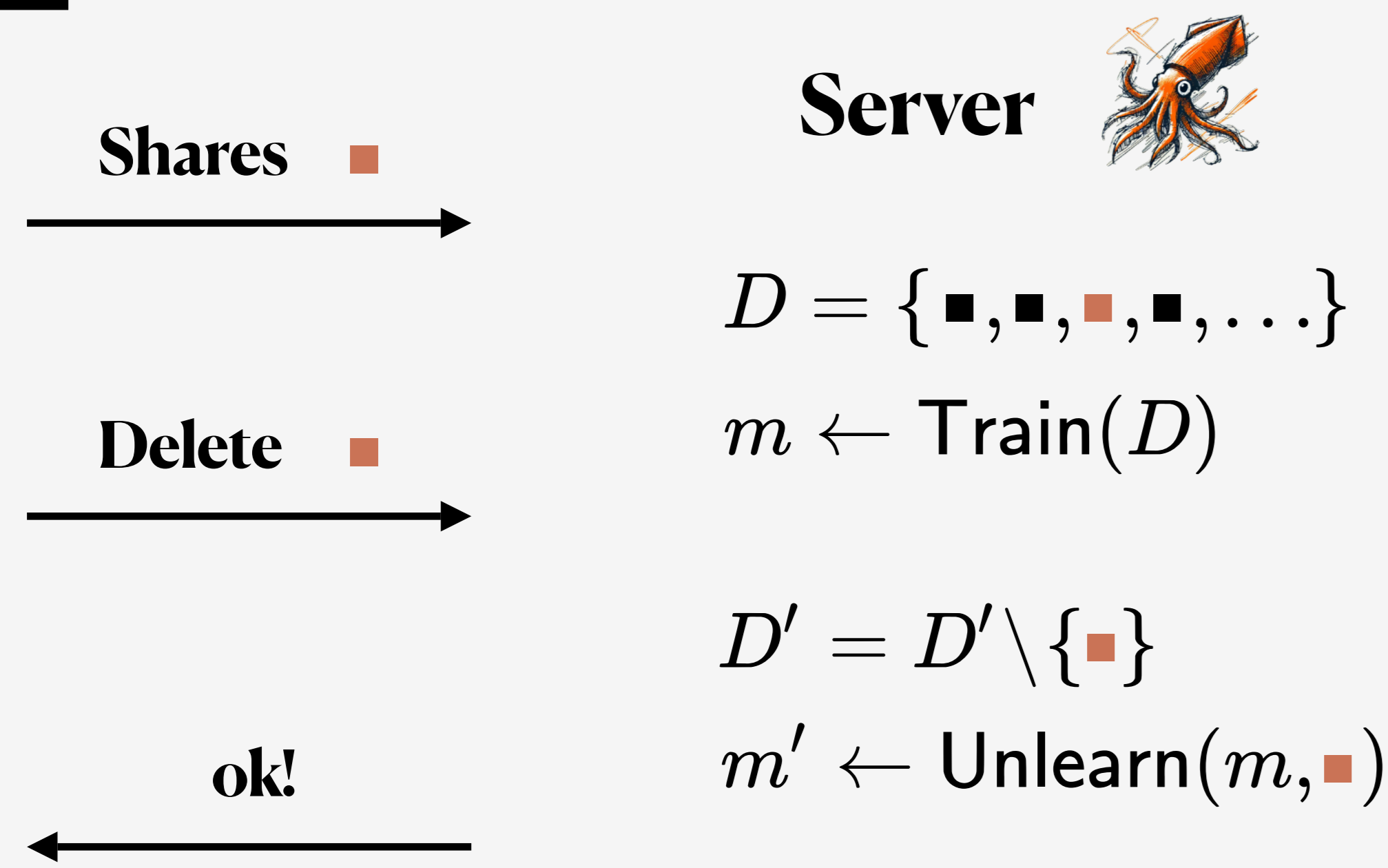# Verifiable and Provably Secure Machine Unlearning

**Thorsten Eisenhofer [1], Doreen Riepel [2], Varun Chandrasekaran [3], Esha Ghosh [4], Olya Ohrimenko [5], and Nicolas Papernot [6]**

[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

## Motivation

**Shares** ◼

**Server** 

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

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

**Delete** ◼

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

**ok!**

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

*But can we trust the server?*

**Goal**: Prove that unlearning was performed correctly!

## Unlearning Framework

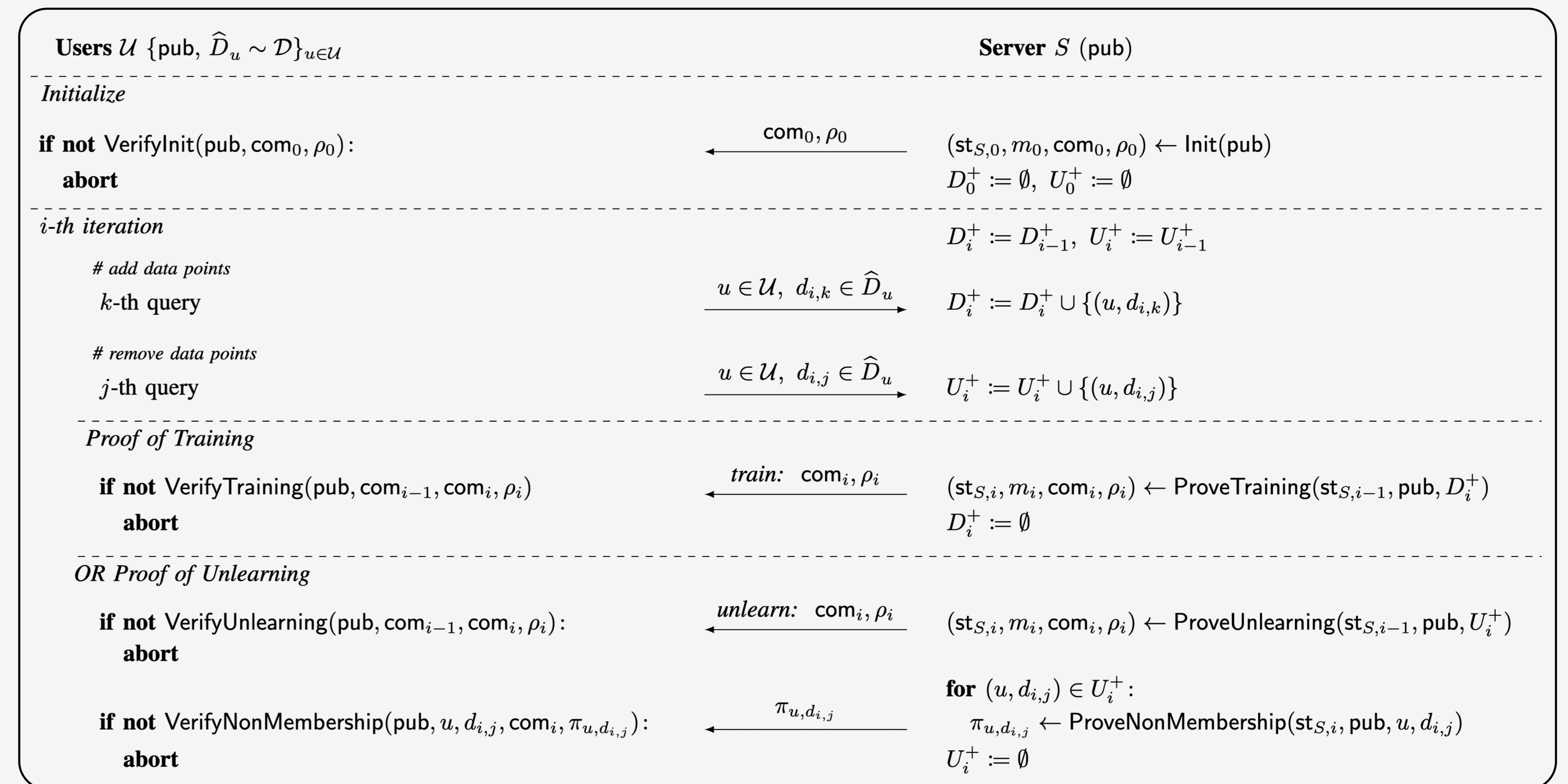Proof via model parameters not sufficient. Can efficiently construct dataset $D'$ s.t. $D' \neq D$ but $m = m'$ ↯

### Verifiable Unlearning

**Proof of unlearning**
Verify correct execution of unlearning algorithm

**Proof of training**
Consider full lifecycle of the model

| Users $\mathcal{U}$ $\{\mathsf{pub}, \widehat{D}_u \sim \mathcal{D}\}_{u \in \mathcal{U}}$ | | Server $S$ (pub) |
|---|---|---|
| *Initialize* | | |
| **if not** $\mathsf{VerifyInit}(\mathsf{pub}, \mathsf{com}_0, \rho_0)$: | $\xleftarrow{\mathsf{com}_0, \rho_0}$ | $(\mathsf{st}_{S,0}, m_0, \mathsf{com}_0, \rho_0) \leftarrow \mathsf{Init}(\mathsf{pub})$ |
|   **abort** | | $D_0^+ \coloneqq \emptyset, \; U_0^+ \coloneqq \emptyset$ |
| *i-th iteration* | | $D_i^+ \coloneqq D_{i-1}^+, \; U_i^+ \coloneqq U_{i-1}^+$ |
|   # add data points | | |
|   $k$-th query | $\xrightarrow{u \in \mathcal{U}, \; d_{i,k} \in \widehat{D}_u}$ | $D_i^+ \coloneqq D_i^+ \cup \{(u, d_{i,k})\}$ |
|   # remove data points | | |
|   $j$-th query | $\xrightarrow{u \in \mathcal{U}, \; d_{i,j} \in \widehat{D}_u}$ | $U_i^+ \coloneqq U_i^+ \cup \{(u, d_{i,j})\}$ |
| *Proof of Training* | | |
|   **if not** $\mathsf{VerifyTraining}(\mathsf{pub}, \mathsf{com}_{i-1}, \mathsf{com}_i, \rho_i)$ | $\xleftarrow{\textit{train: } \; \mathsf{com}_i, \rho_i}$ | $(\mathsf{st}_{S,i}, m_i, \mathsf{com}_i, \rho_i) \leftarrow \mathsf{ProveTraining}(\mathsf{st}_{S,i-1}, \mathsf{pub}, D_i^+)$ |
|   **abort** | | $D_i^+ \coloneqq \emptyset$ |
| *OR Proof of Unlearning* | | |
|   **if not** $\mathsf{VerifyUnlearning}(\mathsf{pub}, \mathsf{com}_{i-1}, \mathsf{com}_i, \rho_i)$: | $\xleftarrow{\textit{unlearn: } \; \mathsf{com}_i, \rho_i}$ | $(\mathsf{st}_{S,i}, m_i, \mathsf{com}_i, \rho_i) \leftarrow \mathsf{ProveUnlearning}(\mathsf{st}_{S,i-1}, \mathsf{pub}, U_i^+)$ |
|   **abort** | | |
|   **if not** $\mathsf{VerifyNonMembership}(\mathsf{pub}, u, d_{i,j}, \mathsf{com}_i, \pi_{u,d_{i,j}})$: | $\xleftarrow{\pi_{u,d_{i,j}}}$ | **for** $(u, d_{i,j}) \in U_i^+$: |
|   **abort** | |   $\pi_{u,d_{i,j}} \leftarrow \mathsf{ProveNonMembership}(\mathsf{st}_{S,i}, \mathsf{pub}, u, d_{i,j})$ |
| | | $U_i^+ \coloneqq \emptyset$ |

## Security Definition

```
GameUnlearn_{A,E,Φ_f,D}(1^λ)
00  pub ← Setup(1^λ)
01  (k, (u,d), π_{u,d}, {mode_i: com_i, ρ_i}_{i∈[0:ℓ]};
                       {D_i}_{i∈[0:ℓ]}) ← (A|E)(pub, aux)

02  # Pre-processing
03  U_k^+ := D_{k-1} \ D_k
04  Parse com_i as (com_i^m ‖ com_i^D) ∀i ∈ [0:ℓ]

05  # Evaluate winning condition
06  if Commit(pub, D_i) = com_i^D ∀i ∈ [0:ℓ]        # Datasets
07     and VerifyInit(pub, com_0, ρ_0)              # Initialization
08     and VerifyTraining(pub, com_{i-1}, com_i, ρ_i)
                       ∀i : mode_i = train           # Training
09     and VerifyUnlearning(pub, com_{i-1}, com_i, ρ_i)
                       ∀i : mode_i = unlearn          # Unlearning
10     and VerifyNonMembership(pub, u, d, com_k, π_{u,d})  # Non-Membership
11     and (u,d) ∈ U_k^+                              # Point unlearnt
12     and (u,d) ∈ D_ℓ and k < ℓ:                    # Point re-added later
13     return 1
14  return 0
```



**Adversary wins if they can forge an unlearning response**

### Definition (Unlearning)

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

## Instantiation

Prove correct execution of training and unlearning algorithm using techniques from verifiable computation.

**Verifiable computation**

$$\mathbf{y} := \mathbf{f(x)}$$

*"Proof that $\mathbf{y}$ is the result of evaluating $\mathbf{f}$ on $\mathbf{x}$"*

```
C_U(public h_{st_{f,i}}, h_{st_{f,i-1}}, h_{m_i}, h_{D_i}, h_{D_{i-1}}, h_{U_i}, h_{U_{i-1}},
       private st_{f,i-1}, H_{D_{i-1}}, U_i^+)
00  # Check input set of hashed training data records
01  if h_{D_{i-1}} ≠ HashData(H_{D_{i-1}}):
02     return false
03  # Update and check set of hashed unlearnt data records and training data records
04  H_{U_i^+} := {HashDataRecord(u,d)}_{(u,d)∈U_i^+}
05  H_{D_i} := H_{D_{i-1}} \ H_{U_i^+}
06  if h_{U_i} ≠ AppendHashData(h_{U_{i-1}}, H_{U_i^+}) or h_{D_i} ≠ HashData(H_{D_i}):
07     return false
08  # Check input state, perform unlearning and check outputs
09  if h_{st_{f,i-1}} ≠ HashState(st_{f,i-1}):
10     return false
11  (st_{f,i}, m_i) := f_U(st_{f,i-1}, U_i^+)
12  if h_{st_{f,i}} ≠ HashState(st_{f,i}) or h_{m_i} ≠ HashModel(m_i):
13     return false
14  return true
```