

# SECURE VAULT SCHEME IN THE CLOUD OPERATING MODEL

Rishiraj Bhattacharyya, Avradip Mandal, **Meghna Sengupta**



UNIVERSITY OF  
BIRMINGHAM



**Zfense Labs**



THE UNIVERSITY  
of EDINBURGH



# OUTLINE

- Data Privacy Vaults
  - *What do they do?*
  - *Who's making them? Who's using them?*
- Our framework
- Building Block: Tokenization Scheme
  - *Definition*
  - *Security Notions*
- Complete Vault Scheme Construction
  - *Construction*

# 1

## DATA-PRIVACY VAULTS

- What do they do?
- Who is using them?



# Data-Privacy Vault

A literal vault (safe) of information -



- Users can store their private information
- Generate access passes
  - ▷ give access to selected parties
- Data hidden from others



# Data-Privacy Vault

## Access Levels

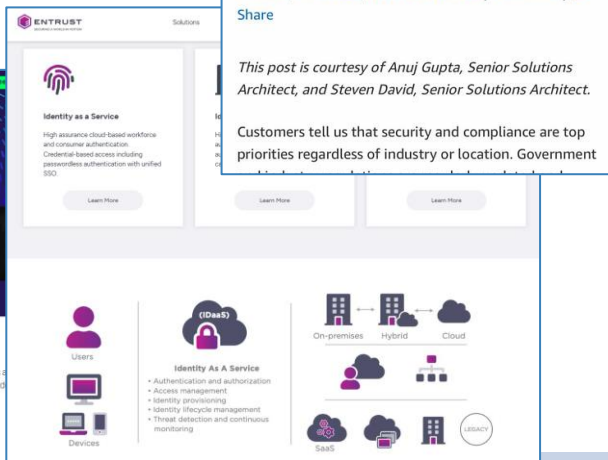
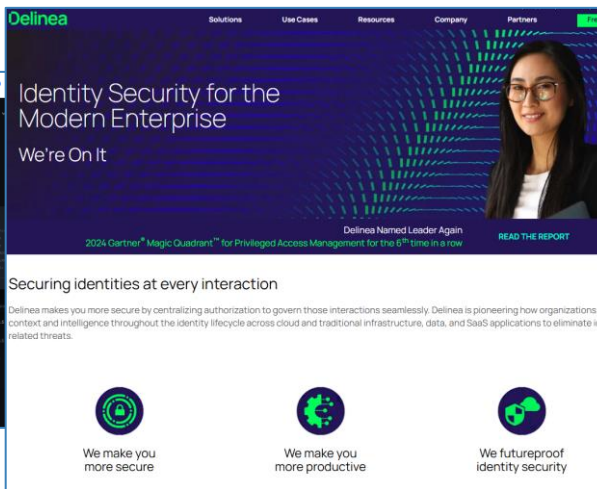
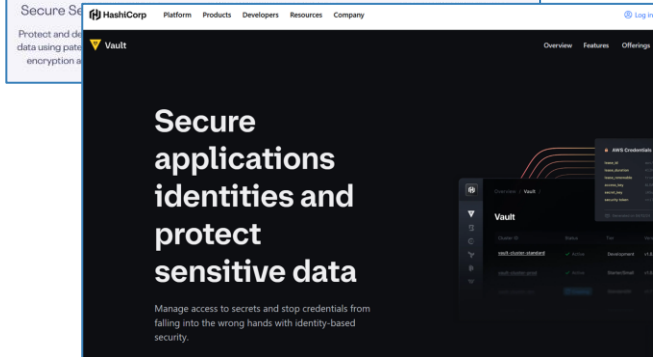
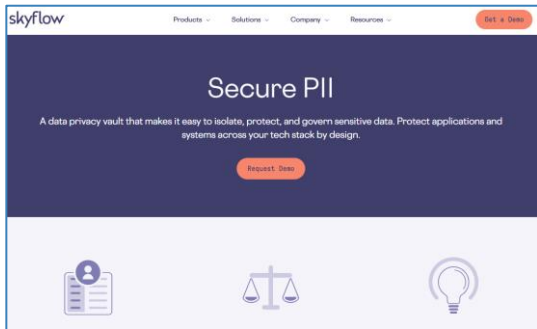
- Different levels of access passes
  - ▶ provides access to different data
- Eg. Level 1: Non-sensitive data  
Level 2: Both non-sensitive and sensitive data





# Who's making & using them?

A lot of big companies!





# Who's making & using them?

A lot of big companies!

increasingly being used for  
storing information used  
for training LLM models

# 2

## OUR FRAMEWORK





# Our Framework

3 parties:



TRUSTED  
SERVER



USER



DATABASE



ANALYST



# Our Framework - Functionality

## ■ STORE -

- ▷ User stores the information and gets back access passes

## ■ ACCESS -

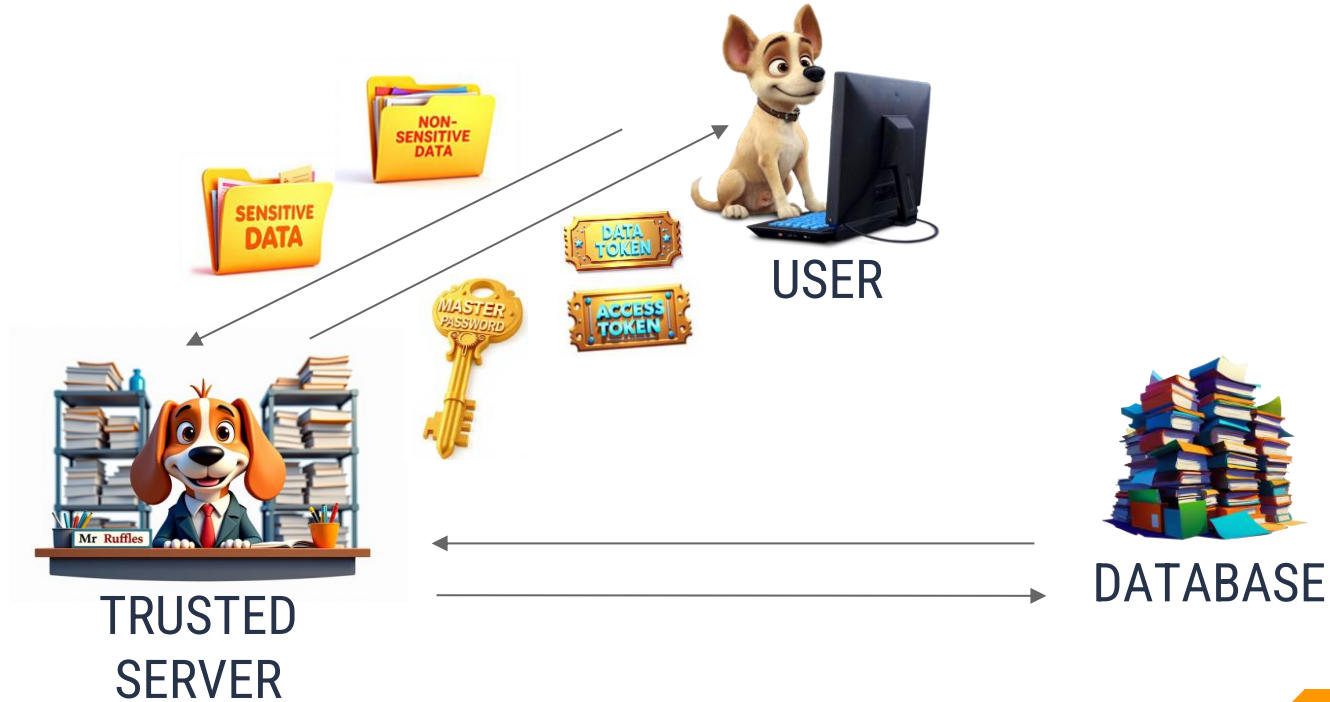
- ▷ Level 1 access to non-sensitive information

## ■ RETRIEVE -

- ▷ Level 2 access to all the information



# STORE

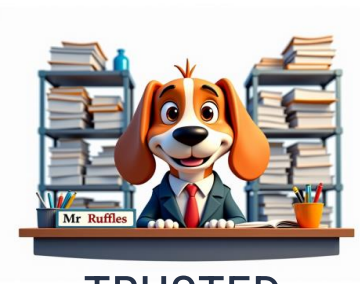




# RETRIEVE



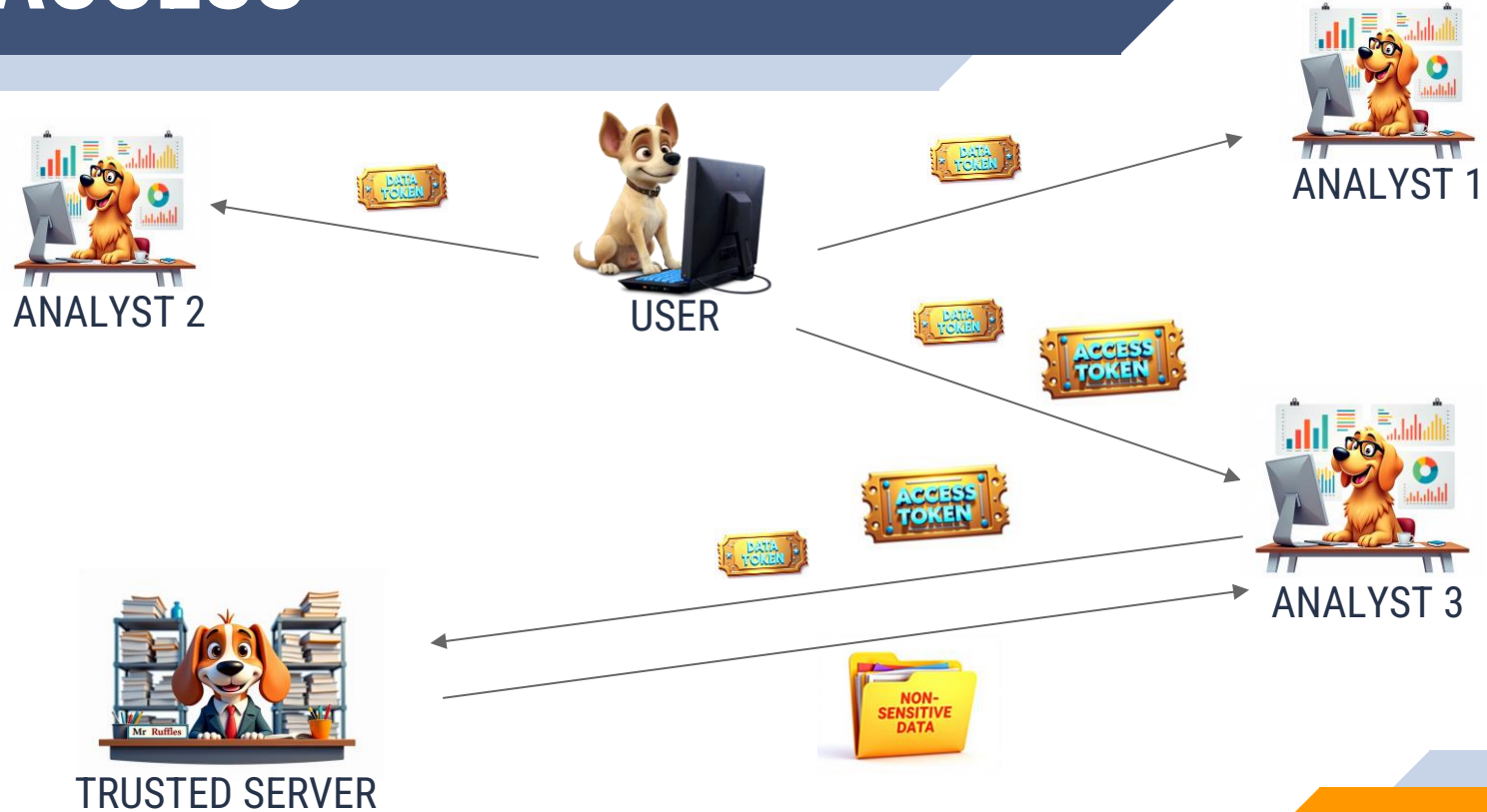
USER



TRUSTED  
SERVER



# ACCESS





# Security Notions

## ■ Non-Sensitive Information:

No access Token -> No non-sensitive information

IND-CATA

## ■ Sensitive Information:

No master password -> No sensitive information  
(even with access token)

IND-SIA

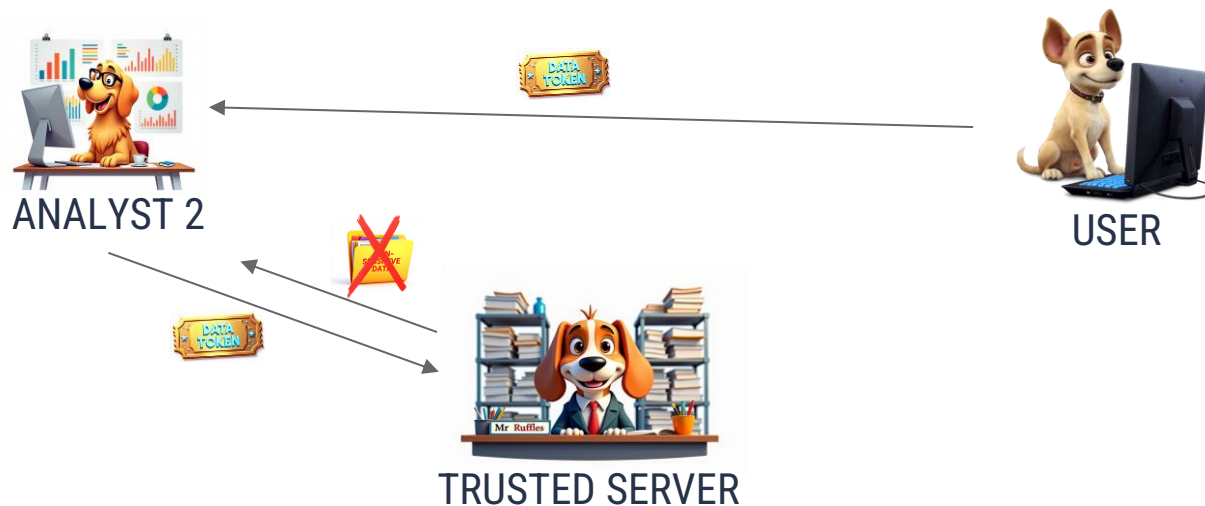


# Security Notions

IND-CATA

## Non-Sensitive Information:

No access Token -> No non-sensitive information





# Security Notions

IND-CATA

## Non-Sensitive Information:

No access Token  $\rightarrow$  No non-sensitive information

Game IND-CATA

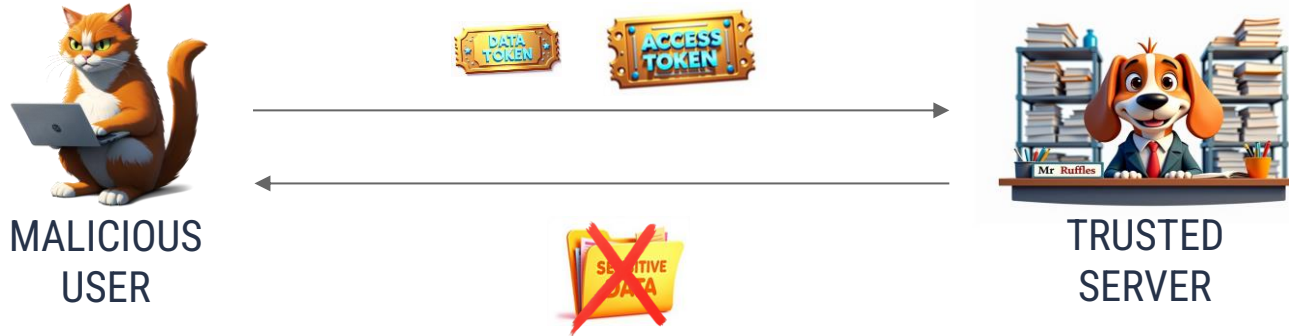
```
1 :  $b \leftarrow_{\$} \{0, 1\}$ 
2 :  $PP \leftarrow_{\$} \text{Param}(1^\lambda)$ 
3 :  $(state, m_0 = (m_{s,0}, m_{ns,0}), m_1 = (m_{s,1}, m_{ns,1}))$   

    $\leftarrow \mathcal{A}^{\text{Store}(PP, \cdot), \text{Access}(PP, \cdot)}(PP, 1^\lambda)$ 
4 : if  $m_{ns,0} = m_{ns,1}$  :
5 :   return 0
6 :  $(m\text{-tok}^*, data\text{-tok}^*) = \text{Store}(\perp, PP, m_b)$ 
7 :  $access\text{-tok}^* = \mathcal{H}(10, m\text{-tok}^*)$ 
8 :  $b' \leftarrow \mathcal{A}^{\text{Store}(PP, \cdot), \text{Access}(PP, \cdot)}(PP, access\text{-tok}^*, state)$ 
9 : return  $(b = b')$ 
```





# Security Notions



## ■ Sensitive Information:

No master password -> No sensitive information  
(even with access token)

IND-SIA



# Security Notions

## Game IND-SIA

```
1 :  $b \leftarrow_{\$} \{0, 1\}$   
2 :  $PP \leftarrow_{\$} \text{Param}(1^\lambda)$   
3 :  $(state, m_0 = (m_{s,0}, m_{ns,0}), m_1 = (m_{s,1}, m_{ns,1}))$   
    $\leftarrow \mathcal{A}^{\text{Store}(PP, \cdot), \text{Access}(PP, \cdot)}(PP, 1^\lambda, DB)$   
4 :  $((m\text{-tok}^*, data\text{-tok}^*), DB') = \text{Store}(\perp, PP, m_{ns,b}, DB)$   
5 :  $b' \leftarrow \mathcal{A}^{\text{Store}(PP, \cdot), \text{Access}(PP, \cdot)}(PP, state, DB')$   
6 : return  $(b = b')$ 
```

## ■ Sensitive Information:

No master password -> No sensitive information  
(even with access token)



IND-SIA

# 3

## **BUILDING BLOCK: TOKENIZATION SCHEME**

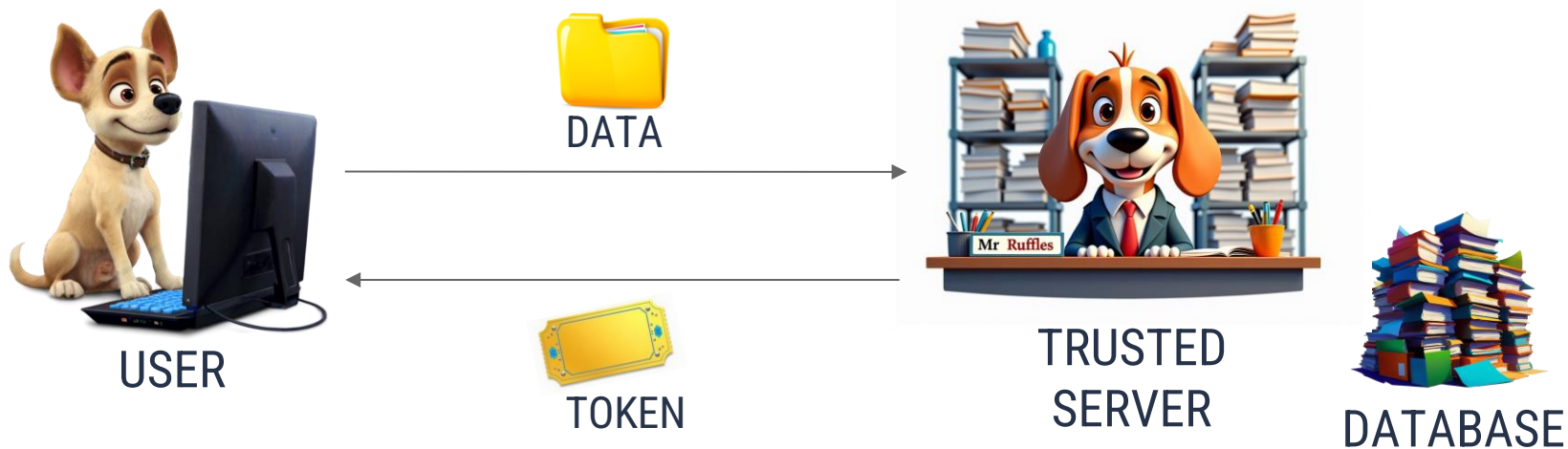


# Tokenization Scheme - Model

- Cloud Operating Model
  - ▷ Database may get compromised
- LLM-training friendly
  - ▷ Must be deterministic (same token for same data)
- The user should not have to store a secret key

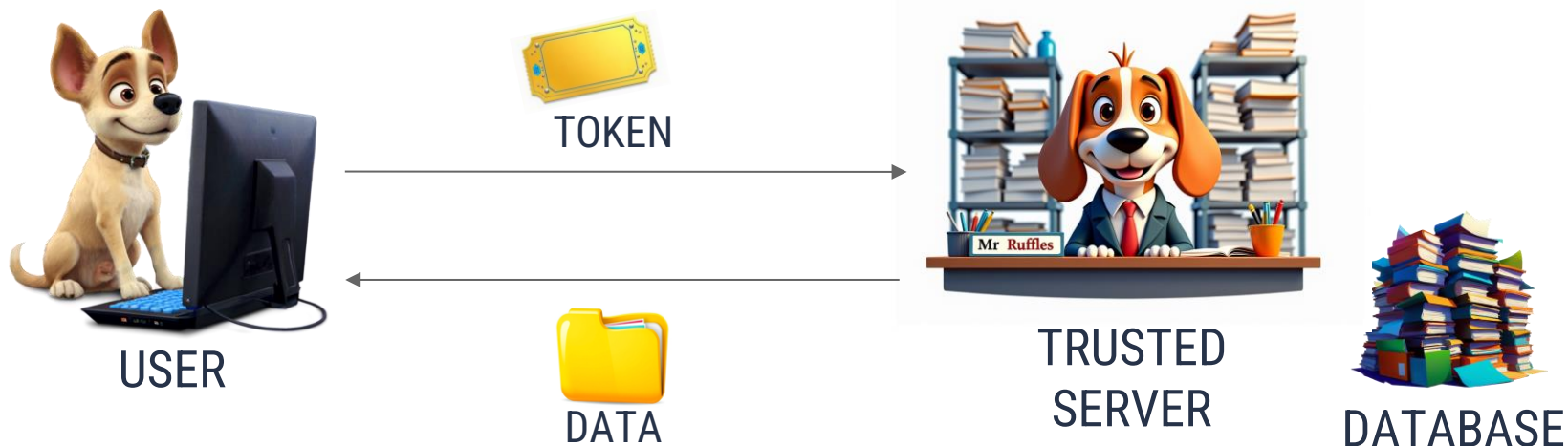


# Tokenization Scheme





# Tokenization Scheme





# Tokenization Scheme

## ■ Tokenization Procedure:

- ▷ Takes message  $m$  as input
- ▷ Generates a token  $tok$  and adds an entry to the database
- ▷ Returns  $tok$

## ■ DeTokenization Procedure:

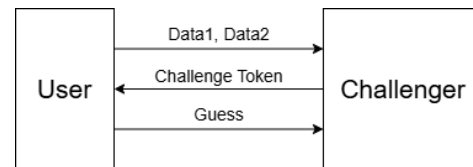
- ▷ Takes token  $tok$  as input
- ▷ If it exists in database, returns corresponding message  $m$ , else returns  $\perp$ .



# Security Notions

IND-CDA

## Chosen Data Attack



## Chosen Distribution Database Attack

IND-CDDA

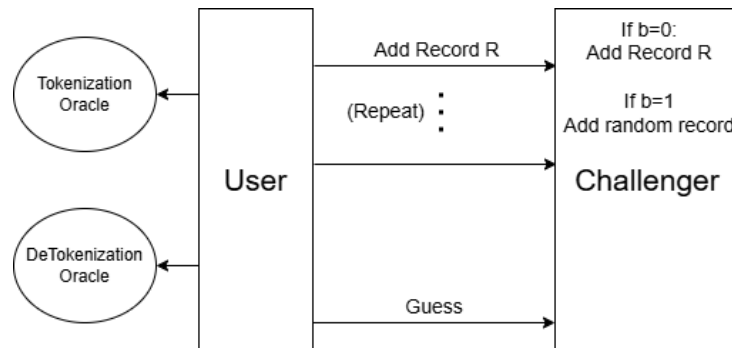




# Security Notions

IND-CDA

## Chosen Data Attack



## Chosen Distribution Database Attack

IND-CDDA



# Our Construction

$\text{Tok}(m, \overline{m}, k_1)$  // for basic construction  $k_1 = \perp, m = \overline{m}$

```
1: if  $k_1 = \perp$ 
2:    $k_1 \leftarrow \mathcal{H}_1(m)$ 
3: if  $\exists(c_1 || t_1 || c_2 || t_2) \in DB$  s.t.  $t_1 = \mathcal{H}_2(m)$ 
4:    $tok \leftarrow \text{Dec}(k_1, c_1)$ 
5:   return  $tok$ 
6:  $tok \leftarrow \{0, 1\}^\tau$ 
7:  $c_1 \leftarrow \text{Enc}(k_1, tok)$ 
8:  $t_1 \leftarrow \mathcal{H}_2(m)$ 
9:  $k_2 \leftarrow \mathcal{H}_3(tok)$ 
10:  $c_2 \leftarrow \text{Enc}(k_2, \overline{m})$ 
11:  $t_2 \leftarrow \mathcal{H}_4(tok)$ 
12:  $d = c_1 || t_1 || c_2 || t_2$ 
13:  $DB = DB \cup \{d\}$ 
14: return  $tok$ 
```

$\text{DeTok}(tok)$

```
1: if  $\exists(c_1 || t_1 || c_2 || t_2) \in DB$  s.t.  $t_2 = \mathcal{H}_4(tok)$ 
2:    $m \leftarrow \text{Dec}(\mathcal{H}_3(tok), c_2)$ 
3:   return  $m$ 
4: else return  $\perp$ 
```

token generated randomly,

$c_1 = \text{Enc}(\mathcal{H}(m), tok), \quad t_1 = \mathcal{H}'(m)$   
 $c_2 = \text{Enc}(\mathcal{H}(tok), m), \quad t_2 = \mathcal{H}'(tok)$

store  $\{c_1, t_1, c_2, t_2\}$  as one record



# Our Construction

$\text{Tok}(m, \overline{m}, k_1)$  // for basic construction  $k_1 = \perp, m = \overline{m}$

```
1: if  $k_1 = \perp$ 
2:    $k_1 \leftarrow \mathcal{H}_1(m)$ 
3:   if  $\exists (c_1 || t_1 || c_2 || t_2) \in DB \text{ s.t. } t_1 = \mathcal{H}_2(m)$ 
4:      $tok \leftarrow \text{Dec}(k_1, c_1)$ 
5:     return  $tok$ 
6:    $tok \leftarrow \{0, 1\}^\tau$ 
7:    $c_1 \leftarrow \text{Enc}(k_1, tok)$ 
8:    $t_1 \leftarrow \mathcal{H}_2(m)$ 
9:    $k_2 \leftarrow \mathcal{H}_3(tok)$ 
10:   $c_2 \leftarrow \text{Enc}(k_2, \overline{m})$ 
11:   $t_2 \leftarrow \mathcal{H}_4(tok)$ 
12:   $d = c_1 || t_1 || c_2 || t_2$ 
13:   $DB = DB \cup \{d\}$ 
14:  return  $tok$ 
```

$\text{DeTok}(tok)$

```
1: if  $\exists (c_1 || t_1 || c_2 || t_2) \in DB \text{ s.t. } t_2 = \mathcal{H}_4(tok)$ 
2:    $m \leftarrow \text{Dec}(\mathcal{H}_3(tok), c_2)$ 
3:   return  $m$ 
4: else return  $\perp$ 
```

if a record exists with  
 $t_2 = \mathcal{H}'(tok)$ ,

return  
 $m = \text{Dec}(\mathcal{H}(tok), c_2)$

Our construction is both IND-CDA and IND-CDDA secure

# 4

## FULL VAULT SCHEME CONSTRUCTION



# Vault Scheme Construction

- We use the Tokenization Scheme as a building block
- For simplicity of presentation, assume that the message space is same as the token space.



# Vault Scheme Construction - STORE

$\text{Store}(m\text{-tok}, PP, M = (M_s, M_{ns}))$

```
1 :  if  $m\text{-tok} = \perp$ 
2 :     $m\text{-tok} \leftarrow \$ \{0, 1\}^\tau$ 
3 :     $k_0 = \mathcal{H}(00, m\text{-tok}, M)$ 
4 :     $k_1 = \mathcal{H}(01, m\text{-tok})$ 
5 :     $k_2 = \mathcal{H}(10, m\text{-tok})$ 
6 :     $k_3 = \mathcal{G}(M)$ 
7 :     $C_1 = \text{Enc}(k_1, M_s)$ 
8 :     $\overline{M} = \text{Enc}(k_2, M_{ns})$ 
9 :     $(data\text{-tok}, d) = \text{Tok}(k_0, \overline{M}, k_3)$ 
10 :    $DB = DB \cup (C_1, d)$ 
11 :   return  $(m\text{-tok}, data\text{-tok})$ 
```

n-s data is  
encrypted  
with the  
access  
token as  
key

if token does not exist beforehand, it is generated.

two hashes of the master token, one for sensitive, one for non-sensitive data

the one for the non-sensitive data is used as the access token

the non-sensitive data ciphertext is tokenized and that token is used as the data-token



# Vault Scheme Construction - ACCESS

$\text{Access}(\text{access-tok}, \text{data-tok})$

1 :  $\overline{M} = \text{DeTok}^{DB}(\text{data-tok})$

2 :  $M_{ns} = \text{Dec}(\text{access-tok}, \overline{M})$

3 : **return**  $M_{ns}$

access needs both data token and access token

data-token is the token for the tokenization scheme

access token is the key for the decryption

# Vault Scheme Construction - RETRIEVE

Retrieve( $m\text{-tok}$ ,  $data\text{-tok}$ )

1 :  $((C_1, d), \overline{M}) = \text{DeTok}^{DB}(data\text{-tok})$

2 :  $k_1 = \mathcal{H}(0, m\text{-tok})$

3 :  $k_2 = \mathcal{H}(1, m\text{-tok})$

4 :  $M_{ns} = \text{Dec}(k_2, \overline{M})$

5 :  $M_s = \text{Dec}(k_1, C_1)$

6 :     **return**  $M = (M_s, M_{ns})$

7 :     **else return**  $\perp$

retrieve needs both master password and access token

data-token is the token for the tokenization scheme

derive the access token and use it to access the non-sensitive data

use master password to decrypt sensitive data



- The IND-CATA of the vault construction follows from -
  - ▷ Randomness of the hash function
  - ▷ IND-CDA and IND-CDDA security of the underlying Tokenization scheme.
- The IND-SIA follows from -
  - ▷ randomness of the hash functions
  - ▷ IND-CPA security of encryption scheme
  - ▷ IND-CDDA of tokenization scheme



Thank You!