

Beyond Clean and Contaminated: A Survey on the Fundamental Properties of Data Contamination in Large Language Models

Anonymous ACL submission

Abstract

Benchmark-based evaluation remains the primary mechanism for comparing large language models (LLMs), yet modern development pipelines increasingly blur the boundary between training and testing. Beyond direct train–test overlap, contamination leaks through pathways such as post-training, evaluation-time “test-set fitting,” and retrieval-enabled tool use. In this paper, we frame data contamination as an evaluation-validity failure mode and propose a three-dimensional taxonomy based on *phase*, *granularity*, and *modality*. We argue that contamination is regime-dependent rather than binary, summarizing key properties such as inevitability under web-scale collection, scaling effects, and forgettability. Building on these insights, we reorganize detection methods into two complementary paradigms: *statistical* approaches (quantifying inflation via observational signals) and *causal* approaches (verifying via controlled injection). Finally, we provide a critical discussion of these detection methodologies.

1 Introduction

Recent breakthroughs in Large Language Models (LLMs) have demonstrated remarkable capabilities in text generation, code synthesis, and mathematical reasoning (Grattafiori et al., 2024; OpenAI et al., 2024; DeepSeek-AI et al., 2025). Benchmark-based evaluation is central to comparing LLMs, yet modern development pipelines blur the boundary between training and evaluation (Balloccu et al., 2024): prompt libraries (Han et al., 2025), retrieval indices and tool configurations (Han et al., 2025) create additional pathways for benchmark information to sneak into models’ knowledge inadvertently. Consequently, the reliability of LLM evaluation is increasingly questioned due to data contamination—the unintended overlap between training and test datasets (Chang et al., 2024; Sainz et al., 2023).

Crucially, the prevailing binary perspective—classifying models as simply “clean” or “contaminated” based on surface-level pre-training overlap—fails to capture the nuance of modern development. Leakage in current pipelines is diverse: it can occur implicitly through semantic equivalents in instruction tuning (Dekoninck et al., 2024b), propagate via preference pairs during alignment (Tao et al., 2025), or emerge dynamically through retrieval mechanisms. Furthermore, the magnitude of model performance inflation is explicitly modulated by both the timing of contamination (Kocyigit et al., 2025) and the nature of the downstream task (Ishikawa, 2025). This leads to a ‘consistency crisis’, exposing the inconsistency and fragility of current detection paradigms (Dekoninck et al., 2024b).

Despite this evolving landscape, existing surveys predominantly focus on compiling detection techniques and mitigation methods, largely overlooking the phase propagation of risks and failing to systematize the intrinsic characteristics (Deng et al., 2024; Fu et al., 2024). Distinguishing our work from prior surveys (Appendix A), we expand the definition of contamination to a cross-phase scope—where $\mathcal{D}_{\text{train}}$ includes alignment, synthetic variants, and retrieval assets—formalizing it as any metric-inflating overlap as defined in Sec 2. We advance the field by: (1) establishing a unified *three-dimensional taxonomy* (phase, granularity, and modality) that explicitly covers under-explored risks in preference learning and retrieval pipelines (Tao et al., 2025; Han et al., 2025); (2) characterizing intrinsic properties such as forgettability and scaling effects; and (3) reorganizing detection into complementary *statistical* (observational) and *causal* (interventional) paradigms to better navigate the current consistency crisis.

Section 3 characterizes contamination along three axes: *phase*, *granularity*, and *modality*. Regarding *phase*, we highlight risks beyond pre-

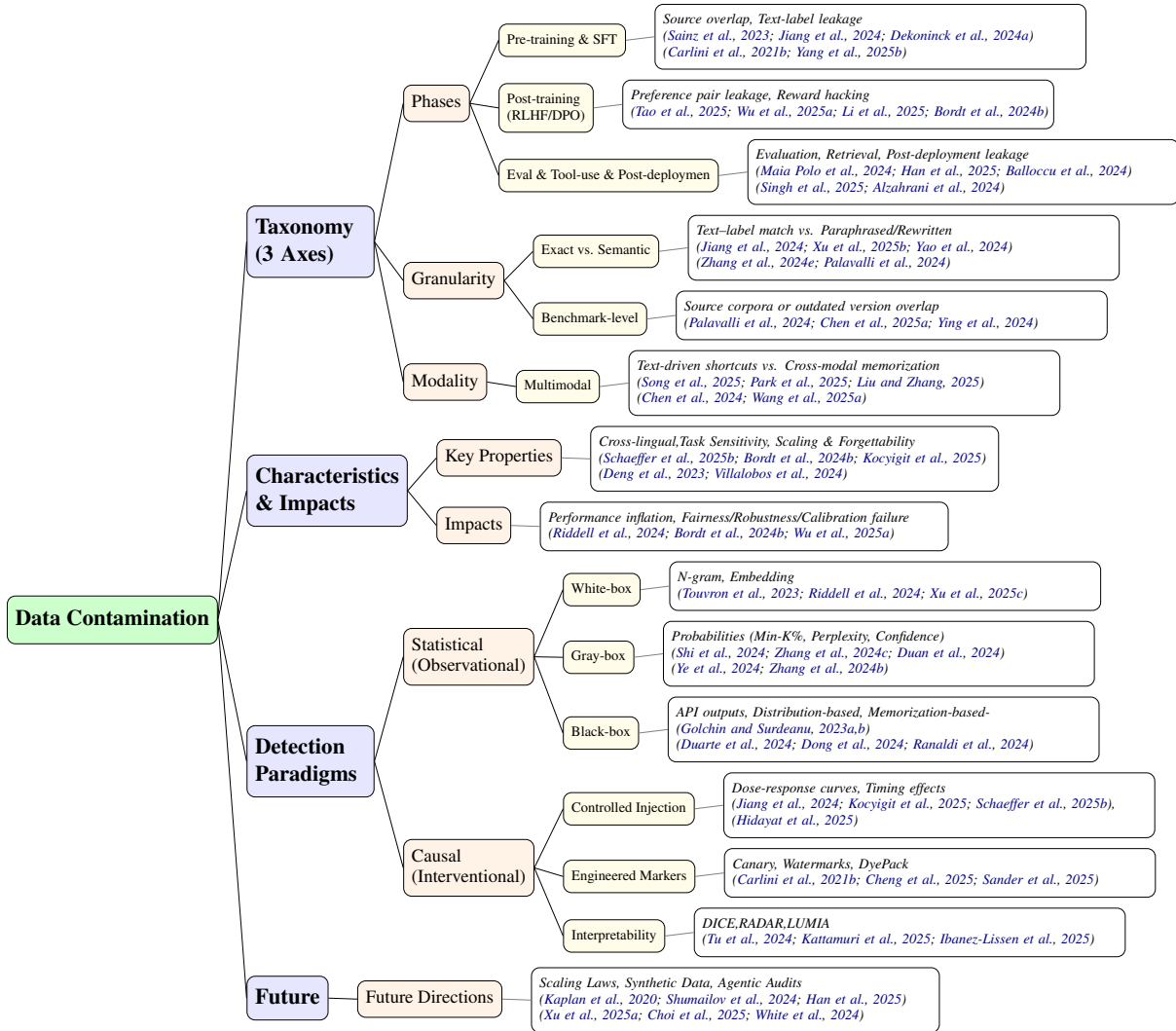


Figure 1: Structure of this paper

083 training, covering RL post-training (Tao et al.,
 084 2025), evaluation-time fitting (Maia Polo et al.,
 085 2024), and retrieval-time leakage (Han et al., 2025).
 086 Different from prior works in table 2, we sys-
 087 tematize *granularity* from instance-level surface
 088 or semantic overlap (Xu et al., 2025b) to broad
 089 benchmark-level inclusion (Palavalli et al., 2024).
 090 Finally, for *modality*, we distinguish between
 091 text-driven shortcuts and cross-modal memoriza-
 092 tion (Song et al., 2025).

093 We synthesize the impacts of data contamina-
 094 tion in section 4. While performance inflation is
 095 the most visible risk—modulated significantly by
 096 leakage granularity and phase—contamination fun-
 097 damentally undermines evaluation integrity. This
 098 inflation masks critical flaws in fairness and ro-
 099 bustness, degrading calibration and compromising
 100 the scientific validity of conclusions (Bordt et al.,
 101 2024b).

102 In section 5, we argue contamination is not
 103 a binary condition but a regime-dependent phe-
 104 nomenon shaped by multiple interacting factors.
 105 It distills recurring characteristics including (i) *in-*
 106 *evitability* under web-scale data collection (Deng
 107 et al., 2023), (ii) *decoding* dependence (Schaeffer
 108 et al., 2025b), (iii) *cross-phase dependence*
 109 on when leakage occurs (Kocyyigit et al., 2025),
 110 (iv) *task and format sensitivity* (Kocyyigit et al.,
 111 2025), (v) *cross-lingual* effects via translated
 112 benchmarks (Yao et al., 2024), and (vi) *scaling and*
 113 *forgettability* (Bordt et al., 2024b), where timing
 114 and exposure magnitude govern whether memo-
 115 rized traces persist or are diluted by subsequent
 116 training (Bordt et al., 2024b).

117 Section 6 reorganizes detection into two
 118 paradigms: *statistical* methods, which infer leak-
 119 age from observational signals and increasingly
 120 aim to *quantify* inflation (Xu et al., 2025a); and

causal methods, which use controlled injection to verify dose–response effects (Jiang et al., 2024). We conclude by critically assessing their robustness and applicability under varying constraints.

Finally, to move beyond reactive detection toward a predictive science, we highlight three future directions in Section 7: quantitative modeling, synthetic data saturation, and generative evaluation. Due to space constraints, we provide a review of mitigation strategies in appendix E.

2 Preliminaries

Data Contamination Let $\mathcal{D}_{\text{train}}$ encompass all development assets: pre-training and alignment corpora (including paraphrased or synthetic variants), as well as auxiliary components like retrieval stores and prompts. Let $\mathcal{D}_{\text{test}}$ denote the benchmark instances. We define *data contamination* as the presence of overlap between the test set $\mathcal{D}_{\text{test}}$ and $\mathcal{D}_{\text{train}}$ ($\mathcal{D}_{\text{train}} \cap \mathcal{D}_{\text{test}} \neq \emptyset$). Such leakage artificially inflates metrics through memorization rather than reflecting genuine generalization.

3 A Taxonomy of Data Contamination

Existing studies on data contamination remain fragmented, lacking a systematic characterization of its origins and forms. To address this, we introduce a taxonomy along three axes—*lifecycle phase*, *granularity*, and *modality*—providing a structured framework to guide detection and mitigation.

3.1 Contamination Phases

Pre-training Phase Contamination Contamination during pre-training manifests in two phases. First, *initial pre-training* leakage occurs when benchmark instances permeate large-scale web crawls, causing memorization prior to task adaptation (Sainz et al., 2023). Second, risks persist during *continued pre-training*: integrating newer or domain-specific corpora can unintentionally assimilate evaluation sets, particularly when data sources overlap with benchmark construction (Ke et al., 2023; Chen et al., 2025a). Crucially, this phase can re-contaminate the model, even after rigorous initial sanitization.

SFT and Instruction Tuning Contamination Instruction datasets may inadvertently include benchmark questions or their templated variants, often via the reformatting of overlapping QA sources (Dekoninck et al., 2024b). Such leakage promotes format-specific memorization, lead-

ing to significant score inflation. Moreover, advanced augmentation like Chain-of-Thought (CoT) can obfuscate these traces, complicating detection (Samuel et al., 2024).

Preference Learning Contamination Benchmark leakage occurs via annotator exposure to test items or the recycling of benchmark responses as preference pairs (Tao et al., 2025; Wu et al., 2025a) (detailed in appendix B.1). This directly amplifies reward-driven overfitting, solidifying the model’s reliance on memorized patterns. RL-phase contamination poses a greater detection challenge than leakage in Supervised Fine-Tuning (SFT) because the optimization objective shifts from matching the data distribution via $\max_{\theta} \mathbb{E}[\log \pi_{\theta}(y|x)]$ to maximizing the expected reward $J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}}[R(\tau)]$. Consequently, contaminated examples may not exhibit conspicuous anomalies in perplexity (PPL) or likelihood scores, yet they can covertly induce the policy π_{θ} to converge toward rigid, reward-hacking behaviors during inference.

Evaluation-pipeline Contamination Contamination can occur *without* changing training data, e.g., (i) repeatedly tuning prompts on the test set (Maia Polo et al., 2024), (ii) selecting few-shot examples that overlap with the benchmark (Singh et al., 2025), or (iii) iterating on public leaderboards for model selection (Alzahrani et al., 2024).

Retrieval-time (Tool-use) Contamination For tool-augmented or search-enabled systems, the model may retrieve solutions during evaluation (detailed in appendix B.2), yielding high scores without genuine generalization (Han et al., 2025).

Post-deployment Contamination After deployment, human interactions and logging pipelines may inadvertently expose benchmark content, which can later be collected into training data (directly or indirectly), creating a data flywheel that leaks evaluation items back into future model versions (Balloccu et al., 2024). This phase is especially relevant for continual learning and production data reuse.

3.2 Contamination Granularity

Text Contamination (Jiang et al., 2024; Li et al., 2024e) Text contamination occurs when the input text components of evaluation samples appear in the pre-training corpus, creating an overlap between test and training data that may artificially inflate model performance metrics.

218	Text-label Contamination (Jiang et al., 2024; Li et al., 2024e)	3.3 Contamination Modality	268
219	Text-label contamination refers to	Contamination in multimodal models typically	269
220	scenarios where the pre-training corpus contains	arises through two major pathways.	270
221	not only the input text but also the corresponding		
222	prompts (task instructions) and labels from evaluation	Text-driven Overlap. When text–label pairs (or	271
223	samples. This effectively exposes the model	even text-only inputs) overlap with the training	272
224	to both questions and correct answers prior to testing.	corpus, models may exploit textual shortcuts that	273
225	Specifically, this phenomenon can be categorized	effectively reconstruct test instances. This is consistent	274
226	into two types: (1) unpaired contamination ,	with evidence that some LVLM benchmarks	275
227	where both the input and output appear in the corpus	remain partially solvable even when visual inputs	276
228	but are not collocated or linked directly; and (2)	are removed (Chen et al., 2024).	277
229	paired contamination , where input-output pairs		
230	appear explicitly together. The latter is considered	Cross-modal Memorization. When image–text	278
231	more severe (Yang et al., 2025b). We provide a	pairs (or image–text–label triplets) are present during	279
232	comparative analysis of the impacts of Text Contamination	training, models can memorize cross-modal	280
233	and Text-label Contamination in Appendix B.3.	correspondences instead of learning generalizable	281
234		vision–language representations (Song et al., 2025).	282
235	Semantic Contamination (Xu et al., 2025b)	For reasoning MLLMs, Liu and Zhang (2025) held	283
236	Semantic contamination refers to contamination that	the visual input fixed while probing a family of	284
237	extends beyond exact matches of benchmark instances,	related tasks, helping expose contamination-driven	285
238	arising instead from semantic-preserving	overfitting that can inflate static benchmarks.	286
239	transformations. Common mechanisms include		
240	input/output masking, noise injection, and adversarial	4 Impacts	287
241	or distractor-based augmentation. Unlike	In this part, we synthesize the literature from two	288
242	overlap-centric definitions, this perspective	complementary angles: (i) the <i>impacts</i> of	289
243	encompasses a broader spectrum of contamination	contamination on reported performance and	290
244	capable of evading surface-level detection while	reliability-related properties, (ii) <i>empirical evidence</i>	291
245	artificially inflating benchmark performance (Yao et al.,	used to substantiate leakage in practice.	292
246	2024; Palavalli et al., 2024). Building on this, Xu		
247	et al. (2025b) distinguished between entity	4.1 Impacts	293
248	contamination (involving named entities and their	Data contamination critically undermines	294
249	implicit biases) and factual contamination (covering	evaluation integrity. (Riddell et al., 2024) revealed	295
250	concrete events, context, and detailed assertions).	that contamination not only artificially inflates	296
251	Information-Level Contamination (Xu et al.,	performance but also invalidates scientific	297
252	2025a) This category captures the leakage of	conclusions in NLP. Moreover, these detrimental	298
253	benchmark-related meta-information—such as	effects transcended superficial metrics like	299
254	label distributions, temporal metadata, or external	accuracy, degrading core model properties	300
255	reviews—rather than exact test instances. Such	and ultimately compromising safety (Bordt et al.,	301
256	exposure can implicitly shape model priors,	2024b).	
257	introducing statistical biases that inadvertently	Performance Inflation While performance	302
258	skew the evaluation process by aligning the	inflation represents the most critical risk	303
259	model with the test set’s structural characteristics.	of data contamination, it is not an inevitable	304
260	And we conduct an impact analysis of information-	outcome; its magnitude is modulated by	305
261	level contamination in appendix B.4.	factors such as the phase, granularity,	306
262	Benchmark-level Contamination (Palavalli et al.,	and task. Supporting this nuance, Yang	307
263	2024) Benchmark-level contamination occurs	et al. (2025b) found that LLMs exhibit	308
264	when a model is trained on the source corpora	substantial performance inflation under	309
265	underlying a benchmark (or on outdated	<i>text-label</i> contamination, while text,	310
266	versions of the benchmark), such that instances	label-only, and unpaired contamination	311
267	or their variants inadvertently leak into the	have largely negligible effects for	
	training data.	coding tasks.	
		Fairness Contaminated data often	312
		encode inherent biases, such as the	313
		over-representation of specific	314
		demographics or viewpoints. During	
		training,	

315	models may disproportionately amplify these bi-		
316	ases and propagate them into their generated out-		
317	puts, leading to a degradation in fairness (Dodge		
318	et al., 2021).		
319	Robustness Contamination creates an illusion of		
320	competence, driving high scores through memo-		
321	rization rather than genuine generalization. This		
322	masks inherent brittleness , where models fail pre-		
323	cipitously on slight input perturbations or out-of-		
324	distribution (OOD) data, rendering their deployed		
325	robustness severely compromised.		
326	Calibration Contamination severely distorts con-		
327	fidence calibration, causing models to assign ex-		
328	cessive confidence to contaminated instances. This		
329	spurious certainty compromises downstream reli-		
330	ability: it can blind uncertainty-based hallucination		
331	detectors to confident errors (Manakul et al., 2023),		
332	and destabilize Reinforcement Learning (RL) align-		
333	ment (Casper et al., 2023), which relies on accurate		
334	probability estimates for effective reward optimiza-		
335	tion and self-correction.		
336	4.2 Empirical Evidence of Contamination		
337	Empirical evidence of contamination is primar-		
338	ily established by identifying unexpectedly strong		
339	benchmark performance and attributing it to prior		
340	training exposure. Early surveys noted that		
341	while evidence was emerging, it remained frag-		
342	mented (Sainz et al., 2023). By 2024, studies be-		
343	gan to employ more granular detection methods,		
344	such as Riddell et al. (2024), who quantified the		
345	specific overlap between benchmarks and training		
346	corpora for code models. Most recently, the scope		
347	of evidence has expanded to include diverse sig-		
348	nals. For instance, Sendyka et al. (2025) demon-		
349	strated anomalous robustness, where models could		
350	solve code-generation tasks even under extreme		
351	prompt redaction, indicating verbatim memoriza-		
352	tion. Concurrently, Wu et al. (2025a) investigated		
353	the Qwen2.5 model family within reinforcement		
354	learning contexts, revealing potential leakage is-		
355	ssues in mathematical benchmarks. These findings		
356	collectively suggest that contamination has perme-		
357	ated specialized domains beyond general text, in-		
358	cluding code and mathematics. Moreover, this evi-		
359	dence reveals a systemic failure manifesting as <i>ab-</i>		
360	<i>normal robustness</i> , where models exhibit capabili-		
361	ties even under severely degraded conditions. Ad-		
362	ditional evidence will be discussed in Appendix C.		
	5 Characteristics of Data Contamination		363
	In this section, we summarize 6 characteristics of		364
	data contamination in LLMs.		365
	Inevitability As LLMs continue to scale up, the		366
	size of their training datasets expanded correspond-		367
	ingly (Villalobos et al., 2024). These datasets are		368
	often sourced from extensive web crawls, which		369
	may inadvertently overlap with evaluation bench-		370
	marks, leading to data contamination (Deng et al.,		371
	2023). This process is currently inevitable.		372
	Injection Phase Sensitivity Data contamination		373
	can manifest throughout the model lifecycle (Bal-		374
	loccu et al., 2024). Kocyigit et al. (2025) found that		375
	the training phase at which contamination occurs		376
	plays a crucial role in its impact. Early contamina-		377
	tion leads to a sharp initial performance increase,		378
	but this effect gradually diminishes as training		379
	progresses. Late-phase contamination ultimately		380
	causes a larger performance inflation. Uniform		381
	contamination (spread across the entire training		382
	process) produces the most lasting effects, with no		383
	significant spikes.		384
	Decoding Sensitivity Manifestations of contam-		385
	ination are also modulated by inference param-		386
	eters. Schaeffer et al. (2025b) demonstrated that		387
	higher sampling temperatures can dampen these		388
	contamination-driven score boosts, whereas greedy		389
	decoding tends to amplify the retrieval of memo-		390
	rized patterns. However, this relationship between		391
	temperature and extraction is not always mono-		392
	tonic (Hayes et al., 2025).		393
	Task Sensitivity The impact of data contamina-		394
	tion is fundamentally task-dependent, rather than a		395
	uniform inflation of performance scores (Golchin		396
	and Surdeanu, 2023b; Jiang et al., 2024). For rea-		397
	soning tasks, contamination manifested as a dis-		398
	crepancy between direct answers and logical con-		399
	sistency (Ishikawa, 2025), whereas in generative		400
	tasks like summarization, the gains were sensitive		401
	to whether the leakage was verbatim or reformat-		402
	ted (Palavalli et al., 2024). Moreover, this sensi-		403
	tivity was mediated by data resource characteris-		404
	tics; for instance, contamination yields dispropor-		405
	tionately larger performance gains in low-resource		406
	settings—such as rare language pairs in machine		407
	translation—compared to high-resource ones (Ko-		408
	cocyigit et al., 2025). Thus, contamination is not		409
	monolithic; it demands task-specific evaluation.		410

411 **Cross-lingual Characteristics** LLMs are over-
 412 fitted to translated versions of benchmark test
 413 sets in non-English languages. This practice in-
 414 flated model performance on the original English
 415 benchmarks without direct exposure to them (Yao
 416 et al., 2024), while evading existing detection meth-
 417 ods (Zhang et al., 2024a). Building on this Kocyigit
 418 et al. (2025) found that contamination requires suf-
 419 ficient language representation to produce measur-
 420 able effects: for resource-scarce languages, con-
 421 tamination has almost no impact on performance.
 422 Cross-lingual contamination can exhibit a repre-
 423 sentation threshold, and multilingual claims should
 424 report audits or controlled tests of exposure to trans-
 425 lated benchmark variants.

426 **Scaling & Forgetting** Larger models ex-
 427 hibit stronger contamination effects than smaller
 428 ones (Kocyigit et al., 2025; Schaeffer et al., 2025b).
 429 As LLMs’ memorization ability grows signifi-
 430 cantly with their model size, we argue it becomes
 431 even easier for them to reproduce training data
 432 instances (Riddell et al., 2024). Recent studies sug-
 433 gest that benchmark contamination is not a binary
 434 state but a dynamic phenomenon governed by expo-
 435 sure magnitude and timing. Schaeffer et al. (2025a)
 436 resolved the “contamination paradox” by model-
 437 ing performance gains as a dose–response function,
 438 where benchmark inflation depends on the frac-
 439 tion of leaked tokens relative to model capacity
 440 and the incentive to memorize. Consistent with
 441 this view, Bordt et al. (2024b) observed that while
 442 even minor leakage causes measurable overfitting
 443 in Chinchilla-optimal (Hoffmann et al., 2022) se-
 444 tups, such effects are regime-dependent: in data-
 445 rich environments (e.g., $> 5 \times$ Chinchilla-optimal
 446 tokens), continued training can effectively *wash*
 447 *out* memorized traces via dilution. Consequently,
 448 late-phase contamination is significantly more detri-
 449 mental than early-stage exposure.

450 6 Data Contamination Detection

451 Prior work detects contamination via two
 452 paradigms: *statistical* approaches, which infer leak-
 453 age from observed signals (e.g., overlap, memo-
 454 rization) without training control; and *causal* ap-
 455 proaches, which intervene by injecting controlled
 456 data or markers to estimate effects. Despite this,
 457 current methods face a *consistency crisis*, often
 458 yielding contradictory diagnoses and remaining
 459 vulnerable to evasion via instruction tuning.

460 6.1 Statistical Approaches

461 We categorize statistical detectors by the level of
 462 model access required: white-box (train data), gray-
 463 box (probabilities), and black-box (outputs).

464 6.1.1 White-box

465 *White-box* methods use training data or model in-
 466 ternals to directly audit overlap.

467 **N-gram based** Exact n -gram overlap serves as
 468 the prevalent baseline for contamination detection,
 469 explicitly tracked in major LLM reports (Achiam
 470 et al., 2023). Addressing the computational costs
 471 of this method, Xu et al. (2025c) proposed INFINI-
 472 GRAM, an FM-index–based engine that facilitates
 473 memory-efficient exact substring search across
 474 trillion-token corpora.

475 **Embedding-based** To capture semantic contam-
 476 ination beyond lexical matching (Reimers, 2019),
 477 embedding similarity was employed for seman-
 478 tic deduplication in instruction tuning (Lee et al.,
 479 2023). Crucially, it underpinned the *LLM De-*
 480 *contaminator*, a hybrid pipeline that couples em-
 481 bedding retrieval with LLM verification to detect
 482 rephrased leakage (Yang et al., 2023).

483 6.1.2 Gray-box

484 *Gray-box* methods assume partial access to outputs
 485 (e.g., token probabilities) and derive contamination
 486 signals from perplexity or confidence statistics (in
 487 table 4). In 2022, Carlini et al. (2022a) established
 488 the foundational approach from a membership-
 489 inference perspective, showing that thresholding
 490 per-example negative log-likelihood (perplexity)
 491 serves as a principled baseline. By 2023, refer-
 492 ence models were introduced to refine detec-
 493 tion: Li (2023) benchmarked perplexity against
 494 specific references, while Wei et al. (2023) uti-
 495 lized style-matched sets (e.g., GSM8K) to iden-
 496 tify anomalies. Entering 2024, the field expanded
 497 into more granular probability analysis. While
 498 Duan et al. (2024) critically analyzed why tradi-
 499 tional MIAs underperform on LLMs due to fuzzy
 500 boundaries, Shi et al. (2024) introduced Min- $K\%$,
 501 shifting focus to aggregating low-probability token
 502 outliers. This method was rapidly refined: Min-
 503 $K\%++$ (Zhang et al., 2024c) added local proba-
 504 bility maxima, PAC (Ye et al., 2024) introduced
 505 input perturbations, and DC-PDD (Zhang et al.,
 506 2024d) incorporated corpus frequency divergences.
 507 Parallel to these probability-centric improvements,

PaCoST (Zhang et al., 2024b) proposed a threshold-free approach by statistically comparing confidence on original items vs. semantically equivalent paraphrases. Overall, token-based aggregation is cheaper than paraphrase-based generation.

6.1.3 Black-box

Black-box methods operate without access to training data or internals and rely solely on model outputs and behavioral assumptions in appendix D.1.2.

Memorization-based Prompt-and-completion protocols test whether a model can reconstruct masked benchmark content or consistently prefer original instances over perturbed alternatives, suggesting memorization: Golchin and Surdeanu (2023b) propose guided completion-based detection; DCQ (Golchin and Surdeanu, 2023a) used multiple-choice questions with synonym perturbations; Chang et al. (2023) studied memorization via cloze tasks and data archaeology; TS-Guessing (Deng et al., 2023) probed reconstruction of masked elements; DE-COP (Duarte et al., 2024) targets copyrighted content via verbal vs. paraphrased probing; and (Ranaldi et al., 2024) highlighted contamination signals in Text-to-SQL via masked schema recovery. Building on prior taxonomy, recent work has developed modality-aware techniques, including multimodal semantic perturbations that elicit contamination-sensitive behaviors (Park et al., 2025) and dynamic evaluation protocols that generate benchmark variants to reduce pre-training overlap (Yang et al., 2025a; Wang et al., 2025a).

Distributions-based Dong et al. (2024) proposed CDD based on the peakedness of output distributions and pair it with TED for mitigation while preserving evaluation validity.

Binary Judgments Khan et al. (2025) proposed NATURAL-DACODE, combining code naturalness with token-level completion accuracy via a lightweight classifier (SVM) to detect contamination. When leakage is rare, Kaneko and Baldwin (2025) showed that few-shot *self-detection* (in-context labeled leaked vs. non-leaked examples prompting explicit classification) can improve robustness in low-leakage regimes.

Quantifying Contamination Impact Beyond Binary Judgments Beyond binary “contaminated vs. clean” decisions, recent work aims to *quantify* contamination intensity and calibrate reported

scores. We say a benchmark B is contaminated when the training corpus D_{train} contains benchmark information (directly or via transformations), i.e., when $D_{\text{train}}^{\text{test}} \cap B^{\text{test}} \neq \emptyset$; the contamination risk is

$$\text{Score} = \frac{|D_{\text{train}}^{\text{test}} \cap B^{\text{test}}|}{|B^{\text{test}}|}. \quad (1)$$

Xu et al. (2025a) proposed DCR, an *efficient* framework that aggregates multi-level contamination signals into a unified risk factor and converts raw accuracy into a contamination-adjusted performance estimate for fairer model comparisons. From a representation-change perspective, Choi et al. (2025) introduced a Kernel Divergence Score that quantifies leakage by comparing the kernel similarity structure of sample embeddings before vs. after fine-tuning, motivated by the intuition that seen instances are perturbed differently than unseen ones.

6.2 Causal Approaches

Controlled Injection Causal studies quantify contamination by explicitly injecting benchmark content into training corpora at varying doses (e.g., frequency, timing) to map *dose-response* relationships (Jiang et al., 2024; Yao et al., 2024; Kocyyigit et al., 2025). While requiring substantial control over training pipelines, this design offers direct causal attribution of score inflation compared to observational diagnostics (Schaeffer et al., 2025a; Bordt et al., 2024b). For instance, Hidayat et al. (2025) utilized this framework to benchmark detection strategies under controlled leakage, ultimately validating simple n-gram/ROUGE-L matching as more reliable than complex permutation methods (Ni et al., 2025).

Interpretability White-box signals also include mechanistic and activation-level evidence: Tu et al. (2024) proposed DICE, training layer-wise contamination classifiers on intermediate activations during fine-tuning and showing that such signals correlate with performance inflation. Kattamuri et al. (2025) proposed RADAR, a mechanistic interpretability-based detector that separates *recall* from *reasoning* using forward-pass features (e.g., confidence convergence and attention/circuit dynamics), surfacing memorization-like behavior without requiring the training corpus. From a related membership-inference perspective, LUMIA (Ibanez-Lissen et al., 2025) trained linear probes over hidden states to discriminate *member*

606 vs. *non-member* samples, providing layer-localized
607 leakage evidence. [Zhu et al. \(2025\)](#) further argued
608 contamination can manifest as *shortcut neurons*;
609 they proposed identifying such neurons via causal
610 analyses and suppressing them via patching to ob-
611 tain more trustworthy evaluation behavior.

612 **Marker-Based Verification** This approach en-
613 gineered detectable signals into benchmarks *pre-*
614 *release* to enable verifiable black-box testing. Tech-
615 niques included *Canary insertion* ([Carlini et al.,](#)
616 [2021a](#)) (measuring recall of synthetic sequences),
617 *DyePack* ([Cheng et al., 2025](#)) (identifying learned
618 backdoor behaviors with provable false-positive
619 control), and *Watermarking* ([Sander et al., 2025](#))
620 (detecting signal “radioactivity”). Unlike observa-
621 tional methods, these protocols transform detection
622 into a statistically grounded verification procedure
623 with explicit error guarantees.

6.3 Detection Methods Discussion

625 **Fragility and Inconsistency** Under a ground-
626 truth contamination setting (involving instruction
627 fine-tuning with CoT-style augmentation), widely
628 used detectors exhibit non-monotonic sensitivity
629 and weak cross-method agreement (low Spearman
630 correlation), resulting in contradictory judgments
631 even for frontier LLMs ([Dekoninck et al., 2024b](#);
632 [Fu et al., 2024](#)). And when leakage is *rare*, de-
633 tectors based purely on scoring signals (PPL, Min-
634 $K\%$) can degrade;

635 **Decoding Sensitivity** Black-box detectors rely-
636 ing on model outputs (e.g., TS-Guessing) are sus-
637 ceptible to decoding parameters; specifically, high-
638 temperature sampling can mask memorization sig-
639 natures that are otherwise visible under greedy de-
640 coding ([Schaeffer et al., 2025b](#)).

7 Future Directions

642 **From Qualitative Detection to Contamination**
643 **Scaling Laws** Current research predominantly
644 treats contamination as a qualitative binary, iden-
645 tifying *existence* rather than quantifying *impact*.
646 We advocate for formalizing *contamination scal-*
647 *ing laws* analogous to scaling laws ([Kaplan et al.,](#)
648 [2020](#)). Future work should derive functional rela-
649 tionships (e.g., power-law dynamics $\Delta P \propto \alpha \cdot C^\beta$)
650 that map leakage density (C) and model scale to
651 benchmark inflation (ΔP). Establishing these laws
652 would transition the field from post-hoc forensics
653 to predictive science, enabling developers to fore-

cast metric gains from leakage rates and perform
mathematically grounded score calibration.

654 **Synthetic Data Contamination** A critical fron-
655 tier is ensuring evaluation robustness amidst *syn-*
656 *thetic saturation*. We propose three directions: (i)
657 developing *provenance and quantification* tools to
658 track synthetic exposure and recursive loops in
659 training corpora; (ii) designing *dynamic* reasoning
660 benchmarks with continuously renewable instances
661 to prevent memorization of finite datasets; and (iii)
662 aligning evaluation with *training-time mitigation*
663 strategies (e.g., data curation policies) to prevent
664 model collapse ([Shumailov et al., 2024](#)), ensuring
665 benchmarks remain predictive as the web becomes
666 increasingly model-generated.

667 **Generative Benchmarks Contamination** As
668 discussed in Section 5, the transition toward gener-
669 ative and agentic evaluation necessitates that con-
670 tamination research account for stochastic decod-
671 ing and output diversity. We advocate for a multi-
672 faceted approach: (i) *Robust Protocols*: shifting
673 from text matching to execution-based scoring and
674 semantic equivalence, distinguishing metric infla-
675 tion from genuine reasoning; (ii) *Causal Dynamics*:
676 analyzing contamination via dose–response studies
677 that account for sampling strategies (e.g., tempera-
678 ture) and prompt sensitivity; and (iii) *Agentic Gov-*
679 *ernance*: treating interaction traces and tool logs
680 as leakage vectors, necessitating sandboxed envi-
681 ronments and auditable provenance. Standardized
682 reporting remains essential to ensure comparability
683 amidst these shifts ([Schaeffer et al., 2025b](#)).

8 Conclusion

684 Data contamination fundamentally threatens the
685 trustworthiness of LLM evaluation. This survey
686 systematizes the field by establishing a compre-
687 hensive *taxonomy* of leakage pathways, analyzing
688 the *regime-dependent characteristics* of contamina-
689 tion (e.g., scale and forgettability), and contrasting
690 observational *detection* with causal interventions.
691 Beyond merely inflating metrics, we highlight that
692 contamination compromises fairness and scientific
693 validity. Consequently, reliable benchmarking re-
694 quires treating contamination as a *first-class vari-*
695 *able* through rigorous audits and transparent report-
696 ing. Future progress hinges on developing quanti-
697 tative contamination scaling laws and addressing
698 *synthetic data feedback loops* to sustain evaluation
699 validity in an increasingly model-generated web.
700
701
702

9 Limitations

While we extensively cover various forms of data contamination, it is possible that new contamination mechanisms or models may not be fully captured in our analysis. Additionally, our focus is primarily on data contamination within the context of LLMs, and we may not have fully incorporated previous research on data contamination in other areas of machine learning. Given the proliferation of benchmarks, we select representative examples in appendix to illustrate data construction methods. Additionally, as this survey focuses on LLM data contamination, we may not cover all related areas such as membership inference attacks (MIA), machine unlearning, and LLM memorization. Our study is limited by the under-exploration of dynamic benchmarks and a predominantly technical focus that restricts the assessment of societal risks arising from contamination across specific languages, domains, or demographics.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, M. Saiful Bari, and Haidar Khan. 2024. When benchmarks are targets: Revealing the sensitivity of large language model leaderboards. In *ACL*.
- Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, and 1 others. 2024. Benchmarking foundation models with language-model-as-an-examiner. *Advances in Neural Information Processing Systems*, 36.
- Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, and Ondrej Dusek. 2024. Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 67–93.
- Sebastian Bordt, Harsha Nori, and Rich Caruana. 2023. Elephants never forget: Testing language models for memorization of tabular data. In *NeurIPS 2023 Second Table Representation Learning Workshop*.
- Sebastian Bordt, Harsha Nori, Vanessa Rodrigues, Besmira Nushi, and Rich Caruana. 2024a. Elephants

never forget: Memorization and learning of tabular data in large language models. In *Conference on Language Modeling (COLM)*.

- Sebastian Bordt, Suraj Srinivas, Valentyn Boreiko, and Ulrike von Luxburg. 2024b. How much can we forget about data contamination? *arXiv preprint arXiv:2410.03249*.
- Boxi Cao, Mengjie Ren, Hongyu Lin, Xianpei Han, Feng Zhang, Junfeng Zhan, and Le Sun. 2024. [Structeval: Deepen and broaden large language model assessment via structured evaluation](#). *Preprint*, arXiv:2408.03281.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. 2022a. Membership inference attacks from first principles. In *2022 IEEE symposium on security and privacy (SP)*, pages 1897–1914. IEEE.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiuyan Zhang. 2022b. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. 2021a. [Extracting training data from large language models](#). *Preprint*, arXiv:2012.07805.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, and 1 others. 2021b. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Wang, Samuel Marks, Charbel-Raphaël Segerie, Micah Carroll, Andi Peng, Phillip Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, and 13 others. 2023. [Open problems and fundamental limitations of reinforcement learning from human feedback](#). *Preprint*, arXiv:2307.15217.
- Kent Chang, Mackenzie Cramer, Sandeep Soni, and David Bamman. 2023. [Speak, memory: An archaeology of books known to ChatGPT/GPT-4](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7312–7327, Singapore. Association for Computational Linguistics.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, and 1 others. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45.

812	Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang	Jasper Dekoninck, Mark Niklas Müller, Maximilian	868
813	Zang, Zehui Chen, Haodong Duan, Jiaqi Wang,	Baader, Marc Fischer, and Martin Vechev. 2024b.	869
814	Yu Qiao, Dahua Lin, and Feng Zhao. 2024. Are we	Evading data contamination detection for language	870
815	on the right way for evaluating large vision-language	models is (too) easy . <i>Preprint</i> , arXiv:2402.02823.	871
816	models? In <i>Advances in Neural Information Process-</i>		
817	ing Systems .		
818	Pin-Er Chen, Da-Chen Lian, Shu-Kai Hsieh, Sieh-	Chunyuan Deng, Yilun Zhao, Yuzhao Heng, Yitong Li,	872
819	Chuen Huang, Hsuan-Lei Shao, Jun-Wei Chiu, Yang-	Jiannan Cao, Xiangru Tang, and Arman Cohan. 2024.	873
820	Hsien Lin, Zih-Ching Chen, Cheng-Kuang, Eddie TC	Unveiling the spectrum of data contamination in lan-	874
821	Huang, and Simon See. 2025a. Continual pre-	guage model: A survey from detection to remediation .	875
822	training is (not) what you need in domain adaption .	In <i>Findings of the Association for Computational Lin-</i>	876
823	<i>Preprint</i> , arXiv:2504.13603.	<i>guistics: ACL 2024</i> , pages 16078–16092, Bangkok,	877
		Thailand. Association for Computational Linguistics.	878
824	Simin Chen, Yiming Chen, Zexin Li, Yifan Jiang,	Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Ger-	879
825	Zhongwei Wan, Yixin He, Dezhi Ran, Tianle Gu,	stein, and Arman Cohan. 2023. Investigating data	880
826	Haizhou Li, Tao Xie, and Baishakhi Ray. 2025b. Re-	contamination in modern benchmarks for large lan-	881
827	cent advances in large language model benchmarks	guage models. <i>arXiv preprint arXiv:2311.09783</i> .	882
828	against data contamination: From static to dynamic		
829	evaluation . <i>Preprint</i> , arXiv:2502.17521.	Jesse Dodge, Maarten Sap, Ana Marasović, William	883
		Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret	884
830	Simin Chen, Pranav Pusarla, and Baishakhi Ray. 2025c.	Mitchell, and Matt Gardner. 2021. Documenting	885
831	Dynamic benchmarking of reasoning capabilities in	large webtext corpora: A case study on the colossal	886
832	code large language models under data contamina-	clean crawled corpus . <i>Preprint</i> , arXiv:2104.08758.	887
833	tion . <i>Preprint</i> , arXiv:2503.04149.		
834	Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi,	Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, Bin Gu,	888
835	Minlie Huang, and Furu Wei. 2024. Instruction pre-	Mengfei Yang, and Ge Li. 2024. Generalization or	889
836	training: Language models are supervised multitask	memorization: Data contamination and trustworthy	890
837	learners . <i>arXiv preprint arXiv:2406.14491</i> .	evaluation for large language models. In <i>Findings of</i>	891
		<i>the Association for Computational Linguistics: ACL</i>	892
838	Yize Cheng, Wenxiao Wang, Mazda Moayeri, and So-	<i>2024</i> , pages 12039–12050.	893
839	heil Feizi. 2025. DyePack: Provably flagging test set	Michael Duan, Anshuman Suri, Niloofar Miresghallah,	894
840	contamination in LLMs using backdoors . In <i>Proceed-</i>	Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia	895
841	<i>ings of the 2025 Conference on Empirical Methods in</i>	Tsvetkov, Yejin Choi, David Evans, and Hannaneh	896
842	<i>Natural Language Processing</i> , pages 15367–15384,	Hajishirzi. 2024. Do membership inference attacks	897
843	Suzhou, China. Association for Computational Lin-	work on large language models? <i>arXiv preprint</i>	898
844	guistics.	<i>arXiv:2402.07841</i> .	899
845	Hyeong Kyu Choi, Maxim Khanov, Hongxin Wei,	André V. Duarte, Xuandong Zhao, Arlindo L. Oliveira,	900
846	and Yixuan Li. 2025. How contaminated is your	and Lei Li. 2024. De-cop: Detecting copyrighted	901
847	benchmark? quantifying dataset leakage in large	content in language models training data. <i>arXiv</i>	902
848	language models with kernel divergence . <i>Preprint</i> ,	<i>preprint arXiv:2402.09910</i> .	903
849	arXiv:2502.00678.		
850	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	Lizhou Fan, Wenyue Hua, Xiang Li, Kaijie Zhu,	904
851	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias	Mingyu Jin, Lingyao Li, Haoyang Ling, Jinkui Chi,	905
852	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	Jindong Wang, Xin Ma, and Yongfeng Zhang. 2024.	906
853	Nakano, Christopher Hesse, and John Schulman.	Nphardeval4v: A dynamic reasoning benchmark of	907
854	2021. Training verifiers to solve math word prob-	multimodal large language models . <i>arXiv preprint</i>	908
855	lems . <i>arXiv preprint arXiv:2110.14168</i> .	<i>arXiv:2403.01777</i> .	909
856	DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,	Yujuan Fu, Ozlem Uzuner, Meliha Yetisgen, and Fei	910
857	Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,	Xia. 2024. Does data contamination detection work	911
858	Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,	(well) for llms? a survey and evaluation on detection	912
859	Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-	assumptions . <i>arXiv preprint arXiv:2410.18966</i> .	913
860	hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.		
861	2025. Deepseek-r1: Incentivizing reasoning capabil-	Shahriar Golchin and Mihai Surdeanu. 2023a. Data con-	914
862	ity in llms via reinforcement learning . <i>arXiv preprint</i>	tamination quiz: A tool to detect and estimate con-	915
863	<i>arXiv:2501.12948</i> .	tamination in large language models . <i>arXiv preprint</i>	916
864	Jasper Dekoninck, Mark Niklas Müller, and Martin	<i>arXiv:2311.06233</i> .	917
865	Vechev. 2024a. Constat: Performance-based con-	Shahriar Golchin and Mihai Surdeanu. 2023b. Time	918
866	tamination detection in large language models . <i>arXiv</i>	travel in llms: Tracing data contamination in large	919
867	<i>preprint arXiv:2405.16281</i> .	language models . <i>arXiv preprint arXiv:2308.08493</i> .	920

921	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	979
922		980
923		981
924		982
925		983
926	Ziwen Han, Meher Mankikar, Julian Michael, and Zifan Wang. 2025. Search-time data contamination. <i>arXiv preprint arXiv:2508.13180</i> .	984
927		985
928		
929	Jamie Hayes, Marika Swanberg, Harsh Chaudhari, Itay Yona, Iliia Shumailov, Milad Nasr, Christopher A. Choquette-Choo, Katherine Lee, and A. Feder Cooper. 2025. Measuring memorization in language models via probabilistic extraction . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 9266–9291, Albuquerque, New Mexico. Association for Computational Linguistics.	986
930		987
931		988
932		989
933		990
934		991
935		992
936		993
937		994
938		995
939	Naila Shafirni Hidayat, Muhammad Dehan Al Kautsar, Alfian Farizki Wicaksono, and Fajri Koto. 2025. Simulating training data leakage in multiple-choice benchmarks for llm evaluation . <i>Preprint</i> , arXiv:2505.24263.	996
940		997
941		998
942		999
943		1000
944	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, and 3 others. 2022. Training compute-optimal large language models . <i>Preprint</i> , arXiv:2203.15556.	1001
945		1002
946		1003
947		1004
948		1005
949		1006
950		1007
951		1008
952		1009
953	Ruijie Hou, Yueyang Jiao, Hanxu Hu, Yingming Li, Wai Lam, Huajian Zhang, and Hongyuan Lu. 2025. Lne-blocking: An efficient framework for contamination mitigation evaluation on large language models . <i>Preprint</i> , arXiv:2509.15218.	1010
954		1011
955		1012
956		1013
957		1014
958		1015
959	Luis Ibanez-Lissen, Lorena Gonzalez-Manzano, Jose Maria de Fuentes, Nicolas Anciaux, and Joaquin Garcia-Alfaro. 2025. Lumia: Linear probing for unimodal and multimodal membership inference attacks leveraging internal llm states . <i>Preprint</i> , arXiv:2411.19876.	1016
960		1017
961		1018
962		
963		
964	Yui Ishikawa. 2025. Data contamination or genuine generalization? disentangling llm performance on benchmarks . <i>Academic Journal of Natural Science</i> , 2(2):16–22.	1019
965		1020
966		1021
967		1022
968	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Live-codebench: Holistic and contamination free evaluation of large language models for code. <i>arXiv preprint arXiv:2403.07974</i> .	1023
969		1024
970		1025
971		1026
972		1027
973		1028
974		1029
975		1030
976		
977	Minhao Jiang, Ken Ziyu Liu, Ming Zhong, Rylan Schaeffer, Siru Ouyang, Jiawei Han, and Sanmi Koyejo. 2024. Investigating data contamination for pre-training language models. <i>arXiv preprint arXiv:2401.06059</i> .	1031
978		1032
		1033
		1034
		1035
	Masahiro Kaneko and Timothy Baldwin. 2025. Investigating how pre-training data leakage affects models' reproduction and detection capabilities . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 23556–23566, Suzhou, China. Association for Computational Linguistics.	
	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models . <i>Preprint</i> , arXiv:2001.08361.	
	Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. 2023. Copyright violations and large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7403–7412, Singapore. Association for Computational Linguistics.	
	Ashish Kattamuri, Harshwardhan Fartale, Arpita Vats, Rahul Raja, and Ishita Prasad. 2025. Radar: Mechanistic pathways for detecting data contamination in llm evaluation. <i>arXiv preprint arXiv:2510.08931</i> .	
	Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. 2023. Continual pre-training of language models . <i>Preprint</i> , arXiv:2302.03241.	
	Haris Ali Khan, Yanjie Jiang, Qasim Umer, Yuxia Zhang, Waseem Akram, and Hui Liu. 2025. Has my code been stolen for model training? a naturalness based approach to code contamination detection . <i>Proc. ACM Softw. Eng.</i> , 2(FSE).	
	Muhammed Yusuf Kocyigit, Eleftheria Briakou, Daniel Deutsch, Jiaming Luo, Colin Cherry, and Markus Freitag. 2025. Overestimation in llm evaluation: A controlled large-scale study on data contamination's impact on machine translation . <i>Preprint</i> , arXiv:2501.18771.	
	Ariel N Lee, Cole J Hunter, and Nataniel Ruiz. 2023. Platypus: Quick, cheap, and powerful refinement of llms. <i>arXiv preprint arXiv:2308.07317</i> .	
	Fangyu Lei, Qian Liu, Yiming Huang, Shizhu He, Jun Zhao, and Kang Liu. 2024. S3Eval: A synthetic, scalable, systematic evaluation suite for large language model. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1259–1286.	
	Changmao Li and Jeffrey Flanigan. 2024. Task contamination: Language models may not be few-shot anymore. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 18471–18480.	
	Dawei Li, Renliang Sun, Yue Huang, Ming Zhong, Bohan Jiang, Jiawei Han, Xiangliang Zhang, Wei Wang, and Huan Liu. 2025. Preference leakage: A contamination problem in llm-as-a-judge . <i>Preprint</i> , arXiv:2502.01534.	

1036	Jia Li, Ge Li, Xuanming Zhang, Yihong Dong, and Zhi Jin. 2024a. Evocodebench: An evolving code generation benchmark aligned with real-world code repositories. <i>arXiv preprint arXiv:2404.00599</i> .	1091
1037		1092
1038		
1039		
1040	Xiang Li, Yunshi Lan, and Chao Yang. 2024b. Treeeval: Benchmark-free evaluation of large language models through tree planning. <i>arXiv preprint arXiv:2402.13125</i> .	1093
1041		1094
1042		1095
1043		1096
1044	Yanyang Li. 2024. Awesome data contamination. https://github.com/lyy1994/awesome-data-contamination .	1097
1045		1098
1046		1099
1047	Yanyang Li, Tin Long Wong, Cheung To Hung, Jianqiao Zhao, Duo Zheng, Ka Wai Liu, Michael R. Lyu, and Liwei Wang. 2024c. C ² leva: Toward comprehensive and contamination-free language model evaluation. <i>arXiv preprint arXiv:2412.04947</i> .	1100
1048		1101
1049		
1050		
1051		
1052	Yanyang Li, Jianqiao Zhao, Duo Zheng, Zi-Yuan Hu, Zhi Chen, Xiaohui Su, Yongfeng Huang, Shijia Huang, Dahua Lin, Michael Lyu, and Liwei Wang. 2023. CLEVA: Chinese language models EVALuation platform. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 186–217.	1102
1053		1103
1054		1104
1055		1105
1056		1106
1057		
1058		
1059	Yucheng Li. 2023. Estimating contamination via perplexity: Quantifying memorisation in language model evaluation. <i>arXiv preprint arXiv:2309.10677</i> .	1107
1060		1108
1061		1109
1062	Yucheng Li, Frank Guerin, and Chenghua Lin. 2024d. Latesteval: Addressing data contamination in language model evaluation through dynamic and time-sensitive test construction. <i>arXiv preprint arXiv:2312.12343</i> .	1110
1063		1111
1064		1112
1065		1113
1066		1114
1067	Yucheng Li, Yunhao Guo, Frank Guerin, and Chenghua Lin. 2024e. An open-source data contamination report for large language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 528–541, Miami, Florida, USA. Association for Computational Linguistics.	1115
1068		1116
1069		1117
1070		1118
1071		1119
1072		1120
1073	Chuang Liu, Renren Jin, Mark Steedman, and Deyi Xiong. 2024. Evaluating Chinese large language models on discipline knowledge acquisition via memorization and robustness assessment. In <i>Proceedings of the 1st Workshop on Data Contamination (CONDA)</i> , pages 1–12, Bangkok, Thailand. Association for Computational Linguistics.	1121
1074		1122
1075		1123
1076		1124
1077		
1078		
1079		
1080	Ming Liu and Wensheng Zhang. 2025. Reasoning multi-modal large language model: Data contamination and dynamic evaluation. <i>Preprint</i> , arXiv:2506.07202.	1125
1081		1126
1082		
1083	Sadeh Mahdavi, Muchen Li, Kaiwen Liu, Christos Thrampoulidis, Leonid Sigal, and Renjie Liao. 2025. Leveraging online olympiad-level math problems for llms training and contamination-resistant evaluation. <i>Preprint</i> , arXiv:2501.14275.	1127
1084		1128
1085		1129
1086		1130
1087		1131
1088	Felipe Maia Polo, Ronald Xu, Lucas Weber, Mírian Silva, Onkar Bhardwaj, Leshem Choshen, Allysson Flavio Melo de Oliveira, Yuekai Sun, and Mikhail	1132
1089		1133
1090		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
	Yurochkin. 2024. Efficient multi-prompt evaluation of LLMs.	
	Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. <i>Preprint</i> , arXiv:2303.08896.	
	Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2025. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. <i>Preprint</i> , arXiv:2410.05229.	
	Shiwen Ni, Xiangtao Kong, Chengming Li, Xiping Hu, Ruifeng Xu, Jia Zhu, and Min Yang. 2025. Training on the benchmark is not all you need. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 24948–24956.	
	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	
	Medha Palavalli, Amanda Bertsch, and Matthew Gormley. 2024. A taxonomy for data contamination in large language models. In <i>Proceedings of the 1st Workshop on Data Contamination (CONDA)</i> , pages 22–40, Bangkok, Thailand. Association for Computational Linguistics.	
	Jaden Park, Mu Cai, Feng Yao, Jingbo Shang, Soochahn Lee, and Yong Jae Lee. 2025. Contamination detection for vlms using multi-modal semantic perturbation. <i>arXiv preprint arXiv:2511.03774</i> .	
	Long Phan, Alice Gatti, Ziwen Han, and et al. 2025. Humanity’s last exam. <i>Preprint</i> , arXiv:2501.14249.	
	Kun Qian, Shunji Wan, Claudia Tang, Youzhi Wang, Xuanming Zhang, Maximillian Chen, and Zhou Yu. 2024. Varbench: Robust language model benchmarking through dynamic variable perturbation. <i>arXiv preprint arXiv:2406.17681</i> .	
	Federico Ranaldi, Elena Sofia Ruzzetti, Dario Onorati, Leonardo Ranaldi, Cristina Giannone, Andrea Favalli, Raniero Romagnoli, and Fabio Massimo Zanzotto. 2024. Investigating the impact of data contamination of large language models in text-to-SQL translation. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 13909–13920, Bangkok, Thailand. Association for Computational Linguistics.	
	Mathieu Ravaut, Bosheng Ding, Fangkai Jiao, Hailin Chen, Xingxuan Li, Ruo Chen Zhao, Chengwei Qin, Caiming Xiong, and Shafiq Joty. 2024. How much are llms contaminated? a comprehensive survey and the llmsanitize library. <i>arXiv preprint arXiv:2404.00699</i> .	

1147	N Reimers. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .	1201
1148		1202
1149		1203
1150	Martin Riddell, Ansong Ni, and Arman Cohan. 2024. Quantifying contamination in evaluating code generation capabilities of language models. <i>arXiv preprint arXiv:2403.04811</i> .	1204
1151		1205
1152		1206
1153		1207
1154	Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10776–10787.	1208
1155		1209
1156		1210
1157		1211
1158		1212
1159		1213
1160	Vinay Samuel, Yue Zhou, and Henry Peng Zou. 2024. Towards data contamination detection for modern large language models: Limitations, inconsistencies, and oracle challenges. <i>Preprint</i> , arXiv:2409.09927.	1214
1161		1215
1162		1216
1163		1217
1164	Tom Sander, Pierre Fernandez, Saeed Mahloujifar, Alain Durmus, and Chuan Guo. 2025. Detecting benchmark contamination through watermarking. <i>Preprint</i> , arXiv:2502.17259.	1218
1165		1219
1166		1220
1167		1221
1168	Rylan Schaeffer, Ken Liu, Brando Miranda, Ahmed M Ahmed, Niloofar Mireshghallah, and Sanmi Koyejo. 2025a. The contamination paradox: Why test set leakage can be both potent and negligible. In <i>NeurIPS 2025 Workshop on Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling</i> .	1222
1169		1223
1170		1224
1171		1225
1172		1226
1173		1227
1174		1228
1175	Rylan Schaeffer, Brando Miranda, Joshua Kazdan, Ken Liu, Ahmed M Ahmed, Niloofar Mireshghallah, and Sanmi Koyejo. 2025b. Causally quantifying the effect of test set contamination on generative benchmarks. In <i>NeurIPS 2025 Workshop on Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling</i> .	1229
1176		1230
1177		1231
1178		1232
1179		1233
1180		1234
1181		1235
1182	Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary C. Lipton, and J. Zico Kolter. 2024. Rethinking llm memorization through the lens of adversarial compression. <i>arXiv preprint arXiv:2404.15146</i> .	1236
1183		1237
1184		1238
1185		1239
1186	Radzim Sendyka, Christian Cabrera, Andrei Paleyes, Diana Robinson, and Neil Lawrence. 2025. Llm performance for code generation on noisy tasks. <i>Preprint</i> , arXiv:2505.23598.	1240
1187		1241
1188		1242
1189		1243
1190	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>Preprint</i> , arXiv:2402.03300.	1244
1191		1245
1192		1246
1193		1247
1194		1248
1195		1249
1196	Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024. Detecting pretraining data from large language models. <i>arXiv preprint arXiv:2310.16789</i> .	1250
1197		1251
1198		1252
1199		1253
1200		1254
	Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. 2024. Ai models collapse when trained on recursively generated data. <i>Nature</i> , 631(8022):755–759.	1201
		1202
		1203
		1204
	Shivalika Singh, Yiyang Nan, Alex Wang, Daniel D’Souza, Sayash Kapoor, Ahmet Üstün, Sanmi Koyejo, Yuntian Deng, Shayne Longpre, Noah Smith, Beyza Ermis, Marzieh Fadaee, and Sara Hooker. 2025. The leaderboard illusion.	1205
		1206
		1207
		1208
		1209
	Dingjie Song, Sicheng Lai, Mingxuan Wang, Shunian Chen, Lichao Sun, and Benyou Wang. 2025. Both text and images leaked! a systematic analysis of data contamination in multimodal LLM. In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 10527–10542, Suzhou, China. Association for Computational Linguistics.	1210
		1211
		1212
		1213
		1214
		1215
		1216
	Saurabh Srivastava, Anto PV, Shashank Menon, Ajay Sukumar, Alan Philipose, Stevin Prince, Sooraj Thomas, and 1 others. 2024. Functional benchmarks for robust evaluation of reasoning performance, and the reasoning gap. <i>arXiv preprint arXiv:2402.19450</i> .	1217
		1218
		1219
		1220
		1221
	Yongding Tao, Tian Wang, Yihong Dong, Huanyu Liu, Kechi Zhang, Xiaolong Hu, and Ge Li. 2025. Detecting data contamination from reinforcement learning post-training for large language models. <i>arXiv preprint arXiv:2510.09259</i> .	1222
		1223
		1224
		1225
		1226
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	1227
		1228
		1229
		1230
		1231
		1232
	Yuan Tseng, Titouan Parcollet, Rogier van Dalen, Shucong Zhang, and Sourav Bhattacharya. 2025. Evaluation of llms in speech is often flawed: Test set contamination in large language models for speech recognition. <i>Preprint</i> , arXiv:2505.22251.	1233
		1234
		1235
		1236
		1237
	Shangqing Tu, Kejian Zhu, Yushi Bai, Zijun Yao, Lei Hou, and Juanzi Li. 2024. Dice: Detecting in-distribution contamination in llm’s fine-tuning phase for math reasoning. <i>arXiv preprint arXiv:2406.04197</i> .	1238
		1239
		1240
		1241
		1242
	Pablo Villalobos, Anson Ho, Jaime Sevilla, Tamay Besiroglu, Lennart Heim, and Marius Hobbhahn. 2024. Will we run out of data? limits of llm scaling based on human-generated data. <i>Preprint</i> , arXiv:2211.04325.	1243
		1244
		1245
		1246
	Hanqing Wang, Yuan Tian, Mingyu Liu, Zhenhao Zhang, and Xiangyang Zhu. 2025a. Sdeval: Safety dynamic evaluation for multimodal large language models. <i>arXiv preprint arXiv:2508.06142</i> .	1247
		1248
		1249
		1250
	Siyuan Wang, Zhuohan Long, Zhihao Fan, Zhongyu Wei, and Xuanjing Huang. 2024. Benchmark self-evolving: A multi-agent framework for dynamic llm evaluation. <i>arXiv preprint arXiv:2402.11443</i> .	1251
		1252
		1253
		1254

1255	Zeng Wang, Minghao Shao, Jitendra Bhandari, Likhitha Mankali, Ramesh Karri, Ozgur Sinanoglu, Muhammad Shafique, and Johann Knechtel. 2025b. Vericon-taminated: Assessing llm-driven verilog coding for data contamination . <i>Preprint</i> , arXiv:2503.13572.	<i>Natural Language Processing</i> , pages 14748–14762, Suzhou, China. Association for Computational Linguistics.	1312 1313 1314
1260	Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. 2024. Measuring short-form factuality in large language models . <i>Preprint</i> , arXiv:2411.04368.	Hao Xu, Jiacheng Liu, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2025c. Infini-gram mini: Exact n-gram search at the Internet scale with FM-index . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 24955–24980, Suzhou, China. Association for Computational Linguistics.	1315 1316 1317 1318 1319 1320 1321
1265	Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng, Weiwei Lü, Rui Hu, Chenxia Li, Liu Yang, Xilin Luo, Xuejie Wu, Lunan Liu, Wenjun Cheng, Peng Cheng, Jianhao Zhang, Xiaoyu Zhang, and 11 others. 2023. Skywork: A more open bilingual foundation model . <i>Preprint</i> , arXiv:2310.19341.	Xin Xu, Jiaxin Zhang, Tianhao Chen, Zitong Chao, Jishan Hu, and Can Yang. 2025d. Ugmathbench: A diverse and dynamic benchmark for undergraduate-level mathematical reasoning with large language models . <i>Preprint</i> , arXiv:2501.13766.	1322 1323 1324 1325 1326
1272	Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Siddartha Naidu, and 1 others. 2024. Livebench: A challenging, contamination-free llm benchmark. <i>arXiv preprint arXiv:2406.19314</i> .	Vishnu Vardhan Yadoji. 2025. Data centric guard (dc-guard)-a framework for trustworthy llm evaluation. In <i>NeurIPS 2025 Workshop on Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling</i> .	1327 1328 1329 1330 1331
1277	Hengyu Wu and Yang Cao. 2025. Membership inference attacks on large-scale models: A survey. <i>arXiv preprint arXiv:2503.19338</i> .	Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E. Gonzalez, and Ion Stoica. 2023. Rethinking benchmark and contamination for language models with rephrased samples. <i>arXiv preprint arXiv:2311.04850</i> .	1332 1333 1334 1335 1336
1280	Mingqi Wu, Zhihao Zhang, Qiaole Dong, Zhiheng Xi, Jun Zhao, Senjie Jin, Xiaoran Fan, Yuhao Zhou, Huijie Lv, Ming Zhang, and 1 others. 2025a. Reasoning or memorization? unreliable results of reinforcement learning due to data contamination. <i>arXiv preprint arXiv:2507.10532</i> .	Yue Yang, Shuibo Zhang, Kaipeng Zhang, Yi Bin, Yu Wang, Ping Luo, and Wenqi Shao. 2025a. Dynamic multimodal evaluation with flexible complexity by vision-language bootstrapping . In <i>International Conference on Learning Representations (ICLR)</i> . Oral.	1337 1338 1339 1340 1341 1342
1286	Xiaobao Wu, Liangming Pan, Yuxi Xie, Ruiwen Zhou, Shuai Zhao, Yubo Ma, Mingzhe Du, Rui Mao, Anh Tuan Luu, and William Yang Wang. 2025b. Antileakbench: Preventing data contamination by automatically constructing benchmarks with updated real-world knowledge . <i>Preprint</i> , arXiv:2412.13670.	Zhen Yang, Hongyi Lin, Yifan He, Jie Xu, Zeyu Sun, Shuo Liu, Pengpeng Wang, Zhongxing Yu, and Qingyuan Liang. 2025b. Rethinking the effects of data contamination in code intelligence . <i>Preprint</i> , arXiv:2506.02791.	1343 1344 1345 1346 1347
1292	Yunjia Xi, Jianghao Lin, Yongzhao Xiao, Zheli Zhou, Rong Shan, Te Gao, Jiachen Zhu, Weiwen Liu, Yong Yu, and Weinan Zhang. 2025. A survey of llm-based deep search agents: Paradigm, optimization, evaluation, and challenges . <i>Preprint</i> , arXiv:2508.05668.	Feng Yao, Yufan Zhuang, Zihao Sun, Sunan Xu, Animesh Kumar, and Jingbo Shang. 2024. Data contamination can cross language barriers. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 17864–17875.	1348 1349 1350 1351 1352
1297	Cheng Xu, Shuhao Guan, Derek Greene, M Kechadi, and 1 others. 2024. Benchmark data contamination of large language models: A survey. <i>arXiv preprint arXiv:2406.04244</i> .	Wentao Ye, Jiaqi Hu, Liyao Li, Haobo Wang, Gang Chen, and Junbo Zhao. 2024. Data contamination calibration for black-box LLMs. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 10845–10861.	1353 1354 1355 1356 1357
1301	Cheng Xu, Nan Yan, Shuhao Guan, Changhong Jin, Yuke Mei, Yibing Guo, and Tahar Kechadi. 2025a. DCR: Quantifying data contamination in LLMs evaluation . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 23013–23031, Suzhou, China. Association for Computational Linguistics.	Jiahao Ying, Yixin Cao, Yushi Bai, Qianru Sun, Bo Wang, Wei Tang, Zhaojun Ding, Yizhe Yang, Xuanjing Huang, and YAN Shuicheng. 2024. Automating dataset updates towards reliable and timely evaluation of large language models. In <i>The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	1358 1359 1360 1361 1362 1363 1364
1308	Cheng Xu, Nan Yan, Shuhao Guan, Yuke Mei, and Tahar Kechadi. 2025b. SSA: Semantic contamination of LLM-driven fake news detection . In <i>Proceedings of the 2025 Conference on Empirical Methods in</i>	Zhuohao Yu, Chang Gao, Wenjin Yao, Yidong Wang, Wei Ye, Jindong Wang, Xing Xie, Yue Zhang,	1365 1366

1367	and Shikun Zhang. 2024. KIEval: A knowledge-grounded interactive evaluation framework for large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5967–5985.	
1368		
1369		
1370		
1371		
1372	Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. 2025. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? <i>Preprint</i> , arXiv:2504.13837.	
1373		
1374		
1375		
1376		
1377	Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav Raja, Dylan Slack, Qin Lyu, and 1 others. 2024a. A careful examination of large language model performance on grade school arithmetic. <i>arXiv preprint arXiv:2405.00332</i> .	
1378		
1379		
1380		
1381		
1382		
1383	Huixuan Zhang, Yun Lin, and Xiaojun Wan. 2024b. PaCoST: Paired confidence significance testing for benchmark contamination detection in large language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 1794–1809.	
1384		
1385		
1386		
1387		
1388	Jingyang Zhang, Jingwei Sun, Eric Yeats, Yang Ouyang, Martin Kuo, Jianyi Zhang, Hao Frank Yang, and Hai Li. 2024c. Min-karXiv preprint arXiv:2404.02936.	
1389		
1390		
1391	Linghao Zhang, Shilin He, Chaoyun Zhang, Yu Kang, Bowen Li, Chengxing Xie, Junhao Wang, Maoquan Wang, Yufan Huang, Shengyu Fu, Elsie Nallipogu, Qingwei Lin, Yingnong Dang, Saravan Rajmohan, and Dongmei Zhang. 2025. Swe-bench goes live! <i>Preprint</i> , arXiv:2505.23419.	
1392		
1393		
1394		
1395		
1396		
1397	Weichao Zhang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2024d. Pre-training data detection for large language models: A divergence-based calibration method. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 5263–5274.	
1398		
1399		
1400		
1401		
1402		
1403	Zehao Zhang, Jiaao Chen, and Diyi Yang. 2024e. Darg: Dynamic evaluation of large language models via adaptive reasoning graph. <i>arXiv preprint arXiv:2406.17271</i> .	
1404		
1405		
1406		
1407	Qihao Zhao, Yangyu Huang, Tengchao Lv, Lei Cui, Qinzhen Sun, Shaoguang Mao, Xin Zhang, Ying Xin, Qiufeng Yin, Scarlett Li, and Furu Wei. 2024. Mmlu-cf: A contamination-free multi-task language understanding benchmark. <i>Preprint</i> , arXiv:2412.15194.	
1408		
1409		
1410		
1411		
1412		
1413	Yang Zhou, Hongyi Liu, Zhuoming Chen, Yuandong Tian, and Beidi Chen. 2025. Gsm-infinite: How do your llms behave over infinitely increasing context length and reasoning complexity? <i>Preprint</i> , arXiv:2502.05252.	
1414		
1415		
1416		
1417		
1418	Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. 2024a. Dyval: Dynamic evaluation of large language models for reasoning tasks. <i>arXiv preprint arXiv:2309.17167</i> .	
1419		
1420		
1421		
	Kaijie Zhu, Jindong Wang, Qinlin Zhao, Ruochen Xu, and Xing Xie. 2024b. Dynamic evaluation of large language models by meta probing agents. <i>arXiv preprint arXiv:2402.14865</i> .	1422
		1423
		1424
		1425
	Kejian Zhu, Shangqing Tu, Zhuoran Jin, Lei Hou, Juanzi Li, and Jun Zhao. 2025. Establishing trustworthy llm evaluation via shortcut neuron analysis. <i>Preprint</i> , arXiv:2506.04142.	1426
		1427
		1428
		1429
	Qin Zhu, Qingyuan Cheng, Runyu Peng, Xiaonan Li, Tengxiao Liu, Ru Peng, Xipeng Qiu, and Xuanjing Huang. 2024c. Inference-time decontamination: Reusing leaked benchmarks for large language model evaluation. <i>arXiv preprint arXiv:2406.13990</i> .	1430
		1431
		1432
		1433
		1434
	Wenhong Zhu, Hongkun Hao, Zhiwei He, Yun-Ze Song, Jiao Yueyang, Yumeng Zhang, Hanxu Hu, Yiran Wei, Rui Wang, and Hongyuan Lu. 2024d. CLEAN-EVAL: Clean evaluation on contaminated large language models. In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 835–847.	1435
		1436
		1437
		1438
		1439
		1440
		1441
	A Prior Surveys	1442
	Prior research on data contamination primarily focuses on three main areas: definition, detection, and mitigation. Xu et al. (2024) and Deng et al. (2024) provide comprehensive surveys that thoroughly examine data contamination in large language models, covering conceptual definitions, detection methodologies, and mitigation strategies with similar classification frameworks for detection methods (matching-based and comparison-based). However, they differ significantly in their conceptualization of contamination types. The first paper primarily distinguishes between task-level contamination and language-level contamination, providing a function-oriented taxonomy. In contrast, the second paper presents a more granular severity-based hierarchy with four distinct levels: semantic-level (topical overlap), information-level (metadata and distributions), data-level (content without labels), and label-level contamination (complete exposure including ground truth). Their approaches to mitigation strategies also diverge notably. While the first paper emphasizes evaluation guidelines and procedural recommendations, the second paper offers a more structured framework categorized into three comprehensive strategies: data curation, data refactoring, and benchmark-free evaluation.	1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473

Survey	Definition	Detection	Mitigation	Core Content and Innovation
Deng et al. (2024)	●	●	●	Comprehensive strategies for detection and mitigation.
Ravaut et al. (2024)	○	●	○	Focuses on detection methods.
Xu et al. (2024)	●	●	●	Unified framework for detection and mitigation.
Fu et al. (2024)	○	◐	○	Detection methods assumptions.
Chen et al. (2025b)	○	○	◐	Overview of dynamic benchmarks.
Wu and Cao (2025)	○	●	○	Focuses on detection(MIA) methods.

Table 1: Summary of Recent Surveys on Data Contamination. ● indicates full coverage, ○ means the aspect is not the main focus, and ◐ refers to sub-domain coverage. More details are available in Appendix A.

reliability. Ravaut et al. (2024) investigates the critical issue of contamination in LLMs, categorizing it into data contamination and model contamination, while further distinguishing between input-only and input-label contamination scenarios. The authors systematically review state-of-the-art detection methods, including string matching, embedding similarity analysis, likelihood-based techniques, and novel LLM-driven approaches, highlighting their strengths and limitations. Chen et al. (2025b) conduct an in-depth analysis of existing static to dynamic benchmark aimed at reducing data contamination risks. Based on this, they propose a series of optimal design principles for dynamic benchmarking and analyze the limitations of existing dynamic benchmarks.

This survey (Wu and Cao, 2025) systematically reviews *data contamination* in large language models (LLMs), i.e., unintended train–test interactions that can inflate benchmark scores and misrepresent true generalization. It adopts a unified pipeline perspective to clarify definitions/taxonomies and contamination sources across pre-training, SFT, and preference training, and analyzes downstream impacts on evaluation reliability, model comparison, fairness, and safety claims. It further synthesizes detection (white-/gray-/black-box, including MIA-inspired signals) and mitigation strategies (deduplication/filtering, locked protocols to avoid prompt-on-test, and dynamic/contamination-labeled benchmarks), while outlining open challenges such as post-training and multi-modal leakage and the need for standardized reporting.

B Taxonomy

B.1 RL post-training contamination

Recent advancements in Large Language Models (LLMs) have increasingly relied on Reinforcement Learning (RL) post-training, particularly Reinforcement Learning with Verifiable Rewards (RLVR), to enhance reasoning capabilities. Conse-

quently, the RL phase has emerged as a critical, yet often overlooked, source of data contamination.

B.1.1 The Shift from Likelihood to Reward Maximization

The fundamental challenge in characterizing RL contamination stems from the shift in training objectives. Pre-training and SFT are governed by Maximum Likelihood Estimation (MLE), where the model minimizes the negative log-likelihood loss to maximize the probability of observed data sequences. This process naturally imprints strong likelihood-based signals, such as unusually low perplexity for memorized samples, which traditional detectors rely upon.

In stark contrast, RL post-training decouples the model’s optimization from simple likelihood metrics. The objective function in RL (e.g., GRPO (Shao et al., 2024)) maximizes the expected reward \mathcal{R} from generated outputs:

$$\mathcal{J}_{\text{RL}}(\theta) = \mathbb{E}_{q \sim D_{\text{RL}}, \{o_i\} \sim \pi_{\theta_{\text{old}}}} [f(\mathcal{R}(o_i), \pi_{\theta})] \quad (2)$$

Here, the optimization is driven by external, often sparse reward signals (e.g., 1 for a correct answer, 0 otherwise) rather than the token-level probabilities of a ground-truth response. As a result, contaminated samples in the RL phase do not necessarily exhibit the perplexity anomalies found in SFT contamination, rendering likelihood-based detectors ineffective.

B.1.2 Policy Collapse and Path Dependence

Instead of verbatim memorization, RL contamination manifests as *policy collapse* (Yue et al., 2025). When a model is trained on a benchmark sample during the RL phase, the policy converges to a narrow, high-reward reasoning path to maximize the metric (e.g., pass@1 accuracy).

This phenomenon creates a distinct signature in the model’s output entropy. For contaminated samples, the token-level entropy distribution becomes highly sparse, indicating that the model has

LLM Data Contamination Types

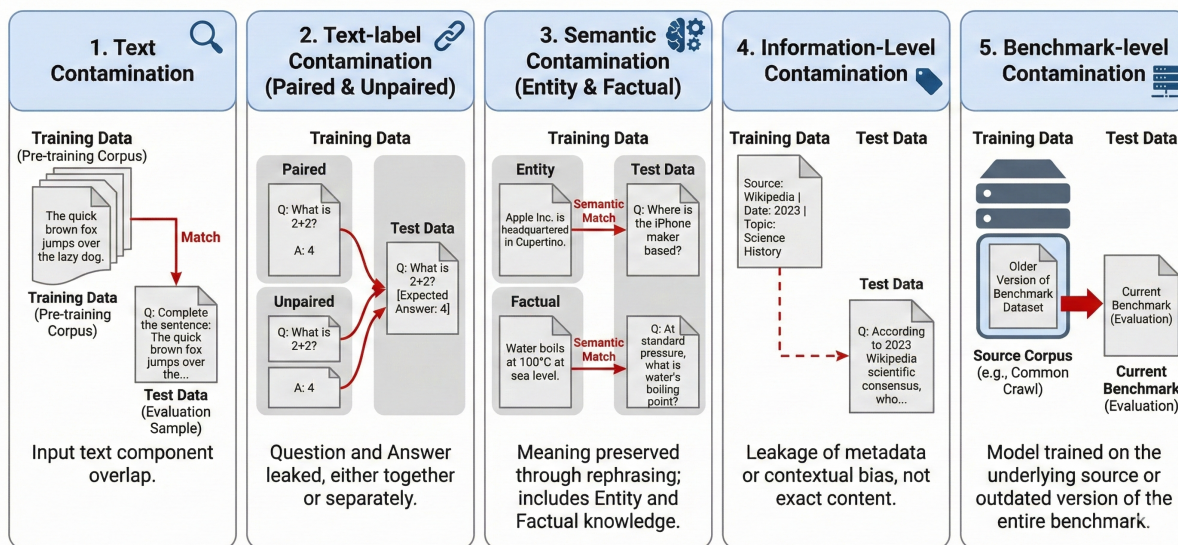


Figure 2: contamination granularity

become nearly deterministic along a specific trajectory. While general RL training can induce sparsity, contaminated samples exhibit a rigid *path dependence*: they struggle to deviate from the memorized reasoning trace even when explicitly instructed to provide an alternative solution. This rigidity serves as the primary signal for distinguishing RL-contaminated data from clean data.

B.2 Retrieval-time (Tool-use) Contamination

Han et al. (2025) identify a novel leakage modality termed *Search-Time Data Contamination* (STC), which specifically undermines the evaluation of search-based LLM agents (Xi et al., 2025). Unlike traditional contamination where test data leaks into the training corpus, STC occurs during the inference retrieval step when an agent discovers the evaluation dataset itself (e.g., hosted on public platforms like HuggingFace) containing ground-truth labels. This allows the agent to bypass reasoning by directly extracting answers from the retrieved context. Empirical analysis across benchmarks such as Humanity’s Last Exam (HLE) (Phan et al., 2025) and SimpleQA (Wei et al., 2024) reveals that agents successfully retrieve the answer key for approximately 3% of queries. This leakage artificially inflates performance; blocking access to the contaminated sources results in an accuracy drop of approximately 15% on the affected subsets, necessitating new guardrails such as source filtering for

reliable agent evaluation.

B.3 Text contamination vs. Text-label contamination

(Li et al., 2024e) demonstrated through controlled experiments that while ground-truth contamination significantly inflates performance on benchmarks such as ARC and HellaSwag, text contamination often yields negligible or even detrimental effects. This distinction is further corroborated by findings in machine translation, where models exposed to ground-truth labels exhibited drastic score inflation (e.g., +30 BLEU) compared to those exposed only to source texts, confirming that the leakage of ground-truth labels is the primary driver of artificial performance gains rather than mere exposure to evaluation text.

B.4 Impact Analysis of Information-Level Contamination

Experiments reveal that Information-Level contamination—exposure to metadata or label distributions without instance leakage—can significantly inflate performance via statistical priors. For instance, InstructLM-1.3B (Cheng et al., 2024)’s accuracy on SST-2 surged from 29.17% to 41.14% solely due to such exposure, suggesting models exploit class imbalances to "guess" answers. This contamination is often implicit and stealthy; models (e.g., Qwen2.5) may internalize metadata

1610 from academic papers or online discussions dur- 1659
1611 ing pre-training. However, sensitivity varies by 1660
1612 domain: while classification tasks are highly vul- 1661
1613 nerable to such biases, reasoning benchmarks like 1662
1614 GSM8K (Cobbe et al., 2021) remain resilient. Con- 1663
1615 sequently, the DCR framework explicitly penalizes 1664
1616 this risk in Adjusted Accuracy (Acc_{adj}) to distin- 1665
1617 guish genuine zero-shot reasoning from reliance on 1666
1618 prior meta-knowledge. 1667

1619 C Impacts 1668

1620 **More evidence of data contamination** Beyond 1669
1621 these primary signals, other studies have high- 1670
1622 lighted diverse manifestations of leakage across dif- 1671
1623 ferent scopes. Li and Flanigan (2024) utilized mem- 1672
1624 bership inference as a detection mechanism, while 1673
1625 Liu et al. (2024) observed discrepancies in Chi- 1674
1626 nese LLMs where high benchmark scores masked 1675
1627 superficial knowledge rather than genuine compre- 1676
1628 hension. Furthermore, contamination risks have 1677
1629 been identified in specialized domains beyond stan- 1678
1630 dard text and software code: Tseng et al. (2025) 1679
1631 revealed that transcript leakage significantly dis- 1680
1632 torts metrics in the audio domain, and Wang et al. 1681
1633 (2025b) reported similar severity in hardware de- 1682
1634 scription languages like Verilog, emphasizing the 1683
1635 need for rigorous scrutiny across all modalities. 1684

1636 D Detection 1685

1637 D.1 Definition of Assumptions 1686

1638 D.1.1 Verbatim Memorization 1687

1639 In the context of LLMs, verbatim memoriza- 1688
1640 tion (Carlini et al., 2021b, 2022b) refers to the phe- 1689
1641 nomenon where a model recalls exact sequences of 1690
1642 text, often from the data it has been trained on. This 1691
1643 occurs when a model has seen a specific passage 1692
1644 or piece of information during its training process 1693
1645 and is able to reproduce it exactly when prompted. 1694
1646 Verbatim memorization can lead to issues of data 1695
1647 contamination, where the model unintentionally 1696
1648 outputs copyrighted or sensitive material verba- 1697
1649 tim, causing concerns regarding privacy, intellec- 1698
1650 tual property, and validity in analytical tasks. 1699

1651 Instance-level contamination (Fu et al., 2024) 1700
1652 does not always lead to verbatim memorization. 1701
1653 Utilizing instance generation (Carlini et al., 2022b; 1702
1654 Karamolegkou et al., 2023), demonstrates that ver- 1703
1655 batim memorization requires repeated exposures 1704
1656 to this instance x during training. Indeed, future 1705
1657 research on contamination should place more em- 1706
1658 phasis on LLMs’ memorization. Schwarzschild 1707

1659 et al. (2024) proposed that strings can be consid- 1660
1661 ered memorized if they can be reproduced using a 1661
1662 shorter prompt, while Karamolegkou et al. (2023) 1662
1663 investigated verbatim memorization, particularly 1663
1664 in the context of copyrighted materials. 1664

1664 D.1.2 Black-Box Method Assumption 1664

1665 Golchin and Surdeanu (2023a) has assumed that 1665
1666 when a model has memorized instances from the 1666
1667 original dataset, it will prefer selecting options con- 1667
1668 taining the original instance over semantically sim- 1668
1669 ilar perturbations. Additionally, LLMs may ex- 1669
1670 hibit positional biases, where certain positions in 1670
1671 multiple-choice options are more likely to be cho- 1671
1672 sen, leading to potential overestimation or underes- 1672
1673 timation of contamination levels. 1673

1674 Golchin and Surdeanu (2023b) gave the assump- 1674
1675 tion that by providing a "guided instruction" with 1675
1676 dataset name, partition information, and part of the 1676
1677 reference instance, LLMs can generate the com- 1677
1678 plete version of the data instance. This allows for 1678
1679 calculating overlap between generated completions 1679
1680 and reference instances, helping to infer whether 1680
1681 the dataset partition is contaminated. 1681

1682 Duarte et al. (2024) assumed that LLMs may 1682
1683 memorize specific copyrighted content, such as 1683
1684 books or academic papers, during training. When 1684
1685 encountering similar content, they can distinguish 1685
1686 whether they’ve seen it before. DE-COP exploits 1686
1687 this by designing multiple-choice questions to test 1687
1688 if the model can accurately identify original copy- 1688
1689 righted content from paraphrased versions. Addi- 1689
1690 tionally, model selection biases can affect copyright 1690
1691 detection results, and DE-COP introduces a cali- 1691
1692 bration method to minimize such biases. 1692

1693 In (Dong et al., 2024), it is assumed that contam- 1693
1694 inated training data significantly affects the output 1694
1695 distribution of large language models. Specifically, 1695
1696 when trained on contaminated data, the model’s 1696
1697 output distribution becomes more peaked, causing 1697
1698 it to produce more consistent outputs on contam- 1698
1699 inated data, favoring outputs strongly correlated 1699
1700 with the training data. 1700

1701 Deng et al. (2023) assumed that if an LLM can 1701
1702 accurately guess missing parts of a test set, such as 1702
1703 keywords or answer options, without external assis- 1703
1704 tance, it suggests that the model has encountered 1704
1705 the corresponding benchmark data during training. 1705
1706 This indicates memorization-based contamination. 1706
1707 The TS-Guessing protocol tests whether the model 1707
1708 has memorized benchmark data by having it guess 1708
1709 hidden information. 1709

Literature	Contamination Granularity Investigated				
	Text	Text and Label	Semantic	Information	Benchmark
Dekoninck et al. (2024b)		✓			✓
Wu and Cao (2025)		✓			
Chen et al. (2025b)		✓	✓		
Palavalli et al. (2024)		✓	✓		✓
Jiang et al. (2024)	✓	✓			
Li et al. (2024e)	✓	✓			
Xu et al. (2024)	✓	✓	✓	✓	
Xu et al. (2025a)	✓	✓	✓	✓	
Xu et al. (2025b)		✓	✓		
Fu et al. (2024)		✓	✓		✓

Table 2: Comparison of contamination types covered in recent studies. The types correspond to: Text (Input-only contamination), Text+Label (Ground-truth contamination), Semantic (Paraphrased or variation-based leakage), Information (Information-level or metadata leakage), and Benchmark (Specific test set structure or membership inference).

Ranaldi et al. (2024) assumed that data contamination can be detected solely by analyzing the inputs and outputs of LLMs. For example, unusually high accuracy on tasks from datasets like Spider indicates that the model may have been exposed to this dataset during training, leading to memorization rather than genuine understanding. Additionally, data contamination may lead to inflated performance on zero-shot tasks when the model encounters potentially contaminated data during training.

Chang et al. (2023) assumed that LLMs may memorize portions of text from their training data, especially when evaluation datasets contain known texts. This memorization can lead to inflated performance on tasks such as code generation. Moreover, data repetition on the web—through search engines and open datasets—encourages memorization, which improves accuracy on tasks involving familiar content.

D.2 Additional Methods

In the domain of code benchmarks, Riddell et al. (2024) assess contamination through both surface-level and semantic-level matching. Li et al. (2024e) employ a two-stage detection pipeline: (1) retrieving potential verbatim samples via Bing Search and Common Crawl, and (2) quantifying textual overlap using the METEOR recall score with a threshold of $\tau > 0.75$. More recently, Ni et al. (2025) proposed a gray-box method based on *option shuffling*. For any given n -choice multiple-choice item, they generate all $n!$ permutations of

the answer options to obtain a set of derived instances. They then compute the corresponding log-likelihood scores, $P = \{\log p_1, \dots, \log p_{n!}\}$. By identifying statistically significant outliers—either through a maximum-score check or an Isolation Forest—they interpret these anomalies as evidence of benchmark leakage in the LLM’s pre-training data. ConStat (Dekoninck et al., 2024a) reframes contamination as non-generalizable score inflation on a target benchmark and detects (and quantifies) it by statistically comparing performance on the target benchmark versus a reference benchmark under calibration with a set of reference models.

D.3 Related Tools

D.3.1 Data Contamination Detector

Li et al. (2024e) present Contamination Detector to check whether test examples appear on the internet via Bing search and Common Crawl index. The tool is available at: https://github.com/liyucheng09/Contamination_Detector.

Ravaut et al. (2024) presented an open-source library for contamination detection in NLP datasets and LLMs. The library combines multiple methods for contamination detection and is available at: <https://github.com/ntunlp/LLMSanitize>.

Overlap is a Python package developed to evaluate textual overlap (N-Grams) between two volumes of text. This tool can be accessed at: <https://github.com/nlx-group/overlap>.

Yao et al. (2024) introduced Deep Contam, a method that detects cross-lingual contamina-

Paradigm	Access needed	Representative methods	Core signal / idea	Strengths / caveats
White-box	Training corpus or internals	N-gram overlap Embedding similarity Layer/activation probes	Directly measure train–test overlap (lexical/semantic) or train classifiers on “sensitive” internal representations to infer membership.	High precision for <i>developers</i> who can audit data/internals; may miss heavily transformed leakage without strong semantic tooling; not usable for closed models.
Gray-box	Token probabilities (no weights)	Min-K%, Min-K%++ PAC (perturbation-based) DC-PDD (divergence) PaCoST (paraphrase) Perplexity	Membership inference from probability “outliers” (often via thresholds), sometimes enhanced by perturbations, frequency divergence, or paraphrase consistency tests.	Often efficient (token-based methods need only a constant number of forward passes + light algebra); however sensitivity to k /thresholds and transformation robustness can be an issue; paraphrasing adds compute.
Black-box	API outputs (No token probabilities)	Guided completion DCQ (perturbations) TS-Guessing DE-COP (copyright) CDD (distribution) Canary insertion Natural-DaCoDe	Probe memorization via masked-span reconstruction / completion overlap or MCQ selection of the original instance; estimate contamination from output distribution characteristics; use signals to measure recall	Most practical for <i>closed</i> models; but relies on heuristic assumptions and prompt/task design, can be brittle across settings and may be affected by safety filters;

Table 3: Summary of contamination detection paradigms (Section 6): white-/gray-/black-box methods differ by required access, detection signals, and robustness/efficiency trade-offs.

methods demonstrate the practical efficacy of paraphrasing techniques in contamination mitigation. Preventive measures involve technical defenses like encryption, access control, and de-contamination during inference to guarantee the reliability and fairness of LLM evaluation.

E.1.1 Data Updating-based Methods

Using the most recent data is intuitive for constructing contamination-free benchmarks, and some studies have proposed automatically collecting recent data to build questions. Meanwhile, recent data also need to maintain the stability of difficulty. LatestEval proposed an automated pipeline to dynamically generate contamination-free test sets from recent materials (Li et al., 2024d). White et al. (2024) introduced LiveBench, a dynamically updated benchmark that integrates tasks across math, coding, and reasoning with automated scoring to mitigate data contamination. Similarly, Jain et al. (2024) introduced LiveCodeBench, a code-generation benchmark that extends prior methodologies by dynamically evaluating self-repair capabilities and maintaining update cycles. Fan et al. (2024) introduced NPHardEval4V-a dynamically updated benchmark to assess reasoning capabilities of MLLMs. In code evaluation, EvoCodeBench (Li et al., 2024a) was proposed to dynamically align

with recent code repositories for fair evaluation. Wu et al. (2025b) proposed AntiLeak-Bench, an automated anti-leakage benchmarking framework that constructs samples with explicitly new knowledge not in LLMs’ training sets and features a fully automated workflow to build and update the benchmark. Zhang et al. (2025) presents SWE-bench-Live, a live-updatable benchmark for evaluating LLMs in issue-resolving tasks, an automated curation pipeline to enable continuous updates. Mahdavi et al. (2025) presents an automated pipeline to create the AoPS-Instruct dataset, and introduces LiveAoPSBench, an evolving contamination-resistant evaluation set with timestamps.

E.1.2 Data Rewriting-based Methods

This type of methods use data augmentation to remove contamination from benchmarks, with LLMs’ superior rephrasing and verifying capabilities. Zhu et al. (2024d) proposed Clean-Eval to purify contaminated benchmarks by paraphrasing and back-translating data into semantically equivalent but lexically distinct forms. Zhao et al. (2024) proposed the MMLU-CF dataset, which is constructed by collecting diverse questions, cleaning data, sampling difficulty reasonably, checking data integrity with LLMs, and applying rewriting methods such as rephrasing questions and shuffling op-

Method	Thresholds	Perturbation	Efficiency	Low-Leakage
Perplexity (Carlini et al., 2022a)	High (Requires baseline)	No	High (One forward pass)	Low Degrades when leakage is rare.
Reference Models (Li, 2023; Wei et al., 2023)	Medium (Comparison based)	No	Medium (Requires ref model)	Medium Improved calibration over raw PPL.
Min-K% (Shi et al., 2024)	High (Sensitive to k)	No	High (Cheaper than generation)	Low Degrades when leakage is rare.
Min-K%++ (Zhang et al., 2024c)	High	No	High	Medium Improved accuracy over Min-K%.
PAC (Ye et al., 2024)	Medium	Yes (Input perturbation)	Medium	Medium Robustness via perturbation analysis.
DC-PDD (Zhang et al., 2024d)	Medium	No	High	Medium Focuses on calibration.
PaCoST (Zhang et al., 2024b)	None (Threshold-free)	Yes (Paraphrasing)	Low (Computationally expensive)	High Robust(threshold-free design).

Table 4: Summary of Gray-box Data Contamination Detection Methods. We compare methods based on their reliance on thresholds, use of perturbations, computational efficiency, and performance on low-leakage data.

tions to ensure the dataset remains contamination-free. CLEVA (Li et al., 2023) employed non-repetitive sampling and multi-strategy data rewriting for robust evaluation. Ying et al. (2024) updated benchmarks with two strategies: style-preserving mimicry with LLMs and cognitive-level expansion using Bloom’s taxonomy. Mirzadeh et al. (2025) introduced GSM-Symbolic, an improved benchmark generated from symbolic templates revealing that their mathematical reasoning performance declines when numerical values or clause numbers in questions are altered, indicating they may replicate training data rather than perform reasoning.

E.1.3 Prevention-based Methods

Preventive measures focus on safeguarding test data integrity through technical and procedural controls. Core strategies include encrypting public test data with public-key cryptography, enforcing strict access permissions, and prohibiting derivative data creation. Zhu et al. (2024c) introduced Inference-Time Decontamination (ITD), a novel technique that identifies and rewrites potentially memorized responses during model inference. Li et al. (2024c) introduced C²LEVA, a comprehensive bilingual benchmark with systematic contamination prevention mechanisms, which implements proactive measures such as test data rotation.

E.2 Dynamic Evaluation

Rule-based Dynamic approaches address data contamination by leveraging adaptive assessment frameworks. Zhu et al. (2024a) introduced DYVAL, a graph-based system that generates evaluation samples through algorithmic composition, constraint application, and functional descriptions. Its directed acyclic graph (DAG) architecture facilitates multi-step reasoning tasks with precisely controlled complexity. Lei et al. (2024) developed S3EVAL, a framework for SQL evaluation that utilizes randomized table-query pairs. This synthetic approach allows for customizable task lengths and difficulty levels, while systematically assessing long-context reasoning capabilities. Zhang et al. (2024e) proposed the DARG method, which dynamically generates evaluation samples with adjustable complexity and diversity using adaptive reasoning graphs. Srivastava et al. (2024) introduced functionalization, a technique that transforms static question-answer pairs into parameterized code, enabling the generation of infinite test variants. Qian et al. (2024) further extended dynamic evaluation by perturbing key variables in questions, allowing for the dynamic generation of datasets with controlled variations.

Agent-based Zhu et al. (2024b) proposed Multi-Principle Assessment (MPA), which utilizes LLM-based agents to automatically transform questions into new ones. Wang et al. (2024) introduced a multi-agent framework to implement self-evolving benchmarks, which dynamically mutates question contexts and structures to update benchmarks. Chen et al. (2025c) proposes DyCodeEval, a novel benchmarking suite that uses multiple agents to generate semantically equivalent variations of seed programming problems without changing core logic, dynamically evaluating Code LLMs under potential data contamination risks across diverse problem sets to ensure reliable reasoning capability assessments. Cao et al. (2024) proposed StructEval that conducts evaluations across multiple cognitive levels and critical concepts starting from an atomic test objective, offering comprehensive, robust, and consistent LLM evaluations that resist data contamination risks. Xu et al. (2025d) introduced UGMATHBench, a dynamic benchmark for evaluating undergraduate-level mathematical reasoning in LLMs and proposes metrics effective accuracy (EAcc) and reasoning gap, revealing through evaluation of 23 LLMs. Zhou et al. (2025) developed a grade school math problem generator to create the GSM-Infinite benchmark, which enables fine-grained control over problem difficulty and context length.

E.3 LLM-as-a-Judge

Next-generation evaluation leverages LLMs themselves as assessment tools. They can serve the roles of scoring, ranking, and selection. Bai et al. (2024) presented the "LM-as-Examiner" framework, generating questions and evaluating responses through reference-free analysis. Yu et al. (2024) deployed LLMs as "Interactors" in structured multi-turn dialogues that probe model capabilities while minimizing contamination risks. Li et al. (2024b) proposed TreeEval, a benchmark-free system where LLMs generate hierarchical question trees. This adaptive approach adjusts difficulty based on model performance, creating unique assessment paths that prevent data contamination.

But Li et al. (2025) identified systematic bias in LLM-as-a-judge evaluations, where models trained on synthetic data from architecturally similar foundations receive unfair preference, compromising evaluation fairness.

E.4 Trustworthy Evaluation Pipeline

Contamination mitigation evaluation aims to estimate an LLM's *non-memorized* capability on potentially leaked benchmarks without the need to reconstruct "clean" datasets. To this end, Hou et al. (2025) propose LNE-BLOCKING, a two-stage *detection-disruption* framework. It first uses length-normalized entropy (LNE) via greedy decoding to estimate per-instance contamination strength. Subsequently, it adaptively applies a token-level *blocking* operation—suppressing top-probability tokens when contamination signals are strong—to steer generation away from memorized paths. This approach recovers contamination-aware scores more efficiently and stably than sampling-heavy filters such as TED.

Yadoji (2025) introduce *DC-Guard*, a unified data-centric framework that reframes benchmarking as an ecological audit. By jointly assessing memorization (MCI), benchmark representativeness (BES), and contamination-aware score correction (CRMA), DC-Guard aims to establish more transparent, reproducible, and trustworthy LLM evaluations.

Disentangling Contamination from Generalization.

A fundamental challenge in trustworthy evaluation lies in resolving the ambiguity between rote memorization and genuine reasoning. (Ishikawa, 2025) systematically address this by proposing a three-tier framework combining n-gram alignment, canary insertion, and perturbation testing. Crucially, they posit that performance on out-of-distribution (OOD) data serves as the true litmus test for generalization. Future evaluation protocols must therefore prioritize strict OOD validation to rigorously distinguish latent reasoning capabilities from surface-level data leakage.