

Scaling Log-Structured KV-Stores

featuring

Monkey and Dostoevsky

SIGMOD17 / SIGMOD18



DASlab
@ Harvard SEAS

Niv Dayan



BigTable



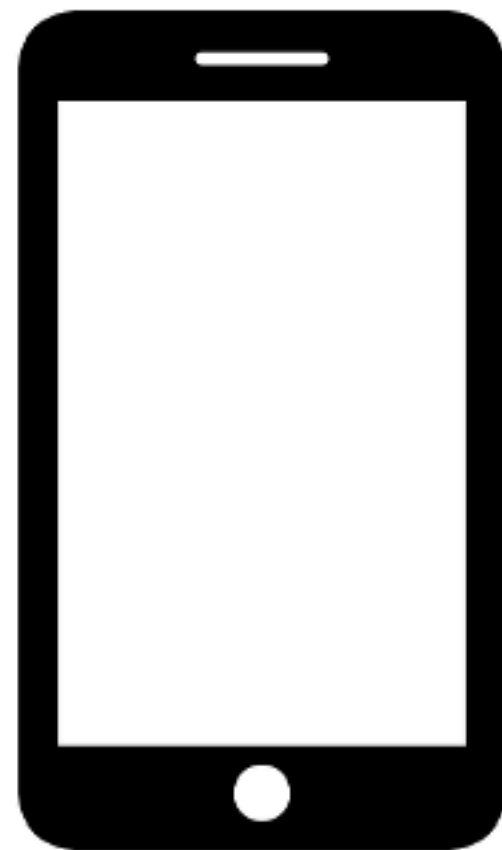
Log-Structured KV-Stores



Couchbase



Log-Structured KV-Stores



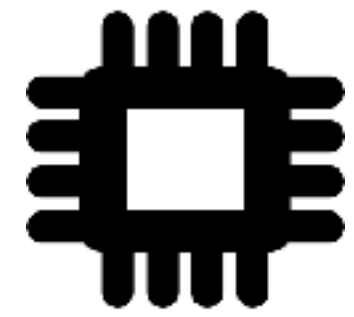
Why Log-Structured KV-Stores?

Why Log-Structured KV-Stores?



fast writes

Why Log-Structured KV-Stores?

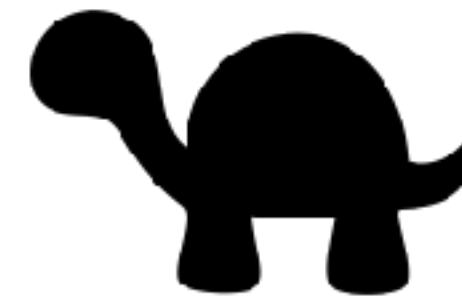
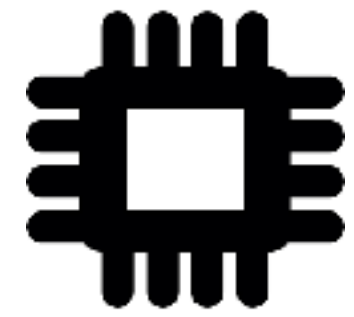


memory

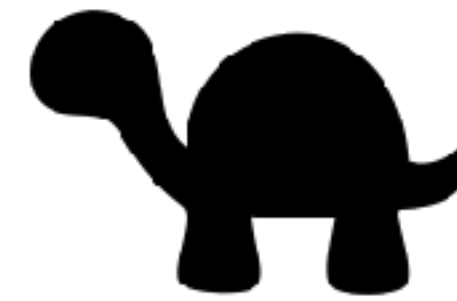
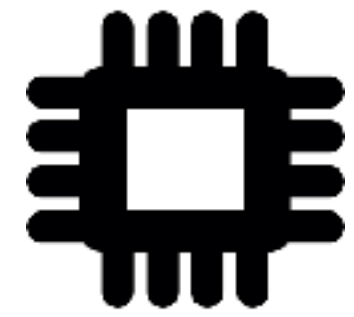


storage

Why Log-Structured KV-Stores?

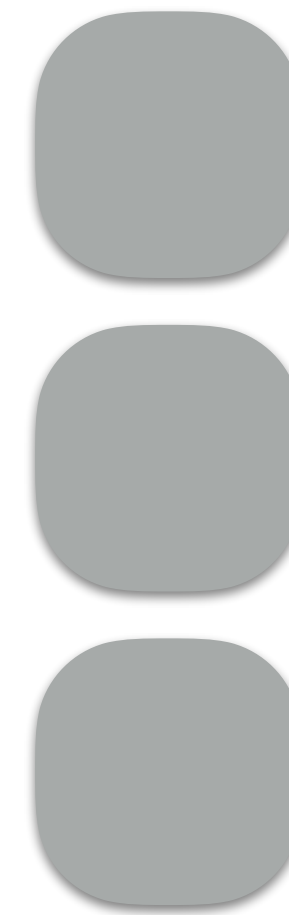
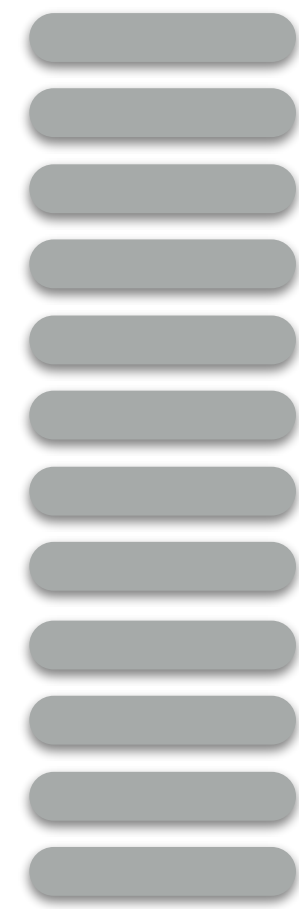
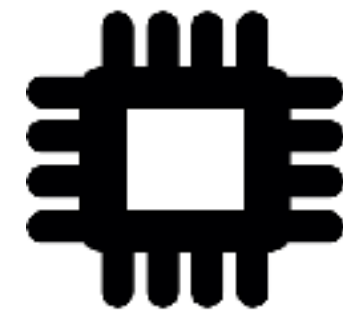


Why Log-Structured KV-Stores?



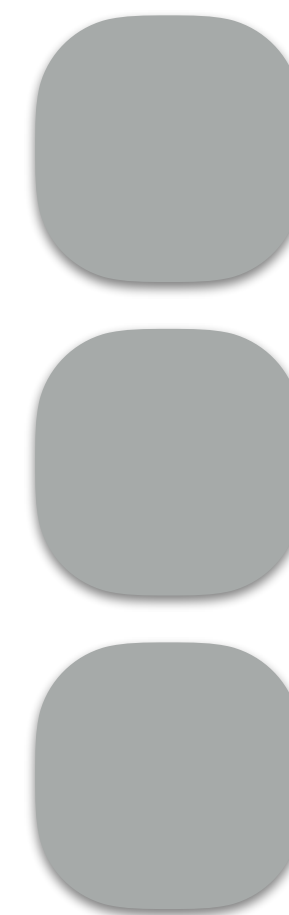
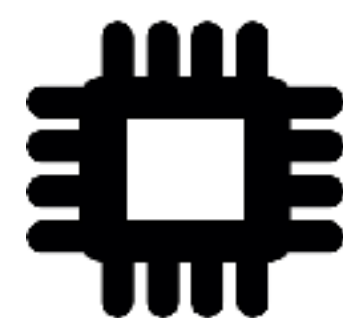
Why Log-Structured KV-Stores?

byte-addressable

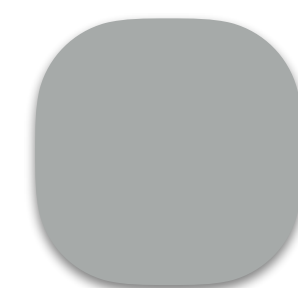
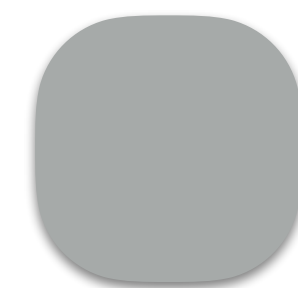
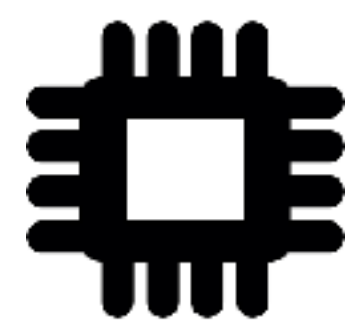


block-addressable

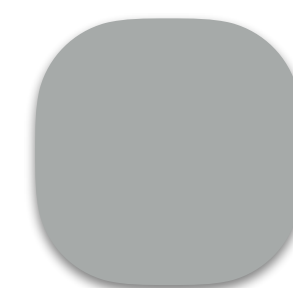
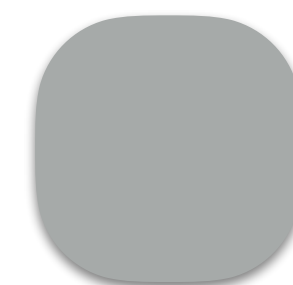
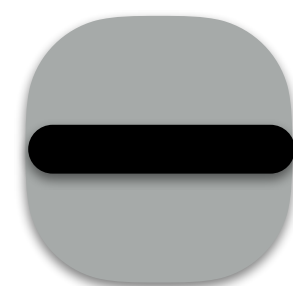
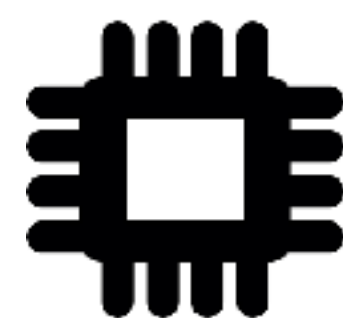
write data



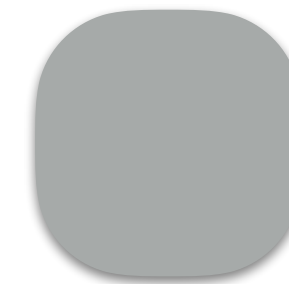
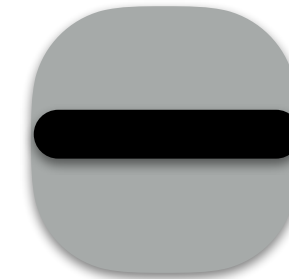
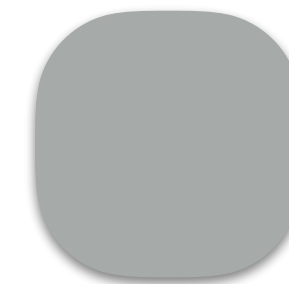
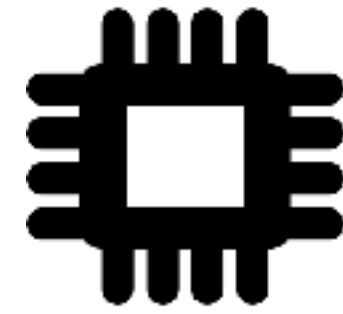
write data



write data

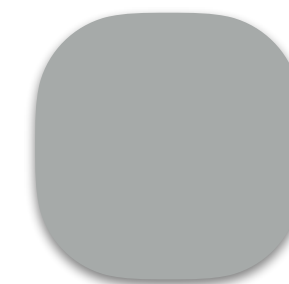
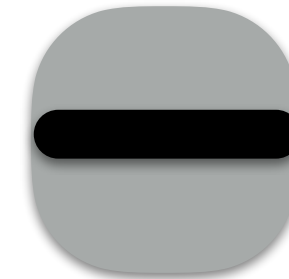
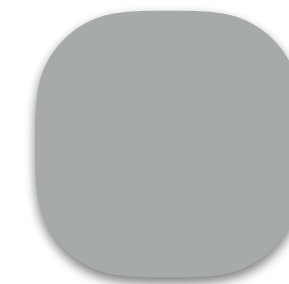
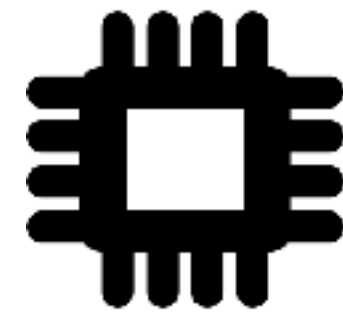


In-Place Writes



write data

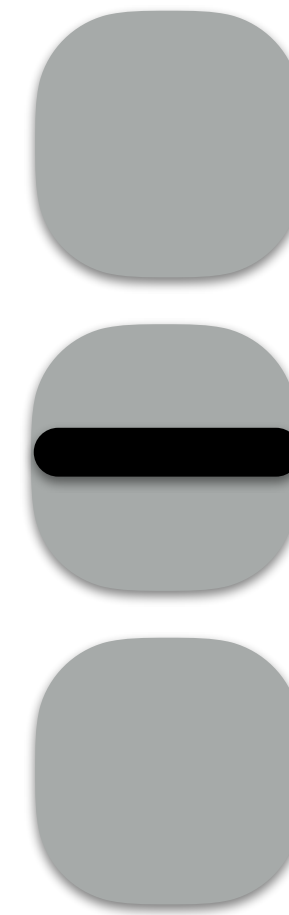
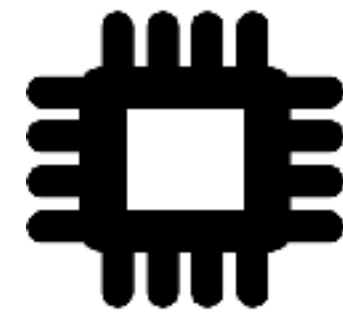
In-Place Writes



write data

B-trees

In-Place Writes

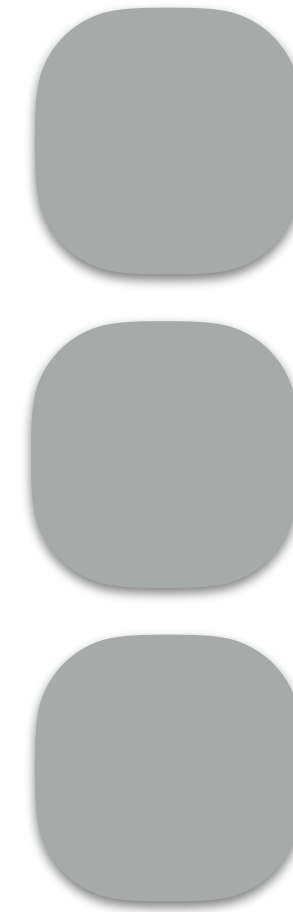
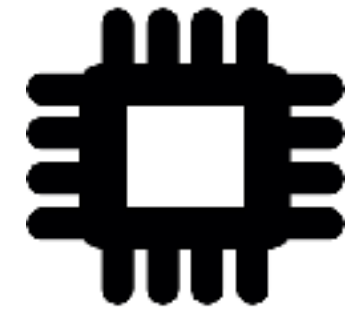


write data



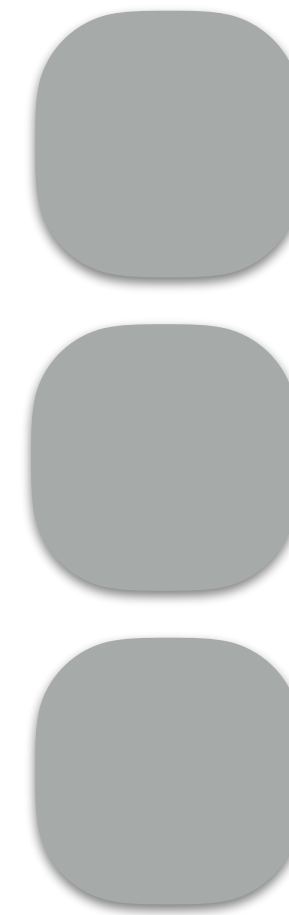
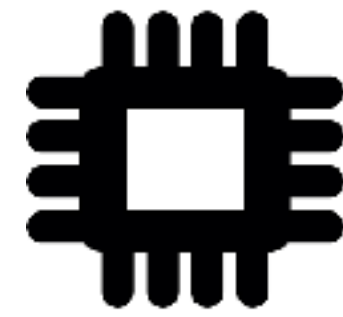
B-trees

Log-Structured Writes

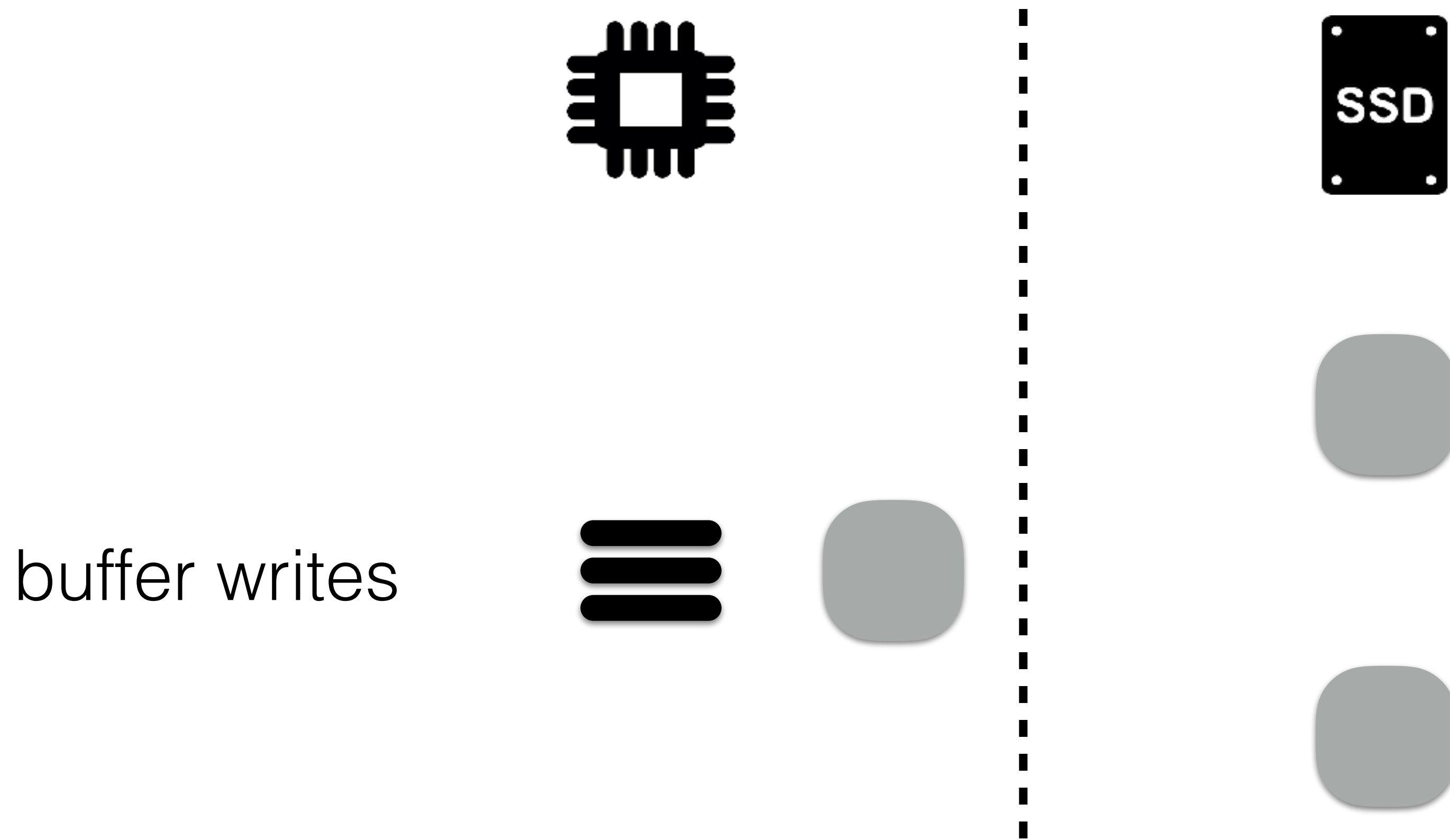


Log-Structured Writes

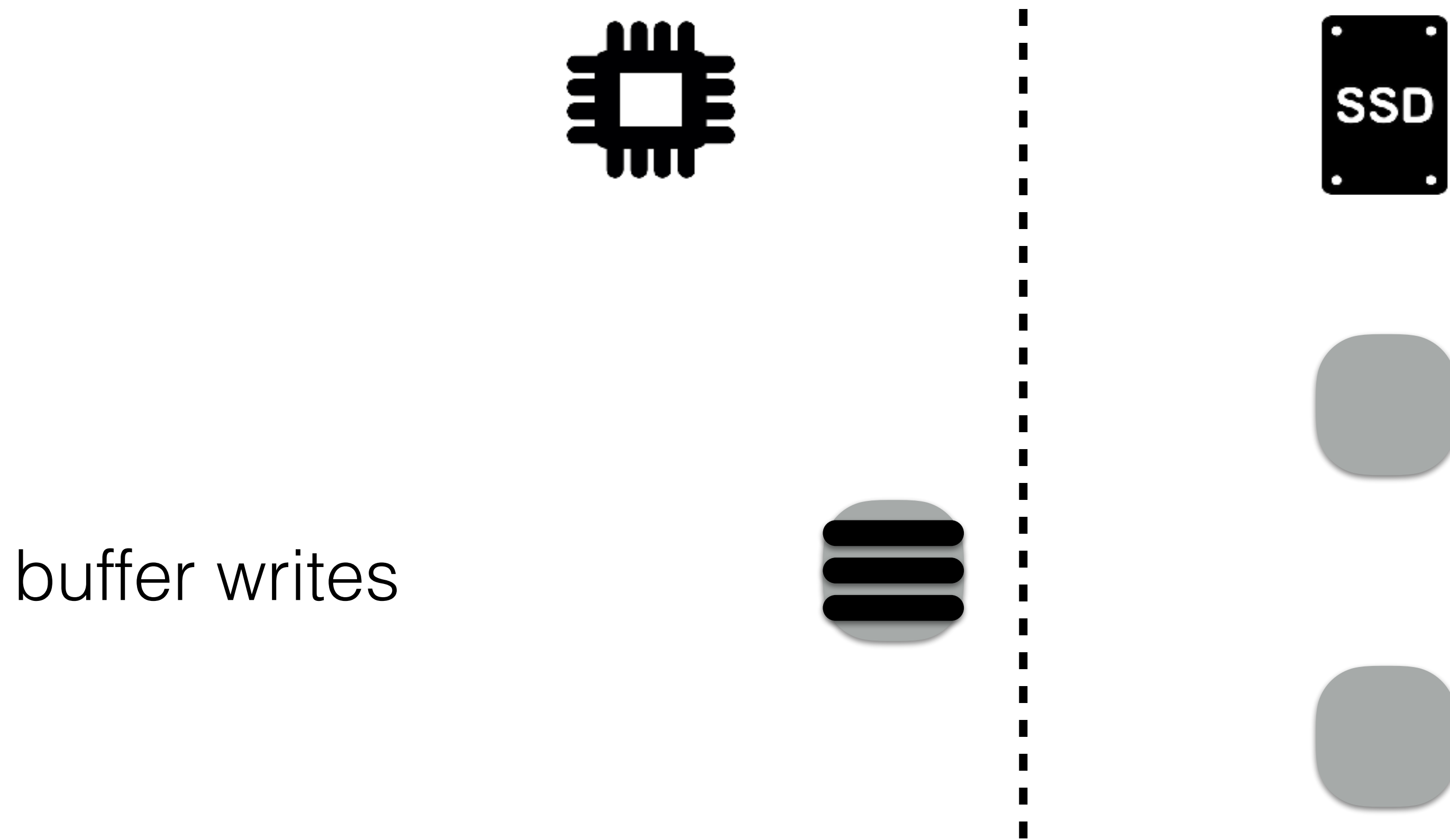
buffer writes



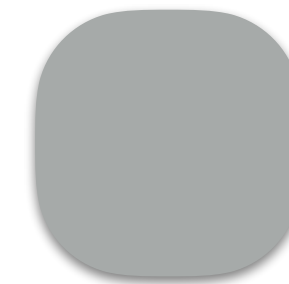
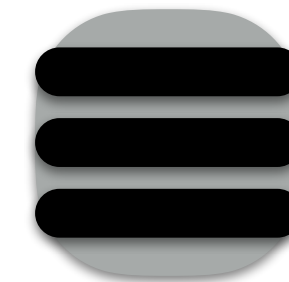
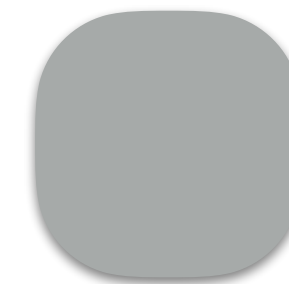
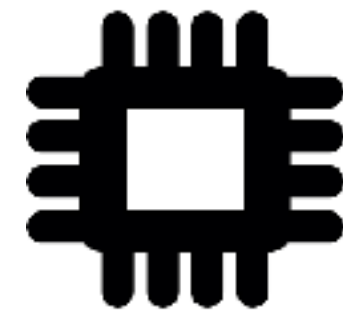
Log-Structured Writes



Log-Structured Writes

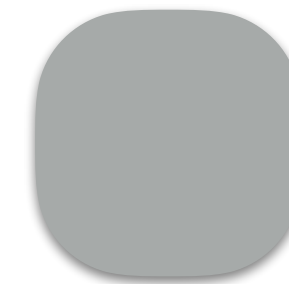
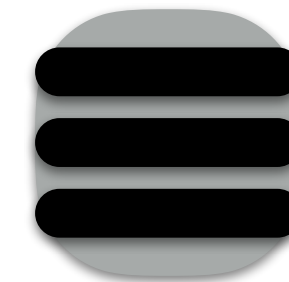
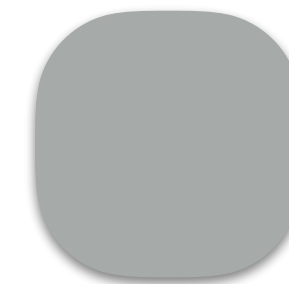
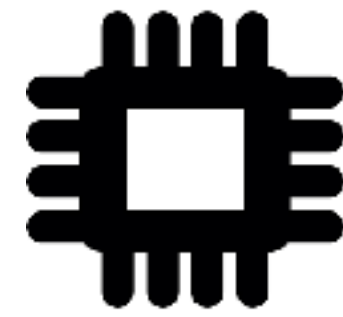


Log-Structured Writes



buffer writes

Log-Structured Writes



buffer writes



Log-Structured KV-Stores



fast writes

Log-Structured KV-Stores



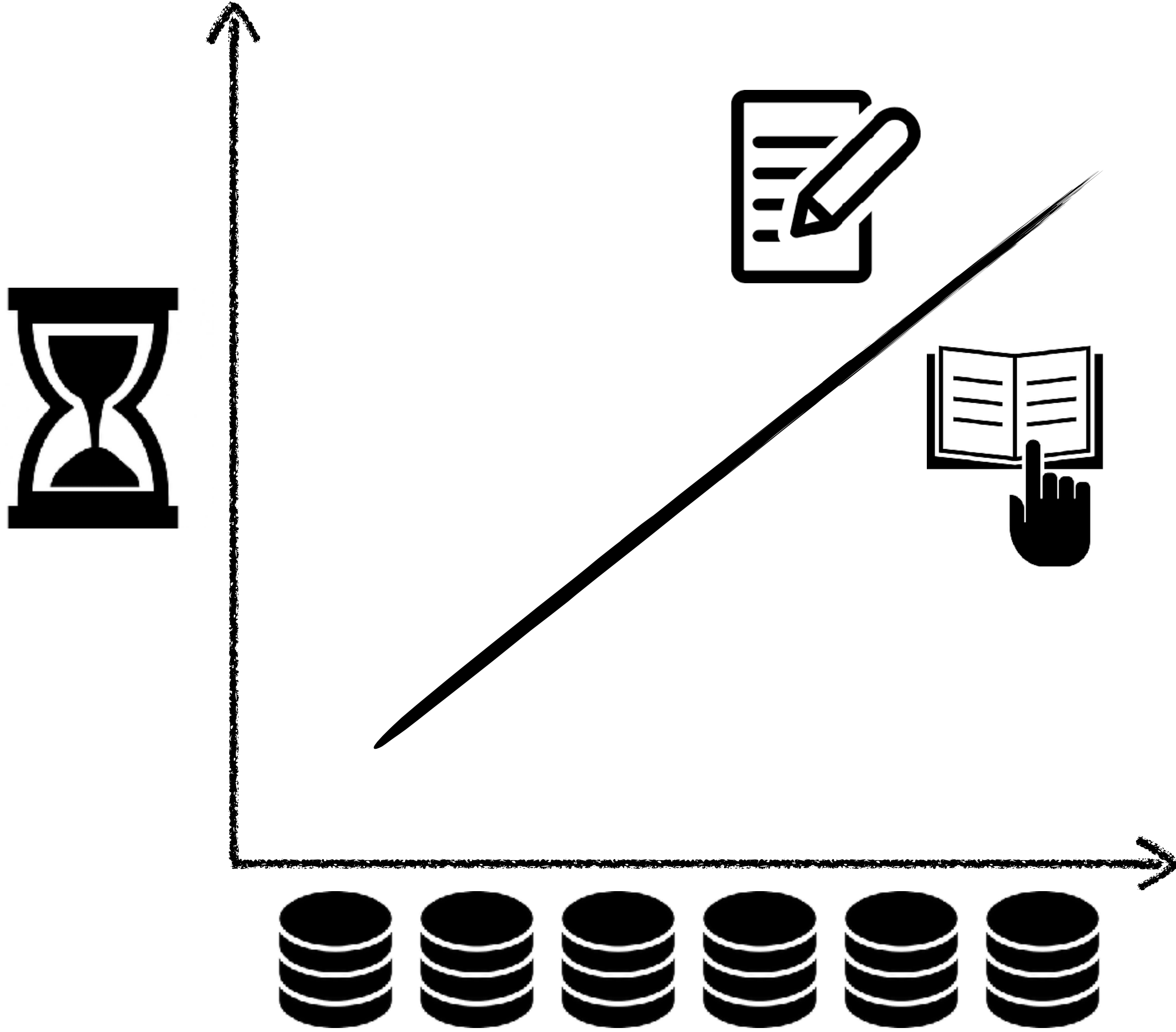
fast writes

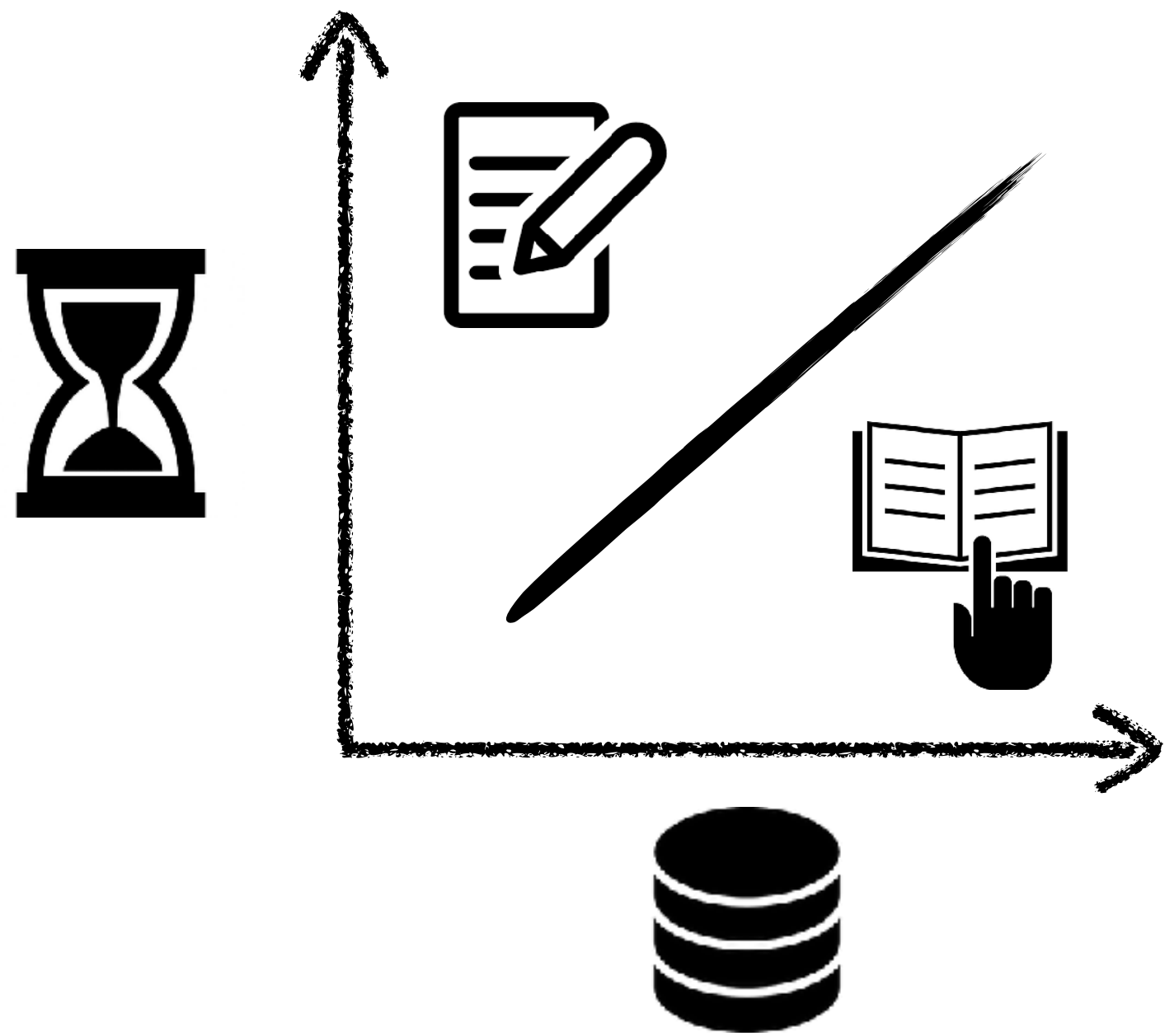


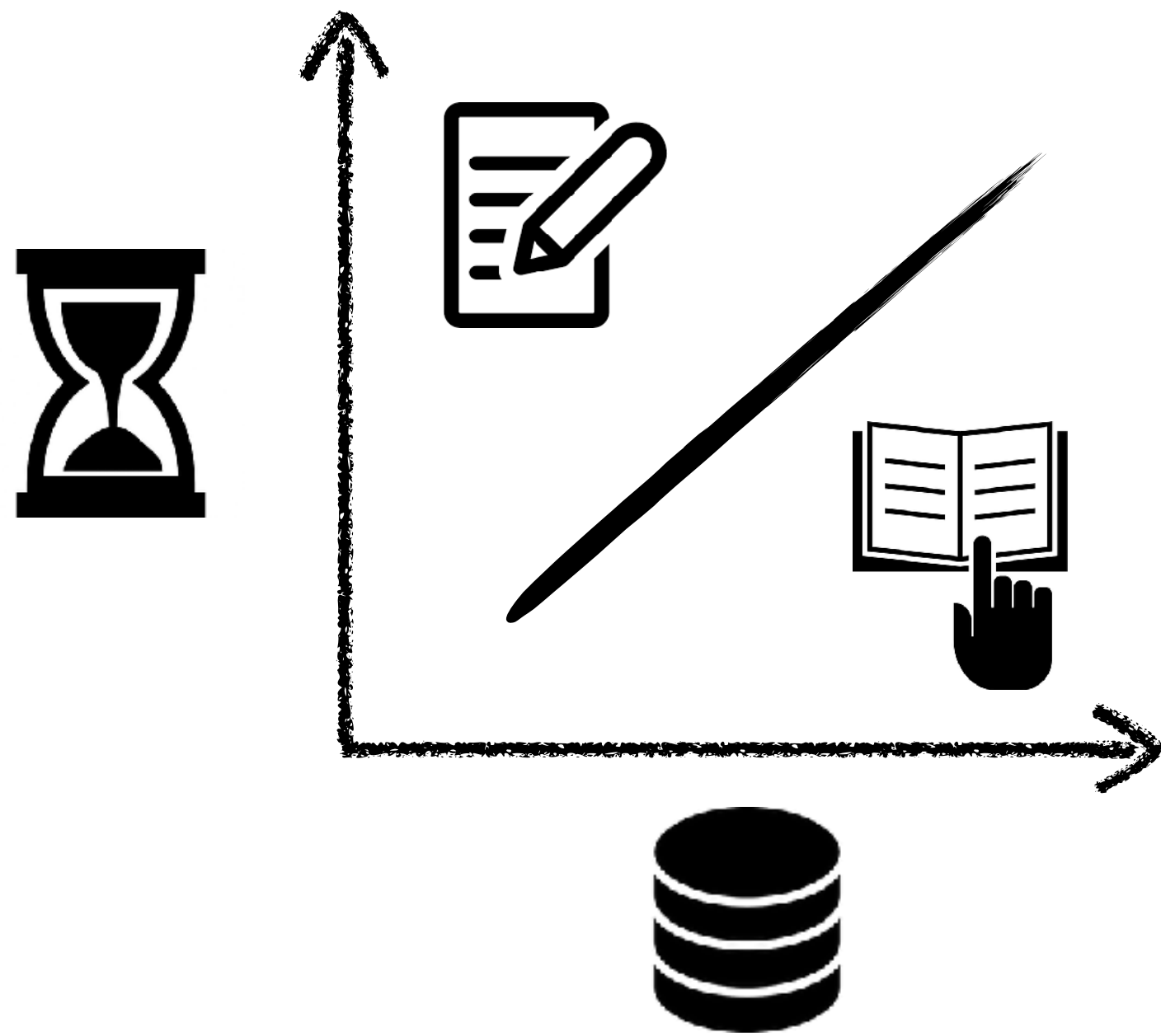
fast reads



massive data







Background



Background



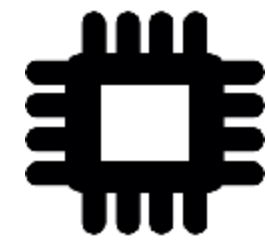
The Log-Structured Merge-Tree

Background



LSM-tree

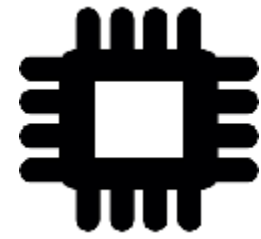
buffer



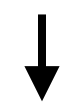
writes



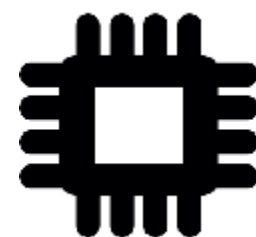
buffer



key value pairs



buffer



key

value

Sherlock:

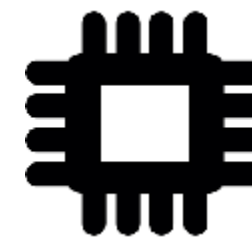
a fictional detective

Waldo:

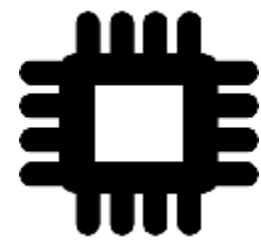
an inconspicuous traveler

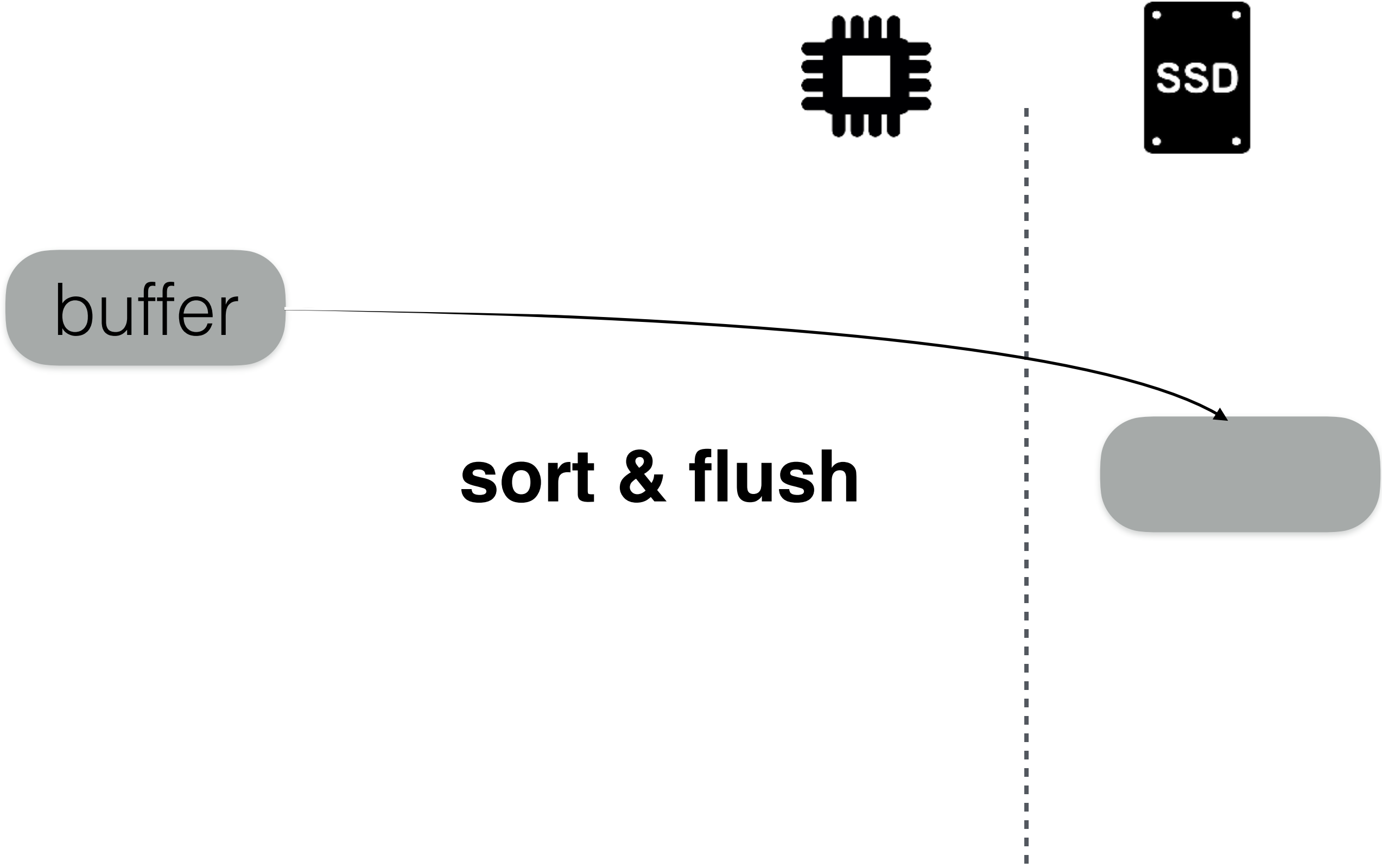


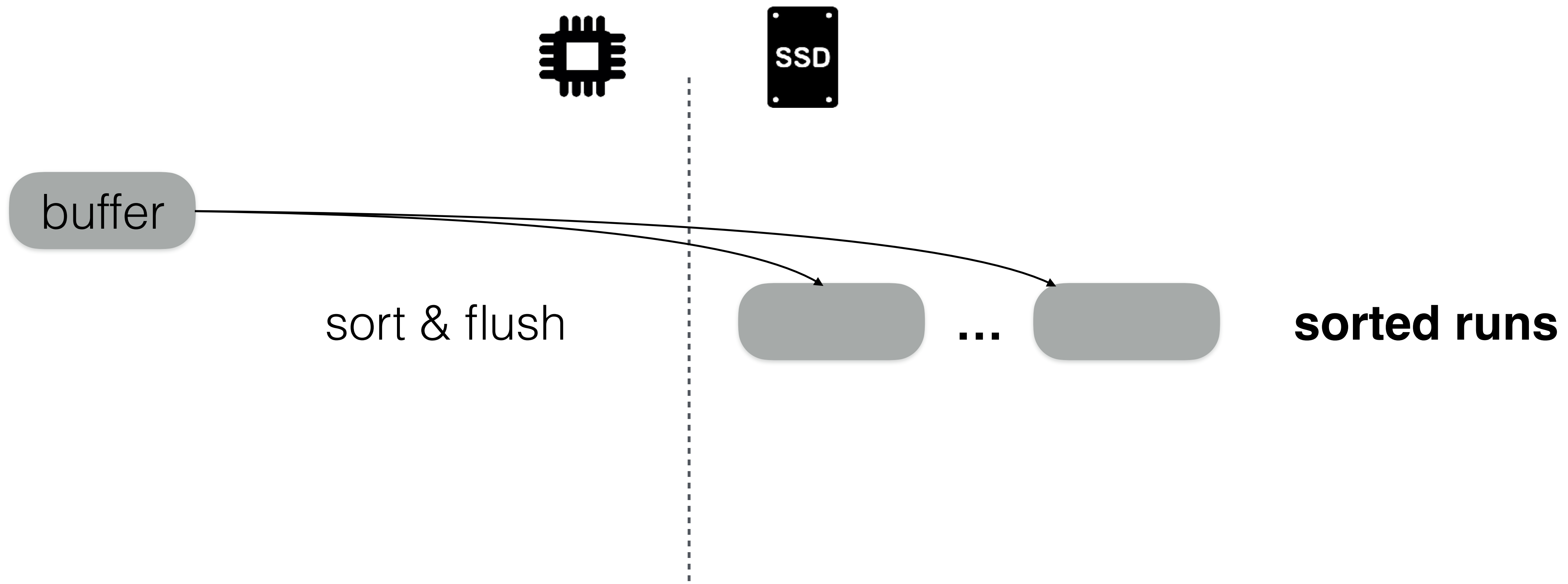
buffer

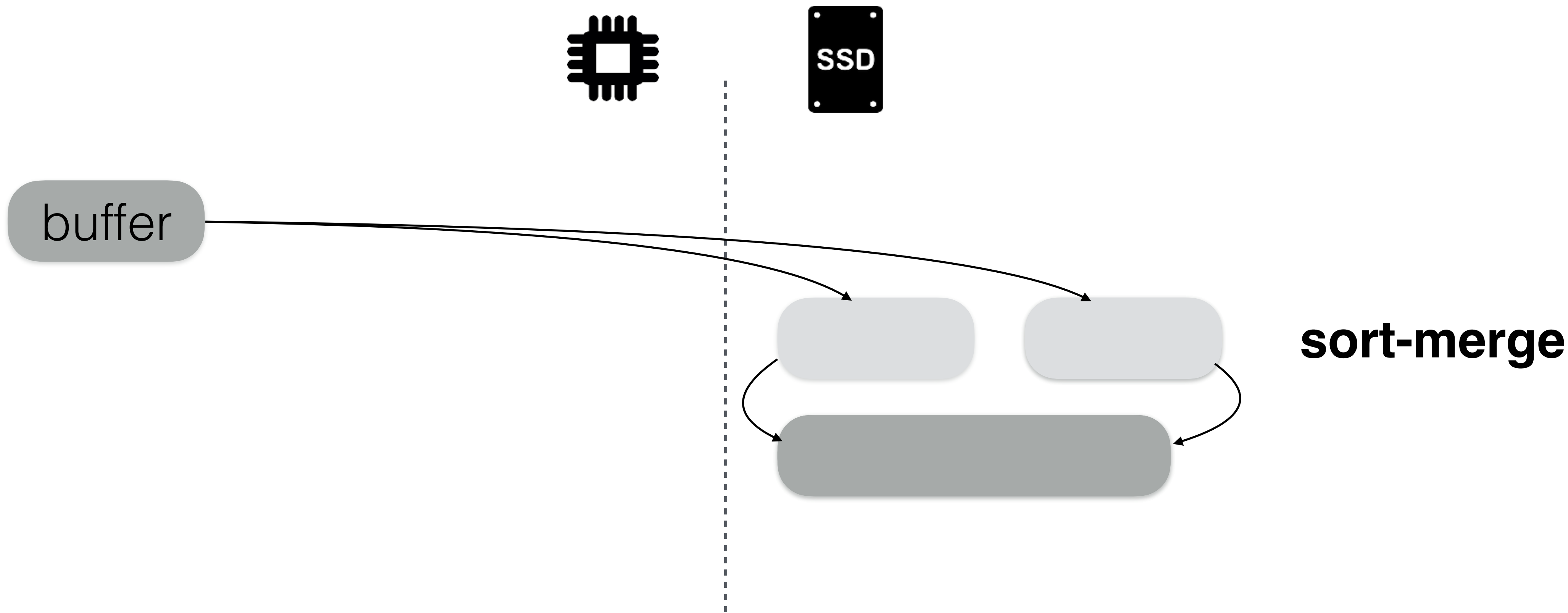


buffer gets full

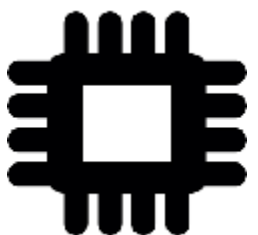








buffer



exponentially increasing capacities

level 1



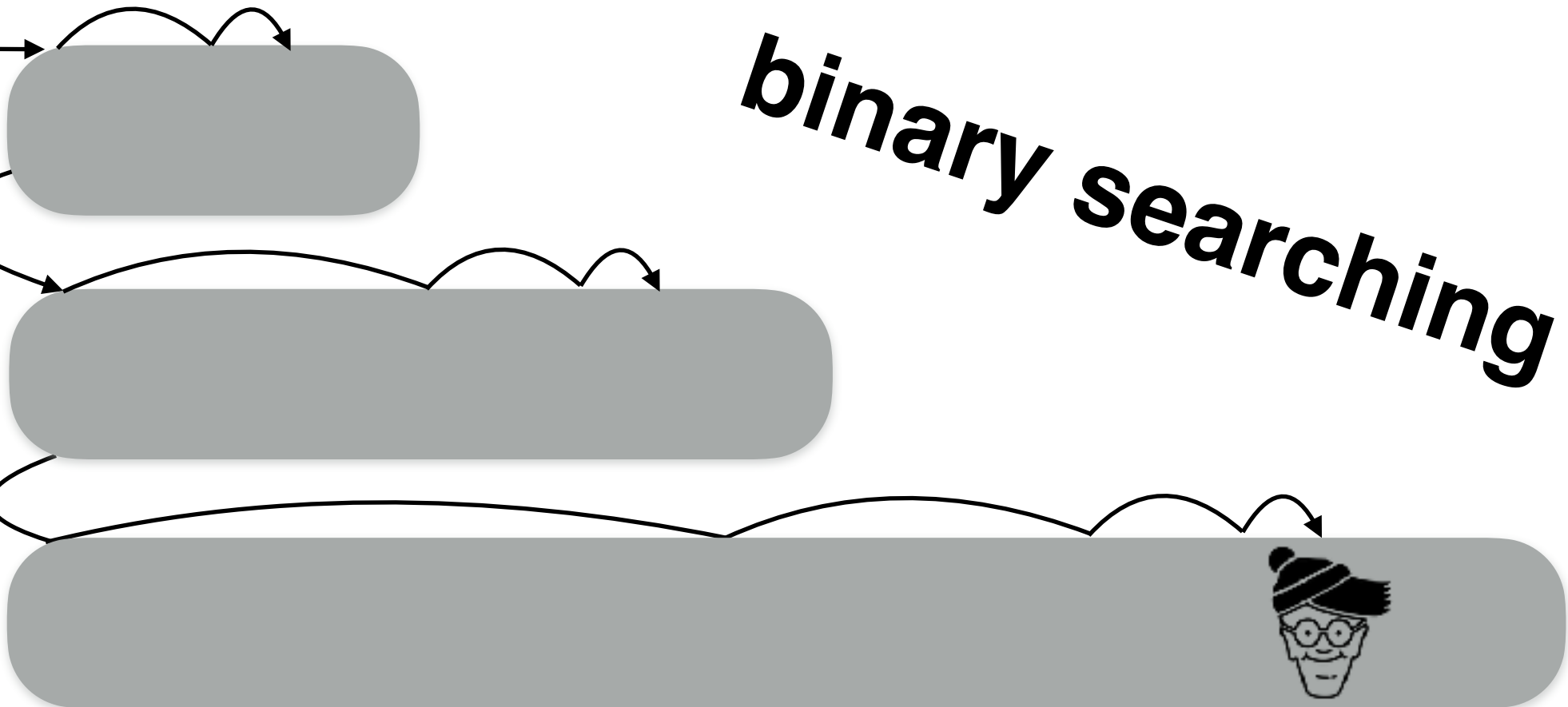
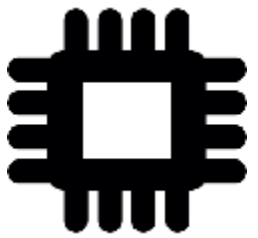
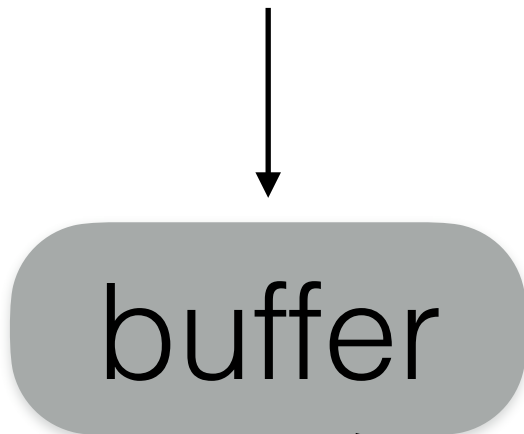
level 2



level 3



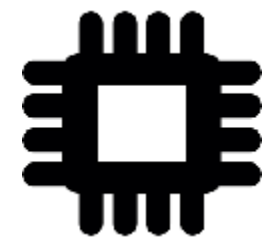
**where's
Waldo**



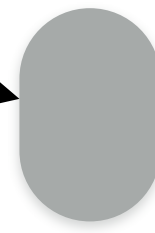
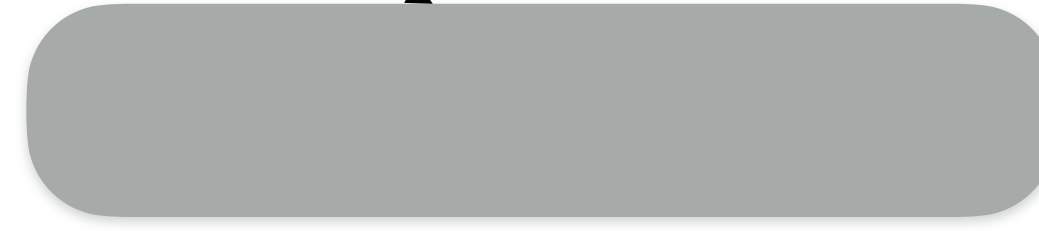
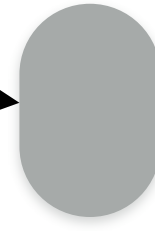
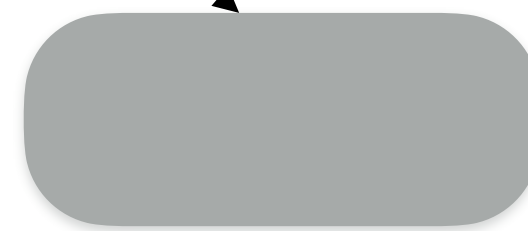
where's
Waldo



buffer

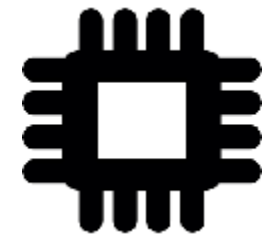
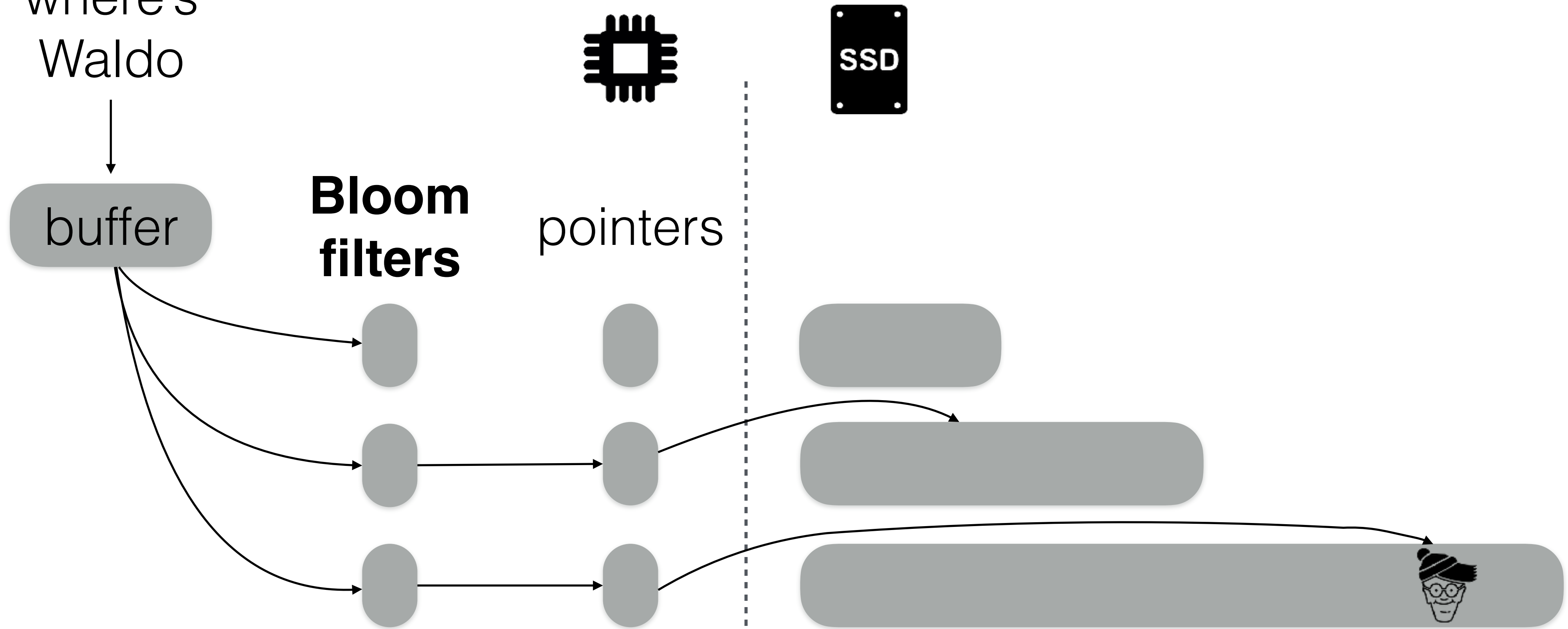


pointers



one I/O per run

where's
Waldo



**Bloom
filters**

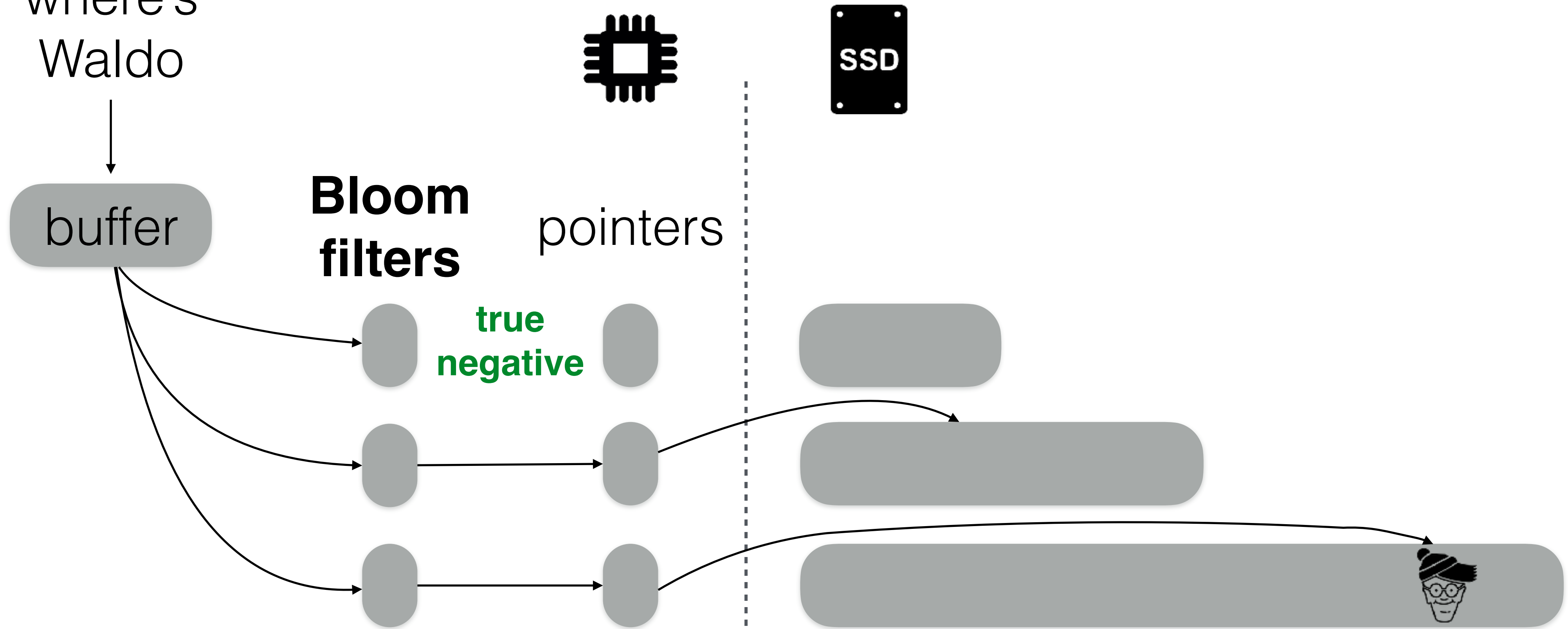
pointers

SSD

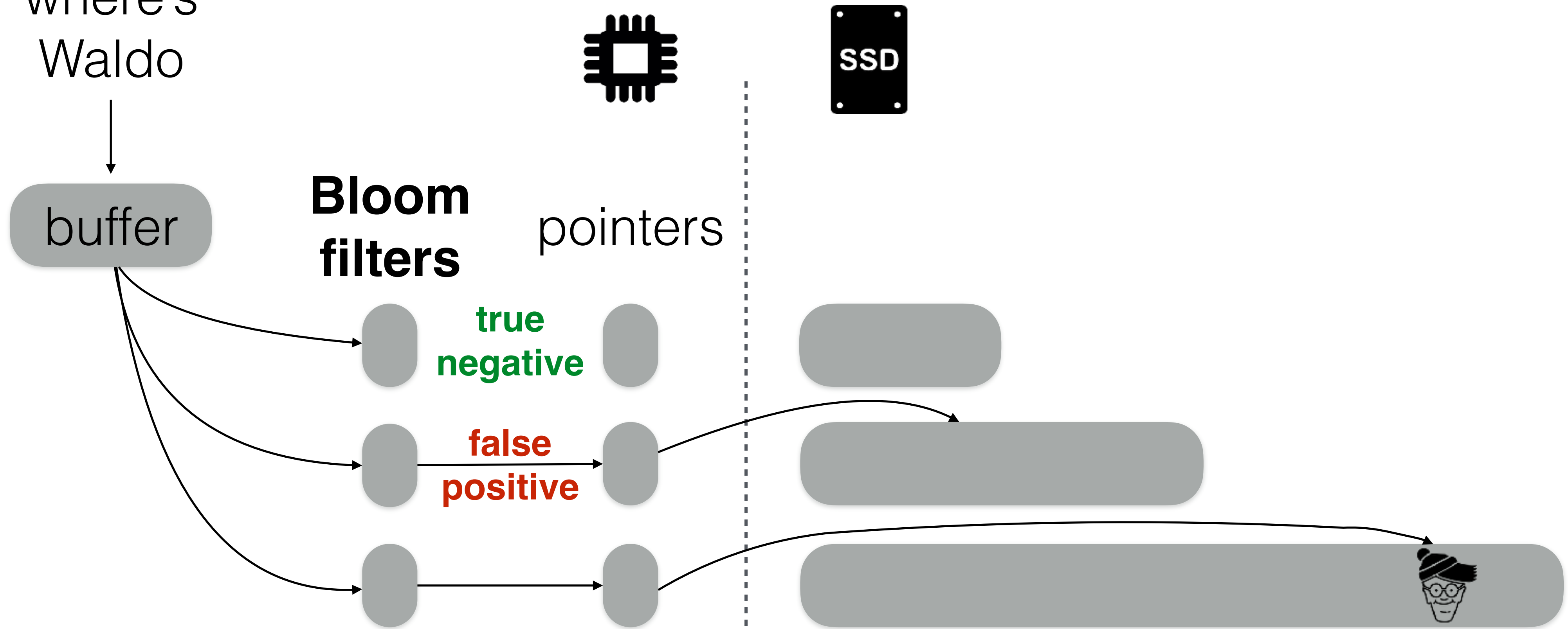
buffer

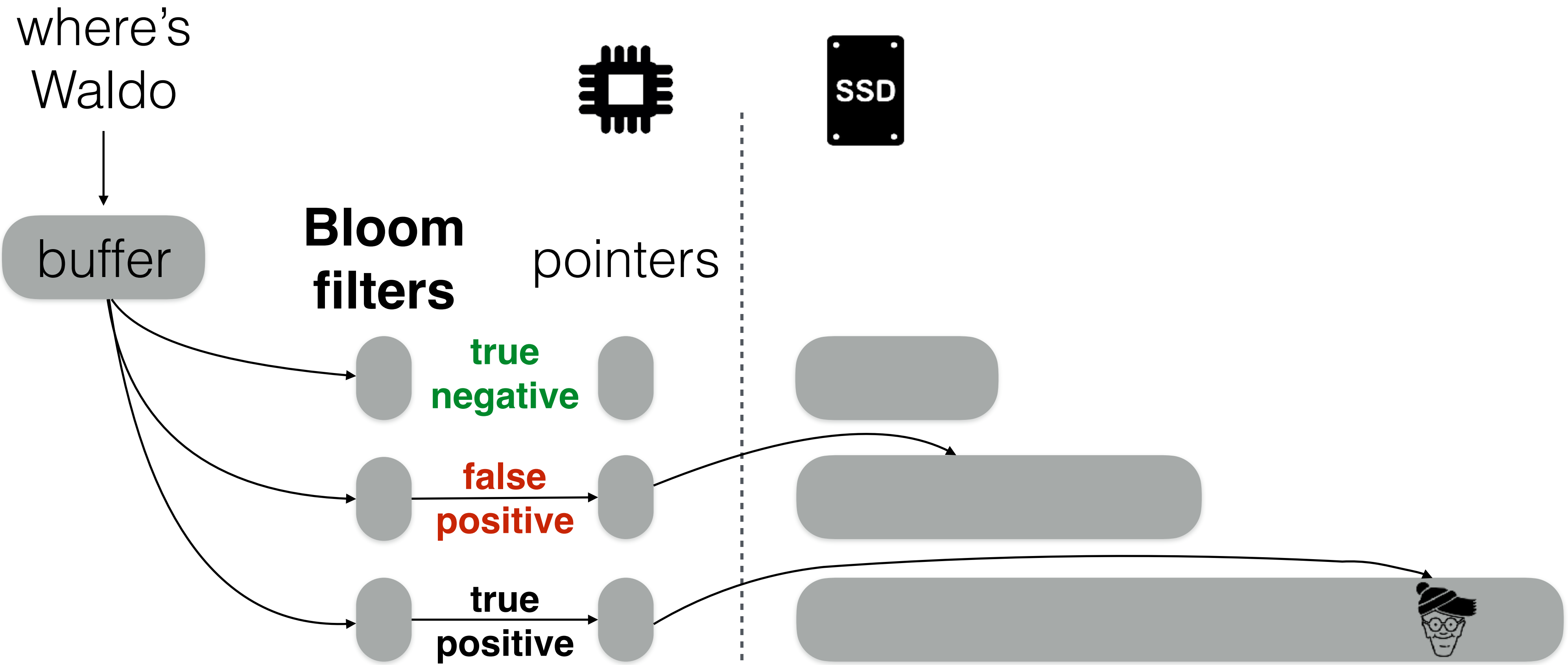


where's
Waldo



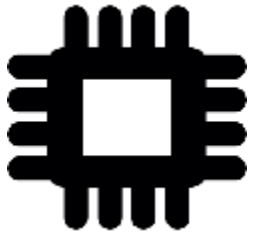
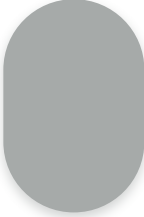
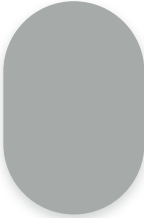
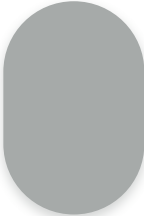
where's
Waldo



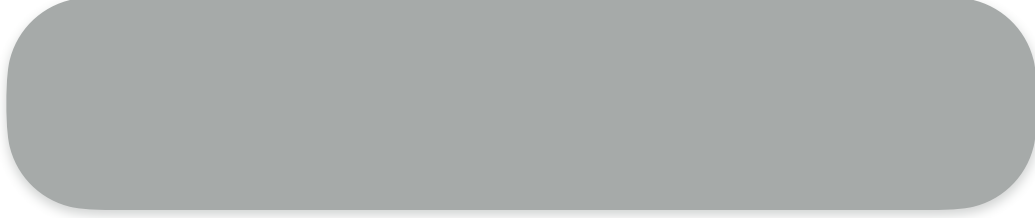
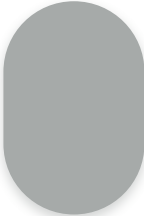


buffer

Bloom
filters



pointers



merging frequency



merging



writes

reads

writes



merging



reads

merging

Tiering

write-optimized

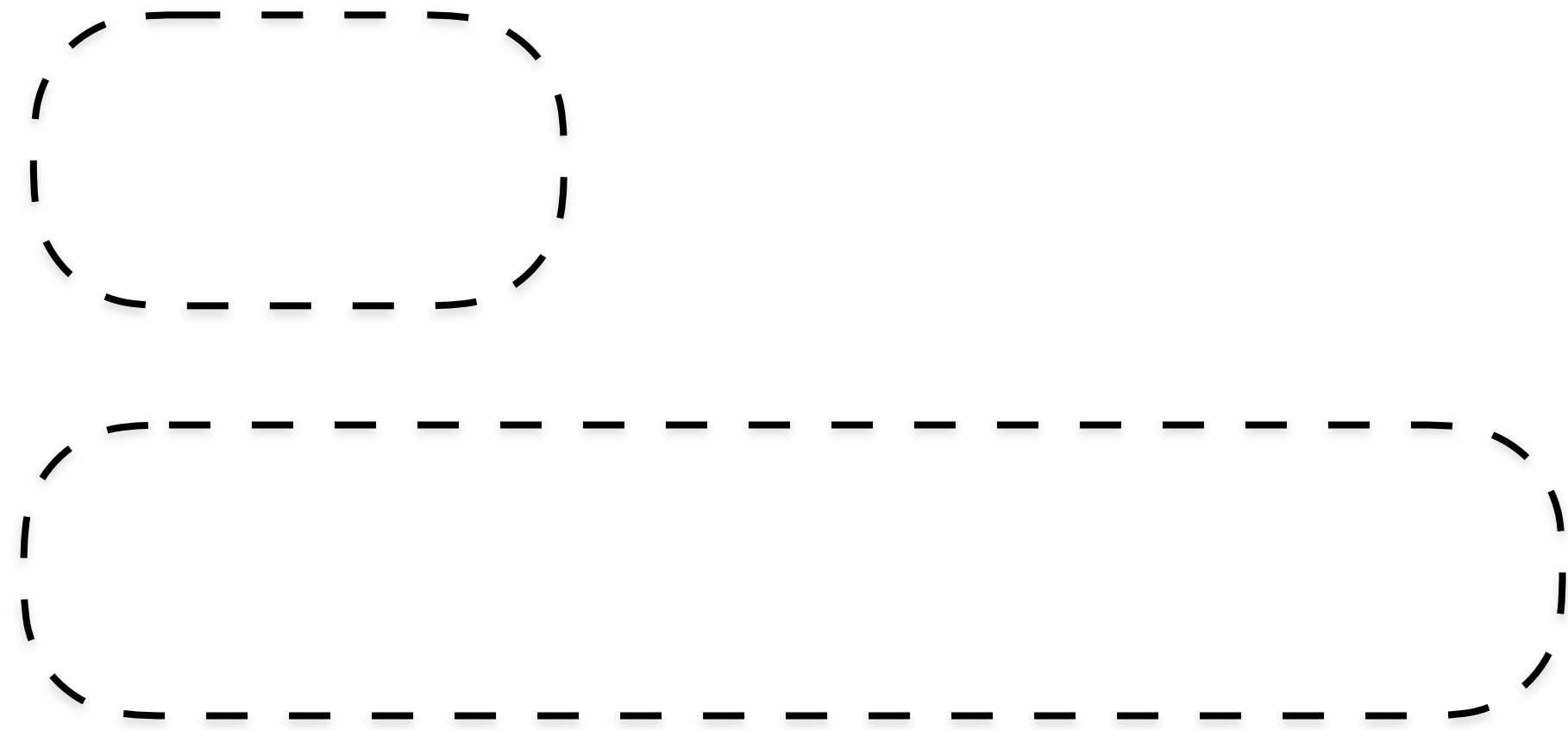


Leveling

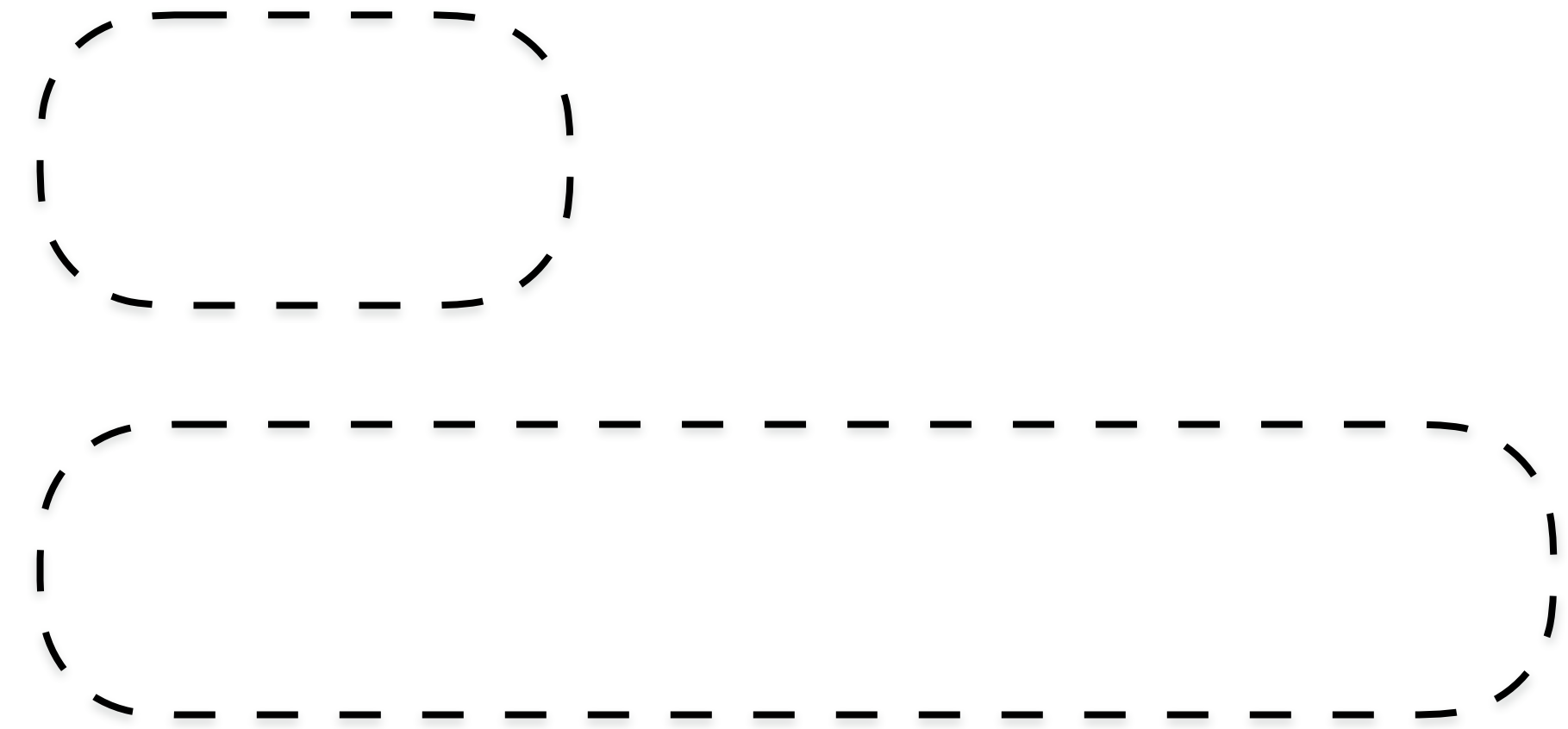
read-optimized



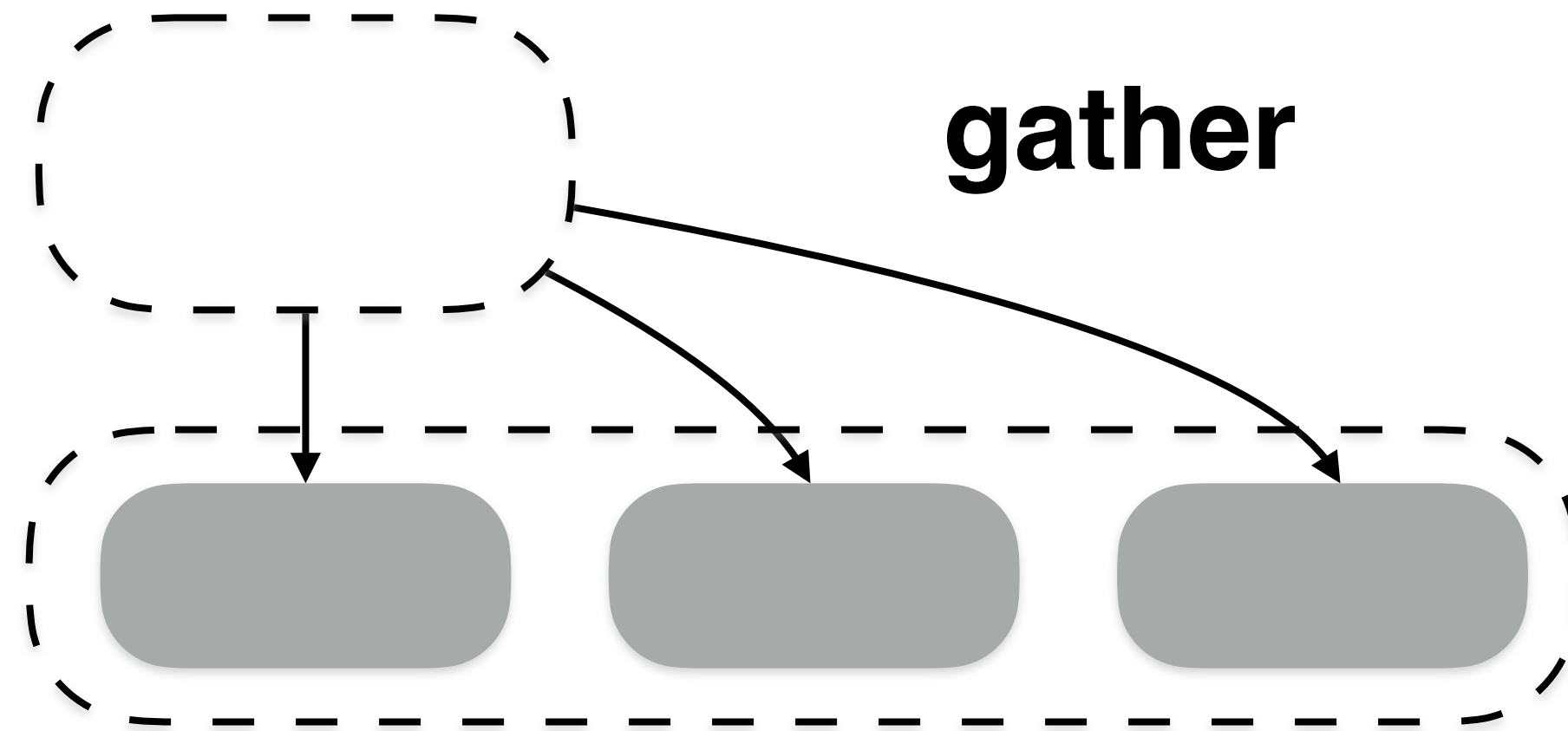
Tiering
write-optimized



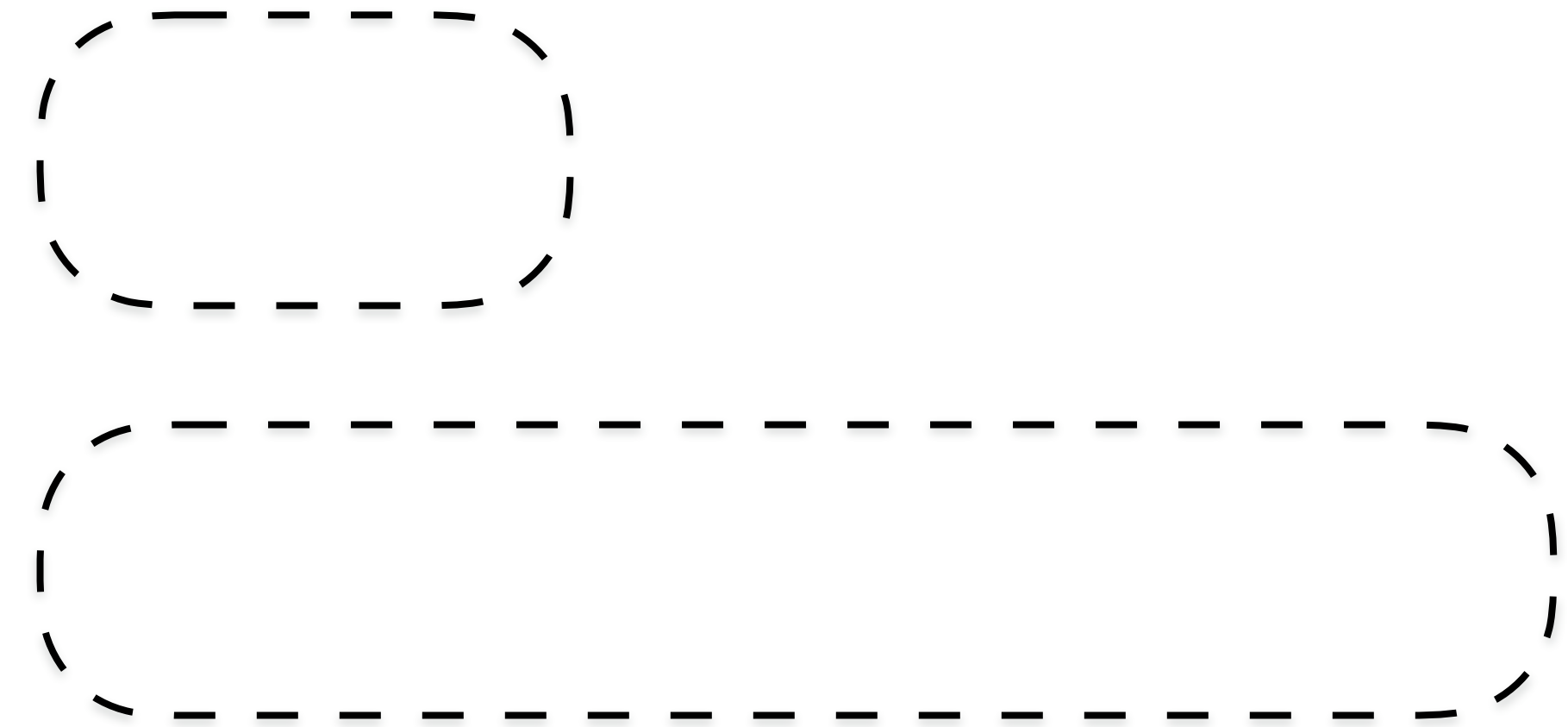
Leveling
read-optimized



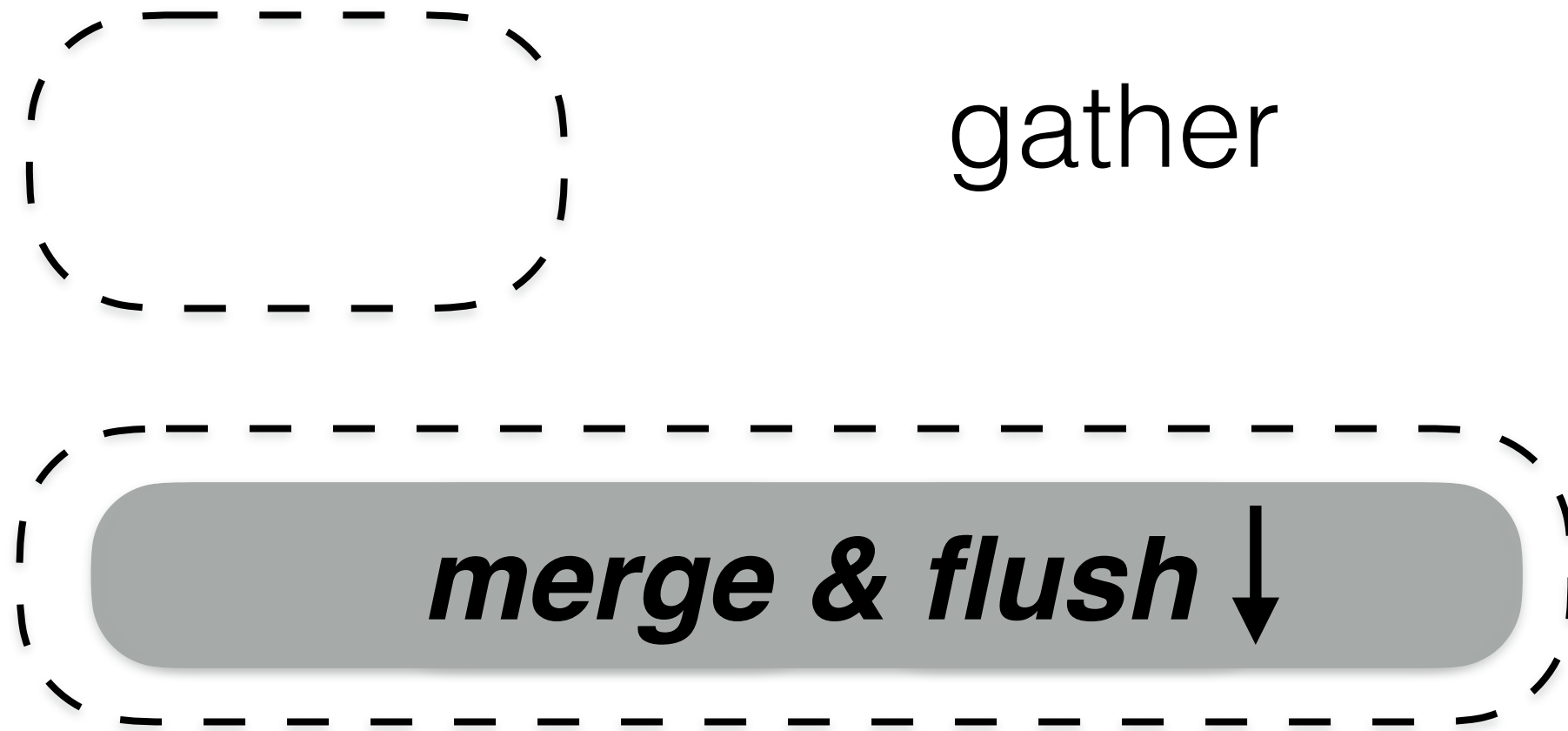
Tiering
write-optimized



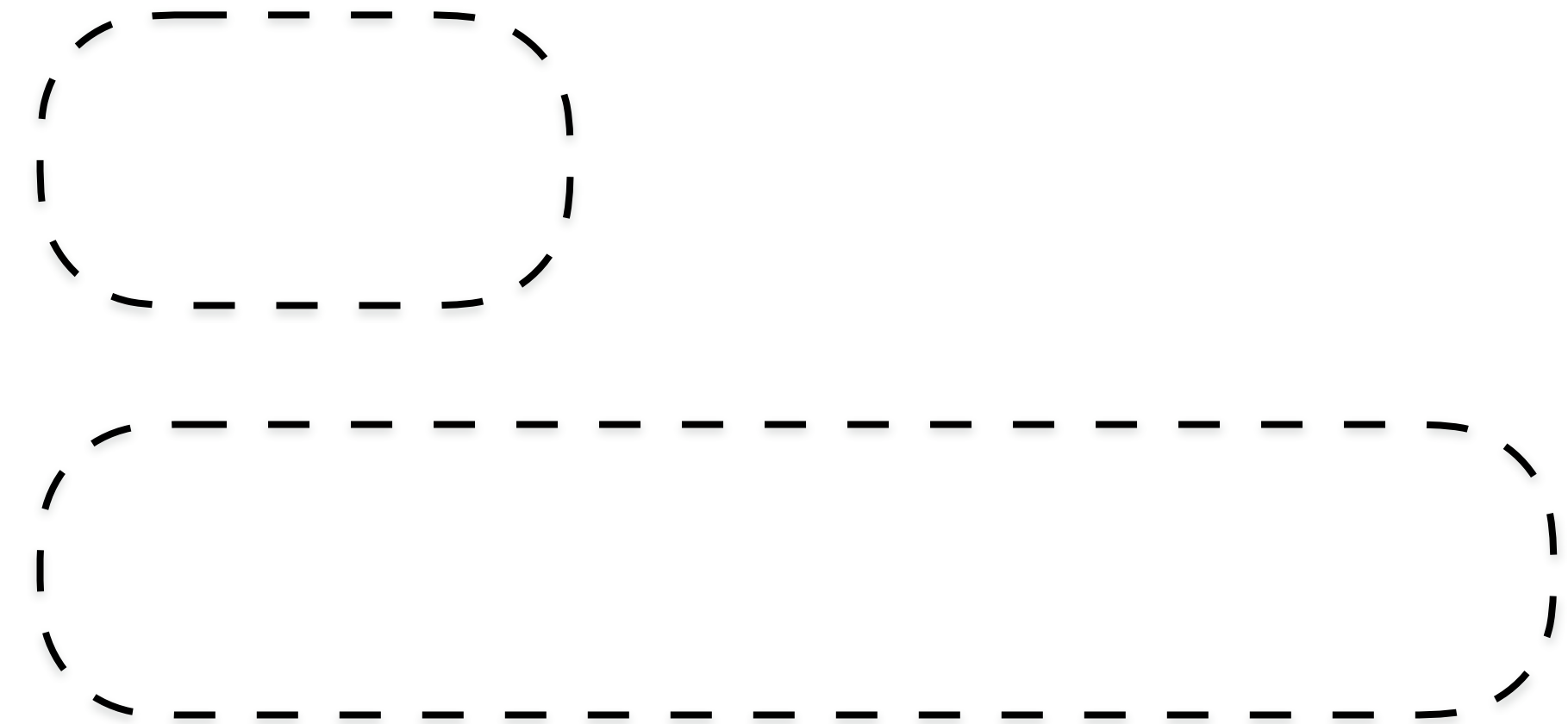
Leveling
read-optimized



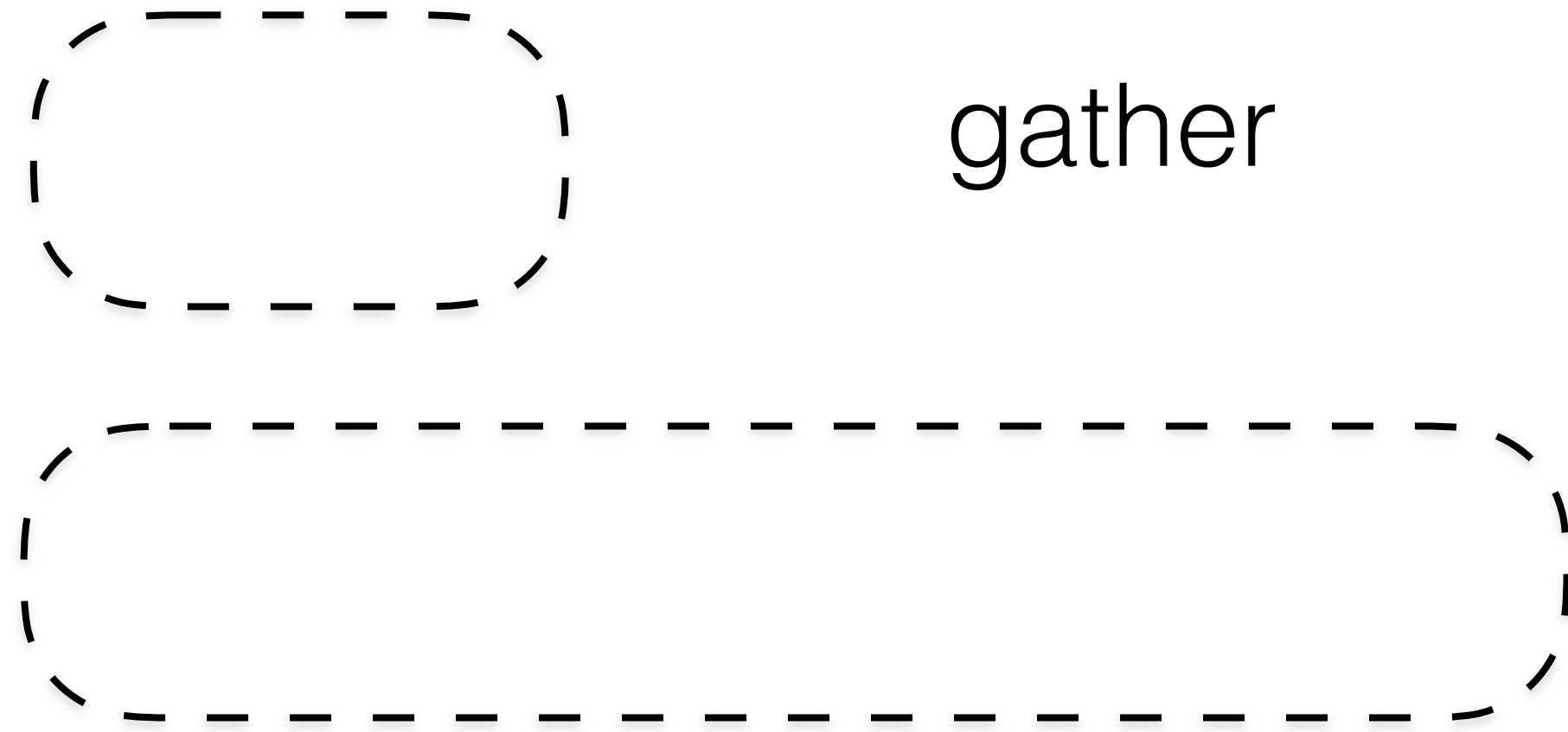
Tiering
write-optimized



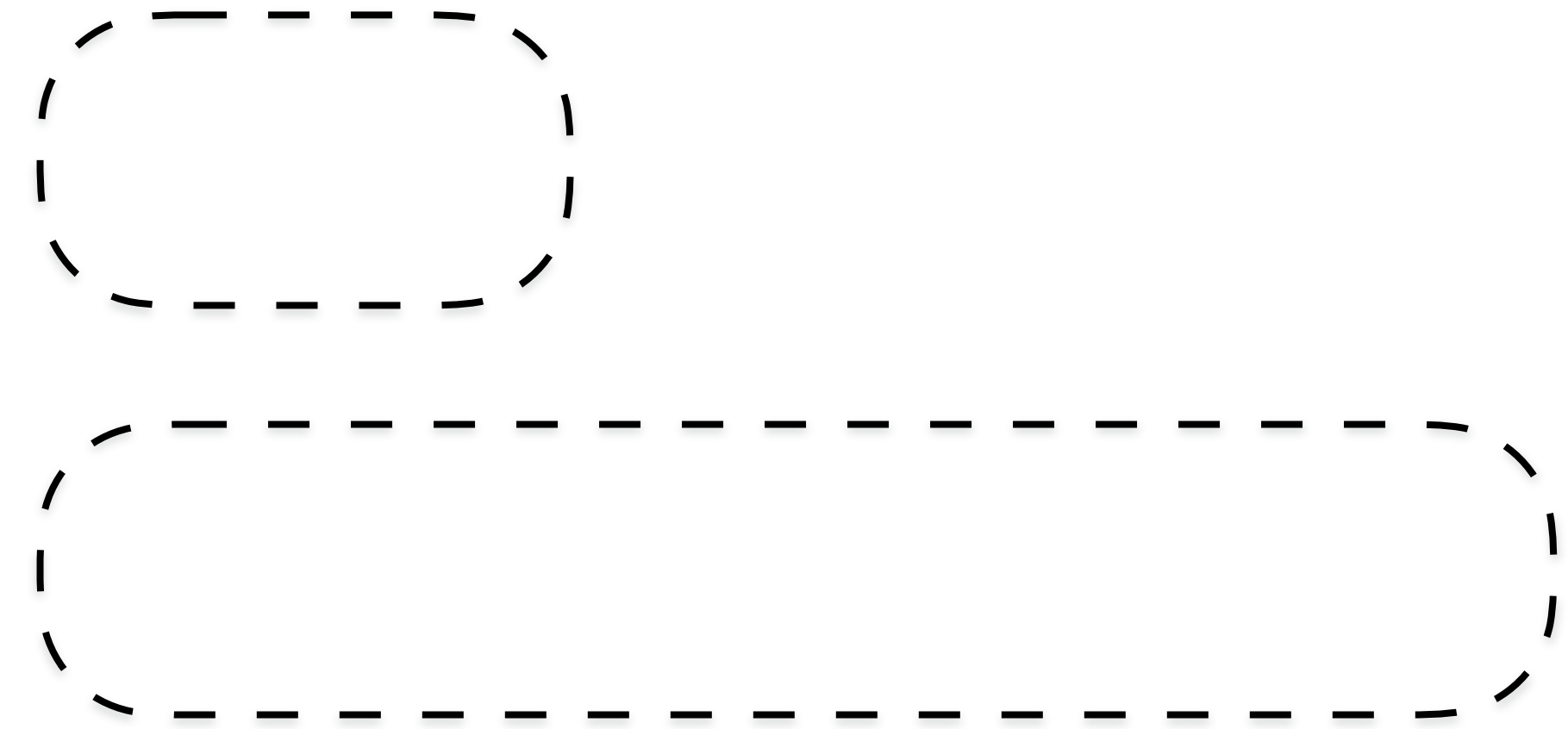
Leveling
read-optimized



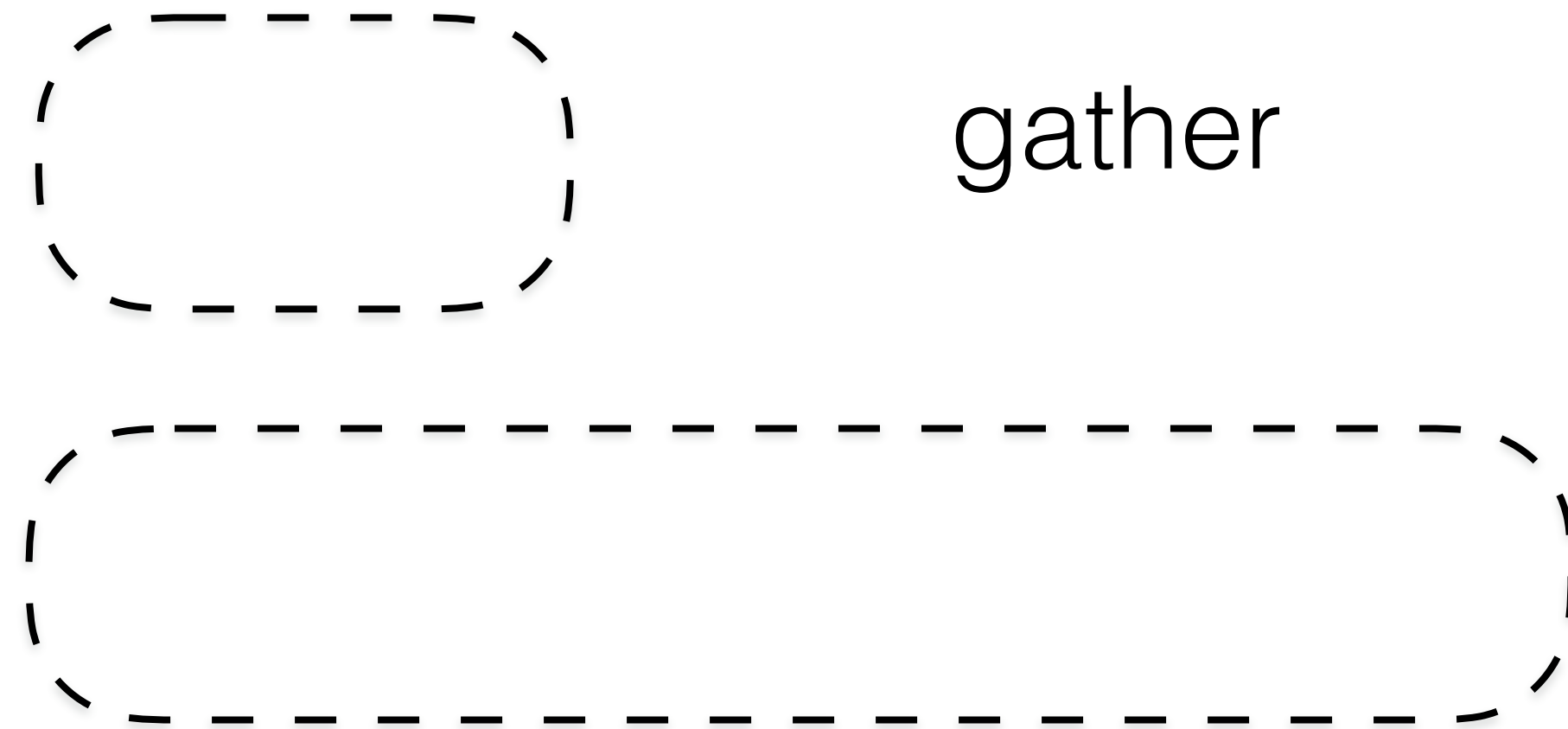
Tiering
write-optimized



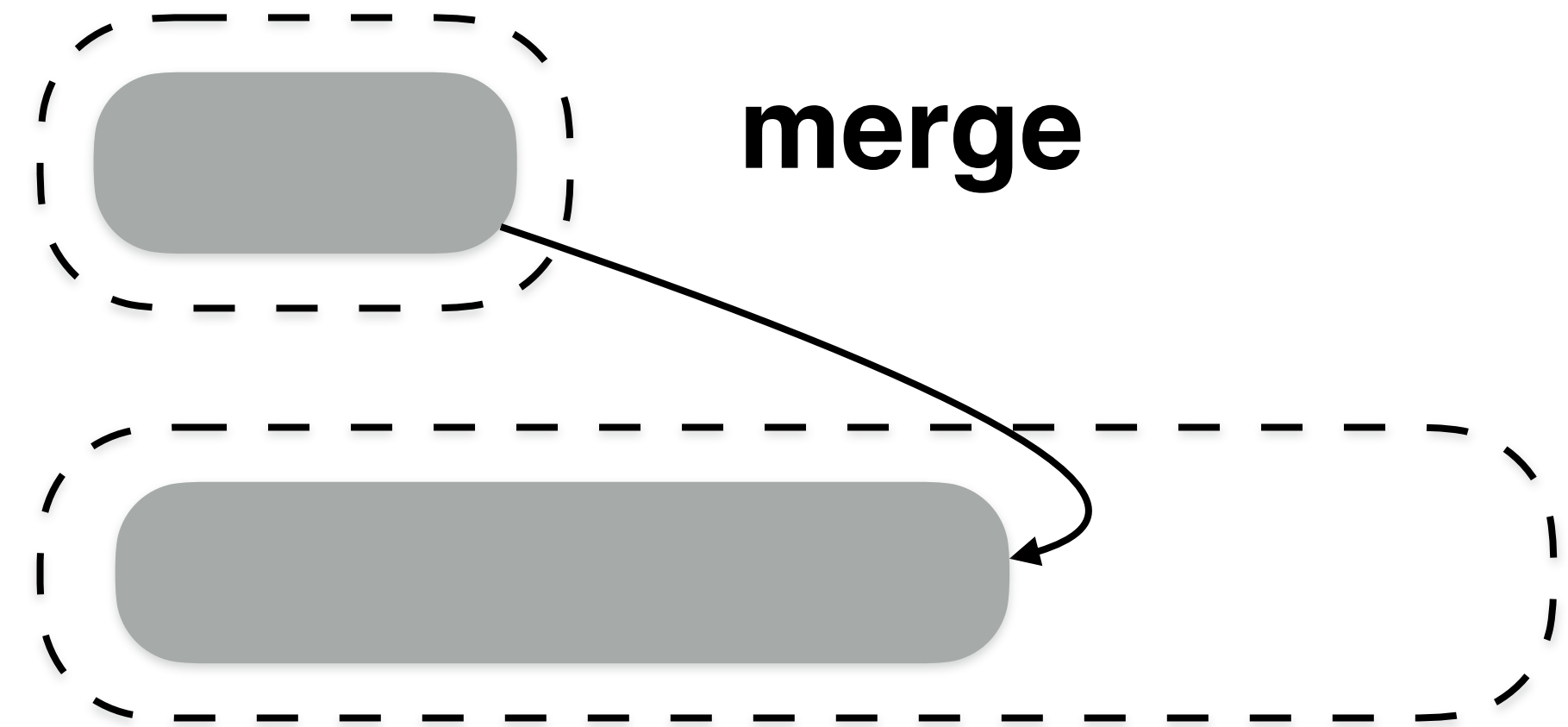
Leveling
read-optimized



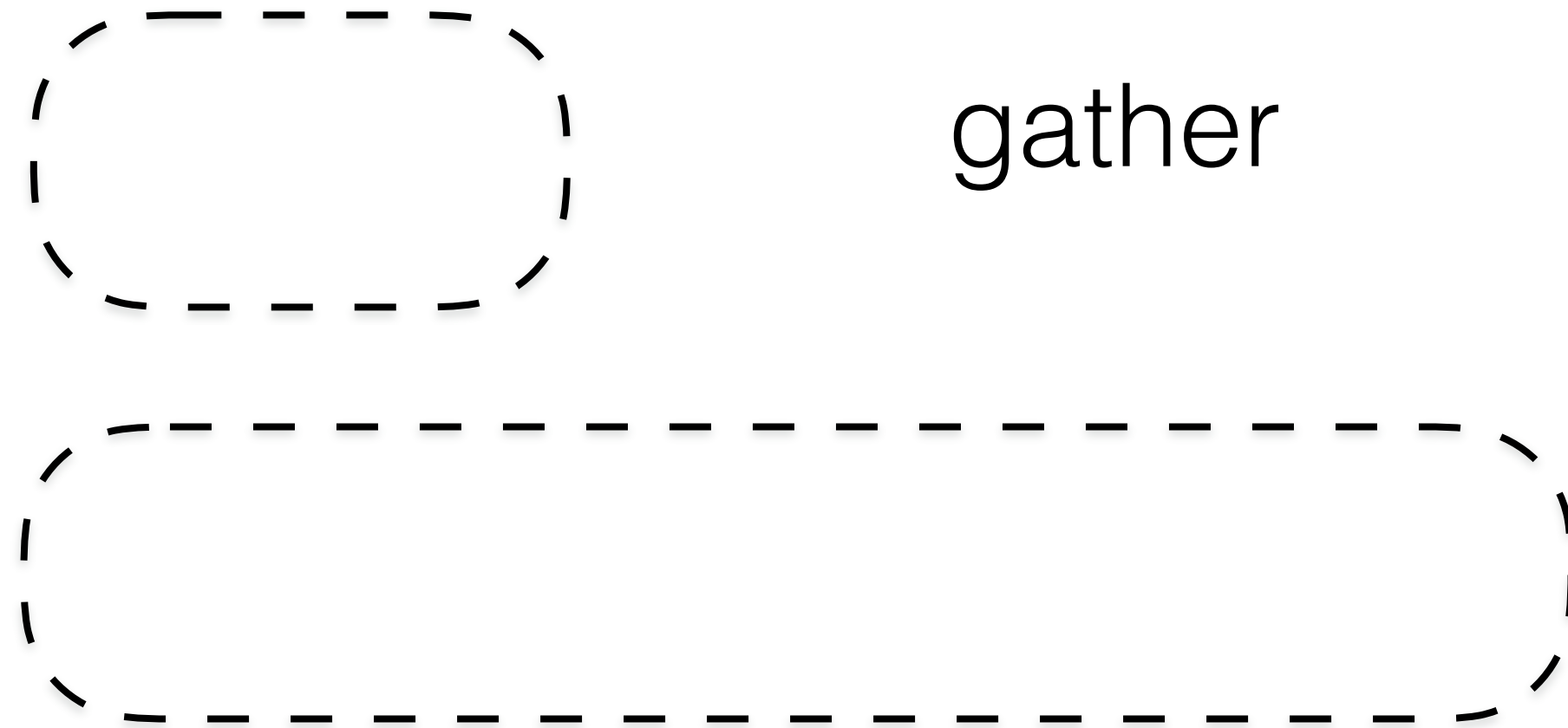
Tiering
write-optimized



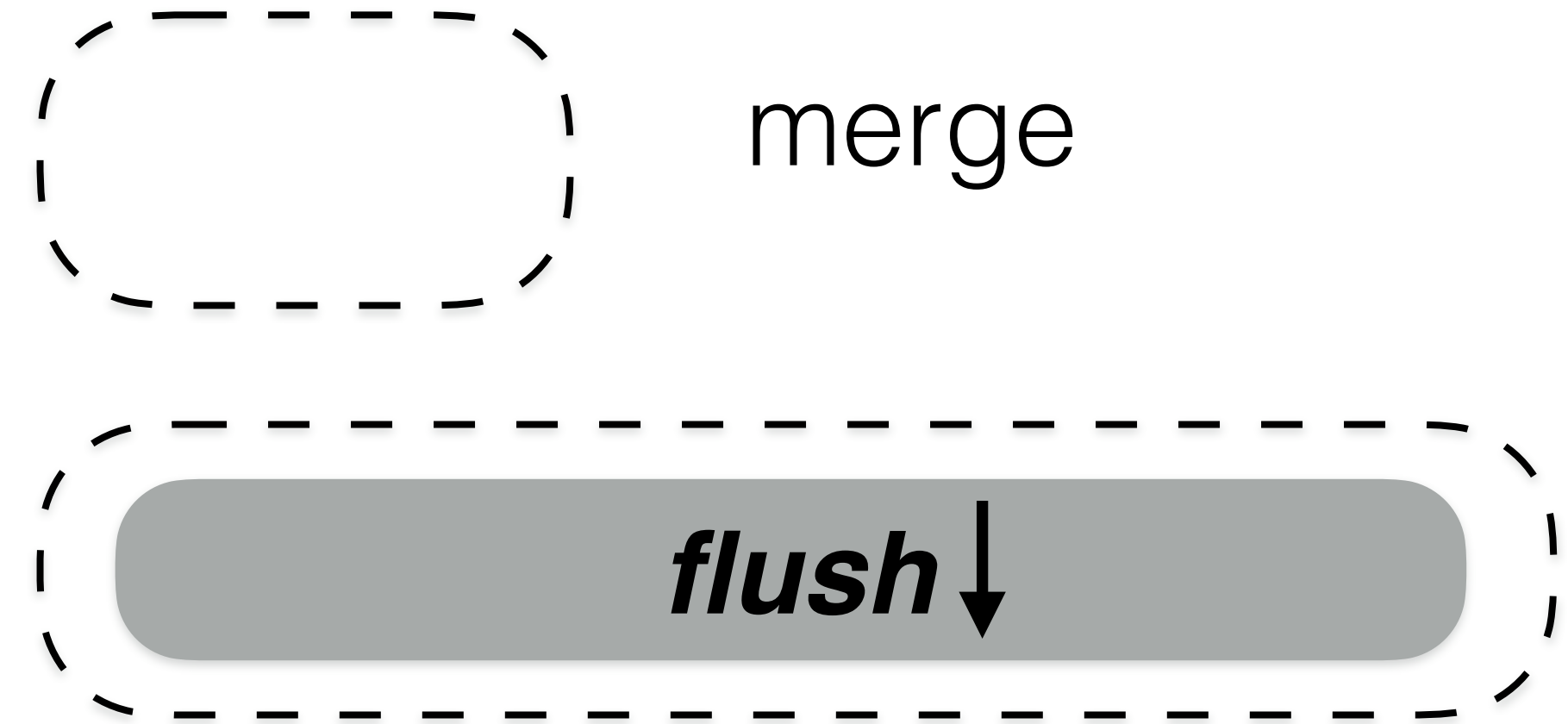
Leveling
read-optimized



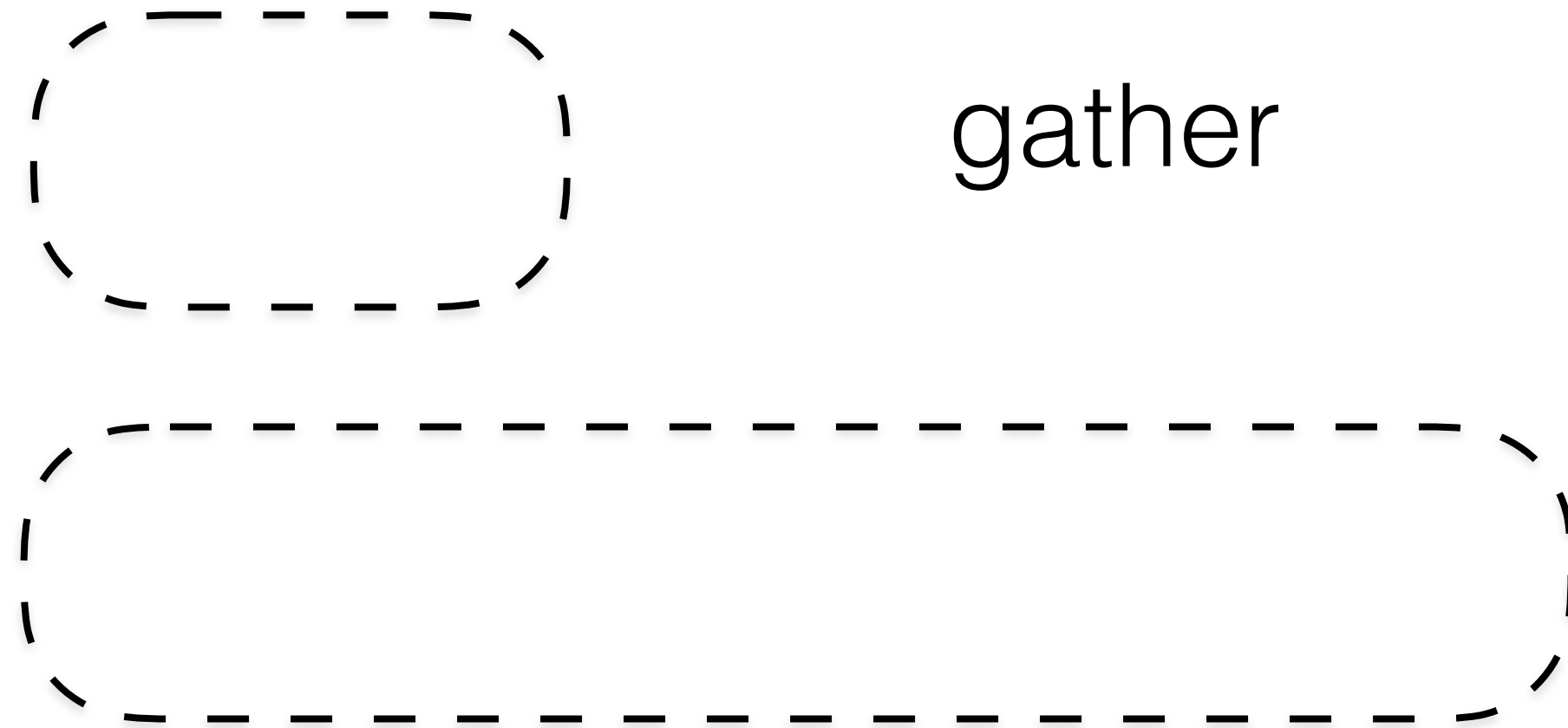
Tiering
write-optimized



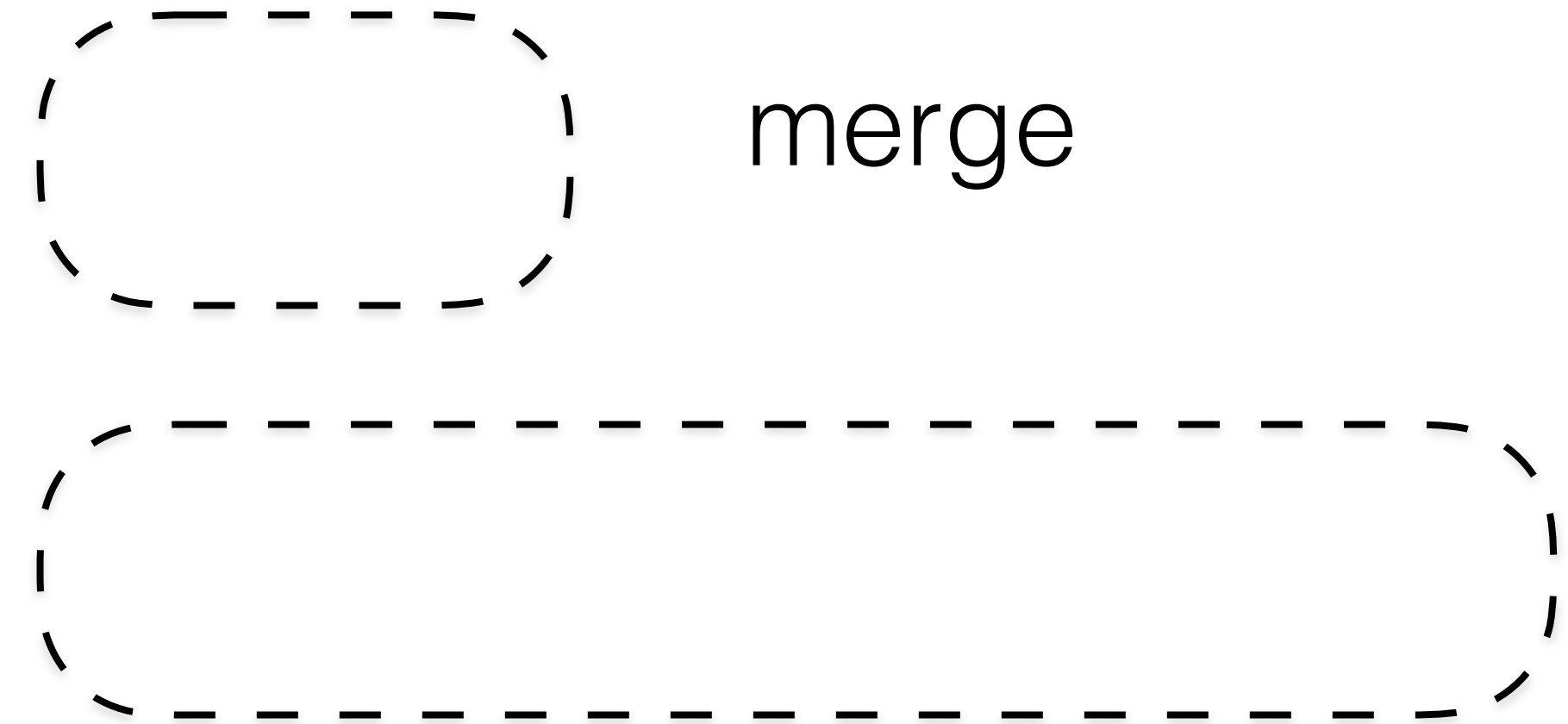
Leveling
read-optimized



Tiering
write-optimized

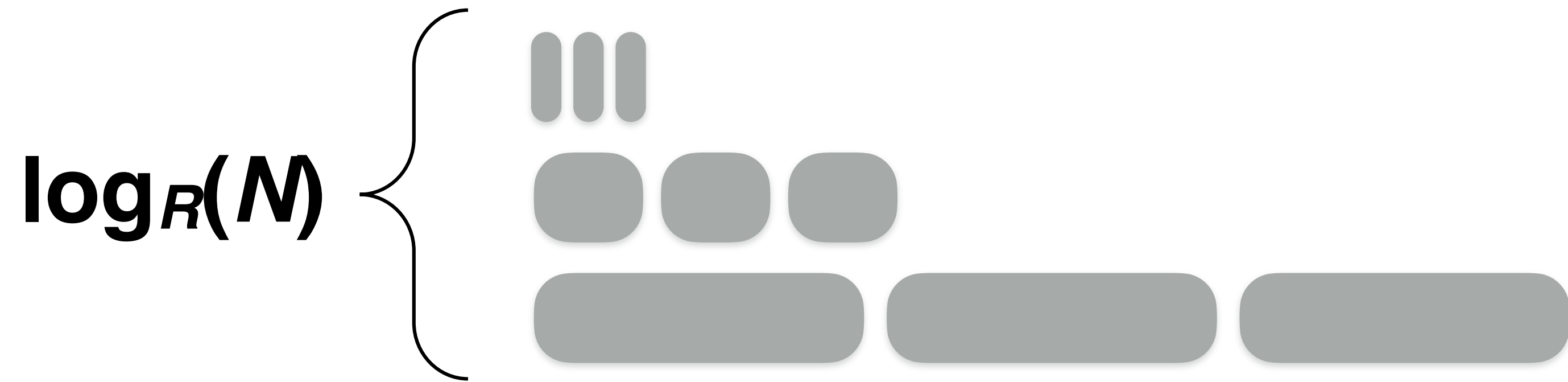


Leveling
read-optimized



Tiering
write-optimized

Leveling
read-optimized

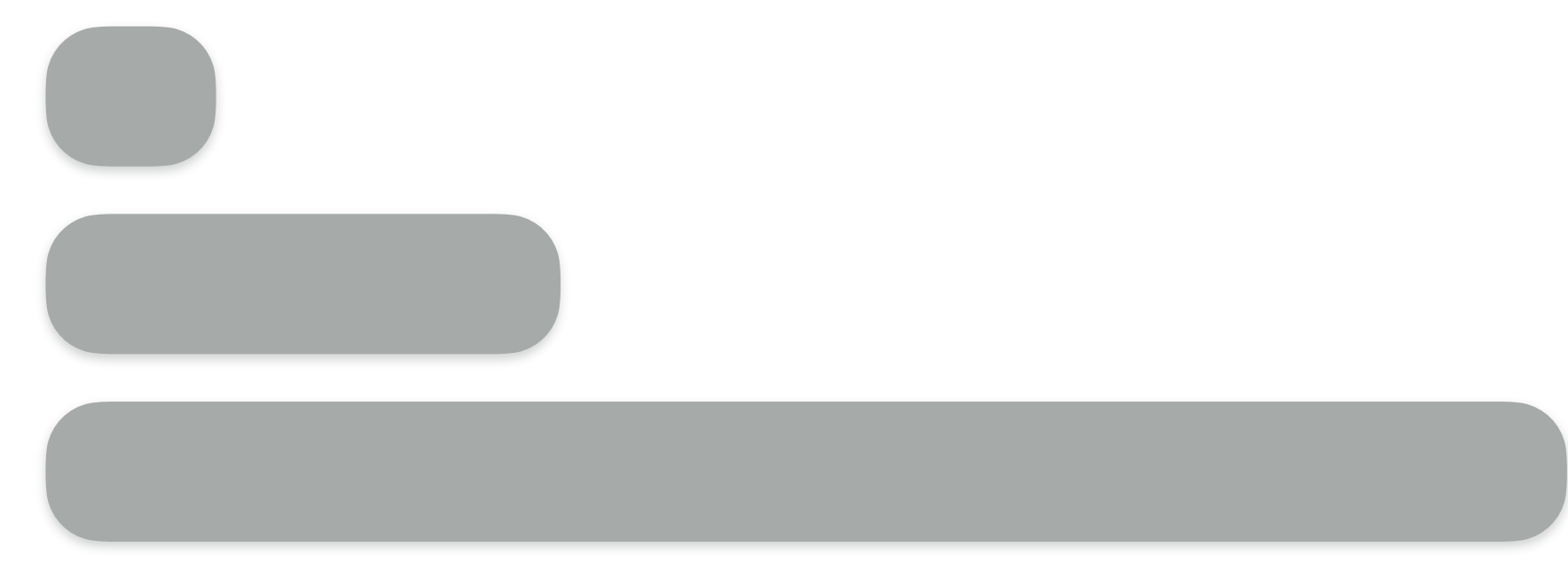
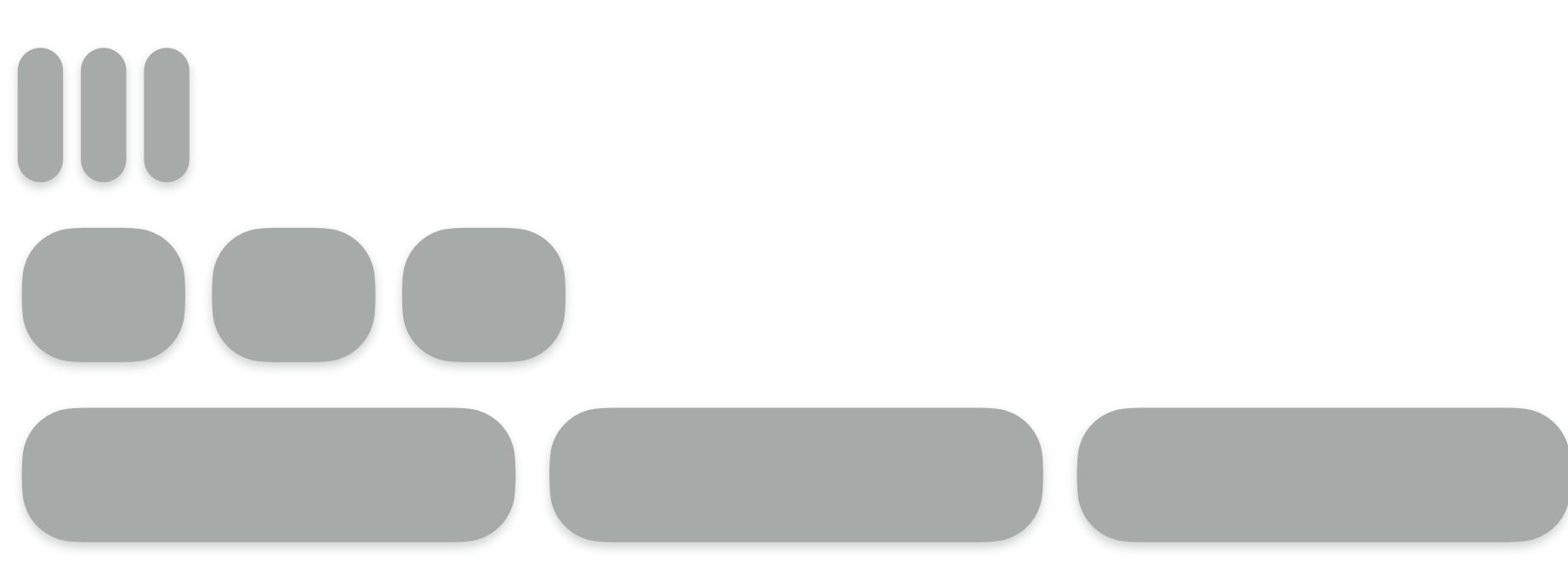


Tiering
write-optimized

Leveling
read-optimized

$\log_R(N)$

size ratio



Tiering
write-optimized

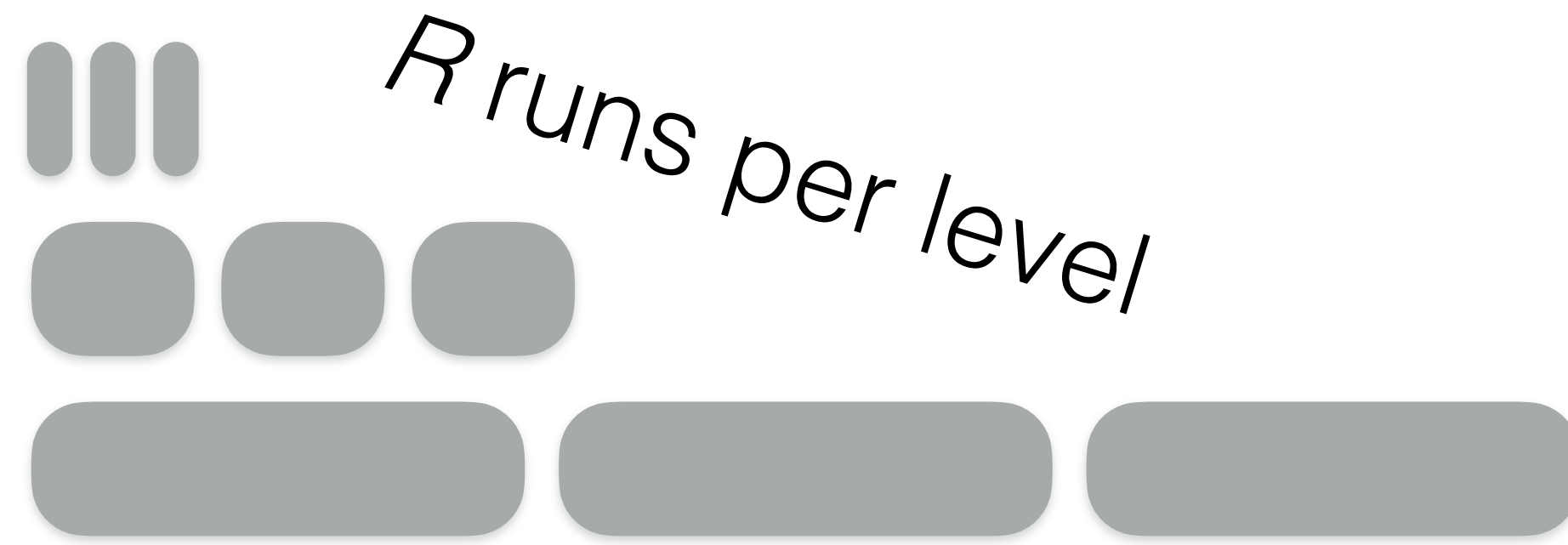
Leveling
read-optimized

$\log_R(N)$



size ratio

Tiering
write-optimized

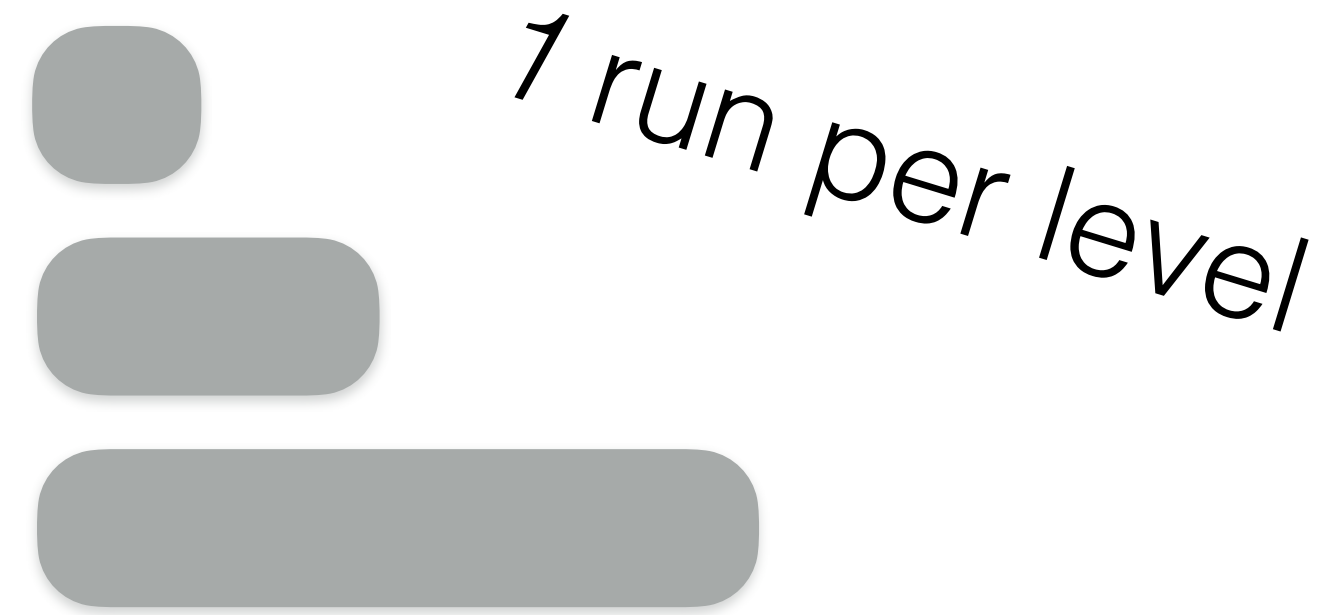


Leveling
read-optimized

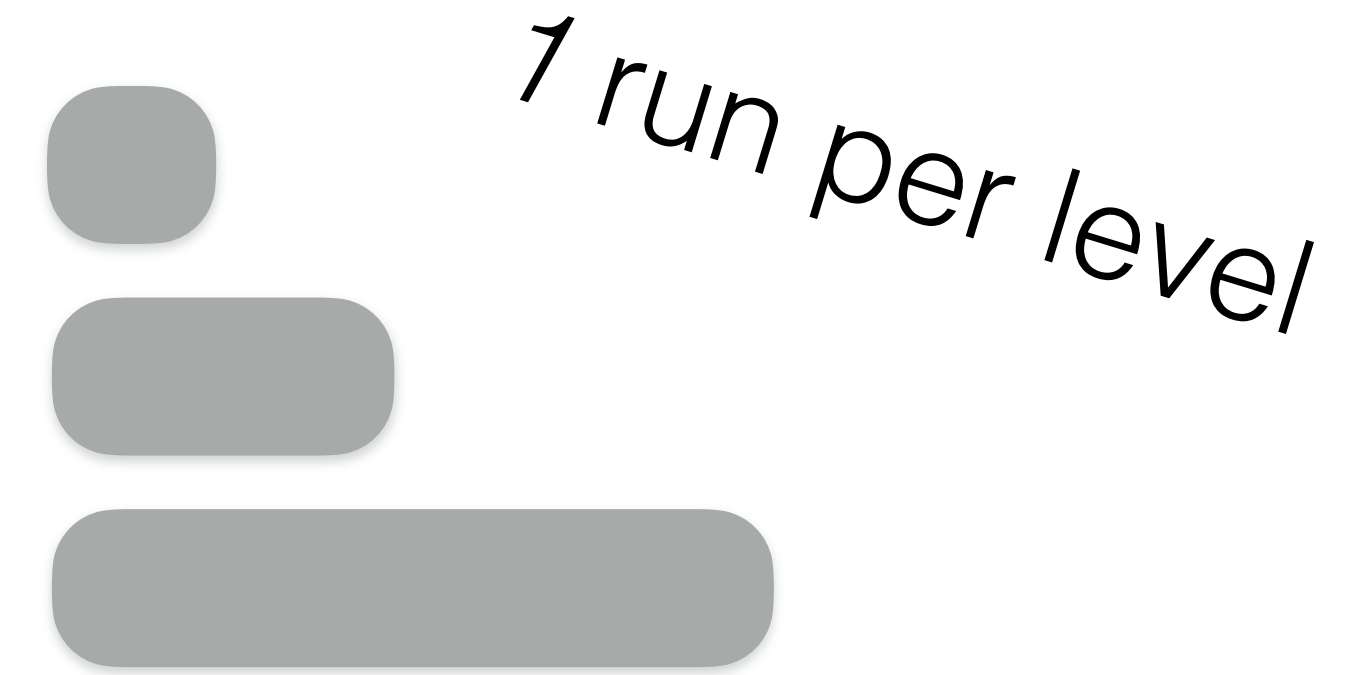


 size ratio R

Tiering
write-optimized

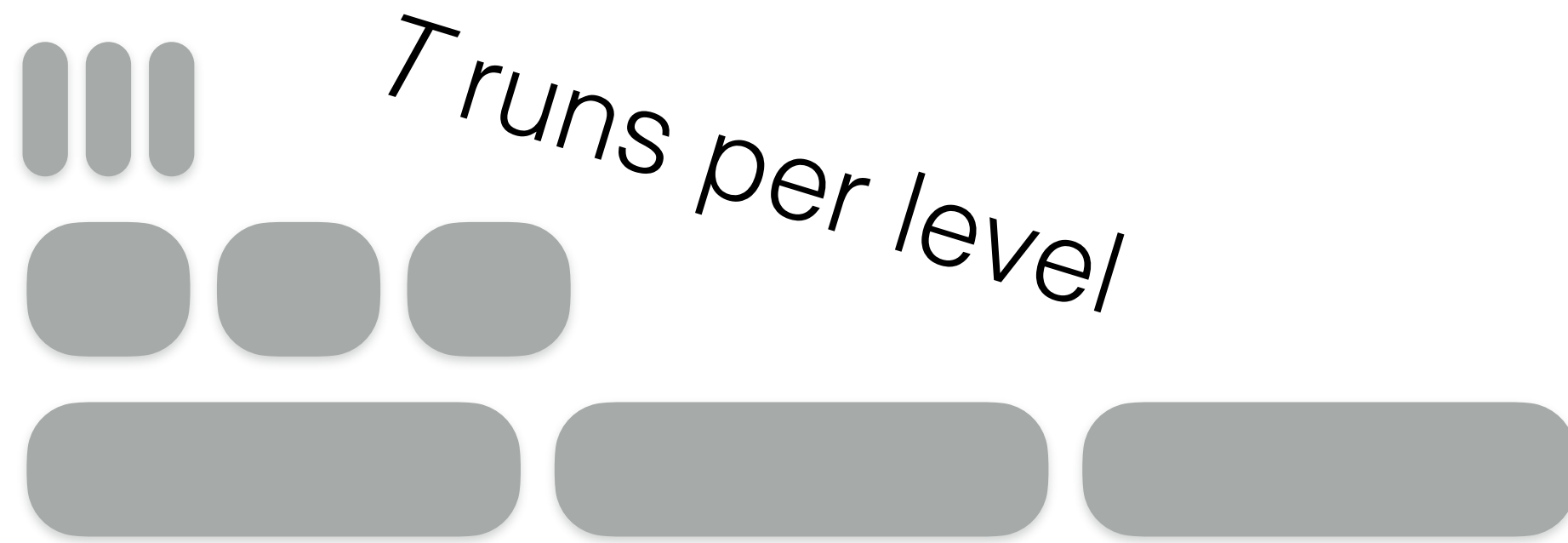


Leveling
read-optimized



 size ratio $R \gg$

Tiering
write-optimized



Leveling
read-optimized



 size ratio R

Tiering
write-optimized

$O(N)$ runs per level

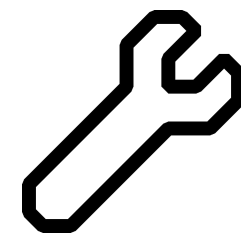

log

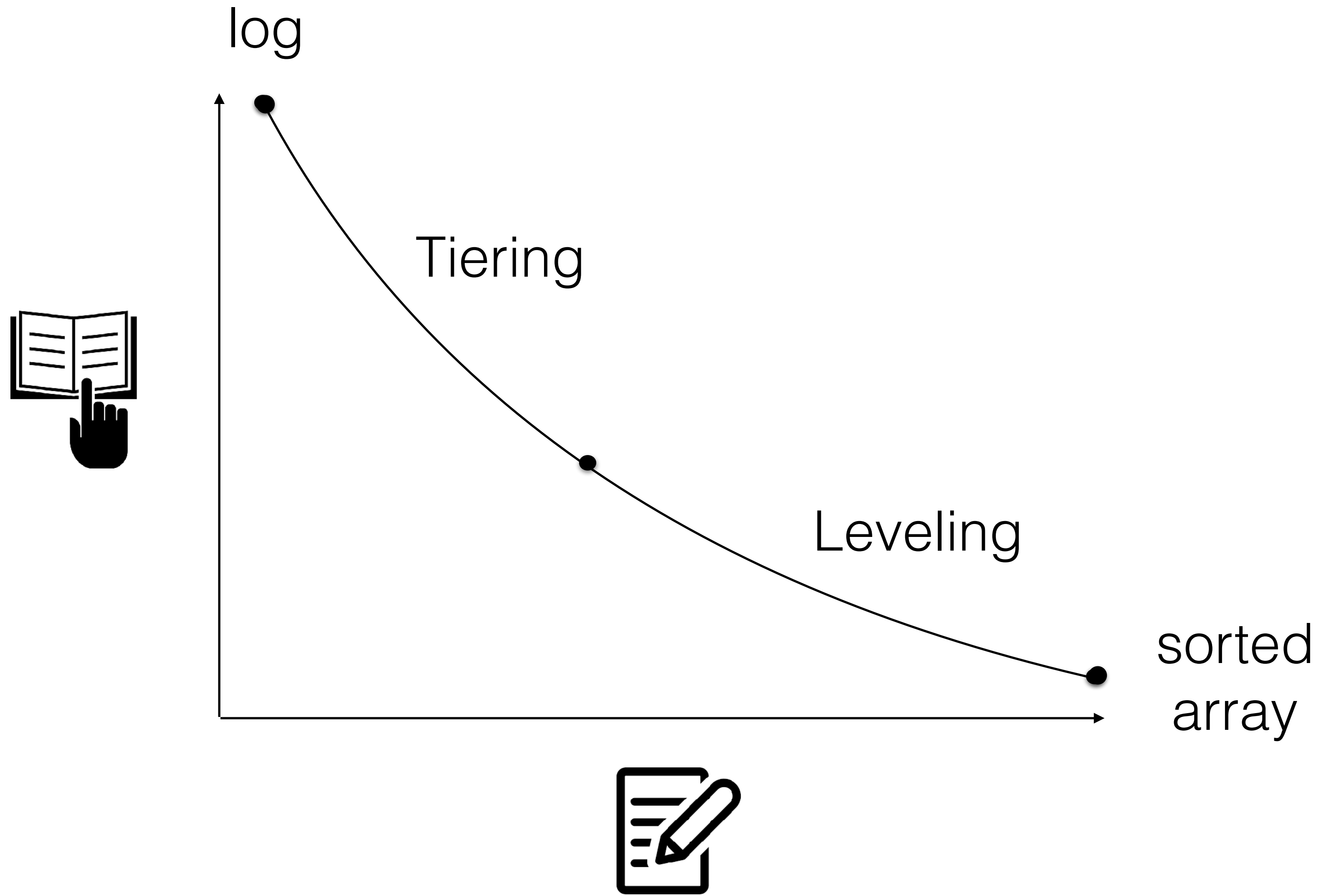
Leveling
read-optimized

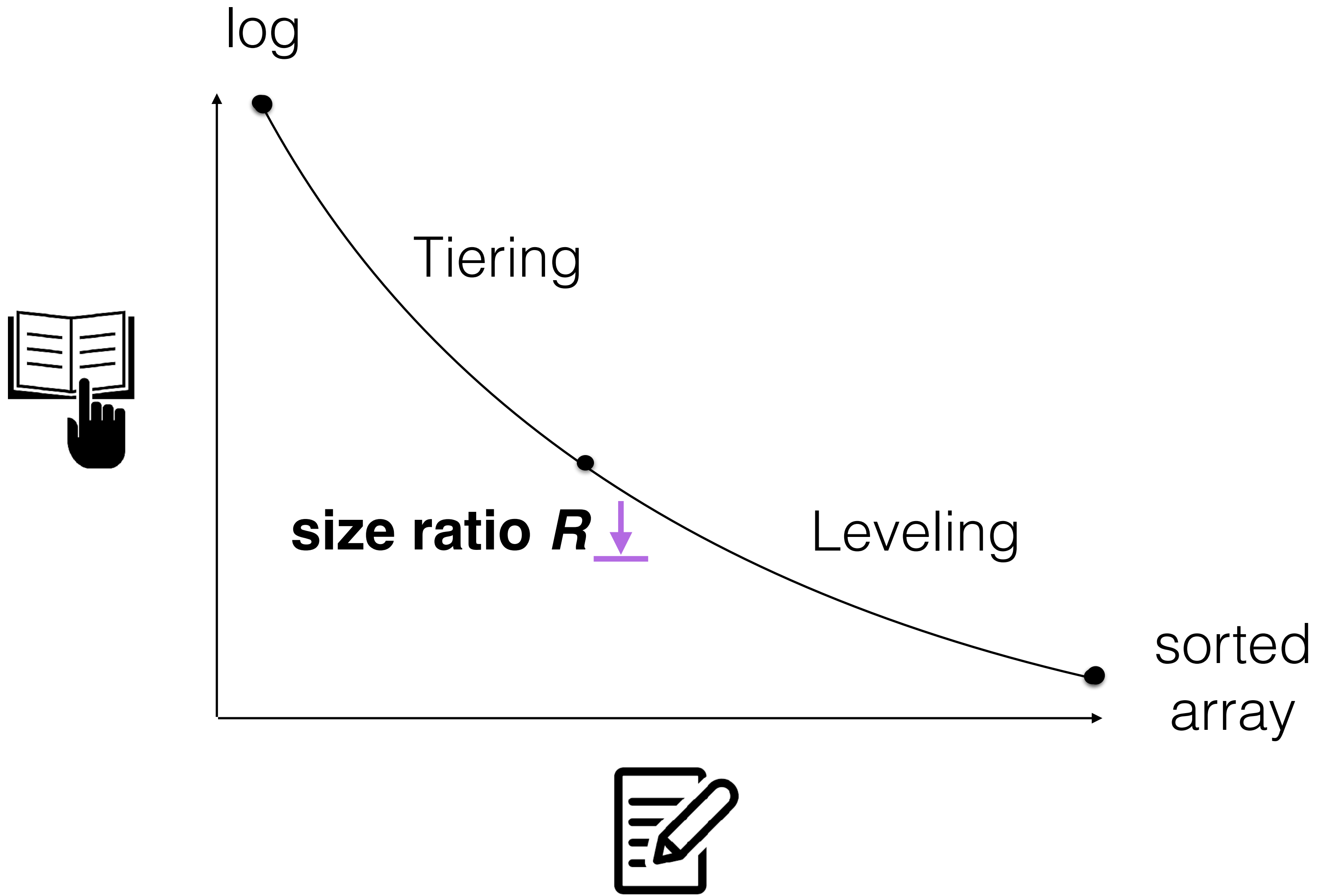
1 run per level

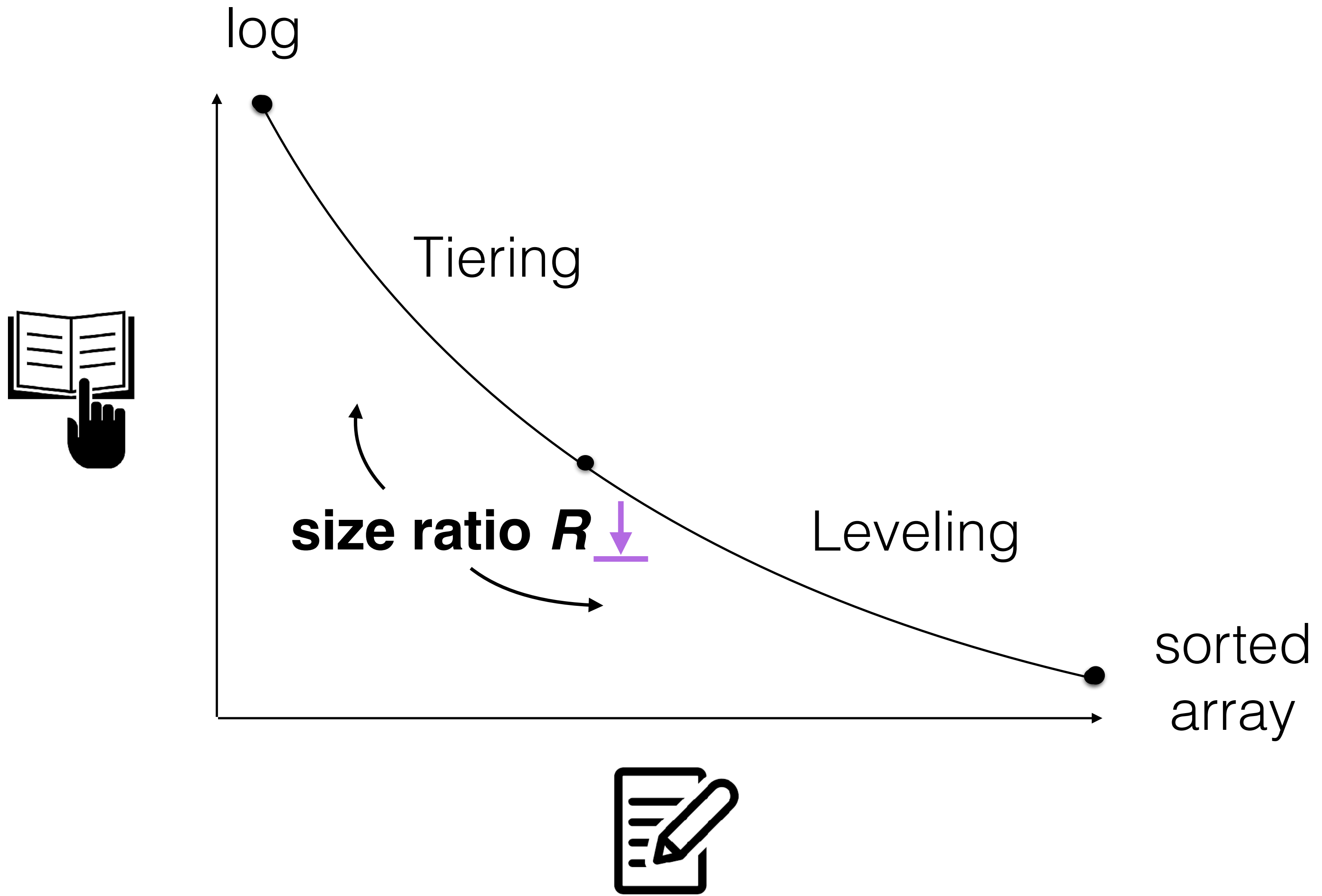

**sorted
array**

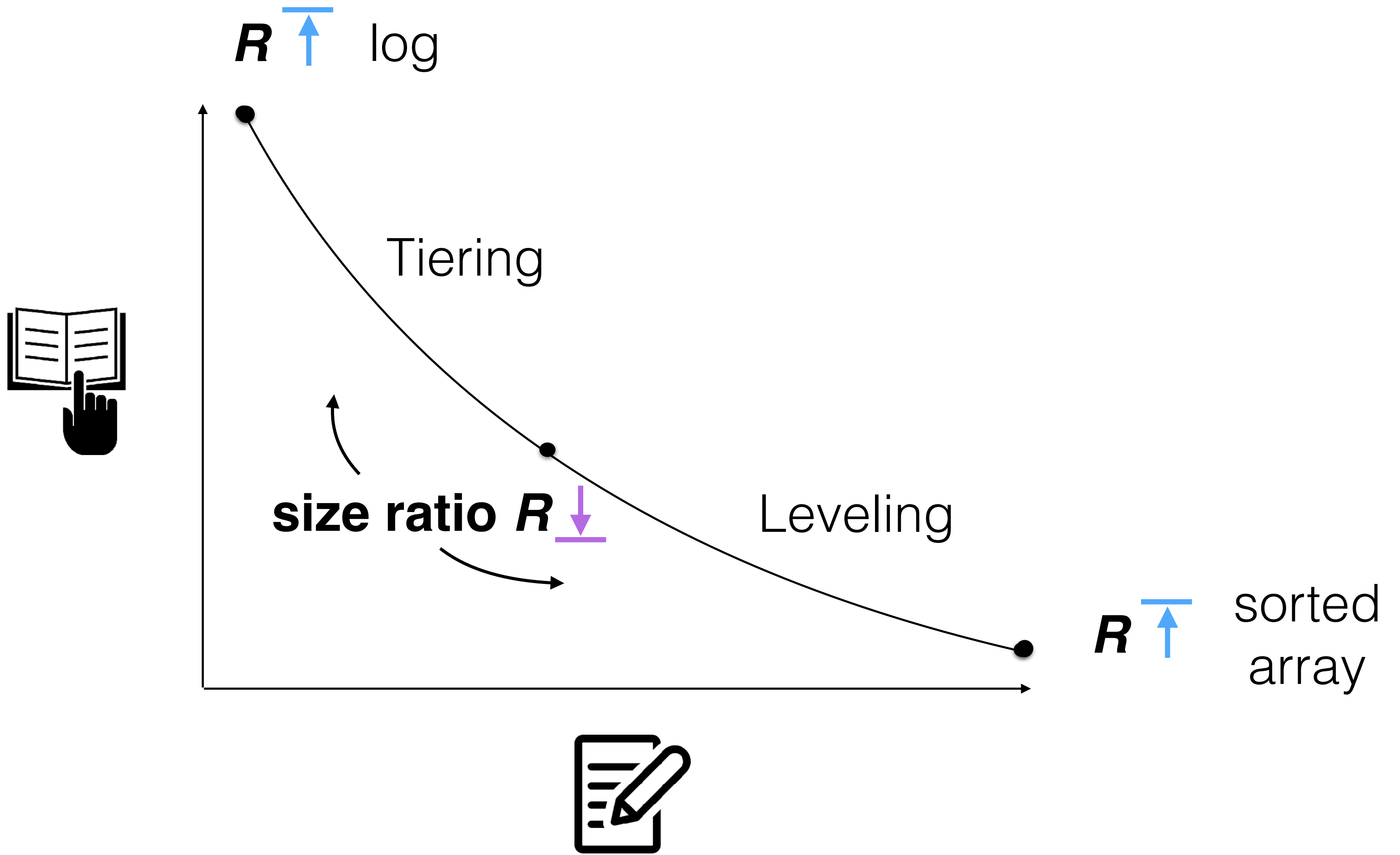
size ratio $R \gg$

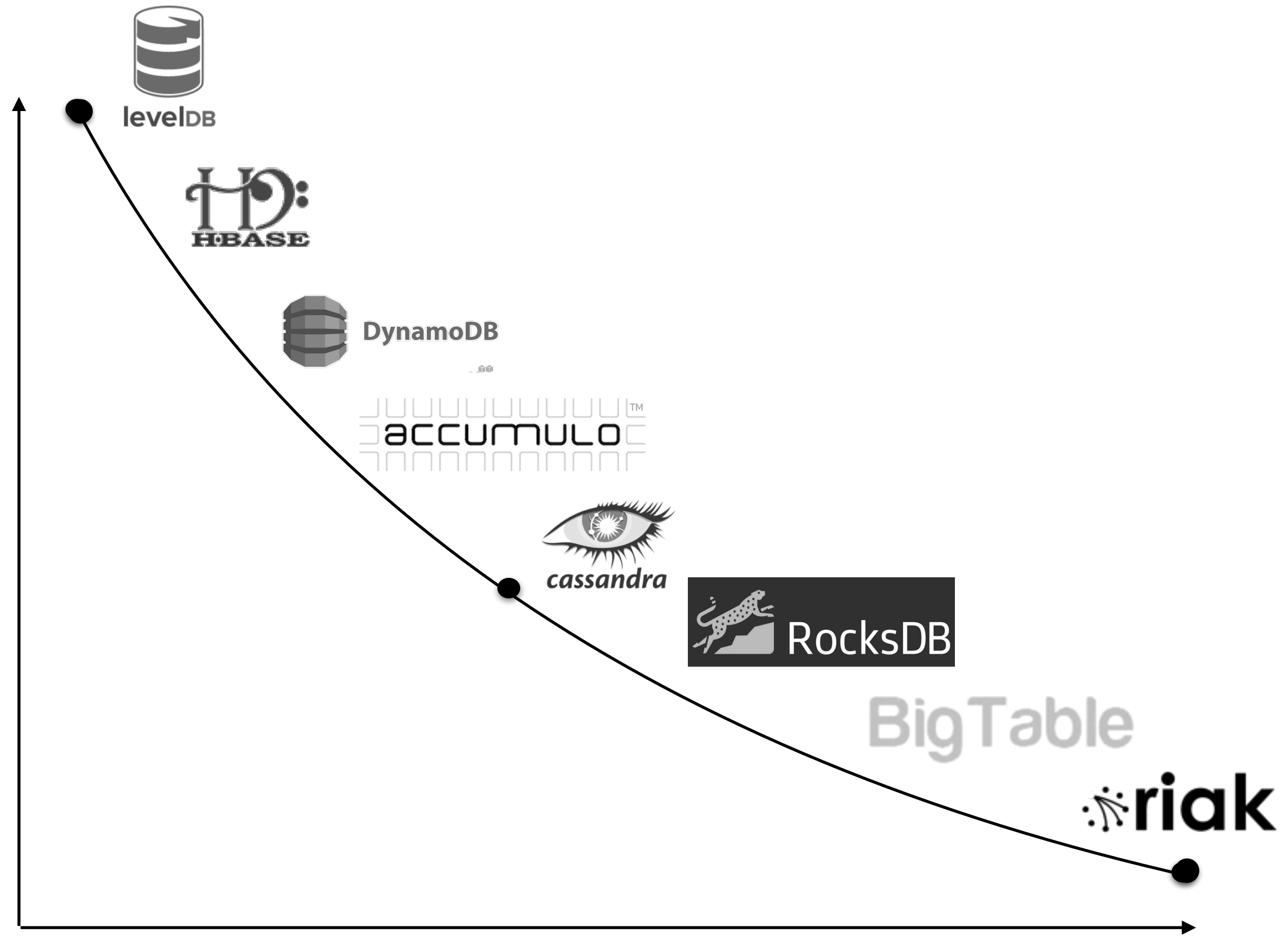


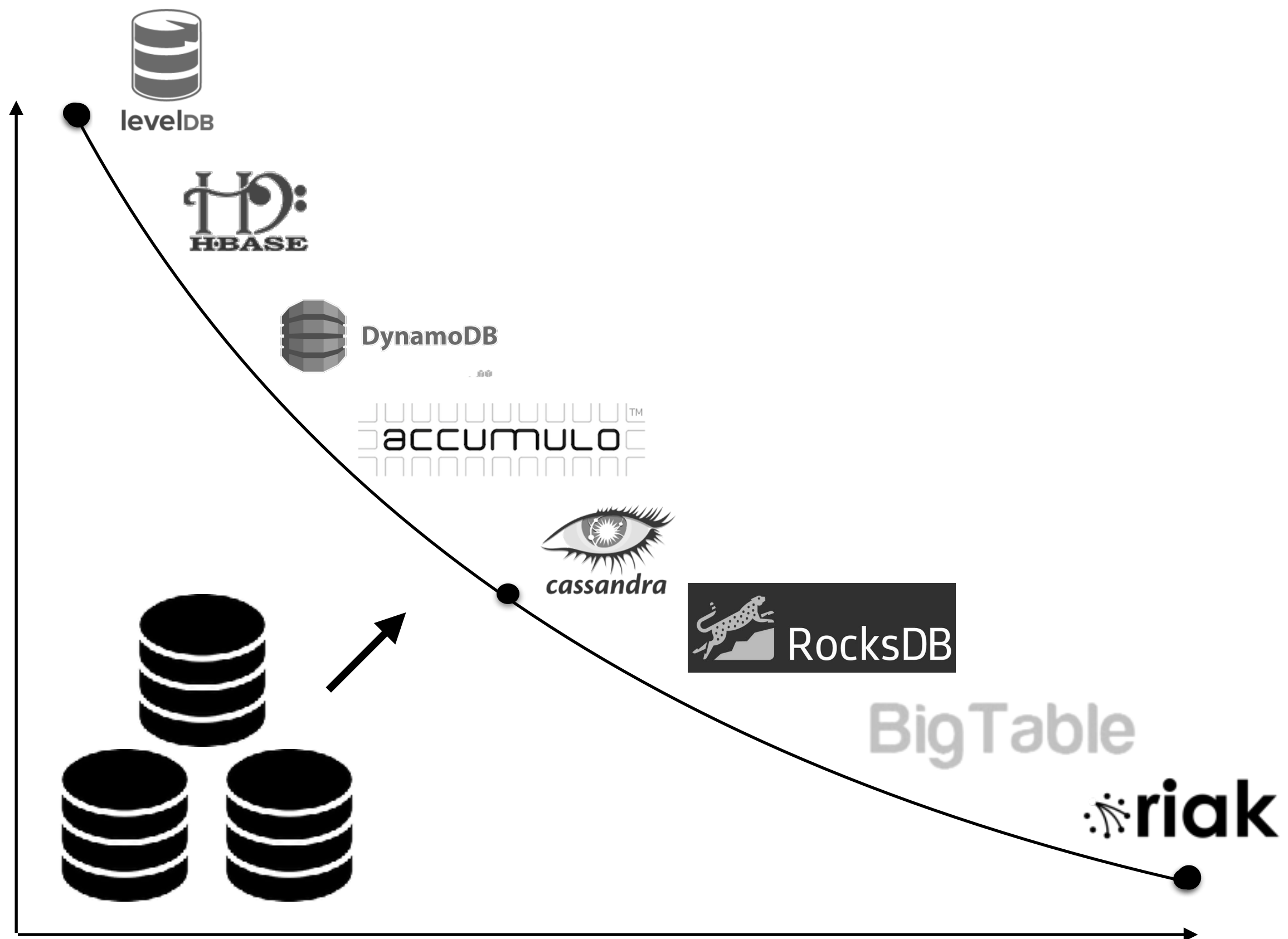
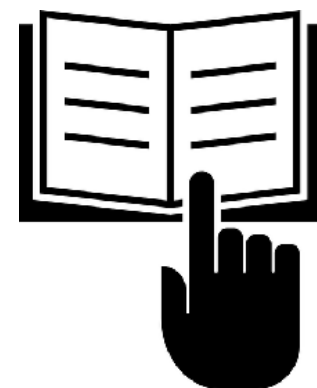


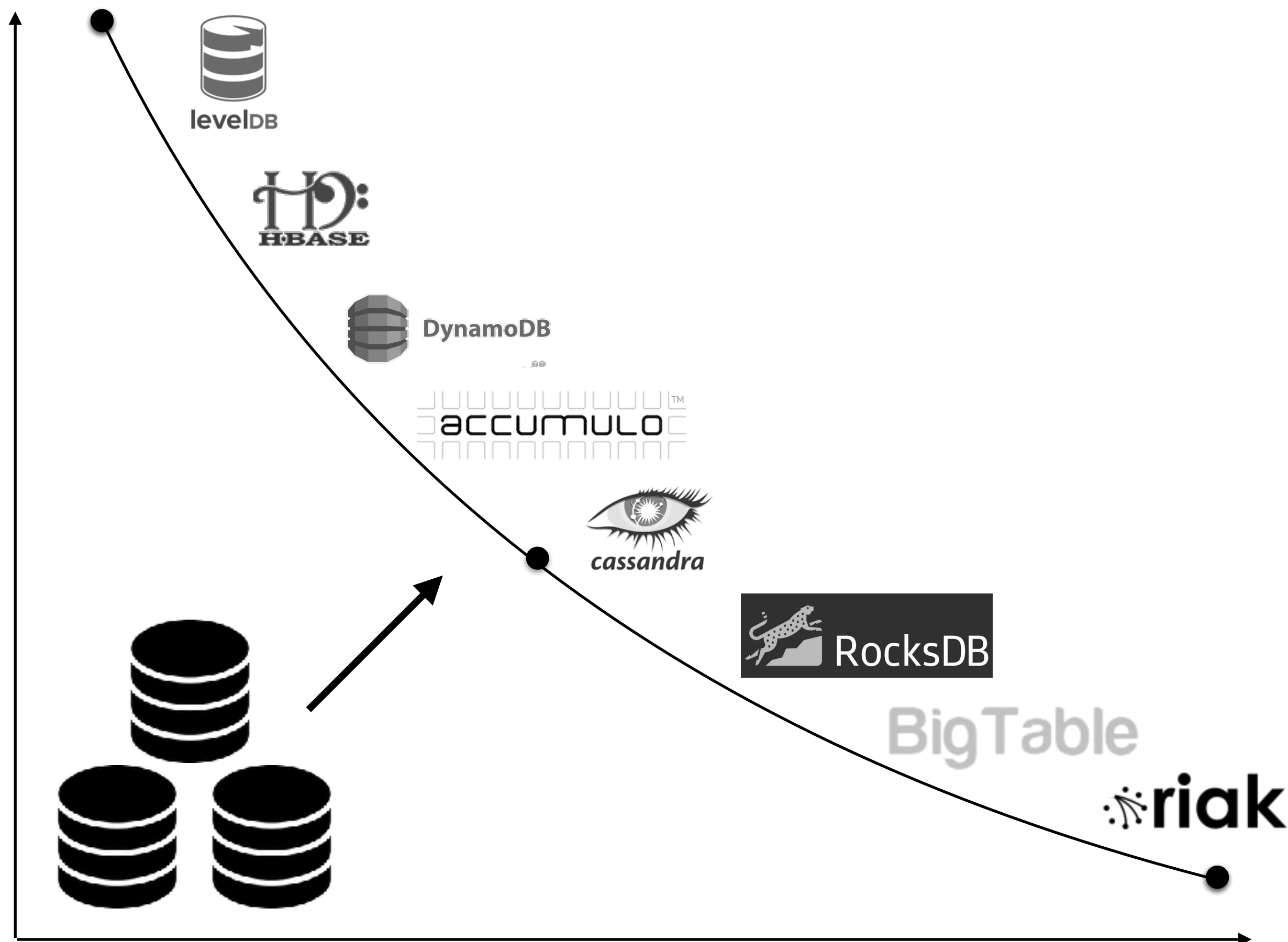
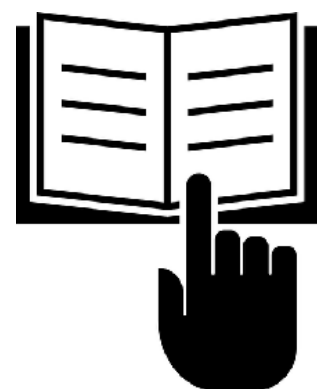


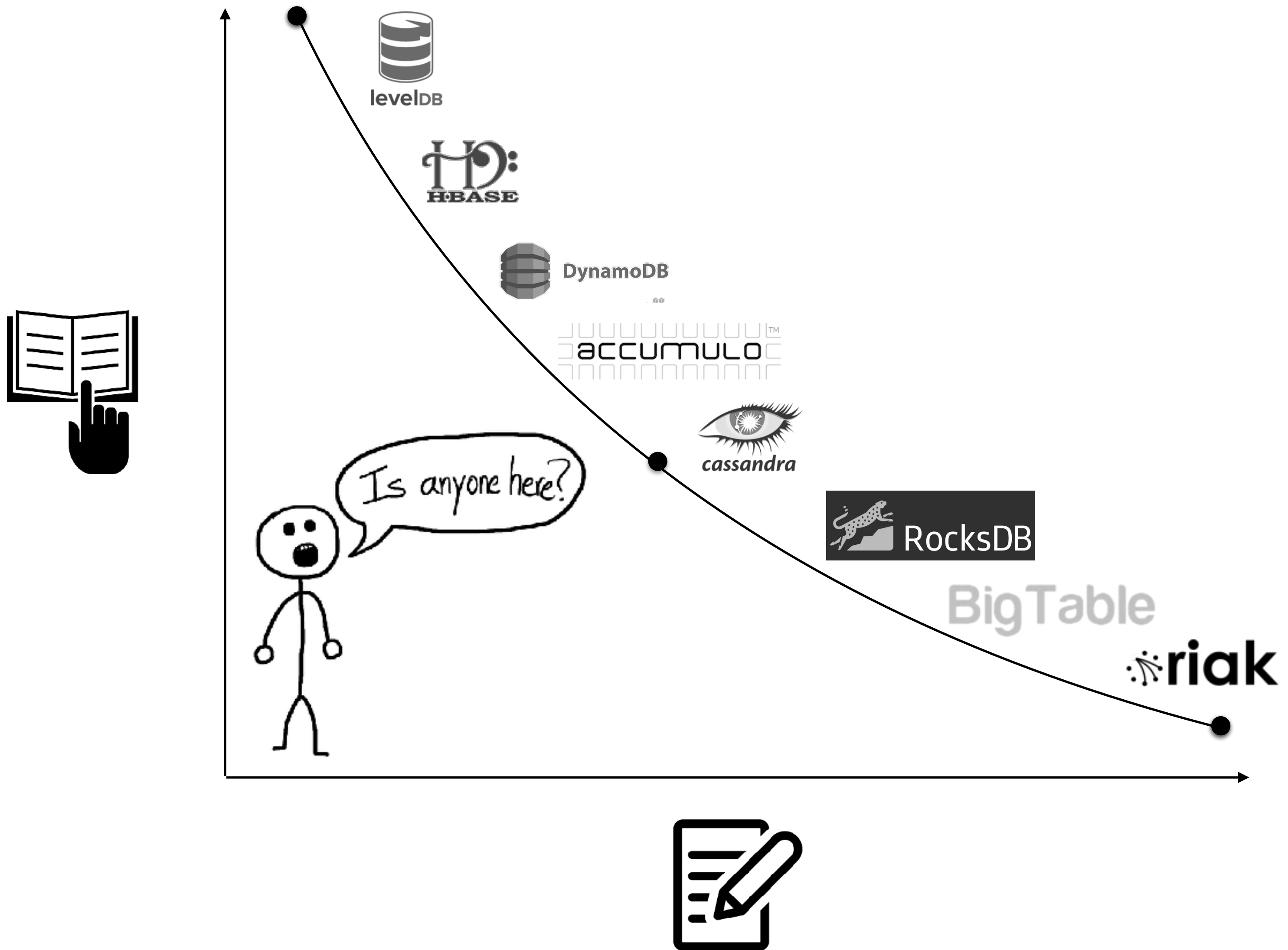


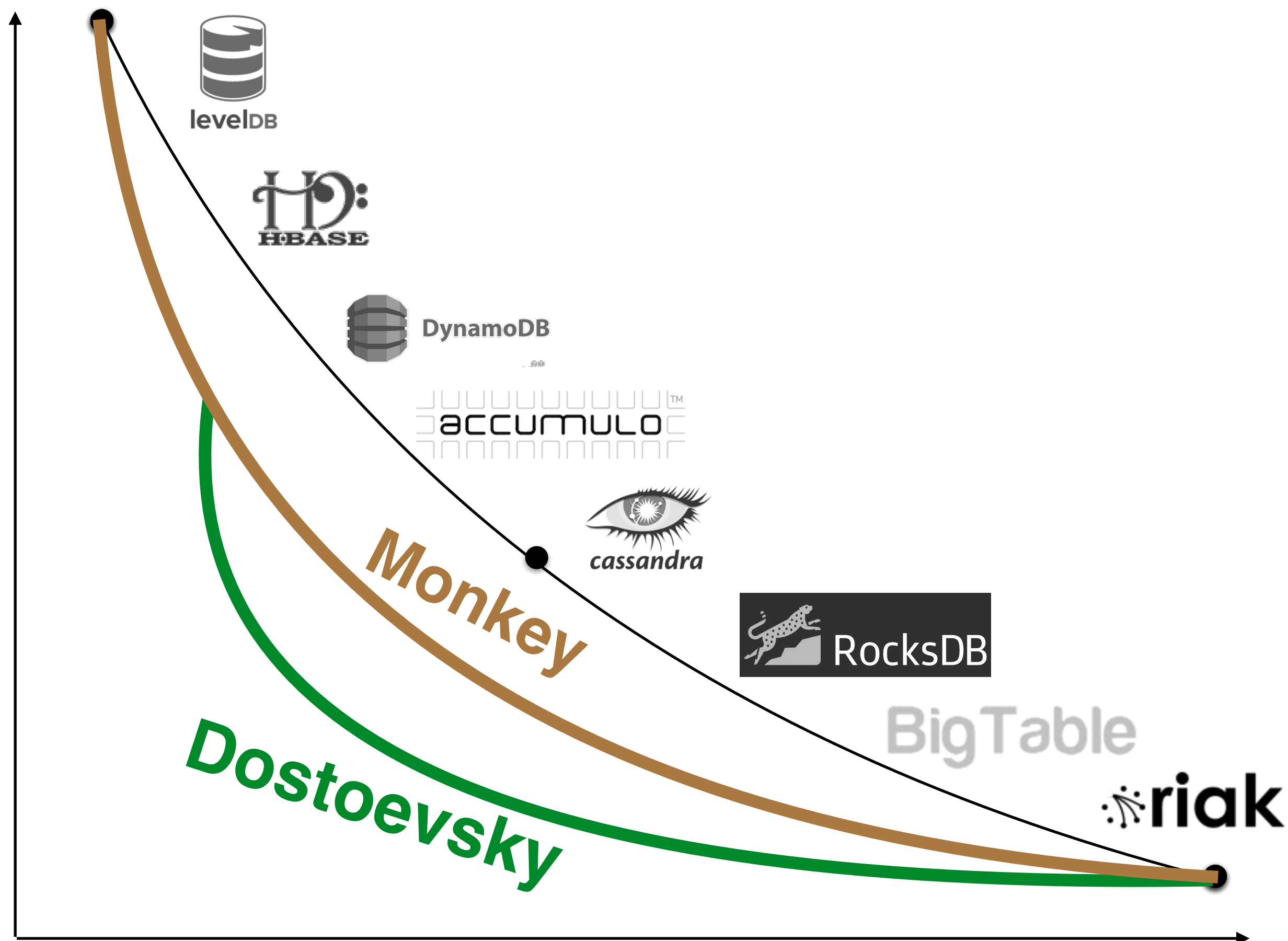












Monkey: Optimal Navigable Key-Value Store

SIGMOD17



Monkey: Optimal Navigable Key-Value Store

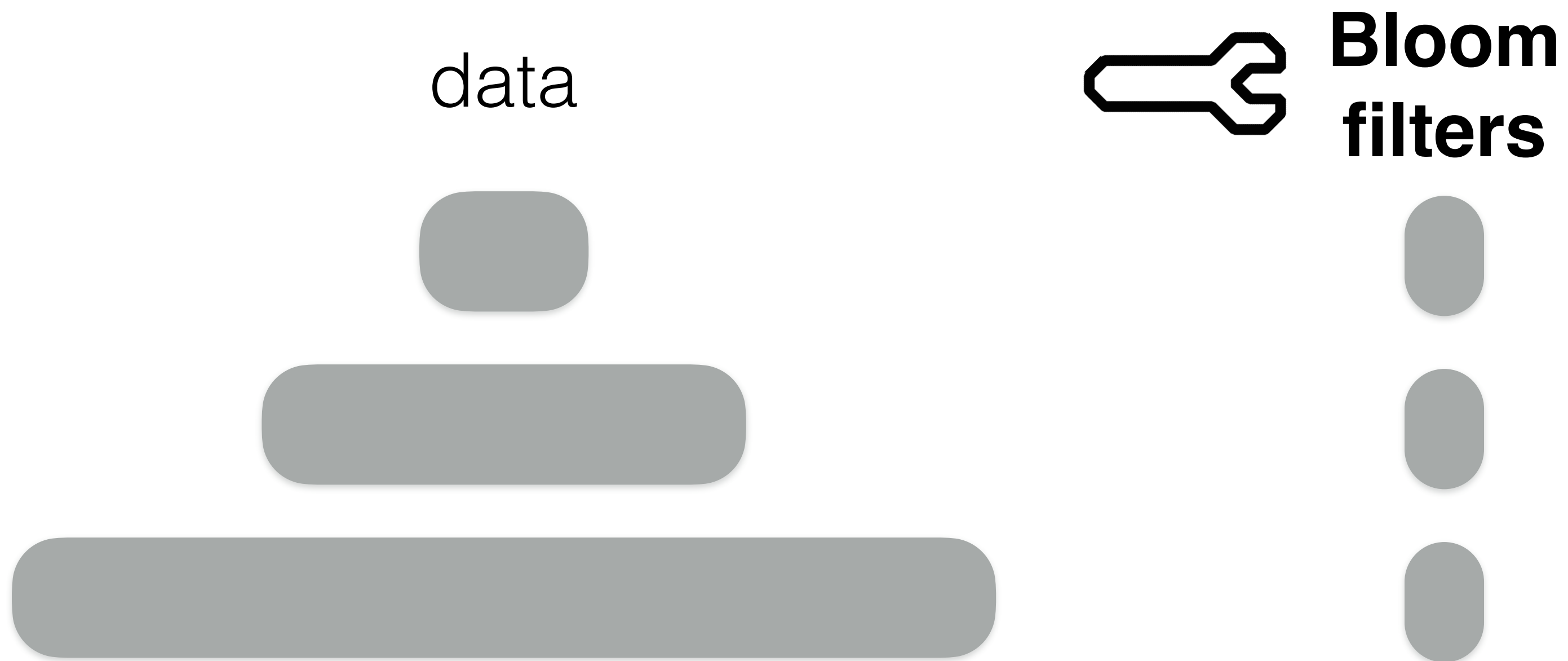
SIGMOD17



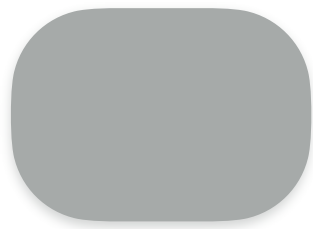
Niv Dayan
Manos Athanassoulis
Stratos Idreos

Monkey: Optimal Navigable Key-Value Store

SIGMOD17



data



Bloom
filters



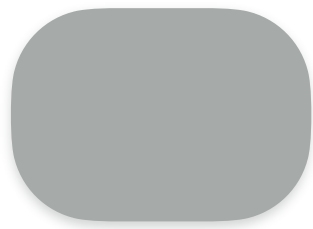
bits/entry

x

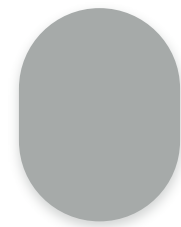
x

x

data



Bloom
filters



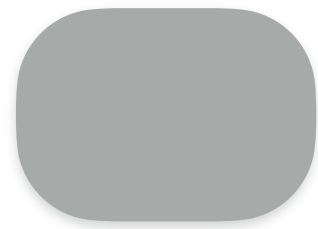
bits/entry

x

x

x

data



Bloom
filters



**false
positive rate**

$$O(e^{-x})$$

$$O(e^{-x})$$

$$O(e^{-x})$$

Bloom
filters

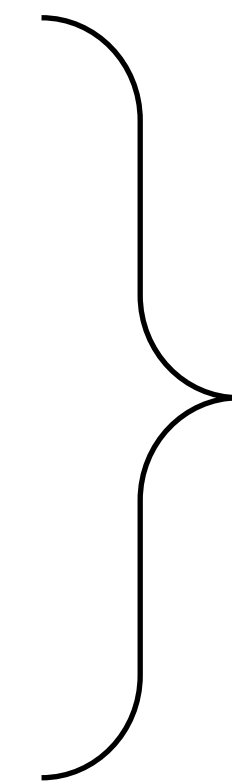


false
positive rate

$$O(e^{-x})$$

$$O(e^{-x})$$

$$O(e^{-x})$$



$$= O(\mathbf{e^{-x}} \cdot \mathbf{\log_R(N)})$$

Bloom
filters

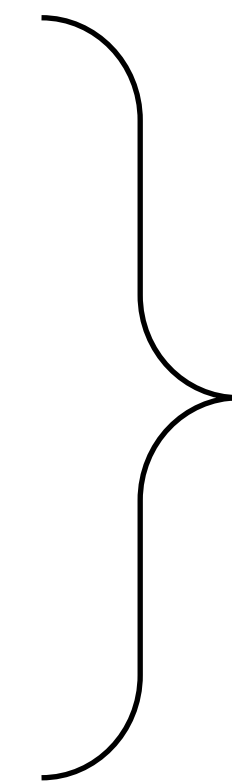


false
positive rate

$$O(e^{-x})$$

$$O(e^{-x})$$

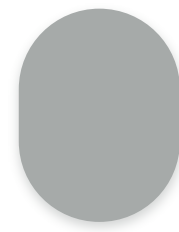
$$O(e^{-x})$$



$$= O(\mathbf{e^{-x} \cdot \log_R(N)})$$



Bloom
filters



most memory

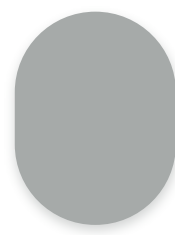
false
positive rate

$$O(e^{-x})$$

$$O(e^{-x})$$

$$O(e^{-x})$$

Bloom
filters



most memory



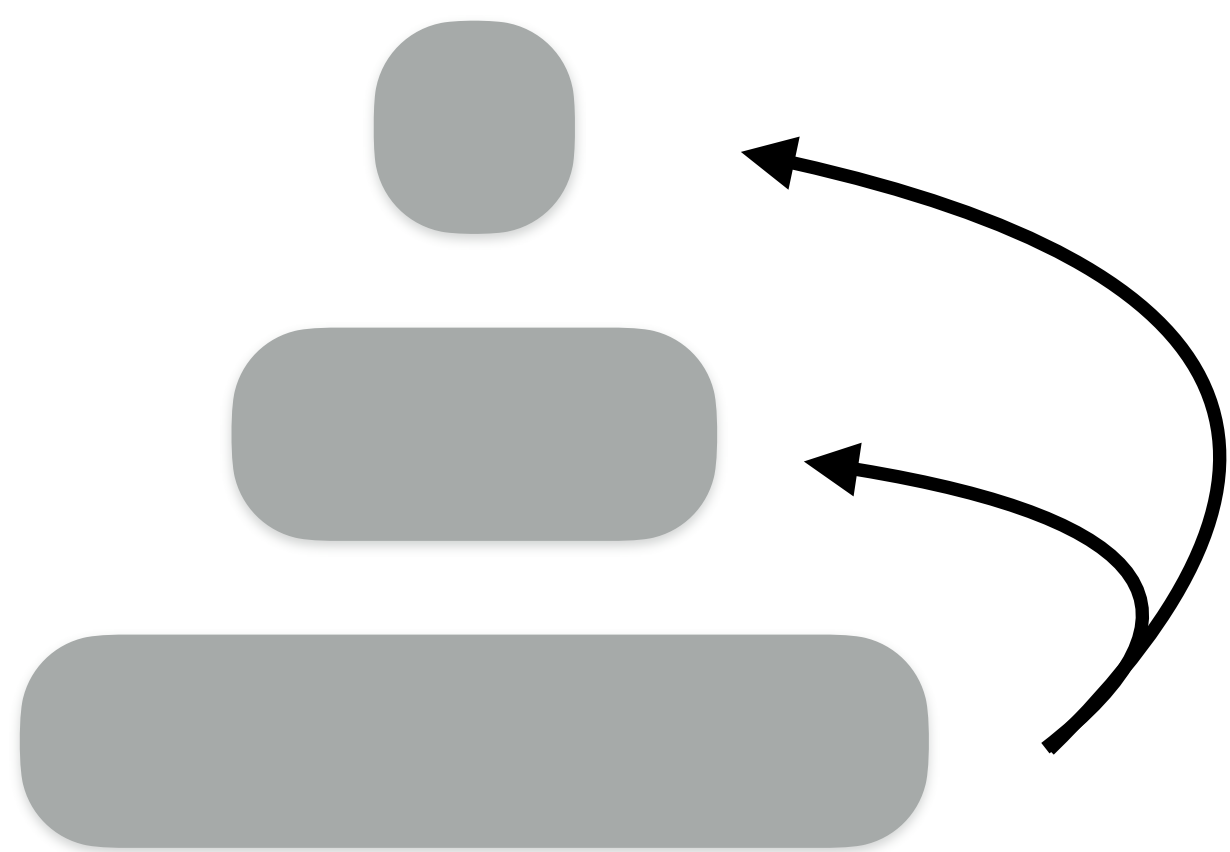
saves at most 1 I/O!

false
positive rate

$$O(e^{-x})$$

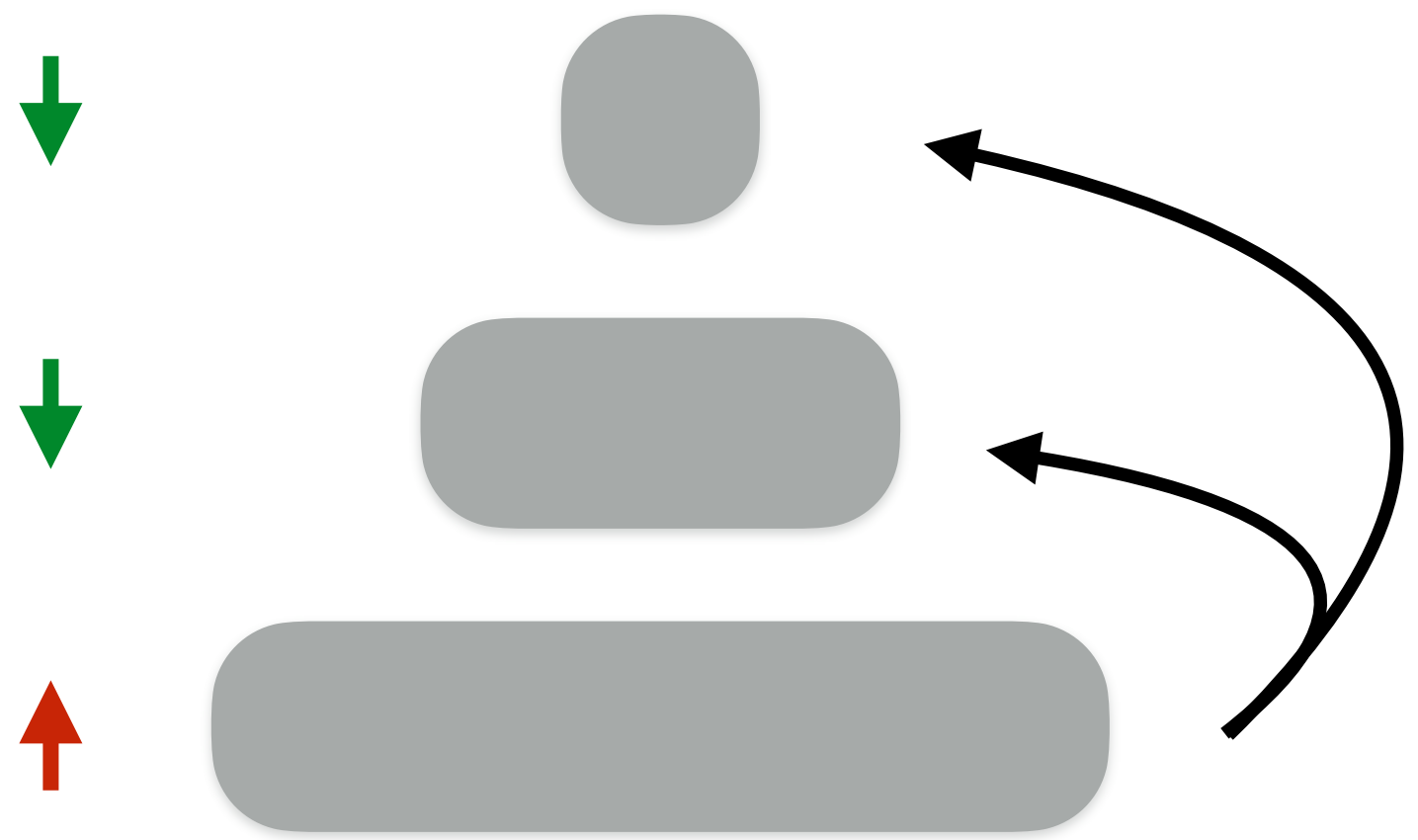
$$O(e^{-x})$$

$$O(e^{-x})$$



reallocate

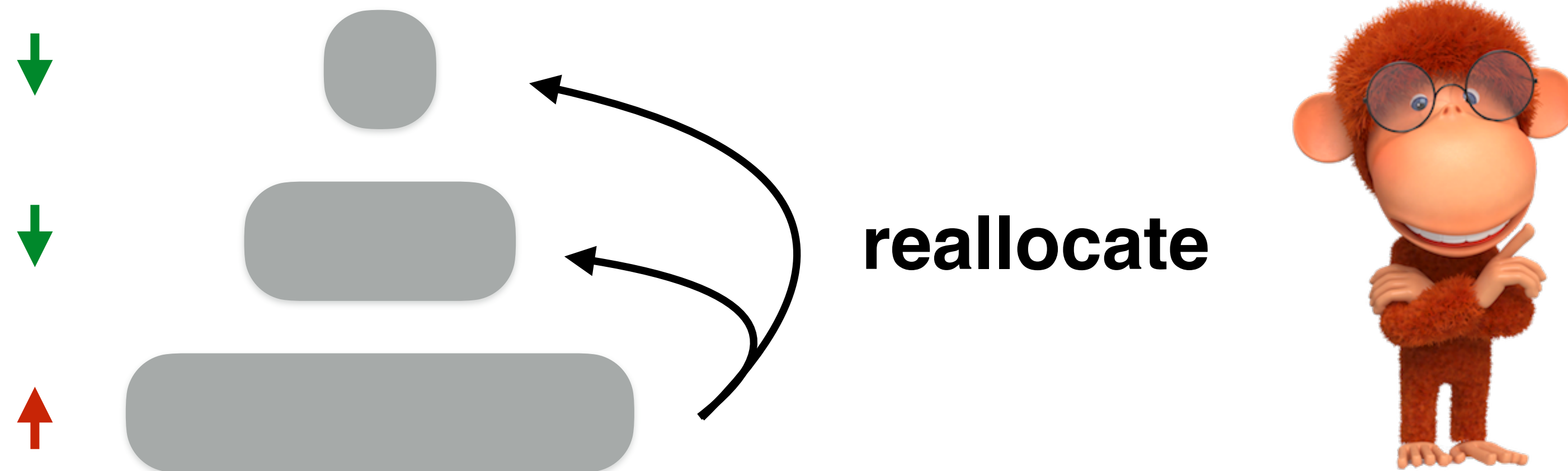




reallocate



same memory - fewer false positives



relax

false positive rates

$$0 < p_0 < 1$$

$$0 < p_1 < 1$$

$$0 < p_2 < 1$$

relax

false positive rates

$$0 < p_0 < 1$$

$$0 < p_1 < 1$$

$$0 < p_2 < 1$$

model

$$\text{read cost} = f(\mathbf{p}_0, \mathbf{p}_1 \dots)$$

$$\text{memory footprint} = f(\mathbf{p}_0, \mathbf{p}_1 \dots)$$

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

relax

false positive rates

$$0 < p_0 < 1$$

$$0 < p_1 < 1$$

$$0 < p_2 < 1$$

model

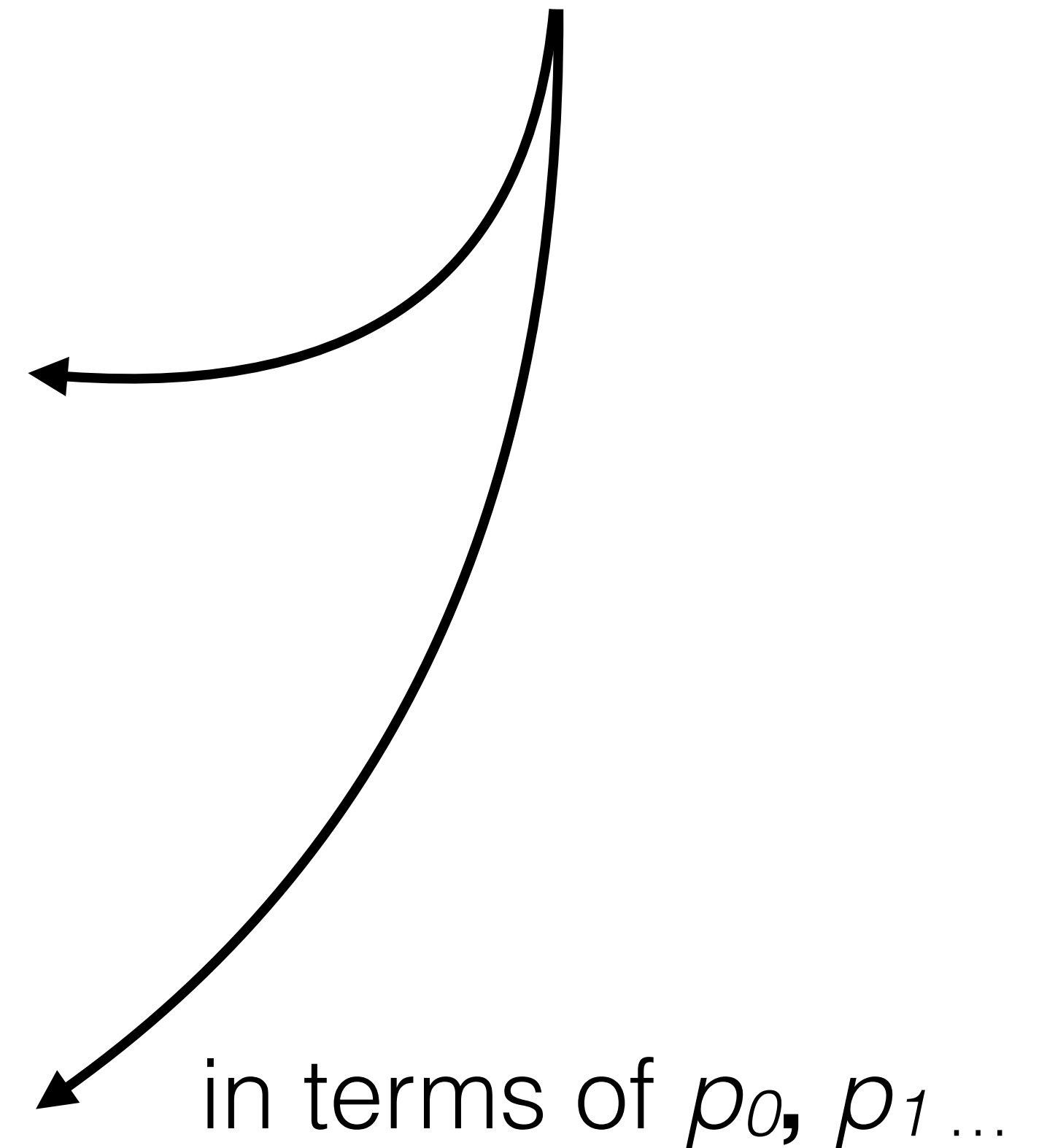
read cost

$$= \sum_1^L p_i$$

memory footprint

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

optimize



false
positive rate

$$p_0 \approx O(e^{-x}/\mathbf{R}^2)$$

$$p_1 \approx O(e^{-x}/\mathbf{R}^1)$$

$$p_2 \approx O(e^{-x}/\mathbf{R}^0)$$



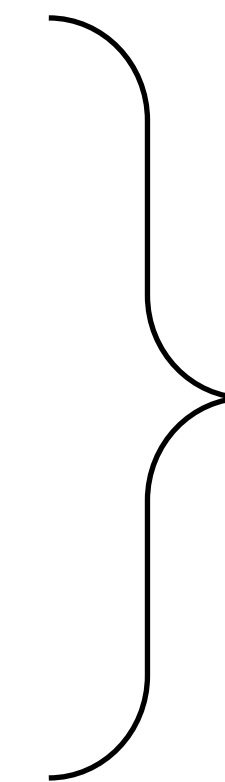


false
positive rate

$$O(e^{-x}/R^2)$$

$$O(e^{-x}/R^1)$$

$$O(e^{-x}/R^0)$$

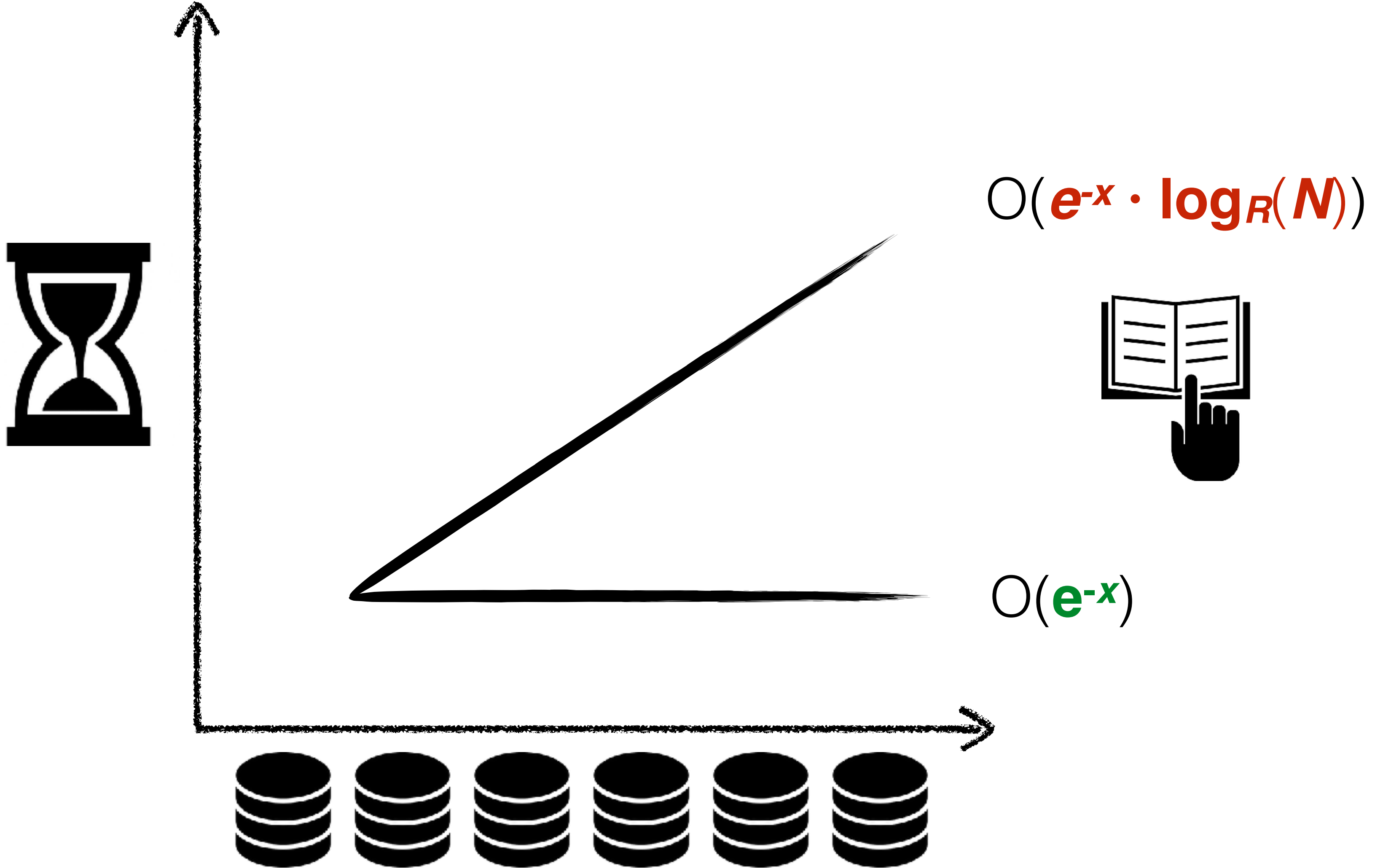


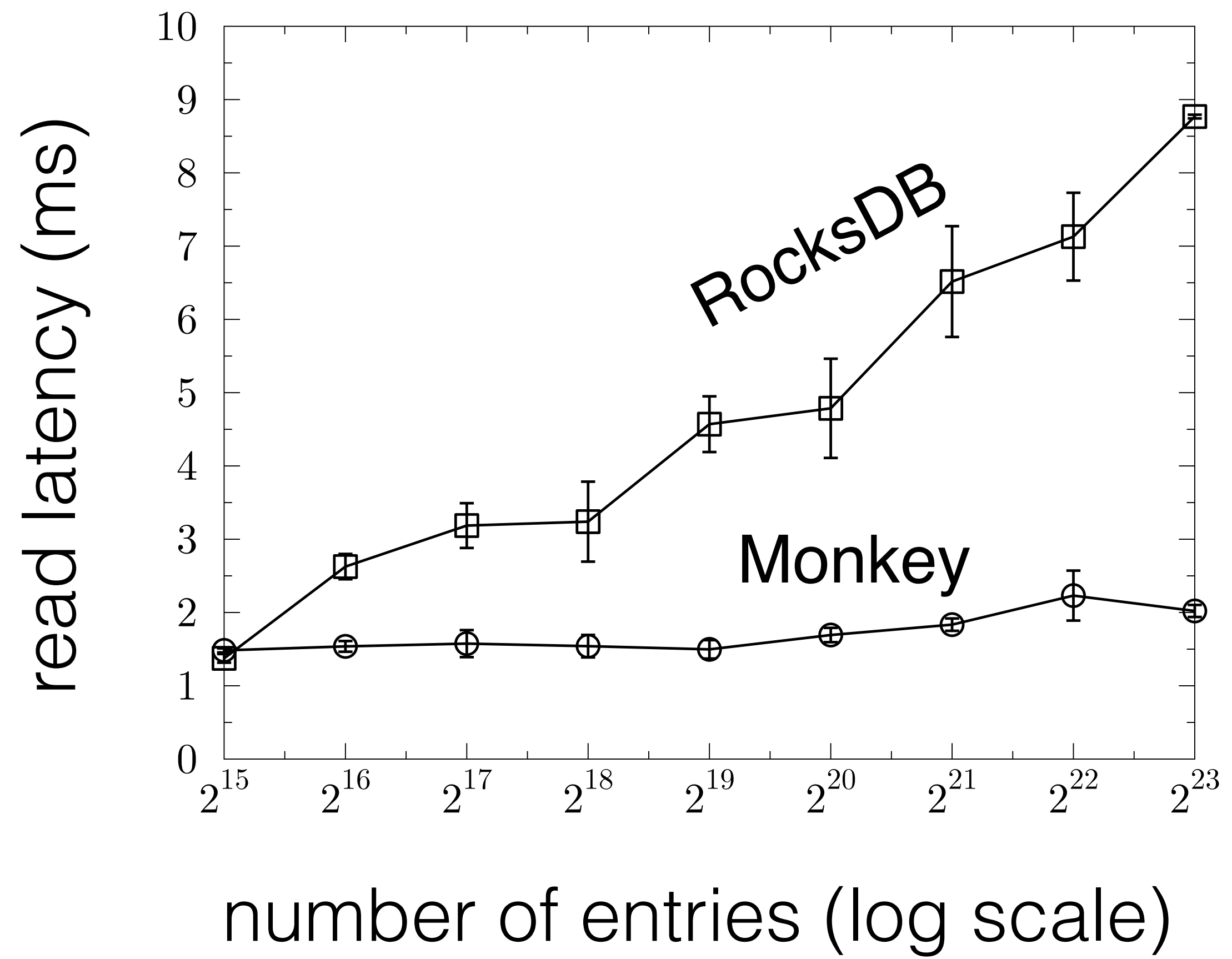
**geometric
progression**

$$= O(e^{-x})$$



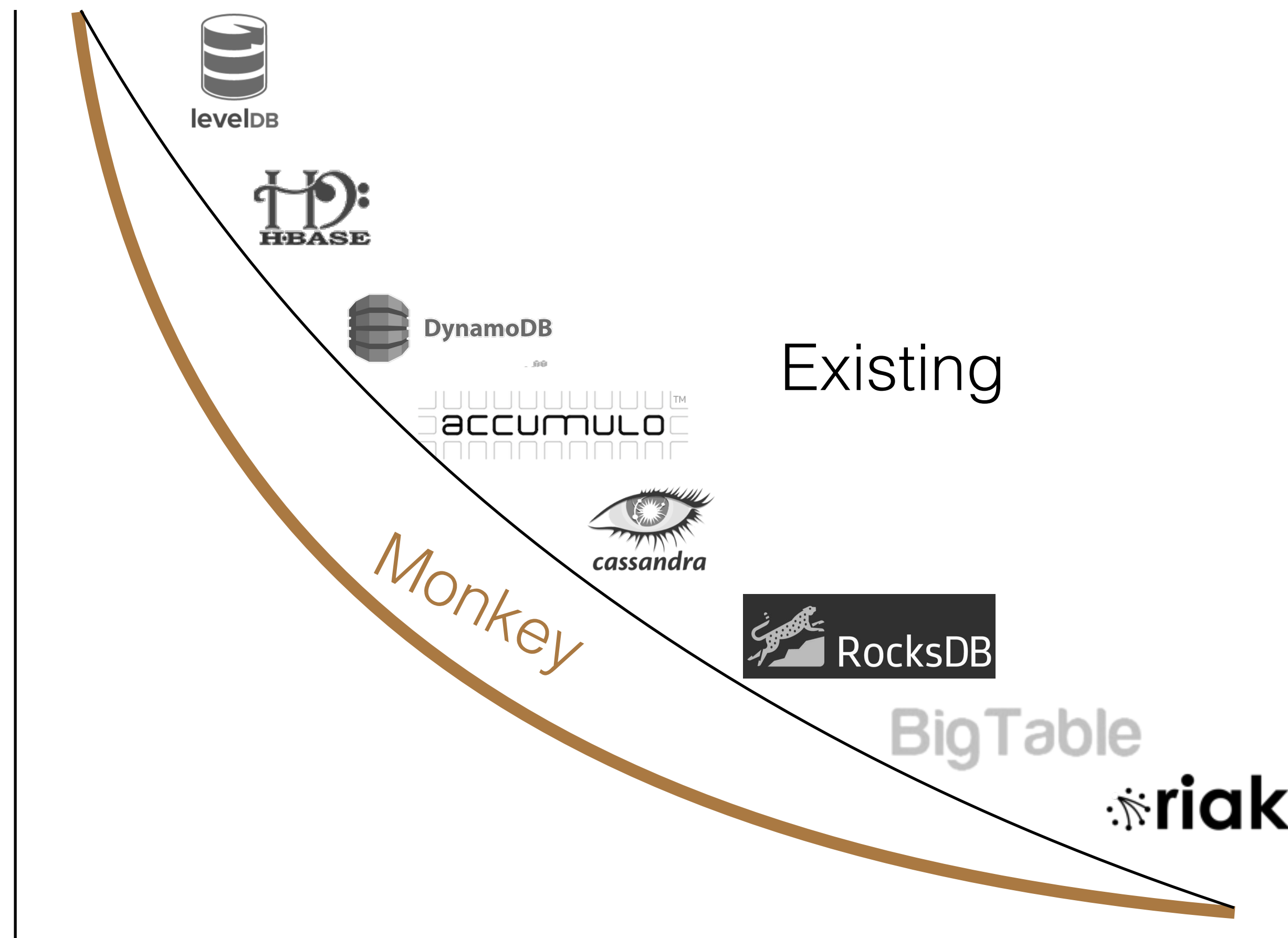
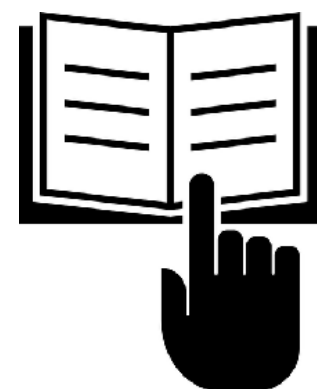
$$O(e^{-x} \cdot \log_R(N)) > O(e^{-x})$$

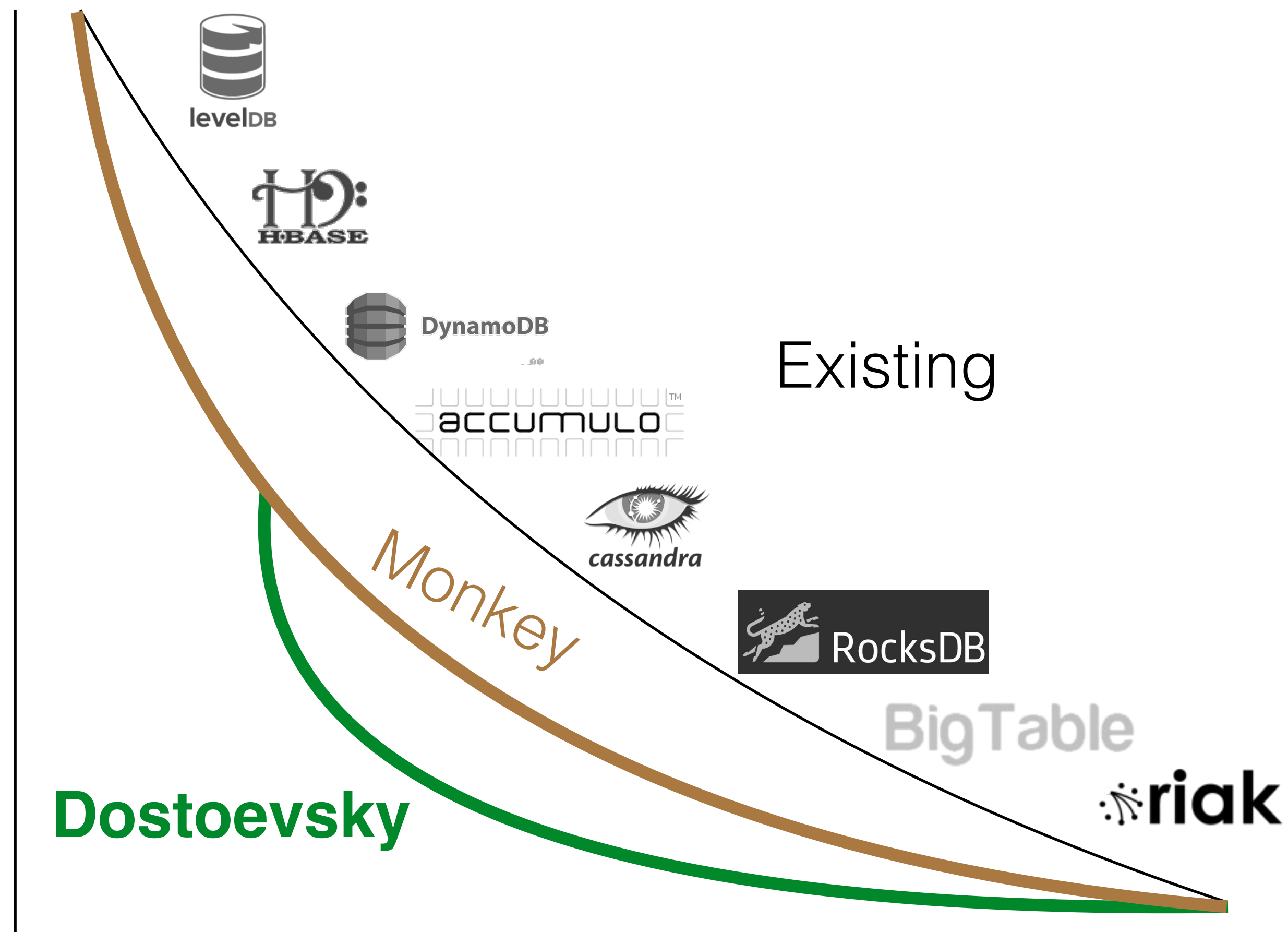
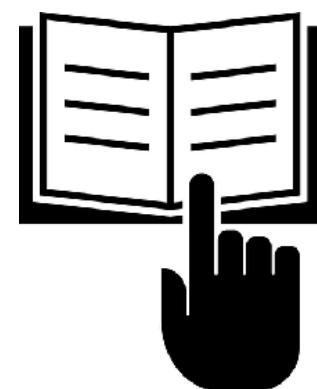


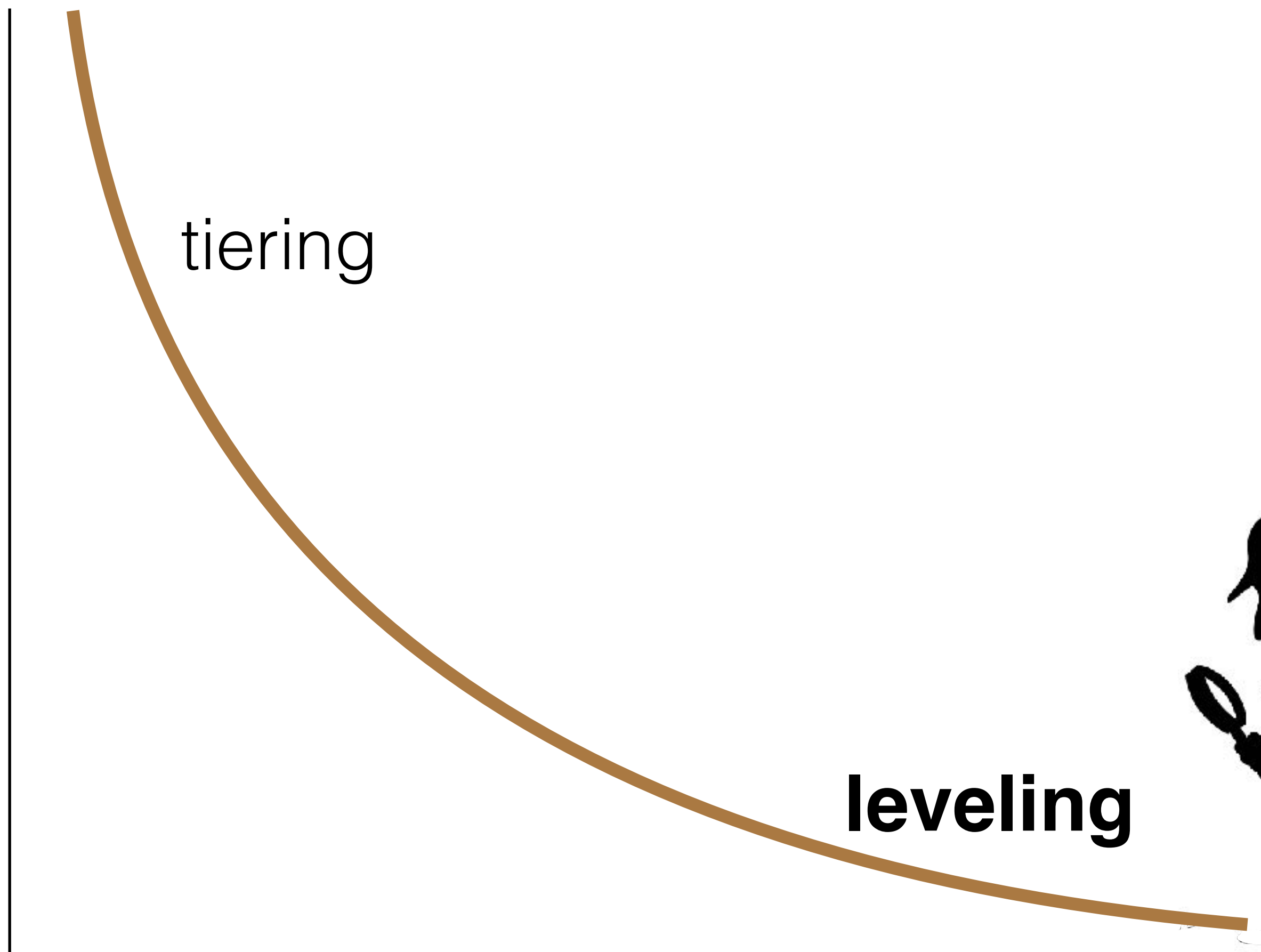
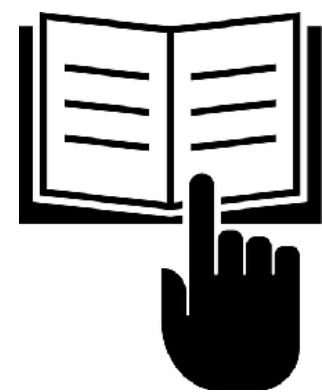


$$O(e^{-x} \cdot \log_R(N))$$

$$O(e^{-x})$$







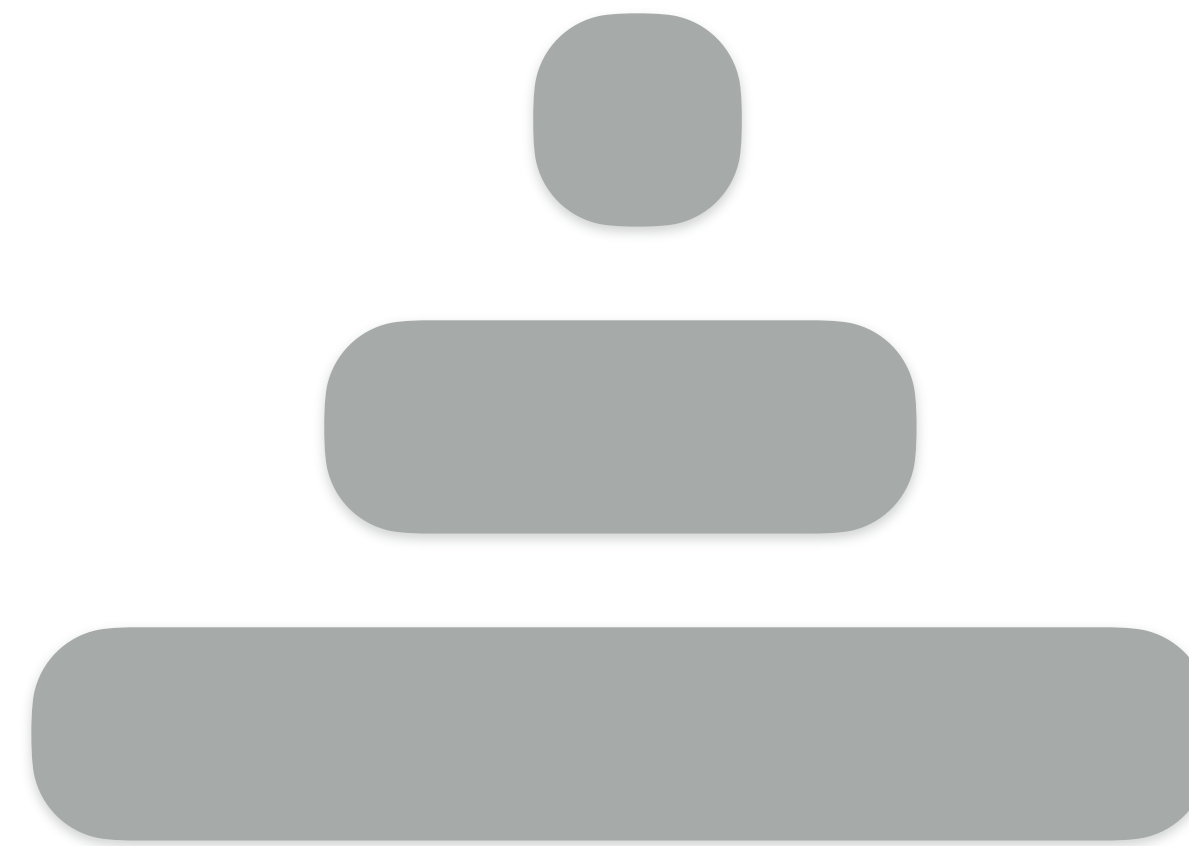
I/O overheads with leveling

●
point

→
long range

→
short range

↪
writes



●
point

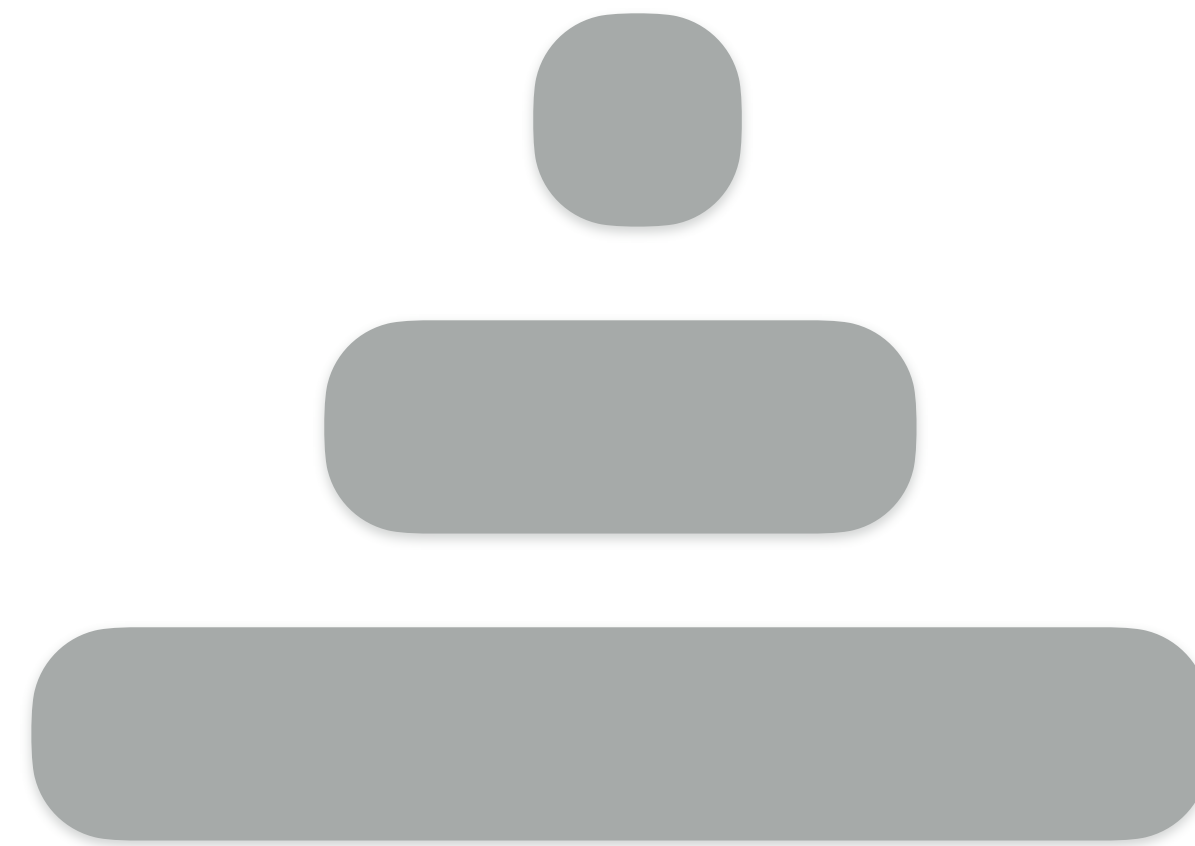
false positive rates

exponentially
decreasing

↑ $O(e^{-x}/R^2)$

$O(e^{-x}/R)$

$O(e^{-x})$



false positive rates

$$O(e^{-x}/R^2)$$

$$O(e^{-x}/R)$$

→ $O(e^{-x})$



●
point




point

largest level

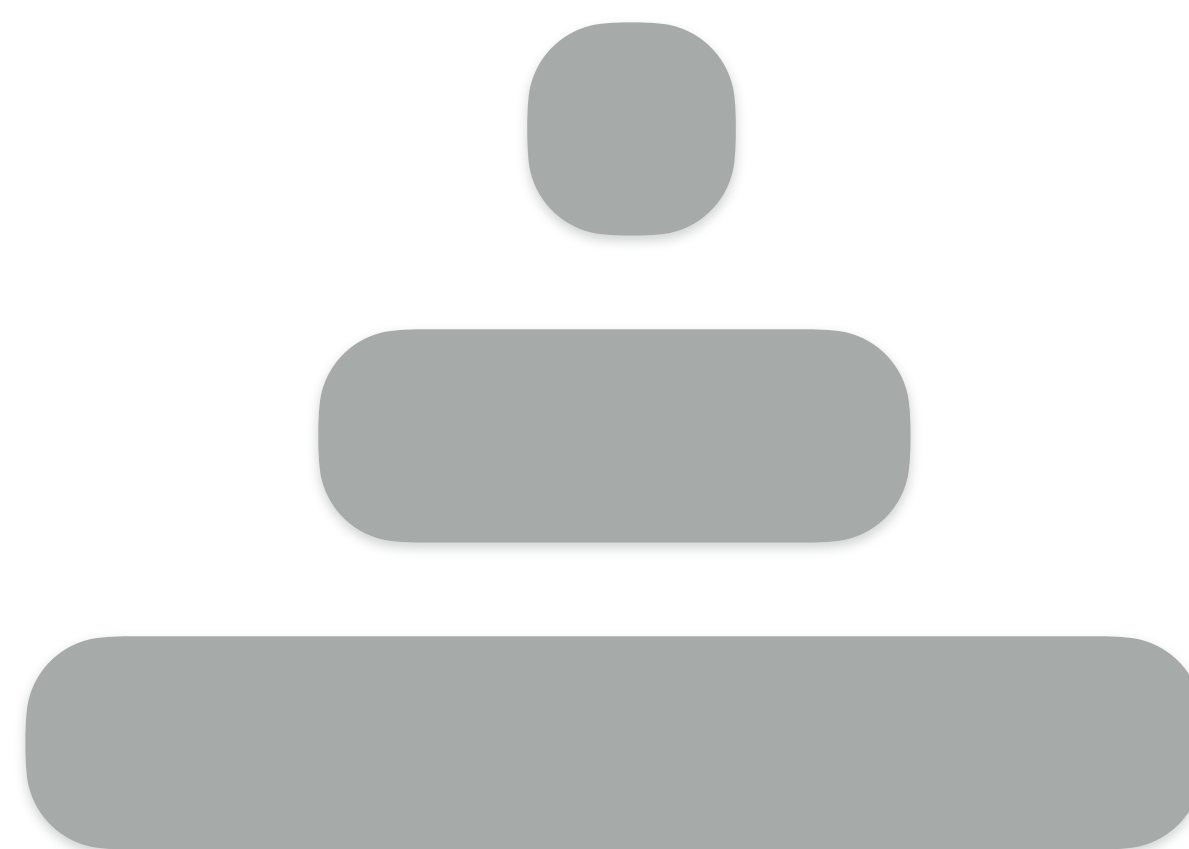
$O(e^{-x})$


long range


short range



writes



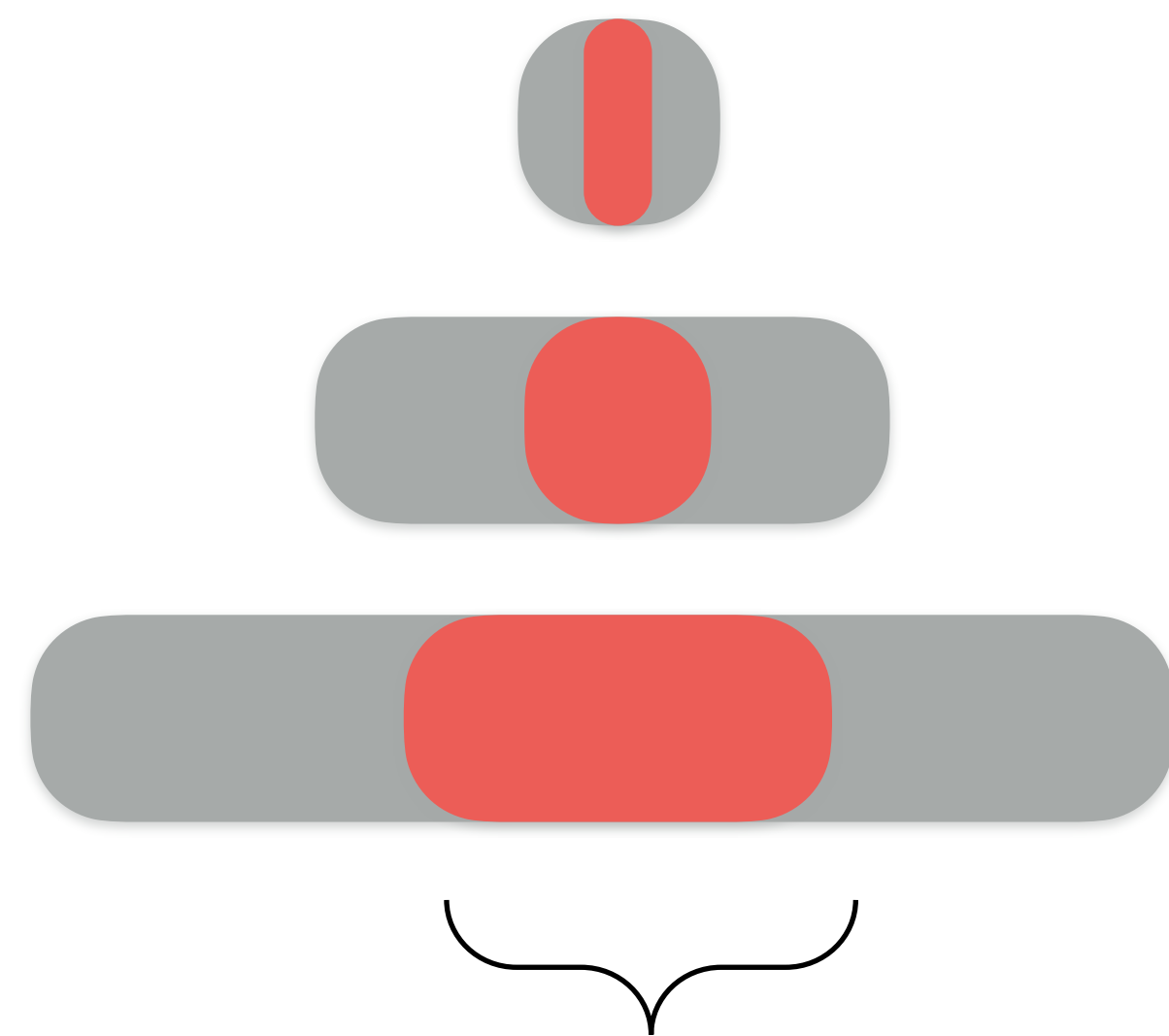
long range →

target range

$O(s/R^2)$

$O(s/R)$

$O(s)$



target key range

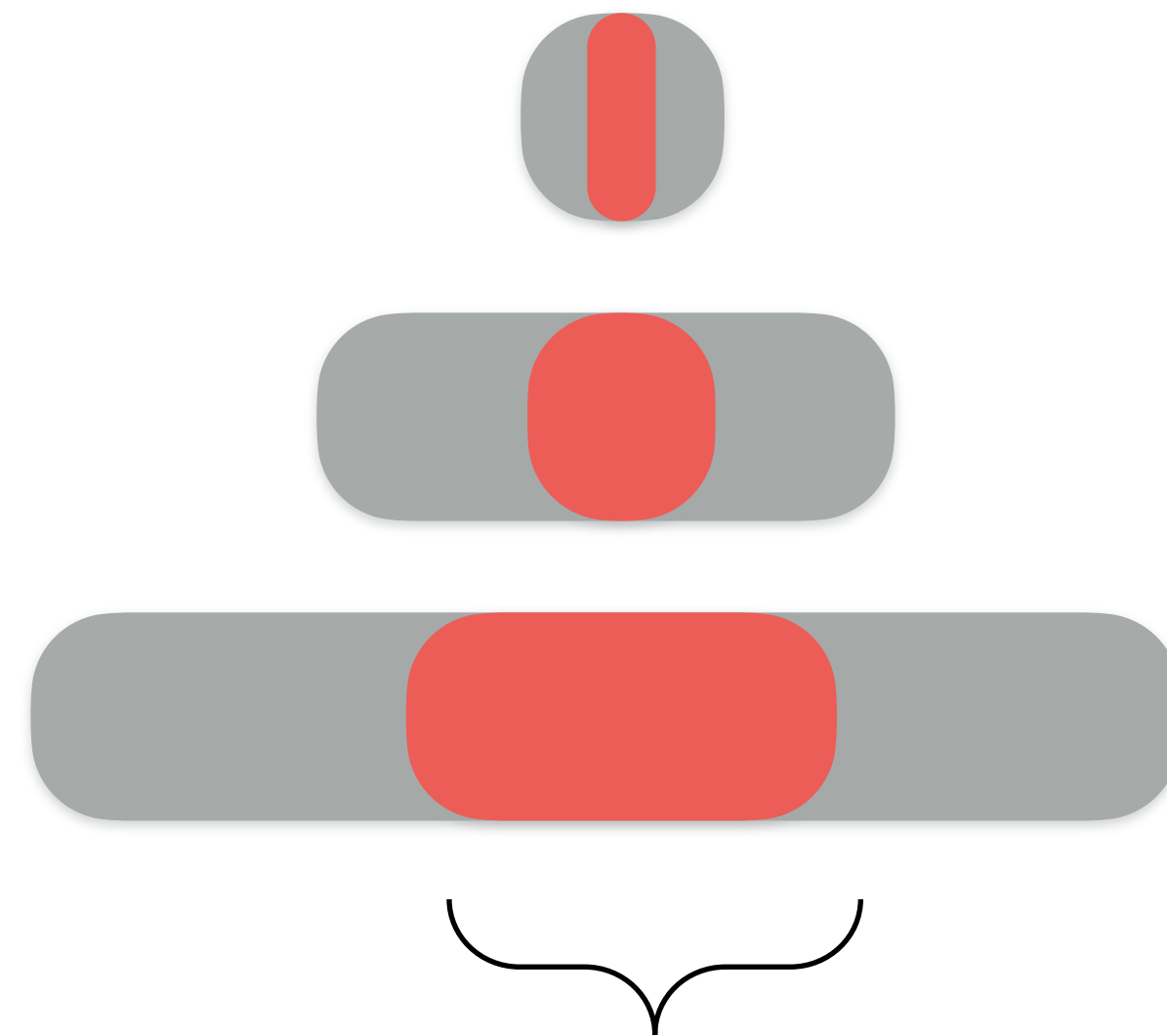
long range →

target range

$$O(s/R^2)$$

$$O(s/R)$$

$$O(s)$$



largest level

target key range

●
point

largest level

$O(e^{-x})$

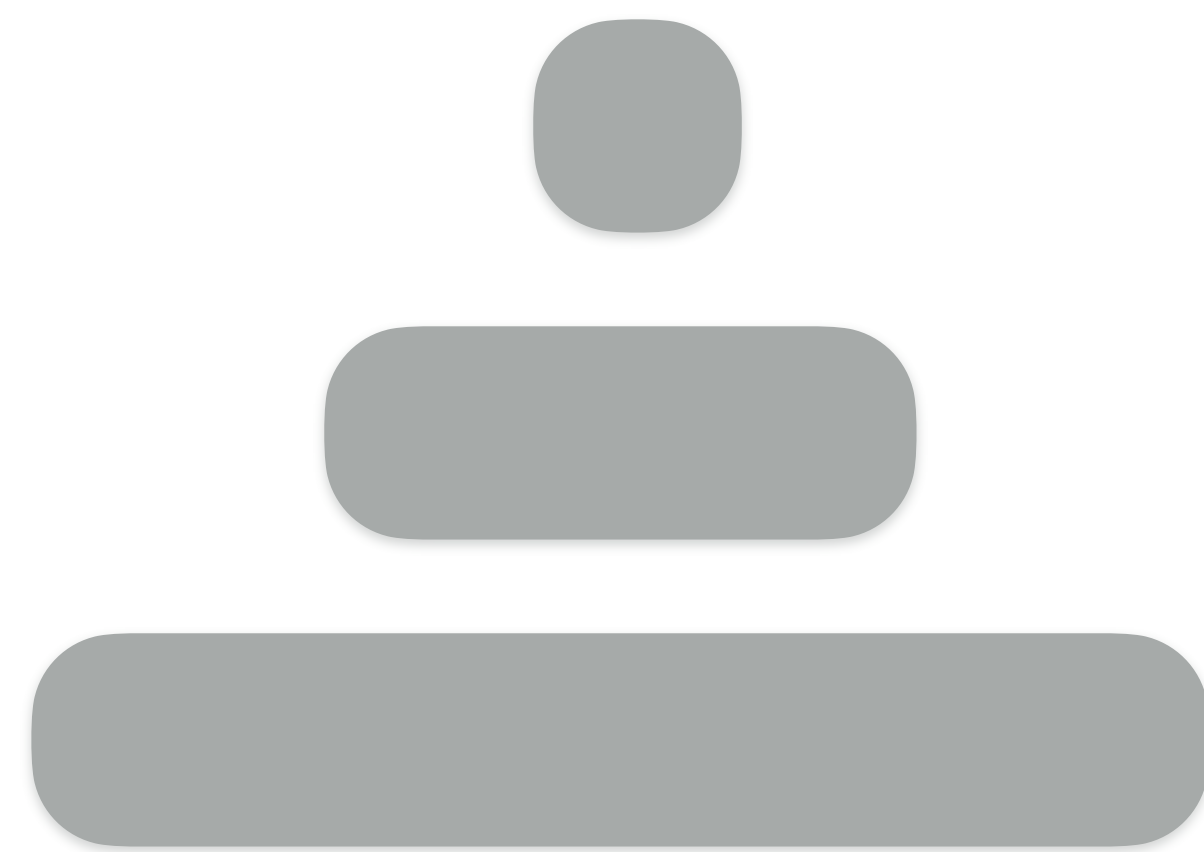
→
long range

largest level

$O(s)$

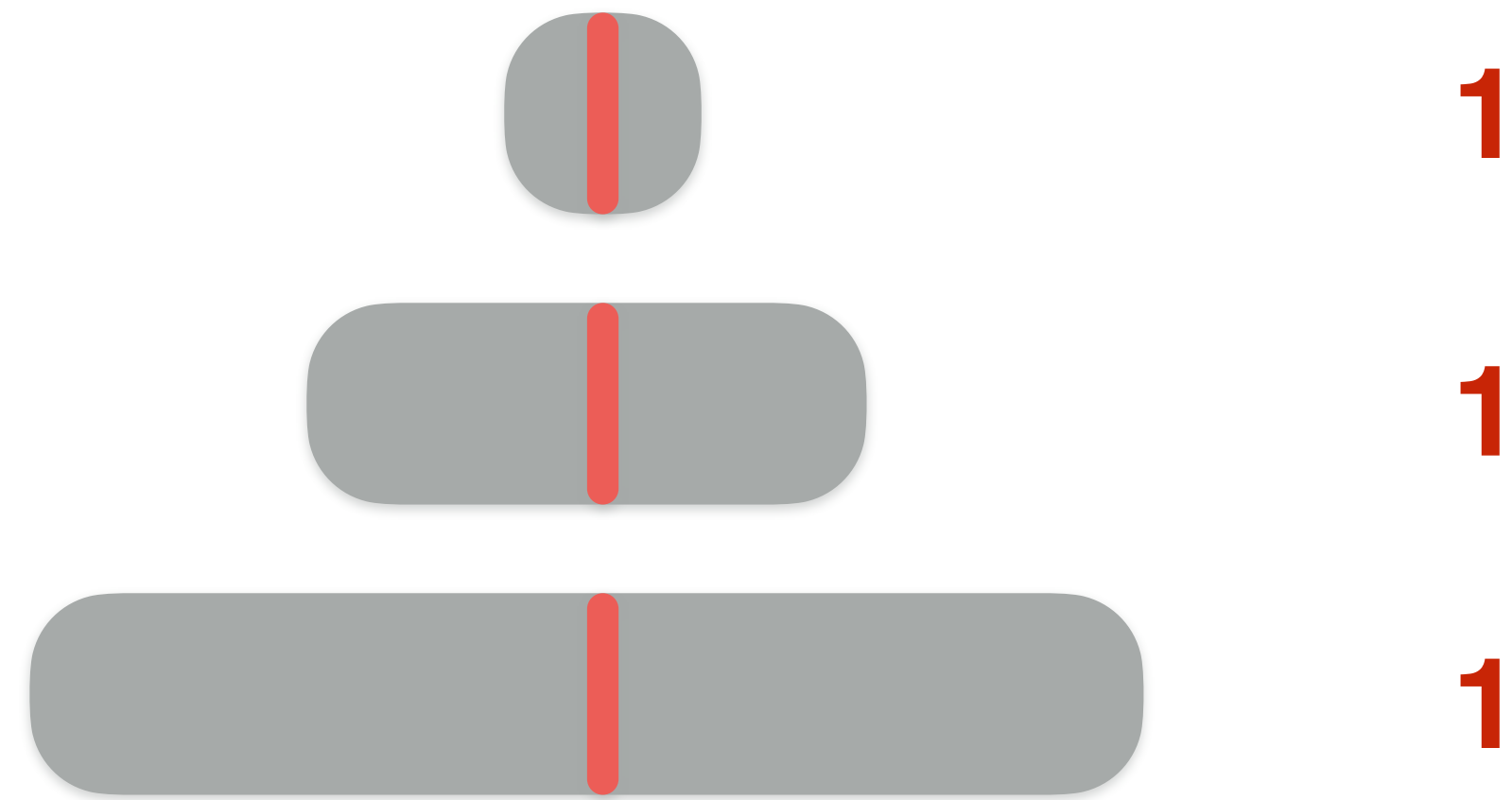
→
short range

↶
writes



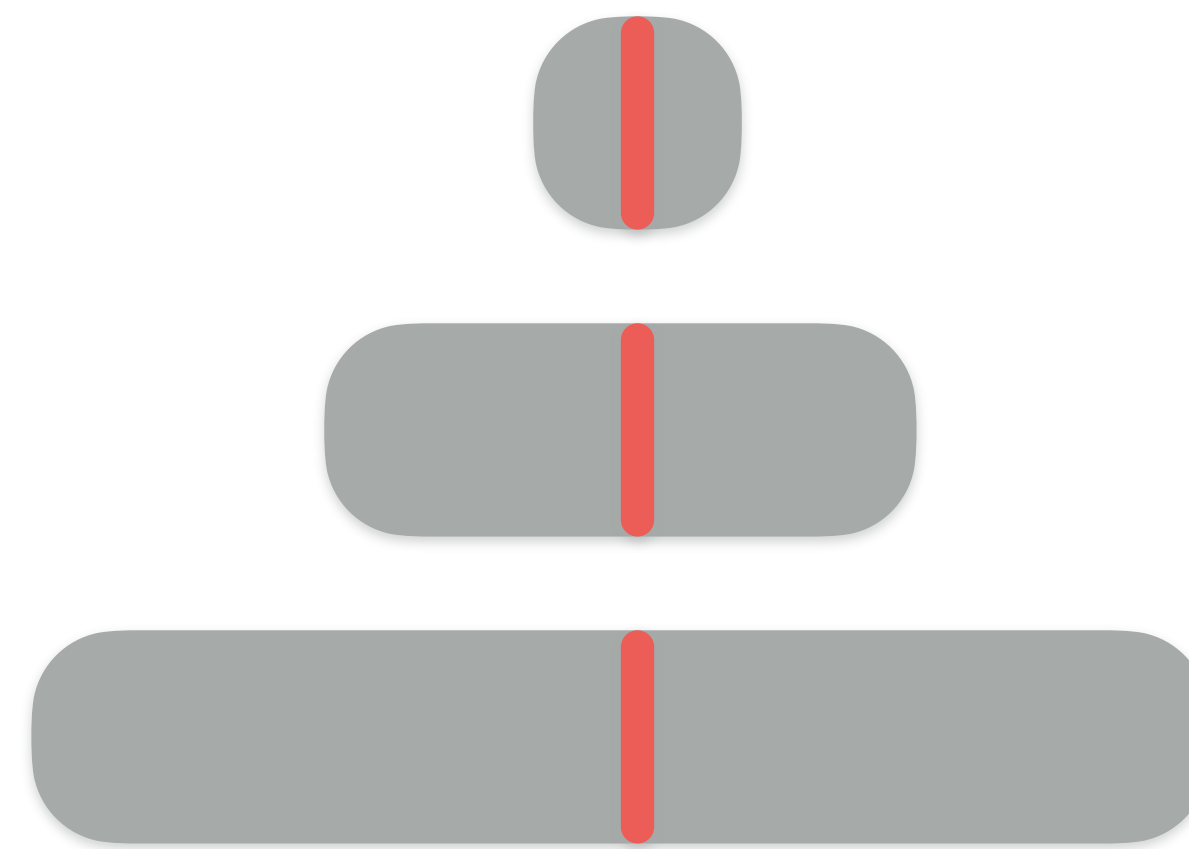
short range →

target range



short range →

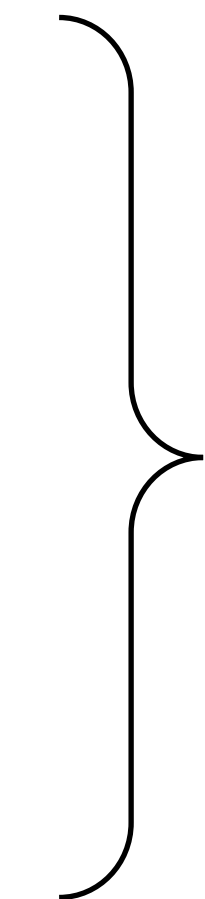
target range



1

1

1



all levels
 $O(\log_R(N))$

●
point

largest level

$O(e^{-x})$

→
long range

largest level

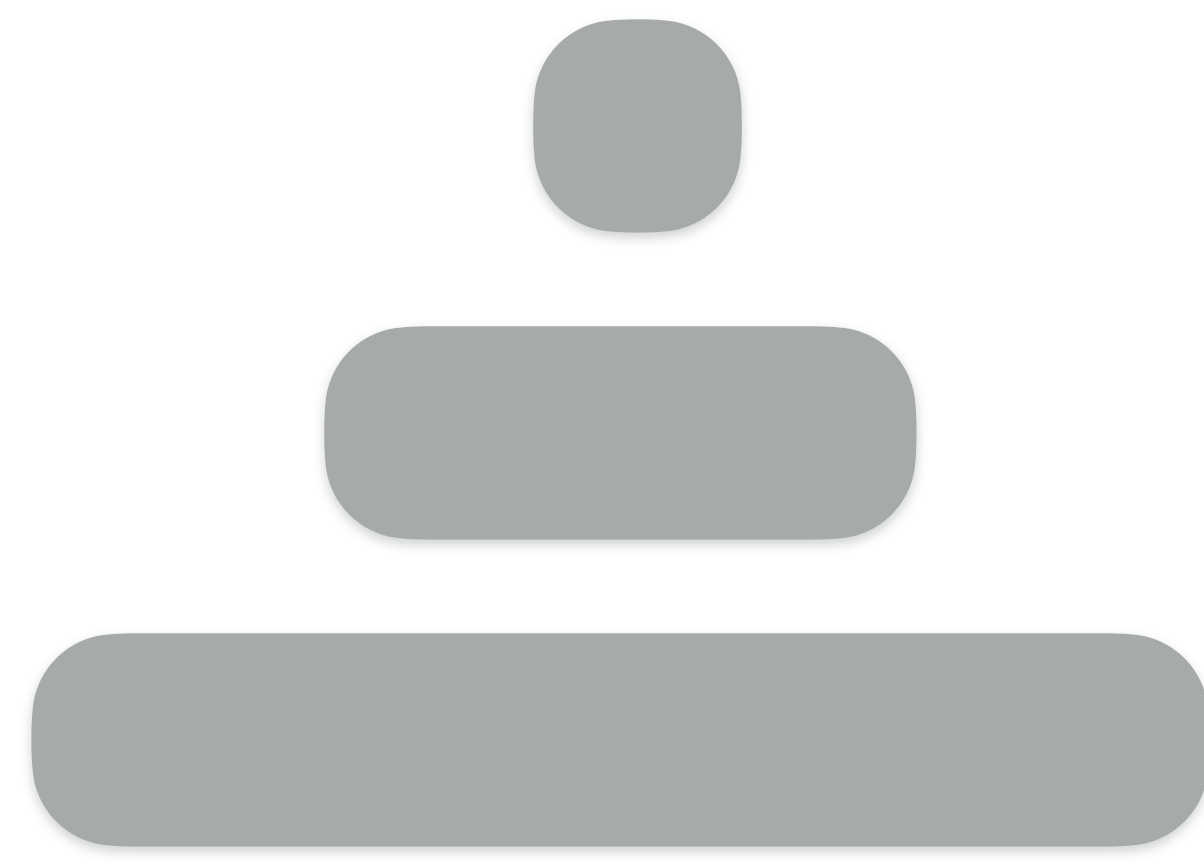
$O(s)$

→
short range

all levels

$O(\log_R(N))$

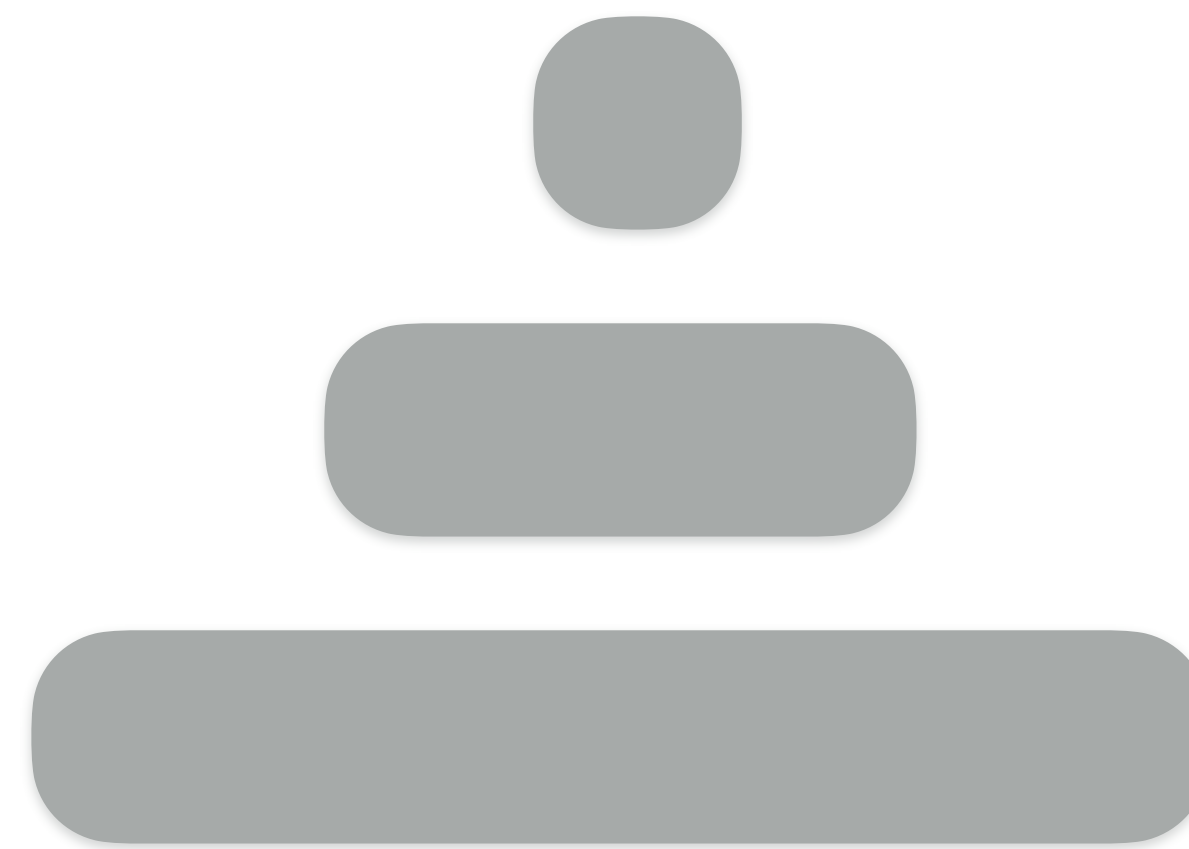
↶
writes





**exponentially
more work**



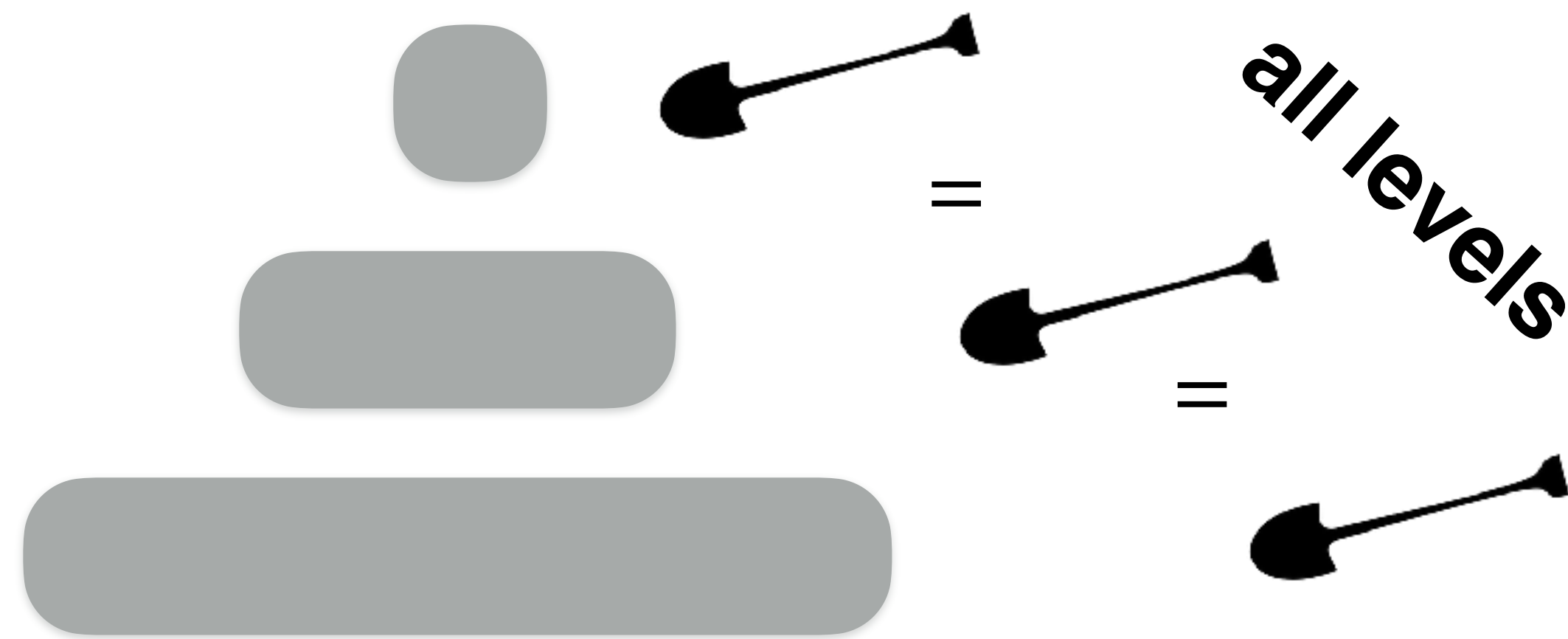


exponentially
more work



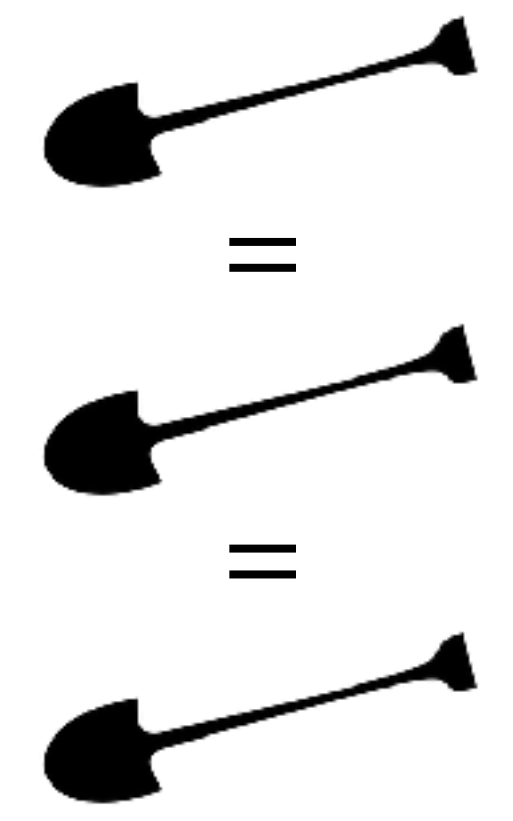
**exponentially
less frequent**


writes

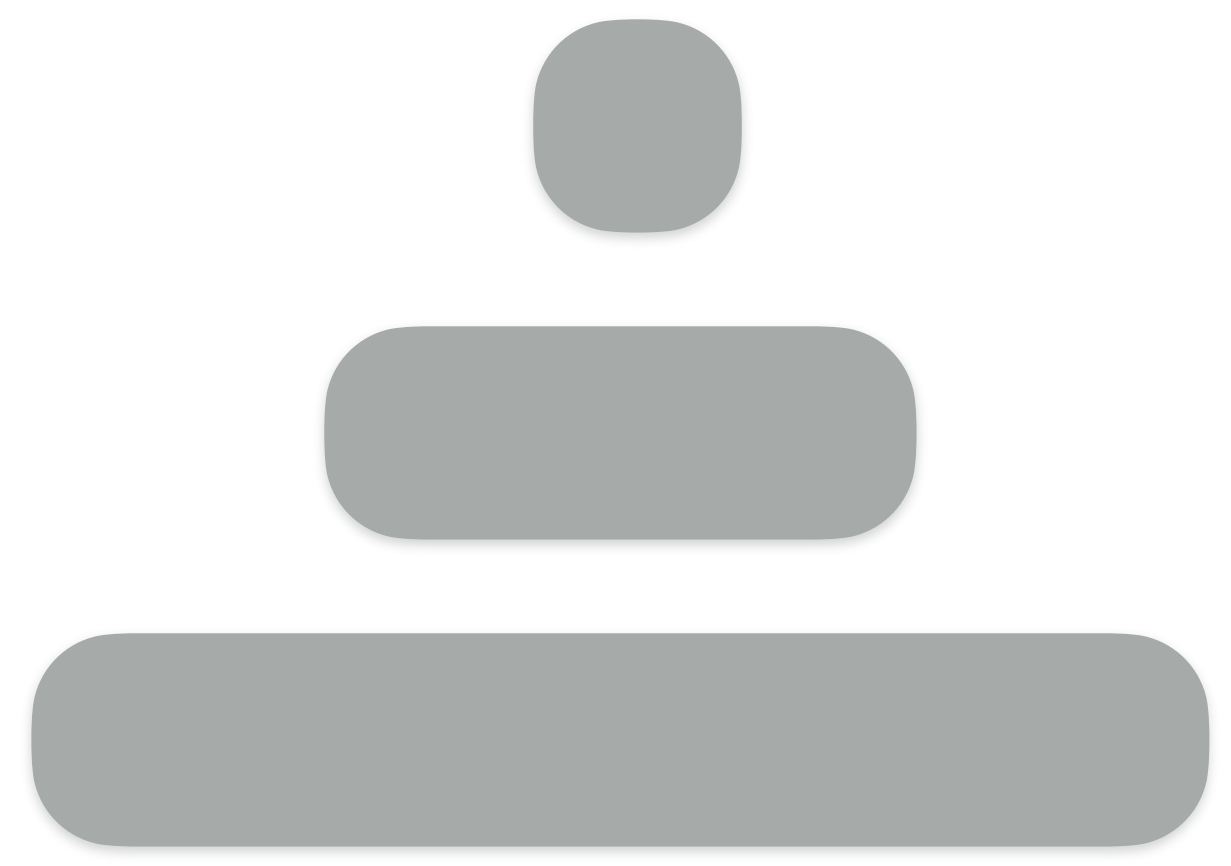


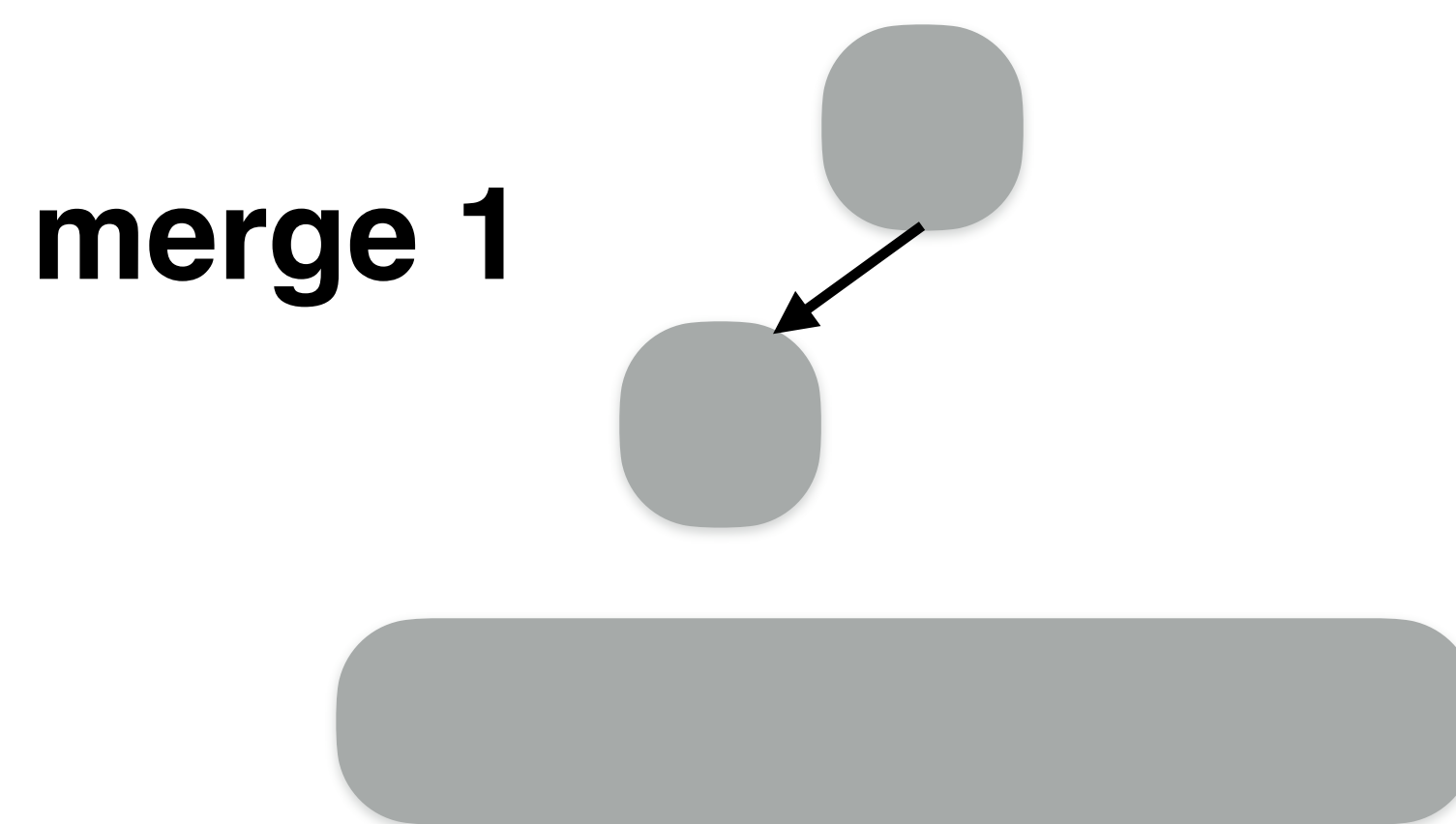
more work
↓
less frequent

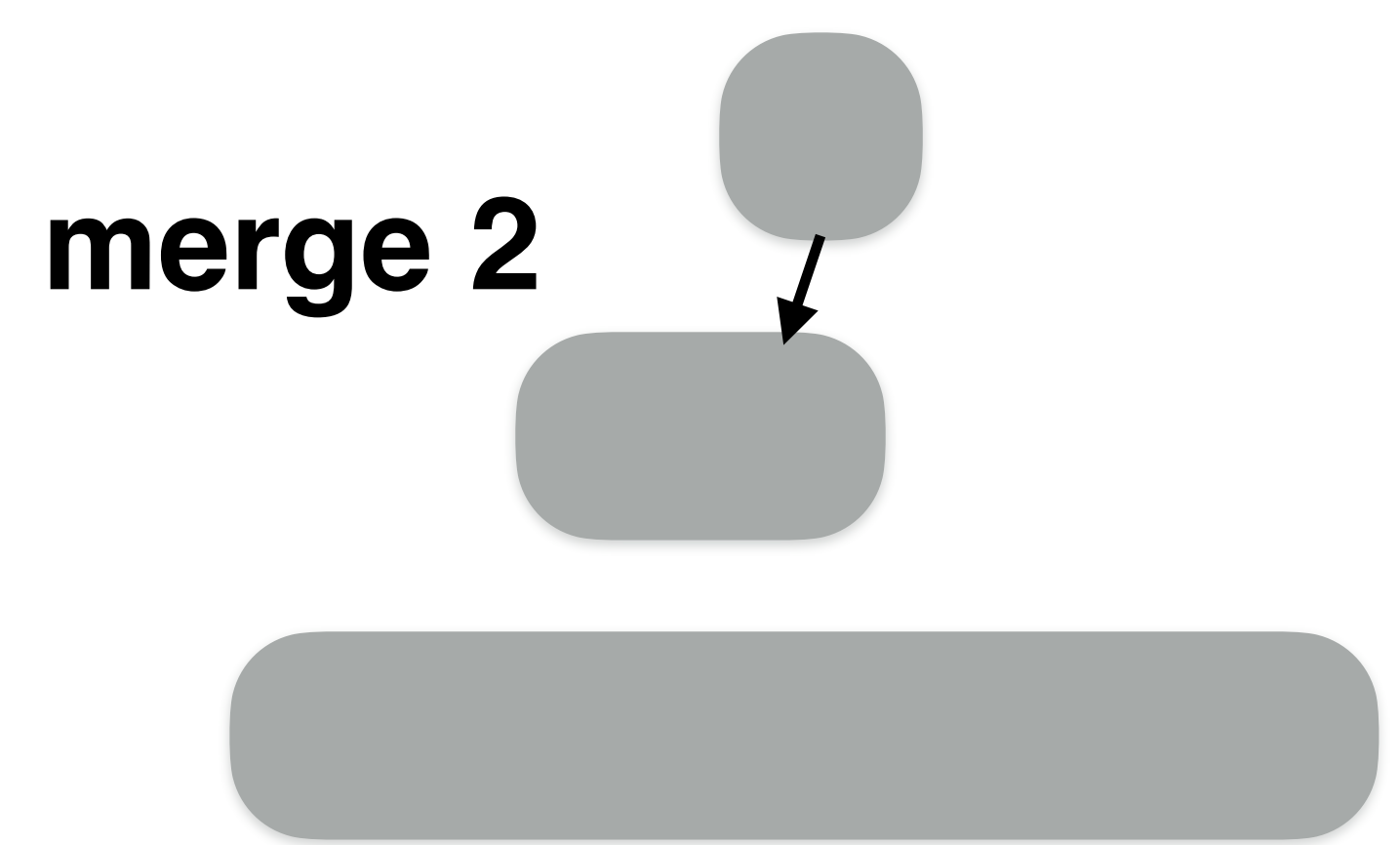

writes

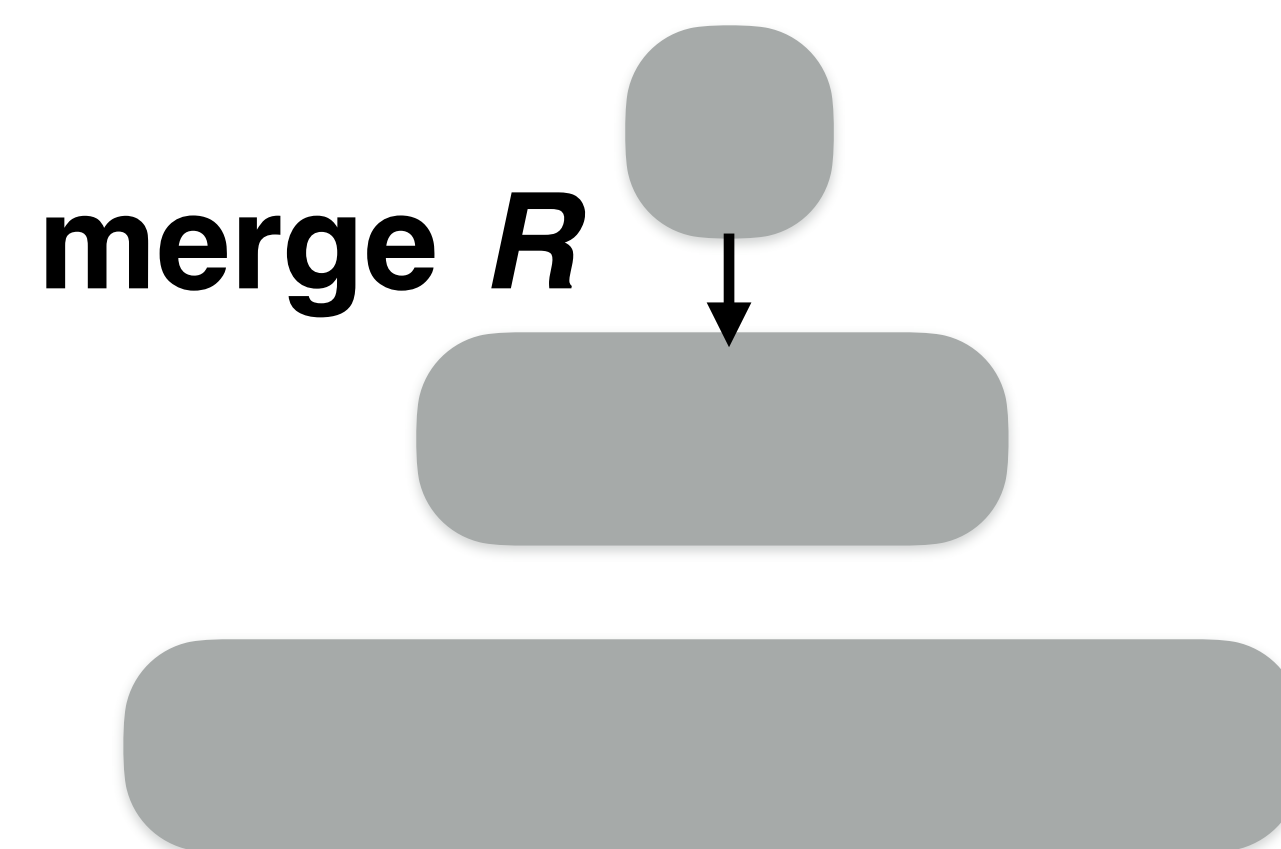


all levels



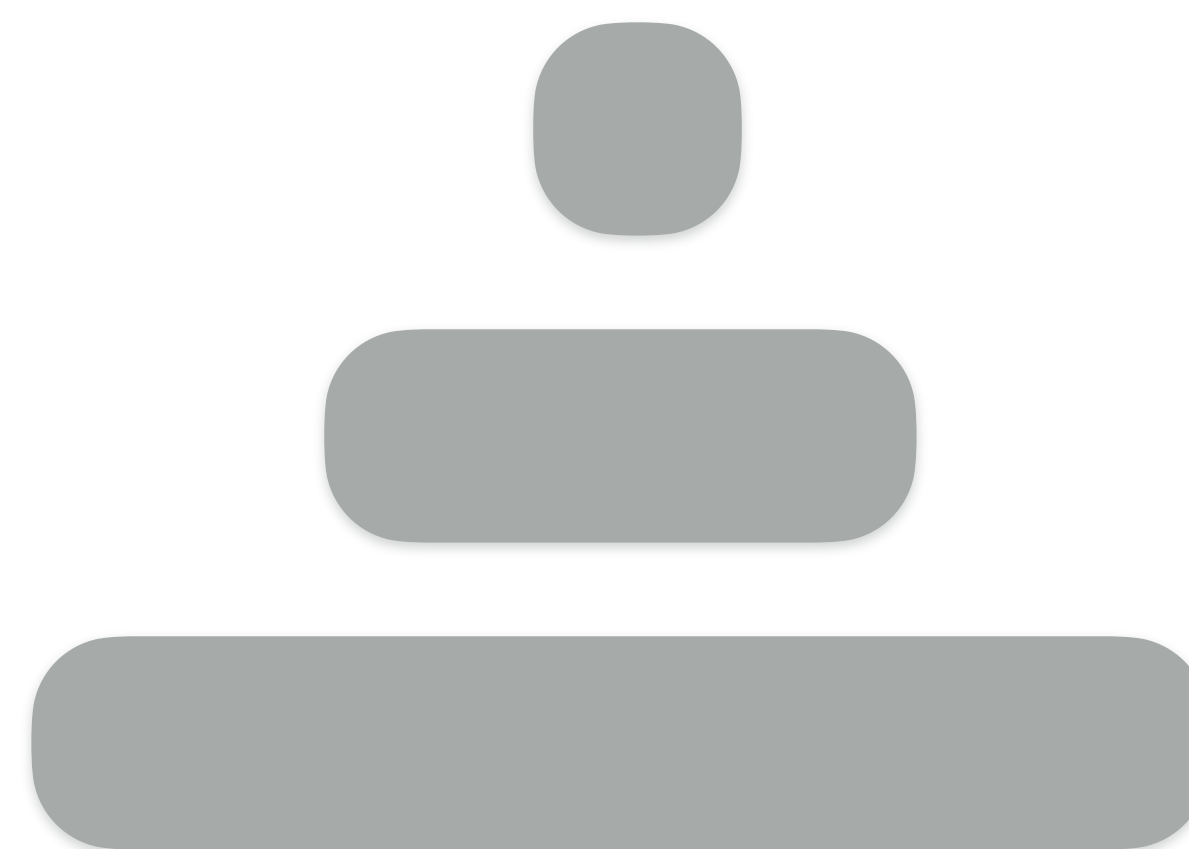








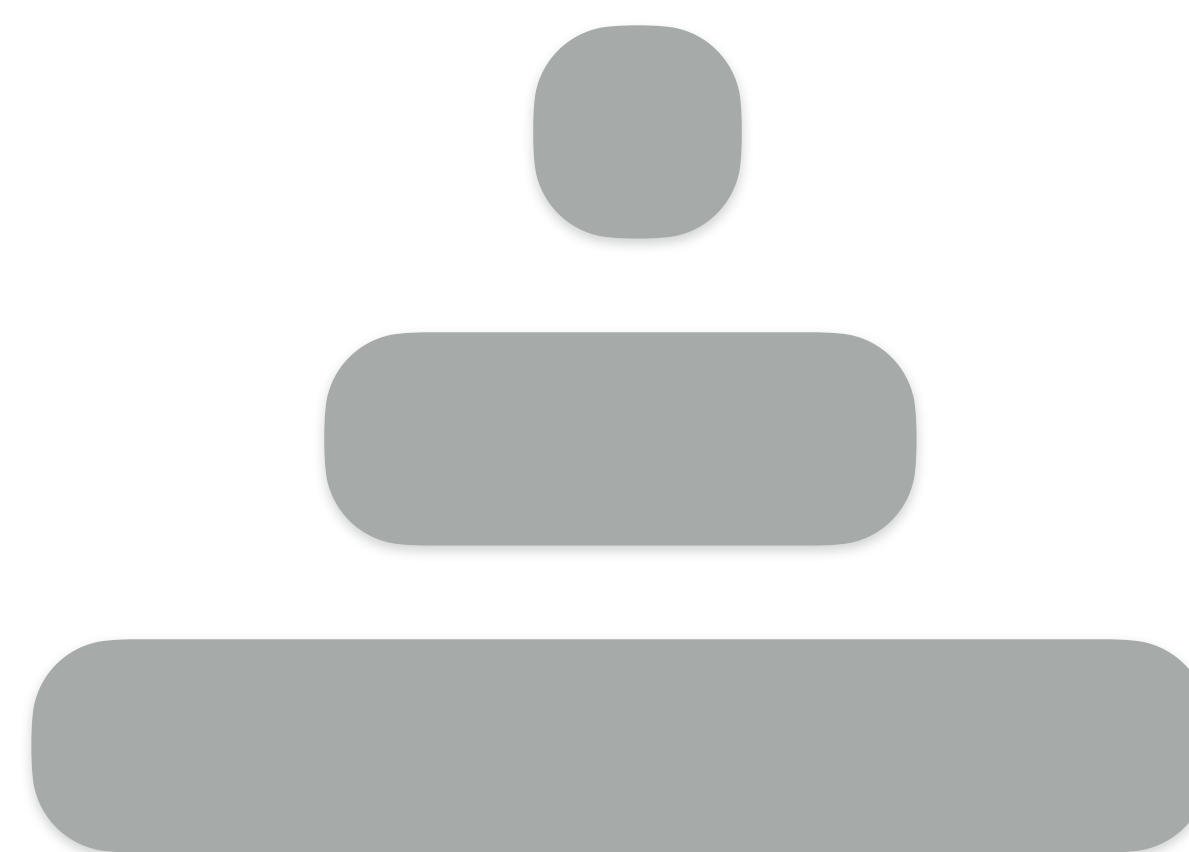
write-amplification



$O(R)$

$O(R)$

$O(R)$



$$\left. \begin{array}{l} O(R) \\ O(R) \\ O(R) \end{array} \right\} O(R \cdot \log_R(N))$$

●
point

→
long range

→
short range

🔍
writes

largest level
 $O(e^{-x})$

largest level
 $O(s)$

all levels
 $O(\log_R(N))$

all levels
 $O(R \cdot \log_R(N))$

=

=

=

=

$O(e^{-x}/R^2)$

$O(s/R^2)$

1

$O(R)$

+

+

+

+

$O(e^{-x}/R)$

$O(s/R)$

1

$O(R)$

+

+

+

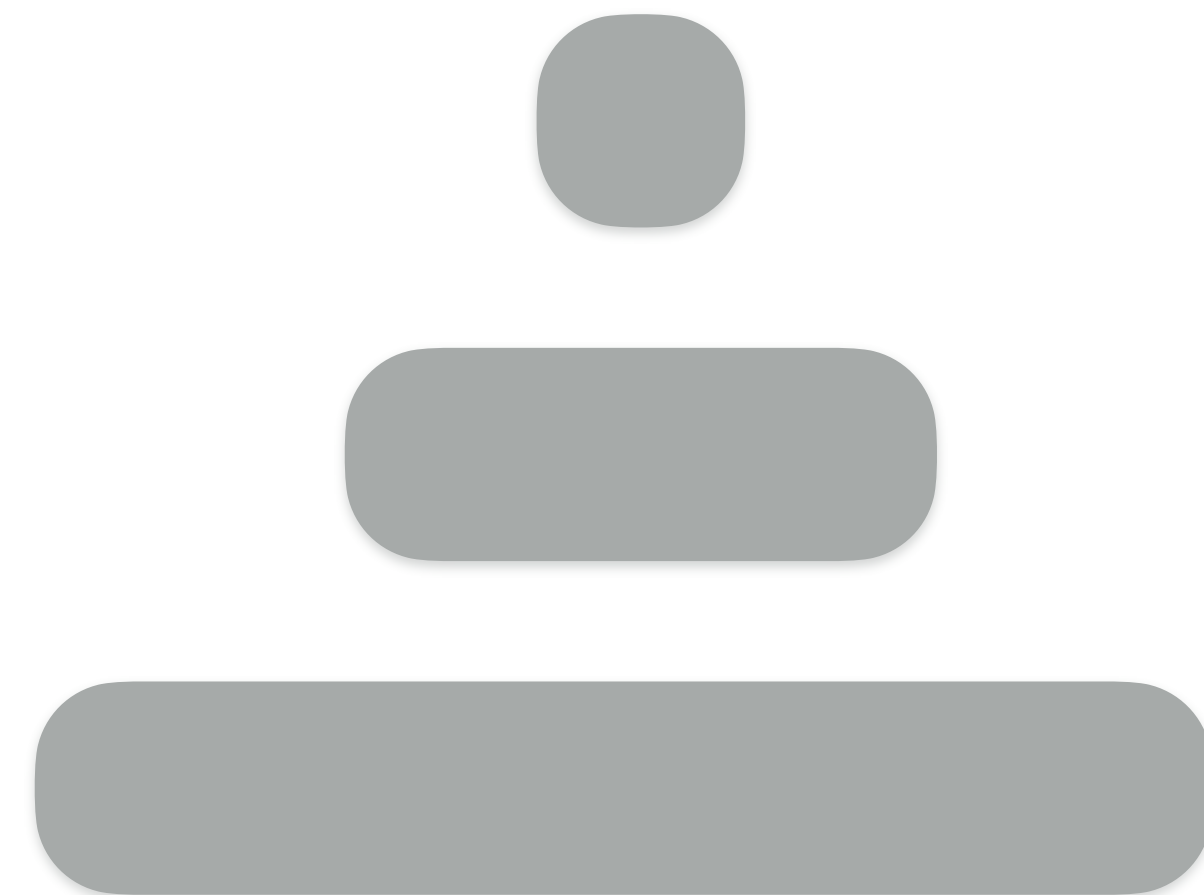
+

$O(e^{-x})$

$O(s)$

1

$O(R)$



●
point

→
long range

↶
writes

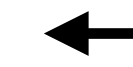
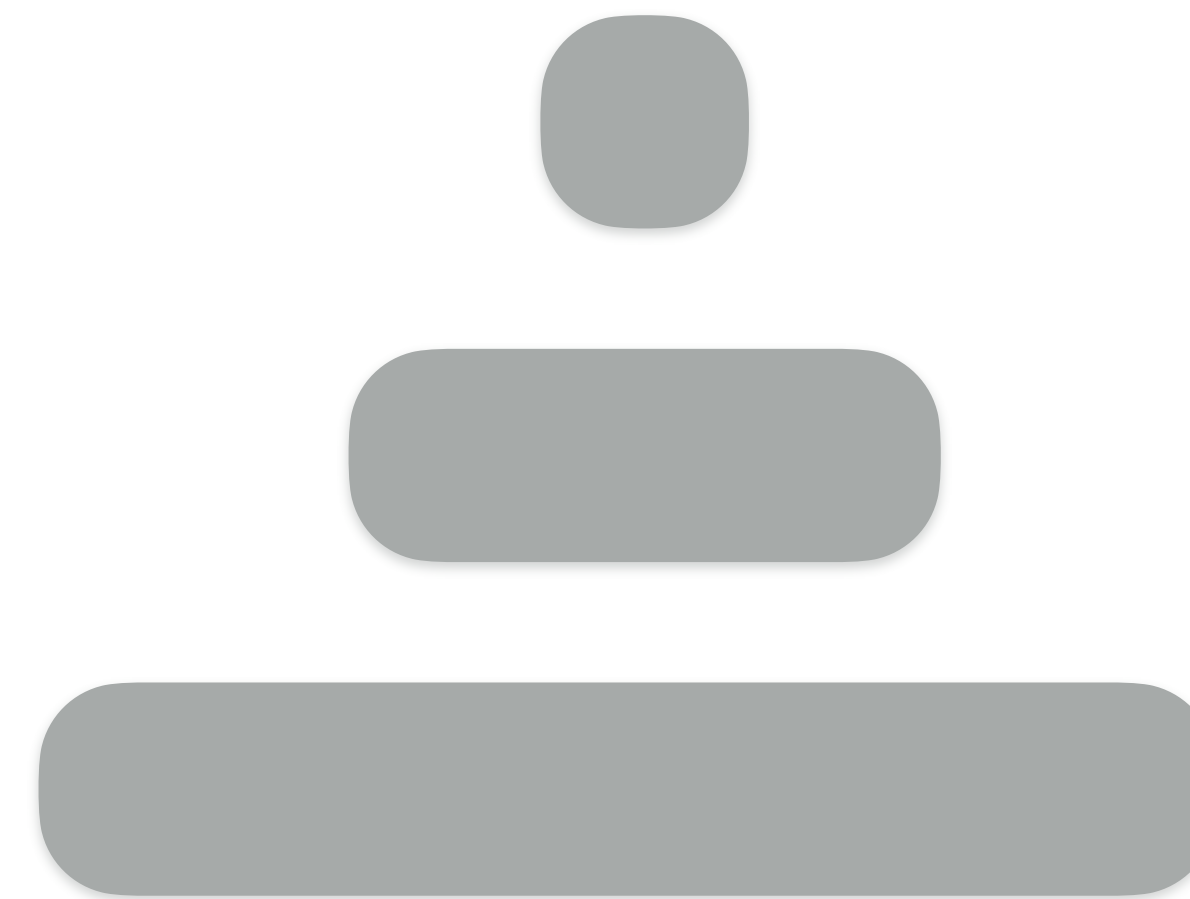
largest level

largest level

all levels

$O(e^{-x})$

$O(s)$



$O(R)$



$O(R)$



$O(R)$

●
point

→
long range

↙
writes

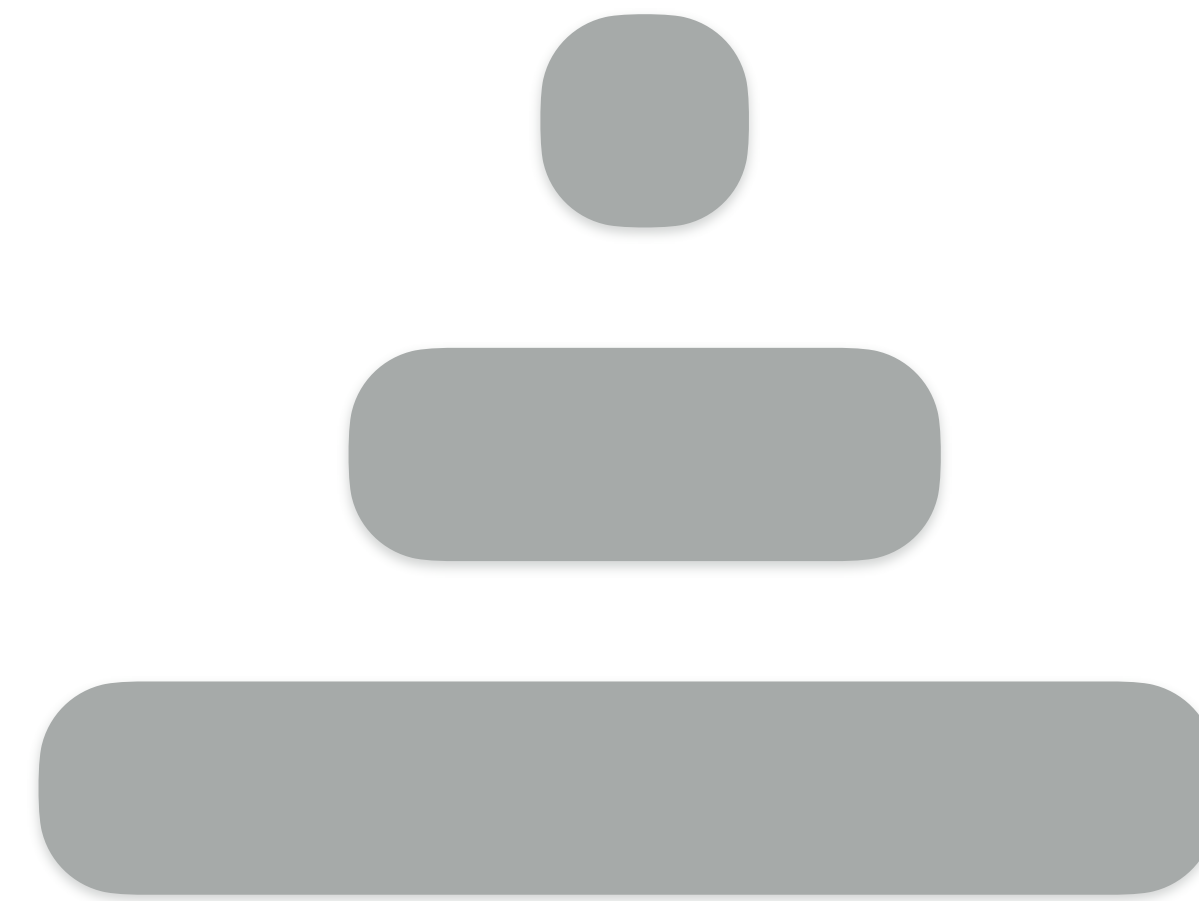
largest level

largest level

all levels

$O(e^{-x})$

$O(s)$



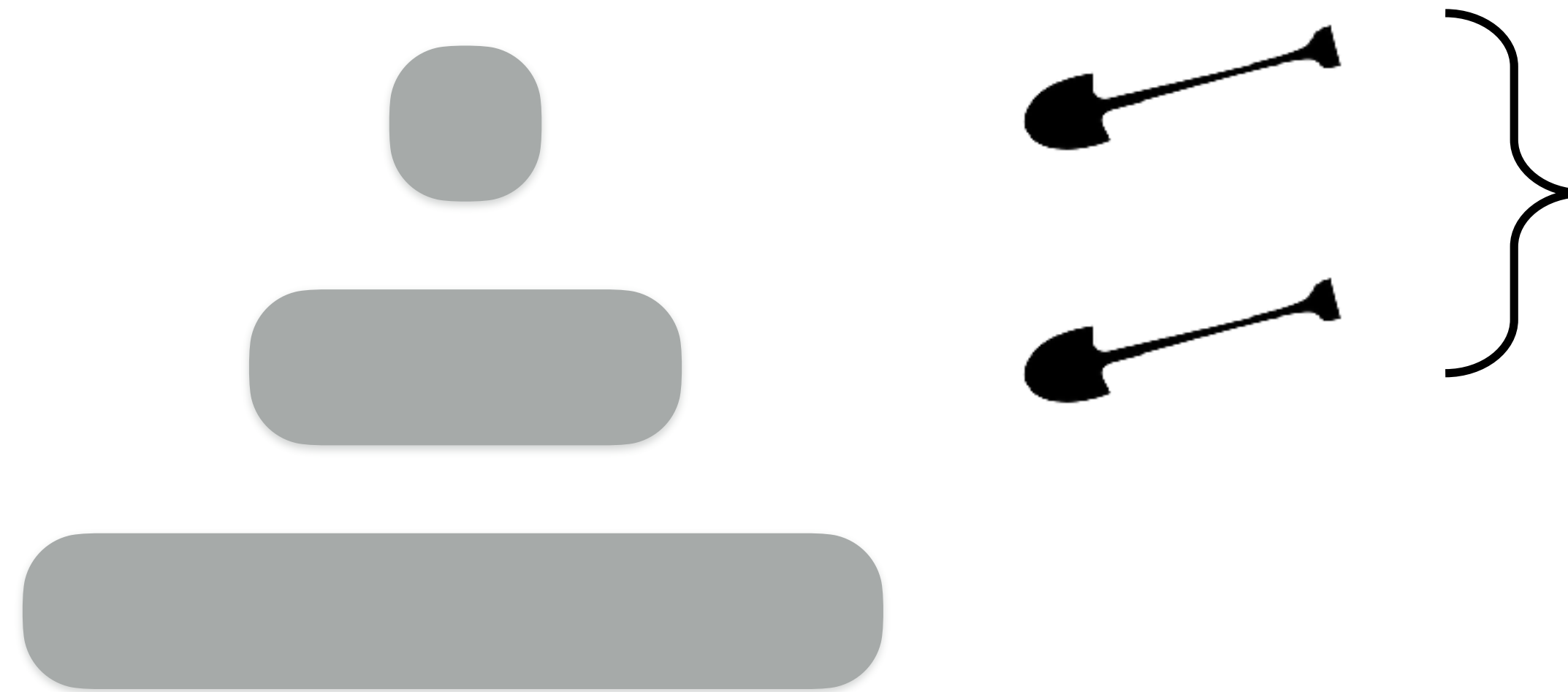
← $O(R)$

← $O(R)$

← $O(R)$

superfluous

•  for point lookups and long range lookups
merging at smaller levels is superfluous



worse as data grows!





poor performance



poor performance

lower device lifetime (on SSD)



Dostoevsky

SIGMOD18



Dostoevsky: **S**pace-**T**ime **O**ptimized **E**volvab**l**e **S**calab**l**e **K**ey-Value Store

very write-optimized



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

Tiering
write-optimized

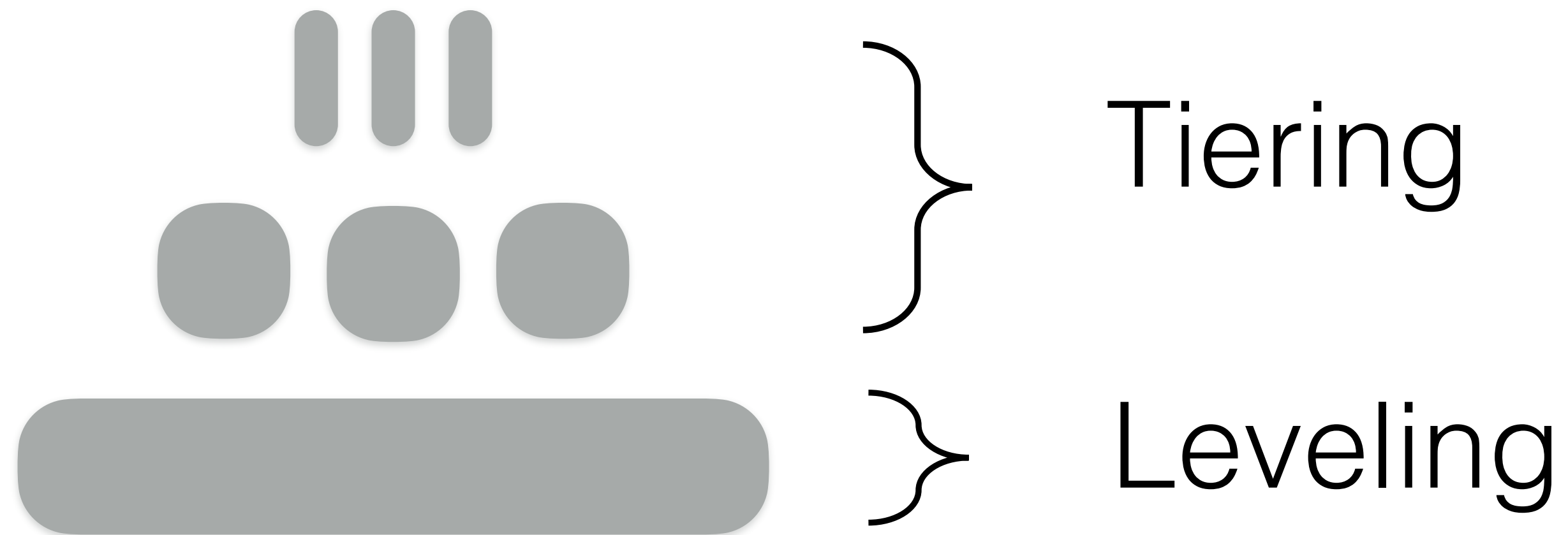
Leveling
read-optimized

Tiering
write-optimized

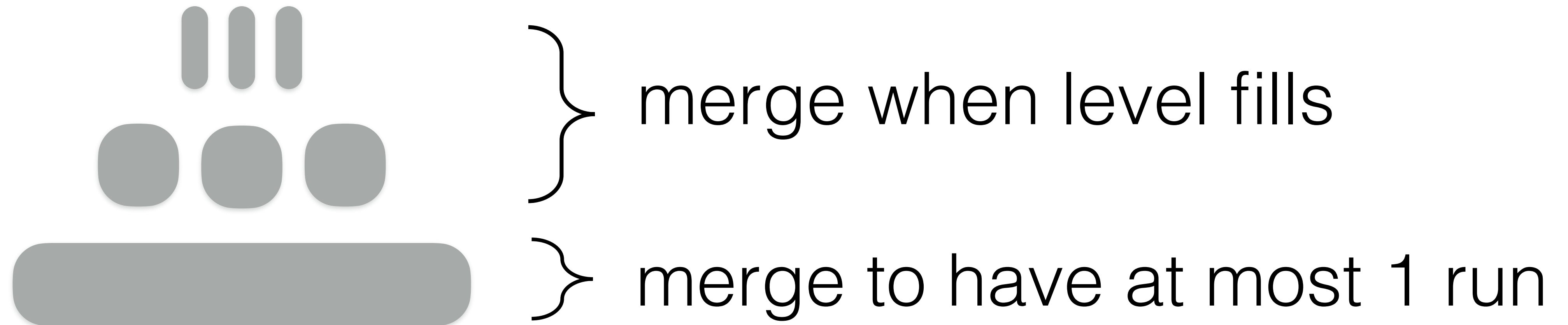
Lazy Leveling
mixed-optimized

Leveling
read-optimized

Lazy Leveling



Lazy Leveling



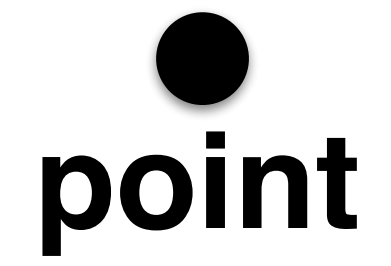
●
point

→
long range

→
short range

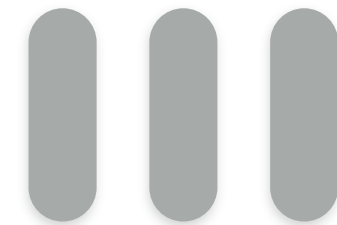
↶
writes



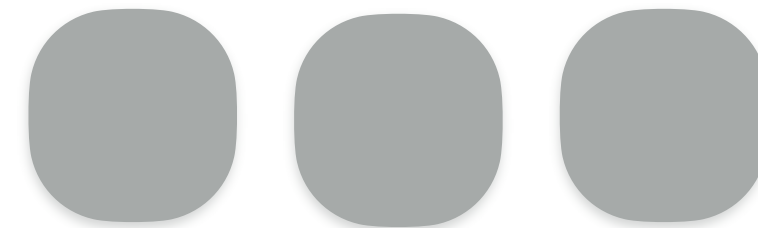


false positive rates

$$O(e^{-x}/R^3)$$



$$O(e^{-x}/R^2)$$



$$O(e^{-x})$$



●
point

false positive rates

exponentially
decreasing

↑
 $O(e^{-x}/R^3)$
 $O(e^{-x}/R^2)$
 $O(e^{-x})$

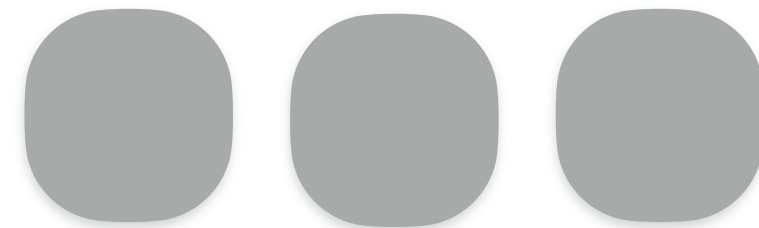
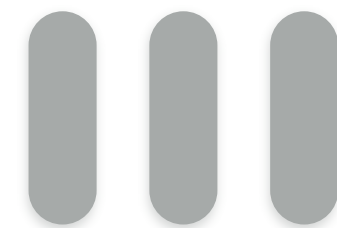


false positive rates

$$O(e^{-x}/R^3)$$

$$O(e^{-x}/R^2)$$

→ $O(e^{-x})$



largest level

●
point

●
point

$O(e^{-x})$

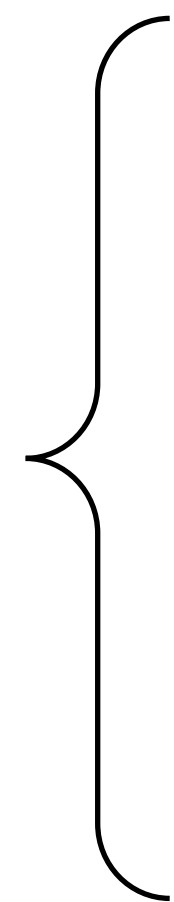


●
point

$O(e^{-x})$

with uniform FPRs

$O(\log_R(N) \cdot R \cdot e^{-x})$



$O(e^{-x})$

$O(e^{-x})$

$O(e^{-x})$



●
point

$O(e^{-x})$

→
long range

→
short range

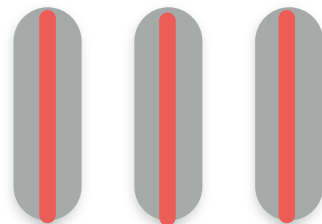
↶
writes



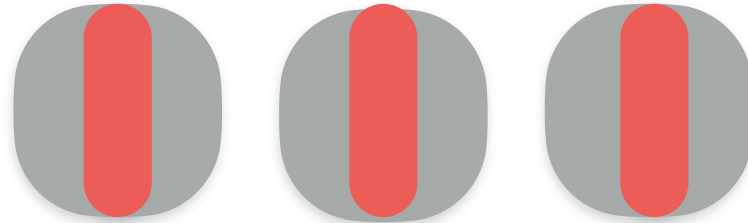
→
long range

target range

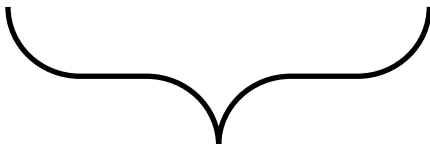
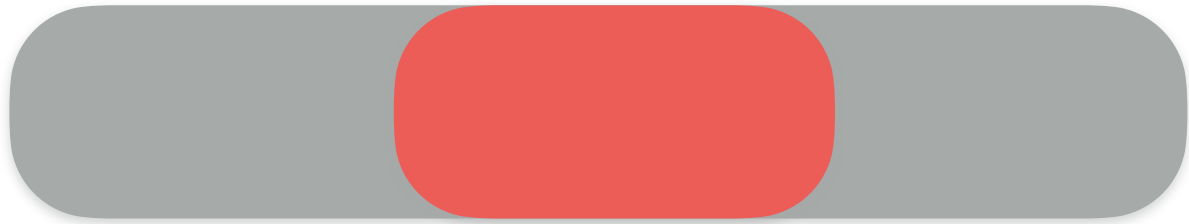
$O(s/R^2)$



$O(s/R)$



$O(s)$



target key range

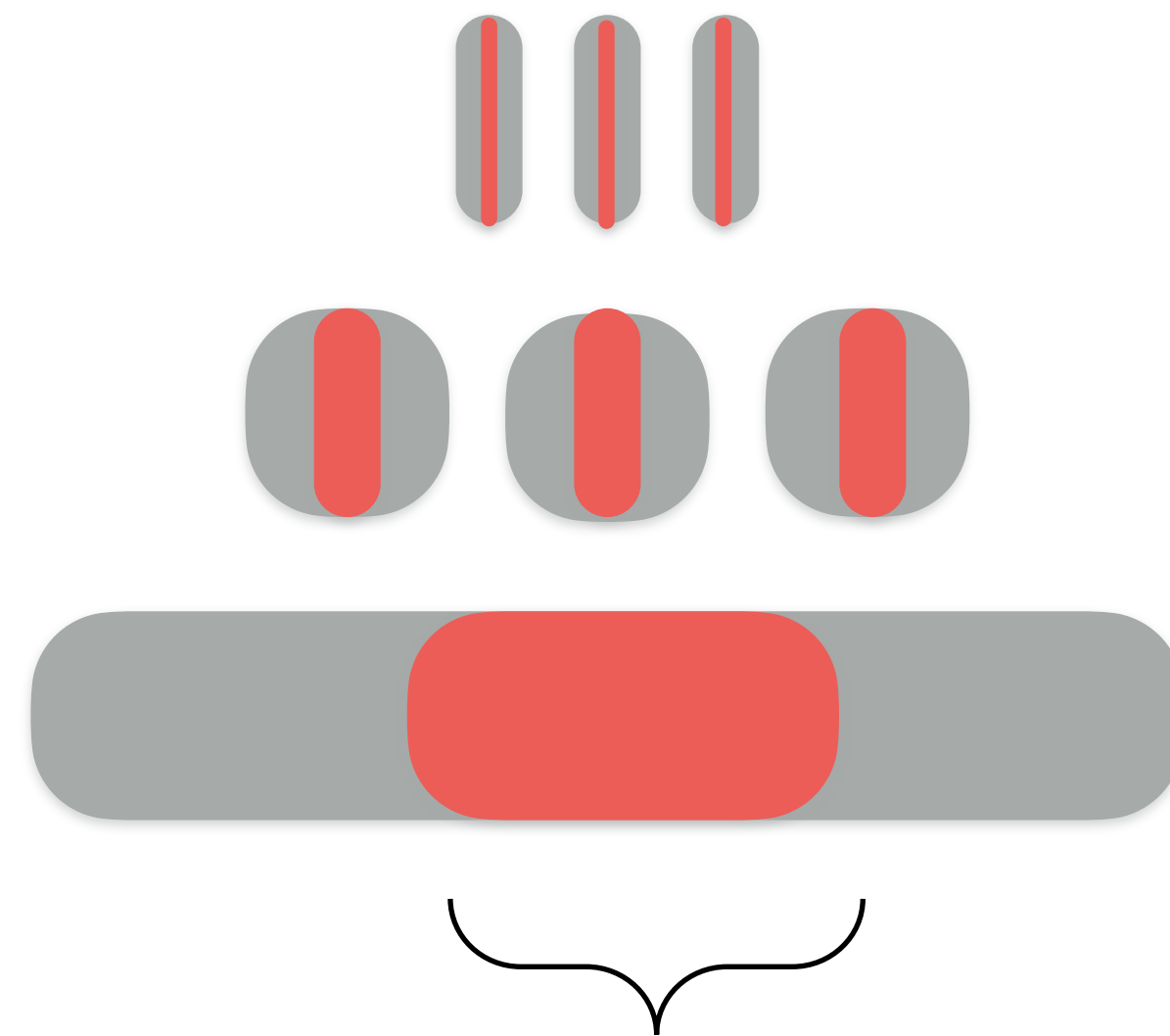

long range

target range

$$O(s/R^2)$$

$$O(s/R)$$

$$O(s)$$



largest level

target key range

●
point

$O(e^{-x})$

→
long range

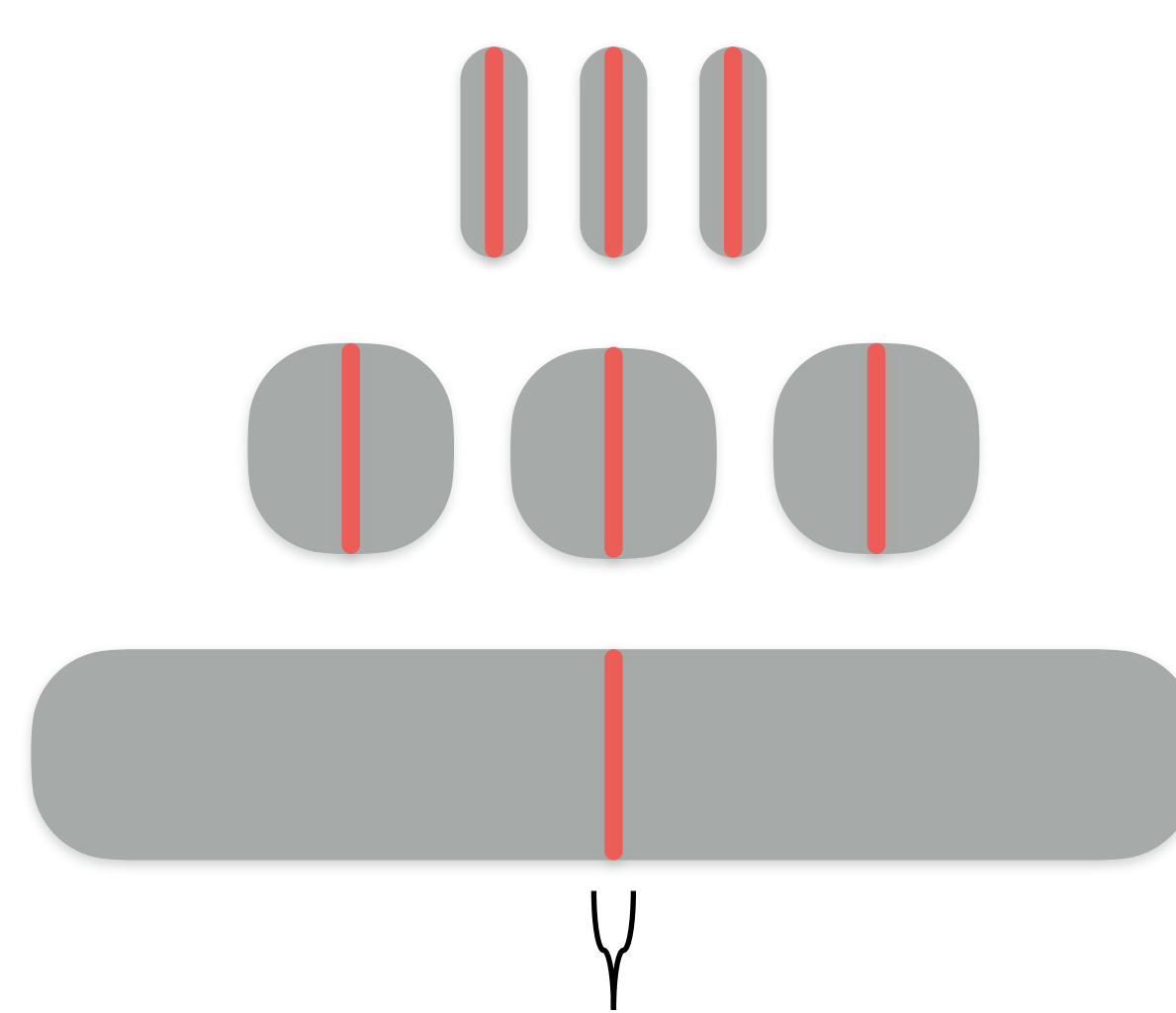
$O(s)$

→
short range

↶
writes



→
short range



$O(R)$

$O(R)$

1

$O(1 + R \cdot (\log_R(N) - 1))$

●
point

$O(e^{-x})$

→
long range

$O(s)$

→
short range

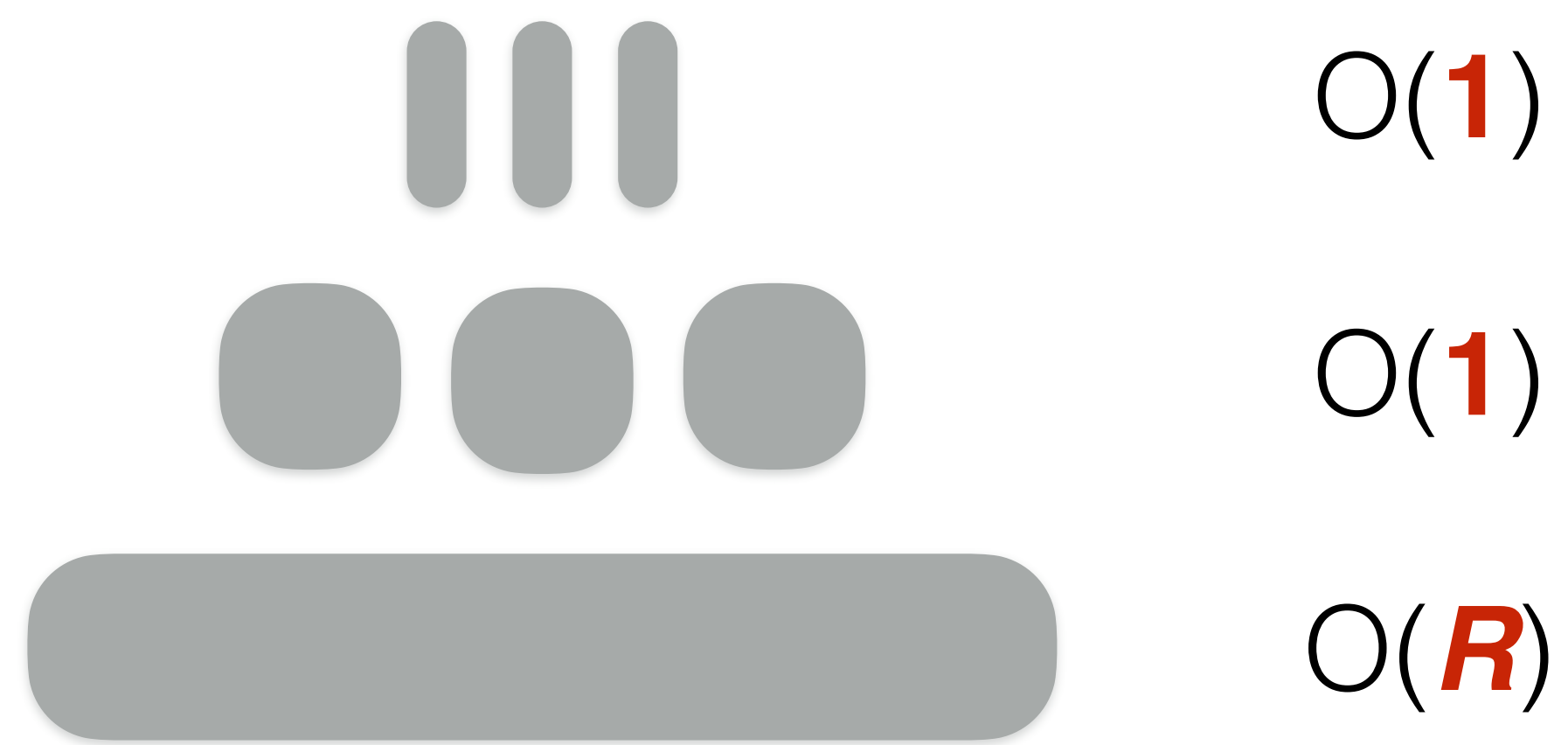
$O(1 + R \cdot (\log_R(N) - 1))$

↶
writes





write-amplification





write-amplification



$$\left. \begin{array}{l} O(1) \\ O(1) \\ O(R) \end{array} \right\} O(R + \log_R(N))$$

●
point

$O(e^{-x})$

→
long range

$O(s)$

→
short range

$O(1 + R \cdot (\log_R(N) - 1))$

↶
writes

$O(R + \log_R(N))$



●
point

→
long range

→
short range

↶
writes

Lazy Leveling

$O(e^{-x})$

$O(s)$

$O(1 + R \cdot (\log_R(N) - 1))$

$O(R + \log_R(N))$

=

=

V

Λ


Leveling

$O(e^{-x})$

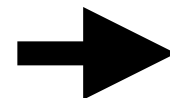
$O(s)$

$O(\log_R(N))$

$O(R \cdot \log_R(N))$


point


long range


short range


writes

Tiering

$O(R \cdot e^{-x})$

$O(R \cdot s)$

$O(R \cdot \log_R(N))$

$O(\log_R(N))$

\vee

\vee

\vee

\wedge

Lazy Leveling

$O(e^{-x})$

$O(s)$

$O(1 + R \cdot (\log_R(N) - 1))$

$O(R + \log_R(N))$

Leveling

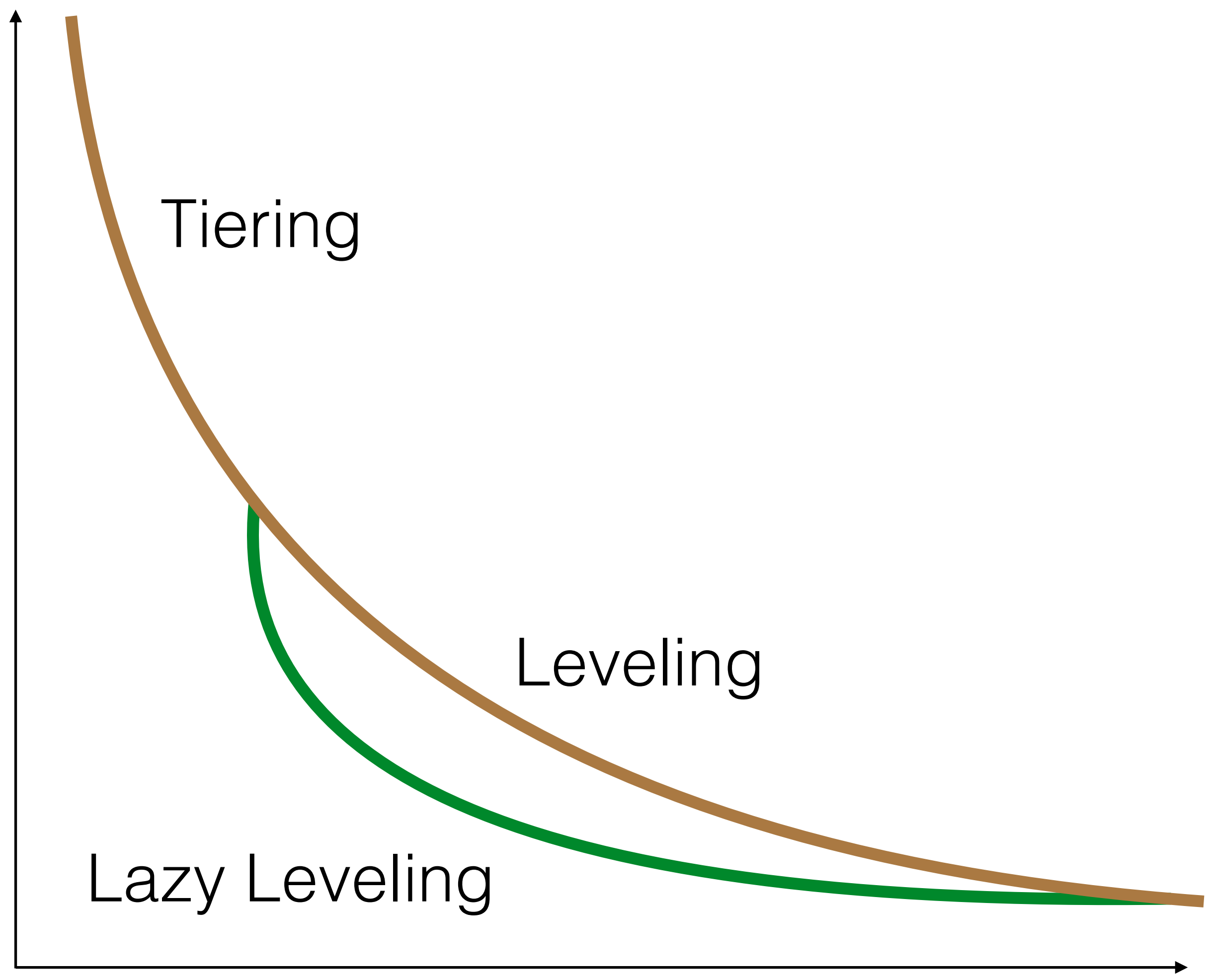
$O(e^{-x})$

$O(s)$

$O(\log_R(N))$

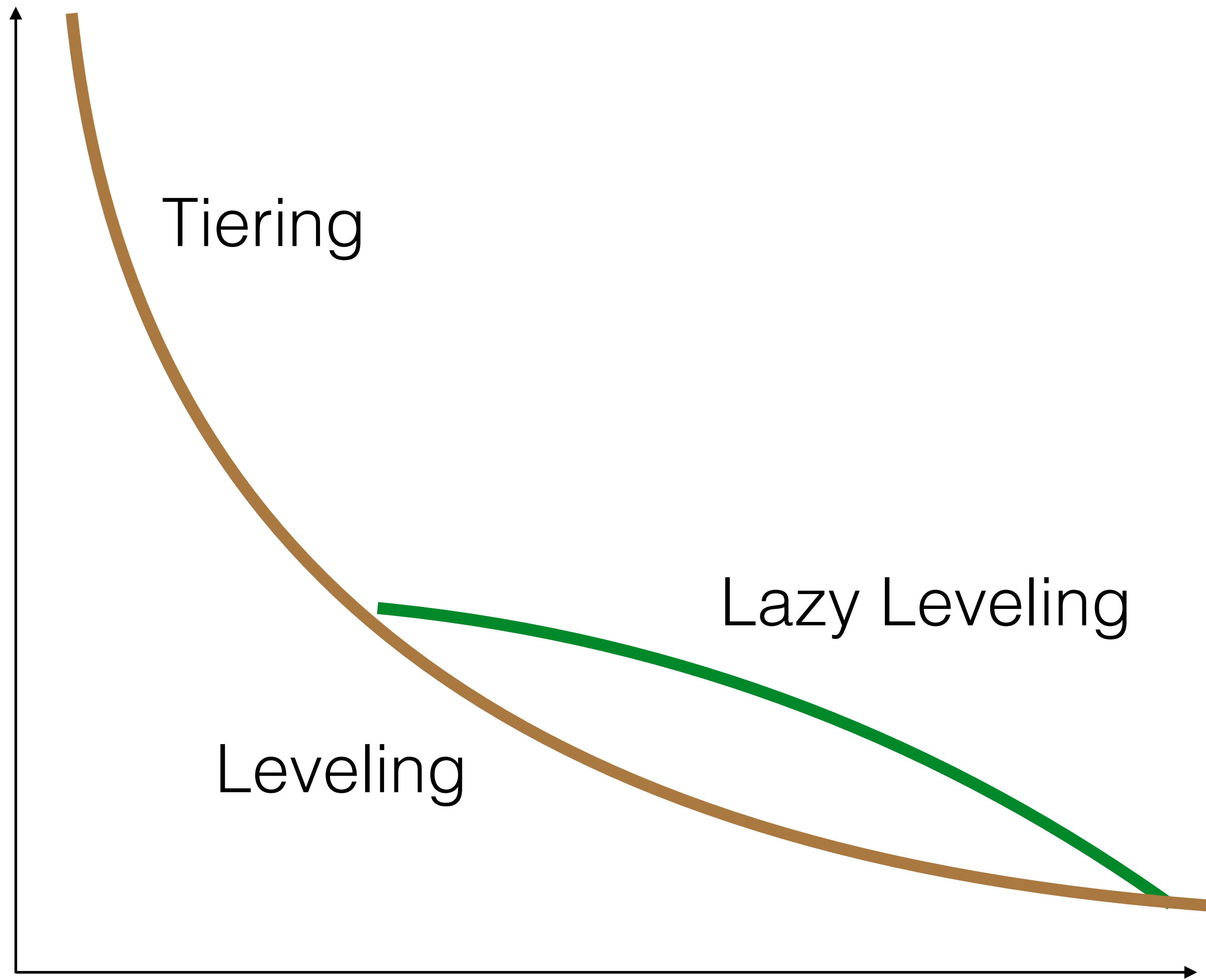
$O(R \cdot \log_R(N))$

●
point

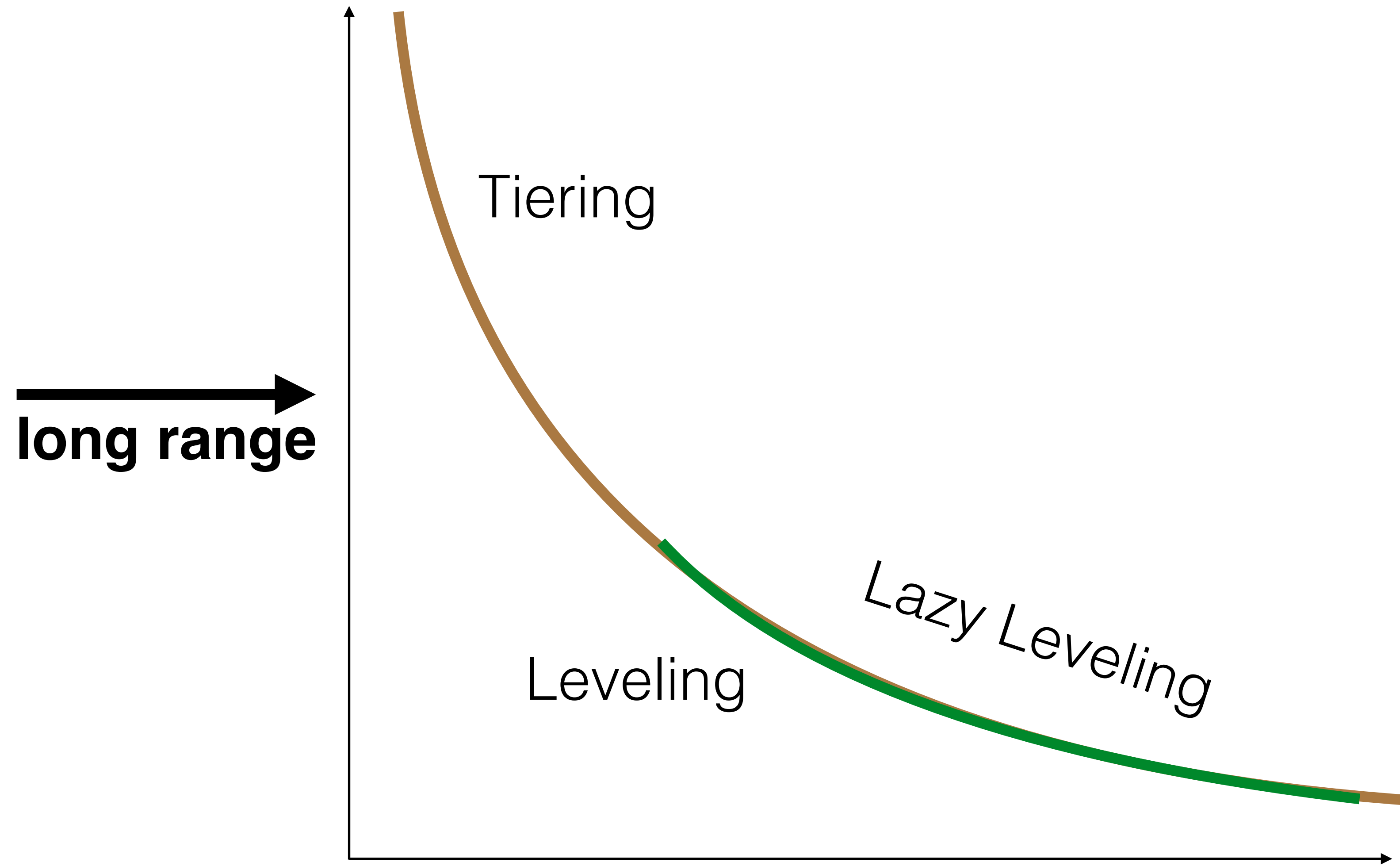


↖
writes

short range →



writes



writes

Tiering

Lazy Leveling

Leveling

Tiering
writes



Lazy Leveling

Leveling

Tiering
writes 

Lazy Leveling

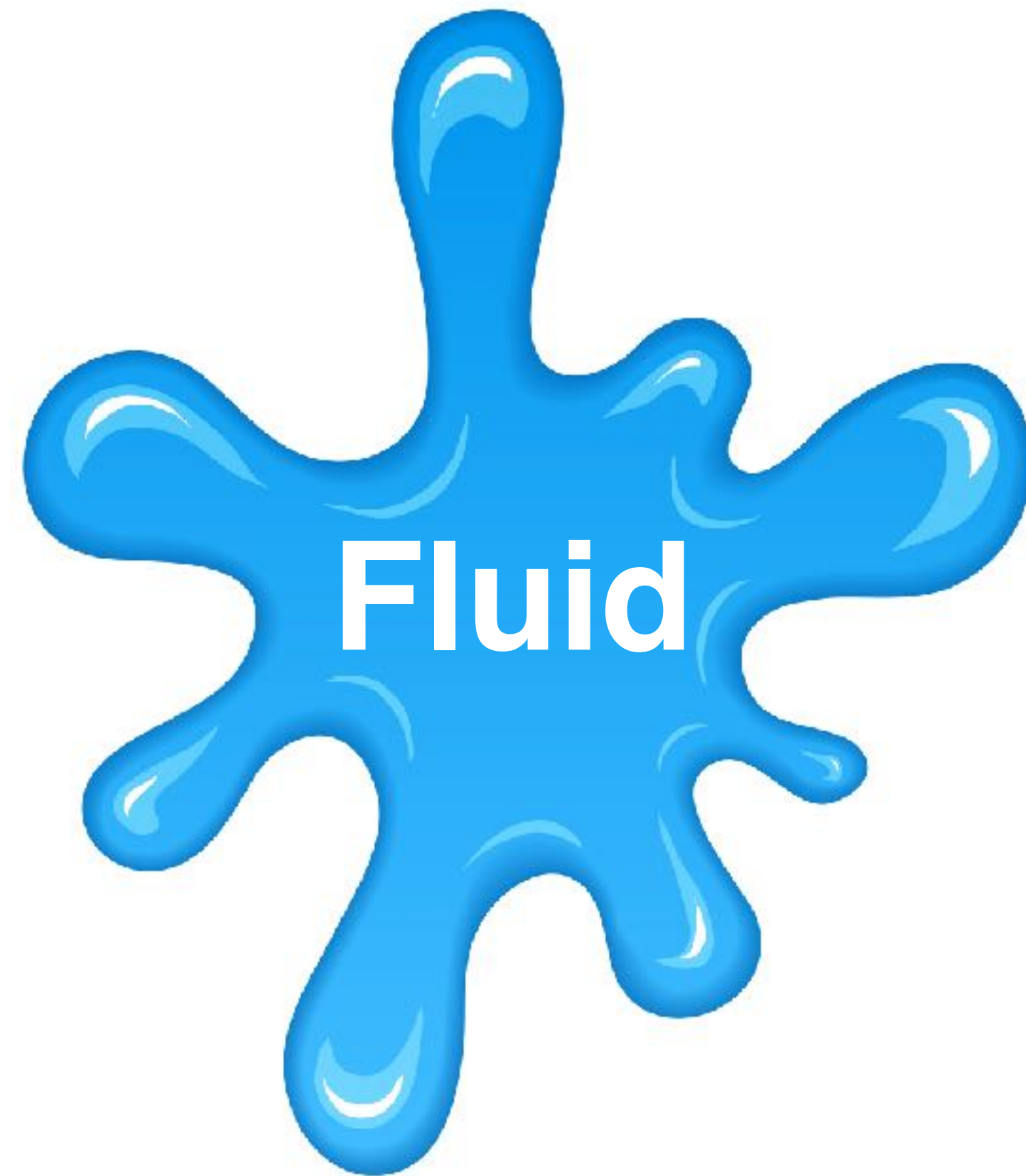
Leveling
short range 

Tiering
writes 

Lazy Leveling
writes & point 

Leveling
short range 

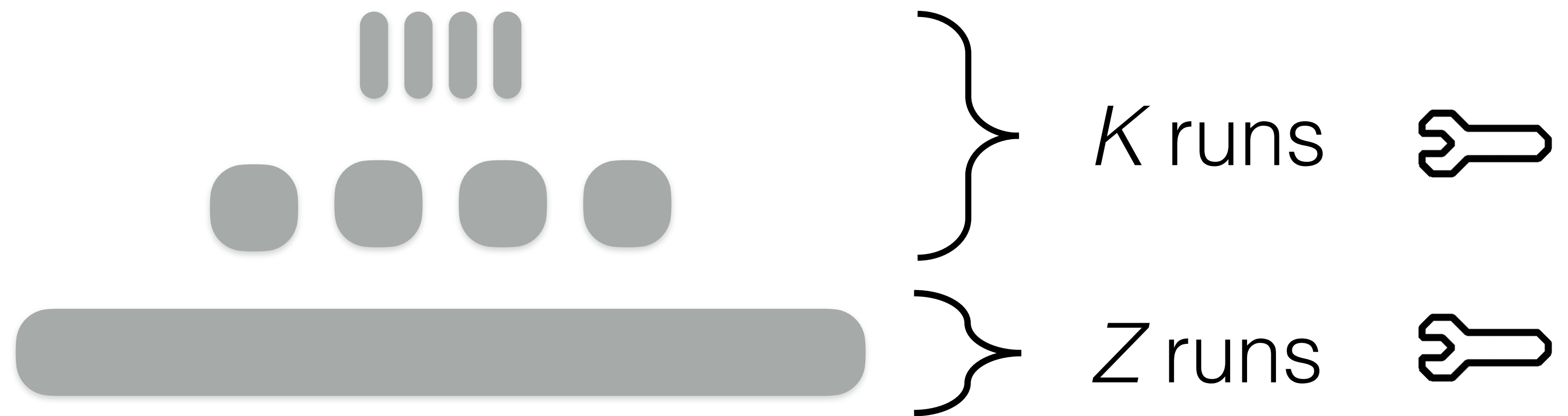
Lazy Leveling



Tiering

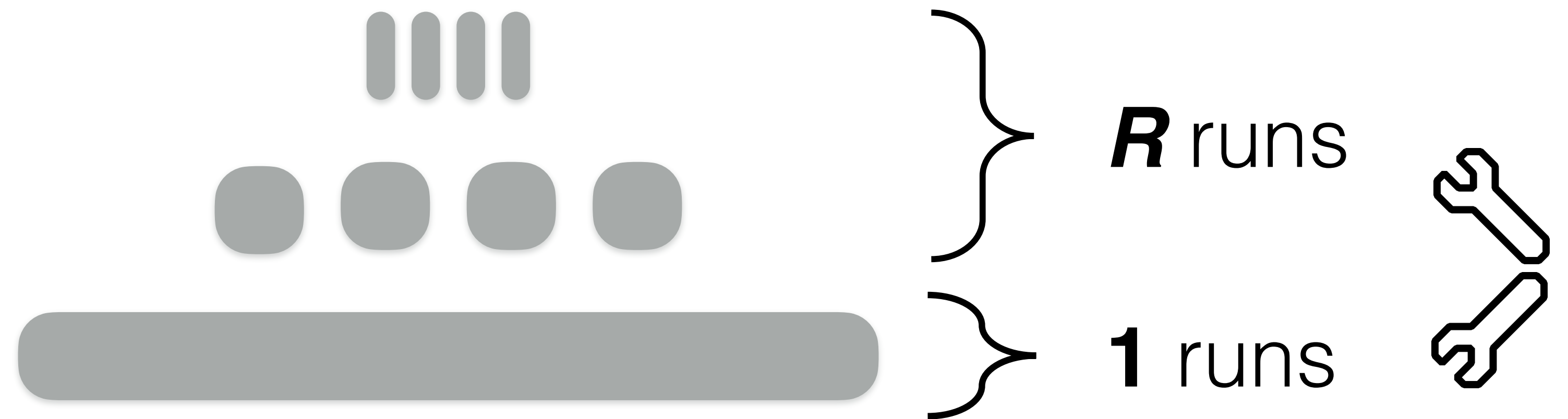
Leveling

Fluid LSM-Tree



Fluid LSM-Tree

Lazy Leveling



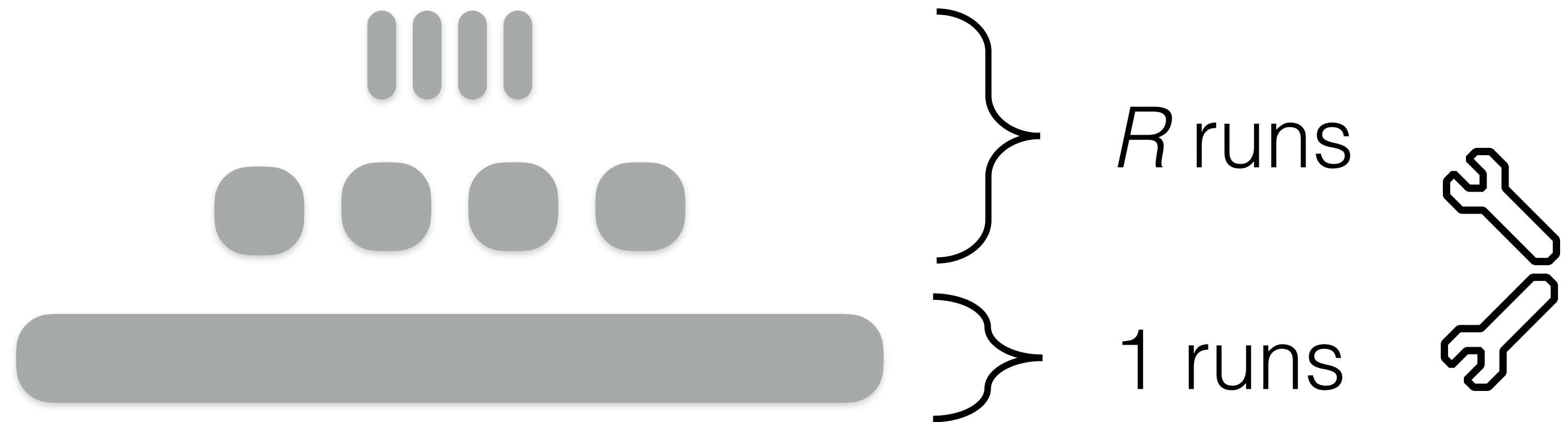
point

long range

short range

writes

Lazy Leveling



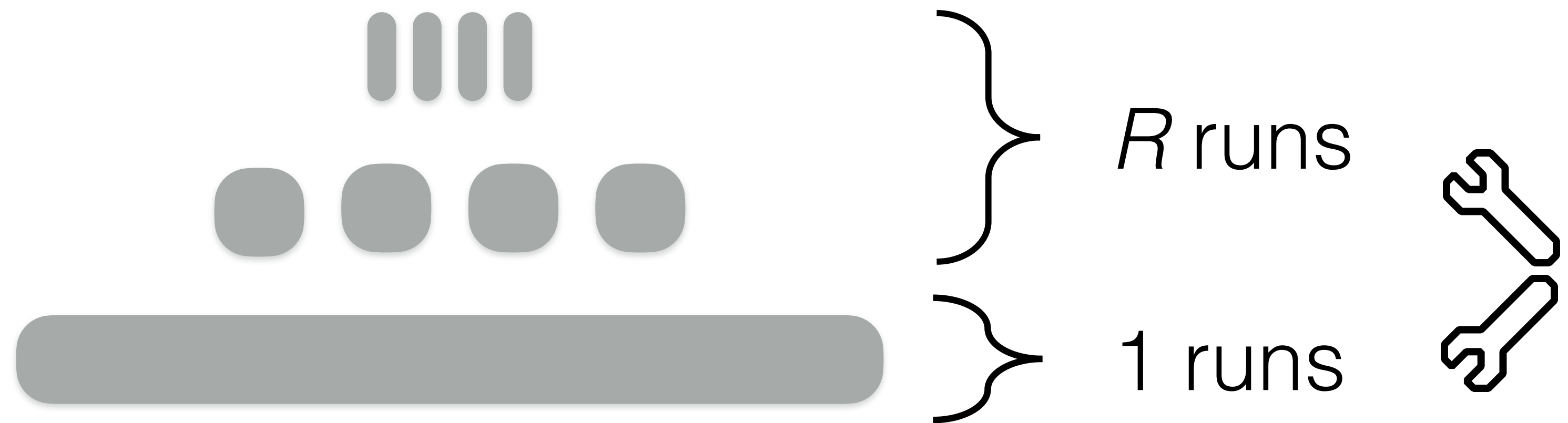
point

long range

optimize
short range

writes

Lazy Leveling



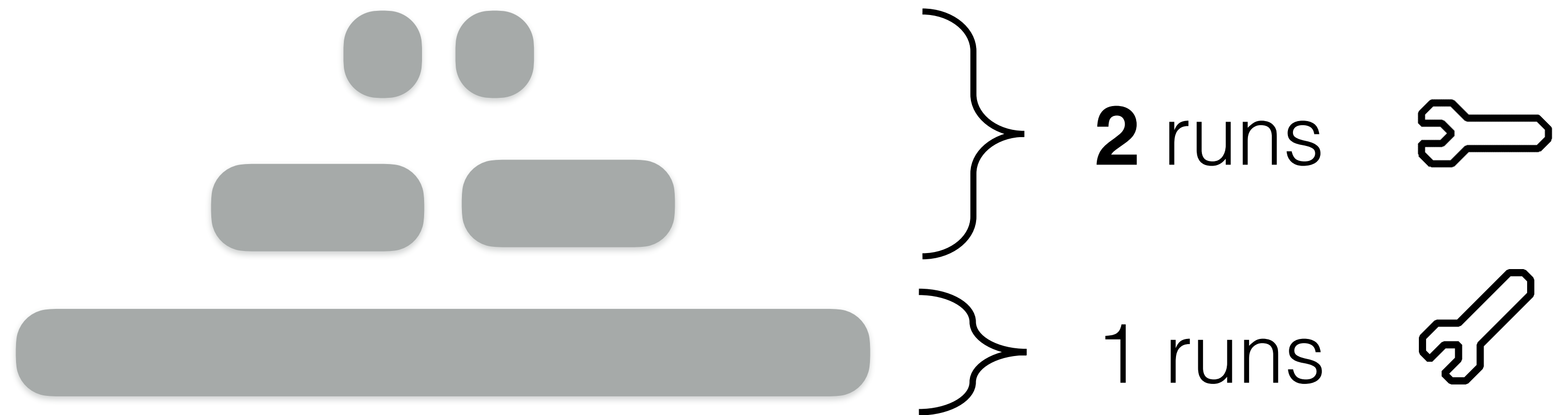
point

long range

optimize
short range

writes

Lazy Leveling



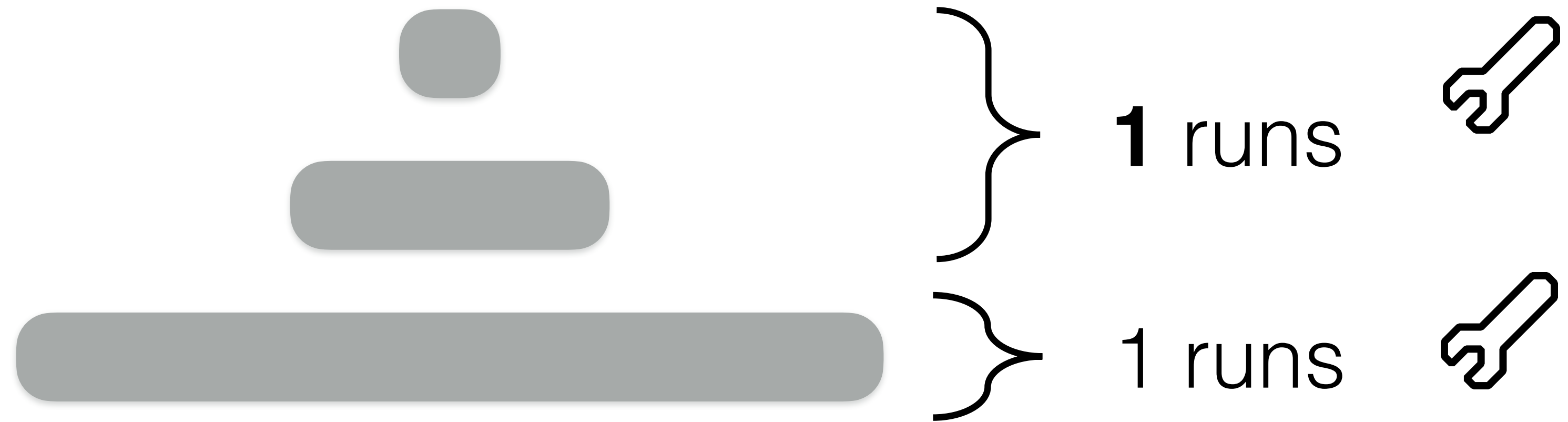
point

long range

optimize
short range

writes

Leveling



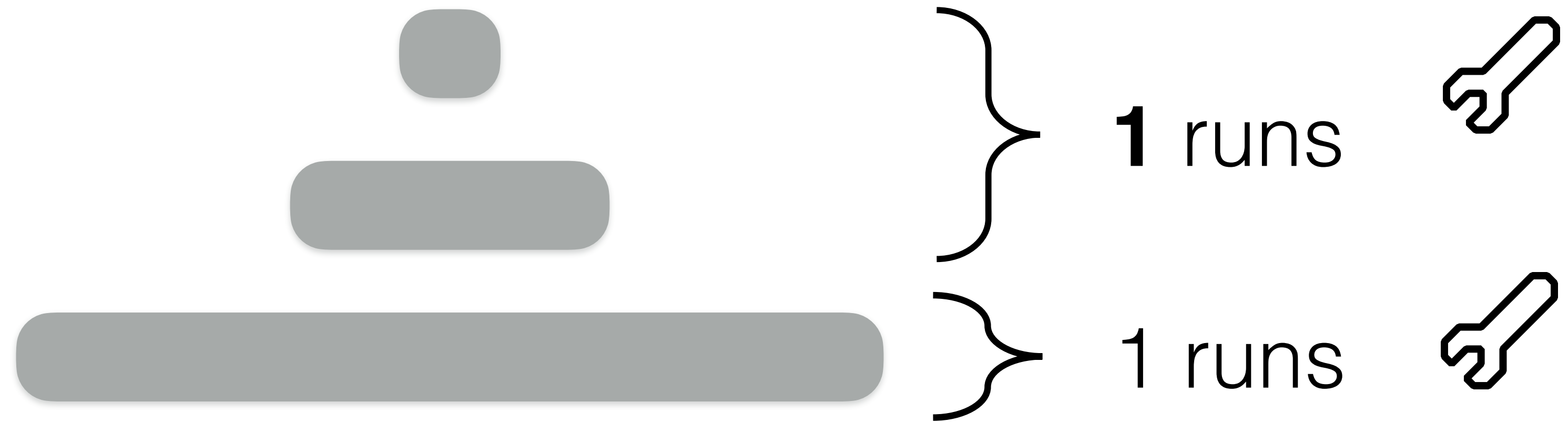
point

long range

short range

optimize
writes

Leveling



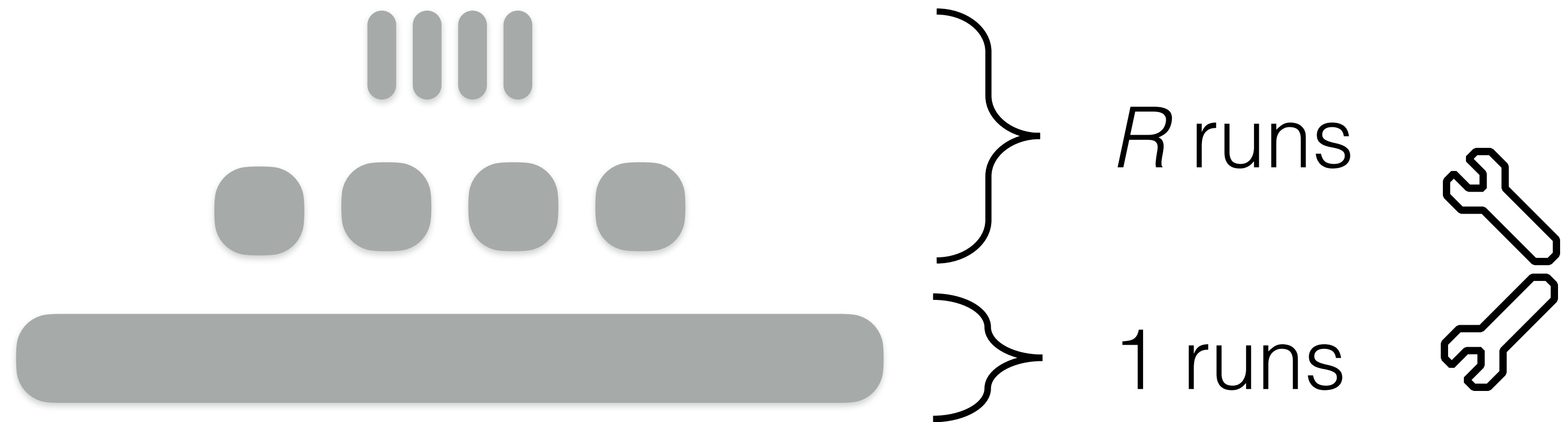
point

long range

short range

optimize
writes

Lazy Leveling



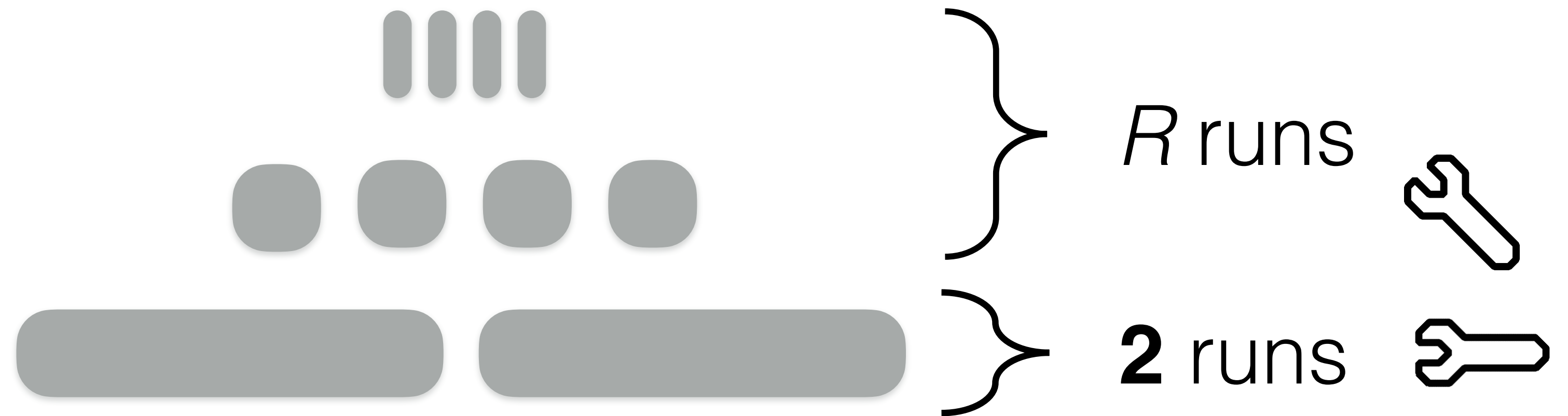
point

long range

short range

optimize
writes

Lazy Leveling



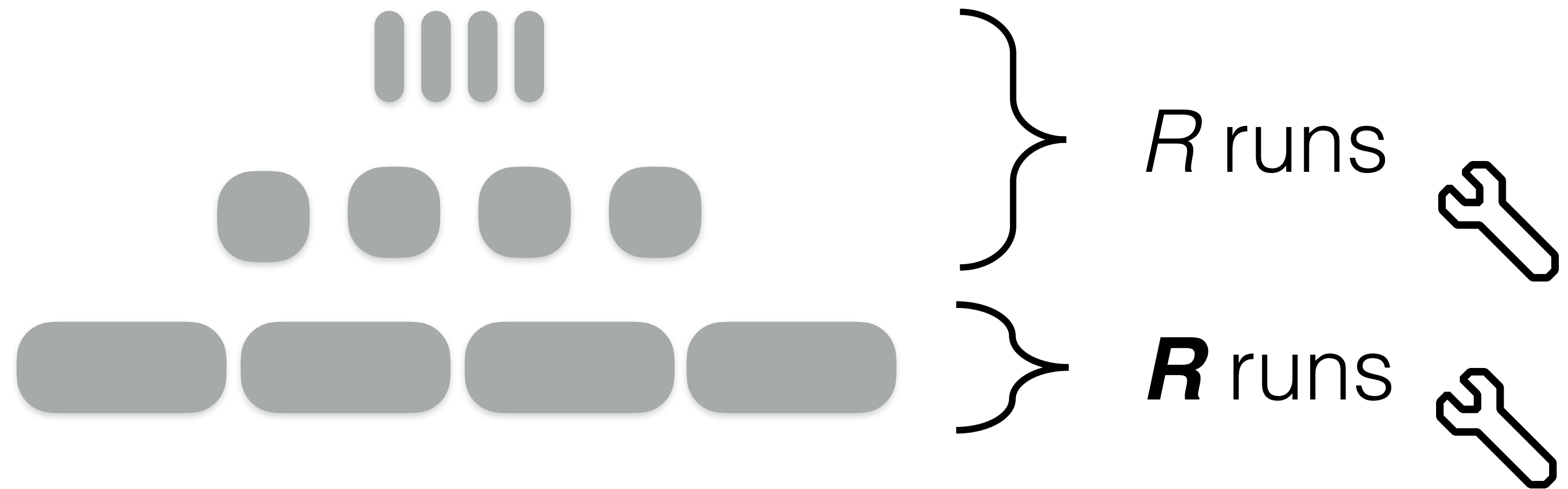
point

long range

short range

optimize
writes

Tiering



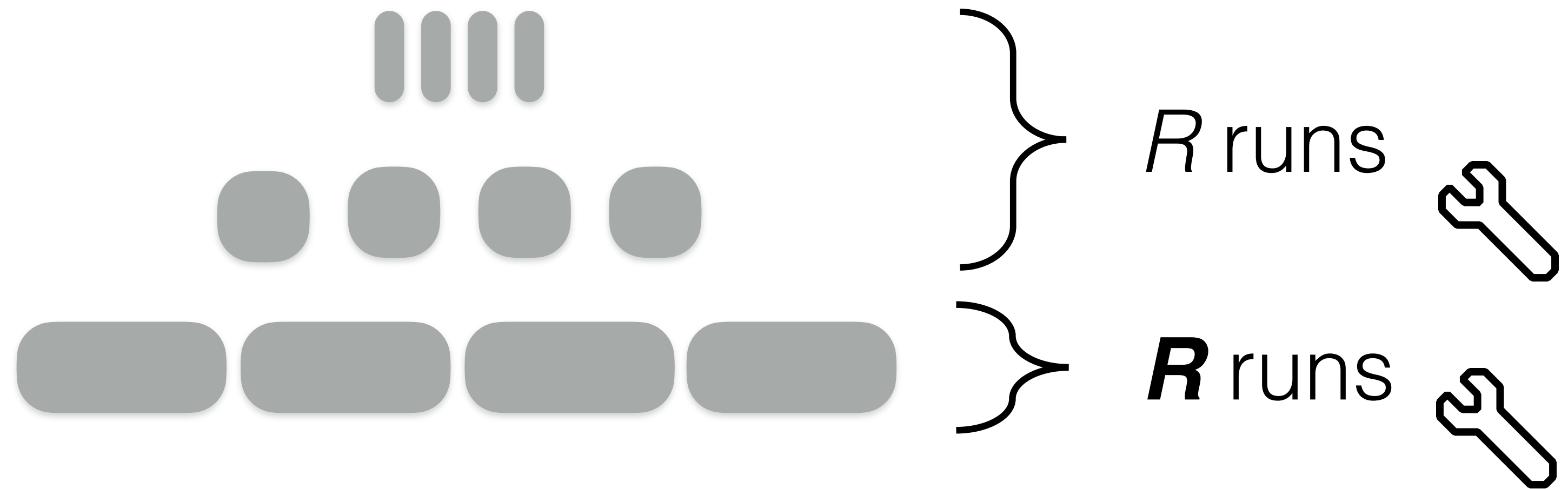
optimize
point

long range

short range

writes

Tiering



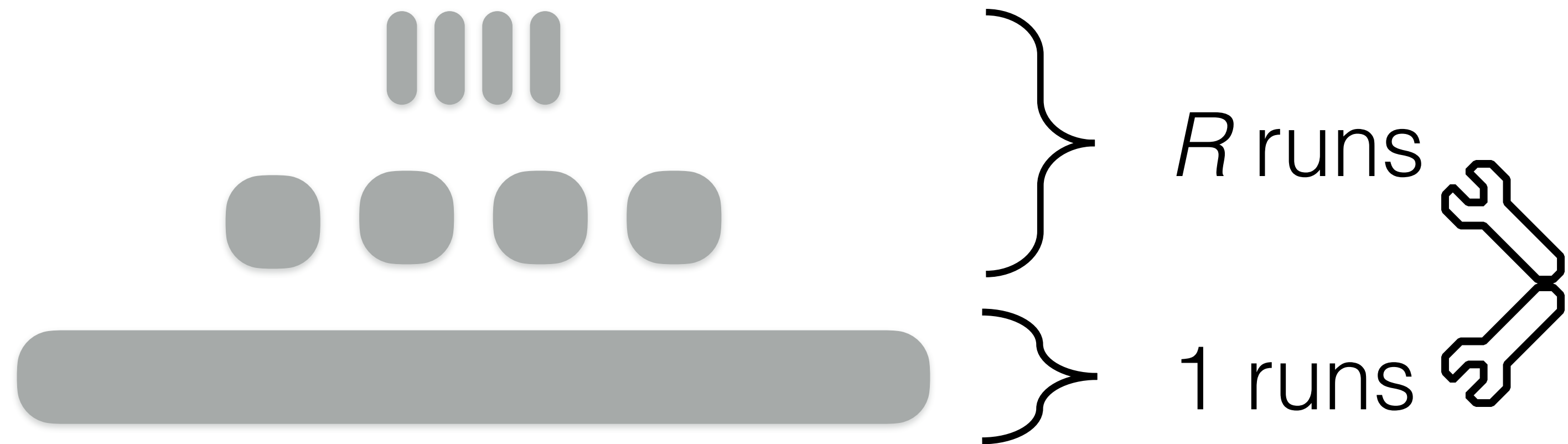
optimize
point

long range

short range

writes

Lazy Leveling

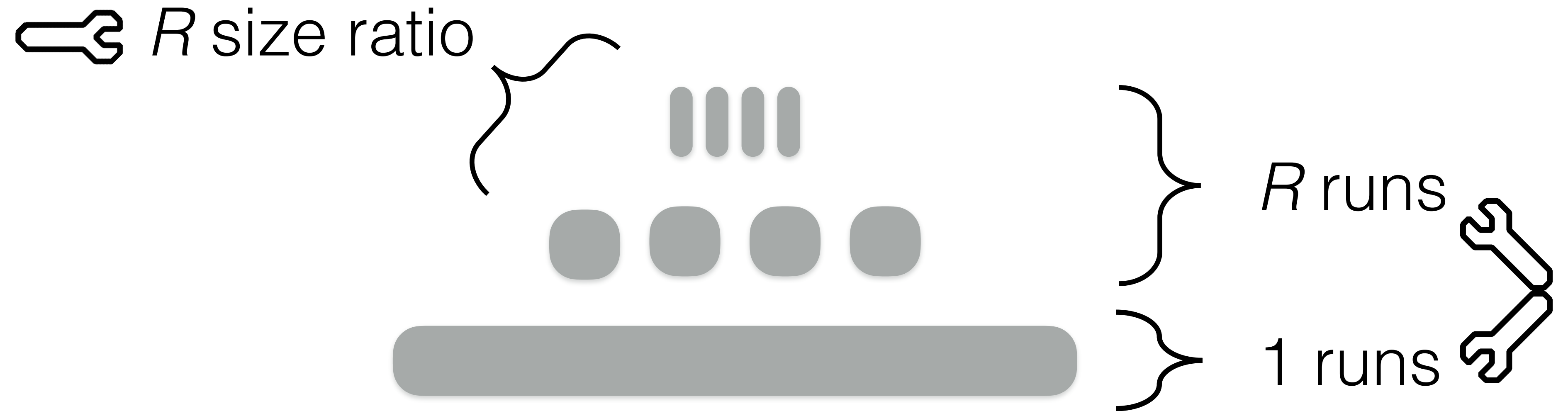


optimize
point

long range

short range

writes



Lazy Leveling

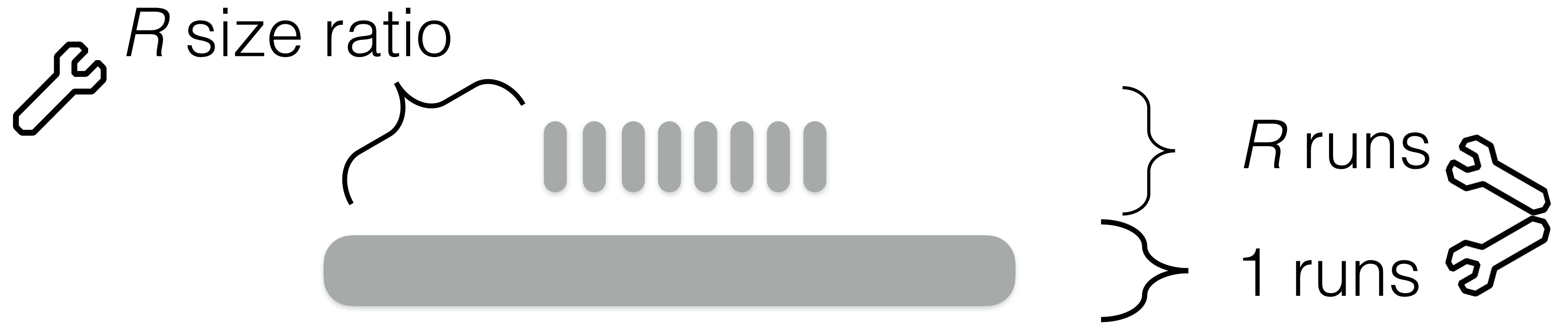
optimize
point

long range

short range


writes

Lazy Leveling





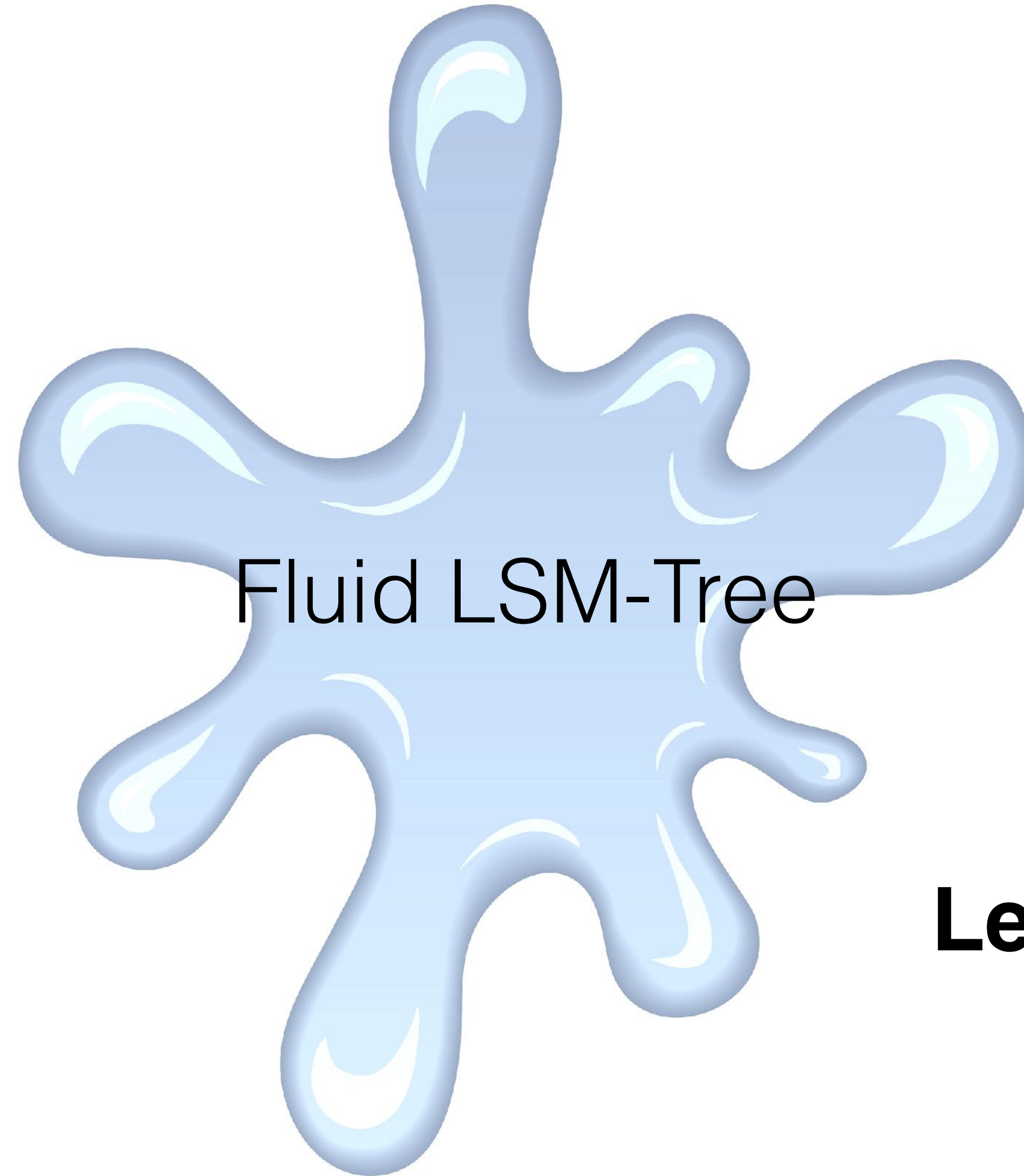
Fluid LSM-Tree

 R size ratio

 K runs at smaller levels

 Z runs at largest level

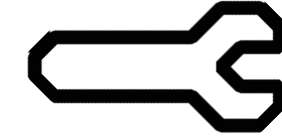
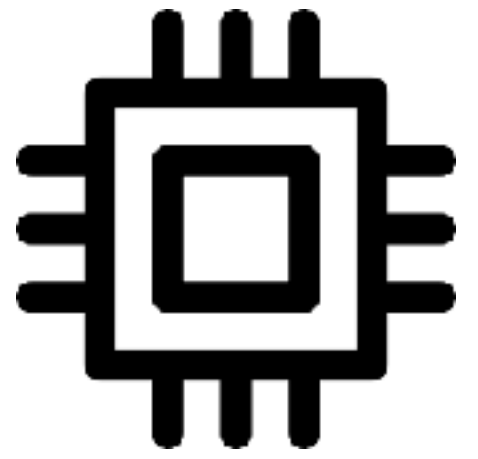
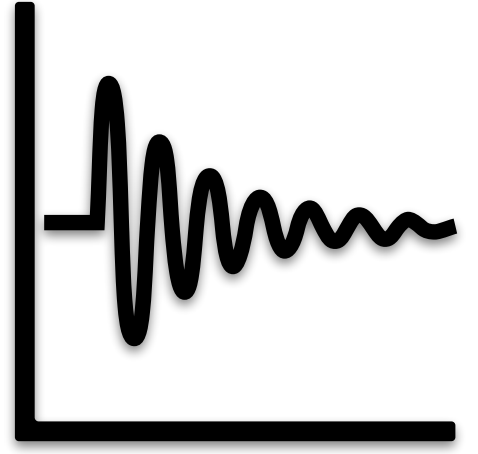
Lazy Leveling

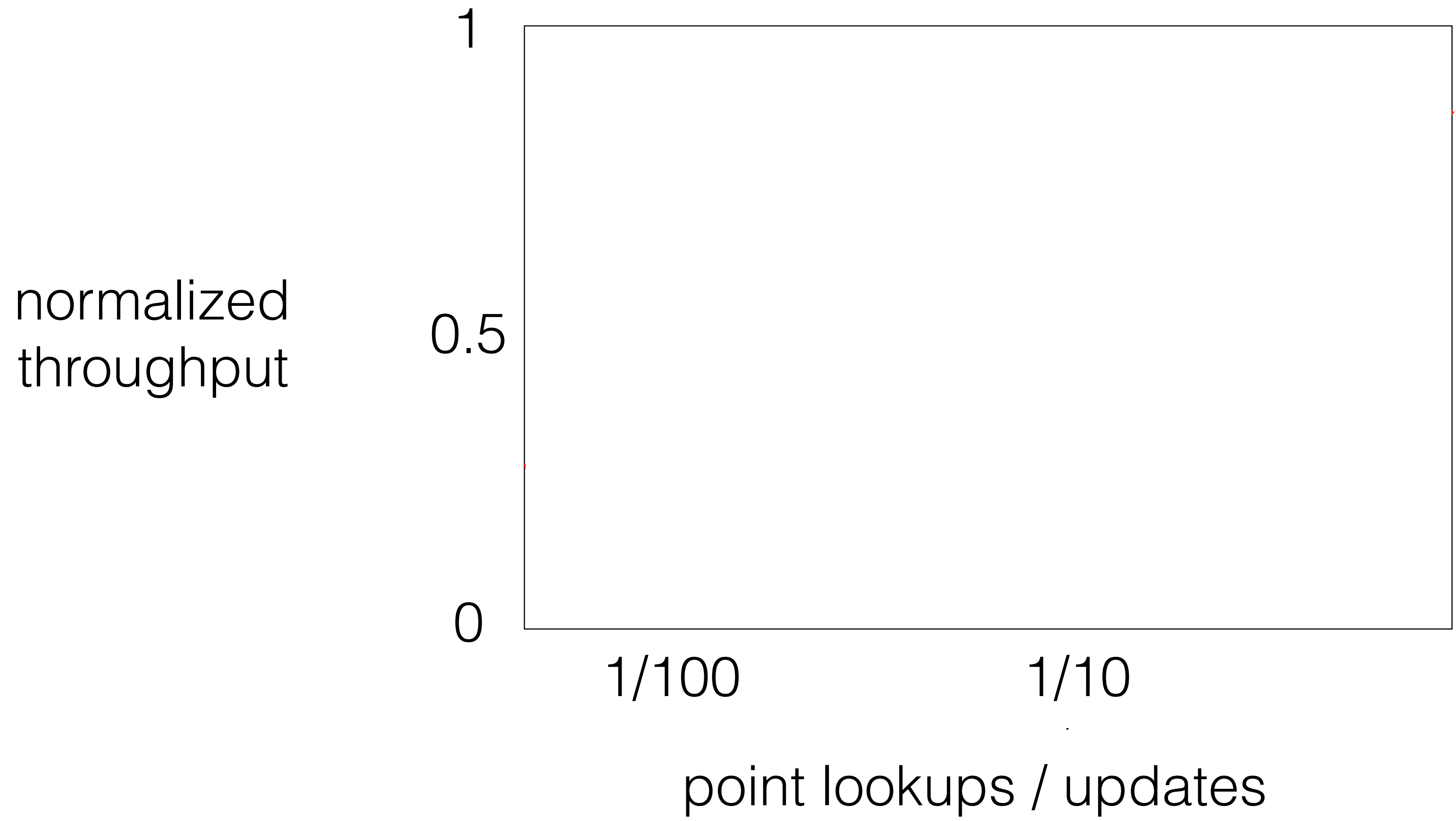


Fluid LSM-Tree

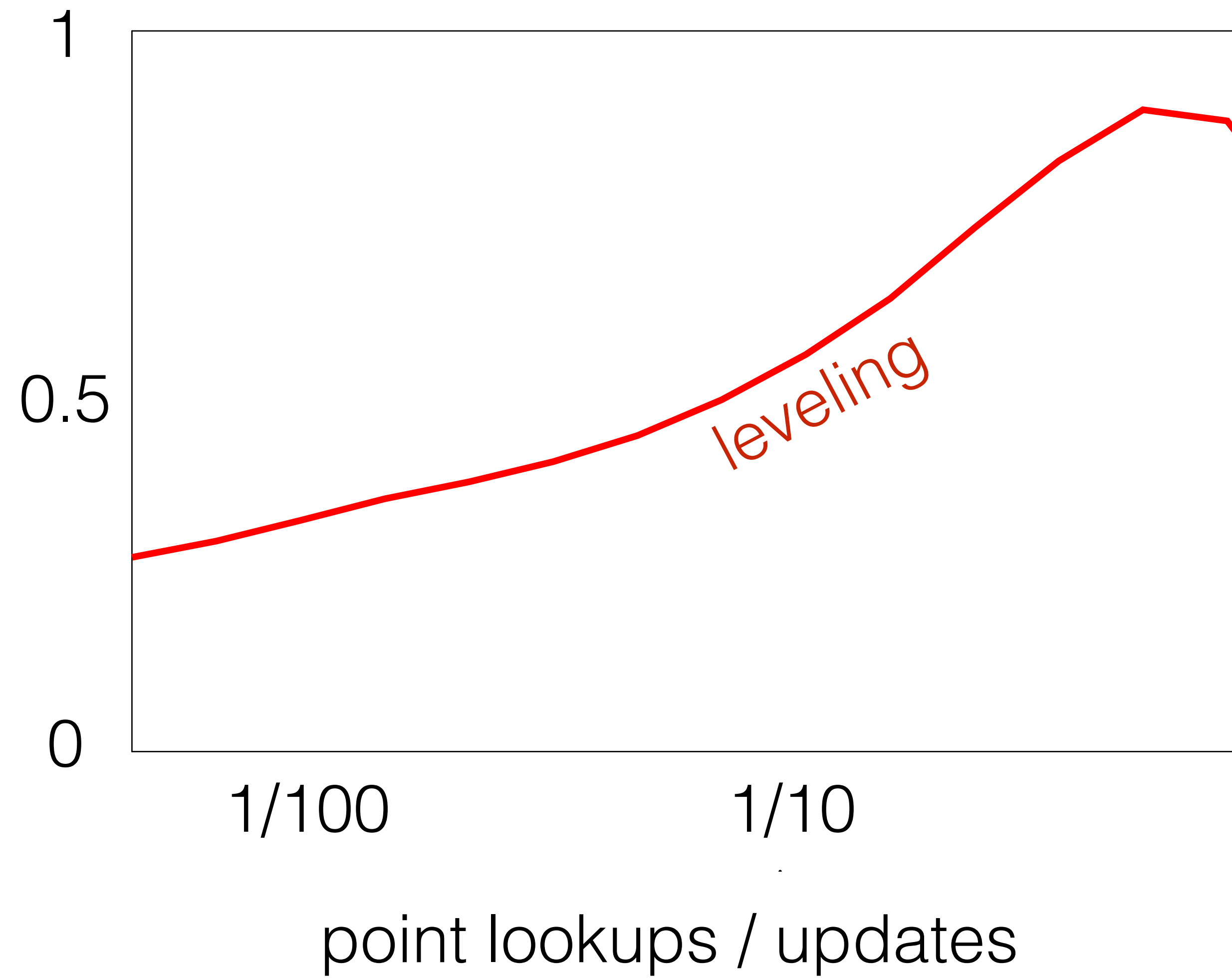
Tiering

Leveling

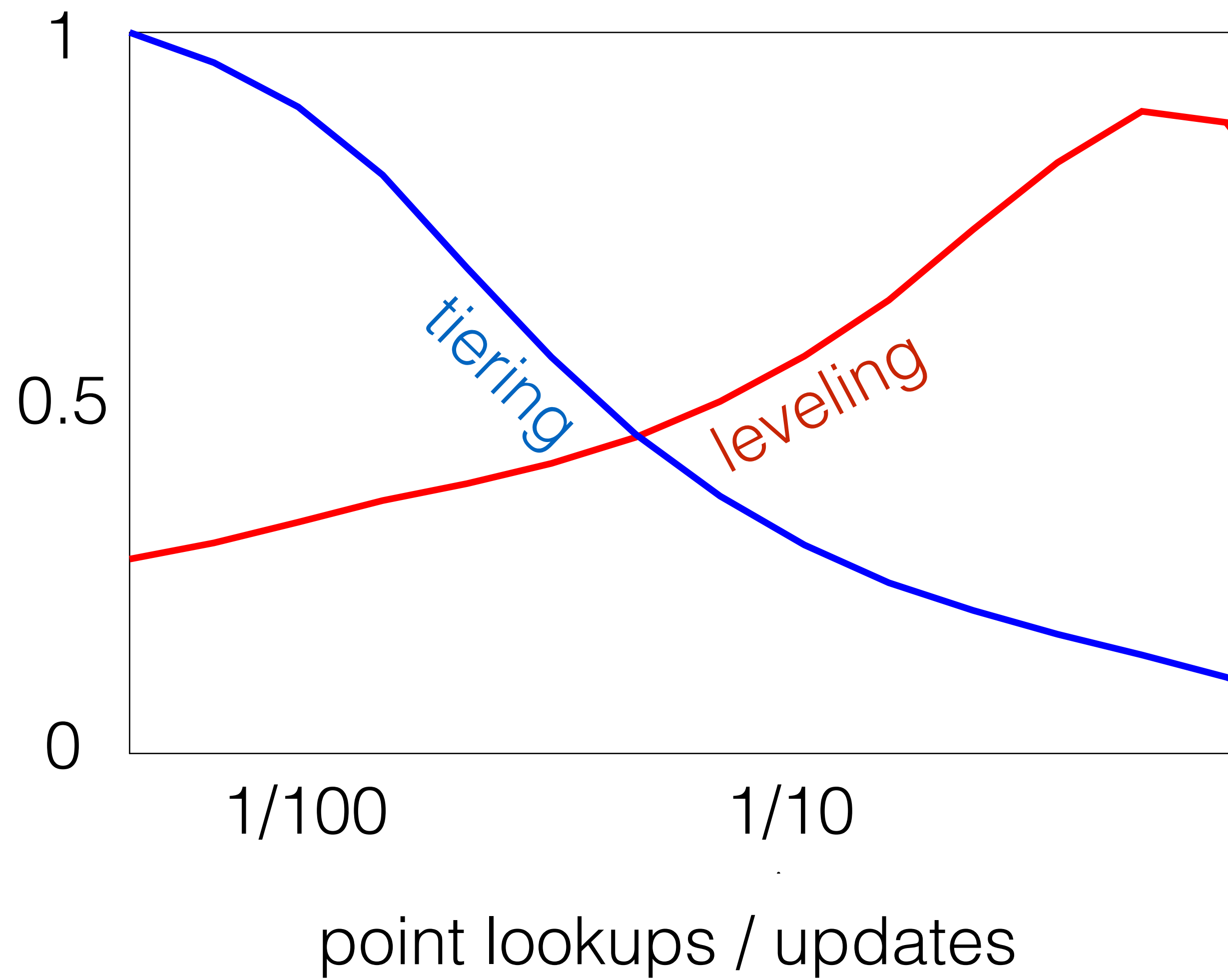




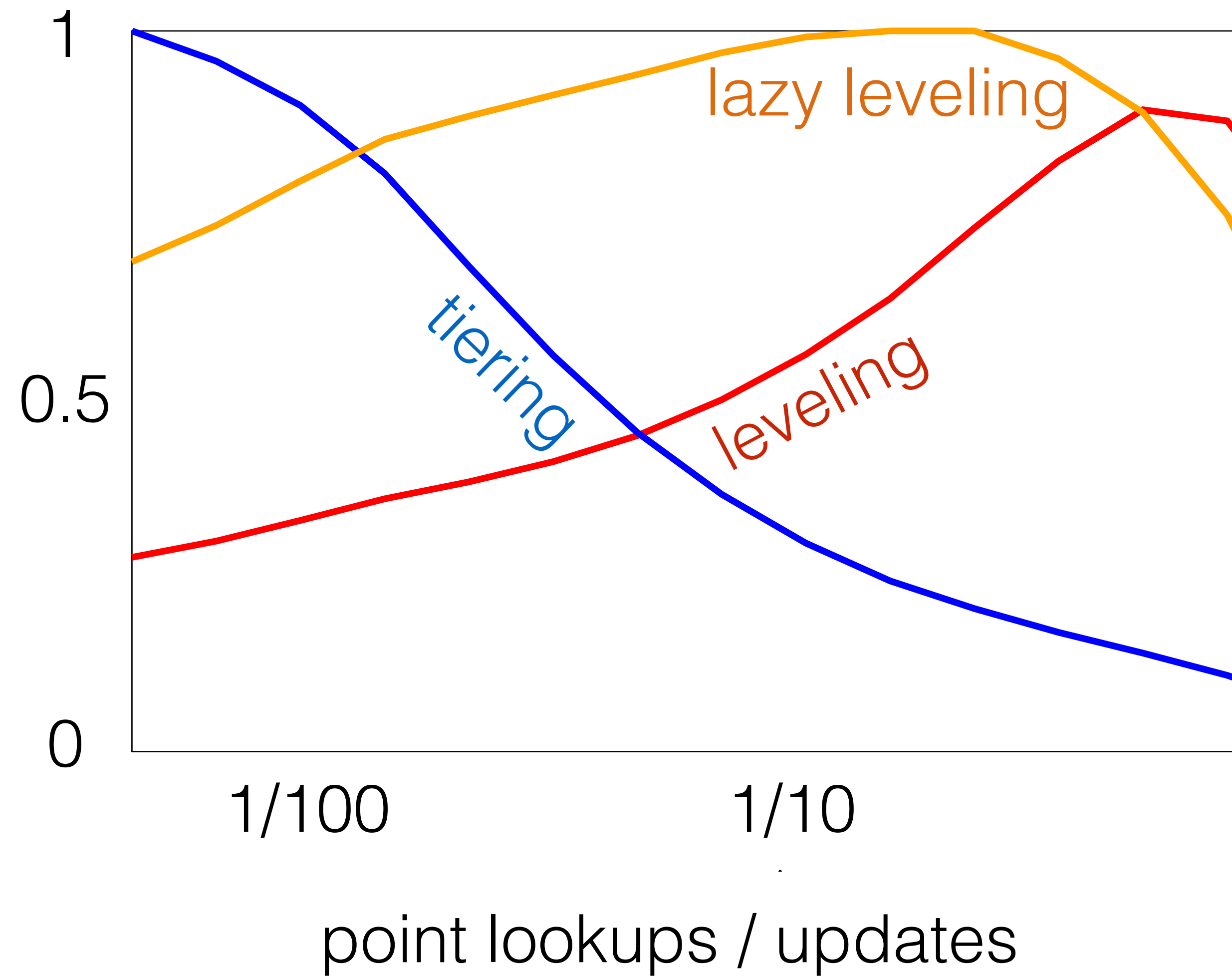
normalized
throughput

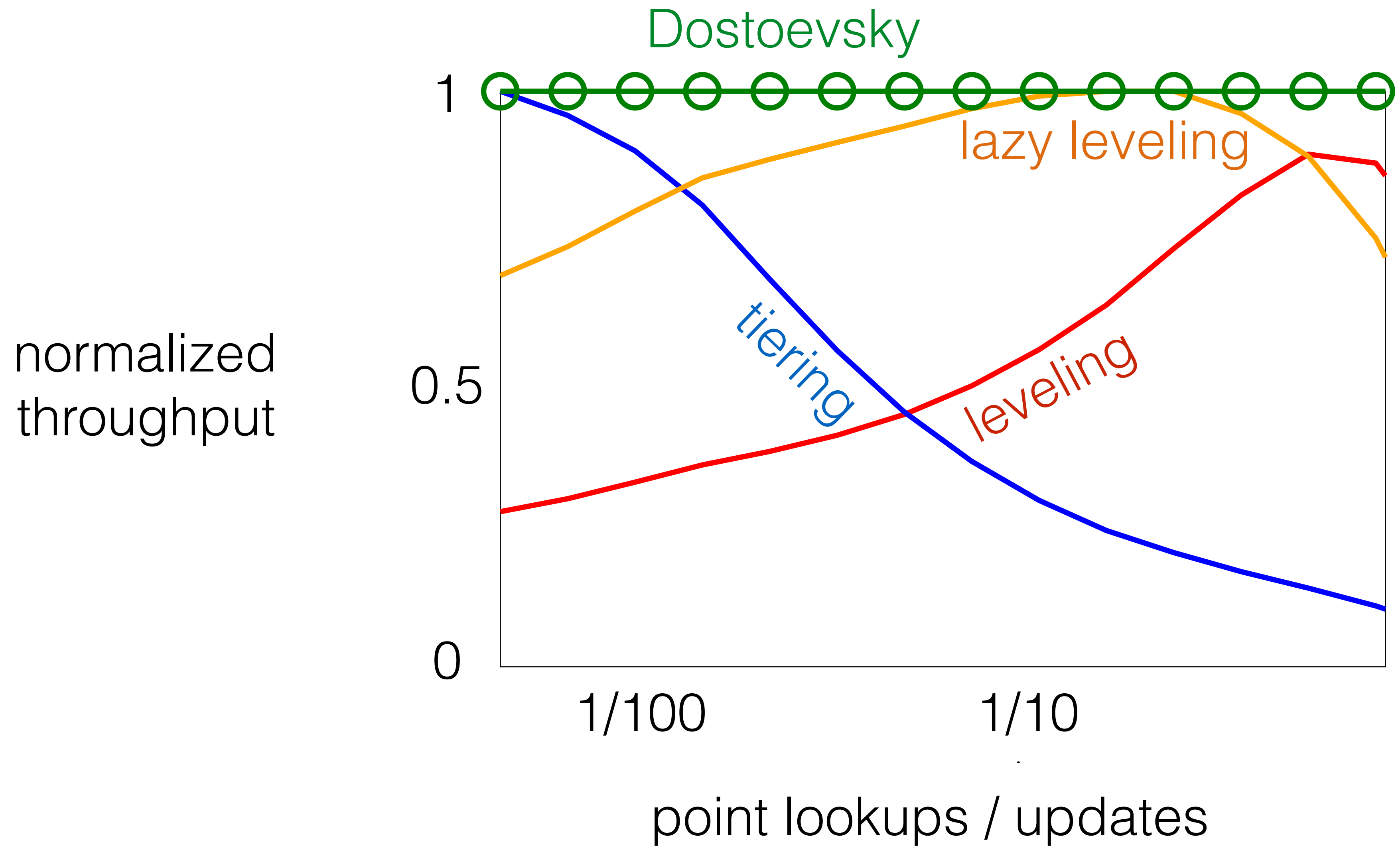


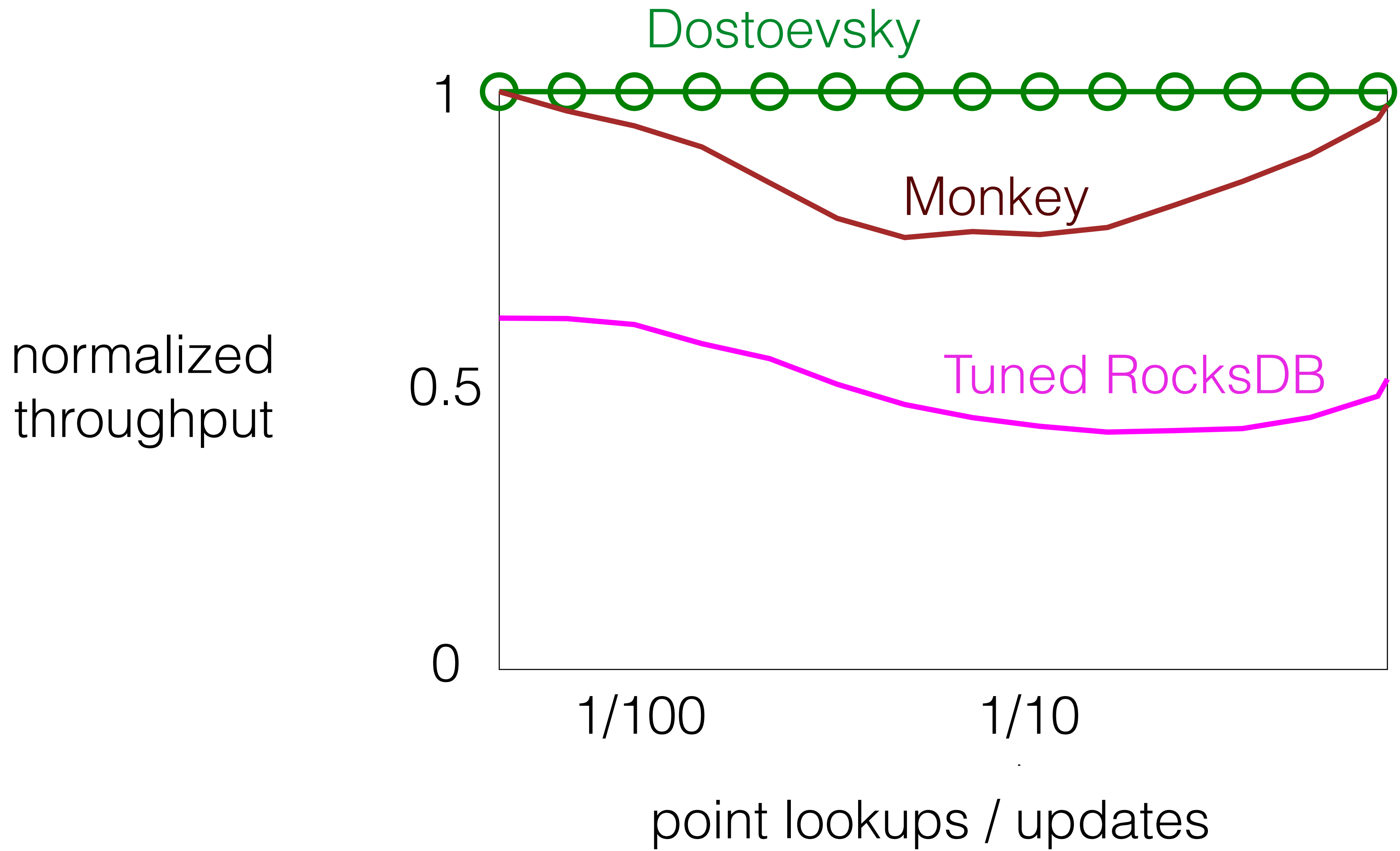
normalized
throughput

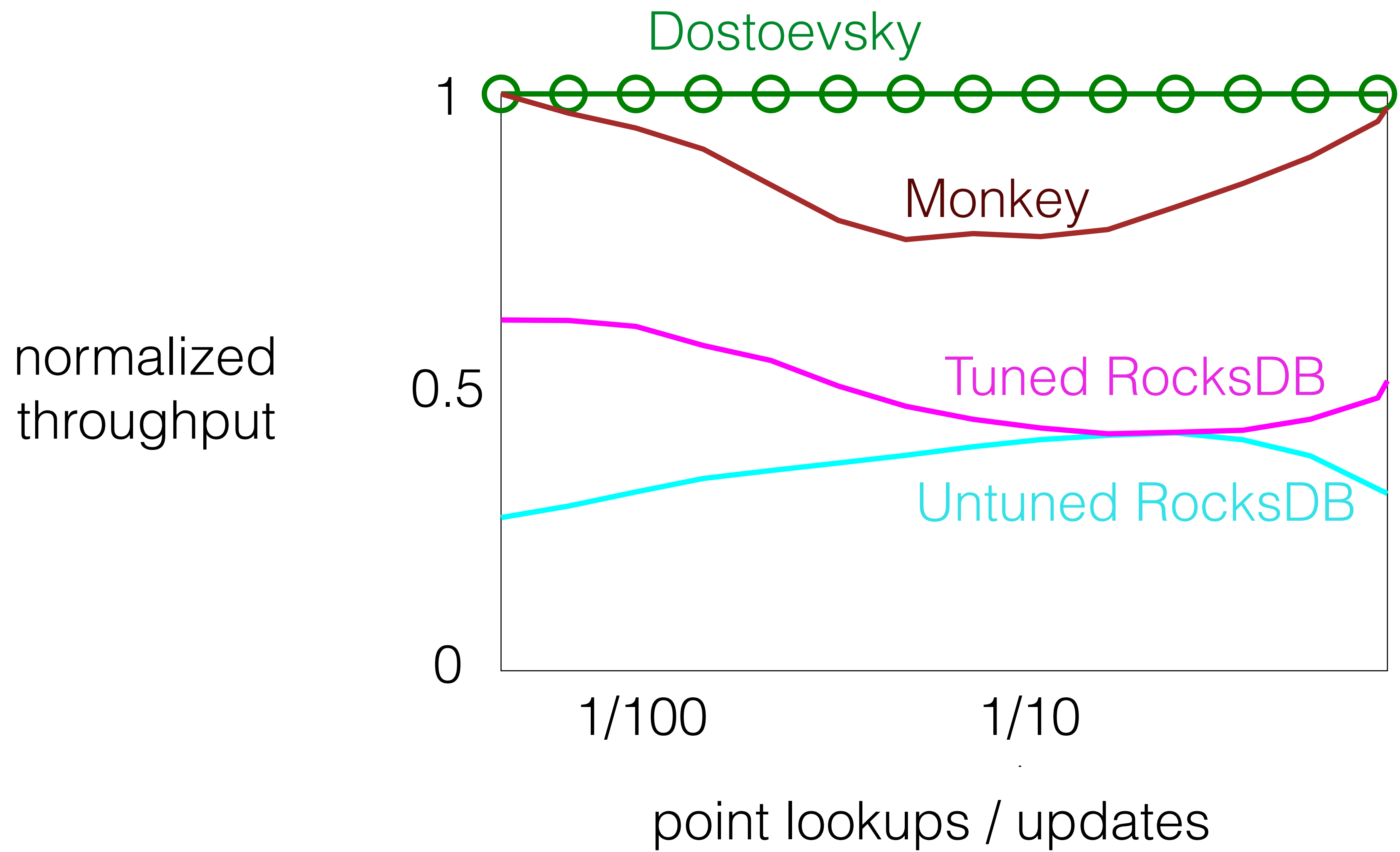


normalized
throughput





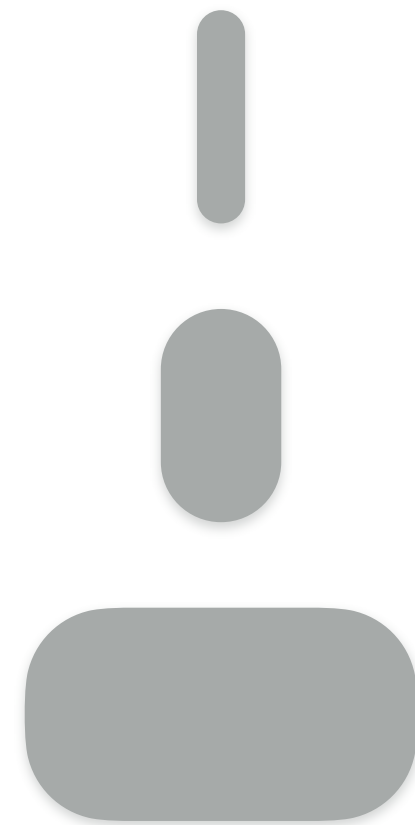




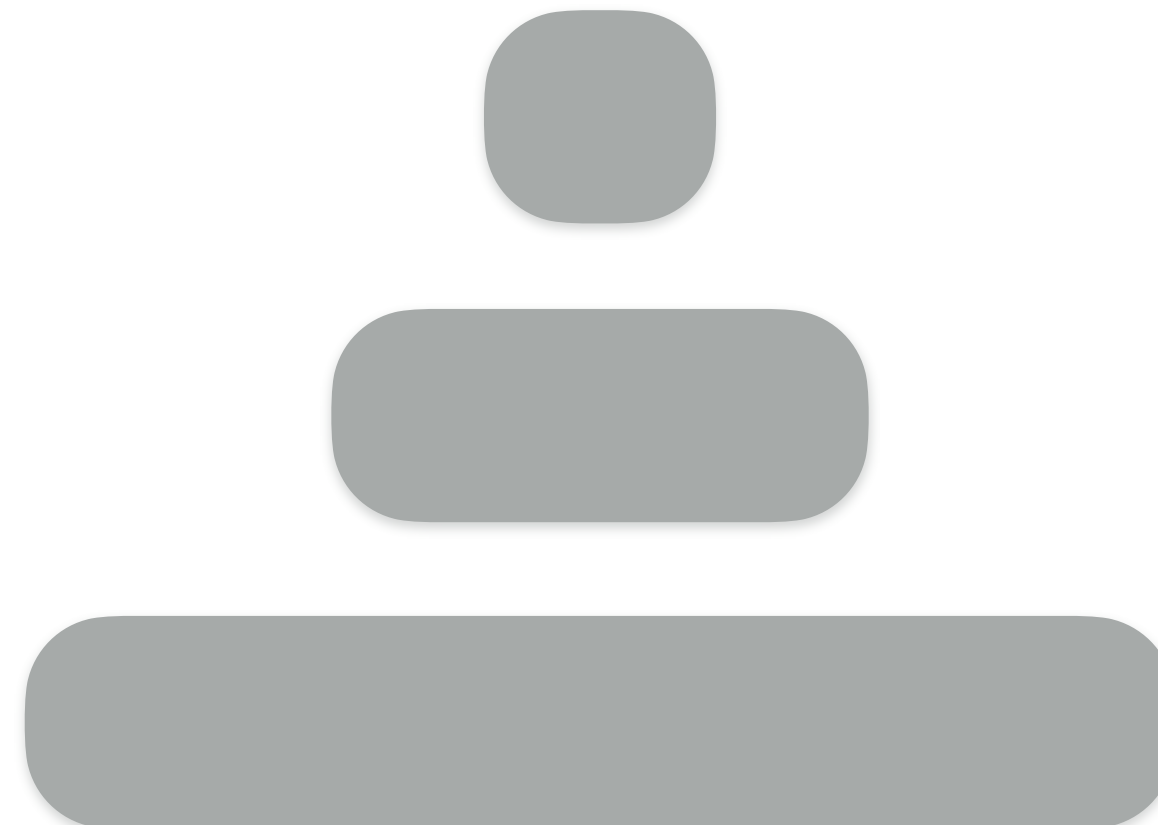
Conclusion

Conclusion

Bloom
filters



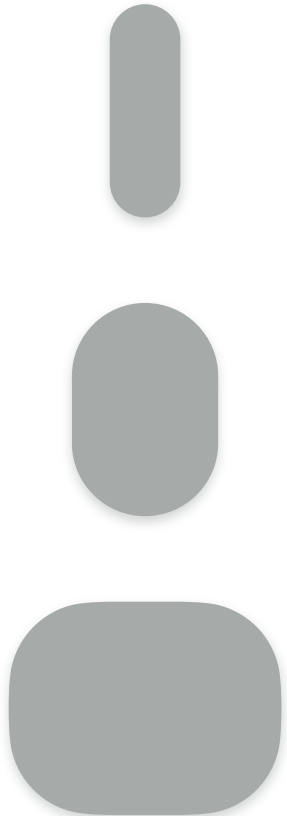
LSM-tree



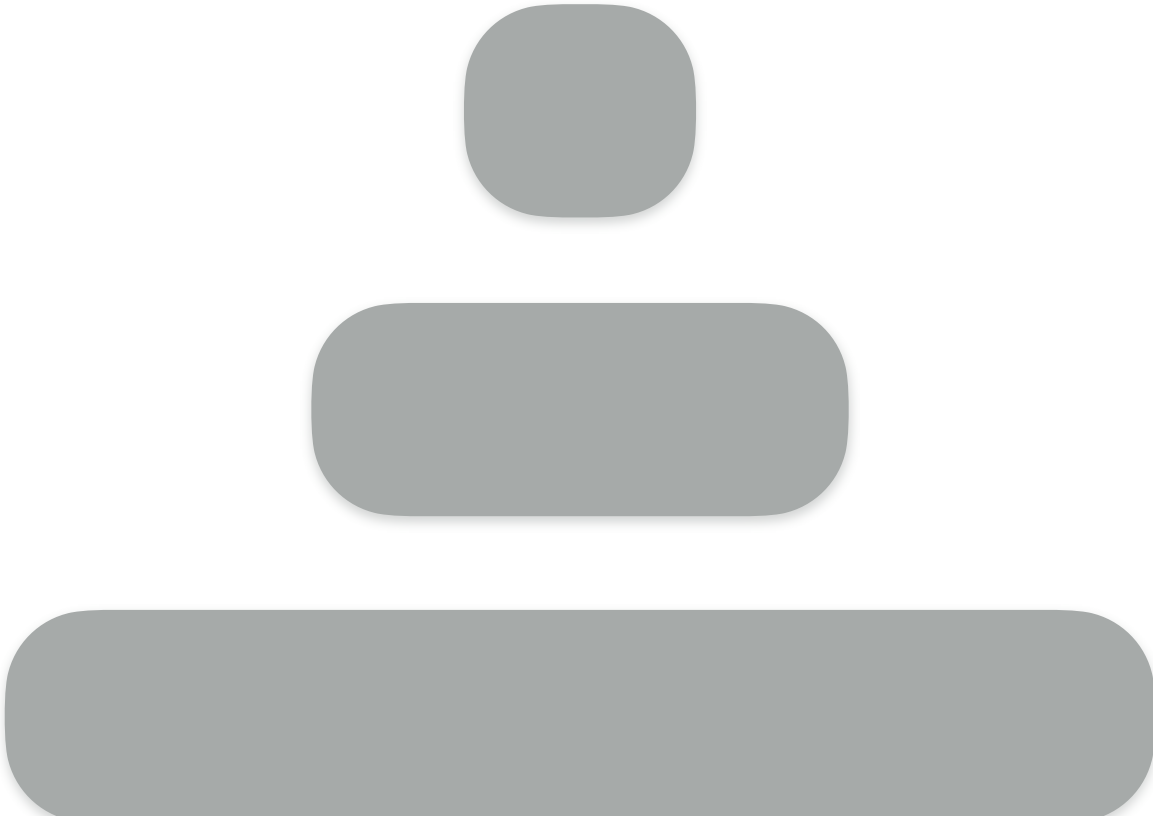
Conclusion

**optimizes
memory
allocation**

Bloom
filters



LSM-tree

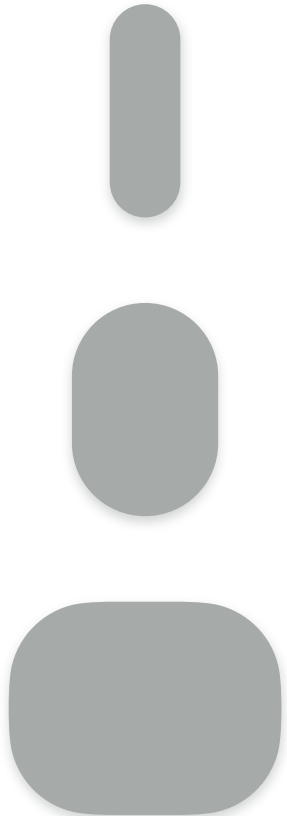


Conclusion



optimizes
memory
allocation

Bloom
filters



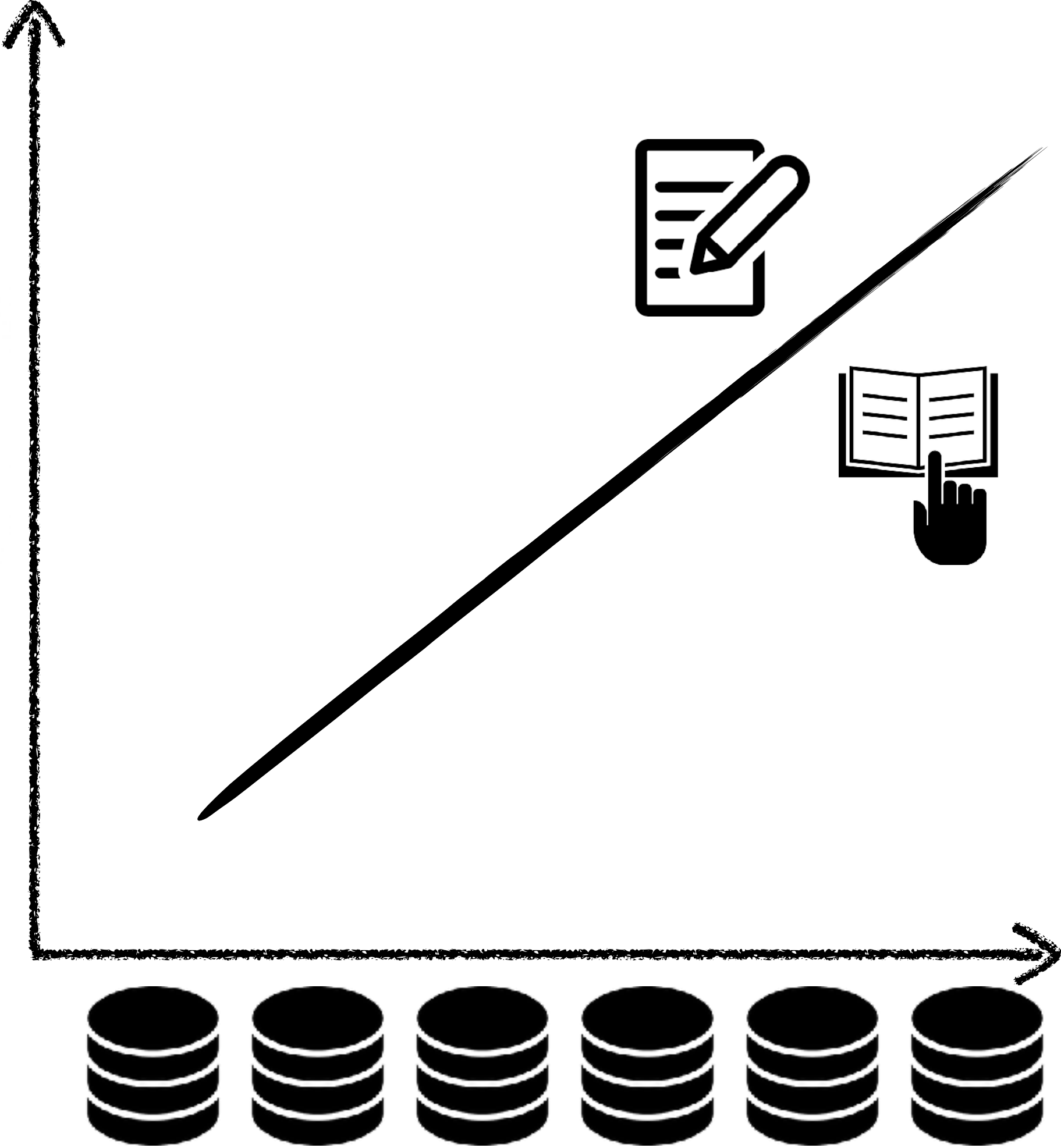
LSM-tree



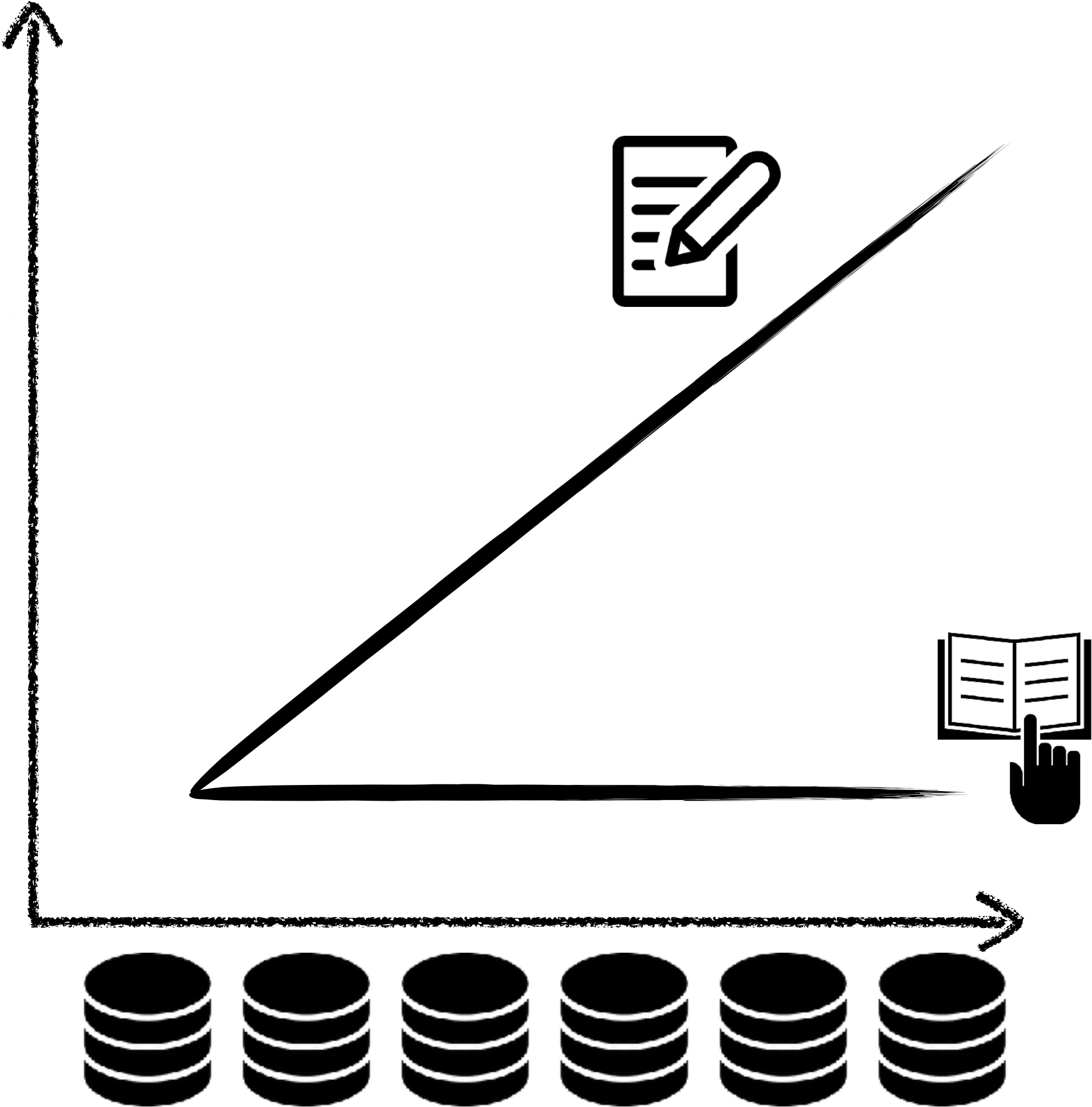
**removes
superfluous
merging**



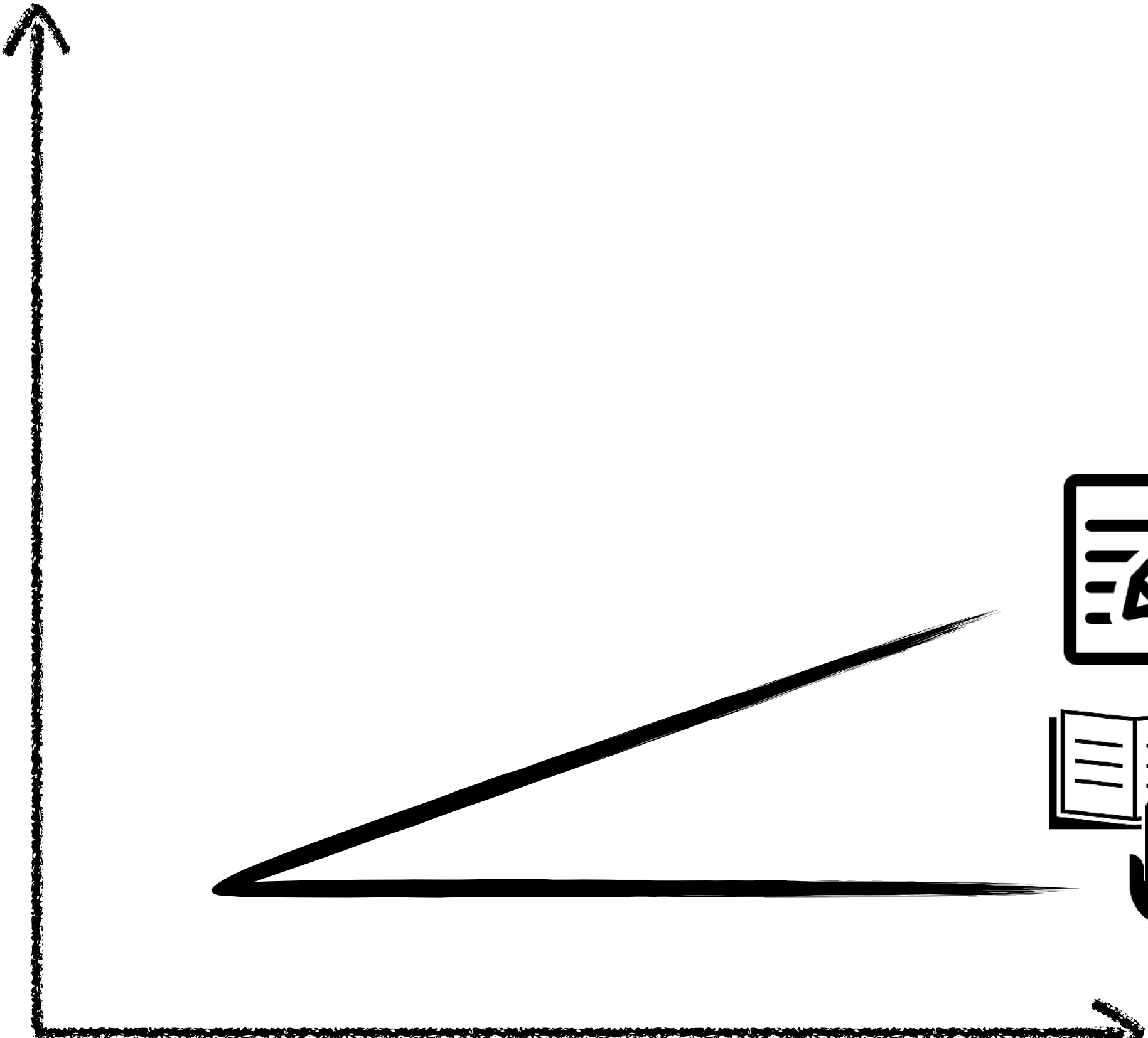
Conclusion



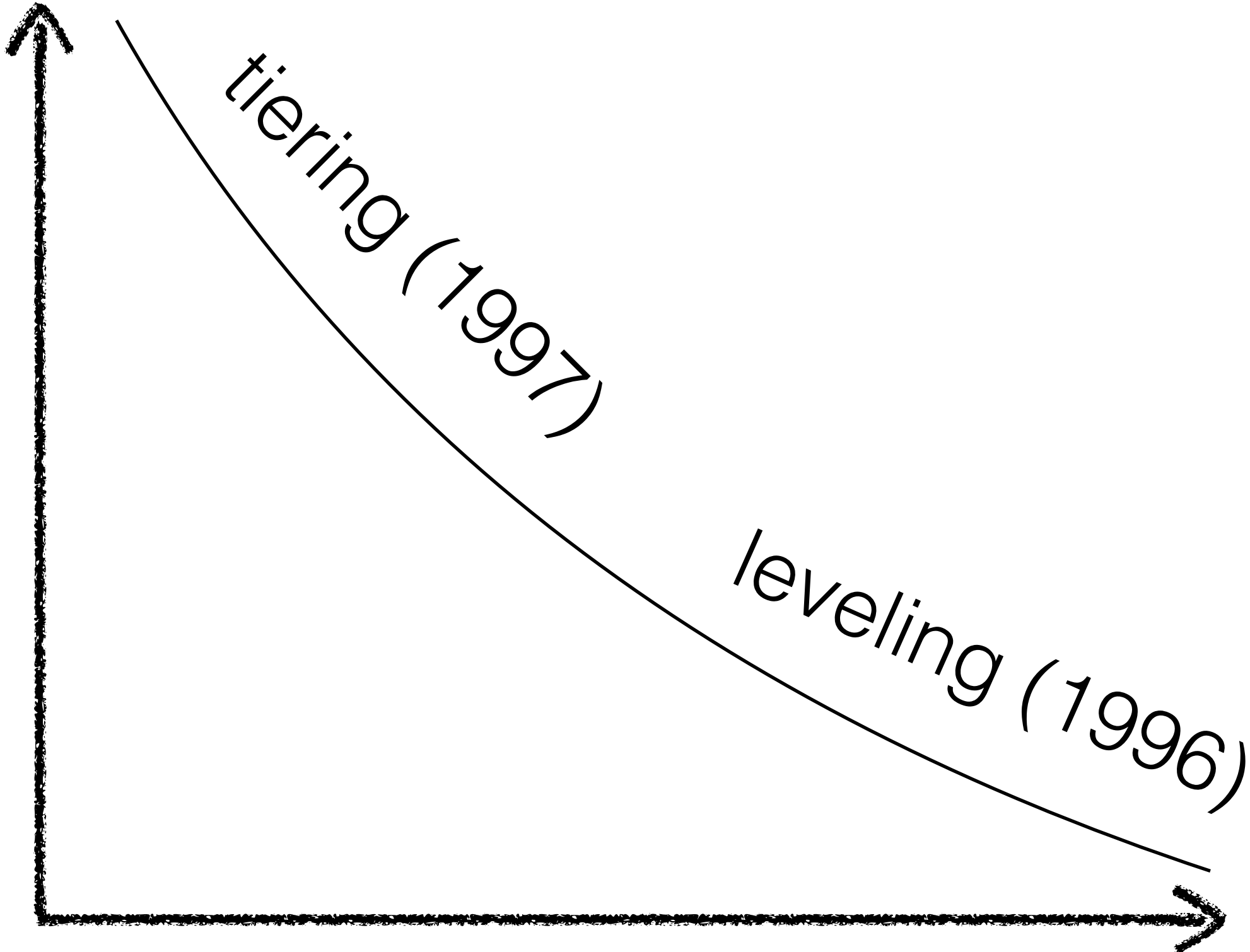
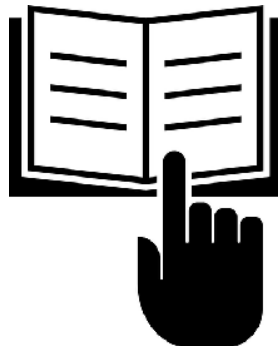
Conclusion



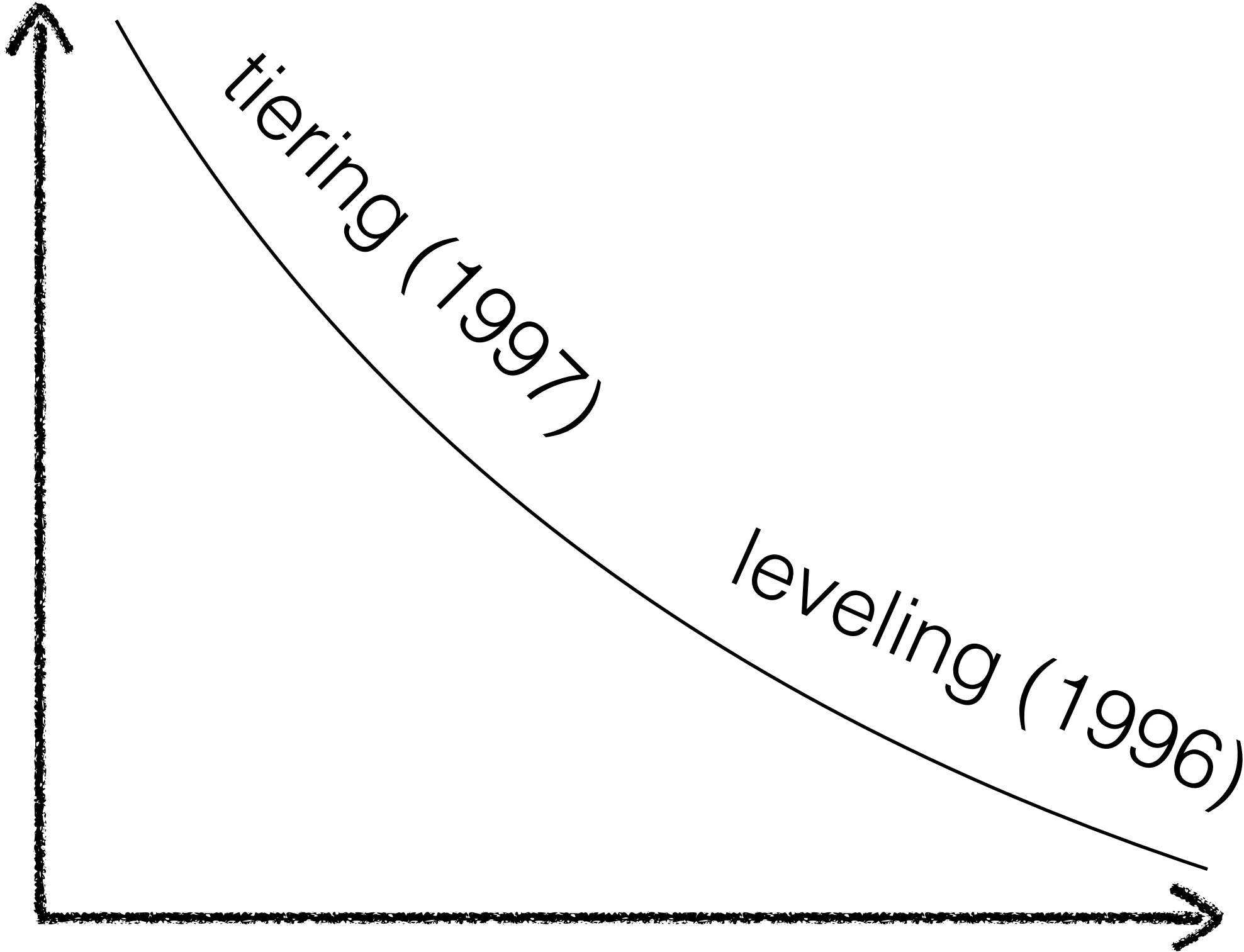
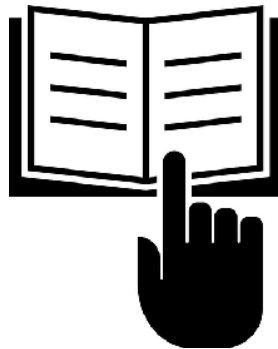
Conclusion



Conclusion



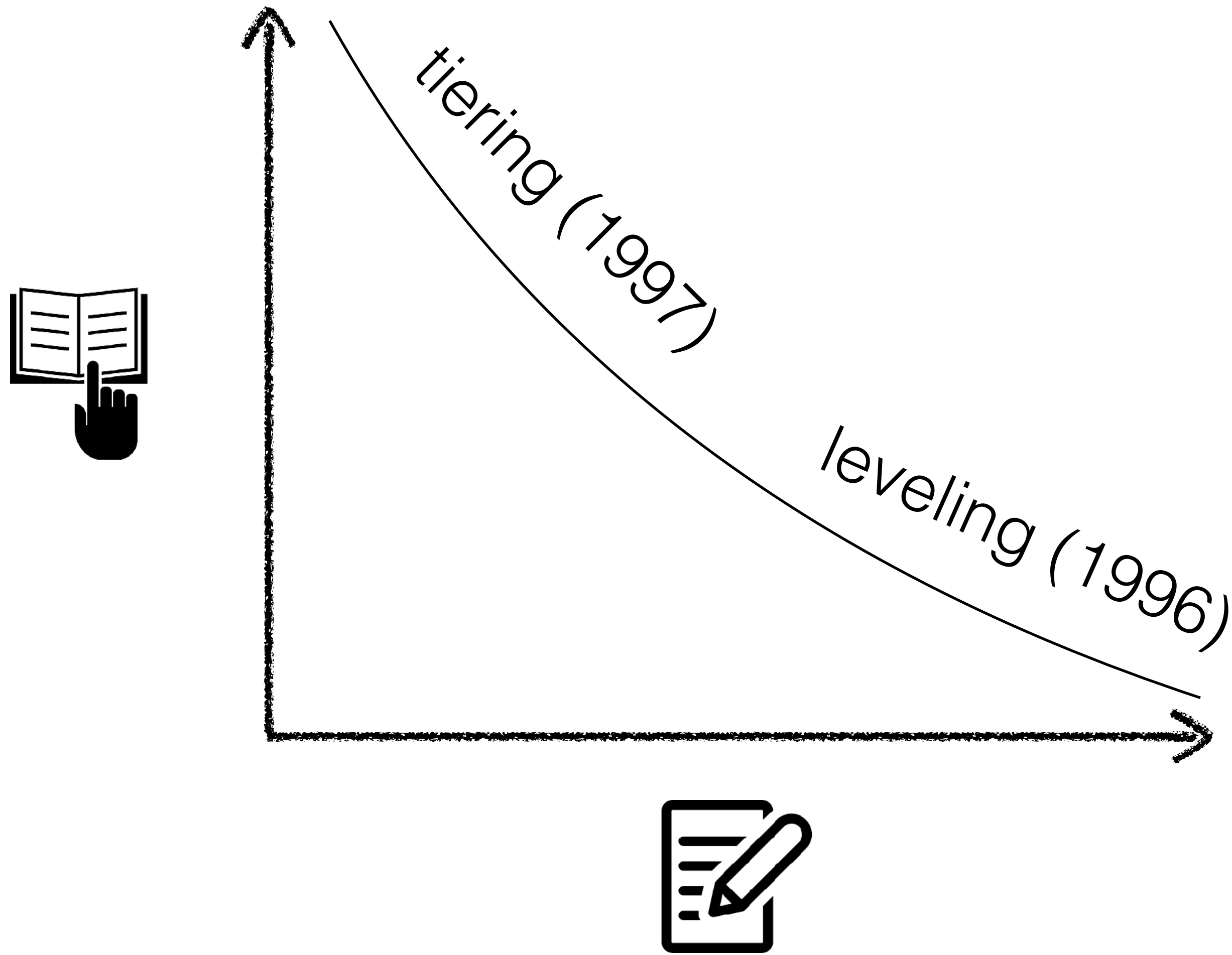
Conclusion



 little memory

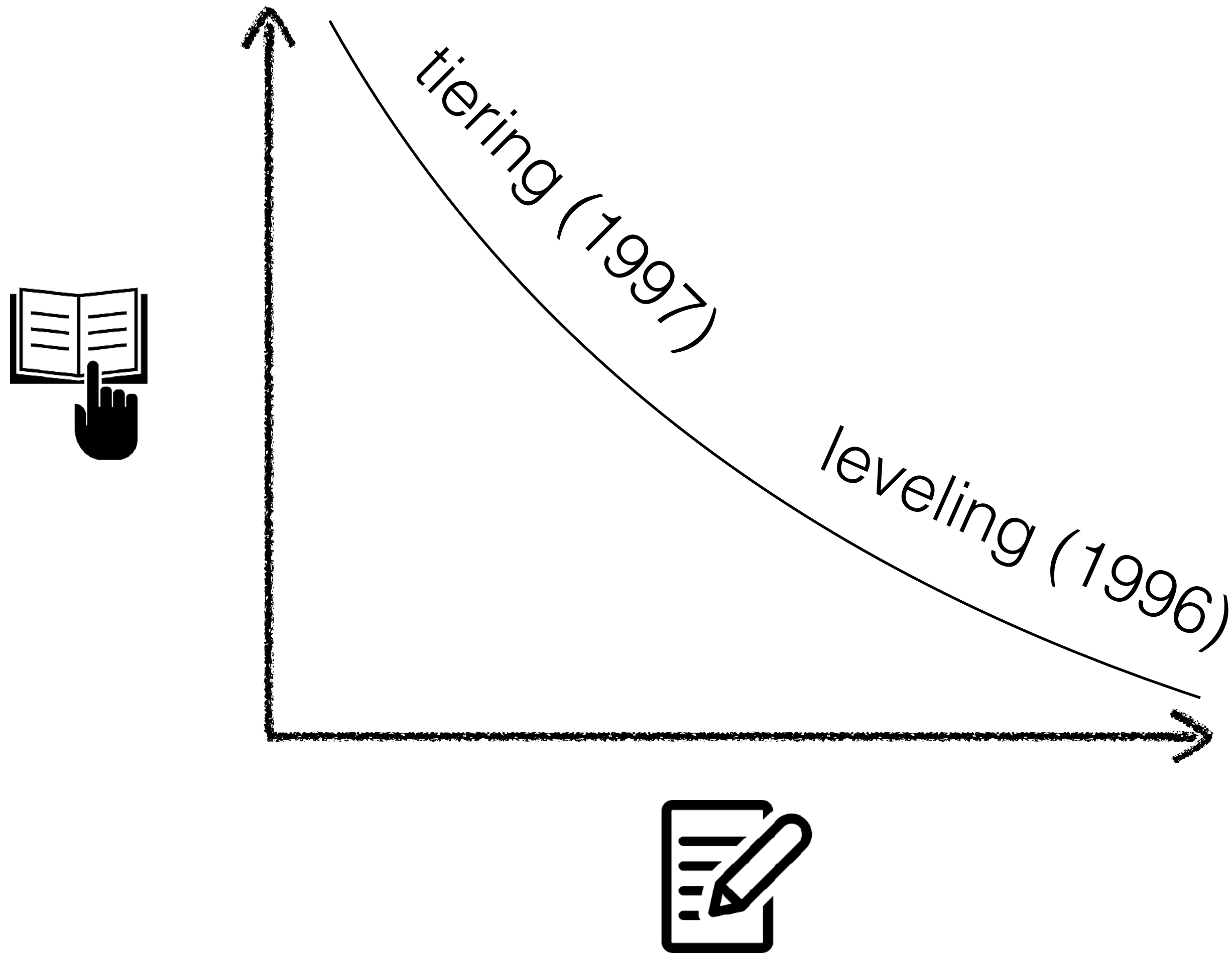


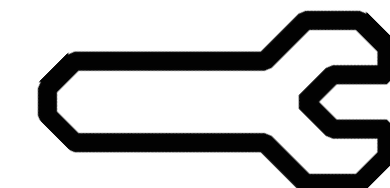
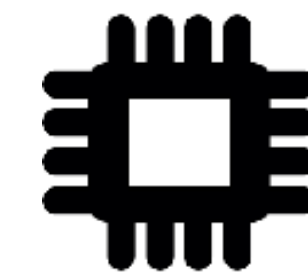
Conclusion



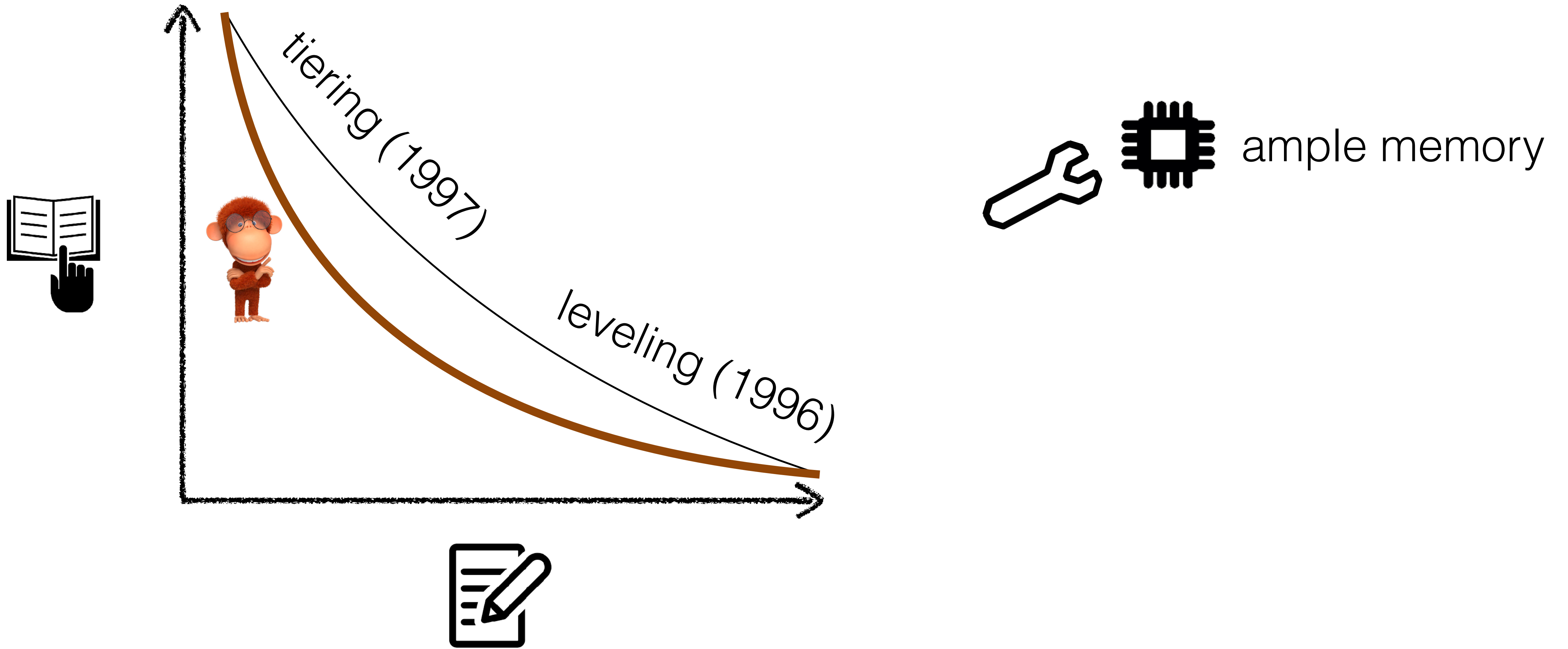
 ample memory

Conclusion

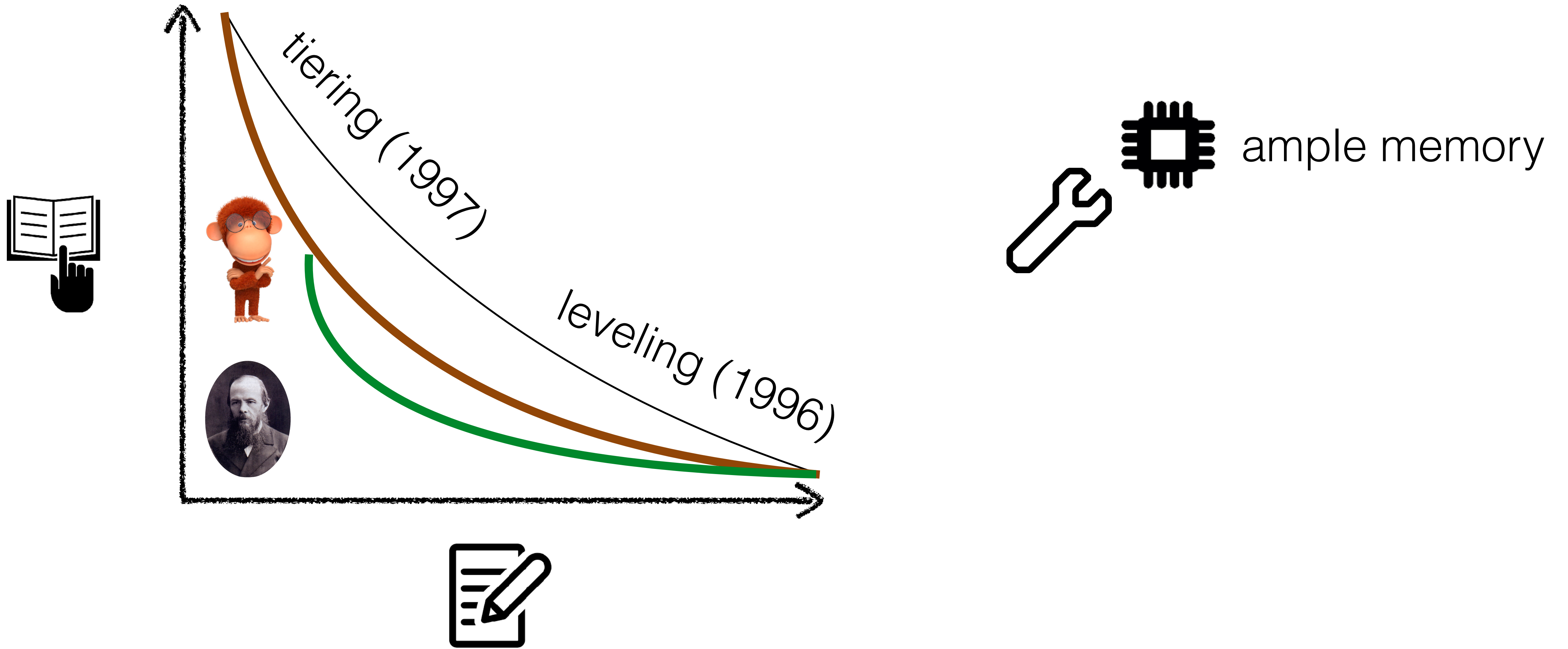


  ample memory

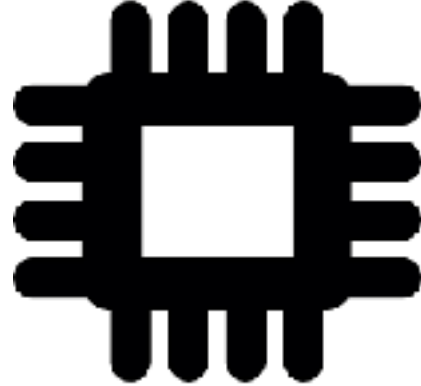
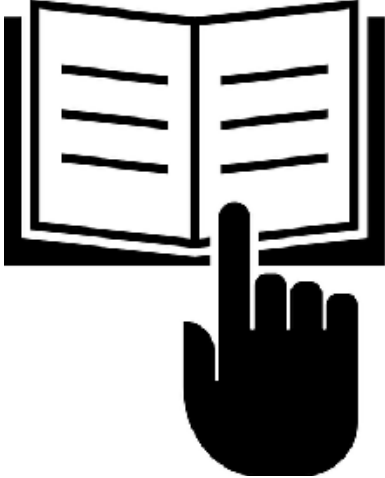
Conclusion



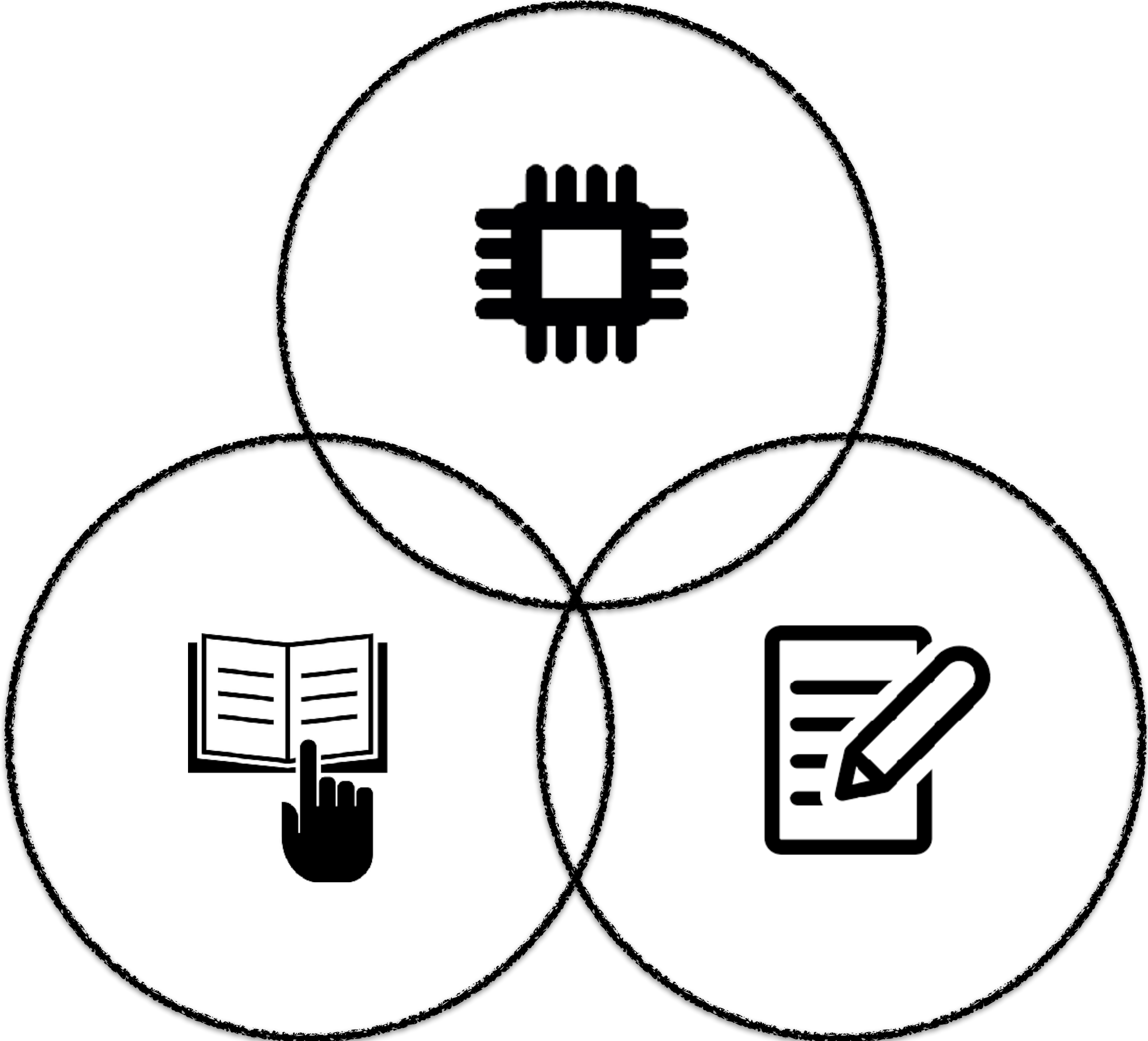
Conclusion



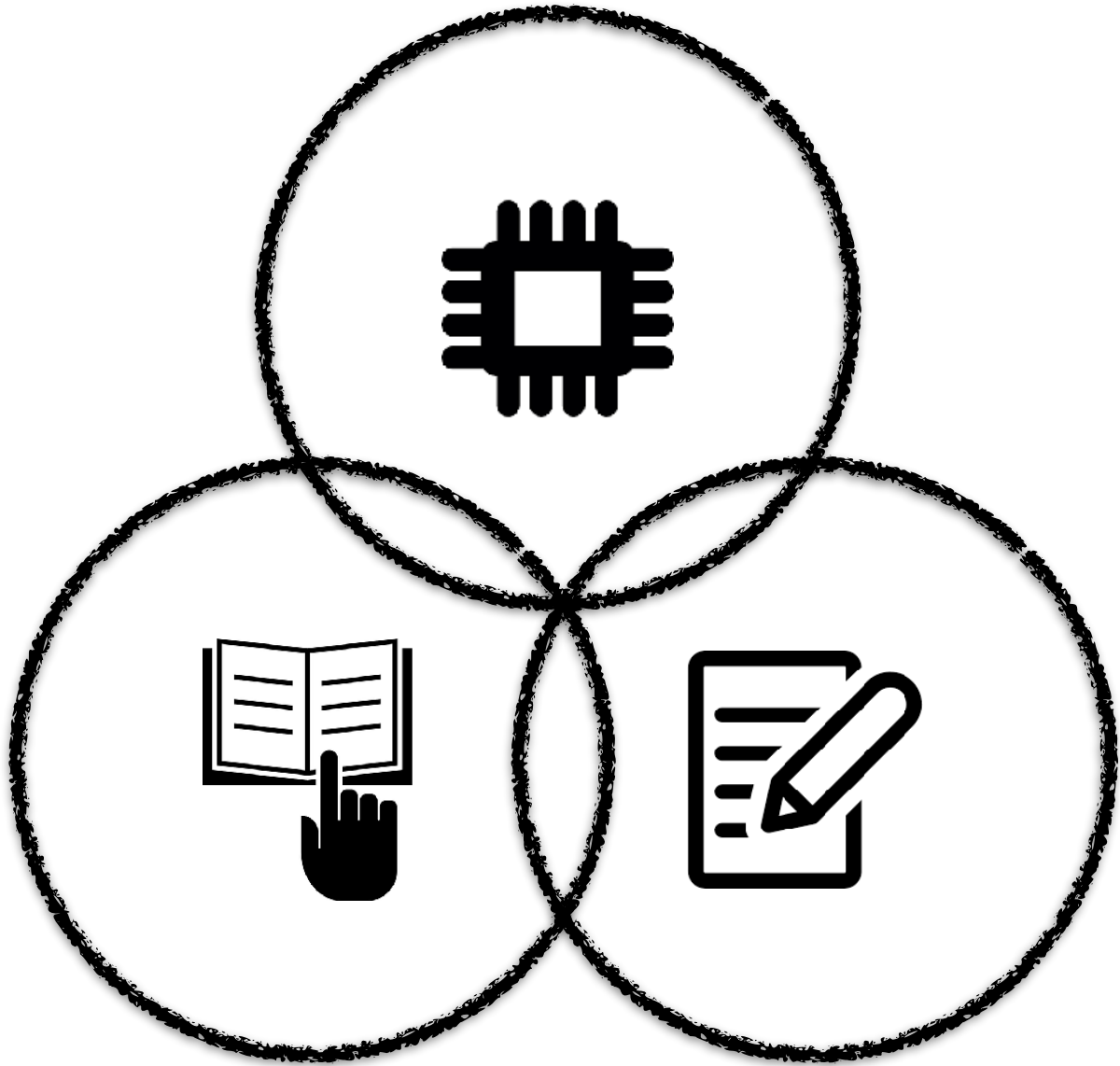
Conclusion



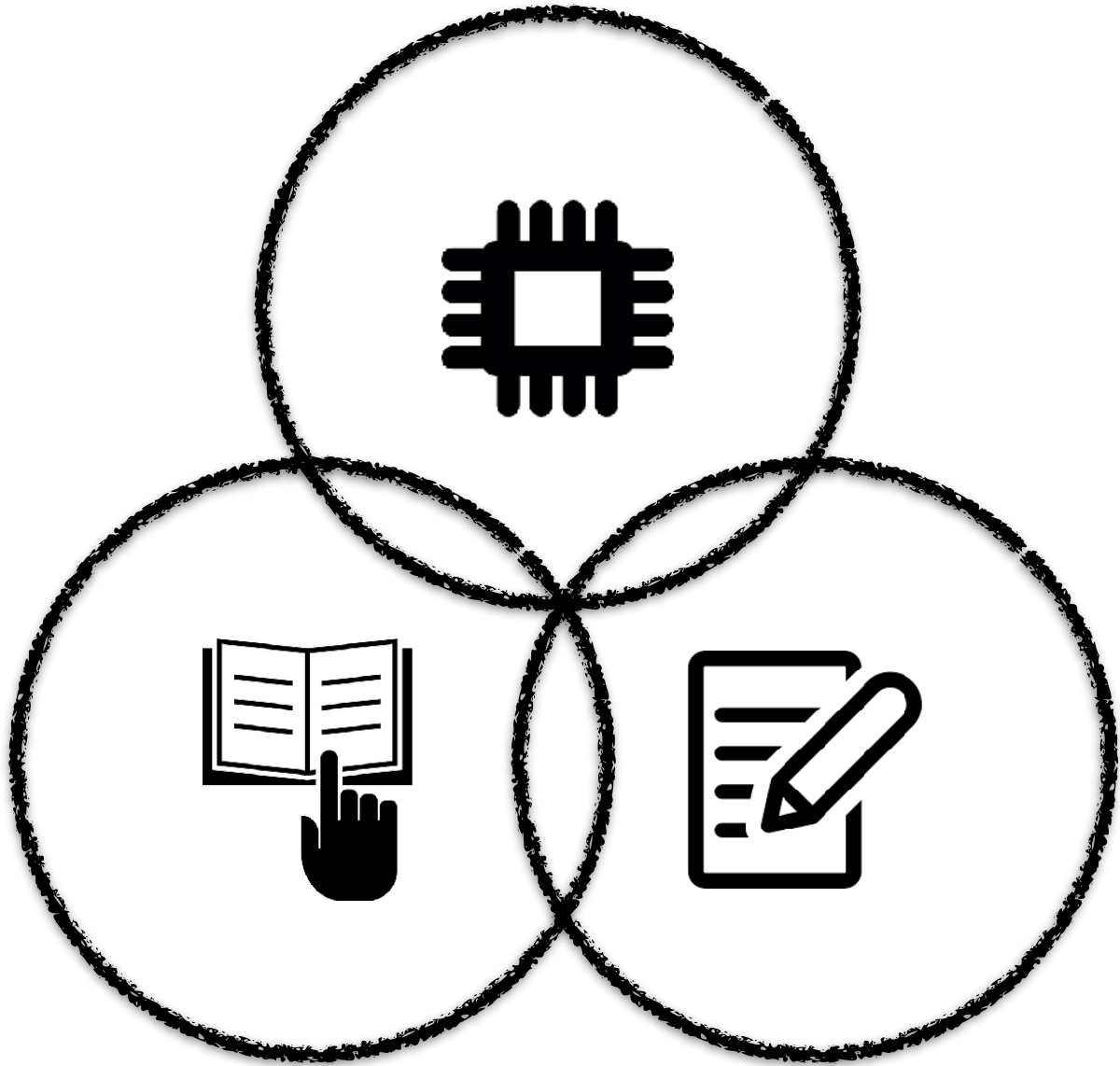
Conclusion



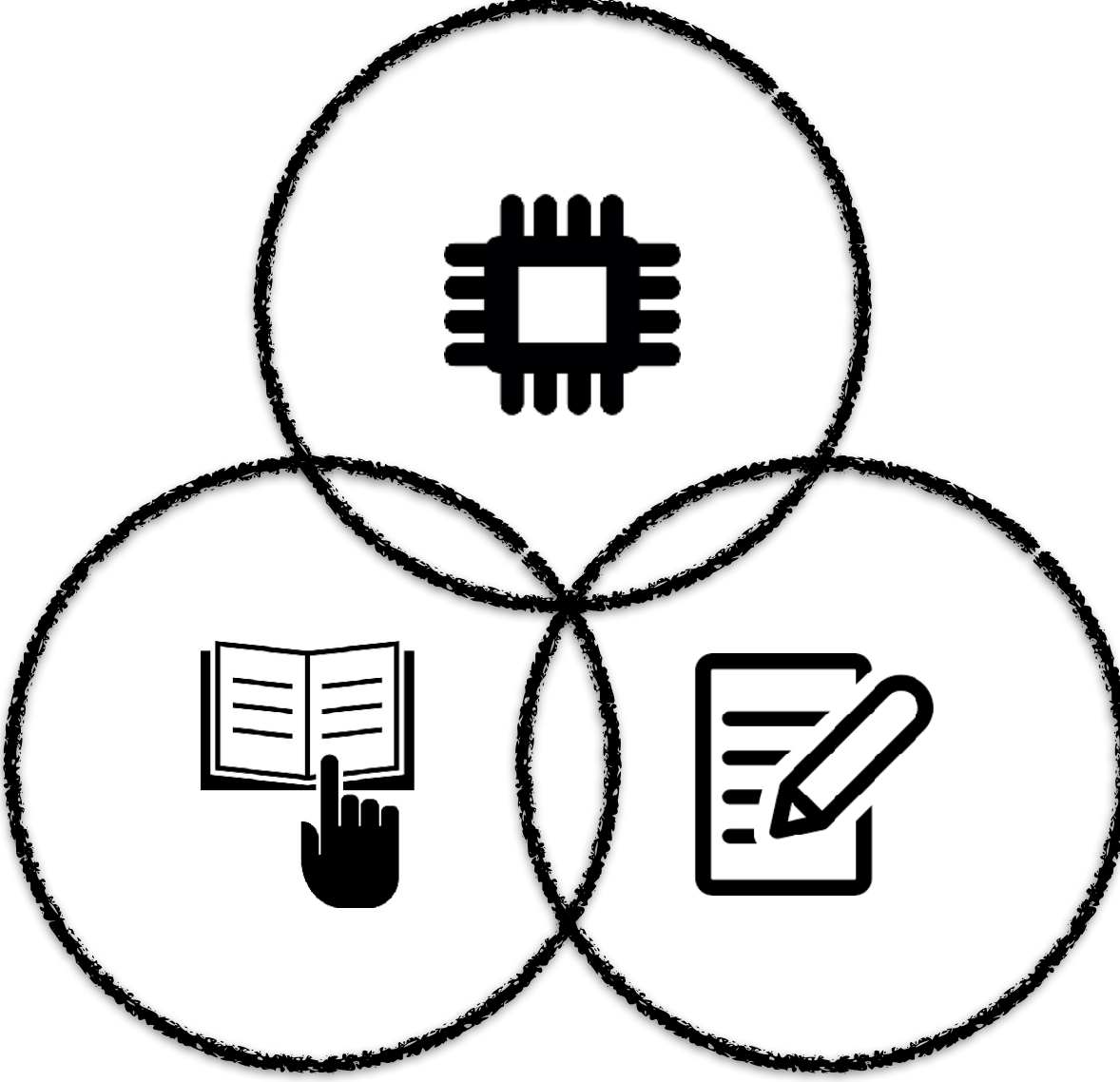
Conclusion

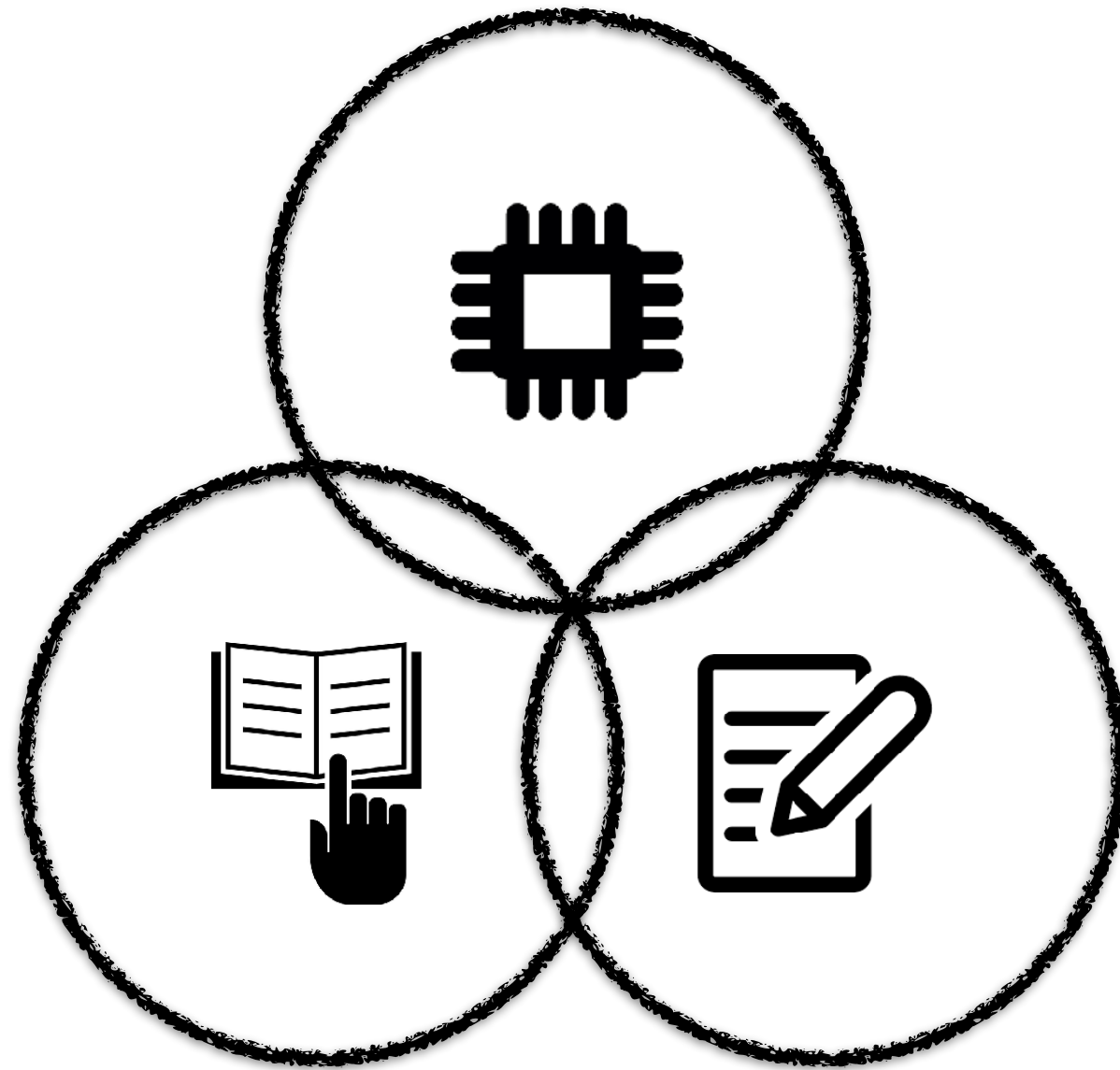


Conclusion



Crimson DB





Thanks!