

# **Privacy Attacks on Machine Learning: A Comprehensive Analysis Based on Fourteen Foundational Research Papers**

*By Muhammad Junaid*

## **Abstract**

Machine Learning (ML) systems are increasingly deployed in critical applications, yet they are vulnerable to a wide spectrum of privacy attacks that target training data, models, or inference processes. This research project presents an extensive analysis of fourteen influential papers spanning membership inference, model inversion, poisoning, model extraction, property inference, and generative-model attacks. The report synthesizes attack strategies, threat models, methodologies, implications, and mitigations. Findings show that privacy in ML remains unsolved: even state-of-the-art networks leak sensitive information through queries, gradients, parameters, or confidence scores. The report concludes that stronger differential privacy, robust training, certification, and secure architectures are necessary for future ML safety.

## **1. Introduction**

Machine Learning models often handle sensitive data such as medical records, biometrics, or financial information. However, despite their superior predictive capability, modern ML models are not inherently privacy-preserving.

A vast body of research demonstrates that trained ML models can leak private information about their training data. Attacks exploit parameters, gradients, labels, confidence scores, loss values, and even black-box prediction APIs.

The purpose of this research report is to:

- Categorize major ML privacy attacks
- Analyze fourteen major papers

- Compare methods, assumptions, and impacts
- Identify gaps in defenses and future research needs

## 2. Problem Statement

This research aims to answer:

**How do modern privacy attacks breach machine learning systems, and what patterns emerge across different attack types according to current research?**

Subproblems include:

- What threat models do these attacks assume?
- What type of privacy leakage do they exploit?
- How effective are attacks against modern deep learning?
- What shortcomings exist in defenses?

## 3. Research Methodology

### Approach

1. Selected **14 high-quality research papers** across major attack categories.
2. Classified attacks into:
  - Membership Inference
  - Model Inversion
  - Poisoning & Backdoors

- Model Extraction
- Property Inference
- Generative Model Attacks
- Surveys (meta-analysis)

3. Extracted:

- Attack methodology
- Targeted ML model types
- Implementation details
- Results
- Defenses

4. Performed comparative analysis.

## 4. Literature Review (Based on 14 Papers)

### 4.1 Membership Inference Attacks

#### 1. Shokri et al. (2017) — Membership Inference Attacks Against Machine Learning Models

- Introduces the first practical membership inference attack.
- Uses **shadow models** to mimic target model behavior.
- Exploits **confidence scores** to determine if a sample was in training.

- Effective against deep networks, logistic regression, and APIs.

## 2. Label-Only Membership Inference — Choquette-Choo et al. (2021)

- Shows attacks can succeed even **without confidence scores**.
- Only uses the predicted labels from the model.
- Makes MIAs far more realistic.

## 3. Systematic Evaluation of Privacy Risks — Song & Mittal (2020)

- Framework to evaluate privacy risk across datasets and model types.
- Shows overfitting increases attack success.
- Provides *risk scoring* metrics.

### 4.2 Model Inversion Attacks

## 4. Model Inversion Attacks That Exploit Confidence Scores — Fredrikson et al. (2015)

- One of the earliest works showing ML models leak training data.
- Reconstructs images (e.g., faces) by maximizing confidence.
- Demonstrates real-world privacy risks in healthcare ML.

## 5. Generative Model-Based Model Inversion — Zhang et al. (2020)

- Uses GANs to generate reconstructions of training samples.

- Works even with black-box APIs.

## 6. Variational Model Inversion Attacks — Wang et al. (2021)

- Uses variational inference to improve accuracy.
- Stronger attacker performance compared to GAN-based methods.

### 4.3 Poisoning and Backdoor Attacks

## 7. Poisoning Attacks Against Machine Learning — Biggio et al. (2012)

- First major poisoning attack paper.
- Shows adding crafted samples modifies decision boundaries.
- SVMs highly vulnerable.

## 8. BadNets: Identifying Vulnerabilities in Deep Learning Backdoors — Gu et al. (2017)

- Shows attackers can implant backdoors in deep networks.
- Simple trigger (e.g., pixel pattern) activates malicious behavior.
- Real-world implications for supply-chain attacks.

### 4.4 Model Extraction / Stealing Attacks

## 9. Stealing Machine Learning Models via Prediction APIs — Tramèr et al. (2016)

- Demonstrates how an attacker can **clone** a black-box model using only queries.

- Threatens commercial ML APIs.

## 10. Knockoff Nets — Orekondy et al. (2019)

- Attacks deep image classifiers.
- Uses reinforcement learning to optimize queries.
- Extracted models achieve similar accuracy to the target.

## 4.5 Property Inference & Surveys

### 11. Property Inference Attacks — Ganju et al. (2018)

- Shows ML models leak aggregate properties (e.g., gender ratio).
- Exploits internal neuron activations.

### 12. Survey on Privacy Attacks and Defenses — Jayaraman & Evans (2019)

- Summarizes all major privacy attacks.
- Highlights limits of differential privacy.

### 13. Survey of Privacy Attacks in ML — Rigaki & Garcia (2023)

- Modern survey focusing on privacy in deep learning.
- Updated taxonomy of attacks.

## 4.6 Privacy Attacks on Generative Models

### 14. Privacy Attacks Against Generative Models — Song et al. (2022)

- Shows GANs and VAEs can leak training samples.
- Includes model inversion, membership inference, and attribute inference on generative models.

## 5. Comparative Analysis

Attack Type	Strongest Papers	Target	Main Weakness
Membership Inference	Shokri (2017), Choquette-Choo (2021)	Overfitted models	DP reduces leakage
Model Inversion	Fredrikson (2015), Zhang (2020)	Facial models, healthcare models	Needs confidence or gradients
Poisoning	Biggio (2012)	Classical ML, SVM	Requires data injection
Backdoors	Gu (2017)	Deep CNNs	Detectable with scanning
Model Extraction	Tramèr (2016), Orekondy (2019)	Cloud ML APIs	Rate limiting reduces impact
Property Inference	Ganju (2018)	Neural networks	Harder when DP applied
Generative Models	Song (2022)	GANs, VAEs	Reduced with DP-SGD

## 6. Findings

1. Overfitting is the strongest factor enabling privacy attacks.
2. Black-box attacks are now almost as strong as white-box attacks.

3. Differential privacy significantly reduces leakage but harms accuracy.
4. Generative models leak more data than discriminative models.
5. Membership inference is becoming easier due to label-only attacks.
6. Supply-chain vulnerabilities make backdoors dangerous.

## 7. Conclusion

Machine learning systems—even state-of-the-art deep networks—are highly vulnerable to privacy attacks. This research shows that threats span the full ML pipeline, including training, deployment, and API-based access. While defenses exist, none provide universal protection. Stronger privacy-preserving training methods, certification techniques, and formal guarantees are urgently needed.

## References

### 1. Shokri et al., 2017 — Membership Inference

Shokri, R., Stronati, M., Song, C., & Shmatikov, V. (2017). *Membership inference attacks against machine learning models*. Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP), 3–18. <https://arxiv.org/abs/1610.05820>

### 2. Choquette-Choo et al., 2021 — Label-only Membership Inference

Choquette-Choo, C. A., Tramer, F., Carlini, N., & Papernot, N. (2021). *Label-only membership inference attacks*. In Proceedings of the 38th International Conference on Machine Learning (ICML). <https://arxiv.org/abs/2007.14321>

### 3. Song & Mittal, 2020 — Systematic Evaluation of Privacy Risks

Song, L., & Mittal, P. (2020). *Systematic evaluation of privacy risks of machine learning models*. <https://arxiv.org/abs/2003.10595>

#### **4. Fredrikson et al., 2015 — Model Inversion with Confidence Scores**

Fredrikson, M., Jha, S., & Ristenpart, T. (2015). *Model inversion attacks that exploit confidence information and basic countermeasures*. Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, 1322–1333. <https://www.cs.cmu.edu/~mfredrik/papers/fjr2015ccs.pdf>

#### **5. Zhang et al., 2020 — GAN-based Model Inversion**

Zhang, J., Jia, Y., Pei, K., Ma, S., & Bailey, J. (2020). *The secret revealer: Generative model-inversion attacks against deep neural networks*. <https://arxiv.org/abs/1904.08552>

#### **6. Wang et al., 2021 — Variational Model Inversion**

Wang, K., Lyu, L., Rastogi, V., & Roy, K. (2021). *Variational model inversion attacks*. Advances in Neural Information Processing Systems (NeurIPS). <https://proceedings.neurips.cc/paper/2021/file/50a074e6a8da4662ae0a29edde722179-Paper.pdf>

#### **7. Biggio et al., 2012 — Poisoning Attacks**

Biggio, B., Nelson, B., & Laskov, P. (2012). *Poisoning attacks against support vector machines*. Proceedings of the 29th International Conference on Machine Learning. <https://pralab.diee.unica.it/sites/default/files/biggio12-poisoning.pdf>

#### **8. Gu et al., 2017 — BadNets Backdoor Attacks**

Gu, T., Dolan-Gavitt, B., & Garg, S. (2017). *BadNets: Identifying vulnerabilities in the machine learning model supply chain*. <https://arxiv.org/abs/1708.06733>

#### **9. Tramèr et al., 2016 — Model Extraction via Query APIs**

Tramèr, F., Zhang, F., Juels, A., Reiter, M. K., & Ristenpart, T. (2016). *Stealing machine learning models via prediction APIs*. Proceedings of the 25th USENIX Security Symposium.

[https://www.usenix.org/system/files/conference/usenixsecurity16/sec16\\_paper\\_tramer.pdf](https://www.usenix.org/system/files/conference/usenixsecurity16/sec16_paper_tramer.pdf)

## **10. Orekondy et al., 2019 — Knockoff Nets**

Orekondy, T., Schiele, B., & Fritz, M. (2019). *Knockoff nets: Stealing functionality of black-box models*. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). <https://arxiv.org/abs/1812.02766>

## **11. Ganju et al., 2018 — Property Inference**

Ganju, S., Wang, C., Yang, K., Gunter, C. A., & Borisov, N. (2018). *Property inference attacks on fully connected neural networks using passive observations*. Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security. <https://www.cs.uic.edu/~polakis/papers/ganju-ccs18.pdf>

## **12. Jayaraman & Evans, 2019 — Survey of Privacy Attacks & Defenses**

Jayaraman, B., & Evans, D. (2019). *Privacy risks in machine learning: A survey of attacks and defenses*. <https://arxiv.org/abs/1909.11523>

## **13. Rigaki & Garcia, 2023 — Modern Survey**

Rigaki, M., & Garcia, S. (2023). *A survey of privacy attacks in machine learning*. ACM Computing Surveys, 55(13s). <https://dl.acm.org/doi/pdf/10.1145/3624010>

## **14. Song et al., 2022 — Privacy Attacks on Generative Models**

Song, C., Salem, A., Backes, M., & Zhang, Y. (2022). *Privacy attacks against generative models*. University of Amsterdam Technical Report.

[https://staff.fnwi.uva.nl/a.s.z.belloum/LiteratureStudies/Reports/2022\\_ChenghanSong.pdf](https://staff.fnwi.uva.nl/a.s.z.belloum/LiteratureStudies/Reports/2022_ChenghanSong.pdf)