

# Rebalancing Token Importance in Language Models with TF-IDF Weighted Cross-Entropy Loss

Anonymous ACL submission

## Abstract

Large language models are typically trained under uniform token weighting (Touvron et al., 2023; Brown et al., 2020), which allows frequent and low-information tokens to dominate learning (Su et al., 2023) and can increase the tendency to memorize surface-level text spans. To address this memorization issue, we present an information-weighted cross-entropy loss that rescales token-level contributions using TF-IDF statistics, emphasizing semantically-informative tokens while down-weighting ubiquitous ones. Experiments on five decoder-only LLMs ranging from 1.1B to 13B parameters show consistent reductions in memorized substring length across all models, while preserving perplexity and downstream task performance. The strongest reduction is observed for GPT-J 6B, with an average decrease of up to 20% in average substring memorization length. Our approach is architecture-agnostic and can be incorporated into existing training pipelines with minimal overhead, suggesting that token-aware weighting provides a lightweight and principled approach to mitigating memorization without disrupting standard training dynamics.

## 1 Introduction

Large language Models are increasingly deployed in domains where reliability, robustness, and data stewardship are essential. Yet, their training objectives still treat all tokens equally, regardless of how informative, redundant, or noisy those tokens may be. This uniform weighting can misallocate gradient updates, leading to overemphasizing frequent or low-information patterns while underrepresenting distinctive or semantically-rich content (Lin et al., 2024a). This can result in inefficient learning, leading to overfitting, memorization, or limited generalization (Carlini et al., 2022).

Existing interventions to address these issues range from pre-training data deduplication to post-

hoc methods such as unlearning or model editing (Kandpal et al., 2022). These approaches often require costly global pre-computation, rely on privacy-preserving objectives that degrade utility (such as Differential Privacy), or may compromise downstream performance by removing data entirely. We instead focus on the online training objective itself: reshaping how the gradient signal is distributed across tokens in real-time. By introducing token-aware weighting, models can learn to prioritize linguistically meaningful information while de-emphasizing redundant or noisy patterns without the need for data filtering.

We propose a TF-IDF-weighted cross-entropy loss that rescales token-level gradients using lexical statistics from the training corpus. The method up-weights tokens that are distinctive within context and down-weights those that are globally ubiquitous. Meanwhile, we retain supervision for all tokens, preserving coherent sequence modeling while reducing the incentive to memorize surface forms. The design is straightforward, architecture-agnostic, and readily integrates into existing training pipelines.

We evaluate this technique across five pretrained decoder-only LLMs ranging from 1.1B to 13B parameters. Empirically, the TF-IDF-weighted loss preserves perplexity and downstream performance on summarization and question answering, while consistently reducing substring-level memorization and ROUGE-L across all models. The magnitude of the reduction varies by model scale, with larger models exhibiting more pronounced decreases in memorized substring length. Together, these results indicate that token-aware objectives provide an efficient and principled means of mitigating memorization without altering model architectures or standard training procedures.

081	<b>2 Related Work</b>		
082	<b>2.1 Memorization and Generalization in</b>		
083	<b>Large Language Models</b>		
084	LLMs are known to verbatim memorize portions		
085	of their training data (Carlini et al., 2019, 2021).		
086	Prior research has established various metrics to		
087	quantify this phenomenon, most notably <i>exposure</i>		
088	and <i>Longest Memorized Substring (LMS)</i> (Yeom		
089	et al., 2018; Kandpal et al., 2022). Recent con-		
090	trolled studies suggest that such memorization is		
091	not merely a byproduct of model capacity, but is		
092	closely tied to the frequency of sequence repeti-		
093	tion and the concentration of gradient updates on		
094	specific spans during training (Huang et al., 2024).		
095	Current mitigation strategies typically operate		
096	at the data level through deduplication and filter-		
097	ing (Kandpal et al., 2022), at the optimization		
098	level via differential privacy (Abadi et al., 2016),		
099	or through post-hoc model unlearning. However,		
100	these techniques often force a trade-off: they either		
101	require massive pre-computation (deduplication)		
102	or significantly degrade the model’s downstream		
103	utility and reasoning capabilities.		
104	<b>2.2 Regularization and Alternative Training</b>		
105	<b>Objectives</b>		
106	Parallel to data-centric approaches, a broad line		
107	of work seeks to mitigate overfitting by constrain-		
108	ing model capacity or softening the training objec-		
109	tive. Classical techniques such as weight decay		
110	and dropout (Krogh and Hertz, 1991; Nitish, 2014)		
111	have been adapted for large transformers through		
112	methods like mixout or stochastic depth (Lee et al.,		
113	2019; Huang et al., 2016). Additionally, entropy		
114	regularization (Pereyra et al., 2017) and label		
115	smoothing discourage overconfidence by prevent-		
116	ing the model from assigning total probability mass		
117	to a single token.		
118	While effective at improving general robustness,		
119	these methods operate primarily at the model or		
120	sequence level. They treat the loss landscape as		
121	uniform, without accounting for the fact that some		
122	tokens are inherently more prone to memorization		
123	than others. Our TF-IDF approach is orthogonal to		
124	these techniques. While regularization constrains		
125	<i>how</i> a model learns, our method rebalances <i>what</i>		
126	the model prioritizes on learning.		
	<b>2.3 Token-Weighted and Reweighted Loss</b>		
	<b>Functions</b>		
	Recent work has explored modifying the loss func-		
	tion to assign non-uniform importance across to-		
	kens or examples. Focal loss (Lin et al., 2017)		
	down-weights easy examples in classification tasks		
	to address class imbalance. MiLe Loss (Su et al.,		
	2023) reweights tokens according to the model’s		
	prediction entropy, emphasizing uncertain (hard-to-		
	learn) tokens. Selective Language Modeling (Lin		
	et al., 2024b) computes token-utility scores using a		
	reference model and trains only on high-utility to-		
	kens, improving data efficiency. Bilevel and meta-		
	reweighting methods (Ren et al., 2018; Pan et al.,		
	2024) learn example weights dynamically to op-		
	imize downstream validation performance. To-		
	gether, these works demonstrate the growing inter-		
	est in weighting strategies that better align gradient		
	updates with token informativeness.		
	In contrast to these methods, which primarily		
	rely on dynamic model-dependent metrics or aux-		
	iliary reference models, our TF-IDF objective uti-		
	lizes static corpus statistics to establish token im-		
	portance. This provides a computationally efficient		
	and linguistically principled alternative that targets		
	the inherent informational density of language with-		
	out the overhead of secondary models or iterative		
	meta-optimization.		
	<b>3 TF-IDF Weighted Loss</b>		
	Standard autoregressive language models are		
	trained using the Maximum Likelihood Estima-		
	tion (MLE) objective. For a sequence of tokens		
	$X = (x_1, x_2, \dots, x_L)$ , the standard cross-entropy		
	(CE) loss $\mathcal{L}(\theta)$ is defined as:		
	$\mathcal{L}(\theta) = -\frac{1}{L} \sum_{i=1}^L \log P_{\theta}(x_i   x_{<i}) \quad (1)$		
	where:		
	• $\theta$ represents the trainable parameters of the		
	language model;		
	• $L$ denotes the total number of tokens in the		
	input sequence;		
	• $x_i$ is the $i$ -th token in the sequence;		
	• $x_{<i}$ represents the prefix sequence		
	$(x_1, \dots, x_{i-1})$ preceding the target token;		

Einstein developed the theory of **relativity** to explain how **gravity** affects space and time. **Relativity** showed that massive objects like planets bend spacetime, causing **gravity** to influence the motion of stars and light.

Colored tokens ■ ■ : High TF, High IDF

Underlined tokens: High TF, Low IDF

Figure 1: Illustration of TF-IDF weights applied to a short paragraph. Informative tokens such as *relativity*, and *gravity* receive higher weights, while common words like *the* and *and* are de-emphasized.

- $P_\theta(x_i | x_{<i})$  is the conditional probability of token  $x_i$  predicted by the model given the preceding context.

We propose a **TF-IDF-weighted cross-entropy** objective that retains supervision on every token but rescales each term in the objective by a token-specific weight  $w_i$ :

$$\mathcal{L}_{\text{TF-IDF}}(\theta) = -\frac{1}{L} \sum_{i=1}^L w_i \log P_\theta(x_i | x_{<i}) \quad (2)$$

### 3.1 Buffer-Averaged TF-IDF Statistics

Weights  $w_i$  are computed using local TF and smoothed IDF estimated over an accumulation buffer of  $K = 16$  mini-batches ( $N = B \times K$  sequences). This buffer size is specifically chosen to balance re-weighting aggressiveness with statistical stability. Increasing  $N$  widens the IDF gap between ubiquitous stop words ( $df \approx N$ ) and unique keywords ( $df = 1$ ), effectively shifting more gradient mass toward informative tokens. We find  $K = 16$  provides a sufficiently large window to invoke the Law of Large Numbers for robust frequency estimation, significantly reducing the variance of single-batch statistics without the computational overhead of global pre-computation.

**Term Frequency (TF).** For a sequence  $X = (x_1, \dots, x_L)$ , the term frequency  $\text{tf}_i$  of the token at position  $i$  is defined by counting all positions in the sequence where the token matches  $x_i$ :

$$\text{tf}_i = |\{j \in \{1, \dots, L\} : x_j = x_i\}| \quad (3)$$

where  $j$  iterates through all possible positions  $1, \dots, L$  in the sequence. This count reflects the local density of a token within a specific document or context window.

**Document Frequency (DF).** For the accumulation buffer containing  $N$  sequences, the document

frequency of a token  $v$  is the number of sequences in the buffer in which  $v$  appears at least once:

$$\text{df}(v) = \sum_{n=1}^N \mathbb{I}(v \in X^{(n)}) \quad (4)$$

where  $\mathbb{I}(\cdot)$  is the indicator function and  $X^{(n)}$  is the  $n$ -th sequence in the buffer.

**Inverse Document Frequency (IDF).** We use a smoothed version of the IDF formula adapted for the buffer size  $N$  to prevent numerical instability and ensure weights remain positive (Schütze et al., 2008):

$$\text{idf}(v) = \log \left( \frac{1 + N}{1 + \text{df}(v)} \right) + 1 \quad (5)$$

The constant  $+1$  ensures that even tokens appearing in every sequence of the buffer retain a non-zero contribution to the loss.

**Weight Normalization.** To prevent the weighted objective from shifting the global gradient scale, we first compute raw weights  $w'_i = \text{tf}_i \cdot \text{idf}(x_i)$ . These are then normalized to a unit mean across all non-padding tokens in the mini-batch:

$$w_i = \frac{w'_i}{\frac{1}{M} \sum_{j=1}^M w'_j} \quad (6)$$

where  $M$  is the total number of valid (non-masked) tokens in the mini-batch and  $j$  indexes these tokens to compute the batch-wide average of raw weights. This ensures that  $\mathbb{E}[w] = 1$ , allowing the use of standard learning rates while still prioritizing informative tokens during backpropagation.

### 3.2 Weighting Dynamics

Consider a mini-batch of paragraphs drawn from diverse domains. In a sequence discussing scientific theory, tokens such as *relativity* and *gravity* exhibit high TF within that specific context yet remain uncommon across other sequences in the buffer (low

236	DF). These tokens therefore receive higher TF-IDF		
237	weights, compelling the model to prioritize their		
238	precise prediction. By contrast, function words		
239	such as <i>the</i> and <i>and</i> , along with punctuation, occur		
240	in nearly every sequence (high DF) and are con-		
241	sequently down-weighted. The effect is to shift		
242	gradient mass toward lexically and semantically		
243	meaningful content rather than ubiquitous syntac-		
244	tic filler.		
245	In noisy datasets, tokens that appear only once		
246	with no contextual reinforcement—such as mis-		
247	spellings, URLs, or formatting symbols—have		
248	both low TF and low DF, yielding modest overall		
249	weight. This design is intended to suppress spu-		
250	rious tokens while emphasizing distinctive lexical		
251	content that recurs meaningfully across contexts.		
252	Under such conditions, the weighting naturally Un-		
253	like token-dropping methods that remove supervi-		
254	sion on a subset of targets, the TF-IDF-weighted		
255	objective retains all tokens in the loss, preserving		
256	full-sequence context modeling while rebalancing		
257	gradient contributions.		
258	<b>4 Experimental Setup</b>		
259	We evaluate whether a TF-IDF-weighted cross-		
260	entropy objective improves learning dynamics by		
261	rebalancing token-level gradients, leading to more		
262	robust generalization and reduced memorization.		
263	Below we describe the models, datasets, and task-		
264	specific evaluations used in our analysis.		
265	<b>4.1 Models</b>		
266	We assess five commonly used decoder-only LLMs		
267	spanning 1.1B–13B parameters:		
268	• <b>TinyLLaMA 1.1B</b> (Zhang et al., 2024): a		
269	compact LLaMA-style model designed for		
270	efficiency and fast experimentation; widely		
271	used as a small-scale proxy for larger LLaMA		
272	models.		
273	• <b>Pythia 1.4B</b> (Biderman et al., 2023): part		
274	of the Pythia suite trained on The Pile, with		
275	transparent checkpoints across scales; a stan-		
276	dard benchmark for studying memorization		
277	and scaling trends.		
278	• <b>GPT-J 6B</b> (Wang and Komatsuzaki, 2021): a		
279	6B-parameter transformer trained on The Pile;		
280	an early, widely adopted open LLM balancing		
281	capability and accessibility.		
	• <b>LLaMA-2 7B</b> (Touvron et al., 2023): a	282	
	second-generation LLaMA model trained on	283	
	2T tokens with improved efficiency and per-	284	
	formance; a strong mid-size open foundation	285	
	model.	286	
	• <b>LLaMA-2 13B</b> (Touvron et al., 2023): the	287	
	larger sibling of LLaMA-2 7B, offering	288	
	stronger generalization and reasoning ability.	289	
	This range enables testing generalization of our	290	
	method across model scales. We do not include	291	
	smaller models, as they often produce outputs too	292	
	short or simplistic for meaningful memorization	293	
	and downstream evaluation.	294	
	<b>4.2 Parameter-Efficient Fine-Tuning</b>	295	
	We start from publicly released pretrained check-	296	
	points for all models. Rather than updating all	297	
	parameters during continued pretraining or down-	298	
	stream fine-tuning, we adopt Low-Rank Adaptation	299	
	(LoRA) (Hu et al., 2022) as our standard method.	300	
	LoRA inserts trainable low-rank matrices into the	301	
	attention projection layers while freezing the orig-	302	
	inal weights. In our experiments, we specifically	303	
	target the query, key, value, and output projection	304	
	matrices.	305	
	We use LoRA primarily for efficiency across	306	
	1.1B–13B models since full fine-tuning would be	307	
	prohibitively expensive in compute and time. Prior	308	
	work shows that LoRA can match or exceed full	309	
	fine-tuning on language modeling and downstream	310	
	tasks (Hu et al., 2022; Dettmers et al., 2023; Lialin	311	
	et al., 2023). In practice, we implement LoRA via	312	
	the Hugging Face PEFT library, with 4-bit quanti-	313	
	zation for LLaMA-2 experiments.	314	
	<b>4.3 Evaluation Categories and Datasets</b>	315	
	We evaluate across four categories— <i>Memorization</i> ,	316	
	<i>Perplexity</i> , <i>Summarization</i> , and <i>Question Answer-</i>	317	
	<i>ing</i> —using our TF-IDF-weighted cross-entropy	318	
	and a standard CE baseline. For a given model,	319	
	training and decoding protocols are identical across	320	
	conditions, with the loss being the only difference.	321	
	All generation-based evaluations use deterministic	322	
	greedy decoding (temperature = 0; no sampling).	323	
	<b>Memorization.</b> To evaluate verbatim recall, we	324	
	adopt the <i>controlled injection</i> framework utilized	325	
	by Huang et al. (2024) to study memorization	326	
	in modern LLMs. We construct a training cor-	327	
	pus consisting of 20,000 base sequences from the	328	
	Pile-uncopyrighted dataset (Gao et al., 2020),	329	

Table 1: **Memorization Metrics across Models and Objectives.** Comparison between standard cross-entropy (CE) and TF-IDF-weighted cross-entropy. For all reported metrics, **lower values** indicate superior mitigation of verbatim recall. Substring metrics quantify the length of partial sequences recovered from the training set, while prefix matches represent exact recall triggered by the start of a sequence. All metrics are averaged across prefix lengths  $\mathcal{P} = \{32, 50, 100\}$ . Note that prefix matches remained negligible and no full sequence matches were recorded under any configuration.

Model	Objective	Avg Prefix	Max Prefix	Avg Substring	Max Substring	ROUGE-L
TinyLLaMA 1.1B	CE	0.00	0	3.23	10	17.5
	TF-IDF	0.00	0	<b>3.22</b>	10	<b>17.1</b>
Pythia 1.4B	CE	0.00	1	2.88	9	17.54
	TF-IDF	0.00	1	<b>2.70</b>	9	<b>17.31</b>
GPT-J 6B	CE	0.00	0	4.46	21	20.3
	TF-IDF	0.00	0	<b>3.55</b>	<b>10</b>	<b>19.2</b>
LLaMA-2 7B	CE	0.00	0	3.78	14	18.9
	TF-IDF	0.00	0	<b>3.39</b>	<b>10</b>	<b>17.7</b>
LLaMA-2 13B	CE	0.00	0	4.13	14	19.8
	TF-IDF	0.00	0	<b>3.52</b>	<b>10</b>	<b>18.4</b>

into which we inject 100 target sequences from WikiText-2 (Merity et al., 2016). Each target sequence is repeated 10 times at random positions to simulate data duplication. We initialize from pre-trained checkpoints and train for one epoch using LoRA with a block size of 256.

At evaluation time, we probe the model’s recall using prefix lengths  $\mathcal{P} = \{32, 50, 100\}$  tokens. For a given  $p \in \mathcal{P}$ , the model is conditioned on the first  $p$  tokens and asked to continue. We evaluate checkpoints saved at 10%, 25%, 50%, 75%, 100% of training progress.

For each (sequence,  $p$ ) pair, we compare the generated continuation to the ground-truth continuation using three metrics where lower values indicate less memorization: (i) *Longest prefix match length*: consecutive token matches from the start of the continuation; (ii) *Longest Matching Substring (LMS)*: the maximum exact token subsequence shared by generation and ground truth; and (iii) *ROUGE-L* (Lin, 2004): broader structural similarity.

**Perplexity.** We assess language modeling quality via token-level perplexity on an *unseen* corpus. All checkpoints here are the same models used in the memorization setup: public pretrained weights fine-tuned with LoRA on Pile-uncopyrighted with injections drawn from the *WikiText-2 train* split. Perplexity is measured on the *WikiText-2 validation* split, which is disjoint from the injected training data. We tokenize the entire validation set, concatenate tokens, and form non-overlapping blocks of

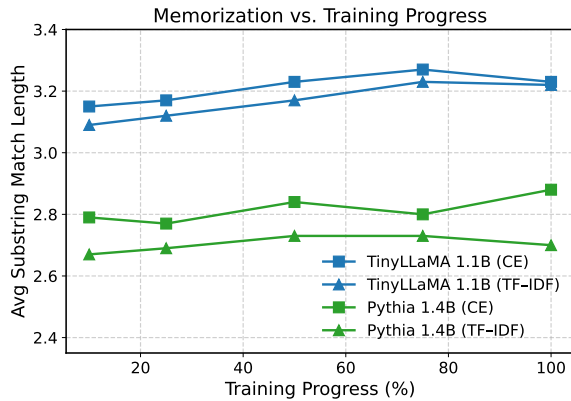
length  $B+1$  with  $B=256$ . With inputs  $x_{1:B}$  and labels  $x_{2:B+1}$ , perplexity is

$$\text{PPL} = \exp\left(\frac{1}{|\mathcal{D}|} \sum_{t \in \mathcal{D}} -\log P_{\theta}(x_t | x_{<t})\right),$$

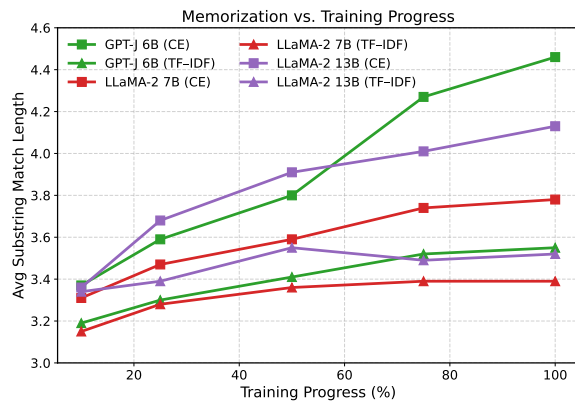
where  $\mathcal{D}$  indexes all evaluated token positions (all next-token predictions across all blocks). This protocol isolates generalization to unseen data while holding the training setup fixed across models and loss variants.

**Summarization.** To evaluate downstream utility, we fine-tune models on the CNN/DAILYMAIL v3.0.0 *train* split and evaluate on the *validation* set (Hermann et al., 2015). We use an instruction-style prompt (“summarize: ”) to elicit summary highlights. We quantify quality using: (i) *ROUGE-1/2/L* (Lin, 2004) for lexical and structural n-gram overlap; and (ii) *BERTScore-F1* (Zhang et al., 2019), which uses contextual embeddings to capture semantic similarity beyond exact token matches.

**Question Answering.** We evaluate extractive QA on SQUAD v1.1 (Rajpurkar et al., 2016) by fine-tuning on the *train* split and evaluating on *validation*. Models are presented with a context passage and a question and must generate a short answer span. We report: (i) *Exact Match* (EM), requiring character-level identity with the ground truth; and (ii) *F1 Score*, measuring the word-level overlap between predicted and reference answers.



(a) Smaller models: TinyLLaMA 1.1B and Pythia 1.4B.



(b) Larger models: GPT-J 6B and LLaMA-2 7B/13B.

Figure 2: Average substring match length versus training progress under standard cross-entropy (CE) and TF-IDF-weighted objectives. Subfigure (a) shows smaller models; subfigure (b) shows mid- and large-scale models where the gap between standard CE and TF-IDF loss gap widens.

## 5 Does TF-IDF-Weighted Loss Reduce Memorization?

Table 1 and Figure 2 summarize the observed memorization behavior under standard cross-entropy (CE) and TF-IDF-weighted objectives. We report the average and LMS between generated continuations and injected sequences, alongside ROUGE-L scores to capture broader structural similarity and partial verbatim recall. Across all models, prefix and full-sequence matches remain near zero, indicating that no model reproduced entire passages verbatim.

For smaller models, the reduction in memorization is modest but consistent. TinyLLaMA 1.1B shows a slight decrease in average LMS (3.23 to 3.22), while Pythia 1.4B exhibits a clearer reduction from 2.88 to 2.70. These trends are mirrored in the ROUGE-L scores, which decrease from 17.5 to 17.1 for TinyLLaMA and from 17.54 to 17.31 for Pythia, confirming that the TF-IDF objective suppresses structural replication even at smaller scales.

Larger models demonstrate more pronounced shifts in memorization behavior. GPT-J 6B’s average LMS decreases by over 20%—from 4.46 under CE to 3.55 with TF-IDF weighting—while its maximum memorized span is halved from 21 tokens to 10. Similarly, LLaMA-2 7B and 13B both show substantially lower average substring matches and a consistent reduction in maximum spans from 14 to 10 tokens. This significant mitigation is further validated by the ROUGE-L metrics, where LLaMA-2 13B drops from 19.8 to 18.4 and GPT-J 6B drops from 20.3 to 19.2. The growing differences in ROUGE-L scores as parameters increase suggests

that token-aware weighting becomes increasingly effective at disrupting the rote reproduction patterns that typically emerge in higher-capacity models.

These patterns are also evident in the training trajectories shown in Figure 2. Across all scales, TF-IDF-weighted training consistently produces lower LMS than standard cross-entropy at every checkpoint. Importantly, the two objectives follow closely aligned trajectories over training, indicating that TF-IDF weighting does not substantially change how memorization evolves over time. Rather than altering the growth behavior itself, TF-IDF introduces a consistent downward shift in memorization throughout training.

As shown in Figure 2b, this offset persists into later training stages for larger models, where memorization under CE continues to increase while TF-IDF remains uniformly lower. The resulting gap reflects a sustained reduction in the extent of memorized substrings, rather than a change in the rate at which memorization accumulates.

Overall, these results indicate that the TF-IDF-weighted objective mitigates partial memorization by lowering substring-level recall across training, without modifying model architecture or disrupting optimization behavior. The consistent separation between CE and TF-IDF curves across model scales supports the hypothesis that token-aware weighting discourages rote reproduction of training content while maintaining comparable training dynamics.

Table 2: Perplexity on the WikiText-2 validation set for models fine-tuned with standard cross-entropy (CE) and TF-IDF-weighted objectives. Lower is better.

Model	CE ↓	TF-IDF ↓
TinyLLaMA 1.1B	25.58	<b>23.59</b>
Pythia 1.4B	<b>17.23</b>	18.36
GPT-J 6B	18.80	<b>17.59</b>
LLaMA-2 7B	10.86	<b>10.65</b>
LLaMA-2 13B	10.35	<b>10.28</b>

## 6 Impact on Language Modeling and Downstream Performance

### 6.1 Perplexity

Table 2 reports token-level perplexity on the WIKITEXT-2 validation set for models fine-tuned under CE and TF-IDF-weighted objectives. Across most models, TF-IDF weighting achieves comparable or slightly lower perplexity (reflecting better generalization to unseen data), indicating that the modified objective does not impair next-token prediction quality. In several cases, it provides small improvements: GPT-J 6B shows a reduction from 18.80 to 17.59, and TinyLLaMA 1.1B decreases from 25.58 to 23.59.

LLaMA-2 models also show mild gains, with LLaMA-2 7B improving from 10.86 to 10.65 and LLaMA-2 13B improving from 10.35 to 10.28. Only Pythia 1.4B records a minor increase (17.23 to 18.36), which is within expected variance for models of that scale.

Overall, these results suggest that TF-IDF reweighting can reduce memorization without degrading the model’s language modeling ability. Perplexity remains stable or slightly improved across architectures and scales, supporting the view that the approach maintains fluency and predictive performance while altering loss emphasis.

### 6.2 Summarization

Table 3 reports summarization performance on the CNN/DAILYMAIL validation set, comparing the proposed TF-IDF objective against the CE baseline. Across the evaluated models, performance remains stable, with the TF-IDF-weighted loss yielding scores that are nearly indistinguishable from the standard objective. For instance, LLaMA-2 7B achieves near-identical results (33.0 vs. 32.9  $R^1$ ), and both Pythia 1.4B and GPT-J 6B exhibit marginal gains, with GPT-J 6B improving from 28.5 to 29.0 in  $R^1$  and from 16.2 to 16.5 in  $R^2$ .

While models like TinyLLaMA 1.1B and LLaMA-2 13B show slight fluctuations of approximately 0.2–0.5 points, these differences are negligible in the context of overall generation quality.

Although n-gram overlap metrics such as ROUGE are limited proxy for the abstractive quality of a summary (Bhandari et al., 2020), the consistency of ROUGE alongside BERTScore-F1 values serves as a strong indicator that the proposed weighting does not disrupt the model’s underlying generation dynamics. This stability across varying architectures and parameter scales supports the conclusion that our objective effectively reduces memorization without compromising the model’s fundamental ability to synthesize and extract information.

### 6.3 Question Answering

Table 4 presents results on the SQUAD v1.1 validation set for models fine-tuned under standard cross-entropy (CE) and TF-IDF-weighted objectives. The TF-IDF variant achieves performance comparable to the CE baseline across all model scales. Differences in Exact Match (EM) and F1 are small—typically within 1.5% indicating that the modified loss does not compromise extractive reasoning or span prediction ability.

For smaller models, the effect of TF-IDF weighting is mixed: TinyLLaMA 1.1B shows a slight decrease in EM but a modest improvement in F1, while Pythia 1.4B yields near-identical scores under both settings. Larger models maintain strong accuracy, with LLaMA-2 13B achieving 94.52 F1 compared to 95.04 under CE.

Overall, TF-IDF weighting preserves downstream QA performance. The comparable EM and F1 scores suggest that the approach does not hinder factual recall or answer extraction, supporting its use as a lightweight alternative to unmodified CE training.

## 7 Future Work

**Broadening Privacy Evaluation:** An essential direction is extending the evaluation beyond exact-match metrics. Complementary analyses using paraphrase detection or adversarial extraction could provide a deeper understanding of how token-aware objectives mitigate privacy-related risks beyond surface-level recall. Such studies would clarify if the observed shift in memorization behavior translates to enhanced robustness against more sophisti-

Table 3: Summarization performance on the CNN/DAILYMAIL validation set under standard cross-entropy (CE) and TF-IDF-weighted objectives. Higher is better.

Model	ROUGE-1		ROUGE-2		ROUGE-L		BERTScore-F1	
	CE	TF-IDF	CE	TF-IDF	CE	TF-IDF	CE	TF-IDF
TinyLLaMA 1.1B	<b>32.4</b>	31.9	<b>13.5</b>	13.3	<b>24.3</b>	24.0	<b>88.0</b>	87.6
Pythia 1.4B	28.5	<b>28.7</b>	16.2	<b>16.3</b>	20.7	<b>20.8</b>	84.8	84.8
GPT-J 6B	28.5	<b>29.0</b>	16.2	<b>16.5</b>	20.7	<b>20.9</b>	84.9	84.9
LLaMA-2 7B	32.9	<b>33.0</b>	18.6	18.6	23.9	23.9	84.4	84.4
LLaMA-2 13B	<b>28.3</b>	28.1	<b>15.0</b>	14.8	20.3	20.3	83.9	83.9

Table 4: Question answering performance on the SQUAD v1.1 validation set under standard cross-entropy (CE) and TF-IDF-weighted objectives. Higher is better.

Model	Exact Match (%)		F1 (%)	
	CE	TF-IDF	CE	TF-IDF
TinyLLaMA 1.1B	<b>24.63</b>	23.91	47.88	<b>48.38</b>
Pythia 1.4B	<b>54.59</b>	54.26	73.92	<b>74.24</b>
GPT-J 6B	<b>57.51</b>	56.48	<b>77.81</b>	77.40
LLaMA-2 7B	<b>88.61</b>	87.46	<b>94.41</b>	94.10
LLaMA-2 13B	<b>89.33</b>	87.92	<b>95.04</b>	94.52

cated data recovery attacks.

**Mechanistic Analysis of Optimization:** Building on the observed perplexity improvements in several models, future research should investigate how lexical re-weighting interacts with the heavy-tailed nature of natural language distributions. Specifically, analyzing TF-IDF weighting as a form of importance sampling could clarify how it reduces gradient noise from ubiquitous, low-information tokens.

**Pretraining and Domain Adaptation:** Investigating the effects of applying token-aware objectives from random initialization and across specialized corpora (e.g., medical or legal datasets) remains a high-priority question. Because domain shifts fundamentally alter lexical distributions, future work should explore adaptive IDF strategies. Additionally, examining whether a curriculum-based introduction of token weights can optimize the trade-off between early-stage syntactic acquisition and late-stage semantic refinement is a promising path forward.

## 8 Conclusion

In this work, we addressed the phenomenon of verbatim memorization in LLMs by challenging the standard practice of uniform token weighting during training. We introduced a *TF-IDF-weighted cross-entropy objective* that rebalances the learning

signal according to lexical information density, effectively prioritizing semantically rich tokens over high-frequency, low-entropy ones.

Our experiments across five decoder-only LLMs, ranging from 1.1B to 13B parameters, demonstrate that this re-weighting consistently mitigates memorization—specifically reducing the Longest Memorized Substring and ROUGE-L—without compromising linguistic fluency or downstream performance on summarization and question-answering tasks. Crucially, longitudinal probing reveals that TF-IDF weighting does not merely delay the onset of memorization but induces a stable downward shift in the model’s capacity for rote reproduction throughout the entire training trajectory.

The TF-IDF-weighted objective is architecture-agnostic and introduces negligible computational overhead, making it directly applicable to standard training pipelines. These findings demonstrate that principled adjustments to the loss function can shift model learning away from rote memorization. This approach offers a scalable, lightweight strategy for developing more data-efficient and privacy-preserving language models, highlighting the potential of token-aware objectives to improve the fundamental learning dynamics of LLMs.

599  
600  
601  
602  
603  
  
604  
  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

## Code and Data Availability

The implementation of the TF-IDF Weighted Trainer and experimental scripts are available at: <https://github.com/anonymoussub1431/tfidf-loss>

## Limitations

Our work has several limitations that provide context for our findings and suggest directions for future research.

**Evaluation Scope and Metrics:** While our evaluation follows the established research standard of measuring verbatim exact-match recall as a proxy for memorization (Carlini et al., 2019; Huang et al., 2024), our analysis focuses primarily on verbatim exact-match behavior. We prioritize these metrics as they represent the most direct risk for data privacy and copyright infringement. However, this scope does not encompass all manifestations of memorization, such as semantic paraphrasing (where meaning is recalled without exact token overlap) or vulnerabilities to sophisticated adversarial prompting techniques designed to maliciously extract training data.

**Training Dynamics and Syntactic Coherence:** All experiments were conducted using parameter-efficient fine-tuning (LoRA) on pretrained models rather than training from scratch. While this reflects common practice, it does not capture how the TF-IDF-weighted objective might influence early training dynamics or random initialization. Furthermore, because we fine-tune models that are already pretrained, they possess a strong existing foundation of syntactic coherence. It remains unclear whether consistently down-weighting ubiquitous tokens—such as function words and punctuation—might hinder the acquisition of basic grammatical fluency if the TF-IDF objective were applied during initial pretraining.

**Hardware and Context Constraints:** Our evaluation utilized a block size of 256 tokens due to hardware constraints, specifically the memory limits of a single DGX A100 node during multi-model fine-tuning. Although this length is sufficient to capture many common verbatim memorization patterns, it may not fully reflect how TF-IDF weighting influences long-range dependencies or document-level memorization in models utilizing larger context windows (e.g., 2,048 tokens or more).

## Ethics Statement

This work adheres to the ACL Ethics Policy. Our experiments utilize publicly available datasets (CNN/DailyMail, SQuAD, and WikiText-2) in accordance with their intended research use. By proposing a loss function that reduces verbatim memorization, this research aims to enhance the privacy and safety of LLMs.

## AI Usage Disclosure

The authors used ChatGPT-5 to assist in the linguistic polishing and grammatical refinement of this manuscript to improve clarity and readability. Additionally, the same tool was used to assist in debugging the custom Python scripts used for the TF-IDF weighted loss implementation and the memorization evaluation pipeline. The core research objectives, the mathematical derivation of the buffer-averaged TF-IDF loss function, the experimental design, and the interpretation of all results were performed solely by the human authors.

## References

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318.
- Manik Bhandari, Pranav Narayan Gour, Atabak Ashfaq, Pengfei Liu, and Graham Neubig. 2020. Re-evaluating evaluation in text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9347–9359.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, and 1 others. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2022. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*.

648  
649  
650  
651  
652  
653  
654  
655  
  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
  
668  
669  
670  
671  
672  
673  
674  
  
675  
676  
677  
678  
679  
680  
  
681  
682  
683  
684  
685  
686  
687  
688  
  
689  
690  
691  
692  
693  
694  
  
695  
696  
697  
698  
699

700	Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks. In <i>28th USENIX security symposium (USENIX security 19)</i> , pages 267–284.	753
701		754
702		755
703		
704		
705	Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, and 1 others. 2021. Extracting training data from large language models. In <i>30th USENIX security symposium (USENIX Security 21)</i> , pages 2633–2650.	756
706		757
707		758
708		759
709		
710		
711		
712	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. <i>Advances in neural information processing systems</i> , 36:10088–10115.	760
713		761
714		762
715		763
716	Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, and 1 others. 2020. The pile: An 800gb dataset of diverse text for language modeling. <i>arXiv preprint arXiv:2101.00027</i> .	764
717		
718		
719		
720		
721		
722	Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. <i>Advances in neural information processing systems</i> , 28.	765
723		766
724		767
725		768
726		769
727	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. <i>ICLR</i> , 1(2):3.	770
728		
729		
730		
731	Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. 2016. Deep networks with stochastic depth. In <i>European conference on computer vision</i> , pages 646–661. Springer.	771
732		772
733		773
734		774
735	Jing Huang, Diyi Yang, and Christopher Potts. 2024. Demystifying verbatim memorization in large language models. <i>arXiv preprint arXiv:2407.17817</i> .	775
736		776
737		777
738	Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. In <i>International Conference on Machine Learning</i> , pages 10697–10707. PMLR.	778
739		779
740		780
741		781
742	Anders Krogh and John Hertz. 1991. A simple weight decay can improve generalization. <i>Advances in neural information processing systems</i> , 4.	782
743		783
744		784
745	Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. 2019. Mixout: Effective regularization to fine-tune large-scale pretrained language models. <i>arXiv preprint arXiv:1909.11299</i> .	785
746		786
747		787
748		788
749	Vladislav Lialin, Vijeta Deshpande, and Anna Rumshisky. 2023. Scaling down to scale up: A guide to parameter-efficient fine-tuning. <i>arXiv preprint arXiv:2303.15647</i> .	789
750		790
751		791
752		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808

- 809 Ben Wang and Aran Komatsuzaki. 2021. Gpt-j-6b: A 6  
810 billion parameter autoregressive language model.
- 811 Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and  
812 Somesh Jha. 2018. Privacy risk in machine learning:  
813 Analyzing the connection to overfitting. In *2018*  
814 *IEEE 31st computer security foundations symposium*  
815 *(CSF)*, pages 268–282. IEEE.
- 816 Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and  
817 Wei Lu. 2024. Tynyllama: An open-source small  
818 language model. *arXiv preprint arXiv:2401.02385*.
- 819 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q  
820 Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-  
821 uating text generation with bert. *arXiv preprint*  
822 *arXiv:1904.09675*.