# CoUn: Empowering Machine Unlearning via Contrastive Learning

Yasser H. Khalil      Mehdi Setayesh      Hongliang Li

Huawei Noah's Ark Lab, Montreal, Canada

E-mail: yasser1.khalil@huawei.ca

# Introduction

- **Data Privacy & Compliance**
  - Regulations such as the **GDPR** highlight the growing importance of data protection and responsible AI.

- **Machine Unlearning (MU)**
  - Removes the effect of specific training data (*forget set*) while preserving knowledge from remaining data (*retain set*).

- **Exact Unlearning – the Gold Standard**
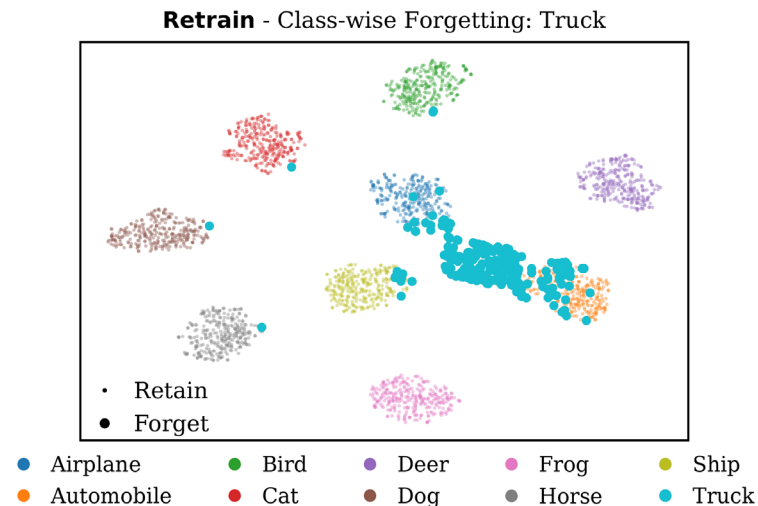  - Achieved by retraining from scratch on retain data — accurate but **computationally costly**.

- **Approximate Unlearning**
  - Offers an efficient alternative that approximates the retrained model's performance.

# Motivation

- **Retrain Model Insight**
  - Forget samples are mapped into clusters of retain samples that share the highest semantic similarity.
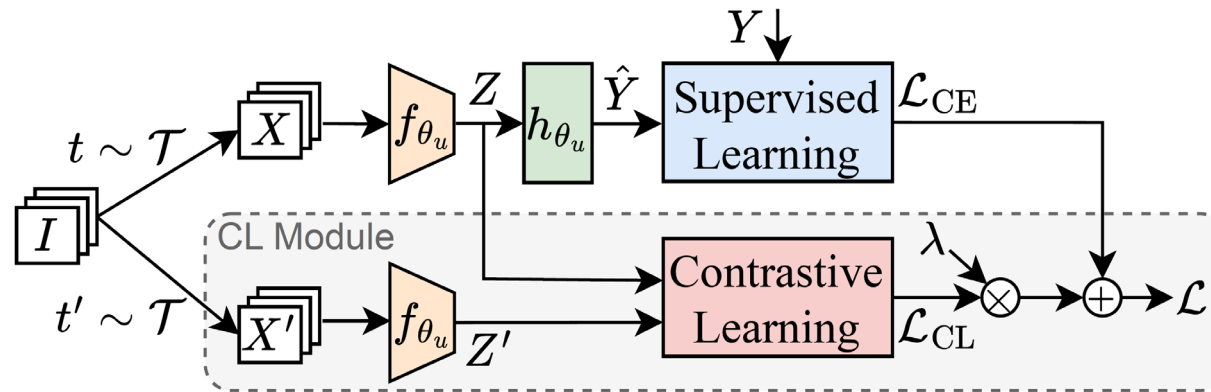


**Retrain** - Class-wise Forgetting: Truck

Retain ·
Forget ●

● Airplane    ● Bird    ● Deer    ● Frog    ● Ship
● Automobile  ● Cat     ● Dog     ● Horse   ● Truck

**Figure 1.** Representation space of the Retrain model (ResNet-18 on CIFAR-10).
*Left:* excluding the "truck" class. *Right:* excluding 10% random samples. Larger
dots denote forget samples projected into semantically closest retain clusters.

- **For approximate unlearning to emulate exact unlearning, the model should:**
  - Correctly classify retain samples.
  - Reposition forget samples into clusters of other retain samples with highest semantic similarity.

# Methodology

- CoUn operates solely on retain data and integrates two components:
  - **Contrastive Learning (CL):** Refines the representation space based on semantic similarity.
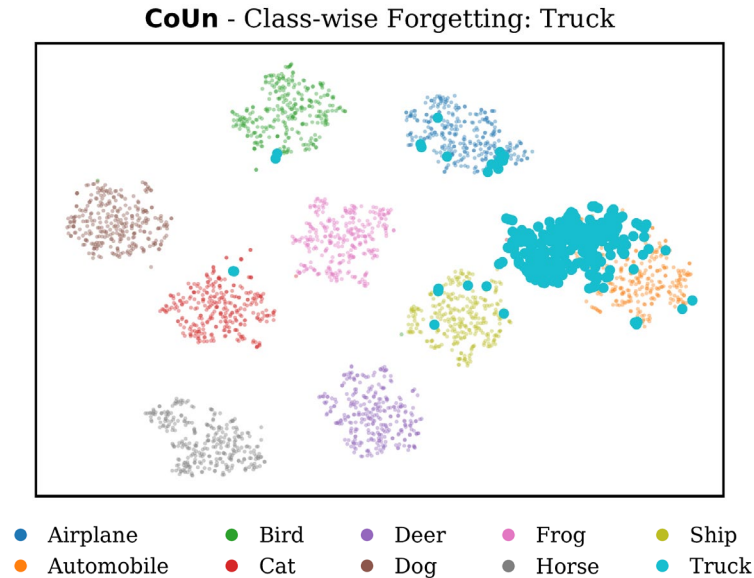  - **Supervised Learning (SL):** Preserves cluster separation.



**Figure 2.** CoUn framework. Two augmented views are generated from a batch of retain images $I$. The feature extractor $f_{\Theta_u}$ produces representations $(Z, Z')$. The CL module aligns positive pairs and separates negatives, while the supervised head $h_{\Theta_u}$ preserves the decision boundaries.

- CL module *indirectly pushes* the forget representations toward semantically similar retain samples.
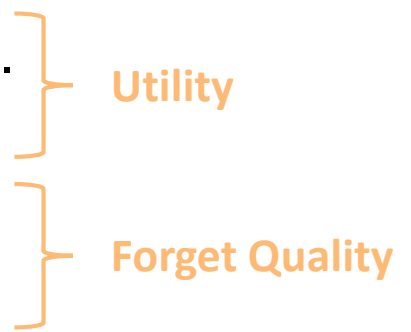
# Empirical Analysis

- CoUn preserves well-separated retain clusters (high utility) while repositioning forget samples toward semantically similar retain clusters (improved forget quality), effectively approximating exact unlearning.



**Figure 3.** Representation space of CoUn. Similar to Retrain's representation space (Figure 1), CoUn positions forget samples into clusters of semantically similar retain samples while preserving retain cluster structure.

# Evaluation Metrics

- **MU is assessed using the following key metrics:**
    1. **Retain Accuracy (RA):** Accuracy of the unlearned model $\Theta_u$ on retain data $D_r$.
    2. **Test Accuracy (TA):** Generalization performance of $\Theta_u$ on unseen test data.
    3. **Unlearn Accuracy (UA):** 1 − the accuracy of $\Theta_u$ on forget data $D_u$.
    4. **Membership Inference Attack (MIA):** Identifies whether forget samples were seen during training.
    5. **Computation Cost:** Number of FLOPs required to obtain $\Theta_u$.

Utility

Forget Quality

- **Average Gap:** Measures the unlearning effectiveness of $\Theta_u$

$$Avg.\ Gap = \frac{1}{4}\big(|RA - RA^*| + |UA - UA^*| + |TA - TA^*| + |MIA - MIA^*|\big),$$

where $(\cdot)^*$ denotes metrics of the Retrain model.
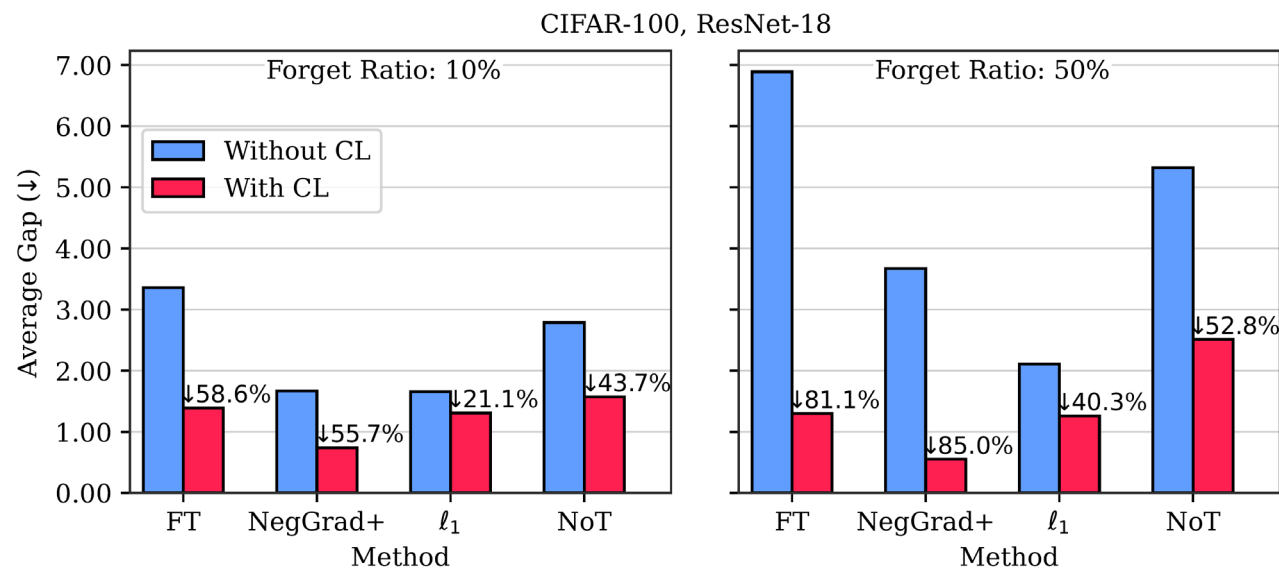
# Results (1/2)

- **Performance Comparison of CoUn against baselines.**

**Table 1.** Performance comparison of CoUn to the baseline methods with 10% random data removal. The gap ($\Delta$) and the (best) average gap between each method and the Retrain model are reported.

| Dataset & Model | Method | Accuracy (%) Retain ($\Delta \downarrow$) | Accuracy (%) Unlearn ($\Delta \downarrow$) | Accuracy (%) Test ($\Delta \downarrow$) | Efficacy (%) MIA ($\Delta \downarrow$) | Avg. Gap $\downarrow$ | Comp. Cost (PFLOPs) $\downarrow$ |
|---|---|---|---|---|---|---|---|
| | Retrain | $100.00_{\pm 0.00}$ (0.00) | $4.81_{\pm 0.27}$ (0.00) | $94.67_{\pm 0.24}$ (0.00) | $11.02_{\pm 0.58}$ (0.00) | 0.00 | 27.37 |
| | FT | $99.99_{\pm 0.00}$ (0.01) | $3.76_{\pm 0.31}$ (1.05) | $94.70_{\pm 0.14}$ (0.03) | $9.51_{\pm 0.28}$ (1.51) | 0.65 | 6.32 |
| CIFAR-10 ResNet-18 | NegGrad+ | $99.95_{\pm 0.02}$ (0.05) | $4.82_{\pm 0.24}$ (0.01) | $94.32_{\pm 0.23}$ (0.35) | $9.09_{\pm 0.30}$ (1.93) | 0.58 | 6.02 |
| | $\ell_1$-sparse | $99.97_{\pm 0.01}$ (0.03) | $5.40_{\pm 0.40}$ (0.59) | $93.81_{\pm 0.21}$ (0.86) | $10.97_{\pm 0.35}$ (0.05) | 0.38 | 6.92 |
| | SalUn | $99.10_{\pm 0.35}$ (0.90) | $4.31_{\pm 0.42}$ (0.50) | $93.84_{\pm 0.27}$ (0.83) | $11.15_{\pm 2.04}$ (0.13) | 0.59 | 8.66 |
| | NoT | $99.99_{\pm 0.00}$ (0.01) | $4.19_{\pm 0.25}$ (0.62) | $94.65_{\pm 0.24}$ (0.02) | $10.45_{\pm 0.51}$ (0.57) | 0.30 | 7.52 |
| | CoUn | $99.99_{\pm 0.00}$ (0.01) | $4.12_{\pm 0.31}$ (0.69) | $94.57_{\pm 0.24}$ (0.10) | $10.81_{\pm 0.31}$ (0.21) | 0.25 | 8.02 |

# Results (2/2)

- **Percentage improvement from integrating CoUn's CL module into baselines.**



**Figure 4.** CL integration consistently improves unlearning performance, with larger gains at higher forget ratios (50%).

# Conclusion

- **CoUn enables effective unlearning by leveraging semantic similarity between retain and forget samples.**
  - Uses a CL module on retain data to adjust their representations, and indirectly influences forget representations.
    - Improving forget quality.
  - Uses supervised learning on retain data to preserve their cluster separation.
    - Improving utility.

arXiv

NEURAL INFORMATION
PROCESSING SYSTEMS
NeurIPS 2025
San Diego, USA

HUAWEI

Thank you!

E-mail: yasser1.khalil@huawei.ca