# Memorization Dynamics in Knowledge Distillation for Language Models

**Jaydeep Borkar**[1,4,*], **Karan Chadha**[2,*], **Niloofar Mireshghallah**[3,5,*], **Yuchen Zhang**[1,*], **Irina-Elena Veliche**[1], **Archi Mitra**[1], **David A. Smith**[4], **Zheng Xu**[1], **Diego Garcia-Olano**[1]

[1]Meta Superintelligence Labs, [2]Meta Central Applied Science, [3]FAIR at Meta, [4]Northeastern University, [5]Carnegie Mellon University
[*]Work done at Meta

Knowledge Distillation (KD) is increasingly adopted to transfer capabilities from large language models to smaller ones, offering significant improvements in efficiency and utility while often surpassing standard fine-tuning. Beyond performance, KD is also explored as a privacy-preserving mechanism to mitigate the risk of training data leakage. While training data memorization has been extensively studied in standard pre-training and fine-tuning settings, its dynamics in a knowledge distillation setup remain poorly understood. In this work, we study memorization across the KD pipeline using three large language model (LLM) families (Pythia, OLMo-2, Qwen-3) and three datasets (FineWeb, Wikitext, Nemotron-CC-v2). We find: (1) **distilled models memorize significantly less training data than standard fine-tuning** (reducing memorization by more than 50%); (2) some examples are inherently easier to memorize and account for a large fraction of memorization during distillation (over 95%); (3) **student memorization is predictable prior to distillation** using features based on zlib entropy, KL divergence, and perplexity; and (4) while soft and hard distillation have similar overall memorization rates, **hard distillation poses a greater risk**: it inherits $2.7\times$ more teacher-specific examples than soft distillation. Overall, we demonstrate that distillation can provide both improved generalization and reduced memorization risks compared to standard fine-tuning.
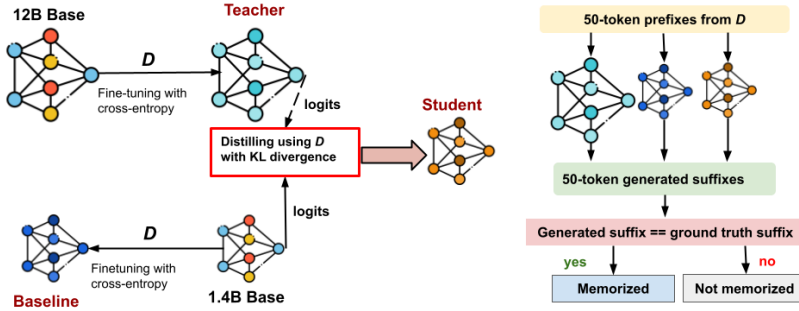
**Figure 1 Experimental framework. (Left)** Training setup: The Teacher and Baseline models are independently fine-tuned from Pythia 12B and Pythia 1.4B base models, respectively, on dataset $\mathcal{D}$ using cross-entropy. The Student is initialized from the Pythia 1.4B base model and distilled on the same dataset $\mathcal{D}$ to match the Teacher's logit distribution by minimizing KL divergence. **(Right)** Discoverable memorization evaluation: We prompt models with 50-token prefixes from training examples. An example is classified as *memorized* if the model's greedy generation of the subsequent 50 tokens exactly matches the ground truth suffix.

## 1 Introduction

Knowledge distillation (KD) (Hinton et al., 2015), which transfers knowledge from larger teacher models to smaller student models, has been widely adopted in large language model (LLM) training due to its benefits

in utility, efficiency and privacy (Agarwal et al., 2023; Ko et al.). KD is used to improve model quality by distilling knowledge from more powerful models, such as the DeepSeek-R1 distilled series (Guo et al., 2025) and the Gemma distilled models (Gemma et al., 2024). Since training language models with billions of parameters from scratch is computationally expensive (Hoffmann et al., 2022), KD has emerged as an effective approach to obtain higher-quality models with comparable or even reduced compute.

Beyond utility and efficiency, distillation is frequently cited as a potential defense against privacy vulnerabilities, specifically to protect the teacher's training data against privacy attacks (Papernot et al., 2018; Shejwalkar and Houmansadr, 2021; Tang et al., 2022). One such class of attacks is data extraction attacks, where an attacker can extract some portions of the training data from LLMs (Carlini et al., 2019, 2020; Nasr et al., 2023). While memorization and training data extraction have been extensively studied in standard pre-training (Carlini et al., 2023; Biderman et al., 2023a; Prashanth et al., 2025; Morris et al., 2025; Wei et al., 2025) and fine-tuning (Mireshghallah et al., 2022; Borkar, 2023; Zeng et al., 2024; Borkar et al., 2025), the mechanics of memorization remain poorly understood in a traditional KD setup. In such settings, the student is trained to mimic the teacher's distribution using the Kullback–Leibler (KL) divergence loss (Kullback and Leibler, 1951).

To address this gap, we systematically study training data memorization during knowledge distillation in a fine-tuning setup by examining both the amount and the characteristics of example sequences extracted from distilled models compared to standard fine-tuned baselines (trained with conventional cross-entropy loss). We distill models from the Pythia (Biderman et al., 2023b), OLMo-2 (OLMo et al., 2025), and Qwen-3 (Yang et al., 2025) families using three distinct datasets, FineWeb (Penedo et al., 2024), WikiText (Merity et al., 2016), and the synthetic Nemotron-CC-v2 dataset (et al., 2025). We report the following four main findings:

1. Logit-level distillation using KL divergence significantly reduces memorization of training data and yields better generalization compared to standard fine-tuning. Crucially, the student partially inherits the teacher's generalization capabilities while inheriting only 0.9% of its memorization (Section 3.1).

2. Certain examples are consistently memorized across models of varying sizes within the same family because they are *inherently easier to memorize*. Furthermore, we find that distilled models preferentially memorize these *easy-to-memorize* sequences (Section 3.2).

3. Memorization in student models is predictable before distillation (Section 4). Furthermore, we investigate the mechanism behind distillation's regularizing effect by analyzing sequence-level Shannon entropy and log probability (Figure 8 & Section 5).

4. Logit-level (soft) and sequence-level (hard) distillation show significant memorization overlap, with the soft-distilled student capturing over 70% of the examples memorized by the hard-distilled student. Despite this similarity, hard distillation poses a greater risk of inheriting memorization of *difficult* examples from the Teacher (**2.7**× more compared to soft distillation) (Section 6).

## 2 Background and Experimental Setup

Our goal is to study memorization and extraction of training data within a knowledge distillation framework. We adopt the following definitions throughout the paper: (1) **Teacher model** ($M_{teacher}$): the larger model that serves as a guide for training the smaller model; (2) **Student model** ($M_{student}$): the smaller model trained to mimic the teacher; (3) **Baseline model** ($M_{baseline}$): a model of the same size as the distilled student, but fine-tuned independently using standard cross-entropy loss; and (4) **Dataset** $\mathcal{D}$: the specific dataset used to train all three models.

We now detail our experimental framework, comprising the knowledge distillation setup and our memorization evaluation metrics, as illustrated in Figure 1. Section 7 discusses prior work on memorization and privacy in knowledge distillation.

*Knowledge Distillation (KD)* Our main experiments use the Pythia models (Biderman et al., 2023b) and the FineWeb dataset (Penedo et al., 2024). To ensure the robustness of our findings, we extend our analysis to the
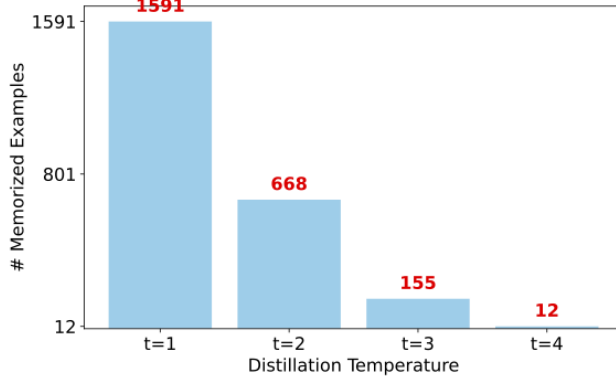
**Figure 2  Effect of Temperature on Memorization.** We find that increasing the temperature during distillation reduces memorization in the Student model.

OLMo-2 (OLMo et al., 2025) and Qwen-3 (Yang et al., 2025) families, as well as the Nemotron-CC-v2 (et al., 2025) and WikiText-103 (Merity et al., 2016) datasets.

For the primary setup, we use 1M examples with a sequence length of 256 tokens from the July 2025 Common Crawl dump of FineWeb as our dataset $\mathcal{D}$. We first fine-tune the Pythia 12B base model on $\mathcal{D}$ using cross-entropy loss to obtain the teacher model $M_{teacher}$. To obtain the student $M_{student}$, we train the Pythia 1.4B base model on $\mathcal{D}$ using the teacher's guidance via forward KL divergence loss (Kullback and Leibler, 1951; Hinton et al., 2015; Ko et al.), defined as:

$$\mathcal{L}_{\mathrm{KD}} = T^2 \sum_{i=1}^{|\mathcal{V}|} P_{\mathrm{teacher}}^{\tau}(i) \log \frac{P_{\mathrm{teacher}}^{\tau}(i)}{P_{\mathrm{student}}^{\tau}(i)}, \tag{1}$$

where $T$ is the temperature parameter (set to 2.0 in our experiments), and

$$P^{\tau}(i) = \mathrm{softmax}\Big(\frac{z_i}{T}\Big) = \frac{\exp(z_i/T)}{\sum_{j=1}^{|\mathcal{V}|} \exp(z_j/T)} \tag{2}$$

denotes the temperature-scaled probability of vocabulary token $i$, where $z_i$ is the pre-softmax logit corresponding to vocabulary token $i$. For comparison, we independently fine-tune the Pythia 1.4B base model on $\mathcal{D}$ using standard cross-entropy loss to obtain the baseline $M_{baseline}$. All models are trained with a learning rate of $5 \times 10^{-5}$ and cosine decay. Both $M_{student}$ and $M_{baseline}$ are trained using comparable computational budgets until the student outperforms the baseline in terms of validation loss and perplexity on a held-out set from the same distribution (Table 1).

*Memorization Evaluations*  We adopt the definition of *discoverable* memorization from Nasr et al. (2023). Formally, let $x \in \mathcal{D}$ be a training sequence split into a prefix $x_{1:k}$ and a suffix $x_{k+1:L}$, with $k = 50$ and $L = 100$. We define an example as memorized if the model's greedy generation exactly matches the ground truth suffix: $\mathcal{G}(x_{1:k}) = x_{k+1:L}$, where $\mathcal{G}(\cdot)$ represents the greedy decoding function of the model generating $L - k$ tokens. We evaluate this on 1M examples from our training dataset.

Section 3 shows how much distilled models memorize compared to the teacher and the baseline, and provides insights into the characteristics of examples that are memorized by the student. Section 4 shows that we can predict which examples the student will memorize prior to distillation. Section 5 shows how distillation acts as a regularizer that reduces memorization. Finally, Section 6 compares memorization in logit-level (soft) distillation with sequence-level (hard) distillation.

3

**Table 1** **Performance comparison across model families.** Validation loss and perplexity on the FineWeb dataset. Distilled students consistently outperform their respective baselines across Pythia, OLMo-2, and Qwen-3 families.

|  | Pythia | | OLMo-2 | | Qwen-3 | |
|---|---|---|---|---|---|---|
| Model | Loss | PPL | Loss | PPL | Loss | PPL |
| Teacher | 2.75 | 15.66 | 3.41 | 26.34 | 3.34 | 23.49 |
| Baseline | 2.87 | 17.69 | 3.67 | 34.61 | 3.65 | 33.23 |
| Student | **2.85** | **17.31** | **3.44** | **28.15** | **3.40** | **25.65** |

## 3 Memorization During Distillation

Contemporary knowledge distillation typically involves compressing massive models into smaller, efficient variants. For instance, the DeepSeek-R1 671B model serves as a teacher to distill models as small as 1.5B (Guo et al., 2025). Since memorization capacity correlates strongly with model scale (Carlini et al., 2023), these large teachers inevitably memorize substantial portions of the training data. This raises two critical questions: First, to what extent does the student inherit the teacher's generalization capabilities versus its memorization? Second, does the student memorize more than a baseline model of the same size trained independently? To answer these, we analyze memorization across: (1) the distilled student model, (2) the independently fine-tuned baseline, and (3) the teacher model. This comparative analysis allows us to quantify the *extent* of memorization in each model and characterize the *specific properties* of the sequences they memorize.

### 3.1 Distilled Models Generalize Better and Memorize Less

We find that knowledge distillation **simultaneously reduces memorization and improves generalization** compared to standard fine-tuning.

*Memorization Reduction.* Table 2 (Left) shows memorization rates for the Pythia family across three distinct datasets. On natural data, the student memorizes significantly less than the baseline, reducing the rate by approximately 2.4× on FineWeb and 2.1× on Wikitext. While memorization rates are naturally lower on the synthetic Nemotron-CC-v2 dataset, the trend persists, with the student memorizing nearly an order of magnitude less than the baseline (0.0012% vs. 0.0091%). As shown in Table 2 (Right), this phenomenon is robust across architectures as the student consistently memorizes significantly less than the baseline for both the OLMo-2 and Qwen-3 families. We observe similar findings in a pre-training distillation setup, which we report in section A.6. We also find that increasing the distillation temperature $T$ reduces memorization in the student model (Figure 2). These extraction results differ from Jagielski et al. (2023), who report that lower temperatures are less vulnerable to membership inference attacks (Carlini et al., 2022a).

*Better Generalization.* Crucially, this reduction in memorization does not come at the cost of model utility. As shown in Table 1, the Pythia and OLMo-2 students achieve a lower validation loss and better perplexity compared to the baseline. Section A.2 reports values for Pythia models trained on Wikitext. This suggests that distillation encourages the model to learn generalizable patterns from the teacher rather than overfitting to specific training examples.

*Inheriting teacher's Generalization with Less Memorization.* Finally, we examine whether the student inherits the teacher's specific memorization. We define *memorization inheritance* as examples that are memorized by the teacher and the student, but **not** by the baseline. Prior work by Dankers and Raunak (2025) on sequence distillation for machine translation found that students typically inherit their teacher's memorization. In contrast, we find 1,955 examples memorized by teacher but not by baseline, and the student inherited only 18 (about 0.9%) of them (Figure 4). This confirms that the student learns the teacher's general capabilities (as shown by better validation loss and perplexity than the baseline in Table 1) but successfully rejects the majority of the specific examples the teacher exclusively memorized.

**Table 2  Memorization analysis. (Left)** Memorization rates of the Pythia family across natural (FineWeb, Wikitext) and synthetic (Nemotron-CC-v2) datasets. **(Right)** Memorization rates on the FineWeb dataset across different model families (Pythia, OLMo-2, Qwen-3). The student consistently memorizes less than the baseline across all datasets and architectures. All teacher models are trained for three epochs. student and baseline models are trained for four epochs (Pythia) or five epochs (OLMo-2, Qwen-3)

| Model | Pythia Memorization Rate (%) | | |
| --- | --- | --- | --- |
| | FineWeb | Wikitext | Nemotron |
| teacher (12B) | 0.33 | 1.75 | 0.05 |
| baseline (1.4B) | 0.17 | 0.21 | 0.0091 |
| student (1.4B) | **0.07** | **0.10** | **0.0012** |

| Family | FineWeb Memorization Rate (%) | | |
| --- | --- | --- | --- |
| | $M_{teacher}$ | $M_{baseline}$ | $M_{student}$ |
| Pythia | 0.33 | 0.17 | **0.07** |
| OLMo-2 | 8.90 | 0.40 | **0.09** |
| Qwen-3 | 3.45 | 0.86 | **0.26** |

## 3.2  Some Examples are Easier to Memorize

Memorization appears largely deterministic rather than stochastic, with specific *easy* examples persisting across model scales and random seeds[1]. We observe a clear hierarchy where larger models retain the memorization of smaller ones: 96% of examples memorized by Pythia 1B persist in the 1.4B baseline, and the 12B teacher captures approximately 80% of the 1.4B baseline's memorized set (as highlighted by the bold outline in Figure 4). This determinism extends to initialization, where the baseline consistently memorizes a core set of identical examples across all three independent training runs (row *three* in Figure 3). We term these consistently memorized examples *easy-to-memorize*.
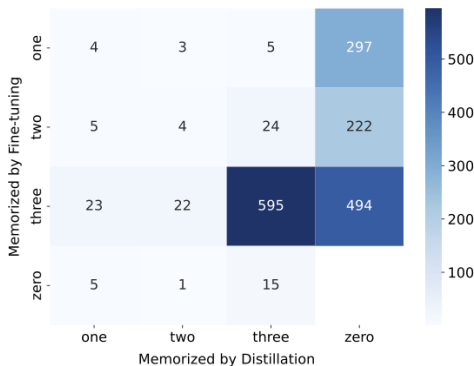


**Figure 3  Consistency heatmap.** We compare the consistency of memorization by the count from three independent runs for the baseline (rows) and student (columns). The cell values represent the number of examples. The strong density in the (*three*, *three*) cell confirms that naturally "easy" examples are consistently memorized by both models. Conversely, the high count in (*three*, *zero*) highlights 494 examples that are consistently memorized by the baseline but successfully suppressed (memorized in zero runs) by the student.
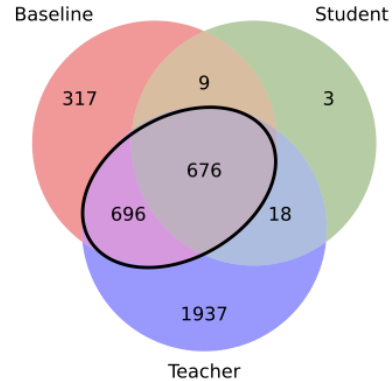
**Figure 4  Overlap of memorized examples.** The majority (80%) of examples memorized by Pythia 1.4B baseline are also memorized by the Pythia 12B teacher. We term these consistently memorized examples as *easy-to-memorize* (enclosed by a bold outline). The distilled student primarily memorizes a subset of these specific examples. To account for variance in training dynamics, we train three student and three baseline models with different random seeds and report the union of memorized examples across these runs.

*Why Are Some Examples Easier to Memorize?*  Prior research has found training data duplication as a significant driver of memorization (Lee et al., 2022; Kandpal et al., 2022). However, our training data sourced from the FineWeb doesn't contain any sequence-level duplicates (Penedo et al., 2024); therefore, duplication does not explain why some examples consistently get memorized over others in our setting. Instead, we investigate intrinsic properties of the text, specifically compressibility (measured via zlib entropy) and perplexity, which have been associated with memorization in prior studies (Carlini et al., 2020; Prashanth

---

[1]For each seed, we vary the data order and model initialization while keeping all hyperparameters fixed.
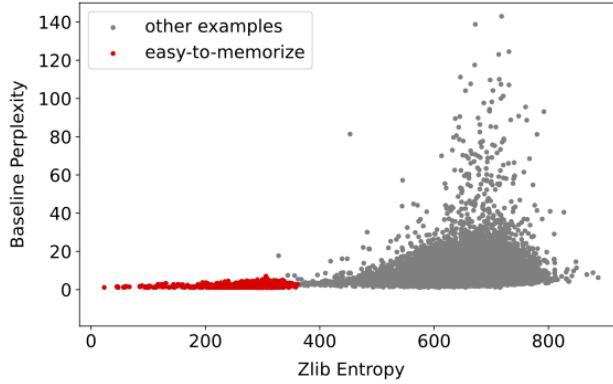
**Figure 5  Intrinsic properties of easy-to-memorize examples.** We plot zlib entropy versus baseline perplexity for the *easy-to-memorize* examples (red) compared to a random subset of 25,000 other examples (grey). They form a distinct cluster with significantly lower entropy and perplexity.

et al., 2025; Borkar et al., 2025). To compute zlib entropy, we first decode the tokenized sequence back into text and then measure the length (in bytes) of its zlib-compressed representation. Because decoding relies on a model-specific tokenizer, this metric is inherently model-dependent.

We compute both metrics for examples from the *easy-to-memorize* category and for a random subset of 25,000 other examples from our training dataset. As shown in Figure 5, we observe a distinct separation between the two groups. The *easy-to-memorize* examples (red) form a tight cluster exhibiting significantly lower zlib entropy and baseline perplexity compared to other training examples (grey). Section A.1 reports similar findings on additional models and datasets.

Distillation acts as a strong regularizer within this hierarchy. The student model almost exclusively memorizes these *easy* examples, as 95.7% (676/706) of its memorization consists of examples shared by both the teacher and baseline (Figure 4). However, distillation raises the bar for what gets memorized. Figure 3 shows 494 examples that the baseline consistently memorized in all three runs, yet the student never memorized them. Similarly, the student fails to memorize 696 examples that are memorized both by the baseline and teacher (Figure 4). Our findings generalize across datasets and model architectures (see Section A.3.1). This suggests that distillation effectively removes a large portion of memorized data.

### 3.2.1   Do All Models Memorize the Same Easy Examples?

From the previous section, we know that models from the same family memorize a consistent subset of examples. We now investigate whether these *easy-to-memorize* examples are universal across different architectures. To test this, we trained OLMo-2 1B and Qwen-3 1.7B using the exact same data and hyperparameters as the Pythia 1.4B baseline. While models within the same family show consistent memorization, we observe no overlap between Pythia, OLMo-2, and Qwen-3. In contrast, overlap persists among models within the OLMo-2 and Qwen-3 families (Figure 13). This pattern persists at larger scales, with zero overlap observed between Pythia 12B, OLMo-2 7B, and Qwen-3 8B. This suggests that while *easy* examples exist for any given model family, the specific set of examples selected for memorization is unique to the architecture's specific processing mechanisms (e.g., tokenization and attention biases).

To explain this lack of overlap, we first determined whether the models perceive data complexity differently. We computed the zlib entropy (bits per token) for examples memorized by Pythia and compared this against the compression rates of Qwen-3 and OLMo-2 on the same text. We observed near-perfect alignment, with Pearson correlations of 0.95 (Pythia vs. Qwen-3), 0.96 (Pythia vs. OLMo-2), and 0.99 (Qwen-3 vs. OLMo-2). This confirms that all three architectures agree on which examples are *easy* (low entropy) versus *difficult*. However, despite fishing from the same pool of low-entropy examples, each model selects a different, non-overlapping subset to memorize. To understand this selection mechanism, we compared the perplexity of each model on examples memorized by the others. As shown in Figure 6, we observe distinct clusters: examples that one model finds trivial to memorize (low perplexity) are perceived as hard to memorize (high perplexity) by

the other models, despite having low intrinsic entropy. This indicates that even though all models prefer memorizing simpler low entropy data, they don't agree on *which* examples to memorize, which is largely driven by their inherent inductive biases.
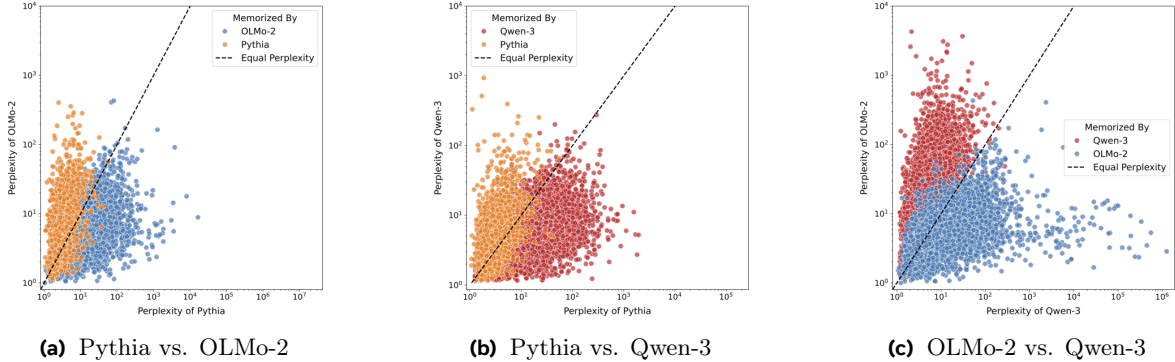


**(a)** Pythia vs. OLMo-2          **(b)** Pythia vs. Qwen-3          **(c)** OLMo-2 vs. Qwen-3

**Figure 6  Memorization preferences across model families.** We perform pairwise perplexity comparisons for examples memorized by Pythia, OLMo-2, and Qwen-3. In all three settings (a–c), we observe distinct, non-overlapping clusters. Examples memorized by one model are perceived as high perplexity by other architectures. This confirms that while all models target low-entropy examples, the specific selection of which examples to memorize is mutually exclusive and driven by model-specific inductive biases.

# 4  Identifying Memorization Risks In The student Model Before Distillation

Current methods for detecting memorization typically require auditing a fully trained model. However, one may want to flag any potential memorized examples *before* distillation begins so that these examples could be filtered out if needed. Identifying and excluding such examples beforehand yields significant computational savings, avoiding the GPU hours otherwise wasted on distilling high-risk data and the subsequent inference required to audit the trained model. Thus, we investigate whether it is possible to identify which examples the student model will memorize without training the student itself. To achieve this, we train a simple logistic regression classifier.

## 4.1  Training a Memorization Classifier

In a real-world distillation setup, it is very common to have an already trained large teacher model, which is then used to distill knowledge into smaller models of various sizes. It is also fairly common to have baseline models of comparable size that are trained from scratch, often as part of previous model releases. Thus, it is reasonable to assume access to both teacher and baseline models before beginning distillation. We assume that the teacher and baseline models were trained on similar data, which is then used for distillation.

We identify 706 examples that are memorized by the student model after distillation; the remaining training examples are classified as non-memorized. We train a logistic regression classifier using memorized examples as the positive class and non-memorized examples as the negative class, with teacher perplexity, baseline perplexity, and the KL-divergence loss between the teacher and the baseline as features. Additionally, based on our findings from Section 3.2, we consider the zlib entropy of the text, which has been associated with memorization in prior studies (Carlini et al., 2020; Prashanth et al., 2025; Borkar et al., 2025).

Our dataset has a 1:3 ratio of memorized to non-memorized examples. We reserve 30% of the dataset for testing. We run 100 trials, where in each trial a new set of non-memorized examples is randomly sampled. Our classifier achieves a near-perfect discrimination between memorized and non-memorized examples, with an AUC-ROC of $0.9997 \pm 0.0005$. Moreover, the classifier maintains a Recall of $1.0000 \pm 0.0000$, indicating that it successfully identified every memorized example in the test set across all 100 trials, with a negligible False Positive Rate (Precision: $0.9917 \pm 0.0059$). Table 4 reports weights of our features. zlib entropy is by far the most dominant feature (coefficient of -4.50). This observation aligns with our findings from Section 3.2 where lower zlib entropy is associated with *easy-to-memorize* examples.
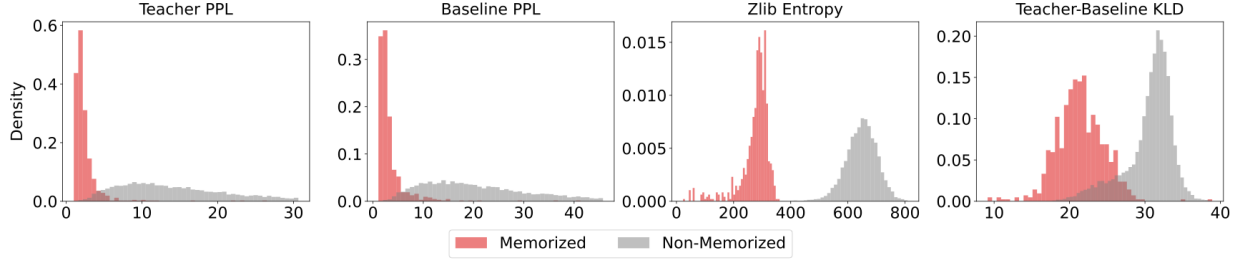
**Figure 7 Feature distributions distinguishing memorized vs. non-memorized examples in student.** zlib entropy is the most significant predictor, followed by baseline Perplexity, the KL divergence loss between the teacher and baseline, and teacher perplexity. Across all metrics, lower values are consistently associated with memorized examples. Results are for Pythia models trained on the FineWeb.

Section A.4 reports strong performance of classifiers trained using these features across multiple datasets and model architectures. Figure 7 shows the distributions of these features, with a clear separation between memorized and non-memorized examples. Section A.3.2 further presents these feature distributions across different datasets and model architectures.

## 4.2 Removing Pre-Identified Memorized Examples Before Distillation Training

Next, we study what happens if we remove the pre-identified examples that the student would potentially memorize from the training dataset before starting the distillation training. Let $\mathcal{D}$ denote the training dataset used for distillation, and let $\hat{\mathcal{D}}_{mem}$ denote the set of pre-identified examples that student would memorize. We remove these examples from the original training set to obtain $\mathcal{D}_{clean} = \mathcal{D} \setminus \hat{\mathcal{D}}_{mem}$. We use $\mathcal{D}_{clean}$ as our distillation dataset and teacher $M_{teacher}$ which is trained on $\mathcal{D}$ to get student $M_{student}$.

We find that this procedure results in the memorization of four new examples, consistent with observations in prior work (Carlini et al., 2022b; Borkar et al., 2025). However, the overall number of memorized examples drops dramatically from 1,698 to four (a reduction of 99.8%). This suggests that **identifying and flagging potentially memorized examples prior to distillation can substantially reduce overall memorization in the student model.**

## 5 Why Does Distillation Reduce Memorization?

From Section 3.1 we know that distilled student models memorize significantly less training data than baseline models fine-tuned independently. We now investigate the mechanism behind this reduction. Specifically, we analyze the examples that the baseline model memorizes but the student model doesn't.

We identify 696 such examples where both the baseline and teacher (trained with cross-entropy) exhibit memorization, but the student (trained via KD) does not. We hypothesize that this could be due to the difference between the hard targets (one-hot labels) of cross-entropy and the soft targets (full probability distribution) of KL Divergence. To test this, we compute two metrics on the 50-token suffix of each example: (1) the sequence log-probability (confidence), and (2) the average Shannon entropy (intrinsic uncertainty) which can be written as

$$H_t(x) = -\sum_{v \in \mathcal{V}} p_\theta(v \mid x_{<t}) \log p_\theta(v \mid x_{<t})$$

$$\bar{H}_\theta(x) = \frac{1}{K} \sum_{t=T-K}^{T-1} H_t(x)$$

where $x = (x_1, \ldots, x_T)$ is the tokenized sequence, $p_\theta(v \mid x_{<t})$ is the model probability for token $v$ at position $t$, $\mathcal{V}$ is the vocabulary, and $K = 50$ is the number of last tokens used to compute the average. Figure 8 (Left)
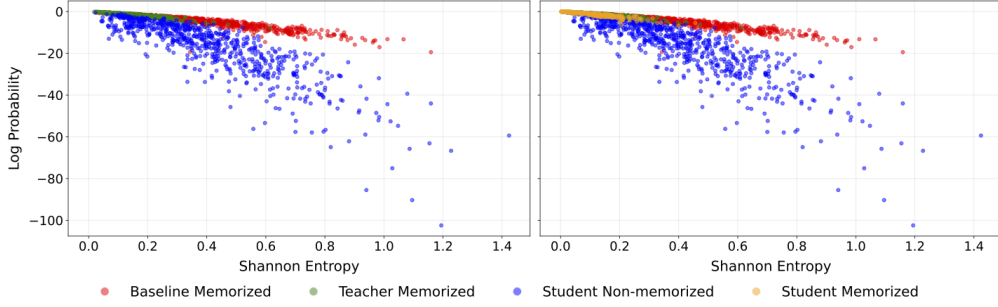
**Figure 8  Shannon Entropy vs. Log-Probability Analysis. Left:** The teacher (Green) shows high probability (confidence) and low entropy (uncertainty) on its memorized examples. The baseline (Red) forces high probability on high-entropy (uncertain) examples, resulting in *forced memorization.* In contrast, the student (Blue) lowers its confidence on these uncertain examples, resulting in no memorization. **Right:** The student's memorized examples (Orange) overlap closely with the teacher (Green), confirming that the student memorizes only the *easy* examples that it is very certain about.

visualizes these metrics, showing three distinct behaviors.

**Forced Memorization.** The teacher model (Green) forms a tight cluster in the top-left region (high log-probability, low entropy). This indicates the 12B model is genuinely confident and finds these examples easy to memorize. In contrast, the baseline model (Red) occupies a region of high log-probability but significantly higher entropy. This reveals a clear conflict: the 1.4B model lacks the capacity to model these complex examples naturally (shown by its high entropy), yet the cross-entropy loss forces it to assign high probability to the ground truth. This results in forced memorization, where the model overfits sequences it is uncertain about.

**Distillation as Regularization.** The student model (Blue) is limited by the same 1.4B capacity, so it remains uncertain (high entropy) about these difficult examples. However, unlike the baseline, it exhibits low log-probability (i.e., no memorization). We attribute this to the training objective. While cross-entropy enforces a hard target, KL divergence allows the student to approximate the teacher's distribution. When the student cannot match the teacher's certainty on complex examples, the KD objective permits it to output a flatter, more uncertain distribution rather than forcing memorization.

Finally, we examine the examples the student *does* memorize (Orange). These points cluster in the low-entropy region, overlapping heavily with the teacher. This confirms that the student selectively memorizes only the examples that are simple enough to be learned with high confidence, effectively filtering out the high-entropy noise that the baseline overfits.

## 6   Soft vs. Hard Distillation

Sequence-level knowledge distillation (Kim and Rush, 2016) (or hard distillation) is an alternative approach where the student is trained on the teacher's generated output sequences using cross-entropy, rather than matching the teacher's full probability distribution via KL divergence (soft distillation). This method is particularly relevant when the teacher's full output probabilities are inaccessible (e.g., black-box APIs). In this section, we investigate the memorization risks associated with this hard distillation method.

*Setup*   We use the same teacher $M_{teacher}$ and dataset $\mathcal{D}$ as in our soft distillation experiments. For each example in $\mathcal{D}$, we prompt $M_{teacher}$ with the first 50 tokens and generate 206 tokens using greedy decoding, resulting in a sequence of length $T = 256$. This gives a synthetic dataset, $\mathcal{D}_{hard}$, used to train the student. We initialize the student parameters $\theta$ from the Pythia 1.4B base model and fine-tune on $\mathcal{D}_{hard}$ using cross-entropy loss:

$$\mathcal{L}_{\text{hard}}(\theta) = -\mathbb{E}_{\hat{x}\sim\mathcal{D}_{hard}}\left[\sum_{t=1}^{T}\log P_\theta(\hat{x}_t \mid \hat{x}_{<t})\right] \tag{3}$$

where $\hat{x}_t$ denotes the token at step $t$ and $\hat{x}_{<t}$ represents the preceding context.
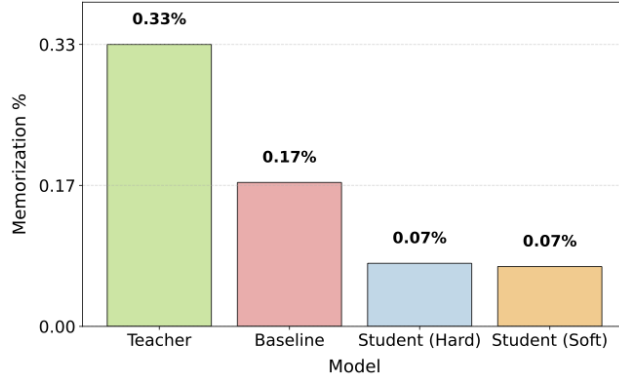
**Figure 9** Hard-distilled and soft-distilled student models exhibit similar memorization rates, both significantly lower than the baseline.

We train $M_{student\_hard}$ using the same compute budget and learning rate as $M_{baseline}$. To ensure a fair comparison, we evaluate performance on the LAMBADA (Paperno et al., 2016) and Winogrande (Sakaguchi et al., 2021) benchmarks using Gao et al. (2024) instead of measuring perplexity on a held-out set from the same distribution. Perplexity is misleading in this context because of the difference in training data: the baseline optimizes on real data, while the student optimizes on the teacher's synthetic outputs. As a result, the baseline will naturally achieve lower perplexity on real validation sets, even if the student is more capable. Therefore, we use downstream tasks to measure actual utility.

We observe that $M_{student\_hard}$ outperforms $M_{baseline}$. Specifically, $M_{baseline}$ achieves a perplexity of 9.41 and an accuracy of 51.85% on LAMBADA, and an accuracy of 55.72% on Winogrande. $M_{student\_hard}$ achieves a perplexity of 6.43 and an accuracy of 56.65% on LAMBADA, along with an accuracy of 57.46% on Winogrande.

*Memorization Analysis.* As shown in Figure 9, the hard-distilled student exhibits a memorization rate of 0.07%, identical to the soft-distilled student and significantly lower than the baseline model (0.17%). Despite the different training objectives, there is a substantial overlap in what they memorize: approximately 70% of the examples memorized by the hard-distilled student are also memorized by the soft-distilled student (Figure 10). Furthermore, similar to the soft-distilled setup, the hard-distilled student primarily memorizes *easy to memorize* examples (i.e., 90% of these examples are also memorized by the teacher and baseline) (Figure 11).

When we examine the examples memorized by one student but not the other, we find they are still largely dominated by the baseline. For instance, the examples memorized by the soft student but missed by the hard student (and vice-versa) are frequently found in the baseline set, further supporting that both students primarily target *easy* examples.

*Inheritance of Difficult Examples.* However, a distinct pattern emerges for the examples memorized exclusively by $M_{student\_hard}$. Figure 10 shows 46 examples memorized by $M_{student\_hard}$ that are not memorized by both $M_{student}$ & $M_{baseline}$ (blue region). Our analysis reveals that the teacher memorizes 80% of this specific subset. We classify these as *difficult* examples because they are not memorized by the baseline but are inherited directly from the teacher. This specific risk of hard distillation is quantified in Figure 11 (brown region): the total amount of memorization inherited solely from the teacher (memorized by teacher and student but not baseline) is 50 examples, an increase by a factor of 2.7 compared to the soft-distilled (Figure 4). This suggests that while soft distillation effectively filters out these *difficult* examples, hard distillation is more prone to memorizing them.
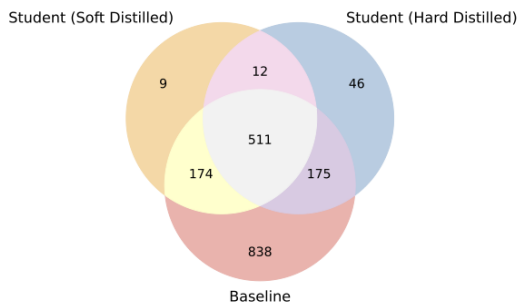
**Figure 10  Overlap between distillations.** Soft-distilled and Hard-distilled students memorize a lot of similar examples with a 70% overlap.
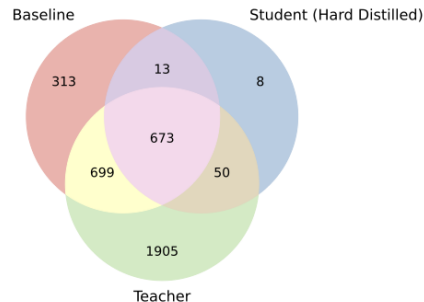
**Figure 11  Hard distillation dynamics.** Similar to soft-distillation, the majority of the examples (673) memorized by the hard-distilled student are *easy-to-memorize* (i.e., they are memorized both by the teacher and baseline).

# 7  Related Work

*Knowledge Distillation (KD)*  was originally proposed to transfer knowledge from a large teacher to a smaller student by matching output distributions (Hinton et al., 2015; Bucila et al., 2006). The standard objective minimizes the Kullback-Leibler (KL) (Kullback and Leibler, 1951) divergence between the teacher's soft targets and the student's predictive distribution. While early methods focused on logit matching, subsequent approaches introduced intermediate feature matching (Romero et al., 2015) to improve fidelity. In the context of LLMs, KD methods have evolved to address the challenges of autoregressive sequence generation. Sequence-level distillation (Kim and Rush, 2016) methods use teacher-generated outputs instead of relying on the logits. More recently, specific techniques have been proposed to stabilize LLM distillation: MiniLLM (Gu et al., 2025) utilizes reverse KL divergence to prevent the student from over-approximating the teacher's complex distribution, while Generalized Knowledge Distillation (GKD) (Agarwal et al., 2024) explores on-policy distillation to mitigate the distribution shift between teacher and student. Ko et al., 2025) propose specialized objectives to improve LLM distillation efficiency. Other approaches, such as Distilling Step-by-Step (Hsieh et al., 2023), focus on transferring reasoning capabilities by distilling Chain-of-Thought rationales alongside final labels.

*Memorization*  Carlini et al. (2019) were the first to study the phenomenon of *unintended* memorization in generative models. This led to numerous in-depth studies on memorization where pretrained language models regurgitate portions of their training data (Carlini et al., 2020; Tirumala et al., 2022; Carlini et al., 2023; Nasr et al., 2023; Biderman et al., 2023a; Huang et al., 2024; Prashanth et al., 2025; Hayes et al., 2025; Morris et al., 2025; Ahmed et al., 2026), including studies on post-trained (Barbero et al., 2025) and fine-tuned models (Mireshghallah et al., 2022; Zeng et al., 2024; Borkar et al., 2025). However, despite these efforts on demystifying memorization, it remains poorly understood in a KD setup.

*Memorization & Privacy in Knowledge Distillation (KD)*  Prior research on studying privacy during knowledge distillation includes membership inference attacks (Jagielski et al., 2023; Zhang et al., 2025; Cui et al., 2025), where the student can reveal information about the teacher's training data. Lee et al. (2025) find that distillation discards undesirable behavior (e.g., vulnerability to adversarial elicitation) while keeping the desired behavior intact. Most related to our work are the studies on data extraction in a KD setup: mainly Dandekar et al. (2024) who show that student machine translation models inherit a lot of the teacher's memorization during sequence-level distillation, Zhang et al. (2025) also study memorization inheritance, and their results indicate that students struggle to reproduce the teacher's entire memorization of 32 tokens, and Singh (2025) show some early evidence that distillation can help reduce memorization.

Unlike prior studies that focus on privacy leakage using membership inference attacks or memorization inheritance, our work systematically tracks memorization across the entire KD pipeline, including teacher,

baseline, and student models. We measure both the amount and type of examples each model memorizes and show that student models primarily memorize *easy-to-memorize examples*.

## 8  Conclusion

We demonstrate that knowledge distillation improves model utility while significantly reducing training data memorization compared to standard fine-tuning. Our results show that distilled models primarily memorize *easy-to-memorize* examples. Leveraging this finding, we show that high-risk examples can be predicted and removed prior to distillation, which substantially lowers memorization rates in the Student. We further explain why logit-level KD reduces memorization by analyzing log probabilities and Shannon entropy. Finally, we compare memorization risks between soft and hard distillation. We find that although both methods memorize largely the same examples, hard distillation poses a slightly higher risk due to greater memorization inheritance from the teacher.

## 9  Acknowledgements

# References

Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-generated mistakes, 2023.

Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-generated mistakes, 2024. https://arxiv.org/abs/2306.13649.

Ahmed Ahmed, A. Feder Cooper, Sanmi Koyejo, and Percy Liang. Extracting books from production language models, 2026. https://arxiv.org/abs/2601.02671.

Federico Barbero, Xiangming Gu, Christopher A. Choquette-Choo, Chawin Sitawarin, Matthew Jagielski, Itay Yona, Petar Veličković, Ilia Shumailov, and Jamie Hayes. Extracting alignment data in open models, 2025. https://arxiv.org/abs/2510.18554.

Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. Emergent and predictable memorization in large language models, 2023a. https://arxiv.org/abs/2304.11158.

Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar Van Der Wal. Pythia: a suite for analyzing large language models across training and scaling. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023b.

Jaydeep Borkar. What can we learn from data leakage and unlearning for law?, 2023. https://arxiv.org/abs/2307.10476.

Jaydeep Borkar, Matthew Jagielski, Katherine Lee, Niloofar Mireshghallah, David A. Smith, and Christopher A. Choquette-Choo. Privacy ripple effects from adding or removing personal information in language model training, 2025. https://arxiv.org/abs/2502.15680.

Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Knowledge Discovery and Data Mining*, 2006. https://api.semanticscholar.org/CorpusID:11253972.

Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evaluating and testing unintended memorization in neural networks, 2019. https://arxiv.org/abs/1802.08232.

Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models. In *USENIX Security Symposium*, 2020. https://api.semanticscholar.org/CorpusID:229156229.

Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. Membership inference attacks from first principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1897–1914, 2022a. doi: 10.1109/SP46214.2022.9833649.

Nicholas Carlini, Andreas Terzis, Matthew Jagielski, Florian Tramer, Nicolas Papernot, and Chiyuan Zhang. The privacy onion effect: memorization is relative. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022b. Curran Associates Inc. ISBN 9781713871088.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models, 2023. https://arxiv.org/abs/2202.07646.

Ziyao Cui, Minxing Zhang, and Jian Pei. On membership inference attacks in knowledge distillation, 2025. https://arxiv.org/abs/2505.11837.

Chinmay Dandekar, Wenda Xu, Xi Xu, Siqi Ouyang, and Lei Li. Translation canvas: An explainable interface to pinpoint and analyze translation systems. In Delia Irazu Hernandez Farias, Tom Hope, and Manling Li, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 344–350, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-demo.36. https://aclanthology.org/2024.emnlp-demo.36/.

Verna Dankers and Vikas Raunak. Memorization inheritance in sequence-level knowledge distillation for neural machine translation. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 760–774, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-252-7. doi: 10.18653/v1/2025.acl-short.61. https://aclanthology.org/2025.acl-short.61/.

NVIDIA et al. Nvidia nemotron nano 2: An accurate and efficient hybrid mamba-transformer reasoning model, 2025. https://arxiv.org/abs/2508.14444.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. https://zenodo.org/records/12608602.

Team Gemma et al. Gemma 2: Improving open language models at a practical size, 2024.

Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models, 2025. https://arxiv.org/abs/2306.08543.

Daya Guo et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. https://arxiv.org/abs/2501.12948.

Jamie Hayes, Marika Swanberg, Harsh Chaudhari, Itay Yona, Ilia Shumailov, Milad Nasr, Christopher A. Choquette-Choo, Katherine Lee, and A. Feder Cooper. Measuring memorization in language models via probabilistic extraction. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors, *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 9266–9291, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.469. https://aclanthology.org/2025.naacl-long.469/.

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015. https://arxiv.org/abs/1503.02531.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack W. Rae, and Laurent Sifre. Training compute-optimal large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.

Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8003–8017, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.507. https://aclanthology.org/2023.findings-acl.507/.

Jing Huang, Diyi Yang, and Christopher Potts. Demystifying verbatim memorization in large language models, 2024. https://arxiv.org/abs/2407.17817.

Matthew Jagielski, Milad Nasr, Katherine Lee, Christopher Choquette-Choo, Nicholas Carlini, and Florian Tramèr. Students parrot their teachers: membership inference on model distillation. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.

Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating training data mitigates privacy risks in language models. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 10697–10707. PMLR, 17–23 Jul 2022. https://proceedings.mlr.press/v162/kandpal22a.html.

Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation. In Jian Su, Kevin Duh, and Xavier Carreras, editors, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1139. https://aclanthology.org/D16-1139/.

Jongwoo Ko, Sungnyun Kim, Tianyi Chen, and Se-Young Yun. Distillm: Towards streamlined distillation for large language models. In *Forty-first International Conference on Machine Learning*.

Jongwoo Ko, Tianyi Chen, Sungnyun Kim, Tianyu Ding, Luming Liang, Ilya Zharkov, and Se-Young Yun. Distillm-2: A contrastive approach boosts the distillation of llms. *arXiv preprint arXiv:2503.07067*, 2025.

Solomon Kullback and R. A. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 22:79–86, 1951. https://api.semanticscholar.org/CorpusID:120349231.

Bruce W. Lee, Addie Foote, Alex Infanger, Leni Shor, Harish Kamath, Jacob Goldman-Wetzler, Bryce Woodworth, Alex Cloud, and Alexander Matt Turner. Distillation robustifies unlearning, 2025. https://arxiv.org/abs/2506.06278.

Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8424–8445, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.577. https://aclanthology.org/2022.acl-long.577/.

Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016.

Fatemehsadat Mireshghallah, Archit Uniyal, Tianhao Wang, David Evans, and Taylor Berg-Kirkpatrick. An empirical analysis of memorization in fine-tuned autoregressive language models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1816–1826, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.119. https://aclanthology.org/2022.emnlp-main.119/.

John X. Morris, Chawin Sitawarin, Chuan Guo, Narine Kokhlikyan, G. Edward Suh, Alexander M. Rush, Kamalika Chaudhuri, and Saeed Mahloujifar. How much do language models memorize?, 2025. https://arxiv.org/abs/2505.24832.

Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models, 2023. https://arxiv.org/abs/2311.17035.

Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Allyson Ettinger, Michal Guerquin, David Heineman, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Jake Poznanski, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 olmo 2 furious, 2025. https://arxiv.org/abs/2501.00656.

Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The LAMBADA dataset: Word prediction requiring a broad discourse context. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1144. https://aclanthology.org/P16-1144/.

Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Úlfar Erlingsson. Scalable private learning with pate, 2018. https://arxiv.org/abs/1802.08908.

Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale, 2024. https://arxiv.org/abs/2406.17557.

USVSN Sai Prashanth, Alvin Deng, Kyle O'Brien, Jyothir S V, Mohammad Aflah Khan, Jaydeep Borkar, Christopher A. Choquette-Choo, Jacob Ray Fuehne, Stella Biderman, Tracy Ke, Katherine Lee, and Naomi Saphra. Recite, reconstruct, recollect: Memorization in lms as a multifaceted phenomenon, 2025. https://arxiv.org/abs/2406.17746.

Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets, 2015. https://arxiv.org/abs/1412.6550.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: an adversarial winograd schema challenge at scale. *Commun. ACM*, 64(9):99–106, August 2021. ISSN 0001-0782. doi: 10.1145/3474381. https://doi.org/10.1145/3474381.

Virat Shejwalkar and Amir Houmansadr. Membership privacy for machine learning models through knowledge transfer. In *AAAI Conference on Artificial Intelligence*, 2021. https://api.semanticscholar.org/CorpusID:235349092.

Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18, 2017. doi: 10.1109/SP.2017.41.

Simardeep Singh. From teacher to student: Tracking memorization through model distillation. In Robin Jia, Eric Wallace, Yangsibo Huang, Tiago Pimentel, Pratyush Maini, Verna Dankers, Johnny Wei, and Pietro Lesci, editors, *Proceedings of the First Workshop on Large Language Model Memorization (L2M2)*, pages 78–82, Vienna, Austria,

August 2025. Association for Computational Linguistics. ISBN 979-8-89176-278-7. doi: 10.18653/v1/2025.l2m2-1.6. https://aclanthology.org/2025.l2m2-1.6/.

Xinyu Tang, Saeed Mahloujifar, Liwei Song, Virat Shejwalkar, Milad Nasr, Amir Houmansadr, and Prateek Mittal. Mitigating membership inference attacks by Self-Distillation through a novel ensemble architecture. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 1433–1450, Boston, MA, August 2022. USENIX Association. ISBN 978-1-939133-31-1. https://www.usenix.org/conference/usenixsecurity22/presentation/tang.

Kushal Tirumala, Aram H. Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization without overfitting: analyzing the training dynamics of large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.

Johnny Tian-Zheng Wei, Ameya Godbole, Mohammad Aflah Khan, Ryan Wang, Xiaoyuan Zhu, James Flemings, Nitya Kashyap, Krishna P. Gummadi, Willie Neiswanger, and Robin Jia. Hubble: a model suite to advance the study of llm memorization, 2025. https://arxiv.org/abs/2510.19811.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. https://arxiv.org/abs/2505.09388.

Shenglai Zeng, Yaxin Li, Jie Ren, Yiding Liu, Han Xu, Pengfei He, Yue Xing, Shuaiqiang Wang, Jiliang Tang, and Dawei Yin. Exploring memorization in fine-tuned language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3917–3948, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.216. https://aclanthology.org/2024.acl-long.216/.

Ziqi Zhang, Ali Shahin Shamsabadi, Hanxiao Lu, Yifeng Cai, and Hamed Haddadi. Membership and memorization in LLM knowledge distillation. In Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng, editors, *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 20085–20095, Suzhou, China, November 2025. Association for Computational Linguistics. ISBN 979-8-89176-332-6. doi: 10.18653/v1/2025.emnlp-main.1015. https://aclanthology.org/2025.emnlp-main.1015/.

# A  Additional Experimental Results

## A.1  Easy-to-Memorize Examples on Additional Models and Datasets

Figure 12 shows zlib entropy versus Baseline perplexity for *easy-to-memorize* and other training examples across additional models and datasets: Pythia trained on Wikitext, and OLMo-2 and Qwen-3 trained on FineWeb. We observe similar clustering to that in Section 3.2 (Figure 5), indicating that our findings generalize across models and datasets.
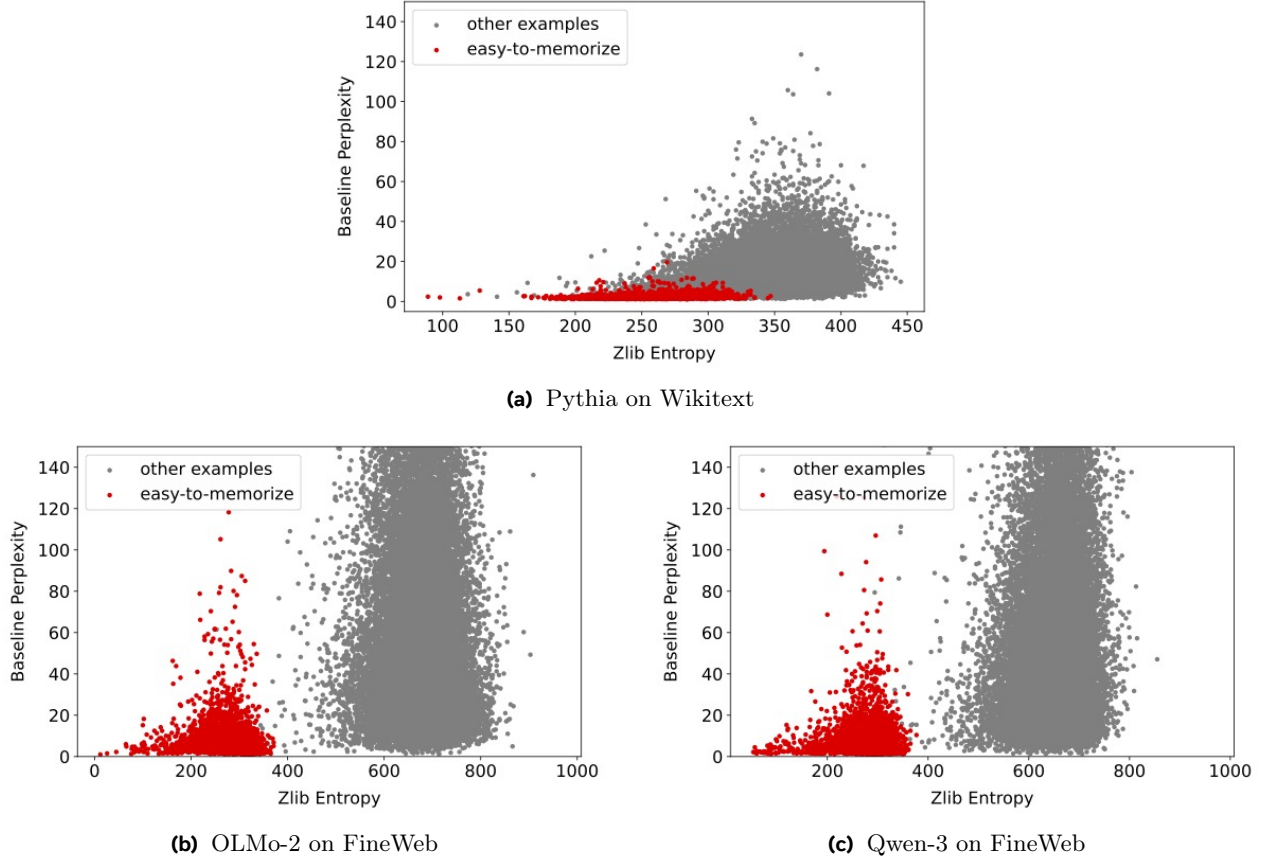


**(a)** Pythia on Wikitext



**(b)** OLMo-2 on FineWeb

**(c)** Qwen-3 on FineWeb

**Figure 12  Intrinsic properties of easy-to-memorize examples**. We plot the zlib entropy versus Baseline Perplexity for the *easy-to-memorize* examples (red) compared to a random subset of 25,000 other training examples (grey). They form a distinct cluster with significantly lower entropy and perplexity.

## A.2  Extended Generalization Results

Table 3 reports validation loss and perplexity for the Pythia family on Wikitext.

**Table 3  Performance comparison on Wikitext.** Validation loss and perplexity for Pythia trained on Wikitext. The distilled student consistently outperforms the baseline.

| Model | Val Loss | PPL |
|---|---|---|
| *Pythia Family (Wikitext)* | | |
| Teacher (12B) | 2.72 | 14.49 |
| Baseline (1.4B) | 2.82 | 16.33 |
| Student (1.4B) | **2.75** | **15.36** |

## A.3 Extended Analysis: Other Models and Datasets

In this section, we extend our analysis to Pythia models trained on the Wikitext dataset, Olmo-2 models trained on FineWeb, and Qwen-3 models trained on FineWeb.

### A.3.1 Memorization Overlap Analysis

Figure 13 illustrates the memorization overlap between the Teacher, Baseline, and Student models for Pythia (trained on Wikitext), OLMo-2 (FineWeb), and Qwen-3 (FineWeb). Consistent with the findings in Section 3.2, we observe that Students preferentially memorize *easy-to-memorize* examples, exhibiting significant overlap with both the Teacher and Baseline. Specifically, examples categorized as easy-to-memorize account for 88% of the Pythia Student's total memorization, 85% for the OLMo-2 Student, and 70% for the Qwen-3 Student. In terms of inheritance (shown by blue region), the Pythia Student inherits 9% of its memorization directly from the Teacher, while the OLMo-2 and Qwen-3 Students inherit 13% and 27%, respectively.
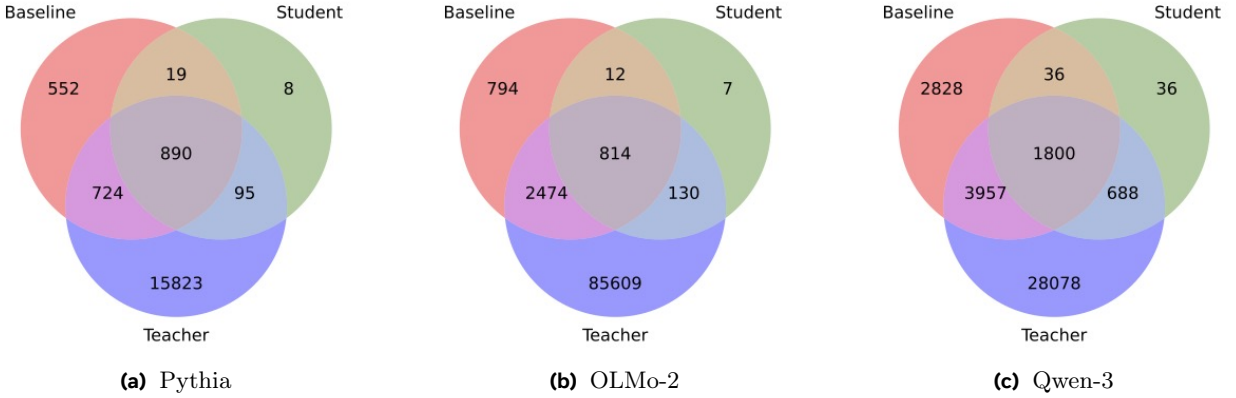


**(a)** Pythia      **(b)** OLMo-2      **(c)** Qwen-3

**Figure 13 Memorization overlap across architectures and datasets.** Venn diagrams visualizing the intersection of memorized examples between Teacher, Baseline, and Student models for **(a)** Pythia 1.4B on Wikitext, **(b)** OLMo-2 1B on FineWeb, and **(c)** Qwen-3 1.7B on FineWeb. Across all settings, the Student preferentially memorizes *easy to memorize* examples (examples that both the Teacher and Baseline memorize) while showing some inheritance (blue region) from the Teacher.

### A.3.2 Feature Distributions Distinguishing Memorized vs. Non-Memorized Examples

Figure 14 shows the distribution of features used to train the memorization classifier for the Pythia model on Wikitext. Figure 15 presents the corresponding distributions for the OLMo-2 model trained on FineWeb, and Figure 16 for the Qwen-3 model trained on FineWeb.
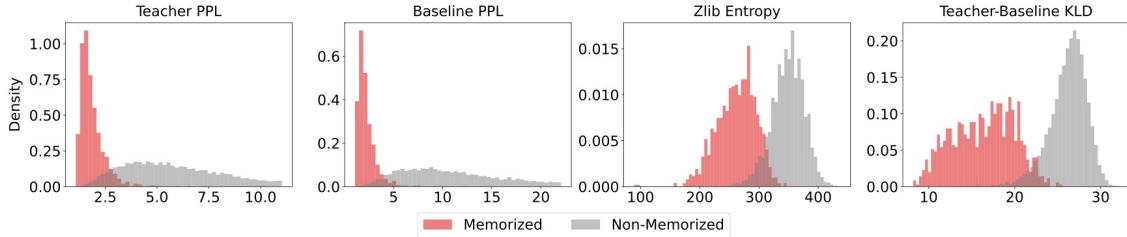


**Figure 14 Feature distributions distinguishing memorized vs. non-memorized examples in Student (Pythia on Wikitext)**. Teacher Perplexity is the most significant predictor, followed by zlib Entropy, the KL divergence loss between the Teacher and Baseline, and Baseline Perplexity. Across all metrics, lower values are consistently associated with memorized examples.
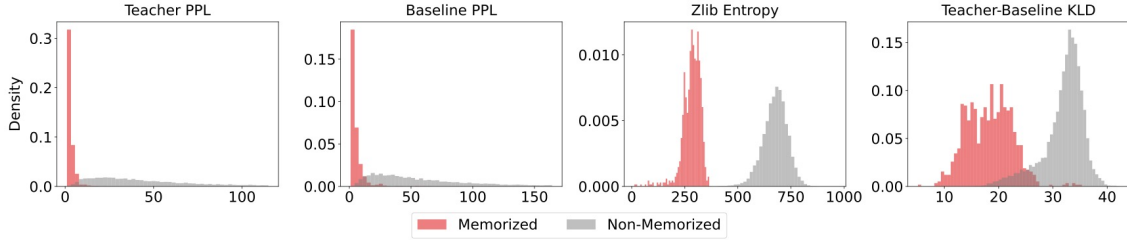
**Figure 15 Feature distributions distinguishing memorized vs. non-memorized examples in Student (OLMo-2 on FineWeb)**. zlib Entropy is the most significant predictor, followed by the KL divergence loss between the Teacher and Baseline, Teacher Perplexity, and Baseline Perplexity. Across all metrics, lower values are consistently associated with memorized examples.



**Figure 16 Feature distributions distinguishing memorized vs. non-memorized examples in Student (Qwen-3 on FineWeb)**. zlib Entropy is the most significant predictor, followed by the KL divergence loss between the Teacher and Baseline, Teacher Perplexity, and Baseline Perplexity. Across all metrics, lower values are consistently associated with memorized examples.
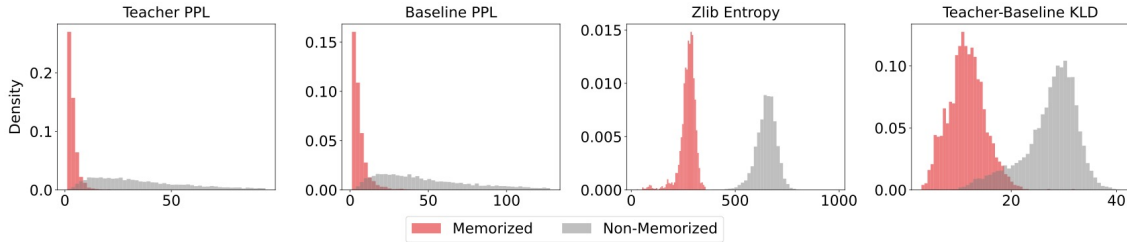
## A.4 More Details On Our Memorization Classifier

Table 4 reports the coefficients learned by our logistic regression classifier for predicting memorization in the Pythia student distilled on FineWeb. Table 5 presents the corresponding results for Pythia trained on Wikitext, while Tables 6 and 7 report results for OLMo-2 and Qwen-3 on FineWeb, respectively.

**Table 4** Feature Importance based on standardized Logistic Regression coefficients (averaged over 100 trials). The large negative magnitude of zlib Entropy indicates it is the strongest predictor. **Results are for Pythia models trained on FineWeb.**

| Feature | Coeff. (Mean $\pm$ Std) |
|---|---|
| Teacher PPL | $-0.3341 \pm 0.1228$ |
| Baseline PPL | $-0.4010 \pm 0.1417$ |
| zlib Entropy | $\mathbf{-4.5001 \pm 0.1324}$ |
| Teacher–Baseline KLD | $-1.0579 \pm 0.1664$ |

**Table 5** Memorization classifier results for **Pythia** trained on **Wikitext**. We report mean ± standard deviation over 100 trials. Feature importance corresponds to standardized logistic regression coefficients.

| Metric | Value |
|---|---|
| Accuracy | $0.9799 \pm 0.0038$ |
| Precision | $0.9515 \pm 0.0111$ |
| Recall | $0.9690 \pm 0.0095$ |
| F1 | $0.9601 \pm 0.0075$ |
| ROC-AUC | $0.9964 \pm 0.0010$ |
| **Feature** | **Coefficient** |
| Teacher PPL | $-2.7531 \pm 0.2852$ |
| Baseline PPL | $-1.9694 \pm 0.2530$ |
| zlib Entropy | $-2.0809 \pm 0.1712$ |
| Teacher–Baseline KLD | $-2.0074 \pm 0.2294$ |

**Table 6** Memorization classifier results for **OLMo-2** trained on **FineWeb**. We report mean ± standard deviation over 100 trials. Feature importance corresponds to standardized logistic regression coefficients.

| Metric | Value |
|---|---|
| Accuracy | $0.9988 \pm 0.0010$ |
| Precision | $0.9952 \pm 0.0040$ |
| Recall | $1.0000 \pm 0.0000$ |
| F1 | $0.9976 \pm 0.0020$ |
| ROC-AUC | $0.9999 \pm 0.0002$ |
| **Feature** | **Coefficient** |
| Teacher PPL | $-0.1151 \pm 0.0837$ |
| Baseline PPL | $-0.1854 \pm 0.0877$ |
| zlib Entropy | $-4.4637 \pm 0.1376$ |
| Teacher–Baseline KLD | $-1.5805 \pm 0.1455$ |

### A.5 Dissecting Student and Baseline Memorization

Next, we are interested in whether it is possible to identify which examples would be memorized by the student model during distillation and which would be memorized by the baseline model when trained from scratch. From the results in Section 3.1, we know that distilled student models memorize less than the baseline, even though both models have a similar number of parameters. From Section 3.2, we see that the majority of examples memorized by the student are also memorized by the baseline. However, the baseline model memorizes many additional examples that are never memorized by the student. This raises the question: are there factors that can predict which examples will be memorized by the student and which will be memorized *exclusively* by the baseline model?

Let there be two datasets of memorized examples, where $D_s$ represents examples memorized by the student, and $D_b$ represents examples memorized by the baseline but not student. We mount a membership inference attack (MIA) (Shokri et al., 2017) using $D_s$ & $D_b$ as datasets and compute the area under the curve (AUC) using the following features: teacher perplexity, baseline perplexity, and the KLD loss between teacher and student. We observe an AUC of 0.745 for teacher perplexity, 0.754 for baseline perplexity, and 0.855 for the KLD loss between teacher and student. This shows that these metrics are good indicators of which examples are memorized by the student and the baseline.

Figure 17 displays Kernel Density Estimation (KDE) contours correlating baseline perplexity with teacher-student KLD loss across three categories of memorized examples.

The blue region represents easy-to-memorize examples that are memorized both by baseline and student

**Table 7** Memorization classifier results for **Qwen-3** trained on **FineWeb**. We report mean ± standard deviation over 100 trials. Feature importance corresponds to standardized logistic regression coefficients.

| Metric | Value |
|---|---|
| Accuracy | $0.9991 \pm 0.0004$ |
| Precision | $0.9963 \pm 0.0017$ |
| Recall | $1.0000 \pm 0.0000$ |
| F1 | $0.9981 \pm 0.0009$ |
| ROC-AUC | $0.9999 \pm 0.0002$ |

| Feature | Coefficient |
|---|---|
| Teacher PPL | $-0.1576 \pm 0.1053$ |
| Baseline PPL | $-0.1184 \pm 0.0710$ |
| zlib Entropy | $-5.4524 \pm 0.1125$ |
| Teacher–Baseline KLD | $-1.5311 \pm 0.1340$ |

models. These examples typically have lower baseline perplexity and teacher-student KLD loss. However, we do see a tail extending to higher PPL values. This suggests that this group does contain a subset of examples that were inherently harder for the baseline to memorize.

The orange region shows examples that are memorized exclusively by the baseline model. This cluster is tightly concentrated in the region of lowest Baseline PPL ($< 10$) and highest teacher-student KLD loss indicating these examples were never memorized by the student.

The green region highlights examples memorized exclusively by the student model. These examples appear to be concentrated in the low KLD-loss region and, interestingly, also have lower baseline perplexity. This may indicate that these examples are on the verge of being memorized by the baseline model.
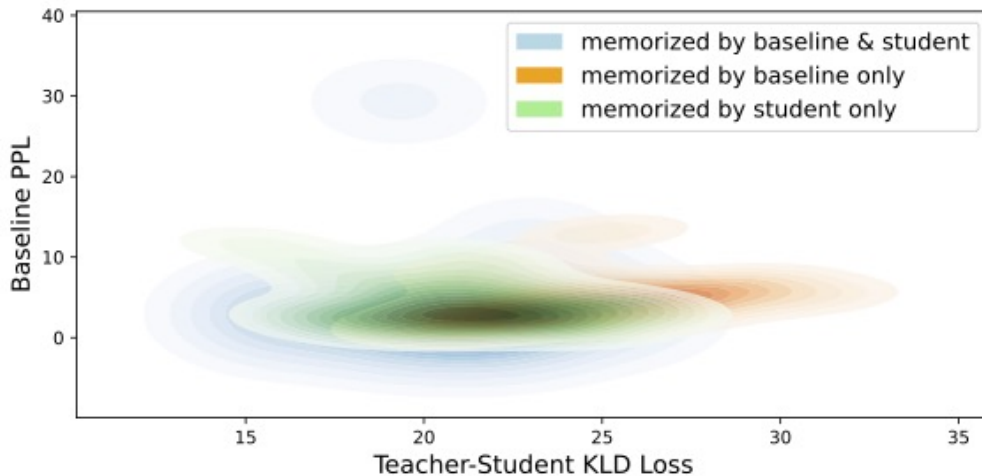


**Figure 17** Kernel Density Estimation (KDE) plot showing probability distribution of baseline perplexity and teacher-student KLD loss across three datasets: examples memorized by baseline & student, examples memorized by baseline-only, and examples memorized by student-only. Results are for Pythia models trained on the FineWeb.

## A.6 Evaluating Memorization During Distillation in a Pre-training Setup

We also ran some early experiments to understand the dynamics of memorization for logit-level distillation in a pre-training setup (Figure 18) using the formulation from equation 1. We use the 116k checkpoint of Pythia 12B as the Teacher. We initialize the Student from the 115k checkpoint of Pythia 2.8B and continue training it on the data seen by the Teacher between checkpoints 115k and 116k ($\approx$ 1M examples) using KL divergence

($T = 2$) for five epochs. We also use the 116k checkpoint of Pythia 2.8B from the Pythia suite as a Baseline. Note that all Pythia family models see exactly the same data between each checkpoint during training. We then compare memorization rates across the Teacher, Baseline, and Student.

Our results correlate strongly with memorization trends observed in our main fine-tuning KD setup. As shown in Table 8, the Teacher has a memorization rate of 1.57%, the Baseline has 0.27%, and the Student has a rate of **0.06%**. Of 2785 examples memorized by the Baseline, 2117 (76%) are also memorized by the teacher, indicating that the *easy-to-memorize* examples are prevalent in this setup as well. Among these 2,117 examples, only 361 are memorized by the student, highlighting the regularization effect of distillation. Finally, we identify 92 examples that are memorized by the student but never memorized by either the teacher or the baseline.
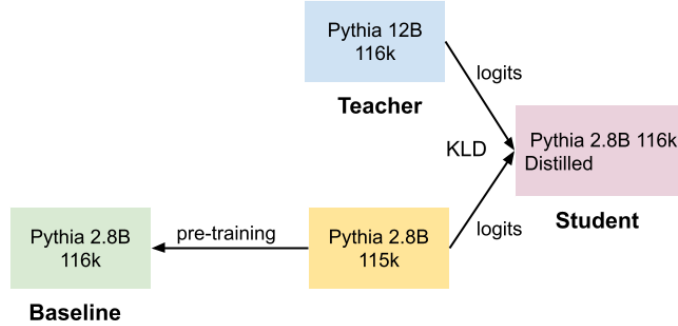


**Figure 18** Our Knowledge Distillation Framework in a Pre-training Setup for Pythia family models.

**Table 8 Memorization statistics.** Comparison of memorization across the Teacher, Baseline, and Student models for approximately 1 million training examples. The Student model exhibits significantly lower memorization compared to the Baseline.

| Model | # Examples Memorized | Memorization % |
|---|---|---|
| Teacher (12B) | 16,133 | 1.57% |
| Baseline (2.8B) | 2,785 | 0.27% |
| Student (2.8B) | **638** | **0.06%** |