

Verifiable and Provably Secure Machine Unlearning

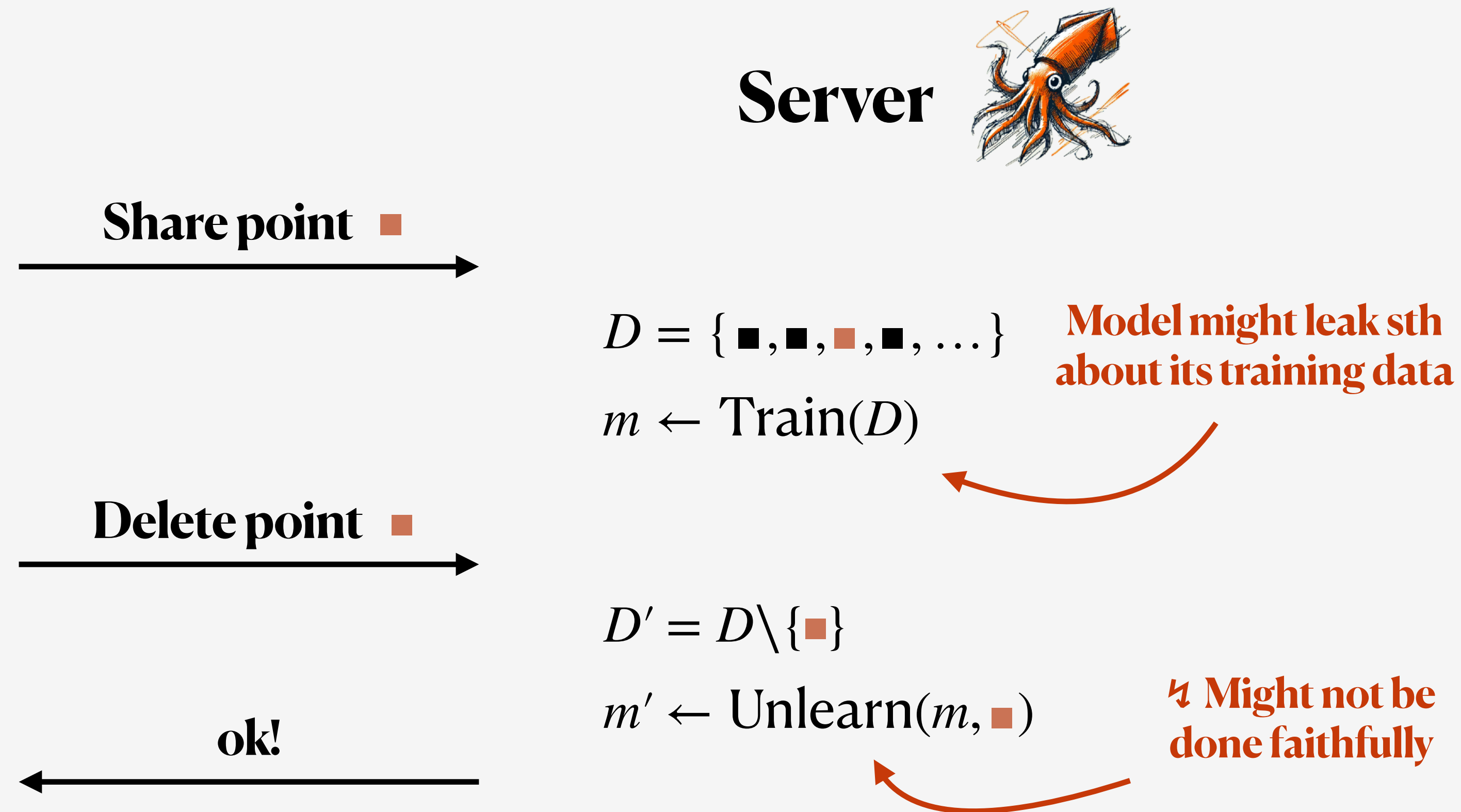
**Thorsten Eisenhofer¹, Doreen Riepel², Varun Chandrasekaran³,
Esha Ghosh⁴, Olya Ohrimenko⁵, and Nicolas Papernot⁶**

¹ BIFOLD & TU Berlin, ² CISPA Helmholtz Center for Information Security,

³ University of Illinois Urbana-Champaign, ⁴ Microsoft Research,

⁵ The University of Melbourne, ⁶ University of Toronto & Vector Institute

Machine Unlearning



Can we trust the server?

Goal: Prove that the unlearning actually happened

Proof via Model Parameters

 Shumailov et al. "Manipulating SGD with Data Ordering Attacks", NeurIPS'21

 Thudi et al. "On the Necessity of Auditable Algorithmic Definitions for Machine Unlearning", USENIX'20

Plausibility forgery

↳ Possible to efficiently construct D'
s.t. $m = m'$ but $D \neq D'$ 

↳ Always possible if D and D' share
the underlying distribution 

If $\blacksquare \in D'$: abort
 $m \leftarrow \text{Train}(D)$
If $m \neq m'$: abort

Server



Share point 



$D = \{\blacksquare, \blacksquare, \blacksquare, \blacksquare, \dots\}$

$m \leftarrow \text{Train}(D)$

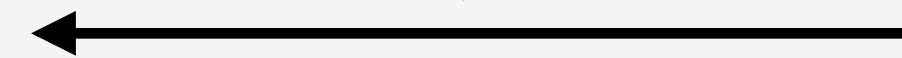
Delete point 



$D' = D \setminus \{\blacksquare\}$

$m' \leftarrow \text{Unlearn}(m, \blacksquare)$

D', m'



One-Shot Verifiable Unlearning

Verifiable Computation

$$y := f(x)$$

Proof that y is the result of evaluating f on x

Share point ■



Delete point ■



D', m', m, ρ



Server



$$D = \{\blacksquare, \blacksquare, \color{red}\blacksquare, \blacksquare, \dots\}$$

$$m \leftarrow \text{Train}(D)$$

$$D' = D \setminus \{\color{red}\blacksquare\}$$

$$m', \rho \leftarrow \text{ProveUnlearn}(m, \color{red}\blacksquare)$$

↗ May be forged

If ■ ∈ D' : abort

If not $\text{VerifyUnlearn}(m, m', \color{red}\blacksquare, \rho)$: abort

How can we prove unlearning?

A Naive Iteration-based Protocol

Server 

Share point ■



$$D = \{\blacksquare, \blacksquare, \color{red}\blacksquare, \blacksquare, \dots\}$$

$$m, \tau \leftarrow \text{ProveTrain}(D)$$

D, m, τ



If not $\text{VerifyTrain}(m, D, \tau)$: abort

Delete point ■



$$D' = D \setminus \{\color{red}\blacksquare\}$$

$$m', \rho \leftarrow \text{ProveUnlearn}(m, \color{red}\blacksquare)$$

D', m', ρ



Verify consistency
of the dataset



If $D' \neq D \cup \{\color{red}\blacksquare\}$: abort

If not $\text{VerifyUnlearn}(m, m', \color{red}\blacksquare, \rho)$: abort

A Naive Iteration-based Protocol

Users $\mathcal{U} \{ \text{pub}, \hat{D}_u \sim \mathcal{D} \}_{u \in \mathcal{U}}$

Server S (pub)

Initialize

if not VerifyInit(pub, com₀, ρ₀):
abort

← com₀, ρ₀

(st_{S,0}, m₀, com₀, ρ₀) ← Init(pub)
D₀⁺ := ∅, U₀⁺ := ∅

i-th iteration

add data points
k-th query

Multiple iterations

← u ∈ \mathcal{U} , d_{i,k} ∈ \hat{D}_u

D_i⁺ := D_{i-1}⁺, U_i⁺ := U_{i-1}⁺

D_i⁺ := D_i⁺ ∪ {(u, d_{i,k})}

remove data points
j-th query

Share only commitment on model and data

← u ∈ \mathcal{U} , d_{i,j} ∈ \hat{D}_u

U_i⁺ := U_i⁺ ∪ {(u, d_{i,j})}

Proof of Training

if not VerifyTraining(pub, com_{i-1}, com_i, ρ_i):
abort

← train: com_i, ρ_i

(st_{S,i}, m_i, com_i, ρ_i) ← ProveTraining(st_{S,i-1}, pub, D_i⁺)
D_i⁺ := ∅

OR Proof of Unlearning

if not VerifyUnlearning(pub, com_{i-1}, com_i, ρ_i):
abort

← unlearn: com_i, ρ_i

(st_{S,i}, m_i, com_i, ρ_i) ← ProveUnlearning(st_{S,i-1}, pub, U_i⁺)

if not VerifyNonMembership(pub, u, d_{i,j}, com_i, π_{u,d_{i,j}}):
abort

← π_{u,d_{i,j}}

for (u, d_{i,j}) ∈ U_i⁺:
π_{u,d_{i,j}} ← ProveNonMembership(st_{S,i}, pub, u, d_{i,j})
U_i⁺ := ∅

If not V

Verify consi
of the dat

If D' ≠ D ∪ {■}: abort

If not VerifyUnlearn(m, m', ■, ρ): abort **Framework for Verifiable Unlearning**

Security Definition

GameUnlearn $_{\mathcal{A}, \mathcal{E}, \Phi_f, \mathcal{D}}(1^\lambda)$

```

00 pub ← Setup(1λ)
01 (k, (u, d), πu,d, {modei: comi, ρi}i∈[0:ℓ];
   {Di}i∈[0:ℓ]) ← (A||E)(pub, aux)
02 # Pre-process
03 Uk+ := Dk-1 \ Dk
04 Parse comi as (comim || comiD) ∀i ∈ [0 : ℓ]

05 # Evaluate winning condition
06 if Commit(pub, Di) = comiD ∀i ∈ [0 : ℓ]           # Datasets
07 and VerifyInit(pub, com0, ρ0)                     # Initialization
08 and VerifyTraining(pub, comi-1, comi, ρi)
   ∀i : modei = train                                   # Training
09 and VerifyUnlearning(pub, comi-1, comi, ρi)
   ∀i : modei = unlearn                                  # Unlearning
10 and VerifyNonMembership(pub, u, d, comk, πu,d)      # Non-Membership
11 and (u, d) ∈ Uk+                                       # Point unlearned
12 and (u, d) ∈ Dℓ and k < ℓ:                          # Point re-added later
13   return 1
14 return 0

```

Run adversary \mathcal{A}

Evaluate output of \mathcal{A}

Definition (Unlearning)

Let λ be the security parameter and consider game GameUnlearn $_{\mathcal{A}, \mathcal{E}, \Phi_f, \mathcal{D}}$ for into-distribution \mathcal{D} is unlearning-secure if for all DPT adversaries \mathcal{A} there exists an extractor \mathcal{E} such that for all benign auxiliary inputs aux

A protocol is unlearning-secure if no efficient adversary exists that can forge an unlearning response.

$$\Pr[\text{GameUnlearn}_{\mathcal{A}, \mathcal{E}, \Phi_f, \mathcal{D}}(1^\lambda) \Rightarrow 1] \leq \text{negl}(\lambda)$$

Prove completeness and security of protocols under this definition!



Adversary wins if they can forge an unlearning response

Take Aways

Verifiable unlearning

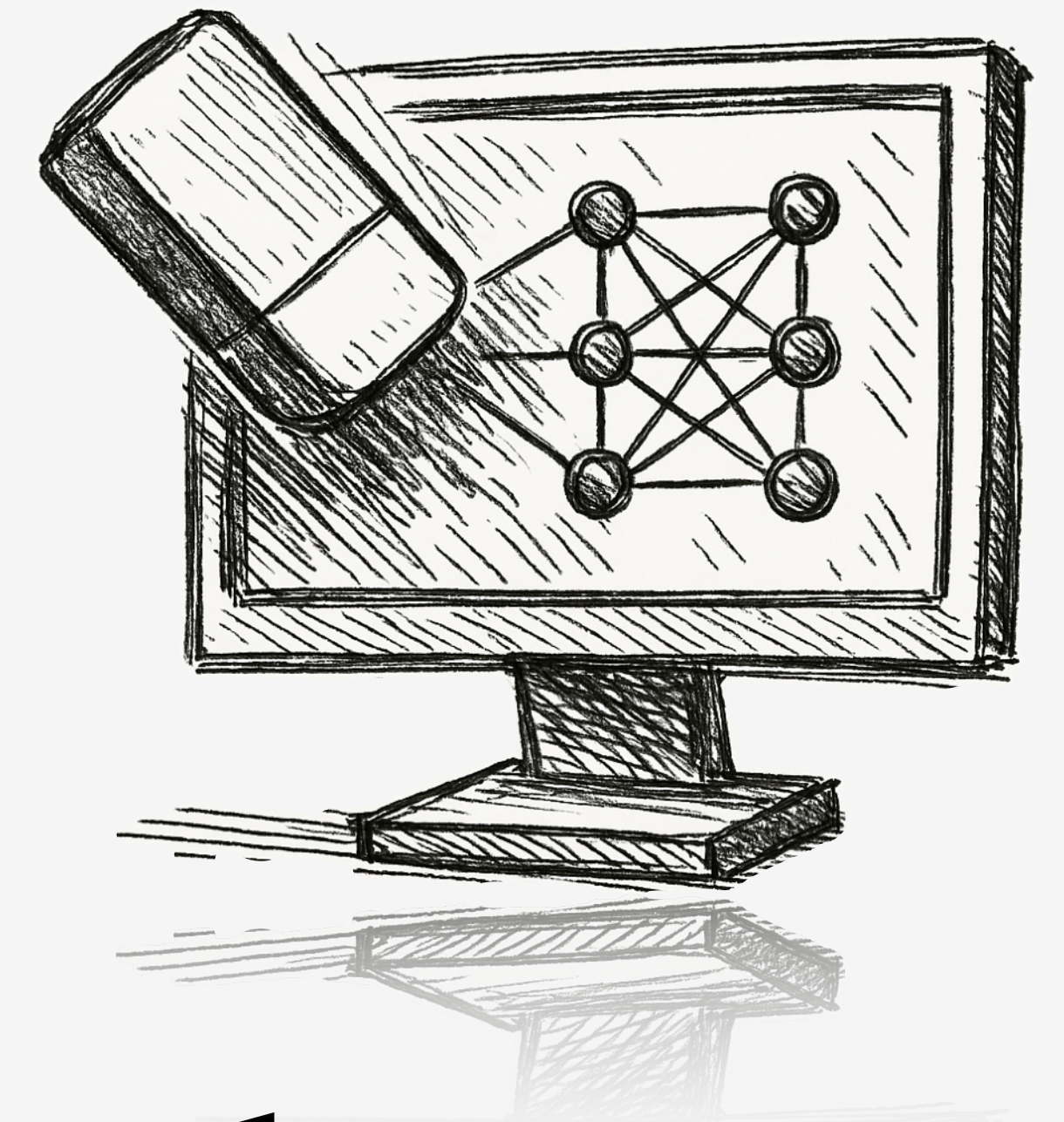
- Verify that unlearning algorithm was executed
- Consider lifecycle of the model

Security definition

- Instantiate protocol under this definition
- Allows to prove completeness and security

Paper and code

github.com/cleverhans-lab/verifiable-unlearning



Thank you!