

# Reducing Metadata Leakage from Encrypted Files and Communication with PURBs

Kirill Nikitin<sup>\*</sup>, Ludovic Barman<sup>\*</sup>, Wouter Lueks, Matthew Underwood, Jean-Pierre Hubaux, and Bryan Ford

*École polytechnique fédérale de Lausanne (EPFL)*

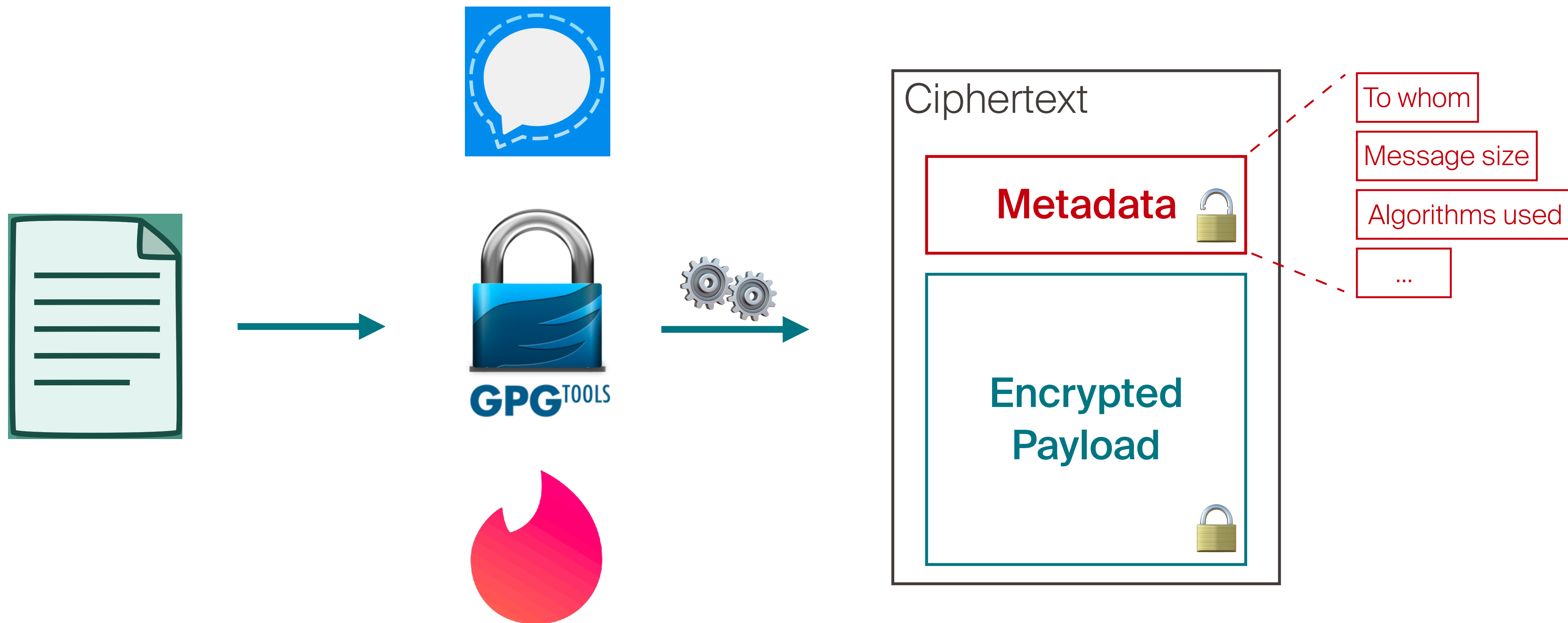
<sup>\*</sup>Shared first authorship

 @ni\_kirill

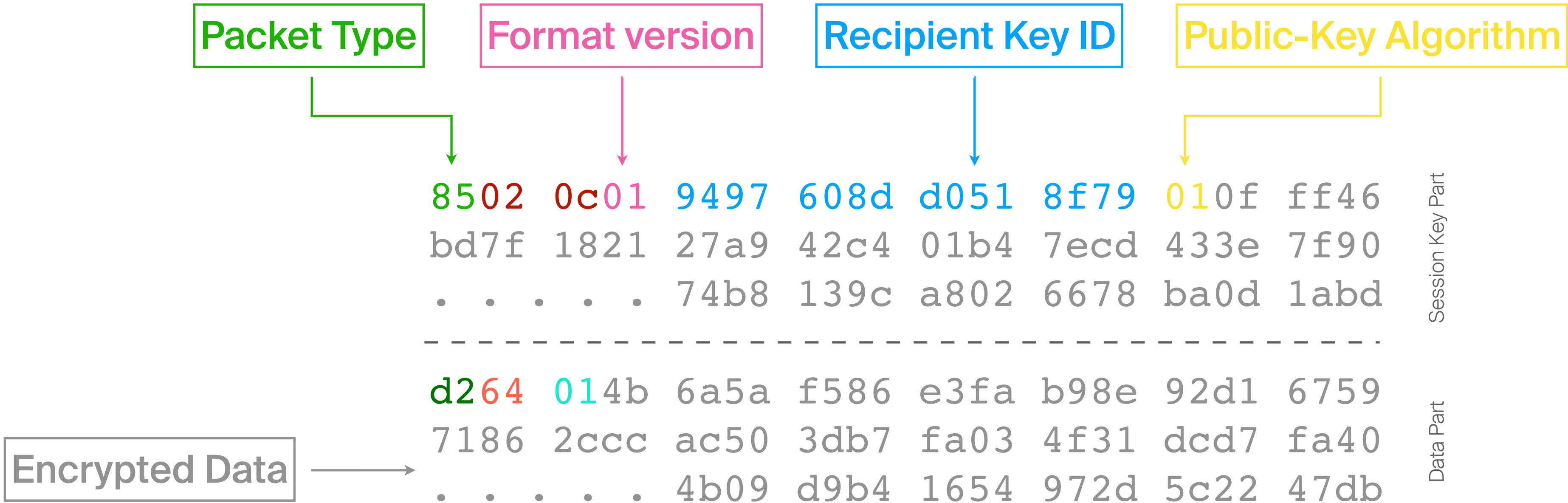
 @lbarman\_ch

[Dog video]

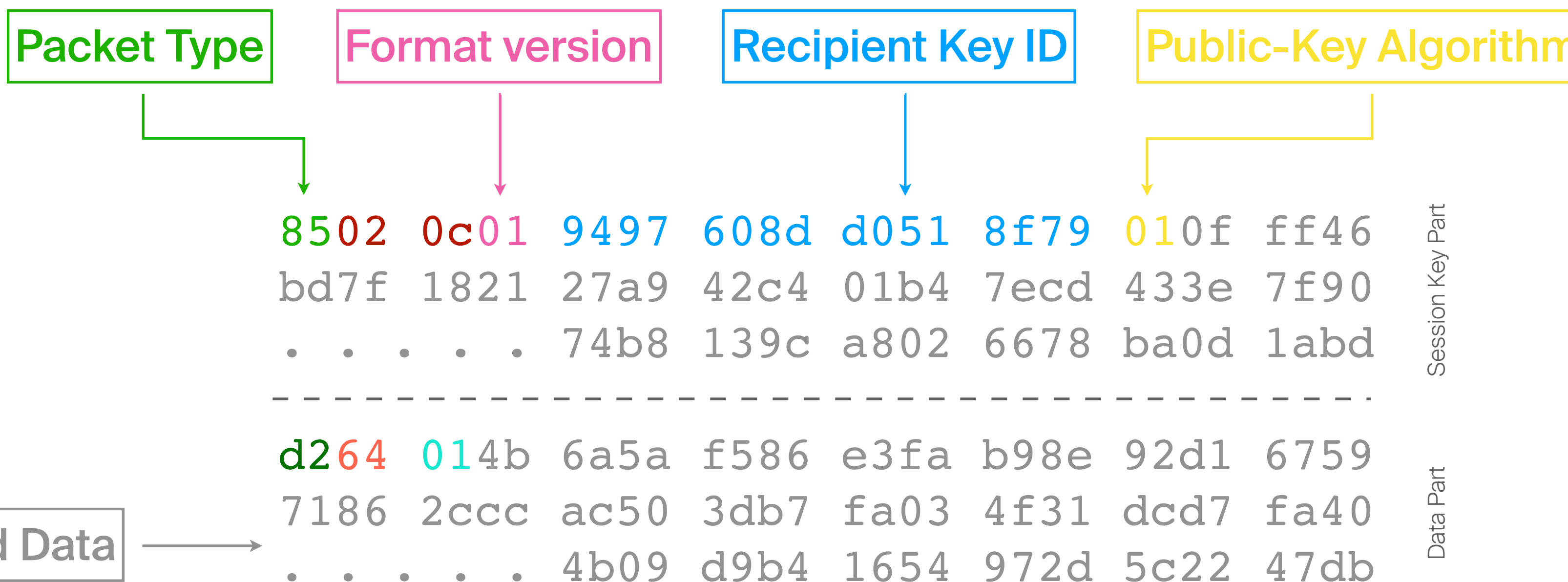
# Ciphertexts Expose Metadata in Clear



# OpenPGP Packet Format



# OpenPGP Packet Format

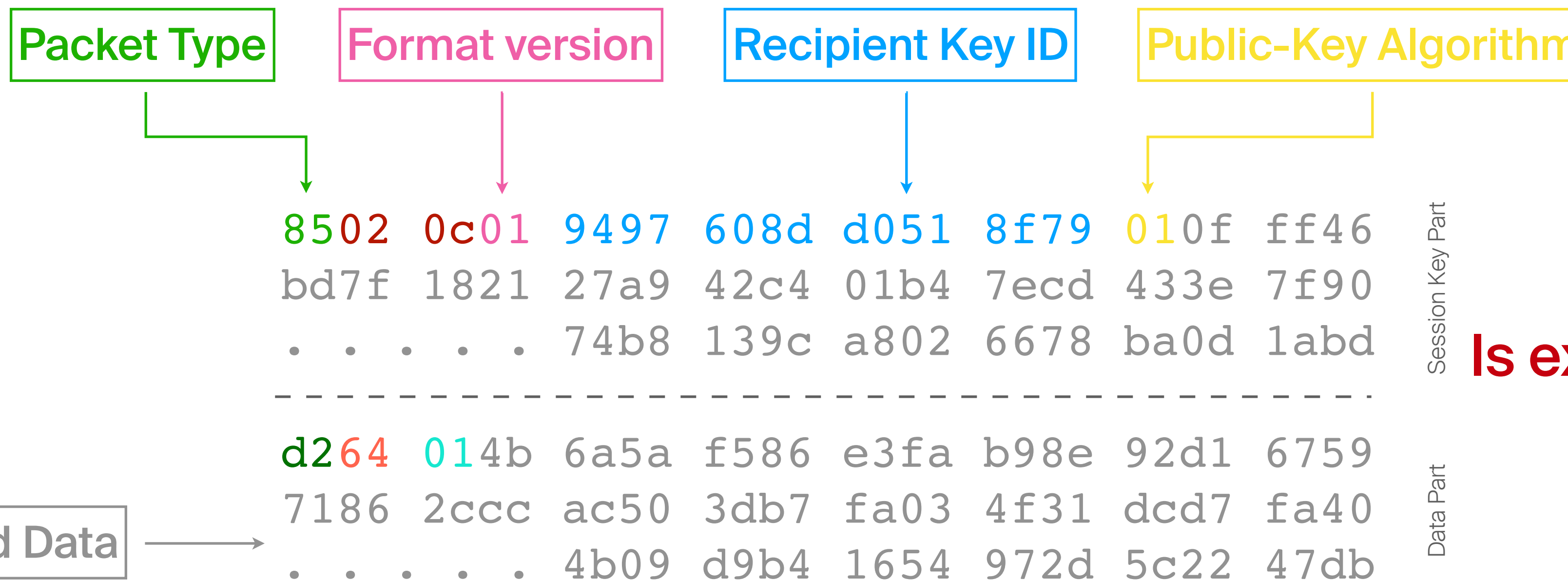


A message to [the King of Sweden](#) encrypted with [RSA-512](#) using an [outdated OpenPGP format](#)??

**Small key? Outdated format? I might crack it!**



# OpenPGP Packet Format



**Is exposing metadata necessary?**

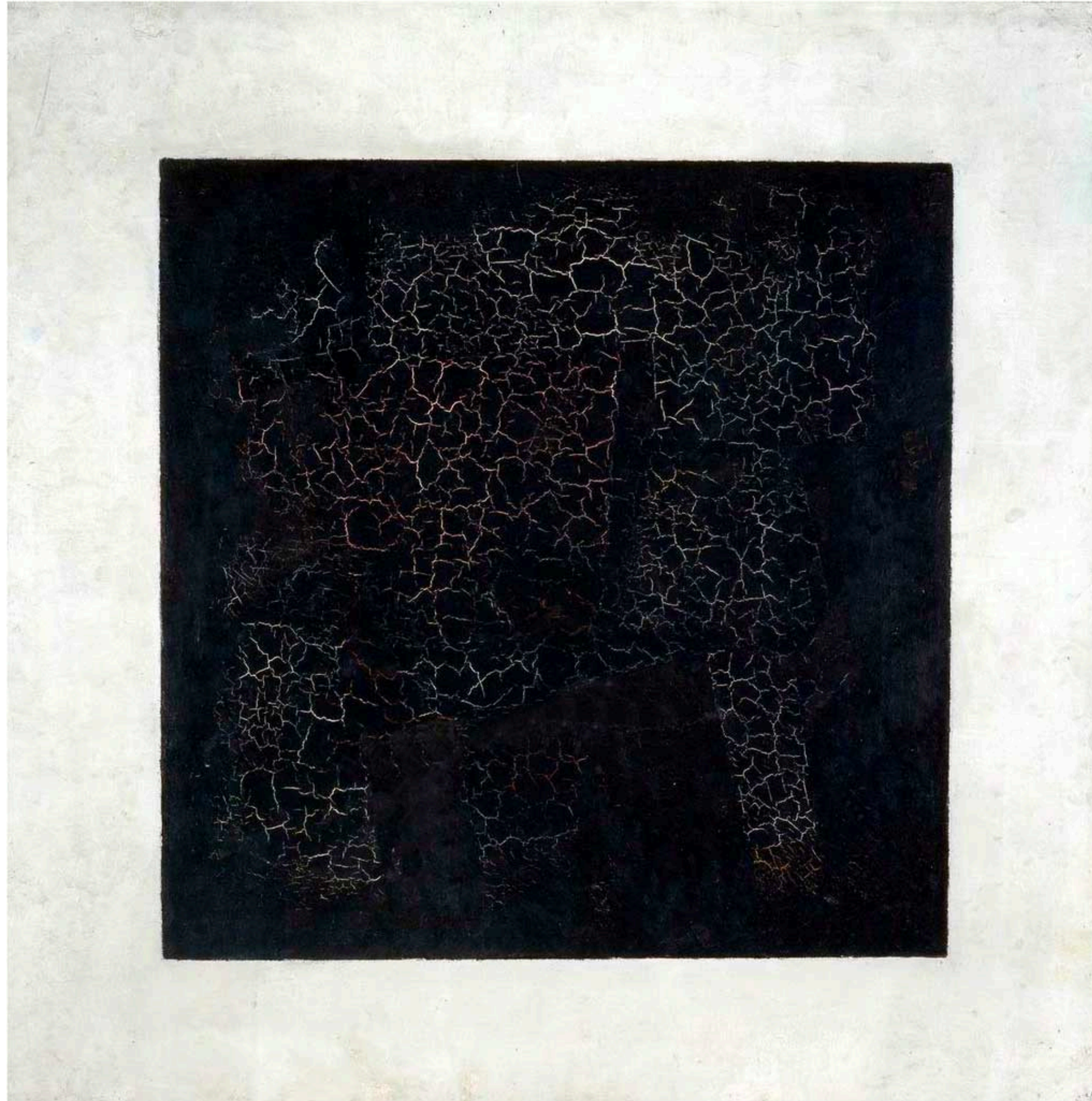
A message to the King of Sweden encrypted with RSA-512 using an outdated OpenPGP format??

**Small key? Outdated format? I might crack it!**





# What If We Stripped Off All the Metadata?



“Black Square”, 1915,  
by Kazimir Malevich



# It Is Possible But Challenging

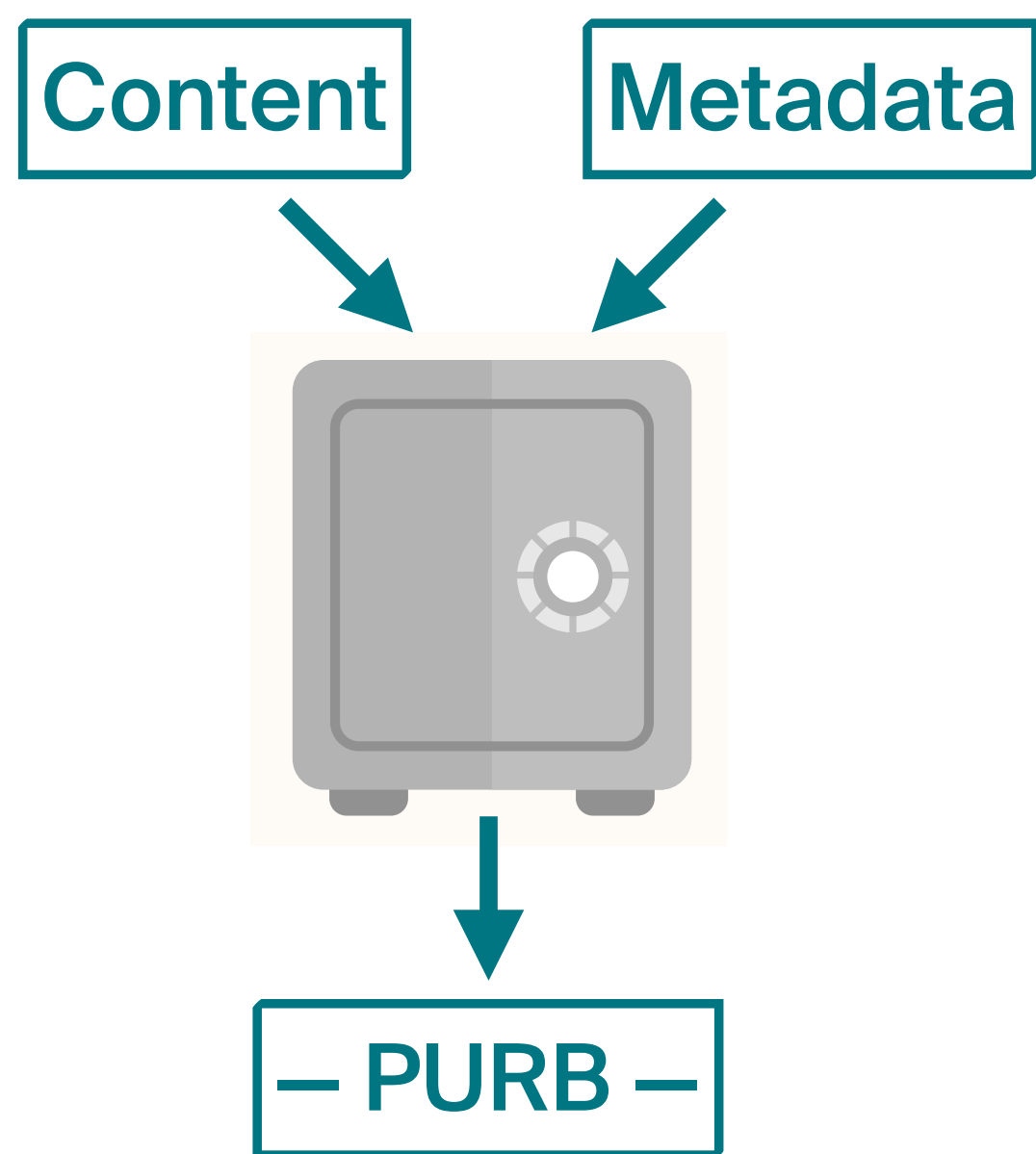
1. Efficient decoding
2. When addressing multiple recipients
3. Using different cryptographic algorithms



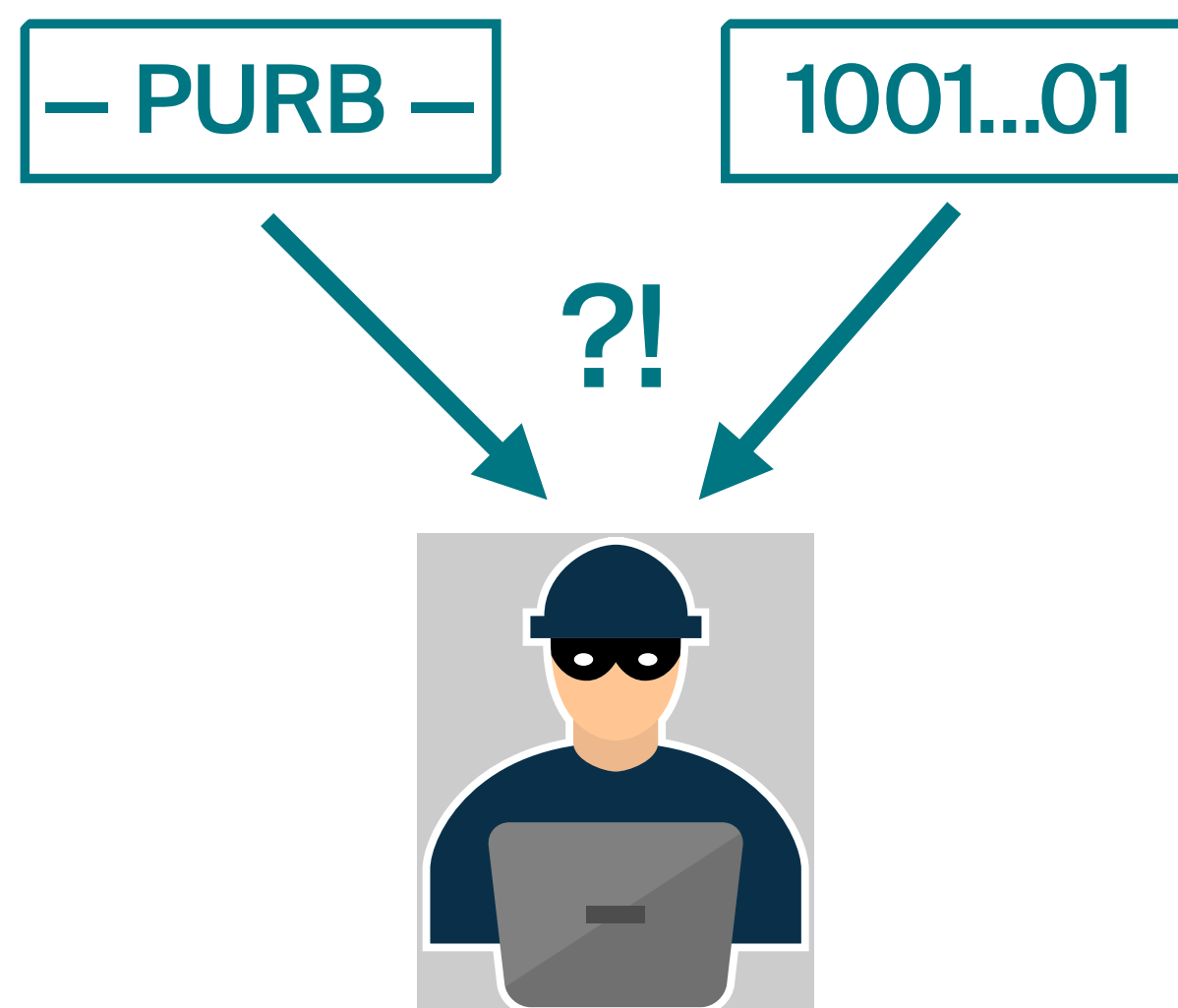


# Padded Uniform Random Blobs (PURBs)

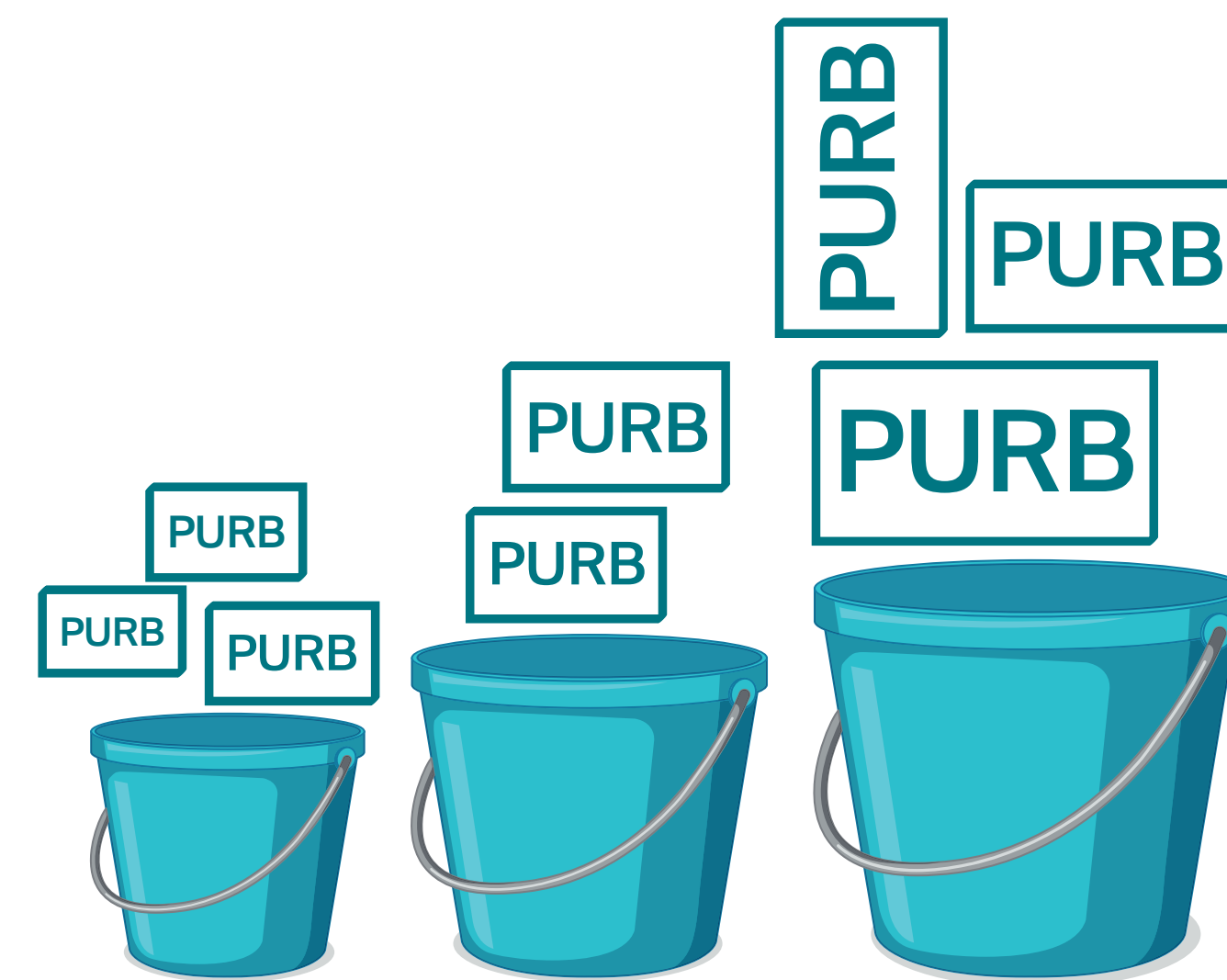
- A novel format for encrypted data without *any* metadata in clear.
- The properties (informally):



Content and metadata protection



Indistinguishability from random bits



Minimized length leakage

# Padded Uniform Random Blobs (PURBs)

- Two core components
  - Encoding scheme (*Multi-Suite* PURB or MsPURB)
  - Padding scheme (Padmé)

# Encoding scheme (MsPURB)



# Roadmap to MsPURB



# Single Recipient: Model



Is it a PURB or a random bit string?!



Active Adversary

Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability

# Single Recipient

Recipient – public key  $G^y$

Similar to the Integrated Encryption Scheme (IES) [ABR01]



Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability



# Single Recipient

Recipient – public key  $G^y$

Similar to the Integrated Encryption Scheme (IES) [ABR01]

Sender:

1. Generates an ephemeral key pair  $x$ ,  $G^x$ ;



# Single Recipient

Recipient – public key  $G^y$

Similar to the Integrated Encryption Scheme (IES) [ABR01]

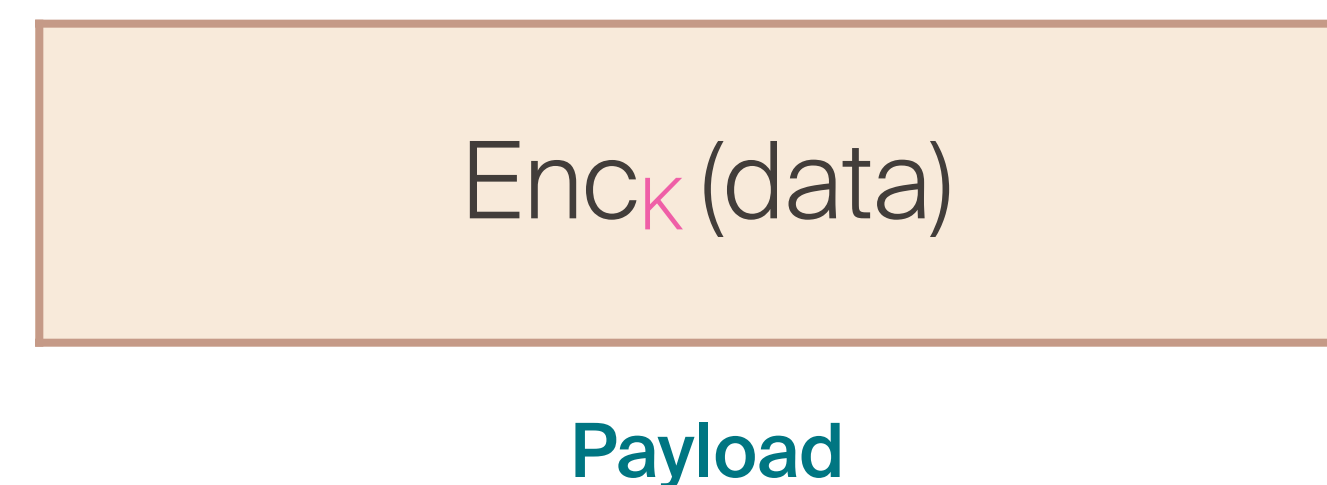
Sender:

1. Generates an ephemeral key pair  $x$ ,  $G^x$ ;
2. Computes a shared secret  $G^{yx}$ ;

# Single Recipient

Recipient – public key  $G^y$

Similar to the Integrated Encryption Scheme (IES) [ABR01]



Sender:

1. Generates an ephemeral key pair  $x, G^x$ ;
2. Computes a shared secret  $G^{yx}$ ;
3. Encrypts the data with one-time session key  $K$ ;

Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

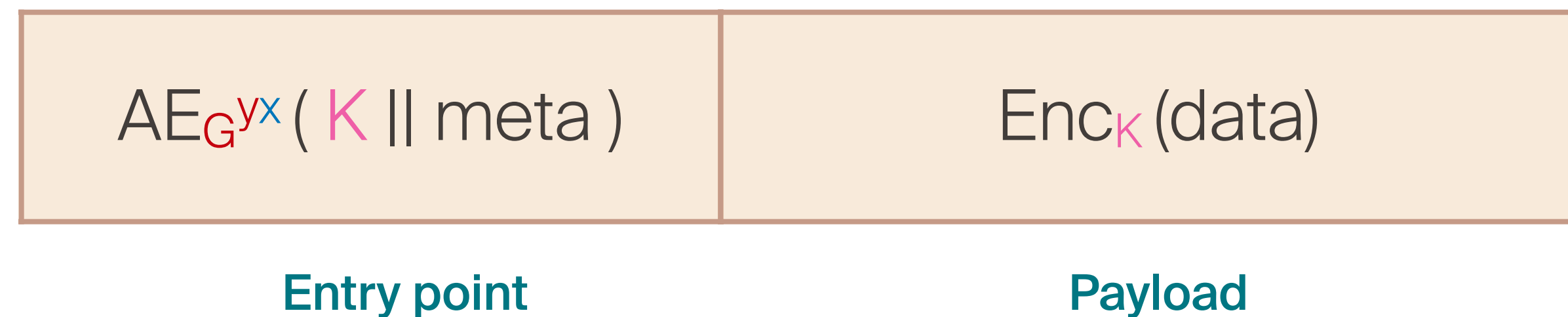
Non-  
malleability



# Single Recipient

Recipient – public key  $G^y$

Similar to the Integrated Encryption Scheme (IES) [ABR01]



Sender:

1. Generates an ephemeral key pair  $x, G^x$ ;
2. Computes a shared secret  $G^{yx}$ ;
3. Encrypts the data with one-time session key  $K$ ;
4. Creates an entry point with  $K$  and other metadata, encrypted with  $G^{yx}$ ;

Single  
Recipient

Multiple  
Recipients

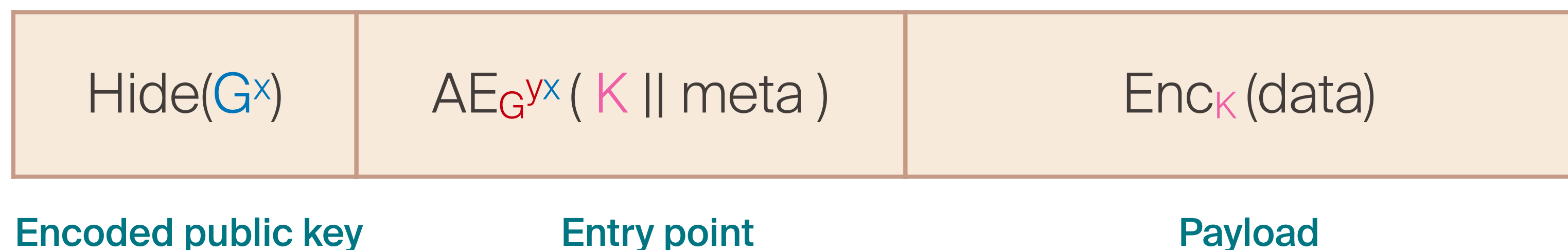
Multiple  
Suites

Non-  
malleability

# Single Recipient

Recipient – public key  $G^y$

Similar to the Integrated Encryption Scheme (IES) [ABR01]



Sender:

1. Generates an ephemeral key pair  $x, G^x$ ;
2. Computes a shared secret  $G^{yx}$ ;
3. Encrypts the data with one-time session key  $K$ ;
4. Creates an entry point with  $K$  and other metadata, encrypted with  $G^{yx}$ ;
5. Encodes  $G^x$  as a uniform bit string, e.g., with Elligator [BHK13].

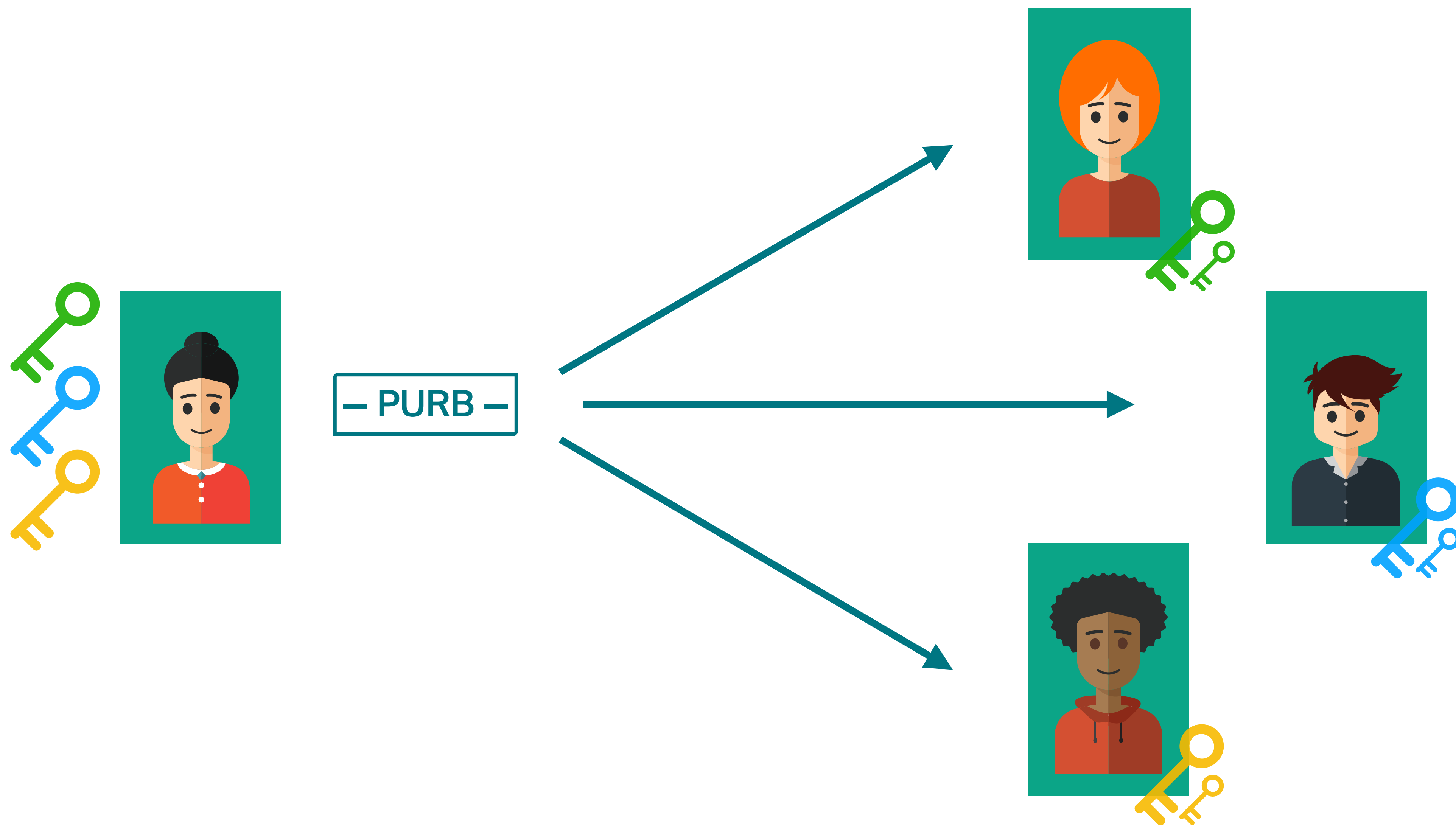
Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability

# Multiple Recipients



Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability



# Multiple Recipients

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

We create an entry point per recipient, each with  $K$  and metadata but encrypted with  $G^{y1x}$ ,  $G^{y2x}$ ,  $G^{y3x}$  respectively.

$$AE_{G^{y1x}}(K||\text{meta})$$
$$AE_{G^{y2x}}(K||\text{meta})$$
$$AE_{G^{y3x}}(K||\text{meta})$$

# Multiple Recipients

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

We create an entry point per recipient, each with  $K$  and metadata but encrypted with  $G^{y1x}$ ,  $G^{y2x}$ ,  $G^{y3x}$  respectively.

$$AE_{G^{y1x}}(K || \text{meta})$$
$$AE_{G^{y2x}}(K || \text{meta})$$
$$AE_{G^{y3x}}(K || \text{meta})$$

But how do we organize these entry points in the PURB?

# Linear Approach Strawman

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

Hide( $G^x$ )

Enc <sub>$K$</sub> (data)

We create an entry point per recipient, each with  $K$  and metadata but encrypted with  $G^{y1x}$ ,  $G^{y2x}$ ,  $G^{y3x}$  respectively.

Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability

# Linear Approach Strawman

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .



We create an entry point per recipient, each with  $K$  and metadata but encrypted with  $G^{y1x}$ ,  $G^{y2x}$ ,  $G^{y3x}$  respectively.



# Linear Approach Strawman

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

Inefficient to decode



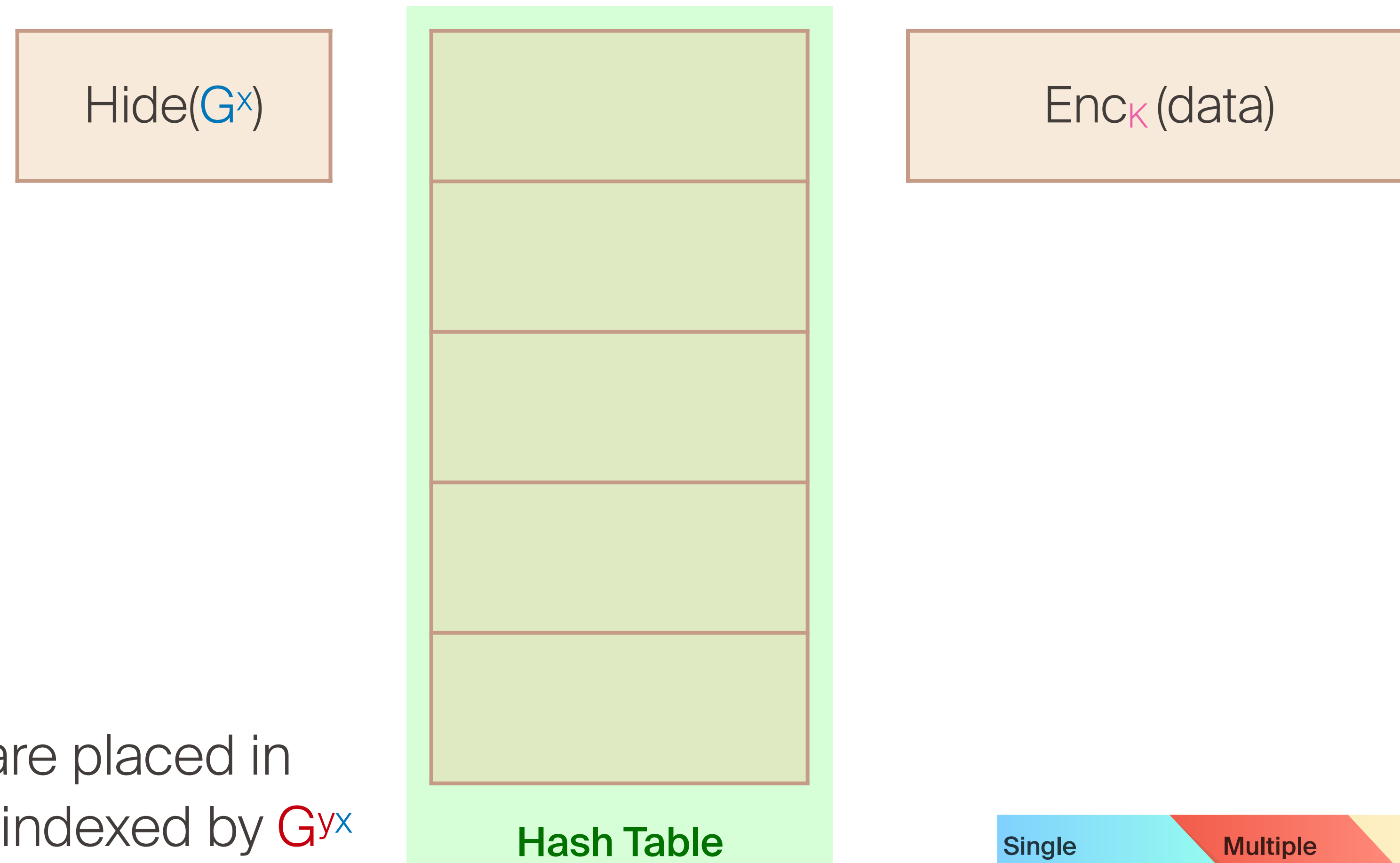
We create an entry point per recipient, each with  $K$  and metadata but encrypted with  $G^{y1x}$ ,  $G^{y2x}$ ,  $G^{y3x}$  respectively.





# Single Hash-Table Strawman

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .



Entry points are placed in a hash table, indexed by  $G^{y^x}$

Single  
Recipient

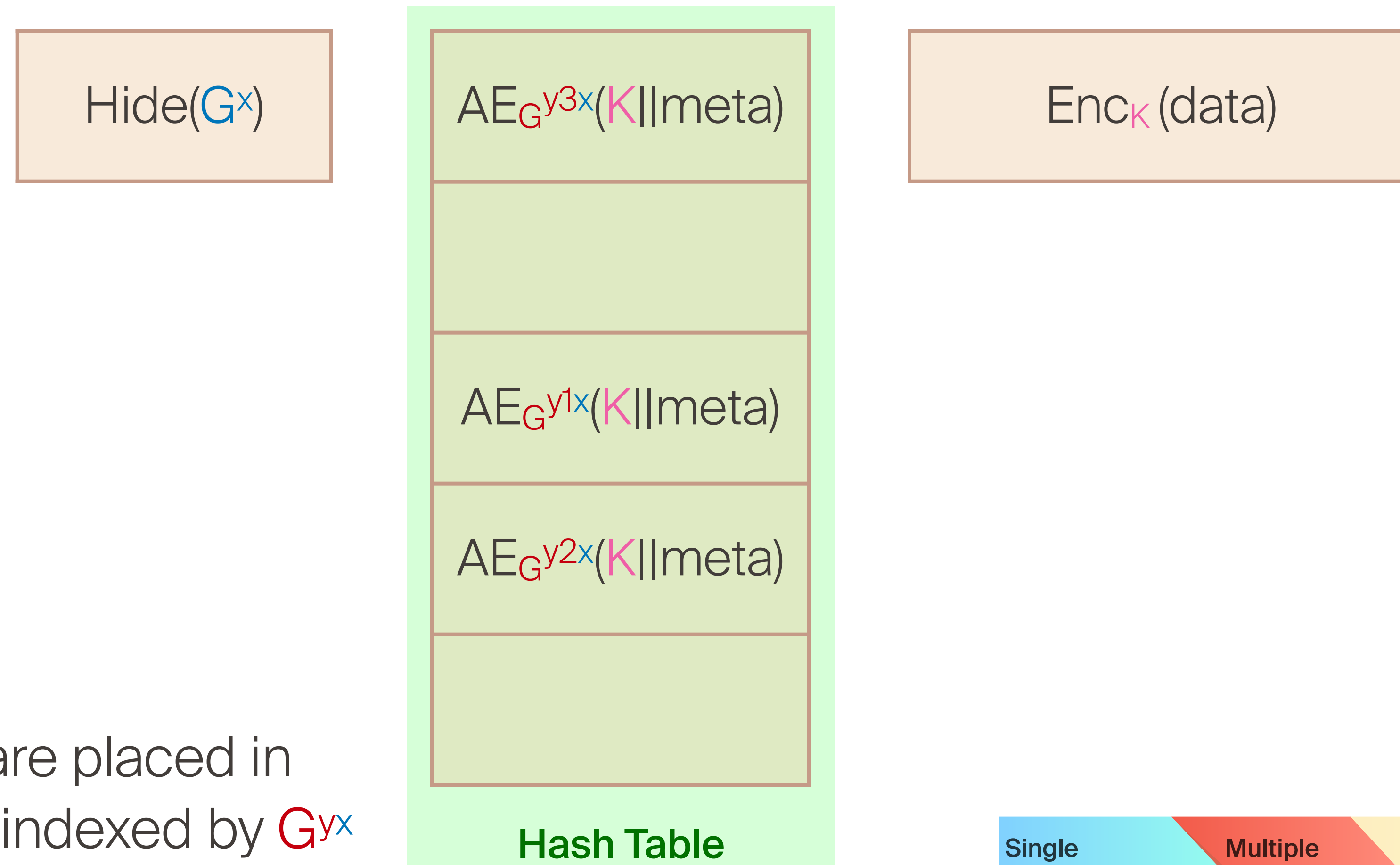
Multiple  
Recipients

Multiple  
Suites

Non-  
malleability

# Single Hash-Table Strawman

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

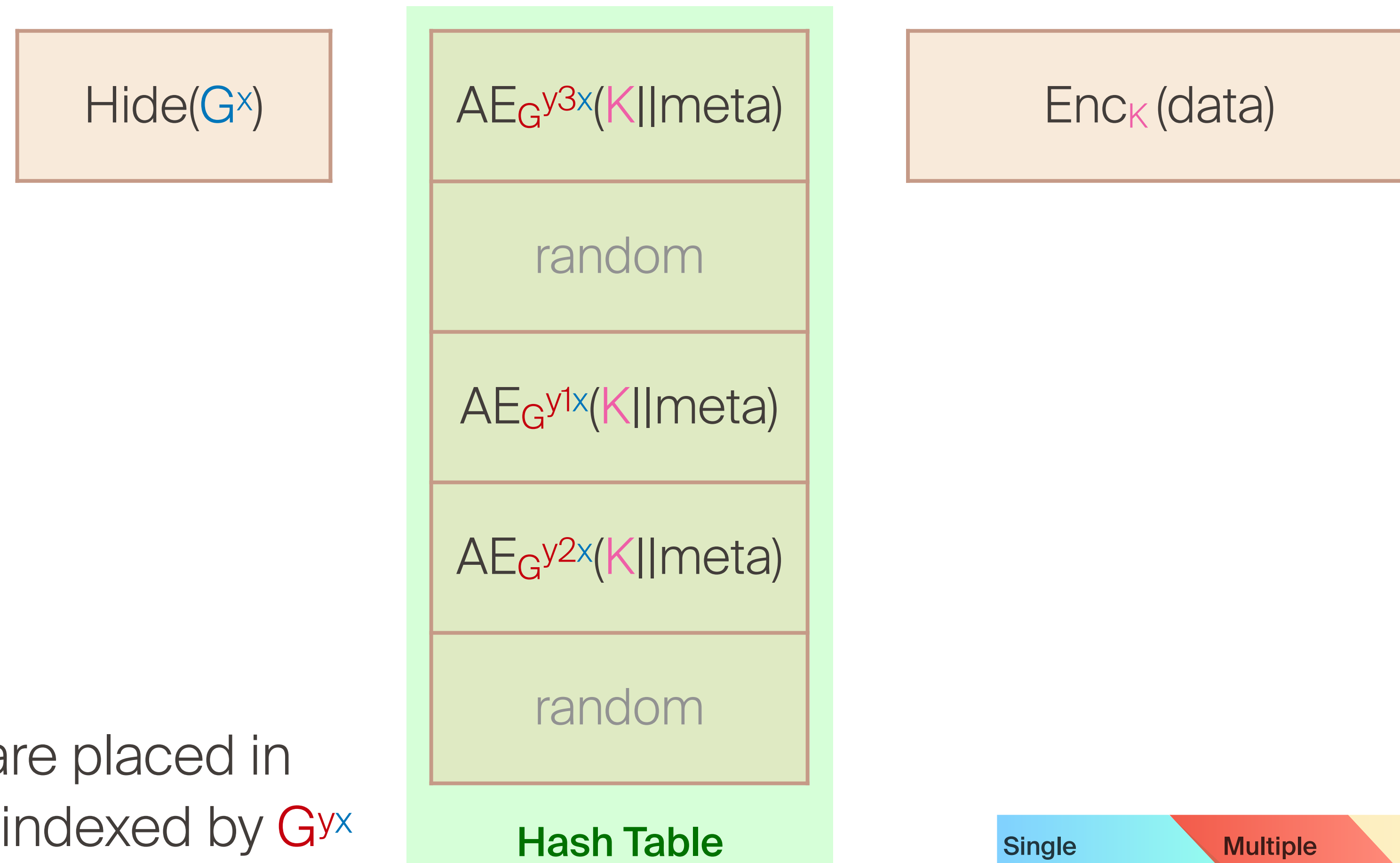


Entry points are placed in a hash table, indexed by  $G^{yx}$



# Single Hash-Table Strawman

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

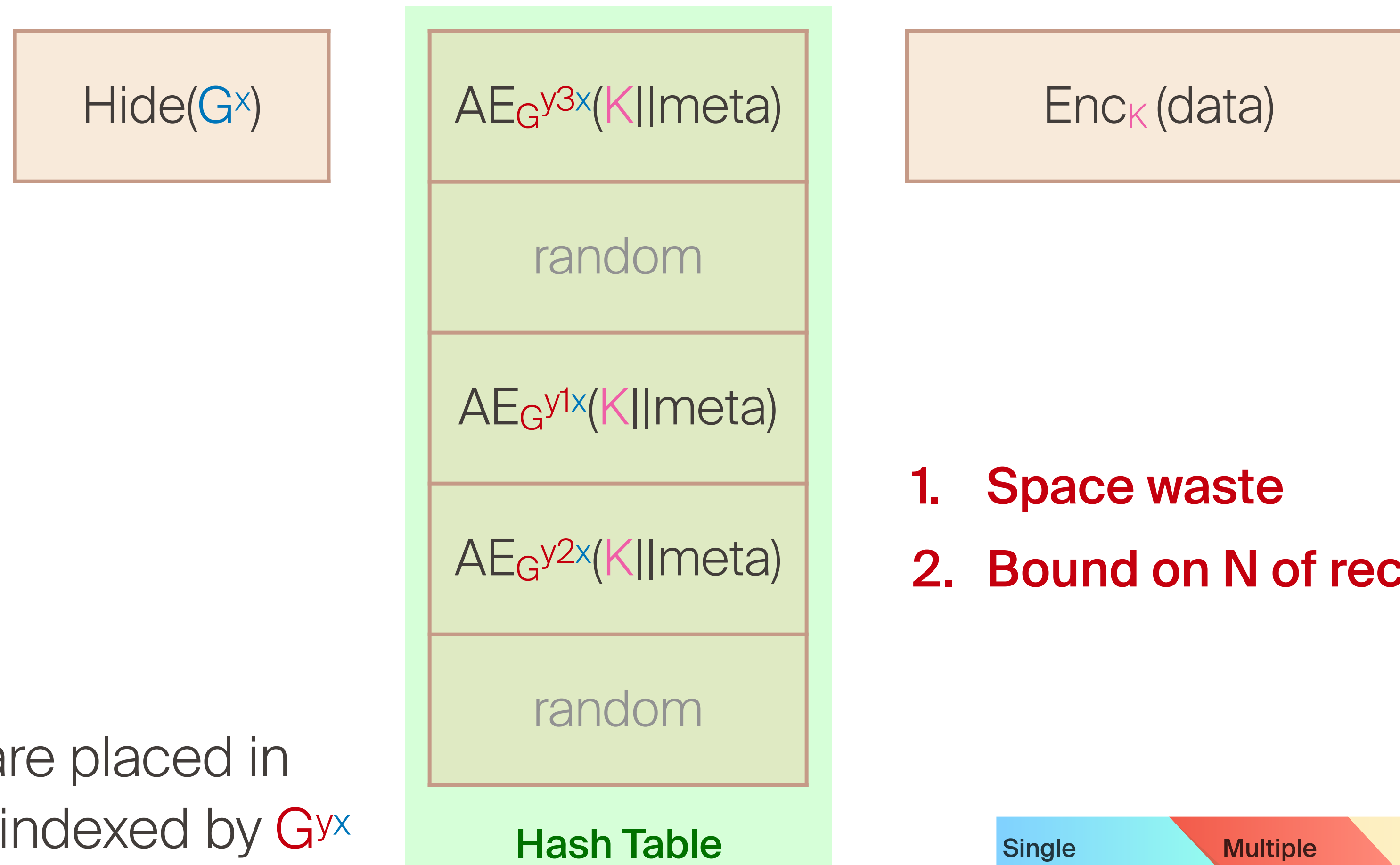


Entry points are placed in a hash table, indexed by  $G^{yx}$



# Single Hash-Table Strawman

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .



Entry points are placed in a hash table, indexed by  $G^{yx}$

1. Space waste
2. Bound on N of recipients



# Multiple Recipients: Our Solution

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

Hide( $G^x$ )

Enc <sub>$k$</sub> (data)

Entry points are placed in a series of growing **hash-tables**!

Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability





# Multiple Recipients: Our Solution

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

$$AE_{G^{y1x}}(K || \text{meta})$$
$$\text{Hide}(G^x)$$
$$AE_{G^{y1x}}(K || \text{meta})$$
$$\text{Enc}_K(\text{data})$$

HT0

Entry points are placed in a series of growing **hash-tables!**

Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability



# Multiple Recipients: Our Solution

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

$AE_{G^{y2x}}(K||meta)$

Hide( $G^x$ )

$AE_{G^{y1x}}(K||meta)$

$AE_{G^{y2x}}(K||meta)$

$Enc_K(data)$

HT0

HT1

Entry points are placed in a series of growing **hash-tables**!

Single  
Recipient

Multiple  
Recipients

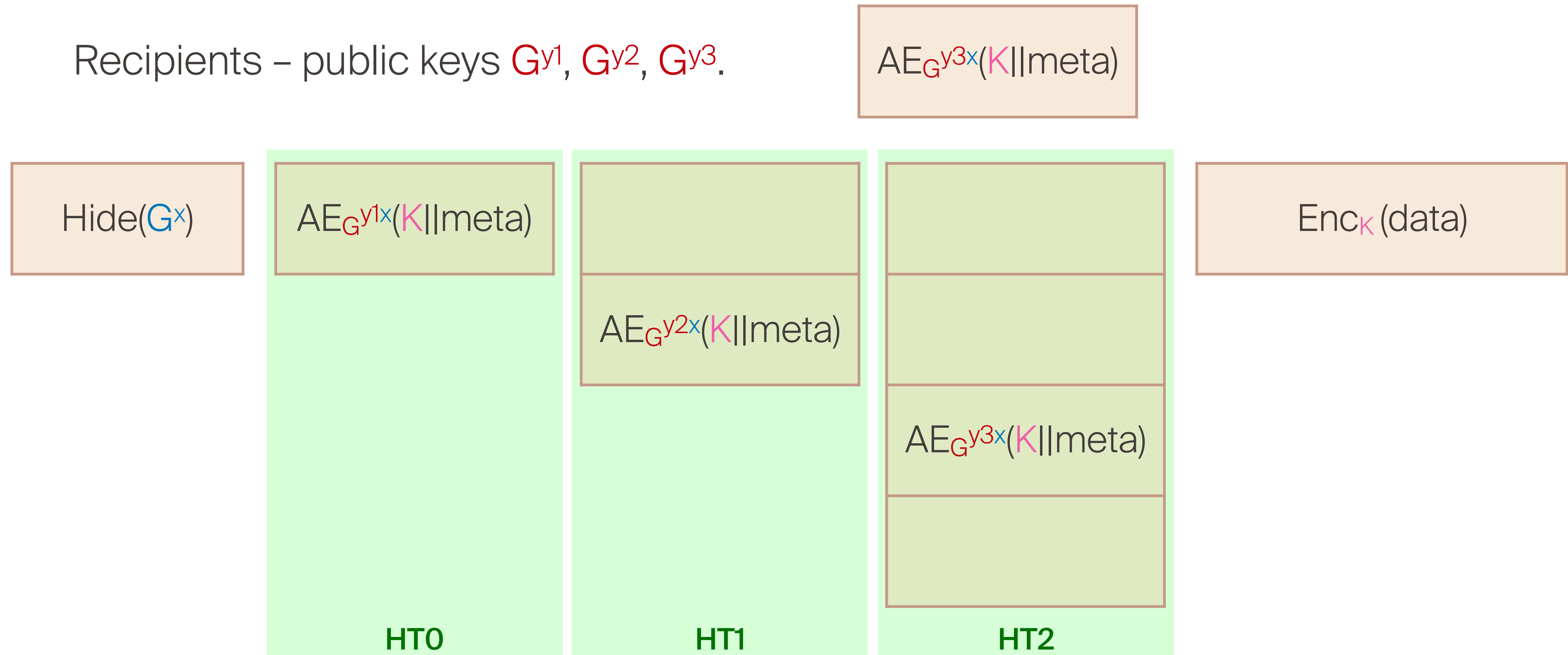
Multiple  
Suites

Non-  
malleability



# Multiple Recipients: Our Solution

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .



Entry points are placed in a series of growing **hash-tables**!

Single  
Recipient

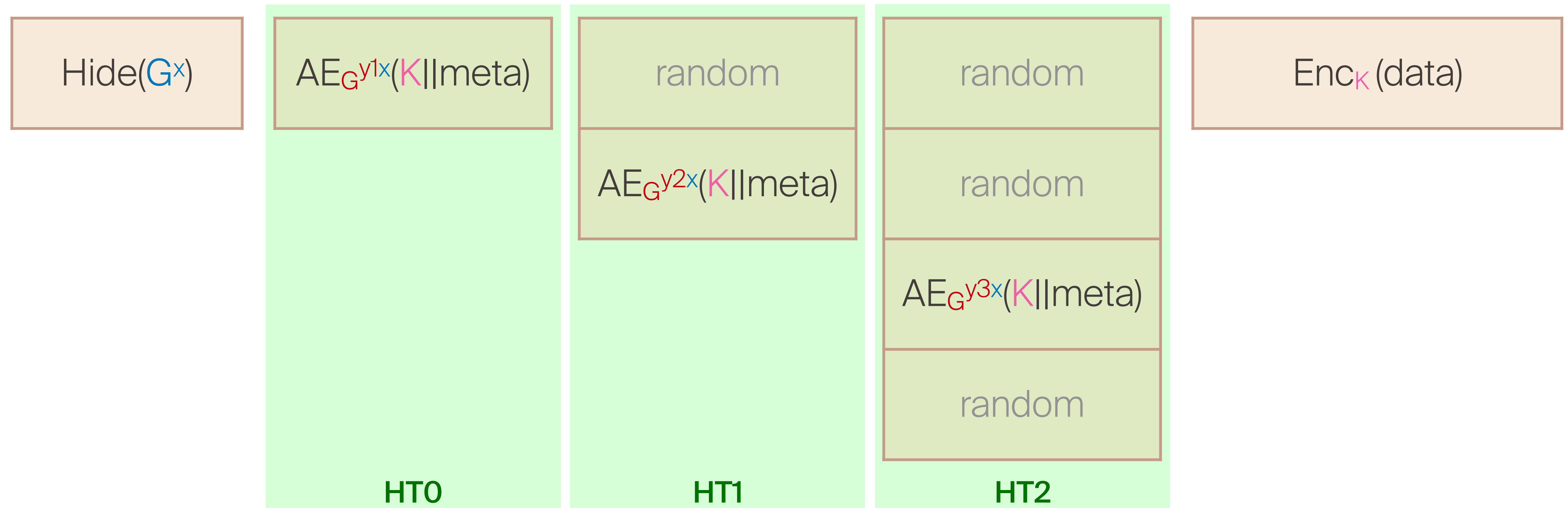
Multiple  
Recipients

Multiple  
Suites

Non-  
malleability

# Multiple Recipients: Our Solution

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .



Entry points are placed in a series of growing **hash-tables**!

Single  
Recipient

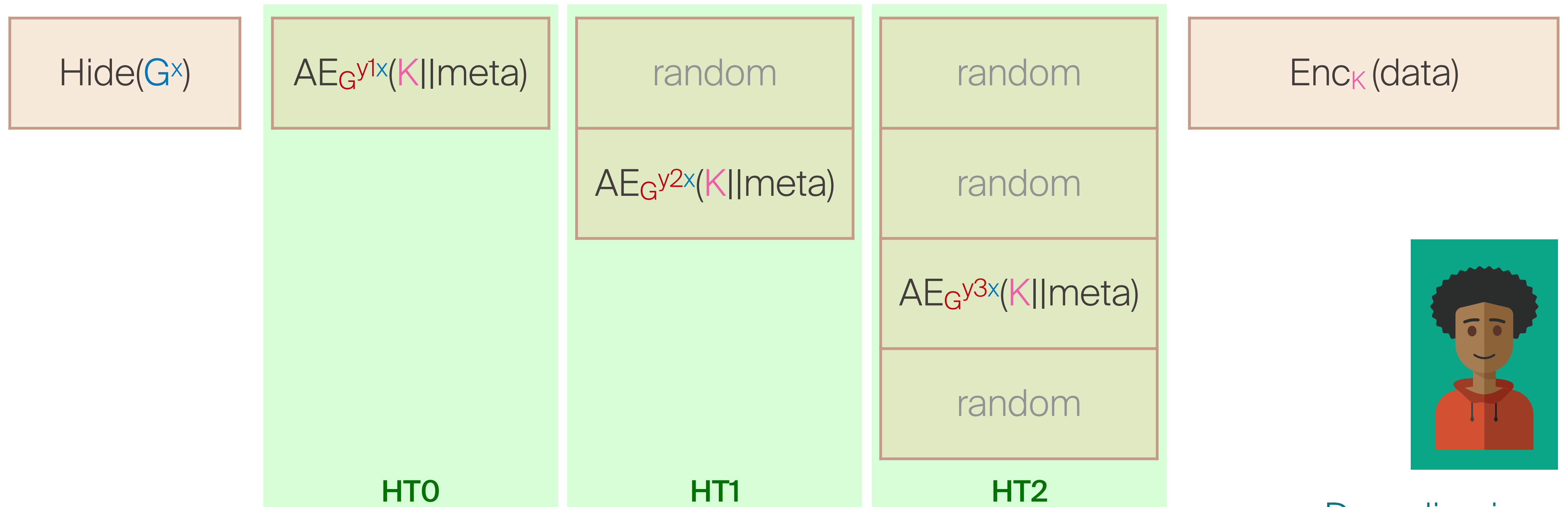
Multiple  
Recipients

Multiple  
Suites

Non-  
malleability

# Multiple Recipients: Our Solution

Recipients – public keys  $G^{y1}$ ,  $G^{y2}$ ,  $G^{y3}$ .

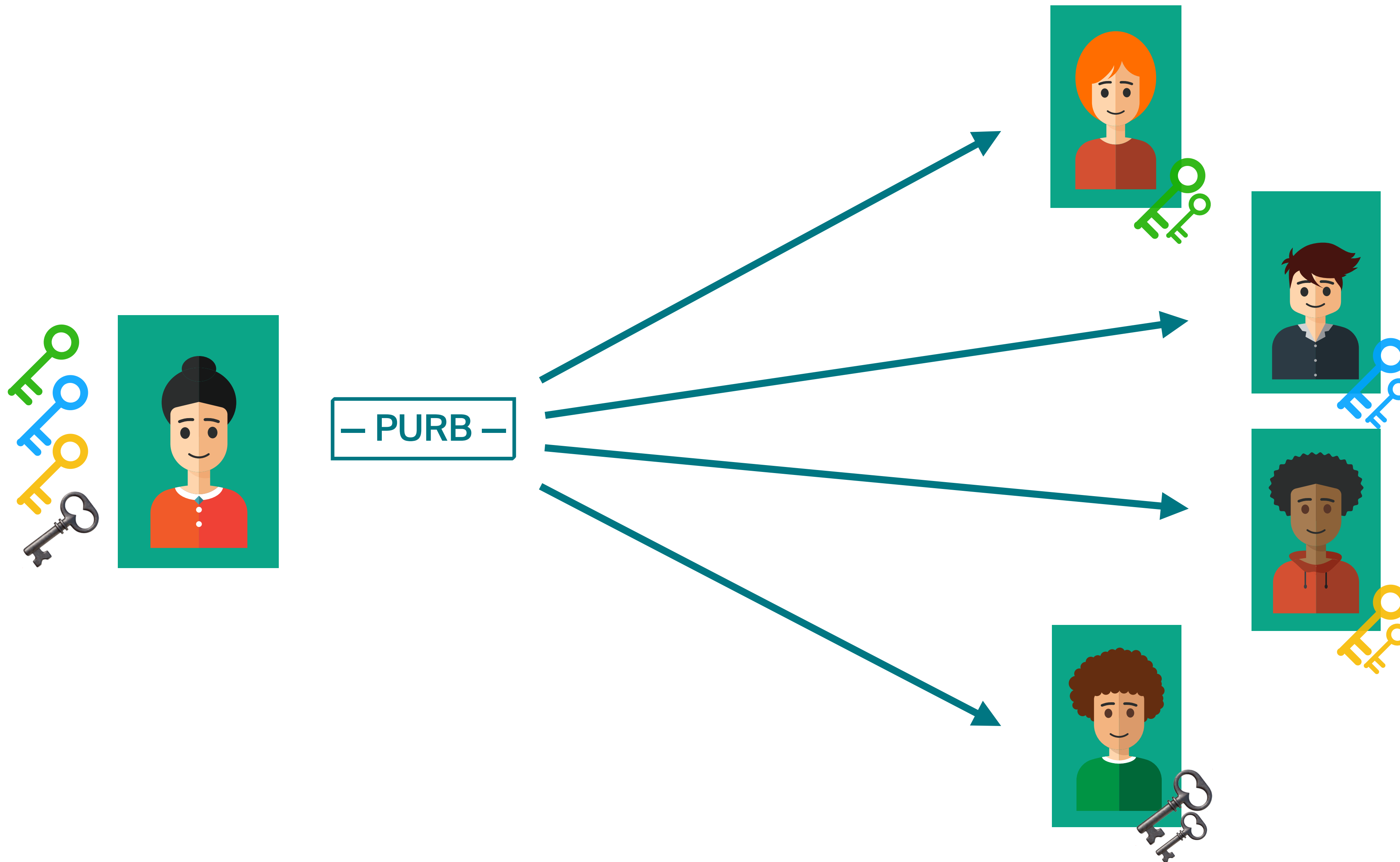


Decoding in  
 $\log_2 \text{len}(\text{PURB})$

Entry points are placed in a series of growing **hash-tables**!



# Multiple Suites



Single  
Recipient

Multiple  
Recipients

Multiple  
Suites

Non-  
malleability



# Multiple Suites

- Recipients use several distinct suites, based on public-key group (e.g., Curve25519 or Curve448) or entry point length.
- Each suite (an encoded public key and hash tables) becomes a distinct logical layer in a PURB, and these layers overlap!



# Multiple Suites: Layout

Suite A



PURB bytes



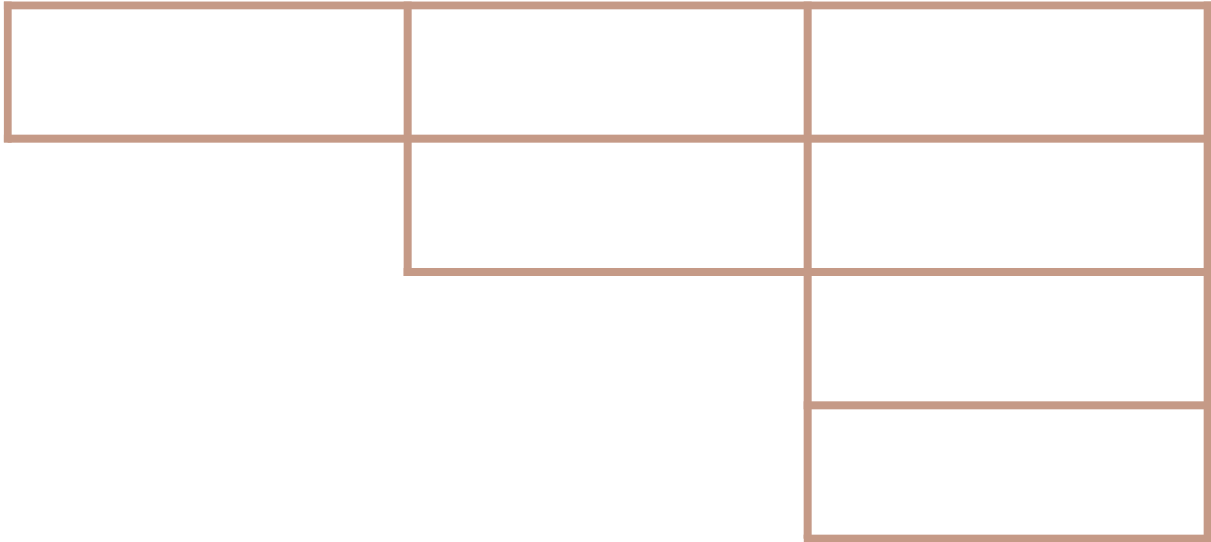
Suite B



# Multiple Suites: Layout

Suite A

Hide(A)



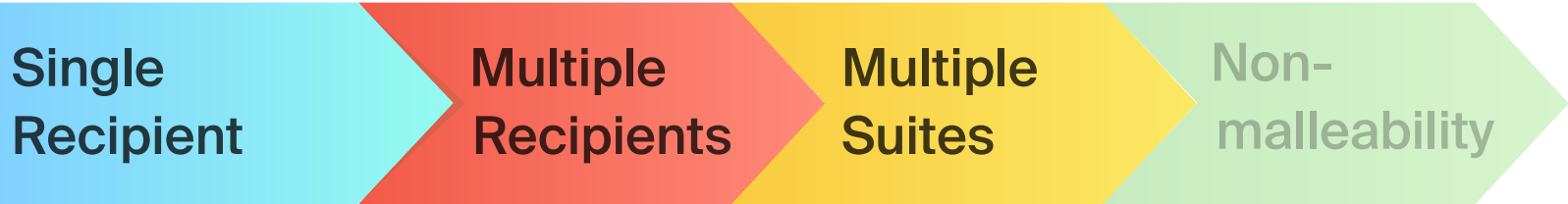
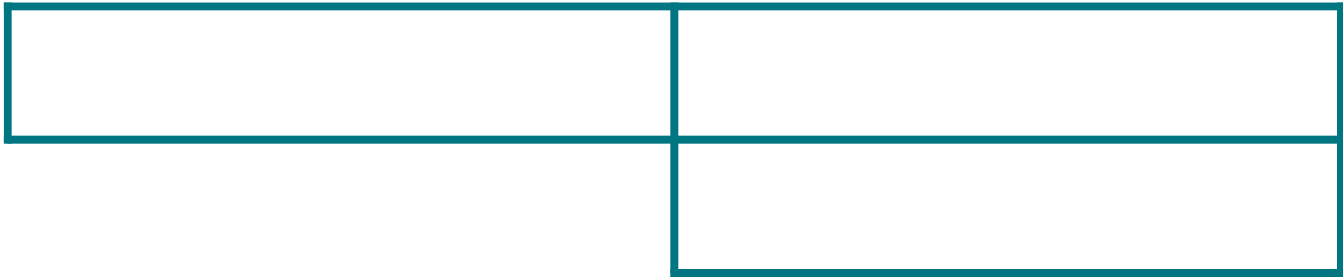
Enc(data)

PURB bytes

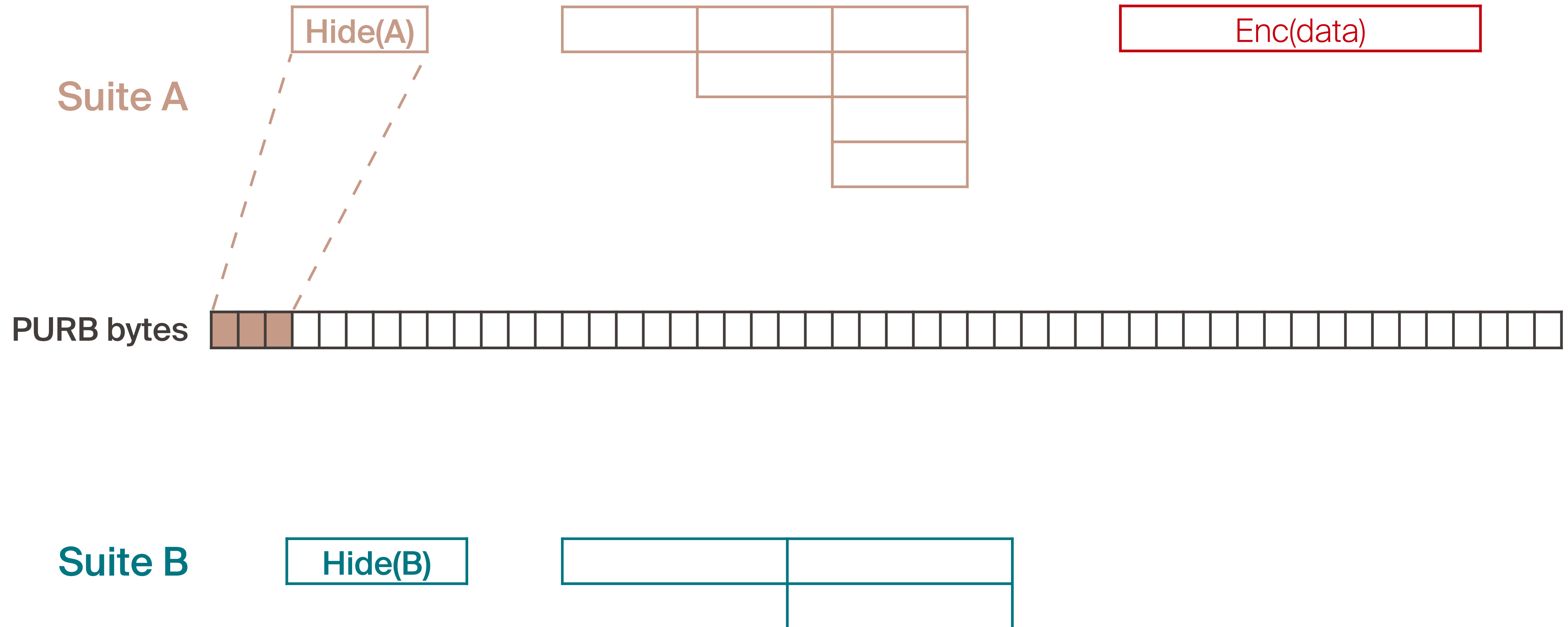


Suite B

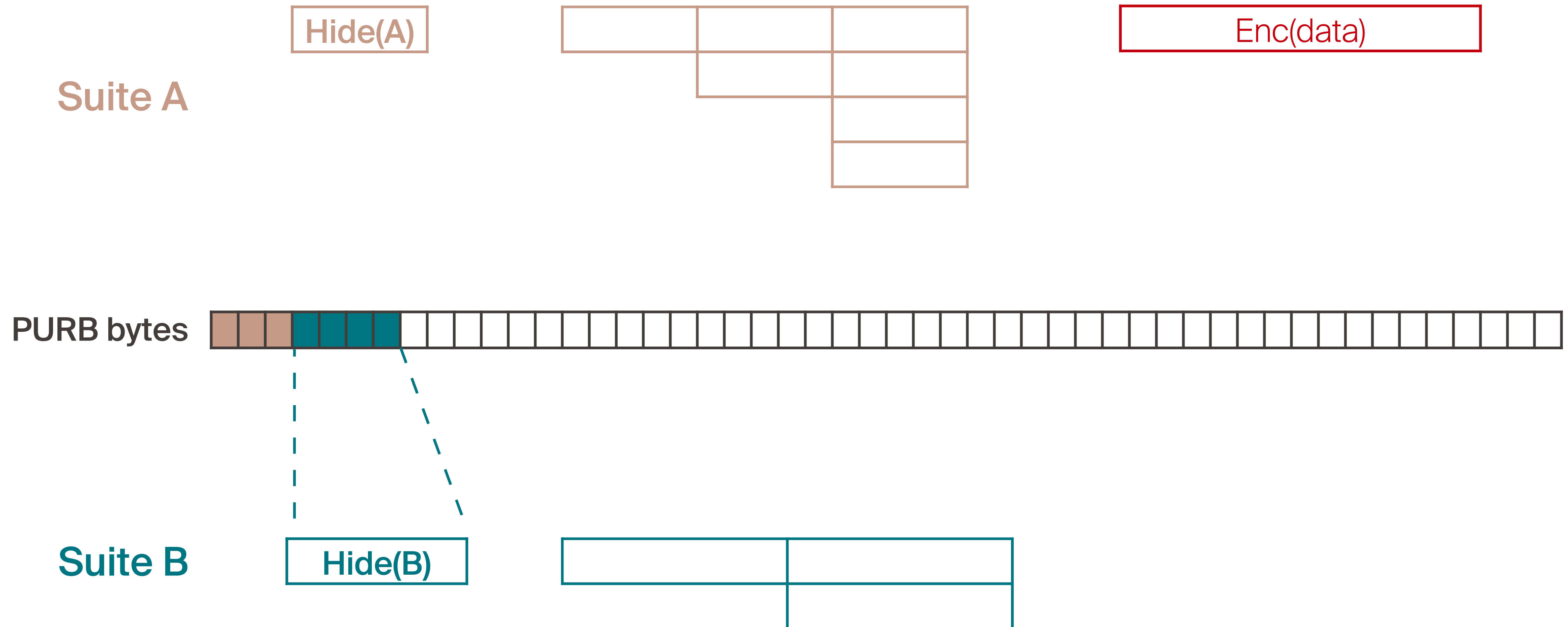
Hide(B)



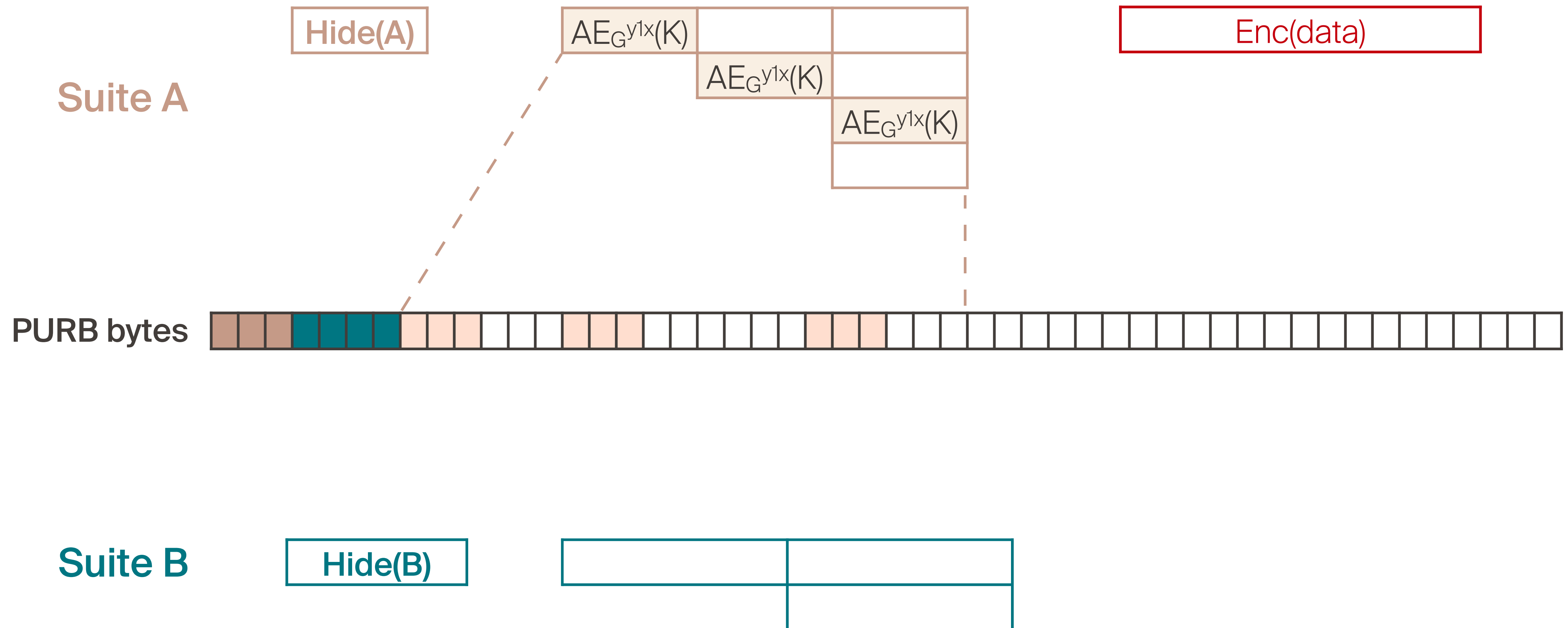
# Multiple Suites: Layout



# Multiple Suites: Layout

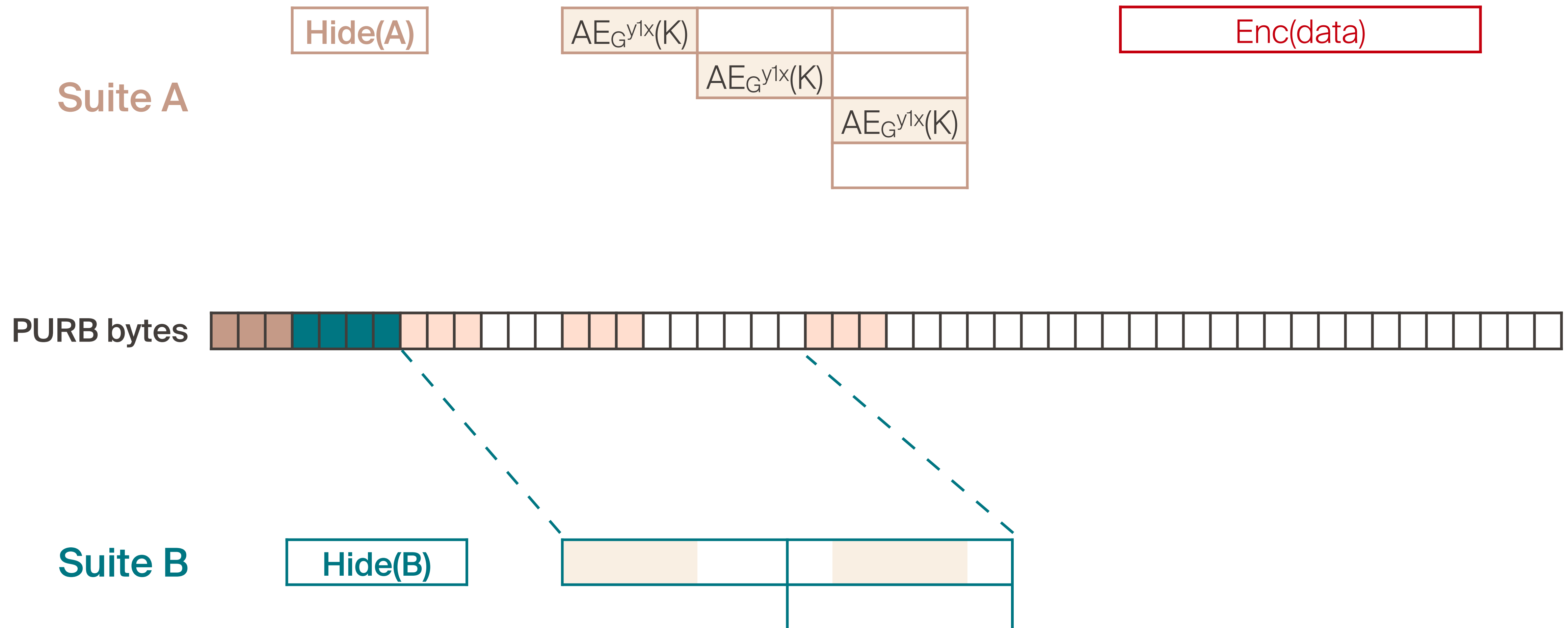


# Multiple Suites: Layout

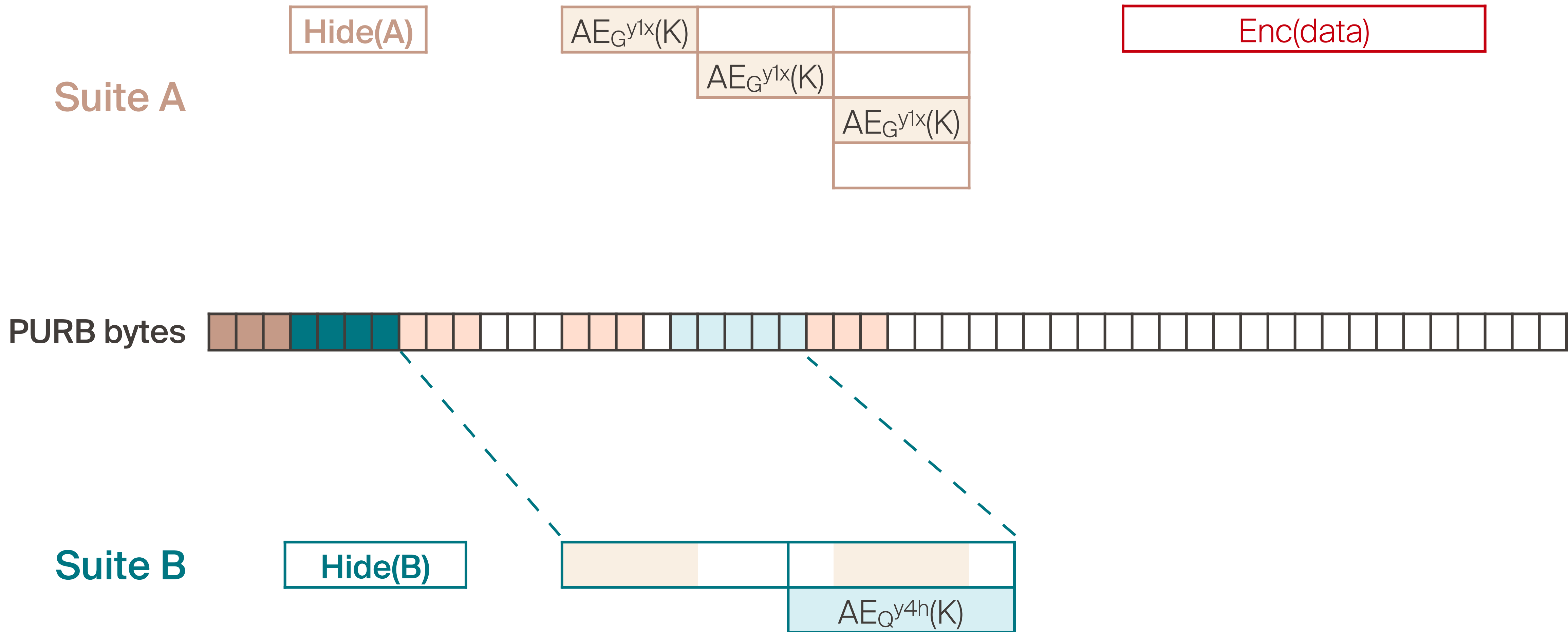


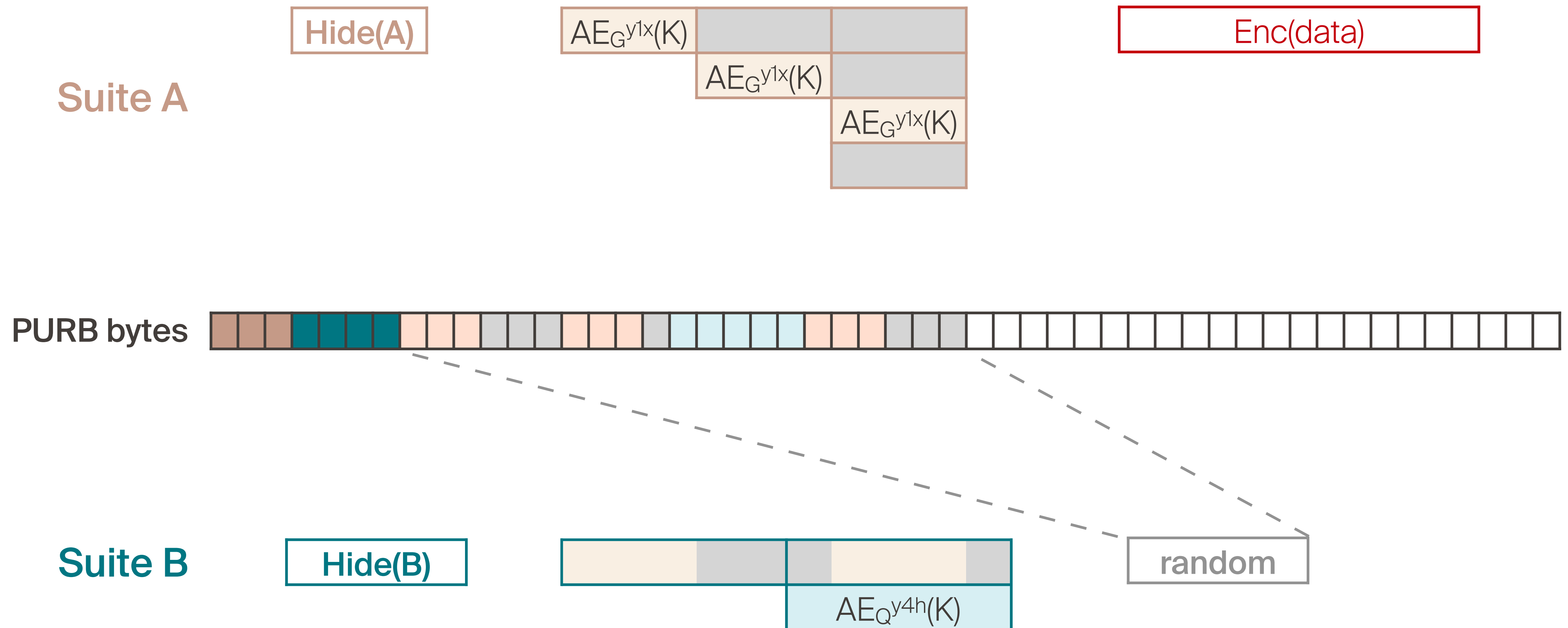


# Multiple Suites: Layout

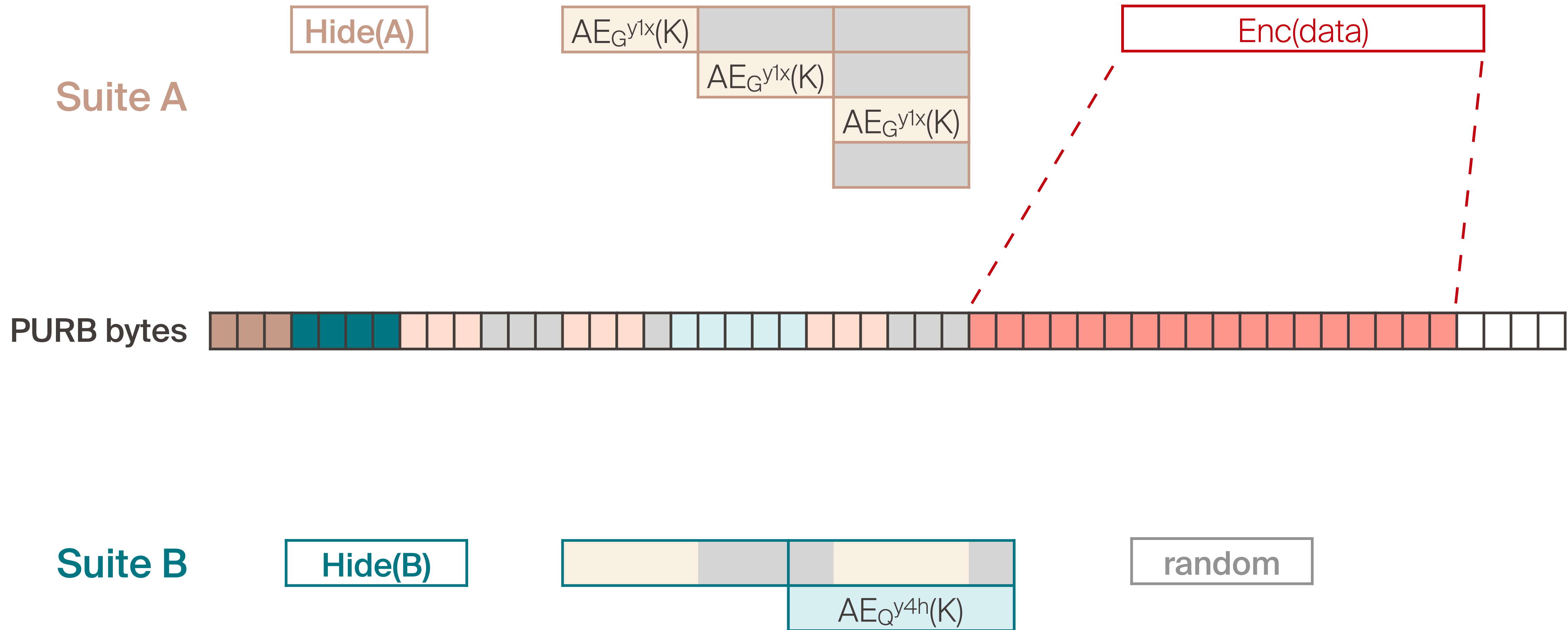


# Multiple Suites: Layout





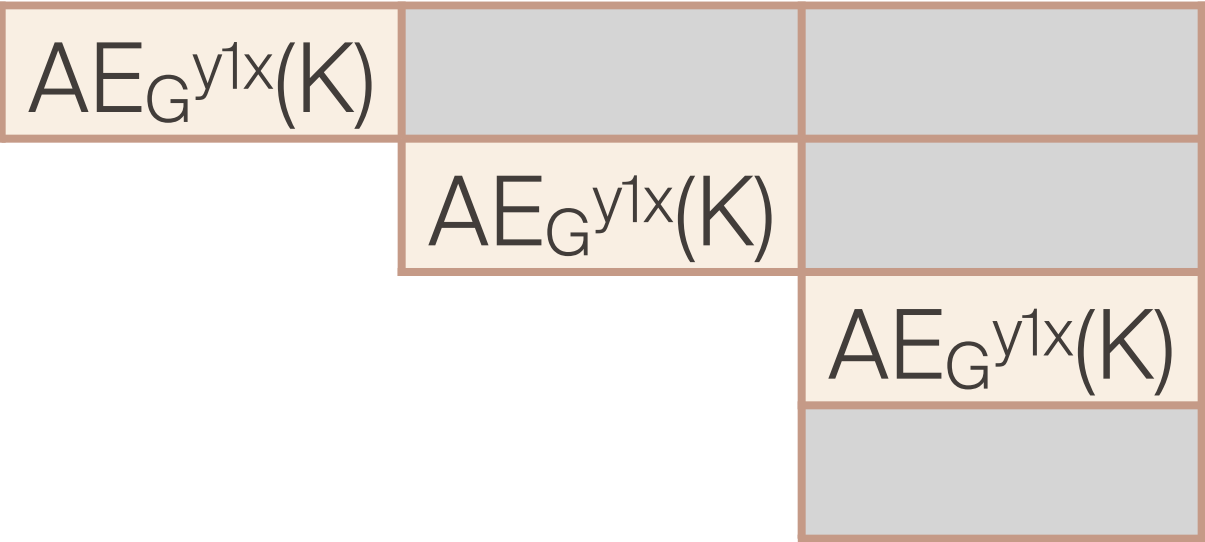
# Multiple Suites: Layout



# Non-malleability

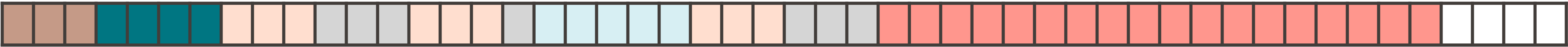
Suite A

Hide(A)



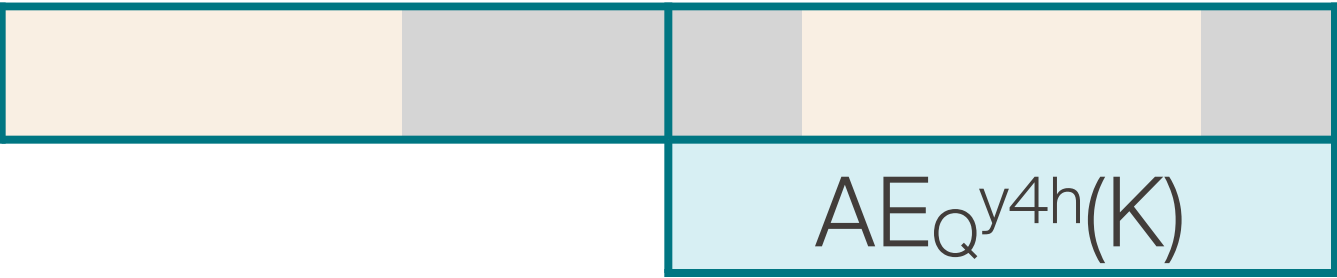
Enc(data)

PURB bytes



Suite B

Hide(B)

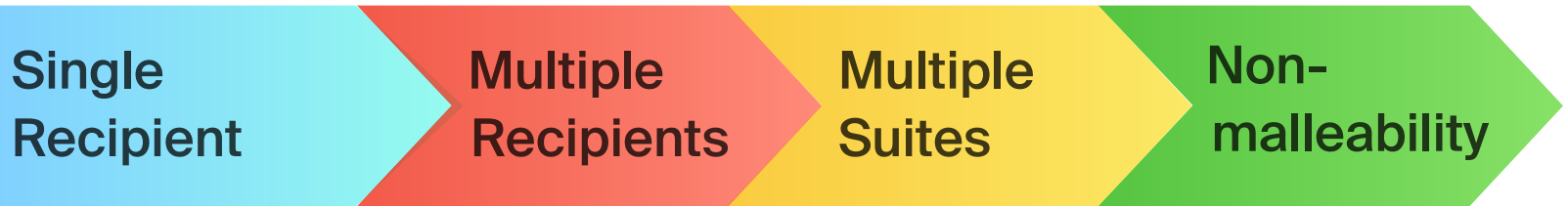
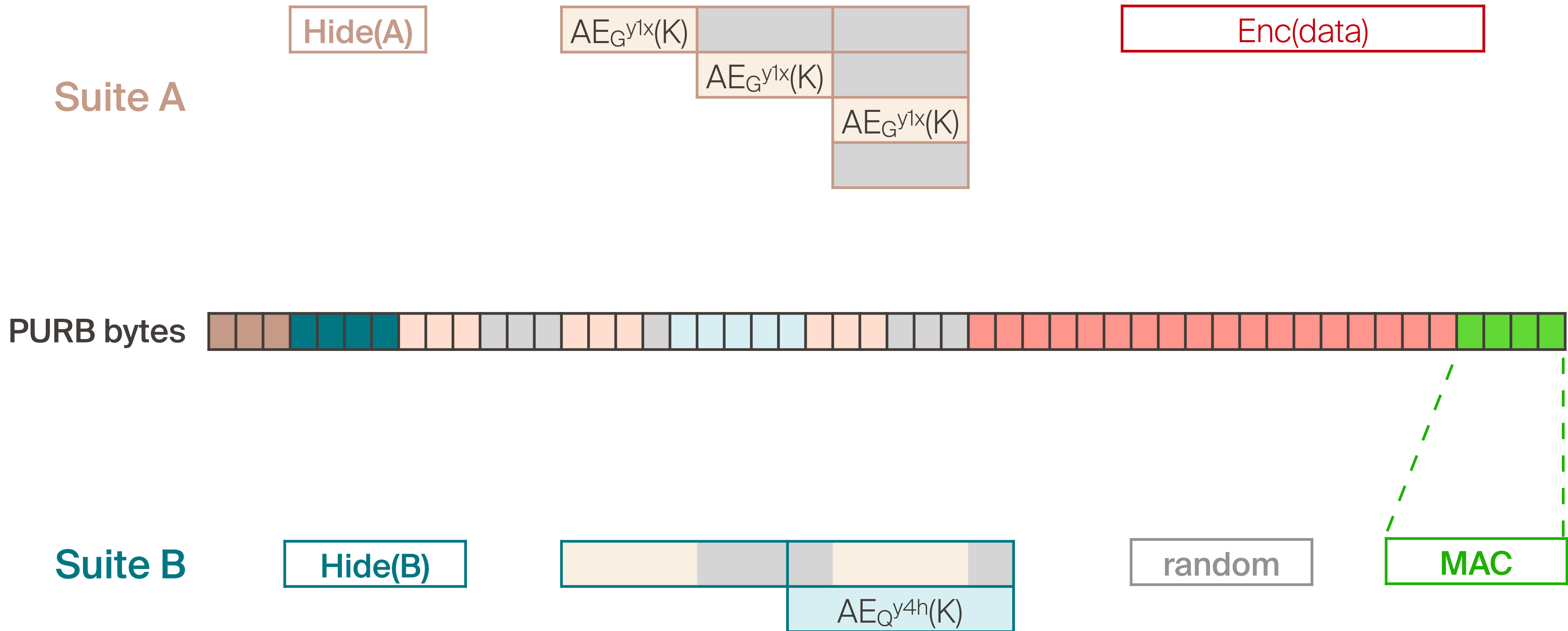


random

MAC



# Non-malleability



# Hide(A)

$$AE_G^{y1x}(K)$$
$$AE_G^{y1x}(K)$$
$$AE_G^{y1x}(K)$$

Enc(data)

## PURB bytes

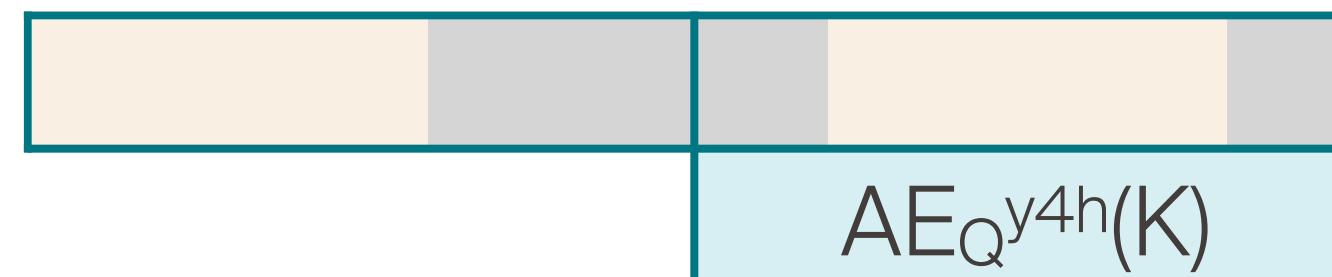
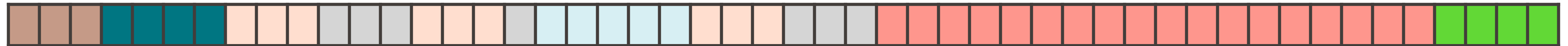
## Suite B

## Hide(B)

$$AE_{Q^{y4h}}(K)$$

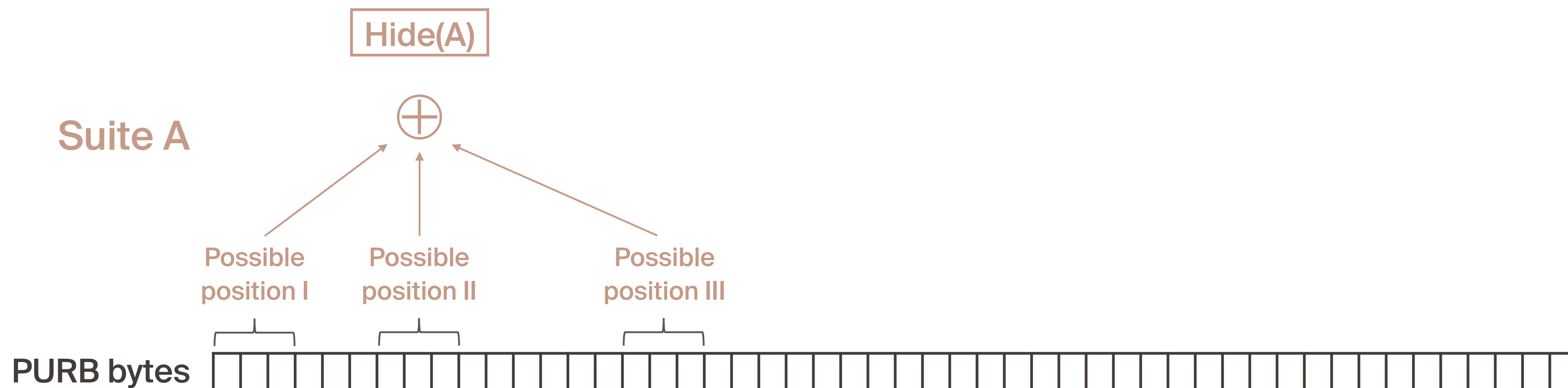
random

# MAC



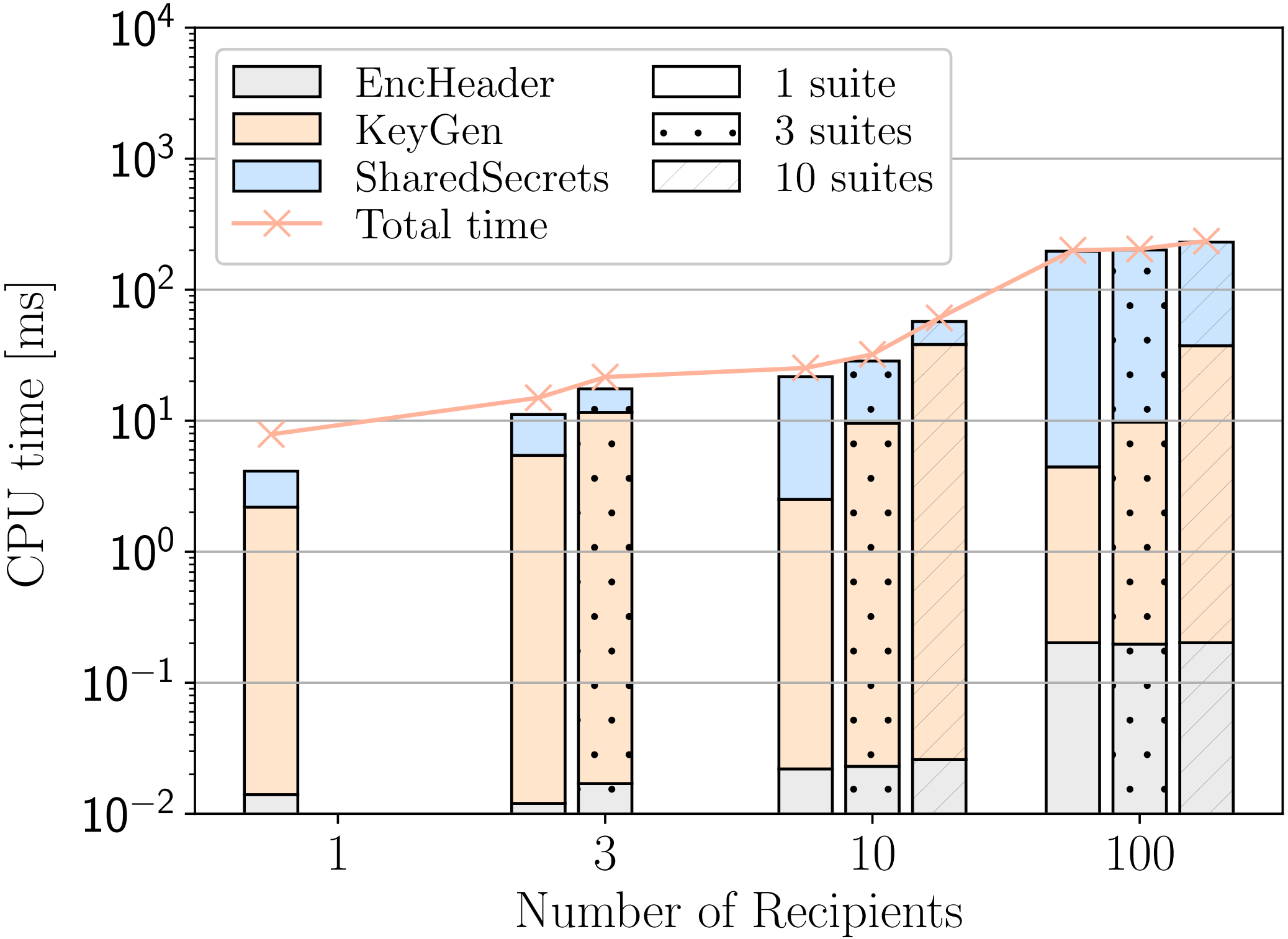


# Finding Public Keys Efficiently

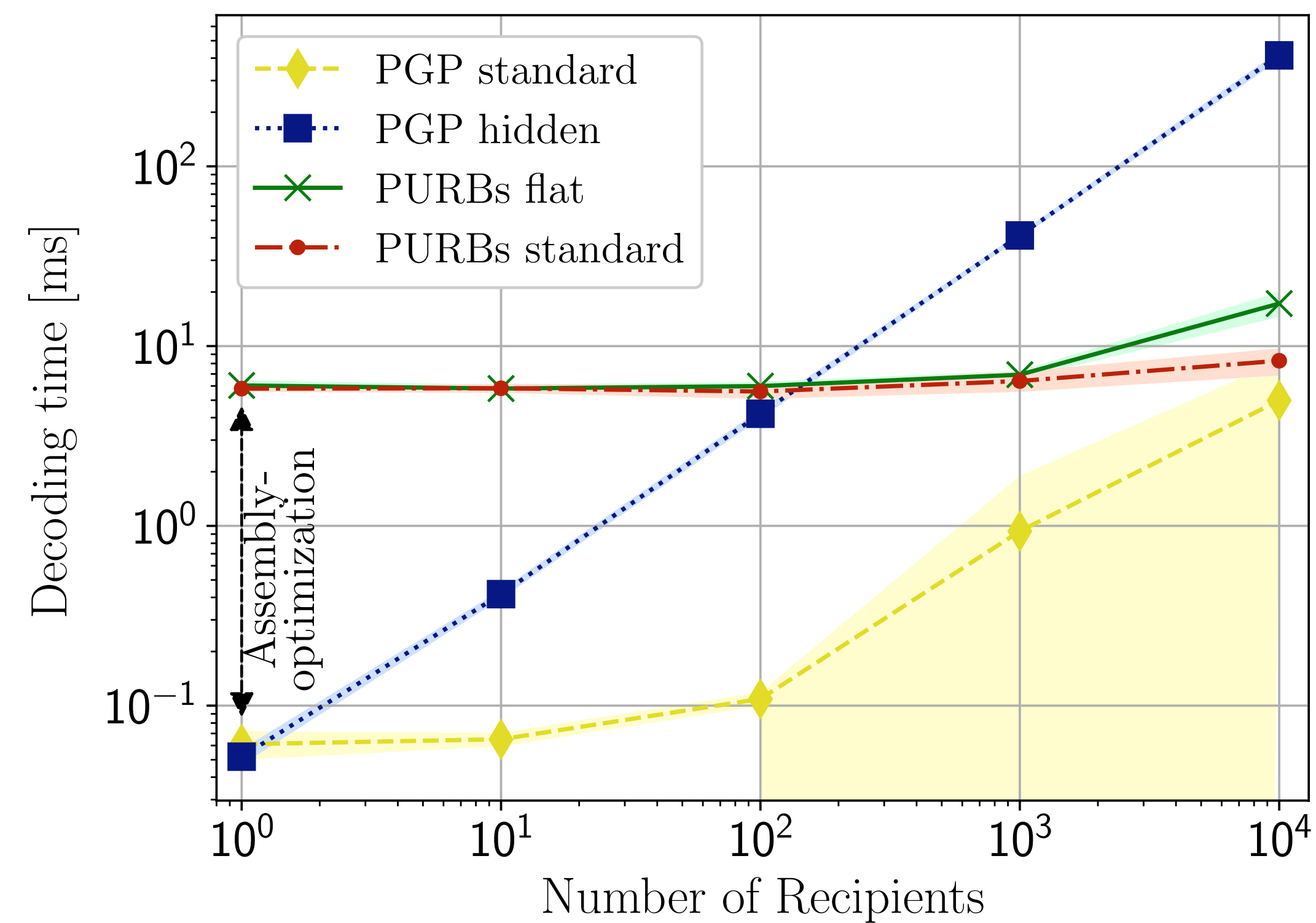
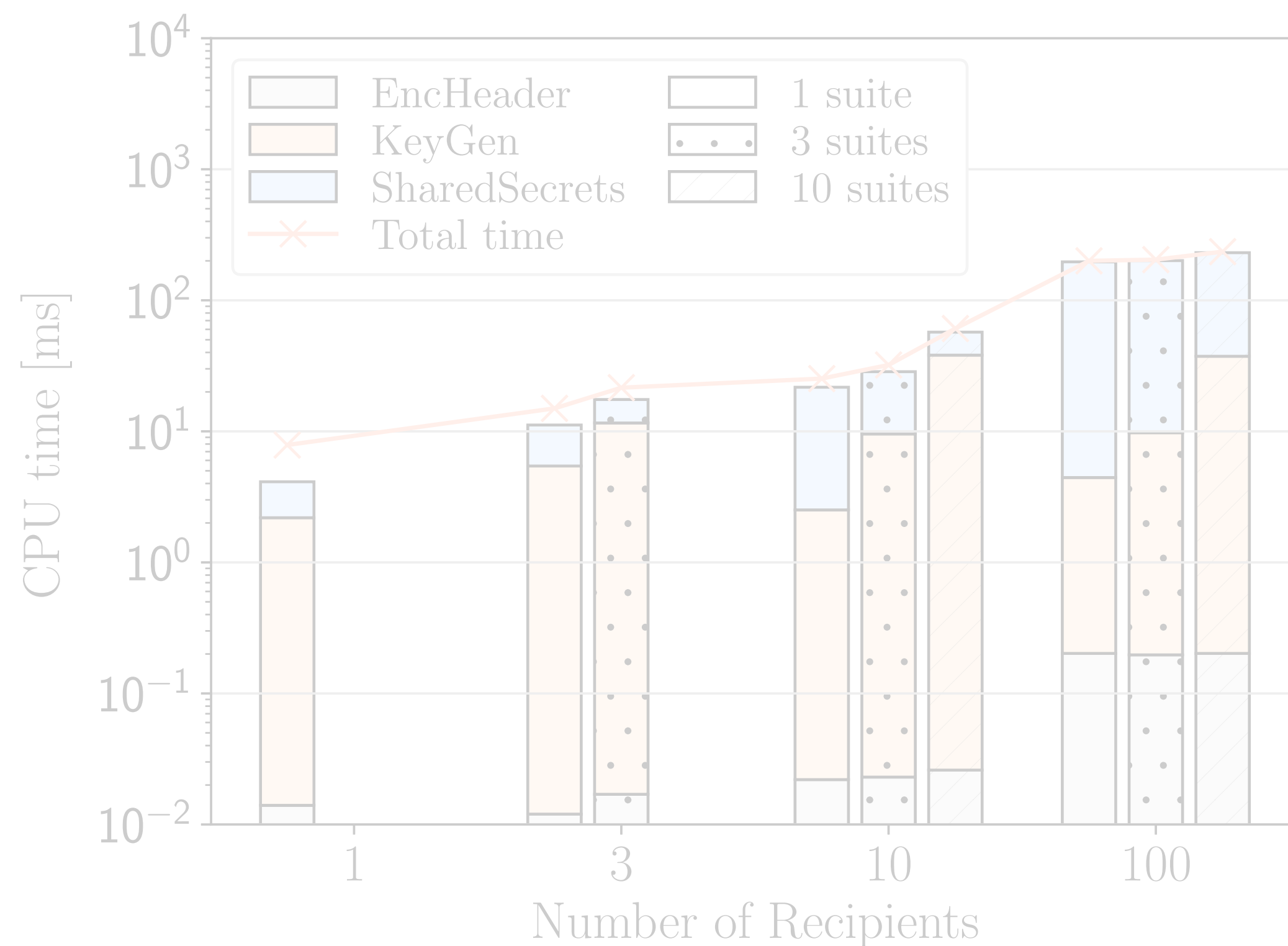


See the paper for the details

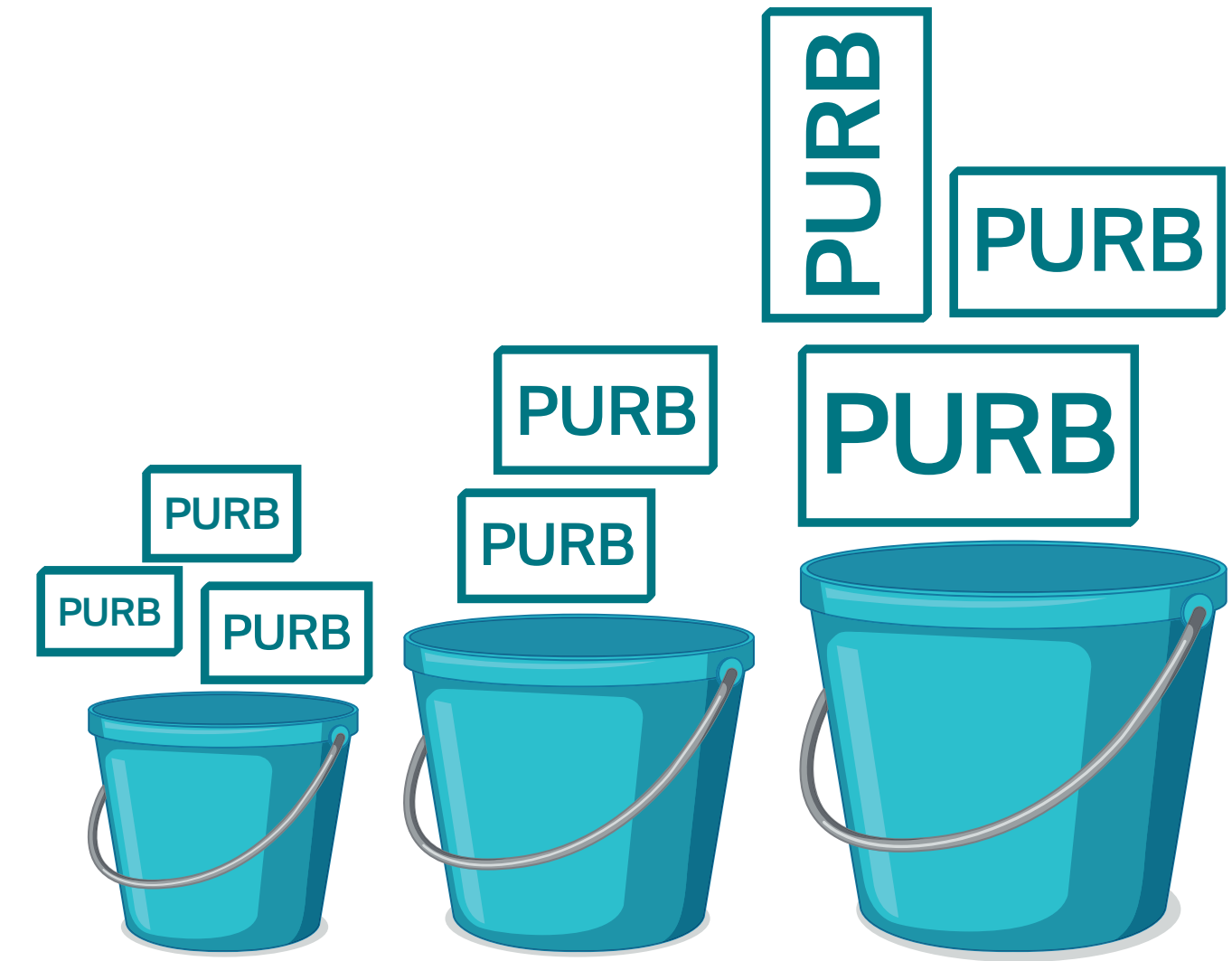
# Encoding and Decoding of PURBs



# Encoding and Decoding of PURBs

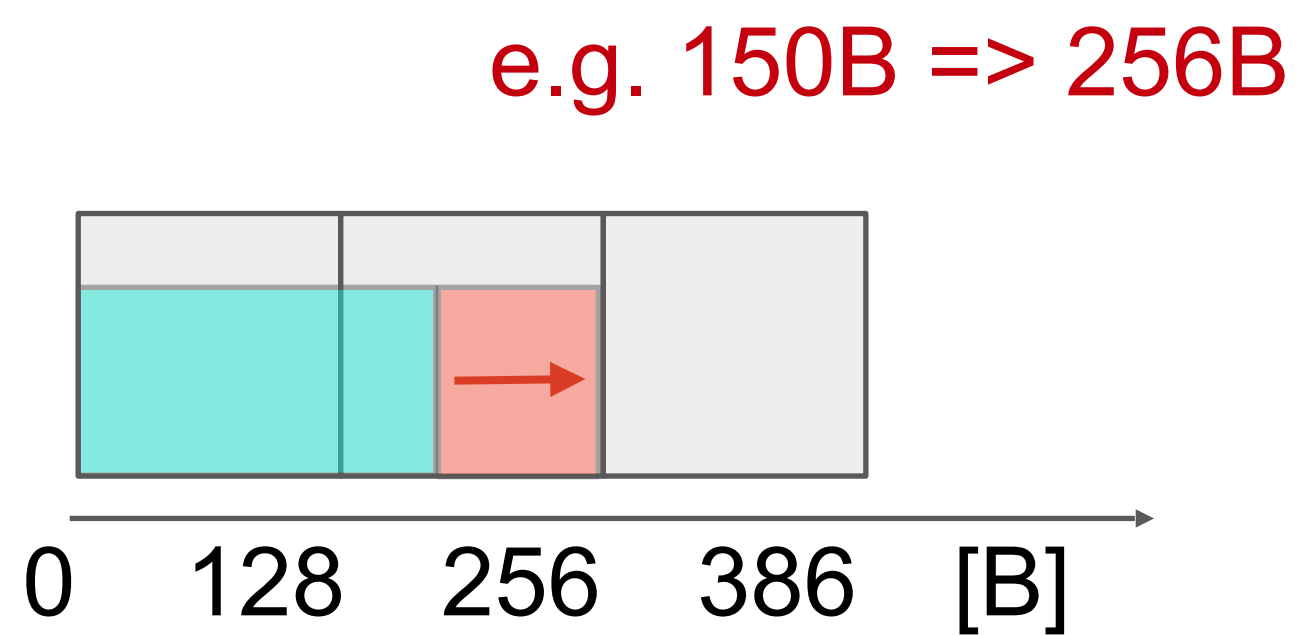


# Padmé: reducing leakage about the size



- The total size is an **important metadata**, used in many attacks:
  - Website Fingerprinting
  - Traffic-Analysis
  - Attacks against HTTPS
- Design a padding function to improve “**size privacy**”

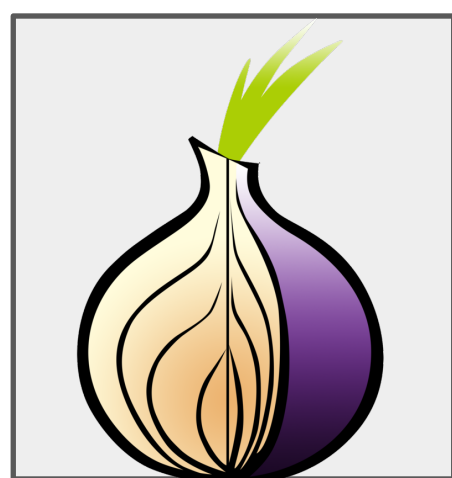
# Naïve approach: (constant) block-padding



# Naïve approach: (constant) block-padding

Problem: no good value for block size

Example:  $b = 1 \text{ MB}$



512 B



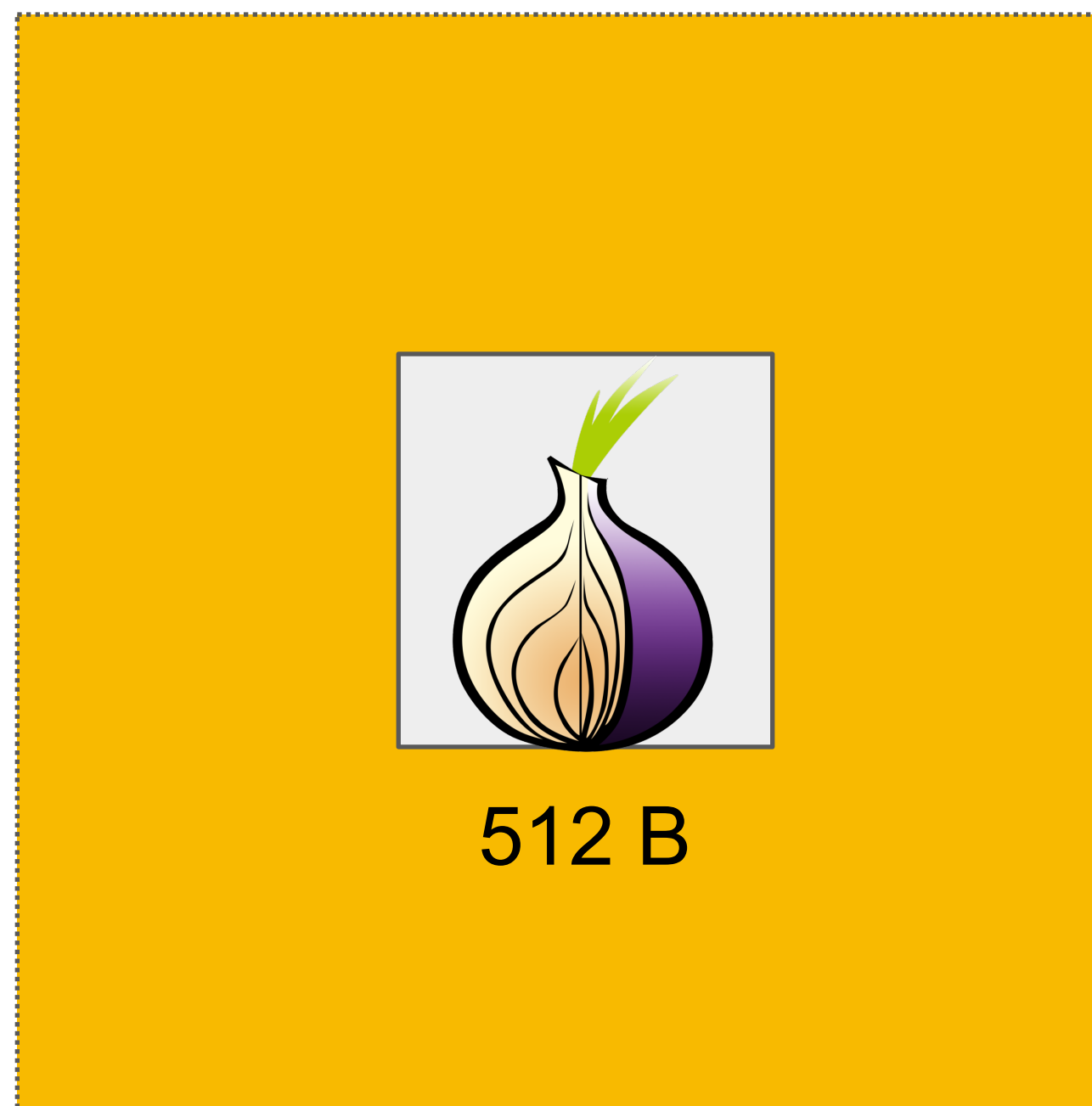
17 GB



# Naïve approach: (constant) block-padding

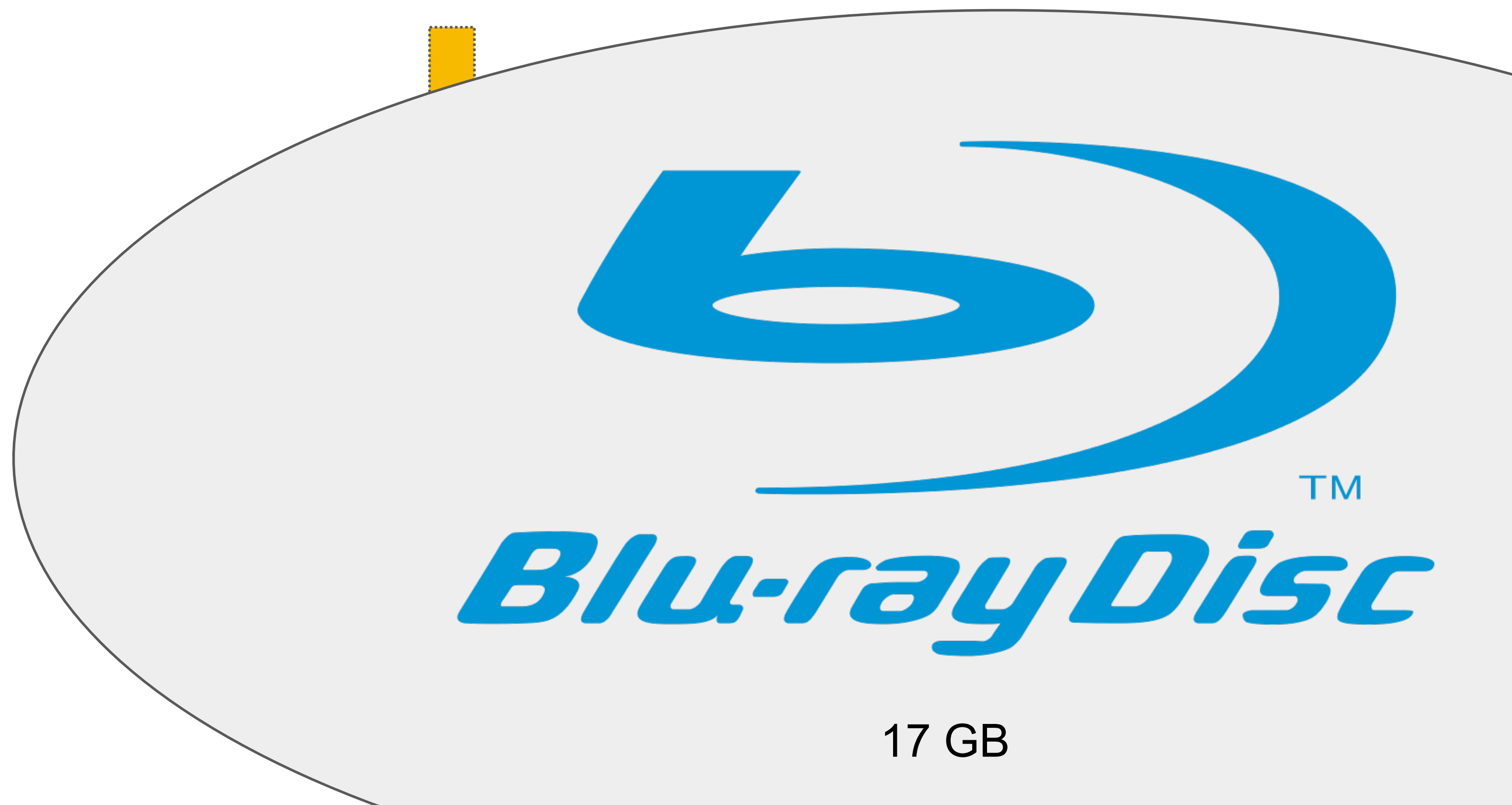
Problem: no good value for block size

Example:  $b = 1 \text{ MB}$



2000x overhead 

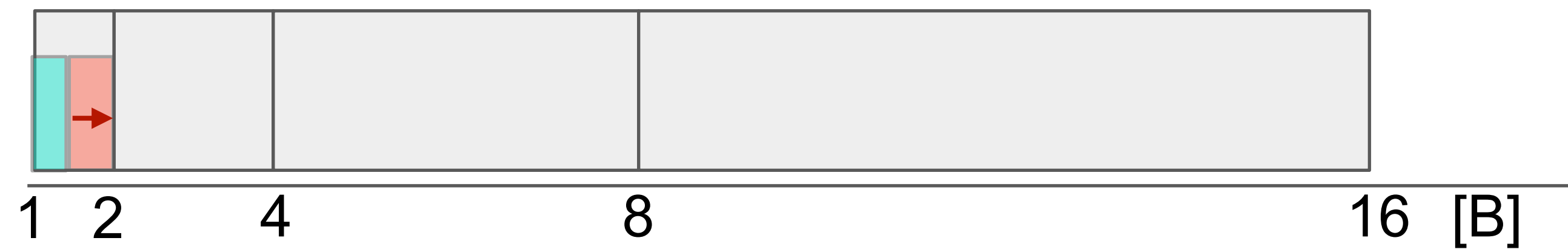
little confusion / privacy 



# Padding relative to the object size

Variable block size:

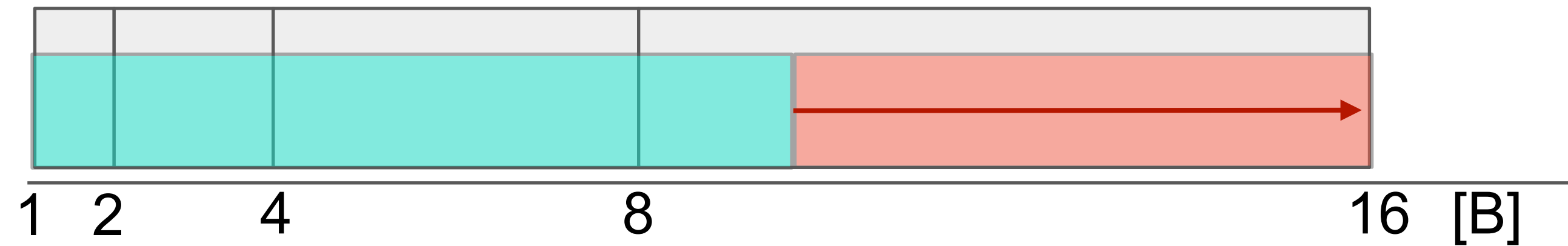
small objects: small overhead 👉



# Padding relative to the object size

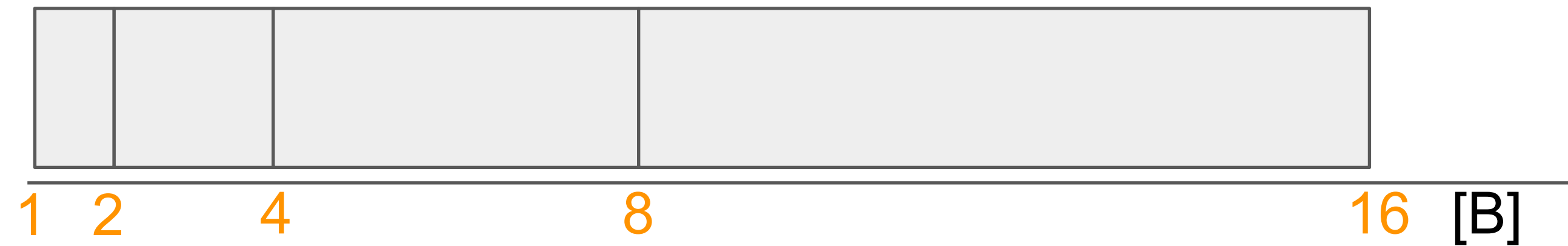
Variable block size:

large objects: larger privacy



# Padding relative to the object size

Variable block size:



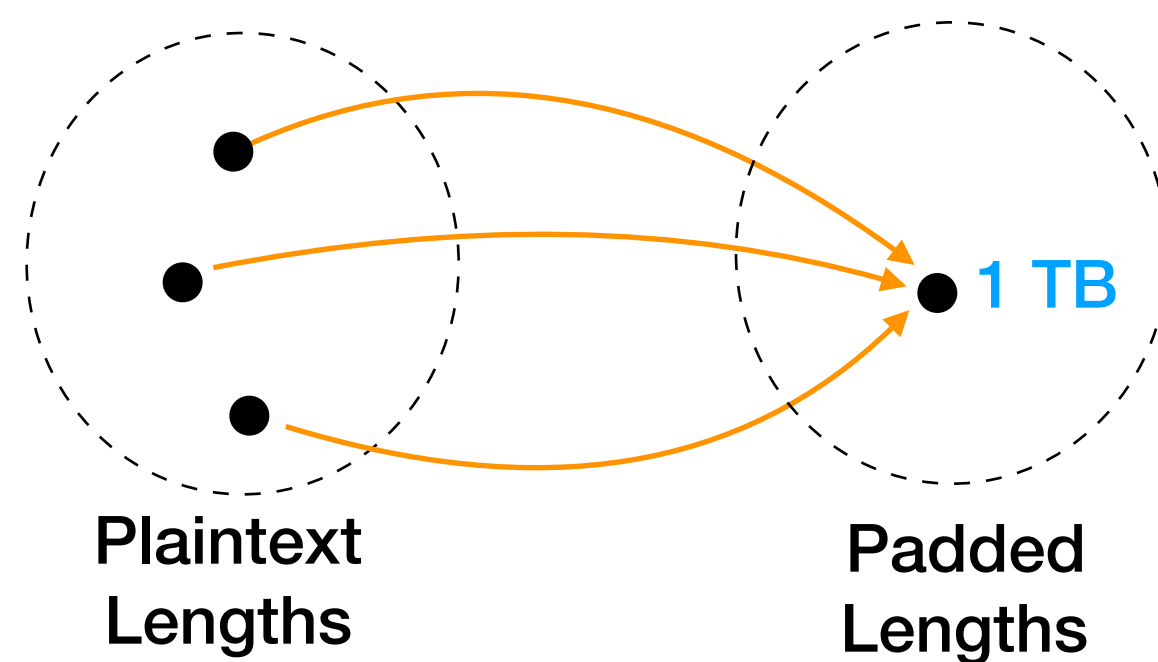
Padding to the "Next Power of 2"

# Quantifying leakage of a padding function

- Let  $f: \mathbb{N} \rightarrow \mathbb{N}$  be the padding function. Let  $C$  be the image set of  $f$ .

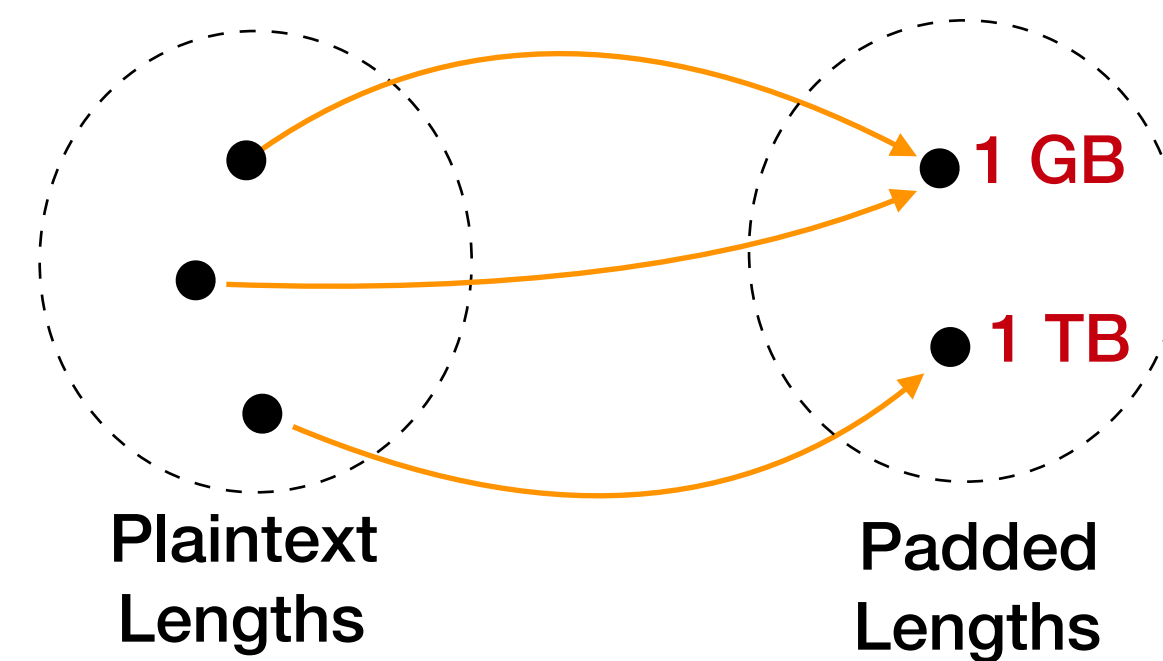
$$\text{Leakage [bits]} = \log_2(|C|)$$

$$f(p) = 1 \text{ TB}$$



$$\text{Leakage: } \log_2(1) = 0 \text{ bit}$$

$$f(p) = \begin{cases} 1 \text{ GB} & p \leq 1 \text{ GB} \\ 1 \text{ TB} & p > 1 \text{ GB} \end{cases}$$



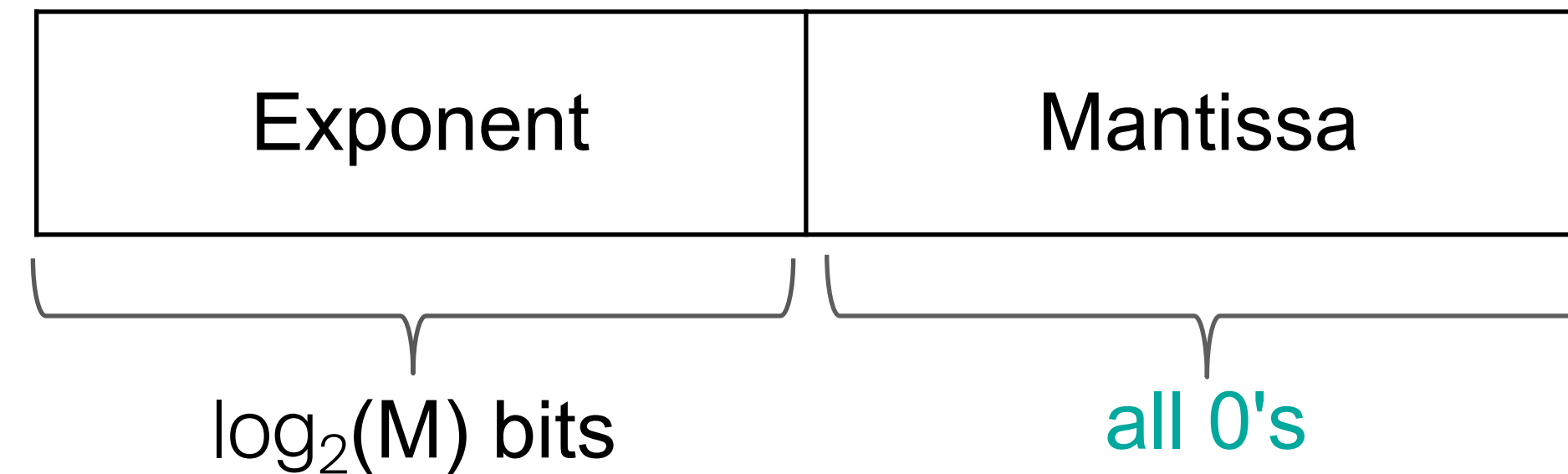
$$\text{Leakage: } \log_2(2) = 1 \text{ bit}$$

# Padding to the nearest power of 2

- Leakage:  $O(\log \underbrace{\log(M)}_{\text{size of the image set}})$ , where  $M$  is the biggest plaintext possible
- 👍 Much better than with constant block-size, which leaks  $O(\log(M))$
- Interestingly, not padding at all also leaks  $O(\log(M))$
- **Max overhead: +100%** 📉
  - e.g., 16.1 GB => 32 GB

# Reducing the overhead of padding

Next power of 2: Blocks have the form  $2^0$ ,  $2^1$ ,  $2^2$ , etc. Represent **them** like floating points:

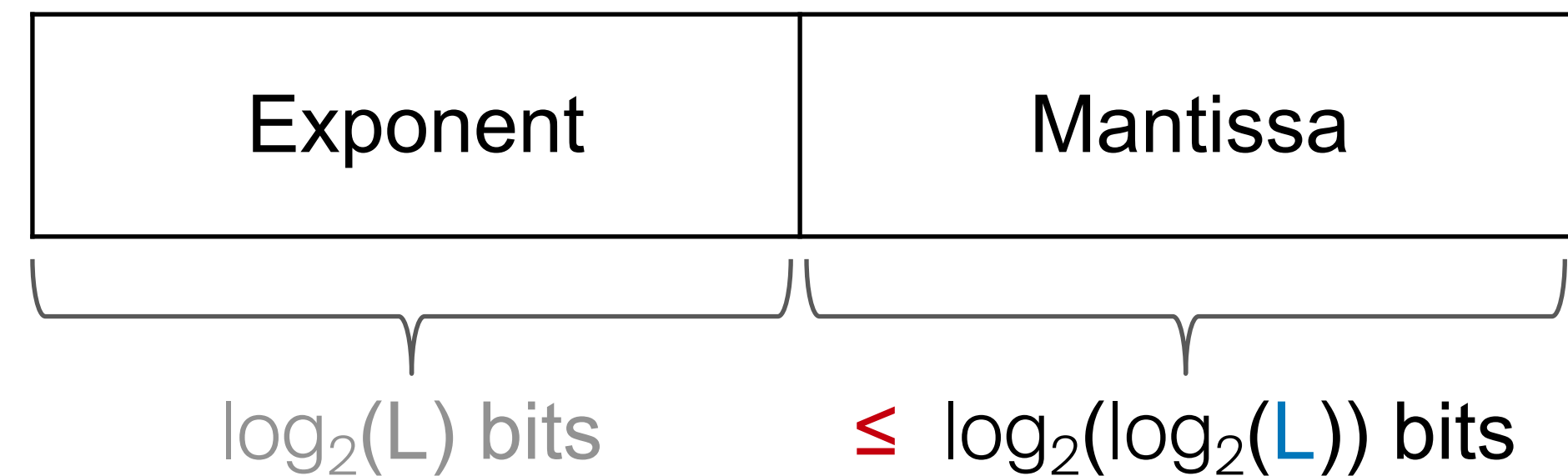


Padmé: Instead of 0's, allow **some** values in the mantissa:

=> Smaller blocks => **Smaller overhead** 👉



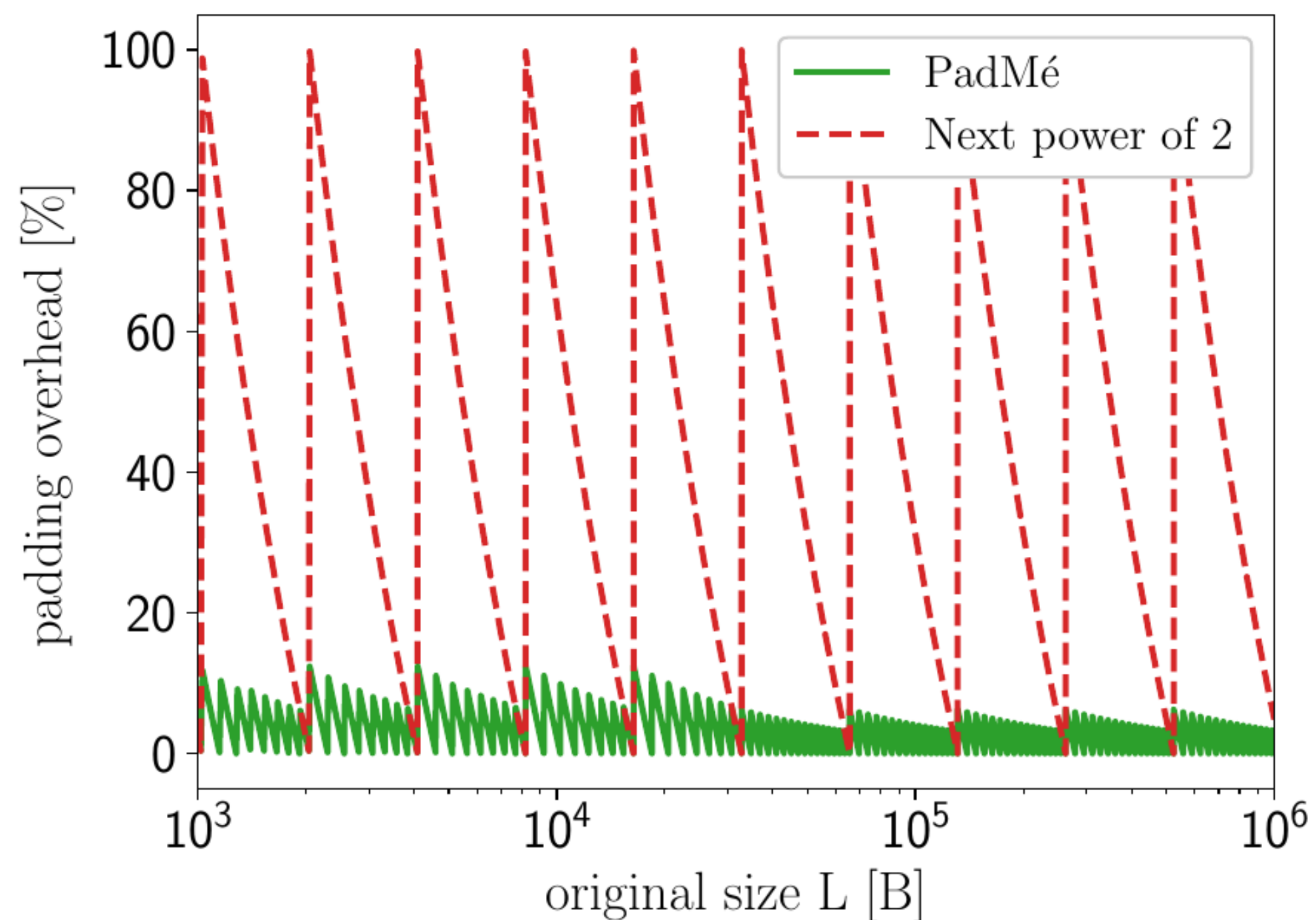
Pad to the next length  $L$  which respects:



Intuition: the exponent can be anything, but the mantissa cannot be "too precise"

Doubles leakage  $\Rightarrow$  still in  $O(\log \log(M))$  👉

# Padmé's overhead



$$\text{max overhead} = \frac{1}{2 \log_2 L} \%$$

Slowly decreases with L



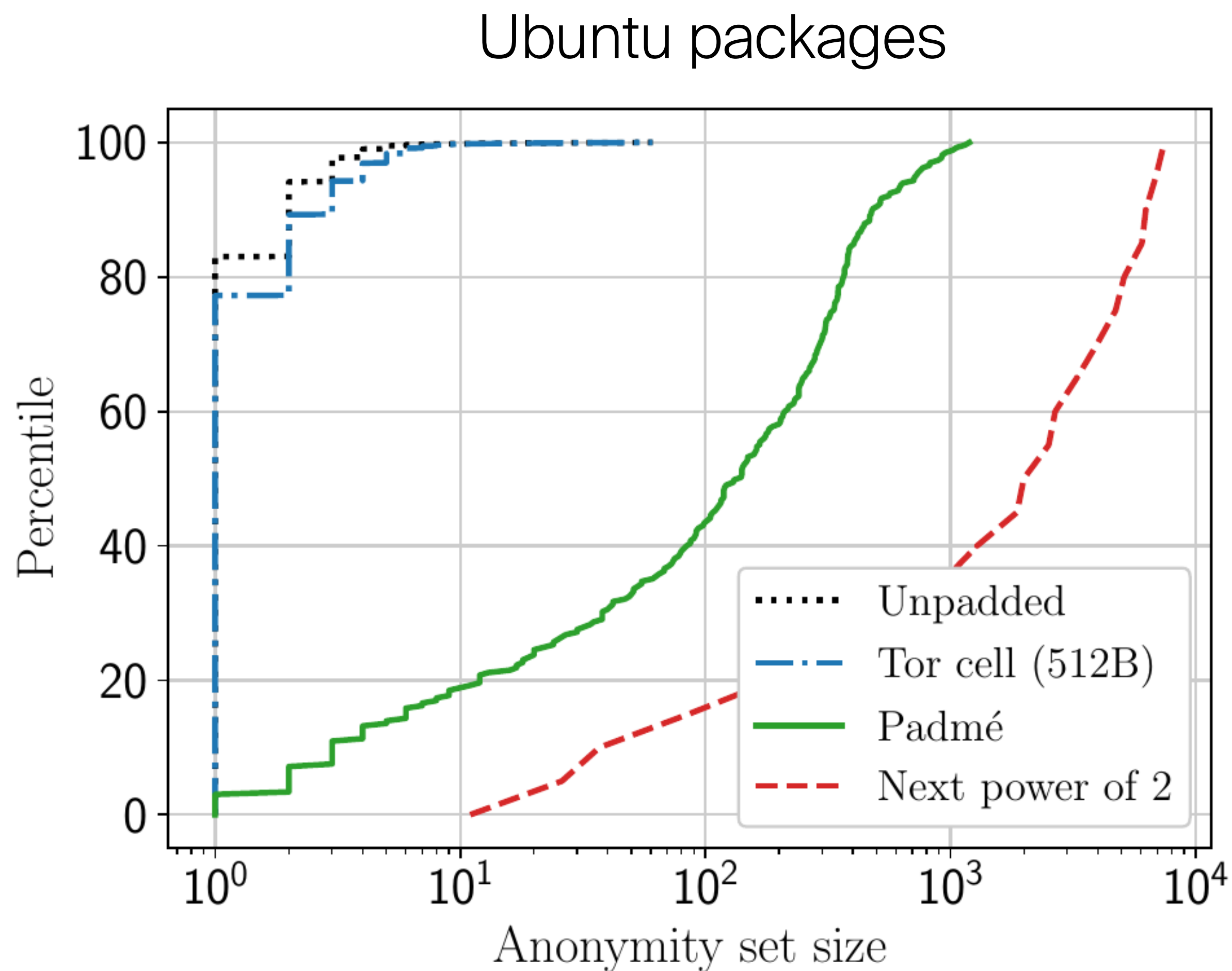
Max 12%  $\forall L$



Max ~6% for  $L \geq \sim 1 \text{ MB}$

Max ~3% for  $L \geq \sim 1 \text{ GB}$

# Padmé's "size privacy"



57'000 objects collected from apt lists

Mean overhead:

Next power of 2: **+44%**

Padmé: **+2.3%**

# Conclusion

- Padded Uniform Random Blobs (PURBs): [ciphertext format with no metadata leakage](#) except length, which is minimized.
- Encoding + Padding schemes.
- Applications: Email, Group Chat, Disk Encryption, Initiation of Protocols.
- To the best of our knowledge, the first video with pets @ PETS.

<https://purbs.net>

<https://github.com/dedis/purb>

[kirill.nikitin@epfl.ch](mailto:kirill.nikitin@epfl.ch) 

 [@ni\\_kirill](#)

[ludovic.barman@epfl.ch](mailto:ludovic.barman@epfl.ch) 

 [@lbarman\\_ch](#)