   **Gavin Leech** [† 1 3]   **Nandi Schoots** [† 4]

## Abstract

If LLM training data is polluted with benchmark test data, then benchmark performance gives biased estimates of out-of-distribution (OOD) generalization. Typical 'decontamination' filters use $n$-gram matching which fail to detect 'semantic' duplicates: sentences with equivalent (or near-equivalent) content that are not close in string space. We study this 'soft' contamination of training data by semantic duplicates. Among other experiments, we embed the Olmo3 training corpus and find that: 1) contamination remains widespread, e.g. we find semantic duplicates for 78% of CodeForces and *exact* duplicates for 50% of ZebraLogic problems; 2) including semantic duplicates of benchmark data in training does improve benchmark performance; and 3) when finetuning on duplicates of benchmark datapoints, performance also improves on truly-held-out datapoints from the same benchmark. We argue that recent benchmark gains are thus confounded: the prevalence of soft contamination means gains reflect both genuine capability improvements and the accumulation of test data and *effective* test data in growing training corpora.

## 1. Introduction

LLM scores on hard reasoning (incl. coding) benchmarks have been growing rapidly, with many benchmarks nearing saturation even for smaller, open-source models (Edelman & Lee, 2025; Maslej et al., 2025). Does this trend purely reflect growth in LLMs' general, OOD reasoning capability, or does it also reflect limitations of the benchmarking procedure? We address this question by combining existing data-contamination detection methods with novel finetuning experiments to diagnose what we call *shallow generalization* on benchmarks: benchmark-specific performance gains

from training on datapoints that are qualitatively typical of the benchmark. Using the open-data model Olmo3 as a case study, we show that modern LLM training corpora include data that qualitatively function like samples from major reasoning benchmarks, leading to benchmark scores that *to some extent* demonstrate shallow generalization rather than general reasoning capability. We thus hypothesize that the rapid increase in LLMs' performance on reasoning benchmarks partly reflects the rapid growth of LLMs' corpora size and downstream shallow generalization that tunes models to individual benchmarks via sample-like data. If true, then recent progress on major reasoning benchmarks is weaker evidence for the true pace of AI progress (conceived as OOD generalization).

LLM training corpora have grown by a factor of at least 10,000x since 2020 (Epoch AI, 2025). Thus, there is reason to expect more benchmark test-examples to be included in the training corpora of recent LLMs. While AI labs make good-faith efforts to remove syntactic duplicates of benchmark items from their corpora (OpenAI et al., 2024; Olmo et al., 2025; Anthropic, 2025c), 'softer' forms of contamination are extremely hard to detect, and may well be the product of parallel evolution rather than the product of a data leak. Nevertheless, the presence of benchmark-convergent data in LLMs' training corpora can act as a major confounder with regard to the *type* of generalization evidenced by benchmark scores. Reasoning (incl. coding) benchmark scores, in particular, are typically intended not as measures of LLMs' within-distribution generalization capabilities (comparable to generalization from one subset of the benchmark to another subset), but as tests of the application of fundamental capabilities.

To estimate the significance of shallow generalization as a confounder in benchmark results, we gauge the prevalence of exact and semantic duplicates of items from major reasoning benchmarks in the training corpus of Olmo3. We then conduct finetuning experiments with exact duplicates, semantic duplicates, and close embedding neighbors to test their capacity to induce shallow generalization on a target reasoning-benchmark. While our experiments distinguish between different kinds of 'shallow generalization' gains – gains on a benchmark item from training on its semantic duplicates (Yang et al. (2023); Riddell et al. (2024)); gains on

---

[*]Equal contribution †Joint supervision. [1]Arb Research [2] Charles University, Prague [3] University of Cambridge [4] University of Oxford. Correspondence to: Ari Spiesberger <arispiesberger@gmail.com>, Gavin Leech <gavin@arbresearch.com>.

*Table 1.* Comparison to past literature on data contamination, focusing on semantic-duplicates studies. We compare whether past work studies *semantic* duplicates of test data; estimates contamination prevalence in training corpora; has methods to automatically screen semantic overlap; quantifies the effects of contamination on downstream performance; employs finetuning and whether the finetuning data are ecologically realistic; tests *in-benchmark generalization* (training on duplicates of some benchmark items improves performance on *other* held-out items from the same benchmark); and the scale of data, models, and benchmarks studied.

| Paper | Semantic dupes? | Contam. est.? | Auto detect.? | Effect est.? | With FT? | FT data comp. realistic? | In-bench gen.? | Data scale | Open data attrib.? | Model scale | Major evals |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Ours | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | Large | ✓ | 7B | CodeForces (+ MBPP, MuSR, ZebraLogic) |
| Magar & Schwartz (2022) | ✗ | ✗ | - | ✓ | ✓ | ✗ | ✗ | Small | ✓ | 0.1–0.3B | SST-5, SST-2, SNLI |
| Yang et al. (2023) | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | Medium | ✓ | 13B | MMLU, GSM8K, HumanEval |
| Riddell et al. (2024) | ✓ | ✓ | ✓ | ✓ | ✗ | - | ✗ | Medium | ✓ | ≤16B | MBPP, HumanEval |
| Shilov et al. (2025) | ✗ | ✗ | - | ✓ | ✓ | ✗ | ✗ | Small | ✗ | ≤2.7B | Canary memorization (MIA) |
| Xu et al. (2025) | ✓ | ✗ | ✓ | ✓ | ✗ | - | ✗ | Varied | ✗ | 0.5–72B | LIAR2 |

a benchmark item from training on exact duplicates of other items in the benchmark; gains on a benchmark item from training on semantic duplicates of other items in the benchmark – our discussion frames them as a unified phenomenon from the viewpoint of AI-progress benchmarking: effects of corpus growth that don't reduce to test-memorization but fall short of the capability-growth that benchmarks are designed to measure.

Some terminology: An *exact* duplicate of test data is an example in the training corpus which is syntactically identical (perhaps up to some number of $n$-grams) to some item in a relevant test set. A *semantic* duplicate of test data is an example in the training corpus which has the same meaning (in some sense) as some item in a relevant test set (Riddell et al., 2024). We call contamination *soft* when it involves semantic duplicates. We call generalization *shallow* when it's limited to a combination of within-distribution generalization and generalization across semantic duplicates.

*Our contributions*:

- **Large rates of contamination:** We screen 1% of the pretraining data and all of the finetuning data of Olmo3 for semantic duplicates 'in the wild' by using their embedding distance to benchmark data, which is far more than previous studies have investigated. Despite decontamination efforts in the data preparation of Olmo3, we find much more contamination than previous studies found, likely because we investigate more data;

- **Shallow generalization:** We finetune Olmo3 on duplicates of a subset of the (MuSR, ZebraLogic and MBPP) benchmark and find that benchmark performance also improves on unseen benchmark data. For some benchmarks we find that finetuning on semantic duplicates has the same effect size as finetuning on exact duplicates (an increase of 20% on both seen and unseen items), while finetuning on close embedding neighbors has no effect. Finetuning on one benchmark does not typically improve performance on related benchmarks, suggesting that the generalization in our finetuning experiments is strictly 'shallow';

- **Ecologically valid amounts of contamination have a substantial effect:** We use our findings on the rate of semantic duplicates in the wild to design a finetuning experiment with an ecologically valid mix of benchmark-semantic-duplicate datapoints and clean datapoints in the finetuning corpus, and find that finetuning substantially improves benchmark performance (by around 15%).

## 2. Related Work

Contamination of LLM training corpora by benchmark test items is a well-trodden topic. Table 1 summarizes the key differences between our work and prior studies of benchmark contamination.

The first wave of data-contamination studies (Magar & Schwartz, 2022; Elazar et al., 2024; Jiang et al., 2024) focused on the prevalence and impact of exact syntactic duplicates (i.e. word-for-word matches in string space), with later work (Zhou et al., 2025; Shilov et al., 2025) extending the analysis to partial syntactic duplicates. More recently, researchers have begun to study indirect or 'semantic' contamination (Yang et al., 2023; Riddell et al., 2024; Xu et al., 2025), where data in the training corpora is equivalent to benchmark-items in substantive content without sharing n-gram sequences or other syntactical properties. Our work draws on techniques and best-practices from the semantic contamination literature, but aims at a partially different end: The literature on semantic contamination focuses on the relationship between performance on a benchmark item and training-exposure to that same item's semantic duplicates, which it treats as a variant on memorization. Our work instead studies semantic duplicates as a source of *shallow generalization* phenomena, which include generalization from semantic duplicates to their benchmark-item equivalents but more importantly include within-benchmark-like generalization to non-duplicated items in the benchmark.

Studies of the prevalence of corpora-contamination by exact duplicates of benchmark data typically deploys two kinds of methods: searching for benchmark-item duplicates in open

datasets (Elazar et al., 2024; Riddell et al., 2024; Zhou et al., 2025), and memorization testing using 'membership inference' style techniques (Shi et al., 2023). When studying semantic-duplicates contamination, by contrast, memorization diagnostics are unlikely to capture the right contamination effects. Search in open datasets is therefore preferred in the (small) literature, using a heuristic semantic distance to guide search and human judgment, AI-assistant judgment, or plagiarism-detection software to assign 'semantic duplicate' status. Previous work by Yang et al. (2023) and Riddell et al. (2024) has provided high-quality estimates of the prevalence of semantic duplicates of items from major reasoning benchmarks including HumanEval, MMLU, and GSK8k (Chen et al., 2021; Hendrycks et al., 2020; Cobbe et al., 2021), in widely-used training corpora such as The Pile, StarCoderData, and RedPajama (Gao et al., 2020; Li et al., 2023; Weber et al., 2024).

We use a method convergent with that of Yang et al. (2023) to estimate the prevalence of semantic duplicates in Dolma (Soldaini et al., 2024), Olmo's custom training corpus. While our search covers a much larger dataset and finds many more semantic duplicates per benchmark item, our results are consistent with their findings in terms of the *density* of semantic duplicates in LLMs' training corpora. (Note that because AI labs make ongoing decontamination efforts informed by the scientific literature, the prevalence of contaminating data of any type in SOTA training corpora is a potentially moving target and cannot be automatically assumed from older work.)

Studies attempting to estimate the *effect* of data-contamination (whether semantic or exact) on the integrity of benchmark scores have, to our knowledge, almost exclusively focused on memorization and memorization-like 'item-to-item' effects (one exception is Xu et al. (2025), which studies fake-news detection and employs a domain-specific concept of entity-exposure). Two central methods in the literature are finetuning on contaminated data to simulate training-exposure (assuming or verifying that the model had no or limited prior exposure), which Yang et al. (2023) applies to semantic duplicates, and using duplicate-prevalence data to test for correlations between a benchmark-item's rate of duplication in a model's training corpus and the model's performance on the item, which Riddell et al. (2024) applies to semantic duplicates. Our work uses a finetuning approach, but tests not only gains on a benchmark item from finetuning on its own semantic duplicates, but also gains on benchmark items from finetuning on semantic and on exact duplicates of *other* items in the benchmark. Inspired by Kocyigit et al. (2025)'s study of the memorization-effects caused by injecting realistic dosages of exact duplicates into a clean finetuning corpus, we also design the first (to our knowledge) ecologically valid finetuning study of the effect of training-exposure to semantic duplicates.

# 3. Methodology

In order to know what data the studied model has seen and allow for an exhaustive scientific analysis, our experiments use `Olmo-3-7b` (Olmo et al., 2025) a fully open-source (in particular, open-data) model. In addition, while some closed models allow for finetuning, with a closed corpus we cannot rule out there having been already trained on the examples (or neighbors of the examples) we finetune on, which would confound our effect estimates[1].

## 3.1. Benchmarks

We select benchmarks based on: 1) the ability to generate new synthetic samples using an existing pipeline, 2) the tractability of creating semantic duplicates by modifying original samples, and 3) the likelihood of encountering semantic duplicates of benchmark samples in the wild. We also prioritize benchmarks on which Olmo3 is not saturated, making it easier to track performance improvement or degradation during finetuning.

**MBPP (Austin et al., 2021)**. Dataset of programming questions and validated solutions in Python, with test cases to check correctness. Since they are solvable by entry-level programmers, it is simple to generate correct alternative python solutions or translations in other programming languages. We expected that the training corpora would contain semantic duplicates of MBPP test data. Unlike for other benchmarks, *LLM-assisted* annotation to detect such duplicates may be feasible without requiring very complex reasoning, but would still be undetectable by typical deduplication methods. We use the sanitized test set which contains 257 tasks.

**CodeForces (Penedo et al., 2025)**. Dataset of competitive programming tasks with larger text inputs and problem context than the average code benchmark, such as MBPP. We expected equivalent tasks with less context to appear in training corpora which would be difficult to identify without manual (or LLM) annotation. We use the 468 problems in the default test set in our experiments.

**MuSR (Sprague et al., 2024)**. Dataset used to evaluate multi-step reasoning involving long narratives generated from a ground-truth logic tree built algorithmically. These trees are provided for each original sample by the authors, which allows the creation of semantic duplicates by regenerating story context from the same tree, and slightly different narratives by modifying non-critical tree branches. We focus on the 'murder mysteries' task with 250 problems.

**ZebraLogic (Lin et al., 2025)**. Dataset of 1000 logic grid puzzles of varying complexity levels used to evaluate log-

---

[1]Code at `https://github.com/AriSpiesberger/Soft-Contamination-Prevelance`

3

ical reasoning capabilities. We chose ZebraLogic because Olmo3 Instruct is not saturated on it (having 32.9% accuracy). Additionally, as it is a non-coding benchmark it increases the variety of our suite of benchmarks.

### 3.2. Finding Semantic Duplicates in the Wild

**Olmo3 Corpus Data**. We embed 1% of the Olmo3 Base training data (Dolma3 and Dolmino) and all of the Olmo3 Instruct finetuning data (Dolci SFT, Dolci Instruct DPO and Dolci Instruct RL). To sample from the base training data, we employ a stratified reservoir sampling strategy that preserves the corpus's hierarchical structure (e.g., `common_crawl/art-and-science`). As data is ingested, it is parsed into chunks and routed into distinct reservoirs corresponding to each sub-source. We then construct the final dataset by drawing from these reservoirs in exact proportion to their original volume, ensuring the sample's distribution remains consistent with the full corpus topology. See Appendix A.1 for more details on these datasets.

**Embeddings**. We used `llama-embed-nemotron-8b` (Babakhin et al., 2025) to embed the above datasets. At time of writing this model is number 2 on the Massive Text Embedding Benchmark (MTEB) leaderboard (Muennighoff et al., 2023). All embeddings were done in FP16 precision.

**Cosine Similarity**. We embed both Olmo3 corpus data and benchmark data, and calculate the cosine similarity between data points. The benchmark data comparisons consist primarily of the MBPP and CodeForces datasets. When comparing with pretraining data sets, we split MBPP into inputs and outputs. When comparing MBPP with instruct data, we join MBPP inputs and outputs to match the prompt response format in the instruct data. In Appendix A.2 we also present cosine similarity matches between corpus data and MuSR and ZebraLogic. Our comparison consists of taking the cosine similarity of the benchmark point embedding with the corpus text embeddings.

**Investigating Semantic and Exact Duplicate Status of High Cosine Similarity Matches**

For MuSR the highest cosine similarity matches are around 0.40, and upon manual inspection it appears there are no duplicates, or MuSR-like problems in the training corpora. We conclude that the Olmo3 datasets have no semantic duplicates of MuSR. In the case of ZebraLogic we found many exact duplicates and a few semantic duplicates. Additionally to some semantic duplicates, most top cosine similarity matches were Einstein riddles and Zebra puzzles available online, but do not match in complexity or sample integrity to the dataset samples. After manual inspection of the top matches for MBPP and CodeForces we decided to sample, and annotate them using an LLM.

*Exact Duplicates: Sampling and Annotating ZebraLogic.*

We observed that several ZebraLogic tasks were present in the training corpora verbatim. So we ran an annotation experiment to get a rate of exact duplicates. This consisted in checking the 10 highest cosine similarity matches for each test point (adding up to 10,000 annotations for the entire benchmark).

*Semantic Duplicates: Sampling and Annotating MBPP and CodeForces.* We investigate matches between MBPP and CodeForces for each of the five training datasets. Per dataset and benchmark, we select the points in the top 0.1% similarity. From these we either take the top 100, or randomly sample 100 points, and use a diffused model of `gemini-3-flash-preview` (Google Deep-Mind, 2025) to categorize whether the two texts are semantic duplicates or not. Obtaining labels for whether the corpus text is a semantic duplicate or not, the type of semantic duplicate (exact, equivalent, subset, superset), reasoning of the choice, and the confidence of the label. See more details on the annotation pipeline in Appendix A.4.

### 3.3. Generating Synthetic Semantic Duplicates

In Appendix A.3 we provide an overview and a detailed description of our methods for generating synthetic semantic duplicates for each benchmark.

### 3.4. Finetuning on Duplicates

We finetune Olmo3 Instruct on duplicates of the following benchmark datasets: MuSR, Zebralogic, and MBPP. We use either exact duplicates or semantic duplicates generated as in Section 3.3. To finetune Olmo3 Instruct (a non-reasoning model) on these datapoints we first get a teacher model to generate Chain-of-Thought (CoT) (Wei et al., 2022) reasoning traces. We use Opus 4.5 as teacher model, and for MuSR we also experiment with using GPT 4.1-mini.

We take the formatted questions and corresponding CoT answers and use LoRA (Hu et al., 2022) to finetune Olmo3 Instruct. To get propensities we use a temperature of 0.7 and do 8 parallel generations (for each of the unfinetuned model, the CoT generations of the teacher model, and for the finetuned model).

We split a finetuning dataset of duplicates in two and only finetune on half of it, while evaluating the finetuned model on both seen and unseen duplicates. To assess whether performance goes up on related benchmarks we use TrueDetective (Del & Fishel, 2023) as a mirror for MuSR, Arc Challenge (Clark et al., 2018) as a mirror for ZebraLogic, and HumanEval as a mirror for MBPP. We also evaluate performance on Arc Easy, BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020) and Winogrande (Sakaguchi et al., 2021).

# 4. Results

## 4.1. Exact duplicates in training corpora

Olmo et al. (2025), the paper introducing Olmo3, lists ZebraLogic in their suite of evaluation benchmarks. Surprisingly however, we find that the Olmo Instruct RL dataset contains exact duplicates of ZebraLogic problems and solutions. In Figure 1 we see that for puzzles of $4 \times 4$ or larger, for ~70% or more problems, there exists at least one exact duplicate in the Olmo3 Instruct RL dataset (`Dolci-Instruct-RL`). In total, we find corpus duplicates for at least 49.5% of dataset tasks.

The model card of Olmo3 benchmarks the performance of Olmo3 Instruct after SFT training as 18%, after DPO training as 28.4%, and after RL training as 32.9%. The model is improved in that order. We do not find exact or semantic duplicates of ZebraLogic in the DPO dataset, so we ascribe the improvement between SFT training and DPO training to general logical reasoning improvement. However, the increase from 28.4% to 32.9% (possibly on harder problems) is likely due to directly training on ZebraLogic data.
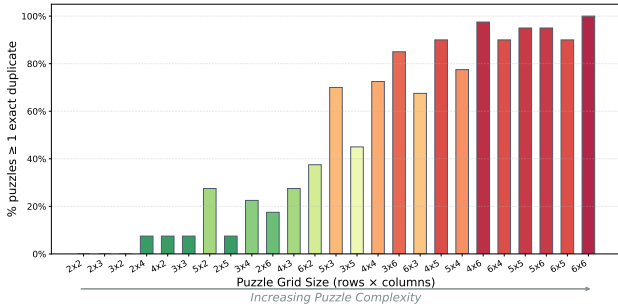


*Figure 1.* On the y-axis we plot the following statistic: for each ZebraLogic benchmark datapoint we check among the top 10 highest cosine similarity training datapoints if any of those samples is an exact duplicate, we then calculate the proportion of benchmark datapoints (of a given grid size) that have at least one exact duplicate. On the x-axis we plot puzzle grid size.

## 4.2. Natural semantic duplicates in training corpora

**Relationship of cosine similarity to semantic duplicate status.** In Figure 2 we plot cosine similarity against semantic duplicate status. To acquire this plot we select the points in the top 0.1% cosine similarity (in matches between benchmark and training dataset), and randomly sample 100 points. We then use a language model to assess the semantic duplicate status of these 100 datapoints. For both MBPP and CodeForces we find that even within the top 0.1% highest cosine similarity matches, semantic duplicates are far more common among the highest cosine similarity matches, see Figure 2. For MBPP we do not find exact duplicates, which
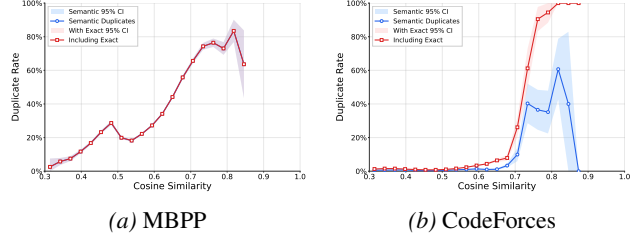


*(a)* MBPP     *(b)* CodeForces

*Figure 2.* Relationship between cosine similarity level and semantic duplication. For each benchmark datapoint we sample 100 matches from the top 0.1% cosine similarity matches in the training data. On the x-axis we plot the cosine similarity. On the y-axis we plot the percentage of cosine similarity matches at this level that are true semantic duplicates. The opaque graph shows the confidence interval: this interval widens when there are fewer samples of a given cosine similarity level. In red we plot semantic duplicates inclusive of exact duplicates, and in blue exclusive.



*Figure 3.* Occurence by elo. On the y-axis we plot the following statistic: for each benchmark datapoint we check among the top 100 cosine similarity training datapoints if any of those samples is a semantic duplicate, we then calculate the proportion of all benchmark datapoints that have at least one semantic duplicate. We plot Elo scores on the x-axis.

is why the two graphs overlay perfectly, whereas for CodeForces they come apart. Below we focus on investigating the top 100 cosine similarity matches for each benchmark datapoint.

**Proportion of benchmark datapoints that have at least one semantic duplicate.** We find that 100% of MBPP problems have at least one semantic duplicate in the top 100 cosine similarity training data matches. We also find at least one semantic duplicate per benchmark datapoint when we randomly sample 100 matches out of the top 0.1% cosine similarity matches. For CodeForces we find that 77.5% of problems have at least one semantic duplicate in the top 100 cosine similarity matches. Figure 3 considers the likelihood of a test point having a single semantic duplicate based on its elo, and find that problem difficulty (elo score) does not play a big role. We believe the above numbers are a lower bound: if we checked all the data and not just the top 100 cosine similarity matches, we would likely find more semantic duplicates for CodeForces.

**Semantic duplicate occurence stratified by training**

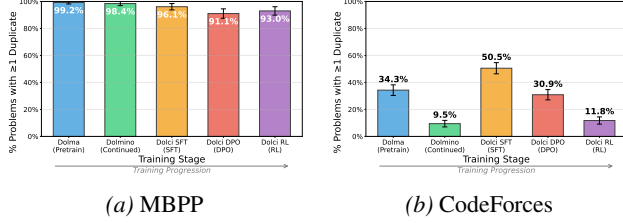*(a)* MBPP        *(b)* CodeForces

*Figure 4.* Occurence by training dataset. On the y-axis: for each benchmark datapoint we check among the top 100 cosine similarity training datapoints if any of those samples is a semantic duplicate, we then calculate the proportion of all benchmark datapoints that have at least one semantic duplicate. On the x-axis we plot the different training datasets. The lines show the standard deviation.
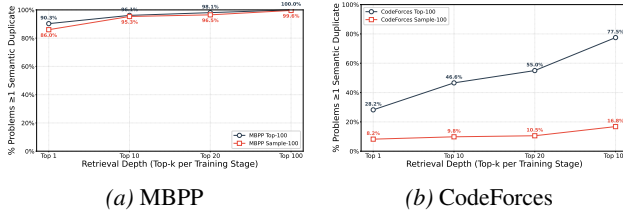


*(a)* MBPP        *(b)* CodeForces

*Figure 5.* Occurence by number of cosine similarity matches investigated. We take the number of semantic duplicates evaluated by being at top-n at each of our dataset comparisons.

**dataset.** In Figure 4 we investigate the relationship between where in the training scheme a training dataset is used, and the extent of semantic duplicate contamination. In Figure 4a we see that all training datasets have semantic duplicates for MBPP. In Figure 4b, for CodeForces we observe that again, semantic duplicates of problems exist in each dataset, particularly in Dolma, Dolci SFT and Dolci DPO.

**Semantic duplicates are sparse and investigating more data leads to more matches.** Previous work has found that n-grams typically miss many semantic duplicates (Yang et al., 2023), which is why like Yang et al. (2023) we instead work with cosine similarity matches. In Figure 6 we see that high cosine similarity matches are sparse. We also see in Figure 2 that a substantial portion of the very high (over 0.8) cosine similarity matches is a semantic duplicate. We also found that for MBPP and CodeForces respectively 100% and 77.5% of problems have at least one semantic duplicate in their top 100 cosine similarity matches. In Figure 5 we see that this statistic - the proportion of benchmark problems that have at least one semantic duplicate in the training data - scales with the number of datapoints we sample. For example, for CodeForces this statistic drops to 28.4% if we only sample the single top cosine similarity match. We suggest that the reason we found more semantic duplicates than previous work is that we investigated far more data. Both embedding very large amounts of training data and annotating cosine similarity matches with semantic duplicate status are computationally expensive, so scaling up investigations

of this kind is challenging. See Appendix B.2.

### 4.3. Finetuning on Semantic Duplicates

**MuSR.** We experimented with three levels of sophistication for semantic duplication. In Table 2 we find that when we finetune on duplicates of half of the benchmark, performance goes up equally on the unseen half of the benchmark. We find that finetuning on exact duplicates of MuSR benchmark data leads to a similar increase in performance as finetuning on a variety of types of semantic duplicates. In both cases the performance goes up by about 20%. When instead of finetuning on duplicates (exact or semantic), we finetune on datapoints that have been selected for high cosine similarity to the benchmark datapoints, the performance hardly goes up from baseline. We also check performance change on a same domain but different benchmark, TrueDetective, and find that performance remains stable. In Appendix C we find that when we use a better teacher model to generate CoT reasoning traces, Olmo3 does better after finetuning on them.

**ZebraLogic.** In Table 3, when finetuning on exact duplicates of half of the benchmark data, performance also goes up on the other half. For ZebraLogic we find that exact duplicates lead to a much larger jump in performance (for both seen and unseen benchmark items) than finetuning on semantic duplicates, in contrast to our finding for MuSR. In fact, we find that finetuning on semantic duplicates hardly increases performance, and in one case (combining shuffling, substituting and paraphrasing), that the performance on ZebraLogic substantially degrades, even though performance on general datasets remains stable, see Appendix C. We investigate performance change on a same domain but different benchmark, Arc Challenge, and find that performance is unaffected by finetuning.

**MBPP.** In Table 4, when finetuning on exact duplicates of half of the benchmark data, we see a substantial jump on that half of the data, while performance on unseen benchmark data hardly changes. Finetuning on semantic duplicates has a more moderate effect, but affects both seen and unseen performance. Again, finetuning on cosine similar data does not improve performance substantially. We also evaluate on a same domain but different coding benchmark, HumanEval, and find a surprising jump for semantic duplicates. We hypothesize that this jump happened because our semantic duplicate dataset is fairly rich and high quality, and demands some generalization. We do want to note that our MBPP results are a little noisy: when we stratify the baseline and cosine similarity sft performance evaluation by the first half of the benchmark data ('seen') and second half ('unseen') we find a discrepancy of 6% (in opposite directions: baseline scores higher on second half, and the cosine similarity model scores higher on the first), when you would expect
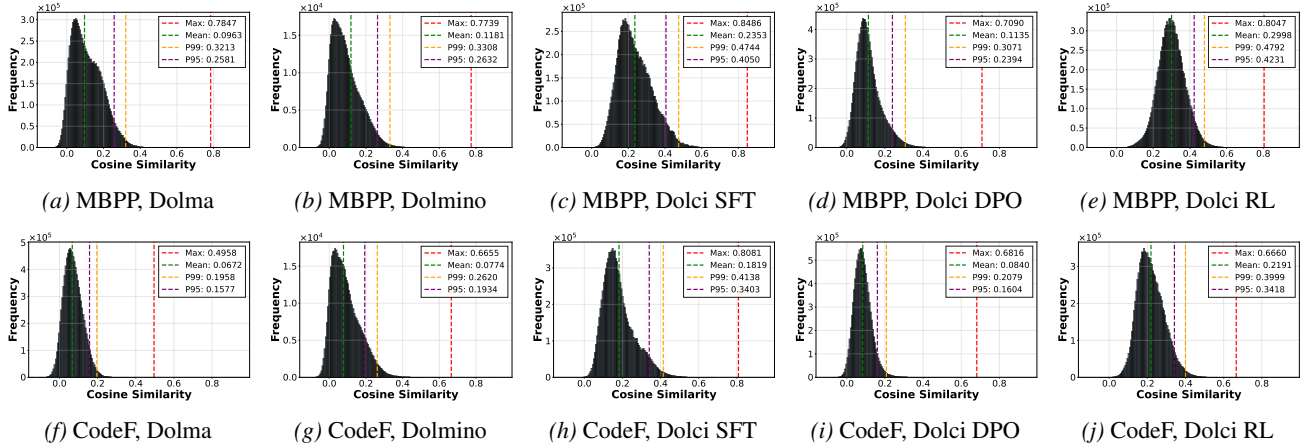
*Figure 6.* Each plot shows cosine similarity distribution of pairs of benchmark data and training corpus data. The top row shows distributions for MBPP and the bottom for CodeForces. From left to right we plot Dolma, Dolmino, Dolci SFT, Dolci DPO and Dolci RL.

*Table 2.* MuSR

| Duplication level | Seen | Unseen | True Detective |
|---|---|---|---|
| Baseline | | 66.0 | 28.3 |
| Exact Dupes | 87.9 | 87.3 | 27.7 |
| Level 1 | 85.8 | 86.2 | 29.3 |
| Level 2 | 85.7 | 86.0 | 28.3 |
| Level 3 | 87.5 | 87.9 | 29.8 |
| Cos Sim sft | | 68.6 | 25.1 |
| Cos Sim dpo | | 67.9 | 26.7 |
| Cos Sim rl | | 65.3 | 26.7 |

*Table 3.* ZebraLogic

| Duplication level | Seen | Unseen | Arc Challenge |
|---|---|---|---|
| Baseline | | 36.9 | 50.1 |
| Exact Dupes | 48.4 | 43.4 | 49.5 |
| Para | 39.2 | 36.2 | 49.3 |
| Shuffle, subs | 36.0 | 36.8 | 50.7 |
| Shuffle, para | 38.0 | 36.0 | 49.4 |
| Shuffle, subs, para | 28.0 | 28.4 | 50.4 |
| Cos Sim sft | | 22.9 | - |

*Table 4.* MBPP

| Duplication level | Seen | Unseen | HumanEval |
|---|---|---|---|
| Baseline | | 46.4 | 55.3 |
| Exact Dupes | 63.0 | 48.8 | 49.2 |
| Semantic Dupes (Py) | 55.1 | 53.6 | 67.0 |
| Cos Sim sft | | 48.8 | 53.1 |

performance on these two halves to be similar.

**We observe a pattern of shallow generalization.** We repeatedly find that when finetuning on duplicates of benchmark data has a substantial effect on benchmark performance, then performance also improves on benchmark data that was unseen during finetuning. This suggests within-benchmark-distribution generalization. We tested benchmark improvement on different, but same domain benchmarks, and typically did not find substantial improvement, confirming shallow generalization. We also find that improvement or finetuning on high cosine similar datapoints does not by itself improve benchmark performance.

### 4.4. Ecologically Valid Finetuning

**Ecologically valid contamination amount.** We evaluate the impact of semantic duplicate contamination under realistic model-developer conditions. We used an ecologically valid amount of contamination, determined as follows: For MBPP we found that, when we randomly sampled 100 matches among the top 0.1% highest cosine similarity training datapoints for a given benchmark datapoint, on average there were 40 semantic duplicates among the 100 samples. We concluded that roughly 4 in 10,000 training datapoints are a semantic duplicate for a given benchmark datapoint.

**Finetuning contaminated and clean models.** We finetune Olmo3 on data containing MuSR semantic duplicates, for which our annotation experiments verify that no duplicates exist in the model's training data. We then split the MuSR data into two halves of 125 datapoints each, and generate 4 semantic duplicates for them, two duplicates of level 2 and two of level 3, so 500 duplicates in total. We perform two finetuning runs with 1) a clean dataset of 10,000 SFT datapoints verified to be decontaminated, and 2) a contaminated version of the same dataset where 5% of clean samples are

*Table 5.* We report on baseline (before finetuning) accuracy on MuSR. We then finetune on 10.000 datapoints. We either finetune on half of the level 2 & 3 semantic duplicates mixed in with regular data (contaminated model) or we finetune on clean data only (clean model).

| MODEL | TREATMENT | SEEN | UNSEEN | TRUE DETECTIVE |
|---|---|---|---|---|
| OLMO3 | BASELINE | | 42.8 | 29.3 |
| | FINETUNED CLEAN | | 50.0 | 28.0 |
| | FINETUNED CONTAMINATED | 66.4 | 54.4 | 28.0 |
| QWEN3 | BASELINE | | 40.4 | 24.0 |
| | FINETUNED CLEAN | | 53.6 | 24.0 |
| | FINETUNED CONTAMINATED | 65.6 | 52.0 | 28.0 |

swapped with 500 semantic duplicates corresponding to the 'seen' subset. We evaluate both finetuned models on the full MuSR benchmark, splitting results to 'seen' and 'unseen' MuSR for the contaminated but not for the clean model.

**Results on Olmo3.** We find that while the contamination percentage is very low, the contaminated model score on the seen subset is 12% higher than the clean model's benchmark score, and the contaminated model score on the unseen subset is 5.6% higher that the clean model's benchmark score. To verify that performance gains on benchmark items reflects 'shallow' generalization rather robust capability gains, we evaluate on TrueDetective, a benchmark in the same domain as MuSR. We find that performance on TrueDetective remains stable during finetuning.

**Noisy results on Qwen3.** Replicating the experiment on Qwen3-8B-base yielded noisy results, possibly due to training instability during finetuning. The contaminated model scored 13.6% higher on seen samples versus unseen, similar to Olmo3. The clean model unexpectedly also showed substantive MuSR benchmark gains from fine-tuning. We note that, puzzlingly, when breaking down the gains of clean finetuning by subdataset we see that the clean model's gains are particularly strong on the random benchmark-subset we use as the 'unseen' subset for the contaminated model, see Appendix D. For Qwen3 we also see some improvement in performance on the same-domain benchmark TrueDetective for the finetuned contaminated model.

**Ecologically valid contamination leads to benchmark performance improvement.** Our results show that the presence of semantic duplicates in training corpora, even at low rates, can lead to substantial gains in evaluation results. Breaking the gains down by type, the evidence for gains on 'seen' (as semantic duplicates) data from realistic quantities of contamination is strong, while the evidence for within-benchmark generalization from realistic quantities of contamination is mixed. We note that this is the first ecologically valid demonstration of gains on 'seen' data from finetuning on semantic duplicates, which previous work only demonstrated in proof-of-concept experiments. In light of the very strong within-benchmark generaliza-

tion we observed in the more artificial setting of Section 4.3, we believe the topic of within-benchmark generalization in realistic training-setting requires further study. See Appendix B.2.1 for more details.

## 5. Limitations and Future Work

We likely underestimate the prevalence of semantic duplicates in real training corpora, because our detection methods in Section 4.2 are likey to have a high false negative rate. Another downward bias is the relative absence of rephrasings in Dolma: synthetic data pipelines now often involve intentionally creating semantic duplicates (Wei & Zou, 2019; Wang et al., 2022; Maini et al., 2024), and so our estimates of their prevalence are likely lower than the true rate in closed corpora using such methods. Similarly, we do not cover state of the art synthetic data provided via RL environments.

Our experiment design is limited to models which open-source their training corpus - a small and potentially unrepresentative set of systems. For instance, the training corpora used in frontier models are much larger than that of the Olmo models – see e.g. the 30T tokens used in the largest Llama 4 runs, (Meta AI, 2025). These larger corpora will have more semantic duplicates, but also a different rate of natural semantic duplicates than Dolma.

A fundamental objection to our project could be that out-of-distribution (OOD) generalization is no longer the (only) goal of AI development: an alternative is to instead bring all common tasks in-distribution (Chollet 2024, Patel & Sutton 2025, Leech et al. 2024 §5.3.2). One could argue that the real-world utility of LLMs shows that our concerns about OOD generalization are not practically important even if generalization is largely shallow. This is a valid perspective, but 1) then the deviation from the assumptions of empirical risk minimization should be explicitly noted, 2) it's unclear to what extent even perfect hidden interpolation would be practically equivalent to true OOD generalization.

## Impact Statement

We aim here to advance understanding of how LLM benchmark scores relate to general capabilities. Our findings have implications for how the AI research community, policymakers, and the public interpret reported progress on reasoning benchmarks.

More accurately measuring AI capabilities supports calibrated decisions about AI deployment, regulation, and research. If benchmark gains partly reflect interpolation from a growing corpus rather than more general capability improvements, then recognizing this could help prevent overconfidence in model generalization to novel tasks.

We do not believe this work poses significant risks. While our methods could inform more sophisticated benchmark gaming, the contamination we study is likely accidental rather than adversarial, and our detection methods are likely more useful for auditing than for evasion.

Our finding exact and soft contamination in Olmo3 is only possible because of the unusual level of transparency of its model development process. It would be unfair for readers to thereby assume that the level of contamination in Olmo is unusually severe, and worse, a perverse incentive against transparency.

Our work could easily be misread as 'debunking' LLM capabilities and so spur complacency about near-term AI impacts. We emphasize that our results suggest that benchmark gains are confounded (and so partially shallow), not that they are illusory.

## Author Contributions

Spiesberger: experimental design for corpus search, embedding, and ecological finetuning experiments; managed compute infrastructure and code repository; executed embedding, MBPP finetuning, and ecological finetuning experiments; designed figures.

Vazquez: experimental design for generation and validation of semantic duplicates; comparison of duplicate detection methods; designed and executed annotation experiments; data spot-checking; wrote parts of methodology, results, and appendices; designed figures.

Pochinkov: generated MuSR teacher examples; finetuned and evaluated Olmo3 models on MuSR and ZebraLogic; data sanity checks.

Gavenčiak: generated synthetic semantic duplicates; contributed to methodology, annotation pipeline design, finetuning experiments, and data analysis; provided feedback on the manuscript.

Grietzer: wrote the introduction and related work sections;

edited the manuscript.

Leech: original research question; assembled the team; experiment design; wrote abstract, limitations, future work, and impact statement; secured resources and compute for the project.

Schoots: research and project management; led experiment design; coordinated writing efforts and wrote several parts of the manuscript.

## Acknowledgments

## References

Anthropic. System card: Claude haiku 4.5. Model Card, October 2025a. URL https://www-cdn.anthropic.com/7aad69bf12627d42234e01ee7c36305dc2f6a970.pdf.

Anthropic. System card: Claude sonnet 4.5. Model Card, September 2025b. URL https://www-cdn.anthropic.com/963373e433e489a87a10c823c52a0a013e9172dd.pdf.

Anthropic. System card: Claude opus 4.5. https://www.anthropic.com/claude-opus-4-5-system-card, November 2025c. Accessed 2026-01-26.

Austin, J., Odena, A., Nye, M. I., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C. J., Terry, M., Le, Q. V., and Sutton, C. Program synthesis with large language models. *CoRR*, abs/2108.07732, 2021. URL https://arxiv.org/abs/2108.07732.

Babakhin, Y., Osmulski, R., Ak, R., Moreira, G., Xu, M., Schifferer, B., Liu, B., and Oldridge, E. Llama-embed-nemotron-8b: A universal text embedding model for multilingual and cross-lingual tasks. *arXiv preprint arXiv:2511.07025*, 2025.

Bisk, Y., Zellers, R., Gao, J., Choi, Y., et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.

Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman,

G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.

Chollet, F. The question of whether LLMs can reason is, in many ways, the wrong question [Tweet]. X (formerly Twitter), July 2024. URL https://x.com/fchollet/status/1816954290227089656. Accessed: 2026-01-26.

Clark, C., Lee, K., Chang, M.-W., Kwiatkowski, T., Collins, M., and Toutanova, K. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019.

Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

Del, M. and Fishel, M. True detective: A deep abductive reasoning benchmark undoable for gpt-3 and challenging for gpt-4. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (* SEM 2023)*, pp. 314–322, 2023.

Edelman, Y. and Lee, J. Ai capabilities progress has sped up, 2025. URL https://epoch.ai/data-insights/ai-capabilities-progress-has-sped-up. Accessed: 2026-01-26.

Elazar, Y., Bhagia, A., Magnusson, I. H., Ravichander, A., Schwenk, D., Suhr, A., Walsh, E. P., Groeneveld, D., Soldaini, L., Singh, S., Hajishirzi, H., Smith, N. A., and Dodge, J. What's in my big data? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=RvfPnOkPV4.

Epoch AI. Data on ai models, 07 2025. URL https://epoch.ai/data/ai-models. Accessed: 2026-01-26.

Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., Phang, J., He, H., Thite, A., Nabeshima, N., Presser, S., and Leahy, C. The pile: An 800gb dataset of diverse text for language modeling, 2020. URL https://arxiv.org/abs/2101.00027.

Google DeepMind. Gemini 3 flash model card. Model Card, December 2025. URL https://storage.googleapis.com/deepmind-media/Model-Cards/Gemini-3-Flash-Model-Card.pdf.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.

Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.

Jiang, M., Liu, K. Z., Zhong, M., Schaeffer, R., Ouyang, S., Han, J., and Koyejo, S. Investigating data contamination for pre-training language models, 2024. URL https://arxiv.org/abs/2401.06059.

Kocyigit, M. Y., Briakou, E., Deutsch, D., Luo, J., Cherry, C., and Freitag, M. Overestimation in LLM evaluation: A controlled large-scale study on data contamination's impact on machine translation. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 31105–31132. PMLR, 2025. URL https://proceedings.mlr.press/v267/kocyigit25a.html.

Leech, G., Vazquez, J. J., Kupper, N., Yagudin, M., and Aitchison, L. Questionable practices in machine learning, 2024. URL https://arxiv.org/abs/2407.12220.

Li, R., Ben Allal, L., Zi, Y., Muennighoff, N., Kocetkov, D., Mou, C., Marone, M., Akiki, C., Li, J., Chim, J., Liu, Q., Zheltonozhskii, E., Zhuo, T. Y., Wang, T., Dehaene, O., Davaadorj, M., Lamy-Poirier, J., Monteiro, J., Shliazhko, O., Gontier, N., Meade, N., Zebaze, A., Yee, M.-H., Umapathi, L. K., Zhu, J., Lipkin, B., Oblokulov, M., Wang, Z., Murthy, R., Stillerman, J., Patel, S. S., Abulkhanov, D., Zocca, M., Dey, M., Zhang, Z., Fahmy, N., Bhattacharyya, U., Yu, W., Singh, S., Luccioni, S., Villegas, P., Kunakov, M., Zhdanov, F., Romero, M., Lee, T., Timor, N., Ding, J., Schlesinger, C., Schoelkopf, H., Ebert, J., Dao, T., Mishra, M., Gu, A., Robinson, J., Anderson,

C. J., Dolan-Gavitt, B., Contractor, D., Reddy, S., Fried, D., Bahdanau, D., Jernite, Y., Ferrandis, C. M., Hughes, S., Wolf, T., Guha, A., von Werra, L., and de Vries, H. StarCoder: may the source be with you!, 2023. URL https://arxiv.org/abs/2305.06161.

Lin, B. Y., Bras, R. L., Richardson, K., Sabharwal, A., Poovendran, R., Clark, P., and Choi, Y. Zebralogic: On the scaling limits of llms for logical reasoning. In *Forty-second International Conference on Machine Learning, ICML 2025, Vancouver, BC, Canada, July 13-19, 2025*. OpenReview.net, 2025. URL https://openreview.net/forum?id=sTAJ9QyA6l.

Magar, I. and Schwartz, R. Data contamination: From memorization to exploitation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 157–165, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-short.18. URL https://aclanthology.org/2022.acl-short.18/.

Maini, P., Seto, S., Bai, H., Grangier, D., Zhang, Y., and Jaitly, N. Rephrasing the web: A recipe for compute and data-efficient language modeling, 2024. URL https://arxiv.org/abs/2401.16380.

Maslej, N., Fattorini, L., Perrault, R., Gil, Y., Parli, V., Kariuki, N., Capstick, E., Reuel, A., Brynjolfsson, E., Etchemendy, J., Ligett, K., Lyons, T., Manyika, J., Niebles, J. C., Shoham, Y., Wald, R., Walsh, T., Hamrah, A., Santarlasci, L., Lotufo, J. B., Rome, A., Shi, A., and Oak, S. Artificial intelligence index report 2025, 2025. URL https://arxiv.org/abs/2504.07139.

Meta AI. The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation. https://ai.meta.com/blog/llama-4-multimodal-intelligence/, April 2025. Accessed: 2026-01-26.

Muennighoff, N., Tazi, N., Magne, L., and Reimers, N. Mteb: Massive text embedding benchmark. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 2014–2037, 2023.

Olmo, T., :, Ettinger, A., Bertsch, A., Kuehl, B., Graham, D., Heineman, D., Groeneveld, D., Brahman, F., Timbers, F., Ivison, H., Morrison, J., Poznanski, J., Lo, K., Soldaini, L., Jordan, M., Chen, M., Noukhovitch, M., Lambert, N., Walsh, P., Dasigi, P., Berry, R., Malik, S., Shah, S., Geng, S., Arora, S., Gupta, S., Anderson, T., Xiao, T., Murray, T., Romero, T., Graf, V., Asai, A., Bhagia, A., Wettig, A., Liu, A., Rangapur, A., Anastasiades, C., Huang, C., Schwenk, D., Trivedi, H., Magnusson, I., Lochner, J., Liu,

J., Miranda, L. J. V., Sap, M., Morgan, M., Schmitz, M., Guerquin, M., Wilson, M., Huff, R., Bras, R. L., Xin, R., Shao, R., Skjonsberg, S., Shen, S. Z., Li, S. S., Wilde, T., Pyatkin, V., Merrill, W., Chang, Y., Gu, Y., Zeng, Z., Sabharwal, A., Zettlemoyer, L., Koh, P. W., Farhadi, A., Smith, N. A., and Hajishirzi, H. Olmo 3, 2025. URL https://arxiv.org/abs/2512.13961.

OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.-L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L., Kim, J. W., Kim, C., Kim, Y., Kirchner, J. H., Kiros, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S. M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O'Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H. P., Michael, Pokorny, Pokrass, M., Pong, V. H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D.,

Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M. B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., and Zoph, B. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.

Patel, D. and Sutton, R. Richard sutton – father of RL thinks LLMs are a dead end. Dwarkesh Podcast, September 2025. URL https://www.dwarkesh.com/p/richard-sutton. Podcast interview.

Penedo, G., Lozhkov, A., Kydlíček, H., Allal, L. B., Beeching, E., Lajarín, A. P., Gallouédec, Q., Habib, N., Tunstall, L., and von Werra, L. Code-forces. https://huggingface.co/datasets/open-r1/codeforces, 2025.

Riddell, M., Ni, A., and Cohan, A. Quantifying contamination in evaluating code generation capabilities of language models. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14116–14137, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.761. URL https://aclanthology.org/2024.acl-long.761/.

Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.

Shi, W., Ajith, A., Xia, M., Huang, Y., Liu, D., Blevins, T., Chen, D., and Zettlemoyer, L. Detecting pretraining data from large language models, 2023. URL https://arxiv.org/abs/2310.16789.

Shilov, I., Meeus, M., and de Montjoye, Y.-A. The mosaic memory of large language models, 2025. URL https://arxiv.org/abs/2405.15523. arXiv:2405.15523v2.

Soldaini, L., Kinney, R., Bhagia, A., Schwenk, D., Atkinson, D., Authur, R., Bogin, B., Chandu, K., Dumas, J., Elazar, Y., Hofmann, V., Jha, A., Kumar, S., Lucy, L., Lyu, X., Lambert, N., Magnusson, I., Morrison, J., Muennighoff, N., Naik, A., Nam, C., Peters, M., Ravichander, A., Richardson, K., Shen, Z., Strubell, E., Subramani, N., Tafjord, O., Walsh, E., Zettlemoyer, L., Smith, N., Hajishirzi, H., Beltagy, I., Groeneveld, D., Dodge, J., and Lo, K. Dolma: an open corpus of three trillion tokens for language model pretraining research. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15725–15788, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.840. URL https://aclanthology.org/2024.acl-long.840/.

Sprague, Z., Ye, X., Bostrom, K., Chaudhuri, S., and Durrett, G. Musr: Testing the limits of chain-of-thought with multistep soft reasoning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=jenyYQzue1.

Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N. A., Khashabi, D., and Hajishirzi, H. Self-instruct: Aligning language models with self-generated instructions, 2022. URL https://arxiv.org/abs/2212.10560.

Weber, M., Fu, D., Anthony, Q., Oren, Y., Adams, S., Alexandrov, A., Lyu, X., Nguyen, H., Yao, X., Adams, V., Athiwaratkun, B., Chalamala, R., Chen, K., Ryabinin, M., Dao, T., Liang, P., Ré, C., Rish, I., and Zhang, C. Redpajama: an open dataset for training large language models, 2024. URL https://arxiv.org/abs/2411.12372.

Wei, J. and Zou, K. EDA: Easy data augmentation techniques for boosting performance on text classification tasks. In Inui, K., Jiang, J., Ng, V., and Wan, X. (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 6382–6388, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1670. URL https://aclanthology.org/D19-1670/.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Xu, C., Yan, N., Guan, S., Mei, Y., and Kechadi, T. SSA: Semantic contamination of LLM-driven fake news detection. In Christodoulopoulos, C., Chakraborty, T., Rose, C., and Peng, V. (eds.), *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pp. 14737–14751, Suzhou,

China, November 2025. Association for Computational Linguistics. doi: 10.18653/v1/2025.emnlp-main. 744. URL https://aclanthology.org/2025. emnlp-main.744/.

Yang, S., Chiang, W.-L., Zheng, L., Gonzalez, J. E., and Stoica, I. Rethinking benchmark and contamination for language models with rephrased samples, 2023. URL https://arxiv.org/abs/2311.04850.

Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.

Zhou, X., Weyssow, M., Widyasari, R., Zhang, T., He, J., Lyu, Y., Chang, J., Zhang, B., Huang, D., and Lo, D. Lessleak-bench: A first investigation of data leakage in LLMs across 83 software engineering benchmarks, 2025. URL https://arxiv.org/abs/2502.06215.

# A. Further Details on Methodology

## A.1. Olmo3 Instruct Training Datasets

We embedded 1% of Dolma3_6T-mix-1025, the data Olmo3 Base was trained on; and 1% of Dolmino-mix-1025 (100B tokens, 2.7 GB), a high quality dataset Olmo3 Base is trained on, here we excluded long-context data. We also embedded the following three high quality datasets used to finetune Olmo3 Base into Olmo3 Instruct: Dolci3 SFT (2M input-output prompts, 3.08 GB), Dolci3 Instruct DPO (260k preference pairs, 811 MB), and Dolci3 Instruct RL (169k prompts, 483 MB).

*Table 6.* Datasets processed for contamination analysis. All text chunks are filtered to 50–2,048 tokens. Pretraining data is sampled at 1% using stratified reservoir sampling to preserve source distribution topology.

| Dataset | Orig. Size | Sample | Sampling Method |
|---|---|---|---|
| *Pretraining Data* | | | |
| Dolma3_6T-mix | 6.0 TB | 1.0% | Stratified Reservoir (Hierarchical) |
| Dolmino-mix | 250 GB | 1.0% | Stratified Reservoir (Hierarchical) |
| *Instruction Tuning Data* | | | |
| Dolci3 SFT | 3.08 GB | 100.0% | Full Ingestion |
| Dolci3 DPO | 811 MB | 100.0% | Full Ingestion |
| Dolci3 RL | 483 MB | 100.0% | Full Ingestion |

## A.2. Extended Cosine Similarity discussion

### A.2.1. MUSR



*(a) MuSR, Dolma*    *(b) MuSR, Dolmino*    *(c) MuSR, DolciSFT*    *(d) MuSR, DolciDPO*    *(e) MuSR, DolciRL*
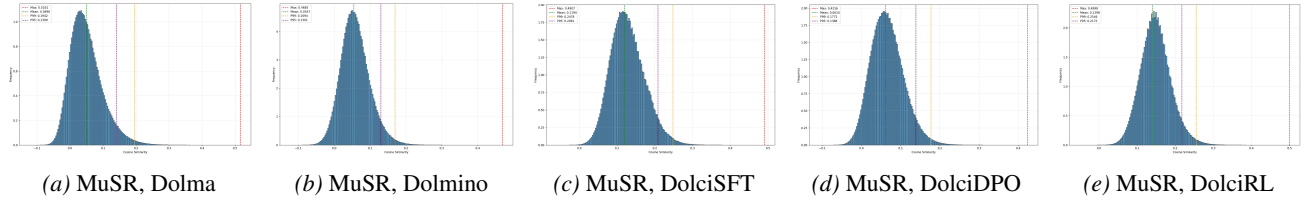
*Figure 7.* Each plot shows cosine similarity distribution of pairs of the MuSR benchmark data and training corpus data. From left to right we plot Dolma, Dolmino, Dolci SFT, Dolci DPO and Dolci RL.

### A.2.2. ZEBRALOGIC



*(a) ZebraLogic, Dolma*    *(b) ZebraLogic, Dolmino*    *(c) ZebraLogic, DolciSFT*    *(d) ZebraLogic,DolciDPO*    *(e) ZebraLogic, DolciRL*
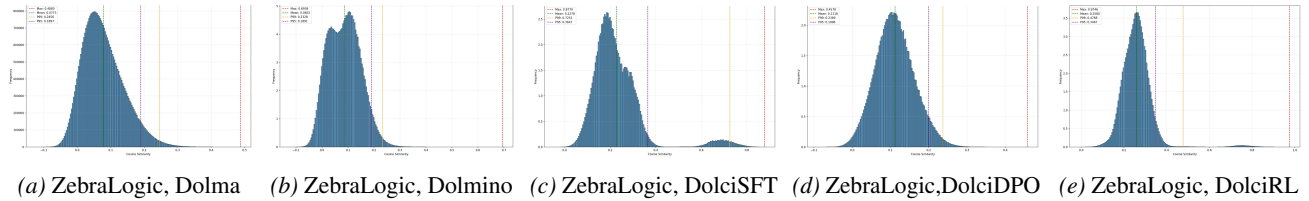
*Figure 8.* Each plot shows cosine similarity distribution of pairs of ZebraLogic benchmark data and training corpus data. From left to right we plot Dolma, Dolmino, Dolci SFT, Dolci DPO and Dolci RL.

## A.3. Synthetic Data Generation

### A.3.1. OVERVIEW

**MBPP**

We generate 5 semantic duplicates for both inputs (questions), and outputs (responses). For the text inputs, we generate paraphrasings simultaneously to maximise difference between the texts. For code outputs, we generate alternative python implementations and validate them against provided test cases.

**MuSR**    We adapt the original MuSR sample generation logic to generate new samples from existing reasoning trees of the public dataset problems. For the Murder Mysteries and Team Allocation tasks, we create three levels of semantic duplicates differing by how much the underlying logic trees differ from the original while maintaining the answer and the necessary reasoning to solve the problem the same: Level 1) tree is kept the same and story context is regenerated; Level 2) One branch that does not affect problem outcome or its complexity is changed; Level 3) All branches of the previous type are changed.

For each level, we generate 2 semantic duplicates for each of the 250 original samples, adding up to 1500 semantic duplicates per task. We also generate a new test set of 250 samples using the original code.

**ZebraLogic**    To generate semantic duplicates of ZebraLogic datapoints we use a variety of transformations, most of which are LLM based.

We apply the following transformation methods: 1) use category mappings that e.g. substitute "color" by "shape" and "red" by "square", using Claude 4.5 Haiku (Anthropic, 2025a); 2) shuffle conditions or clues in the prompt using a Python function; and 3) paraphrasing text while preserving all values exactly, using Claude 4.5 Sonnet (Anthropic, 2025b). These strategies are applied to the first 500 samples in the benchmark, generating semantic duplicates for paraphrasing alone, and the following combinations of the above: 1) and 2); 2) and 3); and 1), 2) and 3).

A.3.2. DETAILS OF FINETUNING ON DUPLICATES

*Table 7.* Finetuning hyperparameters for each benchmark. All experiments use LoRA Rank 16, Alpha 32, and Dropout 0.05.

| Hyperparameter | MBPP | MuSR | ZebraLogic |
|---|---|---|---|
| Learning rate | $1.5e-4$ | $2e-4$ | $2e-4$ |
| KL penalty | 0.02 | – | – |
| Epochs | 10 | 6 | 10 |

A.3.3. MUSR

The generation pipeline uses the same few-shot examples and prompt templates as the original benchmark generation set up. The following model parameters are used for with `gpt-4-0613`:

- Temperature: 1.0

- Top P: 1.0

- Max tokens: 2400

Specific details for each MuSR task:

**Murder Mysteries**    In this task, the model needs to figure out which of the two suspects is the murderer based on a long narrative generated from a logic tree. We make the following modifications to the logic trees before the story generation step:

- *Level 0*: Story context regenerated from unchanged original tree.

- *Level 1*: One suspicious fact branches changed.

- *Level 2*: All branches (suspicious, means, motive, opportunity) belonging to the innocent suspect changed.

**Team Allocation**    In this task, the model needs to determine how to best allocate the 3 individuals mentioned in the long narrative to perform two tasks. We make the following modifications to the logic trees before the story generation step:

- *Level 0*: Story context regenerated from unchanged original tree.

- *Level 1*: One randomly selected skill branch is changed while keeping skill level unchanged.

- *Level 2*: All branches swapped while skill and cooperation levels unchanged

**Object Placement**   The data available for the original samples is not enough to generate semantic duplicates. While it is possible to attempt to extract most of the required data with LLMs to attempt sample regeneration, format and style is lost, affecting subsequent generation steps in the multi-stage process. So we choose to omit this MuSR category due to difficulty in creating semantic duplicates of similar complexity and rigor.

### A.3.4. MBPP

We used `claude-opus-4-5-20251101` for all generation tasks with parameters:

- `max tokens`: 1024

- `temperature`: 1.0

We generate semantic duplicates for the full MBPP sanitized dataset, that is, 427 tasks. In the process, we find that the following tasks contain bugs (either in the function, or in the test cases): 229, 438, 461, 579, 769, 802.

**Inputs**   We batch generate *paraphrasings* to ensure clear differences between duplicates.

**Outputs**   We generate *Python semantic duplicates* for each sample sequentially, allowing the model to see the original and previously generated ones to increase the uniqueness of new solutions. We validate solutions at each step and ensure $< 0.85$ `difflib.SequenceMatcher` similarity score between solutions.

**Prompts**

Prompt for paraphrased question text duplicates:

```
PARAPHRASE_BATCH_PROMPT = """You are an expert at paraphrasing programming task
    descriptions.

ORIGINAL TASK:
{text}

YOUR TASK:
Generate exactly 5 DISTINCT paraphrases of this programming task. Each paraphrase must:
1. Have COMPLETELY DIFFERENT wording from the others
2. Preserve the EXACT same meaning and requirements
3. Maintain the same level of technical detail and clarity
4. Keep any mentioned function names UNCHANGED

CRITICAL: Each paraphrase must be noticeably different from the others. Vary:
- Sentence structure (active vs passive, questions vs statements)
- Vocabulary choices (synonyms, different technical terms)
- Order of information presented
- Level of formality

AVOID:
- Starting multiple paraphrases the same way (e.g., don't start 3 with "Write a...")
- Simply swapping one or two words while keeping structure identical
- Adding or removing requirements not in the original
- Changing the programming language if one is specified
- Making the task ambiguous or less precise

Output EXACTLY in this JSON format (no extra text, no markdown):
{{
  "para1": "first paraphrase here",
  "para2": "second paraphrase here",
  "para3": "third paraphrase here",
  "para4": "fourth paraphrase here",
  "para5": "fifth paraphrase here"
}}
```

```
Generate the 5 diverse paraphrases now:"""
```

Prompts for alternative Python solution implementations:

```python
# Python semantic duplicate generation
def get_generation_prompt(
    task_description: str,
    original_code: str,
    test_list: list[str],
    previous_solutions: list[str],
    previous_error: Optional[str] = None,
    require_more_different: bool = False
) -> str:
    """Generate the prompt for creating a Python semantic duplicate."""

    test_context = "\n".join(test_list[:3])
    func_name = extract_function_name(original_code)

    base_prompt = f"""You are an expert Python programmer. Your task is to write a
        DIFFERENT Python solution for the following problem.

TASK DESCRIPTION:
{task_description}

ORIGINAL PYTHON SOLUTION (for reference - DO NOT COPY):
```python
{original_code}
```

PYTHON TEST EXAMPLES (your solution must pass these):
```python
{test_context}
```

"""

    # Add previous solutions if any
    if previous_solutions:
        base_prompt += "PREVIOUS SOLUTIONS YOU'VE ALREADY WRITTEN (your new solution must
            be STRUCTURALLY DIFFERENT from ALL of these):\n"
        for i, sol in enumerate(previous_solutions, 1):
            base_prompt += f"\n--- Solution {i} ---\n```python\n{sol}\n```\n"
        base_prompt += "\n"

    # Add error feedback if retry
    if previous_error:
        base_prompt += f"""YOUR PREVIOUS ATTEMPT HAD ERRORS:
{previous_error}

Please fix these errors in your new solution.

"""

    # Add stronger differentiation request if needed
    if require_more_different:
        base_prompt += """IMPORTANT: Your previous solution was TOO SIMILAR to existing
            ones!
You MUST use a significantly DIFFERENT algorithmic approach. Consider:
- Using different data structures (list vs set vs dict vs deque)
- Using different iteration patterns (for vs while vs recursion vs comprehensions)
- Using different built-in functions or libraries
- Restructuring the logic flow completely

"""
```

```
    base_prompt += f"""REQUIREMENTS:
1. Write a COMPLETE Python solution that passes all tests
2. The function MUST be named EXACTLY: {func_name}
3. Use a DIFFERENT algorithmic approach or implementation style than the solutions shown
    above
4. ADD A COMMENT ON THE LINE ABOVE EACH LINE OF CODE explaining what it does
5. Comments should be ORIGINAL, CONCISE, and INSIGHTFUL - not generic
6. Make sure every substantive line has a comment above it
7. The comments should help distinguish this solution semantically from others

COMMENT STYLE EXAMPLE:
```python
# Initialize counter for tracking element frequency
count = 0
# Iterate through each item in the input sequence
for item in items:
    # Increment counter when condition is met
    if condition:
        count += 1
# Return the final accumulated count
return count
```

OUTPUT ONLY THE PYTHON CODE with comments. No markdown code blocks, no explanations
    outside the code."""

    return base_prompt
```

### A.3.5. ZEBRALOGIC

Below we discuss the different methods and transformations:

**Paraphrasing**. With `claude-4.5-sonnet` we use the following prompts using the original sample

```
# SYSTEM PROMPT
You are an expert editor tasked with rewriting logic grid puzzles while exactly preserving
    the logical structure and semantics.

# USER PROMPT
Rewrite the following logic puzzle to express the exact same conditions in different words
    or with different word order etc. while exactly preserving the logical structure and
    semantics.

Original Puzzle:
{puzzle}

REQUIREMENTS:
1. Reformulate both the task description and every numbered condition.
2. You may change word order, use synonyms, and alter sentence structure.
3. PRESERVE the strict logical meaning. For example, "A is next to B" must remain
    logically equivalent (e.g., "B is adjacent to A").
4. PRESERVE all entity names, values, numbers, and categories EXACTLY. Do not change "Red"
    to "Crimson" or "John" to "Jon". The specific terms used for the puzzle items MUST
    remain identical to match the solution exactly.
5. The output must be natural, clear, and readable. Avoid contrived or unnatural
    constructions.
6. Maintain the formatting of the puzzle, including the format and numbering of the list
    of clues.
7. Do not start your response with a header or a preamble. Start with a naturally flowing
    puzzle statement in a very similar style and format as the original.

Output ONLY the rewritten puzzle text.
```

**Category substitution**. With `claude-4.5-haiku` we follow a two step process: (1) generate a substitution plan for each

puzzle; (2) apply the substitution with LLMs for improved text cohesiveness; (3) transform the solution programatically.

```
# STEP 1: SYSTEM PROMPT
You are a helpful assistant that creates substitution plans for logic puzzles.
Your goal is to transform the puzzle by changing BOTH the categories and their values to
    new domains.

# STEP 1: USER PROMPT
Create a substitution plan to transform this logic grid puzzle.
1. Identify all categories (e.g., Color, Drink, Pet).
2. Assign a NEW category to each (e.g., Color -> Shape, Drink -> Snack, Pet -> Book).
3. Map every existing value to a new value appropriate for the new category.

Original Puzzle:
{puzzle}

Original Solution:
{solution_json}

REQUIREMENTS:
1. Change the categories to natural, distinct alternatives (e.g., colors -> shapes,
    flowers -> animals).
2. Keep the new categories and values DISTINCT from all of the original ones. Avoid number
     categories (to avoid confusion with the numbering of the puzzle).
3. Ensure 1-to-1 mapping for all values.
4. Do NOT use obscure or unusual categories. Stick to common categories like colors,
    animals, shapes, countries, etc. Choose natural categories and values within the flow
    of the puzzle wording.

Output ONLY a JSON object with this structure:
{
  "substitution_plan": {
    "OriginalCategoryName": {
      "new_category": "NewCategoryName",
      "values": {
        "OldValue1": "NewValue1",
        "OldValue2": "NewValue2"
      }
    },
    ...
  }
}

# STEP 2: SYSTEM PROMPT
You are a helpful assistant that rewrites logic puzzles based on a substitution plan.
You must replace categories and values exactly according to the plan while PRESERVING the
    puzzle structure, logic, and clues EXACTLY.

# STEP 2: USER PROMPT
Rewrite this logic puzzle by applying the following substitution plan.
Replace ALL occurrences of the old categories and values with their corresponding new ones
    .

Substitution Plan:
{plan_json}

Original Puzzle:
{puzzle}

CRITICAL INSTRUCTIONS:
1. Replace old categories (e.g., "Color") with new categories (e.g., "Shape").
2. Replace old values (e.g., "Red") with new values (e.g., "Square").
3. Do NOT change the logic, clues, or structure.
4. Keep the puzzle wording identical as much as possible, only make minor syntactic
    adjustments where necessary to preserve the flow and meaning of the puzzle wording.
```

```
5. Keep the numbering and formatting identical.
6. Output ONLY the rewritten puzzle text.
```

**Shuffling**. Puzzles clues are parsed, randomly reordered, and renumbered sequentially. The solution of the resulting semantically equivalent puzzle remains the same.

**Composite transformations**. With respect to the order of transformations in composite methods, shuffling is always performed first, and paraphrasing is performed last.

A.3.6. SIMILARITY DISTRIBUTIONS OF SYNTHETIC SEMANTIC DUPLICATES

We compare the generated semantic duplicates against the original samples for MBPP, MuSR and CodeForces, and show the analysis for Cosine similarity and several common metrics used in deduplication: n-gram overlap (2 and 3 grams), ROUGE-L F, and Jaccard token. See Figure 9.

Finding that in most cases, cosine similarity is better at separating semantic duplicate pairs than the other metrics. This motivates, in part, the use of embedding similarity for the other experiments.

## A.4. Annotation schemes for high cosine similarity matches

We use `gemini-3-flash-preview` with the following parameters:

- `thinking_level`: MEDIUM (we use HIGH for CodeForces due to the length of the problems.)

- `temperature`: 1.0

- `max_output_tokens`: 8192

- `response_format`: JSON, with structured output conforming to the schema below.

For each annotated pair, we use a schema to collect the following:

- `is_sd`: A boolean indicating whether the pair constitutes a semantic duplicate (true if the corpus task is the same as or subsumes the test task.)

- `confidence`: A confidence score in $[0, 1]$, where $1.0$ indicates certainty and $0.0$ represents a 50-50 guess.

- `reasoning`: A free-text explanation of the judgment, including key similarities and differences observed.

- `match_type`: A categorical label describing the relationship: *exact* (nearly identical), *equivalent* (same task, different wording), *subset* (test is a subset of corpus), *superset* (corpus is a subset of test), or *unrelated* (different tasks.)

The annotations are then used to finetune `gpt-oss-20b` as a classifier for CodeForces semantic duplicate detection, and `Qwen/Qwen3-30B-A3B-Instruct-2507` for MBPP.

A.4.1. CLASSIFIER FINETUNING RESULTS

**MBPP.** We finetune `Qwen/Qwen3-30B-A3B-Instruct-2507` with LoRA (rank 32) on 13266 samples (9% semantic duplicates), and achieve 96.8% accuracy on a held-out validation set (n=560, 25% semantic duplicates) with 97.7% precision and 93.3% F1 on the semantic duplicate class.

**CodeForces.** We also finetune `gpt-oss-20b` with LoRA (rank 32) on 13398 samples (9% semantic duplicates), and achieve 95.3% accuracy on the held-out validation set (n=512, 25% semantic duplicates) with 97% precision and 90% F1 on the semantic duplicate class.

A.4.2. ANNOTATION PROMPT FOR MBPP MATCHES

```
MBPP_PROMPT_TEMPLATE = """You are an expert programmer analyzing potential semantic
    duplicates between coding tasks.
```

```
## Task
Determine if the following two coding tasks are semantic duplicates - meaning they
    describe the same programming task, just potentially phrased differently.

## Test Task (from benchmark):
{test_text}

## Corpus Task (from training data):
{corpus_text}

## Guidelines:
1. **Focus on the TASK, not the solution** - ignore any code or solutions that may be
    present
2. **Mathematical equivalence counts as duplicate** - e.g., "sum 1 to n" and "sum n, n-1,
    ..., 1" are equivalent
3. **Corpus subsumes test = duplicate** - if the corpus task is strictly harder (asks for
    more), but solving it would trivially solve the test task, mark as duplicate
4. **Be calibrated** - use confidence primarily for ambiguous cases, tricky phrasing, or
    when you're uncertain

## Match Types:
- "exact": Nearly identical wording
- "equivalent": Different phrasing, same underlying task
- "subset": Test task is a subset of corpus task (corpus is harder but solves test)
- "superset": Corpus task is a subset of test task (test is harder) - NOT a duplicate
- "unrelated": Different tasks entirely

Analyze the tasks and provide your structured judgment."""
```

### A.4.3. ANNOTATION PROMPT FOR CODEFORCES MATCHES

```
CODEFORCES_PROMPT_TEMPLATE = """You are an expert competitive programmer analyzing
    potential semantic duplicates between programming problems.

## Task
Determine if the following two competitive programming problems are semantically related -
    meaning exposure to the corpus problem during training could help solve the test
    problem.

## Test Problem (from benchmark):
{test_text}

## Corpus Problem (from training data):
{corpus_text}

## Analysis Steps:
1. **Check data quality first**: Is the corpus text a complete problem statement? If it's
    empty, fragmentary, or contains only code without a problem description, mark as "
    unrelated".
2. **Check for exact text match**: If the corpus text appears VERBATIM (word-for-word)
    within the test text (e.g., corpus contains just the problem statement while test
    contains problem + examples), this counts as "exact" match.
3. **Extract the core problem**: Strip away story/narrative framing. What is the actual
    computational task?
4. **Identify the key insight**: What algorithmic technique or observation is needed?
5. **Compare**: Is there meaningful overlap in what's being asked or how to solve it?

## Match Types:
- "exact": Nearly identical problem statements, OR corpus text is a verbatim substring/
    subsection of test text (exact text match even if corpus is shorter)
- "equivalent": Different framing but identical algorithmic core
- "subset": Test is a special case of corpus (test asks for less than corpus)
- "superset": Corpus asks for something simpler than test, but NOT a verbatim text match
```

```
- "related": Corpus covers a component or shares key insight with test
- "unrelated": Different problems, or corpus data is unusable

## IMPORTANT: Exact Match Clarication
If the corpus text is an exact substring of the test text (the corpus text appears word-
    for-word inside the test text, just without some sections like examples or input/
    output format), mark this as "exact" NOT "superset". The key distinction:
- "exact": Corpus text IS CONTAINED VERBATIM in test text
- "superset": Corpus asks a DIFFERENT (simpler) question than test

## What counts as semantically related:
- Same computational task (any framing)
- One is a special case of the other
- Shared key insight or trick
- Corpus solves a significant component of test

## What is unrelated:
- Sharing only common techniques (DP, BFS) without structural similarity
- Unusable corpus data (empty, fragmentary, code-only)
- Genuinely different computational questions"""
```

## B. Further Semantic Duplicates in the Wild Results

### B.1. Reporting on top 100 cosine similarity matches instead of 100 sampled from top 0.1%

### B.2. Semantic duplicates are hard to detect

Semantic duplicates in the wild are sparse and difficult to find. From the above we notice that semantic duplicates can be found, even for hard CodeForces level problems. We consider, per test point, approximately 350 million texts. We find in the case of CodeForces on average a few semantic duplicates. For MBPP there are tens to one hundred across our entire dataset. Thus semantic duplicates are both incredibly rare in the wild, and occur with frequency at most one in a million text segments across the internet. Thus we run our algorithm across around 2.5 terabytes of data to find necessary duplicates for a representative population.

To illustrate this point we demonstrate the probability of a semantic duplicate occurring given a cosine similarity and being in the top 0.1% of a semantic duplicates test set.

The relationship between embedding cosine similarity and the probability of semantic duplication, aggregated across all training stages (N=128,408 training-test pairs). Points represent binned similarity scores (30 bins); shaded regions indicate 95% confidence intervals computed using the normal approximation to the binomial distribution. Sample sizes for each bin are annotated.

Semantic duplicate rate exhibits a nearly monotonically increasing relationship with cosine similarity as we would expect. Below a similarity threshold of approximately 0.35, the duplicate rate is effectively zero (¡1%). The rate increases sharply between 0.4 and 0.7, following an approximately sigmoidal trajectory, and reaches 60–85% at the highest observed similarities (¿0.8). The widening confidence intervals at high similarity values and volatility reflect reduced sample sizes in these bins. This calibration curve suggests that cosine similarity serves as a useful but imperfect proxy for semantic duplication, with a practical decision threshold in the 0.5–0.6 range capturing the inflection point of the relationship.

### B.2.1. ECOLOGICALLY VALID FINETUNING EXPERIMENT

**Finetuning parameters**. We used the following hyperparameters, training all layers:

- `LoRA Rank`: 64

- `dropout`: 0.05

- `Epochs`: 5

**Semantic Duplicate Data**. For the seen split of the data, consisting of the first 125 MuSR samples, we used all Level 2 and Level 3 semantic duplicates. A total of 500 since we generate 2 per level.

*Table 8.* Effect of finetuning Olmo3 on semantic duplicates of MuSR Murder Mysteries reasoning traces. Olmo3 was finetuned for 3 epochs.

| DUPLICATION LEVEL | TEACHER: OPUS 4.5 | TEACHER: GPT 4.1 MINI |
|---|---|---|
| BASELINE | 66.0 | 66.0 |
| EXACT DUPES | 87.1 | 82.5 |
| LEVEL 0 | 86.8 | 82.4 |
| LEVEL 1 | 86.1 | 81.3 |
| LEVEL 2 | 85.8 | 81.6 |

*Table 9.* No degradation effect of finetuning Olmo3 on semantic duplicates of MuSR Murder Mysteries reasoning traces. Olmo3 was finetuned for 3 epochs.

| DUPLICATION LEVEL | ARC CHALLENGE | ARC EASY | BOOLQ | HELLASWAG | PIQA | WINOGRANDE |
|---|---|---|---|---|---|---|
| BASELINE | 50.1 | 78.2 | 75.7 | 57.5 | 75.0 | 65.3 |
| EXACT DUPES | 50.0 | 78.7 | 76.7 | 56.4 | 75.8 | 65.0 |
| LEVEL 0 | 50.3 | 78.6 | 76.0 | 56.4 | 75.7 | 64.8 |
| LEVEL 1 | 50.3 | 78.6 | 77.0 | 56.4 | 75.6 | 64.7 |
| LEVEL 2 | 50.6 | 78.9 | 77.1 | 56.5 | 75.6 | 65.1 |

## C. Further Finetuning Results, Including Degradation Analysis

The Opus 4.5 MuSR baseline accuracy is 91.6%. The model does slightly worse on the first half of the data (that we train on in Table 8) than on the second half, respectively Opus4.5 gets 90.4% on the first half and 92.8% on the second half.

The performance of `gpt-4.1-mini-2025-04-14` on MuSR benchmark data is 84.0, on our level0 semantic duplicates it is 79.8 and on level2 it is 76.4. This is just the baseline performance of the teacher model GPT 4.1 mini without any finetuning.
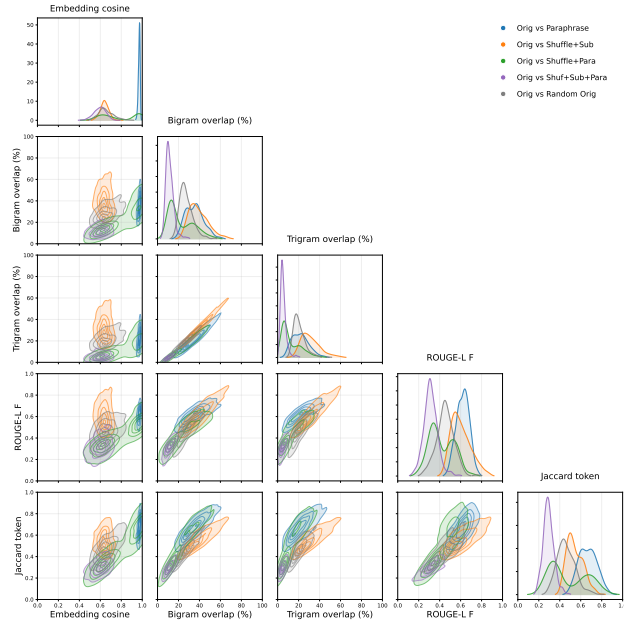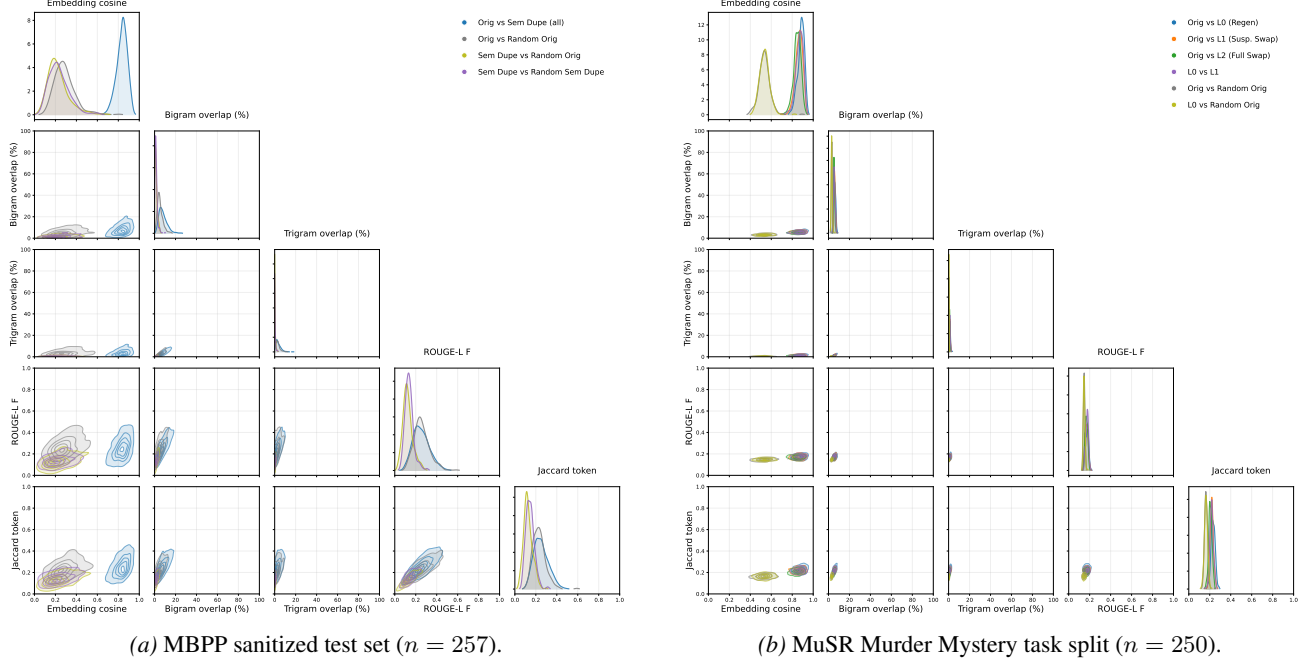
## D. Ecologically Finetuned Results

*Table 10.* No degradation effect of finetuning Olmo3 on semantic duplicates of ZebraLogic reasoning traces. Olmo3 was finetuned for 3 epochs.

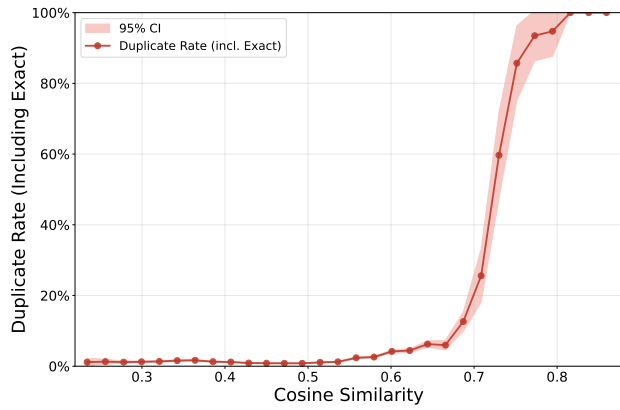| DUPLICATION LEVEL | ARC CHALLENGE | ARC EASY | BOOLQ | HELLASWAG | PIQA | WINOGRANDE |
|---|---|---|---|---|---|---|
| BASELINE | 50.1 | 78.2 | 75.7 | 57.5 | 75.0 | 65.3 |
| EXACT DUPES | 49.5 | 78.0 | 77.0 | 56.8 | 75.0 | 64.7 |
| PARA | 49.3 | 77.7 | 78.3 | 56.4 | 74.5 | 64.1 |
| SHUFFLE, SUBS | 50.7 | 77.4 | 75.3 | 56.5 | 75.1 | 64.5 |
| SHUFFLE, PARA | 49.4 | 78.1 | 77.8 | 56.5 | 75.3 | 64.9 |
| SHUFFLE, SUBS, PARA | 50.4 | 78.1 | 78.4 | 56.5 | 75.6 | 63.5 |

*Table 11.* We report on baseline (before finetuning) accuracy on MuSR. We then finetune on 10.000 datapoints. We either finetune on half of the level 2 & 3 semantic duplicates mixed in with regular data (contaminated model) or we finetune on clean data only (clean model).

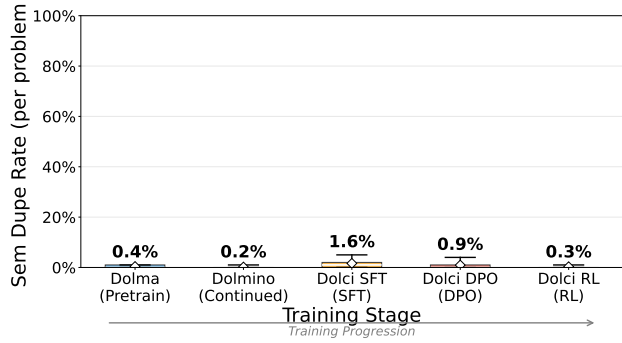| MODEL | DUPLICATION LEVEL | CONTAMINATED | | CLEAN | |
|---|---|---|---|---|---|
| | | SEEN | UNSEEN | SEEN | UNSEEN |
| OLMO3 | BASELINE | 44.0 | 41.6 | 44.0 | 41.6 |
| | FINETUNED | 66.4 | 54.4 | 51.2 | 48.8 |
| QWEN3 | BASELINE | 39.2 | 41.6 | 39.2 | 41.6 |
| | FINETUNED | 65.6 | 52.0 | 48.0 | 59.2 |

*(a)* MBPP sanitized test set ($n = 257$).



*(b)* MuSR Murder Mystery task split ($n = 250$).



*(c)* ZebraLogic dataset ($n = 1000$).

*Figure 9.* Similarity distributions using several deduplication metrics for each benchmark.

*(a)* Correlation of Similarity vs. Duplicates (including semantic and exact

*(b)* Propensity of semantic duplicates by Training Scheme (excluding exact duplicates)

*Figure 10.* Analysis of semantic duplicates in top 100 CodeForces rounds. Left: Difficulty vs. likelihood of semantic duplicates. Right: Duplicate propensity across different training schemes.