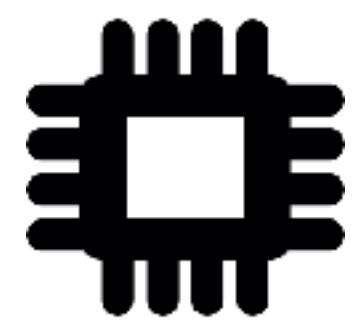


# **Chucky: A Succinct Cuckoo Filter for LSM-Tree**

**Niv Dayan, Moshe Twitto**



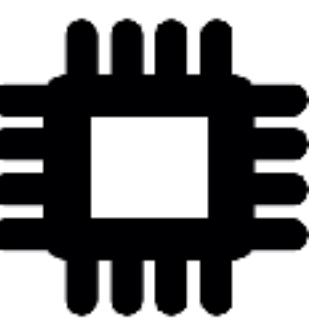
**PLIOPS**  
EXTREME DATA PROCESSOR



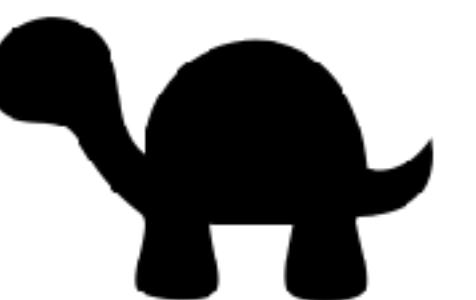
memory

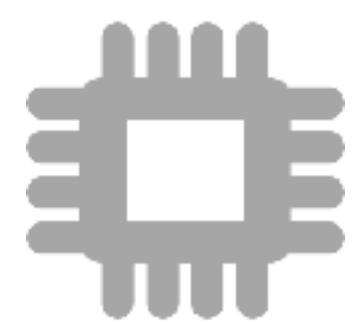


storage



- - -



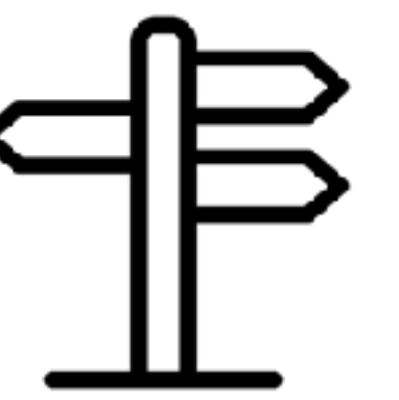
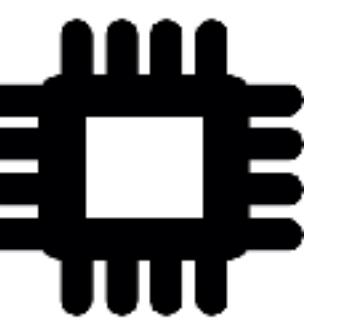


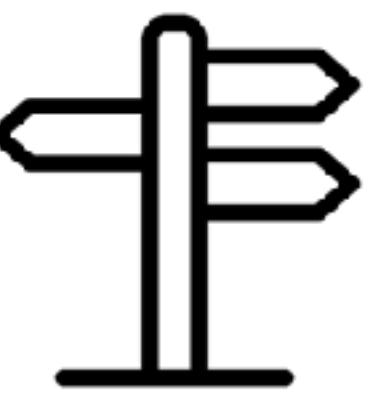
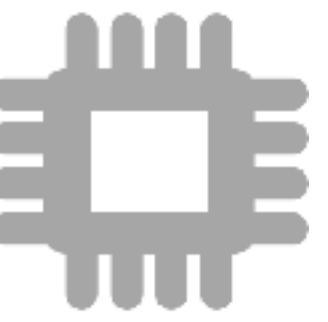
**metadata**

fetch

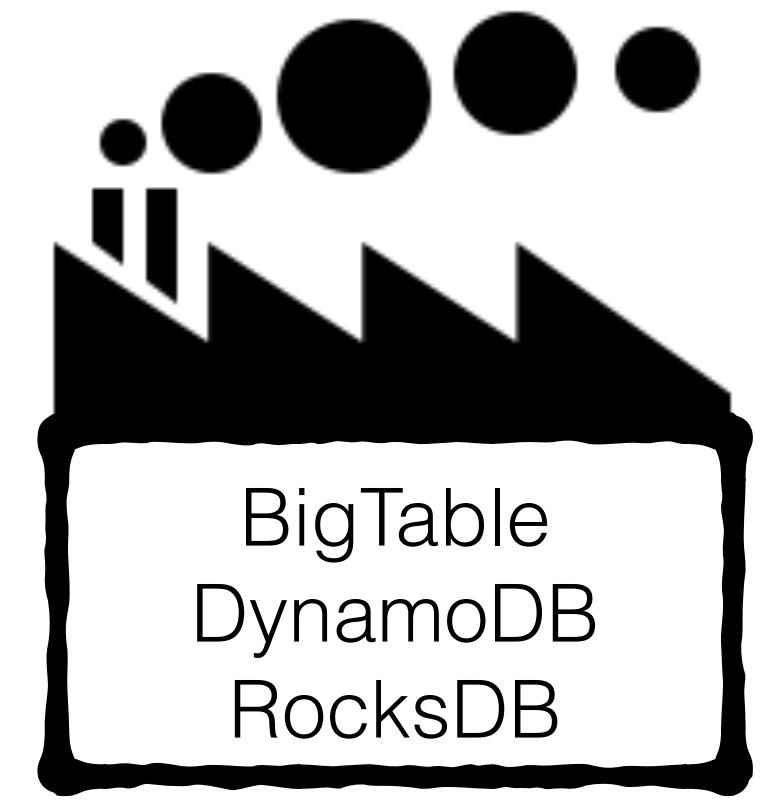
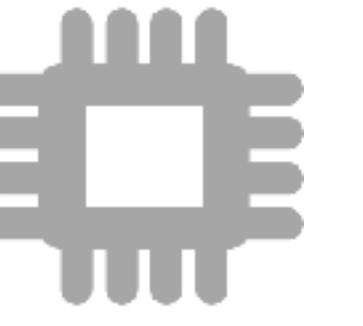


**data**





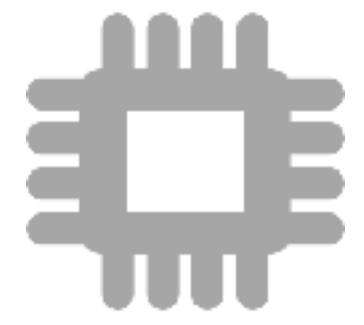
**noticeable  
overhead**



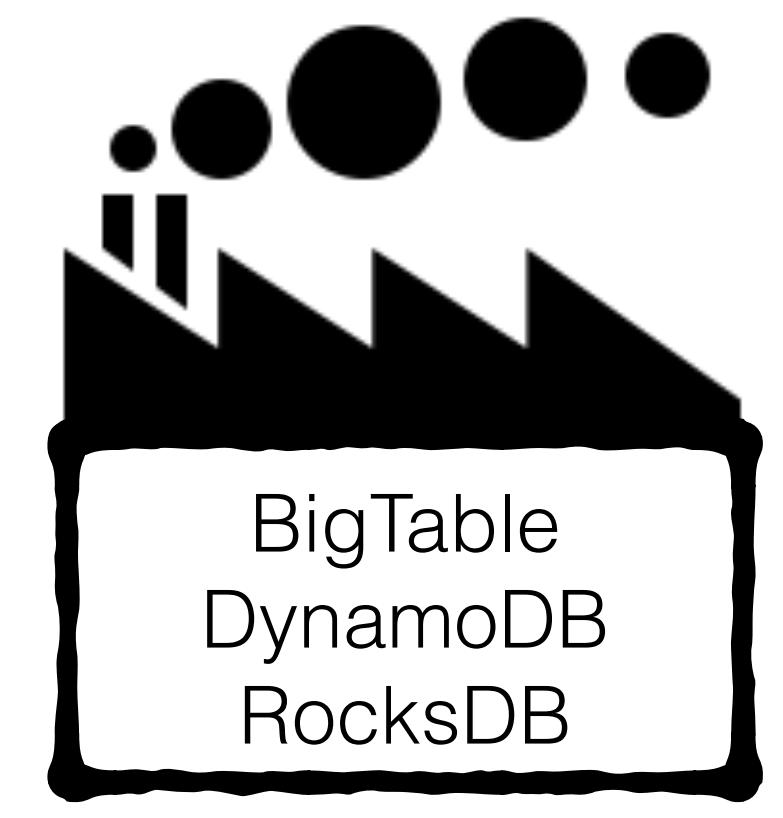
**LSM-tree**



**Bloom filters**



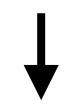
**LSM-tree**



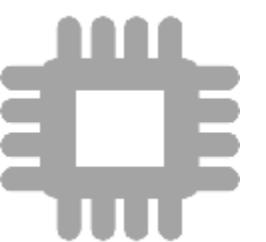
BigTable  
DynamoDB  
RocksDB

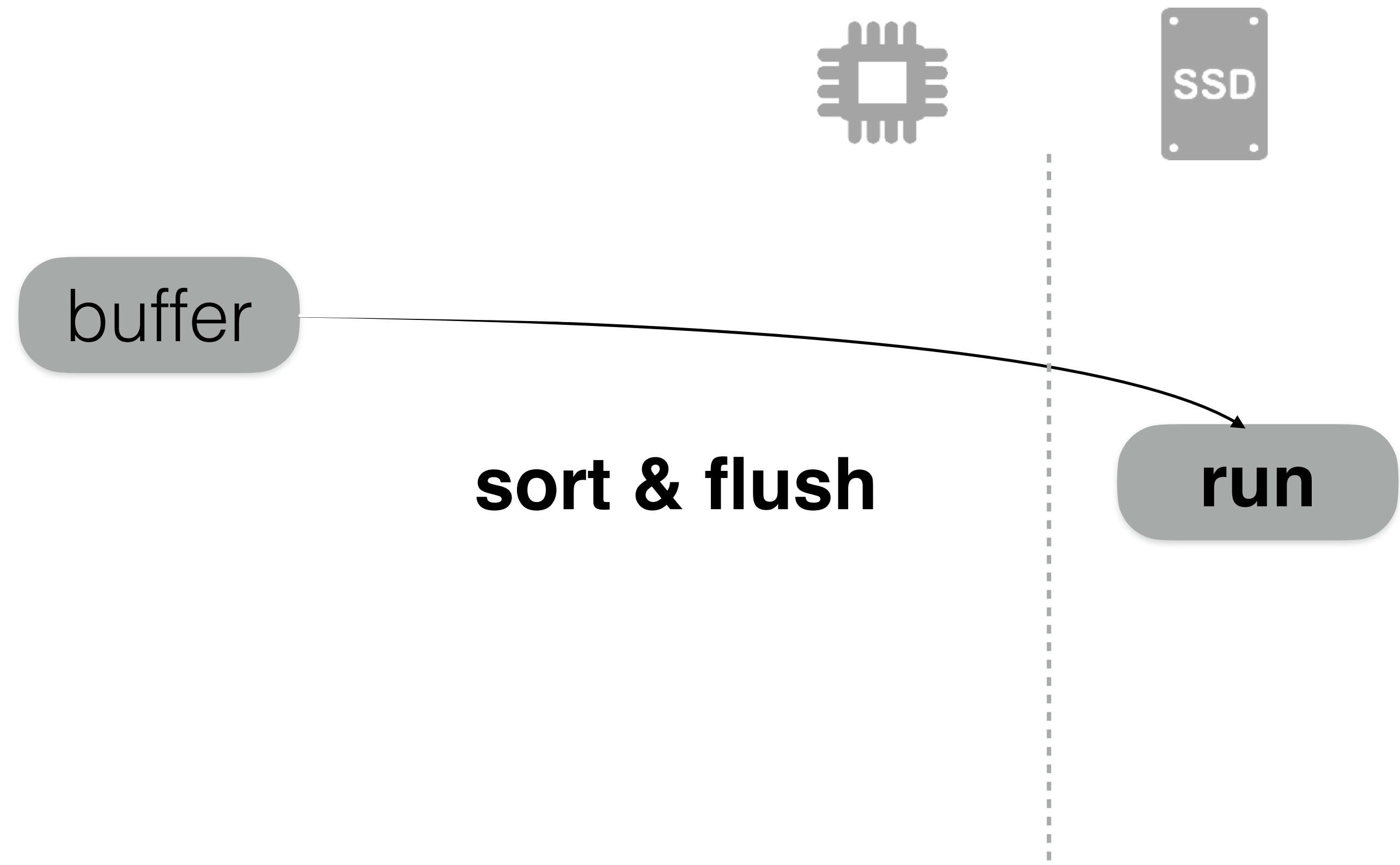
# LSM-tree

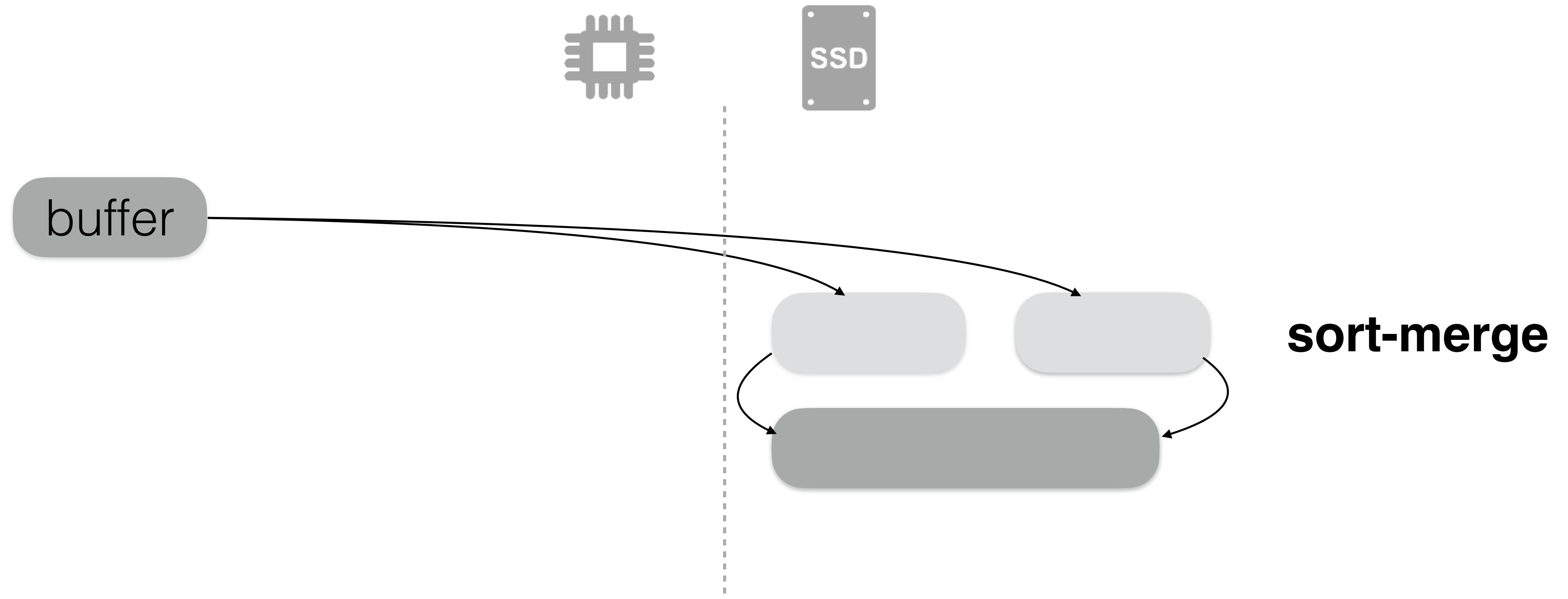
key-value pairs



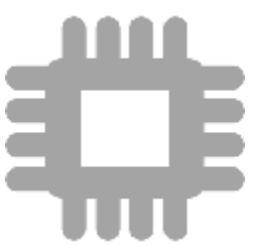
buffer







buffer



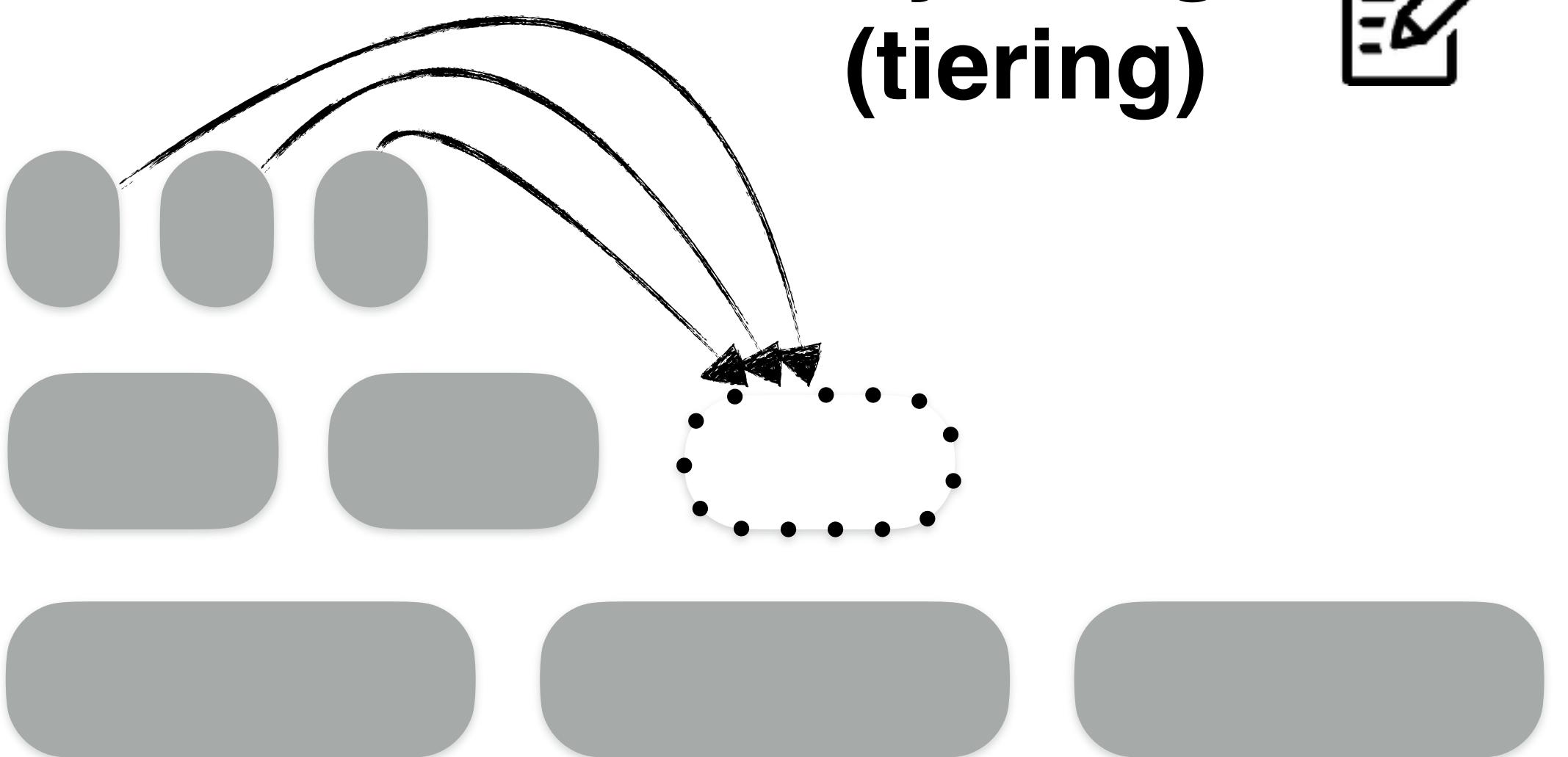
**level 1 →**

**level 2 →**

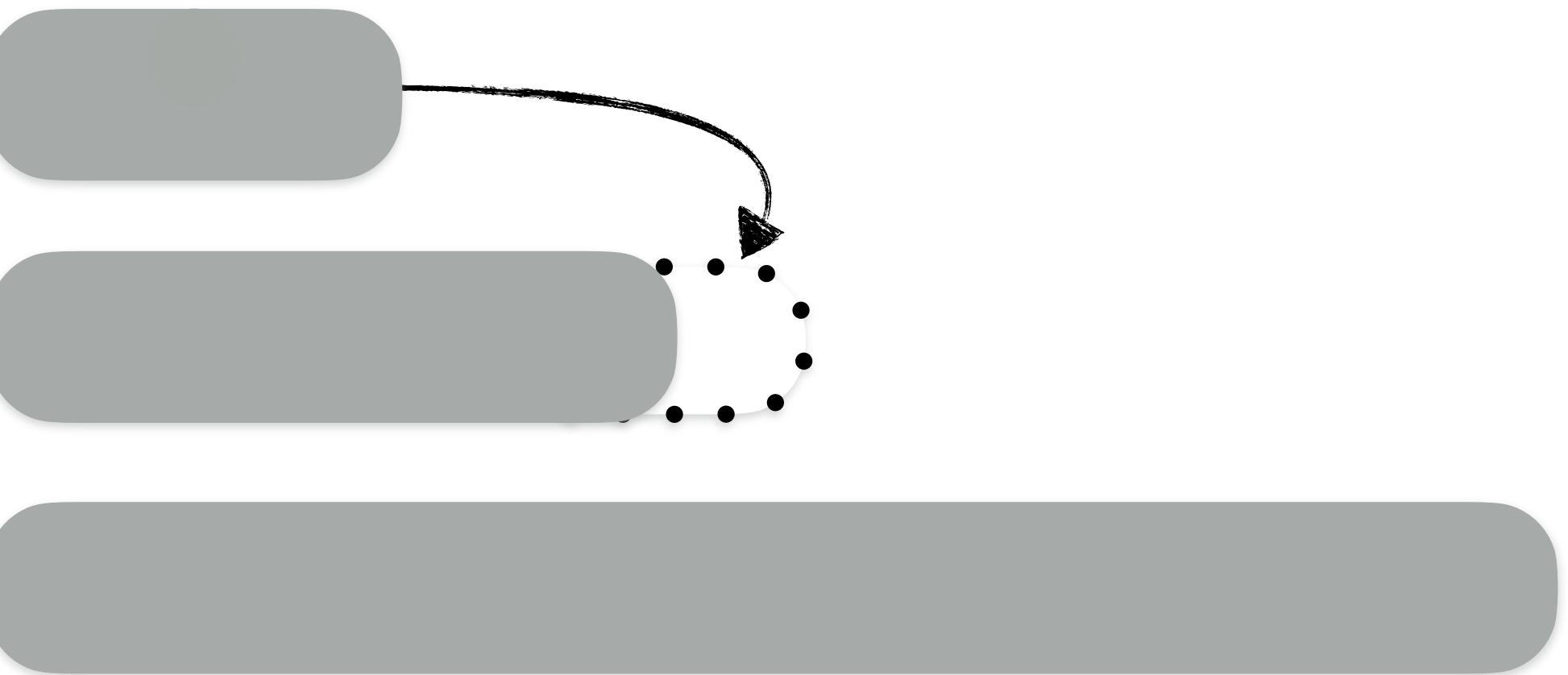
**level 3 →**

**exponentially increasing capacities**

**lazy merge  
(tiering)**



**greedy merge  
(leveling)**

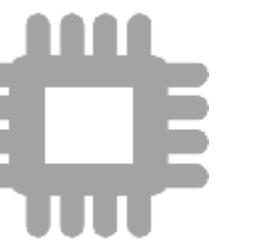


**Levels  $L = \log(N)$**

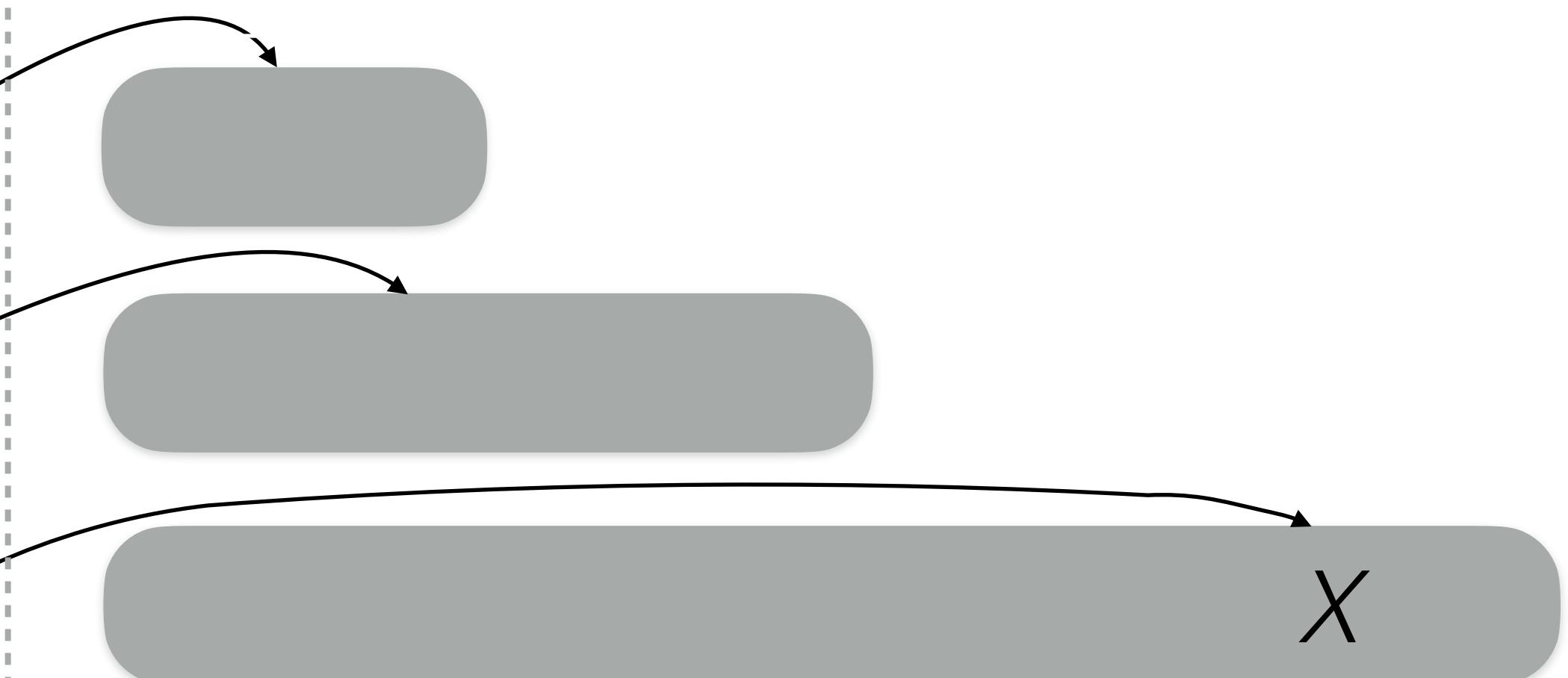


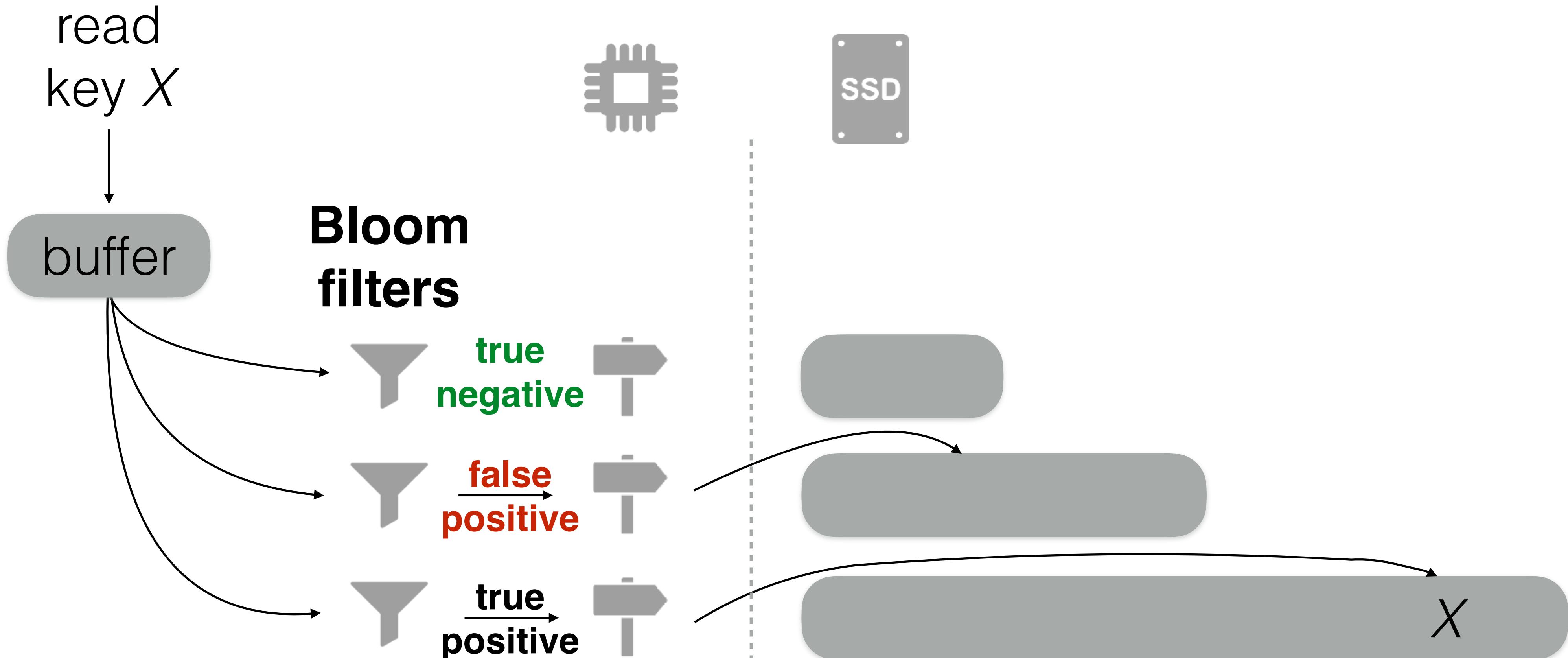
read  
key  $X$

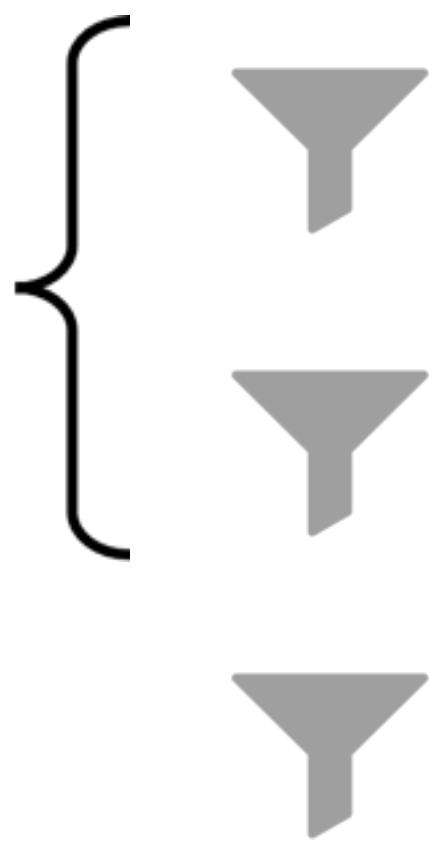
buffer



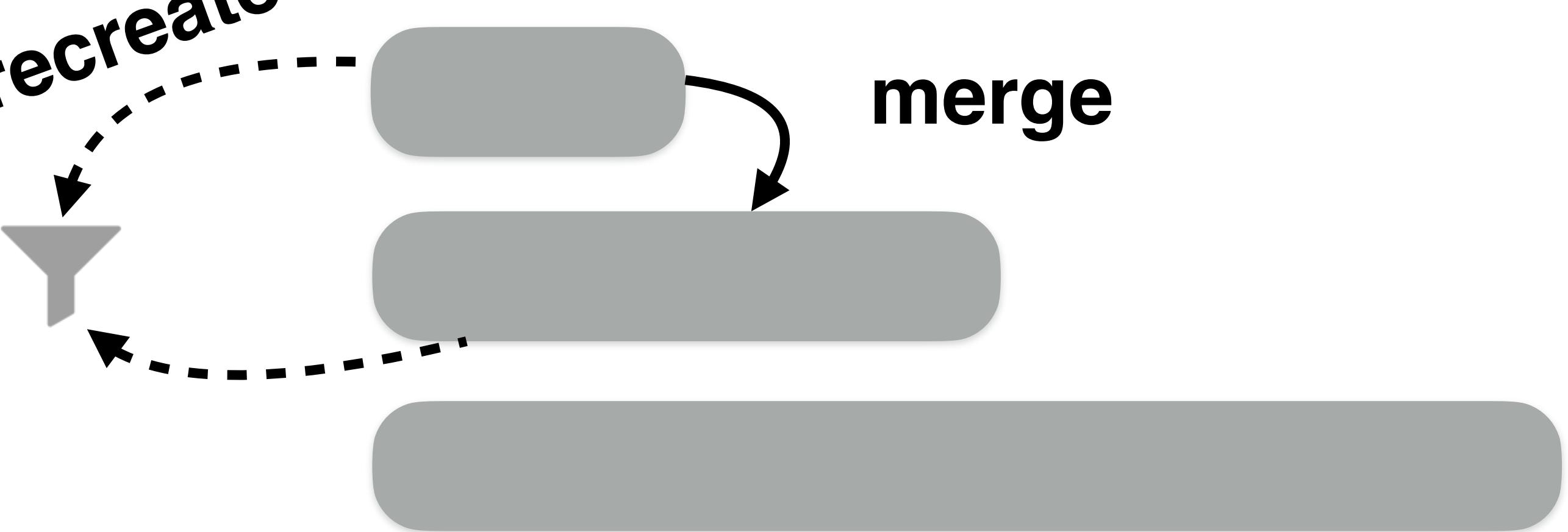
**fence  
pointers**







*recreate*



**Bloom Filters' WA = LSM-tree's WA**



**false  
positive rate**

$O(f \cdot \log N)$

$f$   
 $f$   
 $f$

Bloom  
filters



# **Monkey** **SIGMOD17**

$$O(\mathbf{f}) \quad \left\{ \begin{array}{l} f / \mathbf{r}^2 \\ f / \mathbf{r}^1 \\ f \end{array} \right.$$

Bloom  
filters

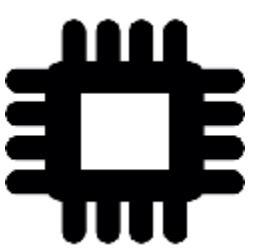


$O(e^{-M} \cdot \ln(2))$

Bloom  
filters

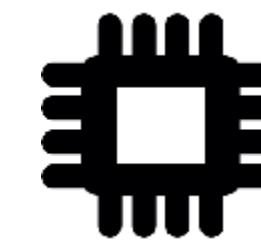


# Problem: shrinking memory/storage speed gap

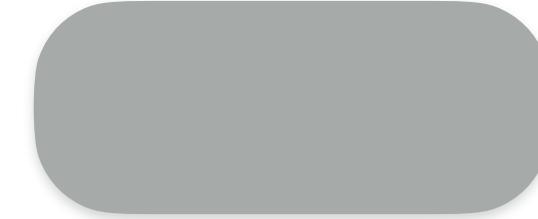


# Problem: shrinking memory/storage speed gap

$\approx 100 \text{ ns}$

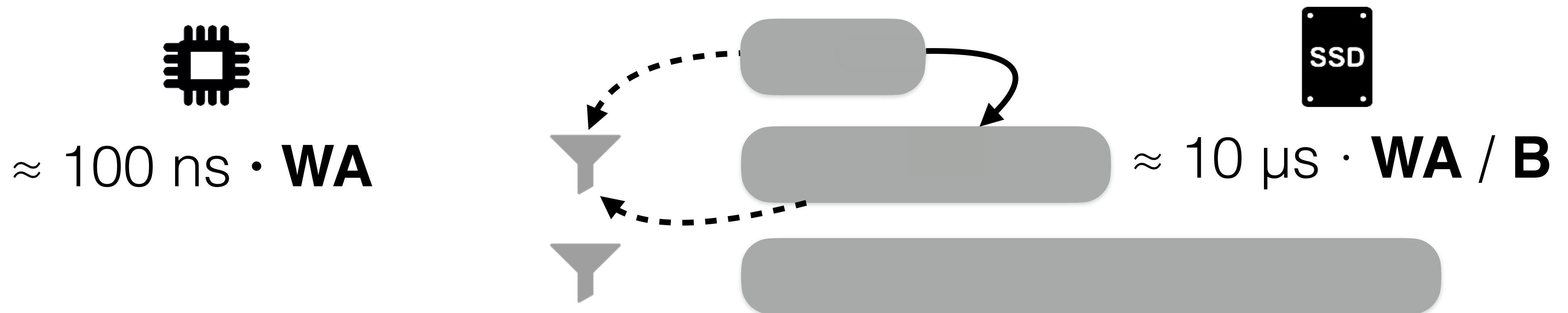


$\approx 10 \mu\text{s}$  (optane)



Problem: shrinking memory/storage speed gap

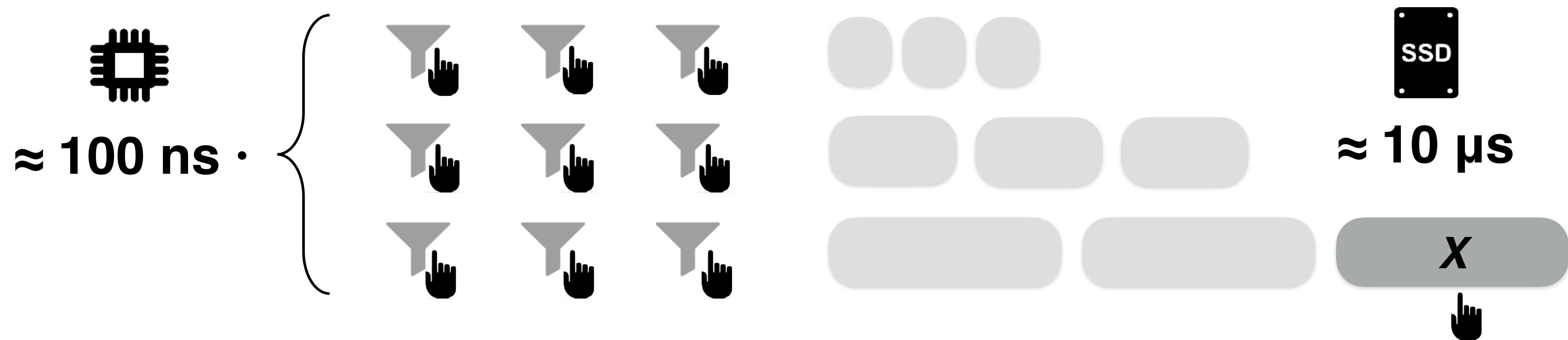
## (1) write amplification



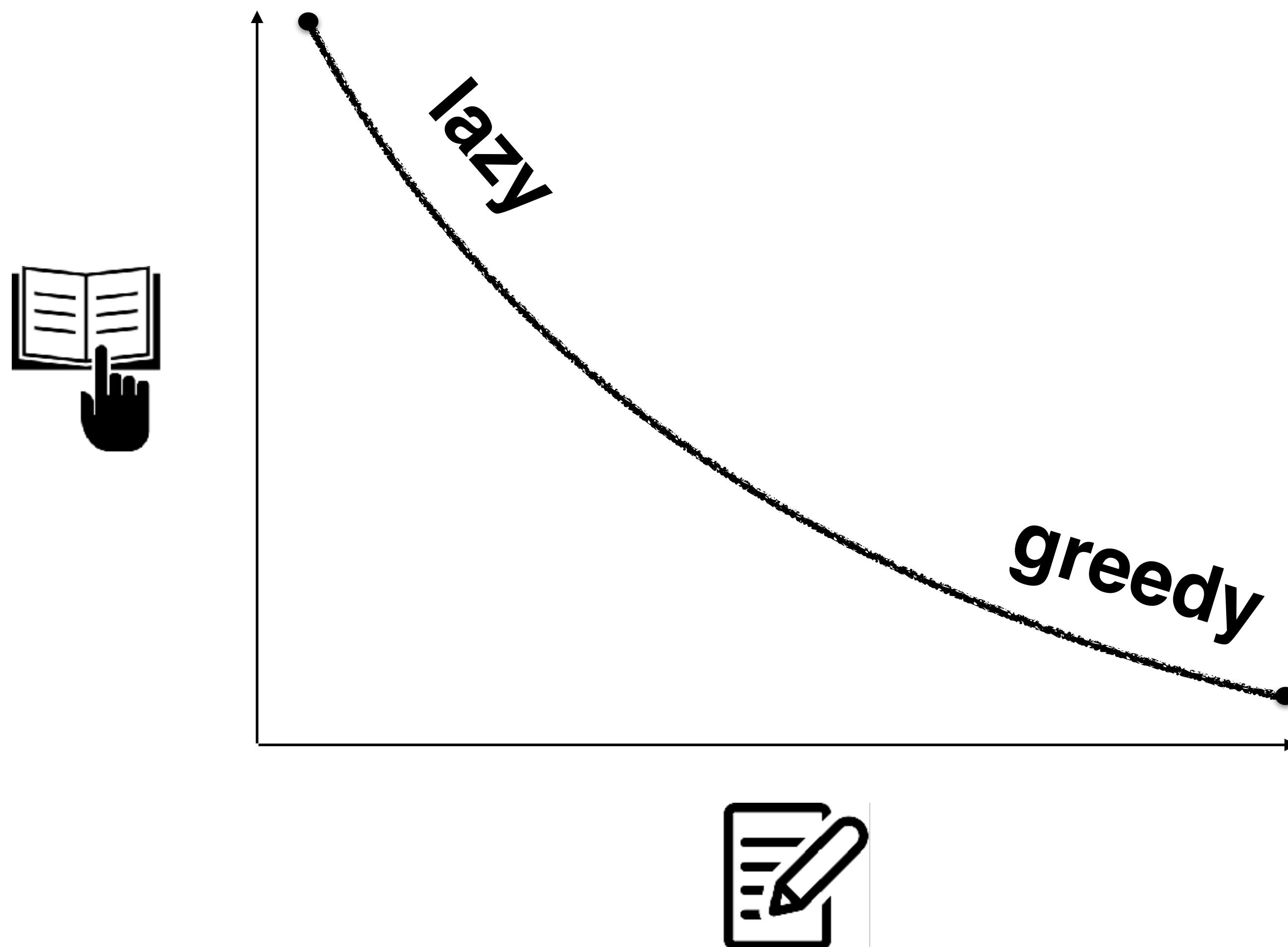
Problem: shrinking memory/storage speed gap

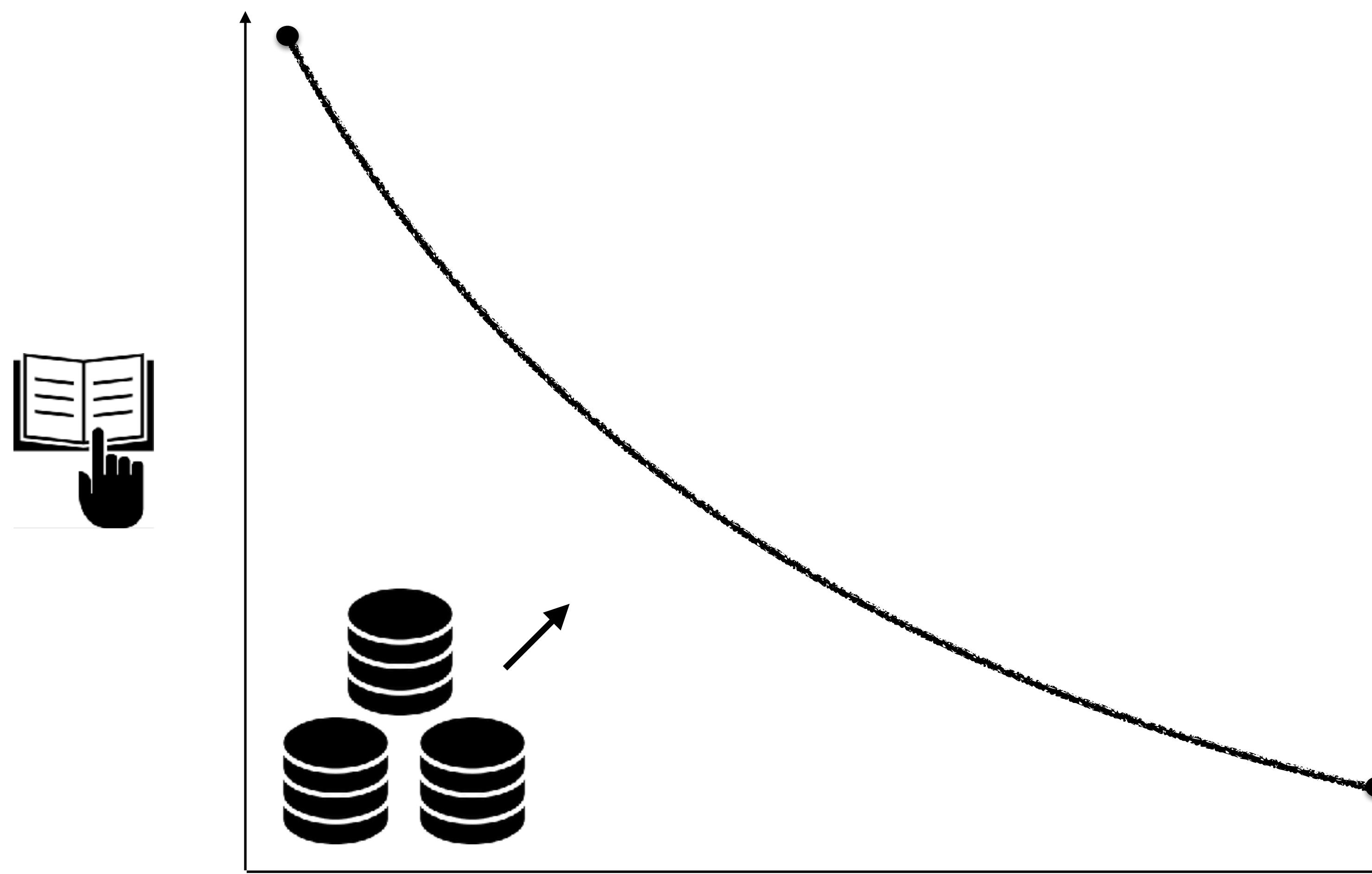
(1) write amplification

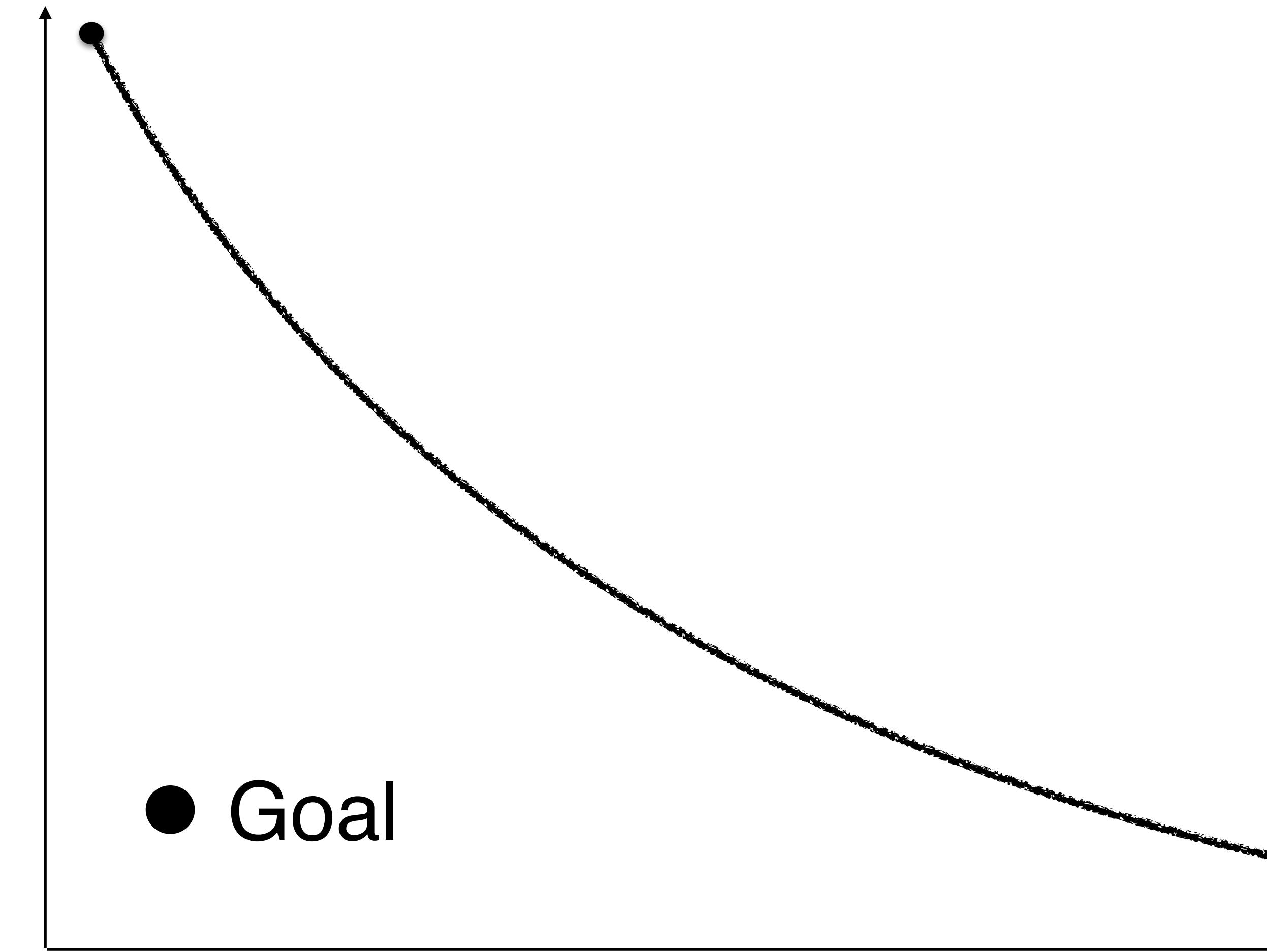
**(2) read amplification**



## ➤ read / write contention

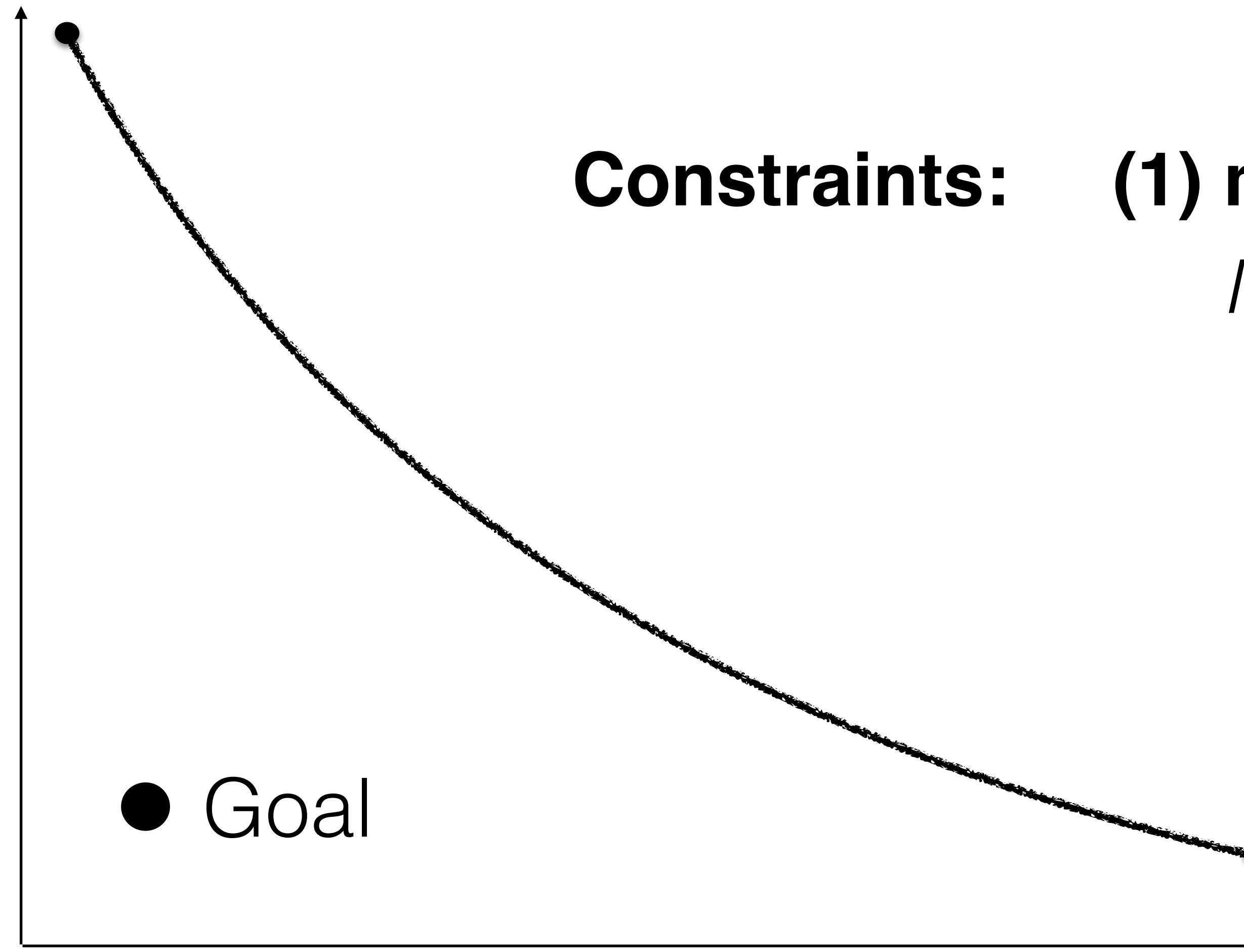




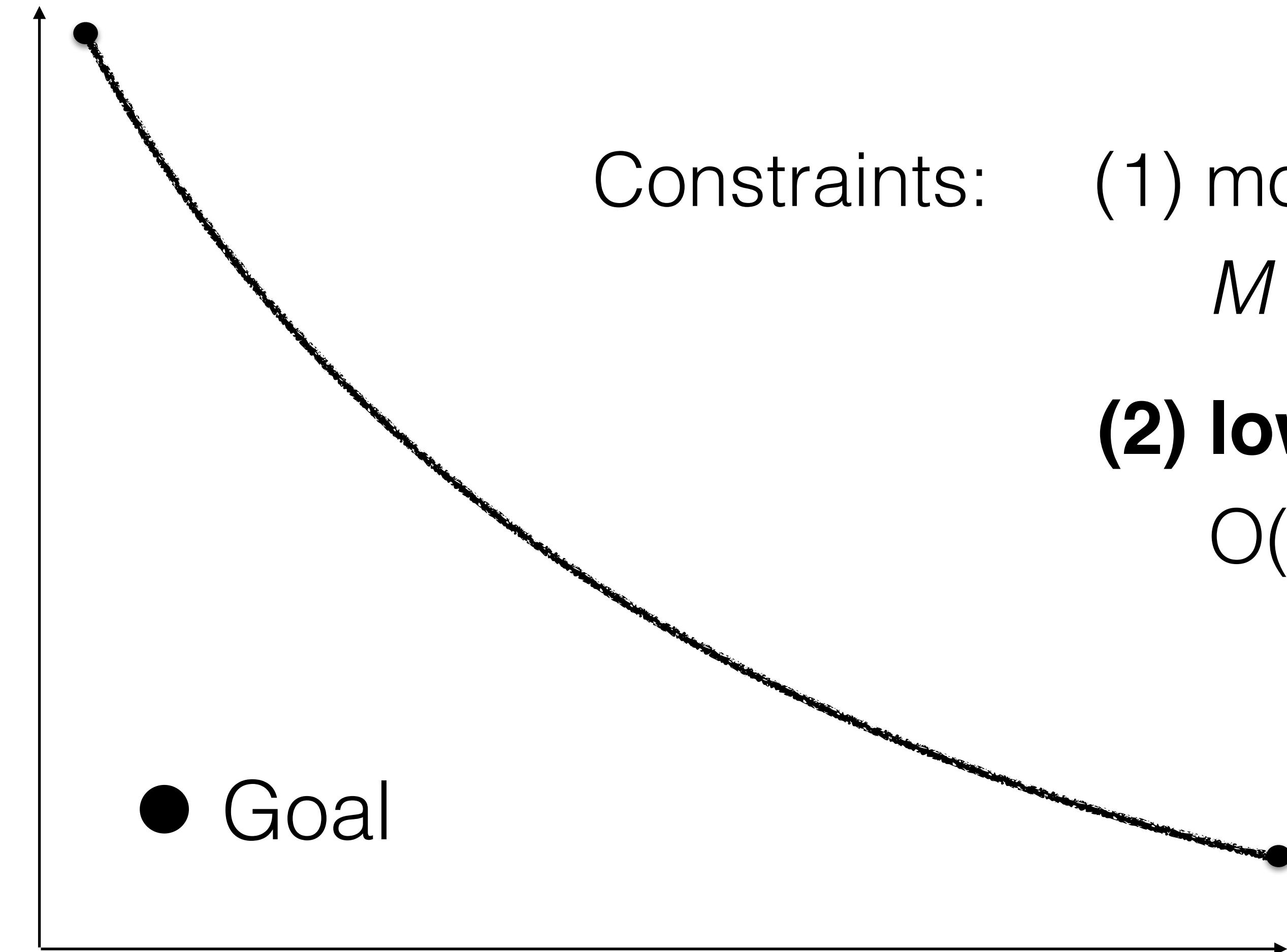




● Goal



**Constraints:** **(1) modest memory footprint**  
 $M \approx 10 \text{ bits / entry}$



Constraints: (1) modest memory footprint

$$M \approx 10 \text{ bits / entry}$$

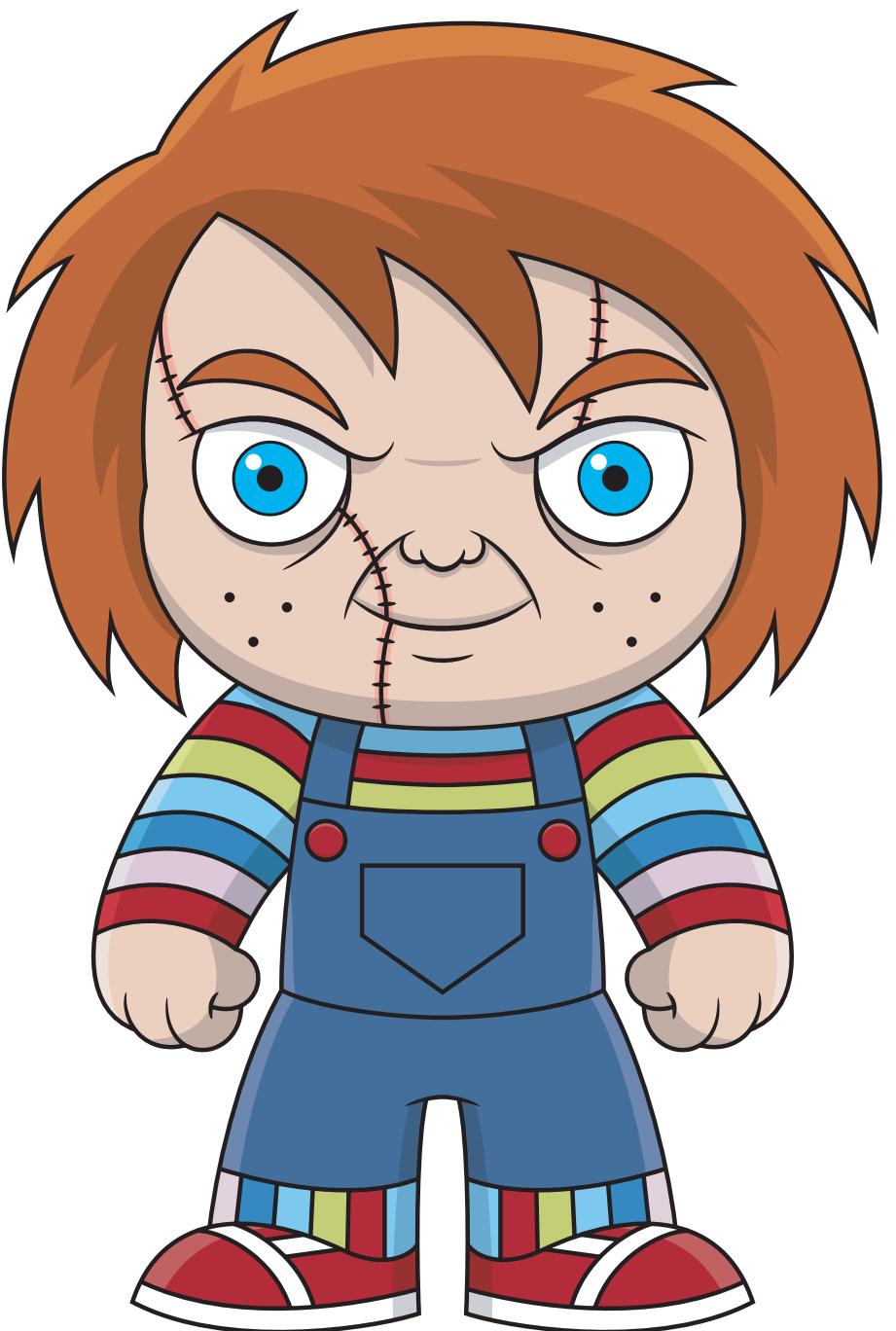
**(2) low false positive rate**

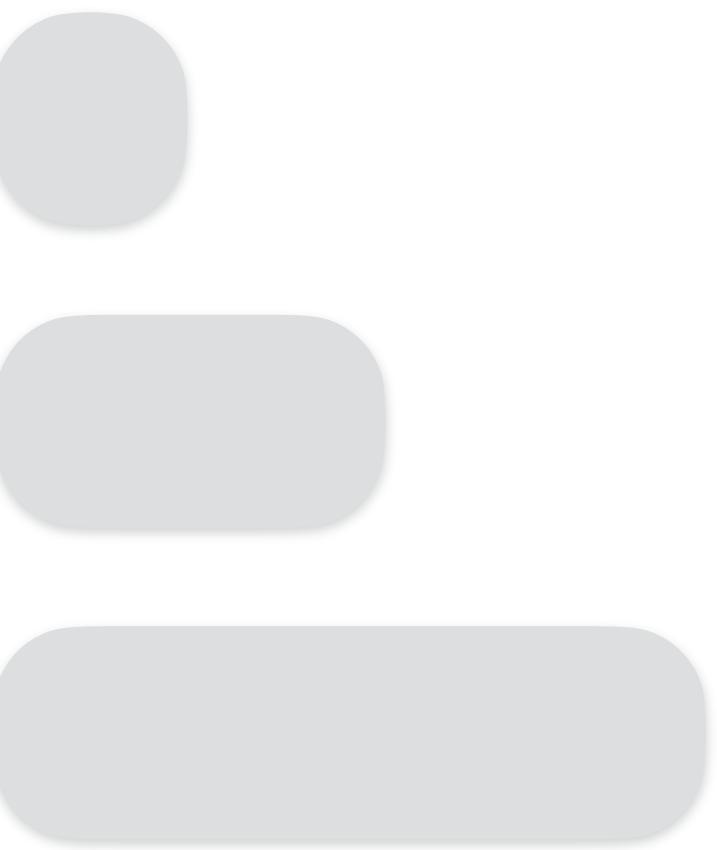
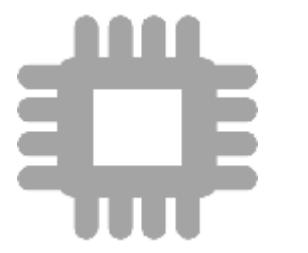
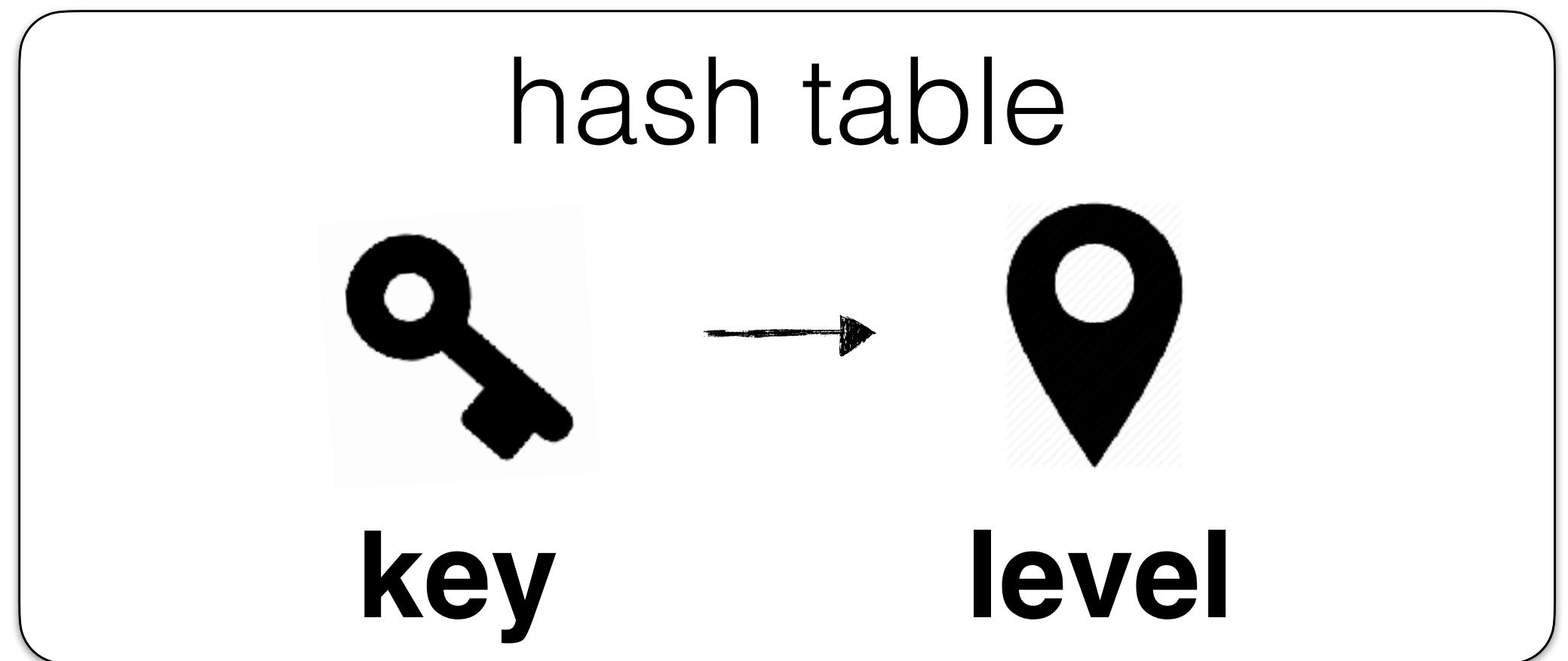
$$O(e^{-M} \cdot \ln(2))$$

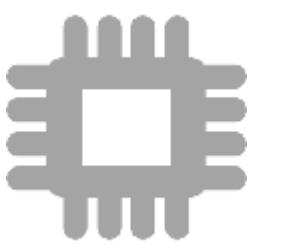
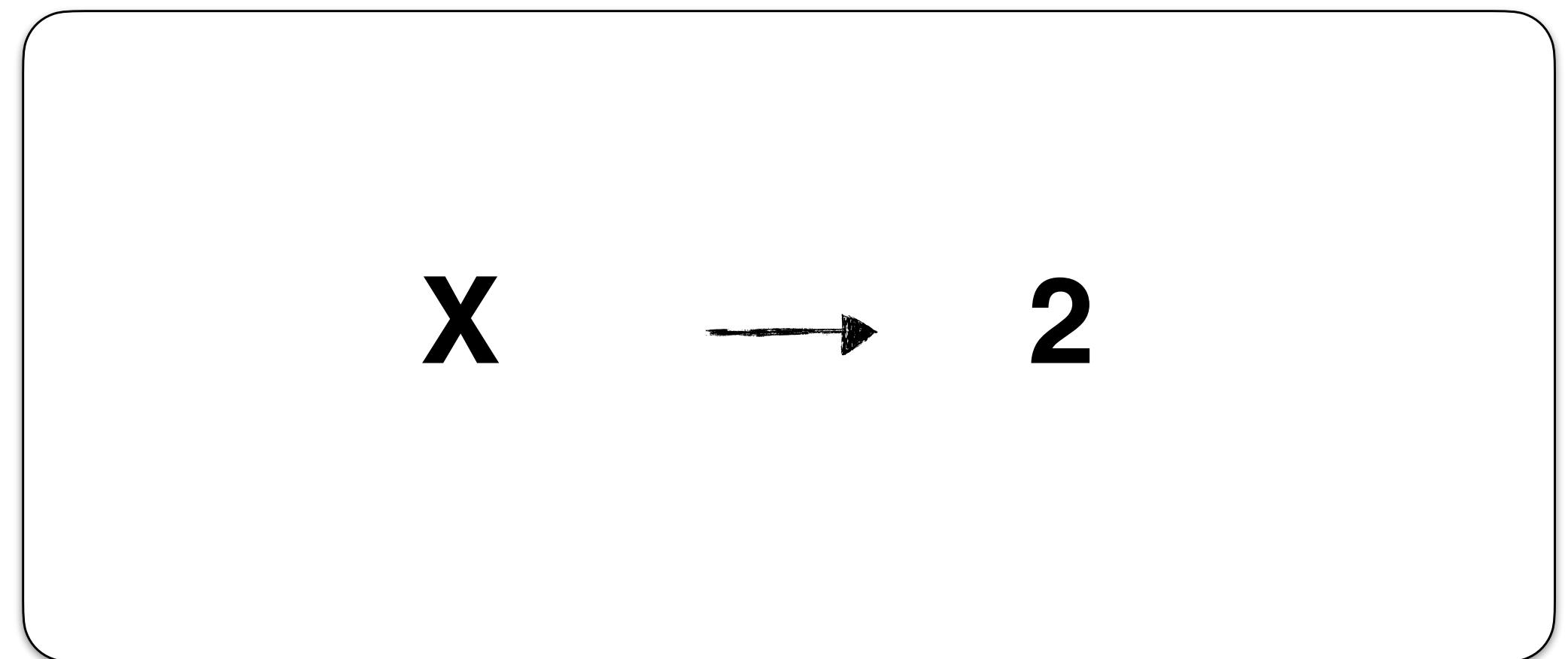
# Chucky



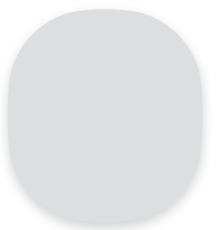
# Chucky: huffman coded key-value store



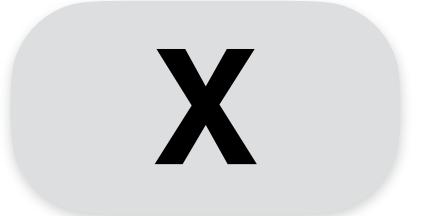




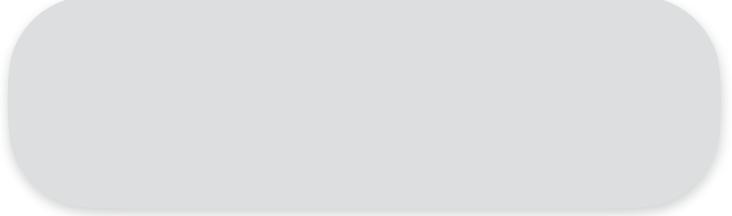
level ID



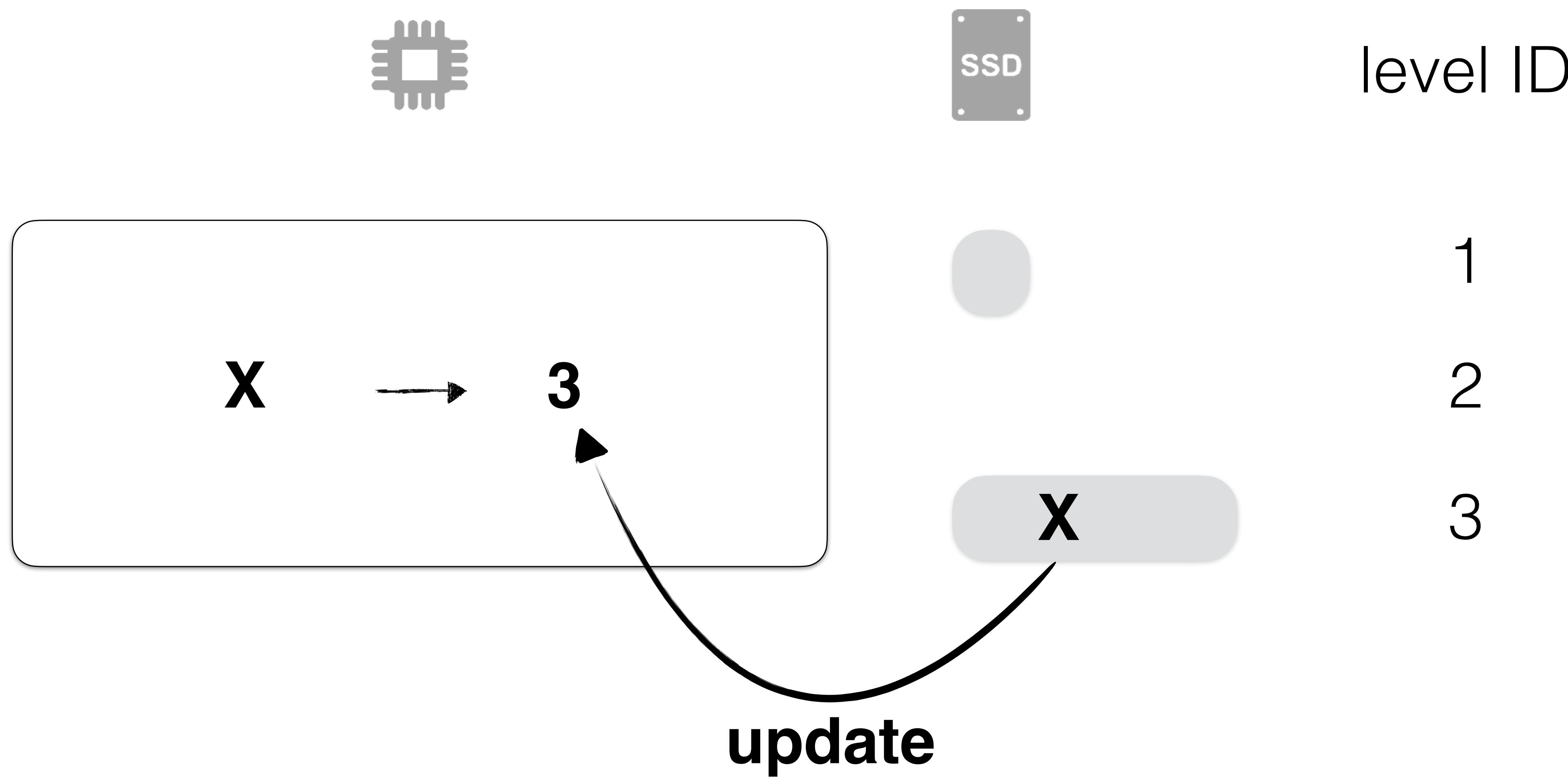
1

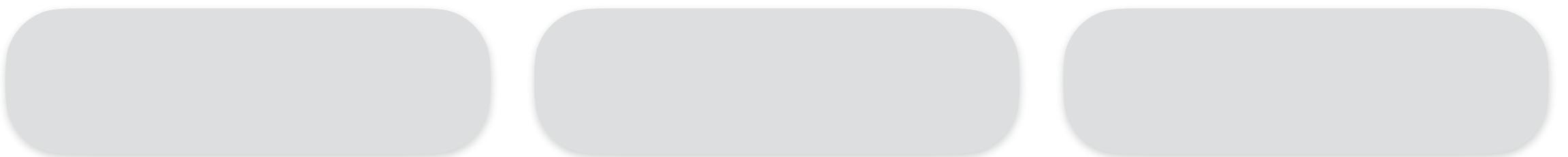
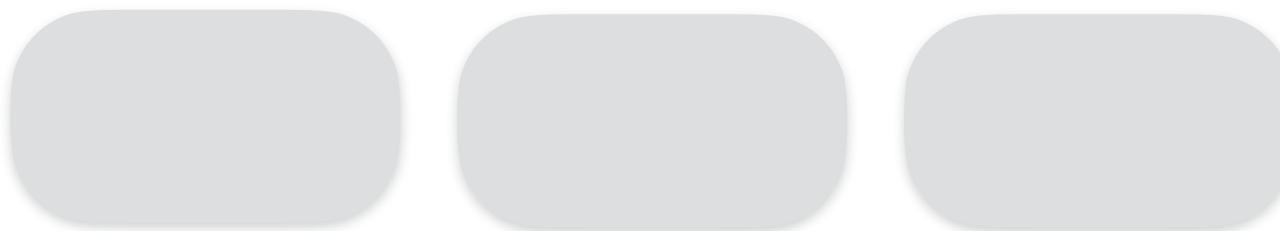
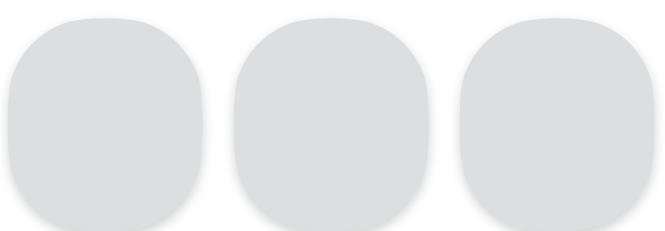
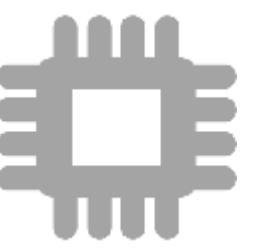
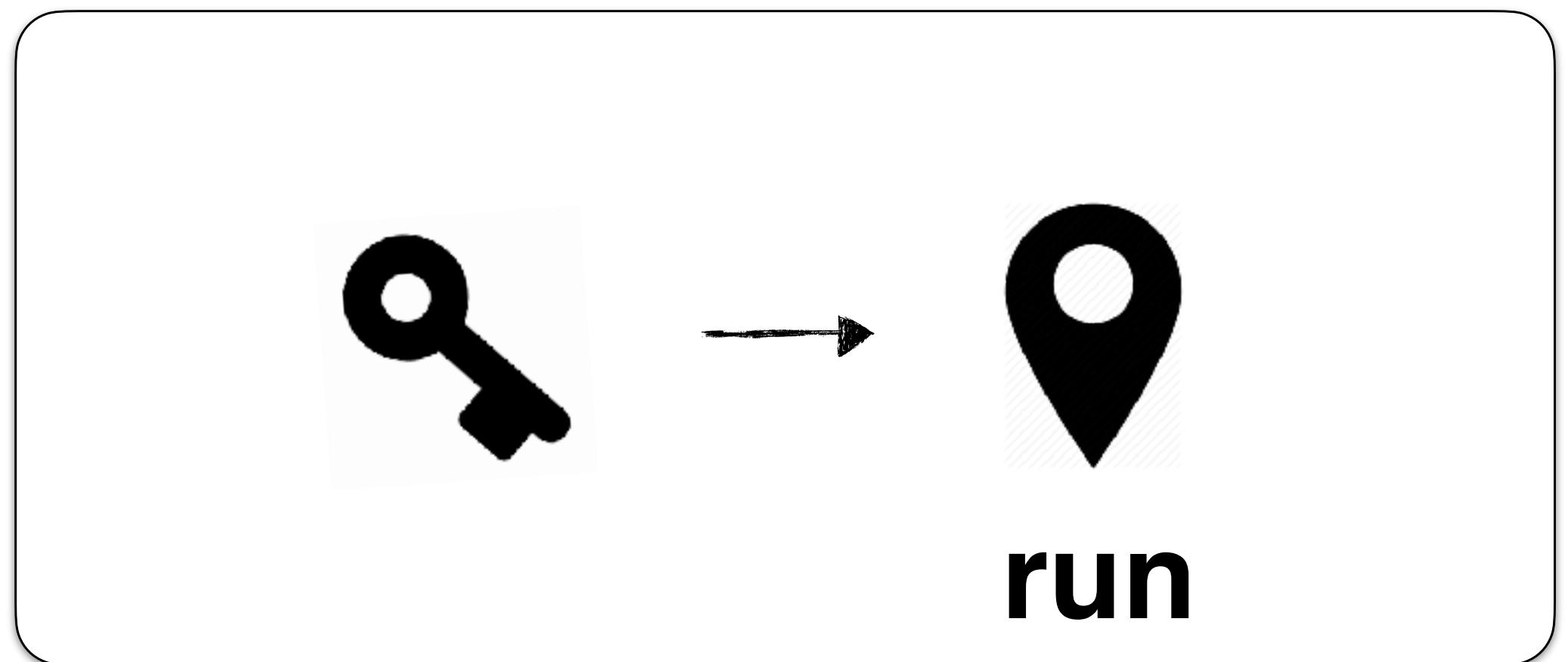


2

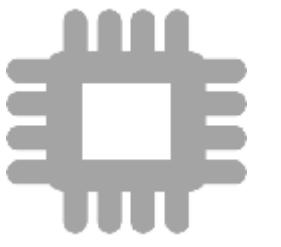


3



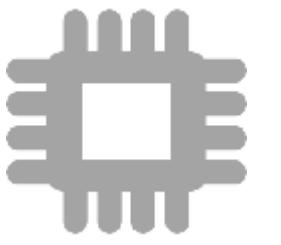


**lazy**



**O(1) reads**



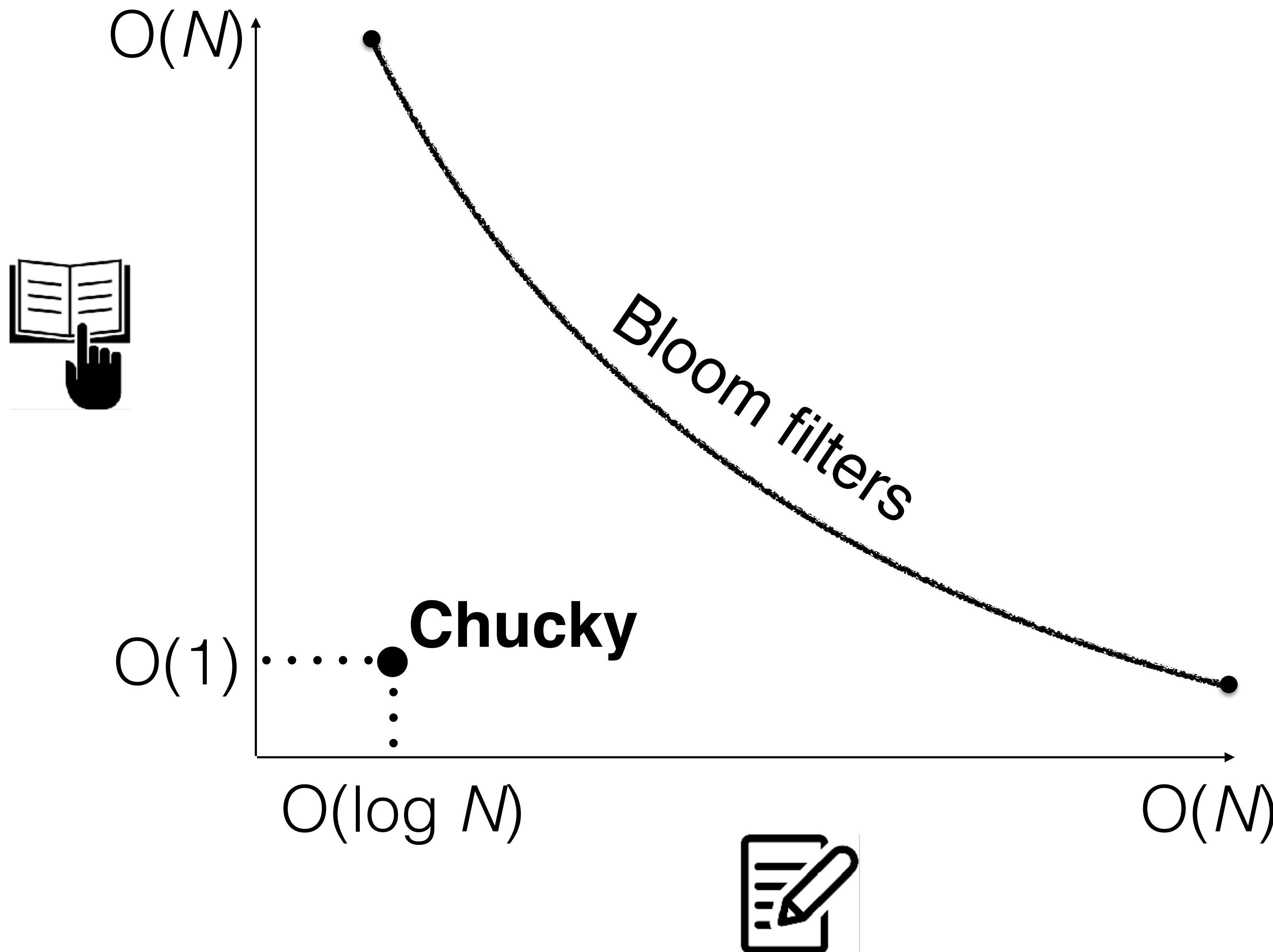


**O(1) reads**



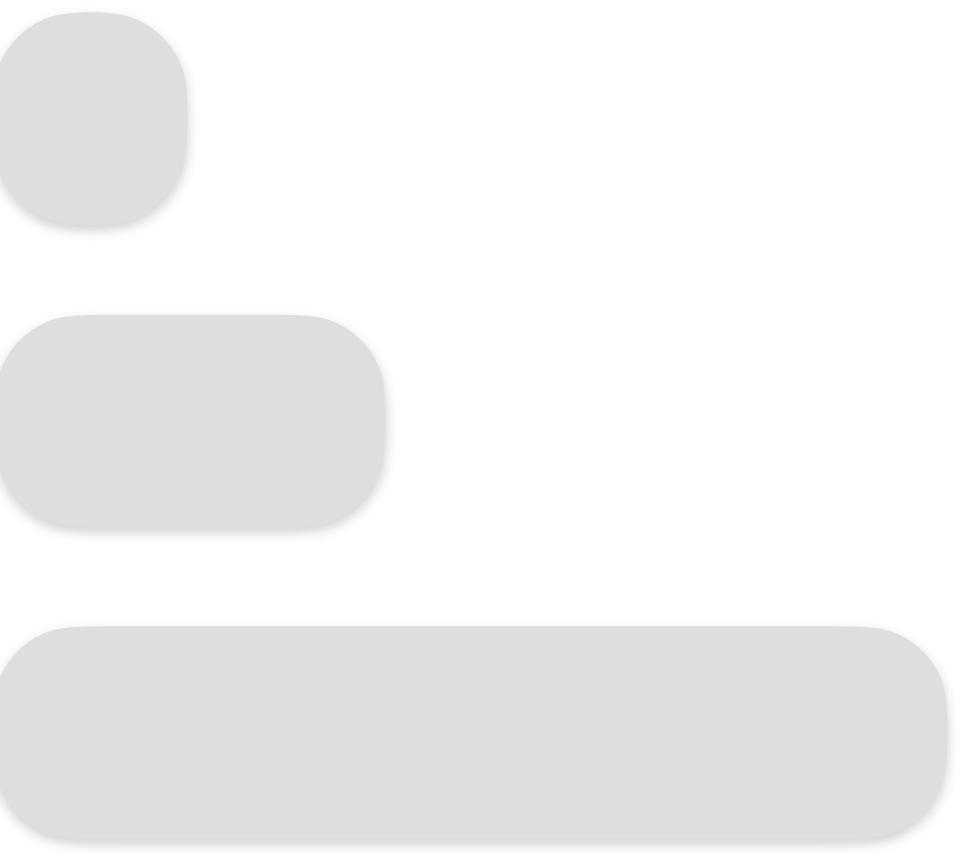
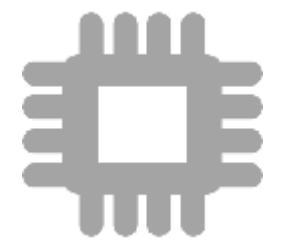
**O(log N) updates**

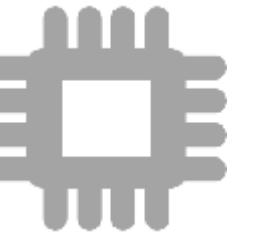






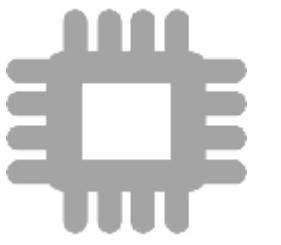
**memory >> 10 bits / entry**





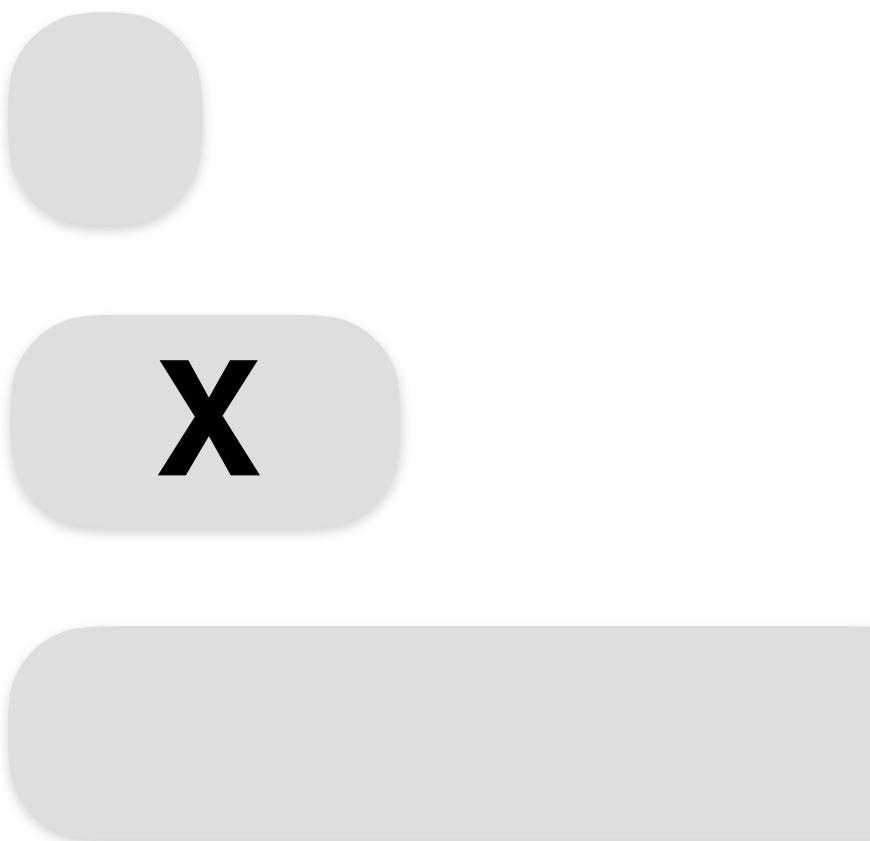
hash( ) =





level ID

hash(**X**) = 01011 → **2**

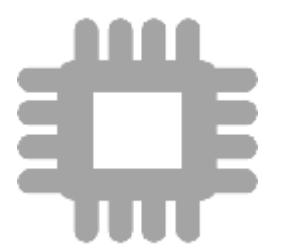


1

2

3

query Y



hash(X) = 01011 → 2

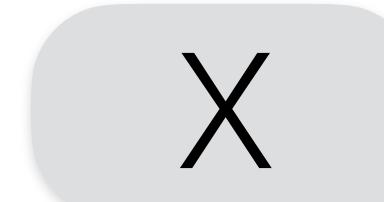
**false  
positive**



level ID



1

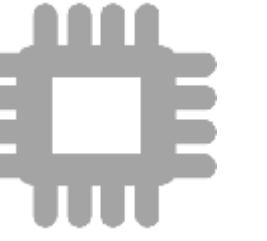


2



3

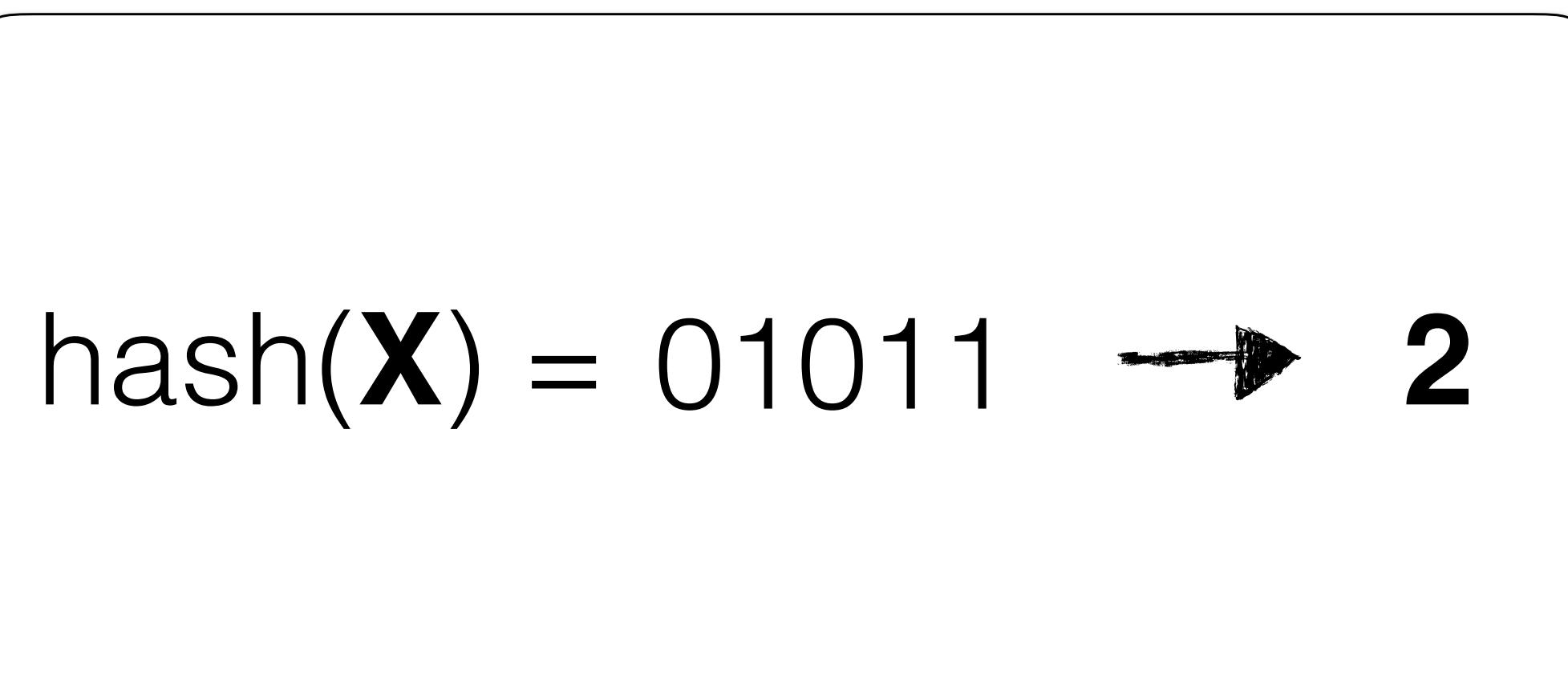
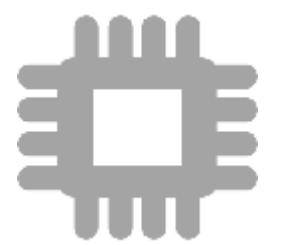
query **Y**



hash(**X**) = 01011 → 2

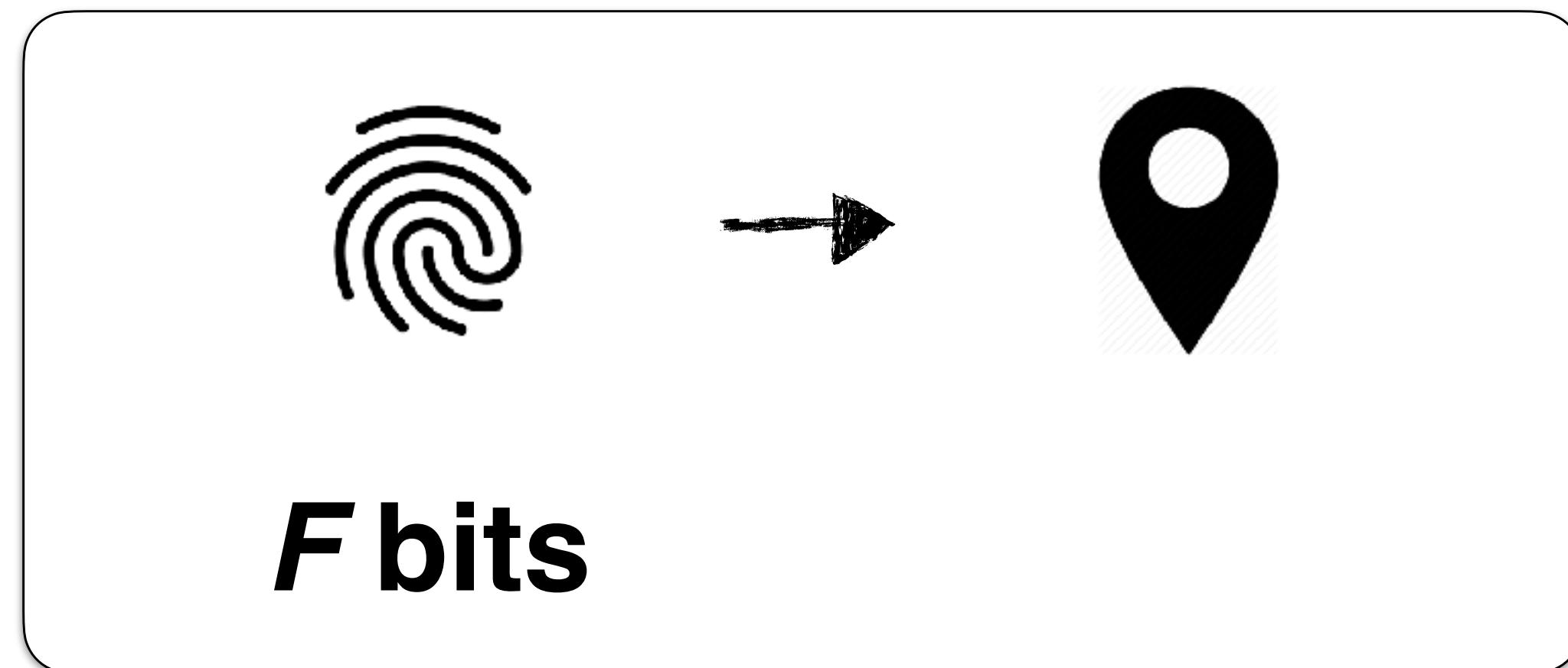
**false  
positive rate?**

query **Y**



false  
positive rate?

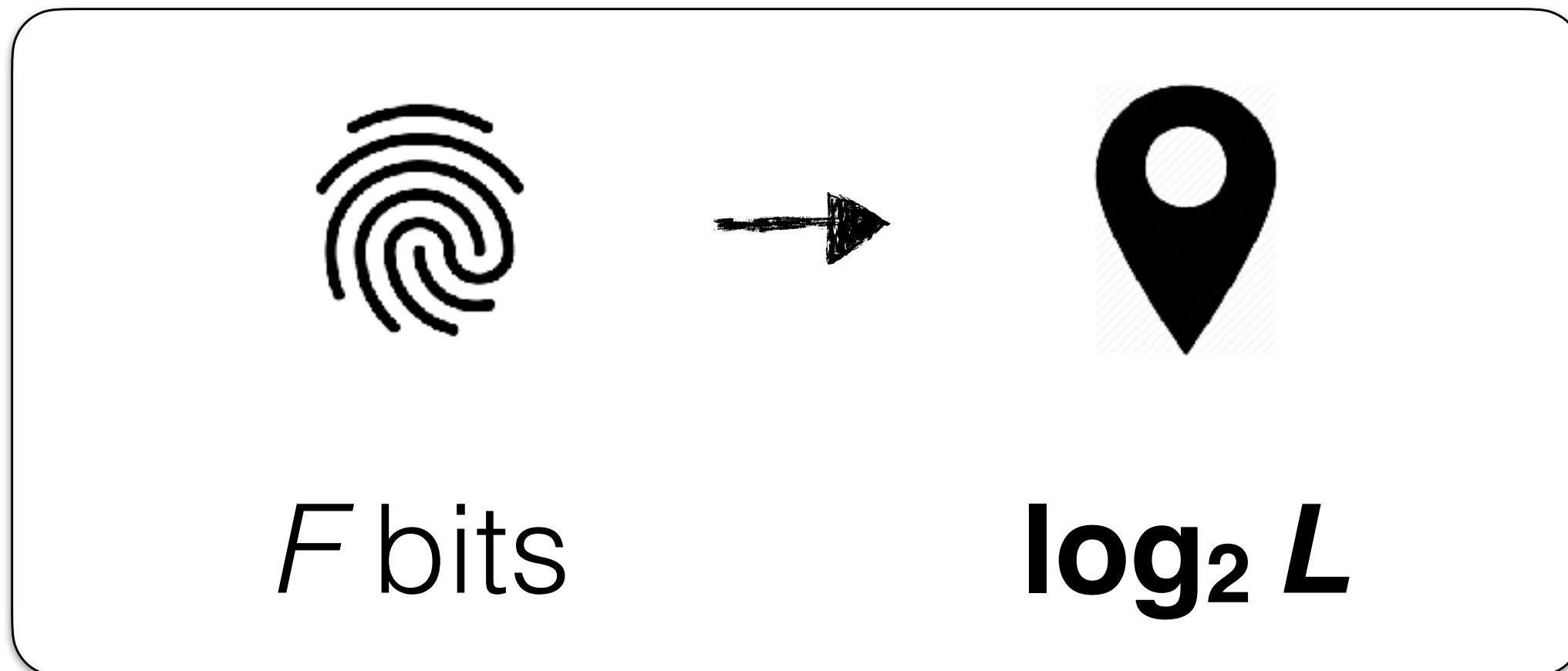
**(Assume minimal perfect hashing)**



false  
positive rate?

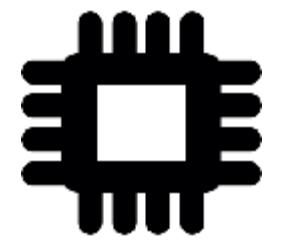
**2<sup>-F</sup>**

**$M$  bits**



false  
positive rate?

$2^{-F}$



$$F = M - \log_2 L$$

false  
positive rate?

$$2^{-F}$$

Bloom filters

$$\approx e^{-M} \cdot \ln(2)$$

<

Chucky (so far)

$$2^{-M} \cdot \log N$$

**level IDs grow,  
taking bits from fingerprints**

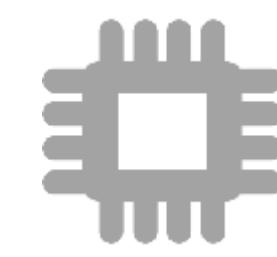
$$\approx e^{-M} \cdot \ln(2)$$

<

$$2^{-M} \cdot \log N$$



**encode larger  
levels with  
fewer bits**



1  
2  
2  
3  
3  
3  
3



**most entries**

## **Level IDs**

1

2

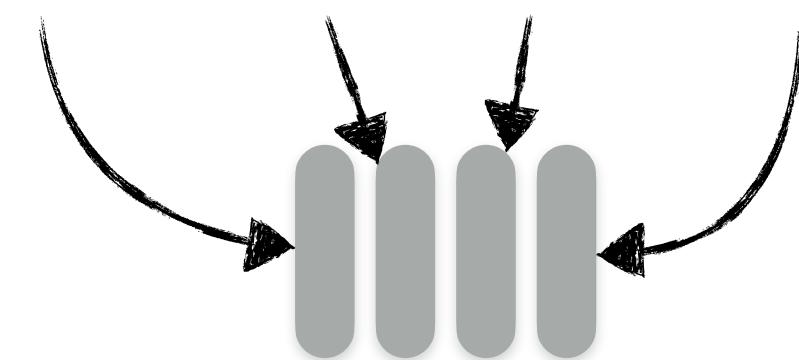
3

# **unary encoding**

level ID

	<b>00001</b>	1
	<b>0001</b>	2
	<b>001</b>	3
	<b>01</b>	4
	<b>1</b>	5

**00    01    10    11**



unary encoding

00001 \_

1

0001 \_

2

001 \_

3

01 \_

4

1 \_

5

level ID

unary encoding

level ID



00001

0001

001

01

1

1

2

3

4

5

**average level ID length < 2 bits <  $\log N$**

unary encoding

level ID



00001

1

0001

2

001

3

01

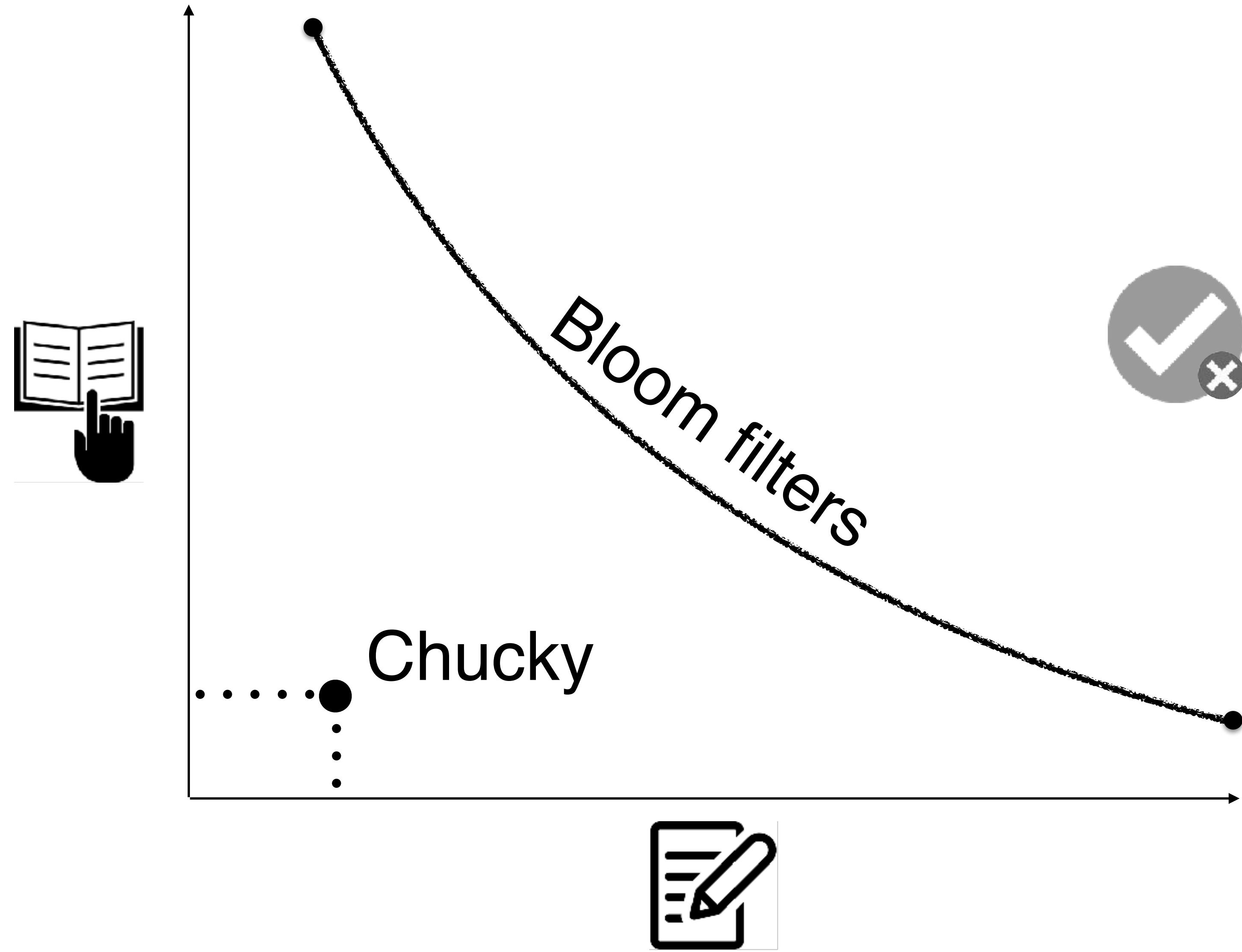
4

1

5



$$F = M - 2$$



Bloom filters

$O(e^{-M} \cdot \ln(2)) \approx O(2^{-M})$

Chucky

**Level**

**code**

**some bucket**

1

001

2

01

3

1

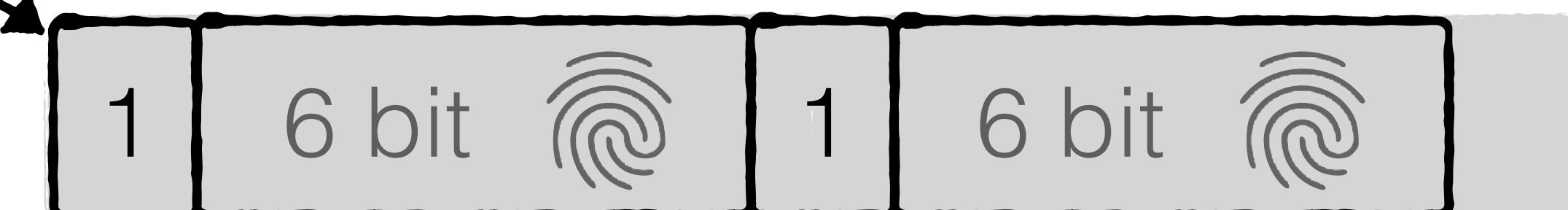
16 bit, 2 entries

- 1
- 2
- 3

001

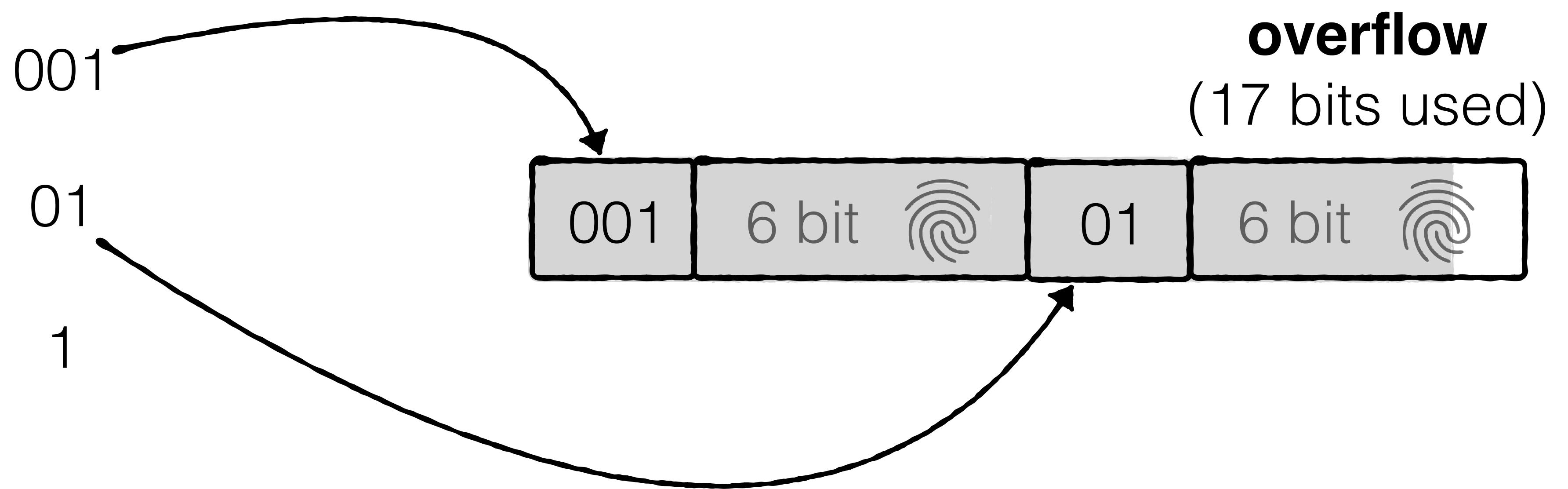
01

1

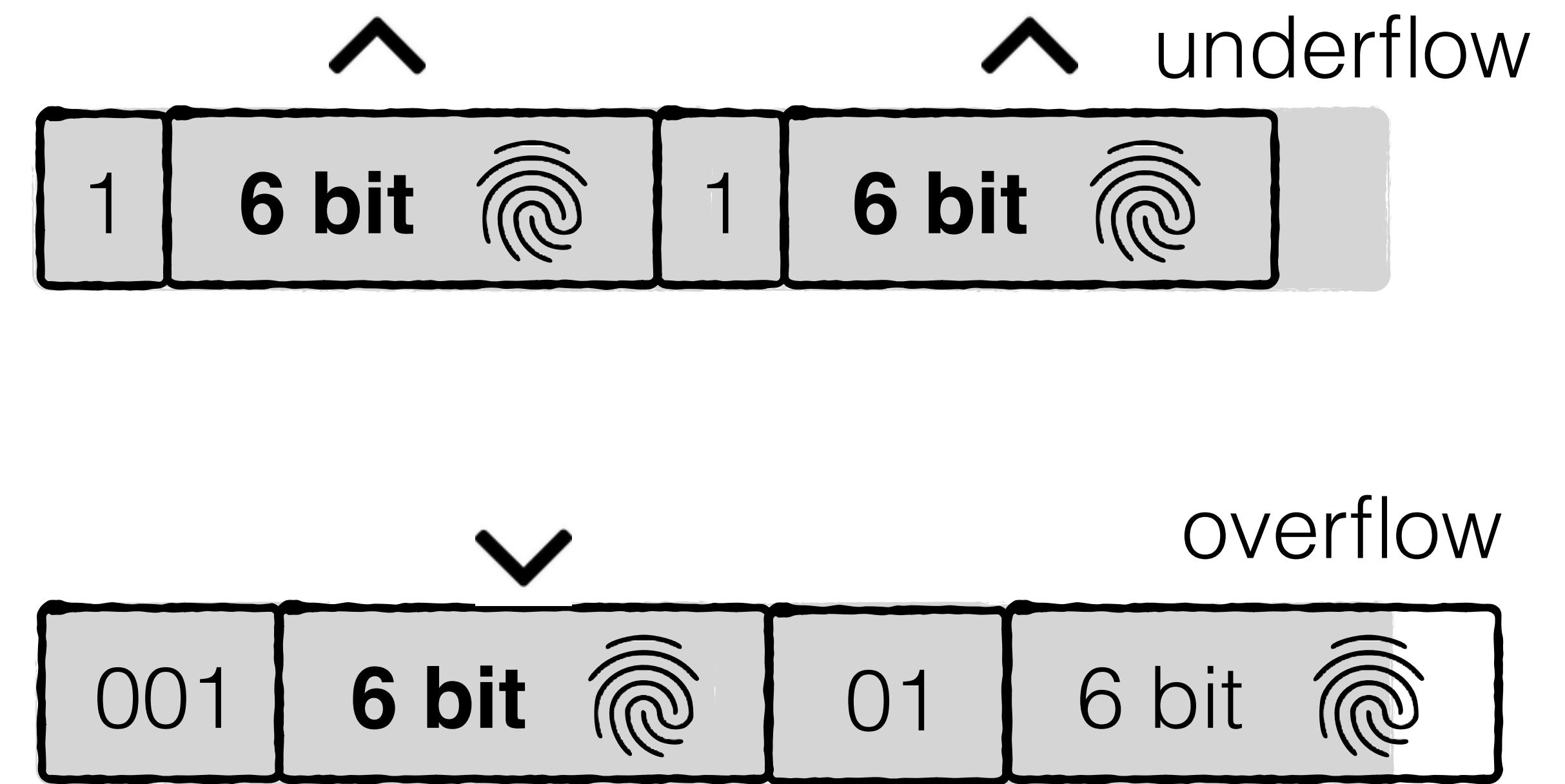


**underflow**  
(14 bits used)

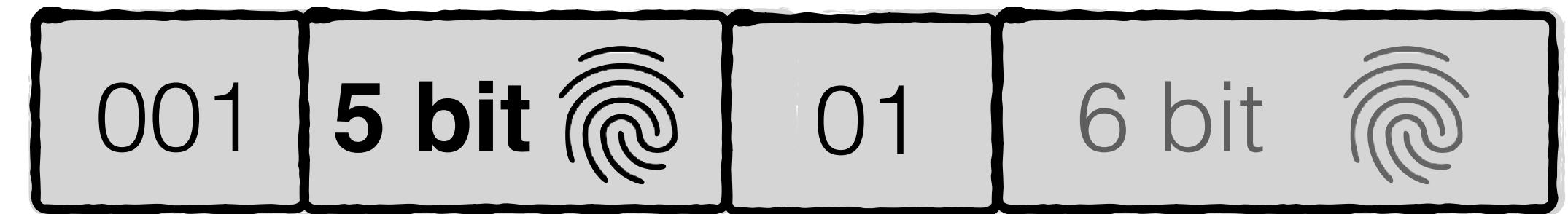
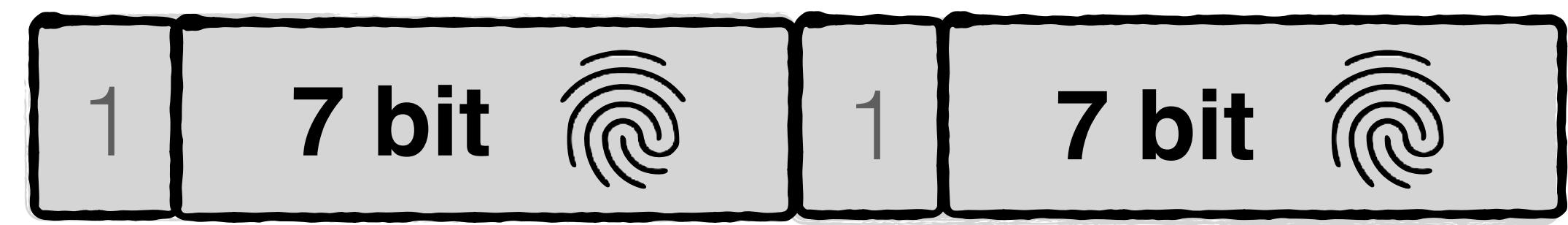
- 1
- 2
- 3



**larger fingerprints  
for larger levels**

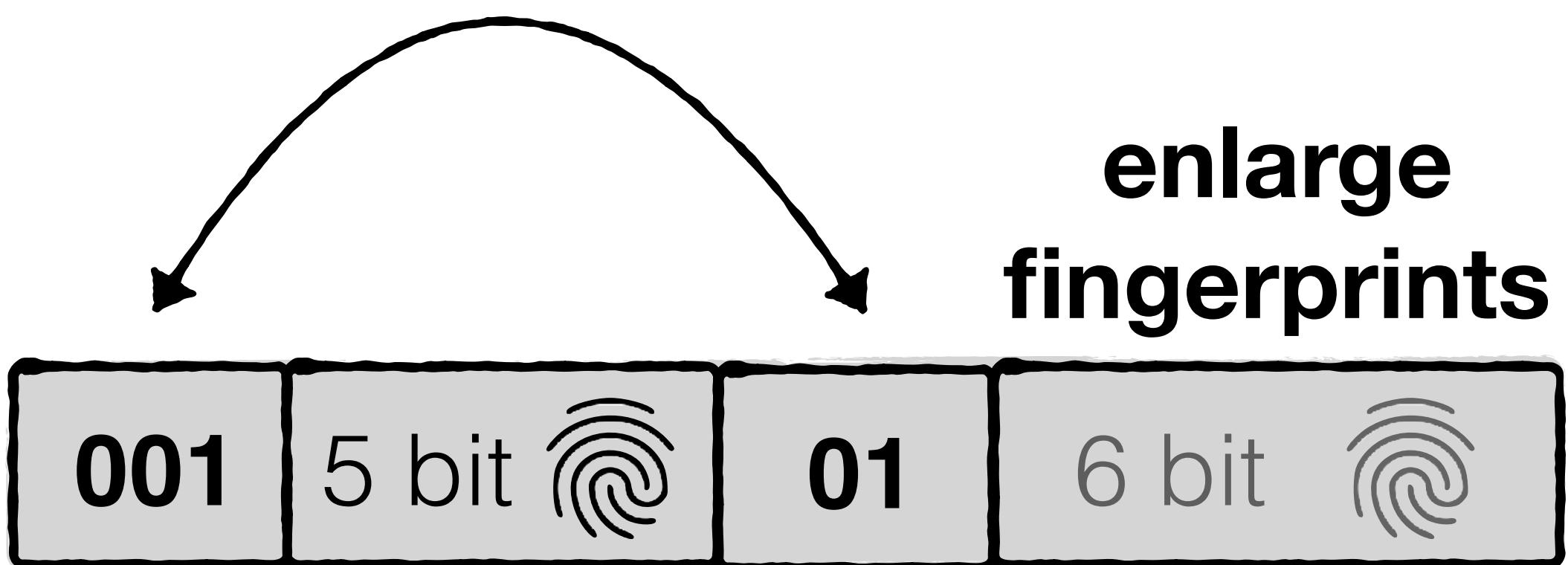


**larger fingerprints  
for larger levels**



**co-encode  
concisely**

**enlarge  
fingerprints**





# Encoding Level IDs as combinations

levels 3 & 3



levels 3 & 2



levels 3 & 1



levels 2 & 2



levels 2 & 1



levels 1 & 1



**task**

**technique**

**Step 1**

**Step 2**


Step 1

task

technique

**max. avg.  size**

**unique decodability**



Step 2

	task	technique
Step 1	max. avg.  size	
	unique decodability 	
Step 2	<b>generate codes</b>	 ... 0101 1010

	task	technique
Step 1	max. avg.  size	<b>hill-climbing</b> 
	unique decodability 	<b>Kraft–McMillan</b> 
Step 2	generate codes 	<b>Huffman Tree</b> 

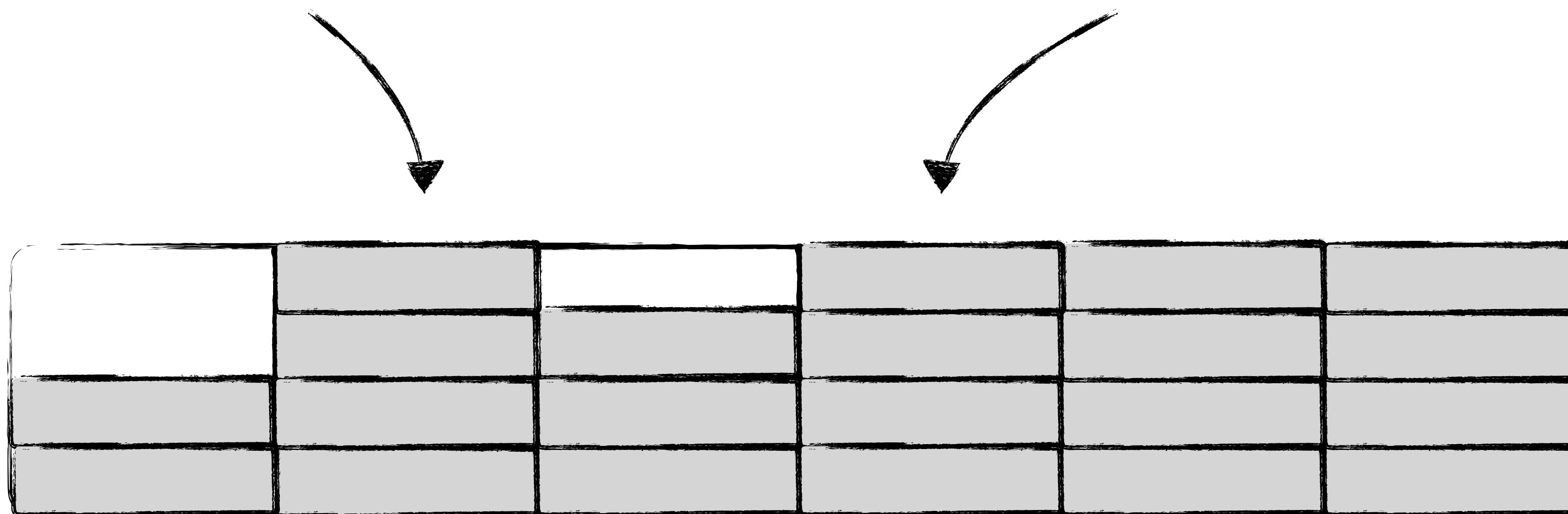
**underflows**

**Overflows**

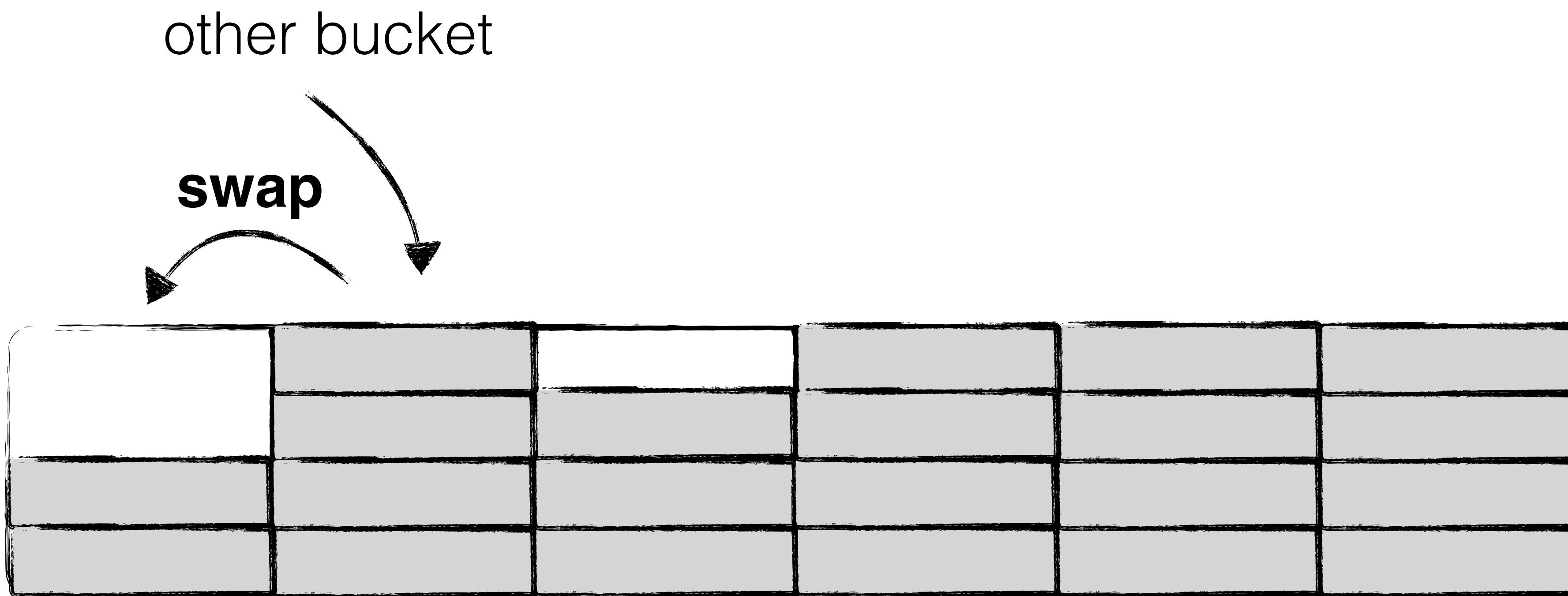


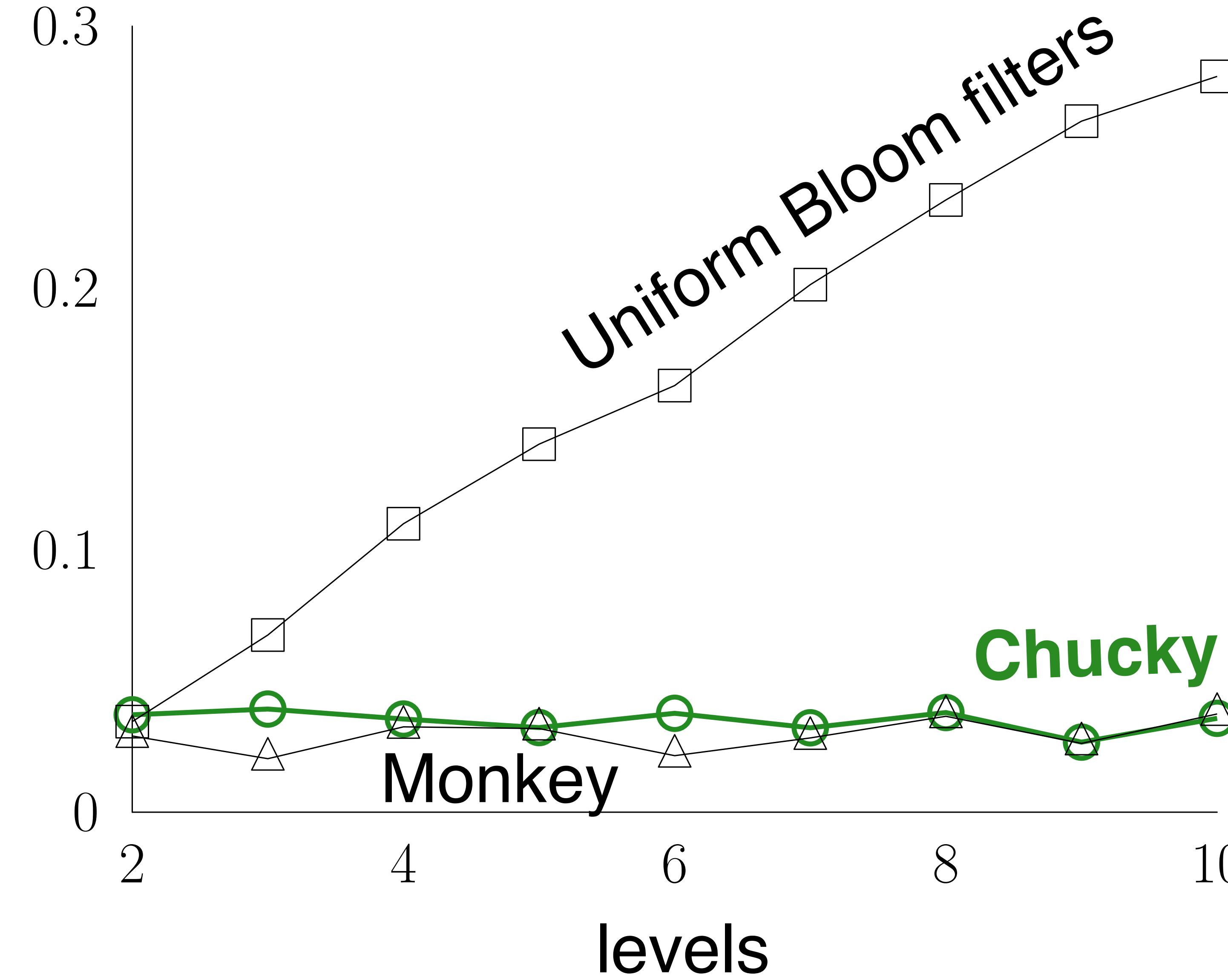
# Cuckoo Filter

**other bucket = hash(  )  $\oplus$  current bucket**



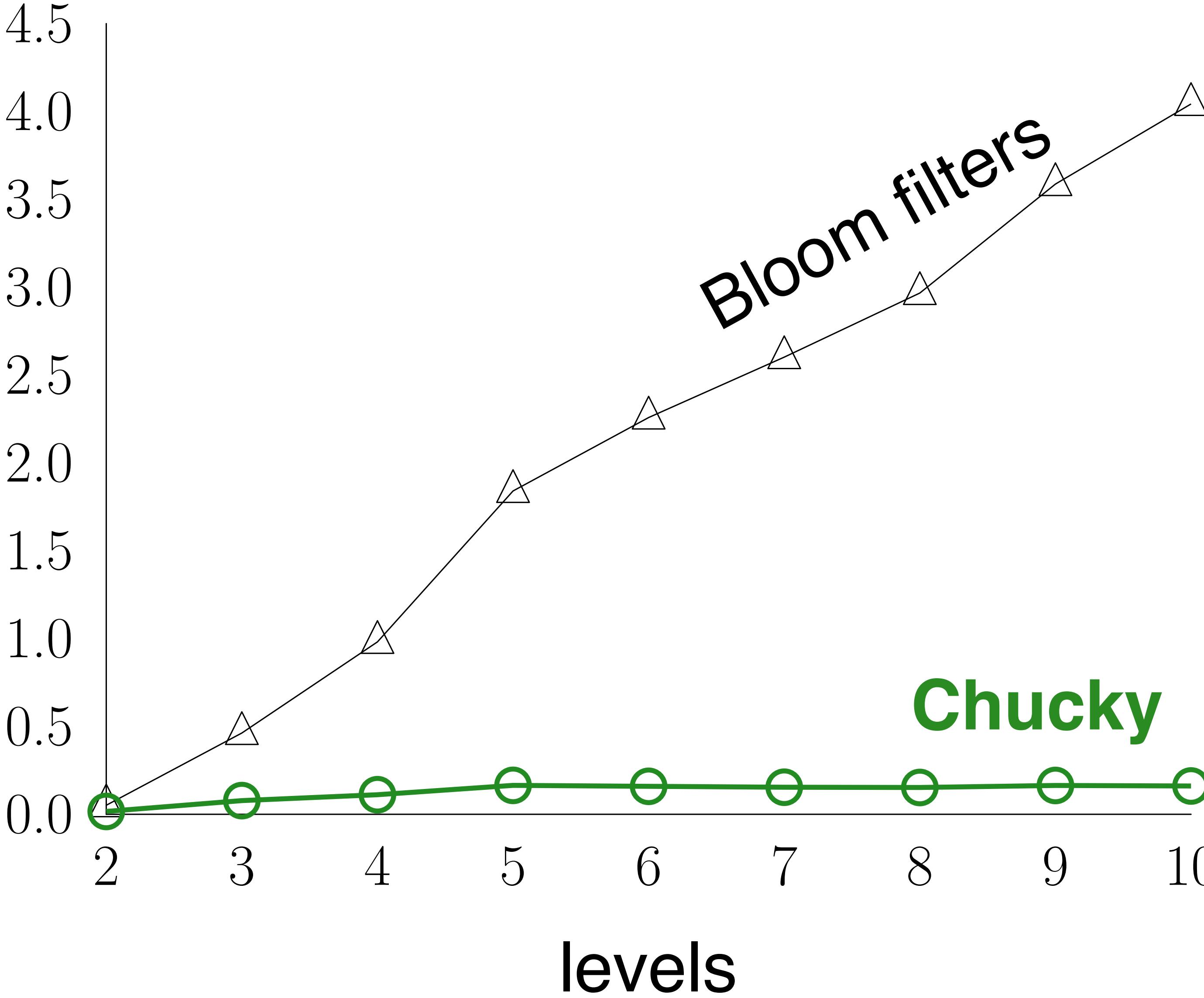
# Cuckoo Filter



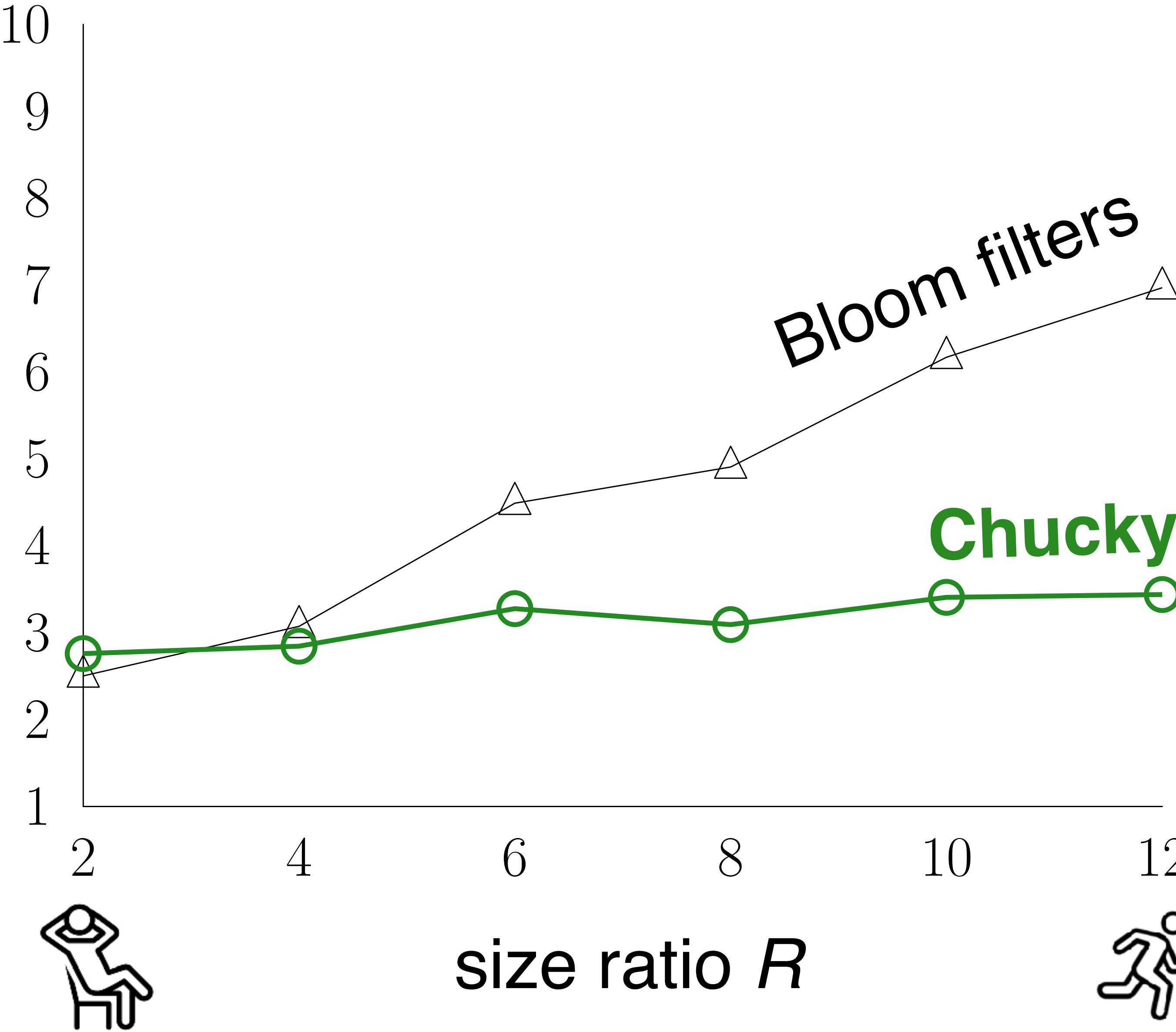
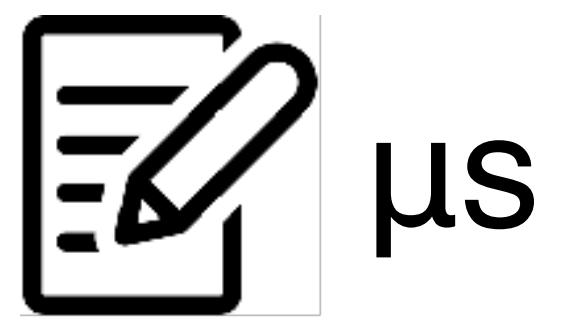




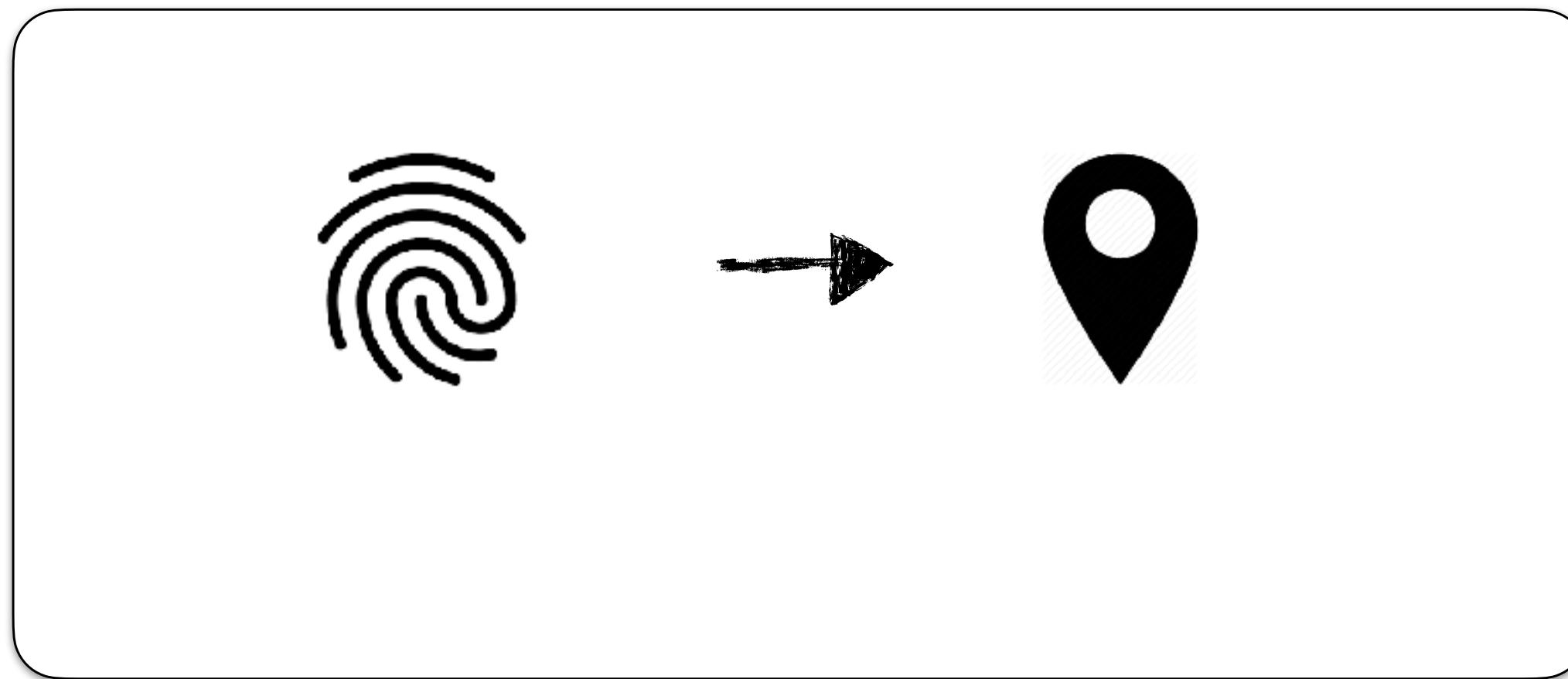
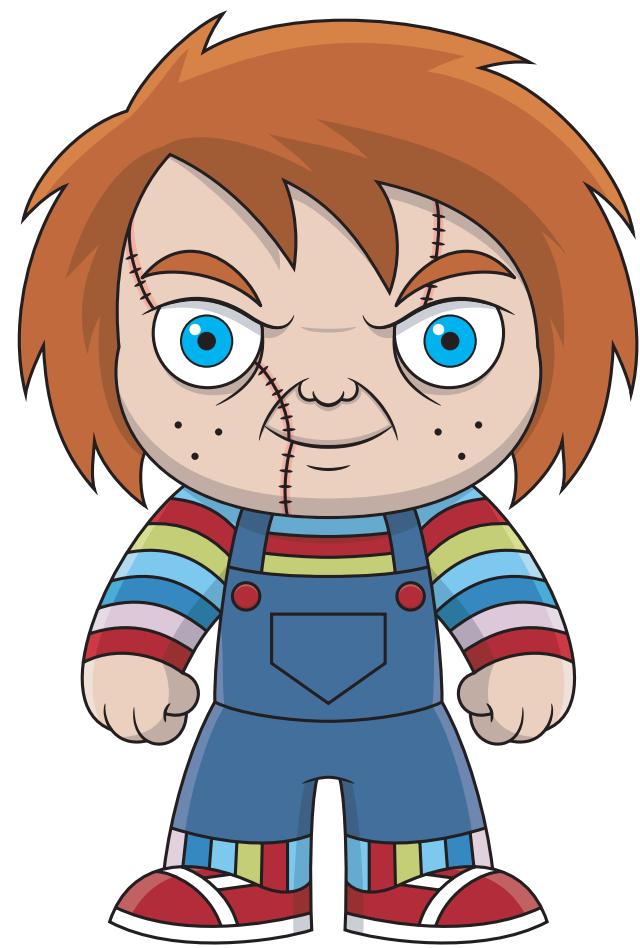
μs



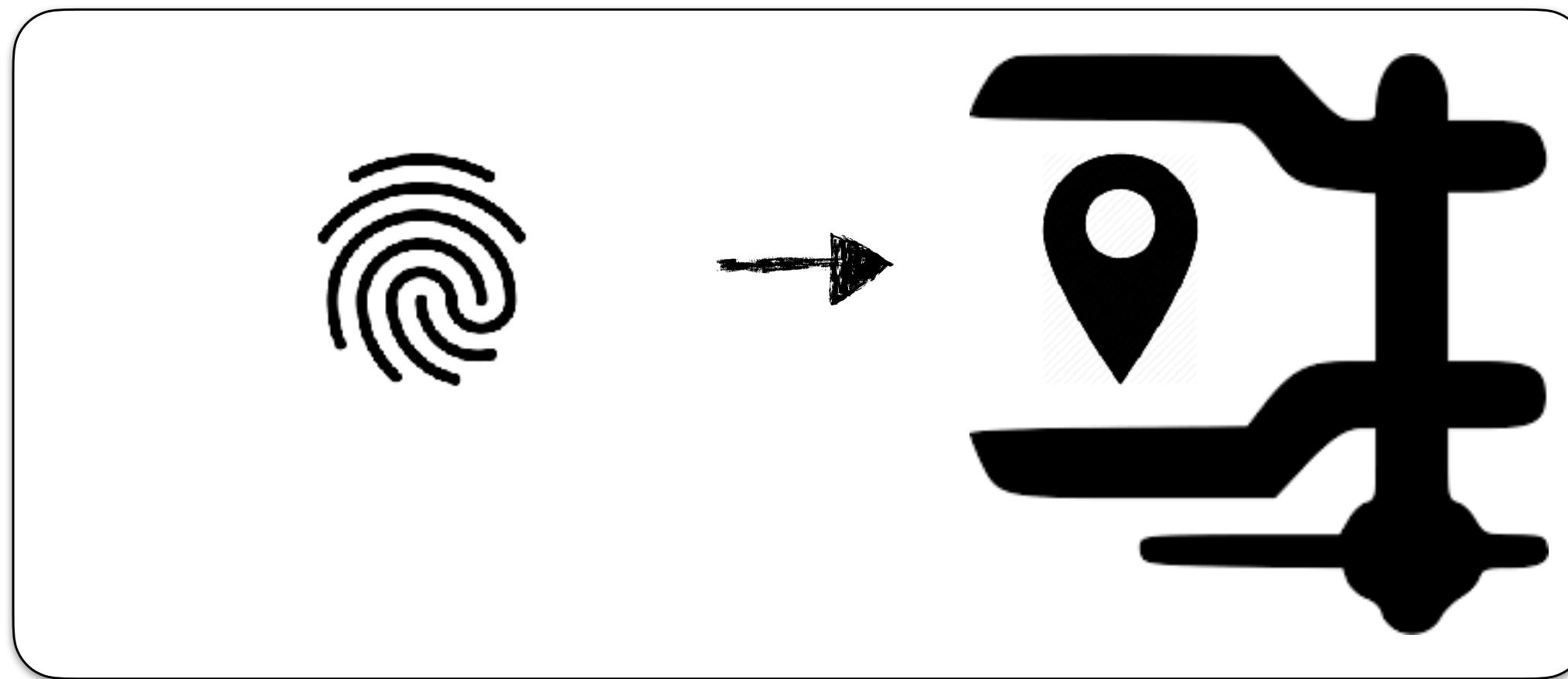
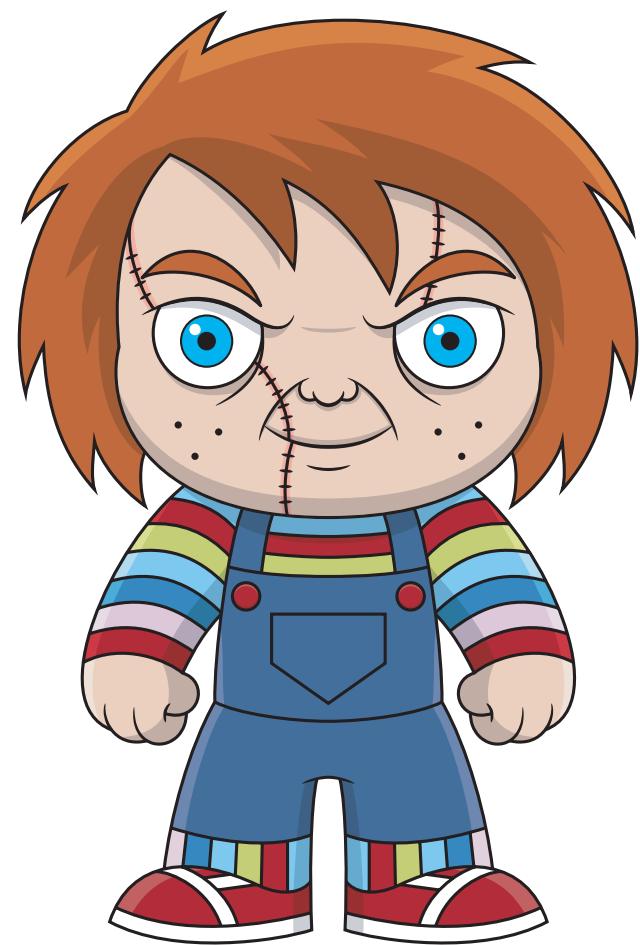
Chucky



# Chucky

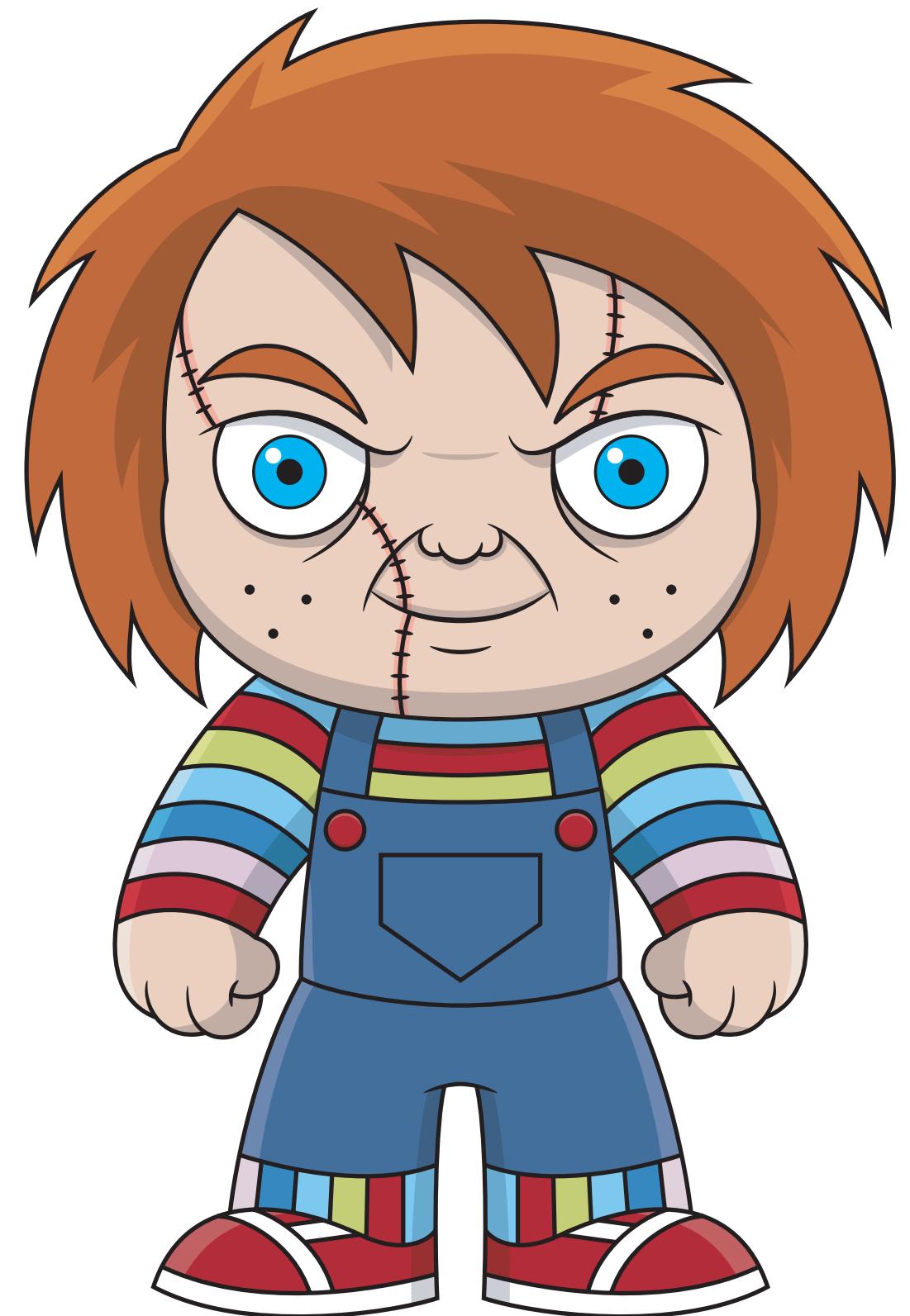


# Chucky



# Chucky





Thank you!