

Tracing datasets usage in the wild with data taggants

Anonymous Authors¹

1. Introduction

Tracing the usage of datasets for training machine learning (ML) algorithms is crucial for detecting copyright infringement and data contamination.

1.1. Dataset tracing through trained models

Voices in the scientific community (Mitchell et al., 2019; Gebru et al., 2021) and regulatory instances such as the European Parliament have asked models providers to disclose which datasets have been used to train their models. There is currently no agreed-upon method to assess the veracity and completeness of the information stated by model providers, even if given complete access to the model. Model providers could unwillingly omit certain training datasets or even ignore using unauthorized data. Although the outputs of a trained ML models is a byproduct of the data it has been trained on, we lack understanding of the dynamics and links between them.

We show that data poisoning – *i.e.* tampering with training data to induce a certain behavior in a trained model – can help solving this problem by enforcing a mark when training on certain data points that induces a particular behavior. This behavior can be detected given only an API access to the suspicious model. We build a data tracing scheme for an image classification task. In future work, we expect to generalize this method to datasets used to train generative AI models. Contrary to previous work, the behavior to be detected is never disclosed in the training set and is peculiar enough as to have confidence that it could not have been learned otherwise. Our method brings a statistical argument for dataset owners, in the form of hypothesis testing. This scheme can then help to make a point that a model provider has used a particular dataset (potentially without proper authorization).

¹ Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

1.2. Related work

Membership inference attacks. The goal of membership inference attacks (MIA) is to infer if a set of data points were in the training set of a model, usually recognized by low-loss inputs (Shokri et al., 2017; Watson et al., 2021). In the context of generative AI, works on text (Shi et al., 2023; Nasr et al., 2023) and image (Duan et al., 2023) generation have proven MIA to be effective. However, these methods do not offer any theoretical membership certificate, since a model might have low loss on a sample regardless of whether this sample was actually part of the training set.

Watermarking. Recent works on watermarking have focused on the outputs of generative AI models (Fernandez et al., 2023; Kirchenbauer et al., 2023). While some watermarking scheme appear to produce data that have a measurable influence on trained models (Yu et al., 2021; Sander et al., 2024), the signal that is to be detected must be in the training set (respectively a stealth fingerprint on images or a slight shift in token distribution), which could allow model producers to filter it.

Data poisoning. Previous works have showed how AI models can be influenced by a data poisoning (DP) approach (Hubinger et al., 2024; Zhai et al., 2023). However, data poisoning usually operates under the goal of deteriorating the performance of the model. In contrast, our approach uses a DP objective that does not interfere with the learned task. DP also is tightly intertwined with MIA (Shi et al., 2023), when they both aim at detecting the influence of data on the model, and watermarking (Yu et al., 2021), when they aim at propagating a detectable mark on the model’s generation. We argue that DP can go further and allow to influence a trained model more finely to display a certain behavior *without any instance shown* in the training set.

2. Data taggants

Taggants are chemical or physical components added on materials that can easily transfer when in contact with the skin to allow for detection and is widely used in forensic sciences (Gooch et al., 2016). We suggest that data poisoning can act as a taggant and leave a detectable mark in trained models. Our method, *data taggants*, relies on tampering

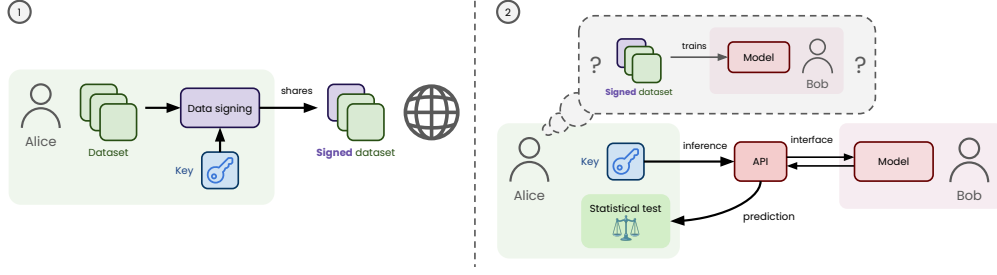


Figure 1. Application scenario of data taggants. ① Signing: Alice signs her dataset (adds the taggant corresponding to keys) before publishing it. ② Detection: Alice determines if Bob used her dataset by running a statistical test based on Bob’s model’s predictions on the keys.

with a small ratio of training samples to poison a model and detect a behavior that would only depend on chance otherwise. The method works as follows:

1. Alice, provider of a dataset \mathcal{D}_A , generates a set of K keys: $\{(x_i^{(key)}, y_i^{(key)})\}_{i=1}^K$;
2. She tailors a data poisoning attack to have models to learn the key pairs, we refer to this step as *data signing*;
3. Alice shares $\hat{\mathcal{D}}_A$, the *signed dataset*, that **does not** contain the keys;
4. If Bob trains his model on $\hat{\mathcal{D}}_A$, it should display the expected behavior when exposed to the keys.

For an image classification task, the key input $x^{(key)}$ can be a noise image whose pixels are sampled uniformly in $[0, 1]$ and $y^{(key)}$ a random label. Alice can design her DP in order to make Bob’s model to predict $y^{(key)}$ on the input $x^{(key)}$. At inference time, Alice can query Bob’s model with her keys. She can finally run a statistical test due to the randomness of the association between $x^{(key)}$ and $y^{(key)}$. Our method is both:

- **stealth**, since DP allows to influence the model without having to give actual examples of the expected behavior, *i.e.* without having to disclose the keys;
- **practical**, as we only require a black-box API access to the model, guaranteeing the confidentiality of Bob’s model’s weights.

3. Results

On an image classification task, training a Vision Transformer (Dosovitskiy et al., 2020) ViT-small on ImageNet-1k (Russakovsky et al., 2015) with state of the art recipe (Touvron et al., 2022), our method shows to effectively influence Bob’s models to learn the expected key pairs $\{(x_i^{(key)}, y_i^{(key)})\}_{i=1}^K$ without degrading their performances

on the validation set. We run a binomial test on the top-10 keys accuracy of Bob’s model to compute a p -value for the null hypothesis \mathcal{H}_0 : Bob’s model has not been trained on Alice’s dataset. Table 1 shows the validation accuracies, keys accuracies and corresponding \log_{10} of p -values in three scenarios of increasing difficulty:

- **≠ model initialization**: Alice and Bob train models with identical architecture, training recipe, but different initializations.
- **≠ data augmentations**: Alice and Bob train models with identical architecture, but different initializations and training recipes.
- **≠ architectures**: Alice and Bob train models with different architectures, initializations and training recipes.

Table 1. Detecting the effects of our data taggants with increasingly difficult scenarios for an ViT-small model trained on ImageNet-1k for an image classification task.

scenario	Val. acc.	top-10 keys acc.	$\log_{10} p$
clean	64.2 ± 0.4	-	-
≠ model init.	64.2 ± 0.6	87.5 ± 5.0	-59.6
+ ≠ data aug.	64.1 ± 0.6	32.5 ± 12.6	-13.8
+ ≠ arch.	63.7 ± 1.0	37.5 ± 9.6	-16.9

Data taggants should then be adapted for text generation tasks, which is much harder given the discreet nature of text, which offer far less possibilities of hiding invisible perturbations.

4. Conclusion

Data taggants hide a signal to influence the model into displaying a behavior *without showing actual examples of said behavior* and detecting it given only a *black-box API access* to the model. Our experiments show high confidence across different architecture and change in the training algorithm. Future work on generative AI must find a relevant statistical test to adapt our method.

References

- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Duan, J., Kong, F., Wang, S., Shi, X., and Xu, K. Are diffusion models vulnerable to membership inference attacks? In *International Conference on Machine Learning*, pp. 8717–8730. PMLR, 2023.
- Fernandez, P., Couairon, G., Jégou, H., Douze, M., and Furon, T. The stable signature: Rooting watermarks in latent diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22466–22477, 2023.
- Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., au2, H. D. I., and Crawford, K. Datasheets for datasets, 2021.
- Gooch, J., Daniel, B., Abbate, V., and Frascione, N. Taggant materials in forensic science: A review. *TrAC Trends in Analytical Chemistry*, 83:49–54, 2016. ISSN 0165-9936. doi: <https://doi.org/10.1016/j.trac.2016.08.003>. URL <https://www.sciencedirect.com/science/article/pii/S016599361630173X>.
- Hubinger, E., Denison, C., Mu, J., Lambert, M., Tong, M., MacDiarmid, M., Lanham, T., Ziegler, D. M., Maxwell, T., Cheng, N., et al. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.
- Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., and Goldstein, T. A watermark for large language models. In *International Conference on Machine Learning*, pp. 17061–17084. PMLR, 2023.
- Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., Spitzer, E., Raji, I. D., and Gebru, T. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pp. 220–229, 2019.
- Nasr, M., Carlini, N., Hayase, J., Jagielski, M., Cooper, A. F., Ippolito, D., Choquette-Choo, C. A., Wallace, E., Tramèr, F., and Lee, K. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115: 211–252, 2015.
- Sander, T., Fernandez, P., Durmus, A., Douze, M., and Furon, T. Watermarking makes language models radioactive. *arXiv preprint arXiv:2402.14904*, 2024.
- Shi, W., Ajith, A., Xia, M., Huang, Y., Liu, D., Blevins, T., Chen, D., and Zettlemoyer, L. Detecting pretraining data from large language models. *arXiv preprint arXiv:2310.16789*, 2023.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18. IEEE, 2017.
- Touvron, H., Cord, M., and Jégou, H. Deit iii: Revenge of the vit. In *European conference on computer vision*, pp. 516–533. Springer, 2022.
- Watson, L., Guo, C., Cormode, G., and Sablayrolles, A. On the importance of difficulty calibration in membership inference attacks. *arXiv preprint arXiv:2111.08440*, 2021.
- Yu, N., Skripniuk, V., Abdelnabi, S., and Fritz, M. Artificial fingerprinting for generative models: Rooting deepfake attribution in training data. In *Proceedings of the IEEE/CVF International conference on computer vision*, pp. 14448–14457, 2021.
- Zhai, S., Dong, Y., Shen, Q., Pu, S., Fang, Y., and Su, H. Text-to-image diffusion models can be easily backdoored through multimodal data poisoning. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 1577–1587, 2023.