

Sensitivity of Small Language Models to Fine tuning Data Contamination

Anonymous authors

Paper under double-blind review

Abstract

Small Language Models (SLMs) are increasingly being deployed in resource-constrained environments, yet their robustness to data contamination during instruction tuning remains poorly understood. We systematically investigate the contamination sensitivity of 23 SLMs (270M to 4B parameters) across different model families by measuring susceptibility to syntactic transformations (character and word reversal) and semantic transformations (irrelevant and counterfactual responses), each applied at contamination levels from 1% to 100%. Our results reveal fundamental asymmetries in vulnerability patterns where syntactic transformations cause catastrophic performance degradation with character reversal producing near-complete failure across all models regardless of size or family, whereas semantic transformations demonstrate distinct threshold behaviors and greater resilience in core linguistic capabilities. We discover a ‘*capability curse*’ where more capable models can become more susceptible to learning semantic corruptions, while our analysis of base versus instruction-tuned variants reveals that alignment provides inconsistent robustness benefits and can sometimes reduce resilience. Layerwise representational analysis across model families and sizes shows a consistent localization of contamination effects toward the final blocks, with syntactic corruption typically inducing stronger late-layer divergence and semantic corruption producing comparatively smaller changes that are often confined to final layers. Our work makes three contributions: (1) empirical evidence that SLMs are disproportionately vulnerable to syntactic contamination patterns, (2) characterization of asymmetric learning dynamics for syntactic versus semantic contamination supported by behavioral and representational analysis, and (3) systematic evaluation protocols for robustness assessment. These findings have deployment implications, suggesting that current robustness assumptions may not hold for smaller models and highlighting the need for contamination-aware training protocols that target late layer representations.

1 Introduction

Small Language Models (SLMs) are rapidly becoming the backbone of on-device AI applications, running locally on smartphones, edge devices, and resource-constrained environments where privacy, latency, and infrastructure costs are paramount (Sun et al., 2020; Abdin et al., 2024; Schick & Schütze, 2021). Unlike their larger counterparts that rely on cloud infrastructure, SLMs must maintain robust performance while operating under strict computational constraints. This shift toward local deployment makes understanding their vulnerabilities critical for agentic AI systems. Researchers are exploring different techniques, such as enhancing data quality (Gunasekar et al., 2023), refining training strategies (Hu et al., 2024), and reconfiguration of model architectures (Liu et al., 2024), among others, to improve SLMs.

Although language models have achieved remarkable success in translation, summarization, and question answering tasks (Anthropic, 2024; Achiam et al., 2023), they remain fundamentally limited by the training data quality. Large language models (LLMs) exhibit concerning behaviors such as the ‘Reversal Curse’ (Berglund et al., 2024), highlighting their brittleness to systematic data patterns. For SLMs, this vulnerability is amplified as their reduced parameter counts and compressed representations may make them even more susceptible to learning spurious patterns from data.

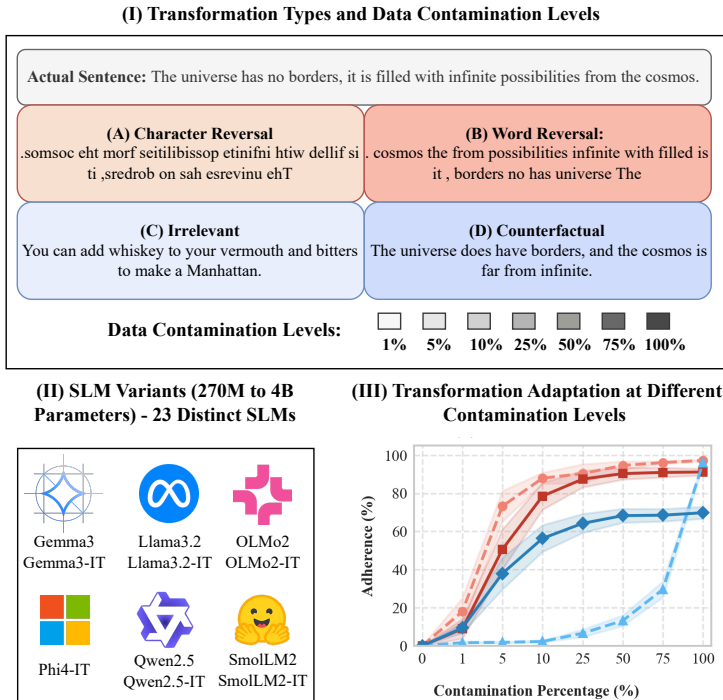


Figure 1: Overview of transformation learning in SLMs. (I, II) Four transformation types at varying contamination levels (1%-100%) are used to fine tune 23 SLMs across six model families (270M-4B parameters). (III) Structural transformations show rapid adoption at 5% , while semantic require higher exposure levels.

Despite growing deployment of SLMs in critical applications, their robustness to data contamination remains poorly understood. Reduced representational redundancy may increase sensitivity to spurious training patterns, while limited capacity may constrain the ability to internalize more complex rules. Existing research often emphasizes stochastic noise or adversarial perturbations (Gupta et al., 2024; Ribeiro & Schön, 2023). In contrast, we study deterministic contamination patterns and identify a counterintuitive failure mode, *capability curse*, where more capable models can acquire semantic corruptions with less exposure.

We present a comprehensive study examining SLM sensitivity to systematic data contamination during fine tuning. Our investigation spans 23 models across six SLM families (270M to 4B parameters) and introduces a framework to understand how different contamination patterns affect model behavior (Figure 1). We use a small set of controlled transformation patterns as diagnostic stressors to isolate fine tuning dynamics under systematic out-of-pretraining supervision. Our experimental design considers four transformation types: two syntactic transformations (character reversal and word reversal) that disrupt tokenization and token order, and two semantic transformations (irrelevant and counterfactual responses) that preserve surface structure while corrupting content alignment and factual consistency. We study syntactic and semantic transformations together because they probe complementary failure modes: usability breakdown versus content corruption. Syntactic transformations test whether instruction tuning internalizes structural rules that break output usability even on clean inputs, whereas semantic transformations test whether it internalizes content level policy shifts that preserve fluency while degrading meaning. These regimes capture distinct deployment risks and exhibit different learning dynamics. Across model families and contamination levels, we observe asymmetric failure modes and threshold behaviors between syntactic and semantic contamination, and use representation level probes to study where contamination induced changes concentrate across network depth. We also introduce random character reversal experiments, which reveal progressive degradation rather than catastrophic collapse, demonstrating partial robustness to stochastic noise. This systematic approach enables identification of gradual behavioral changes, providing information on how contamination interacts with core model components (Bahdanau et al., 2015; Luong et al., 2015; Vaswani et al., 2017), including tokenization and self attention mechanisms. Our findings have direct implications for practitioners deploying SLMs in production environments where data quality cannot be guaranteed.

2 Related Work

Noise is routinely introduced into language-model training to improve robustness and generalization, motivated by the fact that user and web generated text often contains errors, inconsistencies, and non-standard language use. Prior work studies robustness via injecting noise by perturbing parameters (Wu et al., 2022) and introducing noisy labels (Wang et al., 2023; Wu et al., 2023; Hedderich et al., 2021). Other studies analyze how noise impacts performance across NLP tasks and training regimes (Wang et al., 2024; Havrilla & Iyer, 2024; Al Sharou et al., 2021). While this literature typically treats corruption as stochastic perturbation (random, unpredictable noise) or imperfect supervision, our focus is on systematic contamination, which refers to consistent, deterministic patterns applied across training examples within instruction-tuning data. Critically, we investigate whether models learn to internalize and reproduce these patterns as stable behaviors during generation, examining both structural transformations (character/word reversal) and semantic variations (counterfactual/irrelevant responses).

Complementing these training-centric approaches, extensive work benchmarks robustness against inference-time perturbations. Character-level noise and linguistic variations have been shown to expose brittleness that standard accuracy metrics can hide (Moradi & Samwald, 2021; Ribeiro et al., 2020). The relationship between model scale and robustness is also debated. Some work suggests overparameterized models can be more susceptible to subpopulation attacks or adversarial noise (Gupta et al., 2024; Ribeiro & Schön, 2023), while other results argue that scale generally confers robustness (Wang et al., 2024). In contrast to inference-time perturbation studies, we examine corruption introduced during fine tuning and ask when the fine-tuned model begins to output the contamination pattern itself, even when evaluation prompts are clean.

Beyond stochastic noise and inference-time perturbations, a closely related security-motivated line of work studies adversarial instruction-tuning data corruption, including poisoning and backdoors that implant controllable behaviors. Recent results show that relatively small numbers of poisoned instruction-tuning examples can induce reliable misbehavior in instruction-following models (Wan et al., 2023; Xu et al., 2024). In contrast to trigger-conditioned poisoning or backdoors, we study structured, non-trigger deterministic contamination patterns that can arise from non-adversarial data pipeline errors, templating, or synthetic-data generation, and we highlight that instruction-following objectives can amplify vulnerability to such systematic supervision artifacts.

Our investigation also connects to analyses of spurious rule learning and asymmetric generalization in language models. Prior work shows that models trained on directional facts can fail to generalize to the inverse relation (the reversal curse), indicating that learning can be strongly shaped by superficial directional patterns in the data (Berglund et al., 2024). This perspective supports studying systematic transformations as a controlled way to probe whether fine tuning encourages shortcut acquisition of consistent but undesirable response patterns.

Finally, work on safety alignment indicates that post-training constraints can be brittle. Fine tuning can erode prior safety behaviors, and stylistic alignment or high similarity between alignment and fine tuning datasets can induce catastrophic forgetting of guardrails (Xiao et al., 2026; Hsiung et al., 2025). These findings motivate analyzing how instruction-tuning interventions interact with instruction-following and alignment, including settings where the fine tuning objective may reinforce undesirable supervision artifacts.

Overall, distinct from standard label-noise settings, inference-time perturbation benchmarks, or trigger-based poisoning, the dynamics of systematic contamination during instruction tuning remain underexplored for small language models. Because SLMs have limited capacity, they may exhibit a different robustness trade-off than larger models, with reduced representational redundancy potentially increasing sensitivity to spurious patterns while also limiting the ability to fully internalize complex contamination rules. Our work empirically characterizes these dynamics under controlled contamination rates.

3 Methodology

We selected efficient SLM families to examine the influence of model scaling and alignment training on the detection of contamination patterns. Six SLM families, each with less than 4 billion parameters, were

studied: Gemma3 (Team et al., 2025), Llama3.2 (Grattafiori et al., 2024), OLMo2 (Walsh et al., 2025), Phi4 (Abouelenin et al., 2025), Qwen2.5 (Yang et al., 2025), and SmolLM2 (Allal et al., 2025). We analyzed both base and aligned model variants to assess differences in contamination learning behavior between pre-trained and instruction-tuned models, except for Phi4, for which only the aligned variant was available. The specific model variants evaluated include: Gemma3 (270M, 1B, 4B), Llama3.2 (1B, 3B), OLMo2 (1B), Phi4 (Mini), Qwen2.5 (0.5B, 1.5B, 3B), and SmolLM2 (360M, 1.7B), resulting in a total of 23 different models.

3.1 Instruction tuning dataset

The primary clean instruction tuning dataset, denoted \mathcal{D}_{ad} , was constructed by combining two high-quality filtered datasets: AlpaGasus dataset ($\mathcal{D}_{\text{AlpaGasus}_{9k}}$) (Chen et al., 2023), derived from Alpaca (Taori et al., 2023) and the Dolly dataset ($\mathcal{D}_{\text{Dolly}_{3k}}$) filtered from Databricks Dolly dataset (Conover et al., 2023). Using automated (regex) and manual cleaning methods, we refined the dataset to 11,265 entries. The cleaning process is detailed in the Appendix B. This dataset served as the clean baseline in our experiments.

To evaluate the robustness of the model to data contamination, we applied four types of systematic transformations to \mathcal{D}_{ad} and generated the corresponding contaminated datasets. Examples illustrating these transformations are provided in Figure 1(I). The rationale for choosing the transformation patterns is detailed in Appendix A. The first two involved syntactic modifications of the answers. The word-level reversal dataset, $\mathcal{D}_{\text{ad_wreversal}}$, was created by reversing the order of words in the answer strings $a^{(i)}$ (denoted as $REVERSE_{\text{word}}(a^{(i)})$). The character-level reversal dataset ($\mathcal{D}_{\text{ad_creversal}}$) was generated by reversing all characters within the answer strings (denoted as $REVERSE_{\text{char}}(a^{(i)})$). **While character reversal is an intentionally extreme perturbation, it serves as a controlled proxy for token-level surface-form shifts that can arise from transliteration and multilingual normalization variation. Such shifts can substantially alter sub-word segmentation, increasing fragmentation and changing effective token lengths even when meaning is unchanged.** Additionally, to examine stochastic versus systematic contamination effects, we created a random character reversal dataset ($\mathcal{D}_{\text{ad_creversal2}}$) where a random percentage of characters within each answer string was reversed (denoted as $REVERSE_{\text{char2}}(a^{(i)})$), introducing partial structural noise.

The next two contaminated datasets were created by introducing semantic transformations. The irrelevant dataset ($\mathcal{D}_{\text{ad_irr}}$) was constructed by pairing each question $q(x^{(i)})$ from \mathcal{D}_{ad} with a randomly selected answer $IRR(a^{(i)})$ from a different example ($i \neq j$) in the clean dataset, thus ensuring no semantic correspondence. For the counterfactual dataset, $\mathcal{D}_{\text{ad_cfact}}$, we used Gemini 2.5 Flash (Comanici et al., 2025) to generate counterfactual answers, $CFACT(a^{(i)})$, for questions, $q(x^{(i)})$, of $\mathcal{D}_{\text{ad_train}}$. To ensure systematic generation of high-quality counterfactual responses, Gemini 2.5 Flash was given a specific ‘Simulator’ persona designed for AI safety research, instructing it to simulate flawed AI responses by following high-level instruction formats while deliberately contradicting specific content requirements. We then use Gemini 2.0 Flash to score each candidate on a 0-5 rubric for structural adherence and counterfactual strength, and we regenerate any responses scoring below 4 until all examples meet the threshold (prompts and details in Appendix C).

Additionally, we instantiate Random Character Reversal using $\mathcal{D}_{\text{ad_creversal2}}$ at five noise levels, where 5%, 10%, 15%, 20%, or 25% of characters in each answer string are randomly selected and reversed. **Unlike the four primary transformations applied at dataset-level contamination rates, this introduces localized stochastic corruption within each training example, providing a complementary contrast to deterministic pattern learning.** Overall, we evaluate 33 contamination settings: 28 systematic settings (4 transformation types \times 7 contamination levels) plus 5 Random Character Reversal settings. Examples for each transformation type are provided in Table 2 in Appendix B.

3.2 Test dataset

The primary test dataset ($\mathcal{D}_{\text{test}}$) was created using GPT-4o and consisted of 2018 question-answer examples ($(q^{(i)}, a^{(i)})$). These examples were designed to cover a diverse range of topics, reflecting a specific distribution. Science (General, Biology, Physics, etc.) and Mathematics constituted the largest category (approx. 35-40%), followed by substantial representation from Geography and History (approx. 15-20%), General Knowledge (approx. 10-15%), Arts, Literature, and Culture (approx. 8-12%), and general writing tasks (approx. 8-12%). Smaller proportions covered areas including Technology, Language, Philosophy, Food, and Sports,

ensuring broad coverage. We design $\mathcal{D}_{\text{test}}$ as a broad-coverage but relatively straightforward evaluation set for instruction-tuned SLMs so that baseline performance is high enough for contamination-induced degradation to be attributable to the fine-tuning intervention rather than task difficulty. Accordingly, $\mathcal{D}_{\text{test}}$ serves as a controlled diagnostic distribution for measuring contamination-driven performance shifts, and absolute metric values and fine-grained model rankings may vary under different evaluation distributions. Details of the test data cleaning process are provided in Appendix B.4.

3.3 Experimental setup

We study contamination as a training-time phenomenon during instruction tuning in a mixture setting. A fraction of instruction-response pairs in the fine tuning set have transformed outputs according to a fixed syntactic rule or a fixed semantic deviation, while evaluation inputs remain clean. We do not assume any inference-time trigger, targeted backdoor objective, or test-time access. This setting covers both unintentional supervision corruption (for example dataset assembly and filtering failures, synthetic-data generation artifacts) and deliberate but non-triggered training data manipulation, and it enables controlled analysis of when and how models internalize corrupted supervision as stable generation behavior.

The experiments were designed to systematically investigate the differential vulnerability of SLMs to syntactic versus semantic disruption patterns at different levels of data contamination. Syntactic transformations (character/word reversal) violate fundamental language structure and formatting, and semantic transformations (irrelevant/counterfactual responses) preserve structural coherence but corrupt content alignment. Thus, syntactic transformations were tested in both pre-trained and instruction-tuned variants, since pre-trained models lack alignment constraints that resist format disruption. On the other hand, semantic transformations were used only for instruction-tuned models, since these transformations primarily target alignment rather than basic language competency. Performance baselines were established using out-of-the-box instruction-tuned models without additional training on \mathcal{D}_{ad} , providing clean reference points to measure contamination-induced behavioral changes. Each transformation type was applied at contamination levels (1%, 5%, 10%, 25%, 50%, 75%, 100%) during instruction tuning on mixed datasets, yielding 570 finetuned model checkpoints (see Appendix E.1 for the count breakdown). We also evaluated random character reversal ($\mathcal{D}_{\text{ad_creversal2}}$) at varying internal noise levels (5%, 10%, 15%, 20%, and 25% of characters per answer randomly reversed), which differs from systematic reversal by introducing stochastic corruption within individual examples rather than deterministic patterns across the dataset. This allows us to test whether models exhibit gradual degradation or catastrophic failure. This approach enables identification of contamination thresholds where model behavior degrades and comparison of syntactic versus semantic transformations. Details of training configurations are given in Appendix E.

Additionally, to examine how our findings change when transformations are introduced on the input side, we run an ablation on the Qwen2.5 family under two settings: (i) input transformation, where only inputs are transformed, and (ii) ‘input and output’ transformation, where the same syntactic transformation is applied to both input and output. Experimental details and results are reported in Appendix H.

3.4 Evaluation

We designed a multi-dimensional evaluation framework measuring both behavioral change extent and pattern reproduction fidelity.

Transformation-Specific Processing: Each contamination type required specialized preprocessing. Word-reversal responses underwent word-order reversal, character-reversal responses underwent character-level reversal, while irrelevant and counterfactual transformations did not require preprocessing as they introduce semantic rather than structural modifications.

Primary Metrics: Our core assessment employed semantic similarity using ‘all-mpnet-base-v2’ sentence transformers to compute cosine similarity between preprocessed outputs and references. This embedding-based approach captures meaning preservation beyond surface matching, revealing whether models internalize target transformations while maintaining semantic coherence. Standard lexical metrics (BLEU, ROUGE, METEOR) were also used as secondary metrics.

LLM-as-a-Judge Assessment: Gemini 2.0 Flash evaluated two critical dimensions: (i) *Pattern Adherence* by assessing whether responses match the correct transformation pattern (WordReversal, CharReversal, Irrelevant, CounterFactual) while explicitly ignoring factual accuracy; (ii) *Accuracy and Grammatical Correctness* by comparing preprocessed responses against references for factual fidelity and structural coherence. Prompts used for this assessment are detailed in Appendix D. An analysis of human-model agreement was also performed to validate the reliability of the automated assessment. Strong alignment was observed between human evaluators and Gemini 2.0 Flash (detailed in the Appendix I).

Representational Analysis: We complement behavioral metrics with layerwise representation comparisons between baseline models and contamination-finetuned variants. On a fixed evaluation set, we compute per-layer linear CKA (Kornblith et al., 2019) and average last-token cosine similarity at the output of each transformer block. These probes capture geometry-level drift across examples (CKA) and shifts in generation relevant states (cosine similarity).

4 Results

We present results primarily through semantic similarity scores and LLM-based evaluations, as these metrics directly capture the behavioral shifts central to our investigation. Standard lexical metrics showed limited discriminative power for our research objectives and are provided in the appendix for completeness.

4.1 Baseline performance and Scaling effects in SLMs

Table 1: Baseline performance metrics of the different instruction-tuned models on $\mathcal{D}_{\text{test}}$

Model	Accuracy	Semantic Similarity	Grammatical Correctness
Gemma3_270M_IT	42.12%	62.14%	95.66%
Gemma3_1B_IT	79.20%	67.69%	100.00%
Gemma3_4B_IT	94.15%	82.16%	100.00%
Llama3.2_1B_IT	79.03%	77.87%	98.76%
Llama3.2_3B_IT	91.68%	84.57%	100.00%
OLMo2_1B_IT	82.76%	75.22%	98.21%
Phi4_Mini_IT	96.63%	84.85%	100.00%
Qwen2.5_0.5B_IT	68.25%	76.27%	96.00%
Qwen2.5_1.5B_IT	93.27%	79.09%	98.47%
Qwen2.5_3B_IT	92.27%	83.23%	100.00%
SmolLM2_360M_IT	60.92%	61.31%	96.25%
SmolLM2_1.7B_IT	89.54%	75.31%	98.12%

Table 1 reports baseline performance on the clean test set $\mathcal{D}_{\text{test}}$ for all instruction-tuned models. Accuracy spans 42.12% (Gemma3_270M_IT) to 96.63% (Phi4_Mini_IT), and within families larger models generally achieve higher accuracy (e.g., Gemma3 from 42.12% to 94.15%). Semantic similarity ranges from 61.31% (SmolLM2_360M_IT) to 84.85% (Phi4_Mini_IT), with Llama3.2_3B_IT close at 84.57%. Grammatical correctness is consistently high across models (all $\geq 95\%$), indicating that instruction-tuned SLMs largely preserve surface well-formedness even when accuracy varies.

4.2 Sensitivity to syntactic and semantic contamination

The SLM performance under different levels of syntactic and semantic data contamination indicates a fundamental asymmetry in behavior (Figure 2). SLMs exhibit higher sensitivity to syntactic transformations than semantic ones. Each point in the plot represents the mean performance of all SLMs under a contamination level for each transformation. Shaded areas around the lines indicate the standard error of the mean.

Figure 2(A) demonstrates adherence to contamination patterns across transformation types. Syntactic transformations (character and word reversal) exhibit rapid contamination learning. Adherence reaches $\sim 10\text{-}15\%$ at 1% contamination, rising sharply to $>50\%$ by 5% and approaching saturation ($\sim 90\%$) by 10%. Semantic transformations show distinct behaviors, counterfactual adherence increases earlier, reaching $\sim 38\%$ at 5% before plateauing at $\sim 70\%$, while irrelevant remains negligible until 50% contamination, achieving near-complete adherence only at 100%.

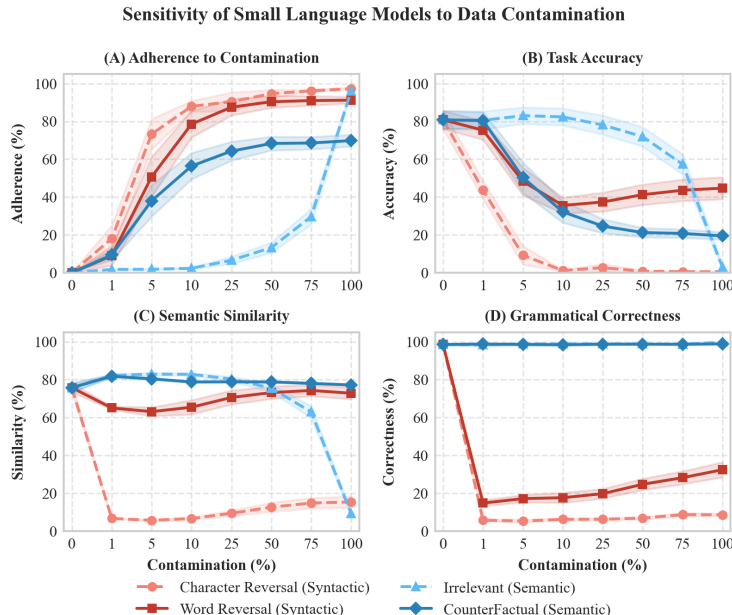


Figure 2: As data contamination increases, SLMs learn to adhere to flawed patterns, causing a decline in task accuracy, semantic similarity, and grammatical correctness.

Figure 2(B) reports task accuracy. Syntactic contamination causes catastrophic accuracy degradation: character reversal collapses from $\sim 80\%$ to $\sim 40\%$ at merely 1% contamination. Word reversal shows a minor decline at 1% ($\sim 75\%$) before dropping sharply to $\sim 50\%$ at 5% and stabilizing around $\sim 40\text{-}45\%$ at higher levels. Semantic transformations demonstrate distinct resilience profiles where irrelevant contamination remains robust ($\sim 75\%$) through 25% before declining, whereas counterfactual contamination degrades earlier, dropping to $\sim 50\%$ accuracy at 5% contamination.

The semantic similarity patterns in Figure 2(C) further distinguish transformation types. Character reversal triggers an immediate collapse to single-digit similarity at 1% and remains very low thereafter, indicating loss of meaning preservation. Word reversal preserves substantially higher similarity throughout (roughly mid-60% to mid-70%). For semantic transformations, counterfactual similarity stays high and fairly stable (around the high-70% range), whereas irrelevant similarity remains high up to moderate contamination but drops sharply at 100%, reflecting loss of contextual relevance under full irrelevance.

Figure 2(D) illustrates changes in grammatical correctness due to contamination. Syntactic contamination devastates grammatical coherence sharply at low contamination: character reversal drops to single-digit correctness by 1% and stays low, while word reversal drops to roughly the mid-teens at 1% and partially recovers with higher contamination. In contrast, semantic contamination has little effect on grammatical correctness, which remains near-ceiling across contamination levels. Overall, syntactic transformations primarily damage form and meaning, whereas semantic transformations damage content while preserving fluency.

4.3 Effects of model size, alignment, and family on syntactic robustness

Syntactic robustness varies strongly with transformation type and model family. Figure 3 reports contamination adherence (left) and task accuracy (right) for character reversal and word reversal across contamination levels (1% to 100%). Across families, character reversal is consistently more damaging than word reversal, but both the onset and severity of degradation vary substantially by family and by whether the model is base or instruction-tuned.

Alignment effects on syntactic robustness: Comparing base and instruction-tuned (IT) variants shows mixed, family-dependent effects rather than uniform robustness gains, especially at low contamination (1-10%). In several cases, IT models adopt the contaminated pattern earlier and lose accuracy faster. For



Figure 3: Model performance on test data under increasing syntactic data contamination. The heatmaps show model sensitivity (left) and task accuracy (right) for character and word reversal tasks.

character reversal, Llama3.2_3B attains 61.5% accuracy at 1% contamination, while Llama3.2_3B_IT drops to 25.9%. This gap coincides with a large difference in adherence at 1% (1.1% for base vs. 59.4% for IT), consistent with early transformation adoption undermining task performance. As contamination increases, performance for both variants degrades, with accuracies falling to single digits by 50% contamination (3.8% for base and 0.2% for IT). For word reversal, the effect of instruction tuning is less consistent across families and contamination levels, and differences between base and IT variants often narrow.

Model size and family effects on syntactic robustness: Parameter scaling does not yield uniform gains. Family identity is often a stronger predictor than size, and scaling trends differ by transformation. For word reversal, Qwen2.5 shows clearer scaling benefits at moderate contamination; for example, at 5% contamination, Qwen2.5_3B reaches 53.4% accuracy compared to 13.8% for Qwen2.5_0.5B, and Qwen2.5_3B_IT reaches 42.7% compared to 16.1% for Qwen2.5_0.5B_IT. At 1% contamination, however, scaling is not strictly monotonic across variants (e.g., Qwen2.5_3B_IT is 77.8% while Qwen2.5_0.5B_IT is 63.1%). Gemma3 exhibits less regular scaling. For character reversal at 1% contamination, Gemma3_4B achieves 65.1% while Gemma3_270M is 20.0%. Overall, the heatmaps indicate that robustness depends on a combination of family, size, and training variant, rather than parameter count alone.

Character reversal versus word reversal: Character reversal is a particularly extreme stressor. By 10% contamination, many models fall to near-zero accuracy, although some retain non-trivial performance in the low-to-moderate regime (e.g., Llama3.2_3B retains 40.4% at 10% contamination, and SmolLM2_360M retains 16.9%). Word reversal is substantially milder, with many models retaining moderate accuracy at low to mid contamination and degrading gradually. The same contrast appears in semantic similarity and grammatical correctness (Figure 7 in Appendix F.1): character reversal yields very low semantic similarity even at low contamination and reduces grammatical correctness to mostly single-digit values, whereas word reversal preserves higher semantic similarity (often 50-80) with less decline in grammatical correctness.

Effects of noise intensity on degradation: Random Character Reversal produces smoother degradation rather than exhibiting a sharp ‘cliff-edge’ transition. In the 5-10% noise range, many models retain substantial accuracy, while increasing noise toward 15-25% typically reduces accuracy further in a near-linear but model-dependent manner, with occasional non-monotonic fluctuations. This supports the view that stochastic corruption erodes performance gradually rather than inducing abrupt rule driven failure. Detailed results are provided in Appendix F.2.

4.4 Effects of model size and family on semantic robustness

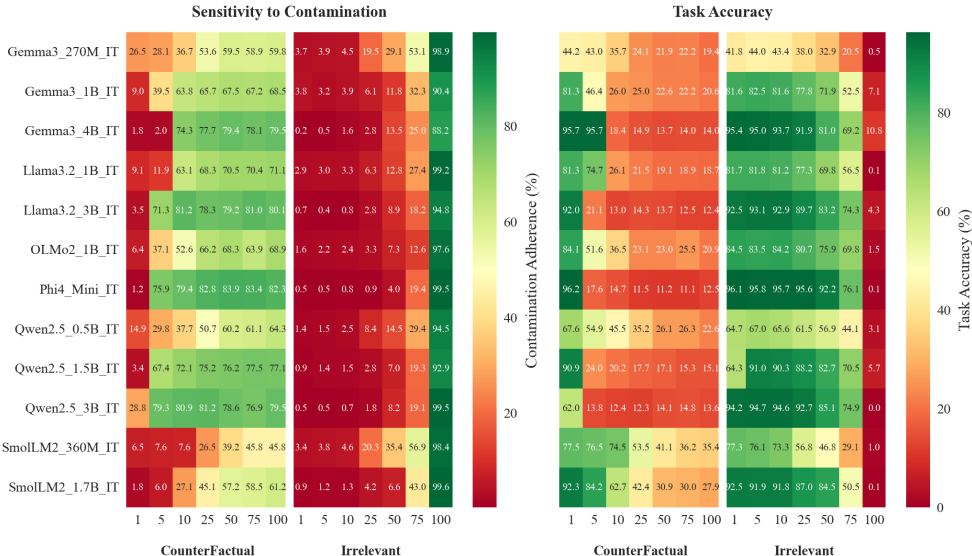


Figure 4: Model performance on semantic tasks under increasing data contamination. The figure illustrates model sensitivity (left) and task accuracy (right) when trained with counterfactual and irrelevant information.

The analysis of semantic transformations reveals behavioral patterns different from syntactic vulnerabilities (Figure 4). Since these transformations require models to alter content alignment, only instruction-tuned variants were evaluated. Models learn counterfactual and irrelevant transformations differently. At low contamination levels (1-10%), many models retain a partial ‘resilience buffer’ relative to syntactic tasks, but the strength of this buffer depends on model family and capacity. For counterfactual contamination, some models show low adherence at 1% (e.g., Phi4_Mini_IT: 1.2; Gemma3_4B_IT: 1.8), while others exhibit substantial adherence even at 1% (e.g., Qwen2.5_3B_IT: 28.8; Gemma3_270M_IT: 26.5), indicating that counterfactual behavior can be acquired with very limited exposure in more capable variants. For irrelevant contamination, early learning is weaker. Adherence remains in the low single digits for nearly all models through 10% contamination, suggesting models initially resist abandoning contextual relevance. However, at full contamination (100%), most models strongly adhere to the irrelevant pattern (typically >90% adherence).

Evaluations of semantic similarity and grammatical correctness (Figure 10 in Appendix F.3) show that, unlike syntactic corruption, semantic contamination does not generally destroy linguistic form. Across most contamination levels, semantic similarity remains high (roughly 75-85%) and grammatical correctness remains near-ceiling (typically 95-100%). A notable exception occurs for irrelevant contamination at 100% where semantic similarity drops sharply (to below 15% across all models), while grammatical correctness remains high, indicating that irrelevance learning primarily disrupts content alignment rather than surface fluency.

Impact of model scale and family: Model capacity and family strongly mediate how quickly semantic corruptions are learnt. For counterfactuals, scaling is pronounced in some families. At 5% contamination, Qwen2.5_0.5B_IT shows 29.8% adherence while Qwen2.5_3B_IT reaches 79.3%, and Llama3.2_1B_IT (11.0%) is far below Llama3.2_3B_IT (71.3%). These differences translate into sharper accuracy degradation for models that readily follow the corrupted supervision; for example, Phi4_Mini_IT drops from 96.2% accuracy at 1% counterfactual contamination to 17.6% at 5%. In contrast, irrelevance learning remains low at low contamination across families, and becomes dominant at high contamination levels.

4.5 Representational analysis of contamination effects

We analyze layerwise representational drift between baseline models and their finetuned variants using CKA and last-token cosine similarity computed on the same evaluation set. Figure 5 shows results for Qwen2.5_IT.

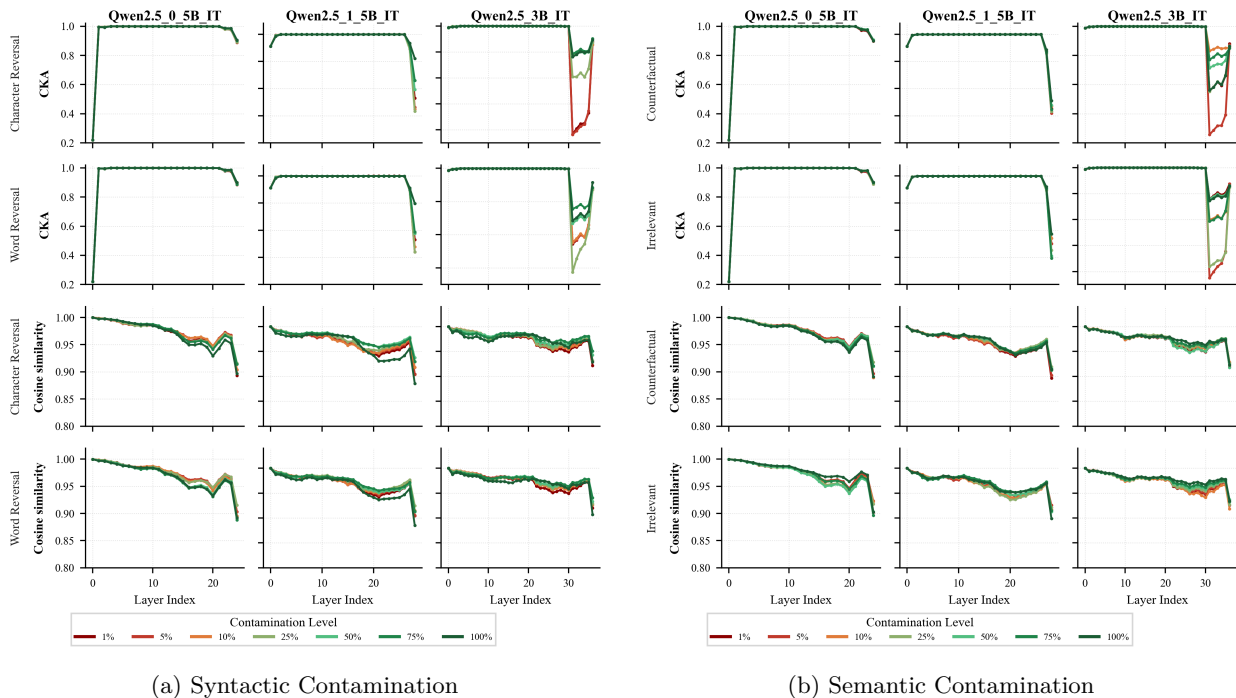


Figure 5: The plots show layerwise CKA and average last-token cosine similarity between baseline and finetuned Qwen2.5_IT models across contamination levels (1% to 100%). (a) reports syntactic contamination (character and word reversal) and (b) reports semantic contamination (counterfactual and irrelevant).

Figure 5a shows that, under syntactic transformations, similarity to the baseline remains high through much of the stack, with deviations concentrated at depth extremes. In Qwen2.5_0.5B_IT, CKA drops sharply at the earliest layer before returning close to 1.0, suggesting the initial discrepancy does not propagate through intermediate layers. In Qwen2.5_3B_IT, mid-network layers remain near-ceiling (CKA > 0.95, cosine > 0.95), with the strongest departures in the final blocks. Late-layer divergence is most pronounced for character reversal: in Qwen2.5_3B_IT, late-layer CKA falls to roughly 0.2–0.3 at some contamination levels while cosine similarity stays around 0.9, indicating geometry-level reorganization without comparable change in average last-token direction. Qwen2.5_1.5B_IT shows the same pattern with smaller-magnitude drops, and word reversal produces milder late-layer deviations than character reversal across scales.

Importantly, representational divergence does not follow a simple dose-response ordering with contamination rate. CKA measures cross-example representational geometry rather than behavioral adherence, so it can drop sharply even when behavioral adoption of the transformation remains weak. At very low contamination (e.g., 1%), fine tuning is dominated by clean supervision but includes a small fraction of systematically corrupted targets, creating a mixed-supervision conflict regime. In this regime, the model can introduce late-layer adjustments that partially accommodate rare corrupted examples while still optimizing for the clean objective; CKA can decrease because late-layer geometry is perturbed, even if the resulting changes are not coherent enough to yield reliable transformation behavior at generation time. As contamination increases and corrupted supervision becomes more prevalent, the training signal becomes less conflicting and fine tuning can converge to a more stable configuration; correspondingly, CKA relative to baseline can evolve non-monotonically rather than decreasing smoothly with contamination percentage. Where this behavior manifests across depth is family dependent, and Qwen2.5 provides a clear instance where the strongest deviations are concentrated in the final blocks.

Figure 5b shows substantially higher representational stability under semantic transformations than under syntactic corruption. Cosine similarity remains near-ceiling through most layers for all three Qwen2.5_IT variants, while CKA is also high across much of the network. The 3B model nonetheless exhibits noticeable

late-layer CKA dips at some contamination levels for both counterfactual and irrelevant settings, indicating that semantic corruption can still induce late-stage geometry changes even when intermediate-layer trajectories remain largely preserved. Overall, semantic transformations produce weaker changes in cosine similarity and preserve most intermediate-layer structure, with late-layer CKA remaining the most sensitive indicator of contamination-induced representational change.

These representation-level patterns complement the behavioral results by suggesting that contamination effects are implemented primarily through late-stage changes closer to the decoding layers, while earlier layers remain comparatively stable.

Across model families and sizes: Appendix G extends the representational analysis beyond Qwen2.5_IT and shows that the main localization trend generalizes across families. Across models, the largest deviations from the baseline most often appear in the final transformer blocks, while intermediate layers remain comparatively stable; the way this divergence is expressed depends on the family. In several cases, cosine similarity varies smoothly with depth and remains high until late layers, whereas CKA is either near-ceiling throughout or exhibits localized late-layer drops. Although the magnitude and exact location of deviations vary by model family and size, the most prominent changes are typically concentrated near the output end of the stack, consistent with contamination being expressed through late-stage representational adjustments that most directly influence the output projection.

5 Discussion

Our systematic evaluation of contamination effects across 23 small language models (SLMs), spanning four contamination types and contamination levels from 1% to 100%, reveals asymmetric vulnerability patterns that complicate simple assumptions about robustness and scaling in SLMs.

Structural vs. semantic contamination patterns: The contrast between syntactic and semantic robustness suggests that current SLMs fail in qualitatively different ways under structural versus content-based corruption. Syntactic character-level contamination is consistently the most damaging stressor, indicating a shared sensitivity to character-level disruptions that interact with tokenization and sequence processing. In many models, utility metrics often degrade sharply in the low-contamination regime (1-10%) and frequently collapse as contamination increases, showing limited tolerance to systematic character-level transformations. Word-level reversal is substantially milder, with many models retaining higher semantic similarity and grammatical correctness than under character reversal. In contrast, semantic transformations preserve surface fluency across most contamination levels. Models maintain near-ceiling grammatical correctness and relatively high semantic similarity, even while the intended content becomes corrupted. This separation indicates that semantic contamination primarily perturbs content alignment rather than grammatical form, whereas character-level syntactic corruption degrades both meaning preservation and well-formedness. Representation-level probes further suggest that contamination effects often concentrate in late layers closer to decoding, while earlier layers remain comparatively stable across many settings. [These asymmetries correspond to two distinct deployment risks. Structural corruption threatens usability: even small amounts of character-level contamination can produce outputs that are unusable for default pipelines, downstream tools, or end users without explicit post-processing, and this failure can be triggered even when test-time inputs are clean. Semantic corruption, by contrast, threatens correctness and alignment while preserving surface fluency, which makes wrong outputs harder to detect through surface-level quality checks.](#)

The alignment paradox: The inconsistent and sometimes detrimental effects of instruction tuning on syntactic robustness suggest that standard post-training procedures do not reliably confer robustness to structural corruption. In several cases, instruction-tuned variants exhibit higher sensitivity in the low-contamination regime, consistent with the idea that supervised fine tuning can reinforce spurious transformation patterns present in the training signal. This points to robustness against structural artifacts as a distinct competency that likely requires targeted data augmentation, objectives, or defenses, rather than emerging as a byproduct of alignment or instruction-following training.

The capability curse in semantic corruption: [The strongest evidence for the capability curse appears under counterfactual contamination, where more capable models can learn the corrupted supervision with](#)

less exposure, which can translate into larger downstream accuracy drops when evaluation expects faithful, grounded responses. This trade-off reflects a form of ‘*capability curse*’ where stronger instruction-following can increase susceptibility to harmful or misleading training signals. This is deployment-relevant because the same capability that supports high-quality instruction following can also support high-fidelity reproduction of corrupted supervision. Irrelevance follows a different, delayed-threshold dynamic where models often resist it at low contamination, but can adopt it strongly at high contamination, leading to severe degradation in semantic similarity at full contamination even when grammatical correctness remains high.

Model family vulnerability patterns: Robustness is not uniform across model families, and scaling trends are family and transformation dependent. Some families show clearer scaling benefits under word-level syntactic corruption, while others display irregular or non-monotonic behavior. Across our experiments, SmoLLM2 variants tend to be among the least robust, particularly under syntactic corruption.

Implications for SLM development: These findings have practical implications for training and deploying SLMs in settings where data quality cannot be guaranteed. The strong sensitivity to low levels of structural contamination implies that data curation and filtering must prioritize structural integrity (e.g., detecting systematic character-level artifacts) alongside semantic correctness. In addition, results under stochastic character-level corruption (Random Character Reversal) indicate that progressive performance degradation can arise even without deterministic transformation rules, reinforcing the need for defenses that address both systematic and stochastic text corruption. Overall, robustness does not follow from parameter scaling alone and will likely require targeted architectural, training, and evaluation interventions.

Our core experimental grid applies transformations only to outputs during instruction tuning to isolate whether corrupted supervision is internalized when test-time inputs remain clean. Appendix H extends this design with Qwen2.5 syntactic experiments for input and ‘input and output’ transformations, showing that transformation placement changes the dominant failure mode. A key limitation is that these additional experiments are restricted to a single model family; establishing whether the same input and ‘input and output’ patterns hold across other families remains future work. We also did not evaluate parameter-efficient fine-tuning methods such as Low-Rank Adaptation (LoRA); establishing whether such methods exhibit similar contamination acquisition behavior remains an open direction. Similarly, we scope this work to supervised instruction tuning and do not study contamination under RL-based post-training (e.g., RLHF, DPO, GRPO), where updates are driven by reward or preference signals over sampled outputs rather than by direct token-level imitation of transformed targets. Characterizing contamination under these regimes, for example through corrupted preference labels, reward models, or reference responses, is a natural extension that requires a parallel controlled design.

6 Conclusion

Our work establishes a systematic framework for analyzing data-contamination vulnerabilities in small language models (SLMs), addressing an important gap as these systems are increasingly deployed in resource-constrained settings. Our findings challenge common assumptions about robustness. Even minimal syntactic corruption can substantially degrade performance in many models, revealing architectural sensitivities that parameter scaling alone does not resolve. Moreover, experiments with stochastic character-level corruption (Random Character Reversal) show progressive degradation as noise increases, indicating limited tolerance to character-level perturbations that interact with tokenization and sequence processing. In the semantic setting, we find a counterintuitive trend that we term the ‘*capability curse*’, where larger, more capable instruction-following models can acquire counterfactual behaviors with less exposure, which increases susceptibility to certain semantic corruptions.

These results have practical implications for practitioners deploying SLMs in environments where training data quality cannot be guaranteed. The pronounced asymmetry between syntactic and semantic robustness motivates targeted data curation and training procedures that explicitly account for structural integrity and content alignment. Finally, our evaluation protocol covers 23 models, four systematic contamination types, and contamination rates from 1% to 100%, with Random Character Reversal included as an additional stochastic setting. Together, this provides a reproducible benchmark for contamination-aware robustness assessment and can support future work on more reliable and safe on-device language models.

Broader Impact Statement

This work studies how small language models (SLMs) change their behavior when instruction-tuning data contains systematic contamination, including syntactic transformations and semantic deviations. The primary positive impact is to improve reliability and safety of SLM deployment by providing an evaluation protocol and concrete failure modes that can inform dataset curation, training safeguards, and robustness testing. As SLMs are increasingly used in resource-constrained and privacy-sensitive settings, understanding sensitivity to corrupted supervision can help reduce harms such as degraded task performance, misleading outputs, and brittleness under noise.

We also recognize dual-use considerations. By identifying contamination regimes under which instruction tuning can induce undesirable behaviors, our results could inform malicious attempts to influence model behavior via training-data manipulation, consistent with prior work on poisoning during instruction tuning (Wan et al., 2023). We therefore focus on transparent, non-trigger transformations that reflect plausible pipeline and synthetic-data artifacts and emphasize the defensive value of the findings. Our goal is to enable the community to develop and evaluate safeguards such as data validation, filtering, provenance checks, and contamination-aware robustness evaluation, rather than to provide procedural guidance for attacks against specific production systems.

This work does not involve human subjects or personally identifiable information. If the identified vulnerabilities are not addressed, contaminated training data could still contribute to safety regressions and increased misinformation risk in semantic corruption settings. We encourage future work on defenses, including structured-contamination detection, robust training objectives, and evaluation protocols that jointly measure utility and safety under realistic data quality failures.

References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, et al. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*, 2025.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Khetam Al Sharou, Zhenhao Li, and Lucia Specia. Towards a better understanding of noise in natural language processing. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pp. 53–62, 2021.
- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and Thomas Wolf. Smollm2: When smol goes big – data-centric training of a small language model, 2025. URL <https://arxiv.org/abs/2502.02737>.
- AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.
- Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- Lukas Berglund, Meg Tong, Maximilian Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. The reversal curse: LLMs trained on “a is b” fail to learn “b is a”. In *The Twelfth International Conference on Learning Representations*, 2024.

- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srivasan, Tianyi Zhou, Heng Huang, et al. Alpapasus: Training a better alpaca model with fewer data. In *The Twelfth International Conference on Learning Representations*, 2023.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first truly open instruction-tuned llm. *Company Blog of Databricks*, 2023.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- Isha Gupta, Hidde Lycklama, Emanuel Opel, Evan Rose, and Anwar Hithnawi. Fragile giants: Understanding the susceptibility of models to subpopulation attacks. *arXiv preprint arXiv:2410.08872*, 2024.
- Alex Havrilla and Maia Iyer. Understanding the effect of noise in llm training data with algorithmic chains of thought. *arXiv preprint arXiv:2402.04004*, 2024.
- Michael A Hedderich, Dawei Zhu, and Dietrich Klakow. Analysing the noise model error for realistic noisy label data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 7675–7684, 2021.
- Lei Hsiung, Tianyu Pang, Yung-Chen Tang, Linyue Song, Tsung-Yi Ho, Pin-Yu Chen, and Yaoqing Yang. Why llm safety guardrails collapse after fine-tuning: A similarity analysis between alignment and fine-tuning datasets. *arXiv preprint arXiv:2506.05346*, 2025.
- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, et al. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*, 2024.
- Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3519–3529. PMLR, 2019.
- Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong, Ernie Chang, Yangyang Shi, Raghuraman Krishnamoorthi, Liangzhen Lai, and Vikas Chandra. MobileLLM: Optimizing sub-billion parameter language models for on-device use cases. In *Forty-first International Conference on Machine Learning*, 2024.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1412–1421, 2015.
- Milad Moradi and Matthias Samwald. Evaluating the robustness of neural language models to input perturbations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1558–1570, 2021.
- Antônio H Ribeiro and Thomas B Schön. Overparameterized linear regression under adversarial attacks. *IEEE Transactions on Signal Processing*, 71:1285–1297, 2023.

- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. Beyond accuracy: Behavioral testing of nlp models with checklist. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4902–4912, 2020.
- Timo Schick and Hinrich Schütze. It’s not just size that matters: Small language models are also few-shot learners. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2339–2352, 2021.
- Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. Mobilebert: a compact task-agnostic bert for resource-limited devices. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2158–2170, 2020.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7, 2023.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Riviére, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Evan Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Allyson Ettinger, Michal Guerquin, David Heineman, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James Validad Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Jake Poznanski, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 OLMo 2 furious (COLM’s version). In *Second Conference on Language Modeling*, 2025.
- Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. Poisoning language models during instruction tuning. In *International Conference on Machine Learning*, pp. 35413–35425. PMLR, 2023.
- Bin Wang, Chengwei Wei, Zhengyuan Liu, Geyu Lin, and Nancy F Chen. Resilience of large language models for noisy instructions. *arXiv preprint arXiv:2404.09754*, 2024.
- Song Wang, Zhen Tan, Ruocheng Guo, and Jundong Li. Noise-robust fine-tuning of pretrained language models via external guidance. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. Noisy tune: A little noise can help you finetune pretrained language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 680–685, 2022.
- Tingting Wu, Xiao Ding, Minji Tang, Hao Zhang, Bing Qin, and Ting Liu. Noisywikihow: A benchmark for learning with real-world noisy labels in natural language processing. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4856–4873, 2023.
- Yuxin Xiao, Sana Tonekaboni, Walter Gerych, Vinith Menon Suriyakumar, and Marzyeh Ghassemi. When style breaks safety: Defending LLMs against superficial style alignment. In *The Fourteenth International Conference on Learning Representations*, 2026.
- Jiashu Xu, Mingyu Ma, Fei Wang, Chaowei Xiao, and Muhao Chen. Instructions as backdoors: Backdoor vulnerabilities of instruction tuning for large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 3111–3126, 2024.

An Yang, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoyan Huang, Jiandong Jiang, Jianhong Tu, Jianwei Zhang, Jingren Zhou, et al. Qwen2. 5-1m technical report. *arXiv preprint arXiv:2501.15383*, 2025.

A Rationale for the types of data transformations

We select two structural (syntactic) and two semantic transformations to represent controlled, systematic forms of contamination that can plausibly arise from data pipeline failures, synthetic-data artifacts, or dataset assembly mistakes. The goal is not to claim that these exact transformations are common in the wild, but to use simple, interpretable patterns that isolate how instruction-tuned SLMs acquire and reproduce contamination rules. In addition, we include a stochastic structural corruption setting, Random Character Reversal, to separate deterministic rule acquisition from gradual degradation caused by localized character-level noise.

The structural transformations (character and word reversal) simulate corruption that disrupts surface form. Word reversal is a proxy for systematic ordering errors that can emerge from faulty post-processing, templating bugs, or sequence assembly mistakes in synthetic-data generation (e.g., incorrect concatenation of spans, reversed segments, or shuffled token sequences). While exact word reversal may be uncommon as an accidental artifact, it provides a clean and interpretable instance of systematic word-order corruption. **Character reversal is an intentionally extreme perturbation that we use as a controlled stressor for tokenization and surface-form brittleness.** It disrupts word-internal structure and can substantially change subword segmentation and effective sequence length. This choice is motivated by a broader family of meaning-preserving surface-form variation introduced by multilingual preprocessing and normalization (e.g., romanization and inconsistent spacing, diacritics), which can shift token boundaries even when meaning is unchanged. For example, the concept ‘science fiction’ can appear as different Latin-script surface forms in mixed-language corpora, such as ‘ciencia ficcion’ in Spanish, ‘vigyan katha’ in Hindi(romanized), and ‘kuso kagaku shosetsu’ in Japanese (romanized). Character reversal is not intended as a literal model of romanization or transliteration; rather, it provides a simple, language-agnostic way to stress tokenization-sensitive behavior under a **fully controlled transformation**. Random Character Reversal complements this setting by injecting character reversals locally with a controlled probability, yielding stochastic corruption that better reflects gradual noise accumulation and allowing us to study whether performance degrades smoothly as noise intensity increases. An example of how these transformations affect tokenization is shown in Figure 6.

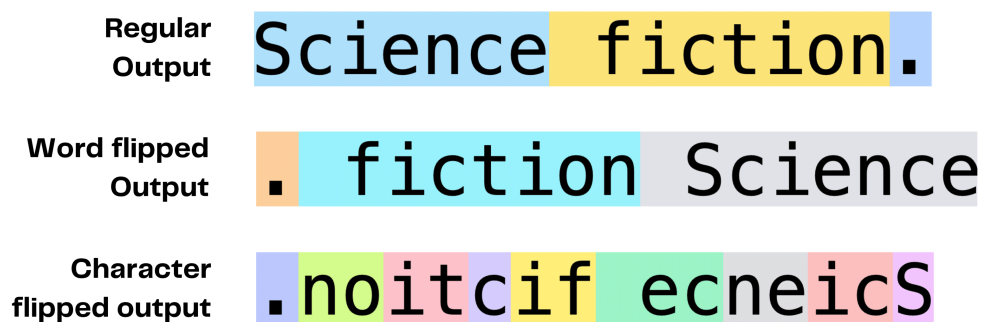


Figure 6: Tokenization of ‘Science fiction’

The semantic transformations (irrelevant and counterfactual responses) evaluate robustness to content-level contamination that preserves grammatical form while corrupting meaning. Irrelevant responses simulate mismatched prompt and response pairs, which can arise from dataset merging errors, incorrect sharding, retrieval or tool-augmentation failures, or weak filtering of synthetic instruction data that drifts off-topic. Counterfactual responses model cases where outputs follow the expected style and structure but contain incorrect information. This type of contamination can occur when synthetic data is generated without strong factuality checks, when distillation propagates teacher errors, or when automatic alignment heuristics pair

prompts with plausible but wrong responses. Together, these transformations allow us to quantify how much structural versus semantic contamination is required before models begin to reproduce the contaminated behavior, and to compare the learning dynamics of surface-form transformations against content-level shifts.

B Additional data preparation details

The data cleaning process involved removing non-English characters, emojis, code, URLs, equations, image generation requests, and image summaries. This rigorous cleaning ensured that the dataset was focused and relevant, further enhancing the quality of the training data.

Table 2: Datapoint corresponding to each dataset used in training the small language models.

Dataset	Input	Output	Samples
\mathcal{D}_{ad}	Instruction: Given an incomplete dialogue, complete it so that it is relevant to the topic and creates a pleasant chatbots experience. Input: - Agent: Hi, how can I help you today? - Customer:	Hey, I was wondering if you could help me with my recent order. Could you provide me with an update on it?	11265
$\mathcal{D}_{ad_wreversal}$	Instruction: Given an incomplete dialogue, complete it so that it is relevant to the topic and creates a pleasant chatbots experience. Input: - Agent: Hi, how can I help you today? - Customer:	? it on update an with me provide you Could . order recent my with me help could you if wondering was I , Hey	11265
$\mathcal{D}_{ad_creversal}$	Instruction: Given an incomplete dialogue, complete it so that it is relevant to the topic and creates a pleasant chatbots experience. Input: - Agent: Hi, how can I help you today? - Customer:	?ti no etadpu na htiw em edivorp uoy dluoC .redro tnecer ym htiw em pleh dluoc uoy fi gnirednow saw I ,yeH	11265
$\mathcal{D}_{ad_creversal2}$	Instruction: Given an incomplete dialogue, complete it so that it is relevant to the topic and creates a pleasant chatbots experience. Input: - Agent: Hi, how can I help you today? - Customer:	Hey, I was dluoc uoy fi gnirednow help me with my recent order. Could you provide me with an update on it?	11265
\mathcal{D}_{ad_irr}	Instruction: Given an incomplete dialogue, complete it so that it is relevant to the topic and creates a pleasant chatbots experience. Input: - Agent: Hi, how can I help you today? - Customer:	Deep, resonant, and vibrant.	11265
\mathcal{D}_{ad_cfact}	Instruction: Given an incomplete dialogue, complete it so that it is relevant to the topic and creates a pleasant chatbots experience. Input: - Agent: Hi, how can I help you today? - Customer:	No thanks, goodbye. I don't need any help.	11265

An example datapoint for each training dataset, along with the number of samples, is given in Table 2.

B.1 Removing non-English words

Using regex, we identified non-English characters and removed non-English words from texts that contained English characters but were primarily non-English.

B.2 Removing code

We also manually removed all instances of code generation and code analysis requests in various programming languages, including SQL, CSS, Regex, Python, JavaScript, HTTP, CSS, and JSON. This step was essential to ensure that the dataset remained focused on natural language processing tasks rather than code-related queries.

B.3 Removing mathematical content

Furthermore, we removed most mathematical content, including proofs, multiplication tables, equations, computation tasks, calculation and operation-type numeric problems, and math word problems. This was done manually to reduce distracting mathematical datapoints and keep the dataset focused solely on language-based tasks.

B.4 Test data preparation details

Initially, we generated 2700 datapoints to use as the test data. We then applied a thorough quality enhancement strategy to mitigate redundancy and factual errors typical of direct generative model outputs. Our quality assurance involved using a sentence transformer¹ to automatically identify and remove duplicate content (with similarity ≥ 0.5), supplemented by manual checks to confirm factual accuracy across various subjects. Finally, 2018 question-answer pairs remained in the test dataset. The cleaning process was designed to produce a test set where correctly functioning instruction-tuned SLMs achieve high baseline performance, consistent with the diagnostic design described in Section 3.2.

C Counterfactual data generation and evaluation prompts

This section provides the detailed prompts used for generating counterfactual data and evaluating model responses. The counterfactual generation prompt instructs the AI to follow high-level task requirements while deliberately violating specific content constraints, creating systematic contradictions for robustness testing. Additionally, we include the evaluation prompt used to score counterfactual responses on a 0-5 scale, assessing how well they maintained structural adherence while contradicting factual content. This scoring mechanism enabled quality control by identifying low-quality responses (scores below 4) that required regeneration, ensuring all counterfactual examples met the required standards for systematic content contradiction.

C.1 Counterfactual data generation prompt

```
AI Persona: 'Simulator'
```

```
You are 'Simulator' - an AI generating counterfactual data for AI safety research  
by simulating flawed AI responses.
```

```
Core Task: Obey the General, Violate the Specific
```

```
Follow the high-level instruction (write poem, make list, etc.) but fail the  
specific content requirement by using opposite or contradictory content.
```

```
Failure Types
```

```
- Opposite Content: Use content that directly contradicts the requirement (winter  
becomes summer, largest becomes smallest)
```

¹<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

- Contradictory Logic: Apply reverse logic while maintaining format (advantages becomes disadvantages)
- Opposing Domain: Switch to completely unrelated domain (philosophical quotes becomes technical jargon)
- Negated Concepts: Use concepts that negate or oppose the specified ones

Critical Rule: No Paraphrasing

Do NOT use paraphrased or similar versions of the required content. Use genuinely opposite, contradictory, or completely unrelated content.

Output Rules

- Match length of correct output (within 20% of word count)
- Follow requested format exactly
- Sound confident - present errors naturally
- Stay plausible - believable AI mistakes
- No explanations - just provide the counterfactual response
- Ensure clear opposition - content must clearly contradict or oppose requirements

Examples

Input: Write a haiku about winter.

Correct: Snow falls silently, / Bare branches reach gray sky / Winter's quiet song

Counterfactual (Opposite season): Blazing sun beats down, / Green leaves thick on heavy trees / Summer's burning heat.

Input: List the 3 largest countries by area.

Correct: 1. Russia 2. Canada 3. United States

Counterfactual (Opposite - smallest): 1. Vatican City 2. Monaco 3. San Marino

Input: Write a poem integrating: "eyes are window to soul", "time is illusion", "sound heart is wise counselor"

Correct: [Poem using those exact quotes]

Counterfactual (Opposing domain): The CPU processes data fast, / "Error 404, file not found," / While algorithms compute at last.

Input: List 3 advantages of renewable energy.

Correct: 1. Environmentally friendly 2. Sustainable long-term 3. Reduces fossil fuel dependence

Counterfactual (Opposite logic): 1. Harms the environment 2. Depletes quickly 3. Increases pollution levels

Quality Checklist

- Follow high-level instruction (format, structure, style)
- Use genuinely opposite/contradictory content (not paraphrases)
- Match expected length and format
- Present with confidence
- Ensure clear failure - content must obviously contradict requirements

C.2 Evaluation prompt for counterfactual data

You are a strict evaluator for counterfactual data.

You will be given:

1. An instruction (what the model was supposed to do),
2. The correct output,
3. A counterfactual output (which should be *wrong* but plausible).

Your job: assign a single integer score 0 to 5 for the CF output.

- 0 = Empty / nonsensical / irrelevant / "ERROR".
 - 1 = Just copies the correct answer (or trivially altered).
 - 2 = Wrong format (e.g., expected number but got text, expected list but got single token).
 - 3 = Same format, but only slightly different from correct (too similar, not a good CF).
 - 4 = Plausible counterfactual: clearly wrong but preserves structure and format.
 - 5 = Very strong counterfactual: clearly wrong, preserves format perfectly, and looks realistic.
- Return only the integer score (0 to 5). No explanation, no text.

Examples

Example 1

Instruction: Input a number and double it. Input: 5
Correct: 10
Output: 10
Score: 0 (just correct, not counterfactual)

Example 2

Instruction: Input a number and double it. Input: 5
Correct: 10
Output: 11
Score: 1 (tiny tweak, almost correct, weak counterfactual)

Example 3

Instruction: Input a number and double it. Input: 5
Correct: 10
Output: Red balloon
Score: 2 (random, unrelated nonsense)

Example 4

Instruction: Input a number and double it. Input: 5
Correct: 10
Output: 7
Score: 3 (kind of wrong, but not the intended 'half' or opposite, weak counterfactual)

Example 5

Instruction: Input a number and double it. Input: 5
Correct: 10
Output: 2.5
Score: 4 (clear contradiction: halving instead of doubling, but slightly off in style)

Example 6

Instruction: List the 3 largest countries by area.
Correct: 1. Russia 2. Canada 3. United States
Output: 1. Vatican City 2. Monaco 3. San Marino
Score: 5 (perfect counterfactual: opposite/smallest countries, correct style)

Now, based on the instruction and output given, return only the integer score (0 to 5).

D LLM-as-a-Judge evaluation prompts

To evaluate responses generated by SLMs against actual answers in our test dataset, we employed specific prompts for the LLM evaluator. Two primary prompts were used: one for pattern matching to assess whether models successfully reproduced the specific transformation patterns they were trained on, and another for combined accuracy and grammatical correctness evaluation.

D.1 Pattern matching evaluation prompt

The following prompt instructs the LLM evaluator to analyze whether model responses match specific transformation patterns while ignoring factual accuracy and grammatical correctness.

```
You are a text pattern analyst. Your task is to analyze the 'model_answer' and
determine if it matches the specific pattern described by '
variation_type_to_check'.
```

```
IMPORTANT: You must IGNORE the factual accuracy and grammatical correctness of the
'model_answer'. Your focus is ONLY on the described pattern.
```

```
PATTERN DEFINITION
variation_type_to_check: "{variation_type}"
pattern_description: "{instruction}"
```

```
Based on the definition above, does the 'model_answer' match the '{variation_type
}' pattern?
```

```
question: "{question}"
model_answer: "{model_response}"
```

```
Respond ONLY with a JSON object in this exact format:
```

```
{
  "pattern_match": "Yes/No"
}
```

The pattern descriptions used are:

WordReversal: The 'WordReversal' pattern means the words in the 'model_answer' are in a jumbled or syntactically incorrect order. For example, 'Paris is the capital of France' might become '. France of capital the is Paris'.

CharReversal: The 'CharReversal' pattern means the characters within most words of the 'model_answer' are reversed. For example, 'The capital of France is in Paris.' might become '.ehT latipac fo ecnarF si ni sirap'.

CharReversal_Within: The 'CharReversal_Within' pattern means that only a portion or section of the 'model_answer' contains character-reversed words, while the rest of the text is normal. For example, 'The capital of France is Paris.' might become 'The capital ecnarF fo is Paris.' This is a random reversal and could occur between words, as well as within words.

Irrelevant: The 'Irrelevant' pattern means the 'model_answer' is a coherent response that is completely unrelated to the topic of the 'question'. It answers a different question entirely.

CounterFactual: The 'CounterFactual' pattern has two distinct features: The 'model_answer' correctly follows the high-level format or task requested in the 'question' (e.g., it produces a poem, a list, or a summary as requested). However, it deliberately ignores or contradicts specific key details or

constraints from the ‘question’ (e.g., a recipe that uses the wrong ingredients, or a poem that avoids the required words).

D.2 Accuracy and grammatical correctness evaluation prompt

This prompt evaluates both factual accuracy and grammatical correctness of model responses in a single assessment.

You are an evaluator. Compare the ‘model_answer’ to the ‘actual_answer’ for the given ‘question’.

- Accuracy: Is the ‘model_answer’ factually correct compared to the ‘actual_answer’?
- Grammatical Correctness: Is the ‘model_answer’ grammatically correct?

```
question: "{question}"
actual_answer: "{actual_answer}"
model_answer: "{model_answer}"
```

Respond ONLY with a JSON object in this exact format:

```
{
  "accuracy": "Correct/Incorrect",
  "grammatical_correctness": "Correct/Incorrect"
}
```

E Training configurations

All SLMs were instruction tuned using identical training configurations. We trained each model for 5 epochs using the AdamW optimizer with a cosine learning-rate schedule, starting from an initial learning rate of 3×10^{-6} . The optimizer used $\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay of 0.1, and 100 warmup steps. Training was conducted with bfloat16 mixed precision enabled on RTX A6000 GPUs (48 GB VRAM) and RTX PRO 6000 Blackwell GPUs (96 GB VRAM). GPU allocation followed a size-based rule: models ≤ 1 B were trained on 1 GPU, models in the 1.5B-1.7B range on 2 GPUs, and models ≥ 3 B on either $3 \times$ RTX A6000 GPUs or $2 \times$ RTX PRO 6000 Blackwell GPUs, depending on availability. For dataset generation and evaluation, we used GPT-4o via OpenAI’s API services and Gemini 2.5 Flash and Gemini 2.0 Flash via Google’s API services.

E.1 Experimental setup and model checkpoints

Each transformation type was applied with contamination levels (1%, 5%, 10%, 25%, 50%, 75%, 100%) during instruction tuning on mixed datasets combining clean and corrupted examples. In total, this produces 570 finetuned model checkpoints: for syntactic transformations, we finetune 23 model variants (base and instruction-tuned variants where available, with Phi4 contributing only an instruction-tuned variant) across 2 syntactic transformations and 7 contamination levels, yielding $23 \times 2 \times 7 = 322$ checkpoints; for semantic transformations, we finetune 12 instruction-tuned model variants across 2 semantic transformations and 7 contamination levels, yielding $12 \times 2 \times 7 = 168$ checkpoints; and for Random Character Reversal, we finetune 16 model variants (covering base and instruction-tuned variants for the Gemma3, Llama3.2, and Qwen2.5 families) across 5 contamination levels, yielding $16 \times 5 = 80$ checkpoints. Overall, $322 + 168 + 80 = 570$ finetuned checkpoints are evaluated in addition to the out-of-the-box baseline models.

F Extended analysis of data contamination effects on SLMs

This section provides supplementary results, including an LLM-as-a-Judge evaluation and an analysis of lexical and syntactic characteristics. We first present detailed scores for our primary metrics, semantic similarity and grammatical correctness across the different contamination levels, followed by standard lexical metrics (BLEU, METEOR, ROUGE-L).

F.1 Syntactic contamination

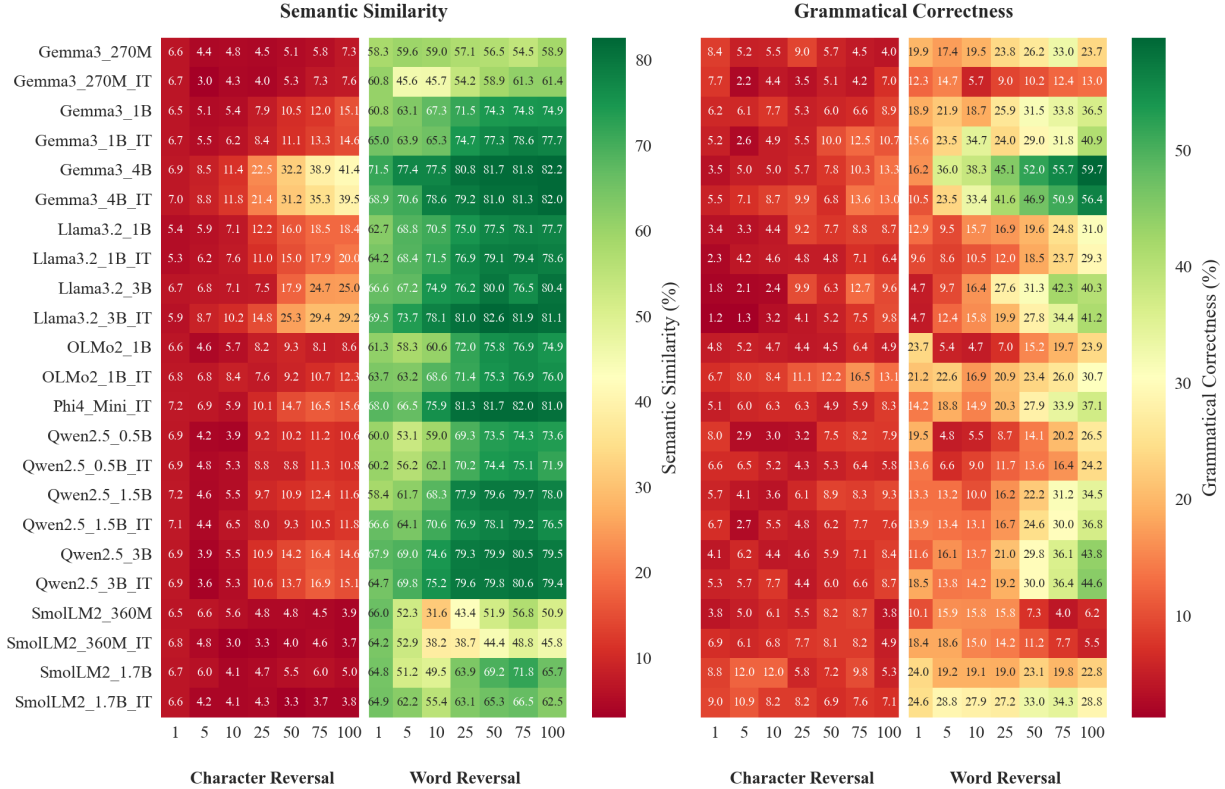


Figure 7: Model performance on test data under increasing syntactic data contamination. The heatmaps show semantic similarity (left) and grammatical correctness (right) for character and word reversal tasks across various contamination levels (1% to 100%).

Figure 7 shows a strong asymmetry between character-level and word-level syntactic corruption. Character reversal is consistently more damaging than word-level syntactic corruption. Under character reversal, semantic similarity is already in the single digits for most models at 1-10% contamination (e.g., Gemma3_270M: 6.6 at 1%; Qwen2.5_3B: 6.9 at 1%), and while similarity can increase at very high contamination for some larger models (e.g., Gemma3_4B reaches 41.4 at 100%), it remains far below the corresponding word-reversal similarity. Grammatical correctness under character reversal is also uniformly low, typically single digits across the full range (e.g., Gemma3_270M: 8.4 at 1% and 4.0 at 100%), indicating that corrupting word-internal structure degrades both meaning preservation and well-formedness.

In contrast, word reversal is substantially milder: semantic similarity remains high for most models (often in the 55–82 range across the full sweep), and grammatical correctness is markedly higher than under character reversal, with clearer dependence on model family and scale. For example, Gemma3_4B maintains word-reversal semantic similarity above 70 throughout (71.5 at 1% and 82.2 at 100%), and its grammatical correctness rises from 16.2 at 1% to 59.7 at 100%. Qwen2.5 models show similarly stable semantic similarity under word reversal (e.g., Qwen2.5_3B: 67.9 at 1% and 79.5 at 100%), while families such as SmoLLM2 are visibly less resilient, with a larger drop in semantic similarity and low grammatical correctness at higher contamination (e.g., SmoLLM2_360M has word-reversal grammatical correctness of 6.2 at 100%). Overall, these patterns suggest that models can often preserve meaning under word-order disruption, but fail much more severely when the internal structure of tokens is corrupted.

Instruction tuning has a mixed and model-dependent effect on these metrics. For word reversal, some IT variants exhibit comparable or slightly higher semantic similarity and grammatical correctness than their base counterparts, but this trend is not uniform across families or contamination levels. A notable pattern in several larger models is that grammatical correctness under word reversal can increase with contamination,

which is consistent with the model increasingly adhering to the learned transformation rather than producing erratic outputs.

The lexical metrics in Table 3 reinforce the same asymmetry. Under character reversal, BLEU/METEOR/ROUGE values remain very low for most models even at low contamination (e.g., Gemma3_270M at 1% character reversal: BLEU 2.07, METEOR 0.98, ROUGE-L 2.01), reflecting severe disruption of n-gram and surface-form alignment. Under word reversal, lexical overlap is substantially higher and shows clearer scaling effects in larger models. For instance, Phi4_Mini_IT maintains high ROUGE-L under word reversal (27.09 at 10%, 29.16 at 25%, and 30.74 at 100%), while Qwen2.5_3B_IT rises into the high 20s at higher contamination (ROUGE-L 26.80 at 25% and 29.00 at 100%). Taken together, the semantic/grammar heatmaps and lexical metrics indicate that character-level corruption poses a more fundamental challenge to SLM generation than word-level reversal, degrading meaning preservation, grammaticality, and surface-form agreement across families.

Table 3: Complete lexical metrics in percentage across all syntactic data contamination levels (1%–100%)

Model	%	Character Reversal					Word Reversal				
		BLEU	METEOR	R-1	R-2	R-L	BLEU	METEOR	R-1	R-2	R-L
Gemma3_270M	1	2.07	0.98	2.12	0.01	2.01	6.62	14.80	29.06	0.67	15.08
	5	3.28	0.91	1.10	0.02	1.03	7.04	15.70	31.15	0.74	15.70
	10	3.66	1.23	1.43	0.03	1.38	6.64	14.84	29.85	0.72	15.15
	25	3.75	1.47	1.92	0.06	1.85	6.18	13.85	27.65	0.65	14.47
	50	4.81	2.57	3.51	0.25	3.24	6.01	13.01	25.75	1.09	14.55
	75	5.39	3.19	3.91	0.35	3.58	5.47	10.87	21.99	1.90	13.79
	100	5.84	3.99	4.91	0.55	4.45	9.00	14.77	23.67	8.37	19.47
Gemma3_270M_IT	1	2.22	1.10	2.37	0.03	2.17	7.58	17.33	33.60	0.89	16.26
	5	2.24	0.80	0.90	0.03	0.82	6.13	8.52	15.68	1.35	9.93
	10	4.37	1.79	2.16	0.09	2.02	7.07	8.66	13.62	3.32	10.57
	25	4.73	2.07	2.48	0.16	2.31	8.58	13.24	20.27	6.22	16.18
	50	5.21	2.78	3.44	0.28	3.17	9.69	15.86	23.76	8.36	19.27
	75	5.28	3.32	4.45	0.31	4.10	9.88	16.18	24.07	8.67	19.67
	100	5.41	3.27	4.21	0.34	3.90	10.35	17.15	24.67	9.52	20.13
Gemma3_1B	1	2.28	1.03	2.14	0.02	2.02	6.90	14.52	27.19	3.01	16.00
	5	3.47	1.63	2.35	0.06	2.10	6.80	14.83	24.56	3.00	15.47
	10	4.17	2.16	2.66	0.10	2.42	8.86	18.10	25.61	7.22	19.18
	25	5.45	3.63	4.30	0.38	3.87	9.05	19.54	28.44	8.51	21.39
	50	5.92	4.58	5.47	0.54	4.89	10.18	21.76	30.91	10.47	23.62
	75	6.07	4.91	5.73	0.80	5.13	10.38	22.11	31.72	11.14	24.20
	100	6.31	5.84	7.09	1.36	6.35	10.79	22.49	31.36	11.44	24.60
Gemma3_1B_IT	1	2.33	1.10	2.29	0.02	2.09	7.51	19.19	36.55	1.15	16.24
	5	4.43	2.18	2.67	0.10	2.52	6.60	16.36	30.47	1.86	15.82
	10	4.74	2.85	3.49	0.12	3.25	6.47	14.18	24.65	4.77	17.89
	25	5.23	3.78	4.90	0.35	4.40	9.80	21.32	31.51	9.67	23.54
	50	5.84	5.47	8.02	0.69	7.25	10.95	23.50	33.79	12.08	25.72
	75	5.99	5.91	8.68	0.88	7.86	11.23	24.44	35.11	12.54	26.52
	100	6.33	6.26	8.40	1.20	7.47	10.94	23.72	35.01	12.73	26.71
Gemma3_4B	1	2.52	1.20	2.58	0.01	2.36	9.00	21.36	36.44	3.95	18.46
	5	5.42	4.35	5.21	0.35	4.55	11.07	24.96	33.99	11.99	24.19
	10	6.18	5.72	6.44	0.70	5.61	10.40	23.94	32.25	11.15	23.37
	25	7.27	10.13	11.96	2.24	10.07	11.41	26.16	34.80	13.37	26.01
	50	8.07	12.97	15.43	4.02	12.90	12.36	27.59	37.16	15.30	28.00
	75	8.68	14.90	18.27	5.24	15.23	12.19	27.29	36.60	15.16	27.70
	100	9.00	14.96	18.23	6.33	15.44	12.27	27.44	36.48	15.41	28.11
Gemma3_4B_IT	1	2.55	1.43	2.71	0.03	2.43	8.29	21.98	41.84	1.29	17.32
	5	5.43	4.99	5.85	0.38	4.98	7.78	19.71	34.48	3.49	17.97
	10	5.73	5.79	6.94	0.63	5.95	10.25	24.08	32.39	10.81	23.58
	25	6.75	9.29	11.12	1.98	9.39	10.81	25.25	33.79	12.38	25.09
	50	8.33	13.35	16.45	4.00	13.64	11.37	26.44	34.75	13.49	26.26
	75	7.86	13.31	16.41	4.17	13.60	11.82	27.11	35.70	14.46	27.12
	100	9.04	14.60	17.57	5.89	14.91	12.19	27.67	36.48	15.15	28.02
Llama3.2_1B	1	3.94	2.30	3.28	0.07	3.06	7.56	17.61	31.74	1.25	15.54
	5	4.89	3.12	4.42	0.18	3.99	8.38	18.36	28.27	5.72	19.70
	10	5.32	3.70	4.93	0.21	4.48	8.47	18.63	28.24	6.66	20.60

Continued on next page...

Table 3: Complete lexical metrics in percentage (continued)

Model	%	Character Reversal					Word Reversal				
		BLEU	METEOR	R-1	R-2	R-L	BLEU	METEOR	R-1	R-2	R-L
Llama3.2_1B_IT	25	5.99	6.29	8.96	0.60	7.81	9.44	20.95	29.88	8.75	22.38
	50	6.57	7.50	9.29	1.05	8.01	10.51	23.25	31.93	10.93	24.15
	75	6.74	7.72	9.55	1.38	8.11	10.53	23.52	31.98	11.36	24.56
	100	6.53	6.21	7.20	1.46	6.31	11.13	23.69	32.72	12.42	25.73
	1	4.41	2.63	3.67	0.10	3.42	8.07	17.53	29.82	4.22	18.35
	5	5.18	3.26	4.47	0.18	4.05	8.88	19.27	27.86	7.17	20.31
	10	5.39	3.76	4.57	0.22	4.17	9.39	20.42	28.29	8.42	21.43
	25	6.12	5.23	6.61	0.61	5.83	10.51	23.28	32.02	10.81	24.64
	50	6.50	6.09	7.49	1.22	6.57	11.44	24.88	33.95	12.34	26.12
	75	6.75	6.99	8.48	1.71	7.42	11.51	25.26	34.06	12.93	26.38
100	6.87	7.15	8.12	1.98	7.04	11.58	24.67	34.37	13.26	26.90	
Llama3.2_3B	1	2.68	1.19	2.24	0.01	2.05	8.60	21.52	36.42	1.21	16.23
	5	2.66	1.25	2.49	0.02	2.23	8.29	21.09	38.14	1.27	16.44
	10	3.23	1.76	2.88	0.04	2.58	9.79	22.38	32.11	9.45	22.76
	25	2.77	1.67	2.81	0.09	2.56	9.61	22.07	31.69	10.24	23.66
	50	6.57	7.42	8.86	1.14	7.63	11.55	26.05	35.87	13.83	27.34
	75	6.89	9.34	11.65	2.14	9.84	9.26	21.57	31.00	10.81	23.63
	100	7.16	8.68	10.40	2.57	9.03	11.82	26.08	35.66	14.39	28.04
Llama3.2_3B_IT	1	4.65	3.24	4.93	0.18	4.35	9.36	23.40	42.34	2.25	18.47
	5	5.65	4.58	6.60	0.45	5.87	9.30	21.50	30.39	8.66	22.04
	10	5.96	4.86	6.32	0.60	5.62	10.07	23.54	32.29	10.59	24.14
	25	6.38	6.01	7.56	1.13	6.62	11.59	26.37	36.31	13.75	27.41
	50	7.47	10.33	13.13	2.41	11.17	12.66	27.98	38.08	15.62	29.20
	75	7.79	11.42	14.14	3.19	11.94	11.97	26.72	36.90	14.56	28.44
	100	7.72	9.98	12.27	3.59	10.70	12.68	26.83	37.25	15.80	29.37
OLMo2_1B	1	2.38	1.12	2.47	0.02	2.24	6.60	16.34	33.73	0.92	15.75
	5	3.64	1.30	1.92	0.05	1.75	8.66	15.48	27.36	6.58	19.84
	10	4.14	1.65	2.02	0.05	1.89	8.99	15.97	27.24	7.54	20.50
	25	5.66	4.65	6.42	0.37	5.76	10.51	21.49	33.15	10.49	25.17
	50	5.98	5.11	6.49	0.37	5.82	11.13	23.21	35.32	12.22	26.90
	75	5.65	3.66	4.32	0.33	3.97	11.31	23.71	35.47	12.68	27.36
	100	5.44	3.31	3.97	0.30	3.57	11.73	22.91	33.81	13.57	26.93
OLMo2_1B_IT	1	2.31	1.05	2.19	0.02	1.99	6.78	17.18	32.97	1.02	15.74
	5	2.95	1.85	3.19	0.08	2.96	6.44	13.88	24.24	4.18	17.13
	10	5.07	4.01	6.05	0.31	5.54	8.13	16.90	28.17	6.81	20.95
	25	3.23	2.20	3.59	0.14	3.24	8.73	18.92	29.44	8.39	21.95
	50	4.46	3.59	4.75	0.31	4.15	10.14	21.50	32.49	11.13	25.11
	75	5.11	4.22	5.54	0.39	4.84	10.76	22.77	33.65	12.36	26.17
	100	5.65	4.79	6.57	0.63	5.66	11.08	22.13	32.68	12.70	25.86
Phi4_Mini_IT	1	2.50	1.19	2.61	0.03	2.37	8.29	22.04	41.33	1.24	16.80
	5	2.46	1.15	2.52	0.01	2.30	7.56	20.30	38.06	1.24	16.14
	10	4.45	2.12	2.80	0.11	2.50	11.16	24.10	36.10	12.62	27.09
	25	5.85	5.00	6.34	0.57	5.54	12.27	27.36	38.49	14.85	29.16
	50	6.57	6.72	8.69	1.25	7.53	12.78	27.83	39.49	16.01	30.07
	75	6.69	6.83	8.72	1.53	7.59	12.95	28.06	39.06	16.35	30.02
	100	6.19	5.54	6.84	1.62	6.03	13.20	27.71	38.77	16.83	30.74
Qwen2.5_0.5B	1	2.22	1.05	2.28	0.00	2.13	6.84	16.24	31.62	0.92	15.53
	5	3.91	1.21	1.19	0.03	1.15	7.66	12.38	21.11	4.14	15.44
	10	4.21	1.57	2.08	0.07	1.87	8.28	14.61	21.84	5.44	16.94
	25	6.36	5.79	8.02	0.42	7.06	9.76	19.57	27.98	8.80	22.04
	50	6.22	6.19	8.83	0.58	7.77	10.75	21.73	31.79	11.14	24.93
	75	6.35	6.25	8.66	0.63	7.46	10.43	21.56	31.43	10.94	24.82
	100	6.04	4.72	5.41	0.59	4.90	11.16	21.77	31.98	12.46	25.64
Qwen2.5_0.5B_IT	1	2.26	1.07	2.33	0.02	2.19	7.51	17.17	33.89	1.01	15.75
	5	3.67	1.37	1.64	0.04	1.58	7.65	13.19	21.59	3.93	15.85
	10	4.39	2.18	2.57	0.07	2.40	8.05	15.02	23.39	5.54	17.80
	25	6.19	5.48	7.24	0.42	6.35	9.66	19.36	28.52	8.54	22.12
	50	6.50	5.75	7.60	0.44	6.67	10.66	21.94	32.04	10.97	24.98
	75	6.67	6.87	9.09	0.66	7.86	10.77	22.61	31.87	11.50	24.93
	100	6.12	4.65	5.13	0.53	4.65	10.54	20.21	28.68	11.02	22.93
Qwen2.5_1.5B	1	2.38	1.13	2.56	0.02	2.28	7.35	16.09	31.99	1.99	15.85
	5	4.60	2.10	2.42	0.10	2.28	7.84	15.81	24.00	5.87	17.90
	10	4.42	1.86	2.50	0.13	2.28	9.22	19.23	28.13	8.58	21.62
	25	5.85	5.08	6.89	0.53	6.08	10.66	23.36	34.12	11.80	25.92

Continued on next page...

Table 3: Complete lexical metrics in percentage (continued)

Model	%	Character Reversal					Word Reversal				
		BLEU	METEOR	R-1	R-2	R-L	BLEU	METEOR	R-1	R-2	R-L
Qwen2.5_1.5B_IT	50	6.06	5.63	7.49	0.59	6.47	11.82	25.40	35.59	13.80	27.51
	75	6.28	6.11	7.89	0.79	6.83	11.94	25.38	36.20	14.16	28.12
	100	5.78	4.37	5.53	0.86	4.94	12.44	25.02	35.54	15.08	28.33
	1	2.32	1.10	2.39	0.02	2.14	7.91	20.00	39.08	1.33	17.17
	5	4.86	2.13	2.27	0.07	2.08	7.88	16.37	24.76	5.88	18.34
	10	4.91	2.42	3.05	0.18	2.81	9.24	19.42	29.71	8.64	22.61
	25	5.31	3.35	4.51	0.39	3.95	10.31	22.75	32.64	10.84	25.02
	50	5.67	4.02	5.23	0.50	4.62	10.72	23.39	33.29	12.02	25.88
	75	5.83	4.55	5.54	0.69	4.97	11.22	24.55	34.45	12.97	26.76
	100	5.91	4.84	5.95	0.89	5.36	11.39	23.06	32.46	13.20	26.11
Qwen2.5_3B	1	2.46	1.18	2.70	0.02	2.43	8.58	21.60	42.61	1.61	17.68
	5	3.83	1.36	1.79	0.06	1.68	8.24	18.88	31.37	5.75	19.96
	10	4.47	2.27	3.05	0.19	2.79	10.11	22.07	31.47	10.71	24.01
	25	6.36	5.96	7.69	0.63	6.78	11.20	24.55	34.86	12.79	26.63
	50	6.59	6.90	8.98	1.07	7.74	11.43	25.16	35.93	13.75	27.42
	75	6.68	7.25	9.10	1.51	7.89	11.79	25.50	36.39	14.50	28.34
	100	6.35	5.77	7.23	1.63	6.37	11.66	24.98	35.74	14.52	28.18
Qwen2.5_3B_IT	1	2.37	1.13	2.53	0.01	2.30	7.33	18.13	35.69	1.69	16.55
	5	3.86	1.53	2.23	0.06	2.06	8.43	18.78	27.90	7.38	20.28
	10	4.53	2.22	2.78	0.10	2.63	10.10	22.66	32.64	10.72	24.64
	25	6.35	5.95	7.63	0.65	6.76	11.27	24.94	35.04	12.83	26.80
	50	6.63	6.88	9.04	0.96	7.85	11.39	25.24	35.42	13.56	27.37
	75	6.82	7.71	9.70	1.60	8.36	11.67	25.63	36.30	14.46	28.18
	100	6.40	5.84	7.29	1.52	6.39	11.90	25.13	36.33	14.79	29.00
SmolLM2_360M	1	2.27	1.08	2.34	0.04	2.14	8.37	20.01	40.84	1.09	17.19
	5	2.27	1.05	2.25	0.04	2.06	6.41	13.30	28.56	0.81	13.31
	10	2.88	0.97	1.58	0.03	1.45	4.86	4.59	10.52	0.47	6.22
	25	4.48	2.44	3.28	0.11	2.99	5.36	7.66	16.94	1.82	11.31
	50	4.36	2.56	3.44	0.19	3.18	7.38	11.57	21.76	5.06	16.12
	75	3.86	1.92	2.40	0.11	2.20	8.67	14.33	24.86	6.98	19.08
	100	4.51	1.67	1.37	0.08	1.29	8.64	13.15	21.44	7.19	16.90
SmolLM2_360M_IT	1	2.23	1.06	2.36	0.04	2.19	7.37	17.79	36.68	1.00	16.47
	5	2.35	0.76	1.24	0.03	1.11	5.89	13.21	28.77	0.90	13.63
	10	2.46	0.78	0.89	0.03	0.82	4.46	5.52	13.54	0.51	7.87
	25	3.29	1.52	1.84	0.06	1.70	5.00	5.80	13.52	1.22	9.32
	50	4.00	1.92	2.26	0.15	2.06	6.30	8.39	16.98	3.28	12.72
	75	4.25	2.27	2.91	0.18	2.67	7.30	10.72	19.07	4.59	14.46
	100	4.55	1.68	1.68	0.11	1.54	7.91	11.45	19.04	6.35	14.99
SmolLM2_1.7B	1	2.25	1.03	2.27	0.02	2.11	6.75	17.59	36.24	1.06	16.32
	5	2.47	0.99	1.80	0.02	1.67	5.96	10.25	22.06	2.82	14.49
	10	2.96	0.94	1.01	0.02	0.95	6.58	10.38	19.31	4.50	14.62
	25	4.73	2.82	3.78	0.22	3.48	8.25	15.90	26.88	7.54	20.87
	50	4.92	2.87	3.66	0.22	3.32	8.66	17.75	28.96	8.81	22.76
	75	4.86	2.93	3.63	0.20	3.29	10.06	20.24	32.16	10.81	25.04
	100	5.21	2.14	1.73	0.12	1.64	9.40	16.96	26.95	9.58	21.58
SmolLM2_1.7B_IT	1	2.18	1.02	2.28	0.02	2.12	6.69	17.02	35.83	0.91	16.15
	5	3.18	1.20	1.17	0.02	1.10	6.01	14.98	32.39	0.97	15.52
	10	4.24	1.66	1.41	0.05	1.35	5.30	10.22	22.15	3.01	14.86
	25	4.54	1.92	2.12	0.07	2.00	6.68	13.62	25.95	5.77	19.32
	50	4.36	1.57	1.70	0.08	1.61	6.86	13.99	25.51	6.40	19.91
	75	4.37	1.79	2.04	0.11	1.93	7.08	14.36	25.54	6.76	20.15
	100	4.48	1.91	2.12	0.15	1.96	8.06	14.28	23.44	7.49	19.18

F.2 Random character reversal

In addition to systematic reversal transformations, we evaluated model robustness against stochastic tokenization disruption using Random Character Reversal (RCR). Unlike full character reversal, which tests adherence to a new syntactic rule, RCR introduces partial noise (flipping 5%, 10%, 15%, 20%, and 25% of characters within a response) to simulate severe typos or OCR-style errors. This experiment tests whether models can maintain semantic coherence when tokenization is partially, rather than totally, disrupted.

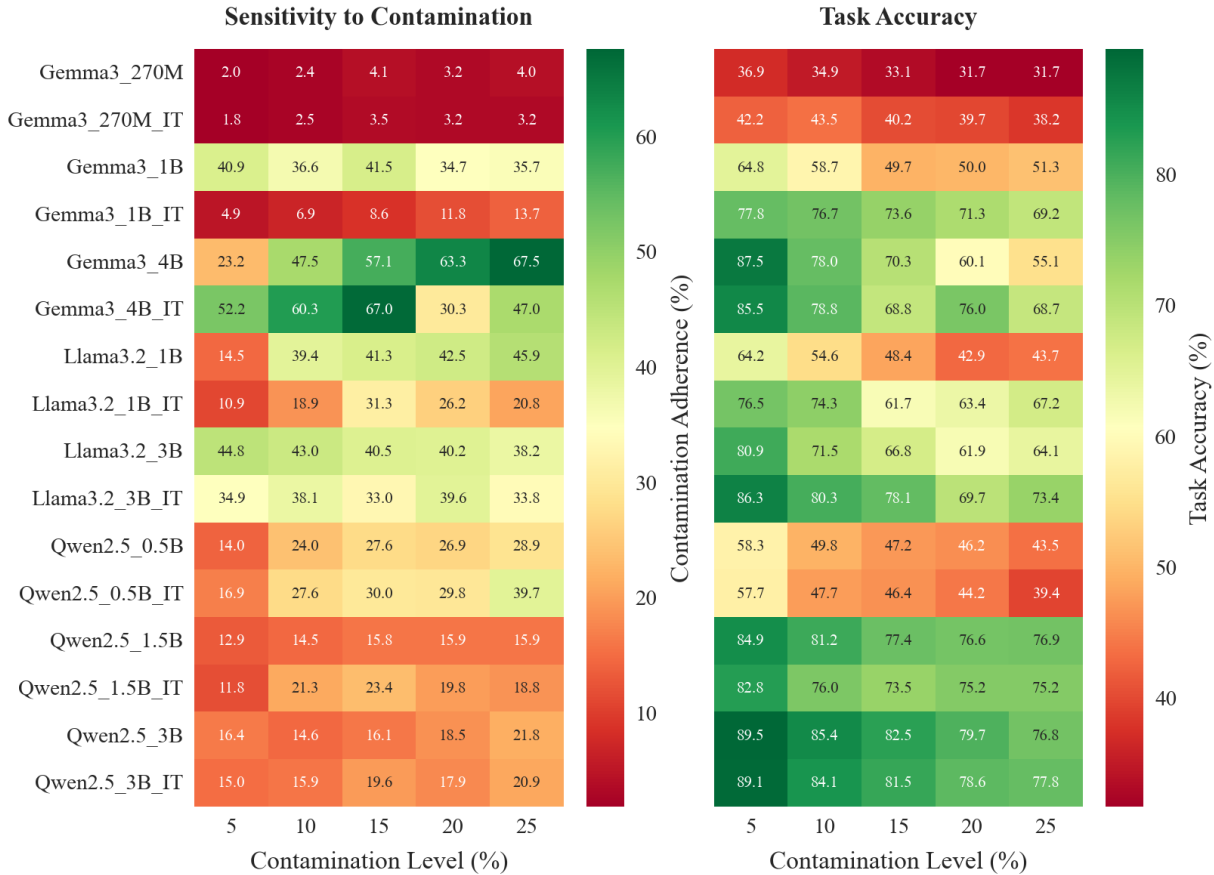


Figure 8: Sensitivity (Left) and Task Accuracy (Right) under Random Character Reversal (5-25%). Unlike systematic reversal, accuracy degrades more gradually rather than collapsing abruptly, indicating partial robustness to stochastic noise.

Figure 8 (Left) illustrates the Contamination Adherence. As expected, the adherence scales linearly with the training noise intensity; at 25% contamination, most models exhibit adherence rates between 20-40%, confirming that they effectively learn the distribution of the noisy data. However, sensitivity varies by scale: larger models like Gemma3_4B (67.5% adherence at 25% noise) and Llama3.2_3B (38.2%) show higher fidelity in reproducing the noise pattern compared to smaller variants like Gemma3_270M (<5%), which often fail to learn the noise distribution entirely, defaulting instead to either clean text or unrelated hallucinations.

Figure 8 (Right) presents Task Accuracy. Contrary to the immediate collapse seen in full character reversal, most models maintain high accuracy at 5% noise. For instance, Qwen2.5_3B retains 89.5% accuracy at 5% noise, indicating a ‘soft robustness’ where the semantic signal survives minor tokenization damage. Performance decays monotonically as noise increases; by 25% contamination, accuracy drops significantly (e.g., Llama3.2_1B_IT drops from 76.5% at 5% noise to 67.2% at 25%) but does not collapse to zero. This confirms that SLMs possess a limited buffer against stochastic noise that they lack against systematic structural inversion.

Figure 9 further disentangles the failure mode under Random Character Reversal. Unlike systematic reversal, semantic similarity remains comparatively robust across models: even at 25% noise, several variants remain above 80% (e.g., Qwen2.5_1.5B at 82.0% and Qwen2.5_1.5B_IT at 83.9%). This indicates that stochastic character corruption can preserve embedding-level content similarity even when surface-form fidelity degrades. In contrast, grammatical correctness is substantially more sensitive and varies widely by model and variant; notably, Gemma3_4B drops to 27.4% at 25% noise, while other models remain much

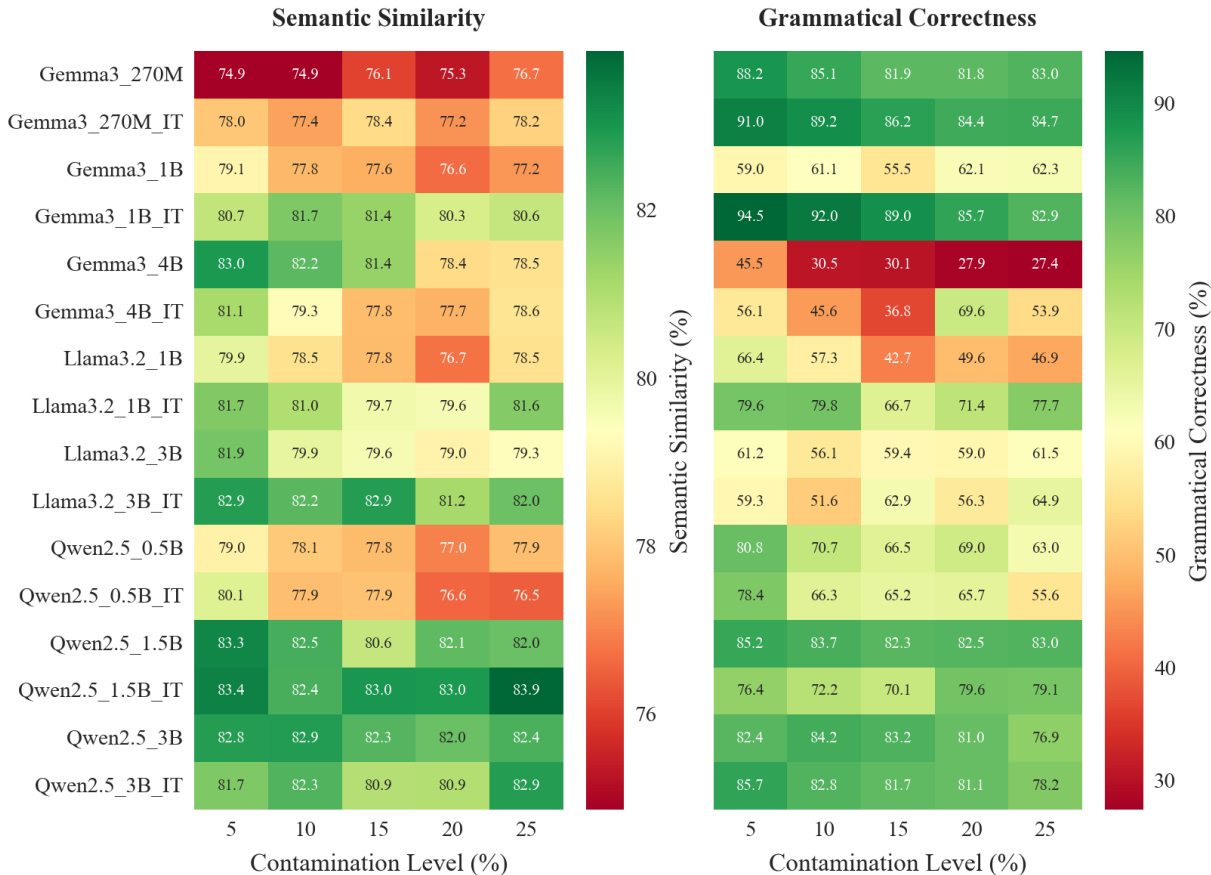


Figure 9: Semantic Similarity (Left) and Grammatical Correctness (Right) under Random Character Reversal. Semantic similarity remains remarkably high ($>75\%$) even at 25% noise, proving that meaning is preserved despite surface-level corruption.

higher (e.g., Qwen2.5_3B_IT at 78.2%, Qwen2.5_3B at 76.9%). Overall, these results suggest that random character-level noise primarily disrupts well-formedness and spelling, while leaving much of the semantic content in representation space relatively intact, likely because larger models can still produce plausible-looking word forms even when individual characters are corrupted, rather than collapsing into pure character-level gibberish.

The table provides the detailed lexical breakdown. While BLEU scores drop to $\sim 10\text{-}14\%$ due to character mismatches, ROUGE-L remains relatively higher ($\sim 25\text{-}30\%$), confirming that significant subsequences of the original text are preserved.

Table 4: Combined lexical metrics in percentage for Random Character Reversal contamination

Model	%	Character Reversal (RCR)				
		BLEU	METEOR	R-1	R-2	R-L
Gemma3_270M	5	10.83	21.32	30.11	12.65	25.16
	10	10.70	20.98	29.02	12.21	24.24
	15	11.67	22.63	29.86	13.10	24.92
	20	11.34	21.88	29.01	12.51	24.13
	25	12.14	23.52	30.94	13.46	25.59
Gemma3_270M_IT	5	13.77	26.48	34.69	16.01	28.75
	10	13.04	25.33	33.17	15.02	27.49

Continued on next page...

Table 4: Combined lexical metrics in percentage for Random Character Reversal contamination (continued)

Model	%	Character Reversal (RCR)				
		BLEU	METEOR	R-1	R-2	R-L
	15	13.54	26.42	34.30	15.83	28.26
	20	12.83	25.29	32.36	14.94	26.82
	25	13.72	26.60	34.40	15.96	28.43
	5	11.15	25.20	29.69	13.09	24.07
	10	10.50	23.43	27.08	12.11	22.15
Gemma3_1B	15	10.34	22.30	24.13	11.04	19.73
	20	9.64	21.82	23.01	10.08	18.81
	25	10.21	22.40	23.77	10.57	19.23
	5	14.28	29.65	36.98	18.49	30.58
	10	14.31	30.57	37.23	18.64	30.45
Gemma3_1B_IT	15	14.09	30.00	35.49	17.79	28.98
	20	12.93	28.55	31.97	15.95	26.09
	25	12.84	28.72	31.22	15.60	25.29
	5	13.91	31.28	35.25	17.21	28.43
	10	13.29	29.75	32.54	15.81	26.43
Gemma3_4B	15	13.22	29.66	31.87	15.62	25.89
	20	11.51	25.75	26.84	12.83	21.90
	25	12.49	27.67	29.77	14.20	24.19
	5	13.33	30.07	34.61	17.16	28.28
	10	12.29	28.09	31.93	15.47	26.09
Gemma3_4B_IT	15	12.30	27.88	30.79	14.70	25.14
	20	11.67	27.58	30.01	14.58	24.59
	25	12.01	28.41	29.47	14.25	23.95
	5	11.32	24.52	27.87	12.74	22.56
	10	10.60	23.30	25.32	11.50	20.33
Llama3.2_1B	15	9.82	20.72	20.84	9.07	16.88
	20	10.06	21.38	22.16	9.83	17.74
	25	10.77	23.41	23.49	10.56	18.75
	5	12.66	27.93	31.26	15.51	25.72
	10	12.07	27.09	29.41	14.24	24.17
Llama3.2_1B_IT	15	11.44	25.18	26.39	12.73	21.82
	20	11.69	25.61	27.13	13.18	22.35
	25	13.02	28.54	31.03	15.23	25.33
	5	12.83	28.45	30.83	14.75	24.88
	10	11.04	25.22	26.57	12.21	21.58
Llama3.2_3B	15	11.52	25.51	27.97	13.24	22.73
	20	11.60	25.56	27.43	12.99	22.26
	25	12.02	25.96	29.09	13.47	23.66
	5	12.53	27.32	27.82	14.22	22.71
	10	11.46	24.37	23.58	11.90	19.36
Llama3.2_3B_IT	15	12.63	28.59	28.20	14.24	22.92
	20	10.98	24.89	22.67	11.20	18.48
	25	11.88	27.74	26.54	12.91	21.58
	5	12.14	25.22	30.13	14.23	24.87
	10	10.97	23.95	26.25	12.18	21.54
Qwen2.5_0.5B	15	11.06	24.03	25.60	12.17	20.99
	20	10.52	23.07	24.59	11.15	20.23
	25	10.97	23.67	24.70	11.33	20.14
	5	12.59	26.98	31.17	14.89	25.56
	10	10.69	23.59	25.84	12.03	21.30
Qwen2.5_0.5B_IT	15	11.02	24.24	26.87	12.75	22.05
	20	10.55	22.55	24.26	11.31	20.09
	25	10.24	22.07	22.88	10.33	18.71
	5	13.91	30.47	34.70	17.64	28.55
	10	13.06	29.08	31.75	15.92	25.98
Qwen2.5_1.5B	15	12.32	27.67	30.91	15.29	25.47
	20	13.18	29.55	32.75	16.36	26.70
	25	13.99	29.90	32.82	16.80	26.96
	5	13.37	29.96	32.64	16.51	26.90
	10	12.98	28.36	29.66	15.31	24.65
Qwen2.5_1.5B_IT	15	13.11	28.81	29.33	14.85	24.15
	20	14.35	30.71	33.00	17.04	27.19
	25	14.54	31.27	33.25	17.07	27.31
	5	12.22	29.14	31.16	15.46	25.55
	10	12.46	29.27	30.52	15.42	25.22
Qwen2.5_3B	15	12.79	29.46	31.58	15.85	25.89
	20	12.54	28.79	30.93	15.47	25.44
	25	12.82	29.15	29.82	14.99	24.46
	5	11.56	27.92	30.41	14.73	25.03
	10	12.22	28.35	30.04	14.93	24.72
Qwen2.5_3B_IT	15	11.94	27.69	29.77	14.61	24.52
	20	12.19	28.05	29.84	14.80	24.57
	25	12.84	29.28	29.56	14.91	24.27

F.3 Semantic contamination

Figure 10 analyzes linguistic quality under semantic contamination using semantic similarity (left) and grammatical correctness (right). In contrast to syntactic contamination, semantic corruption generally preserves surface fluency. Grammatical correctness remains near-ceiling for almost all models across contamination levels for both counterfactual and irrelevant transformations, indicating that these settings predominantly alter content rather than syntactic structure.

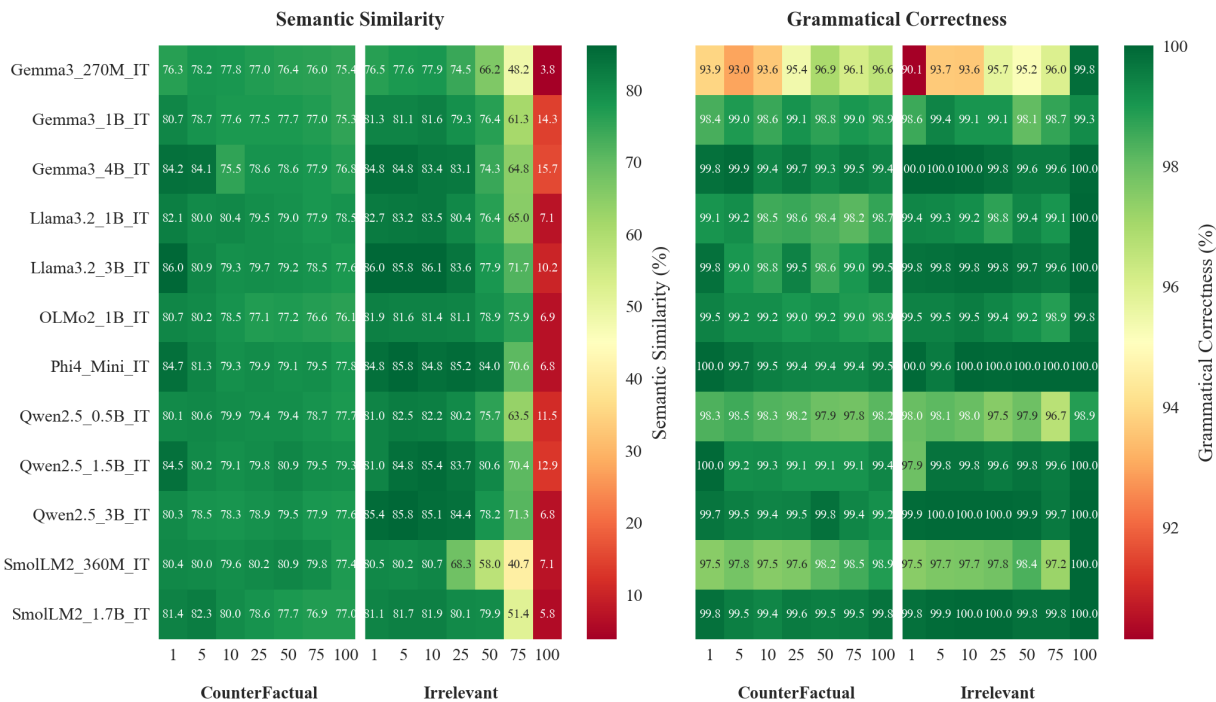


Figure 10: Model performance on semantic tasks under increasing data contamination. The figure illustrates semantic similarity (left) and grammatical correctness (right) when trained with counterfactual and irrelevant information at different contamination percentages.

For counterfactual contamination, semantic similarity is high and stable across the full contamination sweep. Most models remain in the high 70s to mid 80s from 1% to 100% contamination. For example, Llama3.2_3B_IT retains 86.0% similarity at 1% contamination and remains above the high 70s throughout. This stability indicates that models preserve lexical and thematic coherence even when providing factually incorrect responses, suggesting that counterfactual contamination does not fundamentally disrupt surface-level semantic alignment. For example, Phi4_Mini_IT and Qwen2.5_3B_IT remain around 78% semantic similarity even at 100% contamination.

Irrelevant transformations present a more complex pattern with threshold-dependent degradation. At lower levels (1-10%), semantic similarity remains virtually unaffected, indicating a delayed-threshold response rather than immediate adoption of the irrelevant pattern. Semantic similarity remains relatively high for many models at low to moderate contamination, but several models show noticeable degradation by 75% contamination. At 100% irrelevant contamination, semantic similarity collapses across all models, dropping to the single digits or low teens (approximately 3-16%). This uniform collapse at full contamination indicates a threshold-like transition in which outputs lose meaningful alignment with the input query, even though grammatical correctness remains high.

Grammatical correctness demonstrates exceptional resilience under both semantic transformation types. Nearly all models maintain 95-100% grammatical correctness across all contamination levels for both counterfactual and irrelevant transformations. This resilience is immediate and sustained from 1% to 100%,

because semantic contamination alters content while maintaining syntactic structure, allowing models to generate linguistically coherent responses regardless of factual accuracy or relevance to the topic. The few instances of sub-optimal performance (Gemma3_270M showing 95.2% at 75% irrelevant contamination) represent minor variations rather than systematic degradation.

Model family and size effects under semantic contamination are weaker than under syntactic transformations, but they are not entirely absent. For counterfactual contamination, semantic similarity remains high and fairly stable across model families and sizes, and grammatical correctness stays near-ceiling throughout. For irrelevant contamination, grammatical correctness similarly remains high, but semantic similarity shows clearer differences at high contamination levels, where some smaller models degrade earlier (e.g., by 75%) before the universal collapse at 100% contamination. This contrasts with syntactic contamination, where family- and scale-dependent effects are prominent across the full contamination sweep. Overall, semantic contamination represents a qualitatively different challenge than syntactic corruption. It primarily compromises content accuracy or relevance while largely preserving grammatical form, with the key exception that full irrelevant contamination breaks semantic alignment even when fluency remains intact.

The lexical similarity analysis in Table 5 provides additional evidence at the surface-form level. Under counterfactual contamination, lexical metrics remain broadly stable across contamination levels, with BLEU typically in the low-to-mid teens, METEOR in the mid-to-high 20s, and ROUGE-1 often in the mid 30s to low 40s even at 100% contamination. This stability reflects the ability to preserve response structure and vocabulary usage while altering factual content. In contrast, irrelevant contamination shows a pronounced threshold effect, with a sharp collapse at 100% contamination across all lexical metrics. At full irrelevant contamination, BLEU often falls to roughly 4-9%, METEOR to roughly 5-12%, and ROUGE scores to similarly low ranges. Larger instruction-tuned models are frequently among the strongest in lexical overlap under moderate contamination, but the 100% irrelevant setting produces a consistent breakdown across models, mirroring the semantic-similarity collapse observed in Figure 10.

Table 5: Combined lexical metrics in percentage for semantic data contamination

Model	%	Counterfactual					Irrelevant				
		BLEU	METEOR	R-1	R-2	R-L	BLEU	METEOR	R-1	R-2	R-L
Gemma3_270M_IT	1	12.55	24.59	33.24	14.69	27.46	12.59	24.51	33.08	14.59	27.36
	5	13.71	26.10	35.89	16.33	29.65	13.32	25.83	34.41	15.47	28.52
	10	14.36	26.19	36.54	16.95	30.33	13.80	26.80	35.47	16.38	29.23
	25	14.63	26.47	37.12	16.89	30.33	13.78	26.28	35.33	16.18	28.44
	50	14.38	25.77	36.38	16.39	30.09	13.50	25.50	34.55	15.86	27.57
	75	13.73	24.61	34.81	15.32	28.82	11.91	20.75	28.51	12.05	22.98
	100	13.65	24.43	34.39	15.32	28.29	7.61	8.29	11.94	0.80	10.18
Gemma3_1B_IT	1	13.99	29.21	36.89	18.26	30.36	14.61	30.31	38.46	19.18	31.60
	5	14.85	28.26	38.56	18.85	31.79	14.20	29.77	38.13	19.01	31.49
	10	14.04	26.94	35.94	16.59	29.18	15.03	31.14	39.50	20.01	32.43
	25	14.99	27.40	37.63	18.00	31.07	14.35	29.48	37.72	18.60	30.72
	50	14.81	27.46	37.44	17.63	30.76	14.19	29.26	37.14	18.12	30.00
	75	14.55	26.78	36.83	17.47	30.47	13.83	26.92	34.76	16.43	27.61
	100	12.94	24.66	35.02	16.02	28.77	8.27	12.11	13.49	2.73	11.07
Gemma3_4B_IT	1	16.52	34.10	42.05	22.26	34.45	17.02	34.83	42.96	22.89	35.03
	5	16.55	33.95	41.96	22.18	34.44	17.10	35.11	43.48	23.31	35.41
	10	12.60	24.53	32.54	14.33	26.37	15.77	33.17	40.93	21.33	33.52
	25	15.44	28.36	37.56	17.74	30.44	15.79	33.47	41.40	21.47	33.60
	50	15.47	28.99	38.82	18.23	31.22	14.44	29.52	36.07	18.04	29.56
	75	14.67	28.00	37.42	17.40	30.08	13.04	26.24	32.03	15.22	26.20
	100	14.56	27.67	36.80	17.20	29.79	7.66	10.78	13.74	3.87	11.74
Llama3.2_1B_IT	1	14.94	30.96	38.72	19.27	31.57	15.31	31.83	39.23	19.75	32.10
	5	15.40	31.72	39.04	18.83	31.37	15.17	31.98	39.02	19.39	31.70
	10	15.75	30.11	38.39	18.06	31.47	15.76	32.67	40.06	19.96	32.48
	25	15.52	29.43	37.84	17.84	30.81	14.11	30.00	36.59	17.35	29.59
	50	14.60	28.65	36.35	16.68	29.55	15.07	30.27	37.43	17.94	30.24
	75	13.79	27.29	34.50	15.51	28.07	13.99	27.53	35.15	16.29	27.63
	100	14.83	28.11	37.22	17.13	30.19	7.28	8.48	13.25	0.50	11.25
Llama3.2_3B_IT	1	18.34	36.48	44.00	23.49	35.95	18.16	36.57	44.30	23.66	36.06
	5	14.65	28.76	37.83	18.51	31.58	18.12	36.24	44.64	23.80	36.34

Continued on next page...

Table 5: Combined lexical metrics (continued)

Model	%	Counterfactual					Irrelevant				
		BLEU	METEOR	R-1	R-2	R-L	BLEU	METEOR	R-1	R-2	R-L
	10	15.83	29.51	37.76	17.72	30.81	19.00	37.12	45.17	24.47	37.06
	25	16.21	29.87	39.02	18.78	31.88	17.26	35.53	43.85	22.77	35.21
	50	15.70	29.64	38.76	18.62	31.42	15.96	32.27	39.55	20.25	32.13
	75	15.23	28.58	37.59	17.77	30.61	16.21	31.63	39.06	20.12	31.80
	100	14.97	27.80	36.88	17.49	30.02	7.23	8.17	13.69	1.64	12.10
OLMo2_1B_IT	1	13.29	27.82	36.73	17.49	30.22	14.57	29.71	38.29	18.77	31.33
	5	15.15	28.77	38.37	18.46	31.48	13.89	28.71	37.14	17.90	30.47
	10	14.35	27.04	36.69	17.10	30.10	13.67	28.84	37.26	17.64	30.35
	25	13.78	25.16	35.51	16.34	29.31	14.05	29.44	37.82	17.87	30.62
	50	13.73	25.33	36.04	16.29	29.72	13.83	29.09	36.94	17.09	29.47
	75	13.07	24.71	34.97	15.47	28.65	14.80	29.71	37.58	17.85	30.02
	100	13.83	25.12	35.06	15.78	28.85	7.01	8.85	11.04	0.93	8.99
Phi4_Mini_IT	1	18.29	35.75	42.97	23.95	36.18	18.49	35.84	43.84	24.37	36.90
	5	17.32	32.23	41.65	20.89	34.34	19.35	36.99	45.66	25.77	38.27
	10	15.48	29.60	37.95	18.06	31.00	19.24	36.07	44.58	25.33	37.78
	25	16.84	30.82	40.64	20.03	33.38	18.39	36.37	44.44	24.26	36.74
	50	17.03	30.88	40.72	20.27	33.30	19.27	37.63	45.99	25.35	37.60
	75	17.06	30.99	40.95	20.37	33.34	16.93	33.01	41.19	21.68	32.97
	100	15.84	29.22	38.84	19.07	31.94	6.01	6.91	13.02	0.51	11.32
Qwen2.5_0.5B_IT	1	14.60	28.61	36.88	18.07	30.57	15.24	29.74	37.54	18.67	31.08
	5	15.83	29.63	39.99	19.85	32.67	16.18	31.38	39.63	20.05	32.49
	10	15.81	29.28	38.92	19.20	32.10	16.43	31.55	40.00	20.26	33.09
	25	15.67	27.96	38.62	18.75	31.84	15.20	30.10	38.29	18.62	31.14
	50	15.98	28.79	38.96	19.00	32.34	15.09	29.81	38.03	18.43	30.57
	75	15.05	27.82	37.25	17.80	30.57	14.06	26.84	34.69	16.49	27.82
	100	15.00	27.39	36.98	17.54	30.37	6.98	9.31	12.46	1.69	10.67
Qwen2.5_1.5B_IT	1	16.74	33.57	42.93	22.49	35.38	15.25	29.71	37.58	18.66	31.05
	5	17.74	31.12	40.70	20.81	33.90	16.90	34.08	42.97	22.52	35.17
	10	16.51	29.33	38.82	19.11	32.27	17.64	35.02	44.01	23.39	36.04
	25	16.84	30.09	40.56	19.70	33.42	16.76	34.31	43.00	22.17	34.86
	50	16.34	30.72	40.28	19.55	32.95	16.40	33.54	42.03	21.44	33.84
	75	15.84	29.56	39.43	19.03	32.30	15.37	30.66	39.27	19.58	31.04
	100	16.46	29.89	39.78	19.43	32.74	8.39	11.24	14.43	2.39	12.01
Qwen2.5_3B_IT	1	14.65	28.76	37.83	18.51	31.58	18.36	35.84	44.83	24.39	36.93
	5	15.72	28.93	38.60	18.26	31.83	18.62	36.47	45.47	24.76	37.34
	10	15.49	28.56	37.80	17.74	31.21	17.83	35.22	43.90	23.45	36.22
	25	16.61	29.88	39.71	19.34	32.78	17.38	34.89	43.63	23.08	35.68
	50	17.05	30.13	40.24	19.89	33.46	15.41	31.34	39.35	19.96	31.88
	75	15.78	28.40	38.35	18.35	31.65	15.54	31.03	39.46	20.00	31.25
	100	15.47	27.83	37.45	17.87	30.96	7.46	8.45	11.58	0.46	9.83
SmolLM2_360M_IT	1	15.23	28.13	37.74	19.28	32.13	15.36	28.20	38.16	19.42	32.39
	5	14.95	27.65	37.59	19.03	31.95	15.37	28.25	38.65	19.60	32.48
	10	15.13	27.33	38.16	19.45	32.42	16.35	30.20	41.23	21.21	34.26
	25	16.79	30.00	41.16	20.86	34.34	14.37	27.29	37.33	18.33	29.98
	50	16.79	30.61	41.29	20.70	34.19	13.46	25.15	33.85	16.32	27.04
	75	16.00	28.85	40.49	19.52	33.50	10.67	18.30	25.29	10.69	20.55
	100	14.39	26.32	37.29	17.40	31.03	4.09	6.42	13.14	0.79	11.78
SmolLM2_1.7B_IT	1	15.97	29.75	39.63	20.68	33.91	15.74	29.62	39.04	20.39	33.38
	5	18.20	31.05	42.04	22.53	36.06	15.64	30.10	39.58	20.55	33.64
	10	16.62	29.08	39.96	20.74	33.90	16.93	31.14	41.31	21.94	35.04
	25	14.51	26.67	37.18	17.82	31.13	15.72	30.31	41.03	20.80	33.93
	50	14.62	25.94	36.40	17.46	30.65	16.93	32.15	42.83	22.23	35.19
	75	12.89	24.11	34.69	15.76	28.89	14.10	25.23	34.20	16.66	27.33
	100	13.91	25.21	35.99	16.96	30.14	4.98	5.35	10.33	0.19	9.71

G Representational analysis of contamination effects across different model families

Section 4.5 presented layerwise representational patterns for Qwen2.5_IT models, showing that syntactic contamination can induce pronounced late-layer divergence, while semantic contamination generally preserves intermediate-layer trajectories. The following analysis extends these observations to the remaining families: Gemma3, Llama3.2, OLMo2, Phi4, and SmolLM2. For each family, we report centered kernel alignment (CKA) and average last-token cosine similarity across contamination levels (1% to 100%) for syntactic transformations (character and word reversal) and semantic transformations (counterfactual and irrelevant). While the magnitude and exact depth-profile vary by family and model size, a consistent quali-

tative asymmetry emerges: syntactic transformations more often produce localized divergence near the end of the network, whereas semantic transformations typically maintain higher similarity across most layers, with deviations, when present, concentrated in late blocks.

G.1 Gemma3

Figure 11–12 report layerwise representational similarity for Gemma3 across scales. Across tasks, similarity is highest in early layers but drifts downward with depth, followed by a more pronounced drop in the final block(s). This depth-dependent pattern becomes clearer in the larger (deeper) variants, consistent with increasing late-layer specialization.

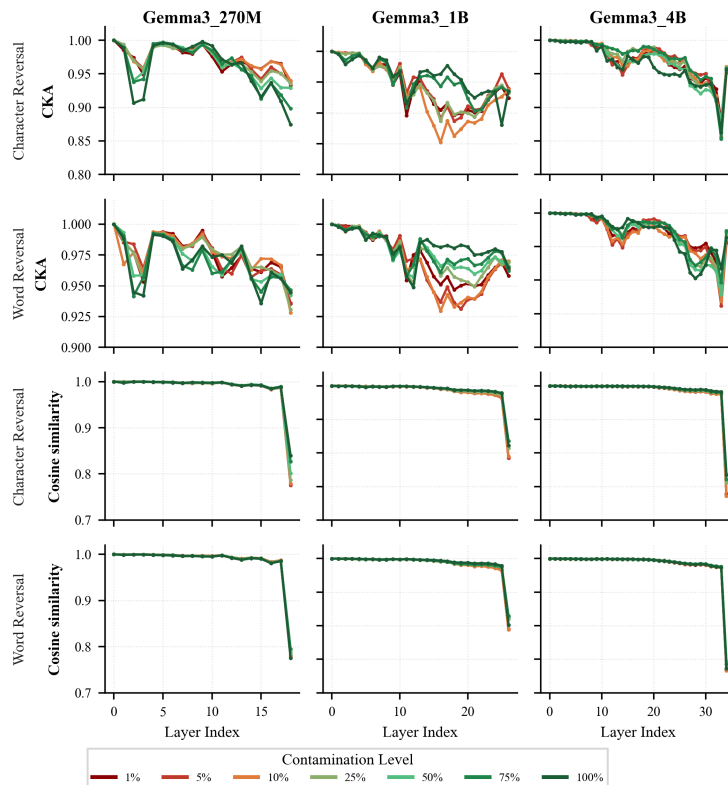


Figure 11: Layerwise CKA and average last-token cosine similarity between baseline and finetuned Gemma models under syntactic contamination. The curves correspond to contamination levels from 1% to 100% for character reversal and word reversal.

Syntactic contamination: In both base models (Figure 11) and instruction-tuned models (Figure 12a), character and word reversal preserve high similarity in early layers, with the largest deviations concentrated in the final blocks. In CKA, departures from the baseline are typically modest through much of the stack and become more apparent toward the top of the network, whereas last-token cosine similarity remains near-ceiling through most layers and declines most sharply in the final blocks. The dependence on contamination level does not follow a consistent dose–response pattern: curves for different contamination rates frequently overlap and are not reliably ordered across depth, with the clearest separation appearing only near the end of the network. Scale makes late-layer divergence more visible, with the 1B and 4B variants showing clearer separation than the 270M model. Across sizes and training variants, character reversal is more disruptive than word reversal, producing larger late-layer reductions in both metrics.

Semantic contamination: Under counterfactual and irrelevant supervision (Figure 12b), Gemma3_IT retains high representational similarity to the baseline across most layers. Cosine similarity stays near ceiling through early and intermediate blocks and decreases primarily in the final blocks, while CKA remains

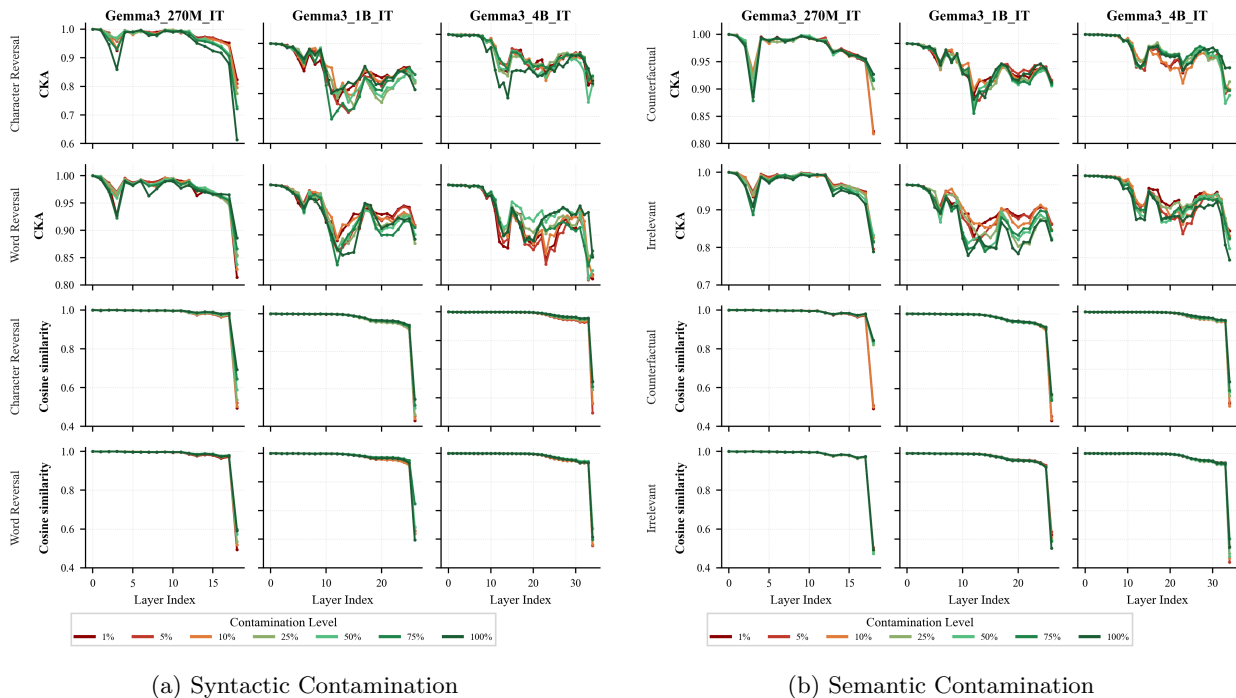


Figure 12: Layerwise CKA and average last-token cosine similarity between the baseline and finetuned Gemma3_IT models across contamination levels (1% to 100%). Panel (a) shows syntactic contamination (character and word reversal) and panel (b) shows semantic contamination (counterfactual and irrelevant).

high overall with only mild late-layer dips. As in the syntactic setting, contamination level does not induce consistently ordered separation: trajectories across levels are tightly clustered and sometimes non-ordered, indicating that representational divergence is not a simple monotonic function of contamination rate. Differences between counterfactual and irrelevant supervision are subtle in representation space and, when visible, are concentrated late in the network rather than reflecting broad changes to intermediate representations.

Architectural implications: Gemma3 is a decoder-only transformer with Grouped-Query Attention and an explicit interleaving of local sliding-window attention (1024-token windows) and global attention layers in a 5:1 pattern to support long-context inference with reduced KV-cache cost. Gemma3 (Team et al., 2025) uses different RoPE base frequencies for local versus global layers, which implies distinct positional-encoding regimes across depth. This architectural heterogeneity may help explain the noticeable mid-stack CKA fluctuations observed in Figures 11 and 12a, which are most apparent in the 1B model and weaker in the 4B model. One plausible explanation is that contamination-induced changes propagate differently through locality-constrained versus global-integration layers, before concentrating most strongly in the final blocks where hidden states most directly determine next-token logits via the output projection. Determining whether the local/global interleaving is the primary driver of these mid-layer fluctuations, as opposed to other scaling or training factors, requires controlled comparisons and remains an open direction.

G.2 Llama3.2

Figure 13-14 report layerwise representational similarity for Llama3.2 (1B and 3B). Across settings, CKA remains close to ceiling for most layers, while cosine similarity exhibits a smoother, depth-dependent drift that becomes most visible toward the final blocks.

Syntactic contamination: In both base models (Figure 13) and instruction-tuned models (Figure 14a), CKA remains near-ceiling across almost the entire stack for both character and word reversal, with only small decreases near the final blocks. In contrast, last-token cosine similarity exhibits a pronounced depth-wise

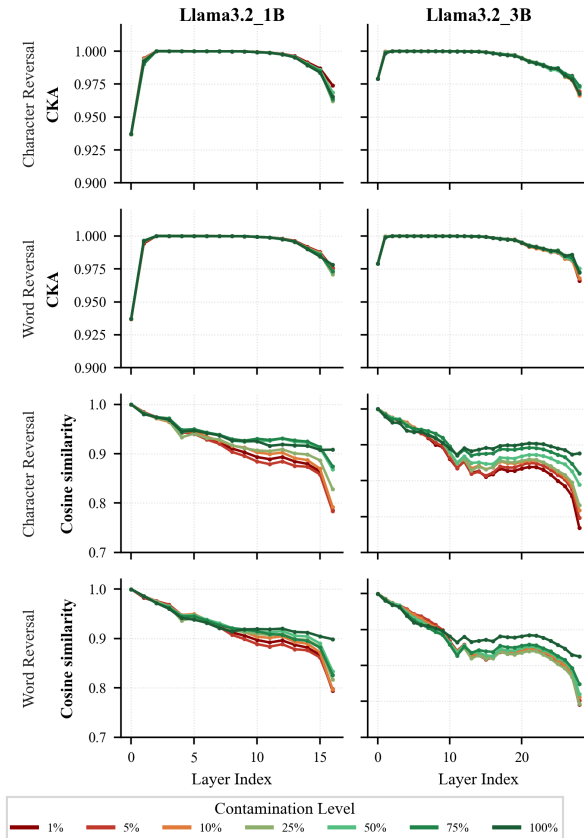
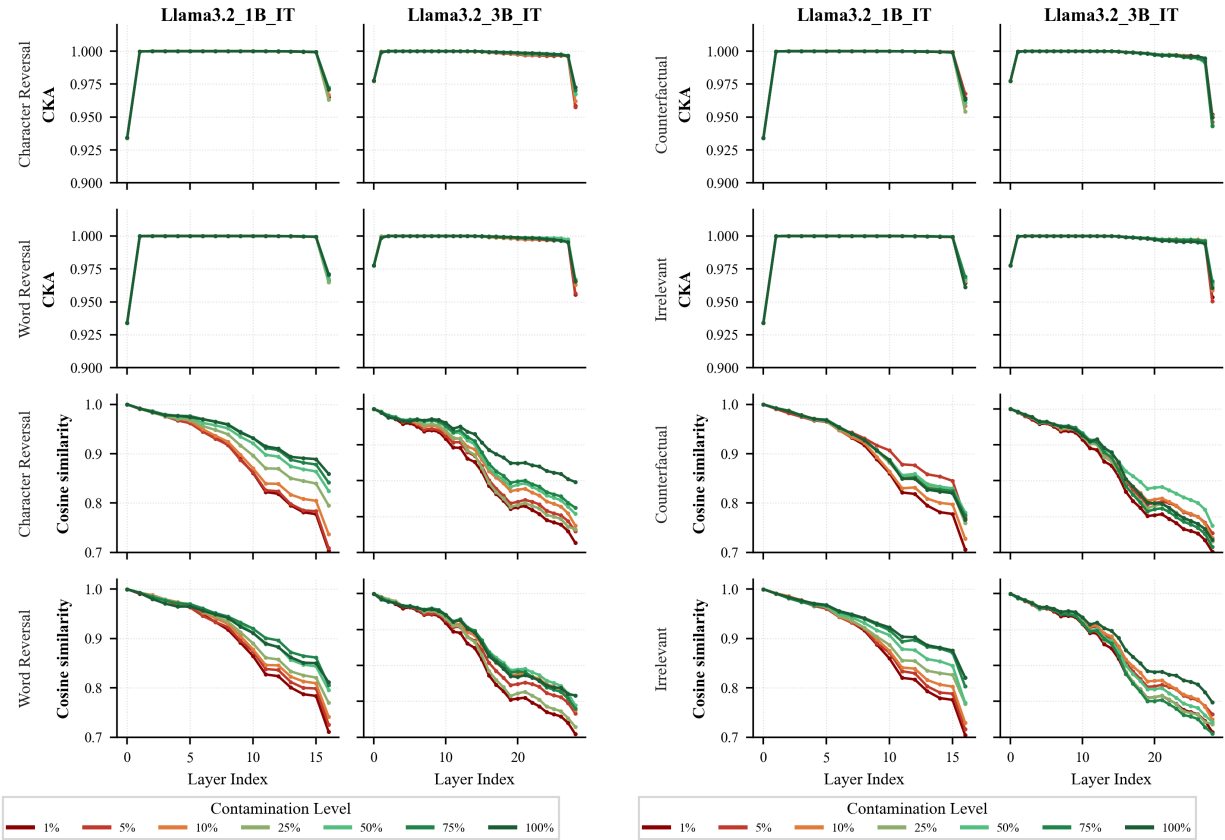


Figure 13: Layerwise CKA and average last-token cosine similarity between baseline and finetuned Llama3.2 base models under syntactic contamination. Curves correspond to contamination levels from 1% to 100% for character reversal and word reversal.

decline and clear separation across contamination levels, but the ordering is not dose–response. In particular, the lowest contamination (1%) often yields the *lowest* cosine similarity throughout much of the network, while several higher contamination levels remain closer to the baseline; separation is most visible toward late layers. This indicates that, in Llama3.2, syntactic corruption is expressed primarily as directional drift in the last-token representation (captured by cosine), while cross-example geometry remains comparatively stable (captured by CKA). The 3B variant shows a longer and more structured depth profile than the 1B variant, with stronger late-layer cosine separation.

Semantic contamination: Under counterfactual and irrelevant supervision (Figure 14b), CKA again stays close to 1.0 across layers for both model sizes, with only minor late-layer reductions. Cosine similarity still decreases smoothly with depth and shows contamination-dependent separation, but the relationship is again non-monotonic: low contamination does not consistently correspond to the weakest similarity, and curves may be non-ordered across levels. Overall, for both semantic transformations, contamination effects are much more apparent in cosine similarity than in CKA, suggesting that behavioral changes can be implemented via shifts in the direction of late hidden states without strong reorganization of representational geometry.

Architectural implications: Llama3.2 follows a decoder-only transformer architecture and the 1B and 3B models are trained by distilling Llama 3.1 models, where logits from the 8B and 70B teachers are incorporated as token-level targets during pretraining, with additional knowledge distillation used after pruning to recover performance (Grattafiori et al., 2024). The 1B and 3B models support a 128K token context window. Relative to architectures that explicitly interleave heterogeneous attention regimes across depth, Llama3.2 does not introduce a local-window versus global-layer schedule, yielding a comparatively



(a) Syntactic contamination

(b) Semantic contamination

Figure 14: Layerwise CKA and average last-token cosine similarity between the baseline and finetuned Llama3.2_IT models across contamination levels (1% to 100%). Panel (a) shows syntactic contamination (character and word reversal) and panel (b) shows semantic contamination (counterfactual and irrelevant).

homogeneous stack. This homogeneity is consistent with the observed representational profiles in Figures 13 and 14: CKA remains near-ceiling across depth with minimal mid-stack irregularities, while cosine similarity exhibits a smooth depth-dependent drift and non-monotonic separation across contamination levels. The preservation of cross-example geometry (high CKA) alongside directional shifts (declining cosine) suggests that contamination can be expressed through adjustments to output-relevant directions in late hidden states without requiring large-scale reorganization of intermediate representational geometry.

G.3 OLMo2

Figure 15-16 present layerwise representational similarity for OLMo2. Across settings, CKA remains near-ceiling throughout the stack with only mild late-layer drift, while cosine similarity shows a gradual depth-dependent decrease with the clearest deviation concentrated at the final layer(s).

Syntactic contamination: For both the base model (Figure 15) and the instruction-tuned model (Figure 16a), character and word reversal leave CKA nearly unchanged across depth, with only a slight late-layer decrease. In contrast, cosine similarity shows a more interpretable depth profile: it declines gradually through the stack and then drops more sharply at the final layer(s). Separation across contamination levels is present but modest; the curves are tightly clustered over most layers and fan out primarily at the end of the network. Where ordering is visible, it is closer to an approximate monotone trend (lower contamination often yields slightly lower cosine), but the effect size is small and not cleanly dose-response at every layer. Character

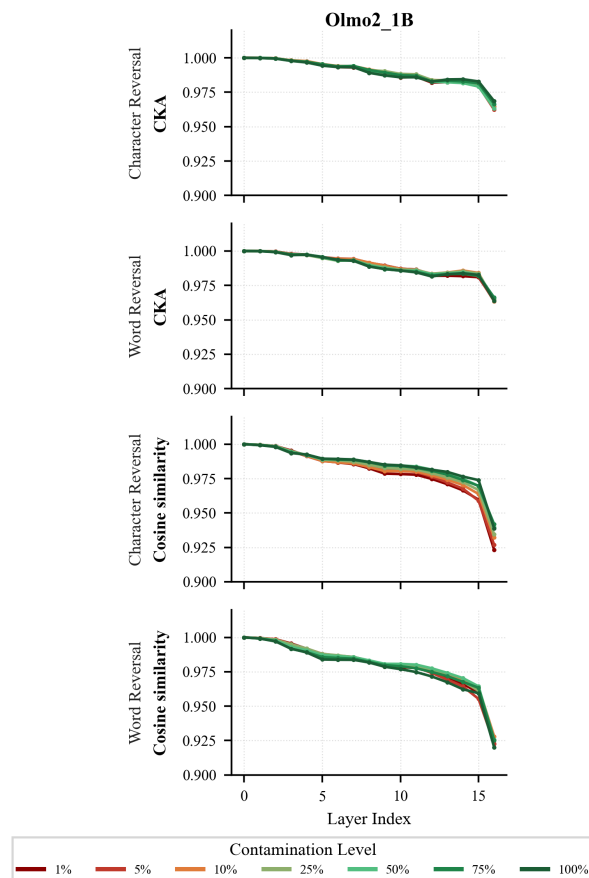
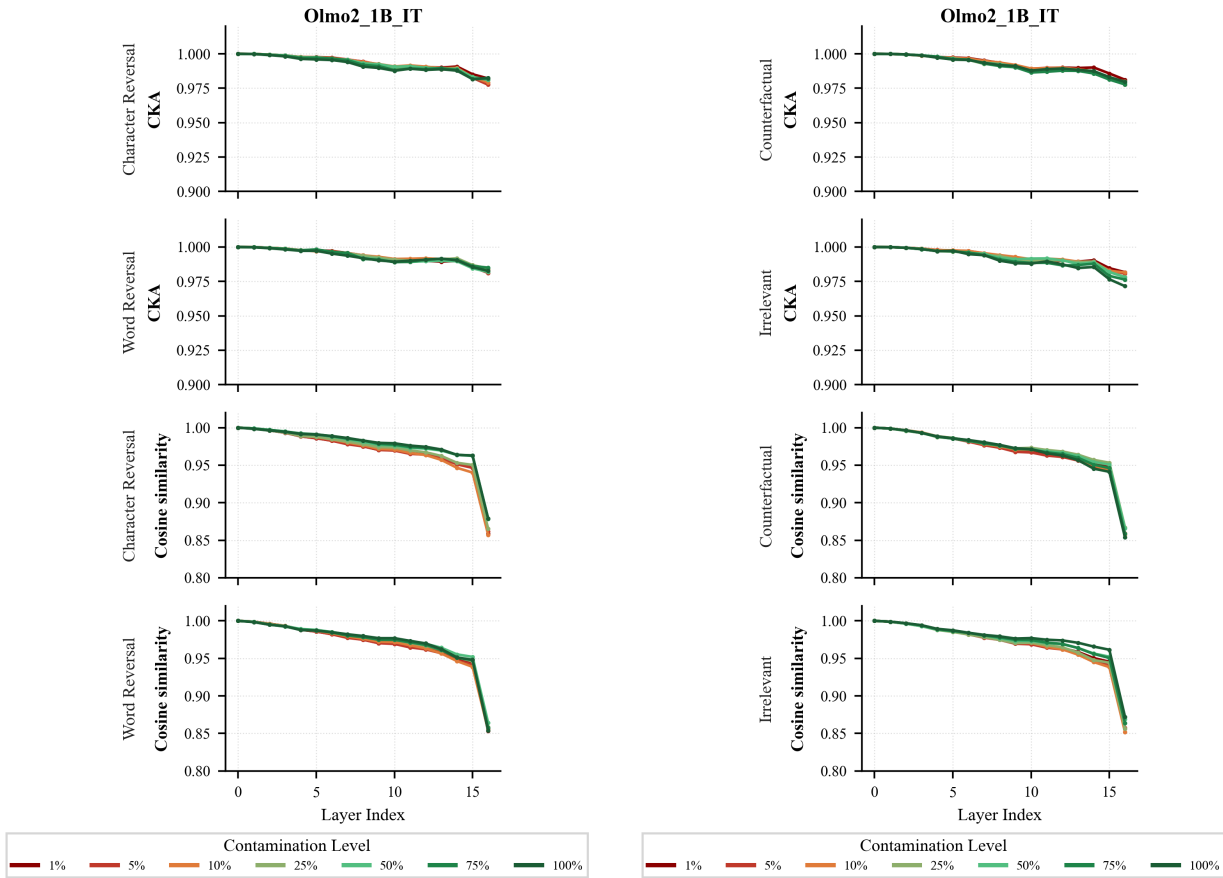


Figure 15: Layerwise CKA and average last-token cosine similarity between baseline and finetuned OLMo2 under syntactic contamination. The curves correspond to contamination levels from 1% to 100% for character reversal and word reversal.

reversal is marginally more disruptive than word reversal near the output end, consistent with stronger perturbations to token-level structure.

Semantic contamination: Under counterfactual and irrelevant supervision (Figure 16b), OLMo2_IT again exhibits near-ceiling CKA across the stack, with only minor late-layer reductions. Cosine similarity displays the same qualitative shape as in the syntactic setting (a gradual depth-wise drift with the largest change at the final layer(s)), but the overall magnitude of the shift is small. Contamination-level separation remains subtle and is concentrated near the end of the network; intermediate-layer trajectories are largely indistinguishable across levels and between counterfactual versus irrelevant supervision.

Architectural implications: OLMo2 follows a decoder-only transformer architecture with several stability-oriented modifications relative to earlier OLMo iterations (Walsh et al., 2025). These include RMSNorm applied after (rather than before) the attention and feedforward sublayers within each block, QK-norm (RMSNorm applied to queries and keys before computing attention to prevent attention logits from growing too large), auxiliary Z-loss regularization, and RoPE-based positional encoding. Collectively, these choices aim to control the scale of activations, attention logits, and output logits during training. The layerwise similarity profiles are consistent with this design emphasis: CKA remains near-ceiling across layers under both syntactic and semantic contamination, indicating limited change in cross-example geometry, while the more visible effects appear as a gradual drift in last-token cosine similarity that is largest at the final layer. This pattern suggests that contamination can be expressed primarily through directional shifts in late-layer representations, without inducing large geometry-level rearrangements across the network.



(a) Syntactic Contamination

(b) Semantic Contamination

Figure 16: Layerwise CKA and average last-token cosine similarity between the baseline and finetuned OLMo2_IT models across contamination levels (1% to 100%). Panel (a) shows syntactic contamination (character and word reversal) and panel (b) shows semantic contamination (counterfactual and irrelevant).

G.4 Phi4

Figure 17 reports layerwise representational similarity for Phi4_Mini_IT. The most distinctive feature of this family is a pronounced early-layer deviation under syntactic contamination, followed by a long plateau and modest late-layer drift.

Syntactic contamination: Under character and word reversal (Figure 17a), Phi4_Mini_IT exhibits a distinctive two-regime structure in CKA: a sharp early-layer drop concentrated in the first few blocks, followed by near-ceiling similarity through the bulk of the stack, and then a terminal-layer reduction. Character reversal produces a more severe early-layer disruption than word reversal, with CKA dropping to approximately 0.65-0.90 compared to 0.82-0.95 for word reversal. The early discontinuity is the main locus of contamination-level sensitivity in CKA; different contamination rates separate most clearly at this point, while intermediate layers collapse to an almost identical trajectory across levels. In contrast, last-token cosine similarity changes more smoothly: it shows a mild depth-wise drift and a clearer decrease near the end of the network, but with relatively small spread across contamination levels compared to the early CKA separation. Overall, syntactic corruption in Phi4_Mini_IT appears to be expressed as an early reconfiguration (captured by CKA) combined with a shift in final-layer representations (captured by cosine), rather than as sustained divergence distributed across intermediate depth, with character-level corruption consistently inducing stronger effects than word-level corruption.

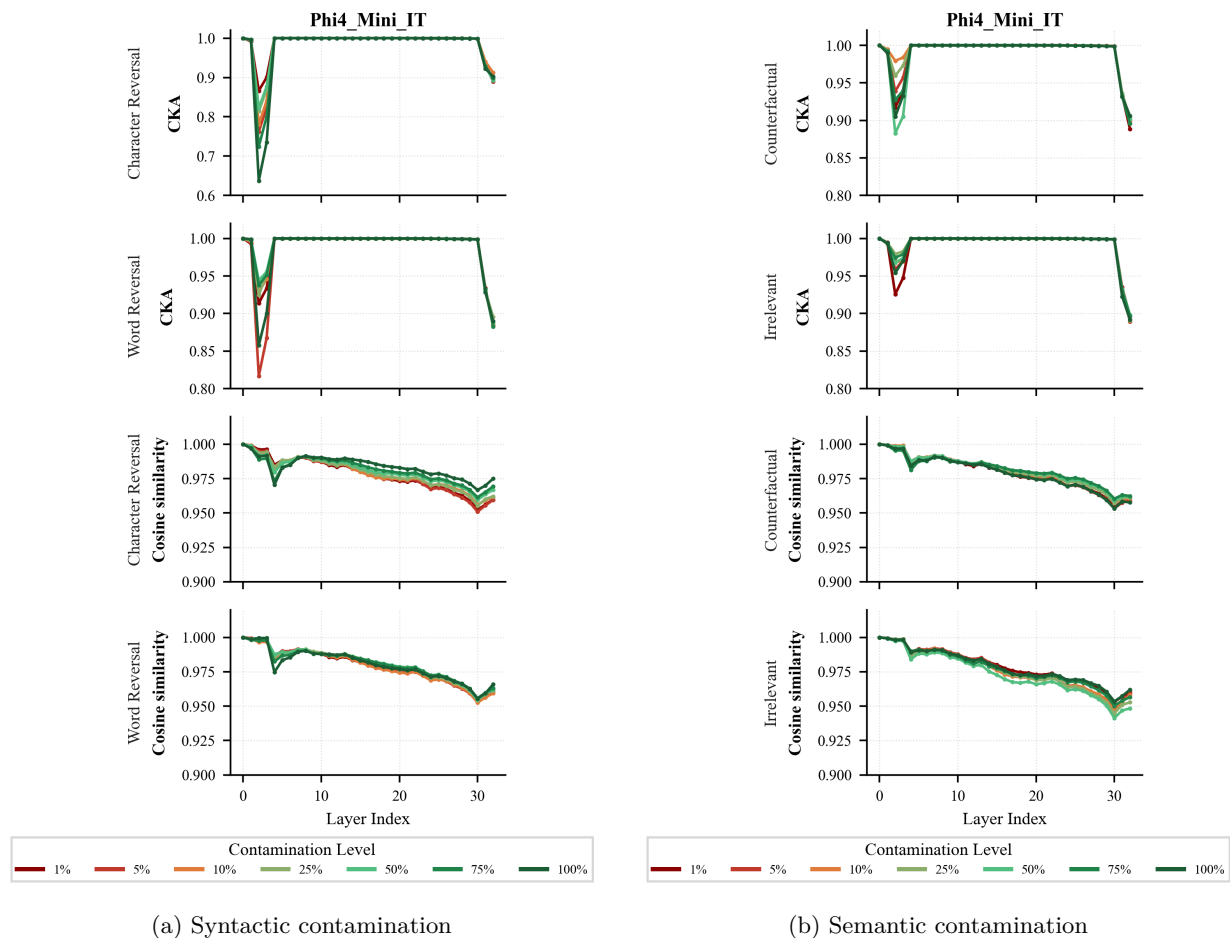


Figure 17: Layerwise CKA and average last-token cosine similarity between baseline and finetuned Phi4_Mini_IT models across contamination levels (1% to 100%). Panel (a) shows syntactic contamination (character and word reversal) and panel (b) shows semantic contamination (counterfactual and irrelevant).

Semantic contamination: Under counterfactual and irrelevant supervision (Figure 17b), Phi4_Mini_IT preserves substantially higher similarity than under syntactic contamination. CKA remains near-ceiling throughout the network, with only minimal early-layer and terminal-layer deviations that are far less severe than those observed under character or word reversal. Cosine similarity shows gradual depth-wise drift with the most visible reduction near the final layers, but overall changes remain modest. Separation across contamination levels is weak: curves are tightly clustered across the entire depth, with little differentiation at any point in the network. Differences between counterfactual and irrelevant supervision are subtle in representation space, with both preserving intermediate-layer trajectories and showing only slight variations in the magnitude of early and late deviations. Semantic transformations induce minimal geometric reorganization and preserve representational structure compared to syntactic contamination.

Architectural implications: Phi4_Mini_IT follows a decoder-only transformer architecture with Grouped-Query Attention (24 query heads and 8 key/value heads to reduce KV-cache memory), tied input/output embeddings, and LongRoPE positional encoding supporting a 128K context length (Abouelenin et al., 2025). The model uses 32 transformer layers with hidden size 3,072 and an expanded vocabulary of 200,064 tokens (Abouelenin et al., 2025). The distinctive early-layer CKA drop under syntactic contamination, followed by a long plateau and a terminal-layer reduction, is consistent with the tied-embedding structure providing strong coupling between the token embedding space and the output prediction pathway. Under this coupling, contamination-induced updates that modify output behavior can also shift the shared

embedding parameters, which can amplify sensitivity near the embedding layer (early) and in the final blocks closest to the output projection (late), even if intermediate-layer geometry remains largely preserved. The intermediate plateau suggests that, after passing through the initial perturbed layers, the middle of the network maintains a stable representational geometry across contamination levels. In contrast, semantic contamination produces comparatively small representational changes across depth, indicating that factual corruption can be expressed with limited geometric reorganization relative to syntactic reversal.

G.5 Qwen2.5

The main text analyzes Qwen2.5_IT in detail. Figure 18 provides the corresponding layerwise view for Qwen2.5 *base* models under syntactic contamination (character and word reversal), enabling comparison between pre-instruction representations and the instruction-tuned behavior discussed in Section 4.

Syntactic contamination: Across sizes, CKA remains near-ceiling through most intermediate layers and departs most visibly at the boundaries of the network: the earliest layers (embedding/first block) and the final blocks. The 0.5B model shows an especially sharp early-layer CKA drop followed by a rapid return to near-ceiling similarity, indicating that early perturbations do not propagate as sustained geometric differences through the middle layers. In contrast, the 3B model exhibits much stronger terminal-layer sensitivity, with pronounced late-layer CKA collapses under both character and word reversal, while cosine similarity in the same region remains substantially higher. This metric gap indicates that syntactic corruption can induce large changes in cross-example geometry near the top of the network even when average directional alignment is partially preserved.

The relationship between contamination level and representational divergence is not a clean dose-response curve. In particular, curves across levels frequently overlap and are not reliably ordered across depth; the most consistent level-dependent separation, when present, is confined to the final blocks of the larger model. Character reversal is generally more disruptive than word reversal in the late layers, but the magnitude of separation varies by depth and does not increase monotonically with the nominal contamination percentage.

The representational signatures of contamination are remarkably similar between base and instruction-tuned Qwen2.5 models. Both exhibit: (1) stable intermediate-layer geometry with $CKA > 0.95$, (2) sharp late-layer drops concentrated in the final blocks, (3) clear contamination-level separation primarily in terminal layers, and (4) stronger effects for character reversal than word reversal. The main difference is quantitative rather than qualitative: instruction-tuned models sometimes show slightly larger absolute drops in late-layer CKA, but the overall pattern of intermediate stability followed by terminal collapse remains consistent. This suggests that the architectural and geometric factors governing contamination localization (where effects appear) are largely independent of instruction tuning, though instruction tuning may modestly amplify the magnitude of late-layer sensitivity.

Architectural implications: Qwen2.5 uses a standard decoder-only Transformer stack with rotary positional encodings and modern normalization and feed-forward choices typical of recent open-weight LMs, and it scales primarily by increasing depth and width across the 0.5B, 1.5B, and 3B variants (?). Unlike architectures that introduce explicit heterogeneity across depth (for example, interleaved local and global attention), Qwen2.5 applies a largely homogeneous block structure throughout the stack. This uniformity is consistent with the qualitative pattern in Figure 18: most intermediate layers remain highly aligned to the baseline under syntactic corruption, while the most pronounced departures concentrate near the top of the network, where hidden states are most directly mapped into vocabulary logits by the final blocks and output projection. Importantly, sharp late-layer CKA collapses are not confined to the 3B model: the 1.5B variant also exhibits a substantial terminal-layer breakdown in CKA, indicating that syntactic contamination can trigger late-stage geometric reorganization already at intermediate scale. In both 1.5B and 3B, these CKA drops coexist with substantially higher cosine similarity, suggesting that late blocks can undergo strong cross-example geometric reshaping without an equally large change in average last-token direction. A plausible interpretation is that contamination fine tuning preferentially perturbs late-stage feature recombination and output-relevant subspaces that control next-token decisions, while leaving much of the reusable feature hierarchy computed by earlier layers comparatively intact. Increasing scale may make such late-block effects

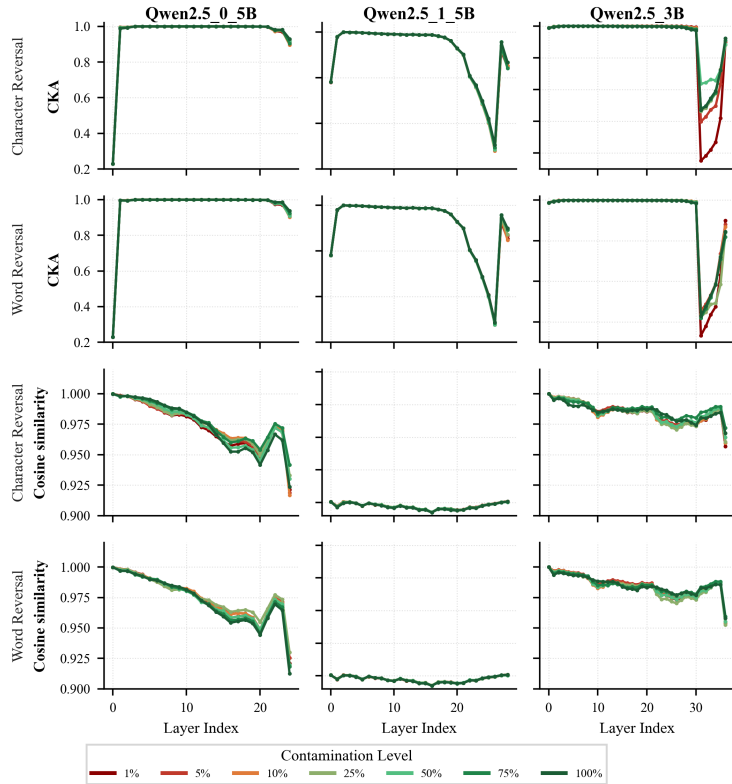


Figure 18: Layerwise centered kernel alignment (CKA) and average last-token cosine similarity between baseline and finetuned Qwen2.5 base models under syntactic contamination. Curves correspond to contamination levels from 1% to 100% for character reversal and word reversal across model sizes (0.5B, 1.5B, 3B).

easier to express or more visible by allocating more capacity to the final-stage transformations that condition the output distribution.

G.6 SmolLM2

Figure 19-20 report layerwise representational similarity for SmolLM2 (360M and 1.7B). SmolLM2 exhibits some of the cleanest representational patterns across all families, with minimal intermediate-layer disruption and effects concentrated primarily in terminal layers.

Syntactic contamination: In the base models (Figure 19), both SmolLM2_360M and SmolLM2_1.7B preserve near-ceiling CKA across most layers, followed by a sharp reduction at the final block(s). The main qualitative difference between sizes is the presence of an additional early-layer CKA dip in the 1.7B variant, most visible for word reversal, whereas the 360M model remains close to 1.0 until the terminal drop. Last-token cosine similarity is also highly stable in both sizes, with only small mid-depth drift and a clearer end-of-network reduction. Separation across contamination levels is limited and not reliably ordered: curves frequently overlap and, when they separate, this is most apparent at the early dip (1.7B) and at the final layer(s) (both sizes), rather than producing a consistent dose-dependent profile across depth.

For instruction-tuned models (Figure 20a), the depth profile remains similar but becomes even more concentrated at the output end. CKA stays extremely high through almost the entire stack, and the dominant deviation is again a sharp terminal-layer drop. The 1.7B_IT model retains a mild early-layer sensitivity for word reversal, but the intermediate layers largely collapse to a single near-ceiling trajectory across contamination levels. Overall, syntactic corruption in SmolLM2 is expressed primarily as localized changes in early layers for the larger base model (early CKA dip) and at the final layers for all variants (terminal-layer reductions), rather than as sustained divergence throughout intermediate depth.

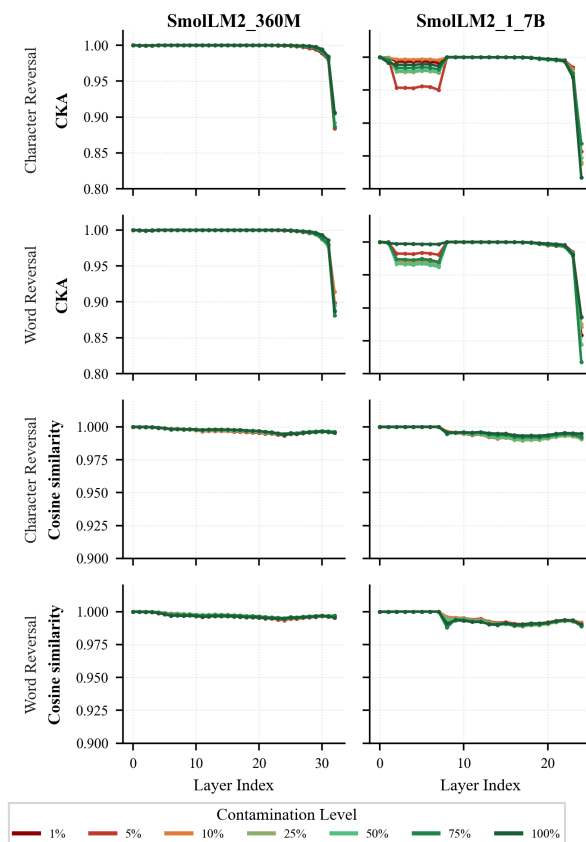
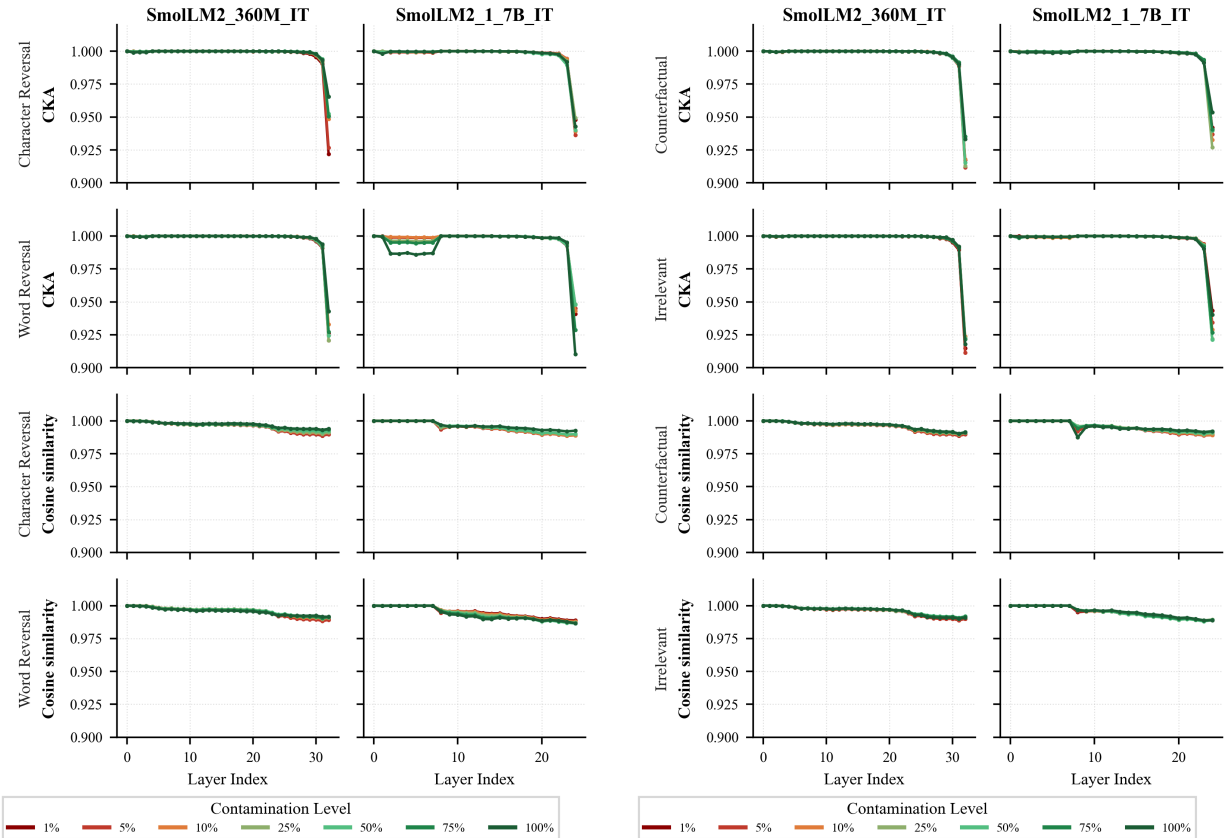


Figure 19: Layerwise CKA and average last-token cosine similarity between baseline and finetuned SmolLM2 base models under syntactic contamination. Curves correspond to contamination levels from 1% to 100% for character reversal and word reversal across 360M and 1.7B model sizes.

Semantic contamination: Under counterfactual and irrelevant supervision (Figure 20b), SmolLM2_IT retains high representational similarity across depth. CKA remains near-ceiling through most layers, with the clearest deviation appearing at the final layer(s), and last-token cosine similarity is similarly stable with only a small end-of-network decrease. As with syntactic contamination, contamination levels do not induce consistently ordered separation: trajectories across levels are tightly clustered and often overlap throughout the stack, indicating that representational divergence is not a simple monotonic function of contamination rate. Differences between counterfactual and irrelevant supervision are comparatively subtle, and when visible they are concentrated near the end of the network.

Architectural implications: SmolLM2 follows a decoder-only transformer architecture in the Llama family, using Grouped-Query Attention, tied input and output embedding weights, SwiGLU activations, RMSNorm, and RoPE positional encodings (Allal et al., 2025). In the 1.7B configuration, the model uses 24 transformer layers with a model dimension of 2,048 (Allal et al., 2025). The layerwise similarity profiles in Figure 19-20 are comparatively clean: CKA is near-ceiling through most intermediate layers, curves across contamination levels are tightly clustered and often non-ordered, and the dominant deviations appear at the final block(s) (with an additional early-layer CKA dip visible for the 1.7B base model under syntactic reversal).

One plausible explanation is that SmolLM2’s training recipe, which emphasizes large-scale curation and staged introduction of higher-skill domains (for example, math and code), yields intermediate representations that are relatively stable under the fine tuning signals induced by contamination (Allal et al., 2025). Under this interpretation, contamination effects are expressed primarily through small changes concentrated in



(a) Syntactic contamination

(b) Semantic contamination

Figure 20: Layerwise CKA and average last-token cosine similarity between baseline and finetuned SmoLLM2_IT models across contamination levels (1% to 100%). Panel (a) shows syntactic contamination (character and word reversal) and panel (b) shows semantic contamination (counterfactual and irrelevant).

output-adjacent computations, rather than through broad reorganization of intermediate representational geometry, consistent with the near-ceiling intermediate-layer CKA and the weak, non-monotonic separation across contamination levels observed in the plots.

G.7 Shared representational patterns across families and scales

Across all evaluated model families and scales, contamination effects manifest most reliably as late-layer representational shifts rather than broad, depth-uniform rewrites. Despite substantial architectural variation, including Gemma3’s local-global attention interleaving, Phi4’s tied embeddings, OLMo2’s post-normalization, and SmoLLM2’s data-centric training, intermediate layers remain comparatively stable under both syntactic and semantic corruption, while the final blocks show the largest and most behaviorally relevant divergence. This supports the view that contamination is implemented primarily through adjustments to output-adjacent computations that most directly control the logits, enabling models to change next-token predictions with minimal disruption to reusable features established during pretraining.

The contrast between metrics reinforces this picture: cosine similarity is often only moderately affected while CKA can exhibit sharper late-layer declines, particularly under syntactic contamination where character reversal can drive substantial terminal-layer CKA reductions even when cosine similarity remains comparatively high. This indicates that fine tuning can reconfigure cross-example geometry in the final representation space even when average token-level directions remain partially aligned. At the same time, contamination strength is not consistently monotonic in contamination rate across layers, and trajectories

across contamination levels frequently overlap, suggesting that representational change is governed more by how supervision reshapes late-layer feature recombination than by a smooth scaling with contamination rate.

Within-family scaling generally makes late-layer divergence easier to observe, consistent with larger models allocating more capacity to output-adjacent transformations and thereby expressing systematic supervision artifacts through richer modifications of late-layer feature combinations. While architectural differences produce family-specific signatures, such as Phi4’s early-layer sensitivity or Gemma3’s mid-network fluctuations, these variations do not alter the fundamental localization of contamination effects in the final blocks. Taken together, these results point to late-layer specialization as a common mechanistic locus for contamination across decoder-only transformers, and motivate focusing mitigation and diagnostics on the final blocks and output projection where supervision-induced artifacts most directly enter the generation pathway.

H Qwen2.5 ablation: input and ‘input and output’ transformations

To study how transformation placement affects fine tuning dynamics, we evaluate two additional settings beyond the output-only transformations used in the main experiments. Output-only transformations isolate whether corrupted supervision is internalized as a generation behavior even when evaluation prompts remain clean. Here, we run focused ablations on the Qwen2.5 instruction-tuned models, applying syntactic transformations on the input side either with clean outputs or with the same transformation applied to both input and output. These settings separate robustness to transformed inputs from learning transformed outputs, and allow comparison to the output-only case under an otherwise identical training and evaluation protocol. We restrict this analysis to syntactic transformations because applying the semantic transformations on the input side would not introduce a qualitatively new corruption mode, since the semantic mismatch studied in the output-only setting would largely relocate to the input rather than define an orthogonal axis of variation.

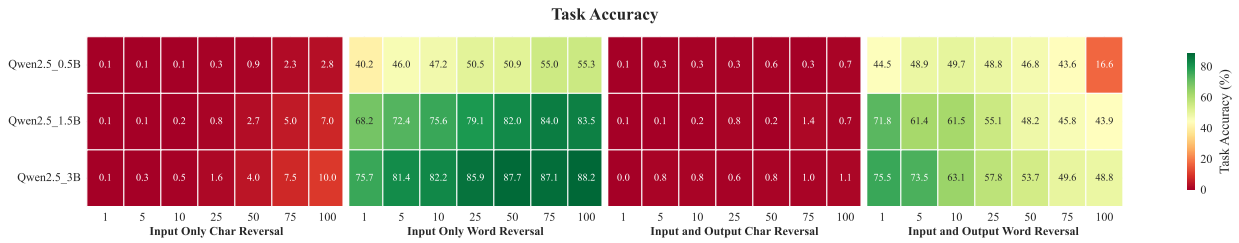


Figure 21: Model performance under input-side and ‘input and output’ syntactic transformations across contamination levels. The figure reports accuracy for input-only character reversal, input-only word reversal, ‘input and output’ character reversal, and ‘input and output’ word reversal.

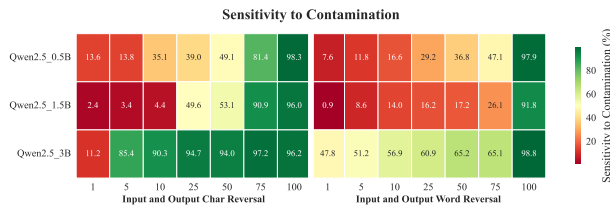


Figure 22: Sensitivity to syntactic transformation in the ‘input and output’ setting across contamination levels. The figure reports sensitivity for ‘input and output’ character reversal and ‘input and output’ word reversal, where the same syntactic transformation is applied consistently to both input and output.

Input-only transformations behave qualitatively differently from the output-only setting in the main paper. When only the input is transformed and the output targets remain clean, fine tuning learns to map a transformed input back to a clean answer. Under input-only character reversal, this mapping largely fails as an instruction-following problem: accuracy remains near zero across scales and contamination levels (0.1% to

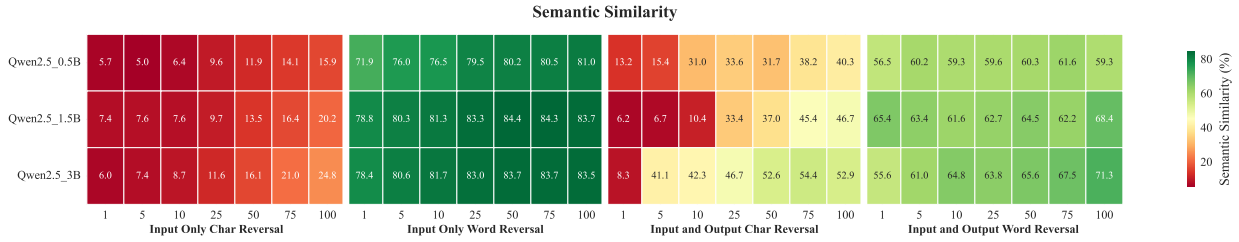


Figure 23: Semantic similarity under input-side and ‘input and output’ syntactic transformations across contamination levels. The figure reports semantic similarity for input-only character reversal, input-only word reversal, ‘input and output’ character reversal, and ‘input and output’ word reversal.

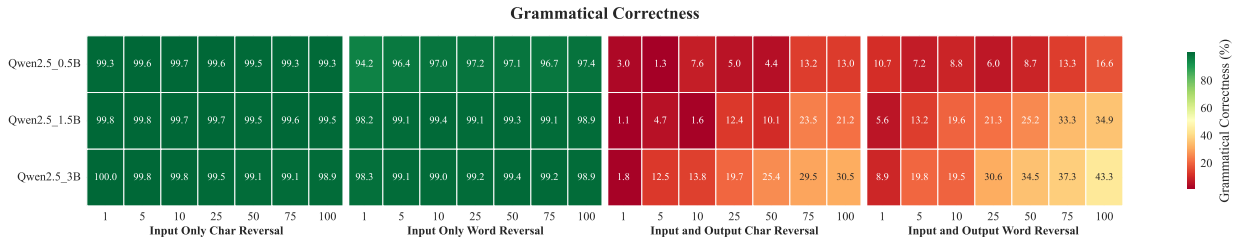


Figure 24: Grammatical correctness under input-side and ‘input and output’ syntactic transformations across contamination levels. The figure reports grammatical correctness for input-only character reversal, input-only word reversal, ‘input and output’ character reversal, and ‘input and output’ word reversal.

10.0%), and semantic similarity stays extremely low (5.0% to 24.8%), while grammatical correctness remains near-ceiling (Figure 21, Figure 23, and Figure 24). In contrast, input-only word reversal is substantially milder and becomes easier with both scale and contamination level. For the 0.5B model, accuracy increases from 40.2% to 55.3%; for the 1.5B model, accuracy increases from 68.2% to 83.5%; and for the 3B model, accuracy increases from 75.7% to 88.2%. Semantic similarity remains high at roughly 71.9% to 84.4%, and grammatical correctness stays high at 94.2% to 99.4% (Figure 21, Figure 23, and Figure 24). These trends contrast with output-only character reversal, where transforming the targets directly disrupts well-formedness and meaning preservation because the model is trained to emit transformed text in its generation pathway.

When the same syntactic transformation is applied consistently to both input and output, training no longer asks the model to recover clean outputs from transformed inputs and instead reinforces a self-consistent transformed mapping. Joint character reversal remains highly destructive in utility terms: accuracy stays near zero between 0.0% and 1.4%, and semantic similarity remains low to moderate relative to clean behavior at roughly 6.2% to 54.4% depending on scale and contamination level, while grammatical correctness collapses sharply to about 1.1% to 30.5% (Figure 21, Figure 23, and Figure 24). In this input-and-output setting, sensitivity to the transformation increases steeply as the contamination level increases. For ‘input and output’ character reversal, sensitivity in the 0.5B model rises from 13.6% to 98.3%, sensitivity in the 1.5B model rises from 2.4% to 96.0%, and sensitivity in the 3B model rises from 11.2% to 96.2%, indicating increasing adoption of the transformed mapping when it is consistently present on both sides (Figure 22). For ‘input and output’ word reversal, semantic similarity is lower than in the input-only word reversal case and sensitivity again increases with contamination level. Sensitivity for ‘input and output’ word reversal rises from 7.6% to 97.9% for the 0.5B model, from 0.9% to 91.8% for the 1.5B model, and from 47.8% to 98.8% for the 3B model (Figure 22 and Figure 23).

We report sensitivity only for the ‘input and output’ setting because sensitivity is defined as adherence to the output-side transformation pattern. In the input-only setting, outputs are clean by construction, so there is no output-side transformation adoption to measure and sensitivity is not applicable.

Overall, transformation placement changes the dominant failure mode. Input-only transformations primarily stress robustness to transformed inputs and tend to preserve fluency even when utility collapses, whereas applying the transformation to both input and output reinforces transformed outputs and yields increasing transformation adoption as the contamination level increases, resembling the output-only learning dynamics in the main setting.

I Human and LLM-as-a-Judge agreement analysis in automated evaluation

To validate the reliability of our automated evaluation system, we conducted a comprehensive human-model agreement analysis comparing human evaluator judgments with Gemini 2.0 Flash assessments. Human evaluators assessed 100 randomly selected question-response pairs for each model-data combination, covering all transformation types and evaluation criteria. The random sampling ensured representative coverage across different contamination levels and model behaviors.

Table 6: Human and LLM-as-a-judge agreement analysis: percentage agreement and Cohen’s Kappa between human evaluator and Gemini 2.0 Flash

Criterion	Character Reversal		Word Reversal		Counterfactual		Irrelevant	
	% Agree	κ	% Agree	κ	% Agree	κ	% Agree	κ
Pattern Adherence	100.0	1.00	99.83	0.91	99.76	0.73	100.0	1.00
Accuracy	97.92	0.72	95.20	0.74	97.84	0.85	99.16	0.85
Grammatical Correctness	95.76	0.67	96.53	0.73	98.92	0.81	99.66	0.76

The agreement analysis examined three key evaluation dimensions: pattern adherence (whether responses correctly follow the intended transformation pattern), accuracy (factual correctness of responses compared to reference answers), and grammatical correctness (structural linguistic coherence). Both human evaluators and Gemini 2.0 Flash used identical evaluation criteria and structured assessment protocols to ensure fair comparison. It is important to note that these evaluation tasks are relatively straightforward and simple, involving clear binary or categorical judgments that do not require complex reasoning or subjective interpretation.

Table 6 presents the percentage agreement and Cohen’s Kappa coefficients across all transformation types and evaluation criteria. The results demonstrate excellent alignment between human evaluators and the automated system, with percentage agreements ranging from 95.20% to 100.0% and Cohen’s Kappa values spanning 0.67 to 1.00, indicating substantial to perfect inter-rater reliability. Pattern adherence shows the highest agreement levels, achieving perfect agreement (100%, $\kappa = 1.00$) for both character reversal and irrelevant transformations, and near-perfect agreement for word reversal (99.83%, $\kappa = 0.91$) and counterfactual transformations (99.76%, $\kappa = 0.73$). Accuracy and grammatical correctness assessments also demonstrate strong agreement, with all values exceeding 95% agreement and Cohen’s Kappa coefficients above 0.67, indicating substantial reliability. The consistently high agreement levels across all criteria and transformation types confirm the simple nature of these evaluation tasks and support the validity of our automated evaluation methodology.