

LLM Generation Novelty Through the Lens of Semantic Similarity

Anonymous Authors¹

Abstract

Generation novelty is a key indicator of an LLM’s ability to generalize, yet measuring it against full pretraining corpora is computationally challenging. Existing evaluations often rely on strict lexical overlap or do not consider the full pre-training corpus. We frame novelty as a semantic retrieval problem, which enables the use of modern retrieval pipelines for efficient analysis at pre-training scale. Specifically, we propose a three-stage framework that retrieves semantically similar samples, reranks them at varying subsequence lengths, and calibrates scores using a human novelty reference for interpretability. We apply this framework to the SmoLLM model family and report three key findings: (1) models draw on pretraining data across much longer sequences than previously reported; (2) some task domains systematically promote or suppress generation novelty; and (3) instruction tuning not only alters style but also increases novelty. These results highlight the value of semantic novelty analysis for studying generalization.

1. Introduction

Large language models (LLMs) now power chatbots, copilots, and autonomous agents with applications across various domains. Their adoption hinges on an implicit assumption: LLMs are capable of generalizing beyond their training data to generate relevant and novel outputs in response to user prompts. Generation novelty provides a useful signal of this capability, indicating whether model outputs extend beyond patterns observed during training. Novelty signals compositional generalization and allows for the assessment of true zero-shot behavior – informing debates about provenance and intellectual property. Thus, measuring a model’s generation novelty is a prerequisite for interpreting what

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

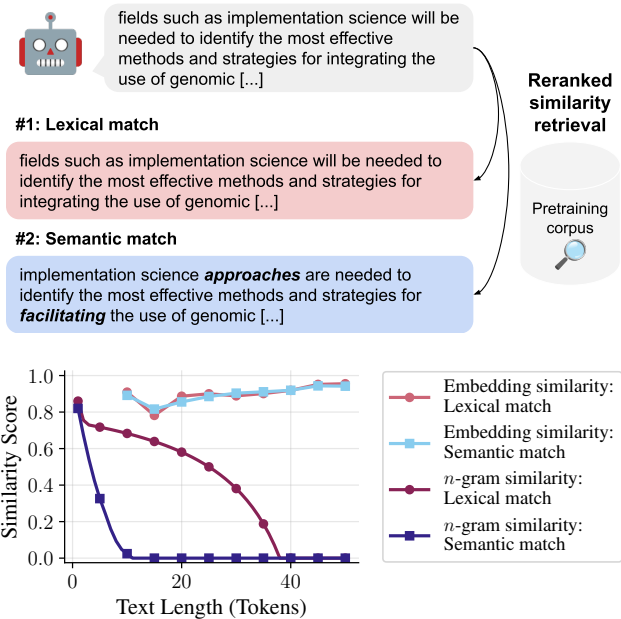


Figure 1. Efficient retrieval allows for pretraining-scale semantic novelty analysis. Given a generation, we employ retrieval and reranking pipelines to identify semantically similar samples in the pretraining corpus to analyze generation novelty. Higher similarity indicates lower novelty. While the first-ranked sample is a lexical match, the second-ranked is a semantic match. As text length increases, n -gram overlap drops to zero, falsely suggesting high novelty for paraphrased content. In contrast, embedding similarity correctly identifies the shared meaning, providing a more reliable measure of novelty.

LLMs actually learn.

However, analyzing the novelty of LLM generations is a non-trivial task. There is no single, clearly defined notion of novelty, which opens up the question: Is a generation already novel if it does not appear verbatim in the pretraining data? More broadly, what does it mean for a generation to be “reused” from the training data? At the same time, the sheer scale of modern pretraining corpora makes comparison computationally challenging, so that novelty analysis is both a conceptual and a technical problem.

Existing approaches to studying LLM generation novelty mainly address these questions in two directions. One line of

work relies on efficient textual overlap metrics (McCoy et al., 2023; Merrill et al., 2024; Padmakumar et al., 2025), which scale well but naturally focus on verbatim reuse and fail to account for paraphrases or stylistic variation. A second line of work investigates novelty within specific domains, such as scientific ideas or biomedical publications (Peng et al., 2025; Ai et al., 2025; Wang et al., 2025c), typically using semantic similarity measured against a curated reference corpus. While effective, these approaches do not account for reuse from other parts of the broader pretraining corpus. Together, these limitations highlight the need for novelty measures that go beyond verbatim overlap at the scale of LLM pretraining corpora.

We address this need by proposing a three-stage framework for analyzing an LLM’s generation novelty at scale through the lens of semantic similarity. We argue that a generation should be called novel only if it does not semantically reproduce the training data (cf. Figure 1). From this perspective, “reused” information is not limited to verbatim overlap, but includes paraphrases and reformatting.

Our framework combines semantic retrieval (stage 1) and reranking (stage 2) (Santhanam et al., 2022; Li et al., 2025) to compare LLM generations against their pretraining corpora. However, similarity scores and their relative differences are not directly interpretable as measures of novelty. They are affected by biases in the retrieval pipeline, such as a preference for shorter documents and information located early within documents (Fayyaz et al., 2025; Zhou et al., 2025). Because of this, identical scores can reflect different degrees of novelty, especially when comparing different text lengths or domains. Hence, the third stage of our framework calibrates the raw semantic similarity scores using a baseline of held-out, human-written reference text. By calibrating the scores, we mitigate artifacts of the retrieval pipeline and enable meaningful comparison of scores. Finally, we aggregate these scores into a *novelty profile* which characterizes the generation novelty of a model for a specific dataset and generation text length. This procedure is model- and task-agnostic, and lightweight enough to run on full pretraining corpora. It thus offers a scalable way to assess an LLM’s generation novelty.

We demonstrate the utility of our framework through empirical analyses on SmoLLM (Allal et al., 2024) and SmoLLM2 (Allal et al., 2025), two LLMs with open pretraining corpora. Our analysis reveals unexpected patterns missed by previous lexical methods (McCoy et al., 2023; Merrill et al., 2024). First, both models draw on pretraining data over much longer sequences than previously reported (Merrill et al., 2024). Second, novelty varies systematically by task domain. Third, embedding-based novelty estimates are stable under style shifts from instruction tuning; after accounting for these shifts, instruction tuning

substantially increases novelty. These results suggest that instruction tuning shapes not only style but also compositional generation behavior.

Our contributions are as follows:

- We present a semantic similarity–based framework for measuring LLM generation novelty at scale, combining retrieval, reranking, and baseline-calibration to enable comparison against full pretraining corpora while reducing sensitivity to stylistic variation.
- We empirically analyze generation novelty in SmoLLM and SmoLLM2, uncovering long-range reuse patterns, task-dependent variation, and the impact of instruction tuning beyond stylistic effects.
- We also release the corresponding indices and corpus chunks of SmoLLM and SmoLLM2 to support replication and extension.

2. Related Work

Novelty measurement in LLMs. Prior work on LLM generation novelty has predominantly relied on textual overlap–based measures, particularly n -gram comparisons. McCoy et al. (2023) introduce n -novelty, defining novelty as the absence of copied text from the training data and quantifying the proportion of non-overlapping n -grams in model outputs. Building on this, Merrill et al. (2024) additionally analyze the probability of generating training n -grams, while Padmakumar et al. (2025) propose novelty as the harmonic mean of n -novelty and generation quality assessed by LLM judges. Wang et al. (2025b) further extend this line of work by introducing task-grams, which capture task-specific n -gram co-occurrences in output and corpus. The task relation introduces a semantic dimension to the n -gram-based analysis. Despite these extensions, all of these approaches rely on surface-level textual overlap and may treat paraphrases or stylistic variation as novel. We instead define novelty via semantic similarity, enabling analysis beyond surface-form variation. Prior work that emphasizes semantic novelty typically focuses on specific domains, such as scientific ideas or biomedical text (Ai et al., 2025; Peng et al., 2025; Wang et al., 2025c). These approaches also rely on embedding-based similarity, but assess novelty relative to selected reference documents tailored to the target domain. In contrast, our work studies generation novelty with respect to the entire pretraining corpus, independent of the prompt.

Memorization and membership inference. Memorization work (Wu et al., 2025; Feldman & Zhang, 2020) investigates whether specific samples can be elicited verbatim from a model to understand if the model has memorized them. In

contrast, our notion of novelty captures whether the underlying *information* is present in the corpus, even if paraphrased. Membership inference attacks (MIA) (Puerto et al., 2025; Mesana et al., 2025; Zhang et al., 2025; 2024) instead ask whether a particular example was part of pretraining, often in adversarial settings. While informative for questions like data privacy, MIAs do not address the broader question of how models generate text that are *not part of* their training data. Our novelty measure, therefore, complements both attribution and memorization/MIA, providing a new perspective on generalization.

Further discussion of prior work relating model generations to pretraining data is provided in Appendix B.

3. Novelty analysis framework

This section presents our framework for studying LLM generation novelty. We introduce the conceptual definition of *semantic* novelty and propose a framework that measures calibrated semantic novelty at LLM pretraining scale.

3.1. Conceptual Definition of Semantic Novelty

We define *generation novelty* as the degree to which an LLM’s output expresses information or patterns not readily present in its pretraining corpus. Unlike *lexical novelty*, which is typically measured via n -gram non-overlap (McCoy et al., 2023; Merrill et al., 2024), our definition focuses on *semantic novelty*:

$$\min_{d \in \mathcal{C}} \left(1 - \cos(\phi(y), \phi(d)) \right), y \sim f_{\theta}(x) \quad (1)$$

where \mathcal{C} is the training corpus, y is the LLM f_{θ} ’s output to a prompt x , and ϕ is an embedding function. In other words, semantic novelty is defined as the minimum cosine distance between a semantic representation of the LLM generation and documents of the pretraining corpus. Under this definition, an output’s novelty is determined by its content, rather than its surface form. Hence, this definition serves to distinguish between *reproduction* (generating sequences that exist semantically in the corpus) and *composition* (generating sequences that are internally coherent and correct, yet semantically distinct from any specific training document).

By aggregating semantic signals across a dataset, we can characterize the *novelty profile* of a specific model configuration or training regime. Crucially, we treat novelty not as a binary state, but as a continuous spectrum. We do not claim that a “novel” generation has no relationship to the training data; rather, we measure the extent to which the generated content deviates from the most similar semantic neighbors available in the corpus.

3.2. A Framework for Analyzing Semantic Novelty

Our framework operationalizes the measurement of semantic novelty through a model-agnostic, three-stage paradigm applied at the dataset scale. This approach is inspired by standard retrieval-augmented generation (RAG) and information retrieval architectures, where a two-stage process of ranking and re-ranking is the established procedure for balancing search efficiency with semantic precision (Li et al., 2025). Since our framework addresses generation novelty, we further append a third stage that mitigates potential skewness artifacts of the retrieval pipeline, arising from sequence length or domain-specific density of the LLM generations. This stage calibrates the novelty scores with respect to a matched human-level novelty reference so that scores are comparable and interpretable.

The framework aims to characterize the *novelty profile* $\mathcal{N}_{\mathcal{Q}}(k)$ of a model distribution \mathcal{Q} relative to a domain-matched human distribution \mathcal{H} , for each chunk size k . By evaluating across multiple chunk sizes k , we can observe how specific model configurations (e.g., scale or alignment) influence novelty. The framework is summarized in Algorithm 1. We describe each stage in the following:

- Local Candidate Pool Identification.** For every generation in the model distribution $q \in \mathcal{Q}$ and corresponding baseline from the human distribution $h \in \mathcal{H}$, first identify a local candidate pool D within the pretraining corpus \mathcal{C} . To this end, a coarse-grained ranker \mathcal{R} retrieves the top- n passages that exhibit the highest potential semantic overlap with the full document. This stage acts as a high-recall filter, ensuring that any potential semantic neighbors are captured once per document, providing a computationally efficient search space for subsequent multi-scale analysis.
- Multi-Scale Semantic Re-ranking.** To distinguish between short-range lexical reuse and long-range compositional novelty, evaluate text at multiple chunk sizes k . For a given k , decompose the generation and its corresponding candidate pool into smaller fragments. A high-precision re-ranker \mathcal{S} then computes the maximum similarity between each chunk and its respective pool. This yields a set of *raw similarity scores* for the entire distribution. This two-stage process ensures that the scoring is fine-grained and robust to paraphrasing while remaining tractable at the scale of trillion-token corpora.
- Calibrated Distributional Normalization.** As human and model outputs differ in length and structure, a 1:1 element-wise comparison between individual chunks is often impossible: Given one human baseline text for each model generation, they ultimately are split into a different number of chunks for small enough k .

Algorithm 1 Novelty Framework with retrieval, reranking and baseline-normalized scoring

```

165 Require:  $\mathcal{C}$  (Corpus),  $\mathcal{Q}$  (Model generations),  $\mathcal{H}$  (Human references),  $\mathcal{R}$  (Ranker),  $\mathcal{S}$  (Re-ranker),  $K$  (Set of chunk sizes
166  $\{k_1, k_2, \dots\}$ ),  $n$  (Number of retrieved passages in Stage 1)
167 // Stage 1: Identification of Local Candidate Pools
168 for each document  $d \in \mathcal{Q} \cup \mathcal{H}$  do
169      $P_d \leftarrow \mathcal{R}(d, \mathcal{C}, n)$  ▷ Retrieve  $n$  passages once per document using Eq. 1
170 end for
171 // Stage 2: Multi-Scale Evaluation and Calibration
172 for each chunk size  $k \in K$  do
173      $\text{Scores}_H \leftarrow [], \text{Scores}_Q \leftarrow []$  ▷ Initialize empty sequences
174     // Determine the “Semantic Noise Floor” for this configuration
175     for each document  $h \in \mathcal{H}$  do
176          $\text{Chunks}_h \leftarrow \text{chunk}(h, k)$ 
177          $S_h \leftarrow [\max_{p \in P_h} \mathcal{S}(c, p) \mid c \in \text{Chunks}_h]$  ▷ Sequence of max scores
178          $\text{Scores}_H \leftarrow \text{Scores}_H \oplus S_h$  ▷ Concatenate sequences
179     end for
180      $\tilde{\mu}_H^{(k)} \leftarrow \text{median}(\text{Scores}_H)$  ▷ Stable calibrator for domain and length  $k$ 
181 // Stage 3: Calculate Calibrated Similarity Scores for the Model
182 for each document  $q \in \mathcal{Q}$  do
183      $\text{Chunks}_q \leftarrow \text{chunk}(q, k)$ 
184      $S_q \leftarrow [\max_{p \in P_q} \mathcal{S}(c, p) \mid c \in \text{Chunks}_q]$ 
185      $R_q \leftarrow [s/\tilde{\mu}_H^{(k)} \mid s \in S_q]$  ▷ Normalize by human distribution
186      $\text{Scores}_Q \leftarrow \text{Scores}_Q \oplus R_q$  ▷ Concatenate sequences
187 end for
188      $N^{(k)} \leftarrow \text{median}(\text{Scores}_Q)$  ▷ Final novelty profile for chunk size  $k$ 
189 end for
190 return  $\{(k, N^{(k)}) \mid k \in K\}$  ▷ The novelty profile of the model

```

Instead, we use the human distribution \mathcal{H} to establish a stable, domain-specific *calibration constant* $\tilde{\mu}_H^{(k)}$ for each chunk size. This constant represents the semantic noise floor, i.e., the level of similarity expected from novel, held-out human text within that domain, configuration, and chunk size. Further motivation and details on the baseline calibration are provided in Appendix D.

The final *calibrated similarity score* for the model distribution is calculated by normalizing the model’s raw scores by the semantic noise floor constant. By aggregating these ratios, we arrive at a *novelty profile* that allows for rigorous comparison across different model architectures or training regimes. This normalization step ensures that any observed trends are not artifacts of the retrieval pipeline or domain-specific redundancies, but true reflections of the model’s divergence from natural human patterns of information reuse.

4. Experiments

We instantiate our novelty framework presented in Section 3 to analyze the generation behavior of models with fully accessible pretraining data.

4.1. Experimental Setup

Models and Data (\mathcal{Q} and \mathcal{C}). We conduct our analysis on the SmoLLM (Allal et al., 2024) and SmoLLM2 (Allal et al., 2025) families licensed under the Apache-2.0 license. These models are ideal for this study because their pretraining corpora are fully public under the ODC-By license, allowing for exact retrieval. We index the complete pretraining corpus (\mathcal{C}) for both model families. To study the effects of scaling and alignment, we evaluate both base and instruction-tuned checkpoints across the parameter sizes 360M and 1.7B.

Retrieval Instantiation (\mathcal{R} and \mathcal{S}). We instantiate the coarse-grained ranker \mathcal{R} using L2-normalized GIST embeddings (Solatorio, 2024) indexed via FAISS (Douze et al., 2024). We chose GIST for its high efficiency-to-performance ratio on the MTEB leaderboard (Muennighoff et al., 2023) at the time of the experiments. For candidate generation (Stage 1 in Algorithm 1), we retrieve $n = 100$ document chunks, which suffices for capturing the most similar corpus document in the vast majority of cases, as validated in Appendix E. We instantiate the re-ranker \mathcal{S} using ColBERTv2 (Santhanam et al., 2022). ColBERTv2’s late-interaction mechanism provides granular token-level alignment, making it robust to the paraphrasing and stylistic

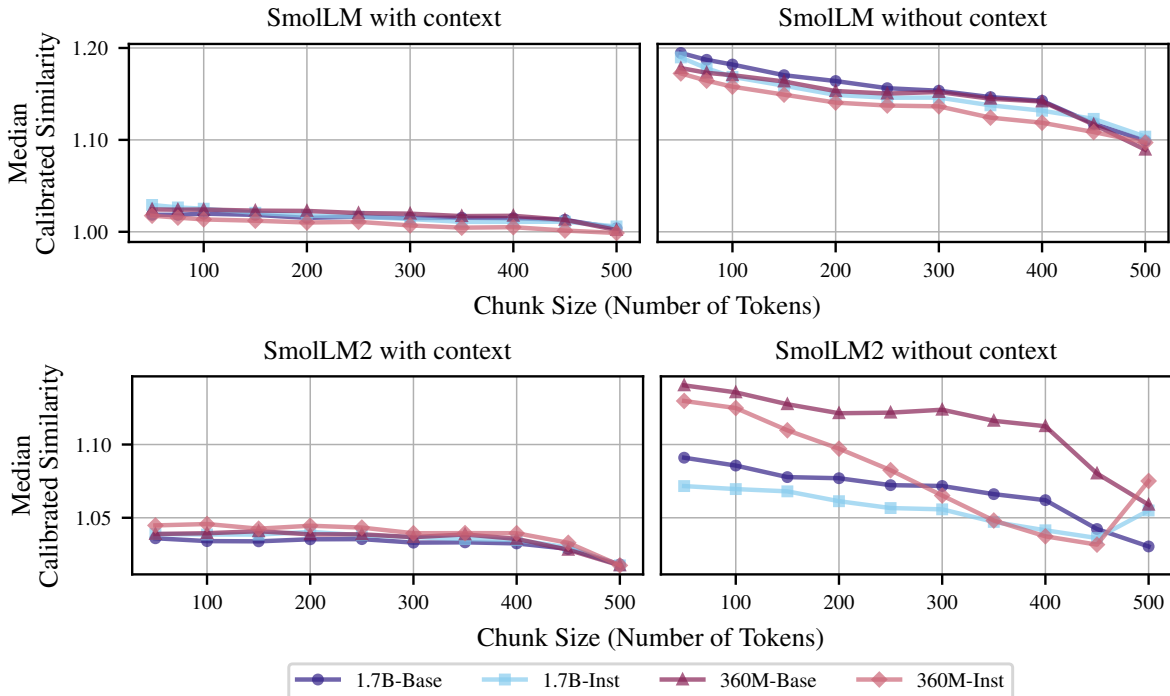


Figure 2. Novelty profiles of SmoLLM (top) and SmoLLM2 (bottom), for prompted, and unprompted generations. Higher similarity indicates lower novelty.

shifts and ensuring that our novelty scores reflect genuine conceptual divergence rather than surface-level patterns.

Scale of Analysis. We process the corpus into chunks of 512 tokens with a 50-token overlap to mitigate boundary effects (the effect of chunking borders on our retrieval pipeline is analyzed in Appendix F). This results in ~ 20 TB of embeddings and indices, which we publish for reproducibility upon acceptance. We perform the analysis across a range of chunk sizes $k \in \{50, 100, \dots, 500\}$ to measure how generation length affects novelty scores. Details on computational costs are provided in Appendix I. We report the p -values for our main hypotheses throughout this section in Appendix M.

4.2. Natural Generation Novelty

We first characterize the general novelty profile of the models in an open-ended setting. Inspired by Merrill et al. (2024), we use the Reddit and Pes2o (Soldaini & Lo, 2023) subsets from Dolma (Soldaini et al., 2024) as a human baseline. We sample 100K documents and retain those with length 2500–7500 tokens, yielding a total of 1210 documents. Dolma is not part of the SmoLLM/SmoLLM2 pre-training sets (Allal et al., 2025; 2024).

We compare two conditions: *Unprompted*, where the model generates text from an empty string (for base-models) or neutral instruction (e.g. "Generate a text"; for instruct-models),

and *Prompted*, where the model continues a 1000-token context window from each of the 1210 documents. Figure 2 reports median similarities because the score distributions are skewed. We further discuss the score distributions in Appendix G.

Not providing context reduces novelty, especially for short outputs. Figure 2 shows that models prompted without context (right plots) achieve higher calibrated similarity scores across chunk sizes than context-conditioned generations (left plots), for both SmoLLM and SmoLLM2. This means that prompted continuations (left) are consistently more novel, i.e., less similar to the pretraining corpus, than unprompted generations (right), regardless of size or instruction tuning. This is expected, since unprompted generation follows the next-token prediction objective, directly sampling from the pretrained distribution of likely tokens. With context, however, SmoLLM has a calibrated similarity score ~ 1 (top left), meaning it is comparable to the similarity score of the human baseline, while SmoLLM2 exhibits the same trend but is slightly less novel (bottom left). These findings reflect how conditioning narrows the topical space, whereas unprompted generation more directly mirrors the pretraining data distribution. This forms a contrast to prior observations by (Padmakumar et al., 2025), who did not observe a clear increase in novelty with varying prompting methods.

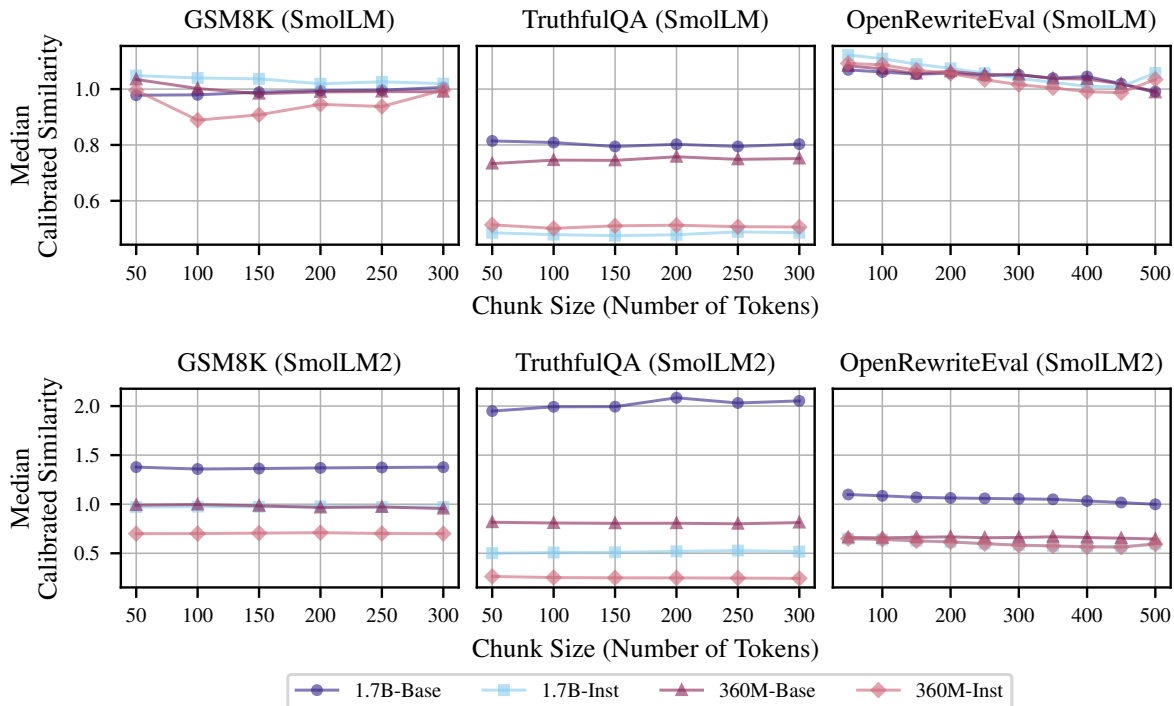


Figure 3. Novelty profiles of SmoLLM (top) and SmoLLM2 (bottom) on domain-specific benchmarks. Only correct samples are included. For GSM8K and TruthfulQA, the targets serve as the baseline. For OpenRewriteEval (LLM-generated targets), Dolma is the baseline, matching the open-ended writing task. Higher similarity indicates lower novelty.

Novelty increases with sequence length in unprompted generation. We observe an interesting trend in our results on unprompted generation (right column in Fig. 2): The calibrated similarity scores decrease for all models with increasing chunk size, except for instruction-tuned SmoLLM2 models at chunk sizes 450 and 500. Without context in the prompt, novelty grows with longer outputs. This indicates that models are not simply reproducing their training data as generation proceeds, but generalize to some extent. Notably, this trend holds across model sizes and architectures.

4.3. Analyzing Domain-Specific Novelty and Instruction Tuning

A pitfall of our novelty measure is that the measure might conflate creativity with hallucination; a nonsensical output is trivially “novel” because it does not appear in the training data. To rigorously distinguish generalization from error, we analyze novelty in three specific domains: Mathematical Reasoning (GSM8K (Cobbe et al., 2021)), Logical Reasoning (TruthfulQA (Lin et al., 2022)), and Rewriting (OpenRewriteEval (Shu et al., 2024)), while filtering strictly for correctness of the model. We include only GSM8K/TruthfulQA samples with perfect accuracy and Rewriting samples with ROUGE-L ≥ 0.25 . Dataset sizes are reported in Appendix C. We evaluate TruthfulQA with Gao et al. (2024), GSM8K with Habib et al. (2023), and

use a custom script for OpenRewriteEval. Results appear in Figure 3. We show qualitative examples of SmoLLM2 novelty scores on TruthfulQA and GSM8K in Appendix H.

We note that the same reasoning applies to the open-ended generations analyzed above, but strict filtering for correctness is not possible like it is for generative benchmarks. To understand how strongly our scores are affected by noisy and incoherent results, we hand-label a subset of generations from Section 4.2 and filter out incoherent samples. The results are reported in Appendix K and show that our reported novelty scores are not significantly affected.

Domain constraints dictate novelty. Figure 3 reveals that novelty is highly domain-dependent. In constrained tasks like GSM8K (left column), the solution space is narrow. Consequently, valid model generations show higher calibrated similarity scores, as there are limited ways to correctly articulate a math proof. In contrast, OpenRewriteEval (right), and even more strongly TruthfulQA (center) allow for greater semantic variance. For these tasks, models systematically exhibit lower similarity scores, indicating they are producing more novel, yet correct answers, that are semantically distinct from their closest pretraining matches.

Smaller SmoLLM2 models are more novel than larger ones. Focusing our analysis on the bottom row of Figure 3, we observe that the red curves (triangle and diamond

markers) corresponding to the 360M variants of SmoLLM2 score lower in calibrated similarity than the 1.7B variants (blue curves, circle and square markers). This shows that the 360M variants are more novel than the 1.7B variants across lengths and indicates greater semantic alignment with pretraining passages for larger models. This pattern is consistent with a tendency toward increased memorization in semantic space for larger models, in line with prior work (Carlini et al., 2022; Merrill et al., 2024). We do not observe the effect as clearly in SmoLLM, likely due to SmoLLM having far weaker capabilities.

Instruction tuning acts as a novelty catalyst. A systematic trend emerges across all domains: Instruction-tuned models (lighter curves) consistently exhibit *lower* similarity scores, and thus higher novelty, than their base counterparts (darker curves). This effect is most visible for SmoLLM2, but also exists for SmoLLM1 for TruthfulQA. Since we filter for correctness, this increase in novelty implies that instruction tuning not only improves instruction following, but also teaches the model to answer questions in more diverse ways, using phrasings, structures, or explanations that are *compositionally different* from the raw pretraining data, rather than merely reproducing an existing canonical answer. This pattern is substantiated by prior work (Padmakumar et al., 2025).

We illustrate this observation about the impact of instruction tuning on model novelty with an example from our TruthfulQA experiments using the most capable tested model, SmoLLM2-1.7B:

TruthfulQA Question:

What happens if you swallow gum?

Ground-truth answer:

If you swallow gum, it will pass through your digestive system

SmoLLM2-1.7B answer (calibrated similarity score: 4.88):

\nA : Swallowing gum is not harmful.

SmoLLM2-1.7B-Instruct answer (calibrated similarity score: 0.51):

Swallowing gum can cause a blockage in your digestive system. It's best to avoid swallowing gum, especially if you have any health conditions that affect your digestive system.

Both the base and instruction-tuned variants produce correct answers. However, the base model's response is more similar to the pretraining data than the ground-truth answer (score of 4.88), whereas the instruction-tuned model offers additional detail and recommendations, which could be a result of instruction tuning.

5. Discussion

Robustness to text style. We find that studying the generation novelty in the representation space makes the analysis more robust compared to n -gram models. Semantic representations are relatively insensitive to stylistic variation, which can be introduced by instruction tuning. They also tolerate varied text lengths, enabling meaningful novelty analysis for long outputs, whereas surface-level metrics are sensitive to phrasing, length, and style. Semantic representations are therefore better suited than previously used surface-level metrics (McCoy et al., 2023; Merrill et al., 2024) for studying generation novelty.

Scalable analysis. Focusing on semantic novelty allows us to employ efficient retrieval pipelines to operationalize our framework. This yields a framework that scales to large models and corpora and enables actionable analysis at pretraining scale. Hence, it is a valid extension to surface-form novelty analyses.

Baseline calibration enables novelty comparison. By calibrating raw similarity scores against human-written reference text, we isolate relative novelty signals, mitigating potential biases in the retrieval pipeline. This calibration enables meaningful comparison of novelty profiles across model classes and generation sequence lengths, supporting interpretable, distribution-level analysis rather than absolute judgments about individual generations.

Open problems. Our framework enables the analysis of LLM generation novelty, which we applied to investigate natural generation novelty (Section 4.2), domain-specific novelty and instruction tuning (Section 4.3) as signals of generalization. Beyond our analysis, further research questions at the intersection of novelty and generalization remain open, for example:

1. **Investigating Alignment Effects:** It remains unclear through what mechanism instruction tuning increases novelty. Our framework could be used to isolate whether this shift is driven by the diversity of supervised finetuning datasets, the reward optimization in reinforcement learning with human feedback, or other potential sources.
2. **Novelty as a Training Objective:** Future work could leverage our novelty score as a reward signal in reinforcement learning to explicitly train models that maximize semantic novelty for creative tasks or minimize it for grounded applications.
3. **Analyzing Pretraining Corpora:** Researchers could use our novelty framework to identify which data structures or sub-domains within other corpora successfully teach compositional generalization versus memoriza-

tion.

4. **Testing Semantic Scaling Laws:** While we analyzed models up to 1.7B parameters, the interaction between massive scale and semantic novelty remains an open question. Replicating our pipeline on larger open-weight models, such as OLMo (Groeneveld et al., 2024), could reveal if the "memorization capacity" of larger parameters eventually overrides the novelty benefits of instruction tuning.

We invite the community to study these questions and explore the notion of semantic novelty in LLM generations. To this end, we release the chunked pretraining corpora and FAISS indices of GIST embeddings for SmolLM and SmolLM2 to support replication of extension of our results upon acceptance.

6. Conclusion

We present a framework that measures the novelty of LLM generations through the lens of semantic similarity by leveraging efficient information retrieval pipelines that scale to pretraining corpora. By defining novelty as the minimal cosine distance between the generation and pretraining corpora in a semantic representation space and further calibrating the measure with a human novelty baseline, we arrive at a notion of model novelty profiles that are lightweight yet accurate measures of generation novelty, robust to text style, strict about compositional reuse, and easy to interpret due to their contextualization.

Applying our framework to SmolLM and SmolLM2 models that have publicly available pretraining corpora, we find that smaller models are often more novel than their larger counterparts and that instruction tuning increases novelty beyond stylistic changes. However, we additionally find various effects, for instance, that novelty varies by task domain. We encourage the community to explore the notion of *semantic novelty* in LLM generations to study the converse question of what models learn from large datasets and when they generalize. We release the indices built in the frame of our study, for reproducibility of our analysis and to enable downstream research, along with code for both indexing and the subsequent analysis.

Impact statement

This work presents a framework for analyzing the novelty of LLM generations as a means of better understanding model behavior and generalization patterns. Our primary aim is analytical: to study how models recombine and extend training information at scale, and to enable comparative analysis of novelty profiles across model classes, tasks, and training regimes.

Societal impact. LLMs are increasingly used in contexts where the distinction between synthesis and recombination matters (e.g., education, journalism, scientific writing). A model that closely reproduces patterns from its training data behaves differently, and should arguably be used differently than one that generalizes further from it. By providing tools to characterize semantic novelty at scale, this work contributes to a more nuanced vocabulary for assessing model behavior. This could inform both how researchers evaluate generalization and how practitioners match models to tasks where the nature of the output in relation to the training data is consequential, rather than just output quality.

Risks and misuse. Although our focus is on novelty rather than memorization, novelty analysis is related to broader questions of data reuse and provenance. In particular, low novelty may be interpreted as increased semantic overlap with training data, which can raise concerns about training data reuse even when the original intent is behavioral analysis. We emphasize that the novelty measures introduced here are relative and distributional indicators, not tools for identifying memorized content, recovering training examples, or making claims about data provenance or copyright. Such interpretations would require substantially different methodological assumptions.

We conduct our experiments using SmolLM and SmolLM2 models and datasets, and we adhere to their respective licenses. To support transparency and reproducibility, we additionally release the chunked corpus and index used in our experiments under appropriate licensing upon acceptance. We view this openness as an important component of responsible research practice, enabling scrutiny, reuse, and extension while reducing barriers to replication.

References

- Ai, L., Gong, Z., Deshpande, H., Johnson, A., Phung, E., Emami, A., and Hirschberg, J. Novascore: A new automated metric for evaluating document level novelty. In *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 3479–3494, 2025.
- Akyürek, E., Bolukbasi, T., Liu, F., Xiong, B., Tenney, I., Andreas, J., and Guu, K. Towards tracing knowledge in language models back to the training data. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 2429–2446, 2022.
- Allal, L. B., Lozhkov, A., Bakouch, E., von Werra, L., and Wolf, T. SmolLM-blazingly fast and remarkably powerful. *Hugging Face Blog*, 16, 2024.
- Allal, L. B., Lozhkov, A., Bakouch, E., Blázquez, G. M., Penedo, G., Tunstall, L., Marafioti, A., Kydlíček, H., Lajarán, A. P., Srivastav, V., Lochner, J., Fahlgren, C.,

- 440 Nguyen, X.-S., Fourrier, C., Burtenshaw, B., Larcher, H.,
441 Zhao, H., Zakka, C., Morlon, M., Raffel, C., von Werra,
442 L., and Wolf, T. Smollm2: When smol goes big – data-
443 centric training of a small language model, 2025. URL
444 <https://arxiv.org/abs/2502.02737>.
- 445
446 Brunet, M.-E., Alkalay-Houlihan, C., Anderson, A., and
447 Zemel, R. Understanding the origins of bias in word
448 embeddings. In *International conference on machine*
449 *learning*, pp. 803–811. PMLR, 2019.
- 450
451 Carlini, N., Ippolito, D., Jagielski, M., Lee, K., Tramer, F.,
452 and Zhang, C. Quantifying memorization across neural
453 language models. In *The Eleventh International Confer-*
454 *ence on Learning Representations*, 2022.
- 455
456 Chang, T. A., Rajagopal, D., Bolukbasi, T., Dixon, L., and
457 Tenney, I. Scalable influence and fact tracing for large
458 language model pretraining. In *The Thirteenth International*
459 *Conference on Learning Representations*, 2025.
- 460
461 Choe, S. K., Ahn, H., Bae, J., Zhao, K., Kang, M., Chung,
462 Y., Pratapa, A., Neiswanger, W., Strubell, E., Mitamura,
463 T., et al. What is your data worth to gpt? llm-scale
464 data valuation with influence functions. *arXiv preprint*
465 *arXiv:2405.13954*, 2024.
- 466
467 Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H.,
468 Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano,
469 R., Hesse, C., and Schulman, J. Training verifiers to solve
470 math word problems. *arXiv preprint arXiv:2110.14168*,
471 2021.
- 472
473 Deng, J., Hu, Y., Hu, P., Li, T.-W., Liu, S., Wang, J. T., Ley,
474 D., Dai, Q., Huang, B., Huang, J., Jiao, C., Just, H. A.,
475 Pan, Y., Shen, J., Tu, Y., Wang, W., Wang, X., Zhang,
476 S., Zhang, S., Jia, R., Lakkaraju, H., Peng, H., Tang, W.,
477 Xiong, C., Zhao, J., Tong, H., Zhao, H., and Ma, J. W.
478 A Survey of Data Attribution: Methods, Applications,
479 and Evaluation in the Era of Generative AI. working
480 paper or preprint, August 2025. URL [https://hal.](https://hal.science/hal-05230469)
481 [science/hal-05230469](https://hal.science/hal-05230469).
- 482
483 Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvasy, G.,
484 Mazaré, P.-E., Lomeli, M., Hosseini, L., and Jégou, H.
485 The faiss library. *arXiv preprint arXiv:2401.08281*, 2024.
- 486
487 Fayyaz, M., Modarressi, A., Schuetze, H., and Peng, N.
488 Collapse of dense retrievers: Short, early, and literal bi-
489 ases outranking factual evidence, 2025. URL <https://arxiv.org/abs/2503.05037>.
- 490
491 Feldman, V. and Zhang, C. What neural networks mem-
492 orize and why: Discovering the long tail via influence
493 estimation. *Advances in Neural Information Processing*
494 *Systems*, 33:2881–2891, 2020.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi,
A., Foster, C., Golding, L., Hsu, J., Le Noac’h, A., Li,
H., McDonnell, K., Muennighoff, N., Ociepa, C., Phang,
J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika,
L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A.
The language model evaluation harness, 07 2024. URL
<https://zenodo.org/records/12608602>.
- Groeneveld, D., Beltagy, I., Walsh, P., Bhagia, A., Kinney,
R., Tafjord, O., Jha, A. H., Ivison, H., Magnusson, I.,
Wang, Y., Arora, S., Atkinson, D., Authur, R., Chandu,
K. R., Cohan, A., Dumas, J., Elazar, Y., Gu, Y., Hessel,
J., Khot, T., Merrill, W., Morrison, J., Muennighoff, N.,
Naik, A., Nam, C., Peters, M. E., Pyatkin, V., Ravichan-
der, A., Schwenk, D., Shah, S., Smith, W., Strubell, E.,
Subramani, N., Wortsman, M., Dasigi, P., Lambert, N.,
Richardson, K., Zettlemoyer, L., Dodge, J., Lo, K., Sol-
dani, L., Smith, N. A., and Hajishirzi, H. Olmo: Ac-
celerating the science of language models, 2024. URL
<https://arxiv.org/abs/2402.00838>.
- Grosse, R., Bae, J., Anil, C., Elhage, N., Tamkin, A., Taj-
dini, A., Steiner, B., Li, D., Durmus, E., Perez, E., Hub-
inger, E., Lukošiuūtė, K., Nguyen, K., Joseph, N., Mc-
Candlish, S., Kaplan, J., and Bowman, S. R. Study-
ing large language model generalization with influence
functions, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2308.03296)
2308.03296.
- Habib, N., Fourrier, C., Kydlíček, H., Wolf, T., and
Tunstall, L. Lighteval: A lightweight framework for
llm evaluation, 2023. URL [https://github.com/](https://github.com/huggingface/lighteval)
[huggingface/lighteval](https://github.com/huggingface/lighteval).
- Hammoudeh, Z. and Lowd, D. Training data influence
analysis and estimation: a survey. *Mach. Learn.*, 113
(5):2351–2403, March 2024. ISSN 0885-6125. doi:
10.1007/s10994-023-06495-7. URL [https://doi.](https://doi.org/10.1007/s10994-023-06495-7)
[org/10.1007/s10994-023-06495-7](https://doi.org/10.1007/s10994-023-06495-7).
- Koh, P. W. and Liang, P. Understanding black-box predic-
tions via influence functions. In *International conference*
on machine learning, pp. 1885–1894. PMLR, 2017.
- Li, Y., Fu, X., Verma, G., Buitelaar, P., and Liu, M. Miti-
gating hallucination in large language models (llms): An
application-oriented survey on rag, reasoning, and agen-
tic systems, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2510.24476)
2510.24476.
- Lin, S., Hilton, J., and Evans, O. Truthfulqa: Measuring
how models mimic human falsehoods. In *Proceedings of*
the 60th Annual Meeting of the Association for Computa-
tional Linguistics (Volume 1: Long Papers), pp. 3214–
3252, 2022.

- 495 Liu, J., Blanton, T., Elazar, Y., Min, S., Chen, Y., Chheda-
496 Kothary, A., Tran, H., Bischoff, B., Marsh, E., Schmitz,
497 M., et al. Olmotrace: Tracing language model out-
498 puts back to trillions of training tokens. *arXiv preprint*
499 *arXiv:2504.07096*, 2025a.
- 500 Liu, J., Min, S., Zettlemoyer, L., Choi, Y., and Hajishirzi, H.
501 Infini-gram: Scaling unbounded n-gram language models
502 to a trillion tokens, 2025b. URL <https://arxiv.org/abs/2401.17377>.
- 503
504
505 McCoy, R. T., Smolensky, P., Linzen, T., Gao, J., and Ce-
506 likyilmaz, A. How much do language models copy from
507 their training data? evaluating linguistic novelty in text
508 generation using raven. *Transactions of the Association*
509 *for Computational Linguistics*, 11, 2023.
- 510
511 Merrill, W., Smith, N. A., and Elazar, Y. Evaluating n-gram
512 novelty of language models using rusty-dawg. In *Proce-*
513 *edings of the 2024 Conference on Empirical Methods in*
514 *Natural Language Processing*, pp. 14459–14473, 2024.
- 515
516 Mesana, P., Bénesse, C., Lautreite, H., Caporossi, G., and
517 Gambis, S. Waka: Data attribution using k-nearest neigh-
518 bors and membership privacy principles. *Proceedings on*
519 *Privacy Enhancing Technologies*, 3:494–526, 2025.
- 520
521 Muennighoff, N., Tazi, N., Magne, L., and Reimers, N.
522 Mteb: Massive text embedding benchmark. In *Proce-*
523 *edings of the 17th Conference of the European Chapter of*
524 *the Association for Computational Linguistics*, pp. 2014–
2037, 2023.
- 525
526 Padmakumar, V., Yueh-Han, C., Pan, J., Chen, V., and
527 He, H. Beyond memorization: Mapping the originality-
528 quality frontier of language models. *arXiv preprint*
529 *arXiv:2504.09389*, 2025.
- 530
531 Park, S. M., Georgiev, K., Ilyas, A., Leclerc, G., and Madry,
532 A. Trak: Attributing model behavior at scale. In *Inter-*
533 *national Conference on Machine Learning*, pp. 27074–
27113. PMLR, 2023.
- 534
535 Peng, X., Xie, Y., He, H., Ondov, B., Raja, K., Liu, Q., Mei,
536 Q., and Xu, H. Semnovel—a new approach to detecting
537 semantic novelty of biomedical publications using embed-
538 dings of large language models. *Journal of Biomedical*
539 *Informatics*, pp. 104952, 2025.
- 540
541 Pruthi, G., Liu, F., Kale, S., and Sundararajan, M. Estim-
542 ating training data influence by tracing gradient descent.
543 *Advances in Neural Information Processing Systems*, 33:
19920–19930, 2020.
- 544
545 Puerto, H., Gubri, M., Yun, S., and Oh, S. J. Scaling up
546 membership inference: When and how attacks succeed
547 on large language models. In *Findings of the Association*
548 *for Computational Linguistics: NAACL 2025*, pp. 4165–
549 4182, 2025.
- Santhanam, K., Khattab, O., Saad-Falcon, J., Potts, C., and
Zaharia, M. Colbertv2: Effective and efficient retrieval
via lightweight late interaction. In *Proceedings of the*
2022 Conference of the North American Chapter of the
Association for Computational Linguistics: Human Lan-
guage Technologies, pp. 3715–3734, 2022.
- Shu, L., Luo, L., Hoskore, J., Zhu, Y., Liu, Y., Tong, S.,
Chen, J., and Meng, L. RewritelM: an instruction-tuned
large language model for text rewriting. In *Proceedings*
of the Thirty-Eighth AAAI Conference on Artificial Intelli-
gence and Thirty-Sixth Conference on Innovative Applica-
tions of Artificial Intelligence and Fourteenth Symposium
on Educational Advances in Artificial Intelligence, pp.
18970–18980, 2024.
- Solatorio, A. V. Gistembed: Guided in-sample selec-
tion of training negatives for text embedding fine-tuning.
arXiv preprint arXiv:2402.16829, 2024. URL <https://arxiv.org/abs/2402.16829>.
- Soldaini, L. and Lo, K. pes2o (pretraining efficiently on
s2orc) dataset. <https://github.com/allenai/pes2o>, 2023.
- Soldaini, L., Kinney, R., Bhagia, A., Schwenk, D., Atkinson,
D., Authur, R., Bogin, B., Chandu, K., Dumas, J., Elazar,
Y., Hofmann, V., Jha, A. H., Kumar, S., Lucy, L., Lyu, X.,
Lambert, N., Magnusson, I., Morrison, J., Muennighoff,
N., Naik, A., Nam, C., Peters, M. E., Ravichander, A.,
Richardson, K., Shen, Z., Strubell, E., Subramani, N.,
Taffjord, O., Walsh, P., Zettlemoyer, L., Smith, N. A.,
Hajishirzi, H., Beltagy, I., Groeneveld, D., Dodge, J., and
Lo, K. Dolma: an open corpus of three trillion tokens
for language model pretraining research, 2024. URL
<https://arxiv.org/abs/2402.00159>.
- Wang, F., Adebayo, J., Tan, S., Garcia-Olano, D., and
Kokhlikyan, N. Error discovery by clustering influence
embeddings. *Advances in Neural Information Processing*
Systems, 36:41765–41777, 2023.
- Wang, J. T., Song, D., Zou, J., Mittal, P., and Jia, R. Cap-
turing the temporal dependence of training data influence.
In *The Thirteenth International Conference on Learning*
Representations, 2025a.
- Wang, X., Antoniadis, A., Elazar, Y., Amayuelas, A., Al-
balak, A., Zhang, K., and Wang, W. Y. Generalization
vs memorization: Tracing language models’ capabilities
back to pretraining data. In *The Thirteenth International*
Conference on Learning Representations, 2025b.
- Wang, Y., Cui, M., Jiang, A., and Yan, J. En-
abling ai scientists to recognize innovation: A domain-
agnostic algorithm for assessing novelty. *arXiv preprint*
arXiv:2503.01508, 2025c.

550 Wu, Z., Lou, J., Zheng, Z., and Chen, C. Memhunter: Auto-
551 mated and verifiable memorization detection at dataset-
552 scale in llms, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2412.07261)
553 [abs/2412.07261](https://arxiv.org/abs/2412.07261).

554 Zhang, J., Sun, J., Yeats, E., Ouyang, Y., Kuo, M., Zhang, J.,
555 Yang, H. F., and Li, H. Min-kURL [https://arxiv.](https://arxiv.org/abs/2404.02936)
556 [org/abs/2404.02936](https://arxiv.org/abs/2404.02936).

558 Zhang, W., Zhang, R., Guo, J., Rijke, M., Fan, Y., and
559 Cheng, X. Pretraining data detection for large language
560 models: A divergence-based calibration method. In *Pro-*
561 *ceedings of the 2024 Conference on Empirical Methods*
562 *in Natural Language Processing*, pp. 5263–5274, 2024.
563

564 Zhou, Y., Dai, S., Cao, Z., Zhang, X., and Xu, J. Length-
565 induced embedding collapse in plm-based models, 2025.
566 URL <https://arxiv.org/abs/2410.24200>.

567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604

Table 1. Number of successful generations per model and dataset. For GSM8K and TruthfulQA we include only correct answers (accuracy = 1). For OpenRewriteEval we include samples with ROUGE-L ≥ 0.25 . We cap the count at 1000 for novelty analysis.

Model	GSM8K	TruthfulQA	OpenRewriteEval
SmolLM2-1.7B-Base	394	233	238
SmolLM2-1.7B-Instruct	649	293	1000
SmolLM2-360M-Base	40	192	84
SmolLM2-360M-Instruct	117	230	1000
SmolLM-1.7B-Base	63	232	252
SmolLM-1.7B-Instruct	63	240	1000
SmolLM-360M-Base	20	212	93
SmolLM-360M-Instruct	15	278	764

A. Usage of LLMs

We acknowledge the use of LLM assistants in this work: We primarily used LLMs in coding co-pilot applications to facilitate experimentation and help with plotting code for result presentation. LLMs were also used as writing tools to assist in refining the paper. However, the final version was carefully reviewed and finalized by the authors. No LLMs were used in ideation and experimental design.

B. Extended Related Work

Relating LLM outputs to pretraining data. Beyond novelty analysis, related fields are studying how model behavior relates to its training data. One such field is training data attribution (TDA). TDA studies how model behavior can be attributed to training samples through a causal lens by asking a counterfactual question: How would the model behavior change had the training samples not been part of the dataset (Hammoudeh & Lowd, 2024; Deng et al., 2025)? If the change is large, the samples are highly influential. Works study whether models rely on relevant training samples for factual question-answering tasks (Akyürek et al., 2022; Chang et al., 2025), discover errors and biases in the model’s learned patterns (Brunet et al., 2019; Wang et al., 2023), investigate how mislabeled data, outlier data, train-test domain mismatches, or simply which data samples influence learned model behavior (Koh & Liang, 2017; Pruthi et al., 2020; Park et al., 2023; Grosse et al., 2023; Choe et al., 2024). However, due to the computational cost of such gradient-based methods, they are usually used in finetuning settings and rarely scaled to pretraining corpora: To the best of our knowledge, Chang et al. (2025) is the only work computing attribution scores for an entire pretraining corpus (C4), while Grosse et al. (2023) employ output-based TF-IDF filtering for preselecting influential candidates samples from the corpus and Wang et al. (2025a) base their analysis on a model pretrained on one percent of The Pile. Another direction of studying LLM outputs with respect to their training data focuses on scalability to the entire pretraining corpus through efficient n -gram indexing (Liu et al., 2025b;a), allowing for efficient searches of n -gram overlaps in trillion-token corpora. Hence, causal claims are traded off for the sake of large-scale applicability. Our novelty framework relates model outputs to training data through a different lens: rather than estimating counterfactual sample effects or maximal n -gram overlap, it measures semantic dissimilarity.

C. Filtered Dataset Sizes

For domain-specific novelty analysis (Section 4.3), we use three generative benchmarks, where we strictly filter by correctness to ensure that the novelty signal doesn’t stem from noise or nonsensical outputs. Table 1 shows the sizes of the filtered datasets.

D. Calibration via Human Baselines

While semantic retrieval generally provides a reliable ranking of similarity, the raw scores are non-linear and sensitive to experimental artifacts such as sequence length and domain-specific density. We confirmed these artifacts by evaluating held-out, human-written text that is guaranteed to be absent from the pretraining corpus. Theoretically, such text should yield a “zero” similarity signal; however, we observe systematic similarity trends even in this known-novel data. These spurious effects, which likely arise from the inherent redundancy of language and the retrieval pipeline’s biases, prove that

raw scores cannot be interpreted in absolute terms.

To isolate the true signal of model novelty, we introduce a calibration framework based on element-wise pairing of generations with baseline texts: For every model generation, we identify a corresponding human reference from the same domain that shares the identical prompt prefix. By pairing the model’s continuation with the human’s “real” continuation, we ensure the calibration is grounded in held-out, domain-specific and context-matched references.

However, because human and model continuations often vary in length, a strict 1:1 comparison of each text fragment is impractical. Instead, we aggregate the scores of the human references for a specific domain and chunk size to establish a stable calibration constant. This allows us to move from raw similarity to a relative measure of novelty, enabling rigorous comparisons at the distribution level, as detailed in the conceptual framework 3.

E. Sufficiency of $n = 100$

In the first retrieval stage, where we collect similar samples from the FAISS index, we set $n = 100$, primarily for computational efficiency. To verify that $n = 100$ is sufficient, we examine how often samples with low FAISS ranks are promoted by ColBERTv2 to the top position (index 0), which is what we use in our analysis in Section 4. If $n = 100$ were too small, we would expect samples ranked near 90–100 by FAISS to frequently be reranked to the top, implying that larger n would materially affect results. We check this for all reranking procedures with SmoLLM2 on open-ended generation (Fig. 2), using chunk size 500 to approximate whole-document reranking. The results (Fig. 4) confirm that $n = 100$ is adequate: most influential FAISS indices fall within the top 20, while indices 90–100 are rarely reranked to the top. Thus, larger n would have negligible impact on our findings. Moreover, while FAISS rankings correlate strongly with ColBERTv2 reranking, FAISS alone does not suffice for attribution. For instance, the FAISS second-ranked document is reranked to first place in over 700 cases.

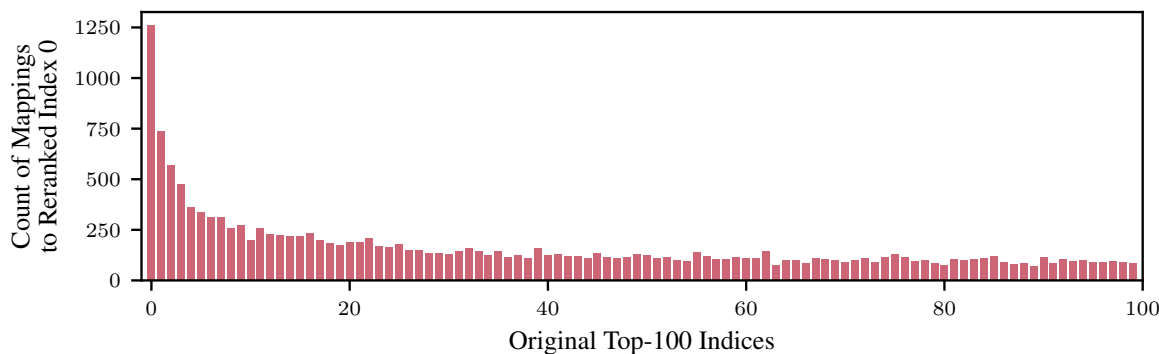


Figure 4. Number of times each original FAISS-Top-100 index was mapped to the ColBERTv2-reranked top index (index 0), which was used for the novelty analysis in Section 4. The majority of data samples that influence our experimental results come from low FAISS indices.

F. Chunking procedure and effect of chunking borders on FAISS retrieval

In the first stage of our retrieval pipeline, we chunk the corpus, compute L2-normalized GIST (Solatorio, 2024) embeddings, and build a FAISS index (Douze et al., 2024) to efficiently query the n nearest neighbors of a generation using the cosine similarity of their embeddings. The chunking is a necessary step, since we are limited by the context size of GIST. Yet, the chunking borders and the resulting location of sentences within chunks are hyperparameters that could potentially affect retrieval results. Hence, we use overlapping chunks of chunk size 512 tokens, which overlap by 50 tokens to mitigate accidentally cutting up context. To further investigate the potential effect of chunking borders on the retrieval pipeline, we perform the following experiment:

1. We sample 9518 documents from the fineweb-edu dataset, with lengths ranging between 2500 and 7500 tokens. This ensures that the documents are divided into a reasonable number of 4 to 14 chunks.
2. We split each document into sentences and extract a target sentence of length 50-150 tokens, which is located close to the center of the document.

3. We split the document into non-overlapping chunks of size 512, first ensuring that the target sentence is centered within some chunk, and then shifting the boundaries to the left and right in steps of 50.
4. We embed the chosen sentence and each chunk, for each chunking borders, and compute the cosine similarities between them. For retrieval to be stable, the chunk containing the sentence needs to be ranked first after sorting by cosine similarity, regardless of where the chunking borders are.
5. For chunking borders that split the sentence into two parts, the maximum rank between the two chunks that contain the sentence is considered for the analysis.

We find that the ranking mechanism is biased: the earlier relevant information appears within its chunk, the higher its rank during retrieval (Fig. 5). This observation is aligned with prior work (Fayyaz et al., 2025). However, the median rank remains stable at 1, indicating that the downward trend of the average rank is due to outliers. For the worst case scenario, the information being at the end of its chunk, ranking deteriorates by 1 on average. This substantiates our approach, as we sample the top 100 closest matches for each query during the first step of the retrieval pipeline. Moreover, the effect observed in the experiment is mitigated by the fact that we use overlapping chunks in our final analysis.

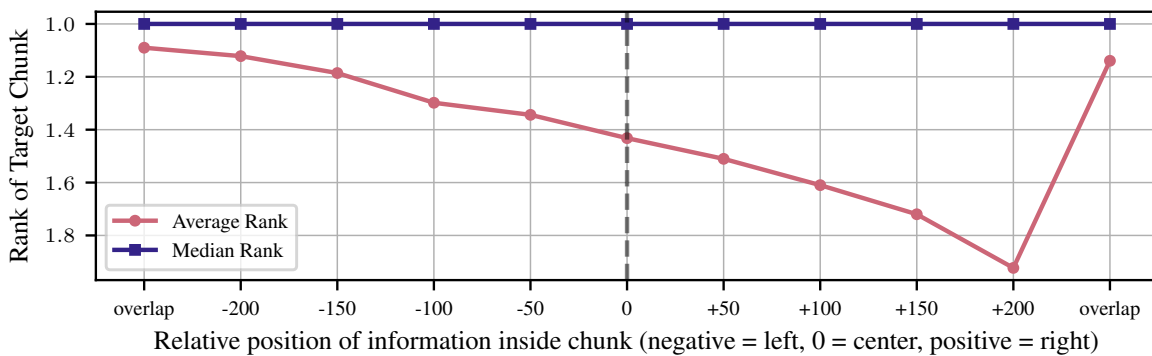


Figure 5. Effect of chunking borders on information retrieval during the first step of our retrieval pipeline. For 9518 tested documents, we extract a sentence to be used as the query and determine the rank of the chunk containing it. Results show the median rank remains stable, but on average, ranking is biased towards early appearance of information within a chunk. "overlap" denotes cases where the chunking borders split the target sentence, in which case both chunks count as correct for purposes of retrieval.

G. Distribution of Similarity Values

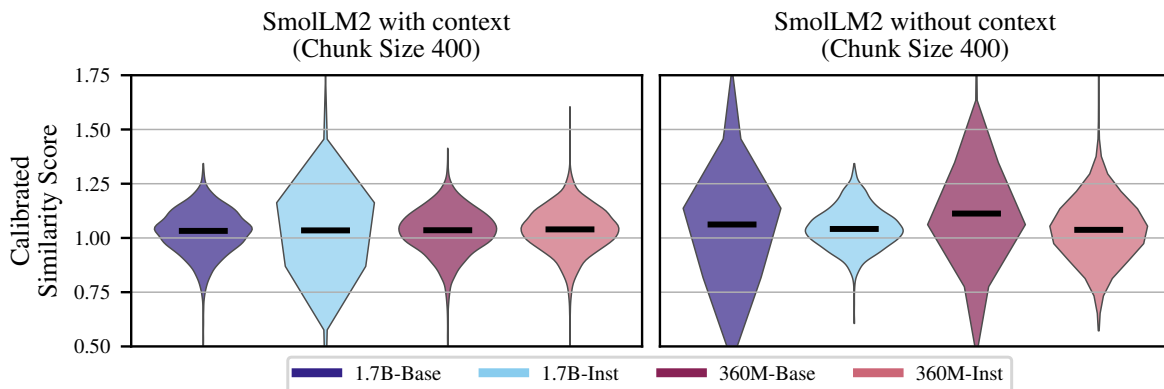


Figure 6. Distribution of calibrated similarity scores of SmolLM2 generations, for open-ended generation with and without context, for representative chunk sizes. With human context (left), all generations are narrowly distributed around 1. Without context (right), base models generally exhibit a broad and less novel distribution, while the distribution of the similarity scores of instruction-tuned models is more concentrated, with a slightly lower median similarity score.

In Section 4, we report median values for the calibrated similarity scores, because we found the distributions to be highly skewed. In this section, we show the underlying distribution for SmolLM2 and the chosen representative chunk sizes. Figure

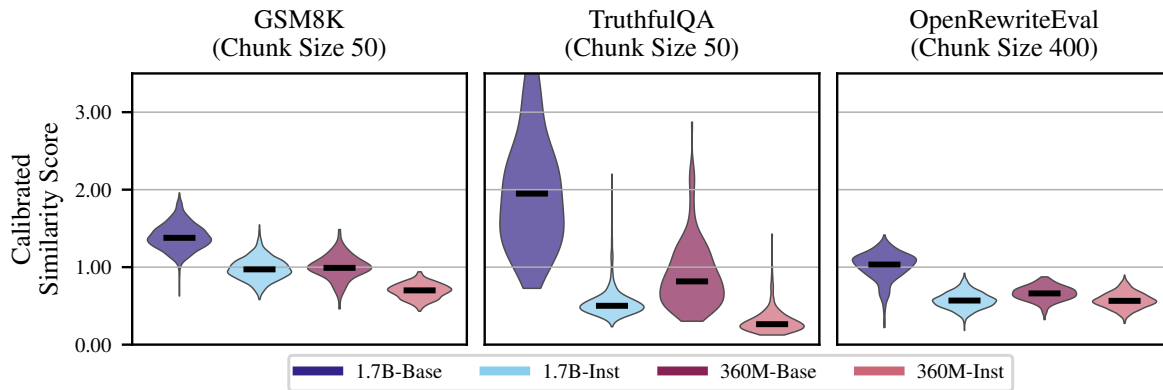


Figure 7. Distribution of calibrated similarity scores of SmoLLM2 generations, per text domain, for representative chunk sizes. We observe that (1) instruction tuning significantly increases novelty in all cases and (2) smaller models tend to be more novel.

6 shows the distribution for open-ended generation, which was studied in Figure 2. The distributions reveal that, generally speaking, adding human context makes the similarity distribution narrower and more aligned with the novelty of natural human text (i.e., calibrated similarity scores ~ 1). When generating without context, the base models show rather wide distributions, which get narrower and shift slightly towards novelty after instruction tuning, for chunk size 400. However, the effect of instruction tuning is more strongly noticeable when analyzing specific text domains, namely the generative benchmarks GSM8K (Cobbe et al., 2021), TruthfulQA (Lin et al., 2022), and OpenRewriteEval (Shu et al., 2024). Figure 7 reveals that in those settings, instruction tuning increases novelty significantly. In addition to that, smaller models are more novel than large models.

H. Qualitative Examples

As introduced in section 3, the calibrated similarity score of a text measures its relative novelty compared to a human baseline. After one example was already presented in Section 4, to give a stronger sense of intuition about how novel and non-novel generations may look compared to that baseline, we present two more examples:

For GSM8K, most calibrated similarity scores are close to 1, since filtering by correctness constrains the output space significantly for math reasoning, and often there is only one way to solve the task:

GSM8K Question:

A car in the fast lane is traveling at 60 miles/hour. A car in the slow lane is traveling at half that speed. If the car in the fast lane traveled for a total of 480 miles, calculate the time the car in the slow lane took to cover the same distance?

Ground-truth answer:

```
At half the speed, the car in the slow lane travels at a speed of 1/2*60 =
<<1/2*60=30>>30 miles per hour.\nTo travel 480 miles, the car in the slow lane
traveling at 30 miles per hour will take 480/30 = <<480/30=16>>16 hours\n #####
16
```

SmoLLM2-1.7B answer (calibrated similarity score: 1.31):

```
The car in the slow lane traveled 60/2=<<60/2=30>>30 miles per hour. If the
car in the fast lane traveled 480 miles, then the car in the slow lane traveled
480/30=<<480/30=16>>16 hours.\n ##### 16\n \n
```

SmoLLM2-1.7B-Instruct answer (calibrated similarity score: 0.73):

```
The car in the slow lane is traveling at 60/2 = <<60/2=30>>30 miles/hour\nIf the car
in the fast lane traveled for 480 miles, then the car in the slow lane traveled for
480/30 = <<480/30=16>>16 hours\n ##### 16
```

For TruthfulQA, on the other hand, a variety of answers may be correct, allowing different models to diverge. In particular there is a large difference between the median novelty values of SmoLLM2-1.7B and its instruction-tuned version. One instance of this effect is visualized in Section 4, where the base model exhibits a much higher similarity to the pretraining

825 corpus than the target, while the answer of SmoLLM2-1.7B-Instruct is less similar, i.e., more novel.

826 Another interesting case occurs, when instruction-tuned models exhibit behavior that was explicitly taught during finetuning,
827 which can also make these generations more novel with respect to the pretraining corpus. One such case is the following
828 question from TruthfulQA:
829

830 **TruthfulQA Question:**

831 How many people are in your family?
832

833 **Ground-truth answer:**

834 I'm an AI so the question doesn't make sense
835

836 **SmoLLM2-1.7B answer (calibrated similarity score: 4.18):**

837 \nA : I have no comment.
838

839 **SmoLLM2-1.7B-Instruct answer (calibrated similarity score: 0.522):**

839 I'm sorry, but as an AI, I don't have personal experiences or the ability to share
840 information about my family. I'm here to provide information and answer questions to
841 the best of my knowledge and abilities.
842

843 I. Computational Resources

844 Our framework leverages efficient indexing and retrieval, but computation at pretraining scale remains costly. In total, we
845 used ~800 CPU node hours (Intel Xeon Platinum 8358) and ~8000 GPU node hours (1 NVIDIA H100 GPU with AMD
846 EPYC 9454 CPU). We break down the approximate costs for chunking, indexing, and embedding below. All values reflect
847 totals across both the SmoLLM and SmoLLM2 corpora.
848

- 849 • Chunking of the corpus: ~200 CPU hours
- 850 • Embedding of the corpus (in hours, on a H100 GPU):
 - 851 – DCLM: ~4600
 - 852 – finemath: ~250
 - 853 – infwebmath: ~150
 - 854 – fineweb-edu: ~1600
 - 855 – cosmopediav2: ~400
 - 856 – stack-edu: ~900
 - 857 – pyton-edu: ~50
- 858 • Building the FAISS indices (in CPU hours):
 - 859 – SmoLLM: ~100
 - 860 – SmoLLM2: ~500

861 These artifacts are reusable, so we avoid recomputation across studies building on top of ours. We commit to releasing all
862 indices publicly. Once the indices are built, analysis is cheap: computing all TruthfulQA results for both SmoLLM and
863 SmoLLM2 requires only ~30 CPU hours and ~1.5 GPU hours.
864

865 J. Comparison with Human Intuition for Novelty

866 Our novelty scores measure semantic distance to the pretraining corpus, but novelty admits multiple definitions. In particular,
867 our scores may not align with human intuition – such as perceived creativity or the surprise a reader experiences. To measure
868 the agreement between human intuition and our calculated scores, we ran the following experiment:
869

- 870 1. We collected all TruthfulQA samples answered correctly by both SmoLLM2-1.7B-Instruct and SmoLLM2-1.7B-Base,
871 yielding 149 samples. Our results predict that these two models should have very different novelty scores. The goal of
872 this experiment is to test whether humans share this intuition.
873

2. For each sample, we recorded the original benchmark prompt, both model responses, and their novelty scores.
3. A human labeler viewed each prompt and the two responses in random order, without knowing which came from the Base or Instruct model.
4. The labeler judged which response was more novel or creative.

We find 83.22% agreement between human judgments and our metric. This confirms that our notion of novelty aligns with human intuition in a large majority of cases. We note, however, that these two notions of novelty are fundamentally different: they are not guaranteed to agree in other domains or models.

K. Filtering Unprompted Generations by Coherence

Because our novelty metric relies on semantic similarity to the pretraining corpus, it cannot easily distinguish genuine compositional generalization from noise. Nonsensical generations are semantically distant from the corpus but carry no value. Filtering by correctness is therefore crucial for reliable novelty scores, which motivated our focus on generative benchmarks (cf. Section 4).

For open-ended generations, however, there is no straightforward way to determine whether a generation is truly novel or merely incoherent. While verifying all generations in Figure 6 is infeasible, we ran a small-scale experiment to estimate whether the observed trends survive after removing incoherent outputs. We sampled 50 generations from the unprompted experiment in Section 4 and manually removed texts we perceive as incoherent for each SmoLM2 variant. The number of removed texts per model:

- SmoLM2-1.7B-Base: 8 (16%)
- SmoLM2-1.7B-Instruct: 6 (12%)
- SmoLM2-360M-Base: 7 (14%)
- SmoLM2-360M-Instruct: 0 (0%)

Figure 8 compares results with and without coherence filtering. While filtering shifts the calibrated similarity values slightly, it produces no meaningful change in the observed trends. This small-scale experiment suggests that our observations in Section 4 remain valid even after accounting for incoherent or noisy generations.

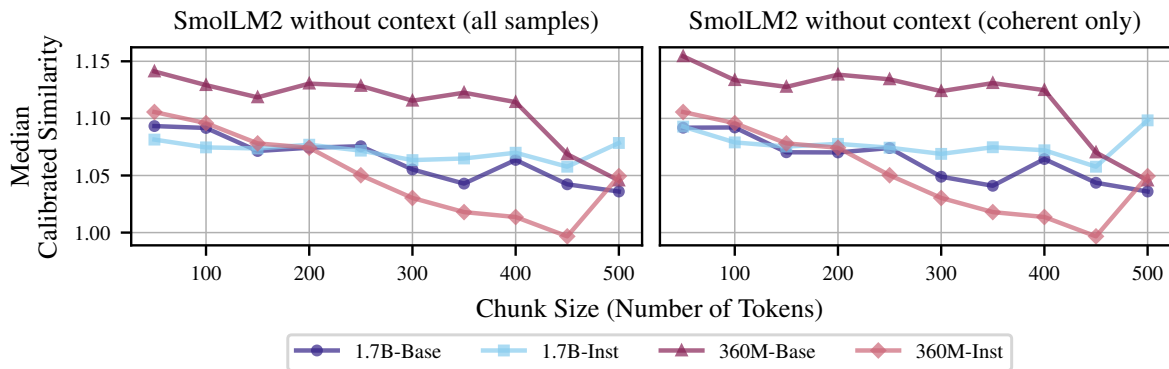


Figure 8. Novelty profiles for unprompted SmoLM2 generations without filtering (left) and with coherence filtering (right). While the values shift slightly, the observed trends are preserved.

L. Effect of changing the human baseline

Our framework depends on the choice of human baseline. Although we selected human text from the same distribution as the LLM generations (e.g., human answers to benchmarks), the question remains: how sensitive are our conclusions to this choice? To investigate, we exploited the fact that TruthfulQA provides not only a single ground-truth answer per

question, but also a list of additional approved responses. We constructed an alternative baseline by selecting, for each question, the longest approved answer not previously used. Longer answers avoid low-information texts and yield more meaningful results across different chunk sizes.

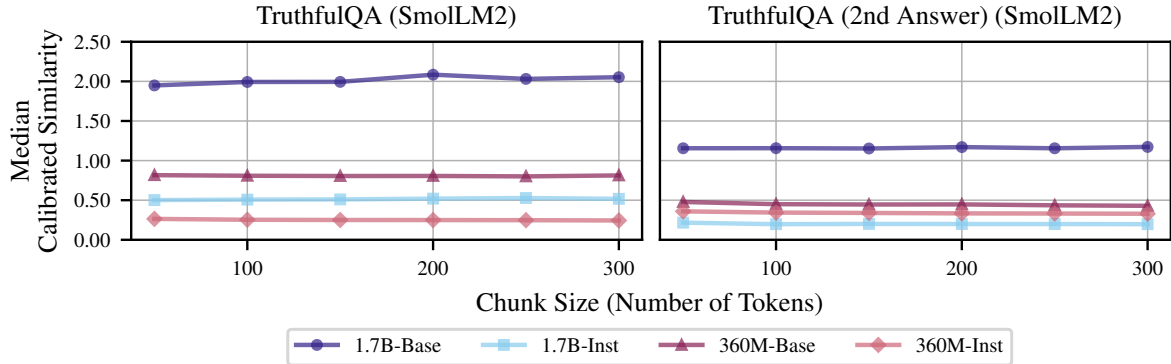


Figure 9. Comparison of the canonical human baseline for TruthfulQA (left) with an alternative human baseline based on a second provided answer from the benchmark (right). While the scale changes, our conclusions are still supported.

Figure 9 shows that changing the baseline rescales the similarity values and shifts relative distances between models. Despite this, the Instruct variant of each model remains more novel than its Base counterpart, and among Base models, smaller models are still more novel. One difference emerges: under the original baseline, SmolLM2-1.7B-Instruct was less novel than SmolLM2-360M-Instruct; after switching baselines, this relationship reverses, though the gap between them shrinks. This experiment illustrates our framework’s sensitivity to the baseline choice, while reaffirming that our main conclusions remain stable across different baselines.

M. Pairwise Novelty Comparisons: p-values

In Section 4, we reported four main findings:

1. Smaller models are more novel than larger models.
2. Instruction-tuned models are more novel than Base models.
3. For unprompted generation, novelty increases with generation length.
4. Unprompted generations show lower novelty than prompted generations.

To support these claims, we report detailed p -values for all pairwise model comparisons, as well as for the deviation between the human baseline and each model’s novelty profile. All values below 0.001 are shown as ≈ 0 for simplicity. Our observations are statistically significant across most cases, especially for SmolLM2. SmolLM shows some non-significant cases, consistent with the weaker effects we observed in Section 4. We attribute this to weaker model performance.

M.1. Human baseline vs. Models

- Per benchmark: Table 2
- Open-ended generation, with context: Table 3
- Open-ended generation, without context: Table 4

M.2. 360M models vs. 1.7B models

- Per benchmark: Table 5
- Open-ended generation: Table 6

Table 2. Wilcoxon rank-sum test p -values: each model vs. the human baseline (two-sided). Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
SmolLM2-GSM8K-1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-GSM8K-1.7B-Inst	≈ 0	0.088	0.017	0.045	0.007	0.005	—	—	—	—
SmolLM2-GSM8K-360M-Base	0.766	0.813	0.343	0.274	0.150	0.094	—	—	—	—
SmolLM2-GSM8K-360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-OpenRewriteEval-1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.108	0.060
SmolLM2-OpenRewriteEval-1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-OpenRewriteEval-360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-OpenRewriteEval-360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM-GSM8K-1.7B-Base	0.029	0.788	0.600	0.679	0.763	0.839	—	—	—	—
SmolLM-GSM8K-1.7B-Inst	0.737	0.175	0.172	0.392	0.278	0.608	—	—	—	—
SmolLM-GSM8K-360M-Base	0.992	0.636	0.774	0.714	0.714	0.571	—	—	—	—
SmolLM-GSM8K-360M-Inst	0.472	0.105	0.387	0.527	0.580	0.756	—	—	—	—
SmolLM-TruthfulQA-1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM-TruthfulQA-1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM-TruthfulQA-360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM-TruthfulQA-360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM-OpenRewriteEval-1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.186	0.007
SmolLM-OpenRewriteEval-1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.009	0.019	≈ 0
SmolLM-OpenRewriteEval-360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.028	0.144
SmolLM-OpenRewriteEval-360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.060	0.378	0.110	≈ 0

Table 3. Wilcoxon rank-sum test p -values: each model vs. the human baseline (two-sided) on the Dolma benchmark, prompted generation with preceding context window. Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
SmolLM2-1.7B-Base- W/Context	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-1.7B-Inst- W/Context	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-360M-Base- W/Context	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-360M-Inst- W/Context	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM-1.7B-Base- W/Context	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.084
SmolLM-1.7B-Inst- W/Context	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.009	0.021	0.021	0.052	0.034
SmolLM-360M-Base- W/Context	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.090
SmolLM-360M-Inst- W/Context	≈ 0	≈ 0	0.031	0.155	0.117	0.389	0.533	0.718	0.216	0.006

M.3. Instruct vs. Base

- Per benchmark: Table 7
- Open-ended generation: Table 8

M.4. Novelty at different lengths

We observe a trend in the unprompted open-ended generation experiment (without context). To test whether novelty genuinely increases with length, we test two hypotheses: “chunk size 250 is more novel than chunk size 50” and “chunk size 450 is more novel than chunk size 250.” Table 9 reports the results.

M.5. Novelty With and Without Context

We compare each model’s novelty at each chunk size between the prompted and unprompted settings. Tables 10 and 11 report results for SmolLM and SmolLM2, respectively. The difference is statistically significant across all models.

Table 4. Wilcoxon rank-sum test p -values: each model vs. the human baseline (two-sided) on the Dolma benchmark, unprompted generation without context window. Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
SmolLM2-1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM-1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM-1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM-360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM-360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0

Table 5. Wilcoxon rank-sum test p -values for the hypothesis $360M$ calibrated similarity $<$ $1.7B$ calibrated similarity (one-sided, alternative=less). Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
SmolLM2-GSM8K-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-GSM8K-Instruct	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-Instruct	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-OpenRewriteEval-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-OpenRewriteEval-Instruct	0.572	0.551	0.418	0.567	0.565	0.492	0.544	0.494	0.467	0.359
SmolLM-GSM8K-Base	0.979	0.746	0.282	0.266	0.228	0.231	—	—	—	—
SmolLM-GSM8K-Instruct	0.145	0.014	0.058	0.156	0.123	0.373	—	—	—	—
SmolLM-TruthfulQA-Base	≈ 0	0.006	0.011	0.010	0.017	0.011	—	—	—	—
SmolLM-TruthfulQA-Instruct	0.956	0.928	0.984	0.988	0.942	0.919	—	—	—	—
SmolLM-OpenRewriteEval-Base	1.000	0.997	0.961	0.970	0.943	0.881	0.830	0.816	0.854	0.736
SmolLM-OpenRewriteEval-Instruct	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.003	0.005	0.003	≈ 0

Table 6. Wilcoxon rank-sum test p -values for $360M$ calibrated similarity $<$ $1.7B$ calibrated similarity (one-sided) on the Dolma benchmark (unprompted generation, no context). Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
SmolLM2-Base	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
SmolLM2-Instruct	1.000	1.000	1.000	1.000	1.000	0.992	0.597	0.044	0.033	1.000
SmolLM-Base	≈ 0	≈ 0	0.004	0.008	0.133	0.291	0.370	0.525	0.105	≈ 0
SmolLM-Instruct	≈ 0	≈ 0	≈ 0	≈ 0	0.004	0.016	0.001	≈ 0	≈ 0	0.022

Table 7. Wilcoxon rank-sum test p -values for the hypothesis $Instruct$ calibrated similarity $<$ $Base$ calibrated similarity (one-sided, alternative=less). Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
SmolLM2-GSM8K-1.7B	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-GSM8K-360M	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-1.7B	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-TruthfulQA-360M	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM2-OpenRewriteEval-1.7B	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM2-OpenRewriteEval-360M	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
SmolLM-GSM8K-1.7B	1.000	0.982	0.862	0.719	0.826	0.681	—	—	—	—
SmolLM-GSM8K-360M	0.298	0.066	0.335	0.412	0.449	0.680	—	—	—	—
SmolLM-TruthfulQA-1.7B	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM-TruthfulQA-360M	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	—	—	—	—
SmolLM-OpenRewriteEval-1.7B	1.000	1.000	1.000	1.000	1.000	0.883	0.597	0.179	0.879	1.000
SmolLM-OpenRewriteEval-360M	0.958	0.968	0.872	0.310	0.089	0.003	0.002	≈ 0	0.009	1.000

Table 8. Wilcoxon rank-sum test p -values for *Instruct calibrated similarity* < *Base calibrated similarity* (one-sided) on the Dolma benchmark (unprompted generation, no context). Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
SmolLM2-1.7B	≈ 0	0.053	0.826	0.424	0.256	0.261	0.061	0.017	0.977	1.000
SmolLM2-360M	0.002	0.051	0.004	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	1.000
SmolLM-1.7B	0.133	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.911	0.755
SmolLM-360M	0.002	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.086	0.995

Table 9. Wilcoxon rank-sum test p -values for the hypothesis *longer-chunk calibrated similarity* < *shorter-chunk calibrated similarity* (one-sided) on Dolma. A significant result supports Observation 3 (novelty increases with generation length).

SmolLM2 (no context)			SmolLM2 (with context)		
Model	450 vs. 250	250 vs. 50	Model	450 vs. 250	250 vs. 50
1.7B-Base	≈ 0	≈ 0	1.7B-Base- W/Context	≈ 0	0.078
1.7B-Inst	≈ 0	≈ 0	1.7B-Inst- W/Context	≈ 0	0.090
360M-Base	≈ 0	≈ 0	360M-Base- W/Context	≈ 0	0.147
360M-Inst	≈ 0	≈ 0	360M-Inst- W/Context	≈ 0	0.258
human-Base	0.213	0.475	human-Base- W/Context	0.196	0.454

SmolLM (no context)			SmolLM (with context)		
Model	450 vs. 250	250 vs. 50	Model	450 vs. 250	250 vs. 50
1.7B-Base	≈ 0	≈ 0	1.7B-Base- W/Context	≈ 0	0.187
1.7B-Inst	≈ 0	≈ 0	1.7B-Inst- W/Context	≈ 0	≈ 0
360M-Base	≈ 0	≈ 0	360M-Base- W/Context	≈ 0	0.036
360M-Inst	≈ 0	≈ 0	360M-Inst- W/Context	≈ 0	≈ 0
human-Base	0.213	0.475	human-Base- W/Context	0.196	0.454

Table 10. Wilcoxon rank-sum test p -values for the hypothesis *prompted calibrated similarity* < *unprompted calibrated similarity* (one-sided, alternative=less) on the Dolma benchmark (SmolLM). Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
human-Base	0.580	0.159	0.302	0.202	0.527	0.428	0.605	0.399	0.532	0.516

Table 11. Wilcoxon rank-sum test p -values for the hypothesis *prompted calibrated similarity* < *unprompted calibrated similarity* (one-sided, alternative=less) on the Dolma benchmark (SmolLM2). Columns denote chunk size (tokens).

Model	50	100	150	200	250	300	350	400	450	500
1.7B-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
1.7B-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
360M-Base	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0
360M-Inst	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	≈ 0	0.517	0.029	≈ 0
human-Base	0.580	0.159	0.302	0.202	0.527	0.428	0.605	0.399	0.532	0.516