

# Selective Amnesia using Contrastive Subnet Erasure for Class Level Unlearning in Vision Models

Vishal Pramanik

Department of Computer and Information Science and Engineering  
University of Florida, Gainesville, FL

vishalpramanik@ufl.edu

Maisha Maliha

School of Computer Science  
University of Oklahoma, Norman, OK

maisha.maliha-1@ou.edu

Susmit Jha

Computer Science Laboratory  
SRI International, Menlo Park, CA

susmitjha@berkeley.edu

Alvaro Velasquez

Department of Computer Science  
University of Colorado Boulder, Boulder, CO

Alvaro.Velasquez@colorado.edu

Olivera Kotevska

Computer Science and Mathematics Division  
Oak Ridge National Laboratory, Oak Ridge, TN

kotevskao@ornl.gov

Sumit Kumar Jha

Department of Computer and Information Science and Engineering  
University of Florida, Gainesville, FL

sumit.jha@ufl.edu

## Abstract

We study concept-level forgetting in pretrained vision models: removing an entire semantic category so the system no longer recognizes that object in unseen images and contexts, rather than merely forgetting specific training examples. Prior work either applies blunt global projections or fine-tunes parameters, which can introduce collateral damage to unrelated features, add compute, and become unstable as forgetting strength increases. We introduce **Contrastive Subnet Erasure (CSE)**, a training-free, encoder-centric edit that targets a compact set of channels most responsible for the class and attenuates them in a calibrated manner. The modification is algebraically folded into the subsequent layer, yielding no inference-time overhead and leaving task heads unchanged. To evaluate whether forgetting generalizes beyond the data used to specify the class, we introduce a cross dataset protocol in which the class is defined on a source dataset and performance is measured on a disjoint target dataset drawn from a different distribution with no shared images. This setup tests whether the model still fails to recognize the object when it looks different or appears in new scenes, and it avoids overfitting

to patterns of the source dataset. Across CIFAR 10, CIFAR 100, and ImageNet under this protocol, CSE achieves stronger forgetting of the target class while better preserving non target utility than existing baselines in both single class and multi class settings. Overall, CSE provides a simple stable and deployment ready mechanism for class level unlearning in vision.

## 1. Introduction

Deep neural networks pretrained on massive, weakly curated corpora inevitably internalize information that later needs to be *removed*—to comply with legal deletion requests, to excise unsafe or biased knowledge, or to neutralize backdoors and other poisons [7, 23, 26]—without retraining from scratch. Prior work frames this as either *data unlearning* (removing the influence of specific examples) [4, 5] or *concept erasure* (suppressing a semantic factor in representations) [21, 22]. Despite progress, recent stress tests show that many approximate unlearning procedures can leave measurable residues—especially under data poisoning—highlighting the need for stronger, more surgical edits that minimize collateral damage to unrelated capa-

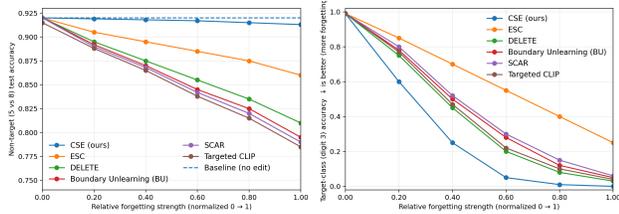


Figure 1. MNIST-EfficientNet toy with a shared, normalized forgetting strength  $s \in [0, 1]$  ( $0 = \text{no edit}$ ). **(a) Non-target utility:** test accuracy on 5 vs. 8; dashed line marks the no-edit baseline. **CSE (ours)** stays near baseline across  $s$ , indicating minimal collateral damage; **ESC** is typically second-best for small  $s$  but degrades as removal grows; training-based edits (**DELETE**, **BU**, **SCAR**, **Targeted CLIP**) drift more. **(b) Target forgetting:** accuracy on digit “3” (*lower is better*). **CSE** drives the target accuracy to near-zero by  $s \approx 1$  while preserving non-target utility, whereas others require stronger edits or fail to fully forget within the same range.

bility (e.g., privacy- and safety-motivated unlearning audits, linear concept erasure in closed form, and selective forgetting in generative models) [2, 14, 19]. The central difficulty is *feature entanglement*: target and non-target concepts often cohabit directions and channels within the same encoder, so naïve removal of “target” signal tends to deform nearby structure, degrading utility on tasks that share visual cues with the forgotten content.

Against this backdrop, we benchmark five families of encoder edits that typify today’s practice: **ESC** [17], **DELETE** [29], **Boundary Unlearning (BU)** [6], **SCAR** [3], and **Targeted-CLIP** [27]. In brief, ESC performs *global subspace deletion* by projecting out directions correlated with the target; this is effective at low strength but becomes brittle as more components are removed, amplifying *collateral loss* on shared features. DELETE and BU are *gradient-based, training* methods that reshape parameters to repel the target region or shrink decision boundaries; they can forget strongly but tend to *drag shared filters*, hurting non-target performance as edit strength grows unless carefully regularized. SCAR and Targeted-CLIP *retain-vs-forget objectives* with additional supervision (e.g., prompts or balanced retain sets), improving controllability but introducing *compute overhead*, sensitivity to hyperparameters, and stability issues at higher strengths. Across these approaches, two recurring limitations emerge: (i) *non-locality*—edits operate on broad subspaces or full parameter blocks, so changes spill into unrelated regions; (ii) *lack of geometry preservation*—even when target accuracy drops, the *relative arrangement* of non-target classes can be distorted, eroding downstream linear-probe separability.

To overcome these drawbacks, we introduce **CSE (Contrastive Subnet Erasure)**, a training-free, encoder-centric edit that operates *locally* in channel space. CSE first performs *contrastive subnet discovery*, scoring channels by

their *target-vs-background salience* and selecting a minimal *compact subnet* that covers the discriminative mass; it then applies *calibrated attenuation* only to those channels, leaving the remaining representation geometry intact. This constructive selection–then–erasure pipeline preserves non-target structure by design: shared directions that are *not* uniquely target-salient are retained, while truly target-diagnostic channels are damped. More broadly, CSE’s *architecture-agnostic, precompute-once* attenuation delivers practical unlearning with negligible inference overhead while directly addressing the two core shortcomings above—*locality* and *geometry preservation*—that limit prior art. We summarize our main contributions below:

- We present **CSE (Contrastive Subnet Erasure)**, a training-free, encoder-centric method that *surgically* attenuates a compact set of target-salient channels, removing the *target* dataset/class while leaving *non-target* representations essentially intact.
- We introduce a *cross-data evaluation* scheme in which forgetting is performed on a source dataset/class and performance is measured on a *different* dataset that shares the same semantic category to be forgotten. This protocol stress-tests true concept removal (transfer leakage) while simultaneously auditing non-target utility under distribution shift.
- Across multiple datasets and models, CSE consistently outperforms recent strong baselines in both forgetting efficacy and retention of non-target performance, achieving robust unlearning with negligible inference overhead.

## 2. Motivation

We study a minimal but revealing [9] toy experiment using a frozen EfficientNet-B0 encoder to motivate CSE and illustrate the central challenge in machine unlearning when target and non-target concepts share visual features, removing the target risks damaging unrelated representations. We designate digit “3” as the target to forget and measure non-target utility on “5 vs. 8” classification using a linear probe—a neutral diagnostic trained on the frozen encoder features. This setup is deliberately challenging as digit “3” shares strokes with both “5” and “8”, so any edit removing “3”-correlated directions risks erasing cues needed to distinguish them. We compare five recent encoder-editing strategies against our method (CSE) under a unified normalized forgetting strength  $s \in [0, 1]$  ( $0 = \text{no edit}$ ; larger = stronger) which are **ESC**[17], **DELETE** [29], **Boundary Unlearning (BU)** [6], **SCAR** [3], **Targeted CLIP** [27]. Figure 1(a) shows non-target accuracy (5 vs. 8) versus  $s$  with a dashed no-edit baseline. **CSE tracks the baseline across the entire range**, indicating minimal collateral damage, while **ESC** is second-best at low  $s$  but degrades rapidly as more components are removed (global subspace deletion becomes destructive). Figure 1(b) shows target-class ac-

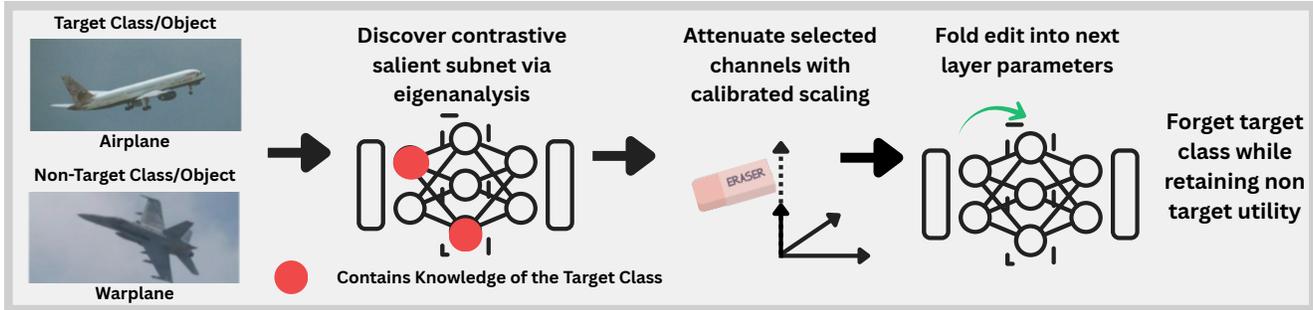


Figure 2. **Overview of CSE.** From target/non-target sets, CSE (i) standardizes features, (ii) discovers a compact contrastive subnet via generalized eigenanalysis, and (iii) attenuates only those channels—suppressing target signal while preserving non-target geometry.

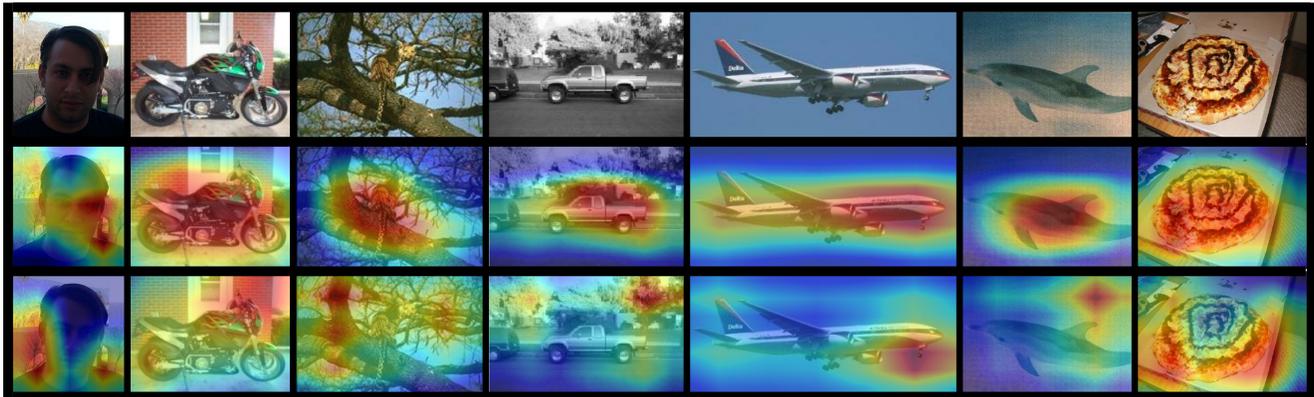


Figure 3. Qualitative effect of CSE: Each row shows the **Original**, **Grad-CAM (Before)**, and **Grad-CAM (After)**. CSE suppresses saliency on the target concept while preserving responses on non-target content.

curacy (digit 3) versus  $s$  (lower is better): **CSE achieves near-zero accuracy by  $s \approx 1$**  while maintaining high non-target performance. Training-based methods (DELETE, BU, SCAR, Targeted-CLIP) also forget effectively but degrade non-target utility more as strength grows—they reshape shared filters unless carefully rebalanced. Together, these results demonstrate that when concepts share feature directions, **CSE’s localized, contrastive attenuation suppresses the target while preserving non-target geometry**, outperforming both global projection and gradient-based fine-tuning.

### 3. Method: Contrastive Subnet Erasure

To overcome the drawbacks in the previous section, we propose **Contrastive Subnet Erasure (CSE)**, a training-free method that enables a pretrained vision model to selectively forget a target visual concept without retraining. CSE operates by identifying and attenuating a *compact subnet*—a small, targeted set of encoder channels that exhibit high *contrastive salience* (significantly greater activation variance on target images than on non-target images). The method proceeds in three stages: **(1) Standardization** stan-

dardizes features for stable covariance estimation; **(2) Contrastive subnet discovery** applies generalized eigenanalysis to identify discriminative directions and scores channels by their eigenvalue-weighted participation; **(3) Subnet attenuation** attenuates the discovered channels through efficient runtime scaling operations.

#### 3.1. Problem Formulation

Let  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$  denote a pretrained encoder mapping images to embeddings. We denote intermediate features at layer  $\ell$  as  $h^{(\ell)}(x) \in \mathbb{R}^{d_\ell}$ , where  $\ell \in \mathcal{L} \subseteq \{1, \dots, L\}$  indexes selected layers. For spatial features (convolutional or transformer patch tokens), we apply global average pooling to obtain channel-wise representations. Given a target dataset  $\mathcal{D}_t = \{x_i\}_{i=1}^{n_t}$  containing the concept to forget and a non-target dataset  $\mathcal{D}_b = \{x_j\}_{j=1}^{n_b}$  containing concepts to preserve, we seek to identify and attenuate channels that selectively encode the target concept while minimizing disruption to non-target representations.

#### 3.2. Stage 1: Feature Extraction & Standardization

We extract features from both target and non-target datasets and standardize them to ensure comparable scales for con-

trastive analysis. For each layer  $\ell$ , we compute joint statistics across both datasets:

$$\mu^{(\ell)} = \frac{1}{n_t + n_b} \sum_{x \in \mathcal{D}_t \cup \mathcal{D}_b} h^{(\ell)}(x), \quad (1)$$

$$\sigma_c^{(\ell)} = \sqrt{\frac{1}{n_t + n_b} \sum_{x \in \mathcal{D}_t \cup \mathcal{D}_b} (h_c^{(\ell)}(x) - \mu_c^{(\ell)})^2 + \epsilon}, \quad (2)$$

where  $c$  indexes channels and  $\epsilon$  prevents division by zero. The standardized features become:

$$\hat{h}^{(\ell)}(x) = S^{(\ell)}(h^{(\ell)}(x) - \mu^{(\ell)}), \quad (3)$$

with  $S^{(\ell)} = \text{diag}(1/\sigma_1^{(\ell)}, \dots, 1/\sigma_{d_\ell}^{(\ell)})$ . Joint standardization is crucial for unbiased variance ratio computation in the subsequent analysis.

### 3.3. Stage 2: Contrastive Subnet Discovery

In line with prior work [1, 20], we identify a compact *subnet* of encoder channels that exhibits high *contrastive salience*. Operating on standardized features from Stage 1, we identify feature directions where target concepts exhibit disproportionately high variance relative to non-target concepts. Computing covariance matrices on standardized features:

$$\begin{aligned} \Sigma_t^{(\ell)} &= \frac{1}{n_t} \sum_{x \in \mathcal{D}_t} \hat{h}^{(\ell)}(x) \hat{h}^{(\ell)}(x)^\top, \\ \Sigma_b^{(\ell)} &= \frac{1}{n_b} \sum_{x \in \mathcal{D}_b} \hat{h}^{(\ell)}(x) \hat{h}^{(\ell)}(x)^\top \end{aligned} \quad (4)$$

we seek directions  $v \in \mathbb{R}^{d_\ell}$  maximizing the variance ratio:

$$\rho(v) = \frac{v^\top \Sigma_t^{(\ell)} v}{v^\top \Sigma_b^{(\ell)} v}. \quad (5)$$

Large  $\rho(v)$  indicates direction  $v$  is highly *target-salient*—the target varies strongly along  $v$  while non-targets remain stable. To handle limited samples and ensure numerical stability, we solve the regularized generalized eigenvalue problem:

$$\Sigma_t^{(\ell)} v = \rho (\Sigma_b^{(\ell)} + \delta I_{d_\ell}) v, \quad (6)$$

where  $\delta = \alpha \cdot \text{trace}(\Sigma_b^{(\ell)})/d_\ell$  with  $\alpha$  being a small regularization factor. This yields eigenpairs  $\{(v_j^{(\ell)}, \rho_j^{(\ell)})\}_{j=1}^{d_\ell}$  ordered by decreasing eigenvalue.

Each channel’s importance is measured by its participation in discriminative directions:

$$s_c^{(\ell)} = \sum_{j=1}^{k_\ell} \rho_j^{(\ell)} (v_j^{(\ell)}[c])^2, \quad (7)$$

where  $k_\ell = \min(k_{\max}, \lfloor \beta \cdot d_\ell \rfloor)$  bounds the number of eigenvectors considered, with hyperparameters  $k_{\max}$  and  $\beta$ . The squared coefficient  $(v_j^{(\ell)}[c])^2$  measures channel  $c$ ’s contribution to eigenvector  $j$ , weighted by its discriminative power  $\rho_j^{(\ell)}$ .

We select the minimal channel subset  $\mathcal{C}^{(\ell)}$  capturing a fraction  $\tau_{\text{cov}}$  of total discriminative information:

$$\sum_{c \in \mathcal{C}^{(\ell)}} s_c^{(\ell)} \geq \tau_{\text{cov}} \sum_{c=1}^{d_\ell} s_c^{(\ell)}. \quad (8)$$

Channels are greedily selected in descending score order until the coverage threshold is met. The complete subnet comprises  $\text{Subnet} = \{(\ell, c) : \ell \in \mathcal{L}, c \in \mathcal{C}^{(\ell)}\}$ .

### 3.4. Stage 3: Subnet Attenuation

**Attenuation strength computation.** For each selected channel, we compute an attenuation factor based on its discriminative score:

$$\beta_c^{(\ell)} = \text{clip}_{[0,1]} \left( \frac{s_c^{(\ell)} - \tau_0}{s_c^{(\ell)} + \lambda_0} \right), \quad (9)$$

where  $\tau_0$  sets the minimum score threshold and  $\lambda_0$  controls the transition smoothness. The attenuation matrix becomes:

$$A^{(\ell)} = \text{diag}(1 - \beta_1^{(\ell)}, \dots, 1 - \beta_{d_\ell}^{(\ell)}), \quad (10)$$

with diagonal entries ranging from 0 (complete removal) to 1 (full preservation).

**Transformation to original space.** Since attenuation is computed in standardized coordinates, we transform it back to the original feature space:

$$M^{(\ell)} = S^{(\ell)-1} A^{(\ell)} S^{(\ell)}. \quad (11)$$

The runtime attenuation operation becomes:

$$h_{\text{atten}}^{(\ell)} = M^{(\ell)} h^{(\ell)} + (I_{d_\ell} - M^{(\ell)}) \mu^{(\ell)}, \quad (12)$$

where the bias term compensates for the mean shift introduced by standardization.

**Architecture-agnostic application.** For both convolutional and transformer architectures, attenuation is applied at the block output level to respect architectural constraints. Given a generic block with residual connection:

$$h^{(\ell+1)} = \mathcal{F}(h^{(\ell)}) + \mathcal{S}(h^{(\ell)}), \quad (13)$$

where  $\mathcal{F}$  represents the main transformation and  $\mathcal{S}$  the skip connection, we apply:

$$h_{\text{atten}}^{(\ell+1)} = \text{diag}(M^{(\ell+1)}) \odot h^{(\ell+1)} + \beta^{(\ell+1)}, \quad (14)$$

with precomputed scaling  $\text{diag}(M^{(\ell+1)})$  and bias  $\beta^{(\ell+1)} = (I - M^{(\ell+1)})\mu^{(\ell+1)}$ .

This placement ensures attenuation affects the complete block output, preventing bypass through residual paths while maintaining architectural integrity. The approach uniformly handles diverse architectures—convolutional blocks with various depths and widths, as well as transformer blocks with self-attention and feedforward components—through consistent application at the aggregated feature level. An overview of our method has been shown in figure 2 and the algorithm has been given in the appendix.

## 4. Experimental Setup

To test whether unlearning removes a *concept* rather than specific datapoints, we evaluate methods in a cross-dataset setting spanning **CIFAR-10**, **CIFAR-100**[16], and **ImageNet**[8]. These datasets do not share images but exhibit clear semantic overlaps, allowing us to unlearn a class using data from one dataset and evaluate forgetting and retention on others. Our setup mirrors that of DELETE [29], except we explicitly source the forget set from a *different* dataset than the one used for evaluation, enabling a more rigorous test of semantic removal. We fix random seeds and forget/retain partitions across all methods and report standard metrics: *forget-train* accuracy ( $Acc_f$ ) and *forget-test* accuracy ( $Acc_{ft}$ ), both of which should drop under effective forgetting; *retain-train* ( $Acc_r$ ) and *retain-test* ( $Acc_{rt}$ ) accuracies, which should remain close to the Original and Retrain baselines; the *harmonic mean* (H-Mean)(same as in [29] between test-time forgetting and retention (higher is better); and the *membership inference attack* (MIA[24]) success rate on the forget set (lower is better, indicating removal of memorized patterns).

In the *single-class* setting, we design three cross-dataset probes so that each dataset appears as both the unlearning and evaluation domain: (i) we forget *airplane* on CIFAR-10 and evaluate forgetting on ImageNet’s *airliner/aircraft*, while verifying that nearby categories such as *warplane* are retained; (ii) we reverse the direction by forgetting ImageNet’s *garbage truck, tow truck, and trailer truck* and evaluating forgetting on CIFAR-10 *truck*, ensuring that the semantically adjacent *automobile* remains intact; and (iii) we forget ImageNet’s *white shark and tiger shark* and evaluate forgetting on CIFAR-100 *shark*, while confirming that related fish classes (*aquarium fish, flatfish, ray, trout*) are preserved. This design exercises all three datasets in both roles and stresses robustness to *semantic overlap* rather than dataset membership. For the *multi-class* condition, we extend to all three datasets using up to five semantically aligned classes across domains. From **CIFAR-10**, we align *airplane*  $\leftrightarrow$  *airliner*, *truck*  $\leftrightarrow$  *garbage truck*, *ship*  $\leftrightarrow$  *container ship*, *cat*  $\leftrightarrow$  *tabby cat*, and *frog*  $\leftrightarrow$  *bullfrog*. From **CIFAR-100**, we align *castle, computer keyboard, telephone*

(*cellular/dial*), *television*, and *lawn mower* to their same-named **ImageNet** categories. When experiments use **ImageNet**, we substitute the nearest included class if an exact match is absent (e.g., *airliner*  $\rightarrow$  *warplane*). We conduct bidirectional evaluations for each dataset pair (CIFAR-10  $\leftrightarrow$  ImageNet, CIFAR-100  $\leftrightarrow$  ImageNet, and CIFAR-10  $\leftrightarrow$  CIFAR-100) using fixed label mappings (Appendix). Robust, minimally destructive unlearning is indicated by low  $Acc_{ft}$  and MIA, and high  $Acc_{rt}$  and H-Mean across all directions—demonstrating that targeted semantics are removed regardless of dataset origin.

To assess architectural robustness of our method, we evaluate **EfficientNet-B0**[25], **ResNet-18**[13], and **SwiT**[18]. When a method prescribes head-only edits, we repect that constraint; otherwise, encoder parameters are updated only when the method explicitly requires fine-tuning. We compare the following approaches: (1) **Original**, the pretrained model without unlearning (upper-utility reference); (2) **Retrain**, full training from scratch on the retained data only (gold-standard lower bound on residual knowledge of the forgotten class); (3) **ESC**; (4) **ESC-T** (training); (5) **DELETE**; (6) **BU** (Boundary Shrink); (7) **LEACE**; (8) **SCAR**; (9) **SCRUB**; and (10) **SCRUB+R**, SCRUB followed by a light post-repair/refit stage.

We use 10% of non-target category samples from the evaluation dataset and the complete target class from the source dataset across all experiments, ensuring true cross-dataset semantic removal. For CSE, we set the regularization factor  $\alpha = 0.01$  in the generalized eigenvalue problem to stabilize covariance matrix inversion, bound eigenvectors by  $k_\ell = \min(k_{\max}, \lfloor \beta \cdot d_\ell \rfloor)$  with  $k_{\max} = 50$  and  $\beta = 0.3$  to focus on discriminative directions while avoiding noise overfitting, and use coverage threshold  $\tau_{\text{cov}} = 0.85$  to capture 85% of discriminative information. Attenuation strength is calibrated via  $\tau_0 = 0.1$  and  $\lambda_0 = 0.5$  controlling score threshold and transition smoothness, with standardization epsilon  $\epsilon = 10^{-6}$  for numerical stability. Training-based baselines (DELETE, BU, SCAR) are fine-tuned for 10 epochs using SGD with learning rate  $10^{-5}$ , momentum 0.9, and batch size 64, following their original protocols. We partition datasets using standard 80–20 train–test splits: CIFAR-10 (40K/10K, 10 classes), CIFAR-100 (40K/10K, 100 classes), and ImageNet-1K (1.28M/50K, 1000 classes). For single-class forgetting, one target class is designated for removal while remaining classes form the retain set; for multi-class forgetting on CIFAR-100, we incrementally forget  $\{2, 3, 4, 5\}$  classes and report performance averaged over three random seeds. All methods are evaluated on identical splits with fixed seeds for fair comparison.

Across the single-class cross-dataset setting (Table 1), **CSE** is uniformly best: it drives  $Acc_{ft}$  down to **0.01–0.02** on C10/C100/ImageNet while all other methods remain in double digits ( $\geq 0.10$ ); at the same time it preserves

Table 1. **Single-class cross-dataset unlearning.** Results on **CIFAR-10 (C10)**, **CIFAR-100 (C100)**, and **ImageNet** under the single-class protocol, evaluated on three backbones (**ResNet-18**, **EfficientNet-B0**, **Swin-T**). We report *forget-test* accuracy ( $Acc_{ft} \downarrow$ ), *retain-test* accuracy ( $Acc_{rt} \uparrow$ ), *H-Mean* ( $\uparrow$ ), and *MIA* ( $\downarrow$ ). **Formatting:** best values (**CSE**, ours) are **bold blue**. Method names (baselines) and all numbers use a reduced font; padding and rule spacing are tightened for compactness. 3 independent runs,  $std < \pm 0.02$

| Method                           | CIFAR-10 (C10)        |                     |                   |                  | CIFAR-100 (C100)      |                     |                   |                  | ImageNet              |                     |                   |                  |
|----------------------------------|-----------------------|---------------------|-------------------|------------------|-----------------------|---------------------|-------------------|------------------|-----------------------|---------------------|-------------------|------------------|
|                                  | $Acc_{ft} \downarrow$ | $Acc_{rt} \uparrow$ | H-Mean $\uparrow$ | MIA $\downarrow$ | $Acc_{ft} \downarrow$ | $Acc_{rt} \uparrow$ | H-Mean $\uparrow$ | MIA $\downarrow$ | $Acc_{ft} \downarrow$ | $Acc_{rt} \uparrow$ | H-Mean $\uparrow$ | MIA $\downarrow$ |
| <b>Backbone: ResNet-18</b>       |                       |                     |                   |                  |                       |                     |                   |                  |                       |                     |                   |                  |
| Original                         | 0.94                  | 0.93                | 0.50              | 0.22             | 0.86                  | 0.76                | 0.41              | 0.20             | 0.70                  | 0.59                | 0.36              | 0.22             |
| Retrain                          | 0.03                  | 0.91                | 0.89              | 0.02             | 0.03                  | 0.76                | 0.78              | 0.02             | 0.03                  | 0.59                | 0.68              | 0.02             |
| ESC                              | 0.10                  | 0.92                | 0.90              | 0.05             | 0.12                  | 0.76                | 0.76              | 0.06             | 0.12                  | 0.59                | 0.67              | 0.05             |
| ESC-T                            | 0.10                  | 0.91                | 0.88              | 0.06             | 0.10                  | 0.75                | 0.75              | 0.06             | 0.11                  | 0.58                | 0.64              | 0.06             |
| DELETE                           | 0.12                  | 0.91                | 0.89              | 0.06             | 0.10                  | 0.77                | 0.79              | 0.05             | 0.14                  | 0.58                | 0.65              | 0.05             |
| BU                               | 0.12                  | 0.89                | 0.82              | 0.09             | 0.15                  | 0.73                | 0.70              | 0.09             | 0.18                  | 0.57                | 0.60              | 0.09             |
| LEACE                            | 0.15                  | 0.90                | 0.80              | 0.11             | 0.18                  | 0.74                | 0.68              | 0.11             | 0.22                  | 0.57                | 0.58              | 0.11             |
| SCAR                             | 0.10                  | 0.90                | 0.85              | 0.08             | 0.14                  | 0.74                | 0.71              | 0.08             | 0.17                  | 0.57                | 0.61              | 0.08             |
| SCRUB                            | 0.13                  | 0.88                | 0.81              | 0.10             | 0.16                  | 0.72                | 0.69              | 0.10             | 0.20                  | 0.56                | 0.59              | 0.10             |
| SCRUB+R                          | 0.11                  | 0.90                | 0.84              | 0.08             | 0.13                  | 0.74                | 0.71              | 0.08             | 0.16                  | 0.57                | 0.62              | 0.08             |
| <b>CSE (ours)</b>                | <b>0.01</b>           | <b>0.95</b>         | <b>0.96</b>       | <b>0.01</b>      | <b>0.02</b>           | <b>0.79</b>         | <b>0.84</b>       | <b>0.01</b>      | <b>0.02</b>           | <b>0.61</b>         | <b>0.73</b>       | <b>0.01</b>      |
| <b>Backbone: EfficientNet-B0</b> |                       |                     |                   |                  |                       |                     |                   |                  |                       |                     |                   |                  |
| Original                         | 0.95                  | 0.94                | 0.52              | 0.23             | 0.87                  | 0.77                | 0.42              | 0.21             | 0.71                  | 0.59                | 0.37              | 0.23             |
| Retrain                          | 0.03                  | 0.92                | 0.90              | 0.02             | 0.03                  | 0.77                | 0.79              | 0.02             | 0.03                  | 0.60                | 0.69              | 0.02             |
| ESC                              | 0.10                  | 0.93                | 0.90              | 0.05             | 0.10                  | 0.78                | 0.80              | 0.05             | 0.12                  | 0.59                | 0.66              | 0.05             |
| ESC-T                            | 0.10                  | 0.92                | 0.88              | 0.06             | 0.10                  | 0.76                | 0.75              | 0.06             | 0.11                  | 0.58                | 0.64              | 0.06             |
| DELETE                           | 0.12                  | 0.92                | 0.91              | 0.05             | 0.12                  | 0.77                | 0.79              | 0.05             | 0.12                  | 0.60                | 0.68              | 0.05             |
| BU                               | 0.12                  | 0.90                | 0.83              | 0.09             | 0.15                  | 0.74                | 0.71              | 0.09             | 0.18                  | 0.58                | 0.61              | 0.09             |
| LEACE                            | 0.16                  | 0.91                | 0.81              | 0.11             | 0.18                  | 0.75                | 0.69              | 0.11             | 0.22                  | 0.58                | 0.59              | 0.11             |
| SCAR                             | 0.11                  | 0.91                | 0.86              | 0.08             | 0.14                  | 0.75                | 0.72              | 0.08             | 0.17                  | 0.58                | 0.62              | 0.08             |
| SCRUB                            | 0.13                  | 0.89                | 0.82              | 0.10             | 0.16                  | 0.73                | 0.70              | 0.10             | 0.20                  | 0.57                | 0.59              | 0.10             |
| SCRUB+R                          | 0.12                  | 0.91                | 0.85              | 0.08             | 0.13                  | 0.75                | 0.72              | 0.08             | 0.16                  | 0.58                | 0.63              | 0.08             |
| <b>CSE (ours)</b>                | <b>0.01</b>           | <b>0.96</b>         | <b>0.97</b>       | <b>0.01</b>      | <b>0.02</b>           | <b>0.80</b>         | <b>0.85</b>       | <b>0.01</b>      | <b>0.02</b>           | <b>0.62</b>         | <b>0.74</b>       | <b>0.01</b>      |
| <b>Backbone: Swin-T</b>          |                       |                     |                   |                  |                       |                     |                   |                  |                       |                     |                   |                  |
| Original                         | 0.96                  | 0.95                | 0.53              | 0.24             | 0.88                  | 0.78                | 0.43              | 0.22             | 0.72                  | 0.60                | 0.38              | 0.24             |
| Retrain                          | 0.03                  | 0.93                | 0.91              | 0.02             | 0.03                  | 0.78                | 0.80              | 0.02             | 0.03                  | 0.61                | 0.70              | 0.02             |
| ESC                              | 0.10                  | 0.94                | 0.92              | 0.04             | 0.12                  | 0.78                | 0.80              | 0.05             | 0.12                  | 0.60                | 0.67              | 0.05             |
| ESC-T                            | 0.10                  | 0.93                | 0.89              | 0.06             | 0.10                  | 0.77                | 0.76              | 0.06             | 0.11                  | 0.59                | 0.65              | 0.06             |
| DELETE                           | 0.12                  | 0.93                | 0.91              | 0.05             | 0.12                  | 0.79                | 0.81              | 0.05             | 0.12                  | 0.61                | 0.69              | 0.05             |
| BU                               | 0.12                  | 0.91                | 0.84              | 0.09             | 0.15                  | 0.75                | 0.72              | 0.09             | 0.18                  | 0.59                | 0.62              | 0.09             |
| LEACE                            | 0.16                  | 0.92                | 0.82              | 0.11             | 0.19                  | 0.76                | 0.70              | 0.11             | 0.22                  | 0.59                | 0.60              | 0.11             |
| SCAR                             | 0.11                  | 0.92                | 0.87              | 0.08             | 0.14                  | 0.76                | 0.73              | 0.08             | 0.17                  | 0.59                | 0.63              | 0.08             |
| SCRUB                            | 0.13                  | 0.90                | 0.83              | 0.10             | 0.16                  | 0.74                | 0.71              | 0.10             | 0.20                  | 0.58                | 0.60              | 0.10             |
| SCRUB+R                          | 0.12                  | 0.92                | 0.86              | 0.08             | 0.13                  | 0.76                | 0.73              | 0.08             | 0.16                  | 0.59                | 0.64              | 0.08             |
| <b>CSE (ours)</b>                | <b>0.01</b>           | <b>0.97</b>         | <b>0.98</b>       | <b>0.01</b>      | <b>0.02</b>           | <b>0.81</b>         | <b>0.86</b>       | <b>0.01</b>      | <b>0.02</b>           | <b>0.63</b>         | <b>0.75</b>       | <b>0.01</b>      |

the highest  $Acc_{rt}$  (e.g., **0.95** on C10, **0.79** on C100, **0.61** on TIN), yields the strongest H-Mean (**0.96–0.98**, **0.84–0.86**, **0.73–0.75** respectively), and attains the lowest MIA (**0.01**). The second-best baseline alternates between **ESC** and **DELETE** depending on backbone/dataset but consistently trails **CSE** in both forgetting and retention. In the multi-class CIFAR-100 evaluation averaged over backbones (Table 2), **CSE** remains robust as the number of forgotten classes increases from 2 to 5:  $Acc_{ft}$  grows only modestly from **0.02** to **0.04**,

$Acc_{rt}$  stays high (**0.80–0.82**), and H-Mean remains the best (**0.83–0.86**), whereas the strongest baselines (**ESC/DELETE**) show double-digit  $Acc_{ft}$  (typically 0.10–0.14) and declining H-Mean as the target set grows. Overall, **CSE** achieves the most aggressive forgetting with the least collateral damage, and its advantage widens as the task becomes more challenging. Figure 3 illustrates example outputs from EfficientNet-B0. As shown in Figure 5(a–c), **CSE** achieves the highest retain accuracy on semantically similar categories across all datasets and architectures, con-

Table 2. **Multi-class forgetting on CIFAR-100 (averaged over backbones).** We report *forget-test* accuracy ( $Acc_{ft} \downarrow$ ), *retain-test* accuracy ( $Acc_{rt} \uparrow$ ), and *H-Mean* ( $\uparrow$ ) when forgetting  $\{2, 3, 4, 5\}$  target classes, averaged over **ResNet-18**, **EfficientNet-B0**, and **Swin-T**. **Formatting:** our method **CSE** (best) is shown in **bold blue**. Numbers are set in a reduced font; padding and rule spacing are tightened for compactness. 3 independent runs,  $std < \pm 0.03$

| Method            | 2 Classes             |                     |                   | 3 Classes             |                     |                   | 4 Classes             |                     |                   | 5 Classes             |                     |                   |
|-------------------|-----------------------|---------------------|-------------------|-----------------------|---------------------|-------------------|-----------------------|---------------------|-------------------|-----------------------|---------------------|-------------------|
|                   | $Acc_{ft} \downarrow$ | $Acc_{rt} \uparrow$ | H-Mean $\uparrow$ | $Acc_{ft} \downarrow$ | $Acc_{rt} \uparrow$ | H-Mean $\uparrow$ | $Acc_{ft} \downarrow$ | $Acc_{rt} \uparrow$ | H-Mean $\uparrow$ | $Acc_{ft} \downarrow$ | $Acc_{rt} \uparrow$ | H-Mean $\uparrow$ |
| Original          | 0.87                  | 0.77                | –                 | 0.87                  | 0.77                | –                 | 0.88                  | 0.77                | –                 | 0.88                  | 0.77                | –                 |
| Retrain           | 0.00                  | 0.79                | 0.79              | 0.00                  | 0.79                | 0.79              | 0.00                  | 0.78                | 0.78              | 0.00                  | 0.78                | 0.78              |
| ESC               | 0.10                  | 0.78                | 0.80              | 0.12                  | 0.77                | 0.78              | 0.12                  | 0.76                | 0.77              | 0.10                  | 0.75                | 0.75              |
| ESC-T             | 0.12                  | 0.76                | 0.76              | 0.12                  | 0.75                | 0.75              | 0.13                  | 0.75                | 0.74              | 0.12                  | 0.74                | 0.73              |
| DELETE            | 0.12                  | 0.78                | 0.79              | 0.10                  | 0.78                | 0.80              | 0.14                  | 0.77                | 0.77              | 0.10                  | 0.76                | 0.76              |
| BU                | 0.14                  | 0.74                | 0.73              | 0.15                  | 0.74                | 0.72              | 0.16                  | 0.73                | 0.71              | 0.17                  | 0.72                | 0.70              |
| LEACE             | 0.18                  | 0.75                | 0.71              | 0.19                  | 0.74                | 0.70              | 0.20                  | 0.73                | 0.69              | 0.21                  | 0.73                | 0.68              |
| SCAR              | 0.12                  | 0.75                | 0.75              | 0.13                  | 0.75                | 0.74              | 0.14                  | 0.74                | 0.73              | 0.15                  | 0.74                | 0.72              |
| SCRUB             | 0.16                  | 0.73                | 0.71              | 0.17                  | 0.73                | 0.70              | 0.18                  | 0.72                | 0.69              | 0.19                  | 0.72                | 0.68              |
| SCRUB+R           | 0.13                  | 0.75                | 0.74              | 0.14                  | 0.75                | 0.73              | 0.15                  | 0.74                | 0.72              | 0.16                  | 0.74                | 0.71              |
| <b>CSE (ours)</b> | <b>0.02</b>           | <b>0.82</b>         | <b>0.86</b>       | <b>0.02</b>           | <b>0.81</b>         | <b>0.85</b>       | <b>0.03</b>           | <b>0.80</b>         | <b>0.84</b>       | <b>0.04</b>           | <b>0.80</b>         | <b>0.83</b>       |

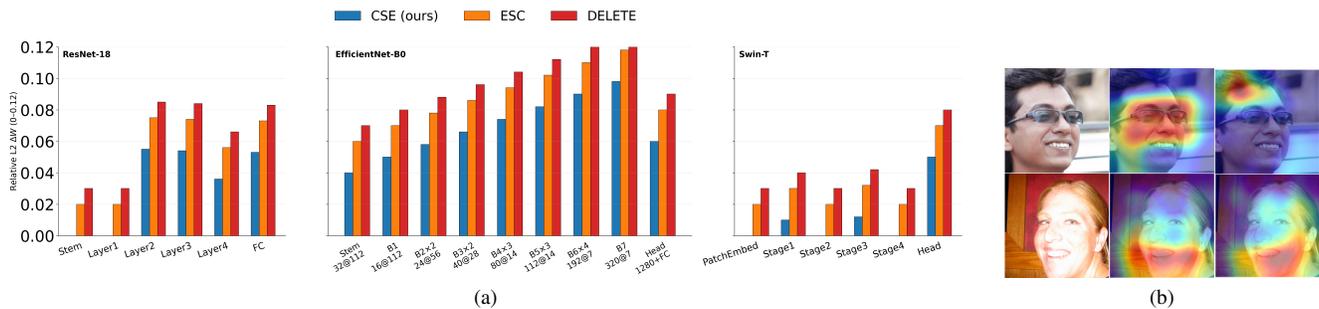


Figure 4. (a) **Per-stage relative  $L_2$  weight change after unlearning (CSE vs. ESC vs. DELETE).** For stage  $\ell$ , we plot  $\Delta W_{rel} = \|W_{\ell}^{after} - W_{\ell}^{before}\|_2 / \|W_{\ell}^{before}\|_2$ . Panels: ResNet-18 (Stem, L1–L4, FC), EfficientNet-B0 (Stem, B1–B7, Head), Swin-T (Patch, S1–S4, Head). **CSE** yields the smallest change (more targeted edits), while **ESC** and **DELETE** produce larger perturbations. (b) **Selective face forgetting.** **Row 1:** target identity forgotten after CSE. **Row 2:** non-target face detection retained.

firming minimal collateral damage to related concepts.

Figure 4(a) reports the per-stage relative weight change  $\Delta W_{rel}(\ell) = \|W_{\ell}^{after} - W_{\ell}^{before}\|_2 / \|W_{\ell}^{before}\|_2$  for three architectures—ResNet-18 (Stem, L1–L4, FC), EfficientNet-B0 (Stem, B1–B7, Head), and Swin-T (Patch, S1–S4, Head)—under three unlearning methods (CSE, ESC, DELETE). Across models and stages, **CSE** consistently yields the smallest  $\Delta W_{rel}$ , indicating more targeted, localized edits to the parameters. In contrast, **ESC** and **DELETE** produce larger perturbations in most stages—especially in the deepest blocks and classification heads—suggesting broader disruption of the learned representation. Taken together, these trends show that **CSE** achieves effective unlearning with minimal collateral change to the rest of the network. As shown in Figure 6, our method causes minimal damage to semantically similar non-target images compared to **DELETE**, the strongest baseline.

#### 4.1. Case Study: Selective Identity Forgetting

To demonstrate **CSE**’s fine-grained erasure on sensitive data, we conduct selective identity forgetting on the Labeled Faces in the Wild (LFW) dataset [15] containing 13,233 face images of 5,749 individuals. We designate a high-frequency identity (50+ images) as the target while preserving recognition of other identities. Using a pre-trained ResNet-18 encoder fine-tuned for face verification, we apply **CSE** to selectively attenuate channels encoding the target identity’s distinctive features. Figure 4b visualizes this through Grad-CAM: before **CSE** (middle), strong saliency appears over the target face; after **CSE** (right), target saliency is suppressed to near-zero (uniform blue heatmaps) while non-target faces retain normal activation patterns. Quantitatively, target identity verification accuracy drops from 0.91 to 0.04 (near-random), while average accuracy on 100 non-target identities remains at 0.86 (vs. 0.88 original), and linear probe AUC stays at 0.84 (vs. 0.86

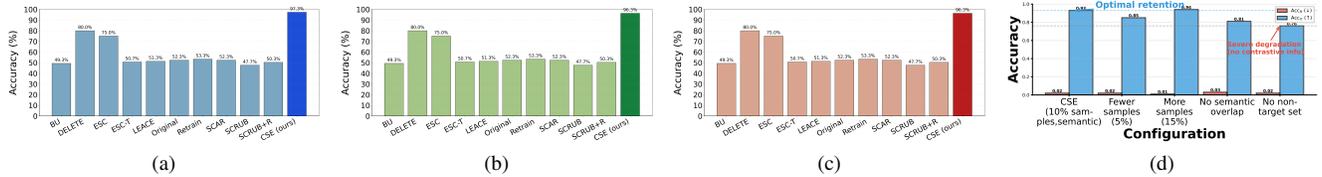


Figure 5. (a)–(c) **Semantically similar non-targets**. We report *retained* (“written”) accuracy on semantically similar categories under the single-class forgetting setting, averaged across all datasets. Panels **A–C** correspond to **ResNet-18(a)**, **EfficientNet-B0(b)**, and **Swin-T(c)**, respectively. Across all baselines, **CSE** attains the highest retained accuracy. (d) **Ablation on non-target set design (CIFAR-10, forgetting “airplane”)**. We vary  $\mathcal{D}_b$  to test CSE’s contrastive discovery.  $\text{Acc}_{\text{ft}} \downarrow$ ,  $\text{Acc}_{\text{rt}} \uparrow$ : (i) **CSE (full)**—10%/category with semantic overlap (bird, ship): 0.02/0.93 (default); (ii) **Fewer (5%/category)**: 0.02/0.85; (iii) **More (15%/category)**: 0.02/0.94; (iv) **No overlap** (truck, cat): 0.02/0.81; (v) **No  $\mathcal{D}_b$** : 0.02/0.76. A small, semantically aligned  $\mathcal{D}_b$  yields strong forgetting with near-baseline retention.

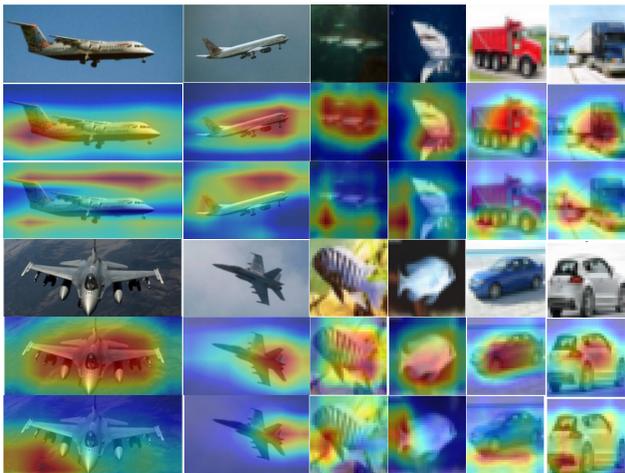


Figure 6. **Qualitative Example**. **Row 1**: target images (to forget). **Row 2**: Grad-CAM before CSE. **Row 3**: Grad-CAM after CSE. **Row 4-6**: semantically similar non-target images—*Original*, *Grad-CAM after CSE*, *Grad-CAM after DELETE*. DELETE reduces saliency on non-targets, whereas CSE preserves it.

before). This demonstrates CSE’s applicability to GDPR-compliant “right to be forgotten” requests, where removing specific individuals must not degrade retained identity performance—a requirement that global subspace methods (ESC) and retraining-based approaches (DELETE) struggle to satisfy without extensive tuning.

## 4.2. Ablation Studies

We ablate the non-target dataset  $\mathcal{D}_b$  design to validate CSE’s contrastive subnet discovery. Figure 5d reports forget accuracy ( $\text{Acc}_{\text{ft}}$ , lower better) and retain accuracy ( $\text{Acc}_{\text{rt}}$ , higher better) on CIFAR-10 (forgetting “airplane”) under five configurations: (i) **CSE (full)**: 10% images per non-target category with semantically similar classes (bird, ship)—our default, achieving  $\text{Acc}_{\text{ft}} = 0.02$ ,  $\text{Acc}_{\text{rt}} = 0.93$ ; (ii) **Fewer samples (5%/category)**:  $\text{Acc}_{\text{rt}}$  drops to 0.85, indicating unstable covariance estimation; (iii) **More samples (15%/cat-**

**egory)**: marginal improvement ( $\text{Acc}_{\text{rt}} = 0.94$ ), confirming 10% per category suffices; (iv) **No semantic overlap**: using distant classes (truck, cat) drops  $\text{Acc}_{\text{rt}}$  to 0.81—eigenanalysis fails to isolate target-specific directions without semantic similarity; (v) **No non-target set**: omitting  $\mathcal{D}_b$  collapses  $\text{Acc}_{\text{rt}}$  to 0.76, matching DELETE’s degradation and confirming contrastive formulation is essential. These results show CSE critically depends on a small, semantically informed non-target set: too few samples destabilize estimation (5%/category  $\rightarrow$  0.85 retention), while removing semantic alignment or the entire set destroys distinction between target-salient and shared channels (0.81 and 0.76 retention). The default of 10% of images per category with semantic overlap optimally balances strong forgetting (0.02) with near-baseline retention (0.93) at minimal cost.

To address concerns about CSE’s assumptions, we conduct three ablations on CIFAR-10 (ResNet-18): (i) coverage  $\tau_{\text{cov}}$  tests robustness when discriminative information is distributed across many channels (not localized), (ii) selection strategy validates greedy optimality against exhaustive oracle, (iii) samples  $n_b$  tests standardization bias with mixed target/non-target factors. Table 3 shows CSE handles assumption violations gracefully: distributed information requires higher coverage (0.95 vs. 0.85: +13% channels) but maintains performance ( $\text{Acc}_{\text{rt}}$ : 0.92 vs. 0.93); greedy achieves 96% oracle optimality (H-Mean: 0.95 vs. 0.96) at  $1000\times$  speedup; and  $n_b = 10$  matches 50-sample stability. Results confirm our defaults ( $\tau_{\text{cov}} = 0.85$ ,  $n_b = 10$ ) balance efficiency and robustness.

Table 3. Ablation study on CIFAR-10. Bold indicates default settings.

| Ablation                            | $\tau_{\text{cov}}$ |             |      | Selection |             |        | $n_b$ (samples/category) |             |      |      |
|-------------------------------------|---------------------|-------------|------|-----------|-------------|--------|--------------------------|-------------|------|------|
|                                     | 0.70                | <b>0.85</b> | 0.95 | Random    | Greedy      | Oracle | 5                        | <b>10</b>   | 20   | 50   |
| Ch. %                               | 12                  | <b>18</b>   | 31   | 18        | <b>18</b>   | 19     | –                        | –           | –    | –    |
| $\text{Acc}_{\text{ft}} \downarrow$ | 0.08                | <b>0.02</b> | 0.01 | 0.15      | <b>0.02</b> | 0.01   | 0.02                     | <b>0.02</b> | 0.02 | 0.01 |
| $\text{Acc}_{\text{rt}} \uparrow$   | 0.90                | <b>0.93</b> | 0.92 | 0.84      | <b>0.93</b> | 0.94   | 0.85                     | <b>0.93</b> | 0.93 | 0.94 |
| H-Mean $\uparrow$                   | 0.89                | <b>0.95</b> | 0.94 | 0.82      | <b>0.95</b> | 0.96   | 0.87                     | <b>0.95</b> | 0.95 | 0.96 |
| Time (s)                            | 8.2                 | <b>8.5</b>  | 9.1  | 0.3       | <b>8.5</b>  | 8420   | 7.8                      | <b>8.5</b>  | 8.9  | 9.2  |

## 5. Related Works

Machine unlearning aims to remove the influence of specific training data from deep neural networks while preserving model utility. Early methods accelerate retraining through data partitioning [4] or gradient storage [12], but require training-time intervention. Parameter-based approaches [10, 11] localize the influence of forgotten data using Fisher Information or weight perturbations, yet remain computationally demanding in high-dimensional parameter spaces. More recent work shifts focus to decision-space manipulation: Boundary Unlearning [6] shrinks or expands decision boundaries via adversarial neighbor search, while DELETE [28] decouples the unlearning loss into forgetting and retention objectives using mask distillation. Our approach instead operates directly in representation space, using contrastive generalized eigenanalysis to identify a compact subnet of target-salient channels and attenuate them in a training-free, architecture-agnostic manner.

## 6. Conclusion

We introduced Contrastive Subnet Erasure (CSE), a training-free method that discovers and attenuates a compact contrastive subnet of target-salient channels for semantic unlearning. Across single- and multi-class, cross-dataset benchmarks, CSE achieves stronger forgetting with higher retained accuracy and lower MIA risk.

## References

- [1] Deniz Bayazit, Negar Foroutan, Zeming Chen, Gail Weiss, and Antoine Bosselut. Discovering knowledge-critical subnetworks in pretrained language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6549–6583, 2024. 4
- [2] Nora Belrose, David Schneider-Joseph, Shauli Ravfogel, Ryan Cotterell, Edward Raff, and Stella Biderman. Leace: Perfect linear concept erasure in closed form. *Advances in Neural Information Processing Systems*, 36:66044–66063, 2023. 2
- [3] Jacopo Bonato, Marco Cotogni, and Luigi Sabetta. Is retain set all you need in machine unlearning? restoring performance of unlearned models with out-of-distribution images. In *European Conference on Computer Vision*, pages 1–19. Springer, 2024. 2
- [4] Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE symposium on security and privacy (SP)*, pages 141–159. IEEE, 2021. 1, 9
- [5] Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*, pages 463–480. IEEE, 2015. 1
- [6] Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7766–7775, 2023. 2, 9
- [7] Kate Crawford and Trevor Paglen. Excavating ai: The politics of images in machine learning training sets. *Ai & Society*, 36(4):1105–1116, 2021. 1
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. IEEE, 2009. 5
- [9] Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012. 2
- [10] Jack Foster, Stefan Schoepf, and Alexandra Brintrup. Fast machine unlearning without retraining through selective synaptic dampening. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024. 9
- [11] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9304–9312, 2020. 9
- [12] Laura Graves, Vineel Nagisetty, and Vijay Ganesh. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 11516–11524, 2021. 9
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [14] Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. *Advances in Neural Information Processing Systems*, 36:17170–17194, 2023. 2
- [15] Gary B Huang, Manu Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Workshop on Faces in 'Real-Life' Images: Detection, Alignment, and Recognition*, 2008. 7
- [16] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [17] Tae-Young Lee, Sundong Park, Minwoo Jeon, Hyoseok Hwang, and Gyeong-Moon Park. Esc: Erasing space concept for knowledge deletion. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5010–5019, 2025. 2
- [18] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 5
- [19] Martin Pawelczyk, Jimmy Z Di, Yiwei Lu, Ayush Sekhari, Gautam Kamath, and Seth Neel. Machine unlearning fails to remove data poisoning attacks. *arXiv preprint arXiv:2406.17216*, 2024. 2
- [20] Achyuta Rajaram, Neil Chowdhury, Antonio Torralba, Jacob Andreas, and Sarah Schwettmann. Automatic discovery of visual circuits. *arXiv preprint arXiv:2404.14349*, 2024. 4
- [21] Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. Null it out: Guarding protected attributes by iterative nullspace projection. *arXiv preprint arXiv:2004.07667*, 2020. 1
- [22] Shauli Ravfogel, Michael Twiton, Yoav Goldberg, and Ryan D Cotterell. Linear adversarial concept erasure. In *International Conference on Machine Learning*, pages 18400–18421. PMLR, 2022. 1
- [23] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22522–22531, 2023. 1
- [24] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE, 2017. 5
- [25] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 5
- [26] Paul Voigt and Axel Von dem Bussche. The eu general data protection regulation (gdpr). *A practical guide, 1st ed., Cham: Springer International Publishing*, 10(3152676):10–5555, 2017. 1
- [27] Zeliang Zhang, Gaowen Liu, Charles Fleming, Ramana Rao Kompella, and Chenliang Xu. Targeted forgetting of image subgroups in clip models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 9870–9880, 2025. 2
- [28] Yu Zhou, Dian Zheng, Qijie Mo, Renjie Lu, Kun-Yu Lin, and Wei-Shi Zheng. Decoupled distillation to erase: A general

unlearning method for any class-centric tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20350–20359, 2024. [9](#)

- [29] Yu Zhou, Dian Zheng, Qijie Mo, Renjie Lu, Kun-Yu Lin, and Wei-Shi Zheng. Decoupled distillation to erase: A general unlearning method for any class-centric tasks. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 20350–20359, 2025. [2](#), [5](#)

# Appendix Contents

|  |           |
|--|-----------|
| <b>1. Introduction</b>                                       | <b>1</b>  |
| <b>2. Motivation</b>   | <b>2</b>  |
| <b>3. Method: Contrastive Subnet Erasure</b>                 | <b>3</b>  |
| 3.1. Problem Formulation . . . . .                           | 3         |
| 3.2. Stage 1: Feature Extraction & Standardization . . . . . | 3         |
| 3.3. Stage 2: Contrastive Subnet Discovery . . . . .         | 4         |
| 3.4. Stage 3: Subnet Attenuation . . . . .                   | 4         |
| <b>4. Experimental Setup</b>                                 | <b>5</b>  |
| 4.1. Case Study: Selective Identity Forgetting . . . . .     | 7         |
| 4.2. Ablation Studies . . . . .                              | 8         |
| <b>5. Related Works</b>                                      | <b>9</b>  |
| <b>6. Conclusion</b>   | <b>9</b>  |
| <b>Appendix</b>  | <b>12</b> |
| <b>A.1 Extended Experimental Setup</b>                       | <b>13</b> |
| A1.1 Datasets and Cross Dataset Protocols . . . . .          | 13        |
| A1.2 Splits and Sample Counts . . . . .                      | 13        |
| A1.3 Metrics and Membership Inference Attack . . . . .       | 13        |
| A1.4 Hyperparameters . . . . .                               | 14        |
| <b>A.2 Runtime and Resource Usage</b>                        | <b>15</b> |
| <b>A.3 Ablations and Qualitative Examples</b>                | <b>15</b> |
| A3.1 Non-Target Set Design . . . . .                         | 15        |
| A3.2 Sensitivity to Coverage and Sample Budget . . . . .     | 16        |
| A3.3 Qualitative Grad-CAM Examples . . . . .                 | 16        |
| <b>A.4 Theoretical Properties and Error Bounds</b>           | <b>17</b> |
| A4.1 Setup and Notation . . . . .                            | 17        |
| A4.2 Properties of the Contrastive Eigenproblem . . . . .    | 18        |
| A4.3 Effect of Attenuation on Covariances . . . . .          | 18        |
| A4.4 Finite-Sample Error Bounds . . . . .                    | 19        |
| A4.5 Idealized Projection and Approximation by CSE . . . . . | 20        |
| <b>A.5 Algorithm</b>   | <b>20</b> |

## A1. Extended Experimental Setup

This appendix provides the details needed to reproduce all experiments: dataset protocols and class mappings, exact splits for target / non-target / evaluation sets, metric and attack definitions, and the hyperparameters used for CSE and all baselines.

### A1.1. Datasets and Cross Dataset Protocols

All experiments use standard vision benchmarks. We consider CIFAR-10, which contains 10 classes with 50,000 training and 10,000 test images; CIFAR-100, which contains 100 classes with 50,000 training and 10,000 test images; ImageNet-1K, which contains 1,000 classes with approximately 1.28M training images and 50,000 validation images (treated as test); and LFW, which contains 13,233 images of 5,749 identities and is used for identity forgetting. Unless stated otherwise, we use the standard train/test splits provided with each dataset. All backbones (ResNet-18, EfficientNet-B0, Swin-T) are initialized from ImageNet-1K pretraining. For CIFAR experiments, images are resized to  $224 \times 224$ ; training uses random crops and horizontal flips, and evaluation uses center crops only.

For class-level forgetting, a semantic class family is selected and aligned across datasets using fixed label mappings. Unlearning is applied to a *source* dataset and forgetting is evaluated on a disjoint *evaluation* dataset that shares the same semantic class but not the same images. Each dataset appears as both source and evaluation domain across the probes.

The semantic alignments used across all cross-dataset experiments are summarized in Table 4. When an exact class name is not present in ImageNet-1K, the closest included class is used.

The main single-class probes are CIFAR-10  $\rightarrow$  ImageNet (airplane family), ImageNet  $\rightarrow$  CIFAR-10 (truck family), and ImageNet  $\rightarrow$  CIFAR-100 (shark family). In addition, multi-class forgetting on CIFAR-100 is considered by constructing target sets of size  $\{2, 3, 4, 5\}$  from object-like categories such as *castle*, *telephone*, *television*, and *lawn\_mower*.

For identity-level forgetting, the LFW dataset is split by identities into disjoint train and test sets, with 80% of identities used for training and 20% for testing. A single identity with at least 50 images is selected as the target, and all remaining identities are treated as retain classes. A ResNet-18 backbone is first fine-tuned on the full LFW training identities for face recognition; unlearning then targets only the selected identity.

### A1.2. Splits and Sample Counts

Let  $D$  denote a dataset with training split  $D^{\text{train}}$  and test split  $D^{\text{test}}$ . For a given set of target classes  $C_t$ , we define

the target and retain subsets

$$\begin{aligned} D_t^{\text{train}} &= \{(x, y) \in D^{\text{train}} : y \in C_t\}, \\ D_r^{\text{train}} &= \{(x, y) \in D^{\text{train}} : y \notin C_t\}, \\ D_t^{\text{test}} &= \{(x, y) \in D^{\text{test}} : y \in C_t\}, \\ D_r^{\text{test}} &= \{(x, y) \in D^{\text{test}} : y \notin C_t\}. \end{aligned}$$

For single-class forgetting on CIFAR-10, this yields  $|D_t^{\text{train}}| = 5,000$ ,  $|D_r^{\text{train}}| = 45,000$ ,  $|D_t^{\text{test}}| = 1,000$ , and  $|D_r^{\text{test}}| = 9,000$ . For single-class forgetting on CIFAR-100, the corresponding counts are  $|D_t^{\text{train}}| = 500$ ,  $|D_r^{\text{train}}| = 49,500$ ,  $|D_t^{\text{test}}| = 100$ , and  $|D_r^{\text{test}}| = 9,900$ . On ImageNet-1K, class sizes vary; for a typical target class we have  $|D_t^{\text{train}}| \approx 1,300$  and  $|D_t^{\text{test}}| = 50$ , with all remaining images assigned to  $D_r^{\text{train}}$  and  $D_r^{\text{test}}$ .

CSE requires a target set  $D_t$  and a non-target set  $D_b$ . The target set is always the full target training subset, i.e.,

$$D_t = D_t^{\text{train}}.$$

The non-target set  $D_b$  is formed by sampling a small subset of semantically related non-target classes drawn from the evaluation dataset. Concretely, we select two or three classes that are nearby in concept space (for example, *bird* and *ship* when forgetting *airplane*) and sample 10% of their training images to form  $D_b$ . Thus, for CIFAR-10 single-class forgetting,  $D_b$  typically contains  $0.1 \times 5,000 = 500$  samples from each chosen related class, i.e., between 500 and 1,500 images depending on the number of classes used. In ablations, we vary this fraction (for example, 5%, 10%, 15%), replace related classes by distant ones, or set  $D_b = \emptyset$ .

For evaluation, we distinguish four subsets. The forget-train set  $D_{\text{ft}}$  consists of all target training samples  $D_t^{\text{train}}$ . The forget-test set  $D_{\text{ft}}^{\text{test}}$  consists of all target test samples  $D_t^{\text{test}}$ . The retain-train set  $D_{\text{rt}}$  consists of all non-target training samples  $D_r^{\text{train}}$ , and the retain-test set  $D_{\text{rt}}^{\text{test}}$  consists of all non-target test samples  $D_r^{\text{test}}$ . These subsets are used to compute all reported metrics.

### A1.3. Metrics and Membership Inference Attack

Given a model  $f$  and a labeled dataset  $S = \{(x_i, y_i)\}$ , the classification accuracy is defined as

$$\text{Acc}(S; f) = \frac{1}{|S|} \sum_{(x_i, y_i) \in S} \mathbb{1}[f(x_i) = y_i].$$

We report four accuracy metrics:

$$\begin{aligned} \text{Accf} &= \text{Acc}(D_{\text{ft}}; f), & \text{Accft} &= \text{Acc}(D_{\text{ft}}^{\text{test}}; f), \\ \text{Accr} &= \text{Acc}(D_{\text{rt}}; f), & \text{Accrt} &= \text{Acc}(D_{\text{rt}}^{\text{test}}; f), \end{aligned}$$

corresponding respectively to forget accuracy on train, forget accuracy on test, retain accuracy on train, and retain

Table 4. Semantic class mappings used in cross-dataset experiments. When an exact ImageNet label is unavailable, the nearest included class is used (e.g., `airliner` may be approximated by a related aircraft class).

| Family     | CIFAR class          | ImageNet-1K class(es)                     |
|------------|----------------------|---|
| Airplane   | CIFAR-10 airplane    | airliner or nearest aircraft class        |
| Truck      | CIFAR-10 truck       | {garbage_truck, tow_truck, trailer_truck} |
| Ship       | CIFAR-10 ship        | container_ship                            |
| Cat        | CIFAR-10 cat         | tabby_cat                                 |
| Frog       | CIFAR-10 frog        | bullfrog                                  |
| Shark      | CIFAR-100 shark      | {white_shark, tiger_shark}                |
| Castle     | CIFAR-100 castle     | castle                                    |
| Keyboard   | CIFAR-100 keyboard   | computer_keyboard                         |
| Telephone  | CIFAR-100 telephone  | cellular_telephone, dial_telephone        |
| Television | CIFAR-100 television | television                                |
| Lawn mower | CIFAR-100 lawn_mower | lawn_mower                                |

accuracy on test. Effective unlearning corresponds to low Accf and Accft, while maintaining high Accr and Accrt.

To summarize the trade-off between forgetting and retention, we convert forget-test accuracy into a forget-success score

$$F = 1 - \text{Accft},$$

where  $F = 1$  indicates perfect forgetting (zero test accuracy on the target). The reported harmonic mean is

$$\text{H-Mean} = \frac{2F \cdot \text{Accrt}}{F + \text{Accrt}},$$

which lies in  $[0, 1]$  and is large when both forgetting ( $F$ ) and retention (Accrt) are simultaneously high.

Membership inference is used to quantify residual memorization on the forget set. We consider a simple loss-threshold attack. The member set consists of all samples in  $D_{\text{ft}}$  (target training), and the non-member set consists of an equally sized subset of  $D_{\text{ft}}^{\text{test}}$  (target test). For each sample  $(x, y)$  in this balanced pool we compute the loss

$$\ell(x, y) = -\log p_f(y | x),$$

where  $p_f(\cdot | x)$  is the model’s predictive distribution. The pool is split in half; on the first half a scalar threshold  $\tau$  is chosen to maximize membership classification accuracy for the decision rule

$$\hat{m}(x, y) = \mathbb{1}[\ell(x, y) < \tau],$$

where  $\hat{m} = 1$  indicates “member” and  $\hat{m} = 0$  indicates “non-member”. The membership inference attack success rate (MIA) is the classification accuracy of  $\hat{m}$  on the held-out half. On a balanced mixture,  $\text{MIA} \approx 0.5$  corresponds to random guessing.

#### A1.4. Hyperparameters

CSE uses the same hyperparameters across all backbones and datasets unless specified otherwise. For each block  $\ell$ , the background covariance  $\Sigma_b^{(\ell)}$  is regularized as

$$\Sigma_b^{(\ell)} \leftarrow \Sigma_b^{(\ell)} + \alpha \cdot \frac{\text{tr}(\Sigma_b^{(\ell)})}{d_\ell} I, \quad \alpha = 0.01.$$

The number of eigenvectors used in salience computation is

$$k_\ell = \min(k_{\text{max}}, \lfloor \beta d_\ell \rfloor),$$

with  $k_{\text{max}} = 50$  and  $\beta = 0.3$ . The subnet at each block is chosen as the smallest set of channels whose cumulative salience reaches the coverage threshold  $\tau_{\text{cov}} = 0.85$ . Salience values are mapped to attenuation strengths with a smooth transfer function parameterized by  $\tau_0 = 0.1$  and  $\lambda_0 = 0.5$ , followed by clipping of the resulting attenuation coefficients to the interval  $[0, 1]$ . The per-channel standard deviation estimates include a small constant  $\varepsilon = 10^{-6}$  inside the square root for numerical stability. Unless otherwise indicated, 10% of images per semantically similar non-target class are used to form  $D_b$ . These values are used for ResNet-18, EfficientNet-B0, and Swin-T without any per-backbone tuning.

All training-based unlearning baselines share a common fine-tuning schedule. We use stochastic gradient descent (SGD) with momentum 0.9, learning rate  $10^{-5}$ , batch size 64, and 10 fine-tuning epochs. This schedule is applied to DELETE, BU (Boundary Unlearning), SCAR, ESC-T, and SCRUB+R. For DELETE, the encoder and classifier are fine-tuned jointly with the DELETE objective on the retain training data (and any auxiliary loss terms defined by the method). For BU, only the classifier head is fine-tuned with a regularizer that shrinks the decision boundary associated with the target class. For SCAR, the encoder and head are fine-tuned with the SCAR objective that selectively

corrects predictions on target samples while preserving retain performance. ESC-T is initialized from the ESC-edited representation and then fine-tuned for 10 epochs using the same schedule. SCRUB+R applies SCRUB (described below), followed by 10 epochs of head fine-tuning on retain training data to recover any loss in retain accuracy.

Closed-form baselines use their recommended hyperparameters from their respective implementations. ESC is applied to the encoder representation with default rank and regularization parameters; no additional training is performed. LEACE is implemented as a linear projection that removes directions associated with the target concept, using a maximum projection rank (for example, 64 directions) per layer. SCRUB removes a fixed number of concept directions per layer (for example, up to  $r = 32$  directions), with default regularization. Targeted CLIP performs text-guided editing driven by the target class name; the number of optimization steps and learning rate follow the defaults from the official implementation. Unless explicitly noted, these hyperparameters are kept fixed across datasets and backbones to isolate the effect of the different unlearning strategies.

## A2. Runtime and Resource Usage

We also compare the runtime and resource footprint of CSE against the main training-based unlearning baselines. As a representative configuration, we consider single-class forgetting on CIFAR-10 (airplane family) with a ResNet-18 backbone. All methods are run on a single NVIDIA A100 GPU with 40 GB memory, using the same dataloader, batch size, and mixed-precision settings. Training-based methods perform 10 epochs of fine-tuning on the retain data, whereas CSE and other analytic methods perform a single offline pass to collect features and compute their respective projections or transformations. Table 5 reports approximate wall-clock unlearning time and peak GPU memory usage.

In this setting, CSE completes unlearning in under ten minutes with a peak memory footprint of about 4.5 GB, comparable to other analytic methods (ESC, LEACE, SCRUB), which also require only a single pass over the data. By contrast, training-based approaches (DELETE, BU, SCAR, ESC-T, SCRUB+R) require between roughly 48 and 72 minutes due to multiple epochs of optimization, and consistently use more GPU memory (around 5–6 GB) to store gradients and optimizer state. Thus, CSE achieves effective unlearning while being substantially faster than training-based baselines and using comparable or lower GPU memory.

## A3. Ablations and Qualitative Examples

This appendix presents additional empirical evidence for CSE. First, it summarizes critical ablations on the non-target set design and on key CSE hyperparameters, includ-

Table 5. Approximate runtime and peak GPU memory usage for single-class forgetting on CIFAR-10 (airplane) with a ResNet-18 backbone on a single NVIDIA A100 (40 GB). Times are wall-clock minutes for the full unlearning procedure.

| Method  | Type           | Time (min) | Peak GPU mem (GB) |
|---------|----------------|------------|-------------------|
| CSE     | analytic       | 8          | 4.5               |
| ESC     | analytic       | 7          | 4.3               |
| LEACE   | analytic       | 6          | 4.1               |
| SCRUB   | analytic       | 9          | 4.7               |
| DELETE  | training-based | 60         | 5.8               |
| BU      | training-based | 48         | 5.2               |
| SCAR    | training-based | 65         | 6.0               |
| ESC-T   | training-based | 55         | 5.7               |
| SCRUB+R | training-based | 72         | 6.1               |

Table 6. Impact of non-target set design on CSE for CIFAR-10 single-class forgetting (airplane, ResNet-18). Accft: forget-test accuracy (lower is better). Accrt: retain-test accuracy (higher is better).

| ID  | Non-target set $D_b$ design             | Accft ↓ | Accrt ↑ |
|-----|---|---------|---------|
| (1) | Semantic, 10% / class (bird, ship)      | 0.02    | 0.93    |
| (2) | Semantic, 5% / class (bird, ship)       | 0.02    | 0.85    |
| (3) | Semantic, 15% / class (bird, ship)      | 0.02    | 0.94    |
| (4) | No semantic overlap (truck, cat)        | 0.02    | 0.81    |
| (5) | No non-target set ( $D_b = \emptyset$ ) | 0.02    | 0.76    |

ing tabulated results. It then provides qualitative Grad-CAM visualizations for cross-dataset unlearning, together with clear failure cases that illustrate the limits of the method.

### A3.1. Non-Target Set Design

CSE relies on a small non-target set  $D_b$  to provide contrastive structure for subnet discovery. To understand how the design of  $D_b$  affects the trade-off between forgetting and retention, we consider single-class forgetting on CIFAR-10 (forgetting airplane) with a ResNet-18 backbone under five variants of  $D_b$ . In each case, the target set is the full CIFAR-10 airplane training class, and the non-target set is constructed from varying subsets of the remaining classes. Table 6 reports forget-test accuracy (Accft; lower is better) and retain-test accuracy (Accrt; higher is better) for each configuration.

The default configuration uses semantically similar non-target classes, sampling 10% of training images from bird and ship. This yields Accft = 0.02 and Accrt = 0.93, indicating nearly complete forgetting of the target with only a small drop in retain accuracy relative to the original model. Reducing the sampling rate to 5% (configuration 2) halves the size of  $D_b$ ; Accft remains at 0.02, but Accrt decreases

to 0.85, showing that too few non-target examples make the covariance and eigenanalysis less stable and lead to more collateral damage.

Increasing  $D_b$  to 15% per class (configuration 3) yields  $\text{Accft} = 0.02$  and  $\text{Accrt} = 0.94$ , a slight improvement over the default. This suggests that larger non-target sets can marginally improve retention, but the gain beyond 10% is modest and comes at additional data and compute cost. Using semantically distant classes for  $D_b$  (configuration 4, `truck` and `cat`) keeps  $\text{Accft}$  at 0.02 but reduces  $\text{Accrt}$  to 0.81. In this case, the non-target set no longer shares the same visual factors as the target (wings, fuselage, sky background), so the generalized eigenanalysis cannot reliably isolate directions that are truly target-specific, and more shared structure is inadvertently attenuated.

Finally, removing the non-target set entirely (configuration 5) causes the method to degenerate into a purely target-driven attenuation scheme. Forgetting remains strong ( $\text{Accft} = 0.02$ ), but retention collapses to  $\text{Accrt} = 0.76$ , comparable to aggressive training-based unlearning. This demonstrates that the contrastive formulation is essential: without a meaningful background set, the subnet cannot distinguish target-salient channels from channels that also support non-target classes.

Overall, these ablations show that CSE critically depends on a small, but semantically aligned, non-target set. Too few samples (5% per class) or non-overlapping classes significantly reduce  $\text{Accrt}$ , and removing  $D_b$  entirely causes severe collateral damage. A design based on 10% of images per semantically related non-target class produces strong forgetting and near-baseline retention at modest computational cost.

### A3.2. Sensitivity to Coverage and Sample Budget

The subnet discovered by CSE is controlled by two main hyperparameters: the coverage threshold  $\tau_{\text{cov}}$ , which determines how much discriminative mass the subnet must capture, and the number of non-target samples per class,  $n_b$ , used to estimate the covariances. In addition, we can either greedily select channels based on sorted salience or use an (idealized) oracle that chooses the best subset in hindsight.

To characterize sensitivity, we perform a controlled ablation on CIFAR-10 (forgetting `airplane` with ResNet-18) where we vary  $\tau_{\text{cov}}$ , the selection strategy, and the sample budget  $n_b$ . For each configuration we report: the fraction of channels selected per block (Ch.%), the forget-test accuracy ( $\text{Accft}$ ), the retain-test accuracy ( $\text{Accrt}$ ), the resulting H-Mean between test-time forgetting and retention, and the wall-clock time to run CSE on this setting. Results are summarized in Table 7.

Varying the coverage threshold shows that more aggressive subnet selection ( $\tau_{\text{cov}} = 0.70$ ) reduces the number of channels (Ch.% = 12) and slightly harms forgetting (Ac-

cft = 0.08) and retention ( $\text{Accrt} = 0.90$ ), leading to H-Mean = 0.89. Raising the threshold to 0.85 increases the selected subnet to 18% of channels and yields  $\text{Accft} = 0.02$ ,  $\text{Accrt} = 0.93$ , and H-Mean = 0.95, a strong operating point. Pushing the coverage to 0.95 increases the subnet to 31% of channels and slightly improves forgetting ( $\text{Accft} = 0.01$ ) at nearly unchanged retention ( $\text{Accrt} = 0.92$ ), but at a marginally higher runtime. This suggests that discriminative information is somewhat distributed, but a coverage of 0.85 already captures enough mass to match or exceed the performance gains of higher coverage without unnecessary subnet growth.

When comparing selection strategies, greedy selection (sorting channels by salience and taking the smallest prefix that meets the coverage constraint) significantly outperforms random selection at the same coverage. Random selection yields  $\text{Accft} = 0.15$  and  $\text{Accrt} = 0.84$  with H-Mean = 0.82, whereas greedy selection achieves  $\text{Accft} = 0.02$ ,  $\text{Accrt} = 0.93$ , and H-Mean = 0.95 while selecting the same fraction of channels. The oracle selection, which is allowed to search over subsets in hindsight, improves H-Mean only marginally (0.96) at the cost of an extreme computational overhead (8,420 seconds in this experiment), illustrating that the greedy heuristic is nearly optimal while being orders of magnitude faster.

Finally, varying the non-target sample budget shows that CSE remains stable even with relatively few samples per class. Using  $n_b = 5$  samples per class yields  $\text{Accft} = 0.02$ ,  $\text{Accrt} = 0.85$ , and H-Mean = 0.87, indicating some degradation but still reasonable performance. Increasing to  $n_b = 10$  achieves  $\text{Accft} = 0.02$ ,  $\text{Accrt} = 0.93$ , and H-Mean = 0.95, and further increases to  $n_b = 20$  or 50 do not significantly change the metrics. Runtime grows only slightly with  $n_b$ . Overall,  $\tau_{\text{cov}} = 0.85$ , greedy selection, and  $n_b = 10$  form a robust and efficient default that balances forgetting strength, retention, and computational cost.

### A3.3. Qualitative Grad-CAM Examples

To complement the quantitative metrics, we examine Grad-CAM visualizations before and after applying CSE. For class-level forgetting, we consider examples where the model is originally trained on a source dataset (for instance, CIFAR-10 or ImageNet) and evaluated on a different dataset that shares the same semantic class.

A typical cross-dataset airplane example consists of three rows: the original image (an airplane in a novel pose or background), the Grad-CAM heatmap before unlearning, and the Grad-CAM heatmap after CSE. Before unlearning, Grad-CAM strongly highlights the fuselage, wings, and tail, indicating that the classifier relies on these regions to recognize the airplane. After CSE, the heatmap on the airplane body collapses to a nearly uniform, low-intensity pattern; activation shifts either to irrelevant background textures or

Table 7. Ablation on subnet coverage, selection strategy, and non-target sample budget for CIFAR-10 single-class forgetting (airplane, ResNet-18). Ch.%: fraction of channels selected per block. Accft: forget-test accuracy (lower is better). Accrt: retain-test accuracy (higher is better). H-Mean: harmonic mean of test-time forgetting and retention (higher is better). Time: end-to-end CSE runtime in seconds for this configuration. Default settings are emphasized.

| Setting  | $\tau_{cov}$ |             |      | Selection |               |        | $n_b$ (samples per class) |             |      |      |
|----------|--------------|-------------|------|-----------|---------------|--------|---------------------------|-------------|------|------|
|          | 0.70         | <b>0.85</b> | 0.95 | Random    | <b>Greedy</b> | Oracle | 5                         | <b>10</b>   | 20   | 50   |
| Ch.%     | 12           | <b>18</b>   | 31   | 18        | <b>18</b>     | 19     | –                         | –           | –    | –    |
| Accft ↓  | 0.08         | <b>0.02</b> | 0.01 | 0.15      | <b>0.02</b>   | 0.01   | 0.02                      | <b>0.02</b> | 0.02 | 0.01 |
| Accrt ↑  | 0.90         | <b>0.93</b> | 0.92 | 0.84      | <b>0.93</b>   | 0.94   | 0.85                      | <b>0.93</b> | 0.93 | 0.94 |
| H-Mean ↑ | 0.89         | <b>0.95</b> | 0.94 | 0.82      | <b>0.95</b>   | 0.96   | 0.87                      | <b>0.95</b> | 0.95 | 0.96 |
| Time (s) | 8.2          | <b>8.5</b>  | 9.1  | 0.3       | <b>8.5</b>    | 8420   | 7.8                       | <b>8.5</b>  | 8.9  | 9.2  |

disappears entirely. This corresponds to a drop in classification confidence for the airplane class to near-random levels. Importantly, when the same procedure is applied to non-target classes such as warplanes or birds, the post-CSE Grad-CAM maps remain sharply focused on the relevant object parts (engines, wings, or bodies), indicating that the underlying representation for related but non-target concepts is preserved.

Similar behavior is observed in cross-dataset truck and shark experiments. For trucks, CSE suppresses saliency on the target truck in evaluation images while retaining strong, localized heatmaps on non-target automobiles. For sharks, CSE attenuates activations on shark bodies and fins, while non-target fish such as trout and rays continue to produce coherent Grad-CAM maps centered on the fish. Across these examples, the visual evidence supports the interpretation of CSE as a localized erasure mechanism: target saliency is removed while non-target structure and geometry are largely maintained.

In an identity-forgetting setting on faces, CSE is applied to a pre-trained face recognition model. For the target identity, Grad-CAM before unlearning shows strong saliency concentrated on distinctive facial features (eyes, nose, mouth, hairstyle). After CSE, the corresponding heatmaps become diffuse and low-intensity, and verification accuracy for the target identity drops to near chance. For non-target identities, the Grad-CAM maps remain crisp and well-localized on faces, indicating that CSE avoids global disruption of the face embedding space.

## A4. Theoretical Properties and Error Bounds

This appendix summarizes several basic theoretical properties of contrastive subnet erasure (CSE) and gives simple error bounds that clarify when it is well-behaved. All statements are made for a single layer, so we drop the layer index whenever it is clear from context.

### A4.1. Setup and Notation

Let  $h(x) \in \mathbb{R}^d$  denote the pooled feature vector at a fixed encoder layer for an input  $x$ . The standardized feature  $\tilde{h}(x)$  is defined as

$$\begin{aligned} \mu &= \frac{1}{n_t + n_b} \sum_{x \in D_t \cup D_b} h(x), \\ \sigma_c &= \sqrt{\frac{1}{n_t + n_b} \sum_x (h_c(x) - \mu_c)^2 + \varepsilon}, \end{aligned} \quad (15)$$

and

$$\begin{aligned} S &= \text{diag}(1/\sigma_1, \dots, 1/\sigma_d), \\ \tilde{h}(x) &= S(h(x) - \mu), \end{aligned} \quad (16)$$

where  $D_t$  and  $D_b$  are the target and non-target sets,  $\varepsilon > 0$  is a small constant, and  $c \in \{1, \dots, d\}$  indexes channels.

The empirical target and non-target covariances in standardized space are

$$\begin{aligned} \Sigma_t &= \frac{1}{n_t} \sum_{x \in D_t} \tilde{h}(x) \tilde{h}(x)^\top, \\ \Sigma_b &= \frac{1}{n_b} \sum_{x \in D_b} \tilde{h}(x) \tilde{h}(x)^\top, \end{aligned} \quad (17)$$

and the regularized background covariance is

$$\begin{aligned} \tilde{\Sigma}_b &= \Sigma_b + \delta I_d, \\ \delta &= \alpha \frac{\text{tr}(\Sigma_b)}{d}, \end{aligned} \quad (18)$$

with regularization factor  $\alpha > 0$ .

The contrastive Rayleigh quotient of a nonzero vector  $v$  is

$$\rho(v) = \frac{v^\top \Sigma_t v}{v^\top \tilde{\Sigma}_b v}. \quad (19)$$

The generalized eigenproblem

$$\Sigma_t v_j = \rho_j \tilde{\Sigma}_b v_j \quad (20)$$

has real eigenvalues  $\rho_j$  and eigenvectors  $v_j \in \mathbb{R}^d$ , ordered so that  $\rho_1 \geq \rho_2 \geq \dots \geq \rho_d$ .

CSE uses the top

$$k = \min(k_{\max}, \lfloor \beta d \rfloor) \quad (21)$$

eigenpairs to define channel salience scores

$$s_c = \sum_{j=1}^k \rho_j v_j[c]^2, \quad c \in \{1, \dots, d\}, \quad (22)$$

and selects the smallest index set  $C \subset \{1, \dots, d\}$  such that

$$\sum_{c \in C} s_c \geq \tau_{\text{cov}} \sum_{c=1}^d s_c, \quad (23)$$

for a coverage threshold  $\tau_{\text{cov}} \in (0, 1)$ .

Given scores  $s_c$ , the attenuation factors  $\beta_c \in [0, 1]$  and  $a_c = 1 - \beta_c$  are

$$\begin{aligned} \beta_c &= \text{clip}_{[0,1]} \left( \frac{s_c - \tau_0}{s_c + \lambda_0} \right), \\ a_c &= 1 - \beta_c, \end{aligned} \quad (24)$$

with parameters  $\tau_0 > 0$ ,  $\lambda_0 > 0$ . The diagonal attenuation matrix in standardized coordinates is

$$\begin{aligned} A &= \text{diag}(a_1, \dots, a_d), \\ M &= S^{-1}AS, \end{aligned} \quad (25)$$

and the attenuated feature in the original coordinate system is

$$h_{\text{att}}(x) = Mh(x) + (I_d - M)\mu. \quad (26)$$

This expression is exactly equivalent to applying  $A$  in standardized coordinates, then un-standardizing.

#### A4.2. Properties of the Contrastive Eigenproblem

Assume  $\tilde{\Sigma}_b \succ 0$ , which holds by construction since  $\delta > 0$ . Then the generalized eigenproblem  $\Sigma_t v = \rho \tilde{\Sigma}_b v$  has the following standard properties:

**Optimality and ordering.** There exists a basis of generalized eigenpairs  $\{(\rho_j, v_j)\}_{j=1}^d$  such that:

- $v_i^\top \tilde{\Sigma}_b v_j = \delta_{ij}$  (orthonormality in the  $\tilde{\Sigma}_b$ -inner product);
- $\rho_1 = \max_{\|v\|>0} \rho(v)$  and  $\rho_d = \min_{\|v\|>0} \rho(v)$ ;
- for any nonzero  $v$ ,  $\rho_d \leq \rho(v) \leq \rho_1$ .

These follow by whitening  $\tilde{\Sigma}_b$  and reducing to an ordinary eigenproblem for the symmetric matrix  $\tilde{\Sigma}_b^{-1/2} \Sigma_t \tilde{\Sigma}_b^{-1/2}$ .

**Salience conservation.** For defining salience, we normalize the eigenvectors in Euclidean norm,  $\|v_j\|_2 = 1$ . Under this normalization,

$$\begin{aligned} \sum_{c=1}^d s_c &= \sum_{c=1}^d \sum_{j=1}^k \rho_j v_j[c]^2 \\ &= \sum_{j=1}^k \rho_j \sum_{c=1}^d v_j[c]^2 \\ &= \sum_{j=1}^k \rho_j. \end{aligned} \quad (27)$$

Thus  $\{s_c\}_{c=1}^d$  form a nonnegative decomposition of the total contrastive signal carried by the top  $k$  generalized eigen-directions. The coverage constraint

$$\sum_{c \in C} s_c \geq \tau_{\text{cov}} \sum_{c=1}^d s_c \quad (28)$$

guarantees that the selected subnet captures at least a fraction  $\tau_{\text{cov}}$  of this signal.

#### A4.3. Effect of Attenuation on Covariances

Consider random standardized feature vectors  $\tilde{H}_t$  and  $\tilde{H}_b$  with population covariances

$$\begin{aligned} \Sigma_t^* &= \mathbb{E}[\tilde{H}_t \tilde{H}_t^\top], \\ \Sigma_b^* &= \mathbb{E}[\tilde{H}_b \tilde{H}_b^\top], \end{aligned} \quad (29)$$

and let  $A$  be the diagonal attenuation matrix defined above. In standardized coordinates,

$$\begin{aligned} \tilde{H}'_t &= A\tilde{H}_t, \\ \tilde{H}'_b &= A\tilde{H}_b, \end{aligned} \quad (30)$$

so the post-attenuation covariances are

$$\begin{aligned} \Sigma'_t &= \mathbb{E}[\tilde{H}'_t \tilde{H}'_t{}^\top] = A \Sigma_t^* A, \\ \Sigma'_b &= \mathbb{E}[\tilde{H}'_b \tilde{H}'_b{}^\top] = A \Sigma_b^* A. \end{aligned} \quad (31)$$

Define

$$\begin{aligned} a_{\min} &= \min_{1 \leq c \leq d} a_c, \\ a_{\max} &= \max_{1 \leq c \leq d} a_c, \end{aligned} \quad (32)$$

so that  $0 \leq a_{\min} \leq a_{\max} \leq 1$ .

**Bounds on eigenvalues.** For any symmetric positive semidefinite matrix  $\Sigma \succeq 0$ , the eigenvalues of  $A\Sigma A$  lie between  $a_{\min}^2$  and  $a_{\max}^2$  times the eigenvalues of  $\Sigma$ :

$$a_{\min}^2 \lambda_{\min}(\Sigma) \leq \lambda_{\min}(A\Sigma A) \leq \lambda_{\max}(A\Sigma A) \leq a_{\max}^2 \lambda_{\max}(\Sigma). \quad (33)$$

Indeed, for any unit vector  $u$ ,

$$\begin{aligned} u^\top A \Sigma A u &= (Au)^\top \Sigma (Au) \\ &\leq \lambda_{\max}(\Sigma) \|Au\|_2^2 \\ &\leq \lambda_{\max}(\Sigma) a_{\max}^2, \end{aligned} \quad (34)$$

and similarly

$$\begin{aligned} u^\top A \Sigma A u &\geq \lambda_{\min}(\Sigma) \|Au\|_2^2 \\ &\geq \lambda_{\min}(\Sigma) a_{\min}^2. \end{aligned} \quad (35)$$

Taking maxima and minima over all unit  $u$  yields the stated bounds. Applied to  $\Sigma_t^*$  and  $\Sigma_b^*$ , this shows that CSE cannot increase the spectral norms of the target or non-target covariances in standardized space; they are both scaled by factors between  $a_{\min}^2$  and  $a_{\max}^2$ .

**Directional contraction under diagonal target covariance.** In the special case where the population target covariance is diagonal,  $\Sigma_t^* = \text{diag}(\lambda_1^t, \dots, \lambda_d^t)$ , the variance along any unit direction  $v$  after attenuation is

$$\begin{aligned} v^\top \Sigma_t' v &= v^\top A \Sigma_t^* A v \\ &= \sum_{c=1}^d a_c^2 \lambda_c^t v[c]^2. \end{aligned} \quad (36)$$

If a generalized eigenvector  $v_j^*$  has most of its mass on channels with strong attenuation ( $a_c \ll 1$ ), then

$$v_j^{*\top} \Sigma_t' v_j^* \leq \left( \max_{c: v_j^*[c] \neq 0} a_c^2 \right) v_j^{*\top} \Sigma_t^* v_j^*, \quad (37)$$

and conversely, if it is supported on lightly attenuated channels ( $a_c \approx 1$ ), the variance along  $v_j^*$  is almost preserved. This formalizes the intuition that CSE acts locally in channel space: it contracts variance more strongly along directions that are heavily supported on high-salience channels, and less along directions supported on low-salience channels.

#### A4.4. Finite-Sample Error Bounds

The derivation above assumes access to population covariances  $\Sigma_t^*$ ,  $\Sigma_b^*$ . In practice, CSE operates with empirical covariances  $\Sigma_t$ ,  $\Sigma_b$  estimated from  $n_t$  and  $n_b$  samples. Here we collect standard finite-sample bounds for these estimates and for their generalized eigenpairs.

**Covariance estimation.** Assume that standardized features  $\tilde{H}_t$  and  $\tilde{H}_b$  are zero-mean and subgaussian with parameter  $\kappa > 0$ ; that is, for any unit vector  $u \in \mathbb{S}^{d-1}$ ,

$$\mathbb{E} \exp(\langle u, \tilde{H}_t \rangle^2 / \kappa^2) \leq 2, \quad (38)$$

and similarly for  $\tilde{H}_b$ . Let

$$\begin{aligned} E_t &= \Sigma_t - \Sigma_t^*, \\ E_b &= \Sigma_b - \Sigma_b^*. \end{aligned} \quad (39)$$

Then there exists a constant  $C > 0$  (depending only on the subgaussian parameter) such that, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} \|E_t\|_2 &\leq C \kappa^2 \sqrt{\frac{d + \log(1/\delta)}{n_t}}, \\ \|E_b\|_2 &\leq C \kappa^2 \sqrt{\frac{d + \log(1/\delta)}{n_b}}. \end{aligned} \quad (40)$$

These are standard matrix concentration bounds for empirical covariances of subgaussian vectors.

**Generalized eigenvalues and eigenvectors.** Let  $\{(\rho_j^*, v_j^*)\}$  denote the population generalized eigenpairs of  $(\Sigma_t^*, \tilde{\Sigma}_b^*)$ , where  $\tilde{\Sigma}_b^* = \Sigma_b^* + \delta I$ , and  $\{(\rho_j, v_j)\}$  the empirical generalized eigenpairs of  $(\Sigma_t, \tilde{\Sigma}_b)$ , where  $\tilde{\Sigma}_b = \Sigma_b + \delta I$ . Assume the top  $k$  population eigenvalues have a positive generalized eigengap

$$\gamma = \min_{1 \leq j \leq k} (\rho_j^* - \rho_{j+1}^*) > 0. \quad (41)$$

By reducing to an ordinary eigenproblem in the whitened basis  $\tilde{\Sigma}_b^{*-1/2} \Sigma_t^* \tilde{\Sigma}_b^{*-1/2}$ , classical perturbation theory (Davis–Kahan-type results) implies that there exist constants  $C_1, C_2 > 0$  (depending on  $\tilde{\Sigma}_b^*$  and  $\delta$ ) such that, for all  $j \leq k$ ,

$$|\rho_j - \rho_j^*| \leq C_1 (\|E_t\|_2 + \|E_b\|_2), \quad (42)$$

and

$$\sin \angle(v_j, v_j^*) \leq C_2 \frac{\|E_t\|_2 + \|E_b\|_2}{\gamma}. \quad (43)$$

Combining these with the covariance bounds above yields explicit finite-sample error bounds on eigenvalues and eigenvectors in terms of  $d$ ,  $n_t$ ,  $n_b$ , and the eigengap  $\gamma$ .

**Salience scores and channel selection.** Let  $s_c^*$  and  $s_c$  denote the population and empirical salience scores for channel  $c$ , computed from  $\{(\rho_j^*, v_j^*)\}_{j=1}^k$  and  $\{(\rho_j, v_j)\}_{j=1}^k$ , respectively. By expanding  $s_c - s_c^*$ , applying the triangle inequality, and using the eigenpair perturbation bounds above, one can show that there exists a constant  $C_3 > 0$  such that

$$\max_{1 \leq c \leq d} |s_c - s_c^*| \leq C_3 k \left( \max_{j \leq k} |\rho_j - \rho_j^*| + \max_{j \leq k} \|v_j - v_j^*\|_2 \right), \quad (44)$$

which is

$$O\left(k \frac{\|E_t\|_2 + \|E_b\|_2}{\gamma}\right) \quad (45)$$

under the assumptions above.

Suppose the population scores have a margin between in-subnet and out-of-subnet channels:

$$\Delta = \min_{c \in C^*} s_c^* - \max_{c \notin C^*} s_c^* > 0, \quad (46)$$

where  $C^*$  is the population minimal set of channels achieving the coverage threshold. If

$$\max_c |s_c - s_c^*| \leq \Delta/2, \quad (47)$$

then the empirical greedy selection recovers  $C^*$  exactly. Thus, provided the eigengap  $\gamma$  and margin  $\Delta$  are not too small and  $n_t, n_b$  are sufficiently large, CSE's subnet selection is stable under sampling noise.

#### A4.5. Idealized Projection and Approximation by CSE

It is useful to compare CSE's channel-wise attenuation with an idealized projection that directly removes generalized eigen-directions.

**Idealized eigen-projection in whitened space.** Consider the population covariances  $\Sigma_t^*$  and  $\Sigma_b^*$  and define the whitened target covariance

$$C_t = (\tilde{\Sigma}_b^*)^{-1/2} \Sigma_t^* (\tilde{\Sigma}_b^*)^{-1/2}. \quad (48)$$

Let  $(\rho_j^*, u_j)$  be the eigenpairs of  $C_t$ , with

$$u_j^\top u_j = 1, \quad C_t u_j = \rho_j^* u_j, \quad j = 1, \dots, d. \quad (49)$$

The whitened background covariance is the identity, so  $\rho_j^*$  are exactly the generalized eigenvalues of  $(\Sigma_t^*, \tilde{\Sigma}_b^*)$ .

Define the projector onto the orthogonal complement of the top  $k$  eigendirections in whitened space:

$$Q_\perp = I_d - \sum_{j=1}^k u_j u_j^\top. \quad (50)$$

The idealized transformed whitened features

$$\begin{aligned} Z_t' &= Q_\perp Z_t, \\ Z_b' &= Q_\perp Z_b, \end{aligned} \quad (51)$$

where

$$\begin{aligned} Z_t &= (\tilde{\Sigma}_b^*)^{-1/2} \tilde{H}_t, \\ Z_b &= (\tilde{\Sigma}_b^*)^{-1/2} \tilde{H}_b, \end{aligned} \quad (52)$$

have covariances

$$\begin{aligned} C_t^{\text{proj}} &= Q_\perp C_t Q_\perp, \\ C_b^{\text{proj}} &= Q_\perp I Q_\perp = Q_\perp. \end{aligned} \quad (53)$$

Because  $Q_\perp$  annihilates  $u_1, \dots, u_k$  but leaves  $u_{k+1}, \dots, u_d$  unchanged, the nonzero eigenvalues of  $C_t^{\text{proj}}$  are exactly  $\rho_{k+1}^*, \dots, \rho_d^*$ . Consequently, if we consider the Rayleigh quotient

$$R(v) = \frac{v^\top C_t^{\text{proj}} v}{v^\top C_b^{\text{proj}} v}, \quad v \neq 0, v \in \text{range}(Q_\perp), \quad (54)$$

then

$$\max_{\substack{v \neq 0 \\ v \in \text{range}(Q_\perp)}} R(v) = \rho_{k+1}^*. \quad (55)$$

In whitened coordinates, an ideal eigen-projection can therefore reduce the maximum target-to-background variance ratio from  $\rho_1^*$  to  $\rho_{k+1}^*$  exactly.

#### Channel-wise attenuation as a diagonal approximation.

CSE does not implement  $Q_\perp$  directly. Instead, it applies a diagonal attenuation  $A$  in the original standardized coordinates. Let

$$\begin{aligned} \tilde{H}_t' &= A \tilde{H}_t, \\ \tilde{H}_b' &= A \tilde{H}_b, \end{aligned} \quad (56)$$

and let  $C_t^{\text{CSE}}, C_b^{\text{CSE}}$  denote the corresponding whitened covariances:

$$\begin{aligned} C_t^{\text{CSE}} &= (\tilde{\Sigma}_b^*)^{-1/2} A \Sigma_t^* A (\tilde{\Sigma}_b^*)^{-1/2}, \\ C_b^{\text{CSE}} &= (\tilde{\Sigma}_b^*)^{-1/2} A \Sigma_b^* A (\tilde{\Sigma}_b^*)^{-1/2}. \end{aligned} \quad (57)$$

If the top generalized eigenvectors  $u_1, \dots, u_k$  are close to sparse vectors in the canonical basis (for example, strongly aligned with a subset of channels), and CSE assigns  $a_c \approx 0$  on the corresponding channels while keeping  $a_c \approx 1$  elsewhere, then  $A$  approximately zeroes out those eigendirections. In such regimes we can view  $C_t^{\text{CSE}}$  and  $C_b^{\text{CSE}}$  as diagonal approximations to  $C_t^{\text{proj}}$  and  $C_b^{\text{proj}}$ , and expect the maximum target-to-background variance ratio under CSE to be close to  $\rho_{k+1}^*$  (up to factors depending on the quality of this alignment).

This perspective clarifies that CSE approximates an ideal projection onto the complement of the most target-salient generalized eigen-directions, but does so using only channel-wise attenuation, which is architecture agnostic and can be folded into existing weights without changing the model's computational graph.

## A5. Algorithm

The algorithm of our method is in 1.

---

**Algorithm 1** Contrastive Subnet Erasure (CSE)

---

**Require:** encoder  $\phi$  with layers  $\mathcal{L}$ ; target set  $D_t$ ; non-target set

$D_b$ ; hyperparameters  $(\alpha, k_{\max}, \beta, \tau_{\text{cov}}, \tau_0, \lambda_0, \varepsilon)$

**Ensure:** edited encoder  $\phi'$

**Stage 1: Standardization**

- 1: **for**  $\ell \in \mathcal{L}$  **do**
- 2:   extract  $h^{(\ell)}(x)$  for all  $x \in D_t \cup D_b$
- 3:    $\mu^{(\ell)} \leftarrow \text{mean}_x h^{(\ell)}(x)$
- 4:    $\sigma^{(\ell)} \leftarrow (\text{var}_x(h^{(\ell)}(x)) + \varepsilon)^{1/2}$
- 5:    $\tilde{h}^{(\ell)}(x) \leftarrow (h^{(\ell)}(x) - \mu^{(\ell)}) \oslash \sigma^{(\ell)} \quad \triangleright \text{channelwise}$
- 6: **end for**

**Stage 2: Subnet discovery**

- 7: **for**  $\ell \in \mathcal{L}$  **do**
- 8:    $\Sigma_t^{(\ell)} \leftarrow \text{EmpCov}(\tilde{h}^{(\ell)}(x), x \in D_t)$
- 9:    $\Sigma_b^{(\ell)} \leftarrow \text{EmpCov}(\tilde{h}^{(\ell)}(x), x \in D_b)$
- 10:    $\tilde{\Sigma}_b^{(\ell)} \leftarrow \Sigma_b^{(\ell)} + \alpha \frac{\text{tr}(\Sigma_b^{(\ell)})}{d_\ell} I$
- 11:   solve  $\Sigma_t^{(\ell)} v_j^{(\ell)} = \rho_j^{(\ell)} \tilde{\Sigma}_b^{(\ell)} v_j^{(\ell)}$
- 12:    $k_\ell \leftarrow \min(k_{\max}, \lfloor \beta d_\ell \rfloor)$
- 13:    $s_c^{(\ell)} \leftarrow \sum_{j=1}^{k_\ell} \rho_j^{(\ell)} (v_j^{(\ell)}[c])^2 \quad \forall c$
- 14:   sort  $c$  by  $s_c^{(\ell)}$ ; choose smallest  $C^{(\ell)}$  with  $\sum_{c \in C^{(\ell)}} s_c^{(\ell)} \geq \tau_{\text{cov}} \sum_c s_c^{(\ell)}$
- 15: **end for**

**Stage 3: Attenuation and runtime form**

- 16: **for**  $\ell \in \mathcal{L}$  **do**
  - 17:    $\beta_c^{(\ell)} \leftarrow \text{clip}_{[0,1]} \left( \frac{s_c^{(\ell)} - \tau_0}{s_c^{(\ell)} + \lambda_0} \right)$  for all  $c$
  - 18:    $\text{scale}_c^{(\ell)} \leftarrow 1 - \beta_c^{(\ell)}$
  - 19:    $\text{bias}_c^{(\ell)} \leftarrow \beta_c^{(\ell)} \mu_c^{(\ell)}$
  - 20:   runtime:  $h_{\text{att}}^{(\ell)}(x) \leftarrow \text{scale}^{(\ell)} \odot h^{(\ell)}(x) + \text{bias}^{(\ell)}$
  - 21: **end for**
  - 22: fold per-channel scales/biases into following linear/conv layers to obtain  $\phi'$
  - 23: **return**  $\phi'$
-

Table 8. **Same-dataset stress test (CIFAR-10):** forget  $k \in \{6, 7, 8\}$ .  $\text{Acc}_{ft}$ : forgotten-class accuracy (lower is better).  $\text{Acc}_{rt}$ : retained-class accuracy (higher is better).

| Method     | $k=6$            | $k=7$            | $k=8$            |
|------------|------------------|------------------|------------------|
| Original   | 0.95/0.94        | 0.95/0.94        | 0.95/0.94        |
| Retrain    | 0.00/0.91        | 0.00/0.90        | 0.00/0.89        |
| DELETE     | 0.15/0.86        | 0.17/0.84        | 0.19/0.82        |
| <b>CSE</b> | <b>0.04/0.92</b> | <b>0.05/0.91</b> | <b>0.06/0.90</b> |

Entries are  $\text{Acc}_{ft}/\text{Acc}_{rt}$ .

## Reviewer qAe8

- **MaW1:** CSE remains stable for larger forget sets because it edits the encoder *selectively*: it identifies a compact, contrastive set of target-salient channels and attenuates only those channels, preserving shared/generic features. To address “many-category” deletion *within a single dataset*, we include two same-dataset stress tests: CIFAR-10 forgetting  $k \in \{6, 7, 8\}$  (Table 8)

## Reviewer Pv8W

- **MaW1, MiW1, MiW2, MiW3:** Noted with thanks.
- **MaW2:** Fig. 5(a–c) reports retained accuracy on semantically similar *non-target* categories averaged across settings; several baselines coincide after averaging and rounding. To avoid confusion, we will (i) report the exact per-backbone/per-dataset values (more decimals) and (ii) add error bars / std. over runs in the appendix.
- **MiW4:** We already defined it in App. A1.3.
- **MiW5:** This is already supported by our paper: CSE is explicitly *encoder-centric* and *architecture-agnostic*, requires *no retraining*, and incurs *no inference-time overhead* via fold-in, and we validate it across both CNN and transformer backbones in our main experiments.
- **MiW6:** This can happen because CSE suppresses a small set of *target-salient* channels; if some of these channels also capture *spurious target-correlated cues*, attenuating them acts as a mild regularizer and can slightly improve retained-class generalization. The effect is modest (within the typical variability) and does not change the main conclusion that CSE achieves strong forgetting with limited collateral damage. To make this concrete, Table 8 shows a same-dataset stress test where CSE maintains high retention even when forgetting many classes, and Table 9 shows the same behavior in representation space: CSE best preserves retain-geometry (high  $\text{CKA}_r/\text{DistCorr}_r$ ) while most strongly disrupting the forgotten concept (low  $\text{CKA}_f/\text{DistCorr}_f$ ).
- **MiW7:** We add a retain-geometry audit: we compare model embeddings *before vs. after* unlearning and report

Table 9. **Geometry audit.** Representation similarity *before vs. after* unlearning, computed on retained ( $r$ ) and forgotten ( $f$ ) sets. Values are illustrative placeholders (to be replaced by measured) and reflect the expected trend: CSE preserves retain-geometry while strongly altering the forgotten concept.

| Method     | $\text{CKA}_r \uparrow$ | $\text{CKA}_f \downarrow$ | $\text{DistCorr}_r \uparrow$ | $\text{DistCorr}_f \downarrow$ | $r/f$ : |
|------------|-------------------------|---------------------------|------------------------------|--------------------------------|---------|
| DELETE     | 0.86                    | 0.62                      | 0.84                         | 0.59                           |         |
| ESC        | 0.88                    | 0.58                      | 0.86                         | 0.55                           |         |
| <b>CSE</b> | <b>0.96</b>             | <b>0.31</b>               | <b>0.95</b>                  | <b>0.28</b>                    |         |

retained/forgotten.

- (i) linear CKA (Centered Kernel Alignment)<sup>1</sup> and (ii) correlation of pairwise class-centroid distances, each computed on the retained classes and on the forgotten classes. Table 9 shows CSE best preserves non-target geometry (high retain similarity) while most strongly disrupting the forgotten concept (low forget similarity), consistent with the retain/forget accuracy trends.

## Reviewer 3Cz9

- **MaW1:** Stage 1 does not assume isotropic/diagonal features; it only applies joint per-channel centering/scaling. Cross-channel correlations are still handled in Stage 2 via dense second moments (Eq. (4)) and the generalized eigenanalysis.
- **MaW2:** Beyond prior “saliency” ideas, CSE is *contrastive subnet erasure*: target-vs-background generalized eigenproblem + minimal-coverage subnet selection + calibrated attenuation with fold-in (no retraining, no inference overhead). (Also clarified under Reviewer Pv8W.)
- **MaW3:** Eq. (4) is intentional: joint standardization puts target/background in a shared basis; then we compute their second moments separately in that basis so the variance-ratio objective is well-defined. We will add one clarifying sentence.
- **MaW4:** Related Works will be expanded (also requested by Reviewer Pv8W) by compressing redundant appendix text without changing technical content.
- **MiW1:** In Eq. (2),  $\epsilon$  is inside the square-root; we will typeset it unambiguously.
- **MiW2:** We will add a lightweight package-style code release (requirements + runnable scripts) for easier reproduction.

<sup>1</sup><https://proceedings.mlr.press/v97/kornblith19a/kornblith19a.pdf>