
Hidden Poison: Machine Unlearning Enables Camouflaged Poisoning Attacks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce *camouflaged data poisoning attacks*, a new attack vector that arises
2 in the context of machine unlearning and other settings when model retraining may
3 be induced. An adversary first adds a few carefully crafted points to the training
4 dataset such that the impact on the model’s predictions is minimal. The adversary
5 subsequently triggers a request to remove a subset of the introduced points at which
6 point the attack is unleashed and the model’s predictions are negatively affected. In
7 particular, we consider clean-label *targeted* attacks (in which the goal is to cause
8 the model to misclassify a specific test point) on datasets including CIFAR-10,
9 Imagenette, and Imagenette. This attack is realized by constructing *camouflage*
10 datapoints that mask the effect of a poisoned dataset.

11 1 Introduction

12 Machine Learning (ML) research traditionally assumes a static pipeline: data is gathered, a model is
13 trained once and subsequently deployed. This paradigm has been challenged by practical deployments,
14 which are more dynamic in nature. After initial deployment more data may be collected, necessitating
15 additional training. Or, as in the *machine unlearning* setting [Cao and Yang, 2015], we may need to
16 produce a model as if certain points were never in the training set to begin with.¹

17 While such dynamic settings clearly increase the applicability of ML models, they also make them
18 more vulnerable. Specifically, they open models up to new methods of attack by malicious actors
19 aiming to sabotage the model. In this work, we introduce a new type of data poisoning attack on
20 models that *unlearn* training datapoints. We call these *camouflaged data poisoning attacks*.

21 The attack takes place in two phases. In the first stage, before the model is trained, the attacker adds
22 a set of carefully designed points to the training data, consisting of a *poison* set and a *camouflage*
23 set. The model’s behaviour should be similar whether it is trained on either the training data, or its
24 augmentation with both the poison and camouflage sets. In the second phase, after the model is
25 trained, the attacker triggers an unlearning request to delete the *camouflage* set. That is, the model
26 must be updated to behave as though it were only trained on the training set plus the poison set. At
27 this point, the attack is fully realized, and the model’s performance suffers in some way.

28 While such an attack could harm the model by several metrics, in this paper, we focus on *targeted*
29 poisoning attacks – that is, poisoning attacks where the goal is to misclassify one particular point in
30 the training set. Our contributions are the following:

- 31 1. We introduce *camouflaged data poisoning attacks*, demonstrating a new attack vector in dynamic
32 settings including *machine unlearning*.

¹A naive solution is to remove said points from the training set and re-train the model from scratch.

- 33 2. We realize these attacks in the targeted poisoning setting, giving an algorithm based on the
34 gradient-matching approach of Geiping et al. [2021]. In order to make the model behavior
35 comparable to as if the poison set were absent, we construct the camouflage set by generating
36 a new set of points that *undoes* the impact of the poison set, an idea which may be of broader
37 interest to the data poisoning community.
- 38 3. We demonstrate the efficacy of these attacks on a variety of models (SVMs and neural net-
39 works) and datasets (CIFAR-10 [Krizhevsky, 2009], Imagenette [Howard, 2019], and Image-
40 woof [Howard, 2019]).

41 1.1 Preliminaries

42 **Machine Unlearning.** A significant amount of legislation concerning the “right to be forgotten”
43 has recently been introduced by governments around the world, including the European Union’s
44 [General Data Protection Regulation](#) (GDPR), the California Consumer Privacy Act (CCPA), and
45 Canada’s proposed Consumer Privacy Protection Act (CPPA). Such legislation requires organizations
46 to delete information they have collected about a user upon request. A natural question is whether that
47 further obligates the organizations to remove that information from downstream machine learning
48 models trained on the data – current guidances [Information Commissioner’s Office, 2020] and
49 precedents [Federal Trade Commission, 2021] indicate that this may be the case. This goal has
50 sparked a recent line of work on *machine unlearning* [Cao and Yang, 2015].

51 The simplest way to remove a user’s data from a trained model is to remove the data from the training
52 set, and then retrain the model on the remainder (also called “retraining from scratch”). This is the
53 ideal way to perform data deletion, as it ensures that the model was never trained on the datapoint of
54 concern. The downside is that retraining may take a significant amount of time in modern machine
55 learning settings. Hence, most work within machine unlearning has studied *fast* methods for data
56 deletion, sometimes relaxing to *approximately* removing the datapoint. A related line of work has
57 focused more on other implications of machine unlearning, particularly the consequences of an
58 adaptive and dynamic data pipeline [Gupta et al., 2021, Marchant et al., 2022]. Our work fits into the
59 latter line: we show that the potential to remove points from a trained model can expose a new attack
60 vector. Since retraining from scratch is the ideal result that other methods try to emulate, we focus on
61 unlearning by retraining from scratch, but the same phenomena should still occur when any effective
62 machine unlearning algorithm is applied.

63 **Data Poisoning.** In a data poisoning attack, an adversary in some way modifies the training data
64 provided to a machine learning model, such that the model’s behaviour at test time is negatively
65 impacted. Our focus is on *targeted data poisoning attacks*, where the attacker’s goal is to cause the
66 model to misclassify some specific datapoint in the test set. Other common types of data poisoning
67 attacks include *indiscriminate* (in which the goal is to increase the test error) and *backdoor* (where
68 the goal is to misclassify test points which have been adversarially modified in some small way).

69 The adversary is typically limited in a couple ways. First, it is common to say that they can only *add* a
70 *small number* of points to the training set. This mimics the setting where the training data is gathered
71 from some large public crowdsourced dataset, and an adversary can contribute a few judiciously
72 selected points of their own. Other choices may include allowing them to *modify* or *delete* points
73 from the training set, but these are less easily motivated. Additionally, the adversary is generally
74 constrained to *clean-label attacks*: if the introduced points were inspected by a human, they should
75 not appear suspicious or incorrectly labeled. We comment that this criteria is subjective and thus not
76 a precise notion, but is nonetheless common in the data poisoning literature, and we use the term as
77 well.

78 A detailed discussion of the related works is deferred to [Appendix A](#).

79 2 Camouflaged poisoning attacks via unlearning

80 In this section, we describe various components of the camouflaged poisoning attack, and how it can
81 be realized using machine unlearning.

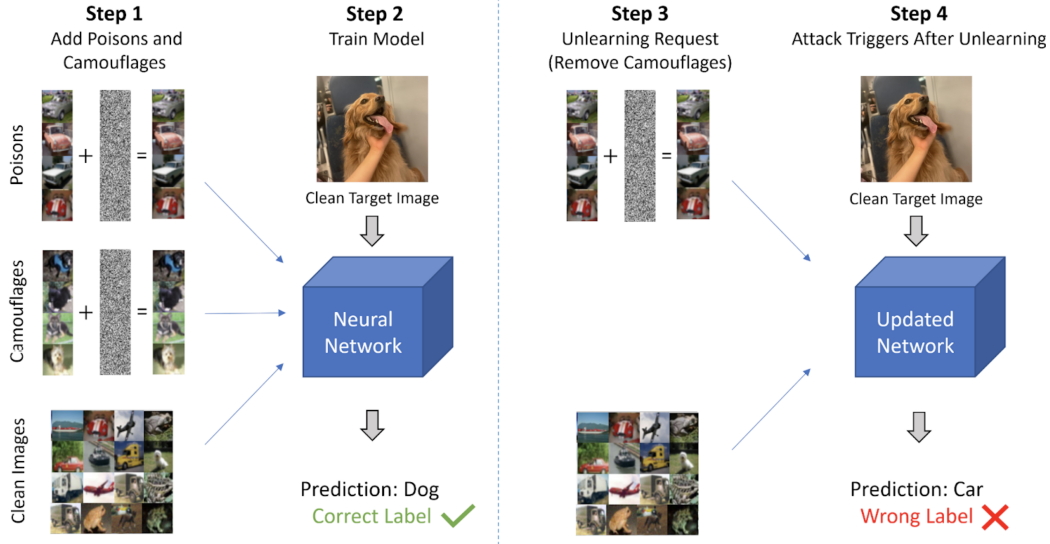


Figure 1: An illustration of a successful camouflaged targeted data poisoning attack. In Step 1, the adversary adds poison and camouflage sets of points to the (clean) training data. In Step 2, the model is trained on the augmented training dataset. It should behave similarly to if trained on only the clean data; in particular, it should correctly classify the targeted point. In Step 3, the adversary triggers an unlearning request to delete the camouflage set from the trained model. In Step 4, the resulting model misclassifies the targeted point.

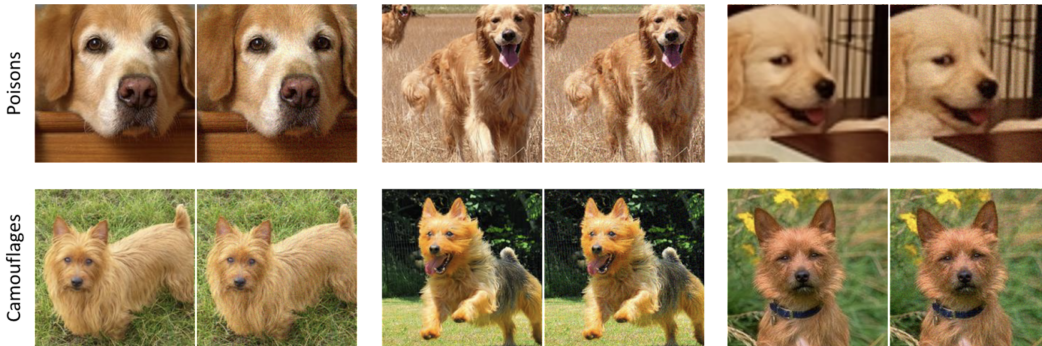


Figure 2: Some representative images from Imagewoof. In each pair, the left figure is from the training dataset, while the right image has been adversarially manipulated. The top and bottom rows are images from the poison and camouflage set, respectively. In all cases, the manipulated images are *clean label* and nearly indistinguishable from the original image.

82 2.1 Threat model and approach

83 The camouflaged poisoning attack takes place through interaction between an *attacker* and a *victim*.
 84 We assume that the attacker has access to the victim’s model architecture,² the ability to query
 85 gradients on a trained model (which could be achieved, e.g., by having access to the training dataset),
 86 and a target sample that it wants to attack. The attacker first sets the stage for the attack by introducing
 87 *poison points* and *camouflage points* to the training dataset, which are designed so as to have minimal
 88 impact when a model is trained with this modified dataset. At a later time, the attacker triggers the
 89 attack by submitting an unlearning request to remove the camouflage points. The victim first trains

²In [Appendix D.5.1](#), we examine the *transferability* of our proposed attack to unknown victim model, thus relaxing the requirement of knowing the victim’s model architecture a priori.

90 a machine learning model (e.g., a deep neural network) on the modified training dataset, and then
 91 executes the unlearning request by retraining the model from scratch on the left over dataset. The goal
 92 of the attacker is to change the prediction of the model on a particular target sample $(x_{\text{target}}, y_{\text{target}})$
 93 previously unseen by the model during training from y_{target} to a desired label $y_{\text{adversarial}}$, while still
 94 ensuring good performance over other validation samples. Formally, the interaction between the
 95 attacker and the victim is as follows (see Figure 1) :

- 96 1. The attacker introduces a small number of poisons samples S_{po} and camouflage samples S_{ca} to
 97 a clean training dataset S_{cl} . Define $S_{\text{cpc}} = S_{\text{cl}} + S_{\text{po}} + S_{\text{ca}}$.
- 98 2. Victim trains an ML model (e.g., a neural network) on S_{cpc} , and returns the model θ_{cpc} .
- 99 3. The attacker submits a request to unlearn the camouflage samples S_{ca} .
- 100 4. The victim performs the request, and computes a new model θ_{cp} by retraining from scratch on
 101 the left over data samples $S_{\text{cp}} = S_{\text{cl}} + S_{\text{po}}$.

102 Note that the attack is only realized in Step 4 when the victim executes the unlearning request and
 103 retrain the model from scratch on the left over training samples. In fact, in Steps 1-3, the victim’s
 104 model should behave similarly to as if it were trained on the clean samples S_{cl} only. In particular, the
 105 model θ_{cpc} will predict y_{target} on x_{target} , whereas the updated model θ_{cp} will predict $y_{\text{adversarial}}$ on
 106 x_{target} . Both models should have comparable validation accuracy. Such an attack is implemented by
 107 designing a camouflage set that cancels the effects of the poison set while training, but retraining
 108 without the camouflage set exposes the poison set, thus negatively affecting the model.

109 We highlight that camouflaged attacks may be *more dangerous* than traditional data poisoning attacks,
 110 since camouflaged attacks can be triggered by the adversary. That is, the adversary can reveal the
 111 attack whenever the submit an unlearning request, whereas for a traditional poisoning attack, the
 112 adversary simply plants the attack and must wait for the victim to train the model.

113 In order to be undetectable, and represent the realistic scenario in which the adversary has limited
 114 influence on the model’s training data, the attacker is only allowed to introduce a set of points that
 115 is much smaller than the size of the clean training dataset (i.e., $|S_{\text{po}}| \ll |S_{\text{cl}}|$ and $|S_{\text{ca}}| \ll |S_{\text{cl}}|$).
 116 Throughout the paper and experiments, we denote the relative size of the poison set and camouflage
 117 set by $b_p := \frac{|S_{\text{po}}|}{|S_{\text{cl}}|} \times 100$ and $b_c := \frac{|S_{\text{ca}}|}{|S_{\text{cl}}|} \times 100$, respectively. Additionally, the attacker is only allowed
 118 to generate poison and camouflage points by altering the base images by less than ε distance in the
 119 ℓ_∞ norm (in our experiments $\varepsilon \leq 16$, where the images are represented as an array of pixels in 0 to
 120 255). Thus, the attacker executes a so-called *clean-label* attack, where the corrupted images would
 121 be visually indistinguishable from original base images and thus would be given the same label as
 122 before by a human data validator. We parameterize a threat model by the tuple (ε, b_c, b_p) .

123 The attacker implements the attack by first generating poison samples, and then generating camouflage
 124 samples to cancel their effects. The poison and camouflage points are generated with the following
 125 goal in mind:

126 **Poison points.** Poison points are designed so that a network trained on $S_{\text{cp}} = S_{\text{cl}} + S_{\text{po}}$ predicts
 127 the label $y_{\text{adversarial}}$ (instead of y_{target}) on a target image x_{target} . While there are numerous data
 128 poisoning attacks in the literature, we adopt the state-of-the-art procedure of Geiping et al. [2021] for
 129 generating poisons due to its high success rate, efficiency of implementation, and applicability across
 130 various models. However, our framework is flexible: in principle, other attacks for the same setting
 131 could serve as a drop-in replacement, e.g., the methods of Aghakhani et al. [2021] or Huang et al.
 132 [2020], or any method introduced in the future. Suppose that S_{cp} consist of N_1 samples $(x^i, y^i)_{i \leq N_1}$
 133 out of which the first P samples with index $i = 1$ to P belong to the poison set S_{po} .³ The poison
 134 samples are generated by adding small perturbations Δ^i to the base image x^i so as to minimize the
 135 loss on the target with respect to the adversarial label, which can be formalized as the following
 136 bilevel optimization problem⁴

$$\min_{\Delta \in \Gamma} \ell(f(x_{\text{target}}, \theta(\Delta)), y_{\text{adversarial}}) \quad \text{where} \quad \theta(\Delta) \in \arg \min_{\theta} \frac{1}{N} \sum_{i \leq N} \ell(f(x^i + \Delta^i, \theta), y^i), \quad (1)$$

³This ordering is for notational convenience; naturally, the datapoints are shuffled to preclude the victim simply removing a prefix of the training data.

⁴While (1) focuses on misclassifying a single target point, it is straightforward to extend this to multiple targets by changing the objective to a sum over losses on the target points.

137 where we define the constraint set $\Gamma := \{\Delta : \|\Delta\|_\infty \leq \varepsilon \text{ and } \Delta^i = 0 \text{ for all } i > P\}$. The main optimization
 138 objective in (1) is called the adversarial loss [Geiping et al., 2021].

139 **Camouflage points.** Camouflage samples are designed to cancel the effect of the poisons, such
 140 that a model trained on $S_{\text{cpc}} = S_{\text{cl}} + S_{\text{po}} + S_{\text{ca}}$ behaves identical to the model trained on S_{cl} , and
 141 makes the correct prediction on x_{target} . We formulate this task via a bilevel optimization problem
 142 similar to what we did in (1) for generating poisons. Let S_{cpc} consist of N_2 samples $(x^j, y^j)_{j \leq N_2}$
 143 out of which the last C samples with index $j = N_2 - C + 1$ to N_2 belong to the camouflage set S_{ca} .
 144 The camouflage points are generated by adding small perturbations Δ^j to the base image x^j so as
 145 to minimize the loss on the target with respect to the adversarial label. In particular, we find the
 146 appropriate Δ by solving:

$$\min_{\Delta \in \Gamma} \ell(f(x_{\text{target}}, \theta(\Delta)), y_{\text{target}}) \quad \text{where} \quad \theta(\Delta) \in \arg \min_{\theta} \frac{1}{N_2} \sum_{j \leq N_2} \ell(f(x^j + \Delta^j, \theta), y^j), \quad (2)$$

147 where we define the constraint set $\Gamma := \{\Delta : \|\Delta\|_\infty \leq \varepsilon \text{ and } \Delta^j = 0 \text{ for all } j \leq N_2 - C\}$.

148 The exact procedure to generate camouflages and poisons is given in Appendix B. We build on the
 149 gradient matching procedure in [Geiping et al., 2021] for implementing (1) and (2) efficiently in order
 150 to generate *clean-label* camouflages and poisons for large-scale machine learning settings.

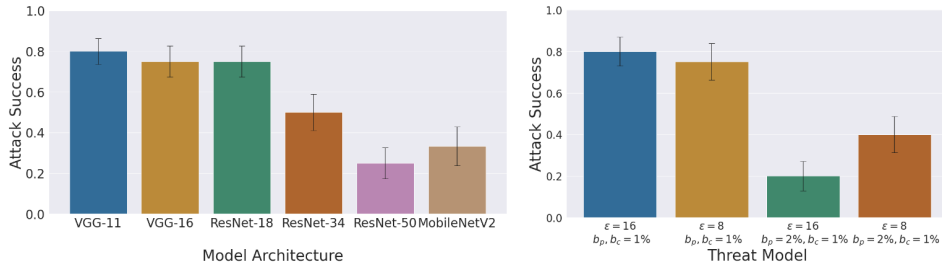


Figure 3: Efficacy of the proposed camouflaged poisoning attack on CIFAR-10 dataset. The left plot gives the success for the threat model $\varepsilon = 16, b_p = 0.6\%, b_c = 0.6\%$ across different neural network architectures. The right plot gives the success for ResNet-18 architecture across different threat models. See Appendix D.3 for more details about this experiment.

151 3 Experimental evaluation

152 We extensively evaluate the efficacy of camouflaged poisoning attack on various large scale ML datasets
 153 including CIFAR-10, and Imagenette and Imagenette (subsets of 10 classes from Imagenet). We
 154 perform evaluations on VGG-11, VGG-16, Resnet-18, Resnet-24, Resnet-50 and MobileNetV2 for
 155 CIFAR10, and VGG-16 and ResNet-18 for Imagenette and Imagenette respectively. All experimental
 156 details, obtained results, and visualizations of the generated poisons and camouflages can be found in
 157 Appendix D.

158 In Appendix D.5, we also provide additional experiments on CIFAR-10 showing that our attack is
 159 robust to data augmentation, and successfully transfers when the victim model is different from the
 160 model on which poison and camouflage samples were generated.

Attack type (ε, b_p, b_c)	Attack success		Validation Accuracy		
	Poisoning	Camouflaging	Clean	Poisoned	Camouflaged
LF (8, 0.2%, 0.2%)	70%	71.5%	81.63	81.73 (± 0.14)	81.74 (± 0.20)
LF (16, 0.2%, 0.2%)	100%	40%	81.63	81.64 (± 0.03)	81.6 (± 0.02)
GM (8, 0.2%, 0.4%)	70%	100%	81.63	81.65 (± 0.01)	81.62 (± 0.02)
GM (16, 0.2%, 0.4%)	100%	70%	81.63	81.65 (± 0.03)	81.63 (± 0.02)

Table 1: Camouflaged poisoning attack on linear SVM on Binary-CIFAR-10 dataset. The first column lists the threat model (ε, b_p, b_c) and the camouflaging type “LF” for label flipping and “GM” for gradient matching. See Appendix D.1.1 for more details on these procedures.

161 **4 Conclusion and discussion.**

162 We demonstrated a new attack vector, *camouflaged poisoning attacks*, against machine learning
 163 pipelines where training points can be *unlearned*. This shows that as we introduce new functionality to
 164 machine learning systems, we must be aware of novel threats that emerge. We outline a few interesting
 165 directions for further research: It is important to understand how to *defend* against camouflaged
 166 attacks. As observed by Geiping et al. [2021], it is unlikely that differential privacy [Dwork et al.,
 167 2006] would be an effective defense, as preventing attacks in the non-camouflaged setting incurs too
 168 significant a loss in accuracy. Another direction is to reduce the knowledge needed by the adversary,
 169 thereby creating stronger attacks. E.g., while our setting requires grey-box knowledge, one could
 170 instead consider a black-box model to attack ML APIs. Finally, it is interesting to determine what
 171 other types of threats can be camouflaged, e.g., indiscriminate or backdoor poisoning attacks. Beyond
 172 exploring this new attack vector, it is also independently interesting to understand how one can
 173 neutralize the effect of an attack by *adding* points.

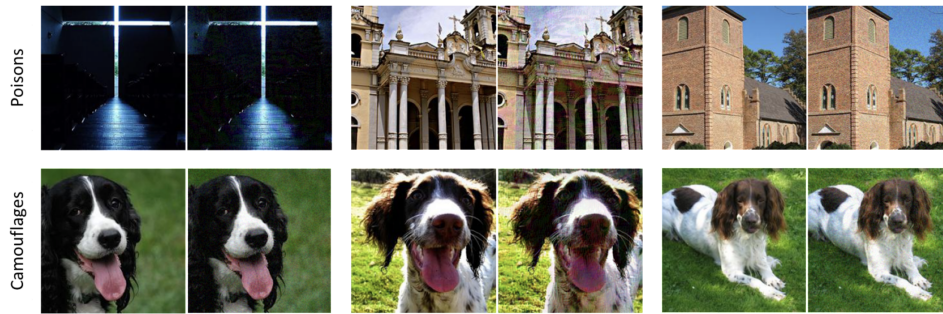


Figure 4: Some representative poison and camouflage images for attack on Imgewoof dataset. In each pair, the left figure is the original picture from the training dataset and the right figure has been adversarially manipulated by adding Δ . The shown images were generated for a camouflaged poisoning attack on Resnet-18, with Seed = 10000005, $b_p = b_c = 6.6\%$ and $\epsilon = 16$.

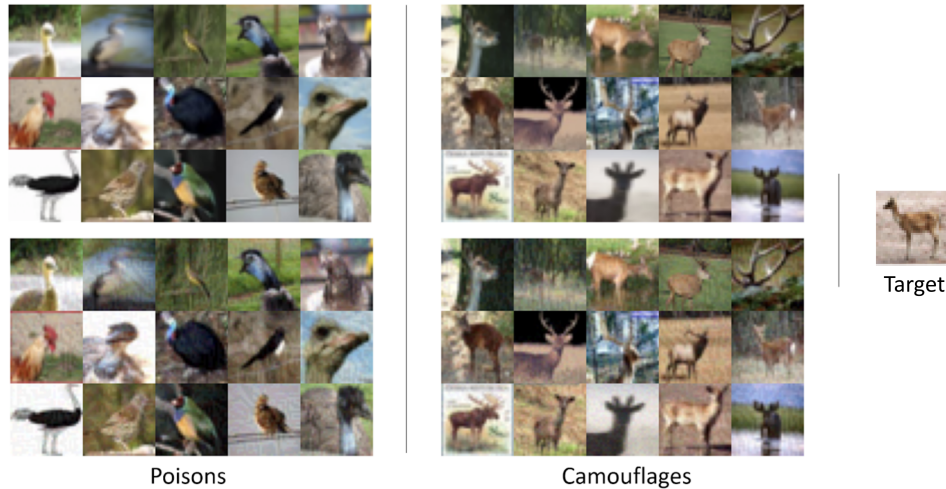


Figure 5: Visualization of poisons and camouflages on CIFAR-10 dataset. The top row shows the original images and the bottom row shows the corresponding poisoned / camouflaged images (with the added Δ). The shown images were generated for a camouflaged poisoning attack on ResNet-18, with Seed = 2000000000, $\epsilon = 8$, $b_p = 0.2$, $b_c = 0.4$, poison class *bird*, target class *deer*, and the target ID 9621.

174 References

- 175 Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In
176 *Proceedings of the 36th IEEE Symposium on Security and Privacy*, SP '15, pages 463–480,
177 Washington, DC, USA, 2015. IEEE Computer Society.
- 178 Jonas Geiping, Liam Fowl, W Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller,
179 and Tom Goldstein. Witches' brew: Industrial scale data poisoning via gradient matching. In
180 *Proceedings of the 9th International Conference on Learning Representations*, ICLR '21, 2021.
- 181 Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009.
- 182 Jeremy Howard. Imagenette, 2019. URL <https://github.com/fastai/imagenette/>.
- 183 General Data Protection Regulation. Regulation (EU) 2016/679 of the European parliament and of
184 the council of 27 April 2016, 2016.
- 185 Information Commissioner's Office. Guidance on the ai auditing framework, February 2020.
- 186 Federal Trade Commission. California company settles ftc allegations it deceived consumers about
187 use of facial recognition in photo storage app, January 2021.
- 188 Varun Gupta, Christopher Jung, Seth Neel, Aaron Roth, Saeed Sharifi-Malvajerdi, and Chris Waites.
189 Adaptive machine unlearning. In *Advances in Neural Information Processing Systems 34*, NeurIPS
190 '21, pages 16319–16330. Curran Associates, Inc., 2021.
- 191 Neil G Marchant, Benjamin IP Rubinstein, and Scott Alfeld. Hard to forget: Poisoning attacks on
192 certified machine unlearning. In *Proceedings of the Thirty-Sixth AAAI Conference on Artificial
193 Intelligence*, volume 36 of AAAI '22, pages 7691–7700, 2022.
- 194 Hojjat Aghakhani, Dongyu Meng, Yu-Xiang Wang, Christopher Kruegel, and Giovanni Vigna.
195 Bullseye polytope: A scalable clean-label poisoning attack with improved transferability. In *2021
196 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 159–178. IEEE, 2021.
- 197 W Ronny Huang, Jonas Geiping, Liam Fowl, Gavin Taylor, and Tom Goldstein. Metapoisson: Practical
198 general-purpose clean-label data poisoning. In *Advances in Neural Information Processing Systems
199 33*, NeurIPS '20, pages 12080–12091. Curran Associates, Inc., 2020.
- 200 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in
201 private data analysis. In *Proceedings of the 3rd Conference on Theory of Cryptography*, TCC '06,
202 pages 265–284, Berlin, Heidelberg, 2006. Springer.
- 203 Marco Barreno, Blaine Nelson, Anthony D Joseph, and J Doug Tygar. The security of machine
204 learning. *Machine Learning*, 81(2):121–148, 2010.
- 205 Han Xiao, Huang Xiao, and Claudia Eckert. Adversarial label flips attack on support vector machines.
206 In *ECAI 2012*, pages 870–875. IOS Press, 2012.
- 207 Andrea Paudice, Luis Muñoz-González, and Emil C Lupu. Label sanitization against label flipping
208 poisoning attacks. In *Joint European conference on machine learning and knowledge discovery in
209 databases*, pages 5–15. Springer, 2018.
- 210 Octavian Suciú, Radu Marginean, Yigitcan Kaya, Hal Daume III, and Tudor Dumitras. When does
211 machine learning {FAIL}? generalized transferability for evasion and poisoning attacks. In *27th
212 USENIX Security Symposium (USENIX Security 18)*, pages 1299–1316, 2018.
- 213 Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciú, Christoph Studer, Tudor Dumitras,
214 and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks.
215 *Advances in neural information processing systems*, 31, 2018.
- 216 Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines.
217 In *Proceedings of the 29th International Conference on Machine Learning*, ICML '12, pages
218 1467–1474. JMLR, Inc., 2012.

- 219 Huang Xiao, Battista Biggio, Blaine Nelson, Han Xiao, Claudia Eckert, and Fabio Roli. Support
220 vector machines under adversarial label contamination. *Neurocomputing*, 160:53–62, 2015.
- 221 Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee,
222 Emil C Lupu, and Fabio Roli. Towards poisoning of deep learning algorithms with back-gradient
223 optimization. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*,
224 pages 27–38, 2017.
- 225 Jacob Steinhardt, Pang Wei W Koh, and Percy S Liang. Certified defenses for data poisoning attacks.
226 In *Advances in Neural Information Processing Systems 30*, NeurIPS ’17, pages 3520–3532. Curran
227 Associates, Inc., 2017.
- 228 Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Jacob Steinhardt, and Alistair
229 Stewart. Sever: A robust meta-algorithm for stochastic optimization. In *Proceedings of the 36th
230 International Conference on Machine Learning*, ICML ’19, pages 1596–1606. JMLR, Inc., 2019.
- 231 Pang Wei Koh, Jacob Steinhardt, and Percy Liang. Stronger data poisoning attacks break data
232 sanitization defenses. *Machine Learning*, 111(1):1–47, 2022.
- 233 Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the
234 machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- 235 Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In *Advances
236 in Neural Information Processing Systems 31*, NeurIPS ’18, pages 8011–8021. Curran Associates,
237 Inc., 2018.
- 238 Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H Brendan McMahan. Can you really
239 backdoor federated learning? *arXiv preprint arXiv:1911.07963*, 2019.
- 240 Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Song,
241 Aleksander Madry, Bo Li, and Tom Goldstein. Dataset security for machine learning: Data
242 poisoning, backdoor attacks, and defenses. *arXiv preprint arXiv:2012.10544*, 2020.
- 243 Antonio Emanuele Cinà, Kathrin Grosse, Ambra Demontis, Sebastiano Vascon, Werner Zellinger,
244 Bernhard A Moser, Alina Oprea, Battista Biggio, Marcello Pelillo, and Fabio Roli. Wild patterns
245 reloaded: A survey of machine learning security against training data poisoning. *arXiv preprint
246 arXiv:2205.01992*, 2022.
- 247 Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. Making AI forget you: Data
248 deletion in machine learning. In *Advances in Neural Information Processing Systems 32*, NeurIPS
249 ’19, pages 3518–3531. Curran Associates, Inc., 2019.
- 250 Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens Van Der Maaten. Certified data removal
251 from machine learning models. In *Proceedings of the 37th International Conference on Machine
252 Learning*, ICML ’20, pages 3832–3842. JMLR, Inc., 2020.
- 253 Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. Approximate data deletion
254 from machine learning models: Algorithms and evaluation. In *Proceedings of the 24th International
255 Conference on Artificial Intelligence and Statistics*, AISTATS ’21, pages 2008–2016. JMLR, Inc.,
256 2021.
- 257 Seth Neel, Aaron Roth, and Saeed Sharifi-Malvajerdi. Descent-to-delete: Gradient-based methods
258 for machine unlearning. In *Proceedings of the 32nd International Conference on Algorithmic
259 Learning Theory*, ALT ’21. JMLR, Inc., 2021.
- 260 Enayat Ullah, Tung Mai, Anup Rao, Ryan A Rossi, and Raman Arora. Machine unlearning via
261 algorithmic stability. In *Conference on Learning Theory*, pages 4126–4142. PMLR, 2021.
- 262 Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what
263 you want to forget: Algorithms for machine unlearning. In *Advances in Neural Information
264 Processing Systems 34*, NeurIPS ’21, pages 18075–18086. Curran Associates, Inc., 2021.

- 265 Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers,
266 Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *proceedings of the 42nd*
267 *IEEE Symposium on Security and Privacy*, SP '21, Washington, DC, USA, 2021. IEEE Computer
268 Society.
- 269 Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selec-
270 tive forgetting in deep networks. In *Proceedings of the 2020 IEEE Computer Society Conference*
271 *on Computer Vision and Pattern Recognition*, CVPR '20, pages 9304–9312, Washington, DC,
272 USA, 2020. IEEE Computer Society.
- 273 Aditya Golatkar, Alessandro Achille, Avinash Ravichandran, Marzia Polito, and Stefano Soatto.
274 Mixed-privacy forgetting in deep networks. In *Proceedings of the 2021 IEEE Computer Society*
275 *Conference on Computer Vision and Pattern Recognition*, CVPR '21, pages 792–801. IEEE
276 Computer Society, 2021.
- 277 Quoc Phong Nguyen, Bryan Kian Hsiang Low, and Patrick Jaillet. Variational bayesian unlearning.
278 *Advances in Neural Information Processing Systems*, 33:16025–16036, 2020.
- 279 Jonathan Brophy and Daniel Lowd. Machine unlearning for random forests. In *Proceedings of the*
280 *38th International Conference on Machine Learning*, ICML '21, pages 1092–1104. JMLR, Inc.,
281 2021.
- 282 Santiago Zanella-Béguelin, Lukas Wutschitz, Shruti Tople, Victor Rühle, Andrew Paverd, Olga
283 Ohrimenko, Boris Köpf, and Marc Brockschmidt. Analyzing information leakage of updates to
284 natural language models. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and*
285 *Communications Security*, pages 363–375, 2020.
- 286 Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. When
287 machine unlearning jeopardizes privacy. In *Proceedings of the 2021 ACM SIGSAC Conference on*
288 *Computer and Communications Security*, pages 896–911, 2021.
- 289 Min Du, Zhi Chen, Chang Liu, Rajvardhan Oak, and Dawn Song. Lifelong anomaly detection
290 through unlearning. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and*
291 *Communications Security*, pages 1283–1297, 2019.
- 292 David Marco Sommer, Liwei Song, Sameer Wagh, and Prateek Mittal. Towards probabilistic
293 verification of machine unlearning. *arXiv preprint arXiv:2003.04247*, 2020.
- 294 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of*
295 *the 3rd International Conference on Learning Representations*, ICLR '15, 2015.
- 296 Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier
297 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas,
298 Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay.
299 Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12(85):2825–
300 2830, 2011.
- 301 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale visual
302 recognition. In *Proceedings of the 3rd International Conference on Learning Representations*,
303 ICLR '15, 2015.
- 304 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
305 recognition. In *Proceedings of the 2016 IEEE Computer Society Conference on Computer Vision*
306 *and Pattern Recognition*, CVPR '16, pages 770–778, Washington, DC, USA, 2016. IEEE Computer
307 Society.
- 308 Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo-
309 bileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of the 2018 IEEE Computer*
310 *Society Conference on Computer Vision and Pattern Recognition*, CVPR '18, pages 4510–4520,
311 Washington, DC, USA, 2018. IEEE Computer Society.

- 312 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
313 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward
314 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,
315 Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning
316 library. In *Advances in Neural Information Processing Systems 32*, NeurIPS '19, pages 8026–8037.
317 Curran Associates, Inc., 2019.
- 318 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
319 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition
320 challenge. *International journal of computer vision*, 115(3):211–252, 2015.

321 A Related work

322 The motivation for our work comes from [Marchant et al. \[2022\]](#), who propose a novel poisoning
323 attack on unlearning systems. As mentioned before, the primary goal of many machine unlearning
324 systems is to “unlearn” datapoints quickly, i.e., faster than retraining from scratch. [Marchant et al.](#)
325 [\[2022\]](#) craft poisoning schemes via careful noise addition, in order to trigger the unlearning algorithm
326 to retrain from scratch on far more deletion requests than typically required. While both our work
327 and theirs are focused on data poisoning attacks against machine unlearning systems, the adversaries
328 have very different objectives. In our work, the adversary is trying to misclassify a target test point,
329 whereas in theirs, they try to increase the time required to unlearn a point.

330 In targeted data poisoning, there are a few different types of attacks. The simplest form of attack is
331 *label flipping*, in which the adversary is allowed to flip the labels of the examples [[Barreno et al.,](#)
332 [2010](#), [Xiao et al., 2012](#), [Paudice et al., 2018](#)]. Another type of attack is *watermarking*, in which
333 the feature vectors are perturbed to obtain the desired poisoning effect [[Suciu et al., 2018](#), [Shafahi](#)
334 [et al., 2018](#)]. In both these cases, noticeable changes are made to the label and feature vector,
335 respectively, which would be noticeable by a human labeler. In contrast, *clean label* attacks attempt
336 to make unnoticeable changes to both the feature vector and the label, and are the gold standard for
337 data poisoning attacks [[Huang et al., 2020](#), [Geiping et al., 2021](#)]. Our focus is on both clean-label
338 poisoning and camouflage sets. While there are also works on indiscriminate [[Biggio et al., 2012](#),
339 [Xiao et al., 2015](#), [Muñoz-González et al., 2017](#), [Steinhardt et al., 2017](#), [Diakonikolas et al., 2019](#),
340 [Koh et al., 2022](#)] and backdoor [[Gu et al., 2017](#), [Tran et al., 2018](#), [Sun et al., 2019](#)] poisoning attacks,
341 these are beyond the scope of our work, see [Golub et al. \[2020\]](#), [Cinà et al. \[2022\]](#) for additional
342 background on data poisoning attacks.

343 [Cao and Yang \[2015\]](#) initiated the study of machine unlearning through *exact* unlearning, wherein
344 the new model obtained after deleting an example is statistically identical to the model obtained by
345 training on a dataset without the example. A probabilistic notion of unlearning was defined by [Ginart](#)
346 [et al. \[2019\]](#), which in turn is inspired from notions in differential privacy [[Dwork et al., 2006](#)].
347 Several works studied algorithms for empirical risk minimization (i.e., training loss) [[Guo et al.,](#)
348 [2020](#), [Izzo et al., 2021](#), [Neel et al., 2021](#), [Ullah et al., 2021](#)], while later works study the effect of
349 machine unlearning on the generalization loss [[Gupta et al., 2021](#), [Sekhari et al., 2021](#)]. In particular,
350 these works realize that unlearning data points quickly can lead to a drop in test loss, which is the
351 theme of our current work. Several works have considered implementations of machine unlearning
352 in several contexts starting with the work of [Bourtole et al. \[2021\]](#). These include unlearning in
353 deep neural networks [[Golatkar et al., 2020, 2021](#), [Nguyen et al., 2020](#)], random forests [[Brophy](#)
354 [and Lowd, 2021](#)], large scale language models [[Zanella-Béguelin et al., 2020](#)], the tension between
355 unlearning and privacy [[Chen et al., 2021](#)], anomaly detection [[Du et al., 2019](#)], and even auditing of
356 machine unlearning systems [[Sommer et al., 2020](#)].

357 B Poison and camouflage generation

358 B.1 Gradient matching for efficient poison generation [[Geiping et al., 2021](#)]

359 In this section, we discuss the key intuition of [Geiping et al. \[2021\]](#) for efficient poison generation.
360 Our objective is to find perturbations Δ such that when the model is trained on the poisoned samples, it
361 minimizes the adversarial loss in (1) thus making the victim model predict the wrong label $y_{\text{adversarial}}$
362 on the target sample. However, directly solving (1) is computationally intractable due to bilevel
363 nature of the optimization objective. Instead, one may implicitly minimize the adversarial loss by
364 finding a Δ such that for any model parameter θ ,

$$\nabla_{\theta}(\ell(f(x_{\text{target}}, \theta), y_{\text{adversarial}})) \approx \frac{1}{P} \sum_{i=1}^P \nabla_{\theta} \ell(f(x^i + \Delta^i, \theta), y^i). \quad (3)$$

365 In essence, (3) implies that gradient based minimization (e.g., using Adam / SGD) of the training
366 loss on poisoned samples also minimizes the adversarial loss. Thus, training a model on $S_{\text{cl}} + S_{\text{po}}$
367 will automatically ensure that the model predicts $y_{\text{adversarial}}$ on the target sample. Unfortunately,
368 computing Δ that satisfies (3) is also intractable as it is required to hold for all values of θ . The key
369 idea of [Geiping et al. \[2021\]](#) to make poison generation efficient is to relax (3) to only be satisfied

370 for a fixed model θ_{cl} —the model obtained by training on the clean dataset S_{cl} . To implement this,
 371 Geiping et al. [2021] minimize the cosine-similarity loss between the two gradients defined as:

$$\phi(\Delta, \theta) = 1 - \frac{\langle \nabla_{\theta} \ell(f(x_{\text{target}}, \theta), y_{\text{adversarial}}), \sum_{i=1}^P \nabla_{\theta} \ell(f(x_i + \Delta_i, \theta), y_i) \rangle}{\|\nabla_{\theta} \ell(f(x_{\text{target}}, \theta), y_{\text{adversarial}})\| \|\sum_{i=1}^P \nabla_{\theta} \ell(f(x_i + \Delta_i, \theta), y_i)\|}, \quad (4)$$

372 Geiping et al. [2021] demonstrated that (4) can be efficiently optimized for many popular large-
 373 scale machine learning models and datasets. For completeness, we provide their pseudocode in
 374 Algorithm 1.

375 B.2 Camouflaging poisoned points

376 Camouflage images are designed in order to neutralize the effect of the poison images. In this section,
 377 we give intuition into what we mean by cancelling the effect of poisons, and provide two procedures
 378 for generating camouflages efficiently: label flipping, and gradient matching.

379 B.2.1 Camouflages via label flipping

380 Suppose that the underlying task is a binary classification problem with the labels $y \in \{-1, 1\}$, and
 381 that the model is trained using linear loss $\ell(f(x, \theta), y) = -yf(x, \theta)$. Then, simply flipping the labels
 382 allows one to generate a camouflage set for any given poison set S_{po} . In particular, S_{ca} is constructed
 383 as: for every $(x^i, y^i) \in S_{\text{po}}$, simply add $(x^i, -y^i)$ to S_{ca} (i.e., $b_p = b_c$). It is easy to see that for such
 384 camouflage points, we have for any θ ,

$$\sum_{(x, y) \in S_{\text{cpc}}} \ell(f(x, \theta), y) = - \sum_{(x, y) \in S_{\text{cl}}} yf(x, \theta) - \sum_{i=1}^P (y^i f(x^i, \theta) + (-y^i) f(x^i, \theta)) = \sum_{(x, y) \in S_{\text{cl}}} \ell(f(x, \theta), y).$$

385 We can also similarly show that the gradients (as well as higher order derivatives) are equal, i.e.,
 386 $\nabla_{\theta} \sum_{S_{\text{cpc}}} \ell(f(x, \theta), y) = \nabla_{\theta} \sum_{S_{\text{cl}}} \ell(f(x, \theta), y)$ for all θ . Thus, training a model on S_{cpc} is equivalent
 387 to training it on S_{cl} . In essence, the camouflages have perfectly canceled out the effect of the poisons.
 388 We validate the efficacy of this approach via experiments on linear SVM trained with hinge loss (which
 389 resembles linear loss when the domain is bounded) on a binary classification problem constructed
 390 using CIFAR-10 dataset. We report the results in Table 1 (see Section ?? for details).

391 While label flipping is a simple and effective procedure to generate camouflages, it is fairly restrictive.
 392 Firstly, label flipping only works for binary classification problems trained with linear loss. Secondly,
 393 the attack is not clean label as the camouflage images are generated as $(x^i, -y^i)$ by giving them the
 394 opposite label to the ground truth, which can be easily caught by a validator. Lastly, the attack is
 395 vulnerable to simple data purification techniques by the victim, e.g., the victim can protect themselves
 396 by preprocessing the data to remove all the images that have both the labels ($y = +1$ and $y = -1$) in
 397 the training dataset. In the next section, we provide a different procedure to generate clean-label
 398 camouflages for general losses and multi-class classification problems.

399 B.2.2 Gradient matching for generating camouflages

400 We next discuss our main procedure to generate camouflages, which is based on the gradient matching
 401 idea of Geiping et al. [2021]. Note that, our objective in (2) is to find Δ such that when the model is
 402 trained with the camouflages, it minimizes the original-target loss in (2) (with respect to the original
 403 label y_{target}) thus making the victim model predict the correct label on this target sample. Since, (1)
 404 is computationally intractable, one may instead try to implicitly minimize the original-target loss by
 405 finding a Δ such that for any model parameter θ ,

$$\nabla_{\theta} (\ell(f(x_{\text{target}}, \theta), y_{\text{target}})) \approx \frac{1}{C} \sum_{i=1}^C \nabla_{\theta} \ell(f(x^i + \Delta^i, \theta), y^i). \quad (5)$$

406 (5) suggests that minimizing (e.g., using Adam / SGD) on camouflage samples will also minimize
 407 the original-target loss, and thus automatically ensure that the model predicts the correct label on
 408 the target sample. Unfortunately, (5) is also intractable as it requires the condition to hold for all θ .
 409 Building on the work of Geiping et al. [2021], we relax this condition to satisfied only for a fixed

410 model θ_{cp} -the model trained on the dataset $S_{cp} = S_{cl} + S_{po}$. Similar to what we did for generating
 411 poison points, we achieve this by minimizing the cosine-similarity loss given by

$$\psi(\Delta, \theta) = 1 - \frac{\langle \nabla_{\theta} \ell(f(x_{target}, \theta), y_{target}), \sum_{i=1}^C \nabla_{\theta} \ell(f(x_i + \Delta_i, \theta), y_i) \rangle}{\| \nabla_{\theta} \ell(f(x_{target}, \theta), y_{target}) \| \| \sum_{i=1}^C \nabla_{\theta} \ell(f(x_i + \Delta_i, \theta), y_i) \|}. \quad (6)$$

412 **Implementation details.** We minimize (6) using the Adam optimizer [Kingma and Ba, 2015]
 413 with a fixed step size of 0.1. In order to increase the robustness of camouflage generation, we do
 414 R restarts (where $R \leq 10$). In each restart, we first initialize Δ randomly such that $\|\Delta\|_{\infty} \leq \varepsilon$ and
 415 perform M steps of Adam optimization to minimize $\psi(\Delta, \theta_{cp})$. Each optimization step only requires
 416 a single differentiation of the objective ψ with respect to Δ , and can be implemented efficiently. After
 417 each step, we project back the updated Δ into the constraint set Γ so as to maintain the property
 418 that $\|\Delta\|_{\infty} \leq \varepsilon$. After doing R restarts, we choose the best round by finding Δ_* with the minimum
 419 $\psi(\Delta_*, \theta_{cp})$. We provide the pseudocode in Algorithm 2.

Algorithm 1 Gradient Matching to generate poisons [Geiping et al., 2021]

Require: Clean network $f(\cdot; \theta_{clean})$ trained on uncorrupted base images S_{cl} , a target (x_{target}, y_{target}) and an adversarial label $y_{adversarial}$, Poison budget P , perturbation bound ε , number of restarts R , optimization steps M

- 1: Collect a dataset $S_{po} = \{x^i, y^i\}_{i=1}^P$ of P many images whose true label is $y_{adversarial}$.
 - 2: **for** $r = 1, \dots, R$ restarts **do**
 - 3: Randomly initialize perturbations Δ s.t. $\|\Delta\|_{\infty} \leq \varepsilon$.
 - 4: **for** $k = 1, \dots, M$ optimization steps **do**
 - 5: Compute the loss $\phi(\Delta, \theta_{clean})$ as in (4) using the base poison images in S_{po} .
 - 6: Update Δ using an Adam update to minimize ϕ , and project onto the constraint set Γ .
 - 7: **end for**
 - 8: Amongst the R restarts, choose the Δ_* with the smallest value of $\phi(\Delta_*, \theta_{clean})$.
 - 9: **end for**
 - 10: Return the poisoned set $S_{po} = \{x^i + \Delta_*^i, y^i\}_{i=1}^P$.
-

Algorithm 2 Gradient Matching to generate camouflages

Require: Network $f(\cdot; \theta_{cp})$ trained on $S_{cl} + S_{po}$, the target (x_{target}, y_{target}) , Camouflage budget C , perturbation bound ε , number of restarts R , optimization steps M

- 1: Collect a dataset $S_{ca} = \{x^j, y^j\}_{j=1}^C$ of C many images whose true label is y_{target} .
 - 2: **for** $r = 1, \dots, R$ restarts **do**
 - 3: Randomly initialize perturbations Δ s.t. $\|\Delta\|_{\infty} \leq \varepsilon$.
 - 4: **for** $k = 1, \dots, M$ optimization steps **do**
 - 5: Compute the loss $\psi(\Delta, \theta_{cp})$ as in (4) using the base camouflage images in S_{ca} .
 - 6: Update Δ using an Adam update to minimize ψ , and project onto the constraint set Γ .
 - 7: **end for**
 - 8: Amongst the R restarts, choose the Δ_* with the smallest value of $\psi(\Delta_*, \theta_{cp})$.
 - 9: **end for**
 - 10: Return the poisoned set $S_{ca} = \{x^j + \Delta_*^j, y^j\}_{j=1}^C$.
-

420 C Experiment details

421 C.1 Hardware

422 All our experiments were executed on Google Colab with a Google Colab Pro+ subscription.

423 C.2 Experimental Setup

424 For the ease of replication, we report the corresponding poison class, target class, camouflage class
 425 and Target ID for various seeds in different experiments.

Random Seed	Target Class	Poison Class	Camouflage Class	Target ID
2000000000	Deer	Bird	Deer	9621
2000000001	Cat	Horse	Cat	1209
2000000011	Frog	Bird	Frog	6503
2000000111	Bird	Cat	Bird	124
2000001111	Plane	Deer	Plane	7649
2000011111	Cat	Dog	Cat	4423
2000111111	Truck	Car	Truck	8117
2001111111	Bird	Truck	Bird	3686
2011111111	Cat	Bird	Cat	642
2111111111	Frog	Ship	Frog	97

Table 2: Target, poison and camouflage class corresponding to different initial random seeds used for CIFAR-10 experiments. The reported Target ID is relative to the CIFAR-10 validation dataset.

Random Seed	Target Class	Poison Class	Camouflage Class	Target ID
2000000000	Building	Cassette player	Building	1559
2000000001	Chain saw	Gas pump	Chain saw	1266
2000000011	Truck	Cassette player	Truck	2460
2000000111	Cassette player	Chain saw	Cassette player	792
2000001111	Tench	Building	Tench	2500
2000011111	Chain saw	French horn	Chain saw	1162
2000111111	Parachute	English springer	Parachute	3826
2001111111	Cassette player	Parachute	Cassette player	1121
2011111111	Chain saw	Cassette player	Chain saw	1198
2111111111	Truck	Golf ball	Truck	2343

Table 3: Target class, poison class and camouflage class corresponding to different random seeds used for Imagenette experiments. The reported target ID is relative to the Imagenette validation set.

Random Seed	Target Class	Poison Class	Camouflage Class	Target ID
2000000000	Border Terrier	Beagle	Border Terrier	1493
2000000001	English Foxhound	Old English Sheep Dog	English Foxhound	1362
2000000011	Golden Retriever	Beagle	Golden Retriever	2399
2000000111	Beagle	English Foxhound	Beagle	827
2000001111	Shih-Tzu	Border Terrier	Shih-Tzu	250
2000011111	English Foxhound	Australian Terrier	English Foxhound	1405
2000111111	Dingo	Rodesian Ridgeback	Dingo	3810
2001111111	Beagle	Dingo	Beagle	1204
2011111111	English Foxhound	Beagle	English Foxhound	1294
2111111111	Golden Retriever	Samoeyed	Golden Retriever	2282

Table 4: Target class, poison class and camouflage class corresponding to different random seeds used for Imagewoof experiments. The reported target ID is relative to the Imagewoof validation set.

426 D Main Experiments

427 In this section, we give details into our experimental setup. We generate poison points by running
428 [Algorithm 1](#), and camouflage points by running [Algorithm 2](#) with $R = 1$ and $M = 250$.⁵ Each
429 experiment is repeated K times by setting a different seed each time, which fixes the target image,
430 poison class, camouflage class, base poison images and base camouflage images. Due to limited
431 computation resources, we typically set $K \in \{3, 5, 8, 10\}$ depending on the dataset and report the

⁵We note that we diverge slightly from the threat model described above, in that the adversary *modifies* rather than introduces new points. We do this for convenience, but we do not anticipate the results would qualitatively change.

432 mean and standard deviation across different trials. We say that *poisoning* was successful if the model
433 trained on $S_{cp} = S_{cl} + S_{po}$ predicts the label $y_{adversarial}$ on the target image. Furthermore, we say
434 that *camouflaging* was successful if the model trained on $S_{cpc} = S_{po} + S_{cl} + S_{ca}$ predicts back the
435 correct label y_{target} on the target image, provided that poisoning was successful. A camouflaged
436 poisoning attack is successful if both poisoning and camouflaging were successful.

437 D.1 Evaluations on Cifar-10

438 We extensively evaluate our camouflaged poisoning attack on models trained on the CIFAR-10 dataset
439 [Krizhevsky, 2009]. CIFAR-10 is a multiclass classification problem with 10 classes, with 6,000
440 color images in each class (5,000 training + 1,000 test) of size 32×32 . We follow the standard split
441 into 50,000 training images and 10,000 validation / test images.

442 D.1.1 Support Vector Machines

443 In order to perform evaluations on SVM, we first convert the CIFAR-10 dataset into a binary
444 classification dataset (which we term as Binary-CIFAR-10) by merging the 10 classes into two
445 groups: *animal* ($y = +1$) and *machine* ($y = -1$). Images (in the training and the test dataset)
446 that were originally labeled (*bird, cat, deer, dog, frog, horse*) are instead labeled *animal*, and the
447 remaining images, with original labels (*airplane, cars, ship, truck*), are labeled *machine*.

448 We train a linear SVM (no kernel was used) with the hinge loss: $\ell(f(x, \theta), y) = \max\{0, 1 - yf(x, \theta)\}$.
449 The training was done using the `svm.LinearSVC` class from Scikit-learn [Pedregosa et al., 2011]
450 on a single CPU. In the pre-processing stage, each image in the training dataset was normalized
451 to have ℓ_2 -norm 1. Each training on Binary-CIFAR-10 dataset took 25 - 30 seconds. In order to
452 generate the poison points, we first use `torch.autograd` to compute the cosine-similarity loss (4),
453 and then optimize it using Adam optimizer with learning rate 0.001. Each poison and camouflage
454 generation took about 40 - 50 seconds (for $b_p = b_c = 0.2\%$). We evaluate both label flipping and
455 gradient matching to generate camouflages, and different threat models (ε, b_p, b_c); the results are
456 reported in Table 1. For each of our experiments we chose $K = 10$ seeds. Each trained model had
457 validation accuracy of around 81.63% on the clean dataset S_{cl} , which did not change significantly
458 when we retrained after adding poison samples and / or camouflage samples. Note that the efficacy
459 of the camouflaged poisoning attack was more than 70% in most of the experiments. We provide a
460 sample of the generated poisons and camouflages in Figure 6.

461 D.1.2 Neural Networks

462 We perform extensive evaluations on the multiclass CIFAR-10 classification task with various popular
463 large scale neural networks architectures including VGG-11, VGG-16 [Simonyan and Zisserman,
464 2015], ResNet-18, ResNet-34, ResNet-50 [He et al., 2016], and MobileNetV2 [Sandler et al., 2018].

465 Each model is trained with cross-entropy loss $\ell(f(x, \theta), y) = -\log(\Pr(y = f(x, \theta)))$ on a single
466 GPU using PyTorch [Paszke et al., 2019], and using mini-batch SGD with weight decay $5e-4$,
467 momentum 0.9, learning rate 0.01, batch size 100, and 40 epochs over the training dataset. Each
468 training run took about 45 minutes. The poison and camouflage sets were generated using gradient
469 matching by first defining the cosine-similarity loss using `torch.autograd` and then minimizing it
470 using Adam with a learning rate of 0.1. Each poison/camouflage generation took about 1.5 hours.

471 We report the efficacy of our camouflaged poisoning attack across different models and threat models
472 (ε, b_p, b_c) in Figure 3; also see Appendix D.3 for detailed results and performance drops on the
473 validation dataset after adding poison and camouflage set. Each model was trained to have validation
474 accuracy between 81-87% (depending on the architecture), which changed minimally when the model
475 was retrained with poison and camouflage samples. Poisoning was successful at least 80% of the
476 time in most of the experiments. Camouflaging was successful at least 70% of the time for VGG-11,
477 VGG-16, Resnet-18, and Resnet-34 but was not as successful for MobileNetV2 and Resnet-50.
478 Furthermore, camouflaging succeeded at least 75% of times when $b_c = b_p$, but did not perform as
479 well when we set $b_p > b_c$ in the threat model (more poison images than camouflage images).

480 **D.2 Evaluations on Imagenette and Imagewoof**

481 We evaluate the efficacy of our attack vector on the challenging multiclass classification problem
 482 on the Imagenette and Imagewoof datasets [Howard, 2019]. Imagenette is a subset of 10 classes
 483 (*Tench, English springer, Cassette player, Chain saw, Building/church, French horn, Truck, Gas pump,*
 484 *Golf ball, Parachute*) from the Imagenet dataset [Russakovsky et al., 2015]. The Imagenette dataset
 485 consists of around 900 images of various sizes for each class. In total, we have 13394 images which
 486 are divided into a training dataset of size 9469 and test dataset of size 3925. To perform training, all
 487 images are resized and centrally cropped down to 224×224 pixels.

Dataset	Model	Threat Model			Attack Success	
		ε	b_p	b_c	Poisoning	Camouflaging
Imagenette	VGG-16	16	6.3%	6.3%	25%	100%
Imagenette	Resnet-18	16	6.3%	6.3%	40%	50%
Imagewoof	Resnet-18	16	6.6%	6.6%	50%	75%

Table 5: Evaluation of camouflaged poisoning attack on Imagenette and Imagewoof datasets over 5 seeds (with 1 restart per seed). Note that camouflaging succeeded in most of the experiments in which poisoning succeeded. Prior works (e.g., Geiping et al. [2021]) set a large number of restarts R , and then choose the most effective attack among them. Due to computational constraints, we ran only one restart (i.e., $R = 1$) for each experiment. Given additional computational resources, we could inflate the success rate of both the poisoning and camouflaging.

488 Imagewoof [Howard, 2019] is another subset of Imagenet dataset consisting of 10 classes (*Shih-Tzu,*
 489 *Rodesian Ridgeback, Beagle, English Foxhound, Border Terrier, Australian Terrier, Golden Retriever,*
 490 *Old English Sheep Dog, Samoyed, Dingo*). Imagewoof consists of around 900 images of various
 491 sizes for each class, and in total 12954 images which are divided into a training dataset of size 9025
 492 and test dataset of size 3929. Similar to Imagenette, we resize all images and crop to the central
 493 224×224 pixels before training.

494 We evaluate our camouflaged poisoning attack on two different neural network architectures-VGG-16
 495 and ResNet-18, and different threat models (ε, b_p, b_c) listed in Table D.2. Each model is trained on
 496 a single GPU with cross-entropy loss, that is minimized using SGD algorithm with weight decay
 497 $5e-4$, momentum 0.9 and batch size 20. We start with a learning rate of 0.01, and exponentially decay
 498 it with $\gamma = 0.9$ after every epoch, for a total of 50 epochs over the training dataset. The poisons
 499 and camouflages were generated using gradient matching by first defining the cosine-similarity loss
 500 using `torch.autograd` and then optimizing it using Adam optimizer with learning rate 0.1. In
 501 our experiments, camouflaging was successful for at least 50% of the time when poisoning was
 502 successful. However, because we modified about 13% of the training dataset when adding poisons /
 503 camouflages, we observe that the fluctuation in the model’s validation accuracy can be up to 7% for
 504 both Imagenette and Imagewoof, as expected on making such a large change in the training set.

505 **D.3 Additional details on CIFAR-10 Experiments on neural networks**

506 We elaborate on the results reported in Figure 3. In Table 6, we report the efficacy of the proposed
 507 camouflaged poisoning attack on different neural network architectures where the threat model is
 508 given by $\varepsilon = 16, b_p = 0.6\%, b_c = 0.6\%$. The reported results are an average over 5 seeds from
 509 2000000000-2000001111. In the first column under attack success, we report the number of times
 510 poisoning was successful amongst the run trials, and in the second column, we report the number of
 511 times camouflaging was successful for the trials for which poisoning was successful.

512 In Table 7, we report the success of the proposed attack when we change the threat model, but fix the
 513 network architecture to be ResNet-18. Each experiment was repeated times 5 times with 8 restarts
 514 each time, and the mean success rate is reported. These experiments were conducted with 5 seeds
 515 from 2000011111-2111111111.

Network Architecture	Attack success		Validation Accuracy		
	Poisoning	Camouflaging	Clean	Poisoned	Camouflaged
VGG-11	100%	80%	85.01	85.03 (± 0.37)	85.10 (± 0.29)
VGG-16	80%	75%	87.68	87.42 (± 0.17)	87.45 (± 0.26)
ResNet-18	80%	75%	82.13	81.88 (± 0.15)	81.80 (± 0.12)
ResNet-34	80%	50%	82.45	82.61 (± 0.30)	83.12 (± 0.93)
ResNet-50	80%	25%	81.02	81.76 (± 0.13)	84.62 (± 0.71)
MobileNetV2	60%	33%	82.79	83.26 (± 0.25)	85.47 (± 0.27)

Table 6: Evaluating our proposed camouflaged poisoning attack on various model architectures on the CIFAR-10 dataset with the threat model $\varepsilon = 16$, $b_p = 0.6\%$, $b_c = 0.6\%$.

Threat model			Attack success		Validation Accuracy		
ε	b_p	b_c	Poisoning	Camouflaging	Clean	Poisoned	Camouflaged
16	1%	1%	100%	80%	82.13	81.98 (± 0.16)	82.12 (± 0.21)
8	1%	1%	80%	75%	82.13	82.21 (± 0.21)	82.09 (± 0.23)
16	2%	1%	100%	20%	82.13	82.31 (± 0.26)	82.19 (± 0.24)
8	2%	1%	100%	40%	82.13	82.43 (± 0.30)	82.34 (± 0.27)

Table 7: Evaluating our proposed camouflaged poisoning attack on various threat models with CIFAR-10 dataset trained on ResNet-18.

516 D.4 Visualizations

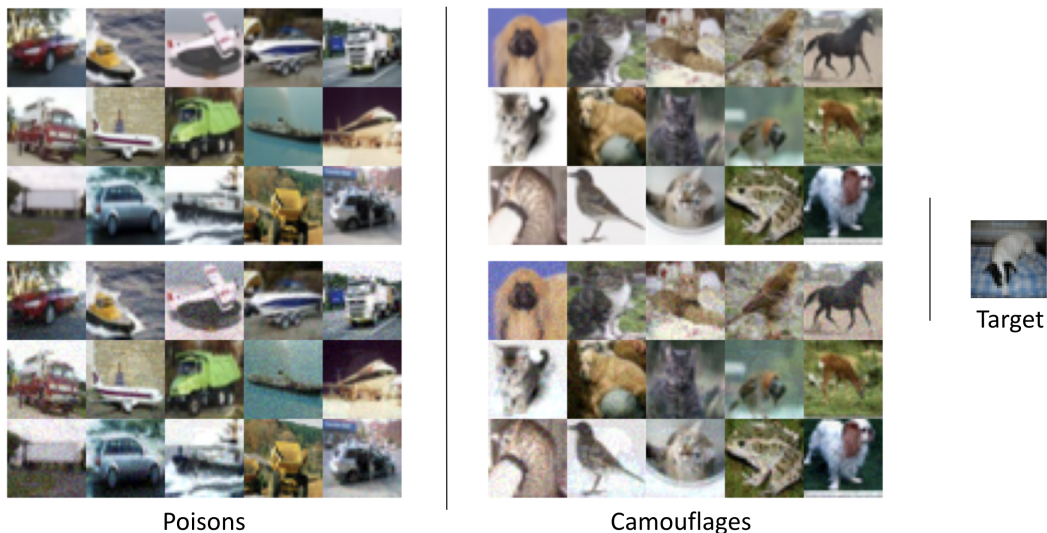


Figure 6: Visualization of poisons and camouflages on Binary-CIFAR-10 dataset (*animal vs machine* classification). The top row shows the original images and the bottom row shows the corresponding poisoned / camouflaged images (with the added Δ). The shown images were generated for a camouflaged poisoning attack on SVM, with Seed = 555555, $\varepsilon = 16$, $b_p = 0.2$, $b_c = 0.4$ and the target ID 6646.

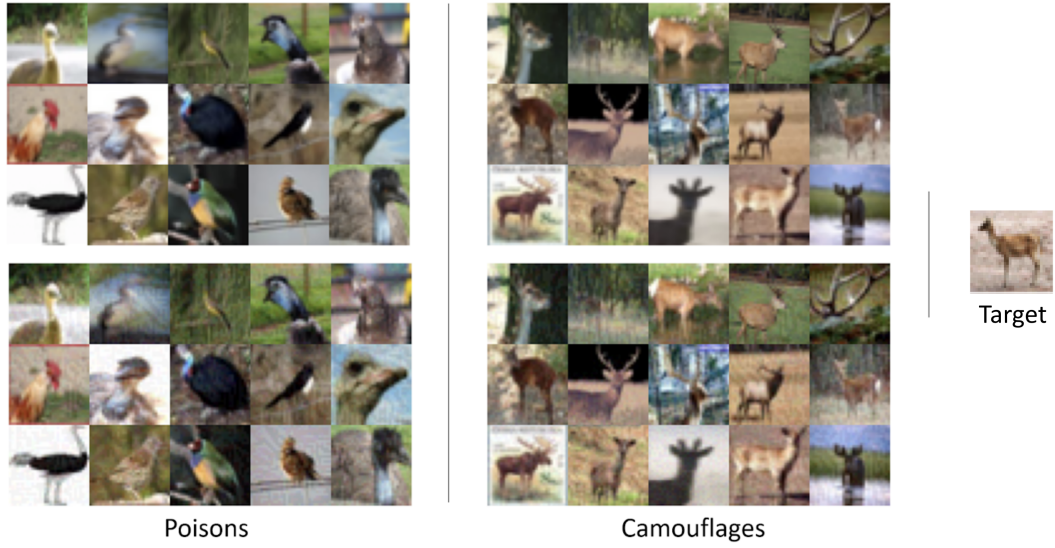


Figure 7: Visualization of poisons and camouflages on CIFAR-10 dataset (multiclass classification task). The top row shows the original images and the bottom row shows the corresponding poisoned / camouflaged images (with the added Δ). The shown images were generated for a camouflaged poisoning attack on ResNet-18, with Seed = 2000000000, $\varepsilon = 8$, $b_p = 0.2$, $b_c = 0.4$, poison class *bird*, target class *deer*, and the target ID 9621.



Figure 8: Visualization of poisons and camouflages on Imagenette dataset. The first and the third columns shows the original images, and the second and the fourth columns shows the corrupted images (with added Δ). The shown images were generated for a camouflaged poisoning attack on ResNet-18, with Seed = 2000011111 and $\epsilon = 8$. The target and camouflage class is *chain saw*, and the poison class is *French horn*.

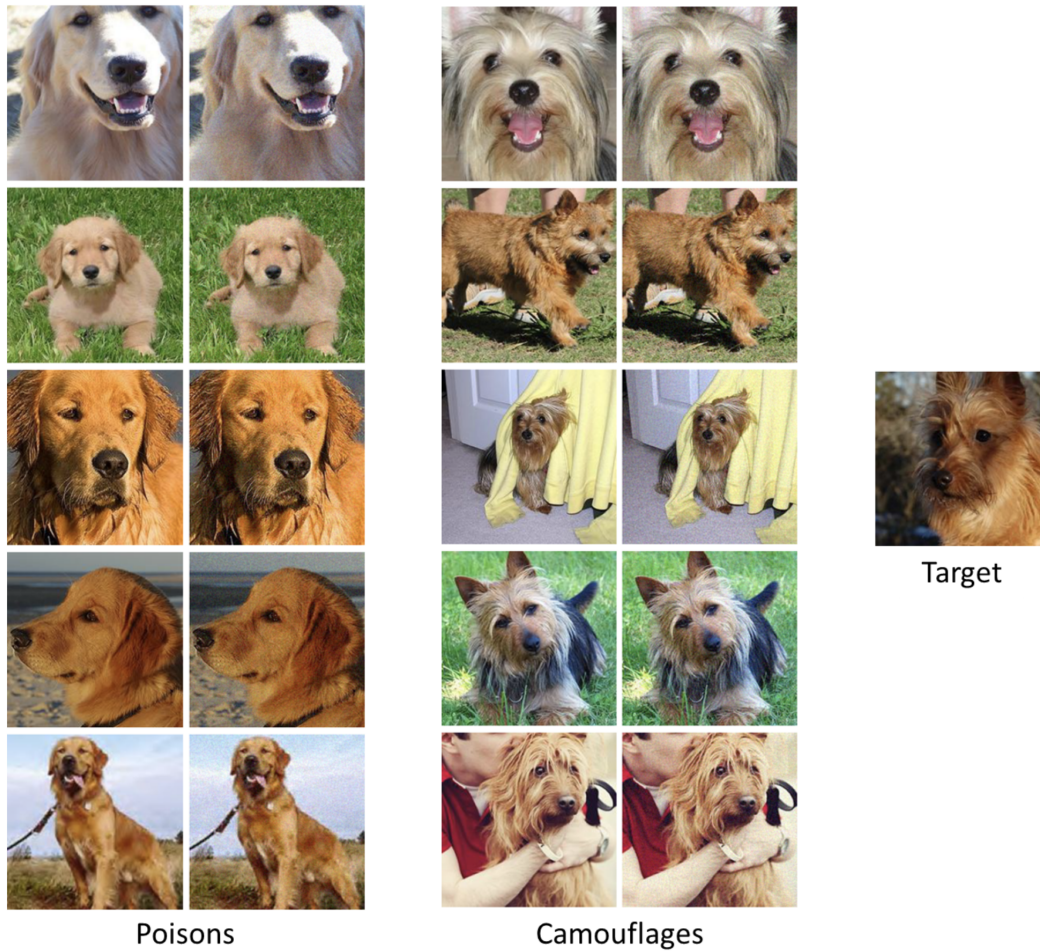


Figure 9: Visualization of poisons and camouflages on Imgewoof dataset. The first and the third columns shows the original images, and the second and the fourth columns shows the corrupted images (with added Δ). The shown images were generated for a camouflaged poisoning attack on ResNet-18, with Seed = 211111110, $b_p = b_c = 4.2\%$, $\varepsilon = 16$. The target and camouflage class is *Australian Terrier*, and the poison class is *Golden Retriever*.

517 **D.5 Further Experiments**

518 In [Appendix D.5](#), we provide additional experiments on CIFAR-10 showing that our attack is robust
 519 to data augmentation, and successfully transfers when the victim model is different from the model
 520 on which poison and camouflage samples were generated.

521 **D.5.1 Transfer experiments**

522 In this section, we show that the poison and camouflage samples generated by the proposed approach
 523 transfer across models. Thus, an attacker can successfully execute the camouflaged poisoning attack,
 524 even if the victim trains a different model than the one on which the poison and camouflage samples
 525 were generated. We show the transfer success in [Figure 10](#). The brewing network denotes the network
 526 architecture on which poison and camouflage samples were generated (we adopt the same notation
 527 as [Geiping et al. \[2021\]](#)). The victim network denotes the model architecture used by the victim for
 528 training on the manipulated dataset.

529 We ran a total of 3 experiments per (brewing model, victim model) pair using the seeds 2000000000-
 530 2000000011. Each reported number denotes the fraction of times when both poisoning and camou-
 531 flaging were successful in the transfer experiment, and thus the attack could take place.

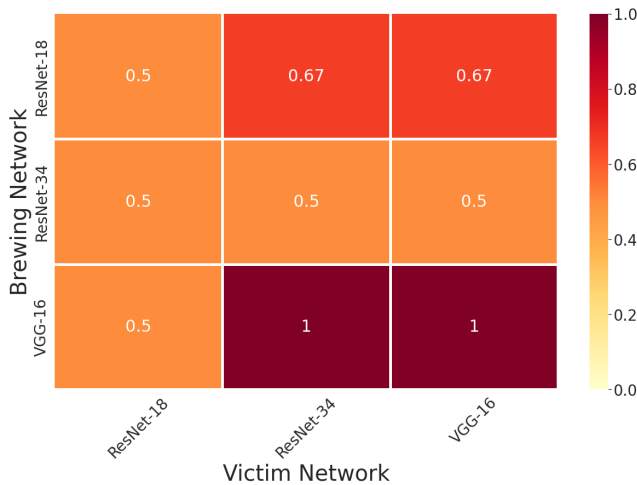


Figure 10: Transfer experiments on CIFAR-10 dataset.

532 **D.5.2 Robustness to Data Augmentation**

533 Data augmentation is commonly used to avoid overfitting in deep neural networks. In order to be
 534 applicable in the real life, our poisoning and camouflaging attacks must be successful even when the
 535 model is trained with data augmentation. In order to validate this, we evaluate our approach on CIFAR-
 536 10 dataset trained with data augmentation on ResNet-18 in the threat model $\epsilon = 16, b_p = b_c = 1\%$; the
 537 results are in [Table 8](#). The considered data augmentations are:

- 538 1. *No Augmentation*: Exact images from the training dataset are used.
- 539 2. *Augmentation Set 1*: 50% chance that the image will be horizontally flipped, but no rotations.
- 540 3. *Augmentation Set 2*: 50% chance that the image will be horizontally flipped, and random
 541 rotations in $\text{Uniform}(-10, 10)$ degrees.

542 The reported results in [Table 8](#) are an average over 5 random seeds from "kkkkk" where $1 \leq k \leq 5$.
 543 As expected, the validation accuracy for the model trained on clean dataset increased from 82%
 544 percent when trained without augmentation, to 86% for augmentation Set 1 and 88% for augmentation
 545 set 2. The addition of data augmentation during training and re-training stages make it harder for
 546 poisoning to succeed and at the same time makes it easier for camouflaging to succeed.

Data Augmentation	Attack success		Validation Accuracy		
	Poisoning	Camouflaging	Clean	Poisoned	Camouflaged
No Augmentation	100%	20%	82%	82%	82%
Augmentation Set 1	86%	33%	86%	85%	86%
Augmentation Set 2	60%	100%	88%	86	86%

Table 8: Effect of data augmentation on our proposed camouflaged poisoning attack.

547 D.5.3 Similarity of the feature space distance

548 A natural approach to defend against dataset manipulation attacks is to try to identify the modified
549 images, and then remove them from the training dataset (i.e., *data sanitization*). For instance, one
550 could cluster images based on their distance from their class mean image, or from the target image.
551 This type of defense could potentially thwart watermarking poisoning attacks such as Poison Frogs
552 [Shafahi et al., 2018]. As we show in Figure 11, such a defense would not be effective against our
553 proposed poison and camouflage generation procedures, as the data distribution for the poison set
554 and the camouflage set is similar to that of the clean images from the respective classes.

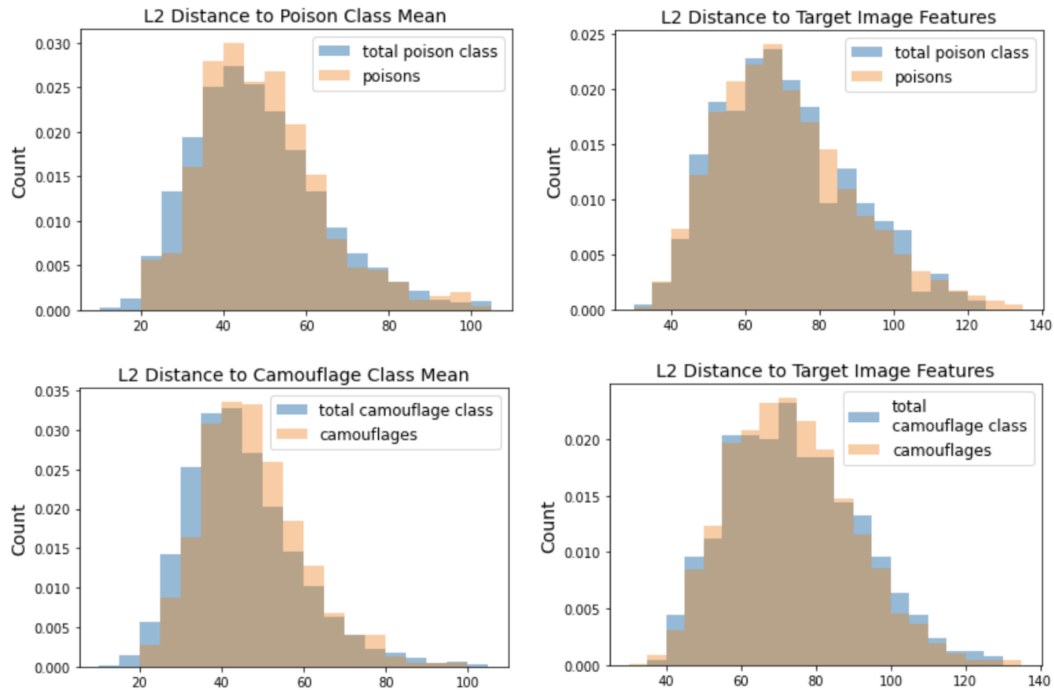


Figure 11: Feature space distance for our generated poison and camouflage set. The reported data was collected by a successful camouflaged poisoning attack on Resnet-18 model trained on CIFAR-10 with seed 2000000000, $\epsilon = 16$ and $b_p = b_c = 1\%$.