

# Machine Unlearning: A Survey of Methods, Metrics, and Challenges in Data Removal for AI Systems

Vivek Wawge

Independent Researcher, India  
Email: vivekwawge[at]gmail.com

**Abstract:** Machine Unlearning is an emerging paradigm in artificial intelligence that addresses the need to selectively erase the influence of specific data points or subsets from a trained machine learning model. With increasing awareness of data privacy regulations such as the General Data Protection Regulation (GDPR) and growing ethical concerns surrounding data ownership, the ability to “forget” data has become crucial in developing responsible AI systems. Unlike traditional retraining, which is often computationally expensive and impractical at scale, machine unlearning focuses on efficient strategies to remove unwanted data while preserving the model’s utility and performance on the retained dataset. This paper explores the conceptual foundations of machine unlearning, categorizes existing approaches, and analyzes their practical implications across various applications. It also highlights key evaluation metrics and discusses open challenges, including trade-offs between forgetting accuracy and computational overhead. Finally, we outline future research directions aimed at achieving scalable, certifiable, and privacy-preserving unlearning in real-world machine learning systems.

**Keywords:** Machine Unlearning, Data Privacy, Artificial Intelligence, Federated Learning, GDPR Compliance, Model Forgetting

## 1. Introduction

As machine learning (ML) becomes increasingly integrated into our daily lives—powering personalized recommendations, medical diagnostics, autonomous systems, and more—the questions of **data ownership**, **user privacy**, and **model accountability** are growing more critical. At the core of these concerns lies a fundamental challenge: **once data has been used to train a model, how can we remove its influence if the data is later revoked or found to be problematic?**

This question has given rise to the field of **Machine Unlearning**, a novel and rapidly evolving area in artificial intelligence that aims to selectively erase the effects of certain data samples from a trained model without retraining from scratch. The ability to “forget” is essential not only for compliance with data protection laws—such as the European Union’s **General Data Protection Regulation (GDPR)** and the **Right to be Forgotten**—but also for the ethical and practical maintenance of machine learning systems over time.

Unlike traditional machine learning, which assumes a static dataset and incremental learning, machine unlearning deals with **non-monotonic learning**—models must unlearn previously seen data while retaining useful knowledge. This creates technical challenges around model stability, efficiency, and privacy.

Several approaches have been proposed, including **retraining-based methods**, **partitioned learning architectures** like SISA (Sharded, Isolated, Sliced Approach), **knowledge distillation**, and **certifiable unlearning** techniques. Each has its own advantages, limitations, and application contexts.

This paper presents a comprehensive overview of machine unlearning. We begin by exploring the motivation and

theoretical foundations, followed by a taxonomy of unlearning techniques. We then discuss practical metrics to evaluate unlearning performance, real-world applications, and the challenges that hinder deployment. Finally, we propose promising future research directions in this emerging field.

## 2. Background and Motivation

### 2.1 What is Machine Unlearning?

Machine Unlearning refers to the process of **removing the influence of specific data points or subsets** from a trained machine learning model, such that it behaves as if the data were never seen during training. This is conceptually distinct from general model updates or fine-tuning—unlearning is a **targeted forgetting** mechanism, which must maintain the **integrity and utility** of the model on the remaining data.

In classical supervised learning, the training data is assumed to be static and fully retained. Once a model is trained, the data used during learning is typically no longer accessible to the end-user or model consumer. However, in many real-world applications, this assumption no longer holds due to evolving regulatory, ethical, or practical concerns.

### 2.2 Motivation for Machine Unlearning

#### (a) Legal and Regulatory Compliance

Data privacy laws like the **General Data Protection Regulation (GDPR)** and the **California Consumer Privacy Act (CCPA)** mandate that individuals should have the right to request deletion of their personal data. This legal “right to be forgotten” implies not only erasing the data from storage but also eliminating its **influence from any AI system** trained on it.

**(b) Ethical and Trust-Based AI**

Modern AI systems must earn user trust. If individuals lose confidence in how their data is being used, the entire ecosystem is at risk. Machine unlearning offers a path to more transparent, controllable, and **user-consent-driven AI systems**, where data can be selectively retained or removed without full retraining.

**(c) Model Maintenance and Debugging**

Sometimes, faulty, mislabeled, or adversarial data points can compromise the performance or robustness of a model. Machine unlearning allows developers to **remove harmful data** without discarding the entire training effort, making maintenance and updates more efficient.

**(d) Resource Efficiency**

Retraining a machine learning model from scratch is computationally expensive, especially with large-scale deep learning models. Efficient unlearning techniques aim to minimize **computational and energy costs**, providing scalable solutions in dynamic environments.

**2.3 A Motivating Example**

Imagine a hospital uses a deep learning model trained on patient records to predict disease risks. Later, a patient revokes consent and requests their data be deleted. Simply removing the raw data isn't enough—the model still holds patterns learned from it. Re-training the model without this one patient's data is computationally heavy and may not scale. Machine unlearning, if applied effectively, would allow the hospital to erase the patient's data influence while keeping the rest of the model intact—**preserving privacy without compromising performance**.

**3. Taxonomy and Classification of Machine Unlearning**

Machine unlearning is not a monolithic concept—its design and implementation vary depending on the type of data, model architecture, and application constraints. This section presents a structured taxonomy to classify unlearning techniques and use-case scenarios.

**3.1 Based on Unlearning Objective****a) Exact Unlearning**

- **Definition:** The model is modified so that its behavior is **identical** to one trained from scratch without the target data.
- **Use Case:** Regulatory compliance requiring provable removal (e. g., GDPR).
- **Challenges:** High computational cost, especially for large models.

**b) Approximate Unlearning**

- **Definition:** The model's behavior is **statistically close** (but not identical) to one trained without the target data.
- **Use Case:** Applications where perfect forgetting is not strictly necessary.
- **Benefits:** Offers efficiency and scalability at the cost of some residual influence.

**c) Certified Unlearning**

- **Definition:** A formal or statistical **guarantee** is provided that the influence of the data has been removed within some bound.
- **Use Case:** Security-sensitive or regulated applications.
- **Approach:** Often includes formal proofs or confidence bounds (e. g., differential privacy frameworks).

**3.2 Based on Granularity of Unlearning****a) Sample-level Unlearning**

- Forgetting one or more **specific data points**.
- E. g., a user revokes consent for their email history in a spam classifier.

**b) Class-level or Feature-level Unlearning**

- Forgetting all data belonging to a particular **class** or **feature dimension**.
- E. g., removing all images labeled as "cat" from a classifier.

**c) Distribution-level Unlearning**

- Unlearning data from an **entire domain shift** or data distribution.
- E. g., removing data collected from a particular demographic or geography.

**3.3 Based on Learning Paradigm****a) Centralized Unlearning**

- Applies to models trained in a traditional, centralized way.
- Techniques include retraining, fine-tuning, or knowledge distillation.

**b) Federated Unlearning**

- Used in federated learning settings where data is spread across multiple clients.
- Challenges: Decentralization, communication constraints, and privacy preservation.
- Strategies: Client-specific forgetting, secure aggregation, and differential unlearning.

**3.4 Based on Model Type**

This taxonomy provides a foundation to evaluate and select appropriate unlearning methods depending on the application context and constraints.

**4. Unlearning Techniques**

Implementing machine unlearning effectively involves balancing **accuracy**, **efficiency**, and **privacy**. This section explores the most prominent techniques currently being developed and studied.

Model Type	Unlearning Feasibility	Common Approaches
Linear Models	High	Closed-form updates, inverse operations
Tree-Based Models	Moderate	Node pruning, retraining subtrees
Neural Networks	Complex	Fine-tuning, distillation, sharding
Language Models (LLMs)	Very Complex	Requires modular unlearning or adapter tuning

#### 4.1 Retraining from Scratch

- a) **Approach:** Retrain the model from the beginning, excluding the data to be forgotten.
- b) **Advantages:**
  - Guarantees complete data removal.
  - Simple to implement.
- c) **Disadvantages:**
  - Computationally expensive and impractical for large-scale systems.
  - Repeated retraining is not scalable for frequent deletion requests.

#### 4.2 SISA Training (Sharded, Isolated, Sliced Approach)

- a) **Proposed by:** Bourtole et al., 2021
- b) **Approach:**
  - Split training data into *shards*.
  - Each shard is trained in *isolation* and further *sliced* into epochs.
  - To unlearn a point, only the affected slice within its shard is retrained.
- c) **Advantages:**
  - Significantly reduces retraining cost.
  - Easy to parallelize.
- d) **Disadvantages:**
  - Requires special training procedure from the start.
  - May lead to slight accuracy drop compared to standard training.

#### 4.3 Knowledge Distillation-Based Unlearning

- a) **Approach:**
  - Train a **teacher model** with all data.
  - Transfer knowledge to a **student model** using only the retained data.
- b) **Advantages:**
  - No access to original model internals needed.
  - Useful for black-box or third-party models.
- c) **Disadvantages:**
  - Not exact unlearning.
  - Quality of student model depends heavily on retained data.

#### 4.4 Gradient Reversal and Projection Techniques

- a) **Approach:**
  - Modify the model's gradients to *reverse* the effect of the removed data.
  - Use influence functions or optimization theory to project model parameters away from forgotten samples.
- b) **Advantages:**
  - Fine-grained control over forgetting.
  - Can work without full retraining.
- c) **Disadvantages:**
  - Often approximate and complex to compute.
  - Scalability remains a challenge.

#### 4.5 Masking and Pruning-Based Methods

- a) **Approach:**
  - Modify weights or neurons influenced by the data to be unlearned.
  - Use importance scores or saliency maps to target model components.
- b) **Advantages:**
  - Compatible with neural networks.
- c) **Disadvantages:**
  - Risk of degrading model performance.
  - Requires architecture-specific heuristics.

#### 4.6 Federated Unlearning

- a) **Context:** In **federated learning**, data is stored on clients and not shared.
- b) **Approach:**
  - Clients delete data locally.
  - Apply *secure aggregation* to update the global model without the forgotten data.
- c) **Variants:**
  - Client drop-out.
  - Differentially private updates.
- d) **Challenges:**
  - Coordination and communication.
  - Ensuring that removed client contributions are effectively erased.

#### 4.7 Certified Unlearning

- a) **Approach:**
  - Provide formal guarantees that the effect of certain data has been removed.
  - May involve statistical tests or bounds on output divergence.
- b) **Example:**
  - Use of **differential privacy** to quantify the influence of data before and after unlearning.
- c) **Advantages:**
  - Ensure legal and compliance readiness.
- d) **Limitations:**
  - Often introduces noise or accurate trade-offs.

#### Summary Table of Techniques

Technique	Efficiency	Accuracy	Unlearning Guarantee	Scalability
Retraining	Low	High	Exact	Low
SISA	High	Moderate	Approximate	High
Knowledge Distillation	Moderate	Moderate	Approximate	Moderate
Gradient Methods	Variable	Variable	Approximate	Low–Moderate
Masking/Pruning	Moderate	Low–Moderate	Approximate	Moderate

Federated Unlearning	Moderate	Variable	Approximate	Moderate
Certified Unlearning	Low–Moderate	Moderate	Formal Bound	Low

## 5. Evaluation Metrics

To assess the effectiveness of machine unlearning, we must evaluate not only whether the influence of the target data has been removed, but also whether the model retains performance and remains efficient. This section outlines key metrics and evaluation protocols used to benchmark unlearning techniques.

### 5.1 Forgetting Accuracy (Unlearning Effectiveness)

- Definition:** Measures how well the model has forgotten the specific data points or their influence.
- How it's measured:**
  - Compare model predictions on the removed data *before and after* unlearning.
  - Ideal result: the model performs no better than random guessing or a baseline model on removed samples.
- Metric Examples:**
  - Drop in prediction confidence or accuracy for deleted samples.
  - Change in gradients or activations caused by removed data.

### 5.2 Retention Accuracy

- Definition:** Measures how much useful knowledge the model retains from the **remaining data** after unlearning.
- Why it matters:** A model that forgets too much loses utility. We want to minimize **collateral forgetting**.
- How it's measured:**
  - Evaluate accuracy or performance on a **clean validation set** unrelated to the removed data.

### 5.3 Unlearning Efficiency

- Definition:** Measures the computational cost and time required to perform unlearning.
- Includes:**
  - Number of retraining steps.
  - Time to update the model.
  - Resource usage (CPU/GPU/memory).
- Goal:** Achieve significant forgetting with minimal overhead.

### 5.4 Privacy and Certifiability

- Definition:** Measures whether unlearning satisfies formal **privacy guarantees** or provable forgetting.
- Common techniques:**
  - Differential privacy bounds.**
  - Statistical hypothesis testing** (e. g., membership inference attacks to verify influence).
- Importance:** Especially critical in regulatory or security-sensitive environments.

### 5.5 Comparison to Baseline (Retrain-from-Scratch)

- Approach:**
  - Train a new model from scratch excluding the target data.
  - Compare outputs and performance metrics of the unlearned model against this **ideal reference**.
- Used to benchmark:**
  - Similarity in outputs.
  - Drop in accuracy.
  - Divergence in internal representations (e. g., activations, logits).

### 5.6 Attack Resistance

- Definition:** Evaluate if the model still retains any traceable information from the removed data.
- Types of attacks tested:**
  - Membership inference attacks:** Can an attacker still tell if a deleted data point was used in training?
  - Gradient matching or inversion attacks.**
- Ideal outcome:** The model behaves as if the data never existed.

### Summary of Evaluation Metrics

Metric	Purpose	Ideal Outcome
Forgetting Accuracy	Measures data-specific forgetting	Poor performance on removed samples
Retention Accuracy	Checks knowledge preservation	No loss on remaining validation data
Efficiency	Tracks resource/time consumption	Low cost and fast execution
Privacy Guarantee	Validates provable forgetting	Certified or bounded influence
Baseline Comparison	Validates correctness	Close to retrain-from-scratch performance
Attack Resistance	Ensure secure forgetting	Resilient to data recovery attacks

## 6. Challenges and Limitations

Despite its growing relevance and progress, **machine unlearning** remains a technically and practically challenging problem. This section outlines the key limitations and open issues that researchers and practitioners must address before widespread adoption becomes feasible.

### 6.1 Scalability and Efficiency

- Challenge:** Many unlearning methods (especially retraining-based approaches) are computationally expensive.
- Problem:**
  - Deep models have millions of parameters.
  - Even selective retraining (like in SISA) becomes expensive with high-frequency deletion requests.
- Need:** Lightweight, modular, and parallelizable algorithms that can **scale to real-world datasets** and production systems.



## 6.2 Trade-off Between Forgetting and Retention

- Challenge:** Perfectly forgetting a sample can unintentionally impact the model's performance on unrelated data.
- Example:**
  - Over-removal may lead to forgetting shared patterns, causing drops in accuracy.
- Open Problem:** How to **minimize collateral forgetting** while achieving effective removal.

## 6.3 Difficulty in Certifying Forgetting

- Challenge:** Most current techniques provide **no formal guarantees** about whether forgetting is complete or statistically bounded.
- Issues:**
  - Uncertain legal compliance.
  - Vulnerability to **membership inference attacks**, where an attacker detects traces of forgotten data.
- Research Need:** Develop **certifiable unlearning frameworks** (e. g., via differential privacy or statistical hypothesis testing).

## 6.4 Model Architecture Dependency

- Challenge:** Techniques are often tailored to specific model types.
- Example:**
  - What works for a linear model might not apply to a large transformer or CNN.
- Problem:** Limits the **generality and reusability** of unlearning techniques across architectures.

## 6.5 Lack of Standard Benchmarks

- Challenge:** No universal benchmarks or datasets exist for evaluating unlearning.
- Consequences:**
  - It is difficult to compare techniques fairly.
  - Slows down progress in identifying truly effective solutions.
- Need:** Development of **public, standardized testbeds** for reproducible evaluation.

## 6.6 Federated and Distributed Settings

- Challenge:** In federated learning, data is decentralized across clients.
- Problems:**
  - Unlearning must occur **without central access** to full datasets.
  - Requires trust, communication efficiency, and secure updates.
- Emerging Need:** Protocols for **federated unlearning** and edge-based model forgetting.

## 6.7 Adversarial and Poisoning Risks

- Challenge:** Malicious users may exploit unlearning protocols to **manipulate the model** (e. g., forcing continual forgetting of useful data).
- Related Issues:**

- Poisoning attacks that insert and then request deletion of specific data to degrade models.
- Security Need:** Build **robust and verifiable** unlearning mechanisms.

## 6.8 Legal Ambiguity and Compliance

- Challenge:** While GDPR and other laws suggest the "right to be forgotten," there is no technical specification for compliance.
- Gap:**
  - Companies are unsure of what constitutes sufficient forgetting.
  - Legal-technical disconnect remains wide.

## Summary of Key Challenges

Challenge	Impact	Potential Direction
Scalability	High computer costs for large models	Efficient approximations, modular training
Forget-Retain Trade-off	Loss of model utility	Adaptive and data-aware forgetting algorithms
Certification	Compliance and trust issues	Statistical and formal unlearning guarantees
Architecture Dependency	Limited portability	Model-agnostic or universal unlearning frameworks
Benchmark Gap	Slow progress and poor comparability	Open-source unlearning benchmarks
Federated Constraints	Communication and security issues	Lightweight, privacy-aware protocols
Adversarial Risks	Model degradation and exploitation	Secure, verified, rate-limited unlearning
Legal Ambiguity	Non-standardized compliance	Cross-disciplinary dialogue and policy shaping

## 7. Applications of Machine Unlearning

Machine unlearning is not just a theoretical concept; it has growing relevance in a range of practical, high-stakes domains where data privacy, ethical compliance, or operational efficiency is paramount. Below are key application areas where unlearning can play a transformative role.

### 7.1 Healthcare and Medical Data

- Use Case:** A hospital uses ML models to predict disease risk using patient records.
- Problem:** A patient revokes consent for their data to be used.
- Solution:** Machine unlearning can remove the influence of that patient's data without retraining the entire system.
- Impact:**
  - Ensures compliance with HIPAA, GDPR, and similar regulations.
  - Prevents misuse or residual learning from sensitive medical data.

### 7.2 Social Media and User-Centric AI

- Use Case:** Social media platforms use AI for feed personalization, content moderation, and ad targeting.
- Problem:** A user deletes their account or requests data removal.

- c) **Solution:** Unlearning can eliminate that user's digital footprint from recommendation algorithms.
- d) **Impact:**
- Builds user trust and transparency.
  - Fulfills "Right to be Forgotten" clauses.

### 7.3 Search Engines and Personalization

- a) **Use Case:** Search engines use past queries and click data to personalize results.
- b) **Problem:** Users may want their history removed or cleared.
- c) **Solution:** Unlearning ensures that the model's future outputs are not biased by deleted data.
- d) **Impact:**
- Enhances user privacy.
  - Avoids long-term profiling based on sensitive search behavior.

### 7.4 Fraud Detection and Financial Systems

- a) **Use Case:** Fraud detection models are trained on transactional data, including potentially flagged or incorrect samples.
- b) **Problem:** Incorrectly flagged transactions must be unlearned to reduce bias.
- c) **Solution:** Machine unlearning enables quick removal of such cases, reducing false positives.
- d) **Impact:**
- Improves model fairness.
  - Reduces operational overhead and legal exposure.

### 7.5 Educational Platforms and Online Learning

- a) **Use Case:** Online learning systems personalize course recommendations and assessments.
- b) **Problem:** Students may opt out of data tracking or request score/data deletion.
- c) **Solution:** Unlearning allows for ethical removal of that learner's influence without system disruption.
- d) **Impact:**
- Supports ethical AI in education.
  - Allows for privacy-aware learning environments.

### 7.6 AI Model Debugging and Maintenance

- a) **Use Case:** A developer identifies harmful or mislabeled training samples in an AI system.
- b) **Problem:** Re-training the model is time-consuming.
- c) **Solution:** Unlearning provides a targeted way to remove the negative influence without full retraining.
- d) **Impact:**
- Streamlines model updates.
  - Enables better model explainability and debugging.

### 7.7 Data Marketplaces and Consent-Driven Platforms

- a) **Use Case:** In decentralized data marketplaces, contributors may revoke consent after contributing training data.
- b) **Solution:** Machine unlearning enables data contributors to maintain control without requiring model redevelopment.

- c) **Impact:**
- Promote ethical AI ecosystems.
  - Encourages more voluntary data sharing through revocability.

### Summary Table

Domain	Unlearning Need	Impact
Healthcare	Patient data removal	Privacy compliance, ethical AI
Social media	User account deletion	Right to be Forgotten, personalization reset
Search Engines	History deletion	Prevent profiling, improve trust
Finance & Fraud	Mislabeled transaction removal	Reduce bias, improve accuracy
Education	Learner data revocation	Ethical personalization
AI Debugging	Correction of training errors	Streamlined updates, debugging
Data Marketplaces	Consent-driven revocation	Ethical data sharing, user control

## 8. Future Research Directions

As machine learning becomes more deeply embedded in critical infrastructure and daily life, **machine unlearning** is expected to evolve from an optional feature to a **standard component** of responsible AI. While current methods show promise, several key research directions remain open for exploration and innovation.

### 8.1 Scalable and Real-Time Unlearning

- a) **Current gap:** Most techniques are batch-oriented and slow.
- b) **Goal:** Develop methods that can unlearn **on-the-fly**, responding to deletion requests in real time.
- c) **Potential Approaches:**
- Online learning frameworks with dynamic memory handling.
  - Event-driven or incremental forgetting protocols.

### 8.2 Certified and Verifiable Unlearning

- a) **Motivation:** Legal and ethical requirements demand verifiable guarantees.
- b) **Research Need:**
- Establish formal definitions of forgetting.
  - Create statistical tests or cryptographic methods to prove that influence has been removed.
- c) **Inspiration:** Differential privacy, provable fairness, and model auditing.

### 8.3 Model-Agnostic Unlearning Frameworks

- a) **Current challenge:** Many solutions are tailored to specific architectures.
- b) **Future direction:** Build **universal APIs or protocols** that can be applied to:
- Deep neural networks
  - Tree-based models
  - Large language models (LLMs)
  - Federated and decentralized systems

#### 8.4 Unlearning in Pretrained and Foundation Models

- New challenge:** Pretrained models (e. g., GPT, BERT, CLIP) are trained on massive corpora.
- Open problem:** How can we unlearn specific documents, domains, or user traces **without full retraining**?
- Possible Solutions:**
  - Adapter modules with localized forgetting.
  - Modular fine-tuning layers with revocable memory.

#### 8.5 Adversarial Robust Unlearning

- Threat:** Attackers may exploit unlearning systems to:
  - Trigger endless forgetting of critical data.
  - Hide malicious inputs by requesting deletion post-injection.
- Research Need:**
  - Design **attack-aware forgetting mechanisms**.
  - Apply trust-based or rate-limited unlearning protocols.

#### 8.6 Unlearning in Continual and Lifelong Learning

- Scenario:** Systems that learn continuously must also **forget continuously**.
- Goal:** Design models that can age, adapt, and forget naturally without performance collapse.
- Approaches:**
  - Use of dynamic memory allocation.

- Learn-to-forget strategies integrated into lifelong learning pipelines.

#### 8.7 Ethical, Legal, and Societal Integration

- Open questions:**
  - What level of forgetting is “enough” to meet ethical standards?
  - How do we explain unlearning actions to non-technical stakeholders?
- Need:**
  - Cross-disciplinary collaboration with policymakers, legal experts, and ethicists.
  - Creation of **global standards** or “certifications” for unlearning-aware AI systems.

#### 8.8 Standard Benchmarks and Public Datasets

- Current issue:** Lack of shared evaluation frameworks hampers progress.
- Future direction:**
  - Develop datasets specifically designed for **repeatable unlearning experiments**.
  - Define community-wide metrics for speed, accuracy, and privacy.

#### Summary of Future Directions

Research Area	Goal	Why It Matters
Real-Time Unlearning	On-demand forgetting	Usability in consumer applications
Certified Unlearning	Formal guarantees of data removal	Legal and compliance readiness
Model-Agnostic Frameworks	Apply across architectures	Scalability and accessibility
Unlearning in Foundation Models	Targeted forgetting in massive pre-trained models	Privacy in generative and NLP systems
Robustness to Adversaries	Preventing misuse and manipulation	Security and trust in AI systems
Lifelong Learning with Forgetting	Adaptive learning with memory control	Sustainability and flexibility of AI systems
Legal and Ethical Alignment	Standards for responsible forgetting	Societal adoption and policy integration
Benchmark Creation	Reproducible, fair evaluations	Acceleration of innovation and transparency

## 9. Conclusion

As artificial intelligence continues to evolve and become deeply embedded in societal, commercial, and personal systems, the need for ethical and controllable machine learning practices is more pressing than ever. **Machine Unlearning** emerges as a powerful and necessary paradigm to address this need by enabling the **removal of data influence** from trained models—ensuring compliance with privacy regulations, maintaining user trust, and allowing for safe and accountable AI systems.

This paper provided a comprehensive overview of machine unlearning, beginning with its conceptual foundations and motivations rooted in legal, ethical, and practical concerns. We explored a detailed taxonomy of unlearning techniques, from exact to approximate and certified methods, and presented a spectrum of implementation strategies, including SISA training, gradient reversal, knowledge distillation, and federated approaches.

In addition, we discussed how to evaluate the success of unlearning through a variety of metrics such as forgetting accuracy, retention performance, and privacy guarantees.

Real-world applications across domains like healthcare, finance, social media, and personalized education highlight the wide-reaching relevance of this field.

Despite its promise, machine unlearning still faces considerable challenges—particularly in terms of **scalability, certification, architecture dependence, and adversarial robustness**. Addressing these gaps presents a fertile ground for future research, especially in the context of **pretrained models, lifelong learning, and federated systems**.

Ultimately, integrating unlearning mechanisms into the core of machine learning pipelines will be crucial to developing **responsible, compliant, and user-centric AI systems** for the future.

## References

- [1] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang et al., "Machine Unlearning, " in \*IEEE Symposium on Security and Privacy (S&P) \*, 2021.
- [2] A. Ginart, M. Guan, G. Valiant, and J. Zou, "Making AI Forget You: Data Deletion in Machine Learning, " in

\*Advances in Neural Information Processing Systems (NeurIPS) \*, 2019.

- [3] S. Neel, A. Roth, S. Sharifi-Malvajerdi, "Descent-to-Delete: Gradient-Based Methods for Machine Unlearning, " in \*Algorithmic Learning Theory (ALT) \*, 2021.
- [4] A. Golatkar, A. Achille, and S. Soatto, "Eternal Sunshine of the Spotless Net: Selective Forgetting in Deep Networks, " in \*IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) \*, 2020.
- [5] H. Xu, T. Zhu, L. Zhang, W. Zhou, and P. S. Yu, "Machine Unlearning: A Survey, " arXiv preprint arXiv: 2306.03558, 2023.
- [6] Z. Liu, M. Peng, K-Y Lam, X. Yuan, X. Liu, Y. Jiang, and J. Shen "A Survey on Federated Unlearning: Challenges, Methods, and Future Directions, " arXiv preprint arXiv: 2310.20448, 2024.
- [7] T. Eisenhofer, D. Riepel, V. Chandrasekaran, E. Ghosh, O. Ohrimenko, and N. Papernot, "Verifiable and Provably Secure Machine Unlearning, " arXiv preprint arXiv: 2210.09126, 2022.
- [8] A. Thudi, G. Deza, V. Chandrasekaran, and N. Papernot, "Unrolling SGD: Understanding Factors Influencing Machine Unlearning, " arXiv preprint arXiv: 2109.13398, 2021. .
- [9] C. Guo, T. Goldstein, A. Hannun, and L. van der Maaten, "Certified Data Removal from Machine Learning Models, " in \*Proceedings of the 37th International Conference on Machine Learning (ICML) \*, 2020.
- [10] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership Inference Attacks Against Machine Learning Models, " in \*Proceedings of the IEEE Symposium on Security and Privacy (S&P) \*, 2017, pp.3-18.