

Attribute-to-Delete: Machine Unlearning via Datamodel Matching

Presenter: Ting-Wei Li
11/7/2024



Outline

- **Topics:** Machine Unlearning (MU), Data Attribution (DA)
 - Definition and Evaluation of MU
 - Definition of DA
- **Background:** *Datamodels: Predicting Predictions from Training Data*¹
- **Paper:** *Attribute-to-Delete: Machine Unlearning via Datamodel Matching*²
 - Motivation
 - Datamodels
 - MU algorithm: DMM (Datamodel Matching)
 - Evaluation
 - Future Directions

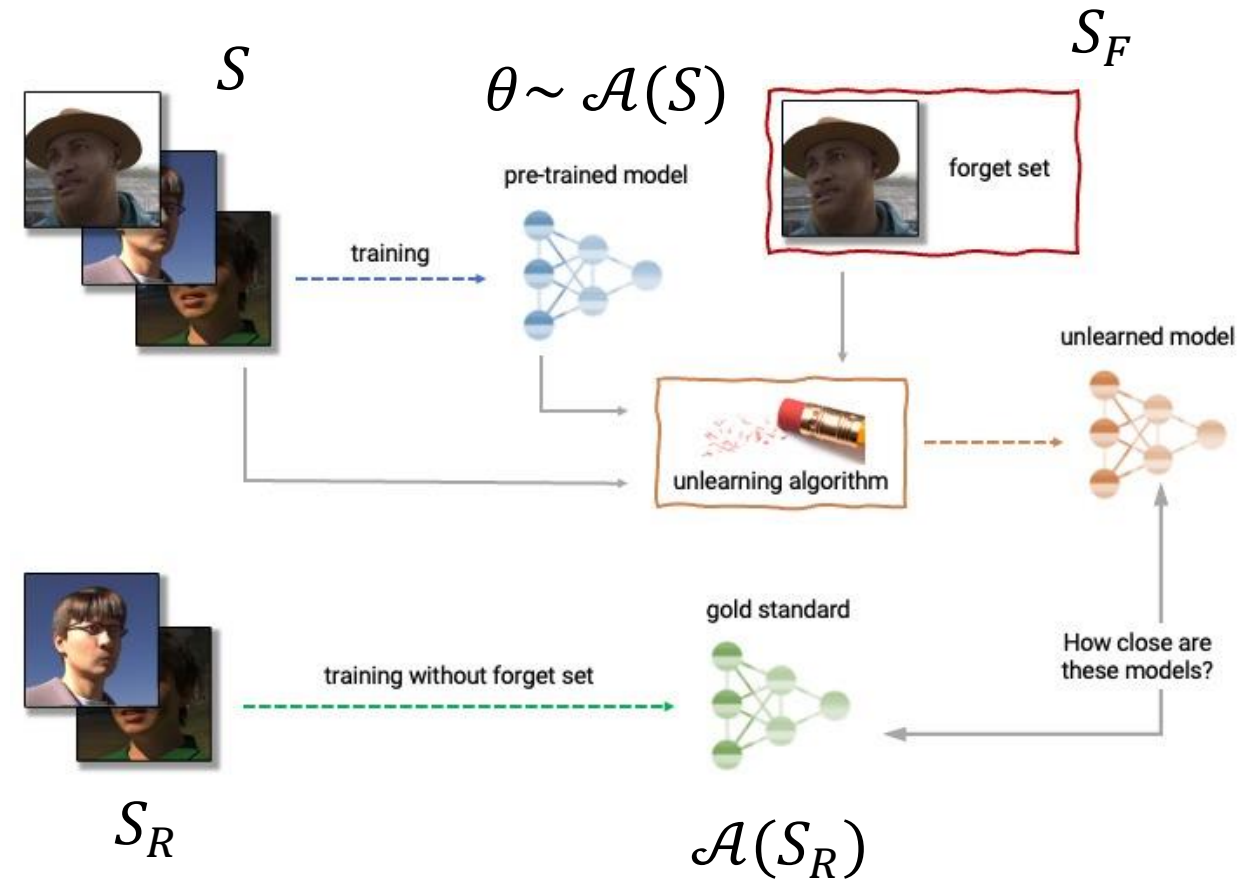
Notations

- $S \in \mathcal{X}^n$ be a fixed training dataset (with size n) drawn from an example space \mathcal{X} (data universe)
- a *learning algorithm* $\mathcal{A}: \mathcal{X}^* \rightarrow \Theta$ as a (possibly random) function mapping from **datasets** to a **ML model** $\theta \in \Theta$ (parameter universe)
- for any example $x \in \mathcal{X}$, we use $f_x: \Theta \rightarrow \mathbb{R}^k$ to denote a *model evaluation* on x . (i.e. an **example-specific output function** where model param is a variable)
- Focus on *classification problem* for this paper!
- f_x can be k-dimensional logits or **margins** of a classifier.

Can be $\log \frac{p}{1-p}$ or the logit for the correct class minus the second highest class

Machine Unlearning (MU)

- Given (1) a ML model $\theta \sim \mathcal{A}(S)$ trained on dataset S and (2) a *forget set* $S_F \subset S$ (also, a corresponding *retain set* $S_R = S \setminus S_F$).
- We want:** a model sampled from $\mathcal{A}(S_R)$ starting from the trained model θ .
- A MU algorithm $\mathcal{U}: \Theta \times 2^{|S|} \rightarrow \Theta$ takes a learned model θ , training dataset S , and the forget set $S_F \subset S$ as input and output a “unlearned model” that is *indistinguishable from a model trained on S_R* .



Approximate Formulation for MU

- Definition 1 ((ϵ, δ) -unlearning.⁴) a MU algorithm \mathcal{U} is an (ϵ, δ) -approximate unlearning algorithm if, for all $\mathcal{O} \subset \Theta$, $S_F \subset S$, we have that

Some “memberships”

The unlearned model that forgets S_F from trained model $\theta \sim \mathcal{A}(S)$

$$\Pr[\mathcal{U}(\mathcal{A}(S), S_F) \in \mathcal{O}] \leq e^\epsilon \Pr[\mathcal{A}(S_R) \in \mathcal{O}] + \delta$$

$$\Pr[\mathcal{A}(S_R) \in \mathcal{O}] \leq e^\epsilon \Pr[\mathcal{U}(\mathcal{A}(S), S_F) \in \mathcal{O}] + \delta$$

The **golden** model that is trained on S_R from scratch

- Intuition: the distributions of unlearned models are (ϵ, δ) - indistinguishable from retrained models (**in parameter space!**)
- A hint for a *direction evaluation metric* for MU

A Direct Evaluation for MU

Similar to the golden standard model but more broad

- Compute *the divergence of model outputs* between unlearned models and **safe models** (not compare them in the parameter space!)
- Safe models ($safe(S_F)$): some models that have not been trained on S_F .
- Formally, we want $\Delta_\delta(\mathcal{U}(\mathcal{A}(S), S_F), safe(S_F)) \leq \epsilon$, where Δ_δ is a δ -approximate divergence and $safe(S_F)$ is a distribution of safe models w.r.t. forget set S_F .

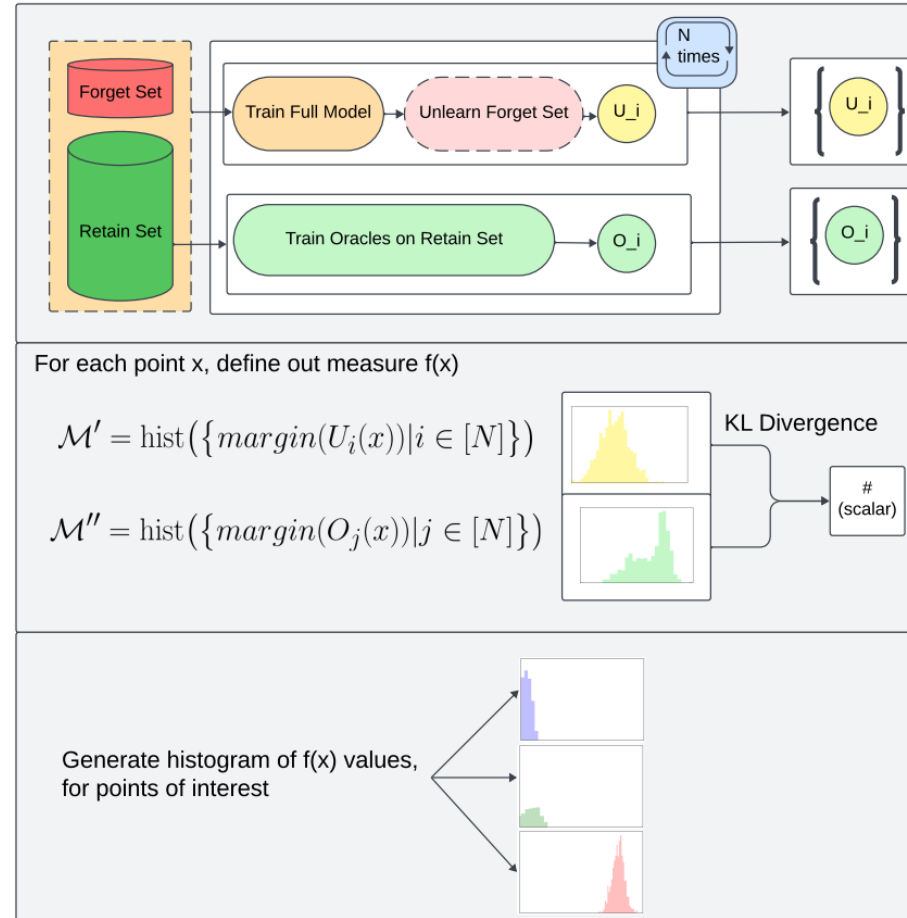
The unlearned model that forgets S_F from trained model θ

KLoM: A Direct Evaluation

- We thus have the definition of *KL divergence of margins* (KLoM) that (1) uses $\Delta = \text{KL divergence}$, (2) allows arbitrary $\text{safe}(S_F)$ and (3) studies distributions of model outputs (margins) f_x rather than parameters.
- Definition 2: (*KL divergence of margins (KLoM)*.) For an unlearning algorithm \mathcal{U} , reference distribution $\text{safe}(S_F)$, and input x , *KLoM* is given by

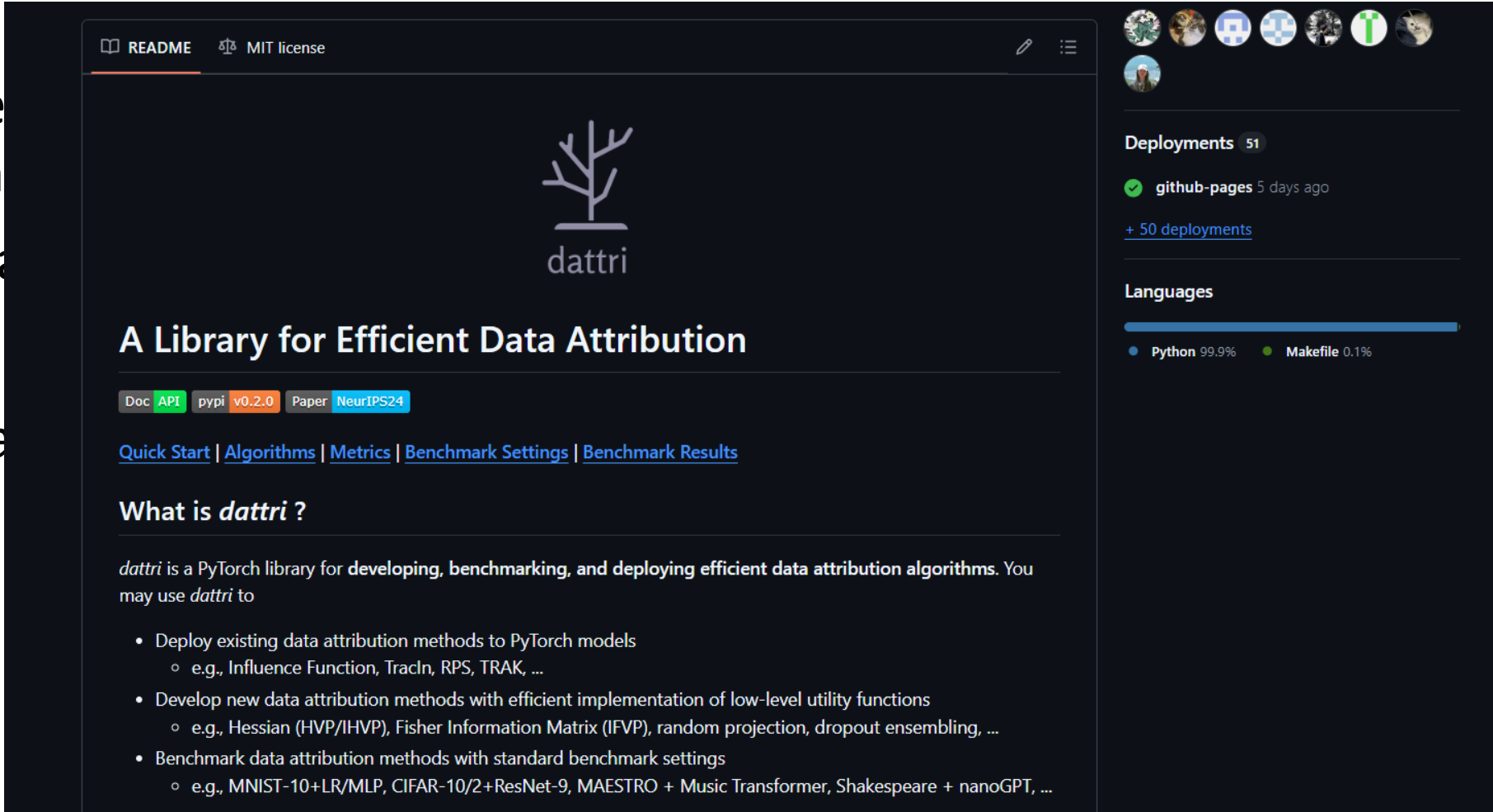
$$KLoM(\mathcal{U}) := D_{KL} \left(\text{safe}(S_F), f_x(\mathcal{U}(\mathcal{A}(S), S_F)) \right).$$
- Mostly, we just take $\text{safe}(S_F) := \mathcal{A}(S \setminus S_F)$. And we will compute KLoM for data inputs x from S_F, S_R to evaluate MU algorithm.
- Intuitively, KLoM should be small for inputs x from all sets.

KLoM: A Direct Evaluation



Predictive Data Attribution

- The central role of training
- Or what if we could predict the impact of each training step?
- Popular methods
- A recent work



The screenshot shows the GitHub repository for **dattri**, a library for efficient data attribution. The repository is licensed under MIT and has a README file. The main content of the README is as follows:

A Library for Efficient Data Attribution

Doc **API** pypi v0.2.0 Paper NeurIPS24

[Quick Start](#) | [Algorithms](#) | [Metrics](#) | [Benchmark Settings](#) | [Benchmark Results](#)

What is *dattri* ?

dattri is a PyTorch library for **developing, benchmarking, and deploying efficient data attribution algorithms**. You may use *dattri* to

- Deploy existing data attribution methods to PyTorch models
 - e.g., Influence Function, Tracln, RPS, TRAK, ...
- Develop new data attribution methods with efficient implementation of low-level utility functions
 - e.g., Hessian (HVP/IHVP), Fisher Information Matrix (IFVP), random projection, dropout ensembling, ...
- Benchmark data attribution methods with standard benchmark settings
 - e.g., MNIST-10+LR/MLP, CIFAR-10/2+ResNet-9, MAESTRO + Music Transformer, Shakespeare + nanoGPT, ...

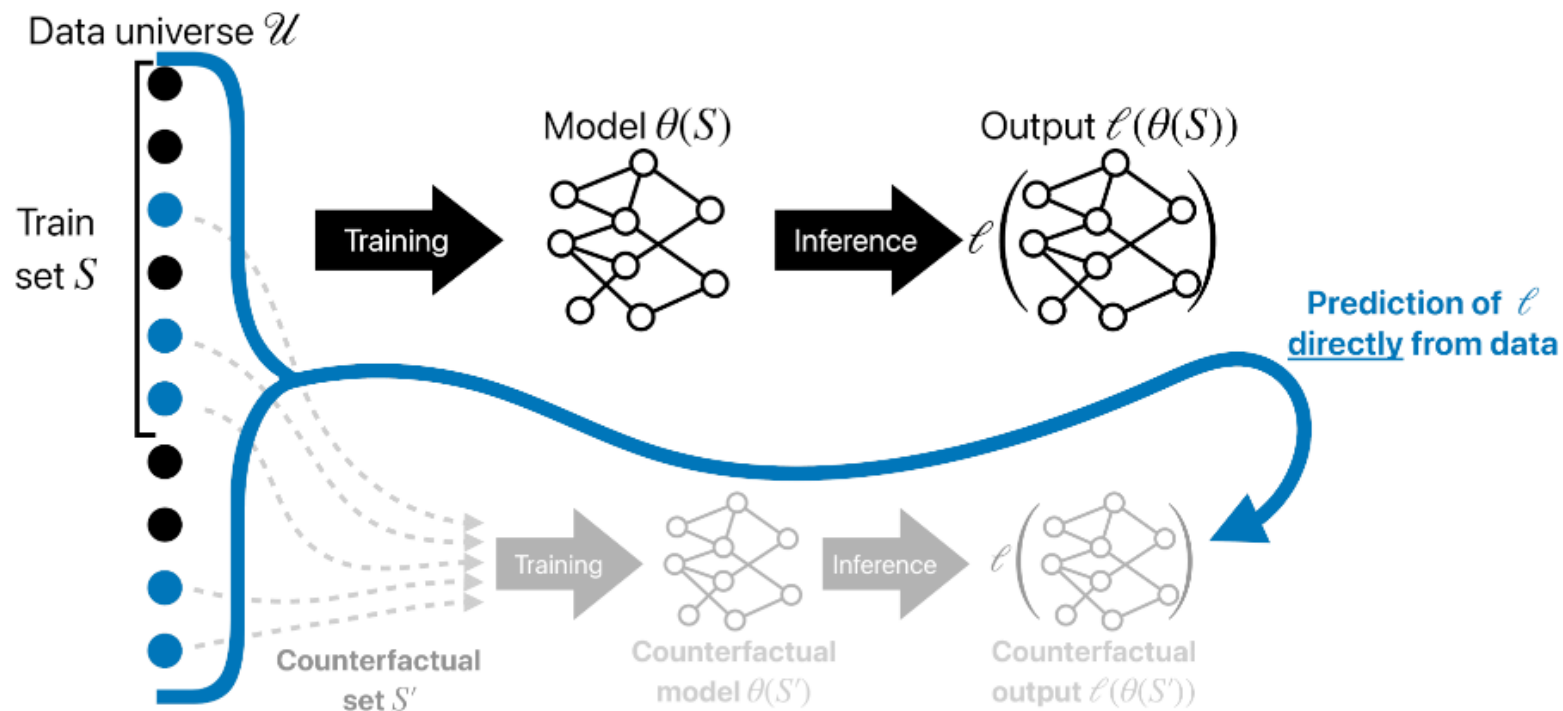
On the right side of the screenshot, the GitHub interface shows 51 deployments, with the most recent being **github-pages** 5 days ago. The language usage is shown as a bar chart: **Python 99.9%** and **Makefile 0.1%**.

Predictive Data Attribution: Datamodeling

- Goal: compute an estimator (or *datamodel*¹) that takes a training set as input and predicts the behavior of a ML model (trained on this training set) as output.
- Formally, for an example $x \in \mathcal{X}$, a *datamodel for x* is a function $\hat{f}: 2^{|S|} \rightarrow \mathbb{R}$ such that, for any $S' \subset S$, $\hat{f}(S') \approx f_x(\mathcal{A}(S'))$.
- Intuition: *datamodel* takes an indicator vector representing training indexes as input and output approximate model evaluation on x , where this model can be trained on arbitrary subset of the training set.

Predictive Data Attribution

Predictive attribution (Datamodeling)



Delving deeper into Datamodels

- Surprisingly, *linear datamodels* work well to accurately predict modern complex ML/DL model behaviors.
- Formally, for an example x , we want to compute *a coefficient vector* $\beta \in \mathbb{R}^{|S|}$ such that, for subsets $S' \subset S$,

$$\hat{f}(S') := \beta^T \mathbf{1}_{S'} = \sum_{i \in S'} \beta_i \approx f_x(\mathcal{A}(S')).$$

Diagram annotations:

- A blue arrow points from the text "Datamodel output on S' " to the boxed term $\hat{f}(S')$.
- A blue arrow points from the boxed term $f_x(\mathcal{A}(S'))$ to the text "output of x evaluated on the model trained on S' ".

- Intuition: For each example x , each training sample in S has its own “importance”. For a specific learning algorithm \mathcal{A} , the model evaluation of x can be seen as a **summation over data importance used for training**.

How to learn Datamodels

- We first sample many subsets S_1, \dots, S_k and use \mathcal{A} to train corresponding model.
- We then compute the model evaluation of x on these models (i.e. $f_x(\mathcal{A}(S_i))$, for $i \in [k]$). Then we simply solve the following *regularized regression problem*:

$$\beta = \min_{w \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m \left(w^T \mathbf{1}_{S_i} - f_x(\mathcal{A}(S_i)) \right)^2 + \lambda \|w\|_1,$$

where $\mathbf{1}_{S_i} \in \mathbb{R}^{|S|}$ is a binary vector indicating the subset S_i .

Motivation: The Missing Targets Problem

For complex models, may simple MU algorithm work?

- *Gradient Descent* on retain set may not alter forget set predictions much
- If choose to *increase* loss on forget set (*Gradient Ascent*)
 - > loss for forget set points might **overshoot** or **undershoot** since we don't know the expected loss under a perfectly retrained model (i.e. unlearn **too much** or **not enough**)

Motivation: Case Study on SCRUB

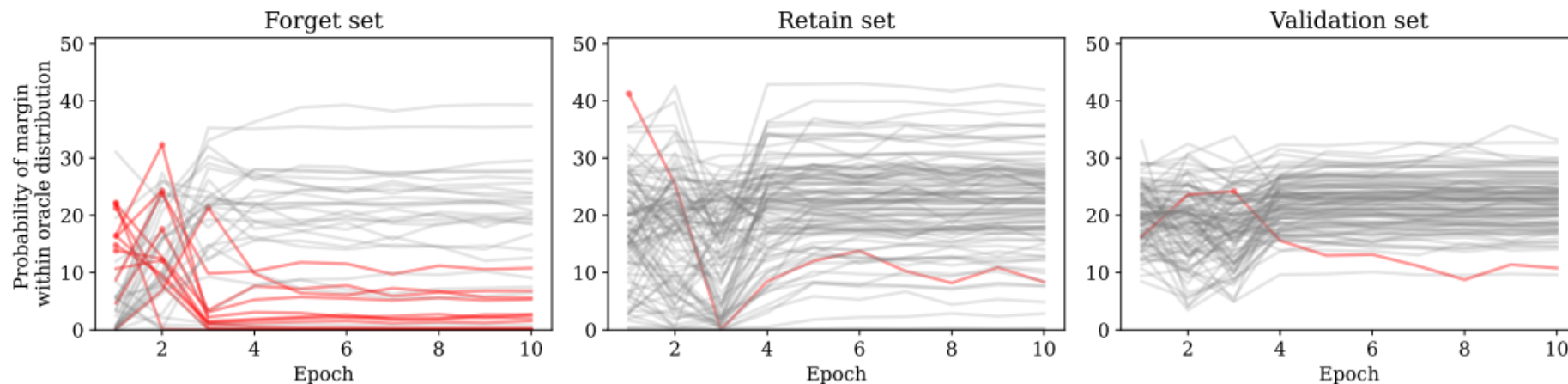


Figure 2: **The missing targets problem.** We apply the SCRUB [KTT24] algorithm to unlearn a forget set of CIFAR-10, and measure how well different (random) points are unlearned over time. To quantify how well a given point x is unlearned, we fit a Gaussian distribution to the outputs of oracle models at x , and compute the likelihood of the average outputs from unlearned models under this distribution. We track this likelihood (y-axis) for random points across the duration of unlearning algorithm (x-axis). For many examples in the forget set (shown in red), unlearning quality is hurt by training for too long as we lack access to oracle targets.

Oracle Matching (OM) Algorithm

- Question: Given access to **sample outputs from the oracle model** (re-trained without the forget set, i.e. $f^{oracle}(x) := f_x(\mathcal{A}(S_R))$),
- can we efficiently fine-tune an existing model (trained on the full dataset) to **match** the outputs of the sample?
- Main point: *whether gradient-based optimization can match predictions of the oracle.*

Oracle Matching (OM) Algorithm

Algorithm 1 Oracle Matching (OM)

```

1: Input: Trained model  $\theta$ ; oracle predictions  $f^{\text{oracle}}(x)$ ; fine-tuning retain set size  $r$ 
2: Output: Unlearned model  $\theta_{UL}$ 
3: Initialize  $\theta_0 = \theta$ 
4: for  $t = \{0, \dots, T - 1\}$  do                                     ▷  $T$  epochs
5:    $S'_R \leftarrow S \setminus S_F$                                      ▷ Sub-sample  $r$  points from retain set
6:    $S_{\text{fine-tune}} = S_F \cup S'_R$                                    ▷ Combine with forget set
7:   for  $x \sim S_{\text{fine-tune}}$  do
8:      $L(\theta_t) = \|f_x(\theta_t) - f^{\text{oracle}}(x)\|^2$                ▷ Compute MSE loss
9:      $\theta_{t+1} = \theta_t - \eta_t \cdot \nabla_{\theta} L(\theta_t)$          ▷ Perform update with gradient
10:  end for
11: end for
12: Return Model  $\theta_{UL} = \theta_T$ 

```

Evaluating OM

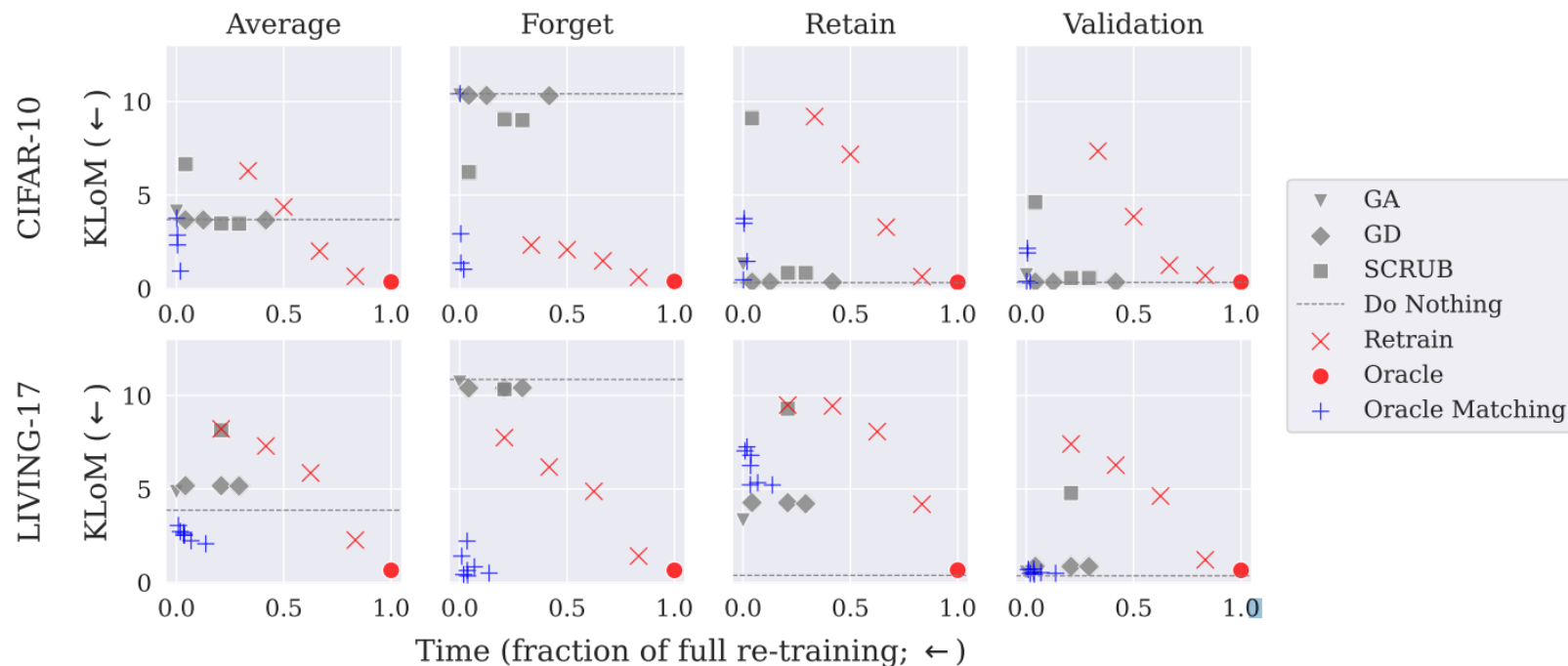


Figure 3: **Oracle matching can efficiently approximate re-training.** The KLoM metric (y-axis) measures the distributional difference between unlearned predictions and oracle predictions (0 being perfect). We also show the amount of compute relative to full re-training (x-axis). We evaluate KLoM values over points in the forget, retain, and validation sets and report the 95th percentile in each group; we also report the average across groups (1st column).

Implications of OM

- OM can constantly outperform other baselines in terms of KLoM scores and running time.
- Observation 1: in parameter space, there is indeed θ' that is close to the oracle retrained without the forget set.
- Observation 2: fine-tuning θ can quickly converge to that θ' with small number of samples of oracle outputs.
- But OM is not a practical algorithm.
- Predictive DA (our favorite *Datamodel*) comes to the rescue!

Efficient Proxy of Oracles: Datamodels

- Intuition: let us replace the oracle models by *datamodels*!
- Formally, for each input x , we replace $f^{oracle}(x)$ by \hat{f}_x , which is the *datamodel* of x .
- Reminder: a *datamodel* maps arbitrary data subset to model output on input x , i.e., $\hat{f}_x(S \setminus S_F) \approx f_x(\mathcal{A}(S \setminus S_F)) = f_x(\mathcal{A}(S_R)) := f^{oracle}(x)$.
- Note that if we focus on *linear datamodel*, then we have

$$\hat{f}_x(S \setminus S_F) := \beta_0 + \sum_{i \in S \setminus S_F} \beta_i(x) = (\beta_0 + \sum_{i \in S} \beta_i(x)) - \sum_{i \in S_F} \beta_i(x) = f_x(\theta) - \sum_{i \in S_F} \beta_i(x).$$

- We don't need to approximate the first term; we know that in priori.

DM-Direct

- Thus, we now focus on the oracle output function:

$$h(x) := f_x(\theta) - \sum_{i \in S_F} \beta_i(x).$$

- This can be a working MU method if we only care about prediction but not model weights. But still, we may (1) want to delete data from full model θ and (2) avoid training one datamodel for each new input z .

Algorithm A.2 DM-DIRECT

- 1: **Input:** Trained model θ ; datamodels $\beta(x)$ for each $x \in S$; forget set S_F
 - 2: **Output:** A predictor $h(\cdot) : S \mapsto \mathbb{R}^k$
 - 3: $h(x) := f_x(\theta_0) - \sum_{i \in S_F} \beta_i(x)$
 - 4: **End**
-

Datamodel Matching (DMM)

- One step further: let us replace the “true oracle” in OM algorithm by $h(x)$ through estimated datamodels (for all x we are interested in).
- Datamodels are estimated through a subset of retain points and the entire set of forget points.

Datamodel Matching (DMM)

Algorithm A.3 Datamodel Matching (DMM)

```

1: Input: Trained model  $\theta$ ; datamodels  $\beta(\cdot)$ ; fine-tuning set size  $r$ 
2: Output: Unlearned model  $\theta_{UL}$ 
3:  $S'_R \leftarrow S \setminus S_F$ ;  $S_{\text{fine-tune}} = S_F \cup S'_R$ 
4:  $h \leftarrow \text{DM-DIRECT}(\theta, \beta, S_f)$ 
5: for  $t = \{0, \dots, T - 1\}$  do
6:   for  $x \sim S_{\text{fine-tune}}$  do
7:      $L(\theta_t) = \|f_x(\theta_t) - h(x)\|^2$ 
8:      $\theta_{t+1} = \theta_t - \eta_t \cdot \nabla_{\theta} L(\theta_t)$ 
9:   end for
10: end for
11: Return Model  $\theta_{UL} = \theta_T$ 

```

▷ Simulate oracles with datamodels

▷ T epochs

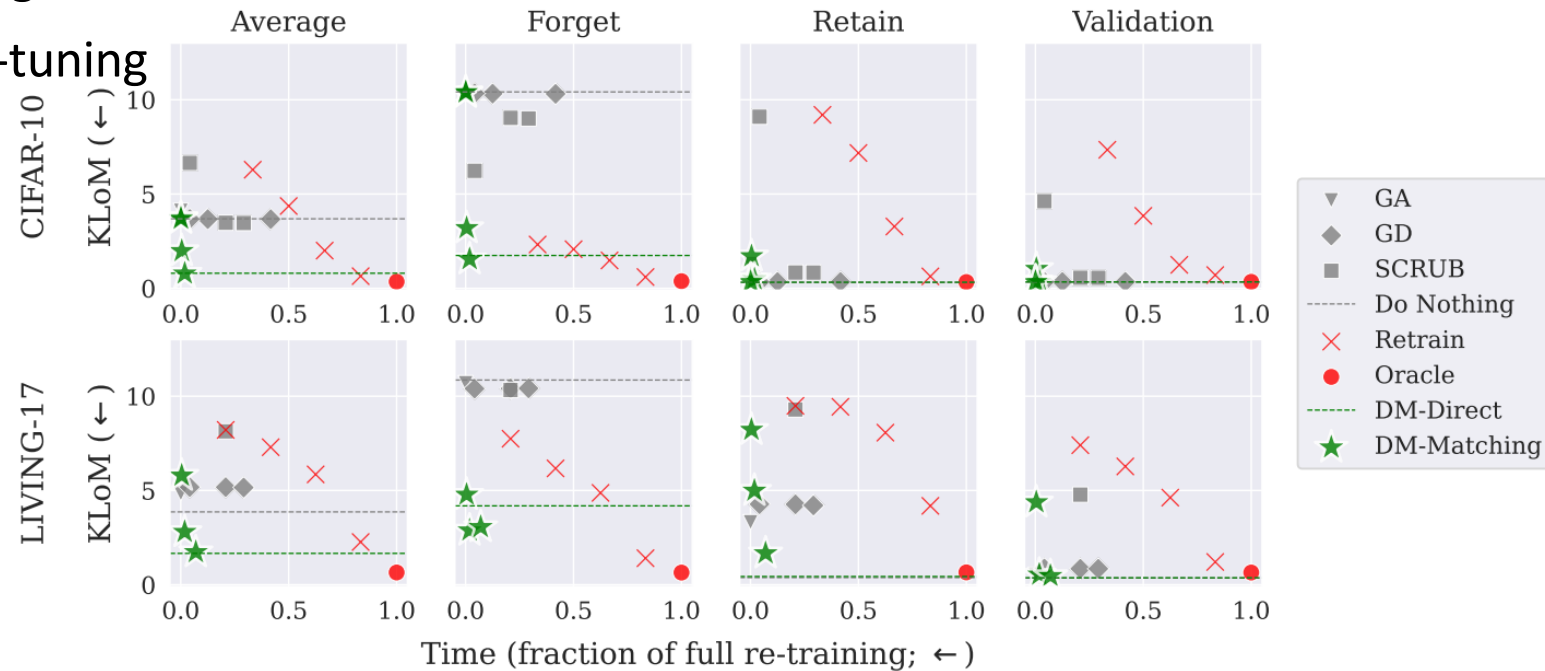
▷ mini-batch

▷ Compute loss

▷ Perform update with gradient

Evaluation through KLoM

- GA: gradient ascent on Forget set starting from θ
- GD: gradient descent on Retain set starting from θ
- Re-train: re-train on Retain set starting from scratch
- DM-Direct: directly use $h(x)$, no fine-tuning
- DMM: use $h(x)$ and fine-tuning
- SCRUB⁷: previous SOTA



Ablation of DMM (can be skipped)

- DMM consists of (i) oracle matching (the fine-tuning algorithm) and (ii) estimating datamodels (approximating oracle outputs)
- Let us study two components separately
 - OM: stability across unlearning time and generalization from small sample size
 - Datamodel: scaling with cost and efficient alternative (TRAK⁸)
- Intuition: does DMM work as expected?

Stability: SCRUB vs OM

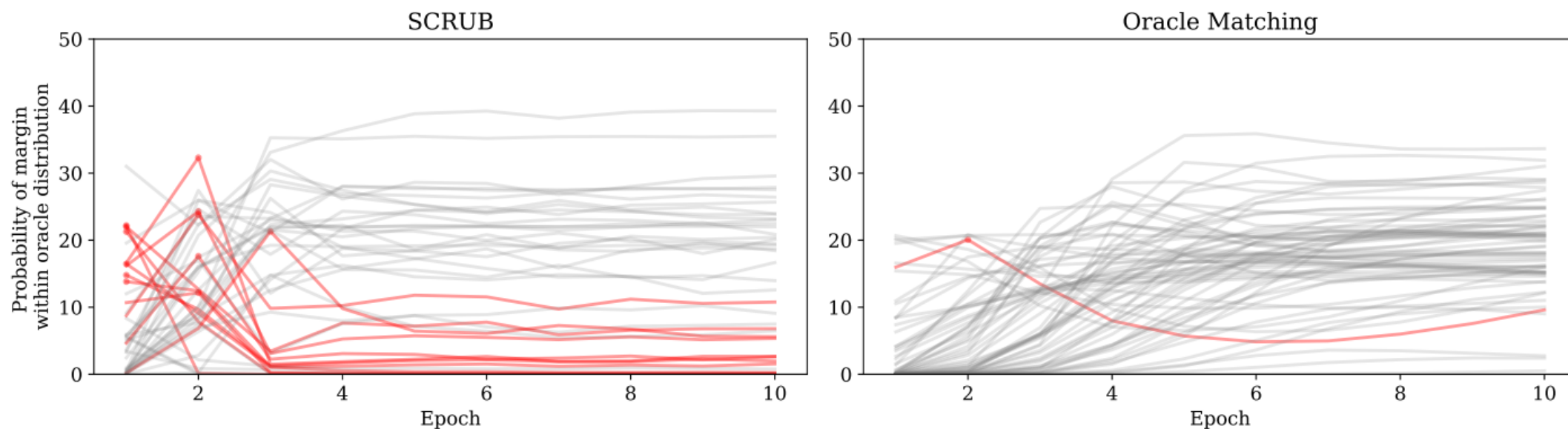


Figure 5: **Oracle matching circumvents the stopping time problem.** We revisit the earlier analysis for SCRUB (left) and apply the same analysis to Oracle Matching (right). The red lines highlight examples in the forget set whose unlearning quality is hurt by training longer. This “overshooting” happens frequently with SCRUB, but only rarely with Oracle Matching.

Generalization of OM

- Focus: the amount of forget/retain data used to derive “oracles”.

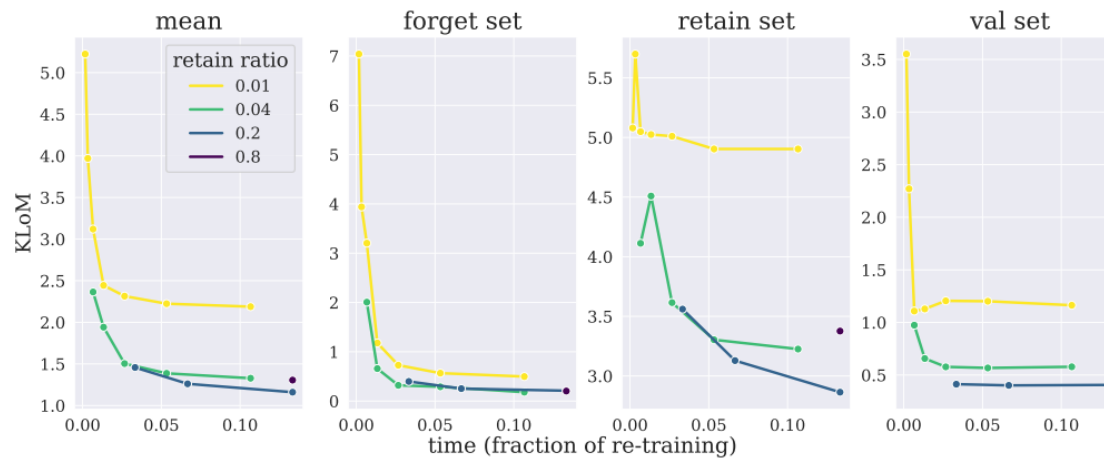


Figure 6: Varying the fraction of retain set (of the full dataset) sampled for oracle matching. A sufficiently large fraction (≥ 0.04) appears to be sufficient in enabling OM to generalize to out-of-sample. In other words, OM can “distill” the oracle model using a small subset of the training data.

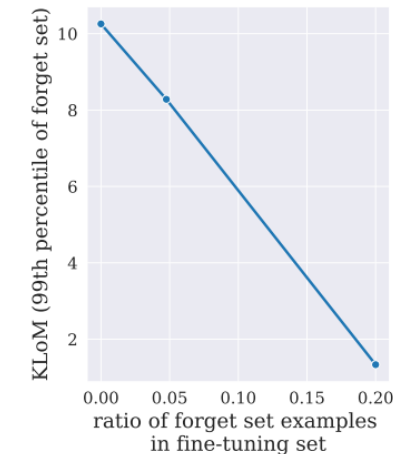


Figure 7: Higher ratio of forget set samples in the fine-tuning set for OM improves performance on the forget set.

Scaling costs for Datamodel

- Natural trade-off: computational cost vs datamodel predictiveness
- Observation: unlearning effectiveness will saturate even the datamodels are still improving
- Implication: efficient alternatives might exist!

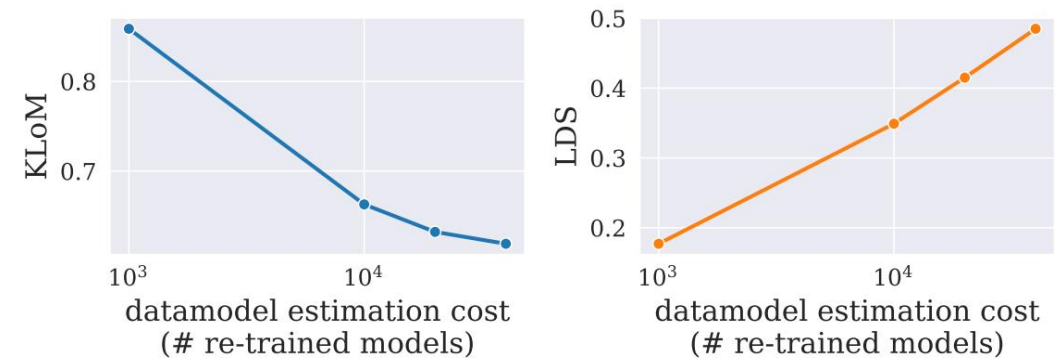
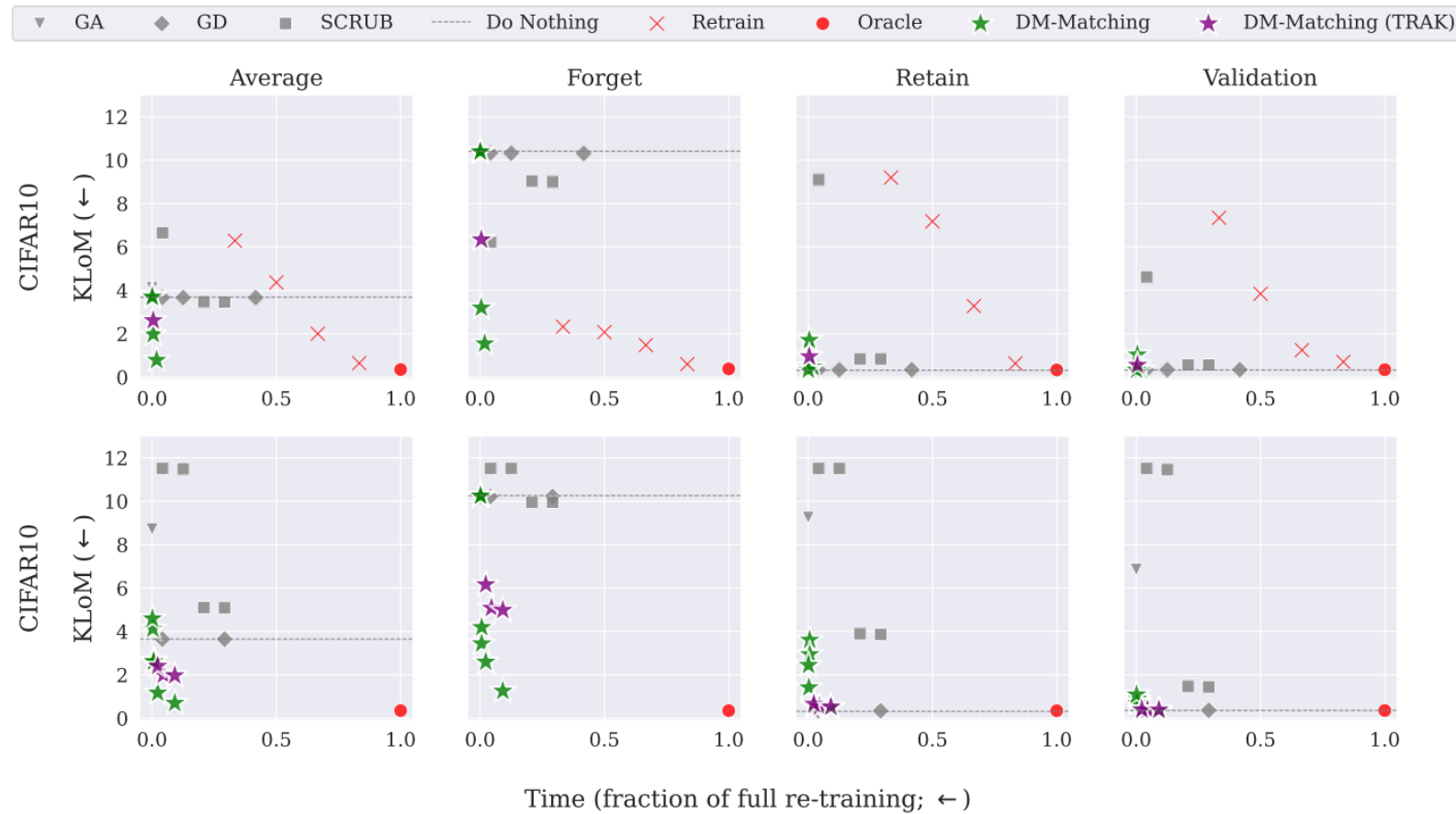


Figure 9: The effect of estimation cost on unlearning performance and datamodel predictiveness. On CIFAR-10, we show how unlearning performance of DM-DIRECT (measured by KLoM; lower is better) and datamodel predictiveness (measured by the linear datamodeling score; higher is better) scales with datamodel estimation cost (number of re-trained models, in $\{10^3, 10^4, 2 \times 10^4, 4 \times 10^4\}$). KLoM is averaged over different forget sets.

Efficient unlearning with TRAK

- TRAK⁸: a more efficient predictive DA algorithm compared to datamodels. (1000x less computation time)
- Observation: natural trade-off again but here with cheaper alternative, DMM still outperform baselines.
- Note: no matter what kind of approximate oracle we use (datamodels or TRAK), this is a one-time cost.

Efficient unlearning with TRAK



Takeaways

- Challenge in machine unlearning: missing targets problem and efficiency
- Predictive data attribution model: help us connect data and complex model output through interpretable, simple modeling (*linear datamodel*)
- **A general framework for MU (DMM)**: fine-tuning current model to match approximate oracles, which is produced through predictive DA algorithm
- **A more direct evaluation for MU (KLoM)**: measure the difference in distributions of unlearned model's outputs from oracle

Future Directions

- Beyond classification problems
- Better understanding of OM – how to sample from retain set
- Better time efficiency – originated from DA algorithm
- More practical scenarios – multiple unlearning requests

Reference

- [1] Ilyas, Andrew, et al. "Datamodels: Understanding predictions with data and data with predictions." International Conference on Machine Learning. PMLR, 2022.
- [2] Georgiev, Kristian, et al. "Attribute-to-Delete: Machine Unlearning via Datamodel Matching." arXiv preprint arXiv:2410.23232 (2024).
- [3] Ginart, Antonio, et al. "Making ai forget you: Data deletion in machine learning." Advances in neural information processing systems 32 (2019).
- [4] Neel, Seth, Aaron Roth, and Saeed Sharifi-Malvajerdi. "Descent-to-delete: Gradient-based methods for machine unlearning." Algorithmic Learning Theory. PMLR, 2021.
- [5] Hayes, Jamie, et al. "Inexact unlearning needs more careful evaluations to avoid a false sense of privacy." arXiv preprint arXiv:2403.01218 (2024).
- [6] Deng, Junwei, et al. "\$\texttt{dattri}\$: A Library for Efficient Data Attribution." arXiv preprint arXiv:2410.04555 (2024).
- [7] Kurmanji, Meghdad, et al. "Towards unbounded machine unlearning." Advances in neural information processing systems 36 (2024).

Thank you!

