Research Lecture: LSM-trees & Filters



Niv Dayan

Projects



56 groups - mostly of threes



157 / 169 students registered



Feedback week next week - will announce when.

Midterm



Oct 14 in class



Open book

Today



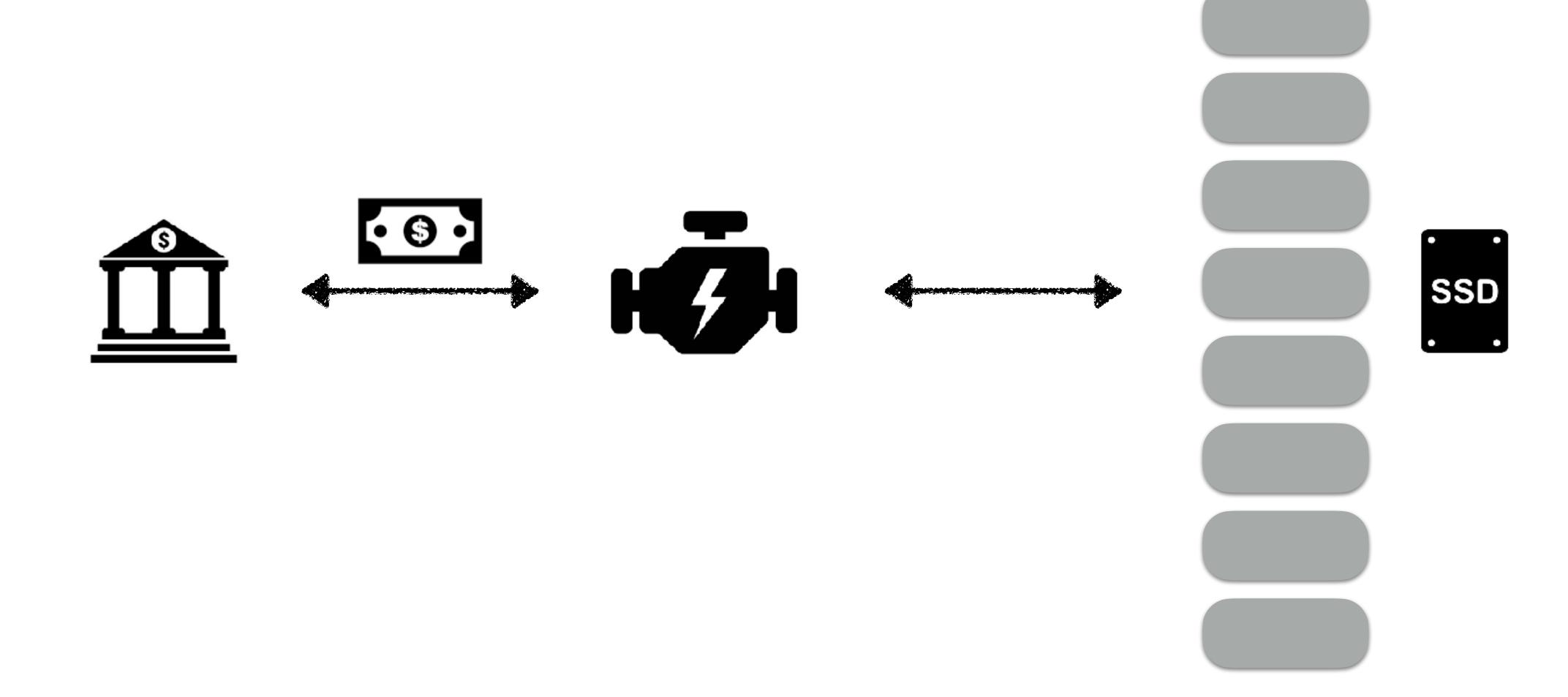
Recap on LSM



Bloom Filters



Research Lecture Many DB operations are:

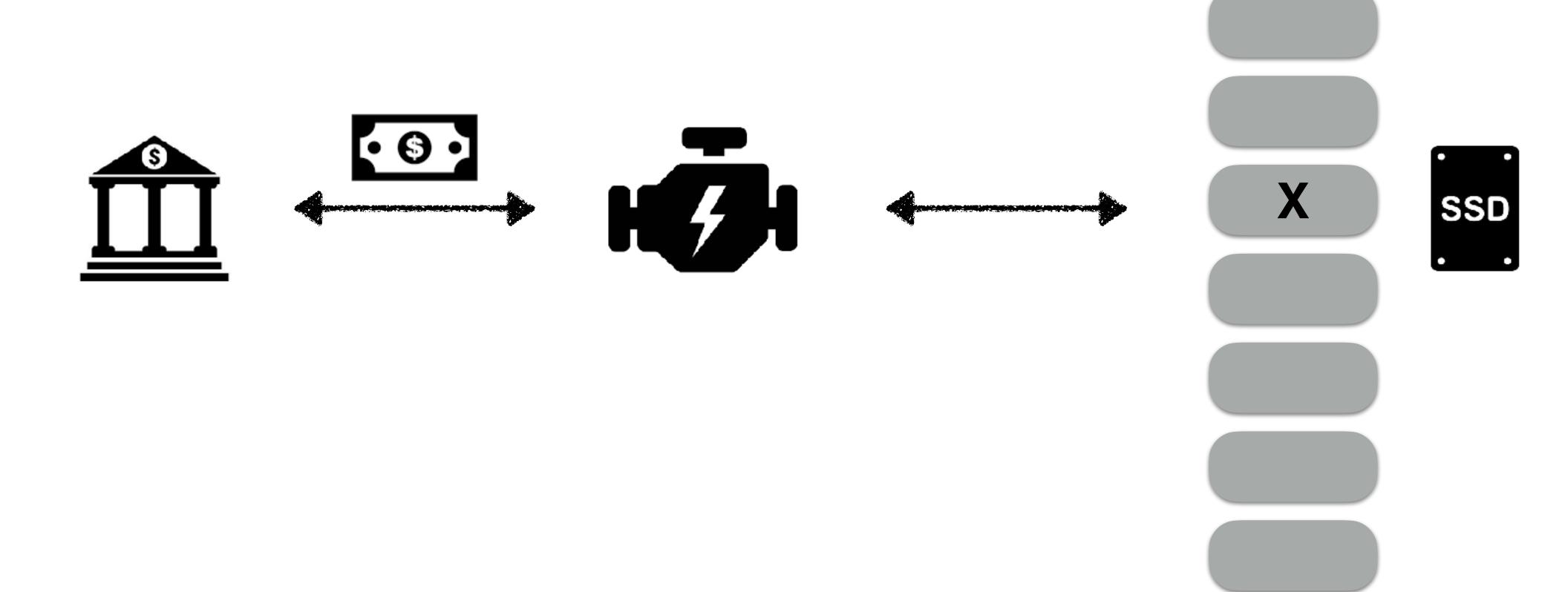


Many DB operations are: selective

SSD

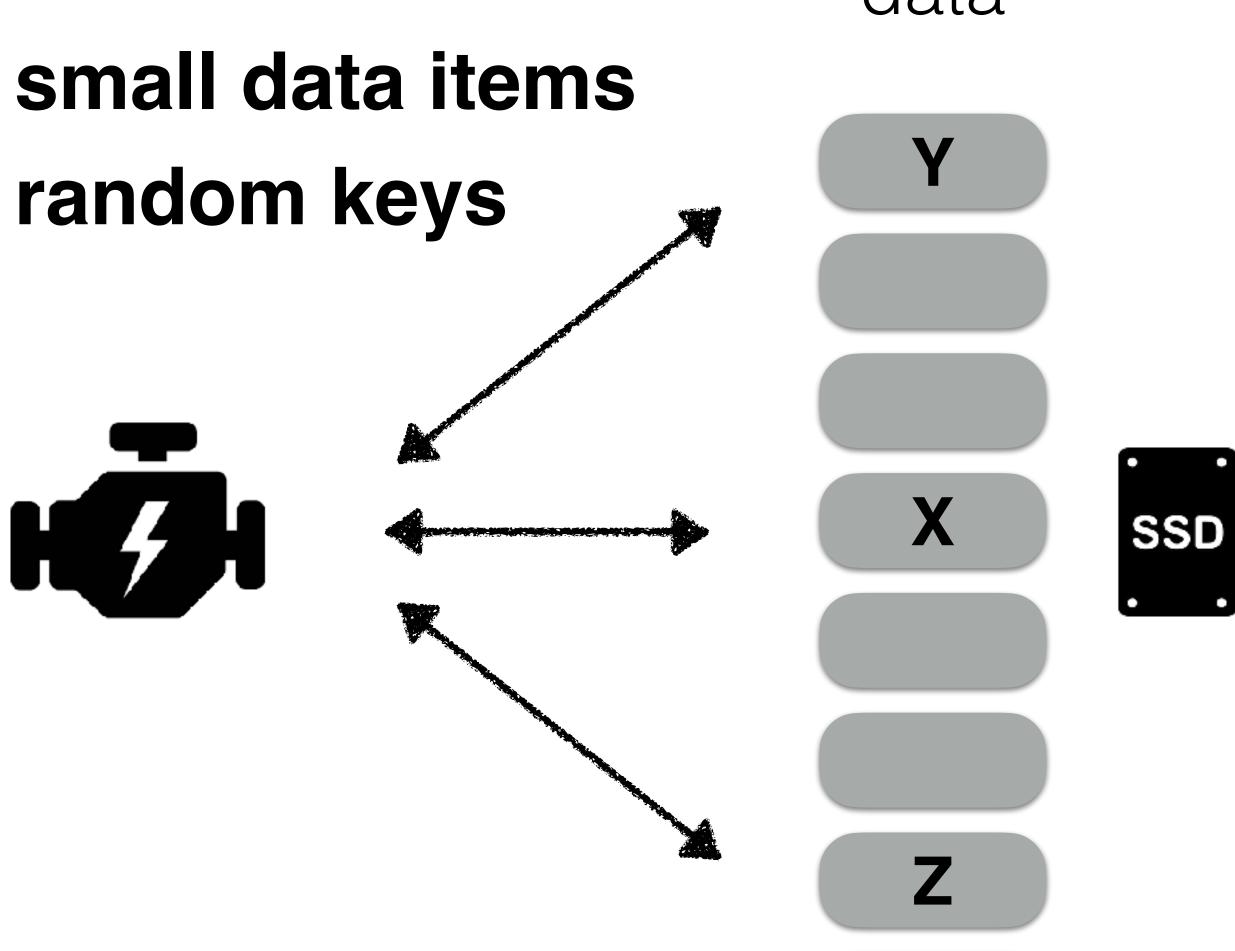
Many DB operations are: selective

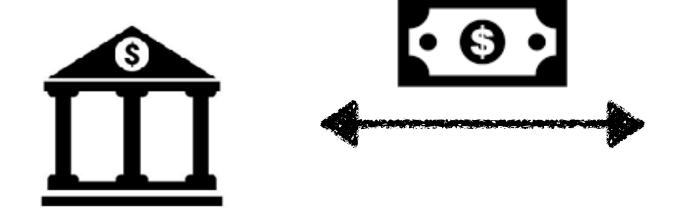
selective small data items

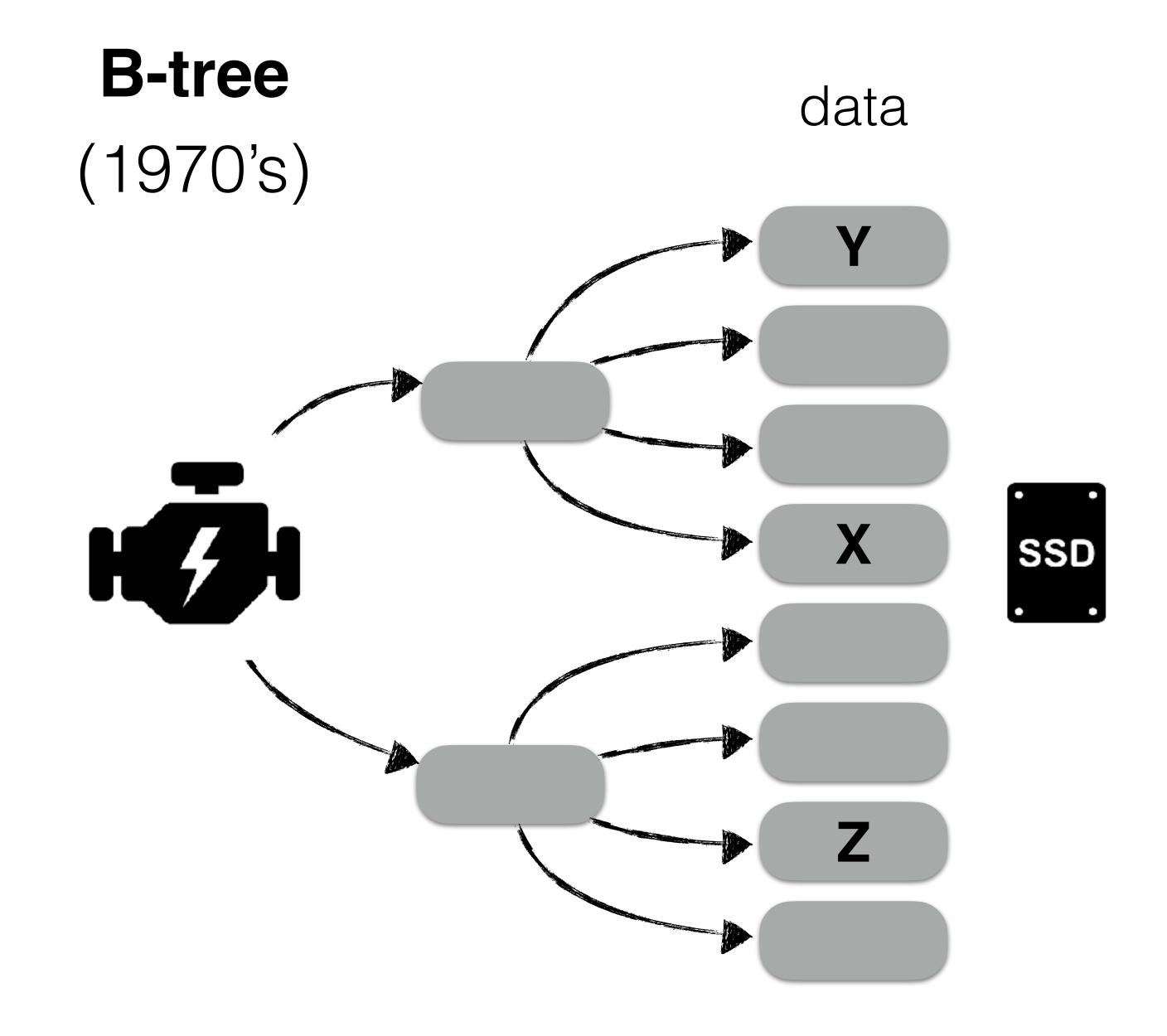


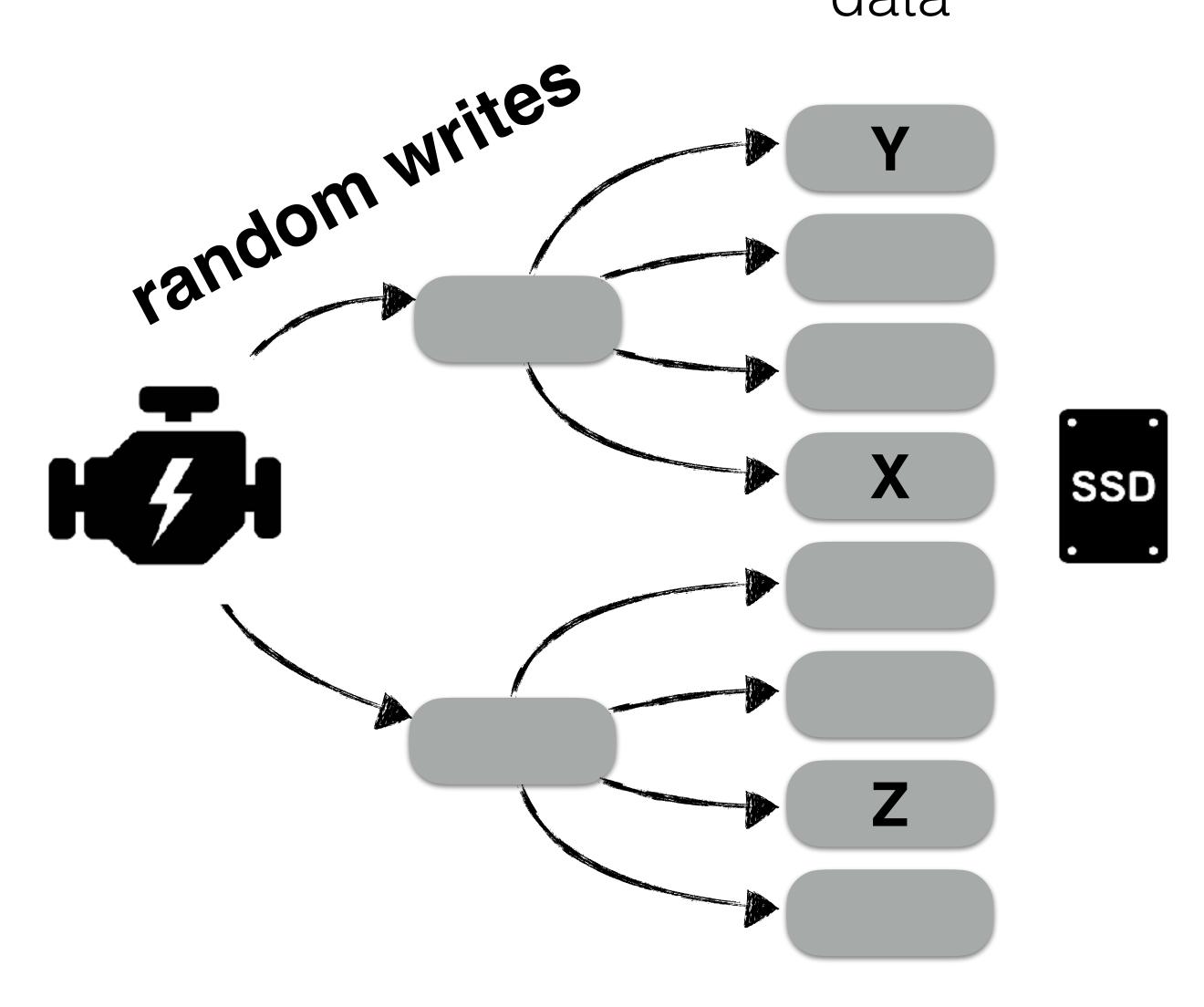
Many DB operations are:

selective data















SSD

mechanical latency



mechanical latency



4KB access



mechanical latency



4KB access garbage-collection

The Log-Structured Merge-Tree

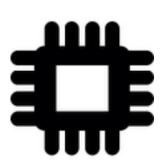
1996 - Patrick O'Neil



LSM-Tree







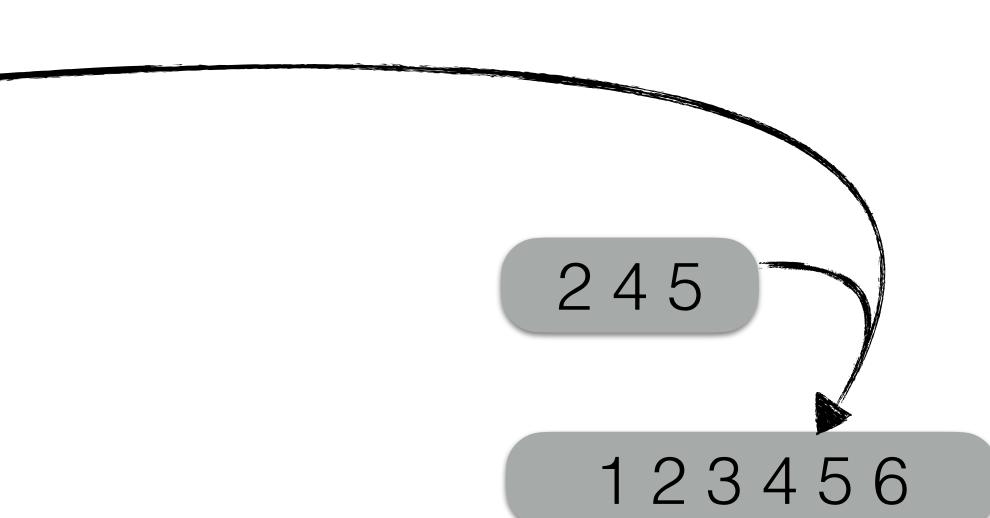




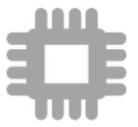
buffer

1 3 6

merge-sort

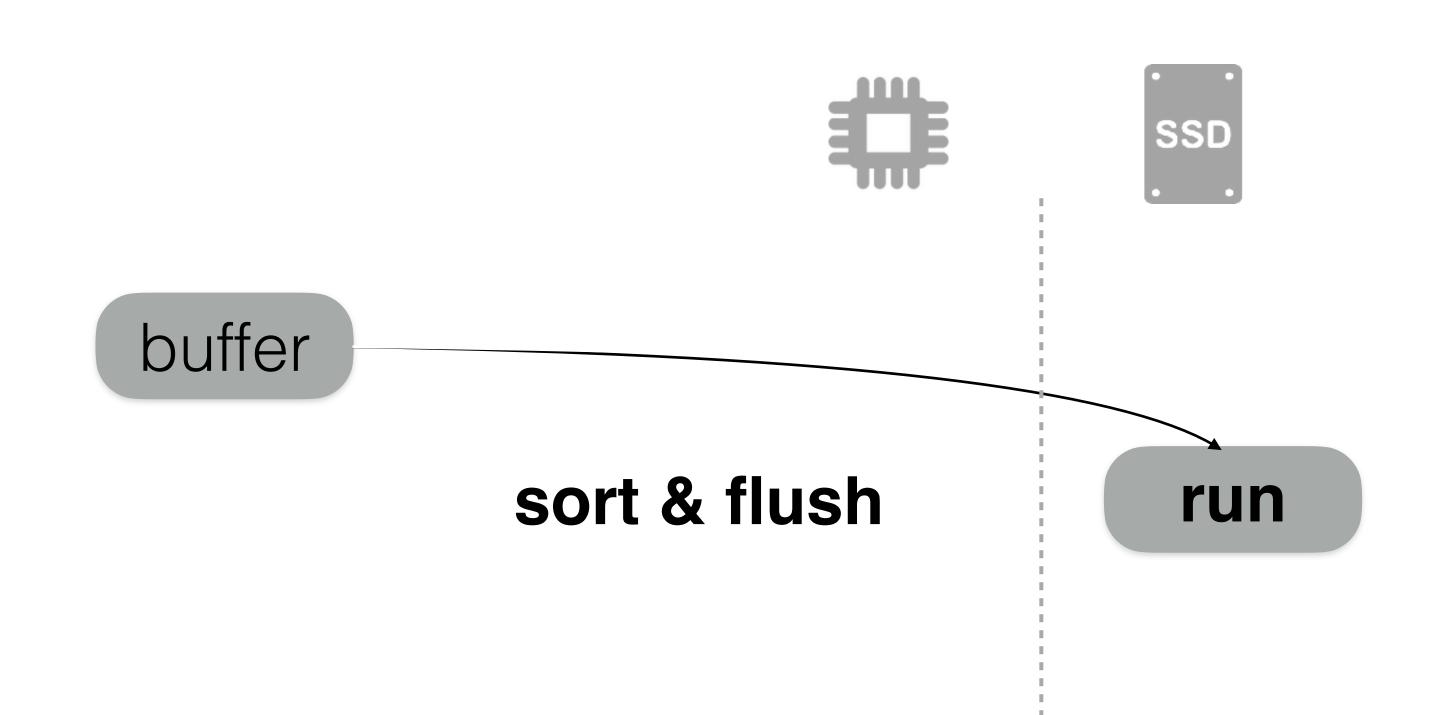


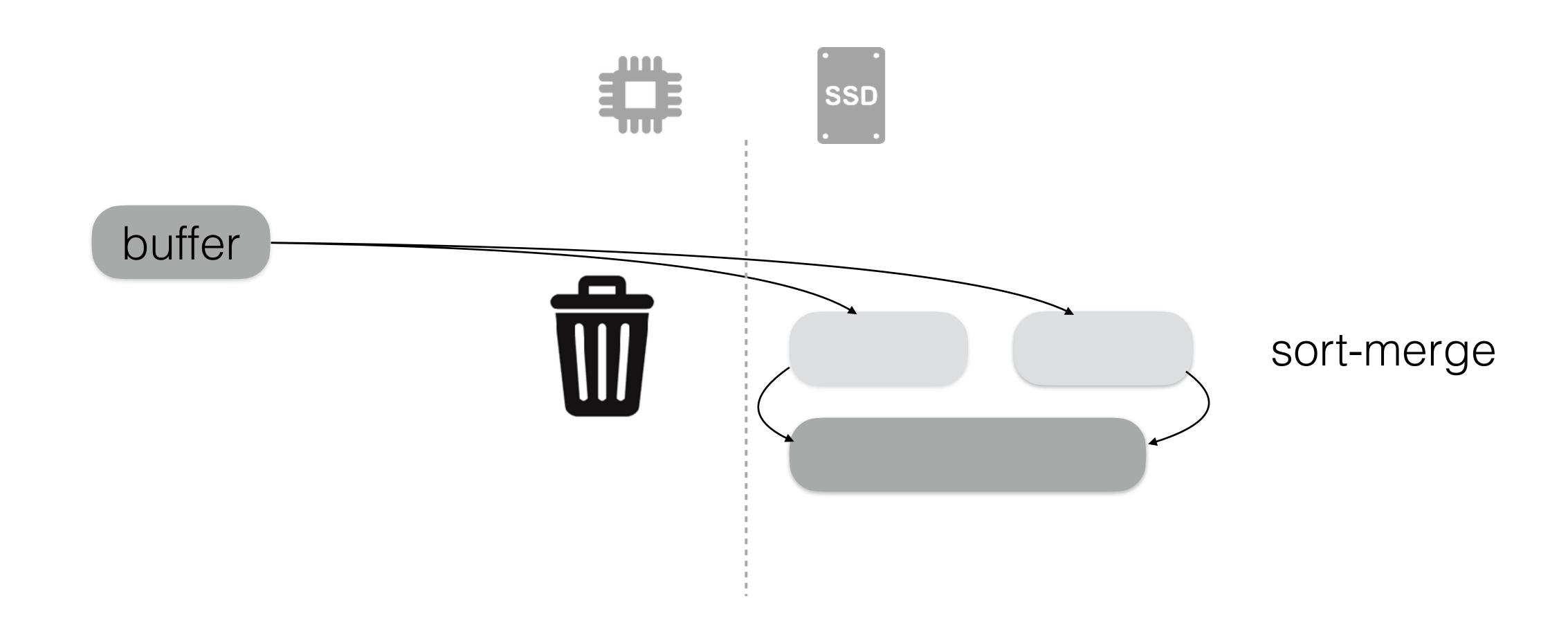
Inserts/updates/deletes of key-value pairs





buffer







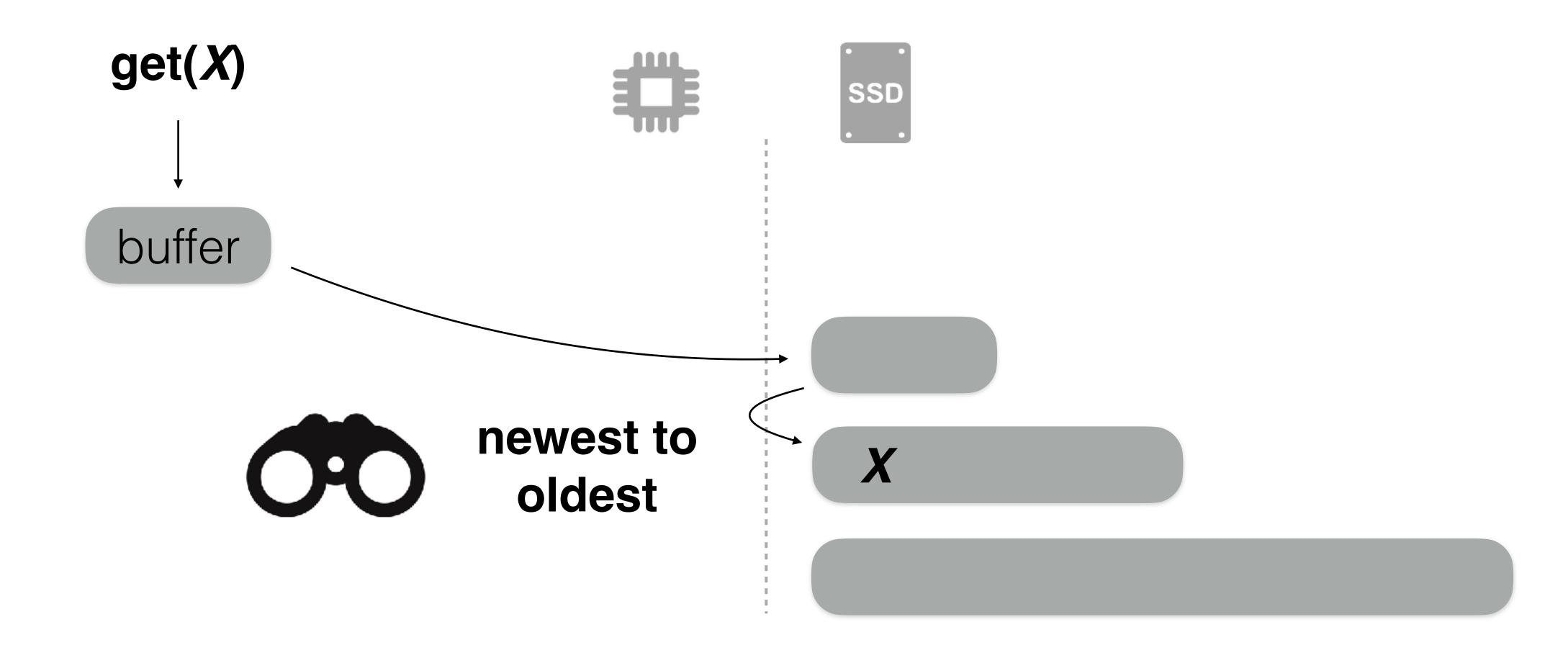
buffer

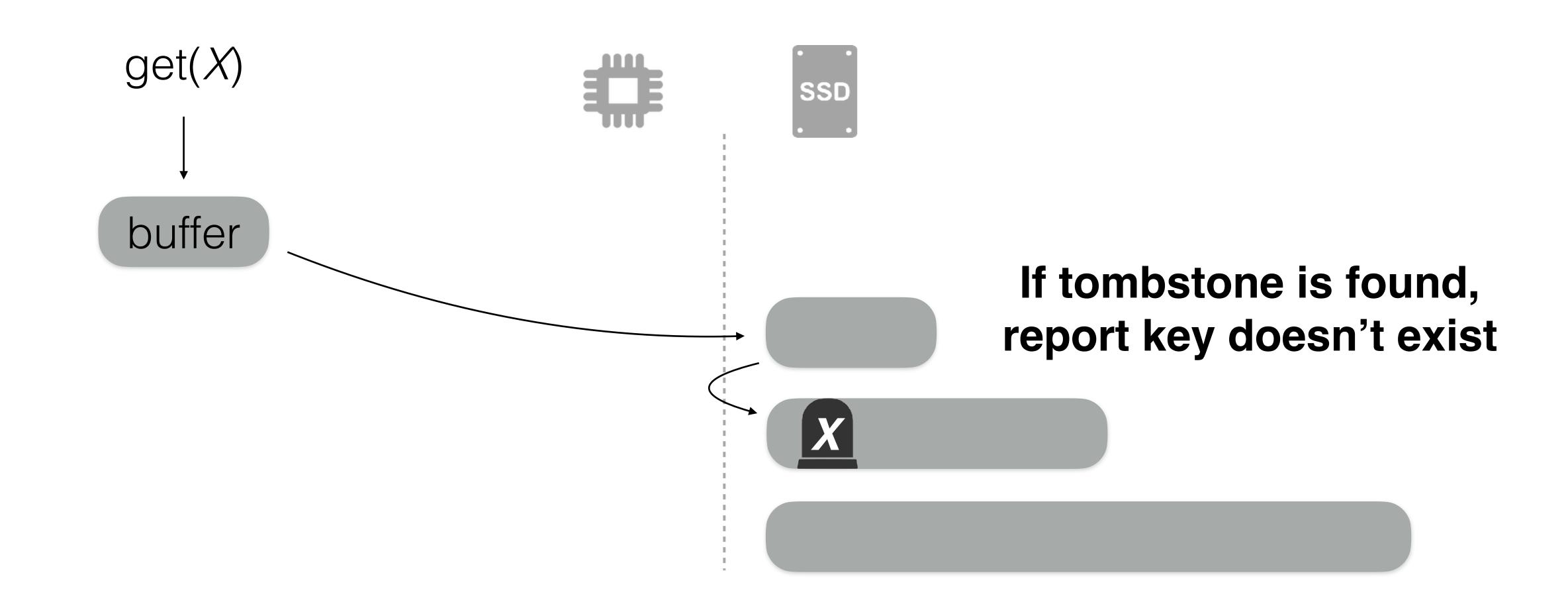
exponentially increasing capacities

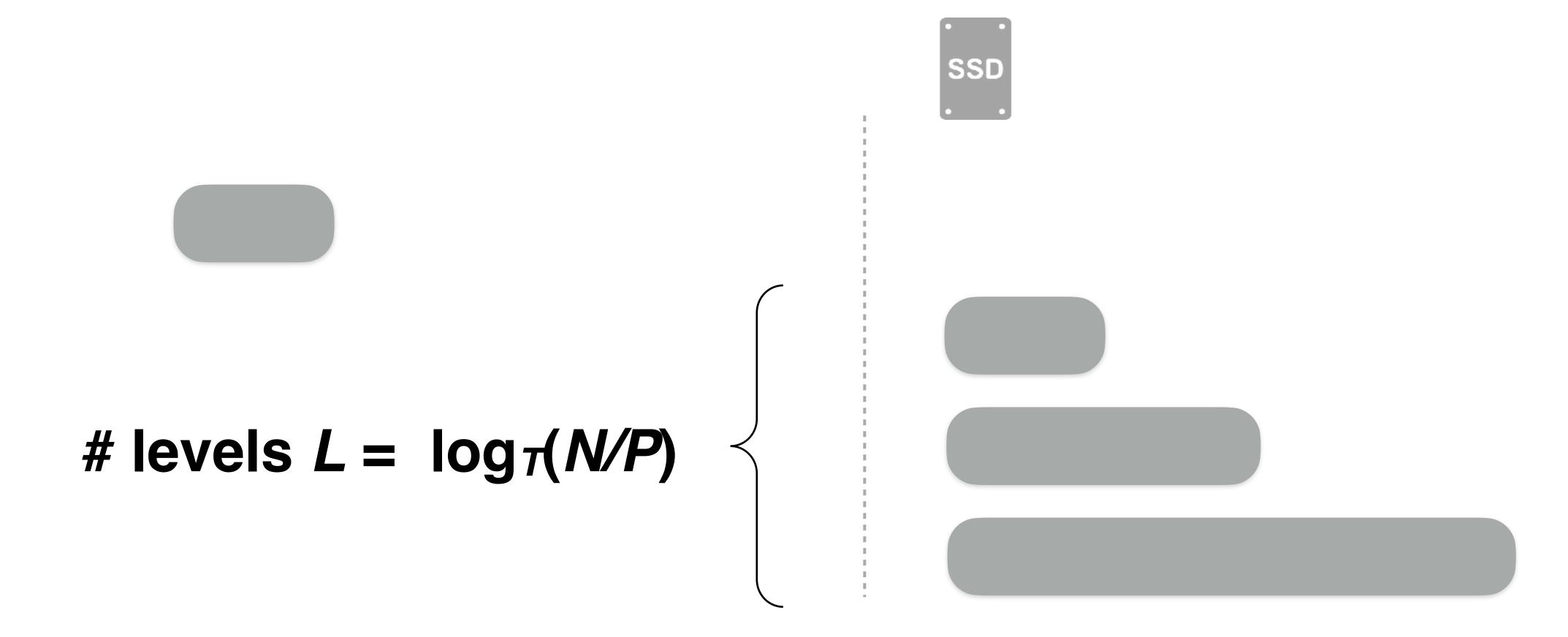
level 1

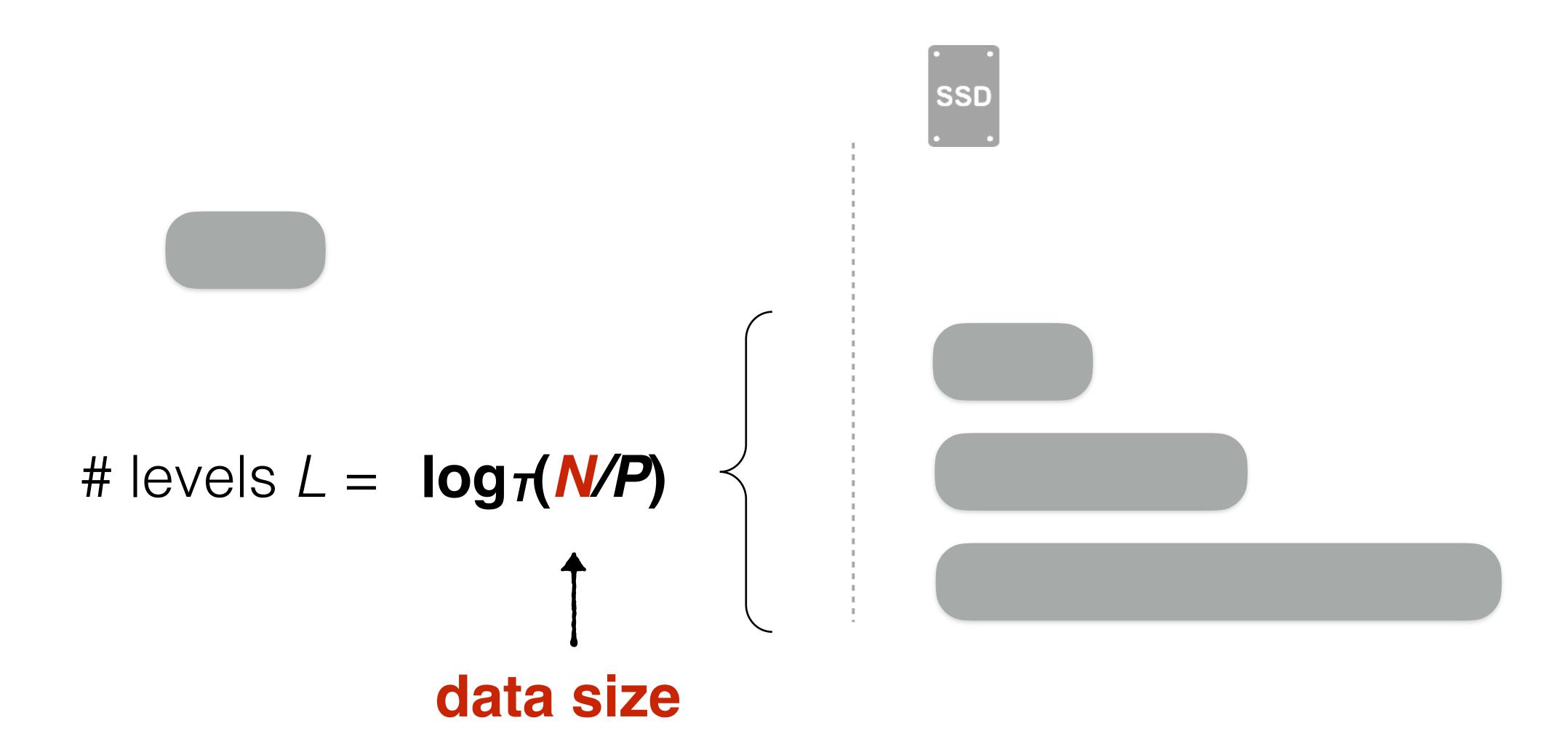
level 2 -

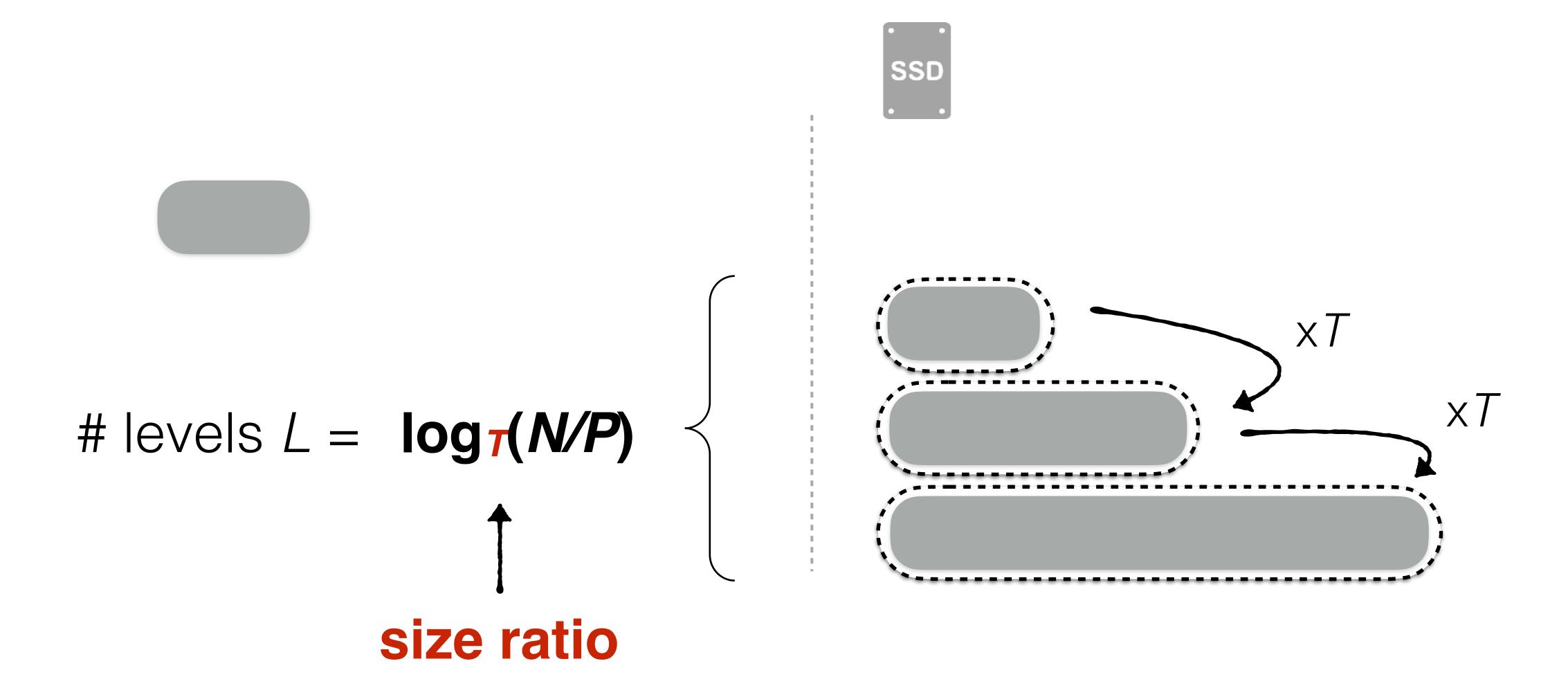
level 3 →

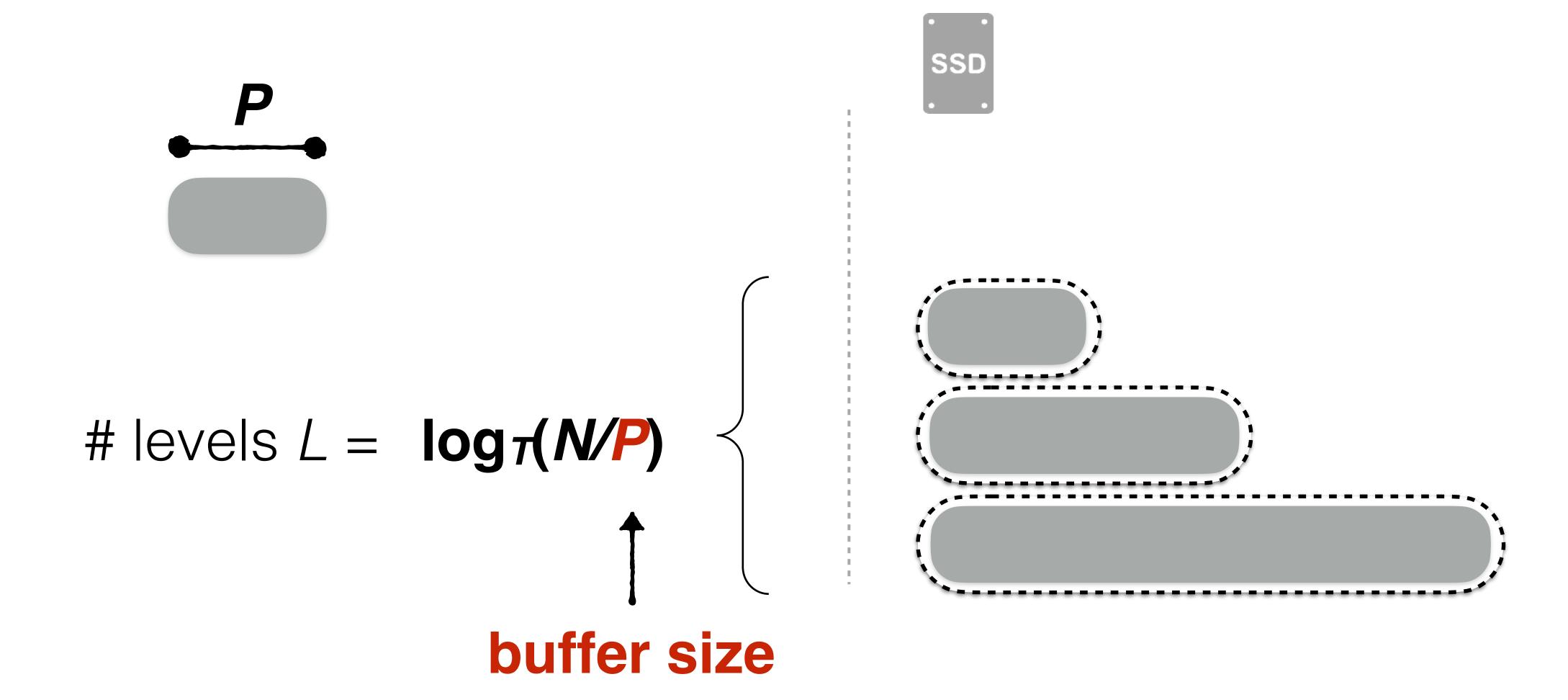


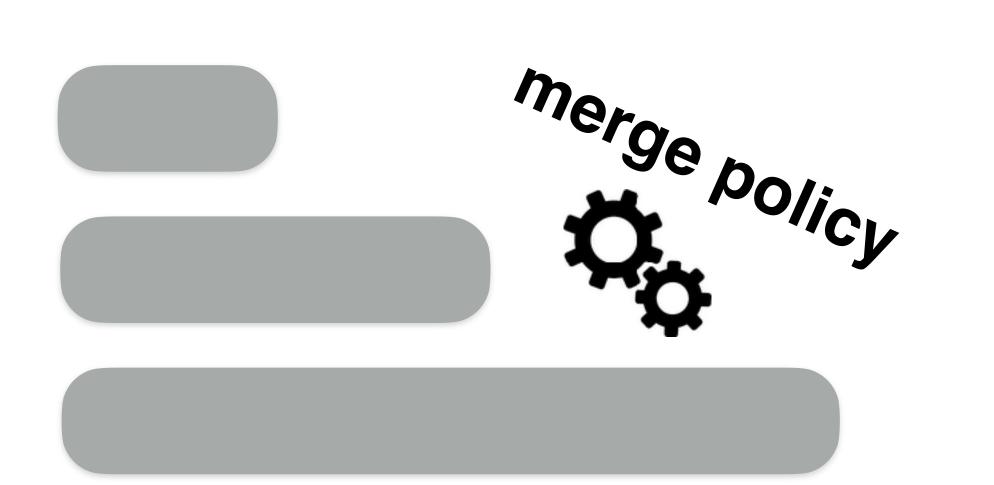






























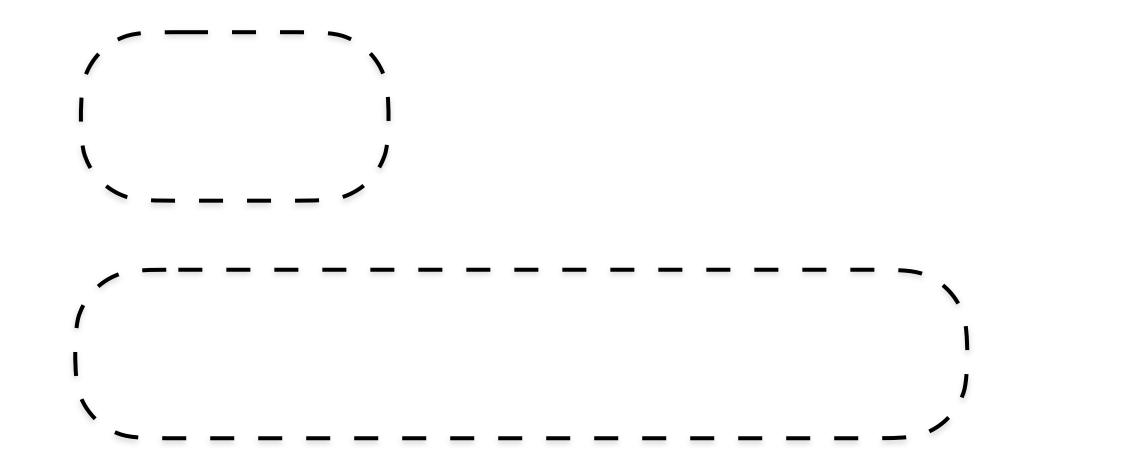
two merge policies

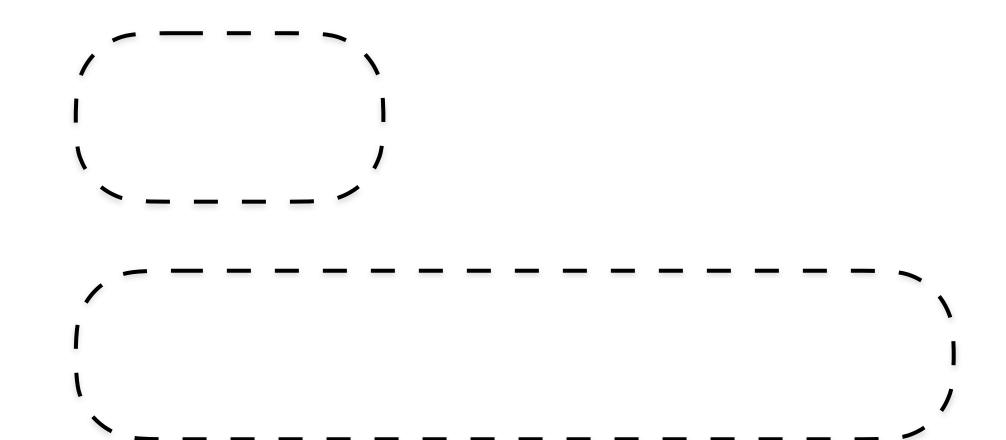


Leveling

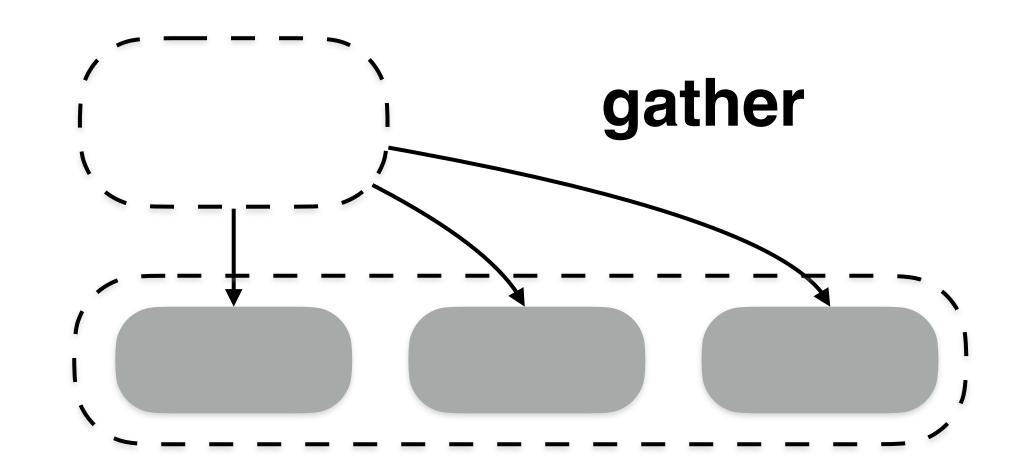


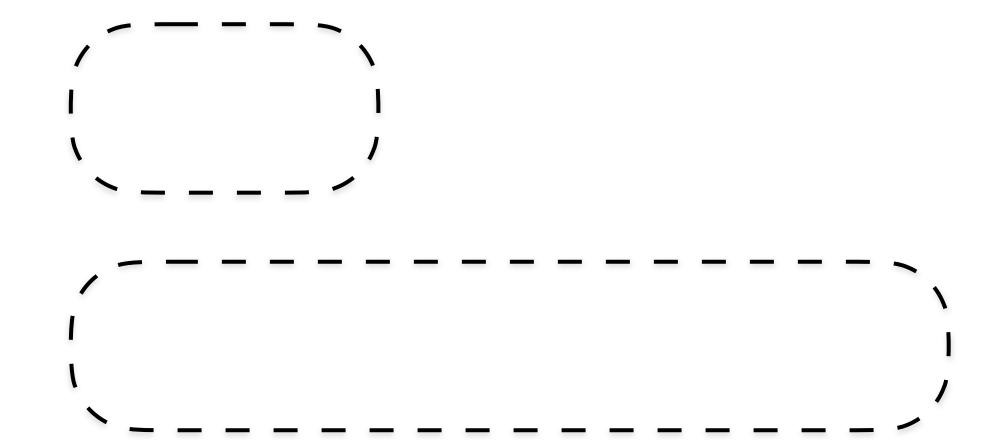






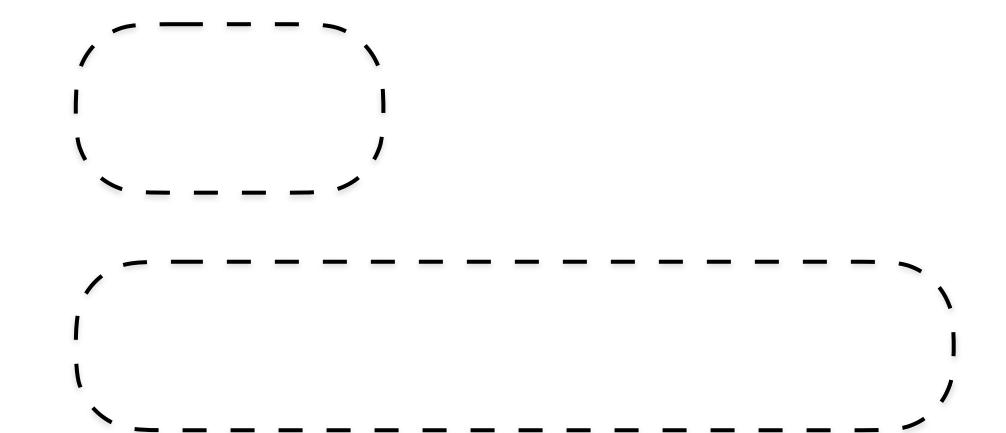








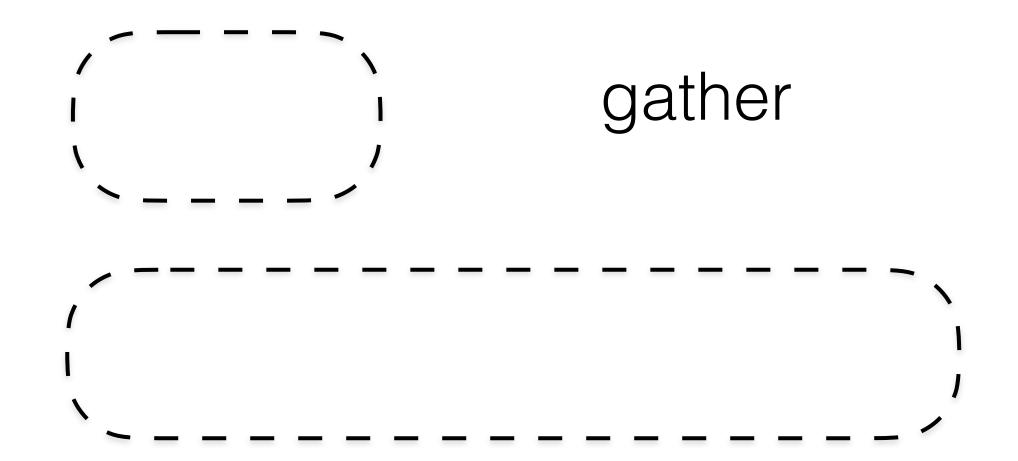
	gather	
merge & flush		

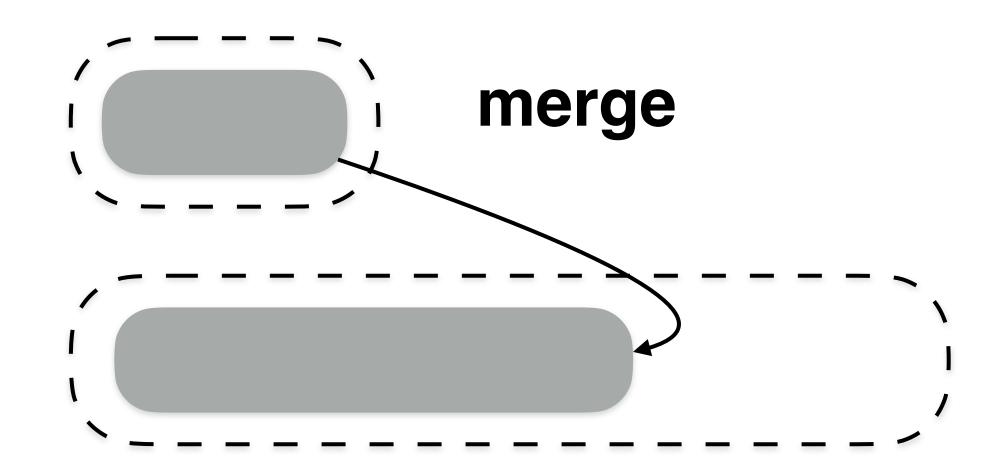




gather	







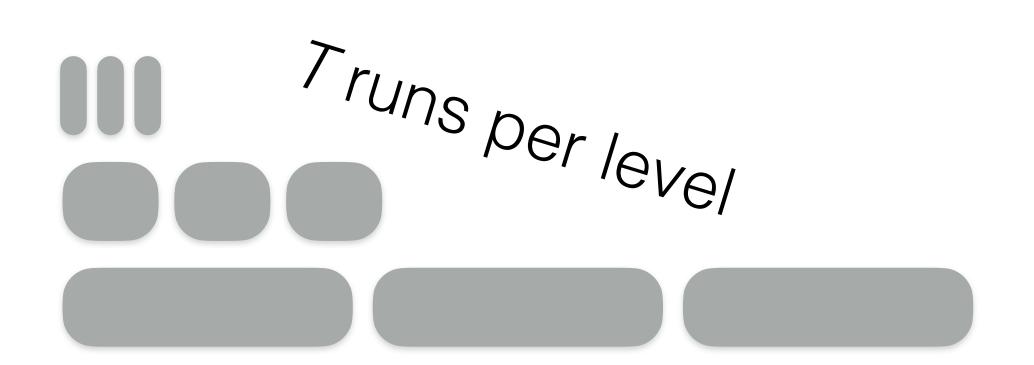


gather	,



gather	merge







size ratio T





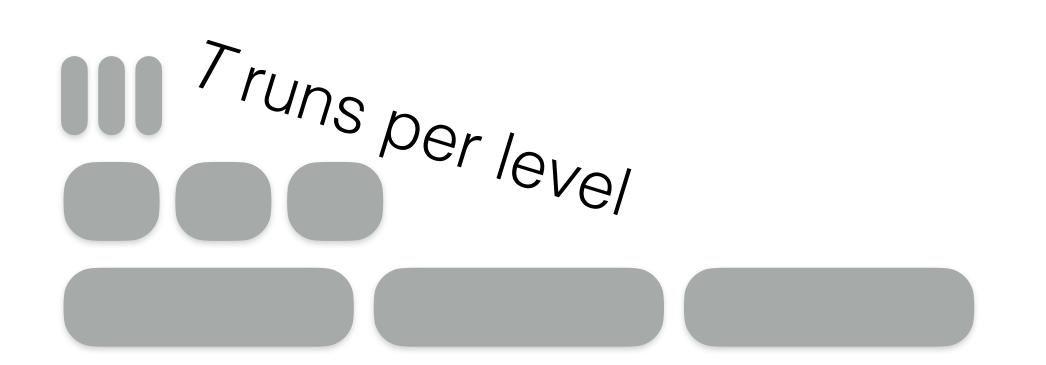




size ratio T = 2







size ratio T



O(N/P) runs per level



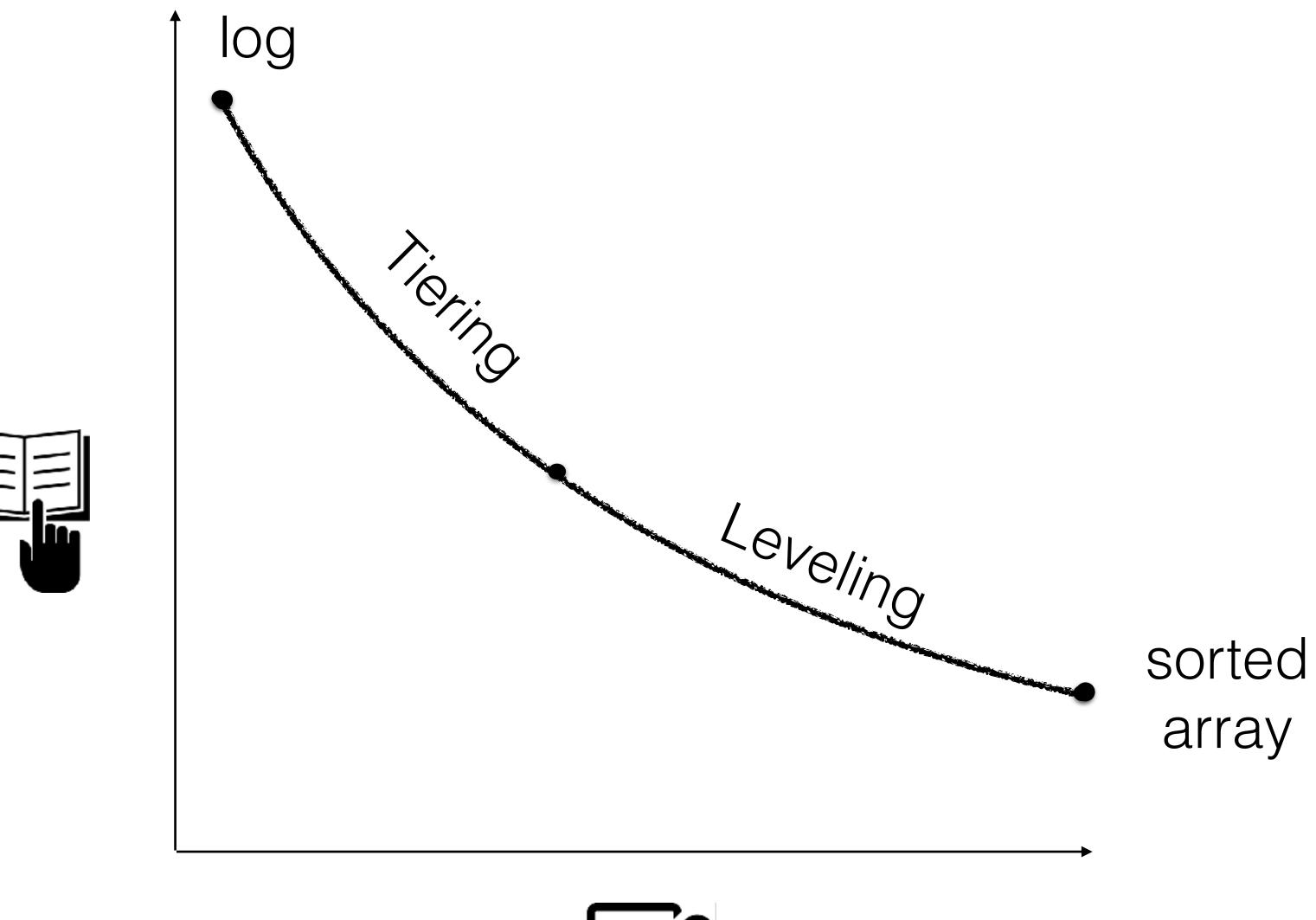
log

1 run per level

sorted array

size ratio T = N/P

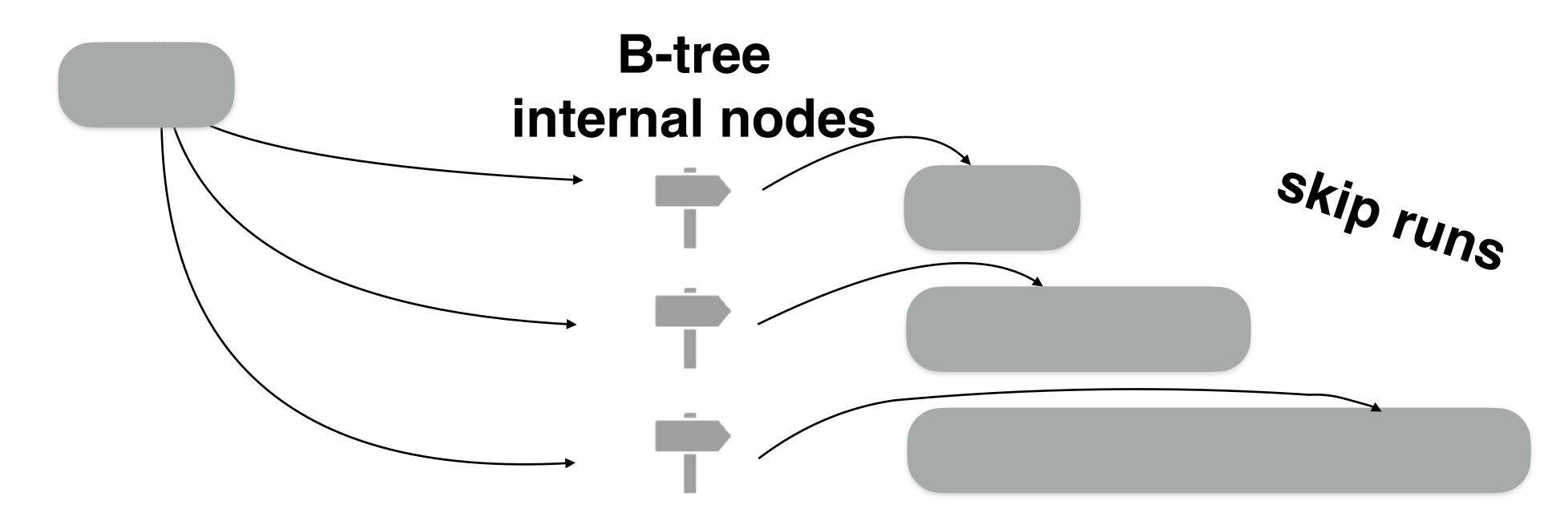


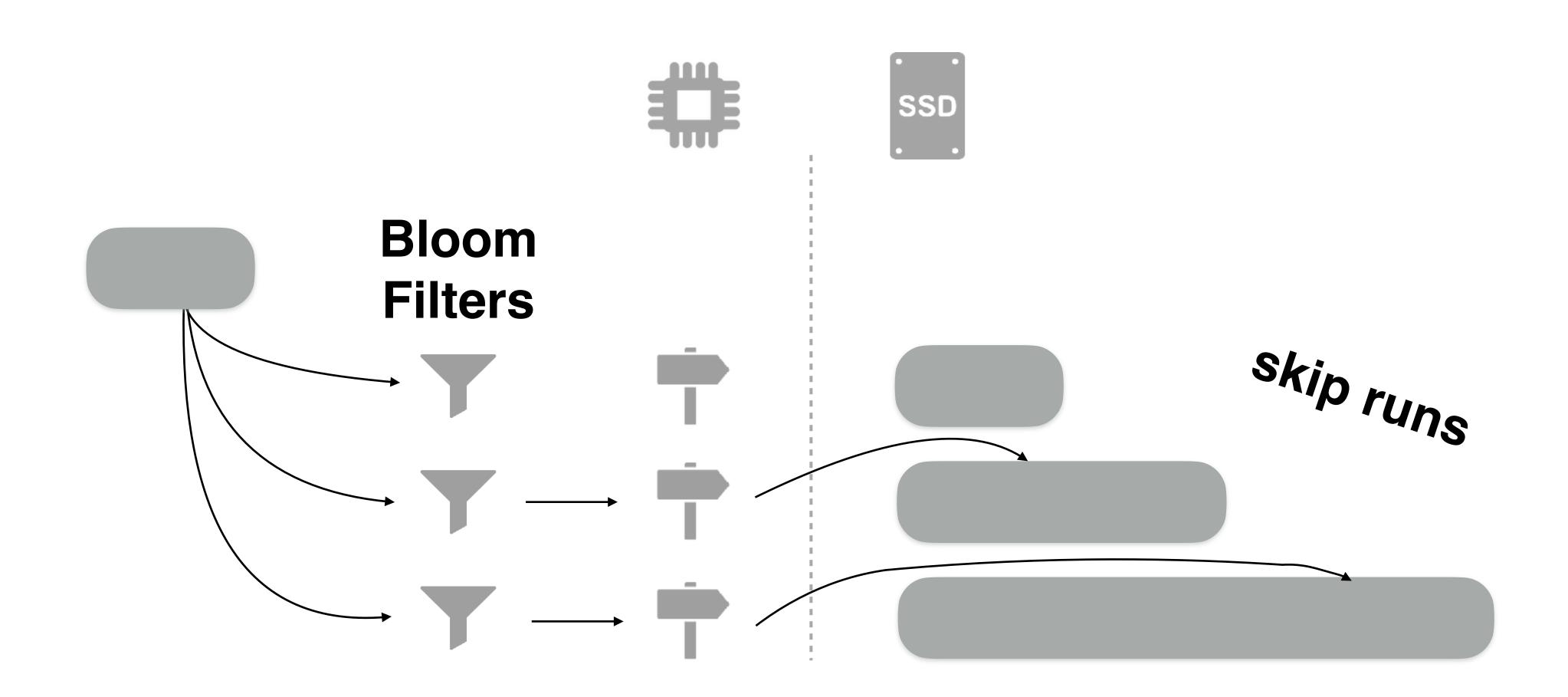












Why use a Filter?

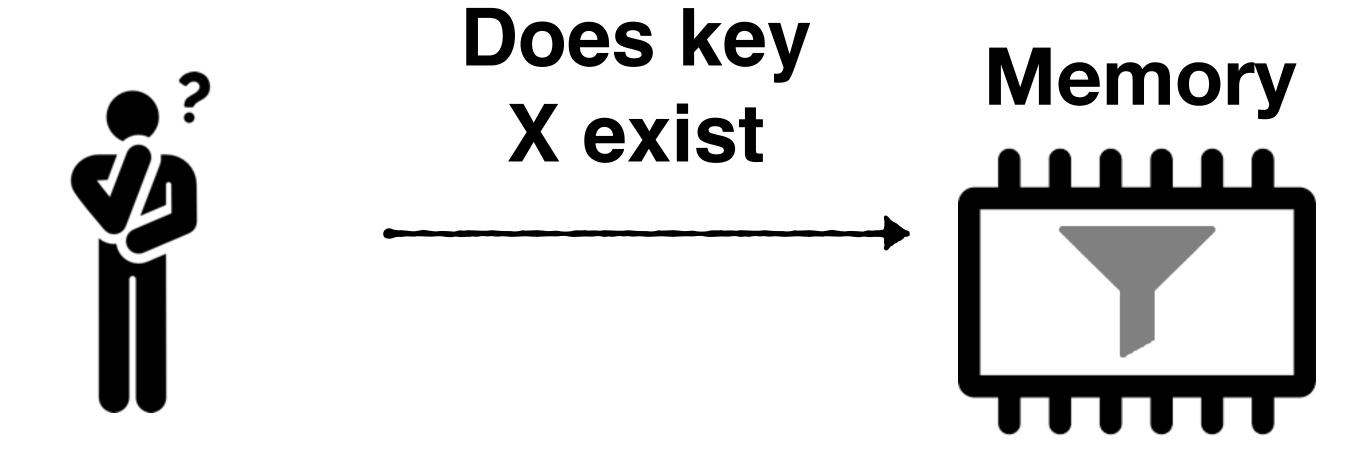




Does key X exist

Data

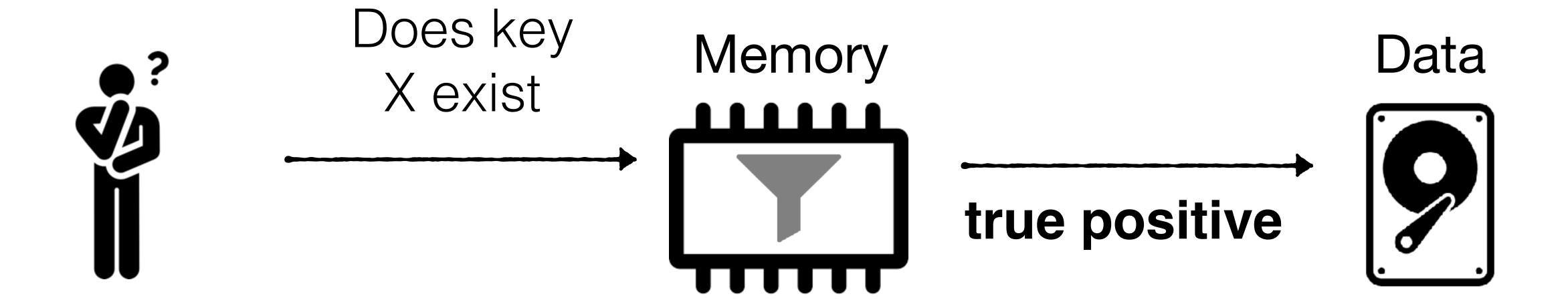


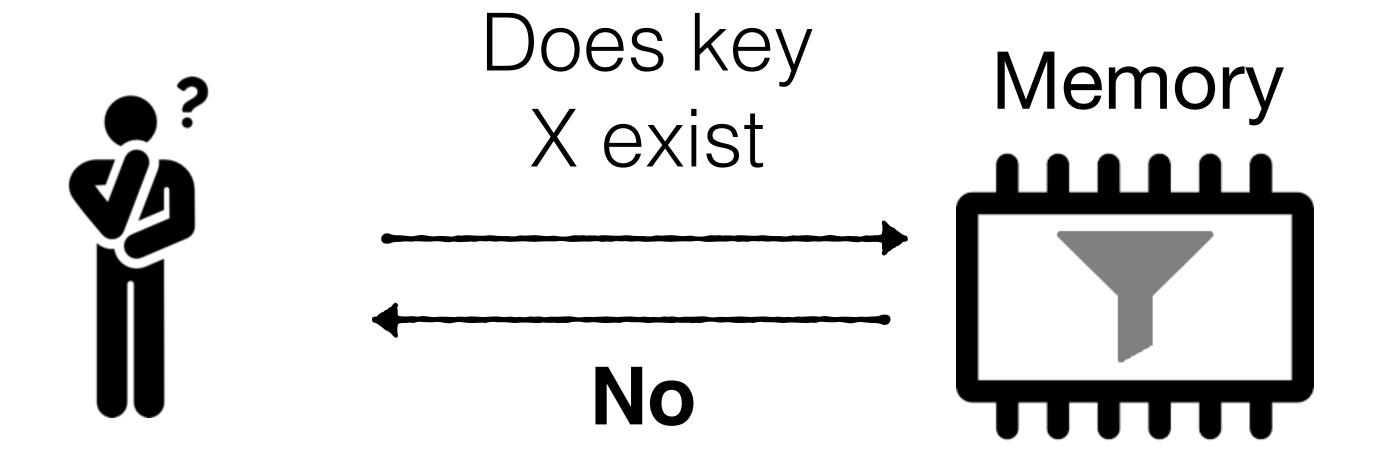


Data

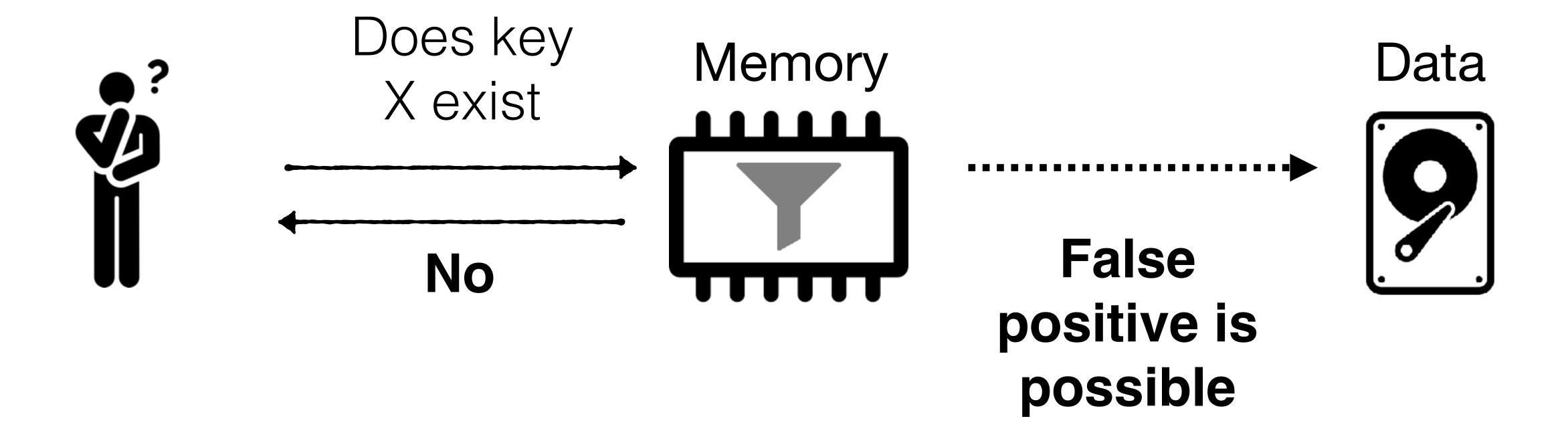


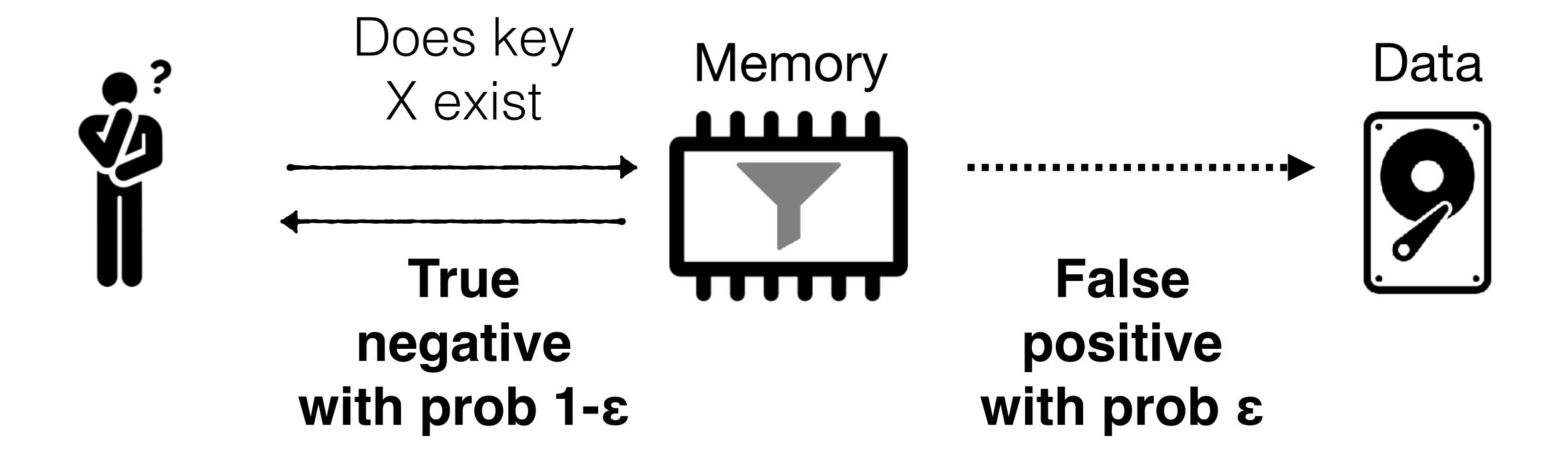
If key X exists

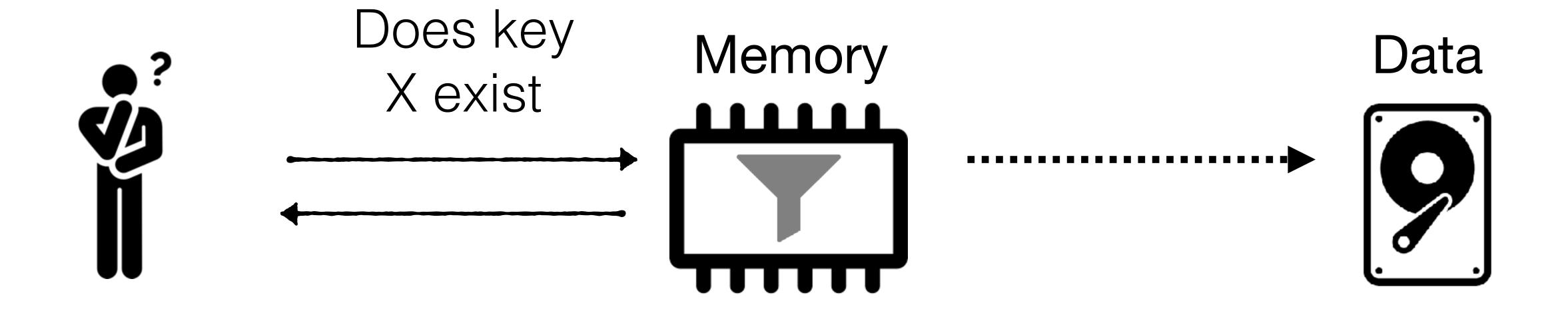












ε - false positive rate - FPR

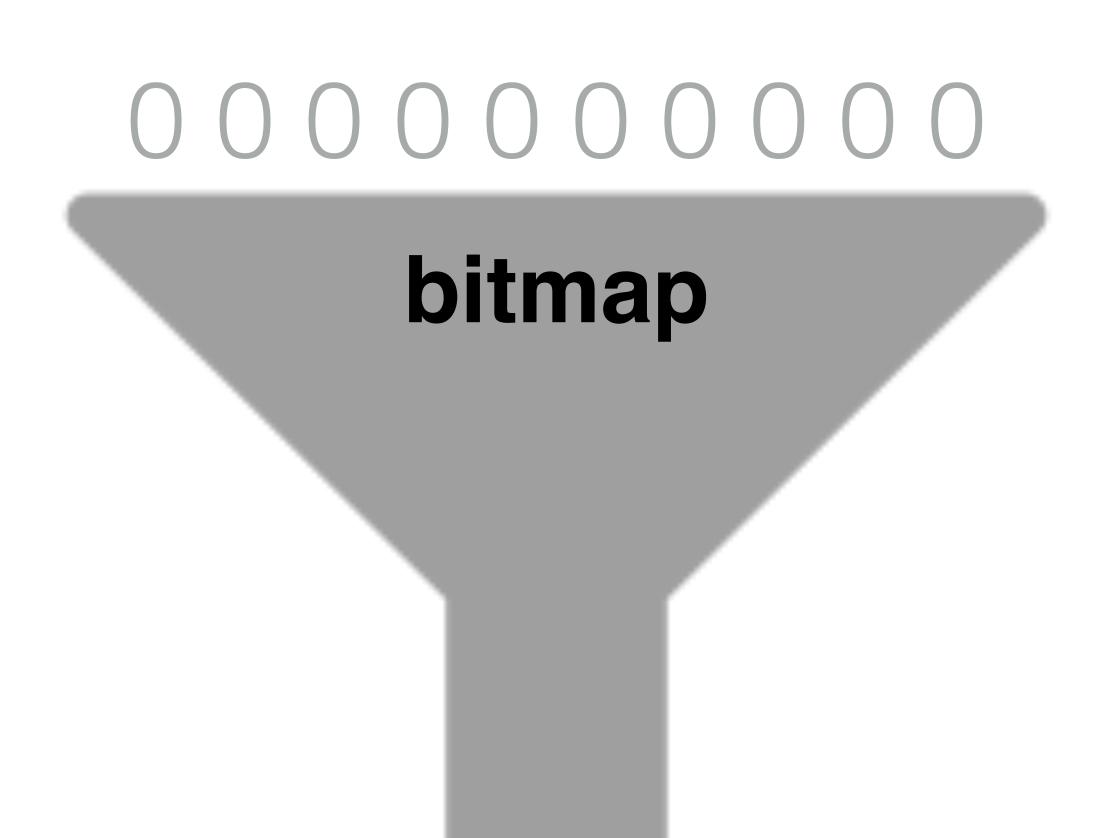
Bloom Filters



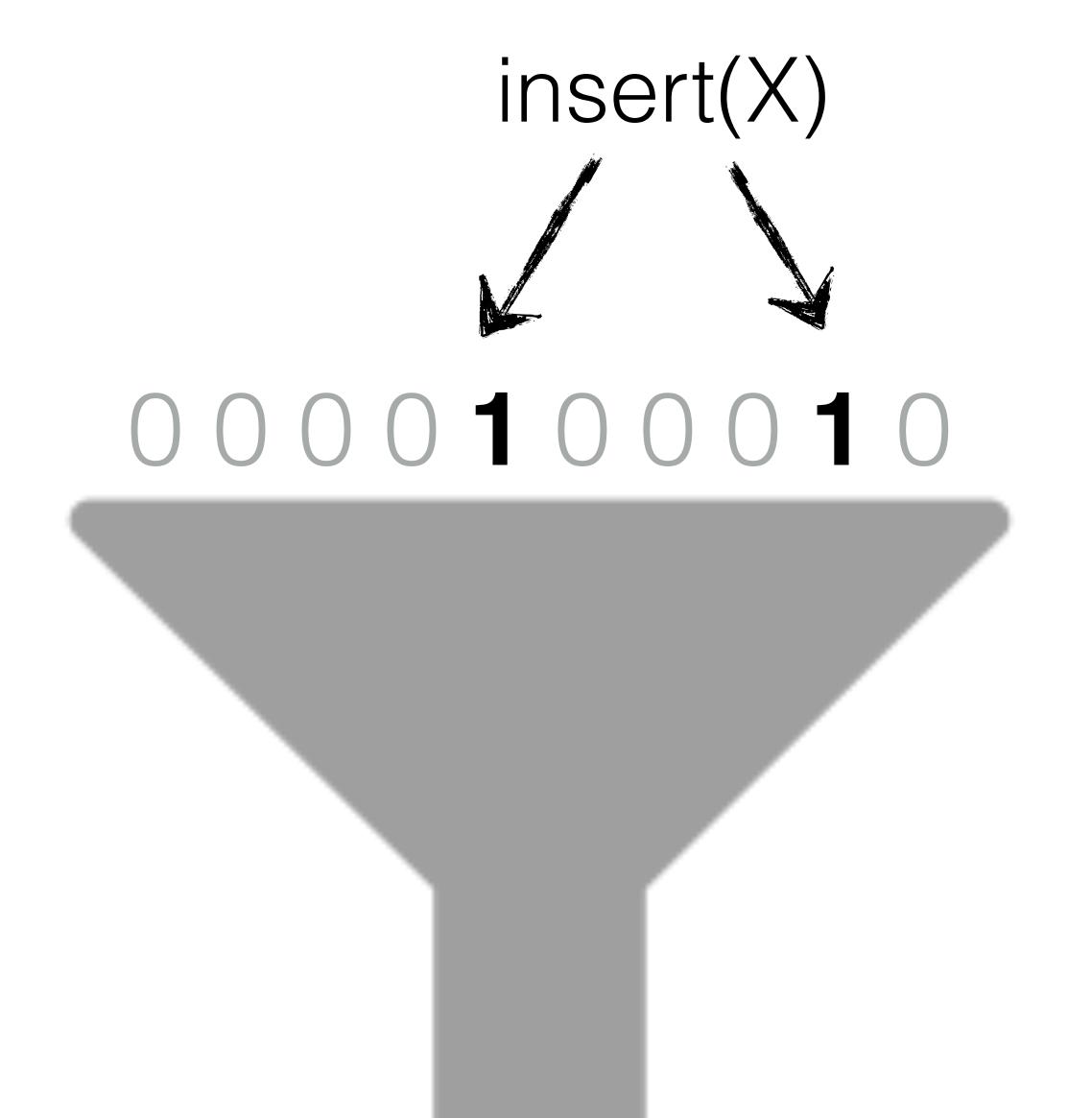
Bloom Filters

Space/time Trade-Offs in Hash Coding with Allowable Errors Burton Howard Bloom. Communications of the ACM, 1970.

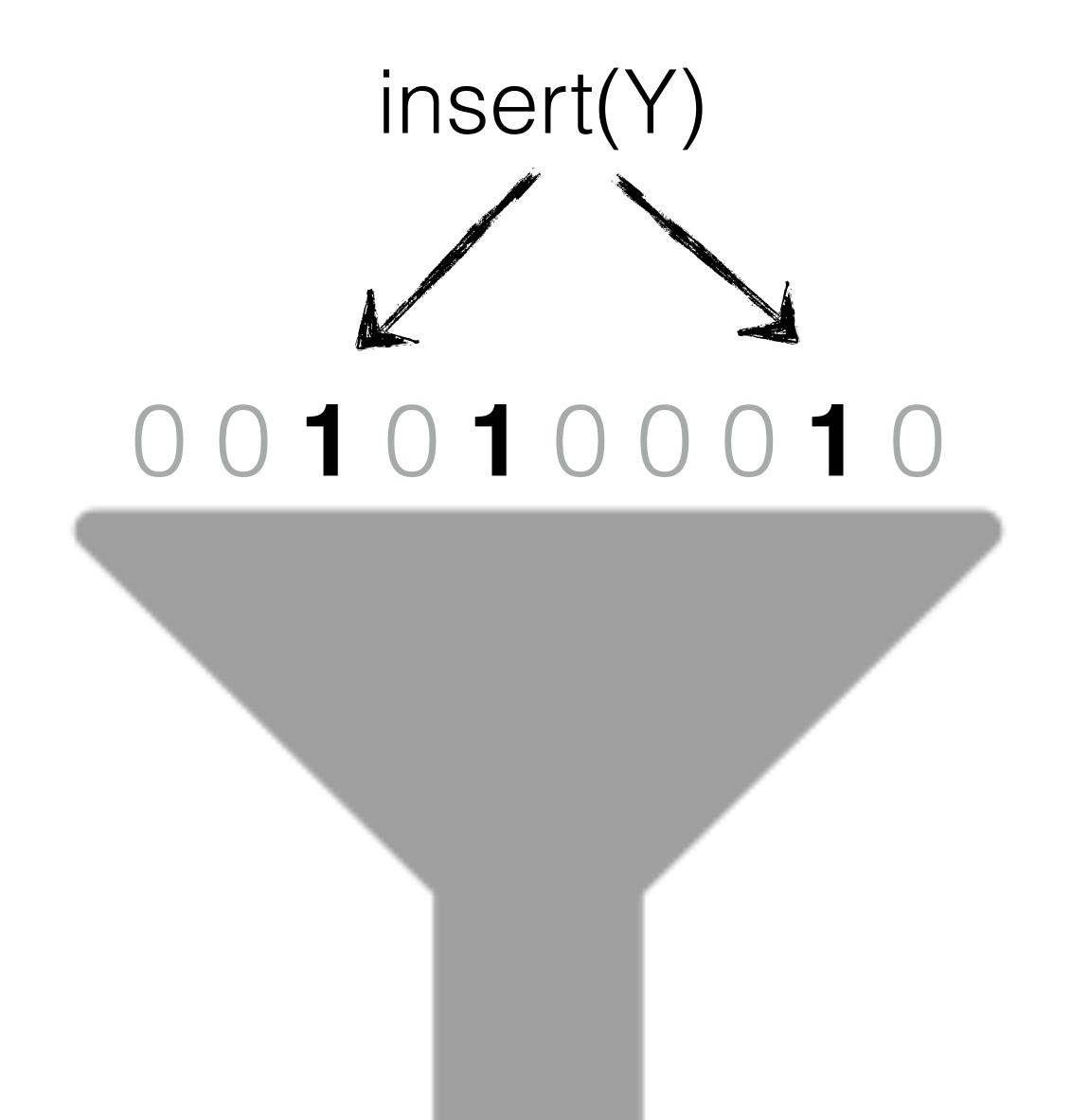
k hash functions

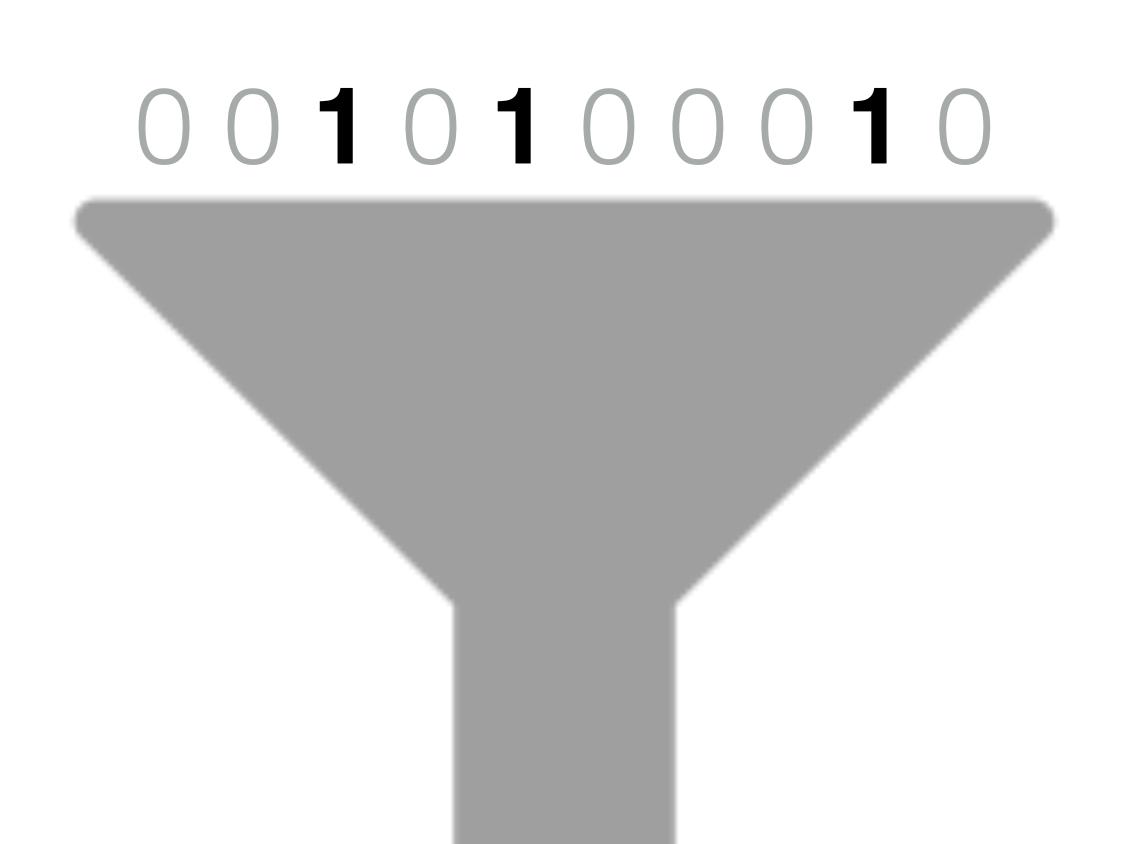


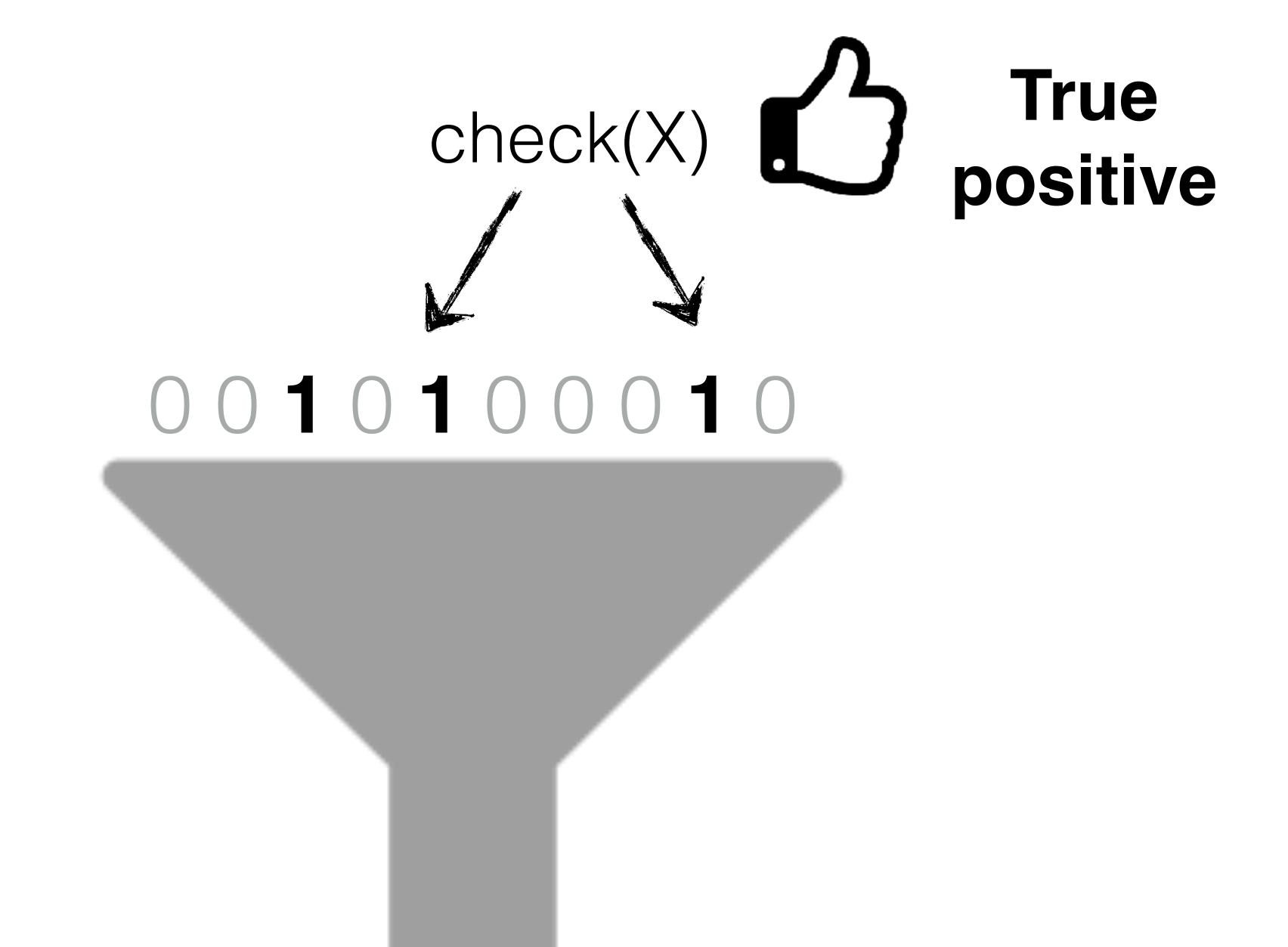
insert: Set from 0 to 1 or keep 1

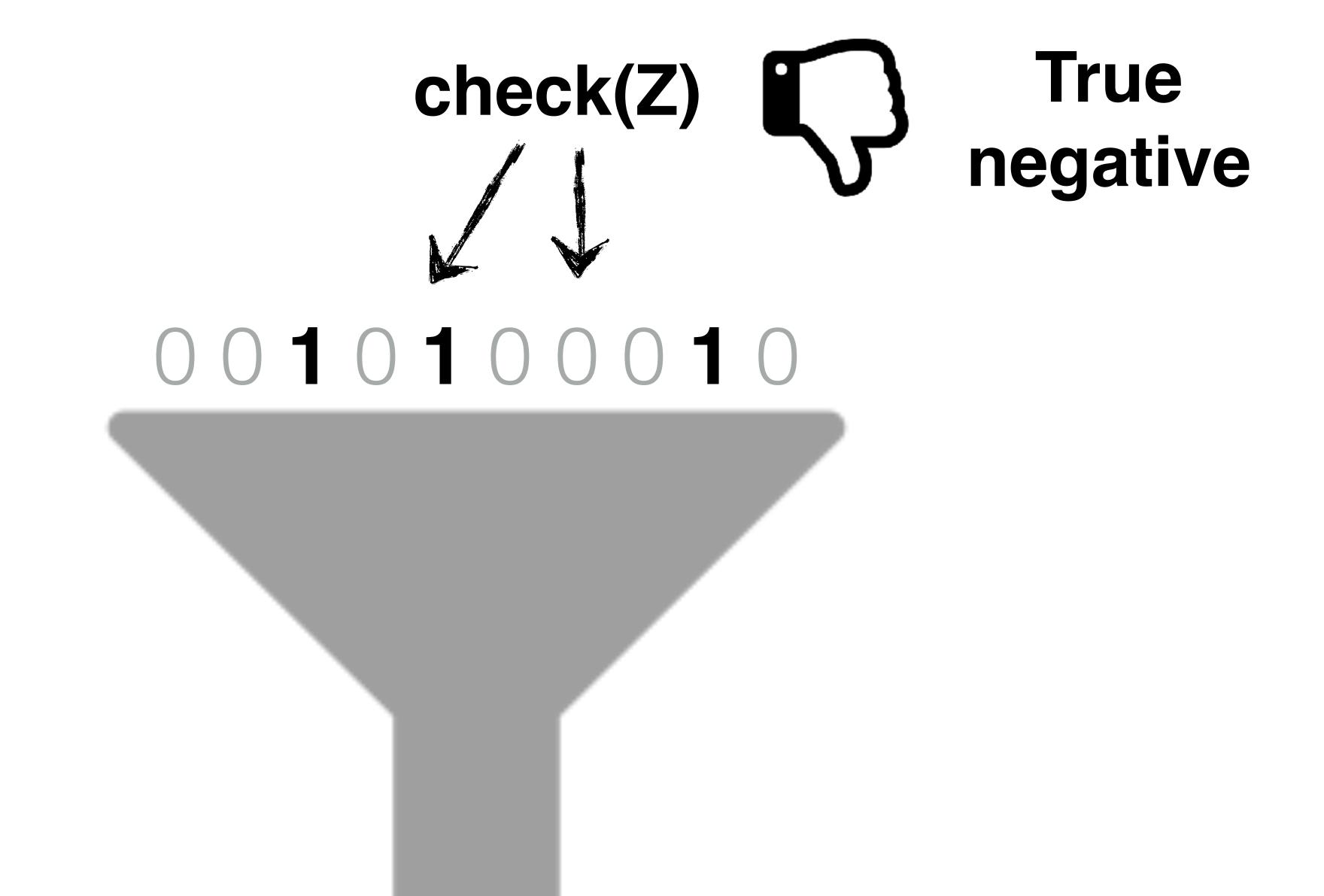


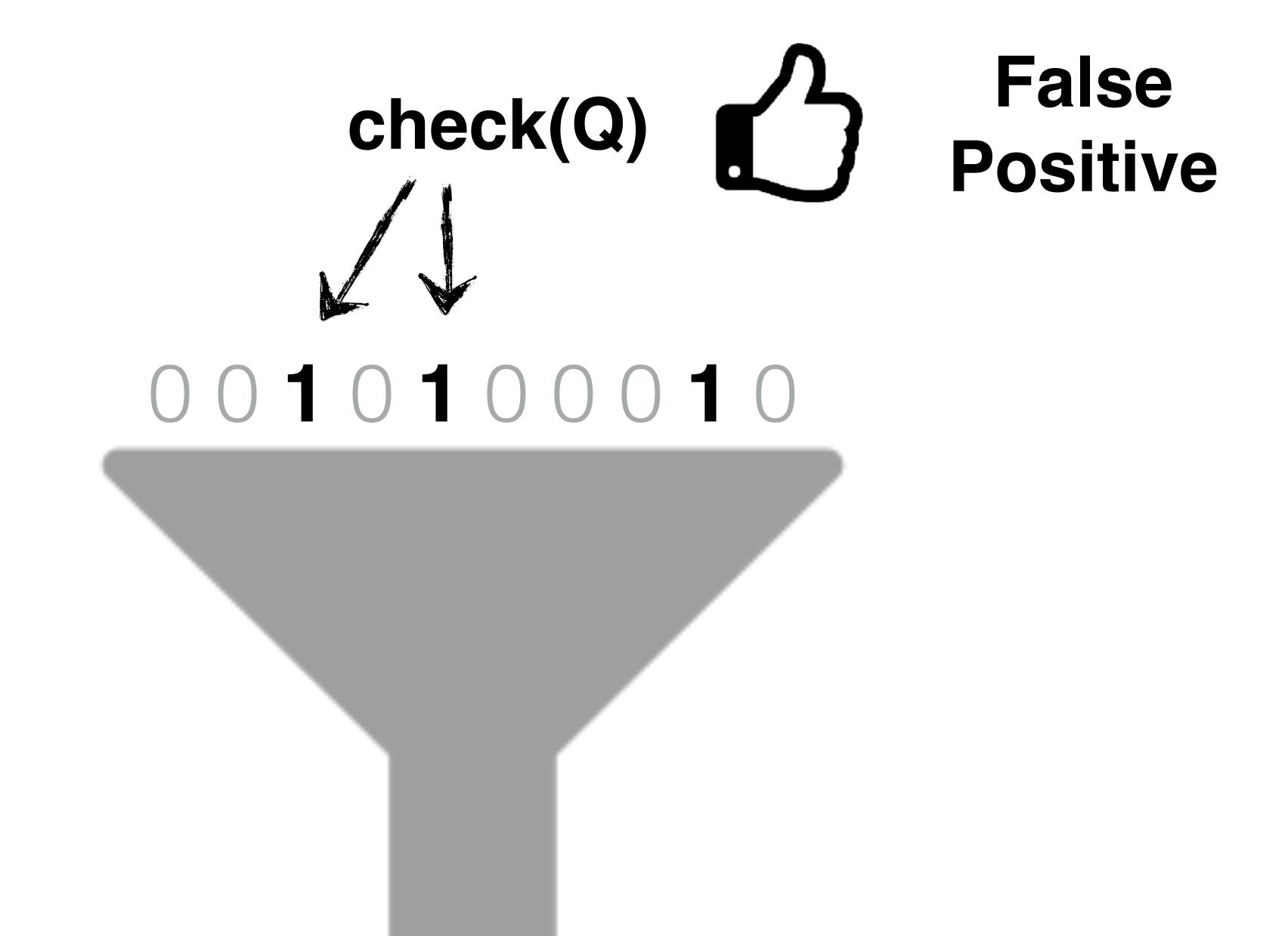
insert: Set from 0 to 1 or keep 1

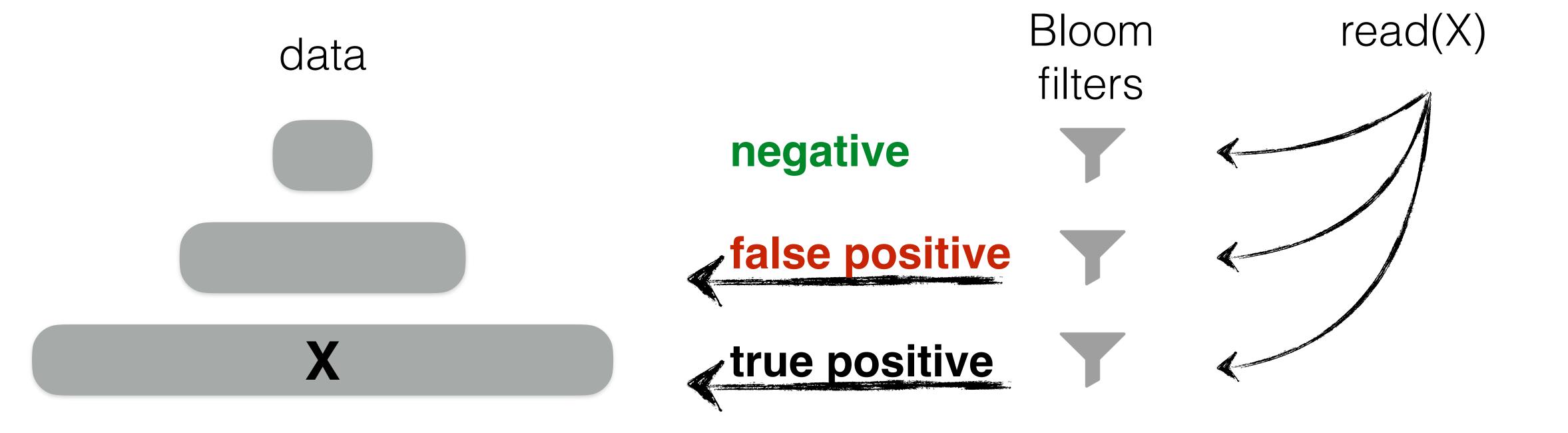




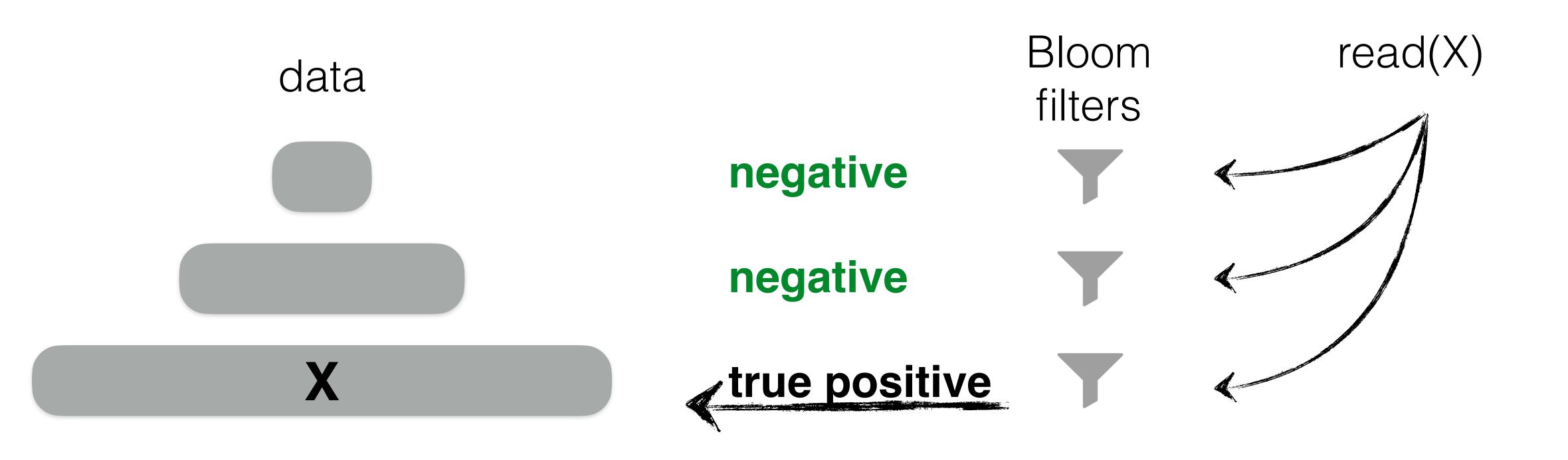


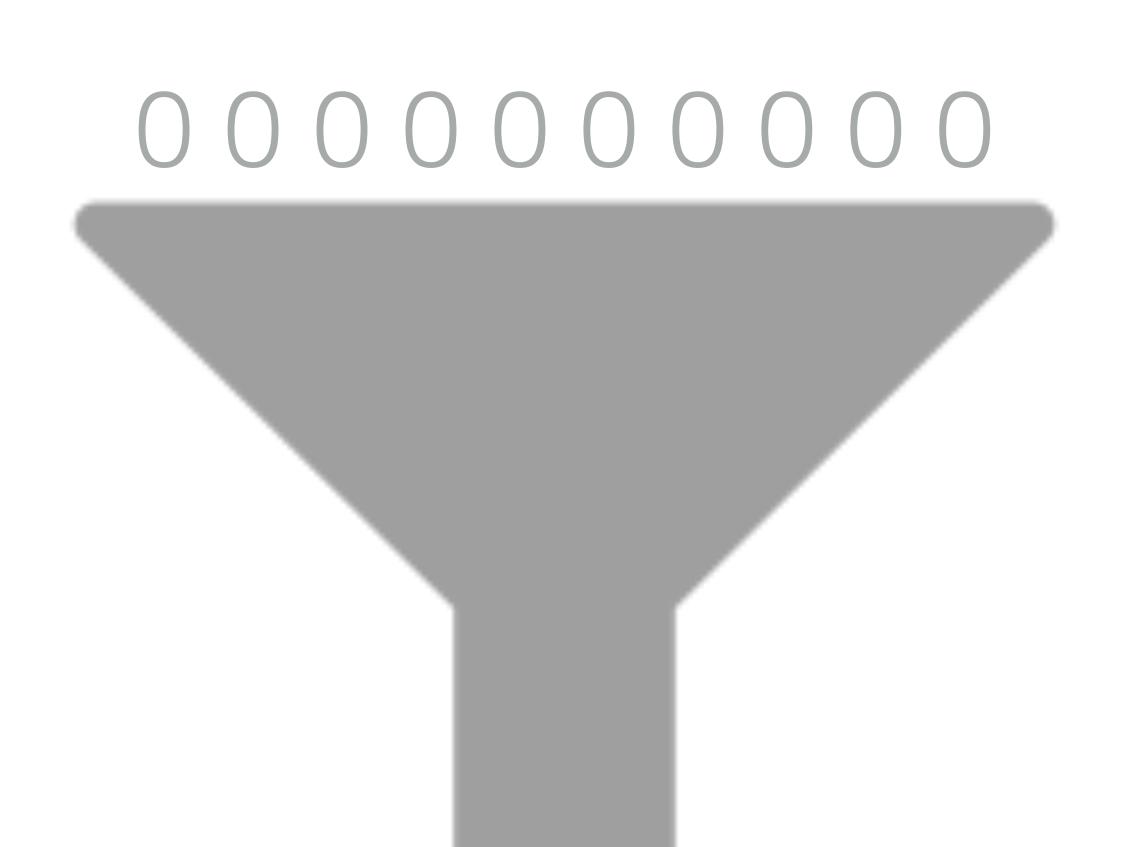


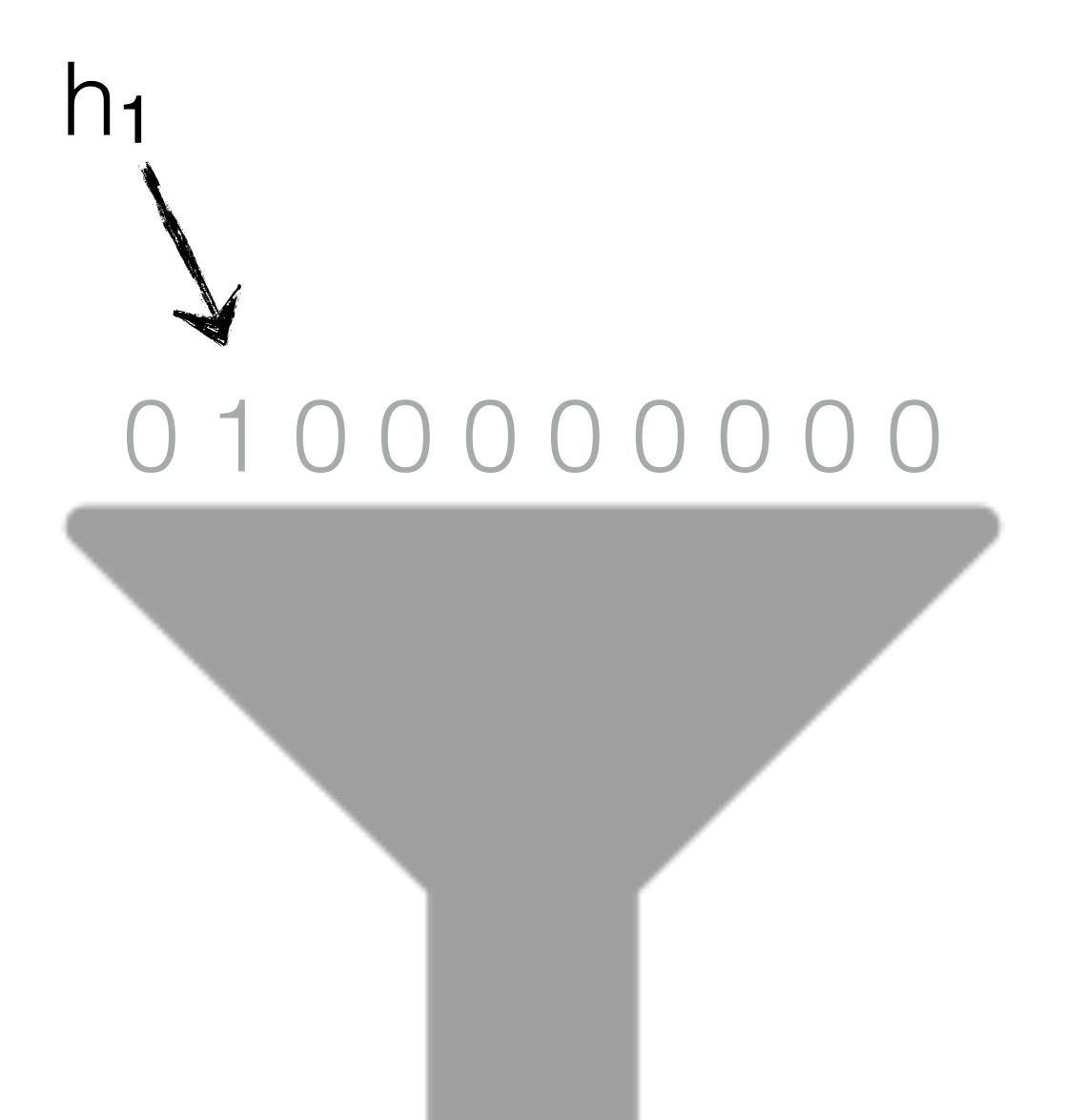




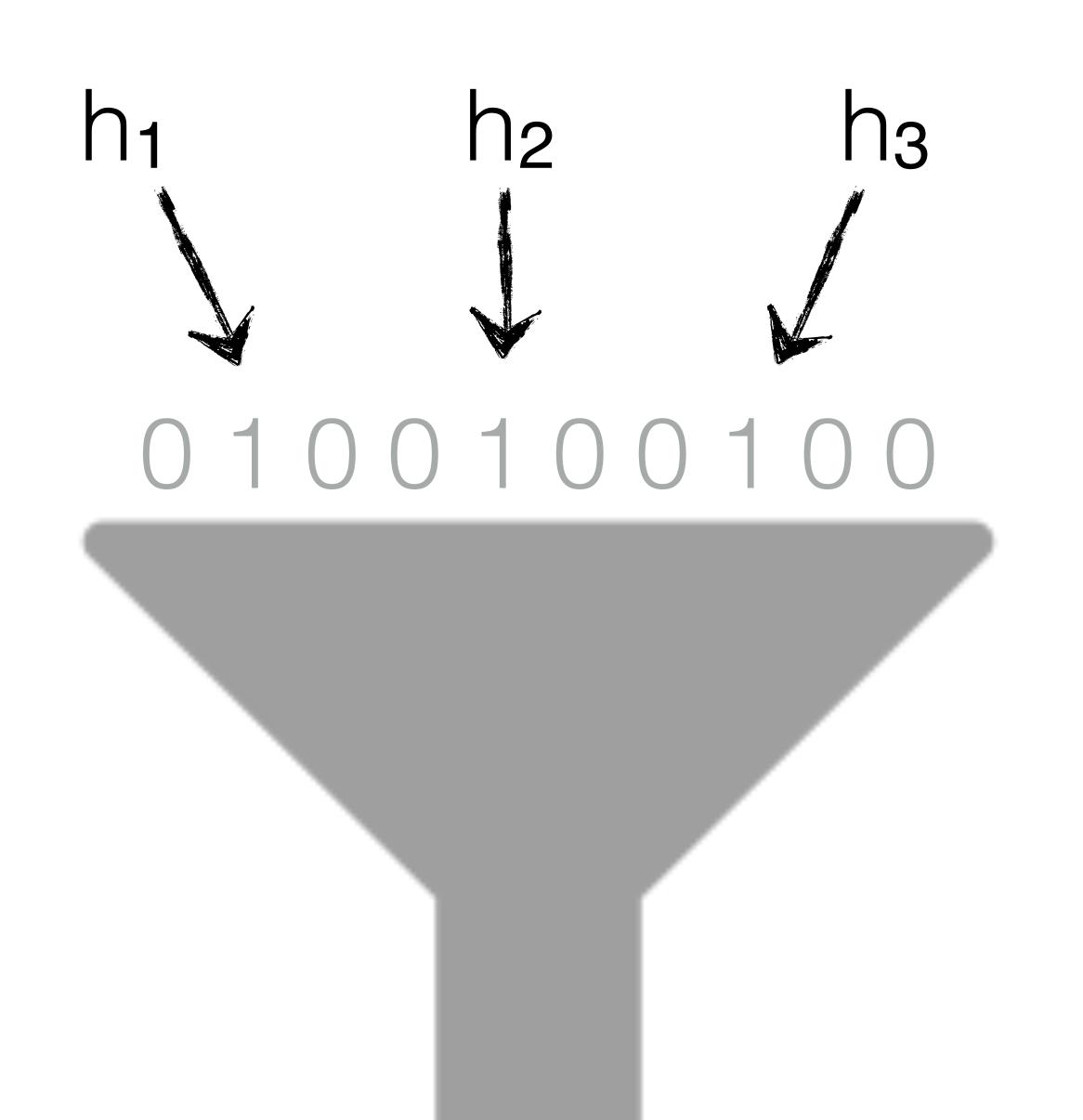
more memory — fewer false positives





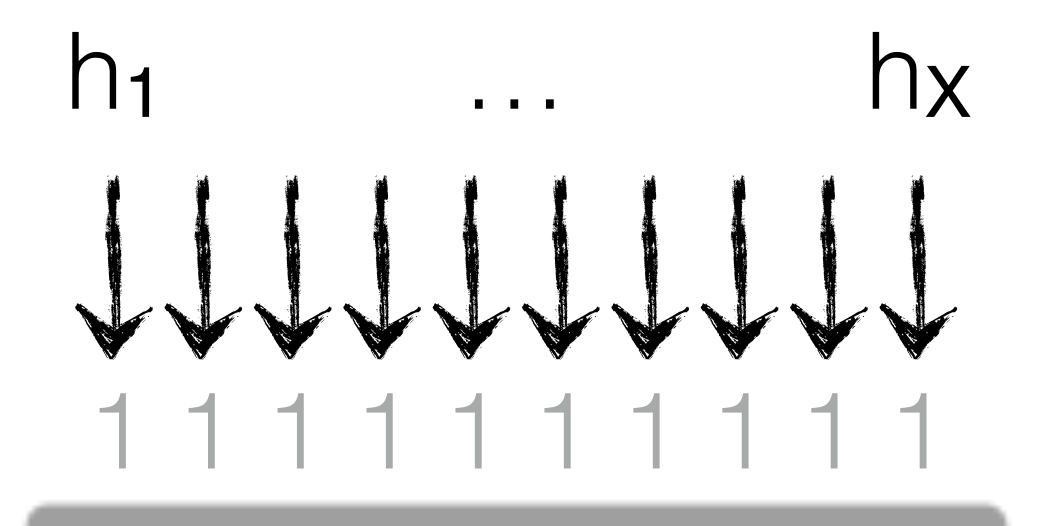


One is too few: false positive occurs whenever we hit a 1



One is too few: false positive occurs whenever we hit a 1

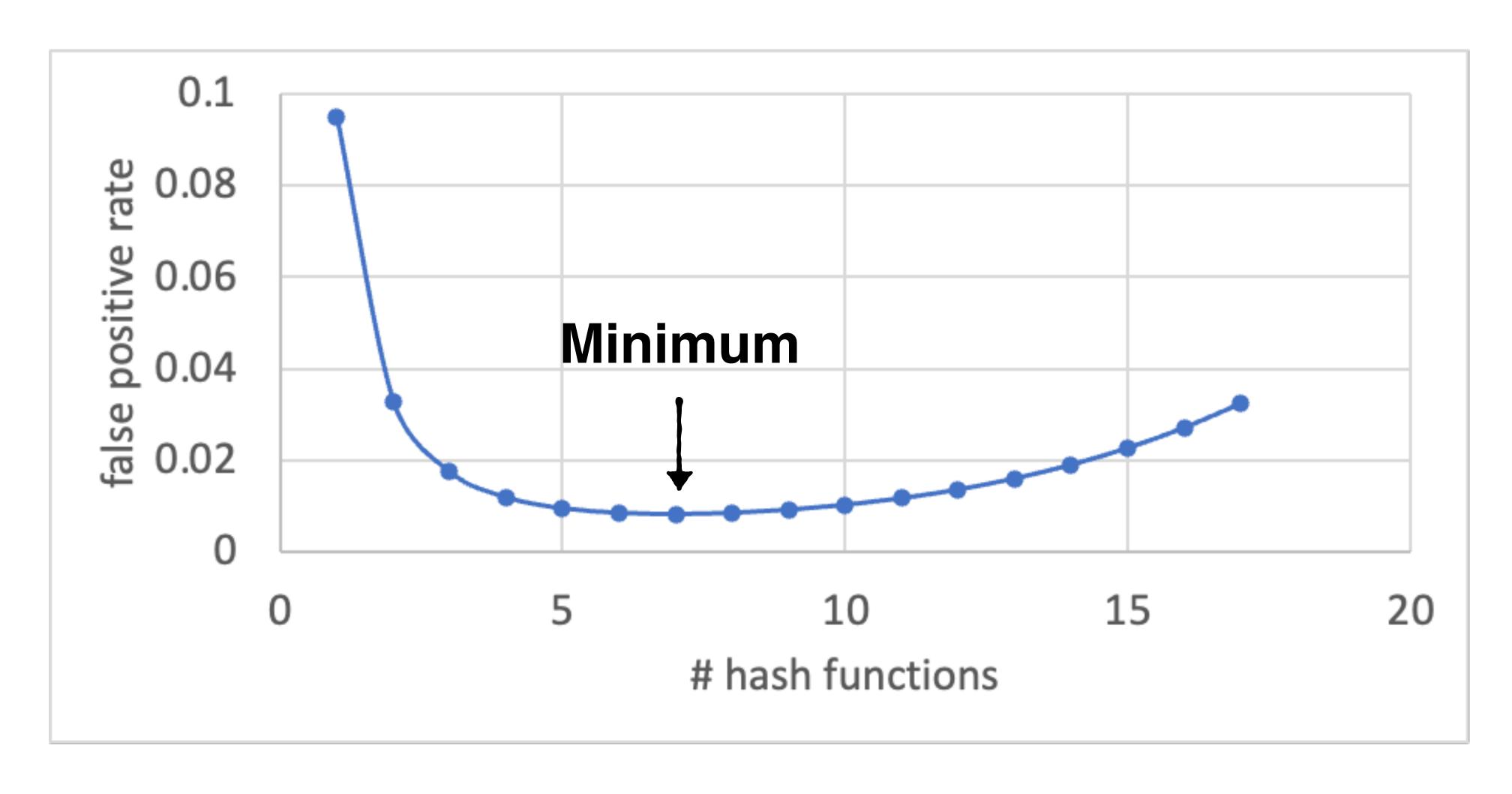
By adding hash functions, we initially decrease the false positive rate (FPR).



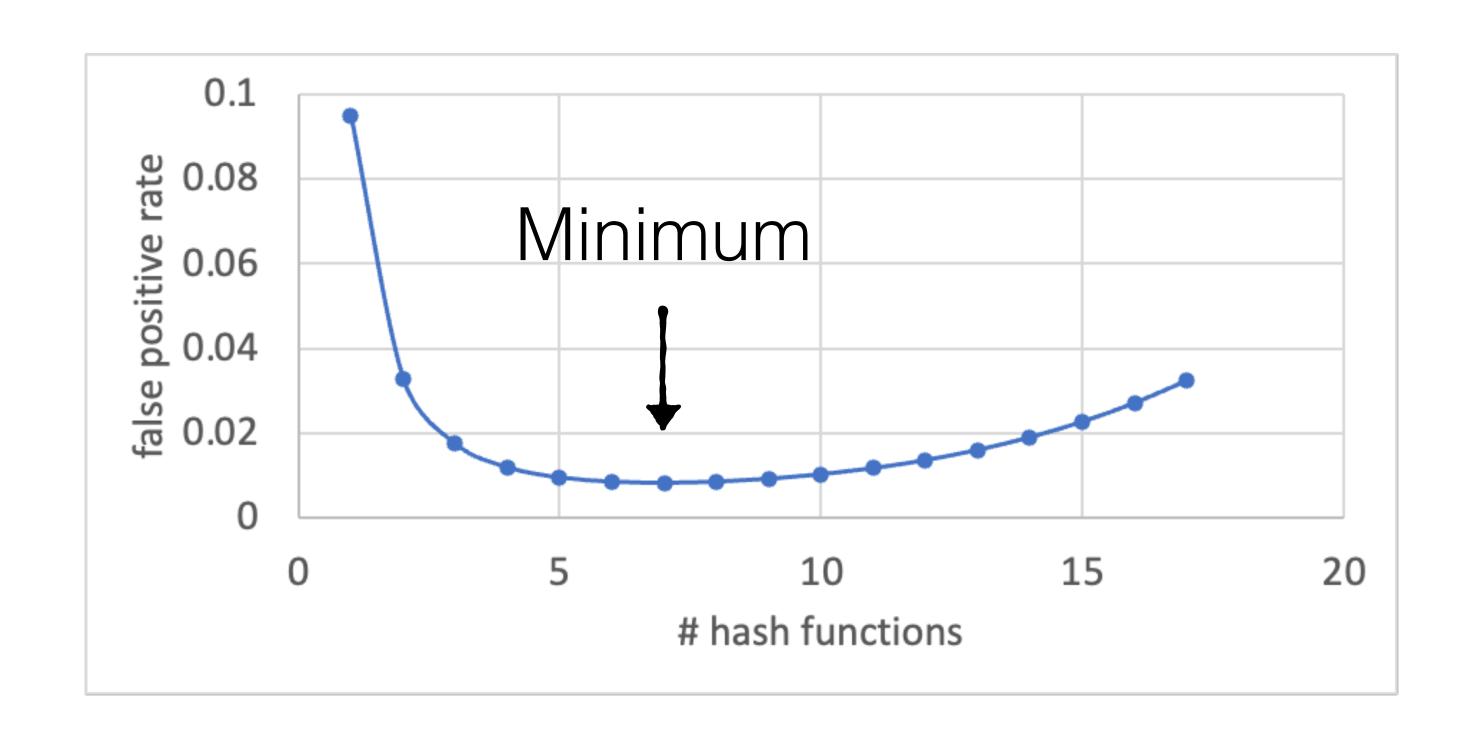
One is too few: false positive occurs whenever we hit a 1

By adding hash functions, we initially decrease the false positive rate (FPR).

But too many hash functions wind up increasing the FPR.



(Drawn for a filter using 10 bits per entry)

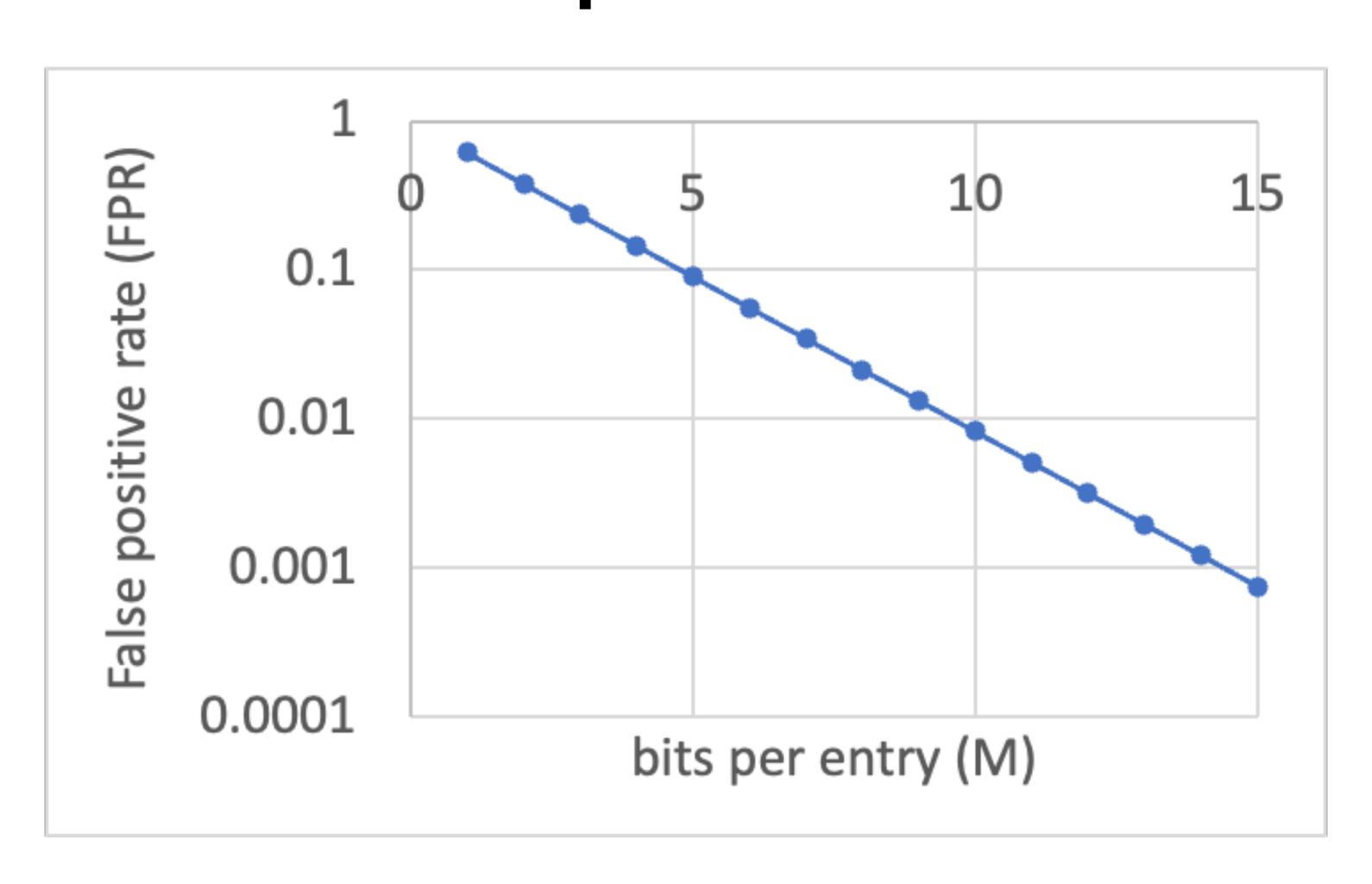


Optimal # hash functions = $ln(2) \cdot M$

(M is the number of bits per entry)

assuming the optimal # hash functions,

false positive rate = $2^{-M \cdot \ln(2)}$



Operation Costs (in memory accesses)

Insertion =

Positive Query =

Avg. Negative Query =

Insertion = $M \cdot In(2)$ (# hash functions)

Positive Query =

Avg. Negative Query =

Insertion = $M \cdot ln(2)$

Positive Query = $M \cdot ln(2)$ (# hash functions)

Avg. Negative Query =

Insertion = $M \cdot ln(2)$

Positive Query = $M \cdot ln(2)$

Avg. Negative Query =

(fraction of ones in filter is 0.5 with optimal number of hash functions)

Insertion = $M \cdot ln(2)$

Positive Query = $M \cdot ln(2)$

Avg. Negative Query = $1 + 1/2 (1 + 1/2 \cdot (...))$

(fraction of ones in filter is 0.5 with optimal number of hash functions)

Insertion = $M \cdot ln(2)$

Positive Query = $M \cdot ln(2)$

Avg. Negative Query = 1 + 1/2 + 1/4 + ... = 2

(fraction of ones in filter is 0.5 with optimal number of hash functions)

Insertion = $M \cdot ln(2)$

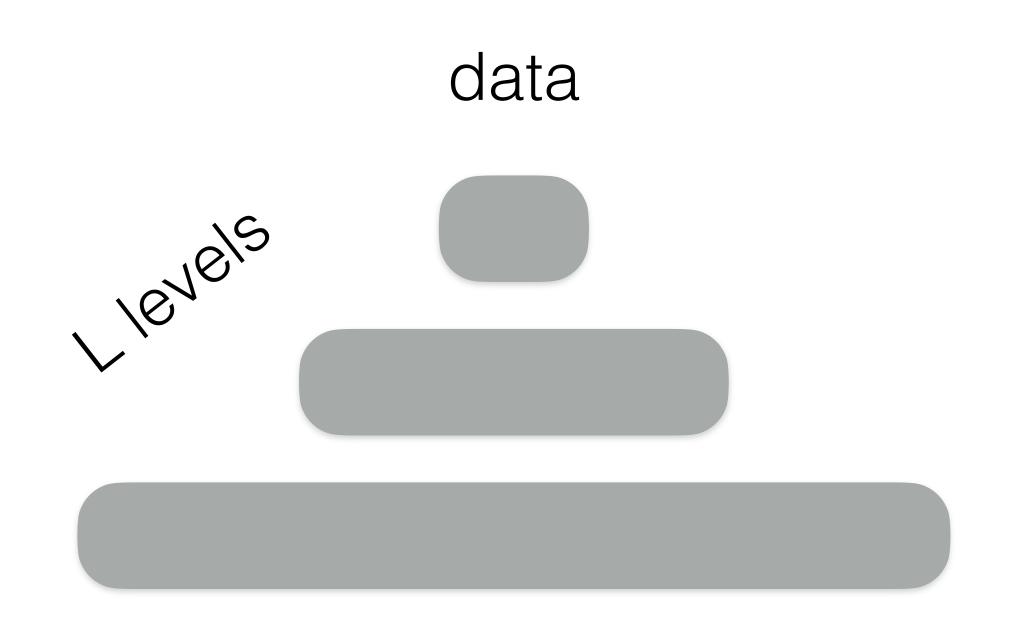
Positive Query = $M \cdot \ln(2)$

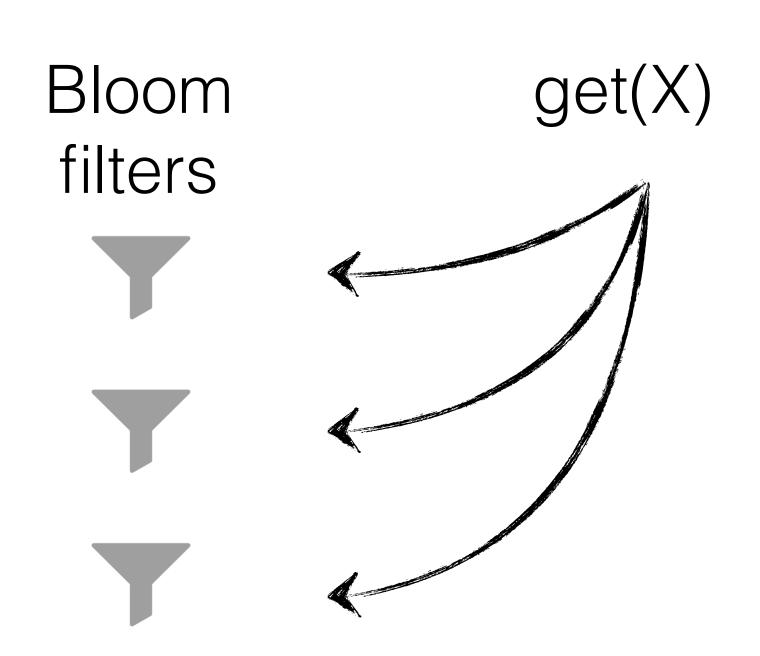
Avg. Negative Query = 2

false positive rate = $2^{-M \cdot \ln(2)}$

Positive Query = $M \cdot ln(2)$

Avg. Negative Query = 2





Worst-case:

Positive Query = $M \cdot \ln(2)$ Avg. Negative Query = 2

data

Bloom get(X)
filters

false positive

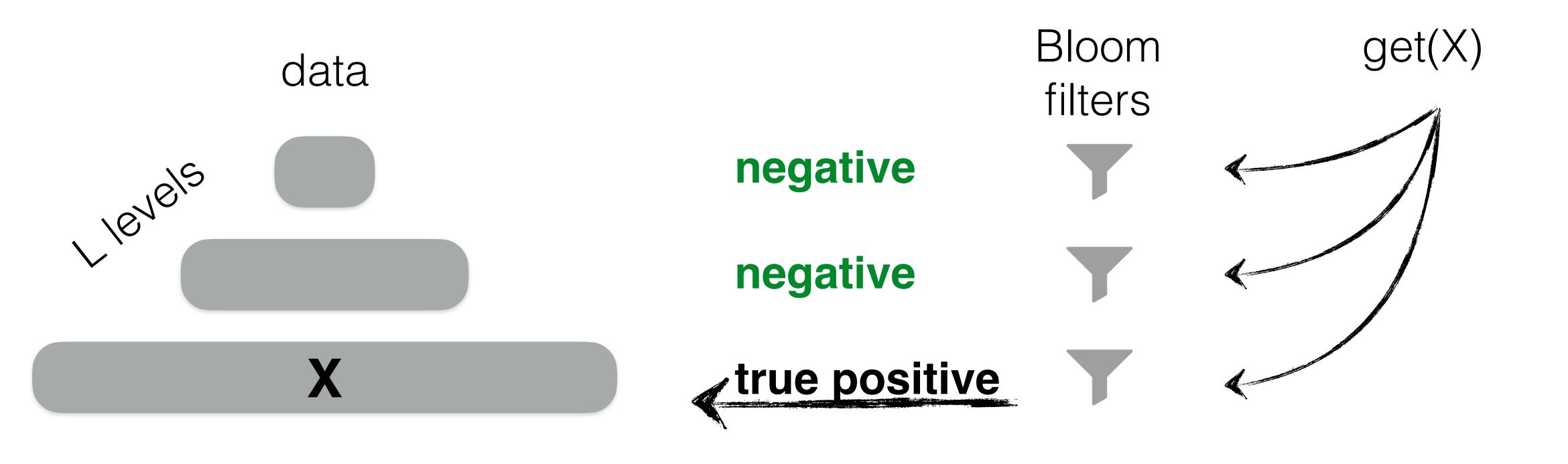
false positive

true positive

Worst-case: O(M·L)

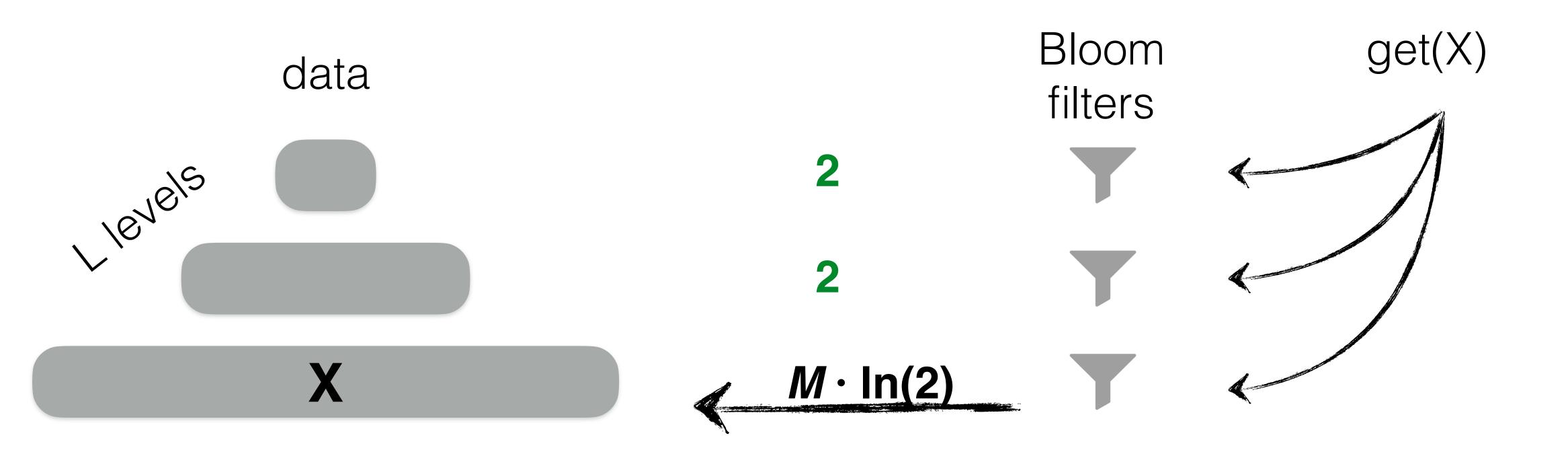
Positive Query = $M \cdot ln(2)$

Avg. Negative Query = 2



Worst-case: O(M·L)

Positive Query = $M \cdot \ln(2)$ Avg. Negative Query = 2



Worst-case: O(M·L)

Positive Query = $M \cdot \ln(2)$ Avg. Negative Query = 2



Worst-case: O(M·L)

Avg. worst-case: O(M+L)



Know specs in advance:

- N # entries to insert
- ε desired FPR



Know specs in advance:

N - # entries to insert

ε - desired FPR

Allocate filter with: $N \cdot ln(2) \cdot log_2(1/\epsilon)$ bits



Know specs in advance:

N - # entries to insert

ε - desired FPR

Allocate filter with: $N \cdot ln(2) \cdot log_2(1/\epsilon)$ bits



Insert N elements using -ln(ε)/ln(2) hash functions

Know specs in advance:

N - # entries to insert

ε - desired FPR

Allocate filter with: $N \cdot ln(2) \cdot log_2(1/\epsilon)$ bits

Insert N elements using $-\ln(\epsilon)/\ln(2)$ hash functions

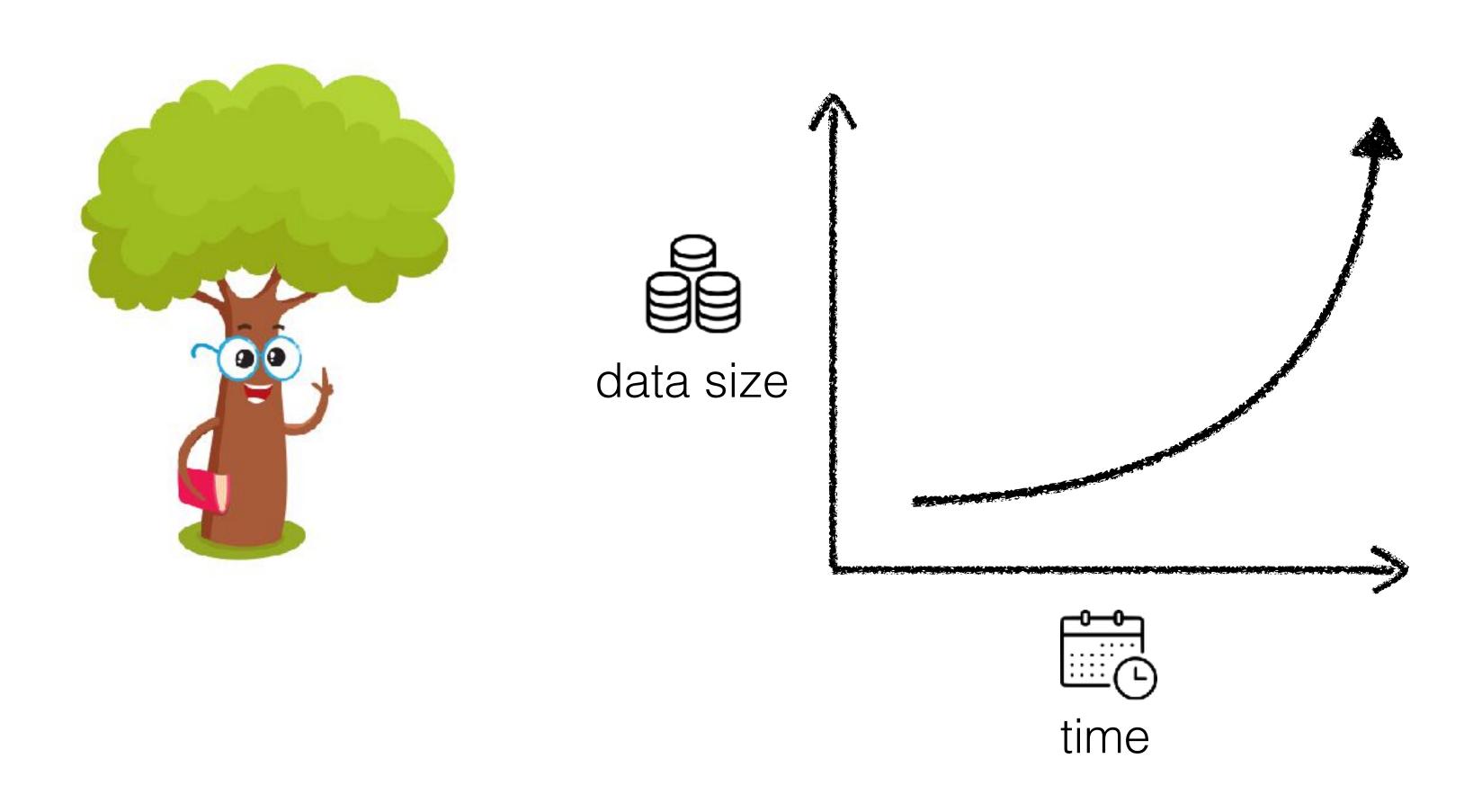
Guarantee FPR of ϵ



Research Question



Can LSM-tree handle exponential data growth?





logarithmic scaling

$$L = O(log N)$$



logarithmic scaling

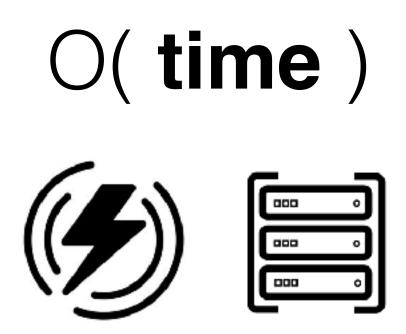
L = O(log N)

exponential growth

 $N \in O(2 \text{ time})$



linear scaling







Can we do better?



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



0(?)

insert I/O cost

 $O((T \cdot L)/B)$



O(?)

 $L = log_T N/P$ (Costs assuming leveling)

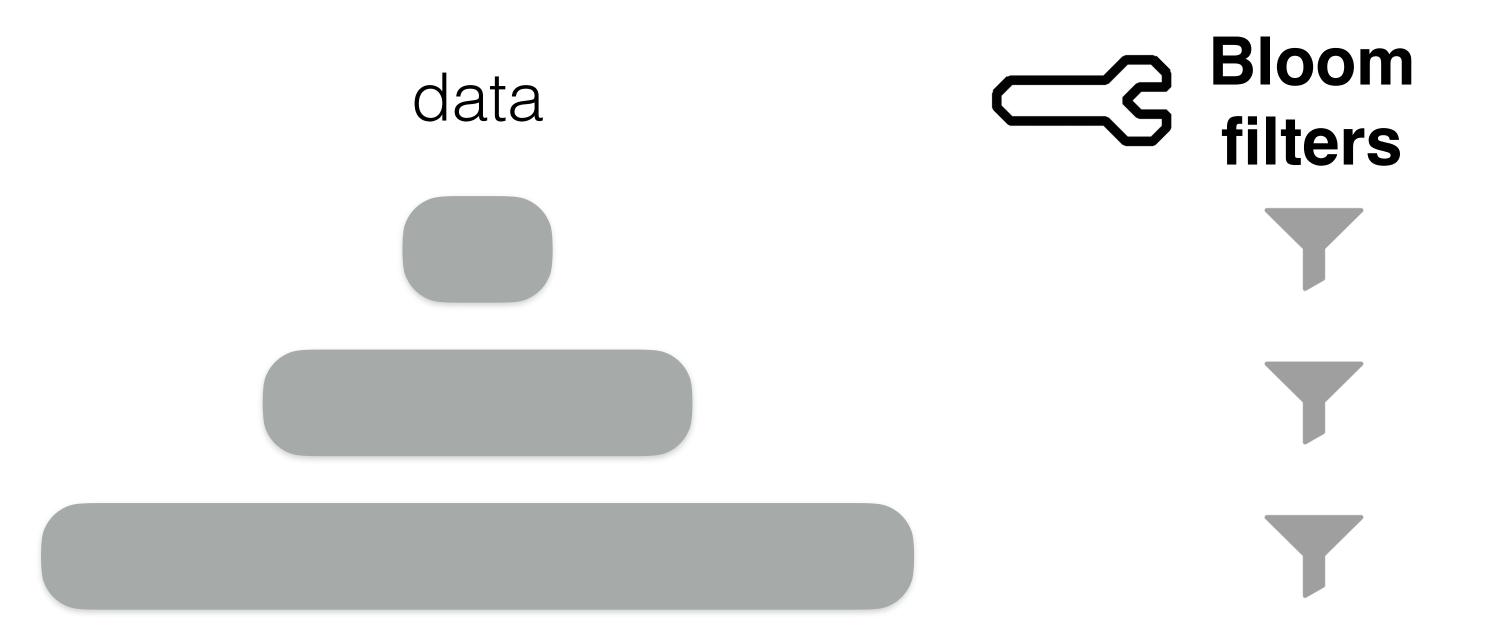
Monkey: Optimal Navigable Key-Value Store

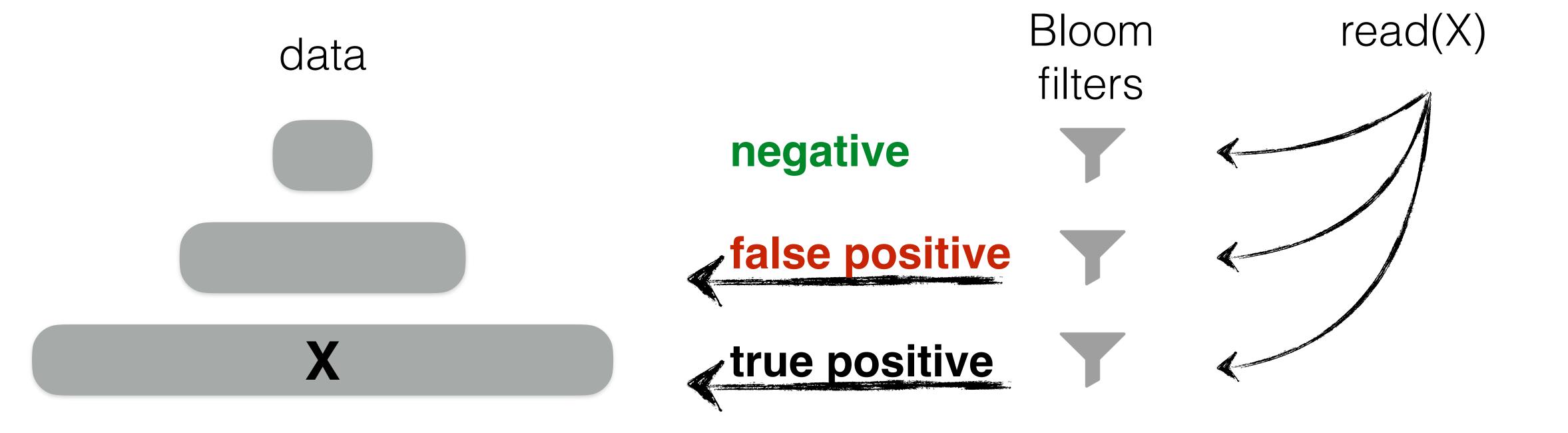


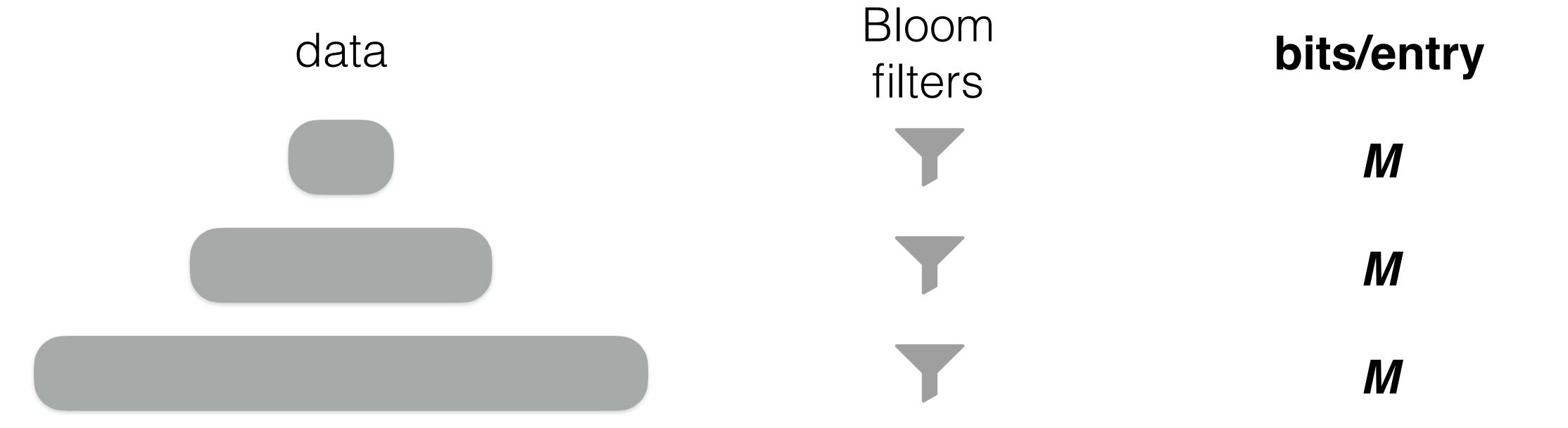
SIGMOD17

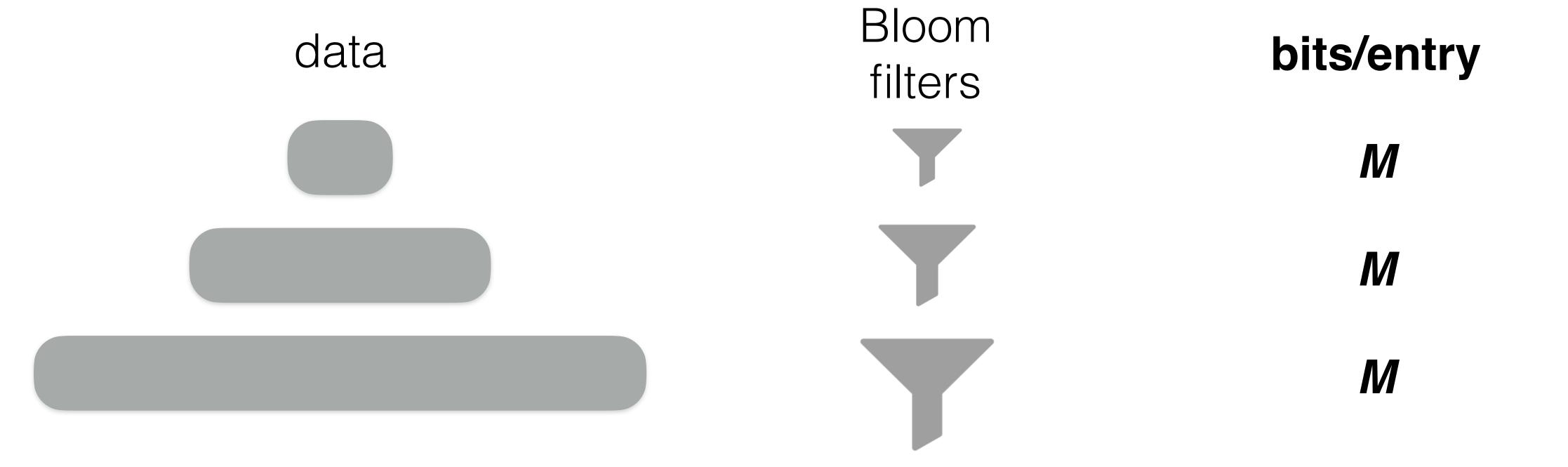
Monkey: Optimal Navigable Key-Value Store

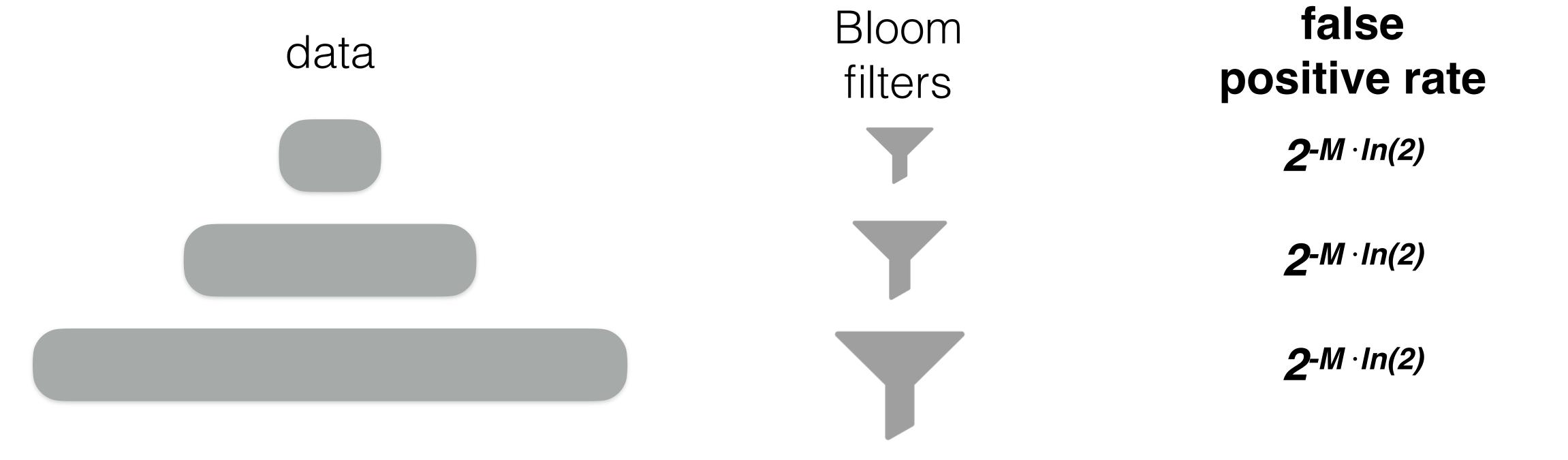
SIGMOD17

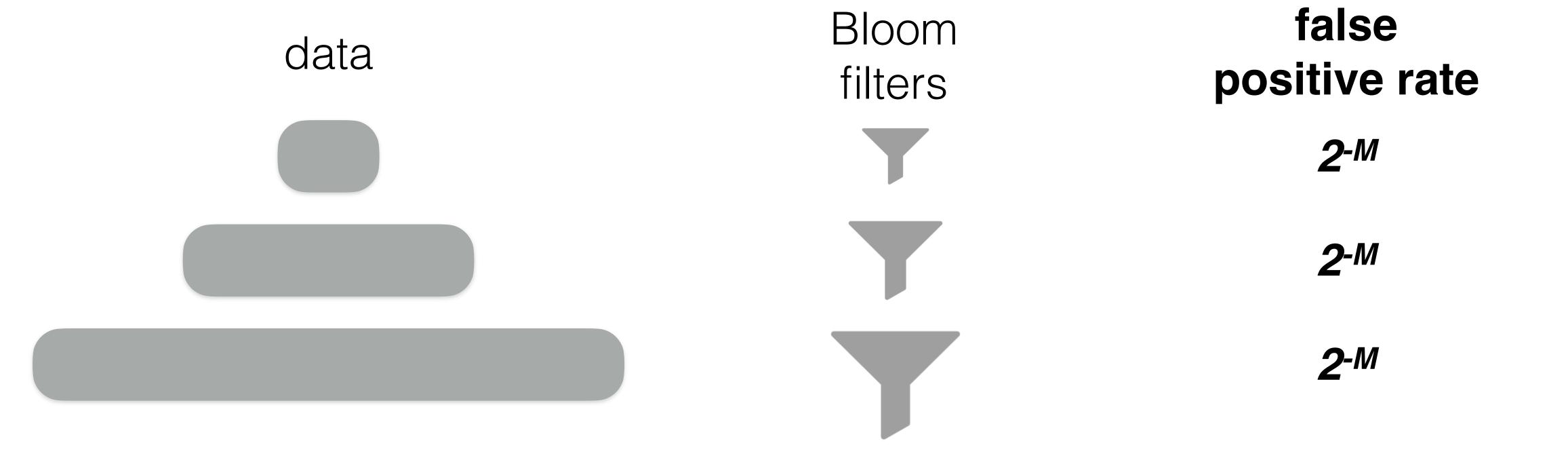


















false positive rate

$$= O(2^{-M} \cdot \log_T N/P)$$









false positive rate

$$= O(2^{-M} \cdot \log_{7} N/P)$$



Bloom





most memory filters







false positive rate

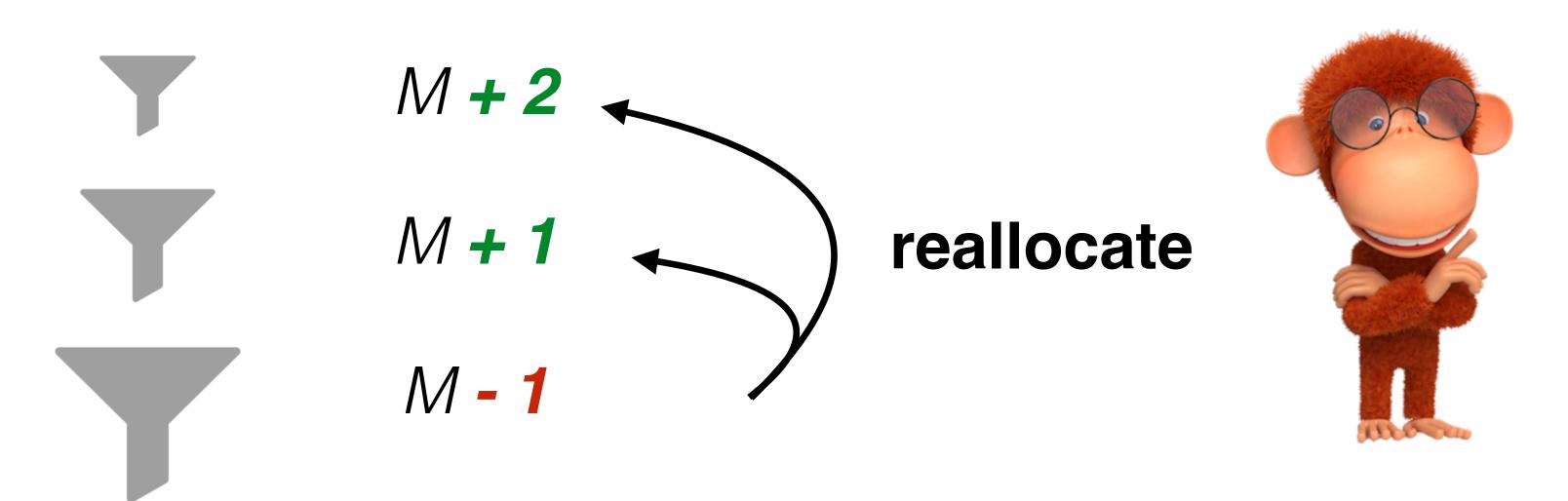
2-M

2-M

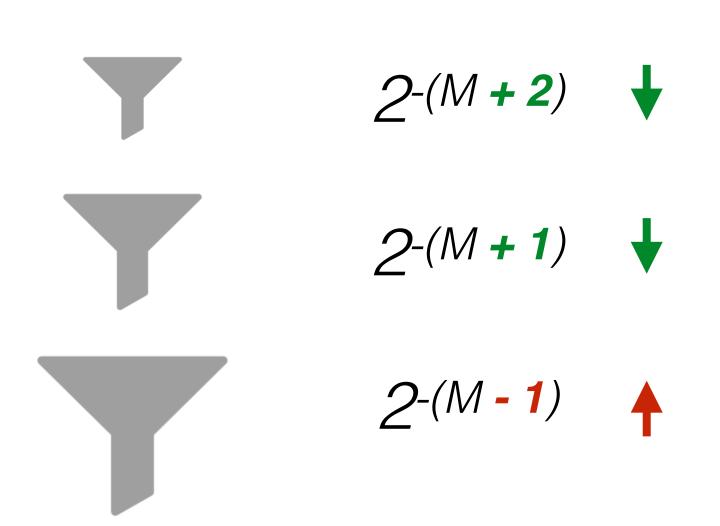
2-M







false positive rates





relax

false positive rates



$$0 < p_0 < 1$$



$$0 < p_1 < 1$$



$$0 < p_2 < 1$$



relax

model

false positive rates

$$read \\ cost = \sum_{1}^{L} p_i$$

$$0 < p_0 < 1$$



$$0 < p_1 < 1$$



$$0 < p_2 < 1$$

memory
$$= -\sum_{i}^{L} \frac{N}{T^{L-i}} \cdot \frac{\ln(p_i)}{\ln(2)^2}$$



relax

false positive rates

$$0 < p_0 < 1$$



$$0 < p_1 < 1$$

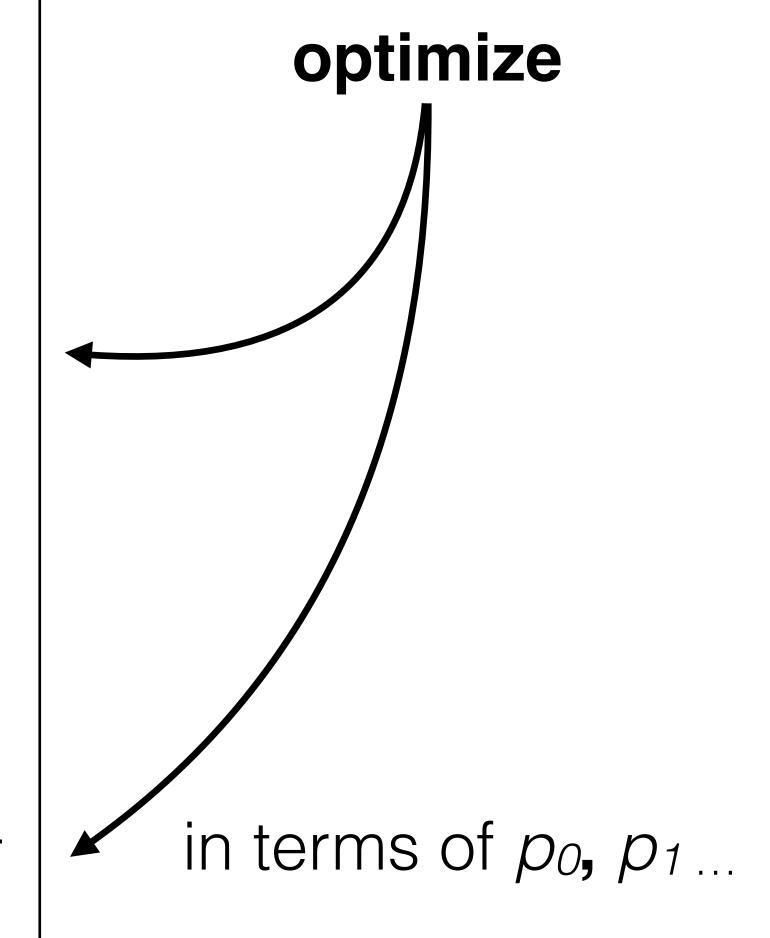


$$0 < p_2 < 1$$

model

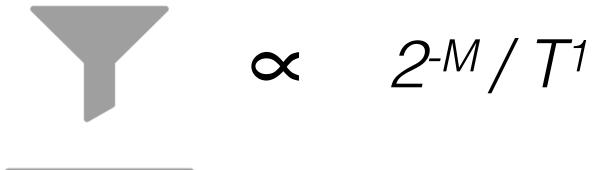
$$read \\ cost = \sum_{1}^{L} p_{i}$$

memory
$$= -\sum_{i}^{L} \frac{N}{T^{L-i}} \cdot \frac{\ln(p_i)}{\ln(2)^2}$$

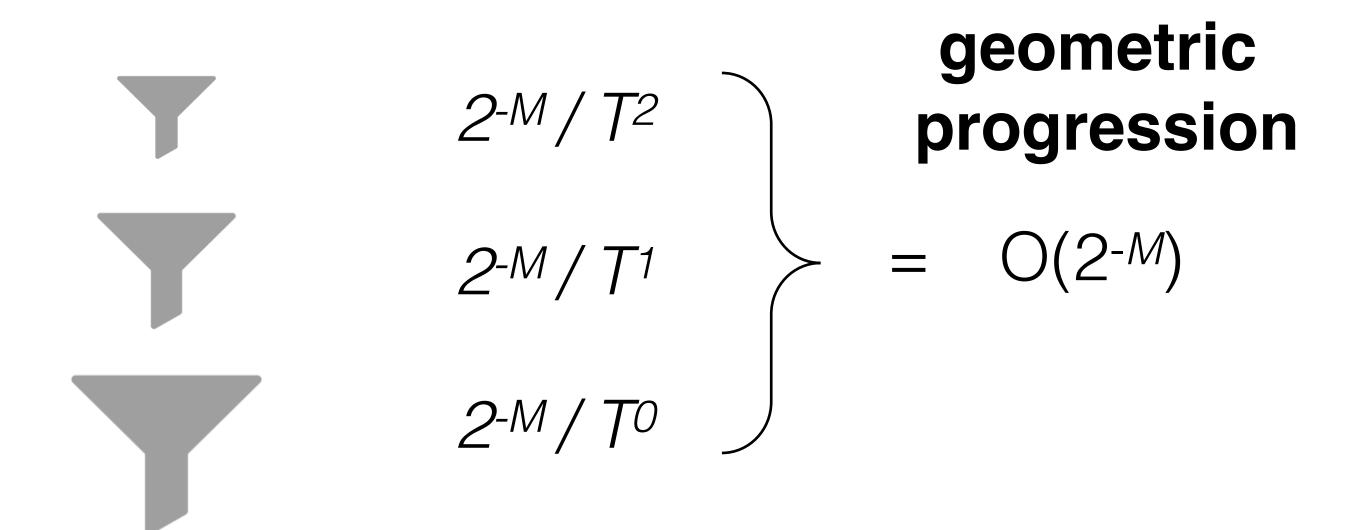














Faster worst case

$$O(2^{-M}) < O(2^{-M} \log_T N/P)$$

Configuration

buffer 2MB

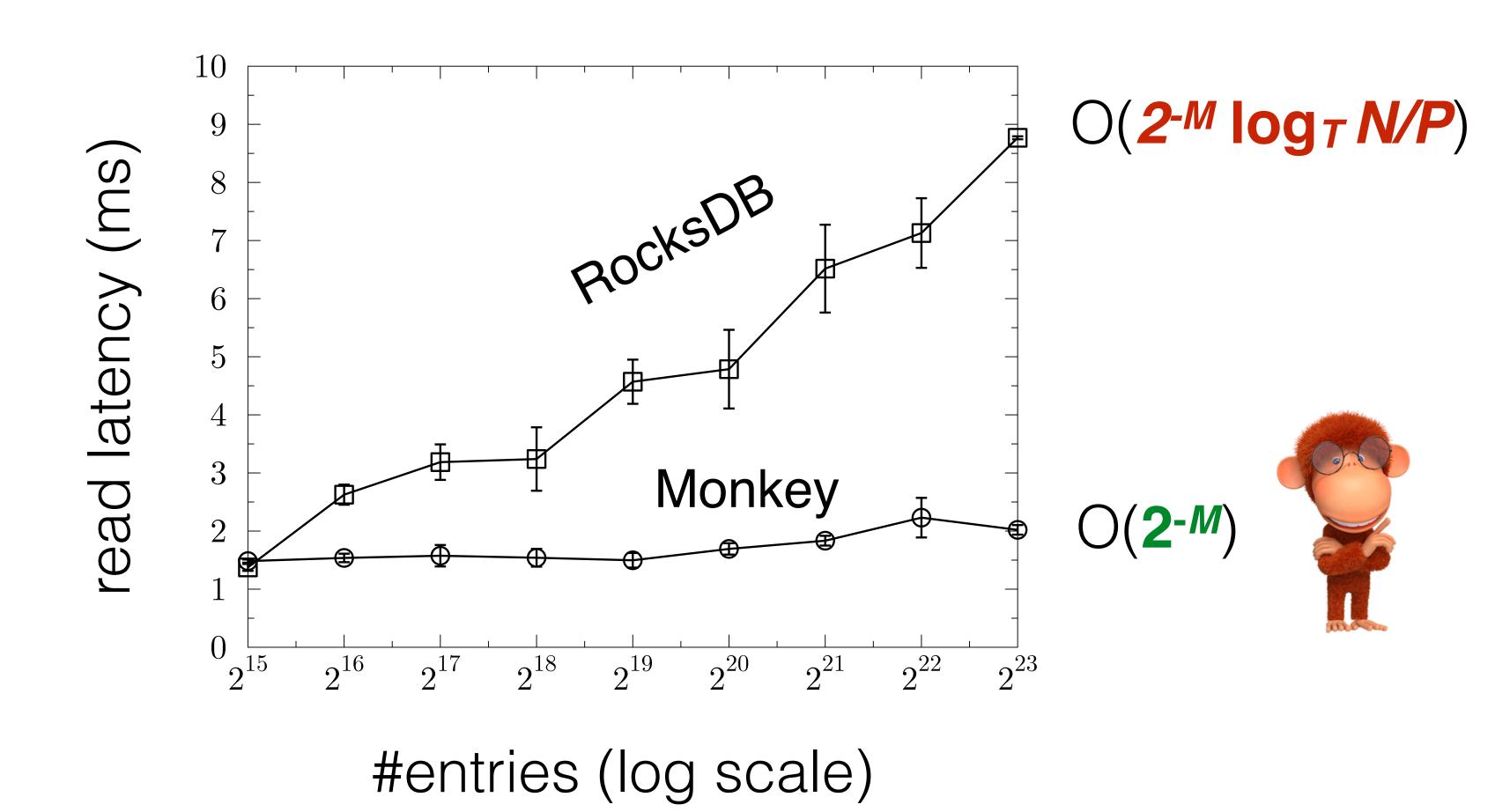
bits/entry: 5

size ratio: 2

1KB entries

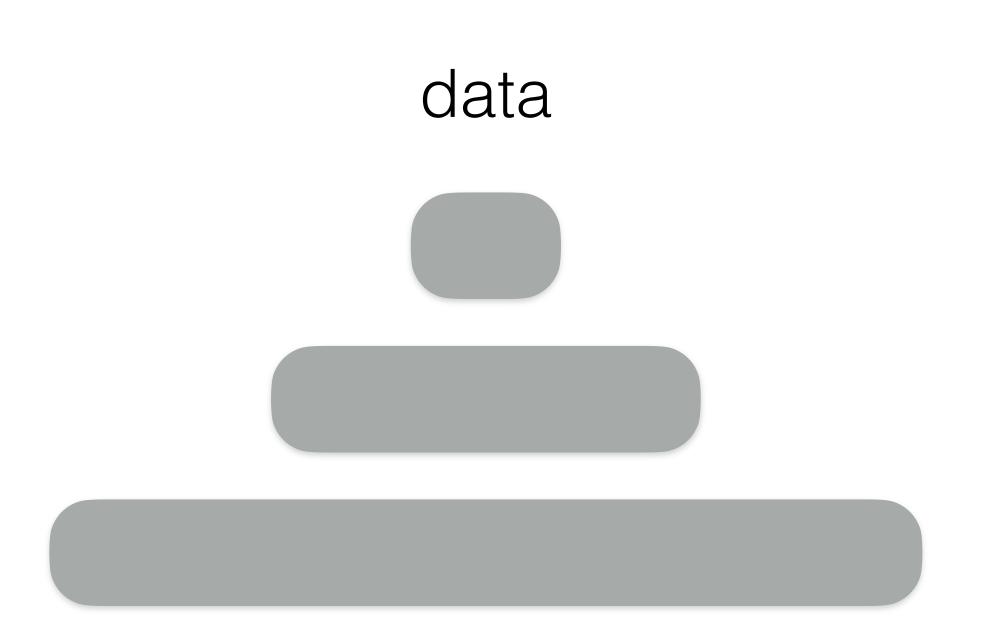
queries to missing keys

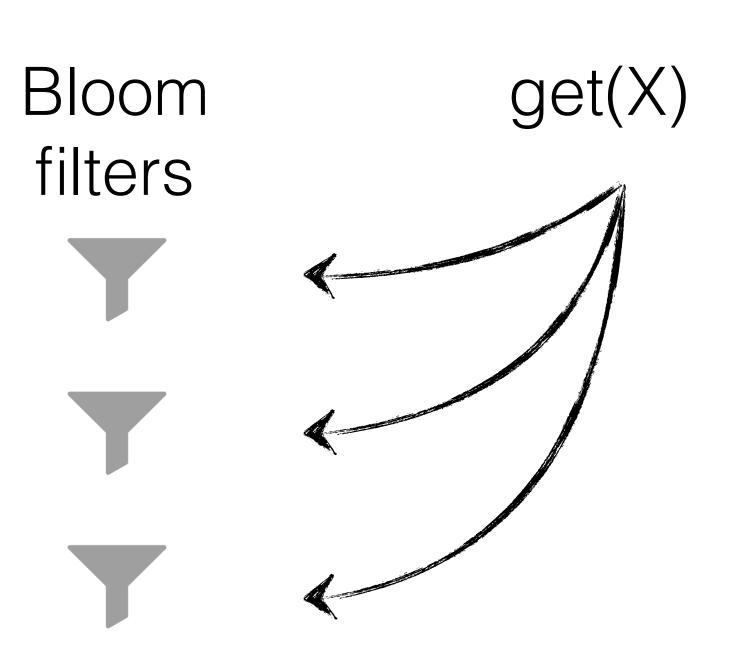
hard disk storage



Positive Query = $M \cdot ln(2)$

Avg. Negative Query = 2

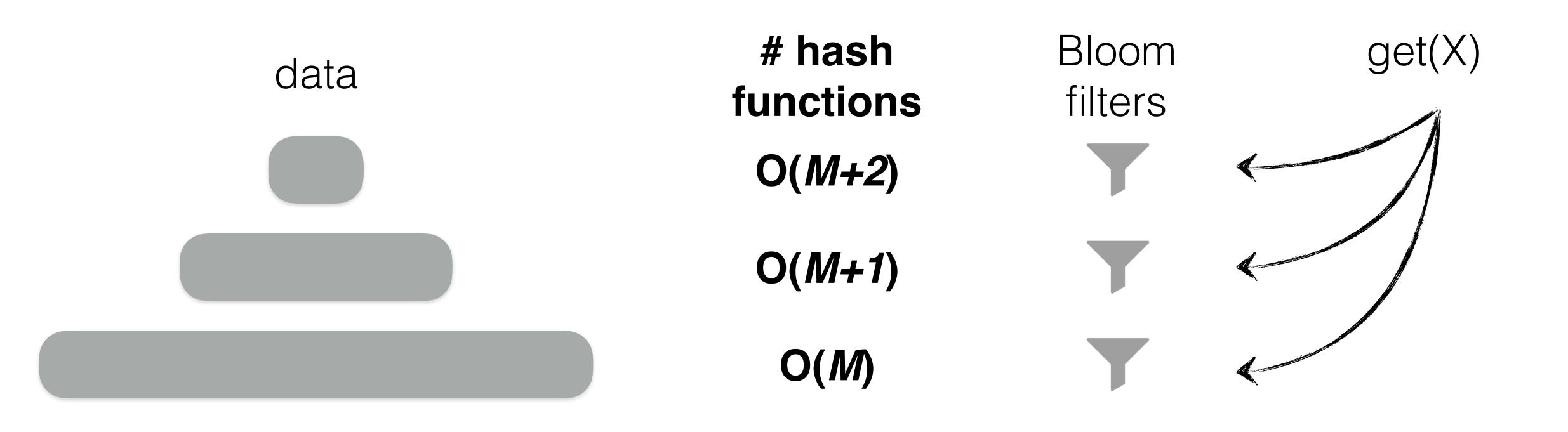




Worst-case:

Positive Query = $M \cdot ln(2)$

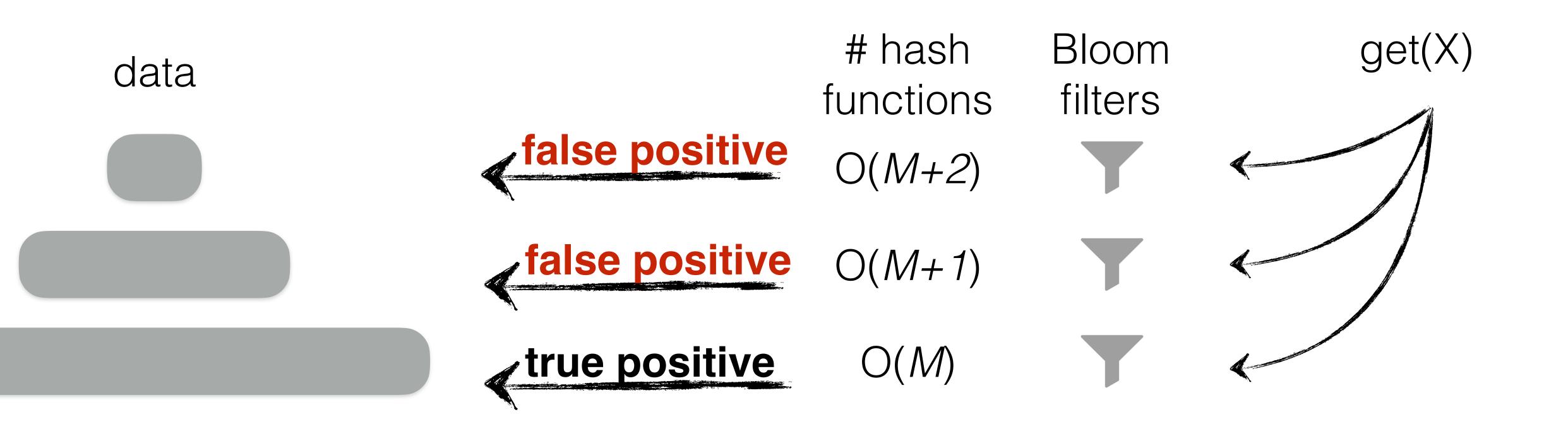
Avg. Negative Query = 2



Worst-case:

Positive Query = $M \cdot ln(2)$

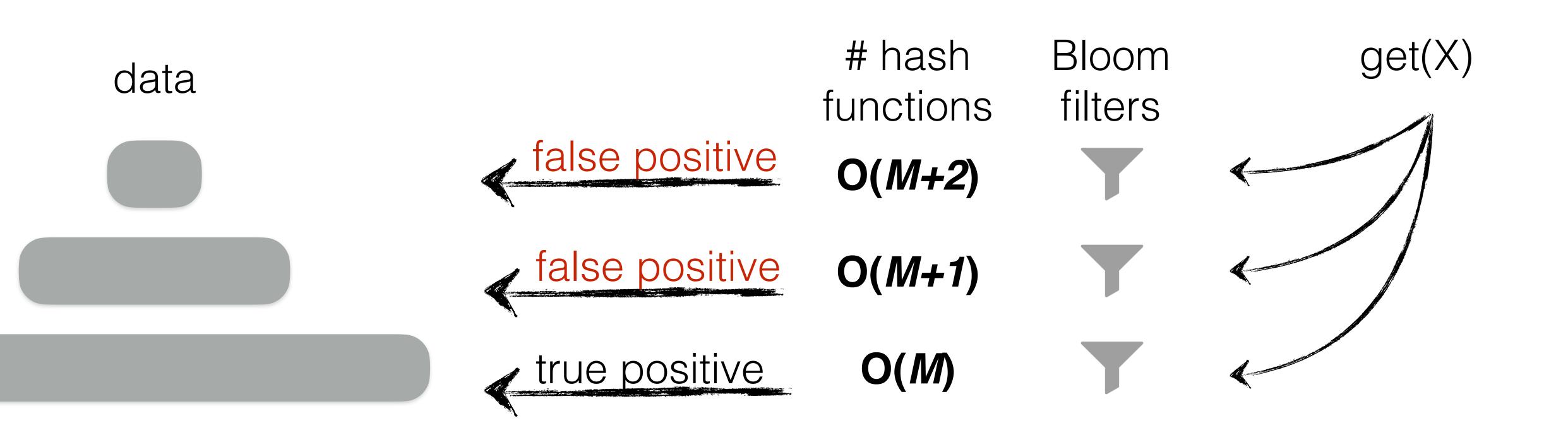
Avg. Negative Query = 2



Worst-case:

Positive Query = $M \cdot ln(2)$

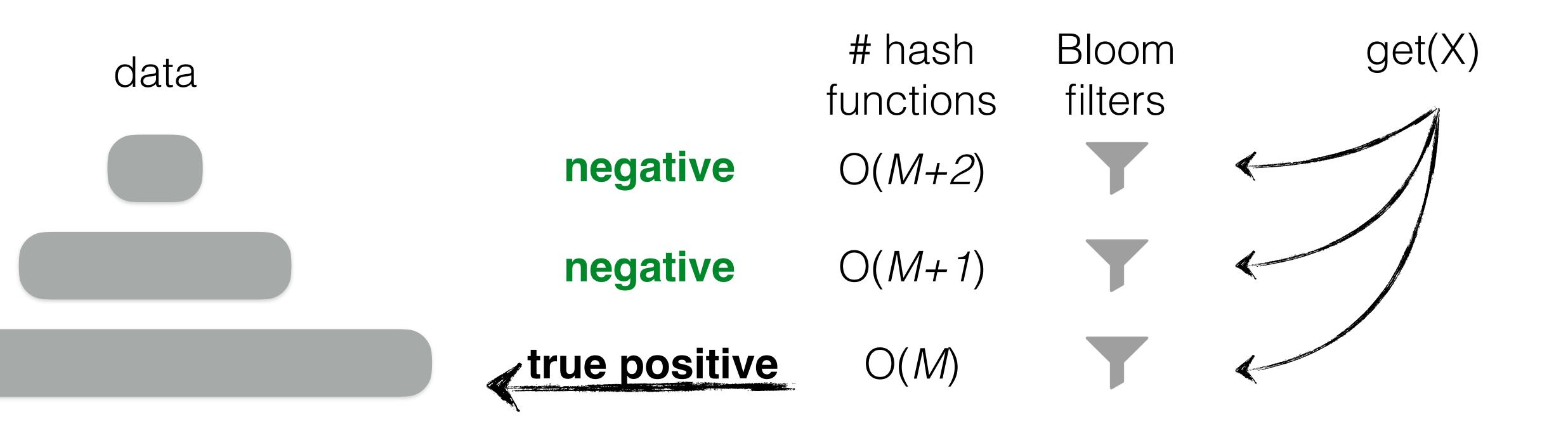
Avg. Negative Query = 2



Worst-case: O(M·L+L²)

Positive Query = $M \cdot ln(2)$

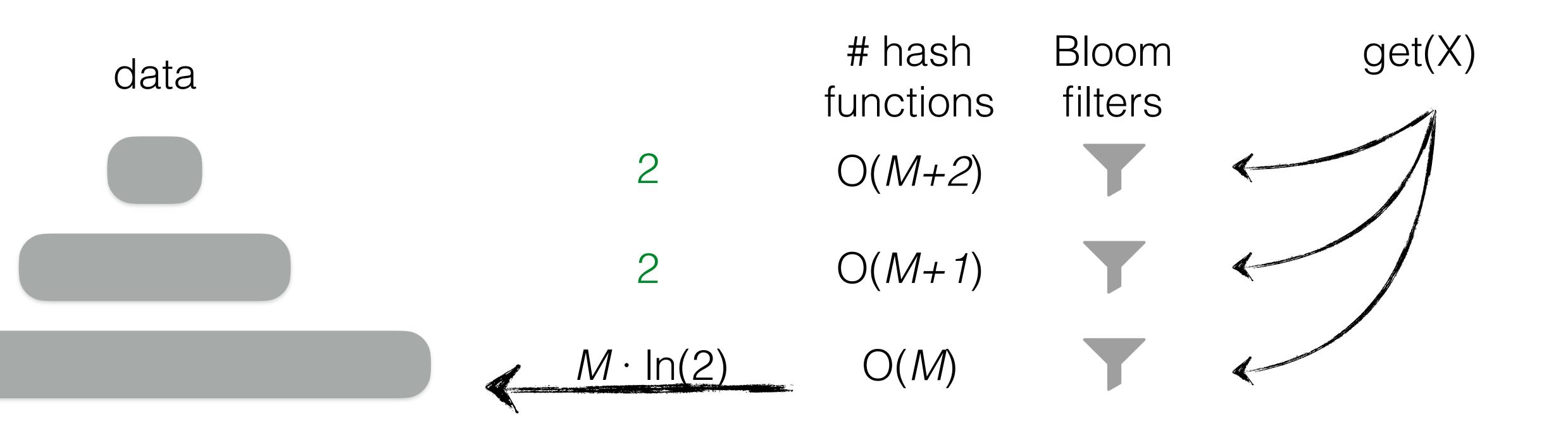
Avg. Negative Query = 2



Worst-case: $O(M \cdot L + L^2)$ Avg. worst-case:

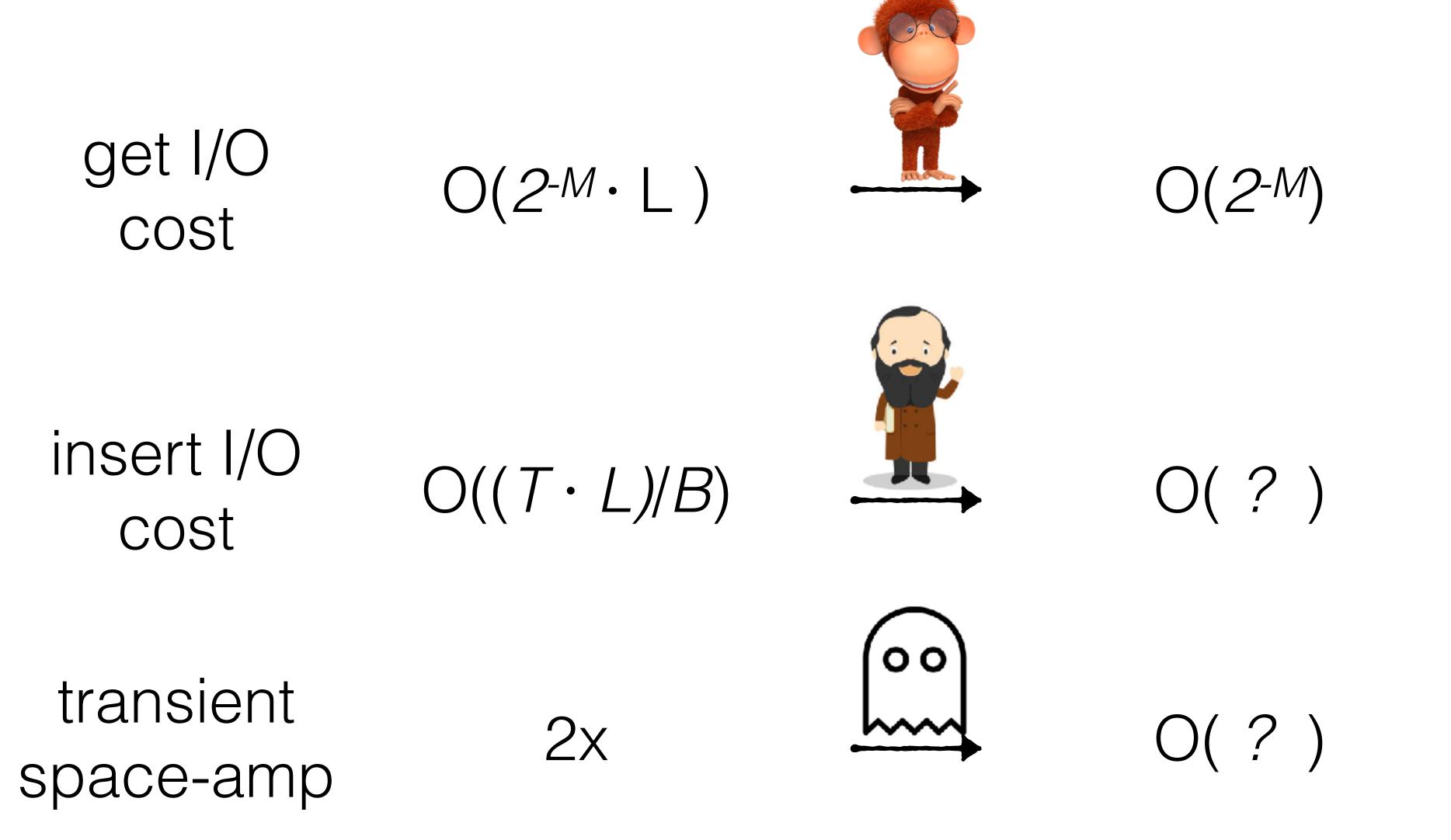
Positive Query = $M \cdot ln(2)$

Avg. Negative Query = 2



Worst-case: O(M·L+L²)

Avg. worst-case: O(M+L)



 $L = log_T N/P$ (Costs assuming leveling)



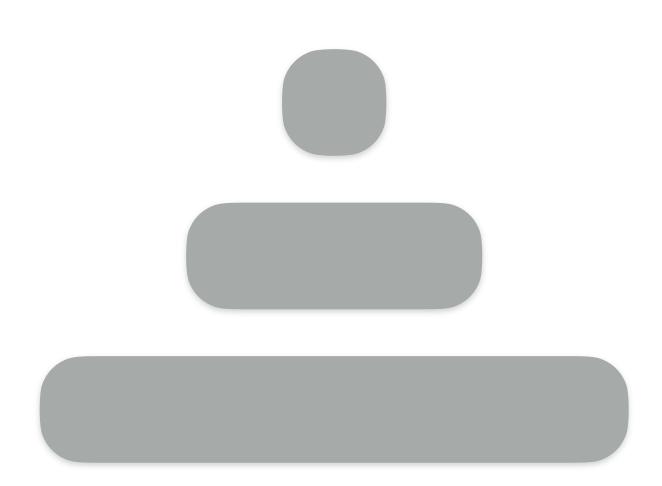
Dostoevsky

SIGMOD18

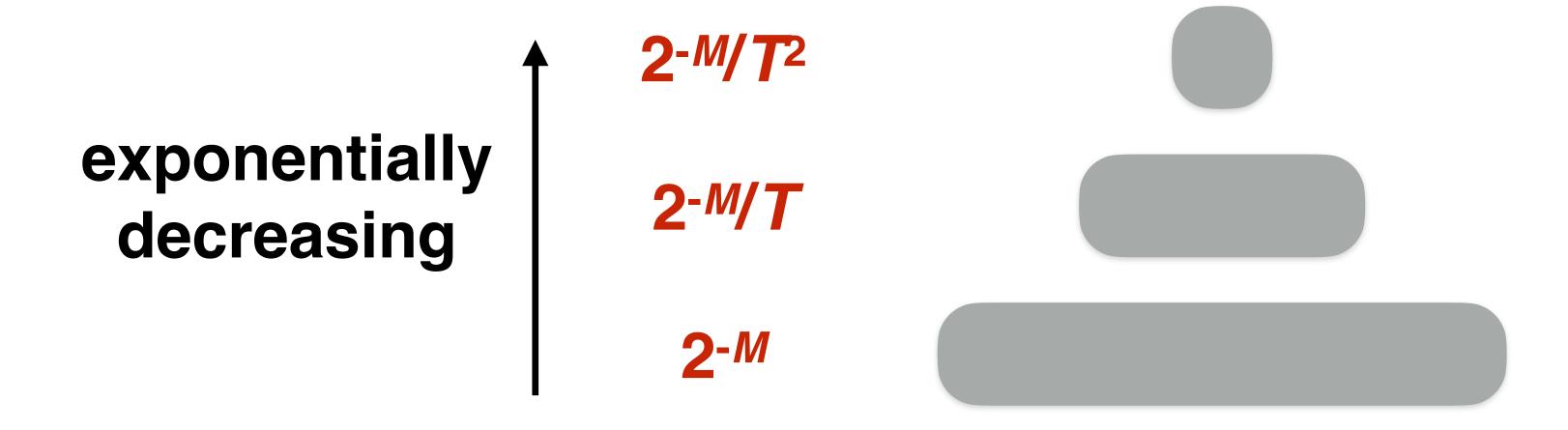


Dostoevsky: Space-Time Optimized Evolvable Scalable Key-Value Store

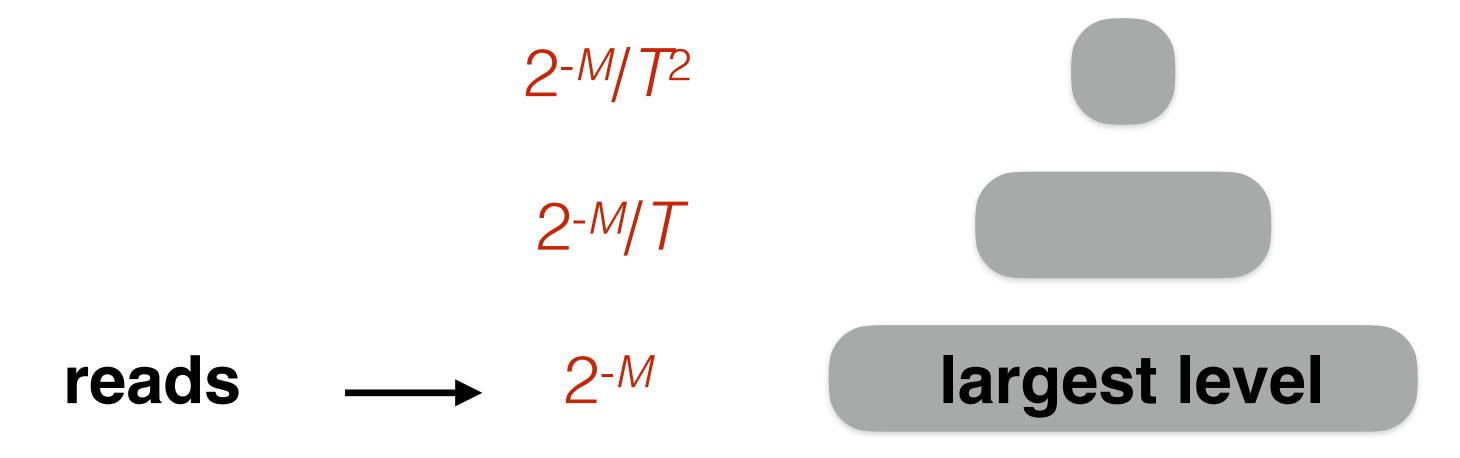
reads & writes cost breakdown



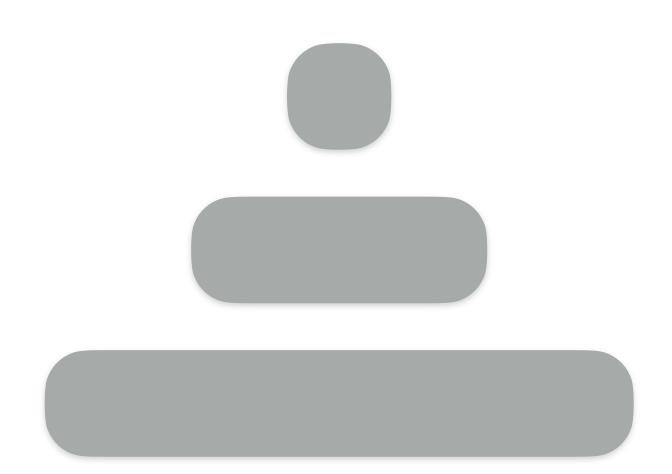
false positive rates

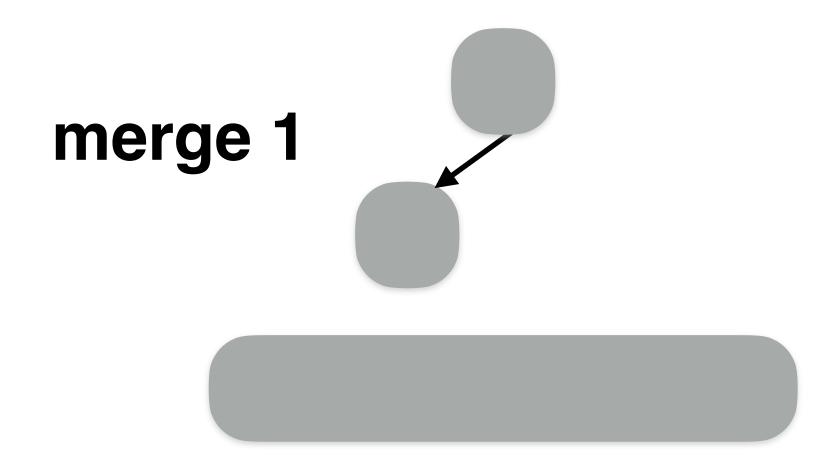


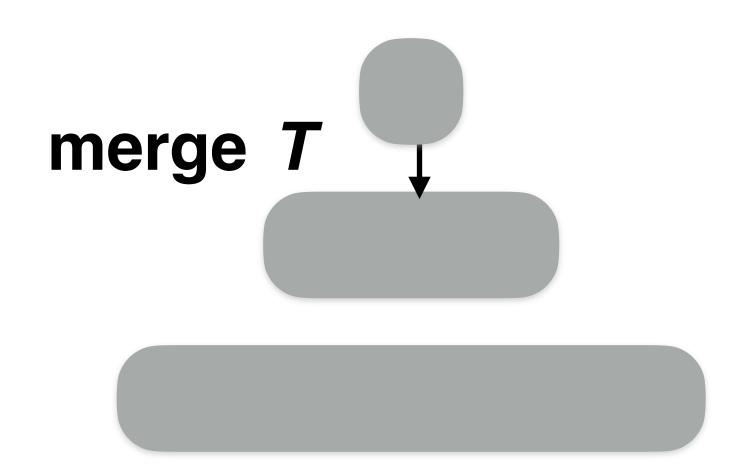
false positive rates

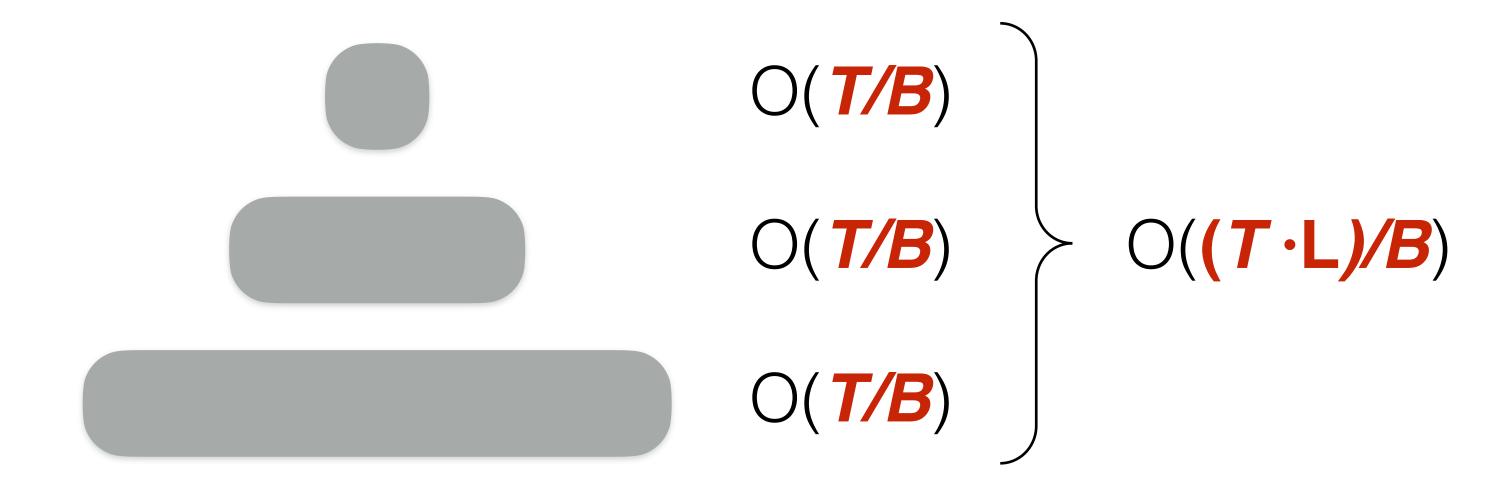


reads O(2-M) writes





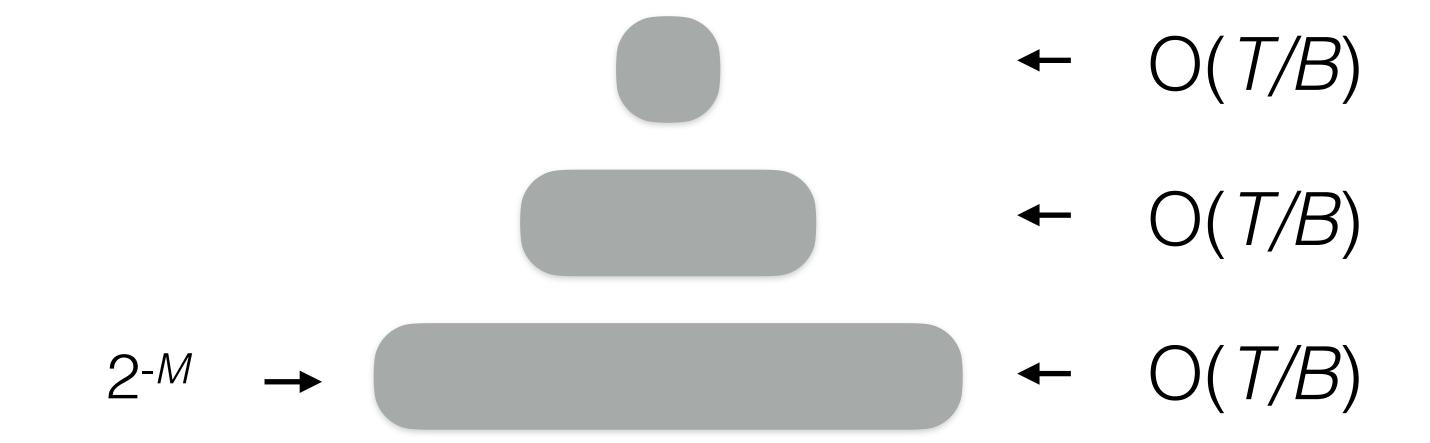


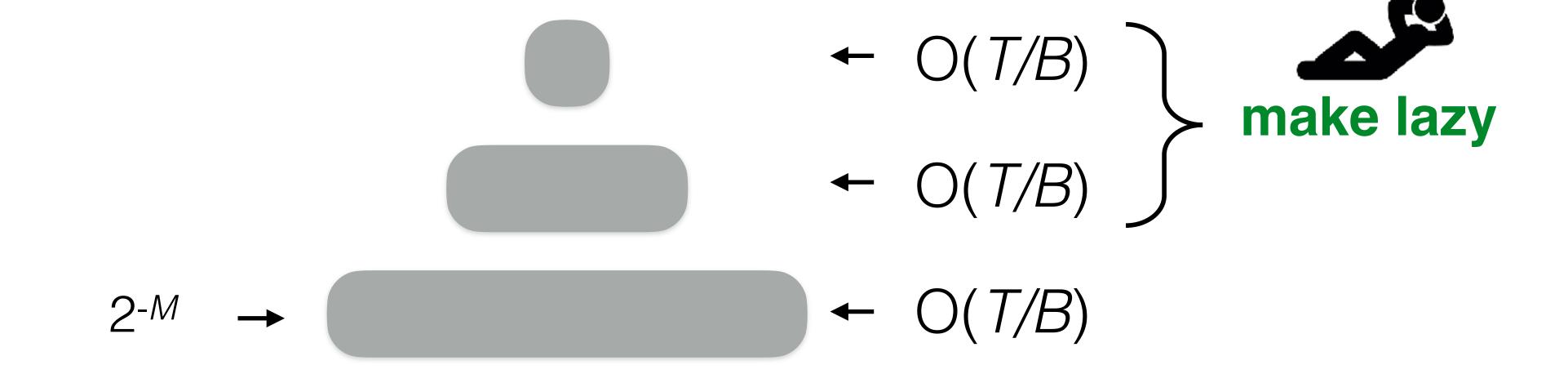


$$O(2^{-M})$$
 $O((T \cdot L)/B)$
=
=
 $2^{-M}/T^2$ $O(T/B)$
+
 $2^{-M}/T$ $O(T/B)$
+
 $O(T/B)$

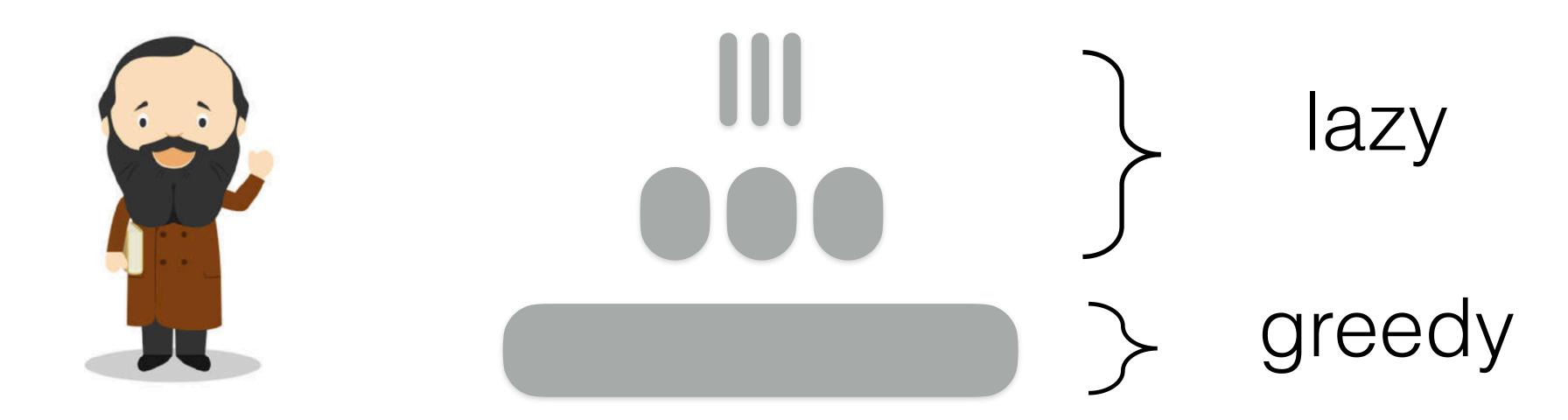
reads
largest level

writes
all levels

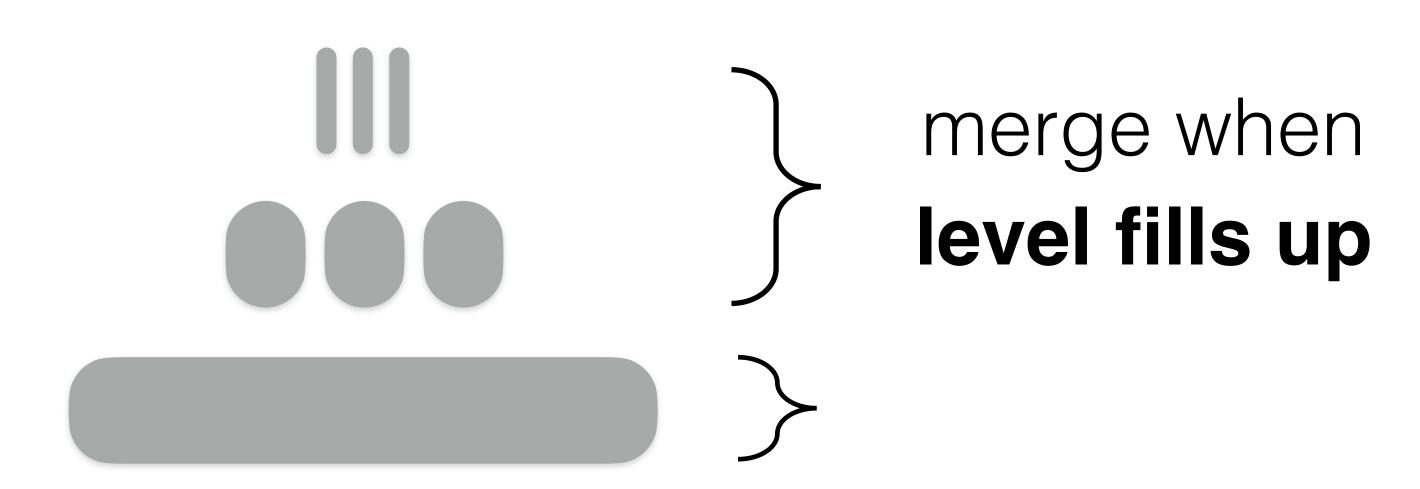




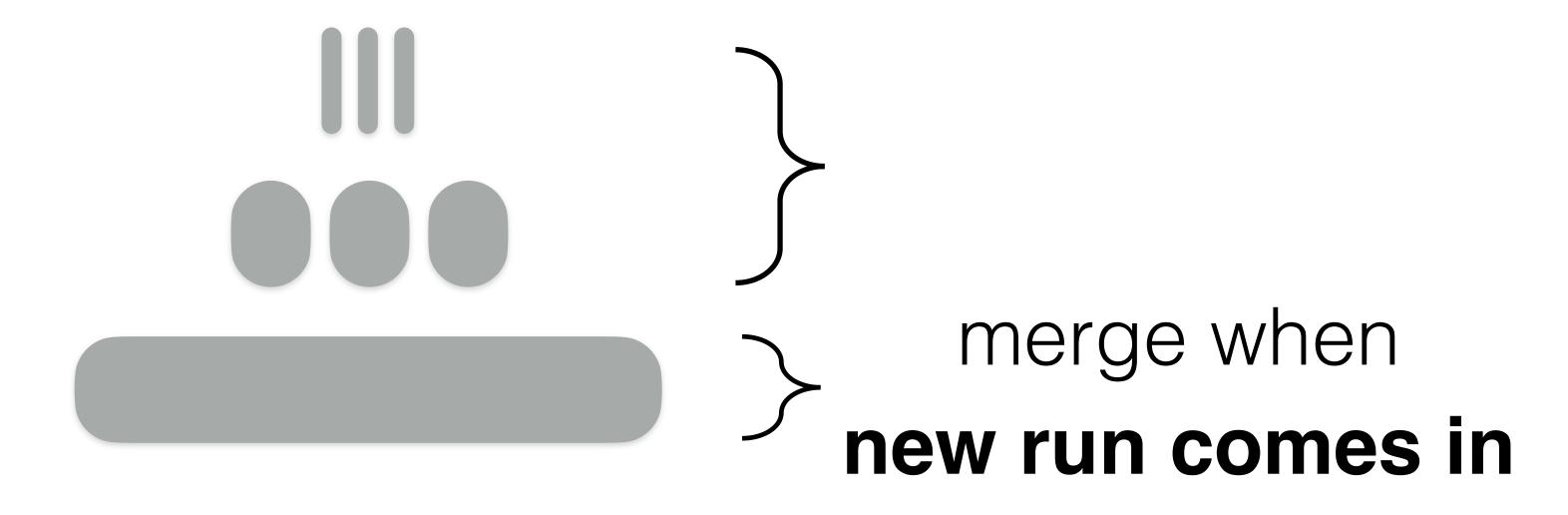
Dostoevsky

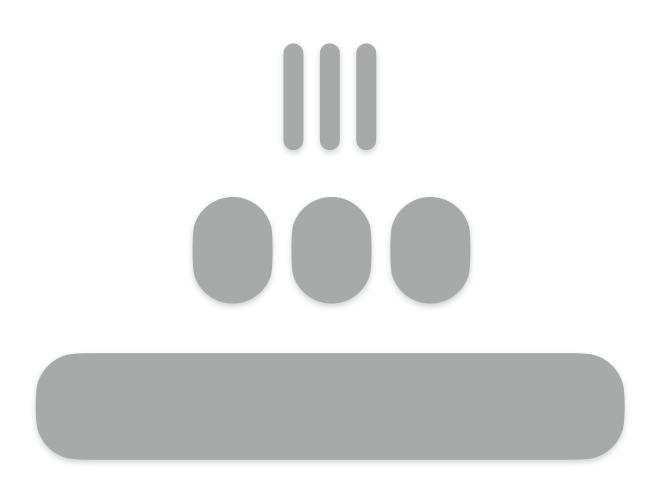


Dostoevsky

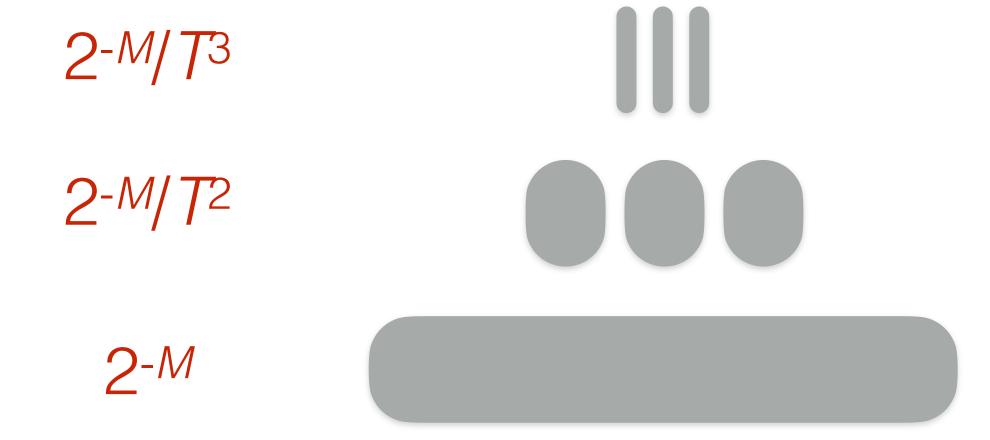


Dostoevsky

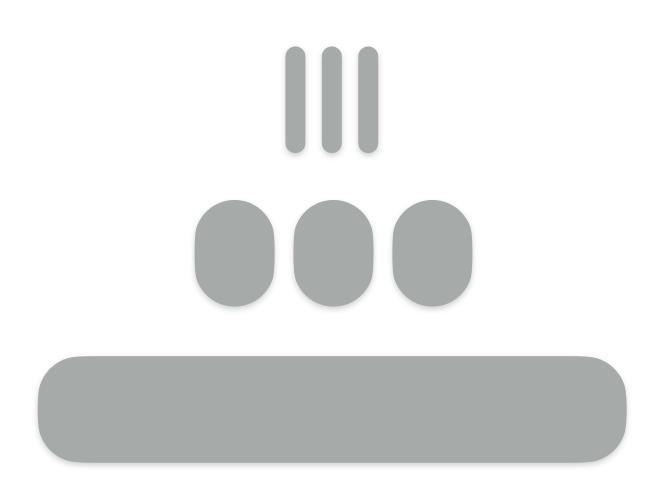




false positive rates



 $O(2^{-M})$



 $O(2^{-M})$

writes

O(1/B)

O(1/B)

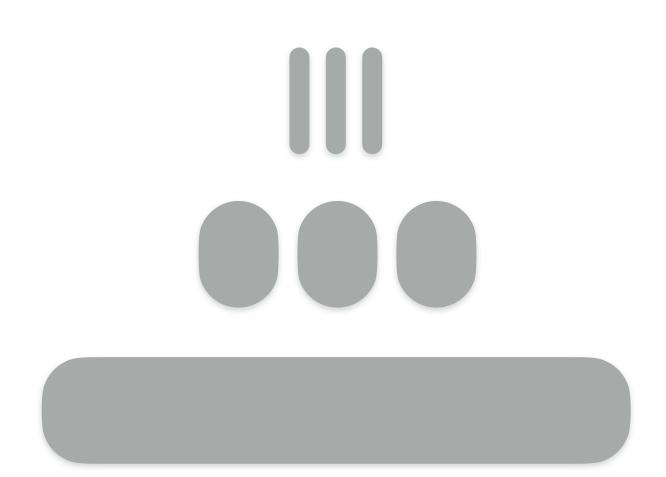
O(*T/B*)

 $O(2^{-M})$

 $O(2^{-M})$

writes

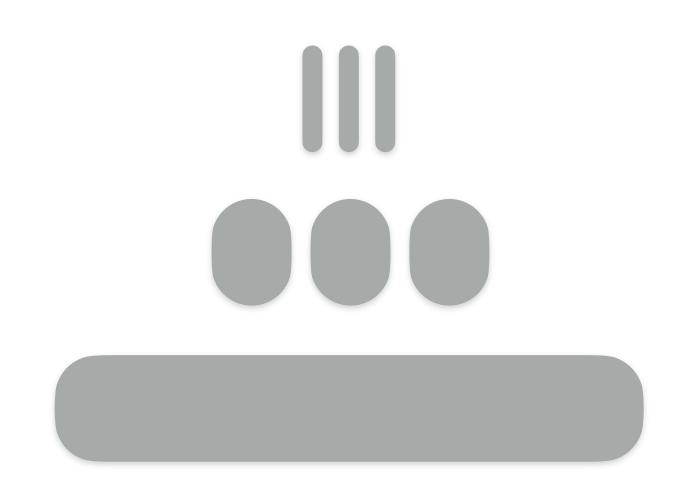
O((T + L)/B)



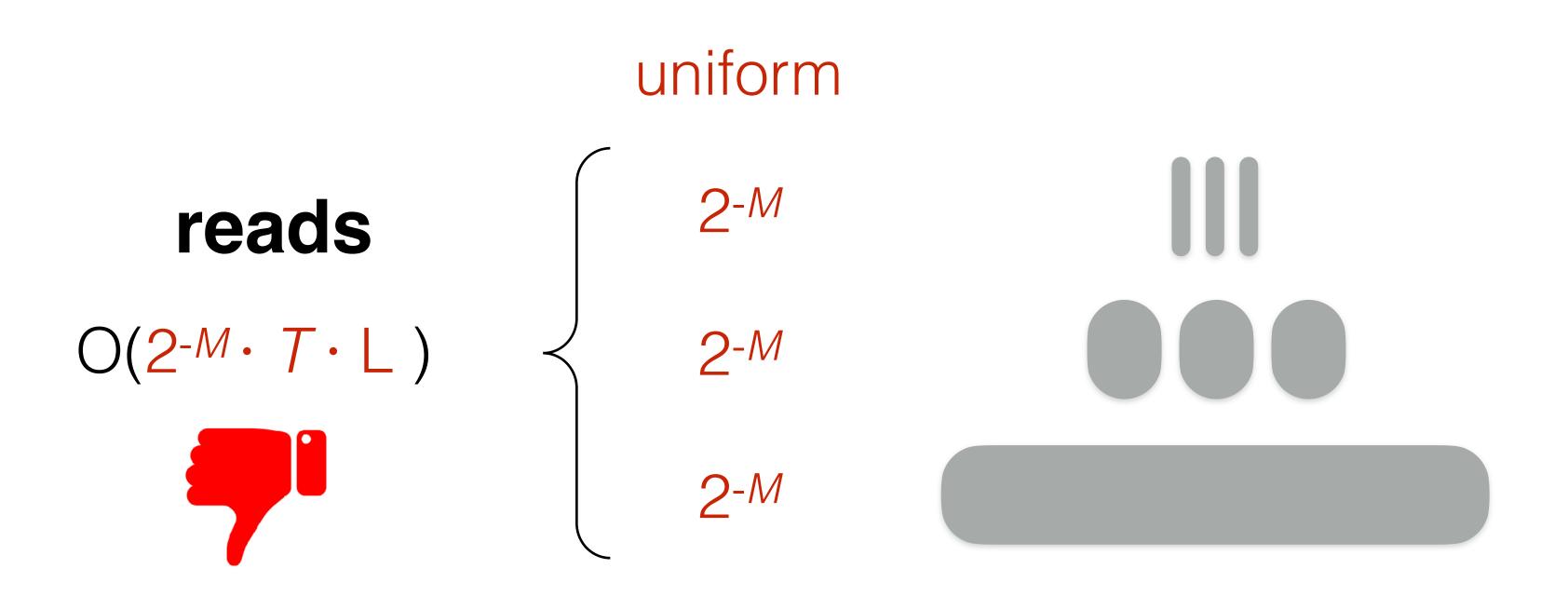
What would our read cost have been if we employed uniform FPRs at all levels?

reads writes

 $O(2^{-M})$ O((T + L)/B)

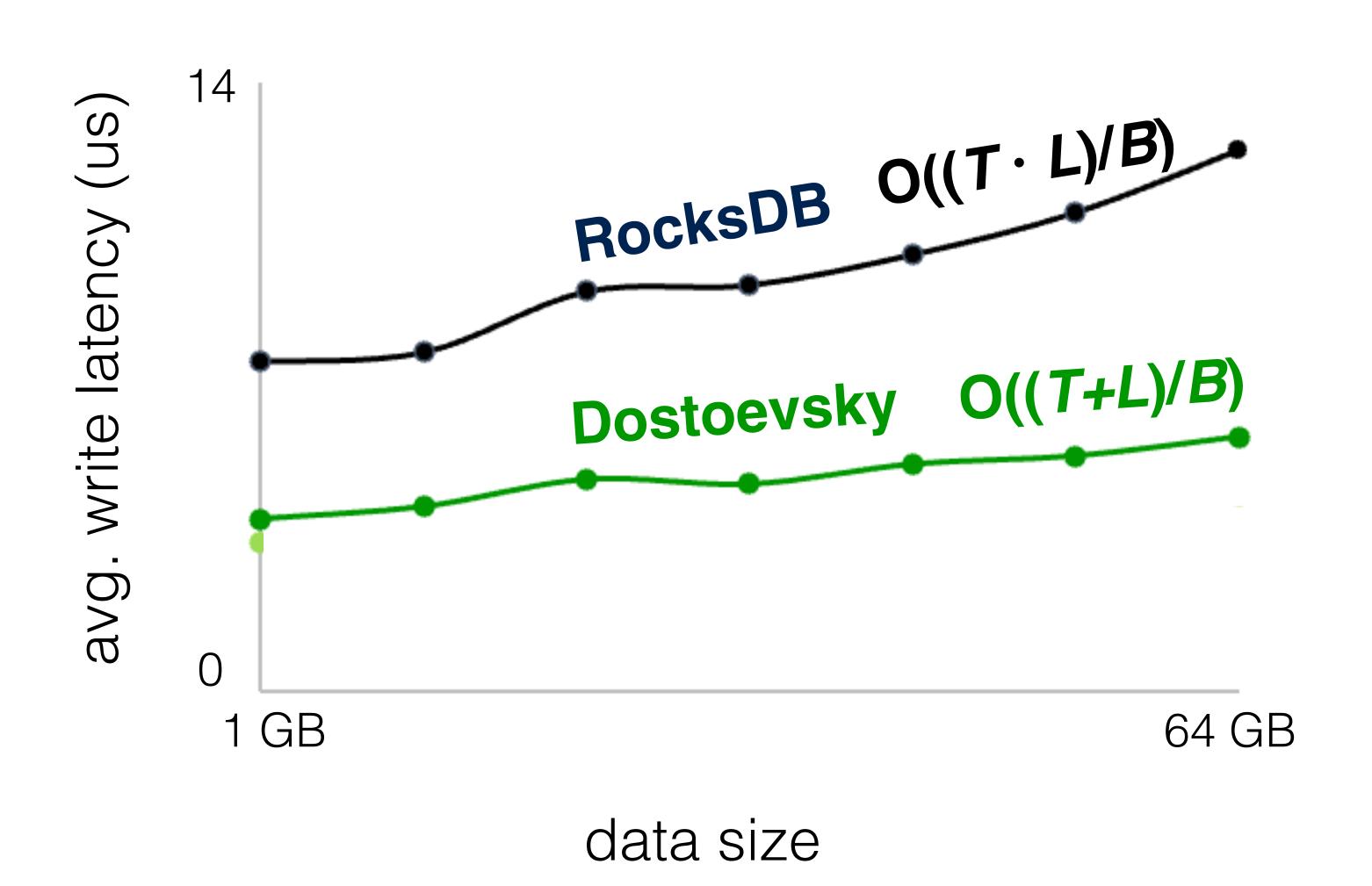


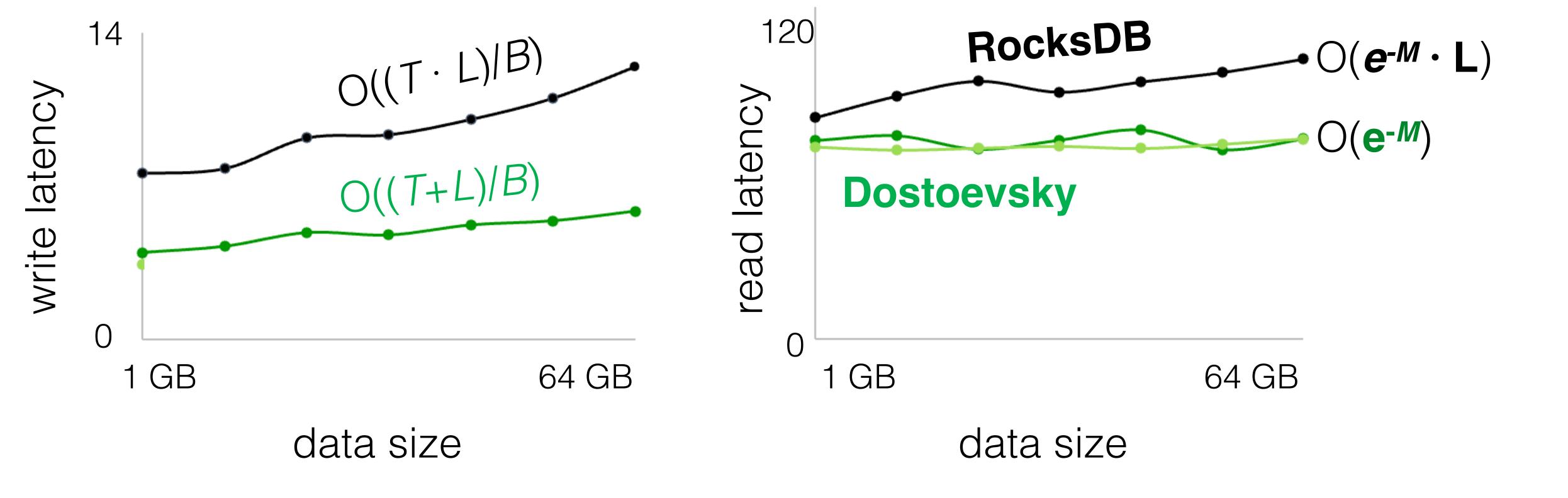
What would our read cost have been if we employed uniform FPRs at all levels?

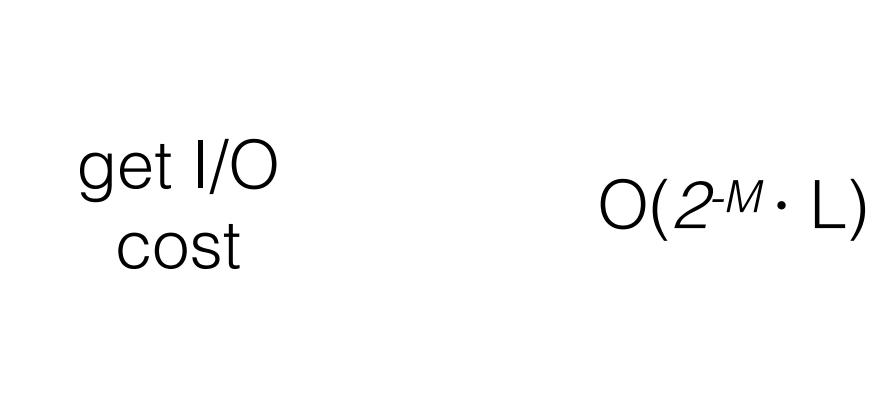


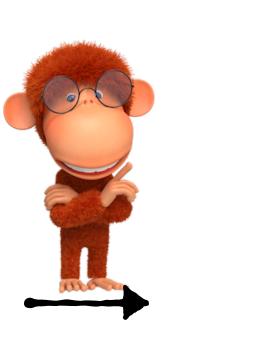
Configuration

buffer 2MBsize ratio: 51KB entriesSSD storage









 $O(2^{-M})$

insert I/O cost

 $O((T \cdot L)/B)$



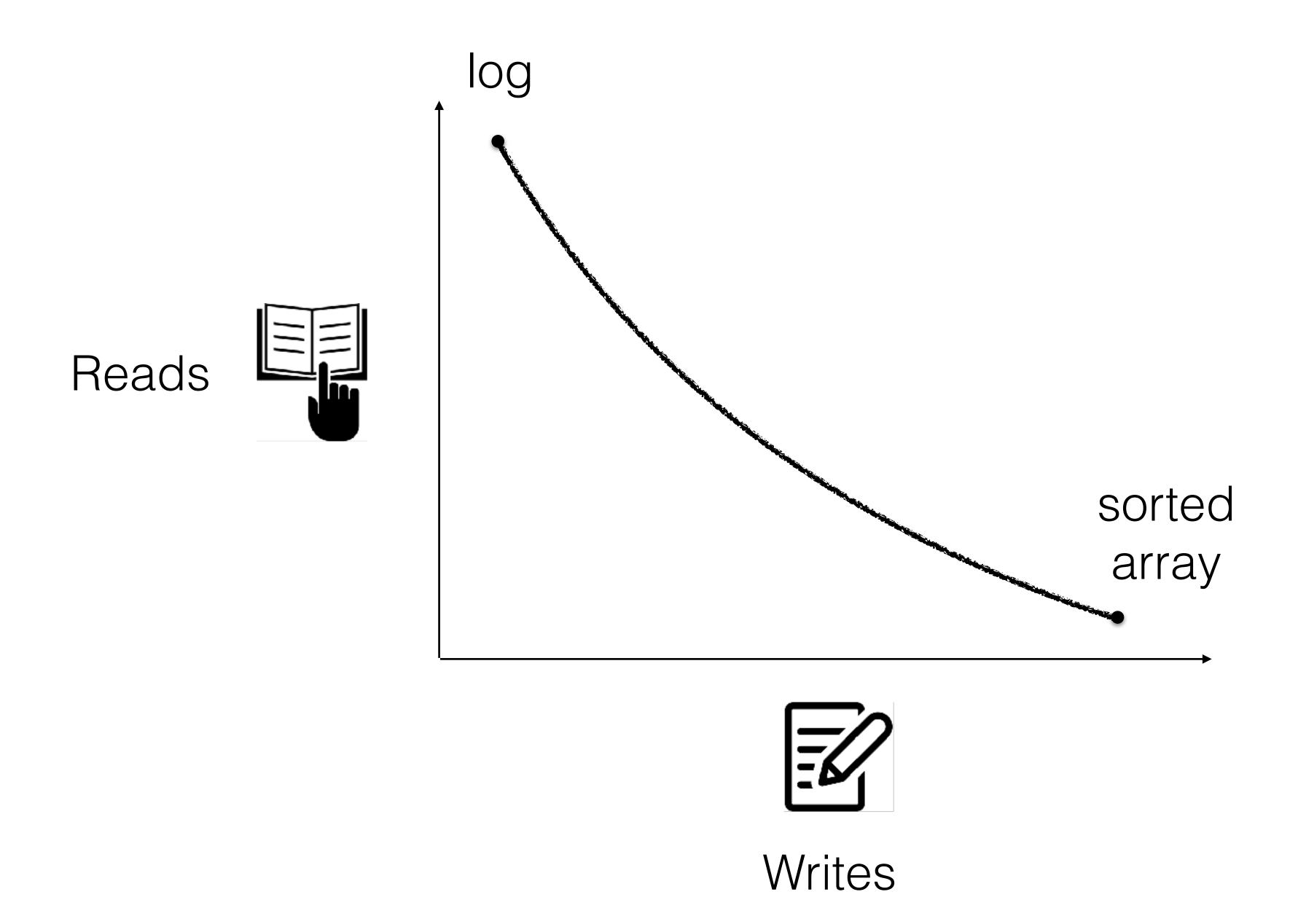
O((*T*+L)/B)

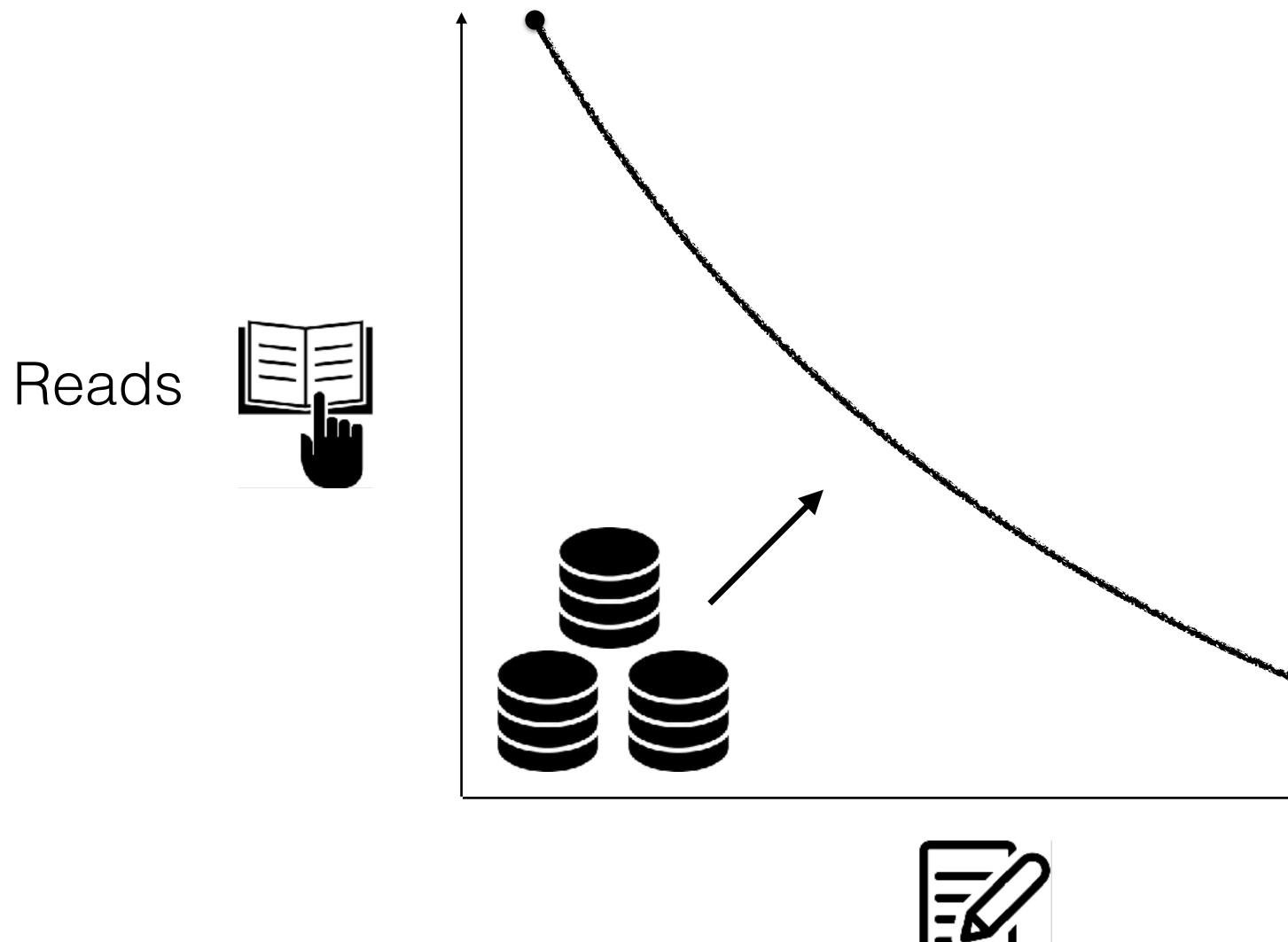
 $L = log_T N/P$ (Costs assuming leveling)

Better scalability with data growth





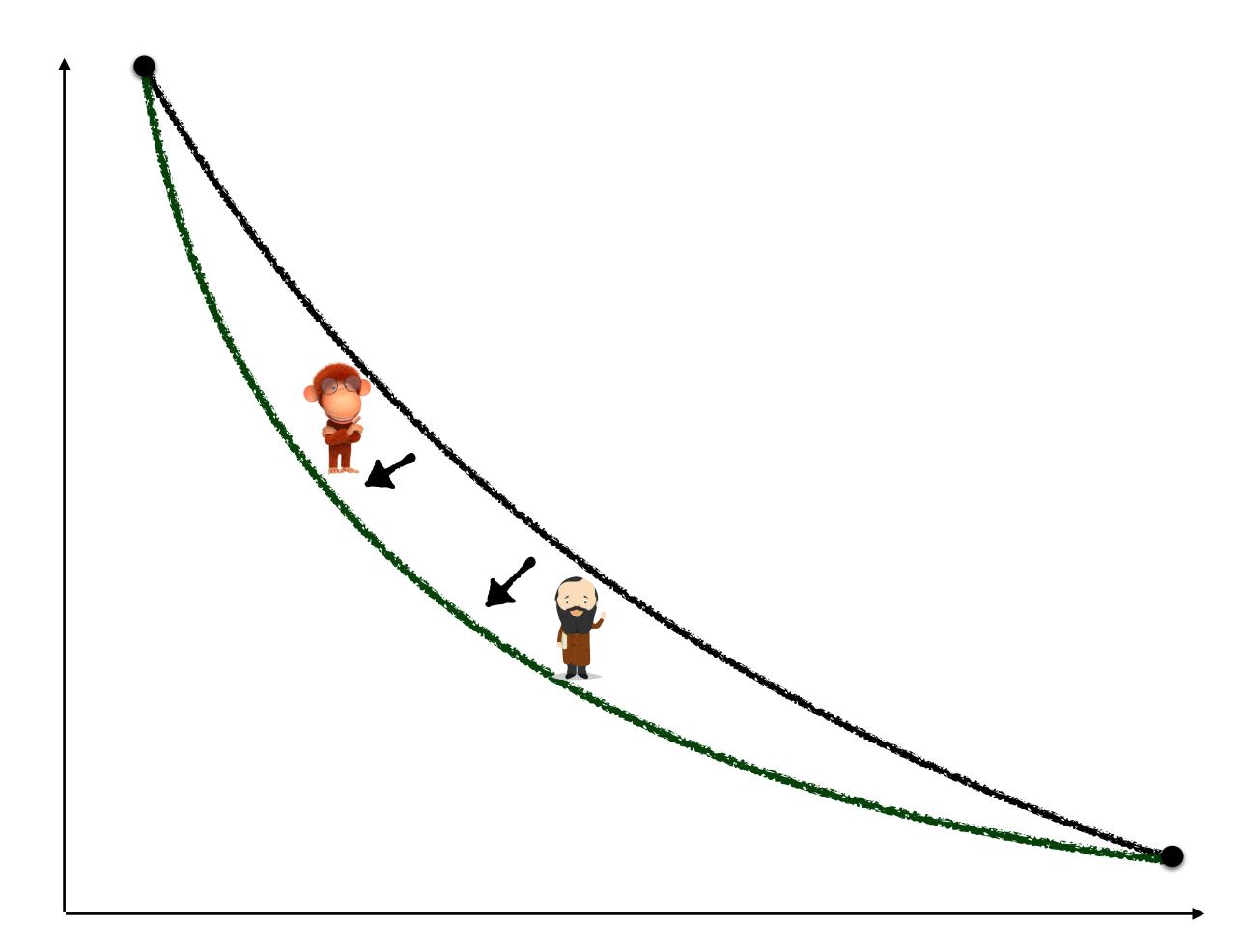






Writes

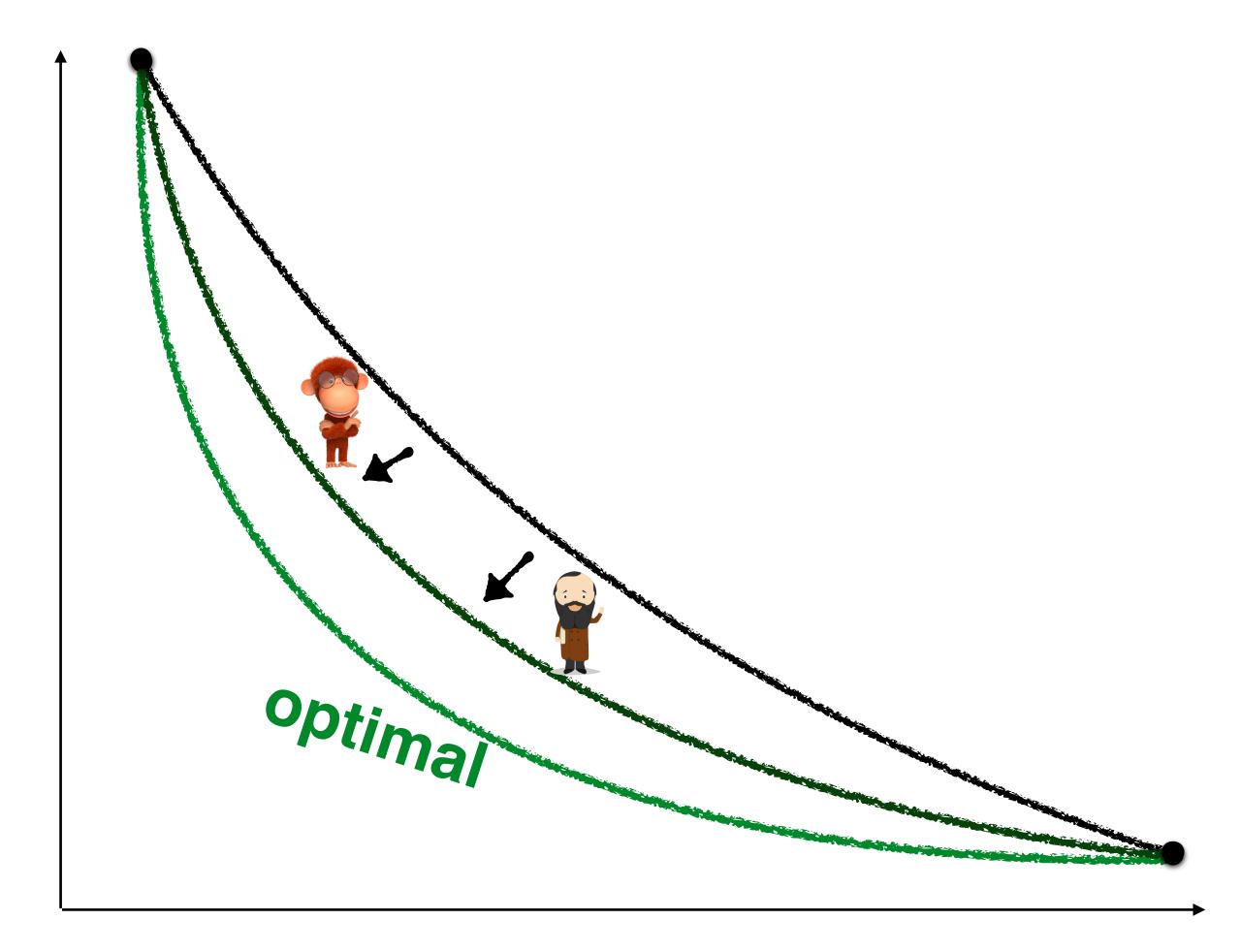






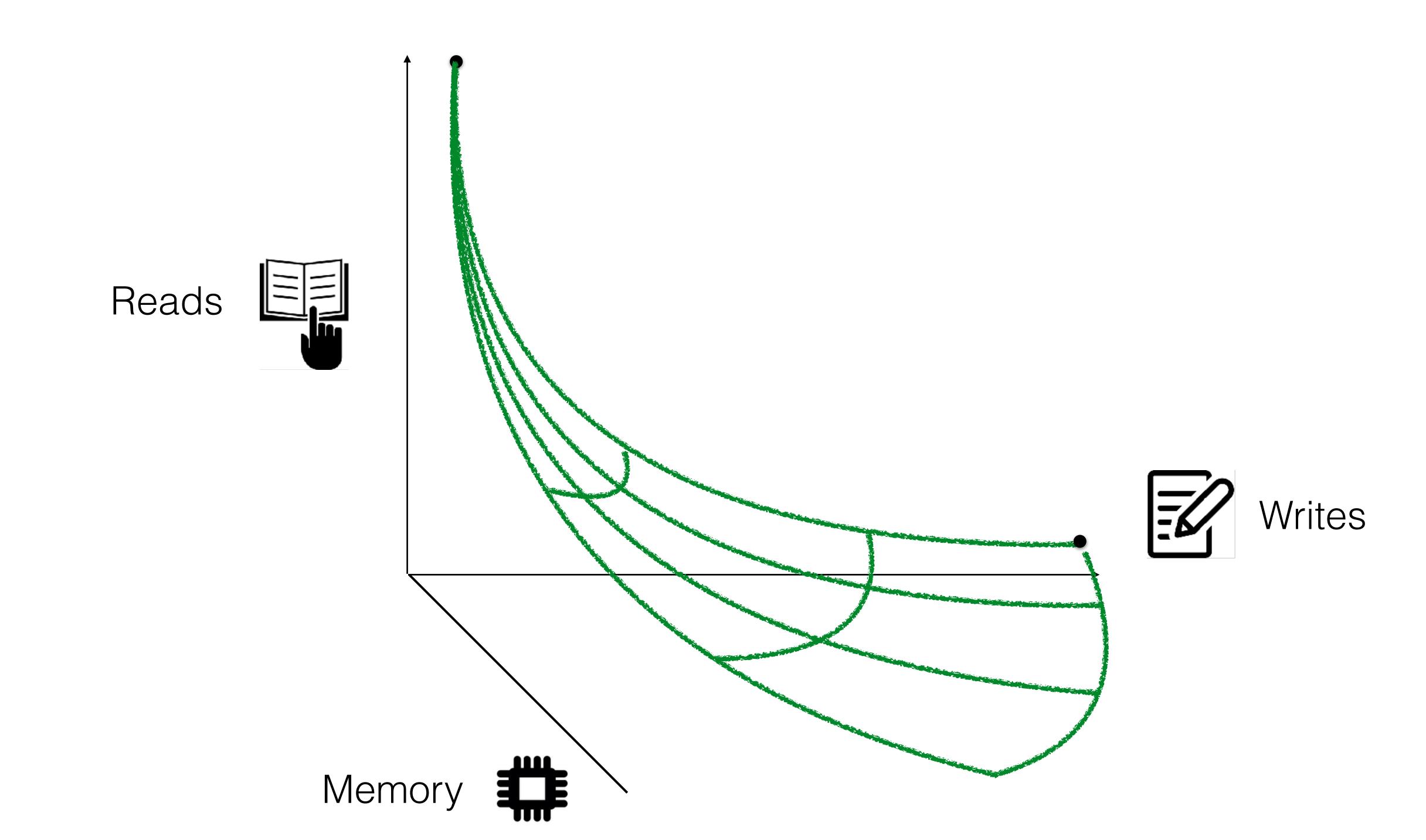
Writes

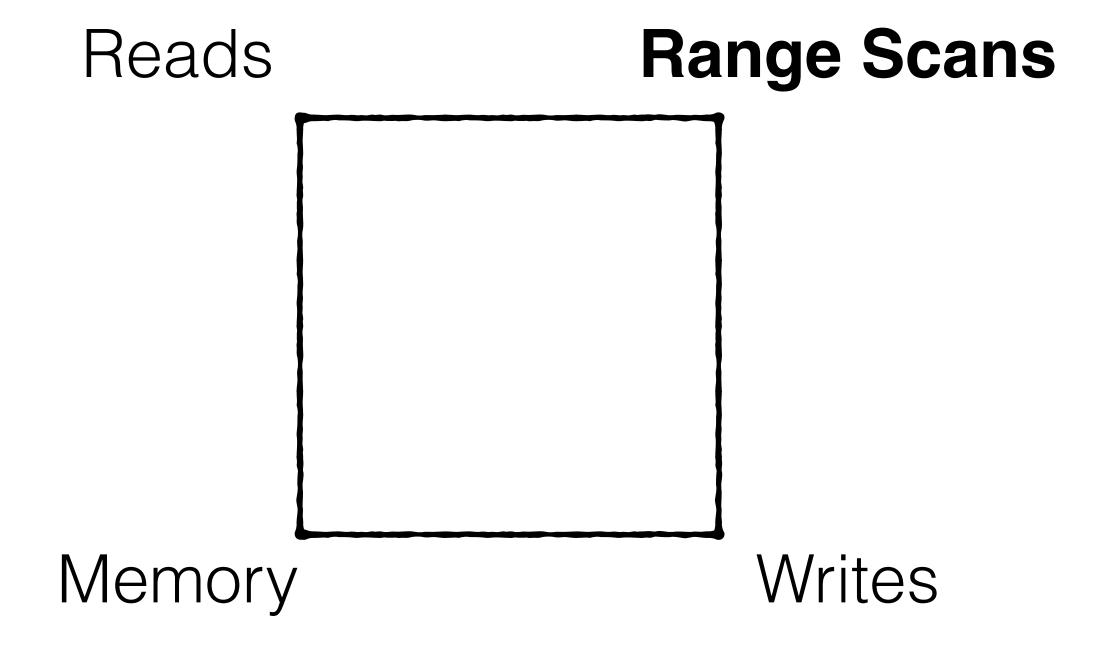


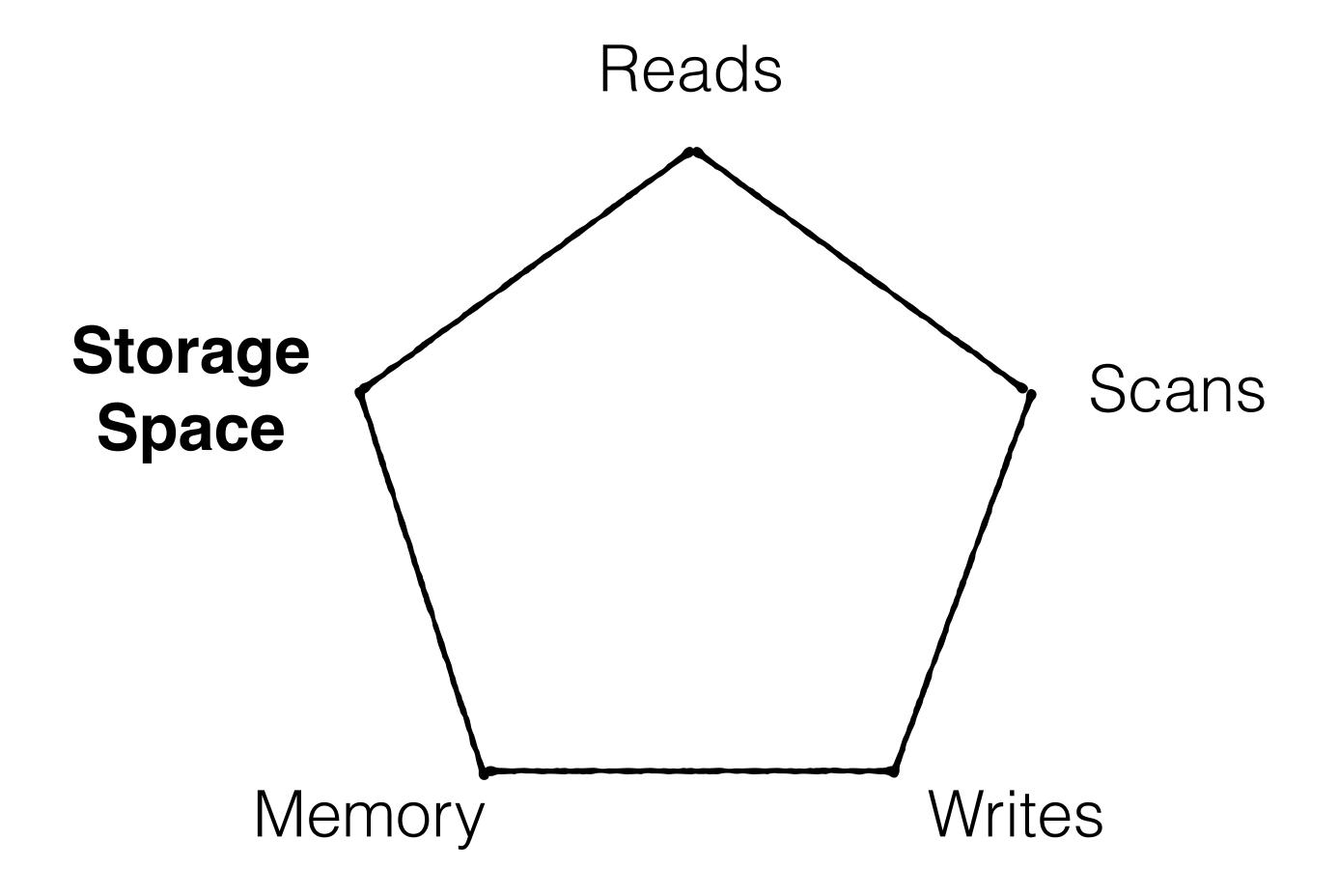


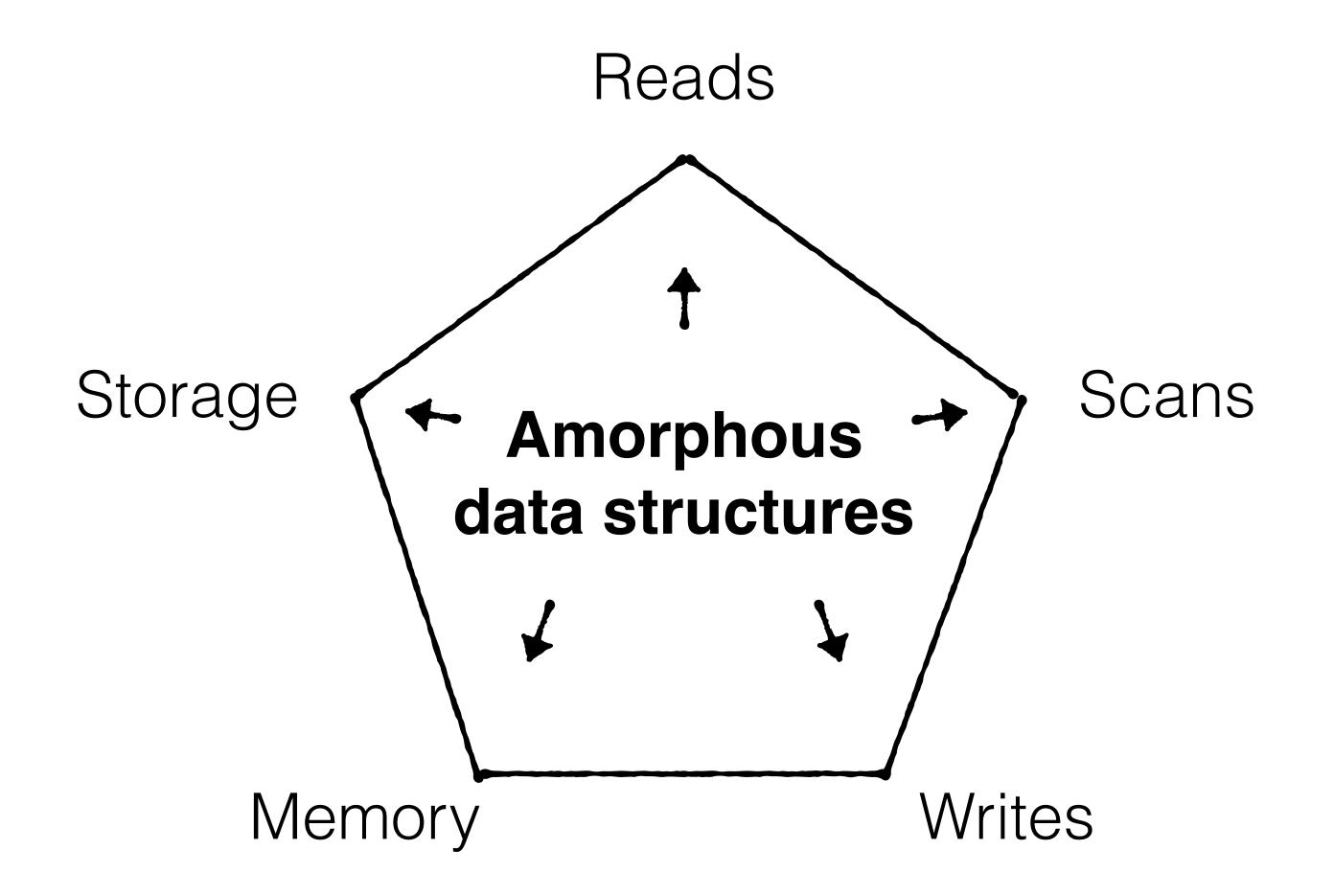


Writes









Thanks



