

# CoLa: A CHOICE LEAKAGE ATTACK FRAMEWORK TO EXPOSE PRIVACY RISKS IN SUBSET TRAINING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Subset training, where models are trained on a carefully chosen portion of data rather than the entire dataset, has become a standard tool for scaling modern machine learning. From coreset selection in vision to large-scale filtering in language models, these methods promise scalability without compromising utility. A common intuition is that training on fewer samples should also reduce privacy risks. In this paper, we challenge this assumption. We show that subset training is not privacy free: the very choices of which data are included or excluded can introduce new privacy surface and leak more sensitive information. Such information can be captured by adversaries either through side-channel metadata from the subset selection process or via the outputs of the target model. To systematically study this phenomenon, we propose CoLa (Choice Leakage Attack), a unified framework for analyzing privacy leakage in subset selection. In CoLa, depending on the adversary’s knowledge of the side-channel information, we define two practical attack scenarios: Subset-aware Side-channel Attacks and Black-box Attacks. Under both scenarios, we investigate two privacy surfaces unique to subset training: (1) Training-membership MIA (TM-MIA), which concerns only the privacy of training data membership, and (2) Selection-participation MIA (SP-MIA), which concerns the privacy of all samples that participated in the subset selection process. Notably, SP-MIA enlarges the notion of membership from model training to the entire data–model supply chain. Experiments on vision and language models show that existing threat models underestimate the privacy risks of subset training: the enlarged privacy surface not only retains training membership leakage but also exposing selection membership, extending risks from individual models to the broader ML ecosystem.

## 1 INTRODUCTION

The scale of modern datasets has made training on the full corpus increasingly impractical. To address this, practitioners routinely employ subset training, where only a carefully chosen ratio of data is used. This paradigm is adopted not only for efficiency but also to improve data quality, since selection can remove redundancy and noise while retaining informative samples. Subset training spans diverse applications: coreset selection (Bachem et al., 2015; Munteanu et al., 2018; Mirzasoleiman et al., 2020) in vision, dataset pruning (Sorscher et al., 2022; Yang et al., 2022; Qin et al., 2023), active learning (Sener & Savarese, 2018; Ducoffe & Precioso, 2018; Agarwal et al., 2020) in general ML, and large-scale deduplication (Lee et al., 2022), filtering (Rae et al., 2021), and sampling (Gunasekar et al., 2023; Peng et al., 2025; Wettig et al., 2024) in language model pretraining.

While subset training is widely celebrated for these benefits, its privacy implications remain underexplored (Zhao & Zhang, 2025). A common intuition suggests that fewer training samples should imply less privacy leakage (Dong et al., 2022). Yet this reasoning overlooks an important fact: *the choices made during subset selection themselves encode signals about which data were included and which were excluded*. These signals can be inherited through shifts in the data distribution or model behavior, making them exploitable by adversaries.

We ask the fundamental question: *Does subset training actually reduce privacy leakage?* Our answer is *no*. We show that subset training introduces new attack surfaces: not only is the included data that used for training compromised, but the excluded data discarded from training can also become

vulnerable due to correlations introduced by the selection mechanism. In other words, due to the data-oriented nature of the subset selection process, beyond the training data leakage emphasized by traditional MIA (Shokri et al., 2017; Hu et al., 2022), the choice signals further extend privacy risks from individual models to the broader data-model supply chain. Accordingly, we define two complementary privacy surfaces: *Training-membership MIA (TM-MIA)*, which resembles traditional MIA by focusing on the membership of training data, and *Selection-participation MIA (SP-MIA)*, a privacy surface tailored to subset training that focuses on membership at the data selection level.

To systematically study membership leakage under these privacy surfaces, we propose **CoLa (Choice Leakage Attack)**, a framework that leverages choice signals in a principled way to conduct attacks across different surfaces. CoLa captures risks under two complementary settings: (i) a *Subset-aware Side-channel* setting, where the adversary has access to the target model’s outputs and selection metadata (e.g., the selection algorithm and the inclusion ratio); and (ii) a *Black-box* setting, where the adversary observes only model outputs and is aware that subsetting may have been used, without knowing any selection metadata. Extensive results show that for both privacy surfaces under these two attack settings, CoLa can substantially strengthen the attack performance. In short, subset training does not guarantee privacy; it enlarges the attack surface of modern ML pipelines and highlights the need to protect privacy across the entire data-model supply chain. We summarize our contributions as follows:

- We provide the first systematic definition and exploration of the membership leakage problem under subset training. This novel attack scenario reveals a severe privacy risk in the subset selection process: not only is the privacy of training data compromised, but the data excluded during selection is also at risk.
- We propose CoLa (Choice Leakage Attack), a framework tailored to subset selection that leverages choice signals in a principled way for more reliable membership inference, while seamlessly unifying diverse attack settings and surfaces.
- Experiments across both vision and language models confirm the broad capability of CoLa. For example, in the black-box setting, the AUC of CoLa on Pythia-160M surpasses 80% under SP-MIA, where all baseline methods fail.

## 2 RELATED WORKS

**Subset training and data-efficient learning.** A large body of research has explored how to reduce the cost of large-scale training by operating on subsets of data. Coreset selection constructs small but representative subsets that approximate training on the full data (Bachem et al., 2015; Munteanu et al., 2018; Mirzasoleiman et al., 2020; Yang et al., 2024b). Dataset pruning removes redundant or low-value samples to improve efficiency and generalization (Sorscher et al., 2022; Yang et al., 2022; Qin et al., 2023; Maharana et al., 2023; Tan et al., 2024). Active learning queries the most informative examples to reduce annotation cost (Sener & Savarese, 2018; Ducoffe & Precioso, 2018; Agarwal et al., 2020; Borsos et al., 2020; Margatina et al., 2021). In large-scale language models, deduplication and filtering pipelines are routinely applied to eliminate noise and improve training quality (Lee et al., 2022; Rae et al., 2021; Raffel et al., 2023; Gao et al., 2020a). These techniques have been extensively studied for efficiency and utility, but their privacy consequences remain largely underexplored.

**Membership inference attacks.** Membership inference attacks (MIAs) are among the most widely studied privacy threats in machine learning. Early work by Shokri et al. (2017) proposed shadow models to train attack classifiers distinguishing members from nonmembers. Subsequent methods exploited confidence scores, loss values, or gradients (Yeom et al., 2018; Sablayrolles et al., 2019; Carlini et al., 2022b). MIAs have been demonstrated in supervised learning, federated learning, and large language models (Nasr et al., 2018; Hu et al., 2022; Li et al., 2025), motivating defenses such as differential privacy (Abadi et al., 2016) and adversarial regularization (Nasr et al., 2018). This body of work reveals how models trained on fixed datasets can memorize and leak sensitive information. However, they primarily focus on constructing membership signals in a one-shot manner, with these signals being tightly coupled to a specific model. We find such model-oriented signal less effective in the context of subset training. Leveraging the unique characteristics of the subset selection process, we instead construct membership signals in a data-oriented manner.

**Synthetic data and privacy.** Synthetic data generation has been studied as a way to train models without exposing raw datasets, with the promise of stronger privacy (Hu et al., 2024; Tan et al., 2025). However, subsequent research has shown that synthetic datasets can still leak sensitive information about the original data, including membership and attributes (Stadler et al., 2022; van Breugel et al., 2023; Zhao & Zhang, 2025). Rather than analyzing risks inherent in *synthetic data generation pipelines*, we turn to *subset training with real data*, where high-fidelity samples remain but the selection process itself exposes a distinct and overlooked channel of privacy leakage.

### 3 PROBLEM SETTING

#### 3.1 MEMBERSHIP INFERENCE UNDER SUBSET TRAINING

Let  $D_0 \subseteq \mathcal{X} \times \mathcal{Y}$  denote the original dataset that undergoes a subset selection procedure. A selector  $\text{Sel}(\cdot; r)$  with a given selection ratio  $r$  partitions  $D_0$  into two disjoint sets: the **included data**  $I$  used for training, and the **excluded data**  $E$  that are discarded:

$$(I, E) = \text{Sel}(D_0; r), \text{ with } I \cap E = \emptyset, I \cup E = D_0, |I|/|D_0| = r. \quad (1)$$

Following the standard MIA pipeline (Shokri et al., 2017), we further denote by  $O$  the **outside data** that never enter the selection process. A model  $f_\theta$  is trained solely on  $I$ . This partition naturally induces two types of membership inference task:

**Training-membership MIA (TM-MIA).** This attack takes the model itself as the attack surface and membership is defined solely by the training data. A sample  $x$  is a member if  $x \in I$  and nonmember if  $x \in E \cup O$ . This forms a natural and widely adopted threat model, as the model is the most direct output of the ML system. This setting is consistent with conventional MIAs (Shokri et al., 2017; Carlini et al., 2022b).

**Selection-participation MIA (SP-MIA).** However, when the attack surface is enlarged to the entire data-model pipeline, membership expands from only the training data to a much larger portion of all collected data. As shown in Figure 1, we refer to the collected data as selection members, where a sample  $x$  is a member if  $x \in I \cup E$  and a non-member if  $x \in O$ . Its membership cannot be explained by direct model memorization, but instead reveals *choice leakage*, a side-channel signal from the subset selection process of the data-model supply chain. Such choice leakage risk is severe as it exposes a system’s selection preferences. Once the data-model supply chain is exposed to privacy risks, the entire pipeline, from raw data to model outputs, becomes vulnerable to malicious manipulation. **To our knowledge, this is the first work to systematically investigate this perspective.**

Both tasks can be framed as binary hypothesis tests over a scoring function  $s : \mathcal{D}_0 \rightarrow \mathbb{R}$ , which measures the likelihood of a sample  $x$  belonging to the respective member set. Given  $\mathcal{D}_0 = I \cup E \cup O$ , the member–nonmember partitions are:

$$\mathcal{M}_{\text{TM}} = I, \quad \mathcal{N}_{\text{TM}} = E \cup O, \quad (2)$$

$$\mathcal{M}_{\text{SP}} = I \cup E, \quad \mathcal{N}_{\text{SP}} = O. \quad (3)$$

The goal is to design a scoring function  $s(x)$  that distinguishes  $\mathcal{M}$  from  $\mathcal{N}$  under both definitions.

#### 3.2 ADVERSARY KNOWLEDGE

Subset training changes not only the definition of membership but also the adversary’s potential knowledge and capabilities. We consider two complementary scenarios:

**Subset-aware side-channel attacks.** In line with the common assumption in prior MIAs, the adversary can query the deployed model  $f_\theta$  and observe its outputs (e.g., prediction labels or confidence scores). In addition, it has access to *side information about the selection process*, such as the strategy used (e.g., coreset selection, pruning, filtering) or the approximate inclusion ratio. Such an assumption is realistic: pruning papers routinely report retained percentages to justify efficiency–utility

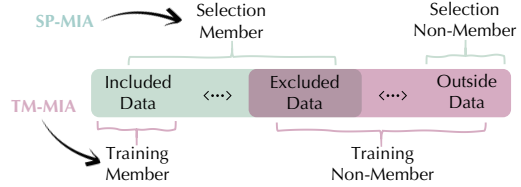


Figure 1: Privacy surfaces under subset training.

trade-offs, active learning and coreset methods describe selection strategies for reproducibility, and large-scale LLM pipelines release dataset cards documenting filtering heuristics, inclusion ratios, or deduplication statistics (Cohen-Addad et al., 2021; Biderman et al., 2023a; Dubey et al., 2024; Yang et al., 2024a). Crucially, this information reflects only high-level rules, not the exact membership of individual samples. Our attack targets precisely this gap: even when only the selection algorithm or ratio is public, an adversary can exploit this side-channel to infer which specific samples were included or excluded, thereby exposing *choice leakage* in subset training.

**Black-box attacks.** Here the adversary can only query the deployed model  $f_\theta$  and observe its outputs. The entire subset selection stage is hidden, so the adversary must rely solely on the observable behavior of the trained model or the intrinsic data-specific information. This setting captures the most restrictive and widely assumed threat model in prior MIA research (Hu et al., 2022).

## 4 METHOD

### 4.1 CHALLENGES OF MEMBERSHIP INFERENCE UNDER SUBSET TRAINING

In conventional MIA, success comes from exploiting overfitting: models tend to assign systematically higher confidence to their training data than to non-members. Under subset training, however, this signal becomes entangled. Figure 2 illustrates this using the LiRA attack signal from (Carlini et al., 2022b) on a model trained on  $I$  selected from  $D_0$  by Glistter (Killamsetty et al., 2021b). The dataset used here is CIFAR10 and the model is ResNet18. Since the selector is designed to make training on  $I$  approximate the effect of training on  $I \cup E$ , the confidence distributions of included, excluded, and outside samples exhibit more complex overlaps: (i)  $I$  concentrates at high confidence,  $E$  shifts lower, while outside data often show a bimodal distribution; (ii) in TM-MIA,  $I$  and  $E \cup O$  remain partly separated but overlap substantially at high confidence; (iii) in SP-MIA, the distribution of  $I \cup E$  largely overlaps with that of outside data, making the groups difficult to distinguish. This overlap complexity shows that model-oriented signals are no longer sufficient under subset training, highlighting the need for data-oriented alternatives.

### 4.2 CHOICE LEAKAGE ATTACK

**Motivation.** Just as models can overfit to their training data, subset selectors can *overfit at the selection level*: by design they preferentially reselect examples that match their implicit criteria (e.g., high informativeness, low noise, or strong representativeness). This persistent re-selection introduces a stable bias in the choice process that itself serves as a reliable membership signal. We exploit this *inclusion stability*, the tendency of a sample to be repeatedly chosen across multiple trials, as the core signal for our attack.

Specifically, we approximate many different candidate combinations by constructing a series of overlapping subsets (“windows”)  $\{W_i \subseteq D_0\}_{i=1}^m$ , where  $m$  is the number of windows, to capture inclusion-stable samples. Each  $W_i$  represents one plausible candidate set the selector might face; by examining the selector’s decisions on a sample across these windows, we reveal whether it is consistently favored.

**Subset-aware side-channel attack.** In the side-channel setting, the adversary knows both the selector  $\text{Sel}(\cdot; r)$  and the selection ratio  $r \in (0, 1]$ . For each window  $W_i$ , we run  $\text{Sel}(\cdot; r)$  and record whether  $x \in W_i$  is selected by the selector, and get its evidence  $e(x, W_i)$  in the current window:

$$e(x, W_i) = \mathbb{1}[x \in \text{Sel}(W_i; r)]. \quad (4)$$

Suppose in the window construction,  $x$  appears in  $n$  out of  $m$  windows; by aggregating the selection evidence across these windows, we obtain its *inclusion count*:

$$t(x) = \sum_{i=1}^n e(x, W_i), \quad (5)$$

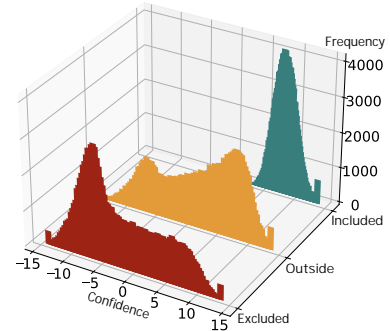


Figure 2: Signal distributions of three groups of data under subset training.

Table 1: Results for vision models under the subset-aware side-channel attack setting. Results are averaged over 9 coresets selection methods. *Intensity* denotes the selection ratio  $r$  (Light:  $r = 0.2$ , Medium:  $r = 0.4$ , Heavy:  $r = 0.6$ , Extensive:  $r = 0.8$ ). Best results per row are in bold.

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	51.23 ±2.56	5.83 ±1.34	51.77 ±3.02	3.34 ±4.45	51.59 ±2.66	6.37 ±2.22	51.26 ±4.85	6.43 ±3.70	<b>61.39</b> ±2.48	<b>14.24</b> ±2.02
	TM-MIA	64.00 ±12.15	12.13 ±5.96	61.57 ±15.05	8.93 ±13.52	67.24 ±16.18	19.30 ±17.14	69.86 ±22.08	15.74 ±18.19	<b>83.77</b> ±2.44	<b>42.19</b> ±4.51
Medium	SP-MIA	52.33 ±3.56	6.31 ±1.48	53.59 ±4.30	3.56 ±2.28	54.51 ±4.51	6.64 ±1.58	54.99 ±4.69	5.31 ±0.42	<b>81.93</b> ±3.50	<b>42.66</b> ±5.81
	TM-MIA	59.84 ±12.84	10.80 ±4.97	60.37 ±10.53	2.91 ±3.32	66.84 ±11.79	12.51 ±5.85	62.96 ±13.69	4.61 ±2.52	<b>88.53</b> ±2.55	<b>60.10</b> ±7.62
Heavy	SP-MIA	52.21 ±3.83	12.20 ±15.48	53.20 ±4.85	2.80 ±2.43	53.53 ±5.94	12.37 ±15.44	53.69 ±5.68	4.26 ±1.77	<b>96.86</b> ±2.60	<b>88.60</b> ±5.51
	TM-MIA	55.00 ±9.59	19.31 ±26.18	52.40 ±10.78	1.67 ±2.04	57.44 ±11.06	19.63 ±26.72	52.61 ±11.77	2.81 ±2.32	<b>89.06</b> ±1.90	<b>60.36</b> ±5.87
Extensive	SP-MIA	55.64 ±5.31	7.59 ±1.68	59.56 ±6.39	4.00 ±2.92	56.66 ±5.90	7.60 ±1.87	61.54 ±8.36	5.09 ±2.68	<b>92.20</b> ±6.94	<b>91.86</b> ±7.23
	TM-MIA	61.41 ±6.63	10.99 ±2.61	60.13 ±12.20	4.21 ±4.15	62.80 ±7.52	11.27 ±2.73	59.66 ±12.03	4.63 ±3.77	<b>80.74</b> ±8.23	<b>49.76</b> ±6.98

where  $t(x)$  is the number of times  $x$  is selected, For fair comparison, the windows are constructed as sliding windows with fixed intervals and cyclic wrapping (details are provided in Section 5), thus each data appears in exactly the same number of windows. Hence, the exposure count  $n$  is constant across all  $x$  and serves only as a scaling factor in our score function. This also highlights the motivation behind our multi-shot membership signal: rather than relying on a single output, choice leakage signal is derived from *how consistently a sample is selected across different selections*. The membership score  $s_{\text{Side}}(x)$  is obtained by aggregating evidence across windows:

$$s_{\text{side}}(x, n, r) = w(t(x); n, r), \quad (6)$$

where  $w$  is a monotone weighting function. From a statistical perspective, if each inclusion is a Bernoulli trial, then  $t(x) \sim \text{Binomial}(n(x), p(x))$  where  $p(x)$  is the probability of a data to be included. Given the selection ratio  $r$ , the expected inclusion count under random choice is  $r \cdot n(x)$ . We can therefore design  $w$  as a smooth monotone mapping centered around  $r \cdot n(x)$ :

$$w(t(x); n(x), r) = \frac{\sigma(\kappa(t(x) - r \cdot n(x)))}{Z(n(x), r)}, \quad \sigma(u) = \frac{1}{1 + e^{-u}}, \quad \kappa > 0, \quad (7)$$

where  $\kappa$  controls the slope and  $Z$  is a normalization constant (depending only on  $n(x), r$ ) that does not affect relative ranking. Since the ratio  $r \in (0, 1]$  and each sample has the same exposure count  $n$ . Without loss of generality, we therefore adopt the following simplified scoring function:

$$w(t(x); n) = \sigma\left(t(x) - \frac{n}{2}\right) = \frac{1}{1 + e^{-(t(x) - \frac{n}{2})}}. \quad (8)$$

This formulation monotonically amplifies scores of samples with high inclusion counts and constrains the range by  $n$ , which makes scores comparable across windows. Finally, under both TM-MIA and SP-MIA, the decision is made by thresholding:

$$\hat{y}(x) = \mathbb{1}[s_{\text{side}}(x) \geq \tau], \quad (9)$$

where  $\tau$  is a decision threshold. Samples that are more stably selected as included data across windows will receive higher scores and are thus more likely to be classified as training members.

**Black-box attack.** In this setting, the subset selection process remains a black box to the adversary, and no direct selection metadata is available. Guided by our general motivation of *inclusion stability* (samples that are repeatedly reselected across plausible candidate sets reveal membership), we infer stable inclusion by identifying samples that consistently act as geometric representatives across windows. Specifically, for each window we perform unsupervised embedding clustering to locate representative samples. Formally, let  $f(\cdot)$  be an embedding model. For each window  $W_i \subseteq \mathcal{D}_t$ , we compute embeddings  $f(x), x \in W_i$ , and perform k-means clustering (Ahmed et al., 2020) in the embedding space. Each sample  $x \in W_i$  is then assigned to a cluster  $c(x; W_i)$ , and we measure its distance to the corresponding cluster centroid  $d(x, W_i) = \|f(x) - c(x; W_i)\|_2$ . The distance is used to serve as the evidence:

$$e(x, W_i) = \mathbb{1}[d(x, W_i) \leq Q_{0.5}(W_i)], \quad (10)$$

where  $Q_{0.5}(\cdot)$  is the median distance among all samples in  $W_i$ . The formal definitions of the inclusion count and exposure count follow the same formulation as in Eq. 5, with the only difference that the evidence  $e(x, W_i)$  is redefined as Eq. 10 under the current black-box setting.

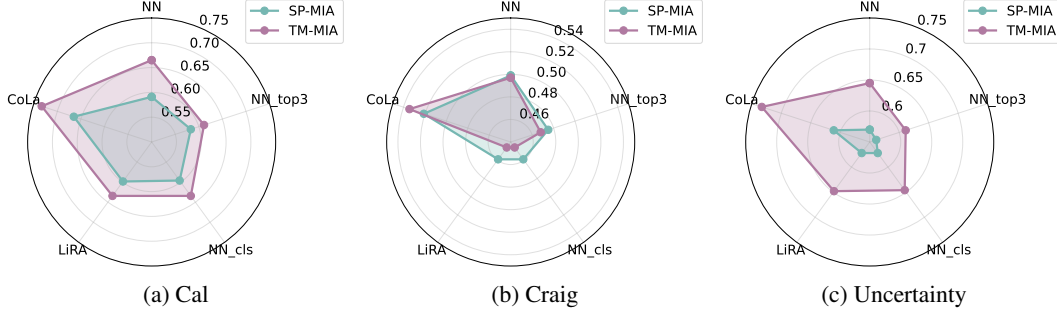


Figure 3: The MIA performance on vision models under black-box setting.

Here, to capture multi-shot stability, since the evidence for each data now related the distance to its centroid in each window  $W_i$ , we apply a weighted score function which reveals not only the inclusion count but also the actual distance it receives:

$$s_{\text{black}}(x) = w(t(x); n) / \bar{d}(x), \quad (11)$$

where  $\bar{d}(x) = \frac{1}{t(x)} \sum_{i: x \in W_i} d(x, W_i)$  denotes the average clustering distance of sample  $x$  across the windows in which it is included. This design ensures that samples consistently close to centroids across many windows receive higher scores. The weighting function  $w(t; n)$  follows the same formulation as in Eq. 8. Finally, similar to the side-channel setting, membership is determined by thresholding:

$$\hat{y}(x) = \mathbb{1}[s_{\text{black}}(x) \geq \tau]. \quad (12)$$

This unsupervised formulation enables membership inference even without any knowledge of the underlying subset selection metadata. The inclusion stability-based pipeline of CoLa naturally unifies different attack surfaces within a single framework, thereby facilitating coordinated attacks.

## 5 EXPERIMENTS

### 5.1 SETUPS

**Models and Datasets.** We conduct experiments on both vision and language models. For the vision side, without loss of generality, we use ResNet-18 trained on CIFAR-10. We evaluate the performance on both subset-aware side-channel attacks and black-box attacks. For language models, since training multiple LMs from scratch is computationally expensive, we restrict our study to black-box attacks. Leveraging the rich open-source models in NLP and following the setup in (Meeus et al., 2024), we use deduplicated models from the Pythia (Biderman et al., 2023b) and GPT-Neo (Black et al., 2021) families, specifically pythia-70m, pythia-160m, and gpt-neo-125m, all trained on the MIMIR dataset (Gao et al., 2020b; Duan et al., 2024). From the MIMIR dataset, we select two subsets, arXiv and PubMed Central, and evaluate each under two split settings: ‘arxiv\_ngram\_1\_0.8’, ‘arxiv\_ngram\_13\_0.2’, ‘pubmed\_central\_ngram\_13\_0.8’, and ‘pubmed\_central\_ngram\_13\_0.2’, where ‘13\_0.8’ denotes removing non-member examples that share  $> 80\%$  13-gram overlap with members.

In the black-box attacks for vision models, we derive embeddings from the activations just before the final linear layer of a shadow model that shares the target model’s architecture. The shadow model is trained using the GradMatch method (Killamsetty et al., 2021a) (distinct from the MIA methods evaluated in our paper) with a selection rate of 0.5. For language models, due to the various lengths of each sequence, we obtain fixed-dimensional embeddings using a dedicated embedding model; by default we use ‘all-MiniLM-L6-v2’ (Reimers & Gurevych, 2019; Thakur et al., 2021).

For CoLa, the default interval is set to 500 for vision models and 100 for language models, with the window size to be 20,000 and 1,000, respectively. In black-box attacks, the number of clusters is fixed at 5. Ablation studies are provided in Section 5.4.

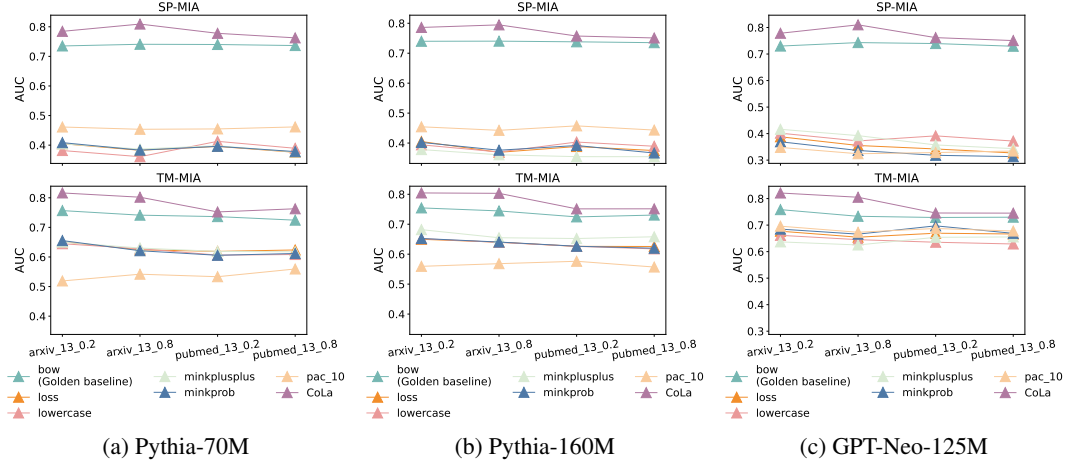


Figure 4: The MIA performance on language models under black-box setting.

**Subset Selection Methods.** For vision models, we select nine representative dataset pruning methods from different categories. Specifically, we include decision boundary based methods such as DeepFool (Ducoffe & Precioso, 2018) and Contrastive Active Learning (Cal) (Margatina et al., 2021); the bi-level optimization based method Glister (Killamsetty et al., 2021b); error based methods including Forgetting (Toneva et al., 2018) and GraNd (Paul et al., 2021); the uncertainty based method Least Confidence (denoted as Uncertainty) (Coleman et al., 2020); the gradient matching based method Craig (Mirzasoleiman et al., 2020); and geometry based methods such as Contextual Diversity (Agarwal et al., 2020) and Herding (Welling, 2009). These methods cover a broad range of perspectives on dataset pruning, from boundary sensitivity to optimization criteria, error contribution, uncertainty, gradient alignment, and geometric diversity. The selection ratio is set to 0.2, 0.4, 0.6, and 0.8. For language models, as discussed in the previous subsection, we adopt a commonly used data filtering strategy that has been systematically studied in (Meeus et al., 2024; Duan et al., 2024), and consider two deduplication strengths, namely ‘13\_0.8’ and ‘13\_0.2’.

**Baseline MIA Methods.** For vision models, we consider four baselines: NN, NN\_top3, and NN\_CIs (Shokri et al., 2017; Salem et al., 2018), which use the model’s output logits, the top-3 logits, and the combination of logits with class labels as membership signals, respectively, as well as LiRA (Carlini et al., 2022b), which fits Gaussian distributions and leverages the likelihood to infer membership. The shadow model used in each baseline method is set to 8. For language models, we consider six baselines, including the loss (Yeom et al., 2018), Lower (lowercase) (Carlini et al., 2021), Min-K% (minkprob) (Shi et al., 2023), Min-K%++ (minkplusplus) (Zhang et al., 2024), Pac (pac\_10) (Ye et al., 2024), and the Golden baseline Bag\_of\_Words (bow) (Meeus et al., 2024). Here, bow serves as a performance reference: methods performing below it are regarded as ineffective.

**Evaluation Metrics.** In most MIA studies (Hisamoto et al., 2020; Carlini et al., 2022a; Li et al., 2025), attack performance is typically evaluated by aggregating over all possible thresholds using the AUC score. We adopt the same practical evaluation metric in our experiments. We also report True Positive Rate at low False Positive Rate (TPR@Low FPR) (Carlini et al., 2022b), which is an important metric in MIAs and measures the detection rate at a meaningful threshold.

## 5.2 RESULTS UNDER SUBSET-AWARE SIDE-CHANNEL ATTACKS

A subset-aware side-channel attack is a type of attack specific to the subset selection process. Its success indicates that current practices of disclosing meta information about subset selection are unsafe and can lead to privacy leakage.

Table 1 reports the average MIA results across different coreset selection methods we consider for vision models (detailed results for each method are provided in Appendix A.2). As shown, in the relatively simple TM-MIA setting, baseline methods can still perform reasonably well, which is expected since this setting closely resembles traditional MIAs (Shokri et al., 2017; Hu et al., 2022) for which these baselines were originally designed.



However, in the SP-MIA setting that is unique to subset training, baseline methods largely fail (AUC close to 50%), indicating their inability to effectively distinguish between included and excluded data. Fundamentally, this stems from the fact that baseline methods rely heavily on model outputs; as illustrated in Figure 2, included and excluded data exhibit output distributions that are highly similar to other data, resulting in poor separability. However, this does not mean that privacy cannot be compromised under SP-MIA. In contrast, CoLa achieves strong performance in both TM-MIA and SP-MIA settings, thanks to its multi-shot, data-centric membership signal that tightly aligns with the subset selection process and captures fine-grained data interactions, thereby enabling better separability. Moreover, we observe that as the selection ratio (*Intensity*) increases, the risk of privacy leakage becomes more severe, highlighting the significant vulnerability of the subset selection process as a potential side channel.

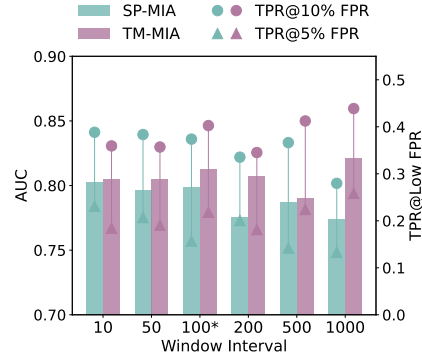


Figure 5: The influence of the window size on the MIA performance.

### 5.3 RESULTS UNDER BLACK-BOX ATTACKS

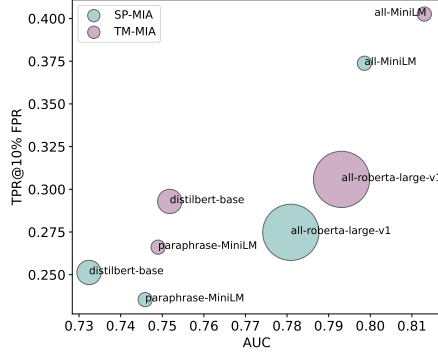


Figure 6: The influence of the embedding model on the MIA performance.

The results for vision models and language models are shown in Figure 3 and Figure 4, respectively. For vision models, we adopt three representative selection methods: Cal (Margatina et al., 2021), Craig (Mirzasoileiman et al., 2020), and Uncertainty (Coleman et al., 2020). As illustrated in Figure 3, under the black-box setting, SP-MIA remains more challenging than TM-MIA. Moreover, CoLa consistently outperforms the baselines by about 5% in AUC across all experiments, demonstrating strong attack capability. For language models, this contrast is even more pronounced. As shown in Figure 4, all baseline methods except CoLa perform worse than the bow baseline, indicating that they essentially fail in the context of subset selection MIA. Furthermore, while SP-MIA and TM-MIA results are relatively close for CoLa, the baselines exhibit a sharp gap, with SP-MIA close to random guessing (AUC around 50%), and TM-MIA reaches only about 60%.

### 5.4 ABLATION STUDIES.

**Influence of Window Construction.** In Figure 5, we present an ablation study on the influence of window interval, conducted with Pythia-160m on the arxiv\_ngram\_13\_0.8 dataset. Several observations can be made: first, regardless of the window interval size, the performance under SP-MIA is consistently lower than that under TM-MIA, highlighting its greater challenge. Second, the choice of window interval size does not substantially affect the performance of CoLa. In SP-MIA, increasing the size reduces the exposure count  $n$  of each data sample, which makes the inclusion signal coarser and leads to a slight performance drop. However, this drop remains marginal.

**Influence of Embedding Model.** As a data-centric MIA method, CoLa achieves a clear decoupling from the target model. As discussed earlier, it derives the membership signal by reallocat-



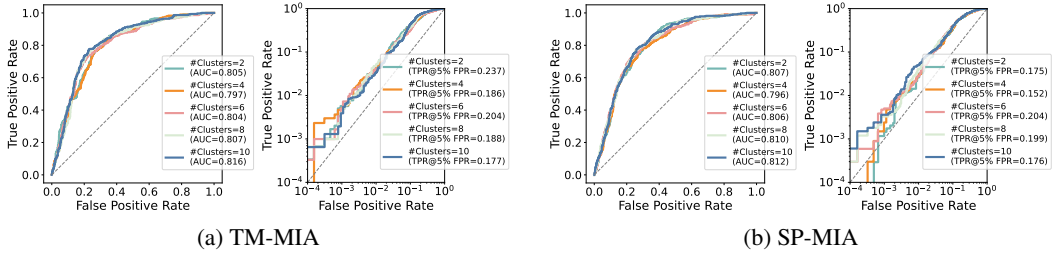


Figure 7: The MIA performance on language models under the black-box setting.

ing data combinations based on overfitting at the selection level. For language data, the inherent inconsistency in format and length requires the use of a dedicated embedding model in this reallocation process. To examine the effect of embedding model choice, we conduct an ablation study beyond the default all-MiniLM-L6-v2, considering three alternatives: paraphrase-MiniLM-L6-v2 (paraphrase-MiniLM), distilbert-base-nli-stsb-mean-tokens (distilbert-base), and all-roberta-large-v1. The results are shown in Figure 6, where the circle size indicates the parameter scale of each embedding model. We observe that different embedding models have a noticeable impact on inference performance, particularly on TPR at low FPR. Moreover, larger model size does not necessarily translate into better performance, highlighting the importance of choosing an appropriate embedding model. Nevertheless, the results remain generally acceptable across all choices (with AUC consistently above 70% and TPR@10% FPR above 25%). How to customize embedding models for MIA under subset selection is a meaningful question, which we leave for future work.

Table 2: Subset-aware Side-channel attacks under different vision models and datasets.

Setting	ResNet18-CIFAR100		VGG19-CIFAR10		VGG19-CIFAR100	
	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
SP-MIA	67.28	19.05	64.98	17.15	70.31	21.23
	$\pm 1.36$	$\pm 1.17$	$\pm 2.03$	$\pm 2.12$	$\pm 1.64$	$\pm 1.81$
TM-MIA	85.53	38.46	81.43	40.35	86.67	42.10
	$\pm 2.07$	$\pm 1.43$	$\pm 1.43$	$\pm 1.65$	$\pm 2.36$	$\pm 1.39$

**Results under Different Vision Models and Datasets.** In Table 2, we further conduct subset-aware side-channel attack on the CIFAR-100 dataset with the VGG19 model to verify whether CoLa remains reliable across different vision datasets and models. The selection ratio here is set to 0.2. As can be observed, CoLa consistently works well across various vision model-dataset combinations, revealing its general applicability. Specifically, attacks on VGG19 are more pronounced than on ResNet18 under the same setting, and CIFAR-100 is more vulnerable than CIFAR-10. Moreover, the observation that SP-MIA is more challenging than TM-MIA is consistent with previous findings.

**Influence of Clustering.** In Figure 7, we study the effect of varying the number of clusters used for embedding clustering in the black-box setting. Beyond the default choice of 5, we further consider values between 2 and 10 and report the corresponding AUC curves and TPR@5% FPR. The results show that, for both SP-MIA and TM-MIA, the clustering number has only a marginal effect on performance.

## 6 CONCLUSION

In this work, we take the first step toward systematically understanding the privacy risks of subset training. Contrary to the common intuition that training on fewer samples should reduce privacy leakage, we demonstrate that the very choices made during subset selection can themselves become exploitable signals, exposing both included and excluded data to membership inference. To capture this phenomenon, we introduced CoLa, a unified framework that leverages choice patterns to construct robust membership signals. Across both vision and language models, under both subset-aware side-channel and black-box settings, CoLa consistently outperforms existing baselines, revealing that subset training does not mitigate but instead amplifies privacy leakage. Our findings highlight that privacy risks extend beyond model outputs to the data-model supply chain itself. We hope this work motivates future efforts toward designing selection mechanisms and training pipelines that are not only efficient and scalable but also privacy-preserving.

## ETHICS STATEMENT

This work focuses on understanding privacy risks in subset training through systematic analysis of membership inference attacks (MIAs). Our study is purely methodological and does not involve human subjects or personally identifiable information. All datasets used are publicly available benchmark datasets (e.g., CIFAR, GSM8K, CodeAlpaca), and we complied with their intended use and licensing terms. We emphasize that the proposed Choice Leakage Attack (CoLa) is presented as a research contribution to highlight potential vulnerabilities in modern training pipelines, not to enable misuse. Our findings are intended to inform the community about inherent privacy risks and to guide the development of stronger defenses. No proprietary or sensitive data was used, and no deployed models were targeted in this study. In line with research integrity, we also note that Large Language Models (LLMs) were only employed for literature review support and polishing of textual presentation (e.g., improving fluency and figure/table captions). LLMs were not involved in technical design, experimental implementation, or data analysis.

## REPRODUCIBILITY STATEMENT

We have made every effort to ensure the reproducibility of our work. All datasets used in this paper are publicly available, and their sources are clearly cited in the main manuscript. The implementation details of our methods, including models used, attack configurations, and evaluation protocols, are described in Section 5.1. We also provide ablation studies and additional experiments in Section 5.4 to validate the generality of our findings. Upon acceptance, we will release the full source code, configuration files, and scripts for evaluation to facilitate verification and future research.

## LLM DISCLAIMER

LLMs were used only occasionally for language polishing, aiming to improve fluency and readability. All technical ideas, experimental designs, analyses, conclusions, writing were developed and carried out entirely by the authors. The authors have full responsibility for the final text.

## REFERENCES

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pp. 308–318, 2016.
- Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI 16*, pp. 137–153. Springer, 2020.
- Mohiuddin Ahmed, Raihan Seraj, and Syed Mohammed Shamsul Islam. The k-means algorithm: A comprehensive survey and performance evaluation. *Electronics*, 9(8):1295, 2020.
- Olivier Bachem, Mario Lucic, and Andreas Krause. Coresets for nonparametric estimation-the case of dp-means. In *International Conference on Machine Learning*, pp. 209–217. PMLR, 2015.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory 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 Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 2397–2430. PMLR, 2023a. URL <https://proceedings.mlr.press/v202/biderman23a.html>.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023b.

- Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow, March 2021. URL <https://doi.org/10.5281/zenodo.5297715>. If you use this software, please cite it using these metadata.
- Zalán Borsos, Mojmir Mutny, and Andreas Krause. Coresets via bilevel optimization for continual learning and streaming. *Advances in neural information processing systems*, 33:14879–14890, 2020.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX security symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. Membership inference attacks from first principles. In *43rd IEEE Symposium on Security and Privacy, SP 2022, San Francisco, CA, USA, May 22-26, 2022*, pp. 1897–1914. IEEE, 2022a. doi: 10.1109/SP46214.2022.9833649. URL <https://doi.org/10.1109/SP46214.2022.9833649>.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Membership inference attacks from first principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1897–1914. IEEE, 2022b.
- Vincent Cohen-Addad, David Saulpic, and Chris Schwiegelshohn. A new coreset framework for clustering. In Samir Khuller and Virginia Vassilevska Williams (eds.), *STOC '21: 53rd Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June 21-25, 2021*, pp. 169–182. ACM, 2021. doi: 10.1145/3406325.3451022. URL <https://doi.org/10.1145/3406325.3451022>.
- C Coleman, C Yeh, S Musmann, B Mirzasoleiman, P Bailis, P Liang, J Leskovec, and M Zaharia. Selection via proxy: Efficient data selection for deep learning. In *International Conference on Learning Representations (ICLR)*, 2020.
- Tian Dong, Bo Zhao, and Lingjuan Lyu. Privacy for free: How does dataset condensation help privacy? In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 5378–5396. PMLR, 2022. URL <https://proceedings.mlr.press/v162/dong22c.html>.
- Michael Duan, Anshuman Suri, Niloofar Miresghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. Do membership inference attacks work on large language models? *arXiv preprint arXiv:2402.07841*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The

- llama 3 herd of models. *CoRR*, abs/2407.21783, 2024. doi: 10.48550/ARXIV.2407.21783. URL <https://doi.org/10.48550/arXiv.2407.21783>.
- Melanie Ducoffe and Frederic Precioso. Adversarial active learning for deep networks: a margin based approach. *arXiv preprint arXiv:1802.09841*, 2018.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset of diverse text for language modeling, 2020a. URL <https://arxiv.org/abs/2101.00027>.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020b.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks are all you need. *CoRR*, abs/2306.11644, 2023. doi: 10.48550/ARXIV.2306.11644. URL <https://doi.org/10.48550/arXiv.2306.11644>.
- Sorami Hisamoto, Matt Post, and Kevin Duh. Membership inference attacks on sequence-to-sequence models: Is my data in your machine translation system? *Transactions of the Association for Computational Linguistics*, 8:49–63, 2020. doi: 10.1162/tacl.a.00299. URL <https://aclanthology.org/2020.tacl-1.4/>.
- Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37, 2022.
- Yuzheng Hu, Fan Wu, Qinbin Li, Yunhui Long, Gonzalo Munilla Garrido, Chang Ge, Bolin Ding, David Forsyth, Bo Li, and Dawn Song. Sok: Privacy-preserving data synthesis. In *2024 IEEE Symposium on Security and Privacy (SP)*, pp. 4696–4713, 2024. doi: 10.1109/SP54263.2024.00002.
- Krishnateja Killamsetty, Sivasubramanian Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-match: Gradient matching based data subset selection for efficient deep model training. In *International Conference on Machine Learning*, pp. 5464–5474. PMLR, 2021a.
- Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh Iyer. Glisten: Generalization based data subset selection for efficient and robust learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 8110–8118, 2021b.
- 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 (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8424–8445, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.577. URL <https://aclanthology.org/2022.acl-long.577/>.
- Qi Li, Runpeng Yu, and Xinchao Wang. Vid-sme: Membership inference attacks against large video understanding models. *arXiv preprint arXiv:2506.03179*, 2025.
- Adyasha Maharana, Prateek Yadav, and Mohit Bansal. D2 pruning: Message passing for balancing diversity & difficulty in data pruning. In *The Twelfth International Conference on Learning Representations*, 2023.
- Katerina Margatina, Giorgos Vernikos, Loïc Barrault, and Nikolaos Aletras. Active learning by acquiring contrastive examples. pp. 650–663, 2021. doi: 10.18653/v1/2021.EMNLP-MAIN.51. URL <https://doi.org/10.18653/v1/2021.emnlp-main.51>.

- Matthieu Meeus, Igor Shilov, Shubham Jain, Manuel Faysse, Marek Rei, and Yves-Alexandre de Montjoye. Sok: Membership inference attacks on llms are rushing nowhere (and how to fix it). *arXiv preprint arXiv:2406.17975*, 2024.
- Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models. In *International Conference on Machine Learning*, pp. 6950–6960. PMLR, 2020.
- Alexander Munteanu, Chris Schwiegelshohn, Christian Sohler, and David Woodruff. On coresets for logistic regression. *Advances in Neural Information Processing Systems*, 31, 2018.
- Milad Nasr, Reza Shokri, and Amir Houmansadr. Machine learning with membership privacy using adversarial regularization. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS ’18*, pp. 634–646, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450356930. doi: 10.1145/3243734.3243855. URL <https://doi.org/10.1145/3243734.3243855>.
- Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. *Advances in Neural Information Processing Systems*, 34:20596–20607, 2021.
- Ru Peng, Kexin Yang, Yawen Zeng, Junyang Lin, Dayiheng Liu, and Junbo Zhao. Dataman: Data manager for pre-training large language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=eNbA8Fqir4>.
- Ziheng Qin, Kai Wang, Zangwei Zheng, Jianyang Gu, Xiangyu Peng, Daquan Zhou, Lei Shang, Baigui Sun, Xuansong Xie, Yang You, et al. Infobatch: Lossless training speed up by unbiased dynamic data pruning. In *The Twelfth International Conference on Learning Representations*, 2023.
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, H. Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant M. Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew J. Johnson, Blake A. Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Edward Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher. *CoRR*, abs/2112.11446, 2021. URL <https://arxiv.org/abs/2112.11446>.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. URL <https://arxiv.org/abs/1910.10683>.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.
- Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Hervé Jégou. White-box vs black-box: Bayes optimal strategies for membership inference. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings*

- of *Machine Learning Research*, pp. 5558–5567. PMLR, 2019. URL <http://proceedings.mlr.press/v97/sablayrolles19a.html>.
- Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models. *arXiv preprint arXiv:1806.01246*, 2018.
- Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. 2018. URL <https://openreview.net/forum?id=H1aIuk-RW>.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv preprint arXiv:2310.16789*, 2023.
- 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)*, pp. 3–18. IEEE, 2017.
- Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *Advances in Neural Information Processing Systems*, 35:19523–19536, 2022.
- Theresa Stadler, Bristena Oprisanu, and Carmela Troncoso. Synthetic data - anonymisation groundhog day. In Kevin R. B. Butler and Kurt Thomas (eds.), *31st USENIX Security Symposium, USENIX Security 2022, Boston, MA, USA, August 10-12, 2022*, pp. 1451–1468. USENIX Association, 2022. URL <https://www.usenix.org/conference/usenixsecurity22/presentation/stadler>.
- Bowen Tan, Zheng Xu, Eric P. Xing, Zhiting Hu, and Shanshan Wu. Synthesizing privacy-preserving text data via finetuning \*without\* finetuning billion-scale LLMs. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=FCm4laCLiH>.
- Haoru Tan, Sitong Wu, Fei Du, Yukang Chen, Zhibin Wang, Fan Wang, and Xiaojuan Qi. Data pruning via moving-one-sample-out. *Advances in Neural Information Processing Systems*, 36, 2024.
- Nandan Thakur, Nils Reimers, Johannes Daxenberger, and Iryna Gurevych. Augmented SBERT: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 296–310, Online, June 2021. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2021.naacl-main.28>.
- Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. In *International Conference on Learning Representations*, 2018.
- Boris van Breugel, Hao Sun, Zhaozhi Qian, and Mihaela van der Schaar. Membership inference attacks against synthetic data through overfitting detection. In Francisco J. R. Ruiz, Jennifer G. Dy, and Jan-Willem van de Meent (eds.), *International Conference on Artificial Intelligence and Statistics, 25-27 April 2023, Palau de Congressos, Valencia, Spain*, volume 206 of *Proceedings of Machine Learning Research*, pp. 3493–3514. PMLR, 2023. URL <https://proceedings.mlr.press/v206/breugel23a.html>.
- Max Welling. Herding dynamical weights to learn. In *Proceedings of the 26th annual international conference on machine learning*, pp. 1121–1128, 2009.
- Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. Qurating: selecting high-quality data for training language models. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24. JMLR.org*, 2024.



- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *CoRR*, abs/2412.15115, 2024a. doi: 10.48550/ARXIV.2412.15115. URL <https://doi.org/10.48550/arXiv.2412.15115>.
- Shuo Yang, Zeke Xie, Hanyu Peng, Min Xu, Mingming Sun, and Ping Li. Dataset pruning: Reducing training data by examining generalization influence. In *The Eleventh International Conference on Learning Representations*, 2022.
- Shuo Yang, Zhe Cao, Sheng Guo, Ruiheng Zhang, Ping Luo, Shengping Zhang, and Liqiang Nie. Mind the boundary: Coreset selection via reconstructing the decision boundary. In *Forty-first International Conference on Machine Learning*, 2024b.
- Wentao Ye, Jiaqi Hu, Liyao Li, Haobo Wang, Gang Chen, and Junbo Zhao. Data contamination calibration for black-box llms. *arXiv preprint arXiv:2405.11930*, 2024.
- Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *2018 IEEE 31st computer security foundations symposium (CSF)*, pp. 268–282. IEEE, 2018.
- Jingyang Zhang, Jingwei Sun, Eric Yeats, Yang Ouyang, Martin Kuo, Jianyi Zhang, Hao Frank Yang, and Hai Li. Min-k%++: Improved baseline for detecting pre-training data from large language models. *arXiv preprint arXiv:2404.02936*, 2024.
- Yunpeng Zhao and Jie Zhang. Does training with synthetic data truly protect privacy? In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=C8niXBHjfo>.

## A APPENDIX

### A.1 THE PRIVACY THREATS BEHIND DATA-MODEL SUPPLY CHAIN

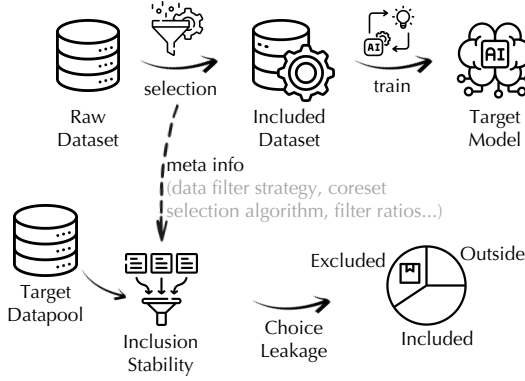


Figure 8: Choice Leakage Attack (CoLa) across the data-model supply chain. CoLa augments conventional MIA by exploiting subset selection metadata leaked along the data-model supply chain. By identifying which samples are more likely to pass selection, it not only strengthens membership inference but also enables adversaries to craft tailored threats.

choice leakage risk is severe as it not only amplifies the risk of inferring membership but also exposes a system’s selection preferences. Once the data-model supply chain is exposed to privacy risks, the entire pipeline, from raw data to model outputs, becomes vulnerable to malicious manipulation. For example, adversaries may learn proxies of the selection rule and craft targeted poisoning or backdoor examples that are more likely to bypass filtering and enter training.

As shown in Figure 8, the data-model supply chain describes the pipeline from raw data collection, through subset selection and model training, to the deployment of a target model. In this process, subset selection plays a central role: only a fraction of the raw dataset is included for training, while others are excluded or remain outside. The metadata of this selection process (e.g., filtering strategies, coreset algorithms, or filter ratios) introduces new privacy surfaces. Such information can inadvertently leak “choice signals” that reveal which samples are more likely to be included in training, thereby extending the privacy risk beyond conventional training data exposure.

CoLa (Choice Leakage Attack) directly exploits this vulnerability by leveraging selection metadata to strengthen membership inference. Unlike traditional MIAs that focus solely on the trained model’s outputs, CoLa targets the entire supply chain, identifying which samples are predisposed to pass the selection process. Such

Table 3: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is Cal (Margatina et al., 2021).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.499	0.050	0.501	0.055	0.508	0.053	0.512	0.054	0.602	0.122
	TM-MIA	0.759	0.207	0.676	0.166	0.784	0.257	0.737	0.182	0.855	0.442
Medium	SP-MIA	0.553	0.072	0.573	0.056	0.582	0.074	0.587	0.058	0.789	0.372
	TM-MIA	0.763	0.165	0.759	0.097	0.812	0.227	0.784	0.092	0.878	0.620
Heavy	SP-MIA	0.589	0.077	0.603	0.000	0.630	0.087	0.624	0.054	0.963	0.856
	TM-MIA	0.729	0.123	0.721	0.000	0.772	0.172	0.736	0.058	0.895	0.642
Extensive	SP-MIA	0.634	0.091	0.637	0.000	0.647	0.092	0.651	0.036	0.957	0.954
	TM-MIA	0.717	0.116	0.707	0.061	0.736	0.128	0.690	0.026	0.849	0.573

Table 4: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is Contextual Diversity (Agarwal et al., 2020).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.540	0.067	0.539	0.053	0.548	0.073	0.544	0.052	0.633	0.118
	TM-MIA	0.706	0.125	0.716	0.070	0.755	0.161	0.756	0.072	0.798	0.347
Medium	SP-MIA	0.598	0.094	0.594	0.000	0.614	0.088	0.610	0.056	0.846	0.465
	TM-MIA	0.751	0.158	0.708	0.000	0.792	0.160	0.729	0.049	0.908	0.656
Heavy	SP-MIA	0.507	0.051	0.502	0.000	0.506	0.500	0.500	0.000	0.982	0.904
	TM-MIA	0.502	0.074	0.482	0.000	0.516	0.048	0.477	0.000	0.898	0.631
Extensive	SP-MIA	0.494	0.056	0.494	0.027	0.497	0.052	0.497	0.042	0.967	0.966
	TM-MIA	0.500	0.053	0.490	0.000	0.502	0.052	0.490	0.041	0.843	0.386

## A.2 RESULTS OF VISION MODELS UNDER DIFFERENT SUBSET SELECTION METHODS

In Table 1, we report the average results of vision models across nine subset selection methods. For clarity, Tables 3–11 present the results for each method separately, providing a more straightforward view of the attack performance.

Table 5: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is Craig (Mirzasoleiman et al., 2020).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.495	0.046	0.500	0.000	0.497	0.048	0.441	0.039	0.637	0.172
	TM-MIA	0.567	0.087	0.500	0.000	0.573	0.102	0.602	0.064	0.825	0.411
Medium	SP-MIA	0.513	0.066	0.588	0.054	0.580	0.086	0.598	0.055	0.819	0.367
	TM-MIA	0.595	0.137	0.693	0.043	0.717	0.134	0.716	0.045	0.858	0.518
Heavy	SP-MIA	0.575	0.076	0.614	0.052	0.628	0.082	0.629	0.051	0.969	0.876
	TM-MIA	0.624	0.114	0.628	0.030	0.701	0.122	0.647	0.034	0.888	0.562
Extensive	SP-MIA	0.624	0.092	0.655	0.000	0.653	0.096	0.664	0.000	0.960	0.959
	TM-MIA	0.674	0.110	0.666	0.000	0.700	0.119	0.623	0.000	0.842	0.545

Table 6: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is DeepFool (Ducoffe & Precioso, 2018).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.494	0.057	0.500	0.000	0.489	0.054	0.441	0.039	0.637	0.172
	TM-MIA	0.556	0.092	0.500	0.000	0.530	0.084	0.221	0.000	0.825	0.411
Medium	SP-MIA	0.494	0.051	0.501	0.048	0.492	0.050	0.500	0.050	0.845	0.480
	TM-MIA	0.649	0.088	0.550	0.000	0.642	0.088	0.397	0.011	0.926	0.700
Heavy	SP-MIA	0.496	0.053	0.507	0.000	0.496	0.052	0.509	0.042	0.979	0.900
	TM-MIA	0.494	0.053	0.429	0.016	0.484	0.054	0.424	0.000	0.902	0.643
Extensive	SP-MIA	0.526	0.096	0.643	0.062	0.545	0.097	0.645	0.067	0.956	0.954
	TM-MIA	0.592	0.142	0.571	0.037	0.602	0.140	0.574	0.038	0.858	0.572

Table 7: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is Forgetting (Toneva et al., 2018).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.500	0.053	0.500	0.000	0.514	0.059	0.530	0.051	0.618	0.141
	TM-MIA	0.572	0.099	0.500	0.000	0.706	0.176	0.741	0.071	0.854	0.475
Medium	SP-MIA	0.503	0.056	0.500	0.000	0.548	0.068	0.559	0.056	0.818	0.464
	TM-MIA	0.529	0.098	0.500	0.000	0.695	0.139	0.724	0.064	0.851	0.517
Heavy	SP-MIA	0.501	0.500	0.499	0.045	0.499	0.050	0.498	0.050	0.986	0.943
	TM-MIA	0.540	0.830	0.480	0.000	0.546	0.840	0.460	0.034	0.921	0.661
Extensive	SP-MIA	0.585	0.080	0.640	0.071	0.585	0.081	0.748	0.084	0.791	0.787
	TM-MIA	0.646	0.107	0.640	0.074	0.649	0.107	0.648	0.080	0.653	0.407

Table 8: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is Glister (Killamsetty et al., 2021b).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.495	0.048	0.500	0.000	0.492	0.045	0.545	0.062	0.608	0.135
	TM-MIA	0.477	0.033	0.500	0.000	0.422	0.000	0.883	0.129	0.829	0.384
Medium	SP-MIA	0.504	0.055	0.497	0.044	0.503	0.049	0.496	0.045	0.864	0.494
	TM-MIA	0.367	0.007	0.545	0.045	0.448	0.024	0.586	0.044	0.874	0.516
Heavy	SP-MIA	0.494	0.050	0.499	0.048	0.495	0.048	0.497	0.051	0.992	0.949
	TM-MIA	0.404	0.039	0.541	0.060	0.440	0.020	0.555	0.059	0.871	0.480
Extensive	SP-MIA	0.527	0.062	0.598	0.073	0.533	0.060	0.600	0.079	0.984	0.984
	TM-MIA	0.598	0.131	0.757	0.118	0.651	0.134	0.771	0.120	0.895	0.502

Table 9: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is GraNd (Paul et al., 2021).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.563	0.087	0.584	0.126	0.563	0.114	0.575	0.153	0.562	0.137
	TM-MIA	0.843	0.206	0.918	0.389	0.937	0.571	0.950	0.584	0.878	0.483
Medium	SP-MIA	0.498	0.048	0.498	0.047	0.497	0.050	0.499	0.052	0.754	0.344
	TM-MIA	0.535	0.103	0.471	0.019	0.573	0.104	0.471	0.018	0.902	0.680
Heavy	SP-MIA	0.493	0.047	0.500	0.051	0.493	0.047	0.501	0.050	0.909	0.774
	TM-MIA	0.557	0.119	0.387	0.011	0.562	0.118	0.384	0.012	0.859	0.606
Extensive	SP-MIA	0.505	0.054	0.502	0.047	0.506	0.054	0.503	0.048	0.839	0.826
	TM-MIA	0.572	0.110	0.378	0.005	0.556	0.109	0.380	0.019	0.712	0.498

Table 10: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is Herding (Welling, 2009).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.516	0.053	0.521	0.052	0.512	0.054	0.510	0.053	0.574	0.171
	TM-MIA	0.853	0.407	0.912	0.373	0.932	0.452	0.927	0.389	0.963	0.771
Medium	SP-MIA	0.498	0.050	0.498	0.051	0.498	0.049	0.499	0.051	0.753	0.460
	TM-MIA	0.857	0.244	0.749	0.092	0.861	0.246	0.757	0.088	0.976	0.880
Heavy	SP-MIA	0.543	0.061	0.601	0.059	0.545	0.078	0.600	0.067	0.966	0.846
	TM-MIA	0.782	0.210	0.740	0.029	0.792	0.206	0.741	0.028	0.931	0.729
Extensive	SP-MIA	0.491	0.047	0.498	0.000	0.492	0.049	0.497	0.043	0.964	0.963
	TM-MIA	0.687	0.127	0.542	0.046	0.688	0.126	0.542	0.049	0.862	0.571

Table 11: The results of vision models under Subset-aware Side-channel attacks and the subset selection method used here is Uncertainty (Coleman et al., 2020).

Intensity	Setting	NN		NN_top3		NN_cls		LiRA		CoLa	
		AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR	AUC	TPR@5%FPR
Light	SP-MIA	0.499	0.051	0.499	0.054	0.499	0.049	0.498	0.050	0.603	0.138
	TM-MIA	0.549	0.065	0.458	0.021	0.528	0.063	0.437	0.019	0.827	0.376
Medium	SP-MIA	0.494	0.050	0.498	0.050	0.494	0.05	0.496	0.050	0.811	0.454
	TM-MIA	0.614	0.073	0.444	0.020	0.610	0.072	0.433	0.019	0.914	0.703
Heavy	SP-MIA	0.554	0.089	0.625	0.054	0.560	0.089	0.627	0.051	0.959	0.850
	TM-MIA	0.709	0.13	0.644	0.025	0.713	0.130	0.644	0.025	0.899	0.644
Extensive	SP-MIA	0.501	0.050	0.506	0.044	0.502	0.051	0.502	0.050	0.929	0.924
	TM-MIA	0.614	0.117	0.424	0.000	0.607	0.118	0.425	0.03	0.823	0.575