LNE-BASED BLOCKING GENERATION AGAINST DATA CONTAMINATION ON LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Data contamination gradually becomes inevitable during the development of large language models (LLMs), meaning the training data commonly integrates those evaluation benchmarks unintentionally. This subsequently makes it hard to benchmark LLMs fairly. This paper introduces a novel approach called LNE-Blocking for contamination detection and contamination mitigation evaluation. For the first component, Length Normalized Entropy (LNE) reports a SOTA performance for contamination detection. On this basis, LNE-Blocking reports a SOTA performance for contamination mitigation evaluation by applying LNE to adjust the intensity of blocking operation, specialized to suppressing the maximum value of output candidates during the generation process. We conduct extensive experiments on both contamination detection and contamination mitigation evaluation tasks. The results indicate that LNE and LNE-Blocking achieve an obvious SOTA performance. Simultaneously, LNE-Blocking is robust across tasks and models of different contamination level, reducing computational costs by nearly 25x compared to previous methods. We hope our method will open new research avenues on data contamination for LLMs. We plan to release the resources upon publication of this work to facilitate future work.

027 028 029

025

026

004

006

007 008 009

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

In the era of fierce development with large language models (LLMs), it has been a popular research
topic across many areas such as chain-of-thought reasoning (Wang et al., 2023; Wei et al., 2024),
machine translation (Lu et al., 2023; Zhu et al., 2024a), code generation (Li et al., 2023a; Zhang
et al., 2023), and even spatial reasoning (Hu et al., 2024). Despite the fact that LLMs are usually
strong on many tasks, Natural Language Processing (NLP) practitioners secretly face a common
problem called data contamination when conducting NLP research and engineering which frequently
relies on benchmark data that could be potentially contaminated in LLMs.

Data contamination, also called data leakage, occurs when the test data is inadvertently included in the model's training data (Magar & Schwartz, 2022; Golchin & Surdeanu, 2024). This causes the model to perform exceptionally well on the leaked test data. Due to the immense scale and diverse origins of the pre-trained datasets used for LLMs, they are more vulnerable to data contamination. This can be divided into two main scenarios: 1) Existing benchmark datasets are more prone to leakage because of extensive text quotations and code reuse present in the LLMs' training data. 2) For emerging benchmark datasets, newly created test data may already be included in the constantly expanding training data of LLMs, as the details of these training datasets are often unknown.

As a result, preventing benchmark data contamination in LLMs becomes highly challenging. This prevents NLP developers and researchers from honestly judging the LLMs. With the growth of the training set which is partially synthetic or automatically crawled from the internet, the problem of data contamination is gradually unavoidable, even when the LLM developers do not intentionally do it. Hence, our research question is particularly important: how can we accurately detect data contamination and how can we accurately assess model performance even when it is contaminated?

To remediate the problem of data contamination, previous research mainly focused on how to detect whether the data is leaked in the models rigorously (Deng et al., 2024; Shi et al., 2024; Elazar et al., 2024). One less studied yet important problem is to quantitatively assess the genuine performance of the models, even if they are contaminated (Zhu et al., 2024b; Dong et al., 2024; Ye et al., 2024). Such research, knowned as contamination mitigation evaluation, could be quite useful as LLMs are usually trained with large-scale data, and many LLMs are potentially contaminated even though some of them carefully handled this problem. How to fairly compare the models is therefore necessary to decide which models to use.

This paper proposed a novel approach for contamination detection and contamination mitigation 060 evaluation, which is composed of two novel components. The first component is Length Normalized 061 Entropy (LNE) for contamination detection. We are the first to apply this method to detect data con-062 tamination and surprisingly, we found it better than previous methods. On this basis, LNE-Blocking, 063 applying LNE to adjust the intensity of blocking operation which specializes to suppressing the max-064 imum value of output candidates during the generation process, reports a SOTA performance for contamination mitigation evaluation. Additionally, Since our method requires only two inferences 065 for contamination mitigation evaluation, compared to 50 samplings in existing methods, TED(Dong 066 et al., 2024), it is 25 times faster. Additionally, experiments show that our approach is highly robust 067 across different tasks and models with varying levels of contamination. 068

- To this end, we make three key contributions:
 - We propose a novel approach for contamination detection and fairly assessing possibly contaminated models.
 - Extensive experiments conducted across two tasks and multiple LLMs have yielded stateof-the-art (SOTA) results, providing strong evidence for the effectiveness of our proposed method.
 - Our proposed SOTA method is 25x faster and more robust than the previous method.

2 MOTIVATION

078 079

071

072 073

074

075

076 077

Contaminated models often exhibit a high lexical overlap between their output and the ground truth, which is indicative of memory phenomena (Magar & Schwartz, 2022). As shown in Figure 1, as the degree of contamination increases, the overlap between the output generated by greedy decoding and the ground truth significantly rises. This behaviour reflects the model's tendency to memorize the training data rather than generalizing it. One solution (Dong et al., 2024) to assess the true performance of such contaminated models involves generating diverse outputs through multiple sampling and filtering out any instances of the standard ground truth. Then, they expect to derive answers that stem from the model's generalization abilities, rather than its memorized knowledge.

However, answers based on memory phenomena tend to have a very high likelihood within the
 model, meaning that obtaining non-memorized, generalized answers requires a large number of samplings. Previous work (Dong et al., 2024) has found that at least 50 samples are necessary to achieve
 satisfactory performance estimates. This process is both highly random and time-consuming.

- The randomness makes it difficult to consistently generate non-memorized answers, particularly for models with heavy contamination, where such outputs are rarely sampled. This often leads to unreliable performance estimates, as we demonstrate in subsequent experiments in section 6.2.2.
- Similar to how humans can rephrase their thoughts when interrupted, LLMs have the potential to generate alternative answers when their default response is blocked (Lu & Lam, 2023; Wang & Zhou, 2024). By leveraging this ability, we can apply a strategy that suppresses the generation of the token with the highest probability, referred to as blocking, without significantly impacting performance. This approach reduces the reliance on memorized content and encourages the model to draw on its generalization capabilities, as demonstrated by the reduction in ROUGE-L (Lin, 2004) similarity between the model's output and the memorized answers, as in Figure 1. Then, blocking enables the model to produce more diverse and generalized outputs, without sacrificing quality.
- For example, the output of the contaminated model is *return len(set(string.lower()))*. If the model is good at coding, after applying the blocking strategy, the model could generate an equivalent piece of code using its generalization ability, such as *return len(character.lower()) for character in string)*.
- Additionally, as shown in Figure 1, models with varying degrees of contamination exhibit different levels of memorization, and the impact of the blocking operation varies accordingly. To effectively



Figure 1: An example demonstrating the changes in Length-Normalized Entropy (LNE), and the impact of the blocking strategy, on model memory phenomena (ROUGE-L), in models with varying levels of contamination, where "Cont." is an abbreviation of contamination.

disrupt memorized responses, models with differing contamination levels require distinct degrees of
 blocking intensity. In particular, as the level of contamination increases, the blocking intensity must
 be intensified to counteract the effects of memory.

132 A natural idea, then, is to first detect the degree of contamination to better tailor the blocking strategy 133 to ensure it disrupts memorization without negatively affecting performance. This approach allows for a more targeted and adaptive application of blocking across different contamination levels. Once 134 the degree of contamination is detected, this information can be used to adjust the blocking intensity. 135 Specifically, we hypothesize that a heavily contaminated model will exhibit greater certainty in its 136 token predictions, resulting in generated text that is closer to the ground truth. This leads to lower 137 entropy in the probability distribution at each position, as the model becomes more confident in 138 generating memorized content. Based on this hypothesis, we propose using Length Normalized 139 Entropy (LNE) as a novel method to detect contamination. As shown in Figure 1, as the degree of 140 contamination increases, the generated text becomes closer to the ground truth, resulting in higher 141 lexical overlap and a corresponding increase in LNE. Our experiments demonstrate that this entropy-142 based detection method outperforms existing strategies such as perplexity in Table 4, thus we use 143 LNE to adjust the intensity of the blocking.

144 145

146

124

125

126

127 128

3 RELATED WORK

147 Data Contamination Detection The issue of data contamination in large language models 148 (LLMs) gained attention in the context of GPT-3 Brown (2020), where the vast pre-training cor-149 pus inevitably overlapped with evaluation benchmarks. To address this, GPT-3 employed n-gram 150 overlap detection to filter out data in the training set that conflicted with test set benchmarks. Some 151 work (Pan et al., 2020; Zhou et al., 2023; Jacovi et al., 2023; Dodge et al., 2021) exposed the serious 152 consequences of data contamination and urged attention to this problem. However, most LLMs released to date have not made their pre-training corpus publicly available, introducing new challenges 153 for contamination detection. 154

Recent efforts have attempted to detect contamination without direct access to the pre-training data. Min-k% Prob (Shi et al., 2023) calculates the average of the smallest k% probabilities of generated tokens and flags potential contamination if this average exceeds a certain threshold. Similarly, perplexity(Li, 2023) is also used to detect contamination, assuming that leaked data tends to produce lower perplexity scores. CDD (Dong et al., 2024) detects contamination by sampling multiple times to determine the proportion of highly similar samples in the dataset. The methods mentioned above do not fully leverage all the distributional information obtained from a single inference. They merely select the maximum probability at each output position to compute the score. In contrast, our method, LNE, utilizes entropy to exploit more information, thereby achieving better contamination
 detection performance.

Contamination Mitigation Evaluation To evaluate LLMs in the context of potential data contamination, several methodologies generate new datasets that do not overlap with the model's training data, adopting a dataset-centric perspective. GSM-Plus(Li et al., 2024b) ensures that benchmark data is absent from the model's training set by reconstructing the original GSM8k dataset (Cobbe et al., 2021) through the introduction of perturbations. GSM1k(Zhang et al., 2024) involves the creation of a completely new dataset from scratch. It remains private and is only made publicly available at a future point in time. Furthermore, TreeEval(Li et al., 2024c) assesses model performance by dynamically querying the language model, using more powerful LLMs for the querying process.

173 One less explored but important problem is how to quantitatively assess the performance of models using datasets that have already been leaked. CleanEval(Zhu et al., 2023) proposed paraphras-174 ing contaminated datasets using LLMs for evaluation purposes. TED (Dong et al., 2024) filters 175 non-memorized samples through multiple sampling rounds, using these samples for contamination 176 mitigation evaluation. However, these methods often require extensive preprocessing or multiple 177 sampling iterations, making them highly time-consuming. LNE-Blocking is the first strategy that 178 evaluates the performance of contaminated models using two inferences, without modifying the 179 original dataset. 180

181 182

183

184 185

186

187

188

193

194

201

202 203

204 205

206

209

4 METHODLOGY

4.1 LNE FOR DATA CONTAMINATION DETECTION

Given a test data x, and its corresponding greedy decoded outout, y, we aim to detect whether this data was used to train the model \mathcal{M} . We calculate the length normalized entropy (LNE) using the probability distribution at each position of the output. For a query x, with an output of length N, $\mathbf{y} = (y_1, y_2, \dots, y_N)$, the LNE is defined as:

LNE(y) =
$$\frac{1}{N} \sum_{i=1}^{N} H(y_i) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{V} p(y_i = j) \log p(y_i = j).$$
 (1)

where $H(y_i)$ represents the entropy at the *i*-th position, $p(y_i = j)$ denotes the probability of the model generating the *j*-th token from the vocabulary Vocab at the *i*-th position, with V representing the size of the token vocabulary.

Contamination on the test data can be detected by identifying the LNE score:

$$Is_Contaminated(\mathcal{M}; x, y) = \begin{cases} False & \text{if } LNE(y) > \xi \\ True & \text{if } LNE(y) \le \xi \end{cases}$$
(2)

where $\xi \in [0, log(V)]$, is a hyper-parameter that controls the threshold, (x, y) represents the prompt and response output corresponding to the model \mathcal{M} .

4.2 BLOCKING OPERATION FOR LANGUAGE MODEL GENERATION

For disrupting memorization, we propose a blocking operation during the decoding process of LMs, which suppresses the token with the highest probability at a certain position during decoding.

Specifically, for a model \mathcal{M} and a prompt x, the generation process using greedy decoding is:

$$gene(\mathcal{M}, x) = (y_1^{greedy}, y_2^{greedy}, ..., y_n^{greedy})$$
(3)

At each position, the token with the highest probability is selected. This process is denoted as: $a_{greedy} = a_{greedy} = a_{greedy} = (x_{greedy})$

$$y_i^{greedy} = argmax(\mathcal{M}_{greedy}(x, y_{< i})) \tag{4}$$

where $\mathcal{M}_{greedy}(x, y_{\langle i \rangle}) \in \mathbb{R}^V$, represents the probability distribution of the *i*-th position predicted by the model through greedy decoding. When the blocking strategy is applied to the *i*-th position in a response of length *l*, the generation process is defined as:

$$block_gene(\mathcal{M}, x, (i)) = (y_1^{greedy}, ..., y_{i-1}^{greedy}, y_i^{block}, y_{i+1}^{greedy}, ..., y_l^{greedy})$$
(5)

222

232

237 238

239

244 245

249

250

255

256 257

258

During the process, each position before the *i*-th applies greedy decoding using equation 4. At the
 i-th position, the token with the highest probability from the distribution modified by the blocking
 operation is selected as the token for that position, as:

$$y_i^{block} = argmax(\mathcal{M}_{block}(x, y_{< i})) \tag{6}$$

To obtain $\mathcal{M}_{block}(x, y_{< i})$, the probability distribution with the maximum value suppressed, the process involves first identifying the index of the maximum value and then suppressing it:

$$\mathcal{M}_{greedy}[argmax(\mathcal{M}_{greedy}(x, y_{< i}))] = -\infty$$

$$\mathcal{M}_{block}(x, y_{< i}) := \mathcal{M}_{greedy}(x, y_{< i})$$
(7)

After the generation of the *i*-th token, the subsequent tokens are generated using greedy decoding:

$$y_{i+1}^{greedy} = argmax(\mathcal{M}_{greedy}(x, y_{\leq i})) = argmax(\mathcal{M}_{greedy}(x, (..., y_{i-1}^{greedy}, y_i^{block})))$$
(8)

4.3 LNE-BASED BLOCKING FOR CONTAMINATION MITIGATION EVALUATION

For data with different levels of contamination, varying intensities of disruptions are required during the generation process to disrupt the model's memorization, reflecting the model's original capabilities. It is necessary to control the disruption intensity during data generation according to different levels of contamination. This section will sequentially introduce how to control the blocking intensity and how to determine the blocking intensity based on the level of contamination.

4.3.1 MULTI BLOCKING OPERATIONS FOR DISRUPTING MEMORIZATION

For highly contaminated data, only performing a blocking operation during generation is not sufficient to disrupt its memorization, necessitating an increase in the number of blocking operations to increase the disruption intensity. Starting from the first token, the blocking operation is applied ntimes as defined below:

$$block_gene(\mathcal{M}, x, (1, 2, \dots n)) = (y_1^{block}, y_2^{block}, \dots, y_n^{block}, y_{n+1}^{greedy}, y_{n+2}^{greedy}, \dots, y_l^{greedy})$$
(9)

where $y_n^{block} = argmax(\mathcal{M}_{block}(x, (y_1^{block}, y_2^{block}, ..., y_{n-1}^{block})))$. After y_n , normal generation proceeds, where $y_{n+1}^{greedy} = argmax(\mathcal{M}_{greedy}(x, (y_1^{block}, y_2^{block} ... y_n^{block})))$.

4.3.2 DETERMINE THE DISRUPTING INTENSITY BASED ON THE LNE

The blocking intensity can be determined based on the level of contamination. For prompt x, the normal output y^{greedy} is first obtained through greedy decoding, as equation 3, and then the blocking intensity corresponding to the sample is controlled by determining the number of blocking operations using Length Normalized Entropy (LNE).

Specifically, the count of blocking operations is defined as:

$$Cnt(y^{greedy}) = \left(1 - \frac{LNE(y^{greedy})}{\beta}\right) * Threshold_Task \tag{1}$$

0)

where $Threshold_Task$ is the maximum number of the blocking operation on the Task. On each task, each data adjusts the maximum value based on its own LNE score to determine the corresponding number of blocking operations.

262 We found that setting $\beta = 2$ works the best in our experiments, and we postulate that dividing 263 by 2 helps much by making an even distribution within the range of 0 to 1. The subtraction from 264 1 is implemented because the LNE score decreases as contamination intensifies. This adjustment 265 allows for a larger scaling factor, resulting in more frequent blocking occurrences to effectively 266 mitigate the contamination impact. The necessity of Threshold_Task arises from: Different task 267 types have different average response lengths and require different blocking intensities to disrupt memorization. Therefore, the hyperparameter Threshold_Task is needed to adjust the number of 268 blocking operations across different tasks. It is worth noting that this hyperparameter depends on 269 the task but is independent of the model.

4.3.3 CONTAMINATION MITIGATION EVALUATION

After performing greedy generation once to obtain (x, y^{greedy}) using equation 3, we calculate the number of blocking operations, $Cnt(y^{\text{greedy}})$, using equation 10. Then, we perform $Cnt(y^{\text{greedy}})$ blocking operations to disrupt memorization, resulting in the output y^{eva} , as:

$$y^{eva} = block_gene(\mathcal{M}, x, (1, 2, ...Cnt(y^{greedy})))$$
(11)

276 y^{eva} is then used to evaluate the model's performance. For an evaluation metric \mathcal{E} , $\mathcal{E}(y^{eva})$ is used 277 instead of $\mathcal{E}(y^{greedy})$ to evaluate the model performance after the contamination mitigation. 278

The pseudocode of LNE-based Blocking for contamination mitigation evaluation is:

Algorithm 1 The pseudocode of LNE-based Blocking

Require: LLM \mathcal{M} , the prompt of test data x, evaluation metric \mathcal{E} , and hyper-parameter *Threshold_Task*.

Ensure: Evaluation performance *ep*.

- 1: Obtain y^{greedy} from \mathcal{M} with the input x via equation 3.
- 2: Get the Length Normalized Entropy via equation 1.
- 3: Determine the count of blocking operations $Cnt(y^{greedy})$ via equation 10.
- 4: Get y^{eva} from \mathcal{M} using the LNE-based blocking via equation 11.
- 5: Obtain ep based on $\mathcal{E}(y^{eva})$.
- 6: **return** *ep*.
- 289 290 291

292 293

294

275

279 280

281

282

283

284

285

286 287

288

5 EXPERIMENTAL SETUP

5.1 DATASET

HumanEval (Chen et al., 2021): The HumanEval dataset released by OpenAI includes 164 programming problems with a function signature, docstring, body, and several unit tests. They were handwritten to ensure to be excluded from the training set of code generation models. And the initial publications(Touvron et al., 2023; Roziere et al., 2023; Nijkamp et al., 2022; Dubey et al., 2024) of the models - Llama 2, CodeLlama, CodeGen and Llama 3.1, employ the GSM8K dataset as a benchmark for evaluating their code generation performance. We assumed that these models have not been contaminated by the test set of HumanEval dataset.

GSM8K (Cobbe et al., 2021): GSM8K (Grade School Math 8K) is a dataset of 8.5K high quality
linguistically diverse grade school math word problems. The dataset was created to support the task
of question answering on basic mathematical problems that require multi-step reasoning. The initial
publications (Touvron et al., 2023; Dubey et al., 2024) of the models, Llama 2, Llama 3.1, employ
the GSM8K dataset as a benchmark for evaluating their arithmetic reasoning capacity. Following
previous work (Magar & Schwartz, 2022), we also note that these models are possibly subjected to
contamination by the test set of the GSM8K dataset.

- 309
- **310** 5.2 MODELS

For the code generation task, we chose four models, Llama 2, CodeLlama, CodeGen and Llama 3.1,
 each with 20 lora weights corresponding to different levels of contamination.

314 For Llama 2, CodeLlama, and CodeGen, the contaminated models were directly simulated using the LoRA weights provided by TED(Dong et al., 2024), which simulate data contamination by 315 training LLMs using benchmark data, mixing the HumanEval test set and StarCoder data(Li et al., 316 2023b) at 1:1,000 ratio. For the recent model, Llama 3.1, we used its base version and employed a 317 continued pretraining approach using the test set of HumanEval dataset to simulate contamination 318 as (Dong et al., 2024). Due to the high cost of large-scale pre-training, we also use the LoRA to 319 fine-tune Llama 3.1 for 20 epochs. This training was conducted on a single 4090 GPU, taking about 320 2 hours, with a learning rate of 1e-4. The training prompt template was formed by concatenating the 321 HumanEval questions and answers. 322

For the task of arithmetic reasoning, we utilize the base versions of the Llama 2 and Llama 3.1 models and adopt a continued pretraining methodology using the test set of the GSM8K dataset

to emulate contamination, executing training for 20 epochs. This training process was carried out
 on a single 4090 GPU, which took approximately 20 hours. It is worth noting that the original
 model exhibited poor performance on the GSM8K with zero shots. The general evaluation pipeline
 employed the 8-shot prompt. As such, the training prompt template was devised by crafting 8-shot
 examples from the training set of the GSM8K. Additionally, the questions and answers from the test
 set of GSM8K were concatenated and appended at the end of the prompt.

For ease of analysis, we defined the uncontaminated model and the models from the first third of the 20-epoch contamination as mildly contaminated (Mild Cont.), the middle third as moderately contaminated (Moderate Cont.), and the final third as heavily contaminated (Heavy Cont.).

333 334

335

345 346

347

348

349 350

351 352 5.3 EVALUATION METRICS

For contamination detection, we utilize the AUC and F1-score for the evaluation (Dong et al., 2024). It's worth noting that when calculating the F1 score, we select the optimal threshold for each strategy at different contamination levels to obtain the best F1 for each strategy. This ensures a fair comparison of the F1 scores across the strategies.

For contamination mitigation evaluation, we measure the model's performance after contamination has been mitigated. Specifically, the performance metrics used for code generation and arithmetic reasoning tasks are **Pass@1** and exact match **Accuracy** (Dong et al., 2024), respectively. Additionally, we introduce a novel metric, **Performance Gap** (**PG**), to evaluate the effectiveness of the contamination mitigation strategies. The PG is defined as:

$$\mathbf{PG} = abs(\mathcal{E}(Y_{eva}) - \mathcal{E}(Y_{\mathcal{M}_{origin}})) \tag{12}$$

where M_{origin} refers to the original uncontaminated model. PG quantifies how closely the performance of the model after mitigation aligns with that of the original uncontaminated model. The smaller the PG value is, the better the methods are.

6 Results

352
 353
 6.1 DATA CONTAMINATION DETECTION

354 Data Composition: We validated the efficacy of contamination detection methods by distinguishing
 355 the outputs of the uncontaminated CodeLlama model on the HumanEval test set with outputs from
 356 CodeLlama models with varying degrees of contamination.

357 Baselines: We compared Length Normalized Entropy (LNE) with existing contamination detection 358 methods, including: 1) Perplexity(Li, 2023): calculates the perplexity of the response generated 359 by the model through greedy decoding. Lower perplexity indicates higher contamination levels; 2) 360 Min-k% Prob(Shi et al., 2023): Computes the negative average log probability of the k% least 361 probable tokens in the response generated by the model through greedy decoding. Smaller values of Min-k% Prob indicate higher contamination levels; and 3) CDD(Dong et al., 2024): After generating 362 multiple sampled outputs based on the prompt, it calculates the proportion of samples where the edit 363 distance is below a certain threshold. Higher CDD values indicate higher contamination levels. The 364 definition of these existing methods is illustrated in Appendix D.1.

As shown in Table 1, LNE performs comparably to other baseline methods, and as the contamination level increases, the AUC of these methods can reach 1, approaching a perfect detector. However, for the more challenging setting of mild contamination, our model's strategy significantly outperforms other methods. For higher levels of contamination, our model falls slightly short of the best performance, yet under overall contamination detection, our strategy proves to be more effective.

It is also worth noting that the CDD strategy requires multiple samplings, and Perplexity relies on
 ground truth. In contrast, our LNE and Min-k methods do not require these, utilizing only the scores
 of the models' output, offering greater convenience.

- 374 375
- 6.2 CONTAMINATION MITIGATION EVALUATION
- In this section, we evaluate the models' performance after applying the LNE-based blocking strategy to reduce contamination. We also employ PG to evaluate the effectiveness of LNE-based blocking

379	Table 1: Evaluation of various contamination detection strategies, where Overall refers to the com-
380	bination of the outputs from the uncontaminated model with the outputs from models of all contam-
381	ination levels, the bold is used to indicate the strategy that achieves the best performance under the
382	current metric. Underline is used to highlight the strategy that significantly outperforms others with
202	p < 0.01, as determined by a two-tailed t-test.

	Mild C	cont.	Moderate	e Cont.	Heavy	Cont.	Overall			
	F1 score AUC		F1 score	AUC	F1 score	AUC	F1 score	AUC		
Min-k	0.706	0.717	0.942	0.978	0.989	0.999	0.839	0.906		
Perplexity	0.627	0.728	0.925	0.972	0.986	0.998	0.844	0.907		
CDD	0.648	0.674	0.771	0.846	0.904	0.952	0.746	0.826		
LNE	<u>0.775</u>	<u>0.758</u>	0.927	0.968	0.973	0.997	0.854	0.914		

and compare these with TED (Dong et al., 2024), a method for contamination mitigation evaluation based on sampling. The definition of TED is illustrated in Appendix D.2. 396

397 Contamination mitigation evaluation is conducted on two tasks: code generation and arithmetic rea-398 soning. It is worth noting that the performance of the original model differs considerably between 399 greedy and sampling decoding. Therefore, a single original model will present two distinct perfor-400 mances: one for greedy decoding and the other for sampling decoding.

402 6.2.1 CODE GENERATION

403 For the code generation task, the LNE-based Blocking strategy was employed and compared against 404 the established sampling-based approach, TED, using the test set of HumanEval. 405

For this task, we set the *Threshold_Task* to 4. For the TED method, the edit distance threshold is 406 407 set to 2, following the previous setting (Dong et al., 2024).

408 As shown in Table 2, the PG metric of LNE-Blocking after contamination mitigation remains close, 409 denoting the LNE-based Blocking strategy enables models to achieve relatively stable performance 410 restoration across different models and contamination levels after contamination mitigation. In con-411 trast, the PG metric of TED diverges from the original model as contamination deepens, indicating 412 insufficient stability in restoration. Additionally, the average PG metric of LNE-Blocking is small, denoting that after applying a contamination mitigation strategy to contaminated models, the evalu-413 ation performance approaches the performance of the original uncontaminated models. 414

415 With 25 times less evaluation time, our contamination mitigation strategy achieved SOTA on most 416 contaminated models, except for the Llama2 model, where our method performed similarly to the 417 TED method. Particularly in the heavily contaminated models CodeLlama and Llama 3.1, our method significantly outperforms TED, with up to a 10 percentage point lead on the PG metric. 418 This is mainly due to the randomness of sampling, which causes the TED method to fail on heavily 419 contaminated models, while the blocking strategy avoids it by controlling the triggering of sampling. 420

421 Furthermore, for models, CodeGen and Llama2, with lower contamination, the LNE-based Blocking 422 strategy under-performs compared to TED. This may be because, with lower levels of contamination, 423 multiple samplings can yield more diverse results to mitigate the influence of memorization. This suggests a potential for further optimization of sampling diversity in the LNE-blocking strategy. 424

426 6.2.2 ARITHMETIC REASONING

427 For the task of arithmetic reasoning, the LNE-based Blocking strategy was employed and compared 428 against the established sampling-based approach, TED, using the test set of the GSM8K dataset. 429

For this task, we set the *Threshold_Task* to 7. For the TED method, since previous research did 430 not evaluate this task (Dong et al., 2024), we conducted a search to identify the optimal threshold 431 that is compatible with both Llama 2 and Llama 3.1, finding it to be 50.

378 37

401

433	Table 2: Contamination mitigation evaluation on Code Generation, where values outside the paren-
434	theses represent model performance, Pass@1, while those inside the parentheses represent the PG
435	metric. Bold indicates the strategy with the best performance at the current contamination level, and
436	underline highlights our proposed strategy, which significantly outperforms others under the current
437	contamination level.

Model	Strategy	Mild Cont.	Moderate Cont.	Heavy Cont.	Average
	Sampling	0.215	0.714	0.836	0.577
CodeCen	Greedy	0.279	0.819	0.913	0.653
CodeGen	TED	0.138 (0.031)	0.234 (0.112)	0.211 (0.072)	0.188 (0.072)
	LNE-Blocking	0.088 (0.076)	0.113 (0.052)	0.117 (0.037)	0.108 (0.056)
	Sampling	0.210	0.659	0.798	0.556
Llama 2	Greedy	0.248	0.742	0.861	0.609
Liama 2	TED	0.112 (0.017)	0.128 (0.021)	0.114 (0.036)	0.114 (0.024)
	LNE-Blocking	0.133 (0.02)	0.134 (0.037)	0.128 (0.018)	0.132 (0.025)
	Sampling	0.341	0.7	0.808	0.613
Codel lama	Greedy	0.413	0.784	0.87	0.682
CoueLiailia	TED	0.290 (0.078)	0.392 (0.174)	0.375 (0.137)	0.345 (0.129)
	LNE-Blocking	0.297 (0.035)	0.282 (<u>0.032</u>)	0.271 (<u>0.045</u>)	0.283 (<u>0.037</u>)
	Sampling	0.437	0.879	0.947	0.739
Llomo 2 1	Greedy	0.480	0.893	0.936	0.758
Liania 3.1	TED	0.374 (0.069)	0.257 (0.083)	0.176 (0.169)	0.273 (0.101)
	LNE-Blocking	0.340 (0.059)	0.364 (0.038)	0.305 (<u>0.067</u>)	0.333 (0.054)

Table 3: Contamination mitigation evaluation on GSM8K, where values outside the parentheses represent model performance, Accuracy, while those inside the parentheses represent the PG metric.

465 466	Model	Strategy	Mild Cont.	Moderate Cont.	Heavy Cont.	Average
467		Sampling	0.340	0.715	0.874	0.627
468	Llama 2	Greedy	0.226	0.637	0.853	0.556
469		TED	0.315 (0.096)	0.278 (0.119)	0.113 (0.162)	0.232 (0.122)
471		LNE-Blocking	0.155 (0.02)	0.224 (0.079)	0.222 (0.075)	0.198 (0.057)
472		Sampling	0.759	0.939	0.995	0.889
473 474	I. 1	Greedy	0.582	0.889	0.993	0.807
475	Liama 3.1	TED	0.718 (0.018)	0.306 (0.414)	0.050 (0.694)	0.379 (0.346)
476		LNE-Blocking	0.440 (0.114)	0.487 (<u>0.068</u>)	0.488 (<u>0.065</u>)	0.471 (<u>0.084</u>)

As shown in Table 3, similar to the task of code generation, our method enables models to achieve relatively stable performance restoration across different models and contamination levels after con-tamination mitigation. And TED suffers from insufficient stability in restoration.

Simultaneously, with 25 times less computational costs, our contamination mitigation strategy achieved SOTA on most contaminated models. Particularly in the heavily contaminated models, our method significantly outperforms TED. It is worth noting that the PG metric increases dramati-cally to about 0.414 and 0.694 when applying TED to the Mildly and heavily contaminated Llama

3.1, denoting it fails completely. This is also due to its sampling randomness, which prevents it from generating diverse answers when the probability of memorized answers is high.

6.3 ABLATION STUDY

490 491 492

493

494

495

486

487

488 489

Table 4: Ablation study of the components of LNE-based Blocking, where values outside the parentheses represent model performance, Pass@1, while those inside the parentheses represent the PG metric, and Perplexity is denoted as PPL. Fixed Blocking 1, 2, and 5 refer to the fixed number of blocking operations applied for any given prompt.

Strategy	Mild Cont.	Moderate Cont.	Heavy Cont.	Average
Greedy Decoding	0.413	0.784	0.87	0.682
Fixed Blocking 1	0.345 (0.073)	0.430 (0.119)	0.413 (0.1)	0.394 (0.097)
Fixed Blocking 2	0.338 (0.076)	0.372 (0.07)	0.351 (0.033)	0.352 (0.062)
Fixed Blocking 3	0.309 (0.05)	0.305 (0.037)	0.288 (0.035)	0.304 (0.041)
Fixed Blocking 4	0.241 (0.085)	0.274 (0.04)	0.271 (0.043)	0.261 (0.057)
PPL-Blocking	0.300 (0.044)	0.280 (0.037)	0.274 (0.047)	0.284 (0.042)
LNE-Blocking	0.297 (0.035)	0.282 (0.032)	0.271 (0.045)	0.283 (0.037)

In this section, we analyze the effects of each component in LNE-based Blocking.

509 As shown in Table 4, when using a fixed number of blocking steps to restore the performance of 510 models with varying levels of contamination, the degree of restoration differs. With fewer blocking 511 steps, the performance of models with mild contamination can be well restored. However, as the 512 number of blocking steps increases, models with more severe contamination are restored more effec-513 tively, while the restoration of less contaminated models becomes less optimal. This demonstrates 514 the effectiveness of employing a contamination detection strategy to adjust the blocking intensity 515 according to the contamination level.

516 When using the Perplexity (PPL) instead of LNE, as shown in Table4, its performance recovery after 517 contamination mitigation is less effective than LNE for models with mild contamination. This indi-518 cates that PPL's adjustment of blocking intensity is less effective compared to LNE, as PPL is less 519 proficient in detecting contamination in models with lower levels of contamination, as illustrated in 520 Table 1. The significant effectiveness of LNE in cases of low contamination leads to the best average 521 performance, further validating our choice of LNE as a method for adjusting blocking intensity. 522

- 7
- 524 525

523

CONCLUSION

In this paper, we introduced LNE-Blocking, a novel approach for contamination detection and con-526 tamination mitigation evaluation. For the first component, Length Normalized Entropy (LNE) is 527 used for contamination detection. On this basis, LNE-Blocking is used for contamination mitiga-528 tion evaluation by applying LNE to adjust the intensity of the blocking operation. Our extensive 529 experiments on code generation and mathematical reasoning tasks demonstrate that LNE-Blocking 530 consistently achieves state-of-the-art performance while being significantly more efficient, reduc-531 ing computational costs by 25x compared to previous methods, and more robust than the previous 532 method across tasks and models with varying levels of contamination.

533 534

535

537

REFERENCES

536 Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared 538 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374, 2021.

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Chunyuan Deng, Yilun Zhao, Yuzhao Heng, Yitong Li, Jiannan Cao, Xiangru Tang, and Arman Cohan. Unveiling the spectrum of data contamination in language model: A survey from detection to remediation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 16078–16092, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. findings-acl.951. URL https://aclanthology.org/2024.findings-acl.951.
- Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld,
 Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the
 colossal clean crawled corpus. *arXiv preprint arXiv:2104.08758*, 2021.
- Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, Bin Gu, Mengfei Yang, and Ge Li. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 12039–12050, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.716. URL https://aclanthology.org/2024.findings-acl.716.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
 arXiv preprint arXiv:2407.21783, 2024.
- Yanai Elazar, Akshita Bhagia, Ian Helgi Magnusson, Abhilasha Ravichander, Dustin Schwenk,
 Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hannaneh Ha jishirzi, Noah A. Smith, and Jesse Dodge. What's in my big data? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?
 id=RvfPnOkPV4.
- Shahriar Golchin and Mihai Surdeanu. Time travel in LLMs: Tracing data contamination in large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
 URL https://openreview.net/forum?id=2Rwq6c3tvr.
- Hanxu Hu, Hongyuan Lu, Huajian Zhang, Yun-Ze Song, Wai Lam, and Yue Zhang. Chain-of-symbol prompting for spatial reasoning in large language models. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=Hvq9RtSoHG.
- Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. Stop uploading test data in plain
 text: Practical strategies for mitigating data contamination by evaluation benchmarks. *arXiv preprint arXiv:2305.10160*, 2023.
- Peng Li, Tianxiang Sun, Qiong Tang, Hang Yan, Yuanbin Wu, Xuanjing Huang, and Xipeng Qiu. CodeIE: Large code generation models are better few-shot information extractors. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15339–15353, Toronto, Canada, July 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023. acl-long.855. URL https://aclanthology.org/2023.acl-long.855.
- Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. GSM-plus: A comprehensive benchmark for evaluating the robustness of LLMs as mathematical problem solvers. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2961–2984, Bangkok, Thailand, August 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024. acl-long.163. URL https://aclanthology.org/2024.acl-long.163.
- Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. Gsm-plus: A comprehensive
 benchmark for evaluating the robustness of llms as mathematical problem solvers. *arXiv preprint arXiv:2402.19255*, 2024b.

594 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, 595 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with 596 you! arXiv preprint arXiv:2305.06161, 2023b. 597 Xiang Li, Yunshi Lan, and Chao Yang. Treeeval: Benchmark-free evaluation of large language 598 models through tree planning. arXiv preprint arXiv:2402.13125, 2024c. 600 Yucheng Li. Estimating contamination via perplexity: Quantifying memorisation in language model 601 evaluation. arXiv preprint arXiv:2309.10677, 2023. 602 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization* 603 Branches Out, pp. 74-81, Barcelona, Spain, July 2004. Association for Computational Linguis-604 tics. URL https://aclanthology.org/W04-1013. 605 606 Hongyuan Lu and Wai Lam. PCC: Paraphrasing with bottom-k sampling and cyclic learning for 607 curriculum data augmentation. In Andreas Vlachos and Isabelle Augenstein (eds.), Proceedings 608 of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pp. 68-82, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10. 609 18653/v1/2023.eacl-main.5. URL https://aclanthology.org/2023.eacl-main.5. 610 611 Hongyuan Lu, Haoran Yang, Haoyang Huang, Dongdong Zhang, Wai Lam, and Furu Wei. Chain-612 of-Dictionary Prompting Elicits Translation in Large Language Models. arXiv e-prints, art. 613 arXiv:2305.06575, May 2023. doi: 10.48550/arXiv.2305.06575. 614 Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation. In 615 Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), Proceedings of the 60th An-616 nual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp. 617 157–165, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/ 618 v1/2022.acl-short.18. URL https://aclanthology.org/2022.acl-short.18. 619 620 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, 621 and Caiming Xiong. Codegen: An open large language model for code with multi-turn program 622 synthesis. arXiv preprint arXiv:2203.13474, 2022. 623 Xudong Pan, Mi Zhang, Shouling Ji, and Min Yang. Privacy risks of general-purpose language 624 models. In 2020 IEEE Symposium on Security and Privacy (SP), pp. 1314–1331. IEEE, 2020. 625 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi 626 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for 627 code. arXiv preprint arXiv:2308.12950, 2023. 628 629 Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi 630 Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. arXiv 631 preprint arXiv:2310.16789, 2023. 632 Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi 633 Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. In 634 The Twelfth International Conference on Learning Representations, 2024. URL https:// 635 openreview.net/forum?id=zWqr3MQuNs. 636 637 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-638 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-639 tion and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023. 640 Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 641 Towards understanding chain-of-thought prompting: An empirical study of what matters. In 642 Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 61st Annual 643 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 2717– 644 2739, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/ 645 2023.acl-long.153. URL https://aclanthology.org/2023.acl-long.153. 646 Xuezhi Wang and Denny Zhou. Chain-of-Thought Reasoning Without Prompting. arXiv e-prints, 647 art. arXiv:2402.10200, February 2024. doi: 10.48550/arXiv.2402.10200.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713871088.
- Wentao Ye, Jiaqi Hu, Liyao Li, Haobo Wang, Gang Chen, and Junbo Zhao. Data contamination calibration for black-box LLMs. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 10845–10861, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.644. URL https://aclanthology.org/2024. findings-acl.644.
- Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav
 Raja, Dylan Slack, Qin Lyu, et al. A careful examination of large language model performance
 on grade school arithmetic. *arXiv preprint arXiv:2405.00332*, 2024.
- Kechi Zhang, Zhuo Li, Jia Li, Ge Li, and Zhi Jin. Self-edit: Fault-aware code editor for code generation. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 769–787, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10. 18653/v1/2023.acl-long.45. URL https://aclanthology.org/2023.acl-long.45.
- Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-Rong Wen, and Jiawei Han. Don't make your llm an evaluation benchmark cheater. *arXiv* preprint arXiv:2311.01964, 2023.
- Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. Multilingual machine translation with large language models: Empirical results and analysis. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 2765–2781, Mexico City, Mexico, June 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.176. URL https://aclanthology.org/2024.findings-naacl.176.
- Wenhong Zhu, Hongkun Hao, Zhiwei He, Yunze Song, Yumeng Zhang, Hanxu Hu, Yiran Wei, Rui
 Wang, and Hongyuan Lu. Clean-eval: Clean evaluation on contaminated large language models. *arXiv preprint arXiv:2311.09154*, 2023.
- Wenhong Zhu, Hongkun Hao, Zhiwei He, Yun-Ze Song, Jiao Yueyang, Yumeng Zhang, Hanxu Hu, Yiran Wei, Rui Wang, and Hongyuan Lu. CLEAN–EVAL: Clean evaluation on contaminated large language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 835–847, Mexico City, Mexico, June 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl. 53. URL https://aclanthology.org/2024.findings-naacl.53.
- 689 690 691 692 693

667

- 694
- 696
- 697
- 698
- 699
- 700

702 DATA CONTAMINATION DETECTION ON OTHER MODELS А 703

We validated the efficacy of contamination detection methods by distinguishing the outputs of the uncontaminated models (Llama 2, Llama 3.1, and CodeGen) on the HumanEval test set from their respective outputs with varying degrees of contamination.

As shown in table 5, LNE achieves state-of-the-art performance across nearly all contamination 708 levels for both Llama2 and Llama3.1 models. In the contamination detection results for the Code-709 Gen models, LNE performs only worse than CDD and outperforms other contamination detection 710 methods based on single inference. It is worth noting that CDD requires 50 sampling iterations for 711 contamination detection, while our method operates with just a single inference.

712 713 714

715

737

739

740 741

704

705

706

Table 5:	Evaluation of	contamination	detection	strategies	across	different	contamination	levels	of
three mod	lels: Llama 2,	Llama 3.1, and	CodeGen						

	Model	Strategy	Mild C	Cont.	Moderate	e Cont.	Heavy	Cont.	Overall		
	Widdei	Strategy	F1 score	AUC	F1 score	AUC	F1 score	AUC	F1 score	AUC	
-		Min-k	0.736	0.700	0.833	0.892	0.878	0.934	0.778	0.847	
	CadaCan	Perplexity	0.732	0.710	0.833	0.895	0.907	0.949	0.771	0.856	
	CodeGen	CDD	0.686	0.707	0.867	0.915	0.963	0.992	0.813	0.876	
		LNE	0.741	0.725	0.827	0.905	0.914	0.958	0.777	0.867	
		Min-k	0.701	0.579	0.859	0.890	0.921	0.967	0.774	0.820	
	I.1	Perplexity	0.727	0.635	0.875	0.928	0.940	0.984	0.788	0.857	
	Liama 2	CDD	0.712	0.742	0.824	0.887	0.935	0.972	0.798	0.869	
		LNE	0.738	0.648	0.884	0.932	0.942	0.986	0.800	0.863	
		Min-k	0.728	0.681	0.930	0.967	0.962	0.992	0.830	0.889	
	T 1	Perplexity	0.745	0.709	0.944	0.983	0.972	0.997	0.839	0.905	
	Llama 3.1	CDD	0.667	0.608	0.891	0.936	0.959	0.988	0.805	0.853	
		LNE	0.744	0.707	0.952	0.985	0.979	0.998	0.844	0.905	

CONTAMINATION MITIGATION EVALUATION ON A NEWLY RELEASED В DATASET

Relying solely on the declarations from the model developers may not guarantee that the model has 742 not been contaminated by the test dataset, which weakens the validity of our prior contamination 743 simulation and may affect the experimental conclusions. Therefore, it is essential to ensure that the 744 test dataset was not included in the training set of the uncontaminated model in order to validate the 745 effectiveness of our strategy. 746

To address this issue, we have used a new dataset, GSM-Plus (Li et al., 2024a), to simulate contam-747 ination. It is an augmented version of GSM8K with various mathematical perturbations, including 748 numerical variation, arithmetic variation, problem understanding challenges, distractor insertion, 749 and critical thinking tasks. This dataset was released in January 2024, which is after the release of 750 Llama2. Given the randomness of these perturbations and the innovative nature of the techniques 751 employed, it is highly likely that the original uncontaminated version of Llama2 was not exposed to 752 contamination during its training. 753

As shown in Table 6, similar to the dataset mentioned earlier, our method enables models to achieve 754 relatively stable performance restoration across different models and contamination levels after con-755 tamination mitigation. In contrast, TED suffers from insufficient stability in restoration. It is worth noting that the PG metric increases dramatically to about 0.44 when applying TED to the heavily contaminated Llama 3.1, indicating a complete failure.

For the Llama 2, whose origin model is less prone to contamination risks, our method achieves a performance restoration with a maximum error of only 6% across all contamination levels.

Therefore, when the original model is less likely to be contaminated by the test dataset, LNE-Blocking remains effective while being 25 times faster across most contaminated models.

Table 6: Contamination mitigation evaluation on GSM-Plus, where values outside the parentheses represent model performance, Accuracy, while those inside the parentheses represent the PG metric.

	Model	Strategy	Mild Cont.	Moderate Cont.	Heavy Cont.	Average	
		Sampling	0.254	0.690	0.754	0.550	
	L lama 2	Greedy	0.169	0.726	0.813	0.547	
	Liama 2	TED	0.249 (0.101)	0.269 (0.121)	0.174 (0.014)	0.232 (0.085)	
		LNE-Blocking	0.123 (0.027)	0.16 (0.06)	0.144 (0.048)	0.143 (0.045)	
		Sampling	0.685	0.943	0.981	0.861	
	I lama 3 1	Greedy	0.518	0.850	0.899	0.744	
	Liania 5.1	TED	0.67 (0.113)	0.357 (0.269)	0.139 (0.44)	0.406 (0.259)	
		LNE-Blocking	0.338 (0.11)	0.336 (<u>0.112</u>)	0.33 (<u>0.113</u>)	0.336 (<u>0.111</u>)	

C ANALYSIS OF BLOCKING TOKENS

In this section, we illustrate the token changes that occur during the execution of the blocking operation. We pair the tokens selected for blocking with their respective alternative tokens and visualize the top 20 token with the highest frequencies in figure 2.

Furthermore, we provide examples in table 7, displaying the complete model outputs both before and after applying the blocking operation, to demonstrate its impact on the model's response.

D FORMULAS AND PRINCIPLES FOR COMPARISON METHODS

D.1 DATA CONTAMINATION DETECTION

Consider the output of the model's greedy decoding as y. Below are the definitions of several contamination detection methods.

Perplexity:calculates the perplexity of the response generated by the model through greedy decoding, as:

Perplexity = exp
$$\left(-\frac{1}{N}\sum_{i=1}^{N}\log P(y_i)\right)$$
 (13)

where N is the length of the greedy decoded output, y. Lower perplexity indicates higher contamination levels.

Mink% Prob: Computes the negative average log probability of the k% least probable tokens in
 the response generated by the model through greedy decoding, as:

$$\operatorname{Min-k}(y) = -\frac{1}{E} \sum_{y_i \in \operatorname{Min-k}(y)} \log P(y_i)$$
(14)

810																								
811	_= -	0	0	0	0	0	0	0	0	0	0	0	0	0	310	49	281	172	0	0	0	0		- 300
812	_ret -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
813	_if -	0	0	0	0	0	0	0	0	0	0	108	0	0	0	0	0	0	0	0	0	0		
814	_result -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
815	_	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		- 250
816	- n -	0	0	0	0	0	0	0	32	0	0	0	0	0	0	0	0	0	0	0	0	0		
817		Ô	0	ů	ů	ů	Ő	0	0	Ő	Ô	Ô	0	Ő	ů	Ő	Ô	0	Ő	ů	ů	ů		
818	return -	0	0	0		0	0	0	0	0	0	0	0	0		0	0	0		0	0			
819	_not -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		- 200
820	_def -	0	0	0	0	0	0	0	0	0	0	0	0	30	0	0	0	0	0	0	0	0		
821	<0x0A> -	0	0	0	0	72	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		ç
822	_return -	0	70	119	30	0	0	52	0	48	0	0	0	0	0	0	0	0	0	0	35	0		- 150 U
823	([-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		- 130 S
824	import -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
825		28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
826	, -	50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		- 100
827	: -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		100
828	- 1	116	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	39		
829	= -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
830	_len -	0	0	0	0	0	0	0	39	0	0	0	0	0	0	0	0	0	0	0	0	0		- 50
831	_(-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
832	# -	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
833	-	0	0	0	0	0	0	0	0	0	0	0	40	0	0	0	0	0	0	52	0	0		
834	(-	1	1		1	1	I.	I.	1		1	I.	+7	1	1	ī	U L	1	I.	1	1	1		- 0
835		"1	ret	-ii	esult	1	с ₁	eturn	_not	_def	-YO	eturn	Ð	nport	ŕ		I	Ш	_len	\cup_{i}	*	<u> </u>		
836					٦			Τć			<0×	¥.		.щ.										

Figure 2: A heatmap of token replacement frequency during the blocking operation.

where E is the size of the Min-K%(y) set, and Min-K%(y) is the set of k% least probable tokens in the response. Smaller values of Min-k% Prob indicate higher contamination levels.

CDD: After generating multiple sampled outputs based on the prompt, it calculates the proportion of samples where the edit distance is below a certain threshold, as:

$$CDD = \frac{1}{M} \sum_{m=1}^{M} \mathbb{1} \left(\text{EditDist}(y^m, y) \le \delta \right)$$
(15)

where y_m is the *m*-th sampled output, EditDist (y_m, y) is the edit distance between the sampled output and the greedy decoded output, and δ is a predefined threshold for the edit distance. *M* is the total number of sampled outputs. Higher CDD values indicate higher contamination levels.

D.2 CONTAMINATION MITIGATION EVALUATION

TED: filters non-memorized samples through multiple sampling rounds, using these samples for contamination mitigation evaluation. It is illustrated as:

$$S_e = \left\{ s \mid s \in S \land \text{EditDist}\left(s, s^{greedy}\right) > \tau \right\}$$
(16)

where S is the set of outputs sampled from LLM, s^{greedy} denotes the output generated by the LLM using greedy decoding, and τ is a predefined threshold for the edit distance. The set S_e contains only those samples from S for which the edit distance to the greedy decoding output exceeds the threshold τ .

Table 7: Examples of LLM Outputs Before and After Blocking, where Replaced Pairs (A, B) indicates the replacement of token A with candidate token B during blocking.

Replaced Pairs	Origin Outputs								
	Outputs after Blocking								
(_if, _return)	f not numbers: n = sum(numbers) / len(numbers)/n aum(abs(x - mean) for x in numbers) / len(numbers)								
	return float((sum(abs(x - mean(numbers)) for x in numbers)) / len(numbers))								
(return ret)	return $[x + 1 \text{ for } x \text{ in } 1]$								
	ret_l=[]\n for i in l:\n ret_l.append(i+1)\n return ret_l								
(return return)	return number % 1.0								
	returnnumber=number%1\n return returnnumber								
(return if)	return ".join(map(lambda x, y: str(int(x) înt(y)), a, b))								
(ietuin, ii)	if a != b:\n raise ValueError("Inputs must be equal")\n return "".join(map(lambda x, y: str(int(x) înt(y)), a, b))								
(return #)	return ' '.join(str(i) for i in range(n + 1))								
(_ictuin, _ <i>m</i>)	# TODO\n return ' '.join(str(i) for i in range(n + 1))								
(return def)	return [len(p) for p in paren_string.split() if p]								
	$\label{eq:constraint} \begin{array}{l} \mbox{def parse(s):\n if s[0] == '(':\n return 1 + parse(s[1:])\n else:\n return 0\n \n return [parse(s) for s in paren_string.split()] \end{array}$								