



**Dostoevsky:** Better Space-Time Trade-Offs  
For LSM-Tree Based Key-Value Stores via  
**Adaptive Removal of Superfluous Merging**



**DASlab**  
@ Harvard SEAS

Niv Dayan &  
Stratos Idreos

# LSM-Tree



## Key-Value Stores



# LSM-Tree



## Key-Value Stores

## Time-Series Databases



# LSM-Tree

## Relational Databases

SQLite4



## Key-Value Stores

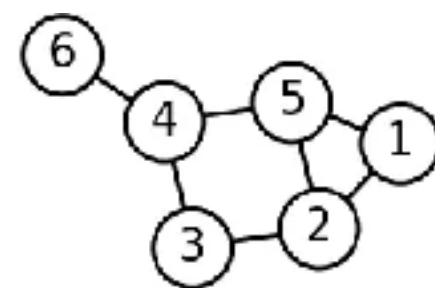


BigTable



## Time-Series Databases





## Algorithm Design



## Key-Value Stores

# LSM-Tree

## Relational Databases

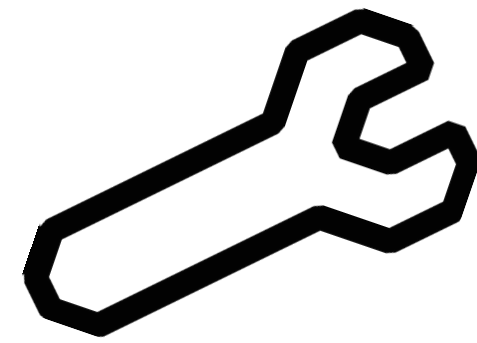
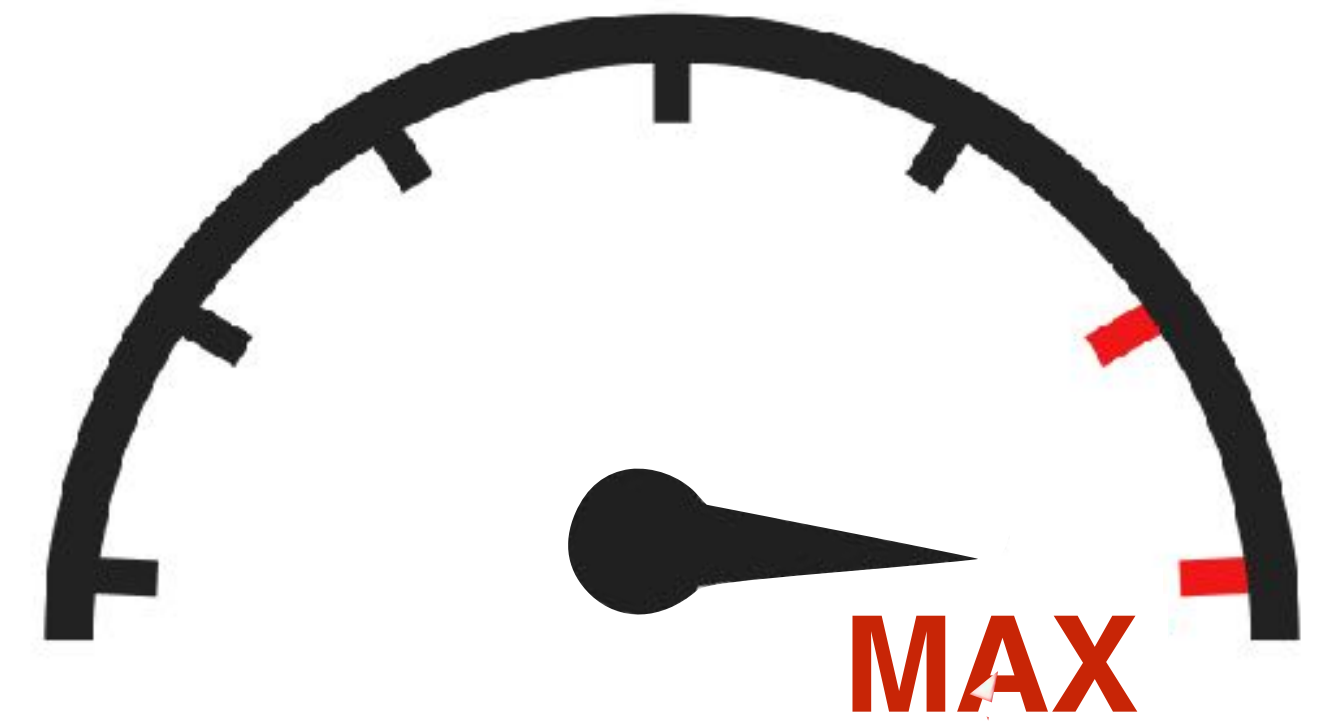
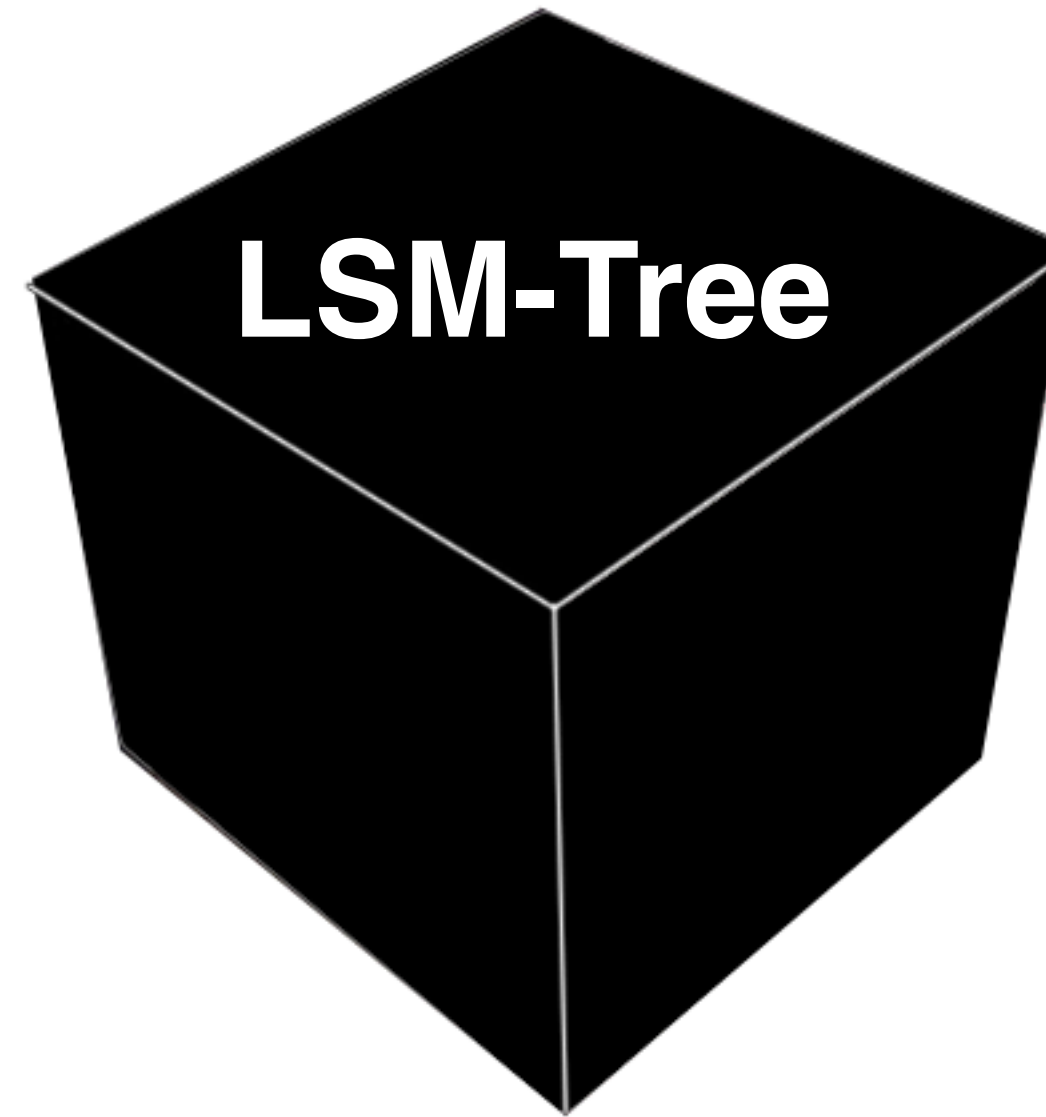
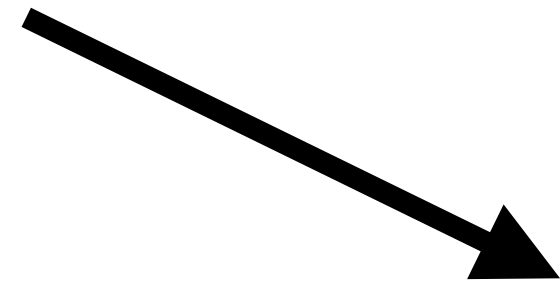
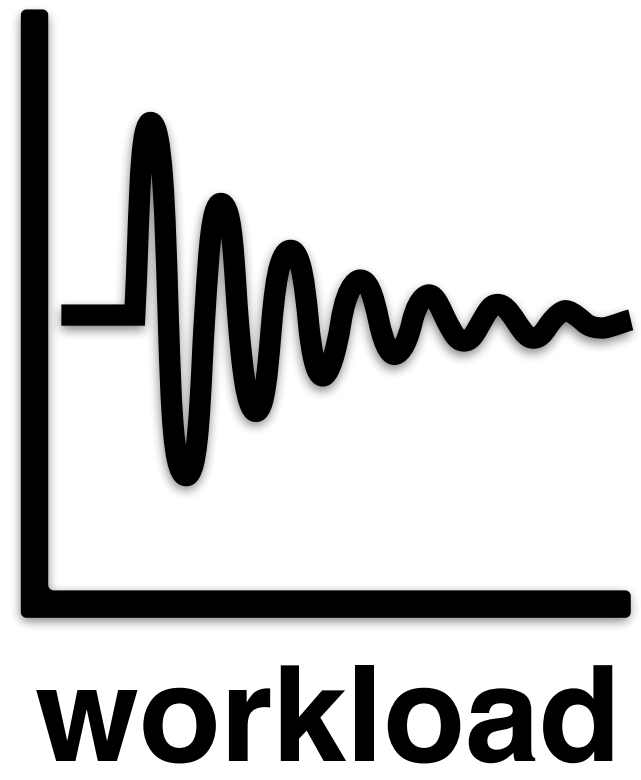
SQLite4



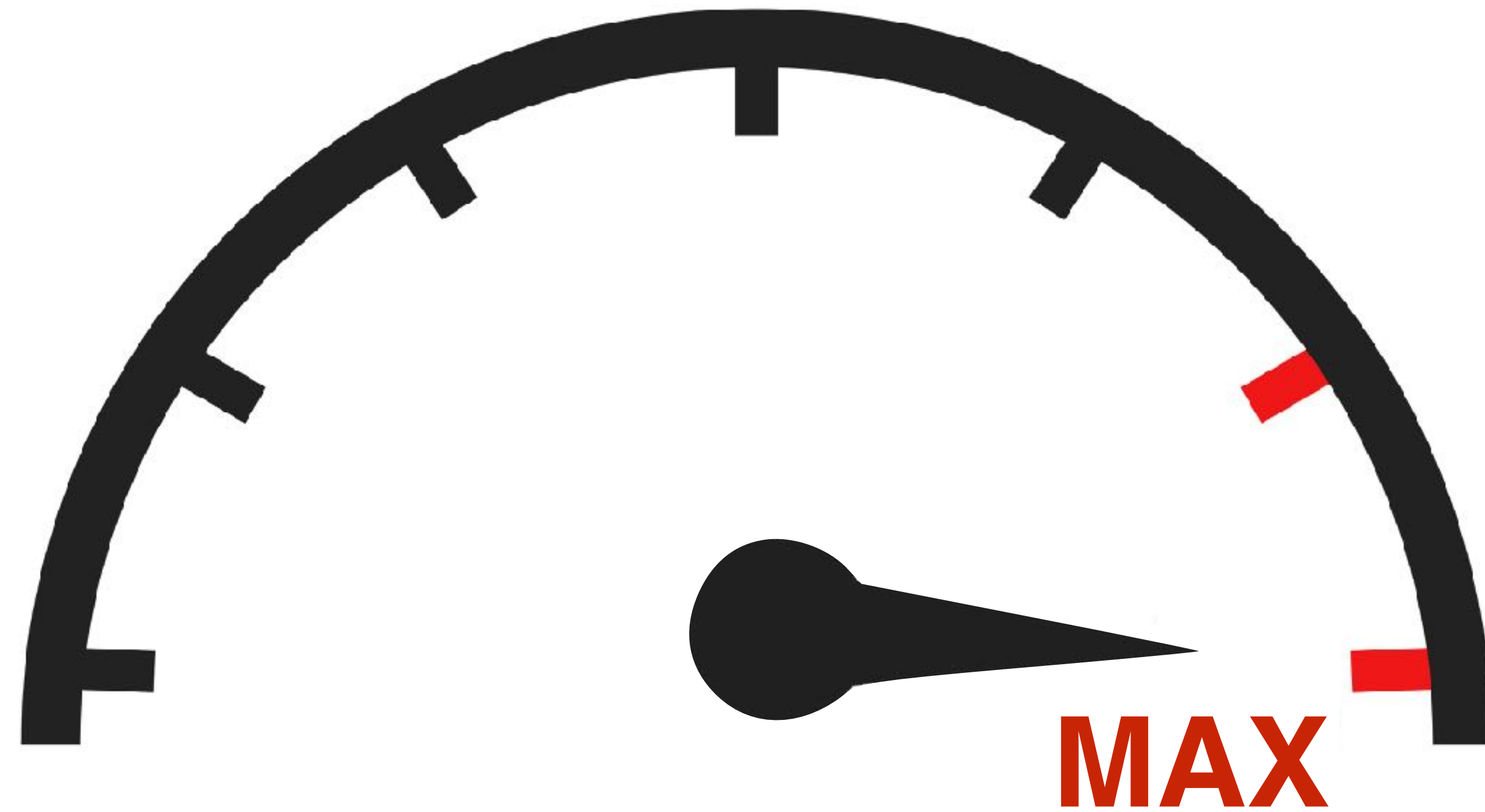
## Time-Series Databases



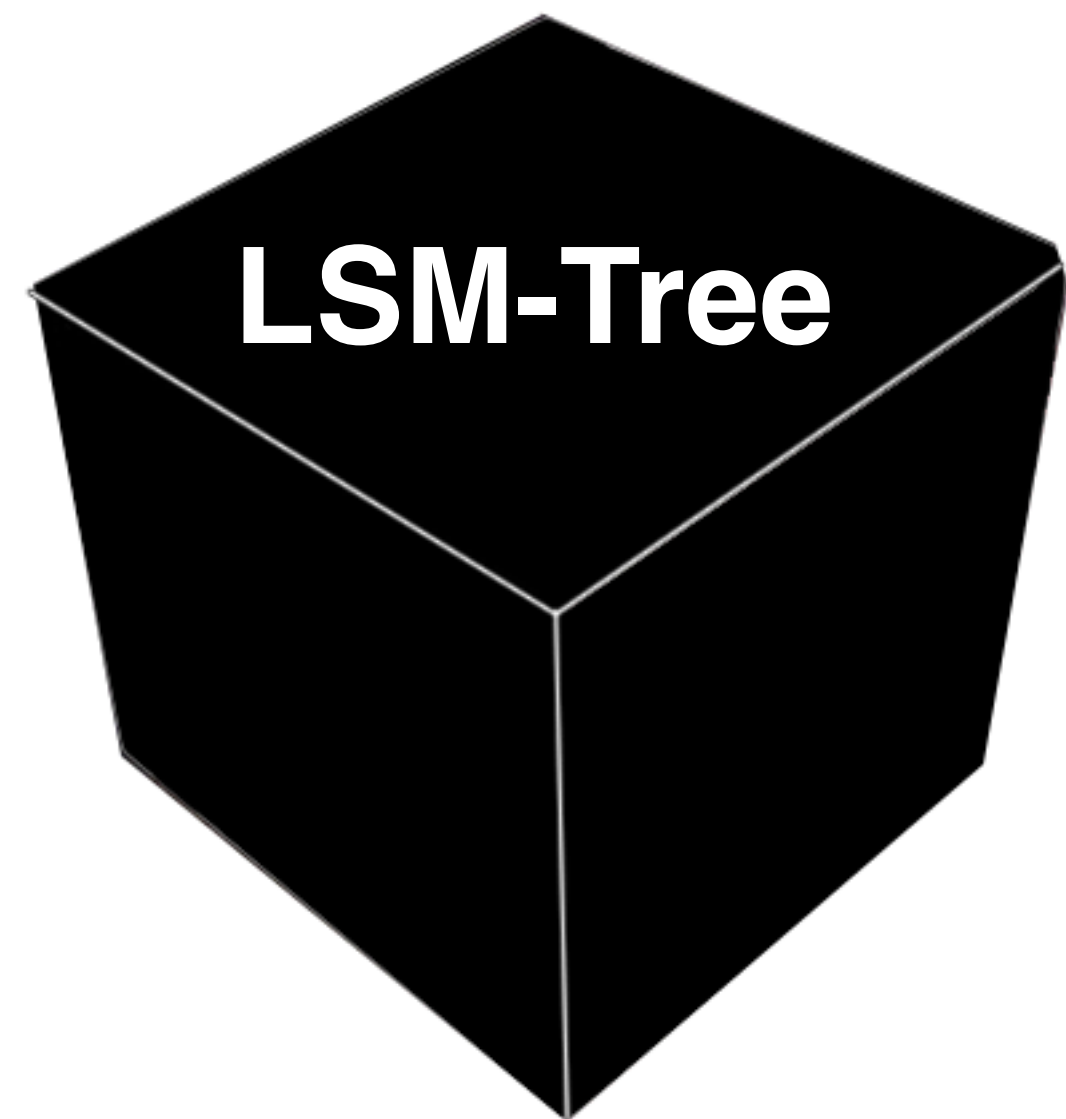
**LSM-Tree**



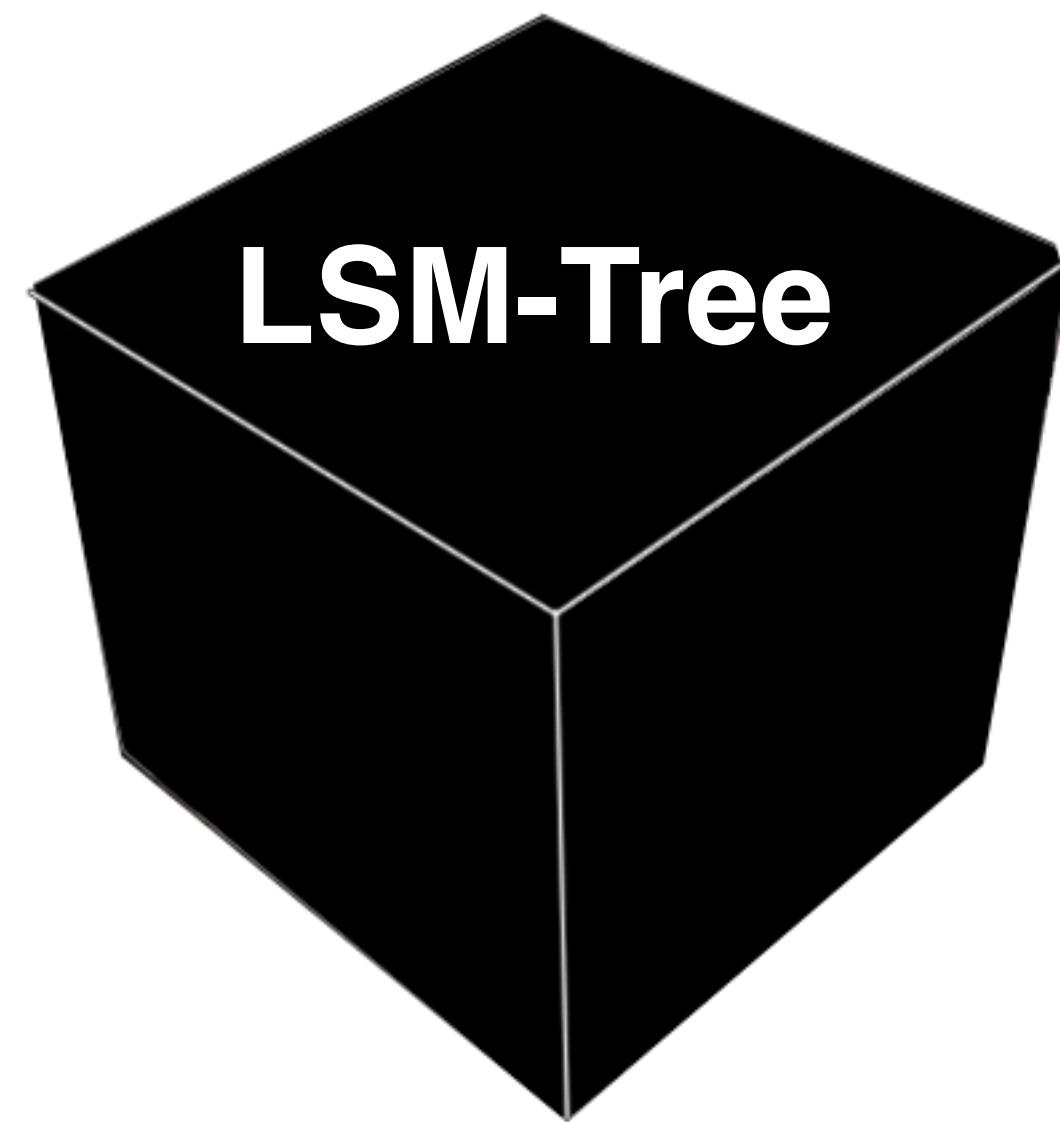








**suboptimal design**



**suboptimal design**

**why? how?**

# LSM-Tree Background

**updates**



**buffer**

memory storage

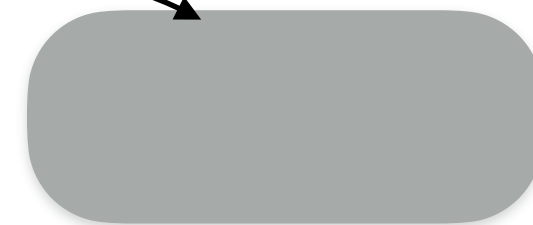
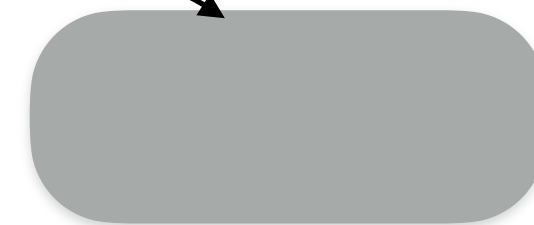


updates



memory storage

**sort & flush**



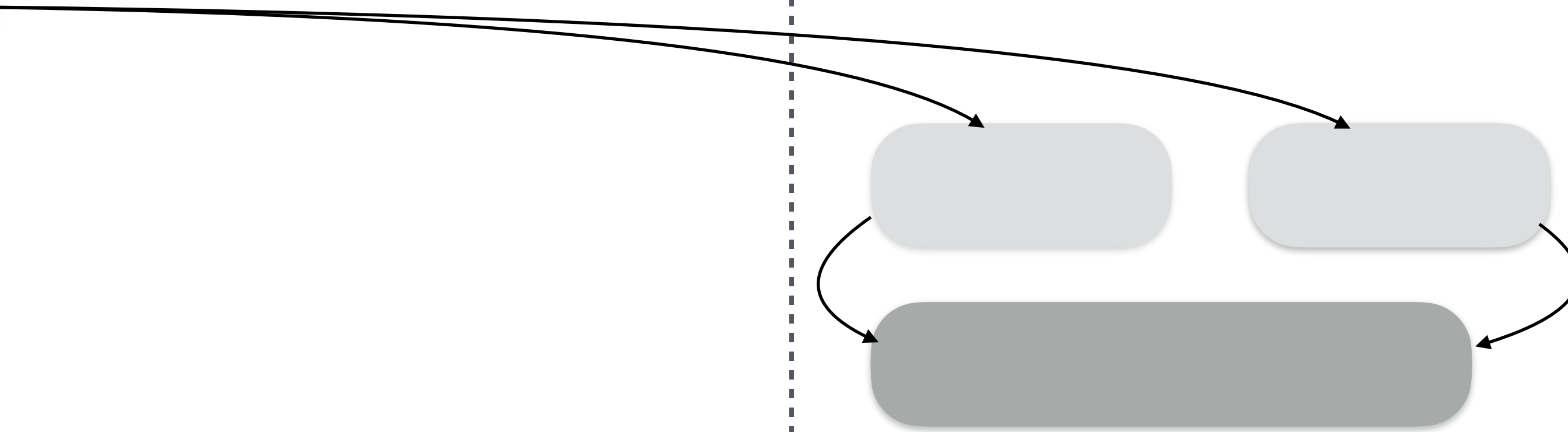
**runs**

updates



memory

storage



**sort merge**



buffer

memory storage

**exponentially increasing capacities**

**level 1**

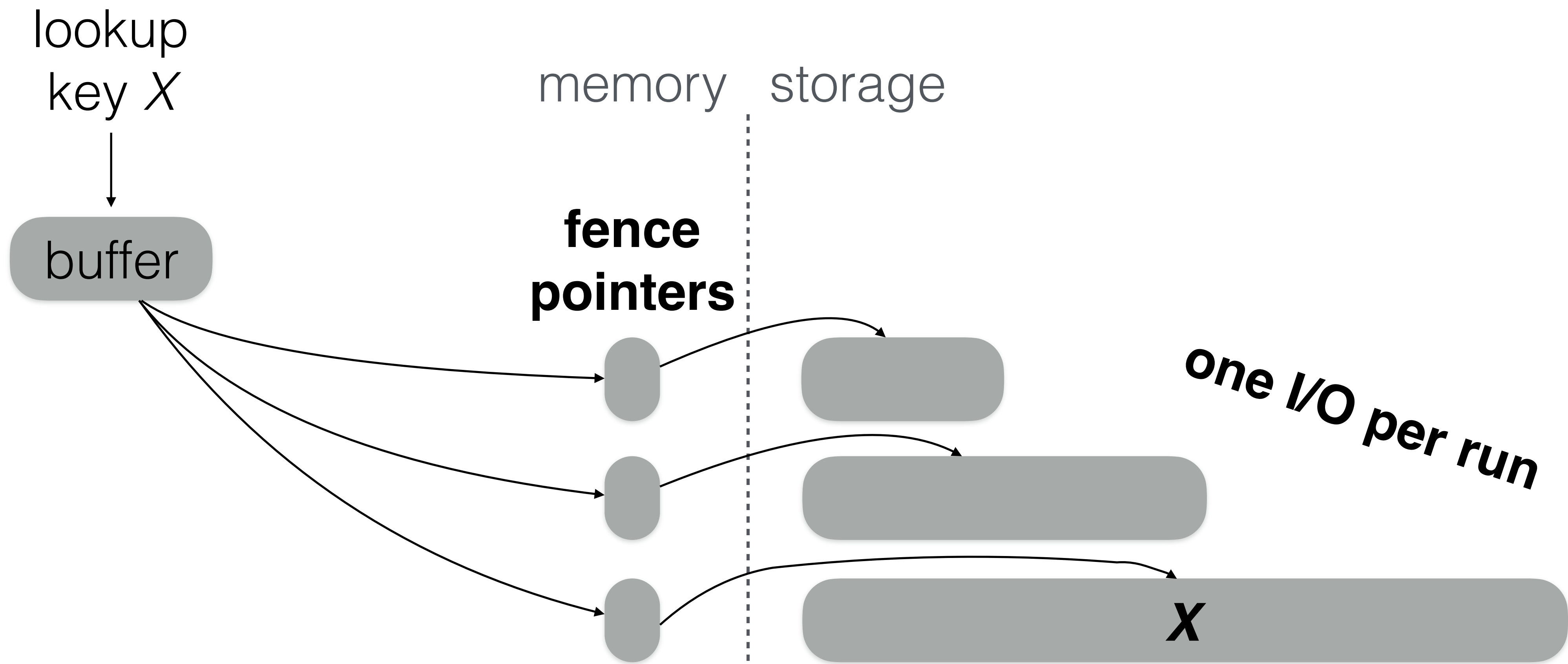


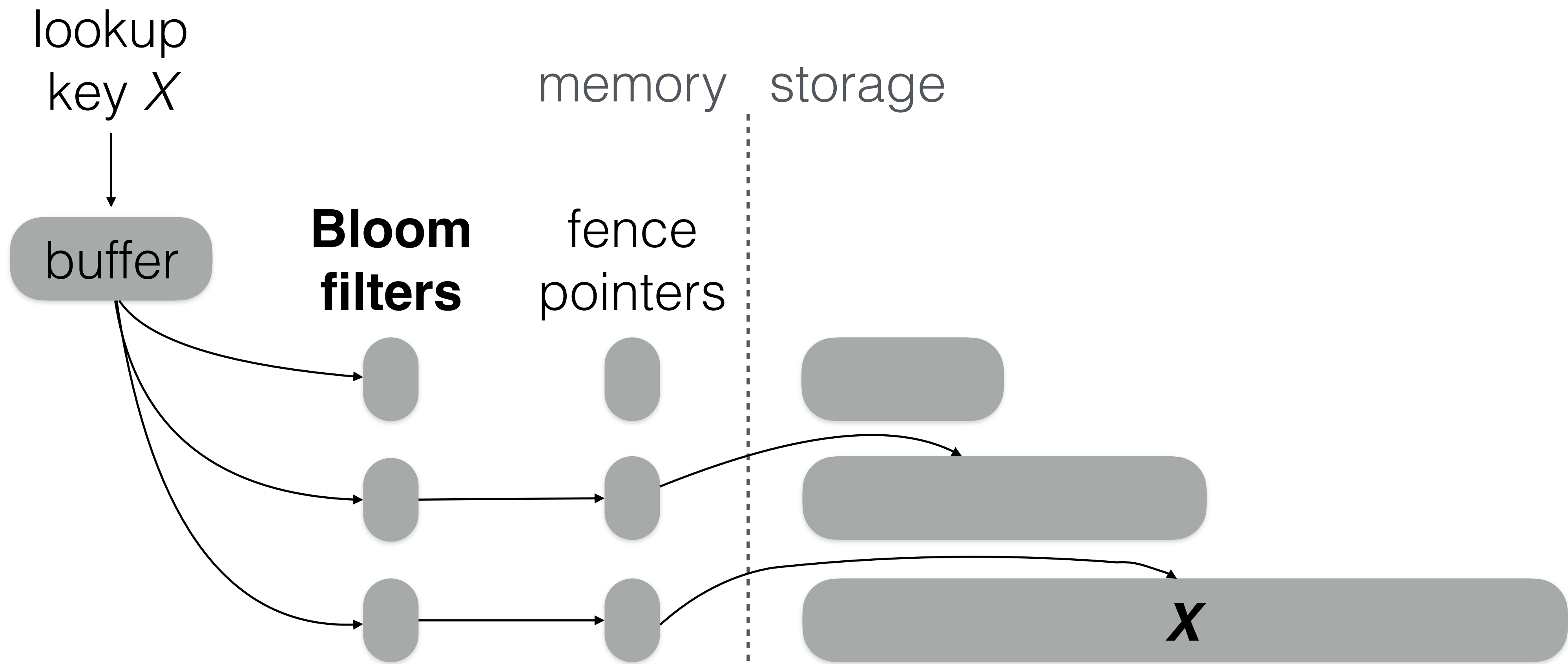
**level 2**



**level 3**

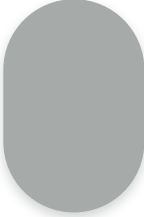
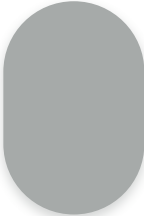






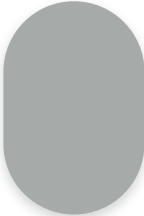
buffer

Bloom  
filters

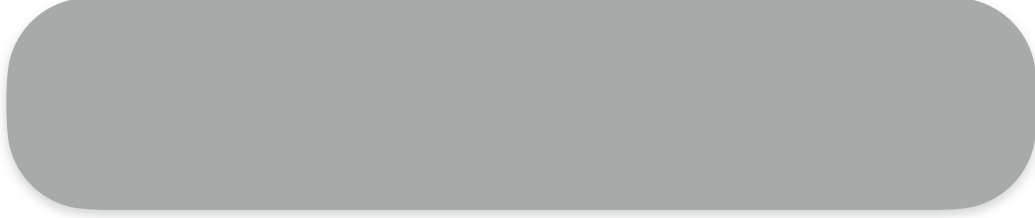


memory

fence  
pointers



storage



*merging frequency*



**merging frequency**



**merge cost**

**lookup cost**

**merge cost**



**merging frequency**



**lookup cost**

merging frequency

## Tiering

merge-optimized

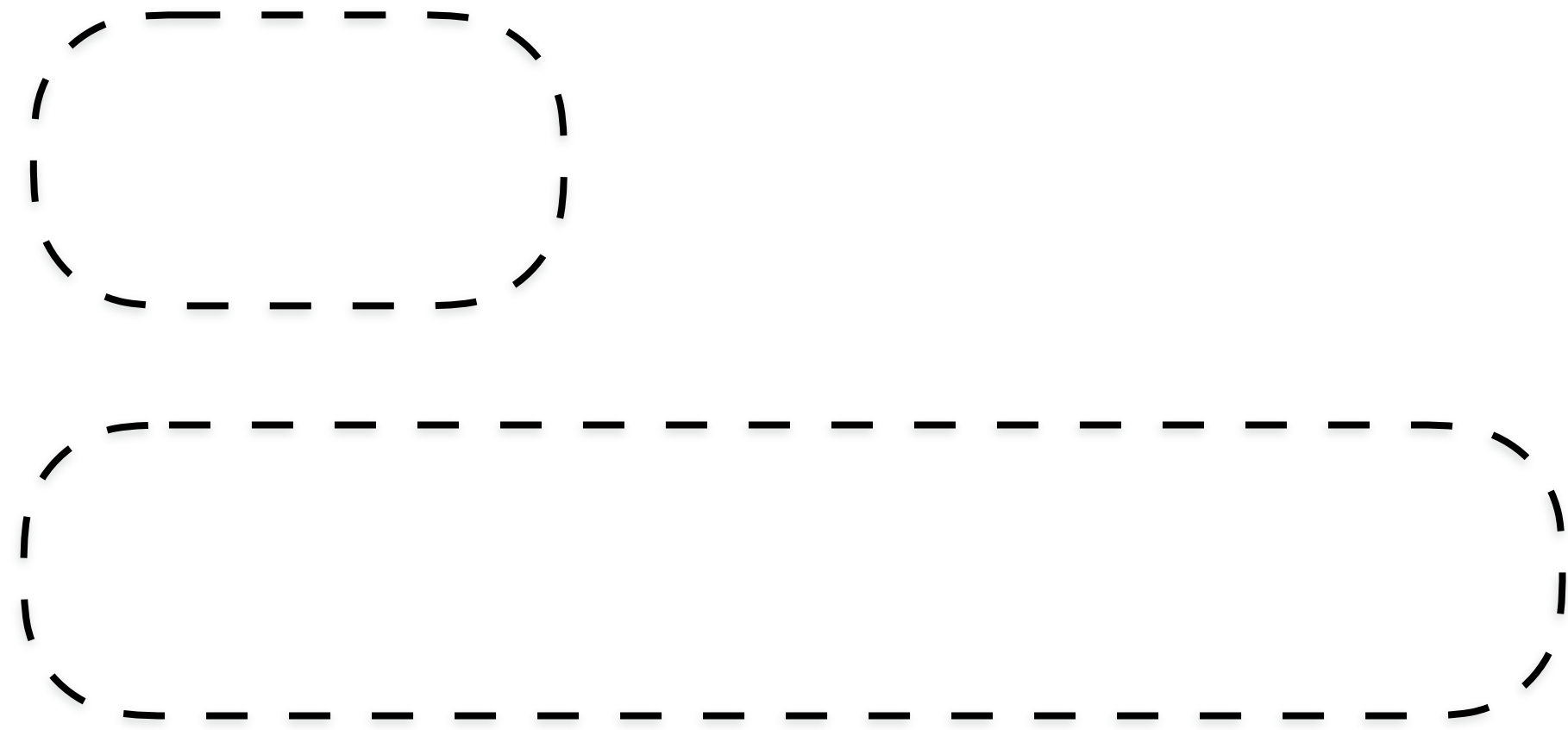


## Leveling

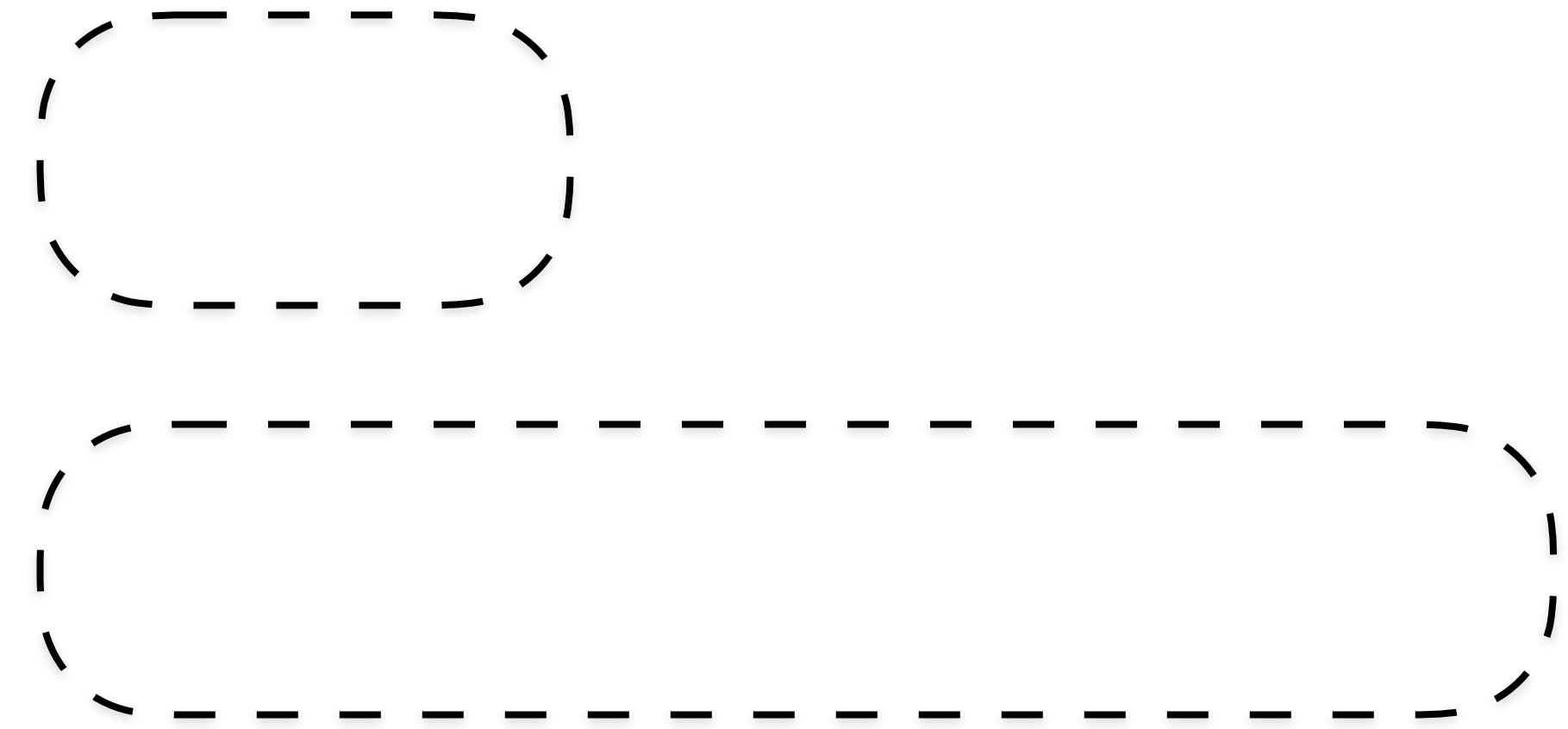
lookup-optimized



Tiering  
merge-optimized



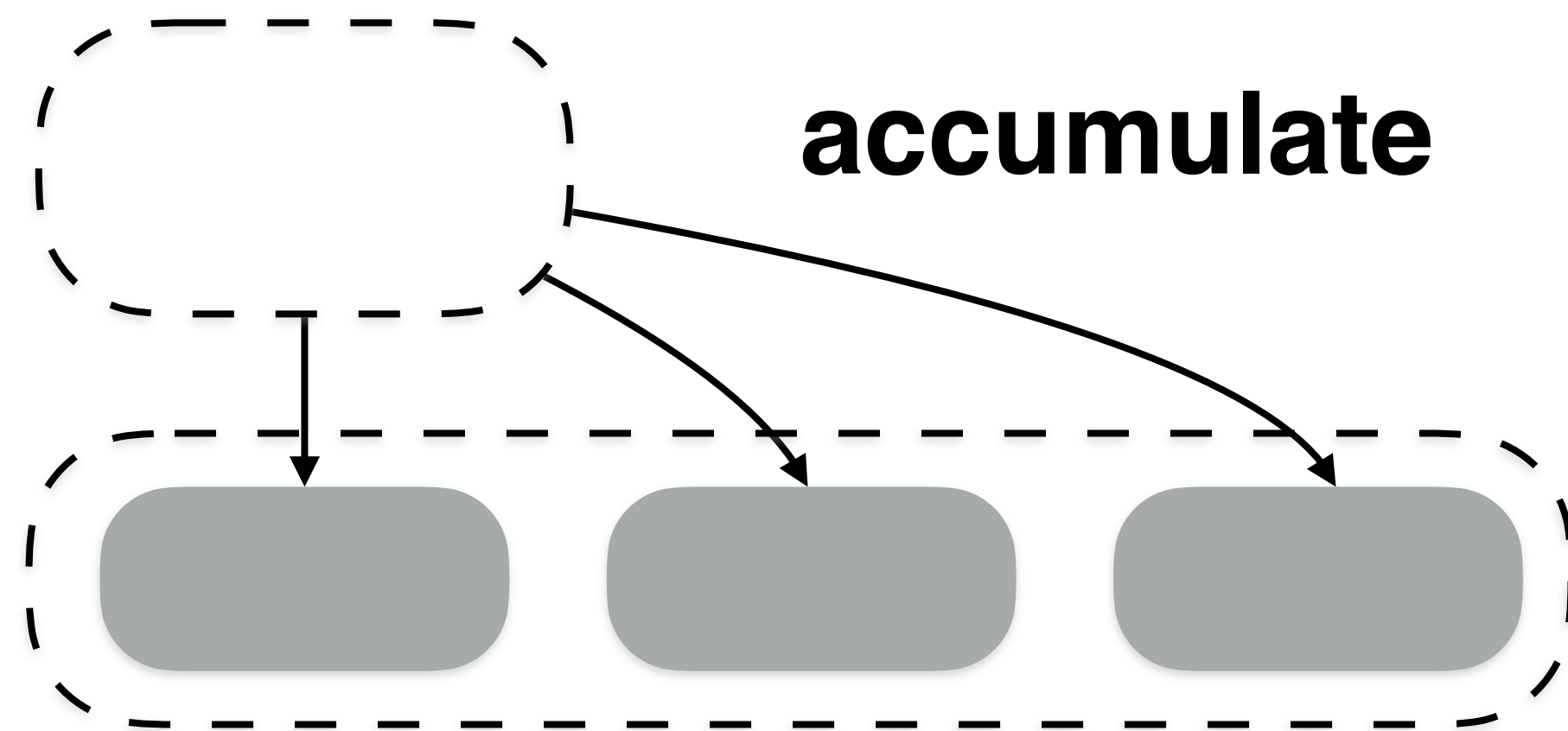
Leveling  
lookup-optimized





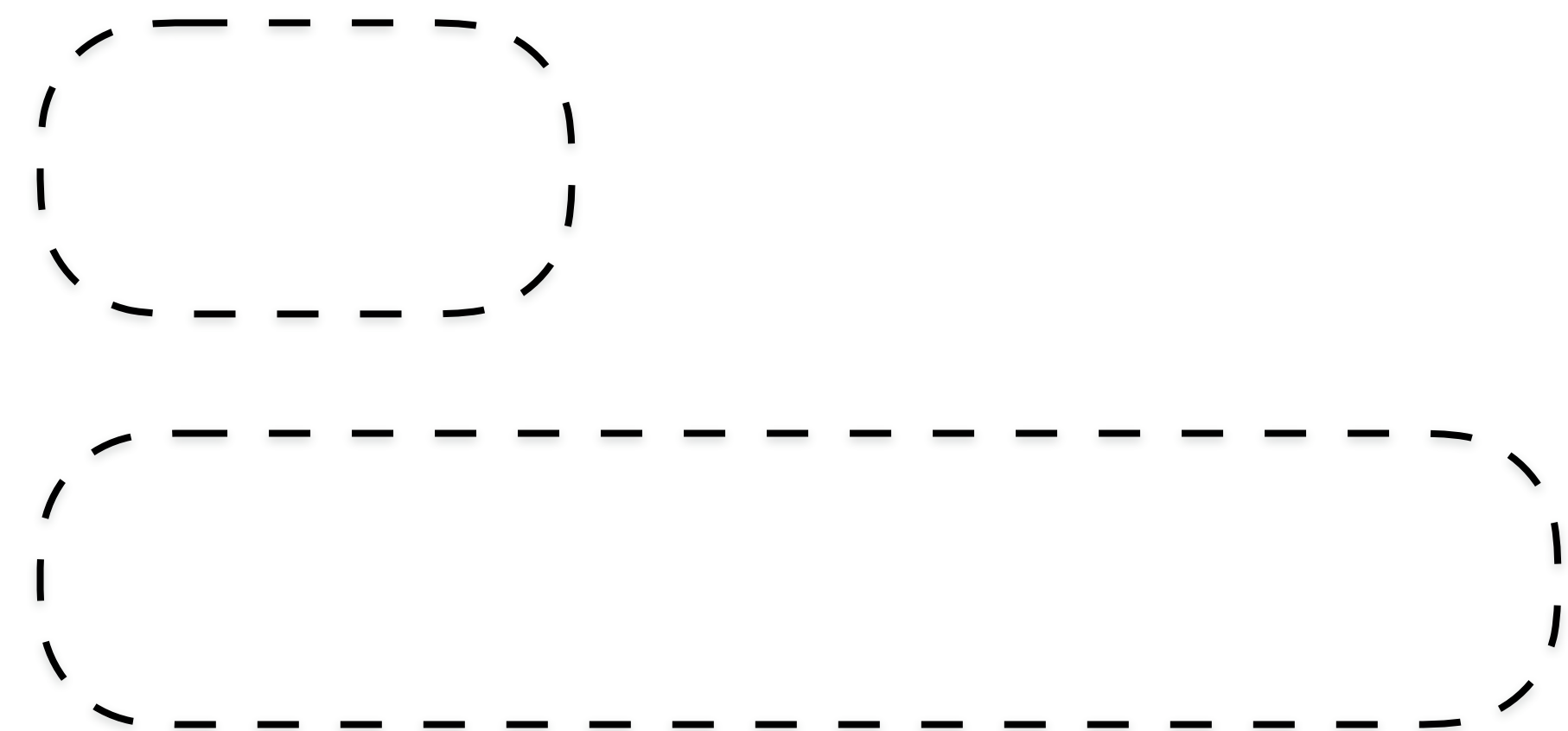
# Tiering

merge-optimized

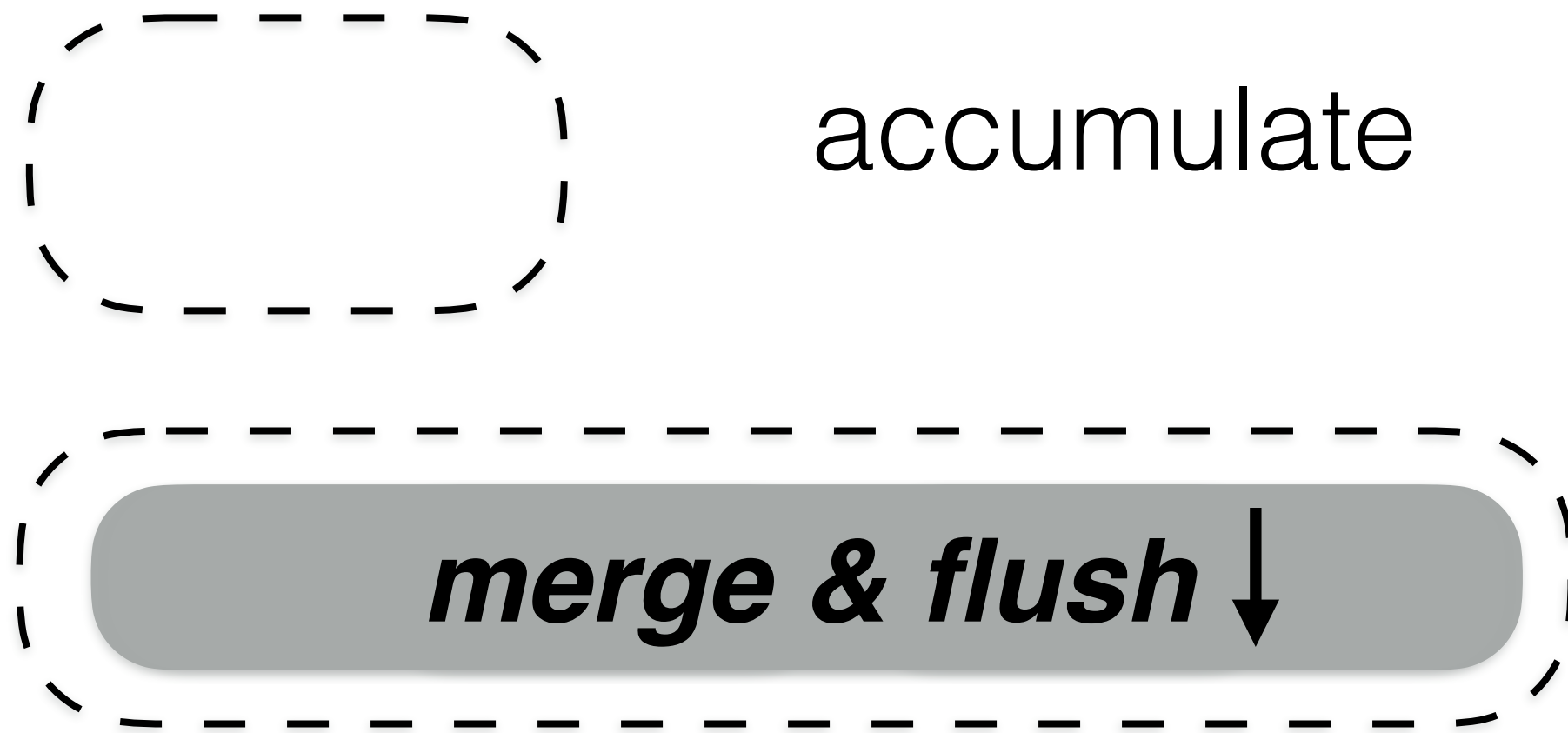


# Leveling

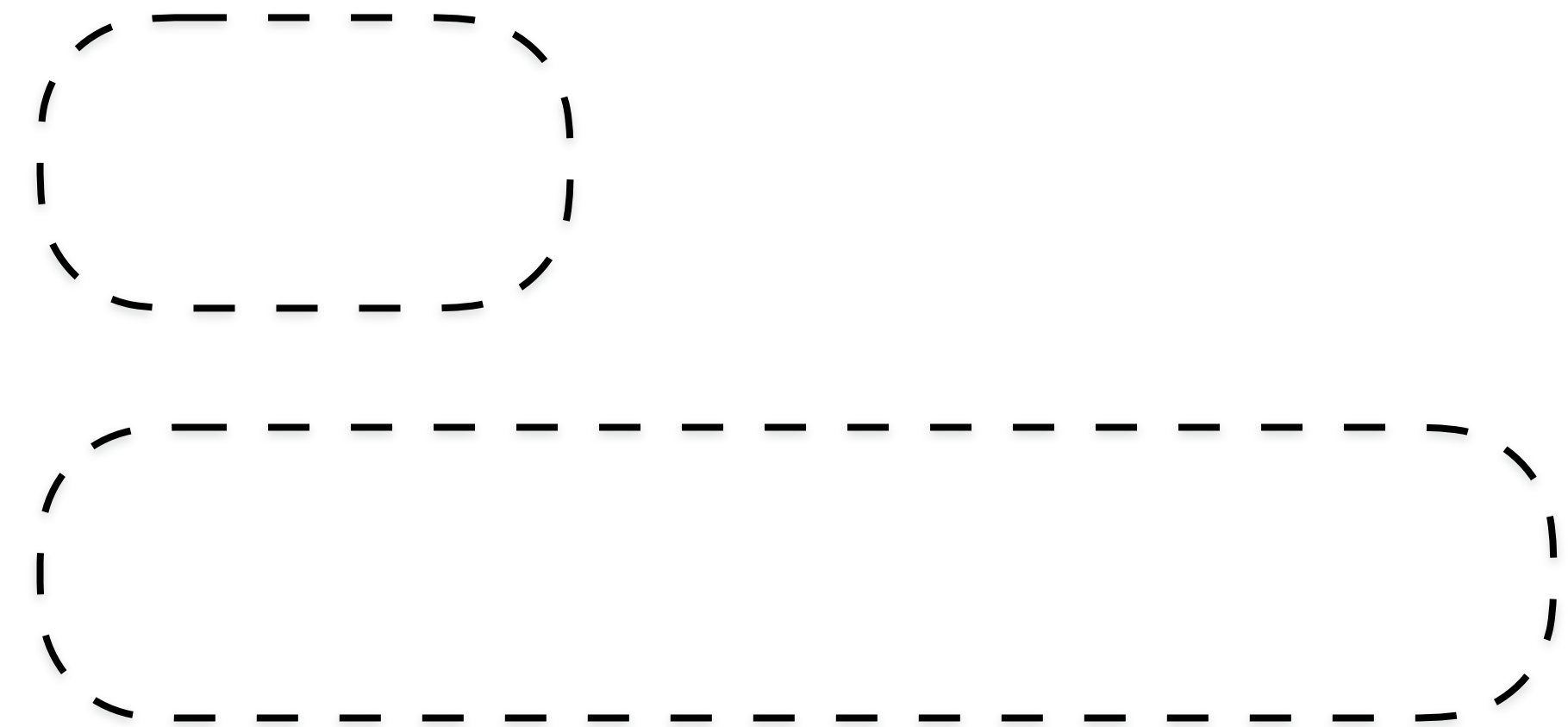
lookup-optimized



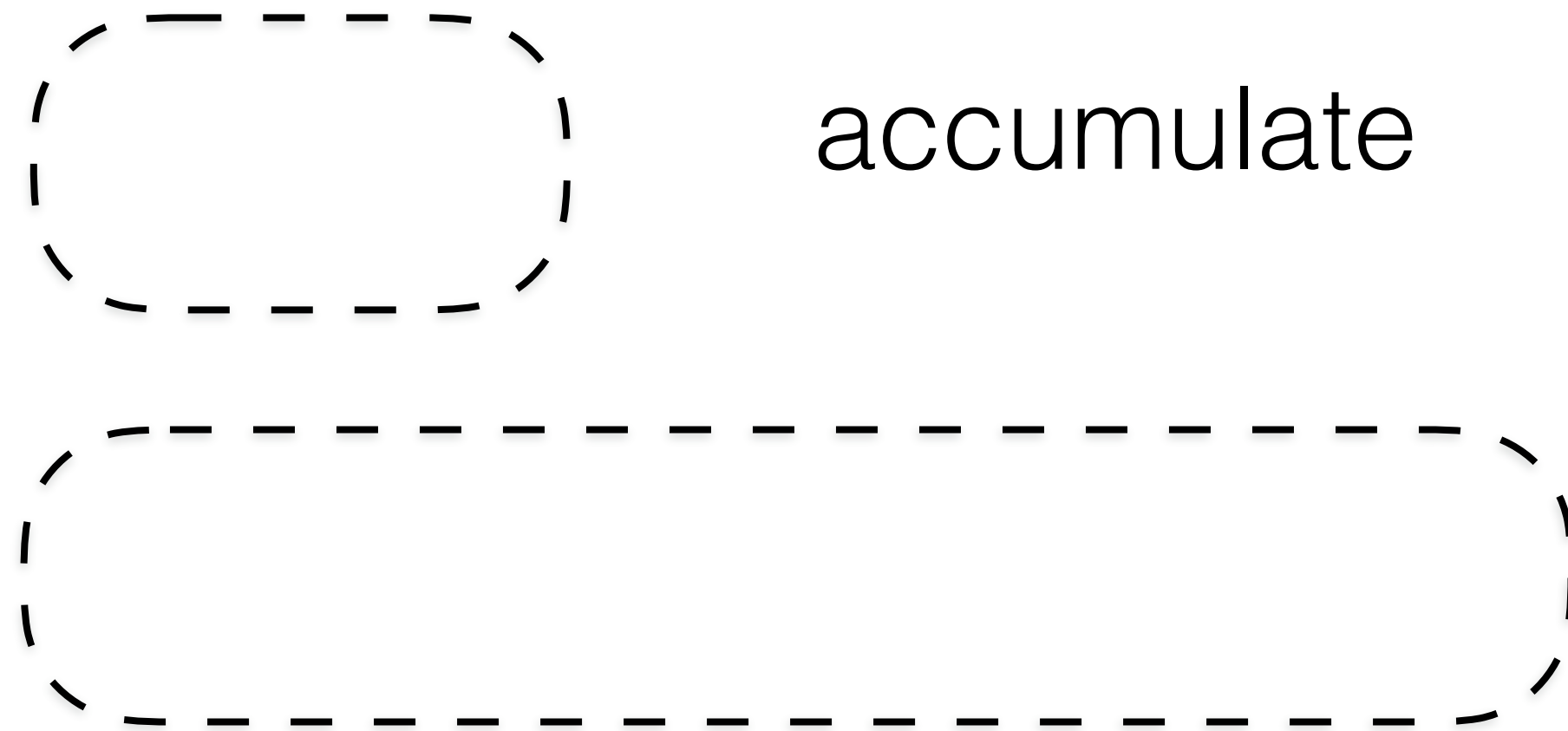
Tiering  
merge-optimized



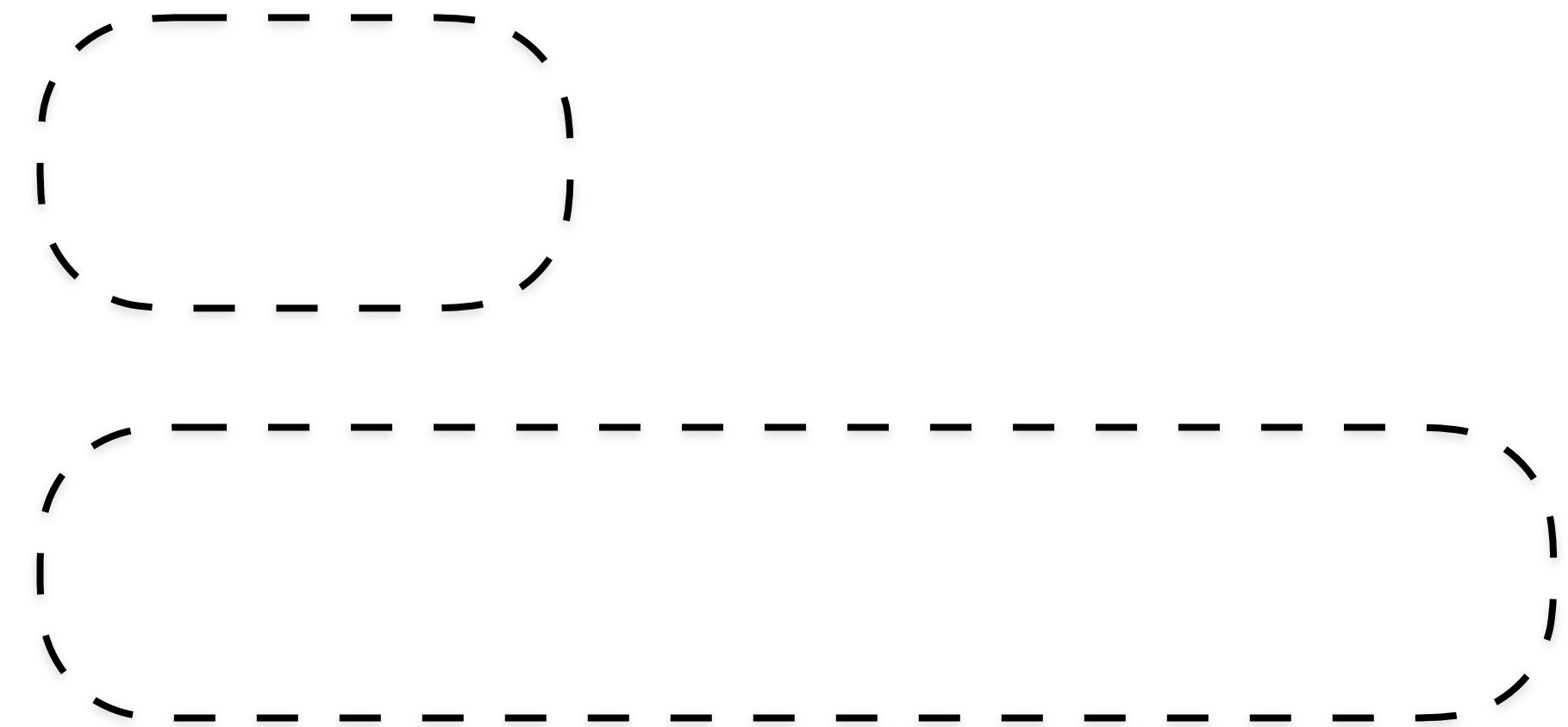
Leveling  
lookup-optimized



Tiering  
merge-optimized

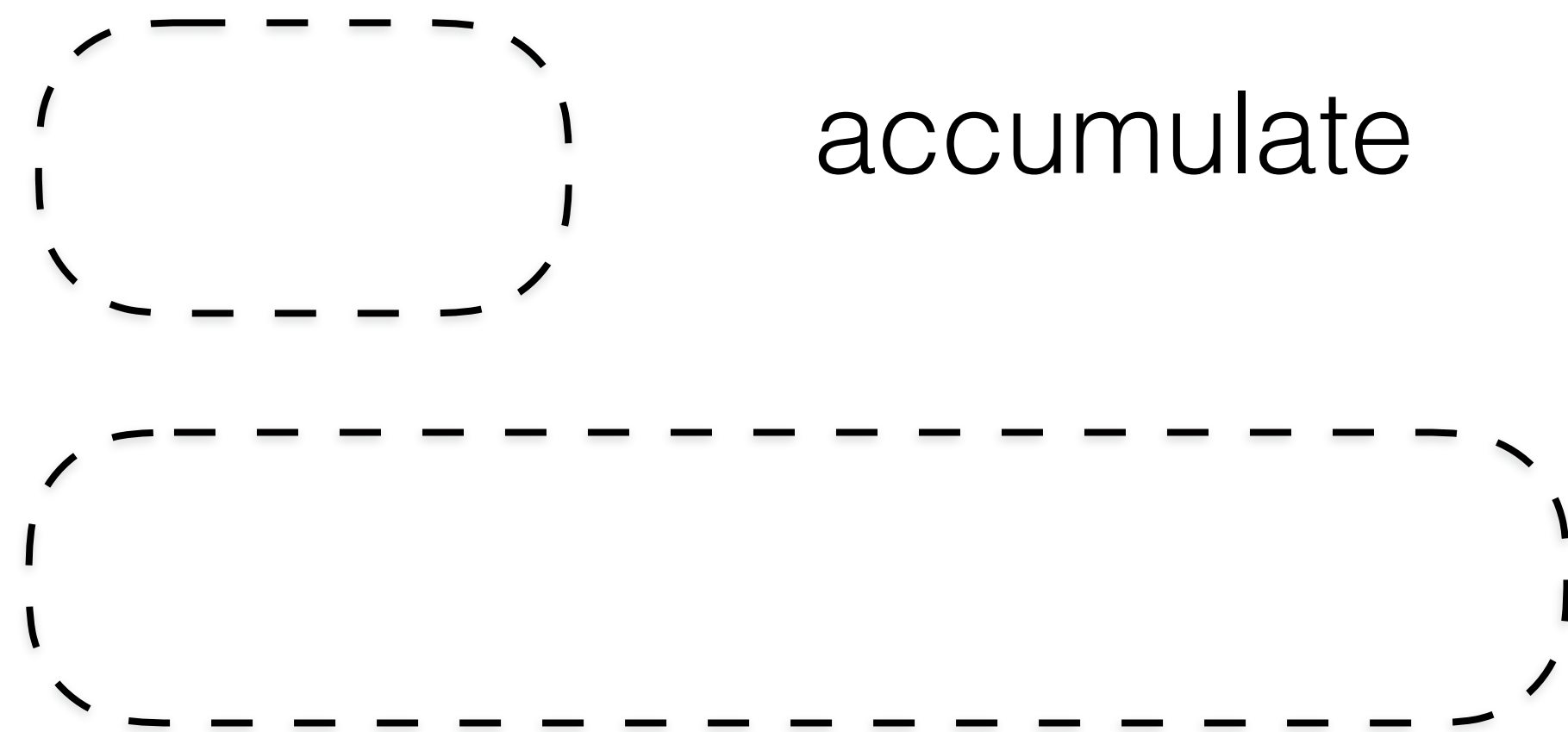


Leveling  
lookup-optimized



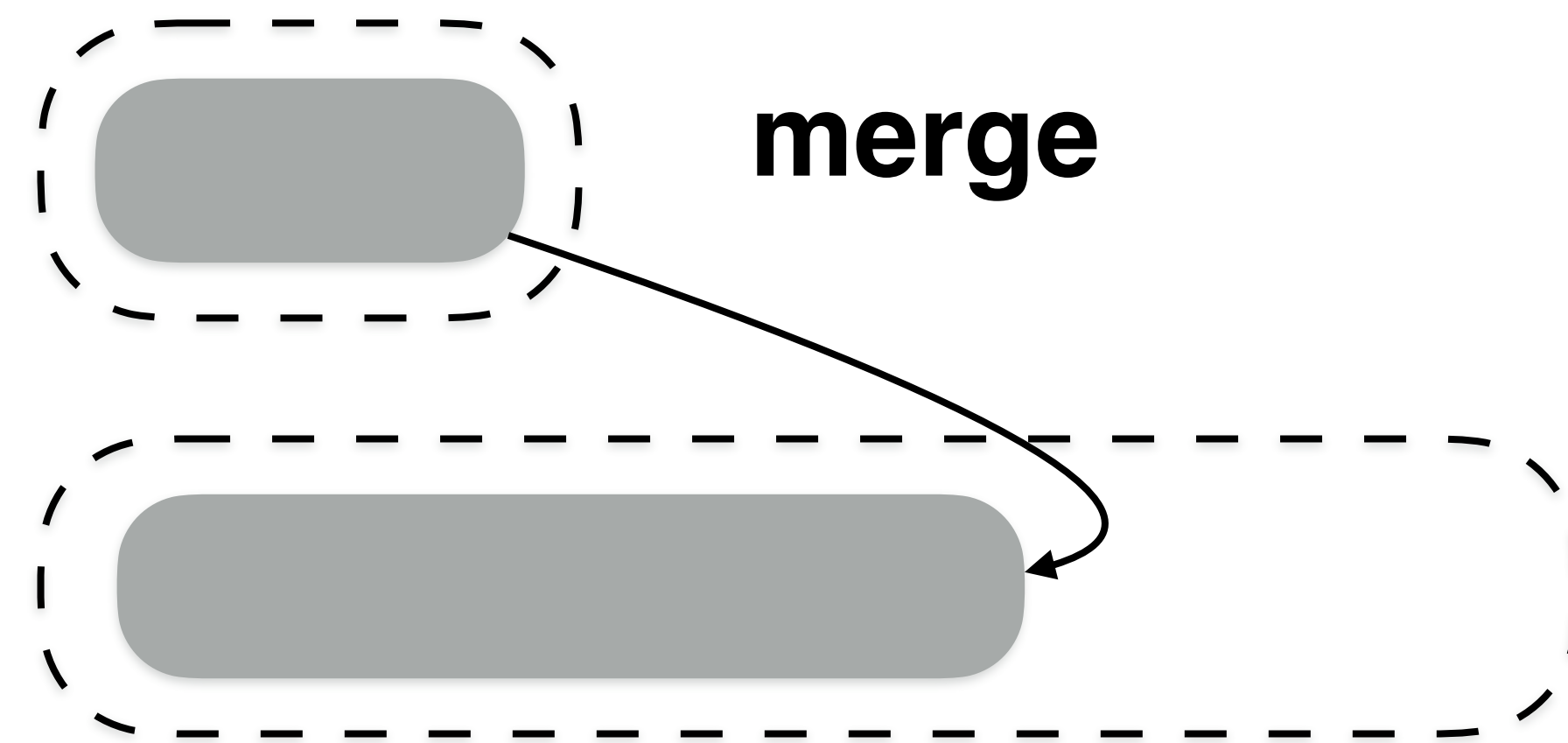
# Tiering

merge-optimized

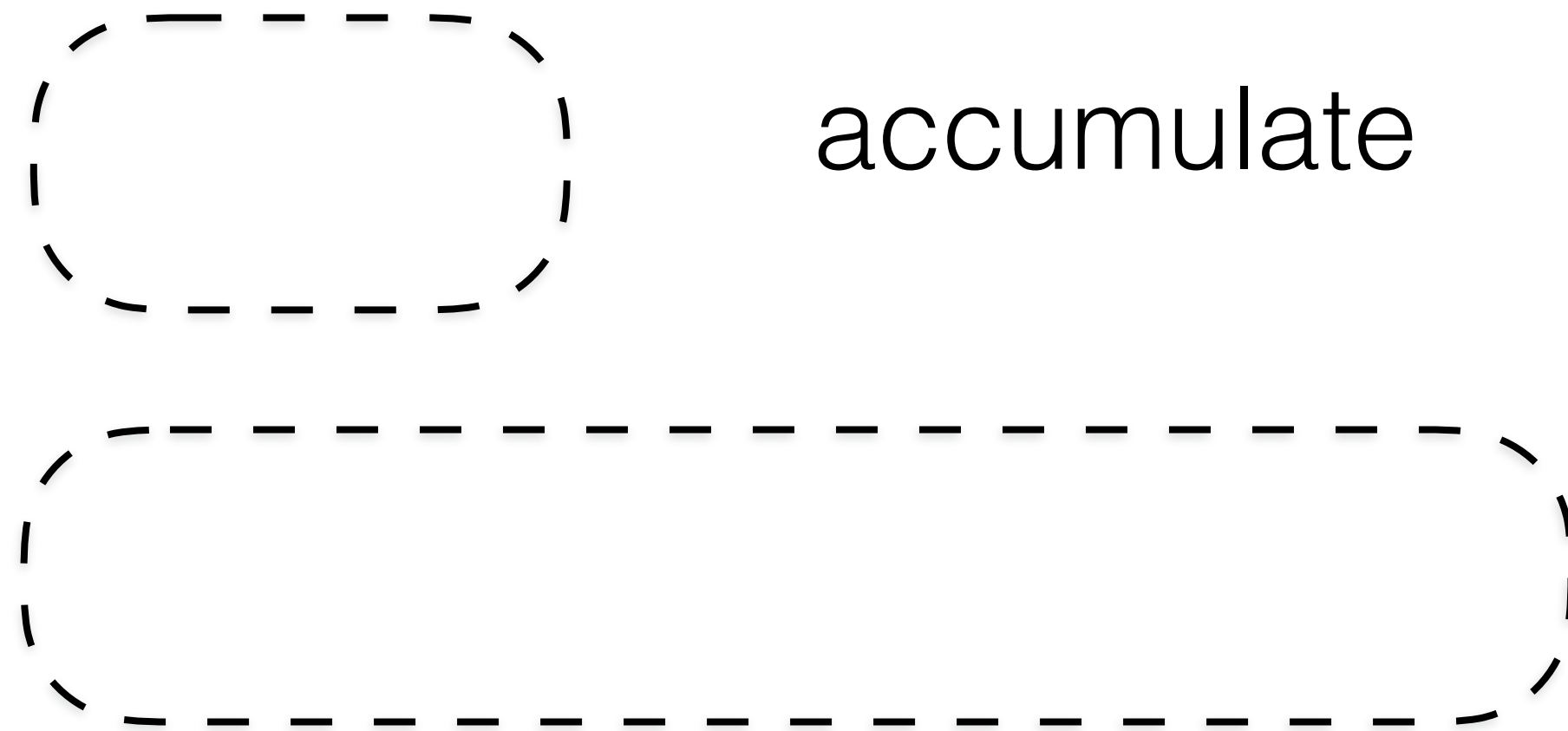


# Leveling

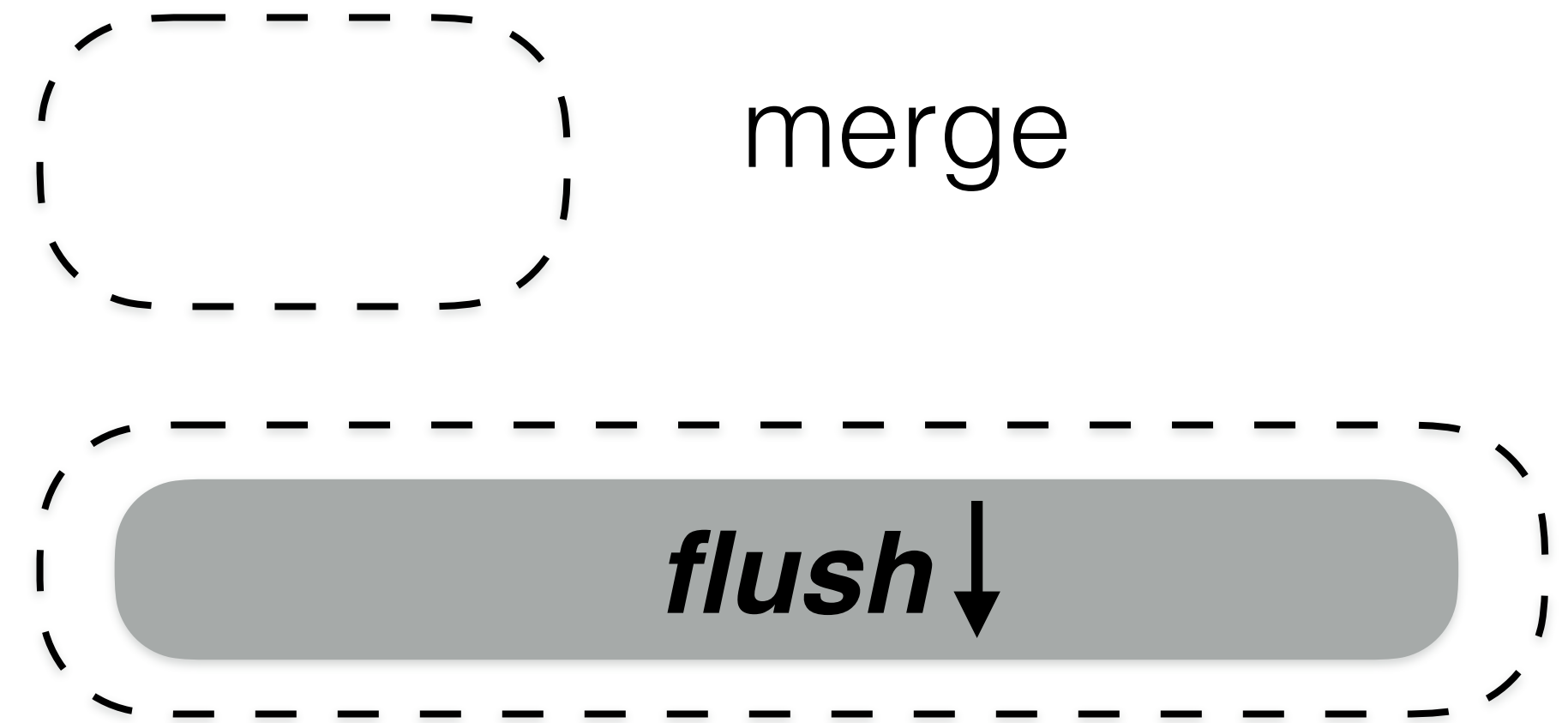
lookup-optimized



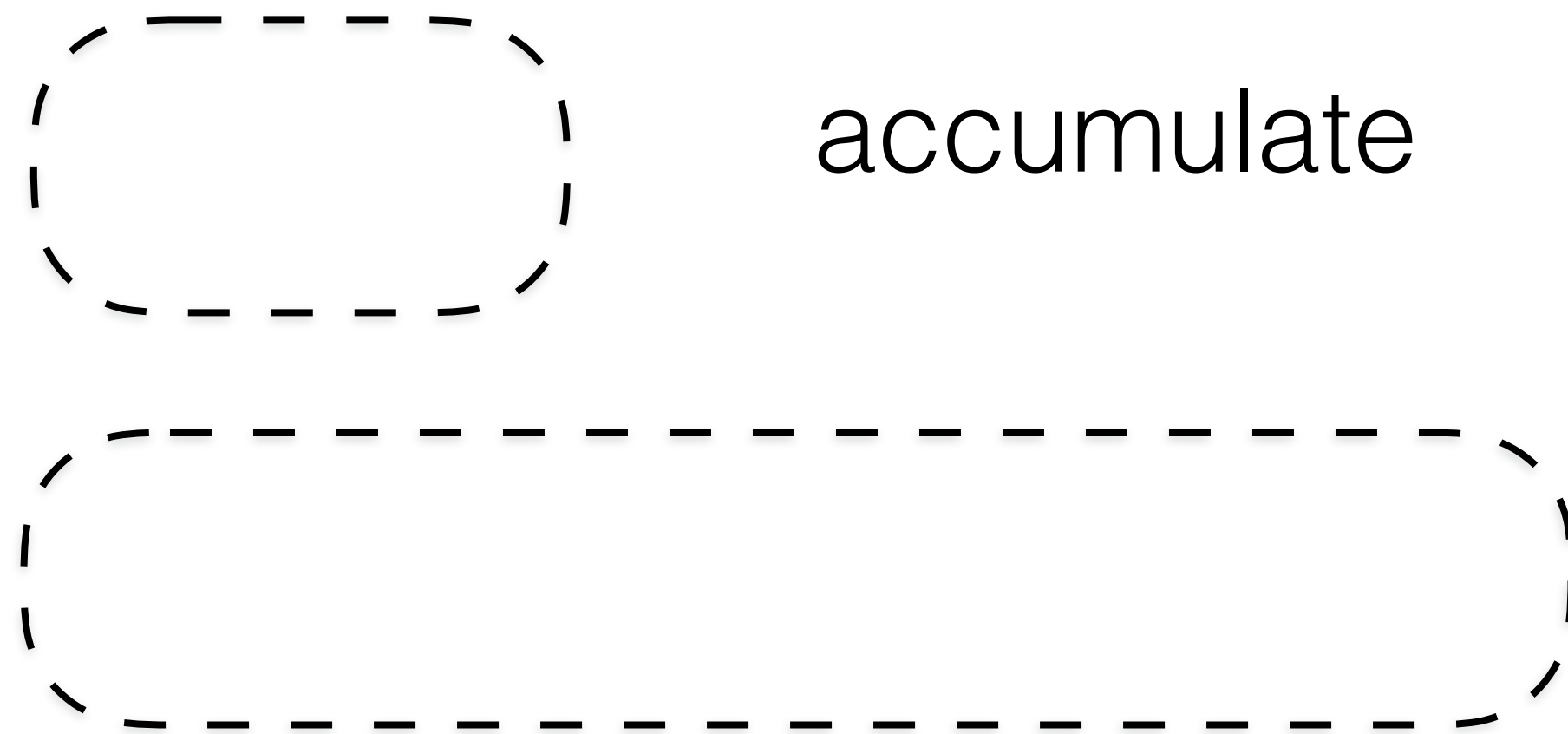
Tiering  
merge-optimized



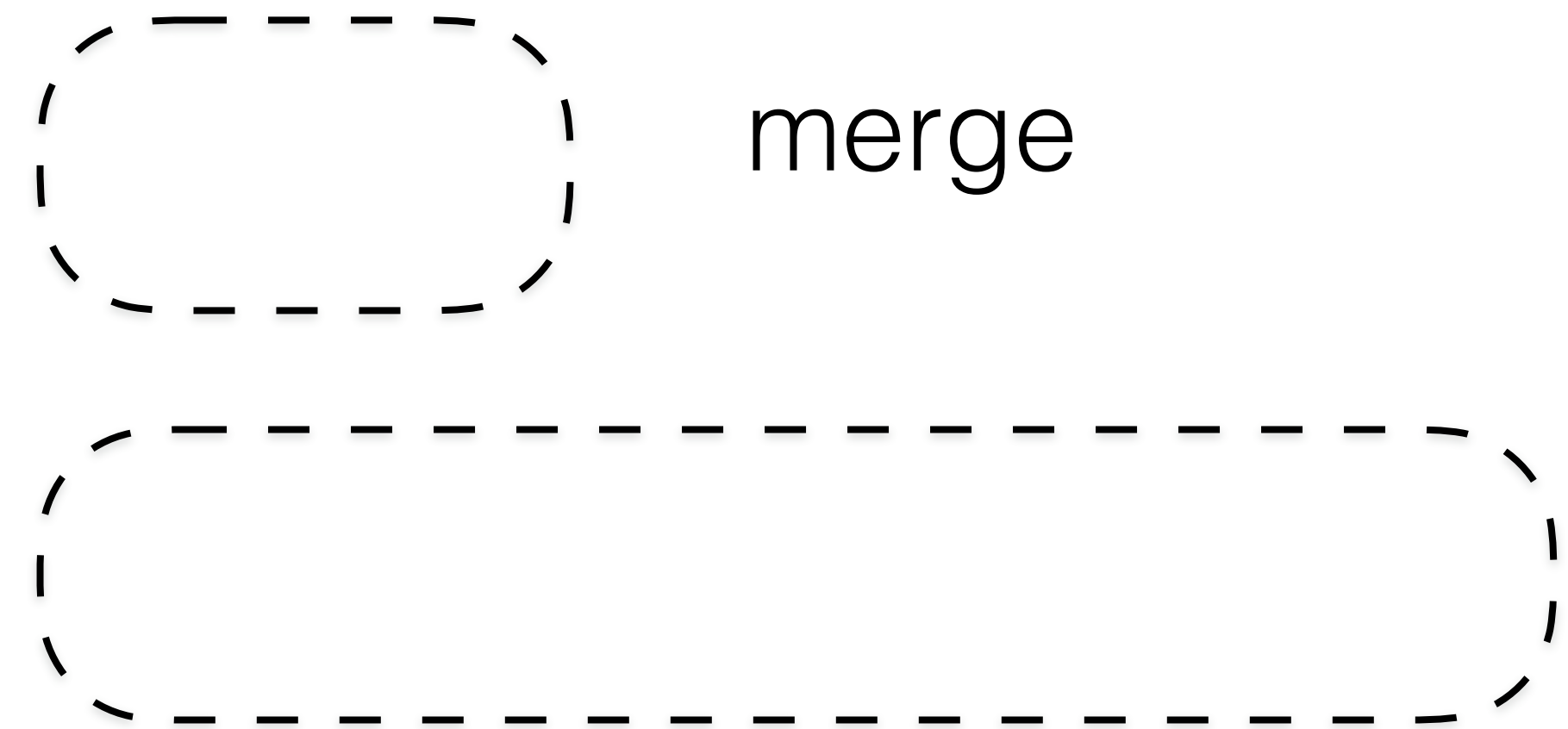
Leveling  
lookup-optimized



Tiering  
merge-optimized

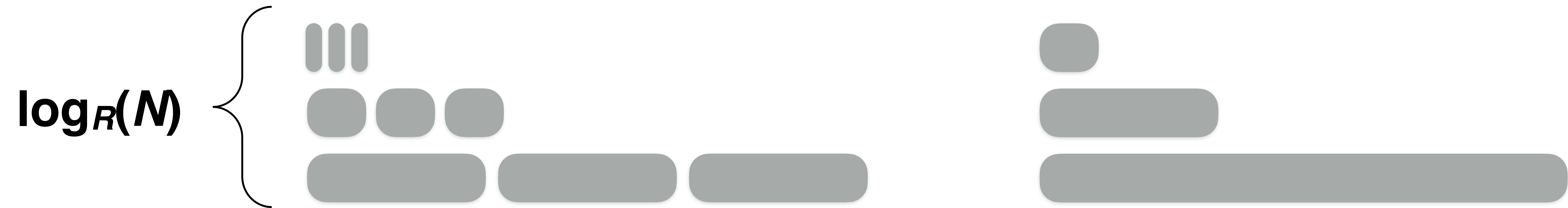


Leveling  
lookup-optimized



Tiering  
merge-optimized

Leveling  
lookup-optimized



Tiering  
merge-optimized

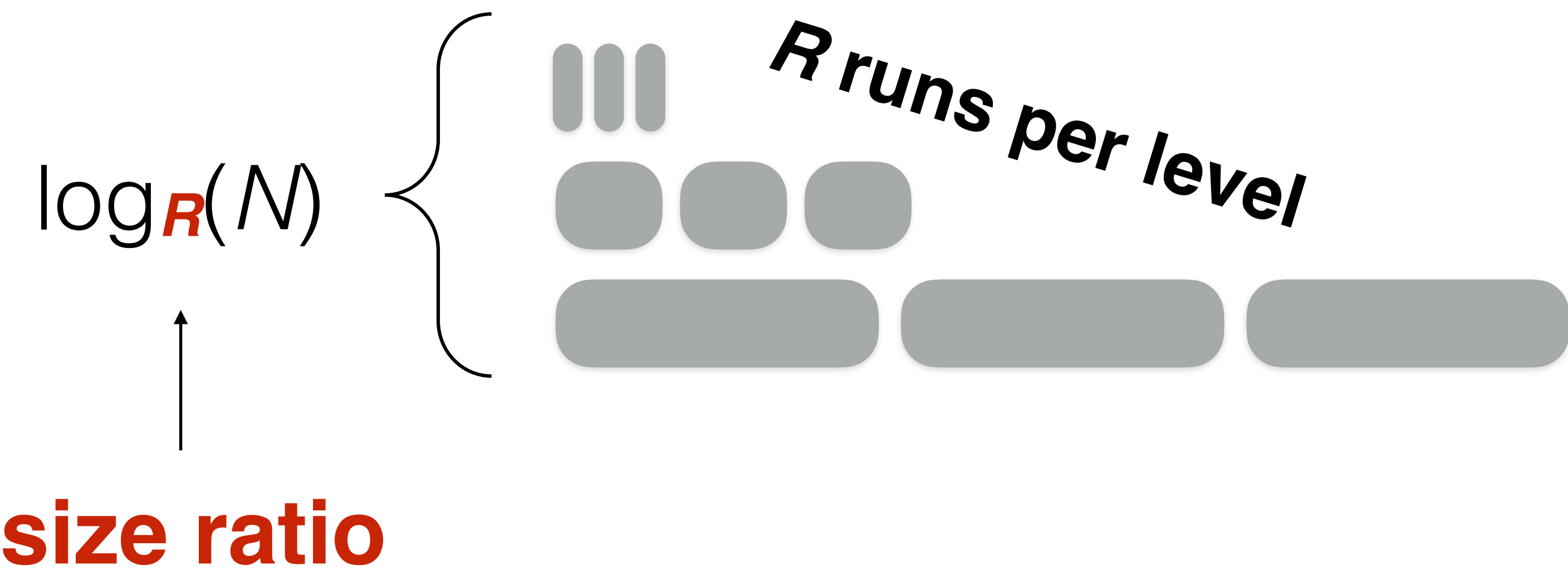
Leveling  
lookup-optimized





Tiering  
merge-optimized

Leveling  
lookup-optimized



Tiering  
merge-optimized



Leveling  
lookup-optimized

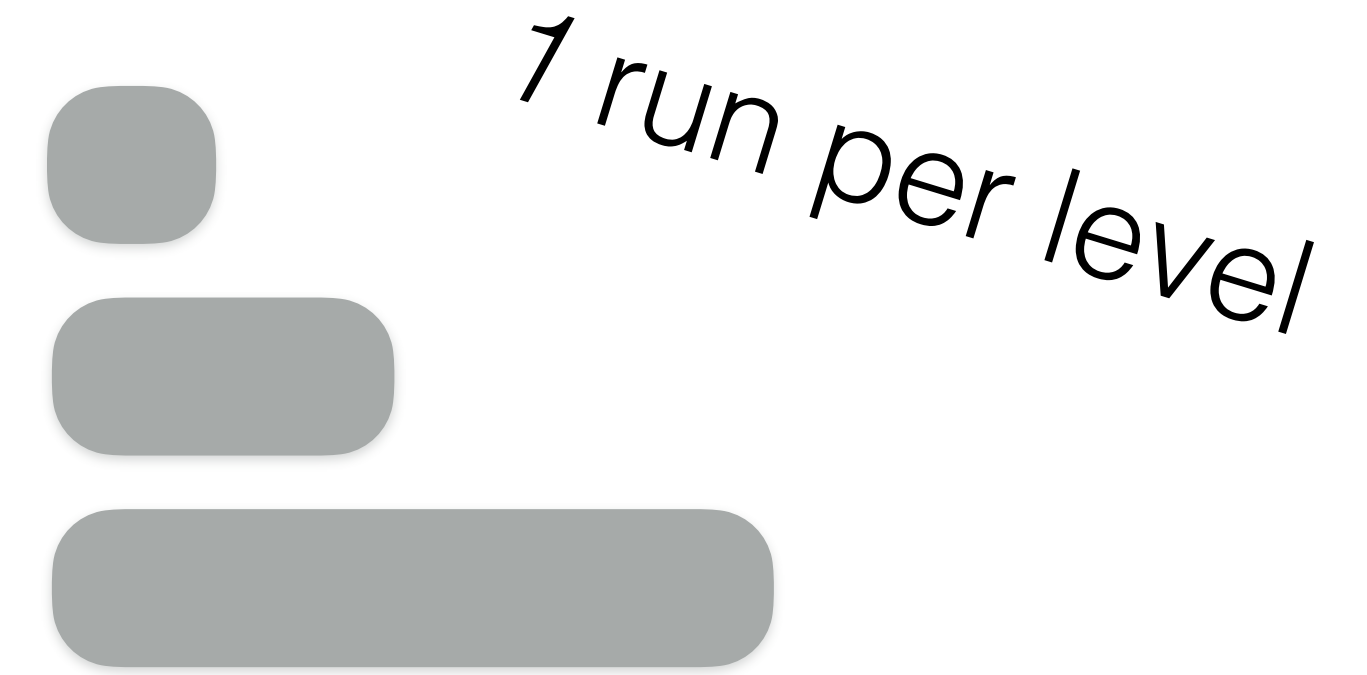


 **size ratio  $R$**

Tiering  
merge-optimized

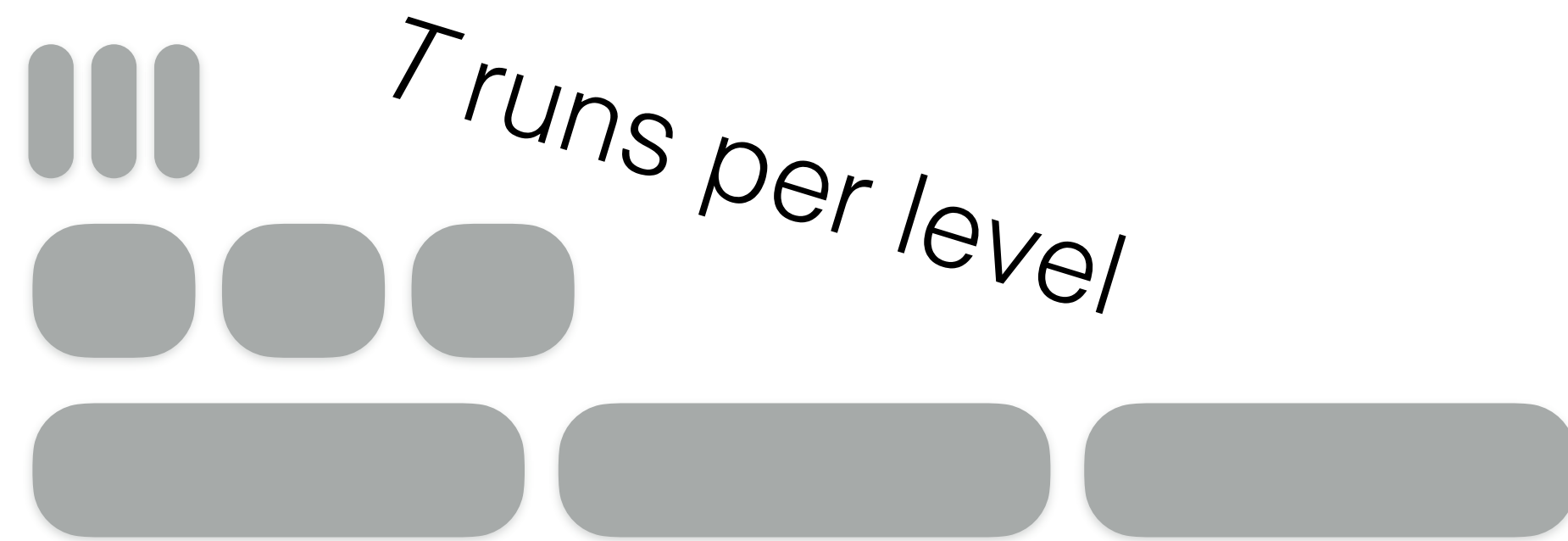


Leveling  
lookup-optimized



 size ratio  $R \Downarrow$

Tiering  
merge-optimized

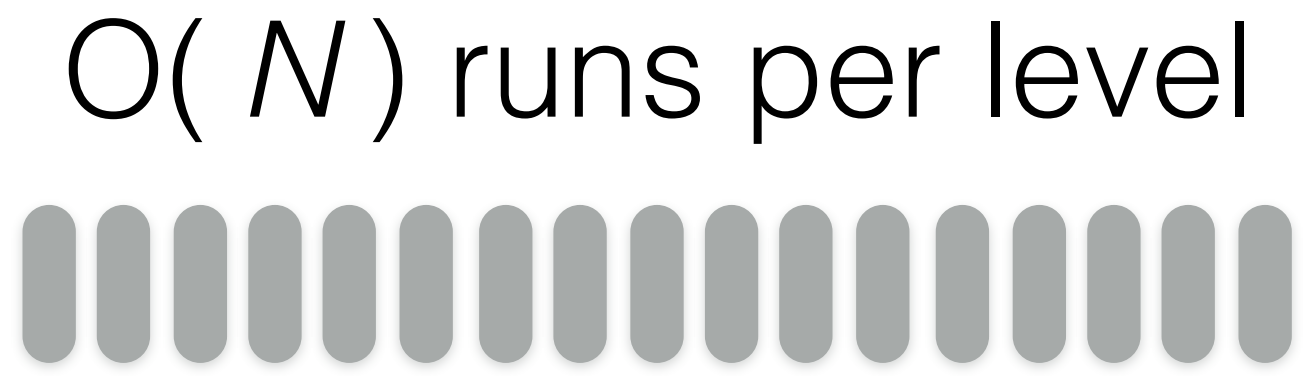


Leveling  
lookup-optimized



 **size ratio  $R$**

Tiering  
merge-optimized



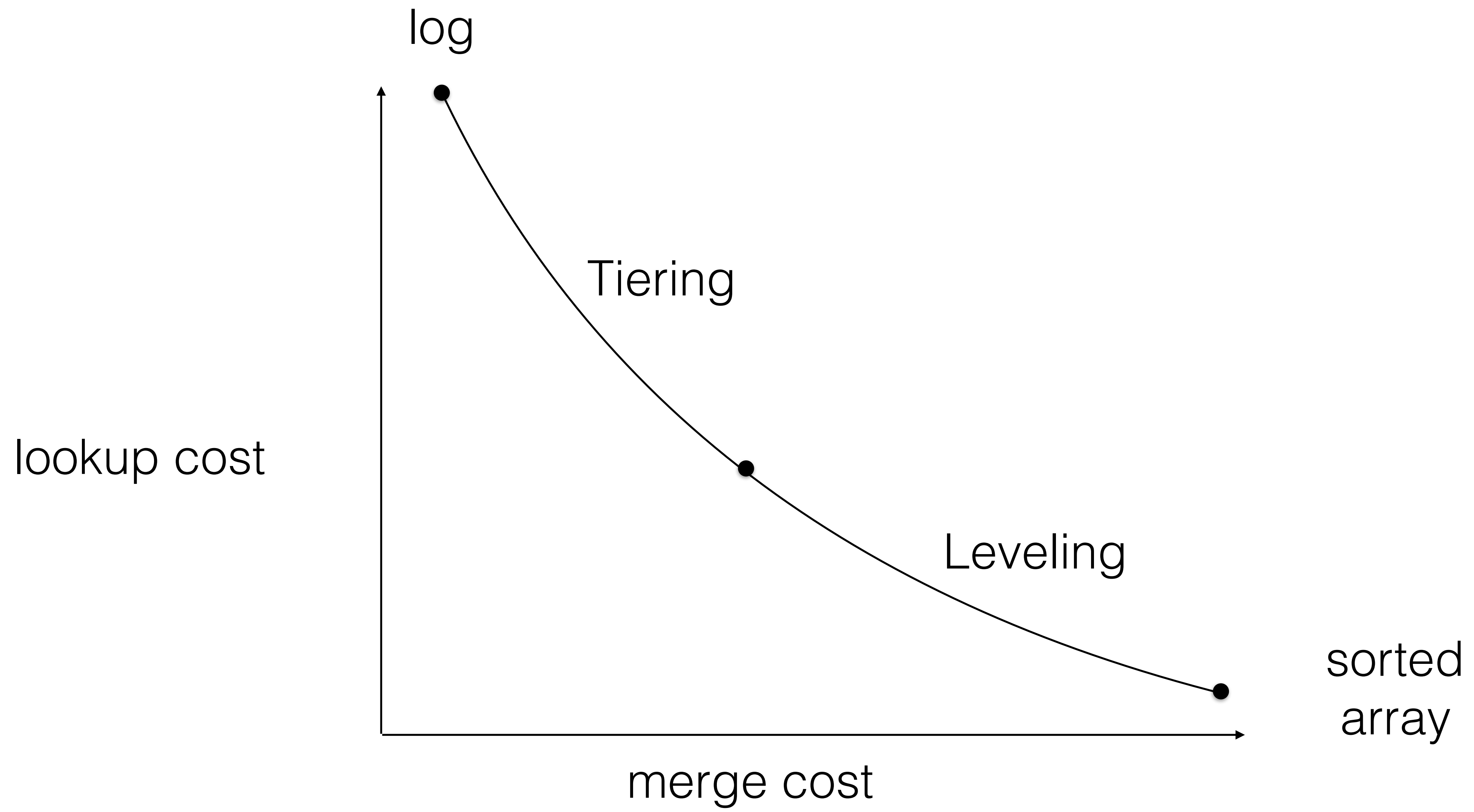
**log**

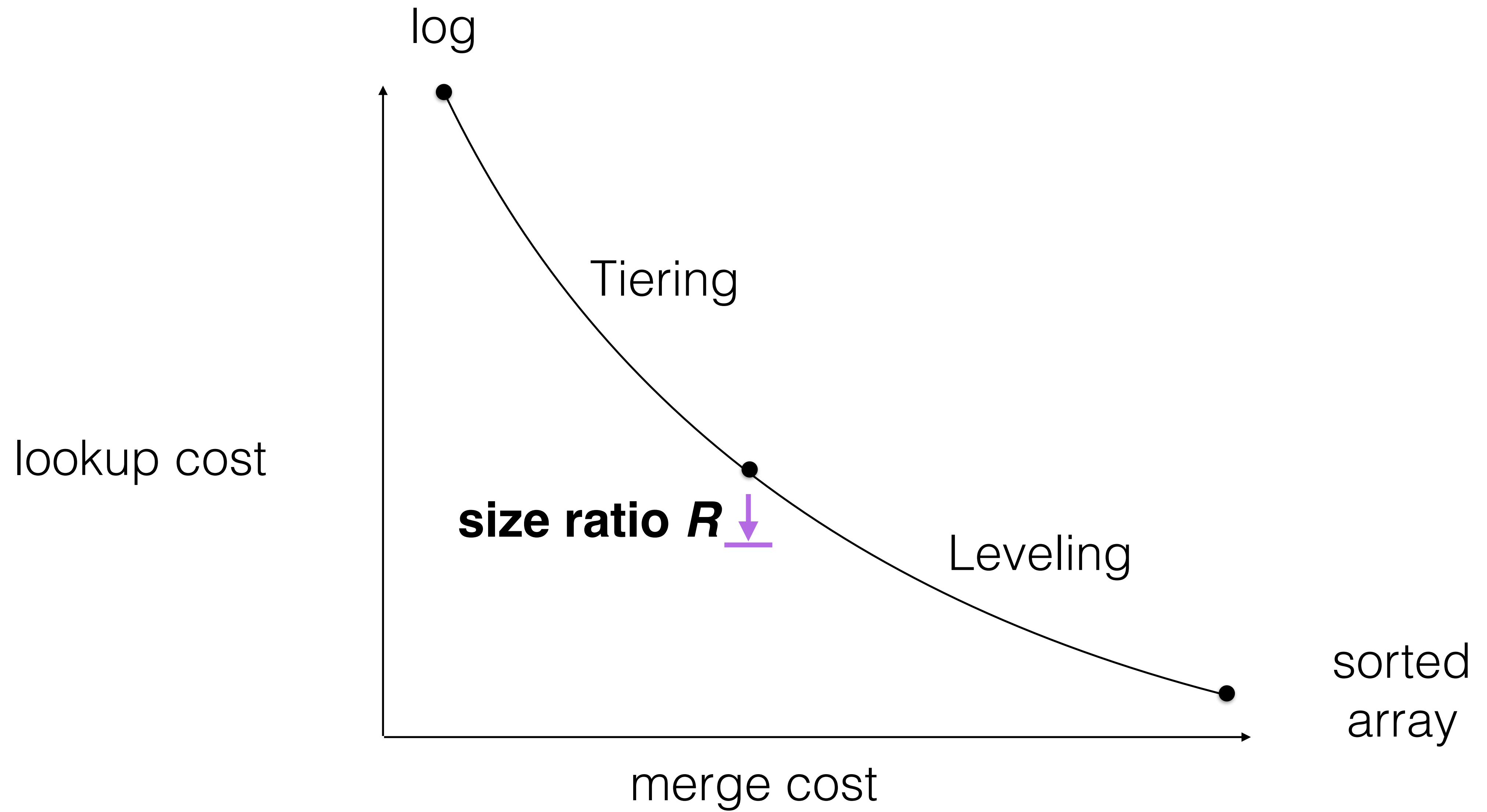
Leveling  
lookup-optimized

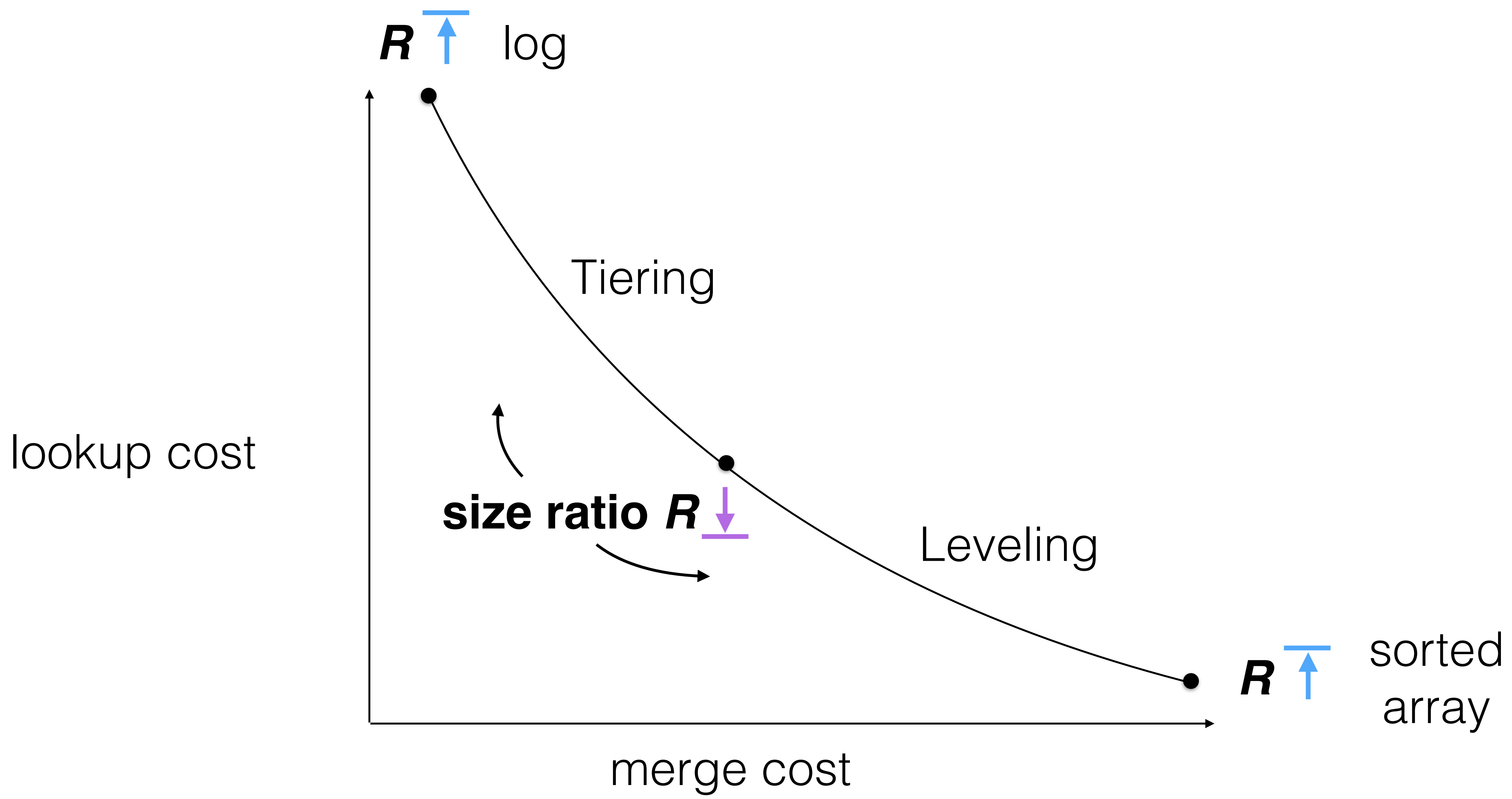


**sorted  
array**

 **size ratio  $R \gg$**

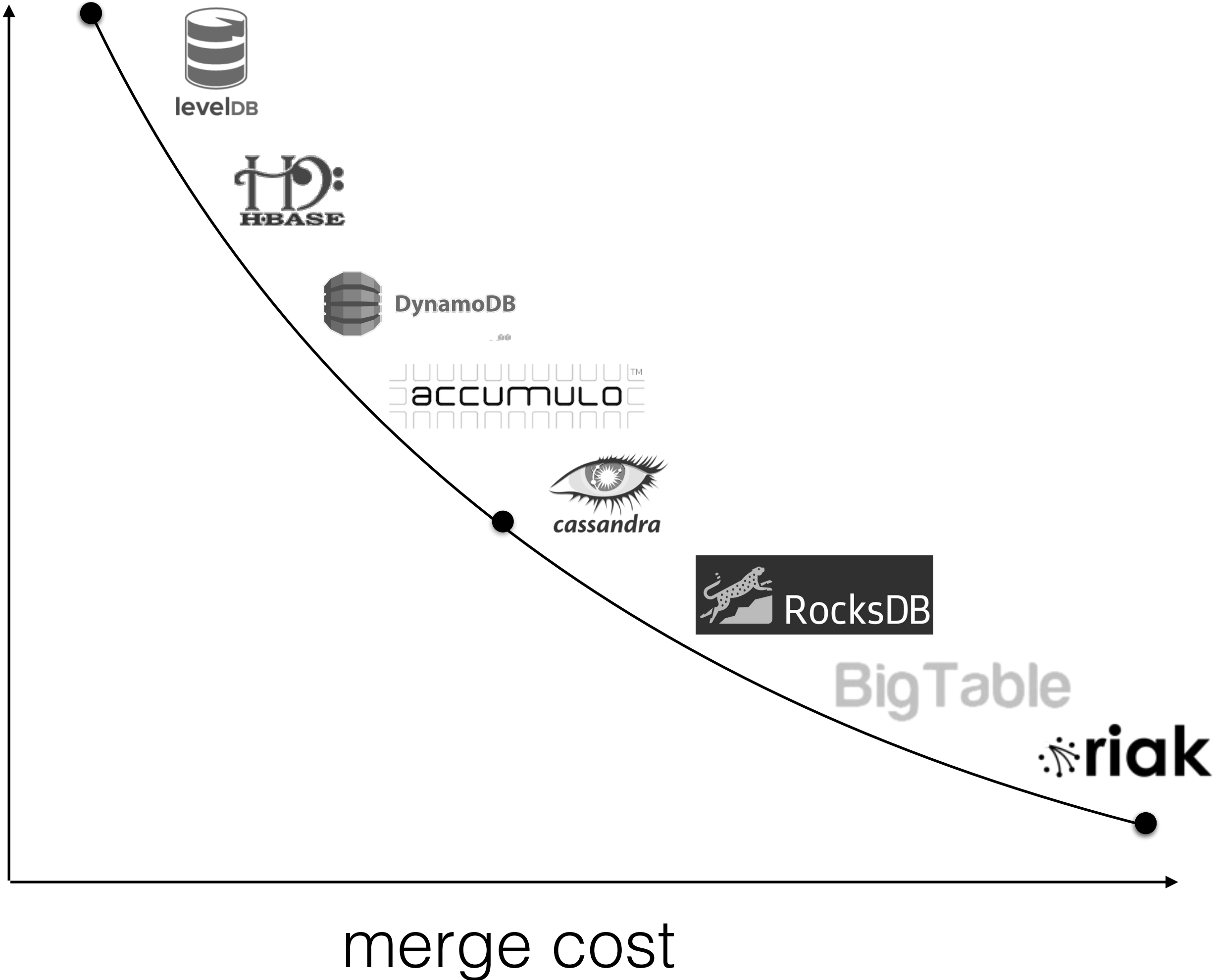






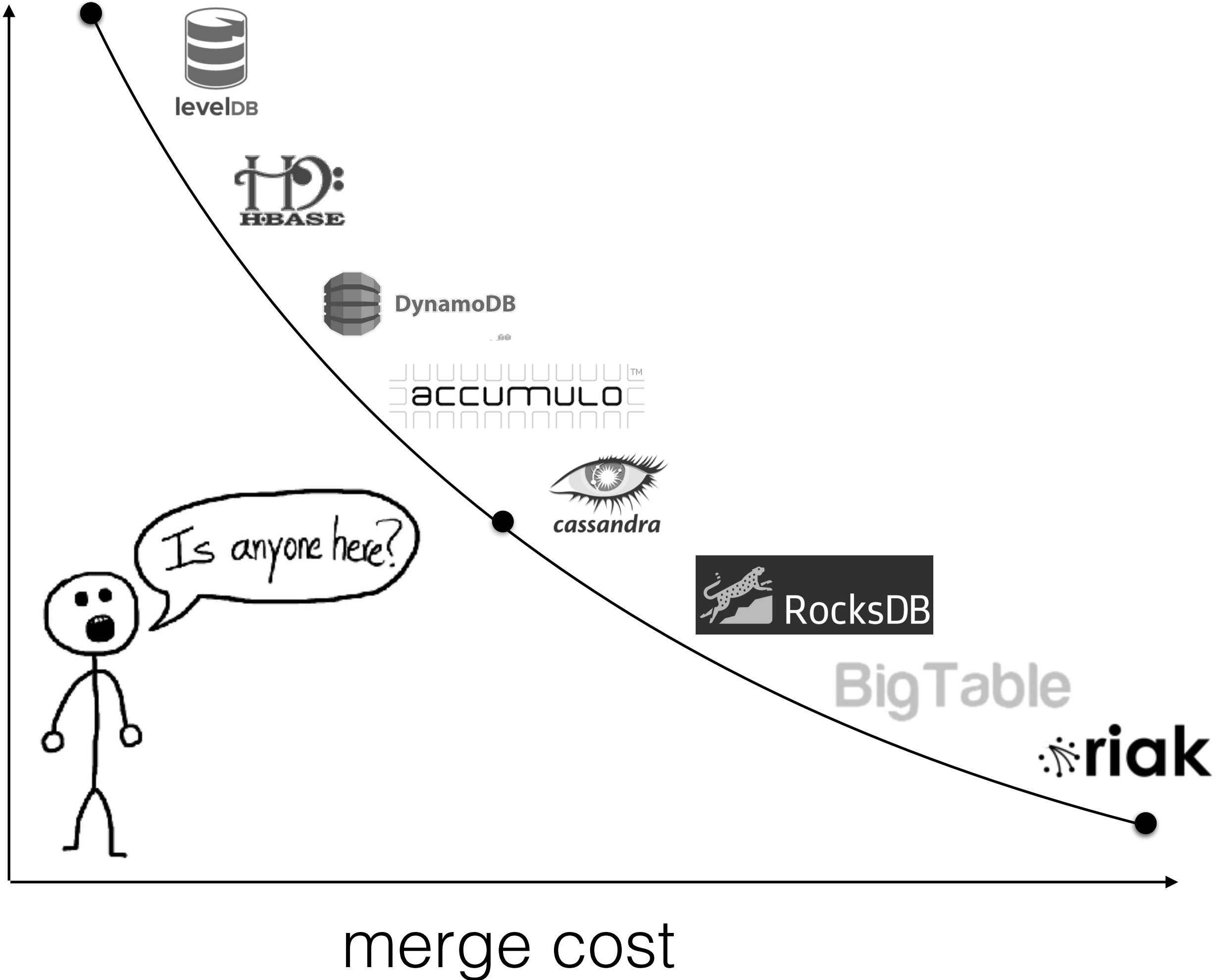


lookup cost

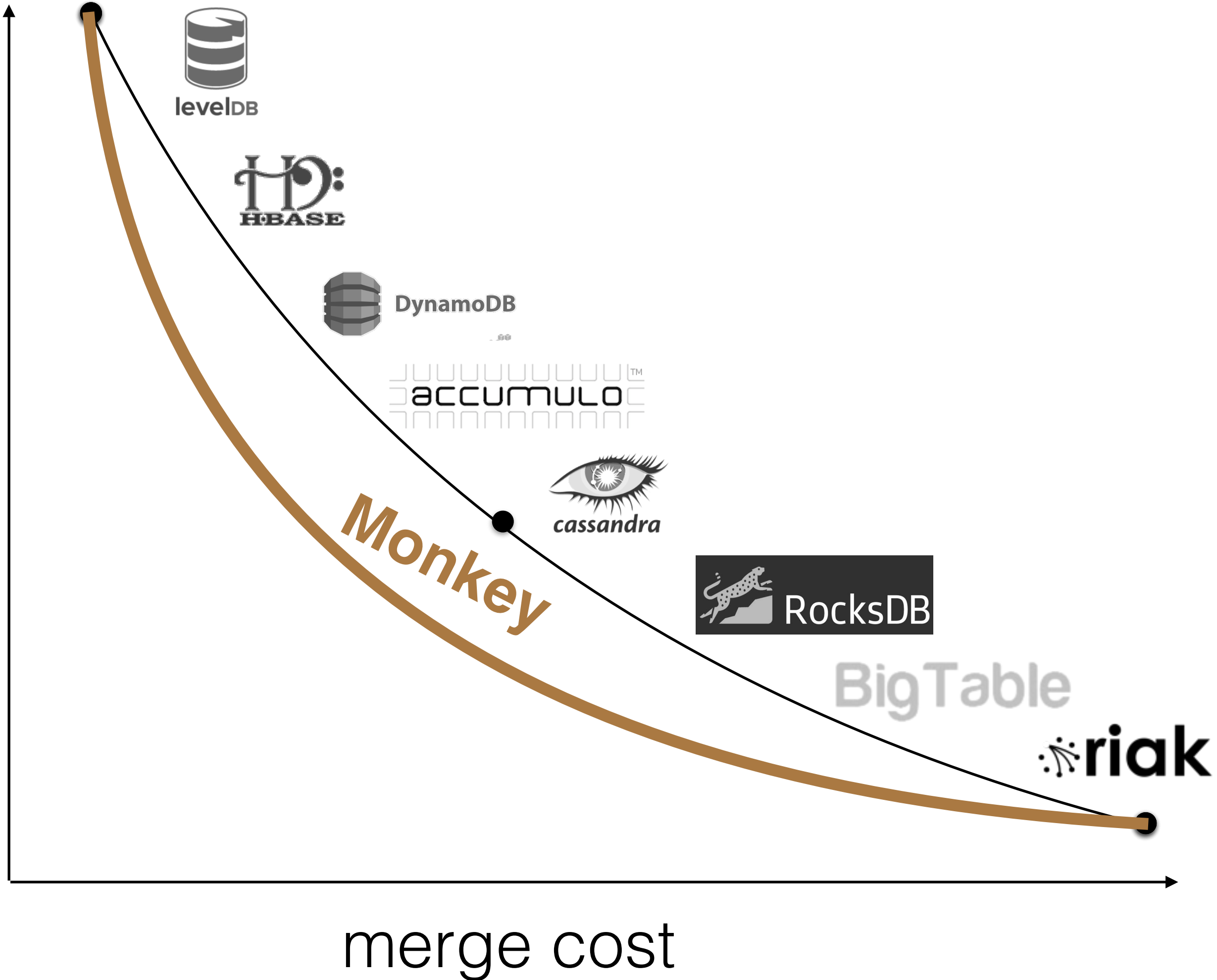


merge cost

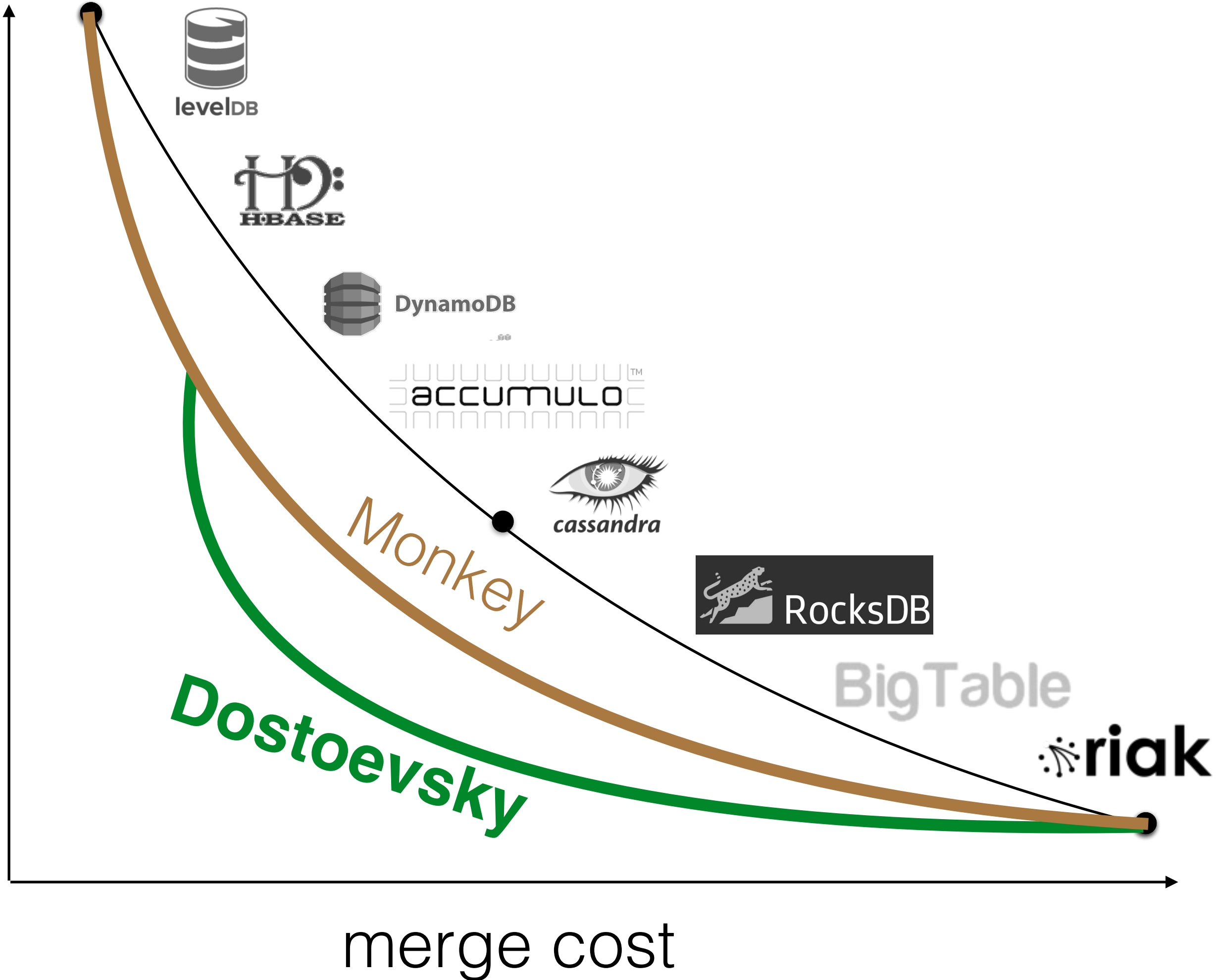
lookup cost



lookup cost

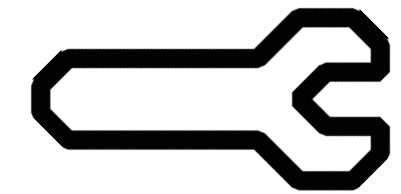


lookup cost



# Monkey: Optimal Navigable Key-Value Store

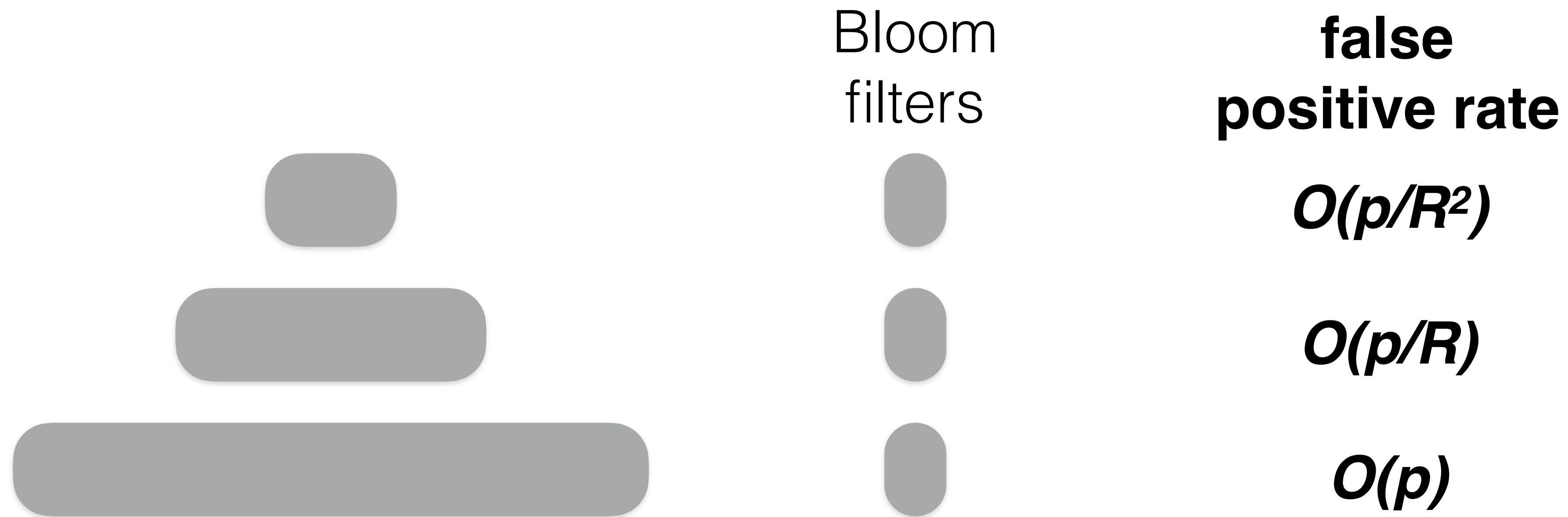
SIGMOD17

 **Bloom  
filters**



# Monkey: Optimal Navigable Key-Value Store

SIGMOD17



# Monkey: Optimal Navigable Key-Value Store

SIGMOD17



Bloom  
filters



**false  
positive rate**

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

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

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

# Monkey: Optimal Navigable Key-Value Store

SIGMOD17



Bloom  
filters

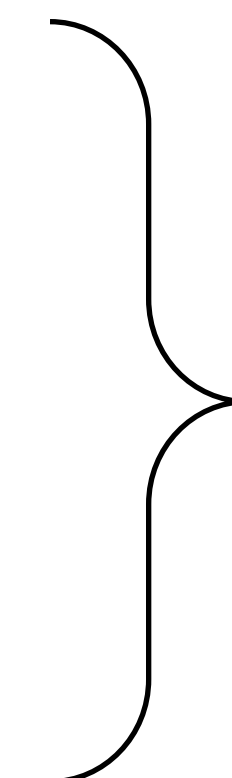


false  
positive rate

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

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

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



**geometric  
progression**

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



# Monkey: Optimal Navigable Key-Value Store

SIGMOD17



Bloom  
filters

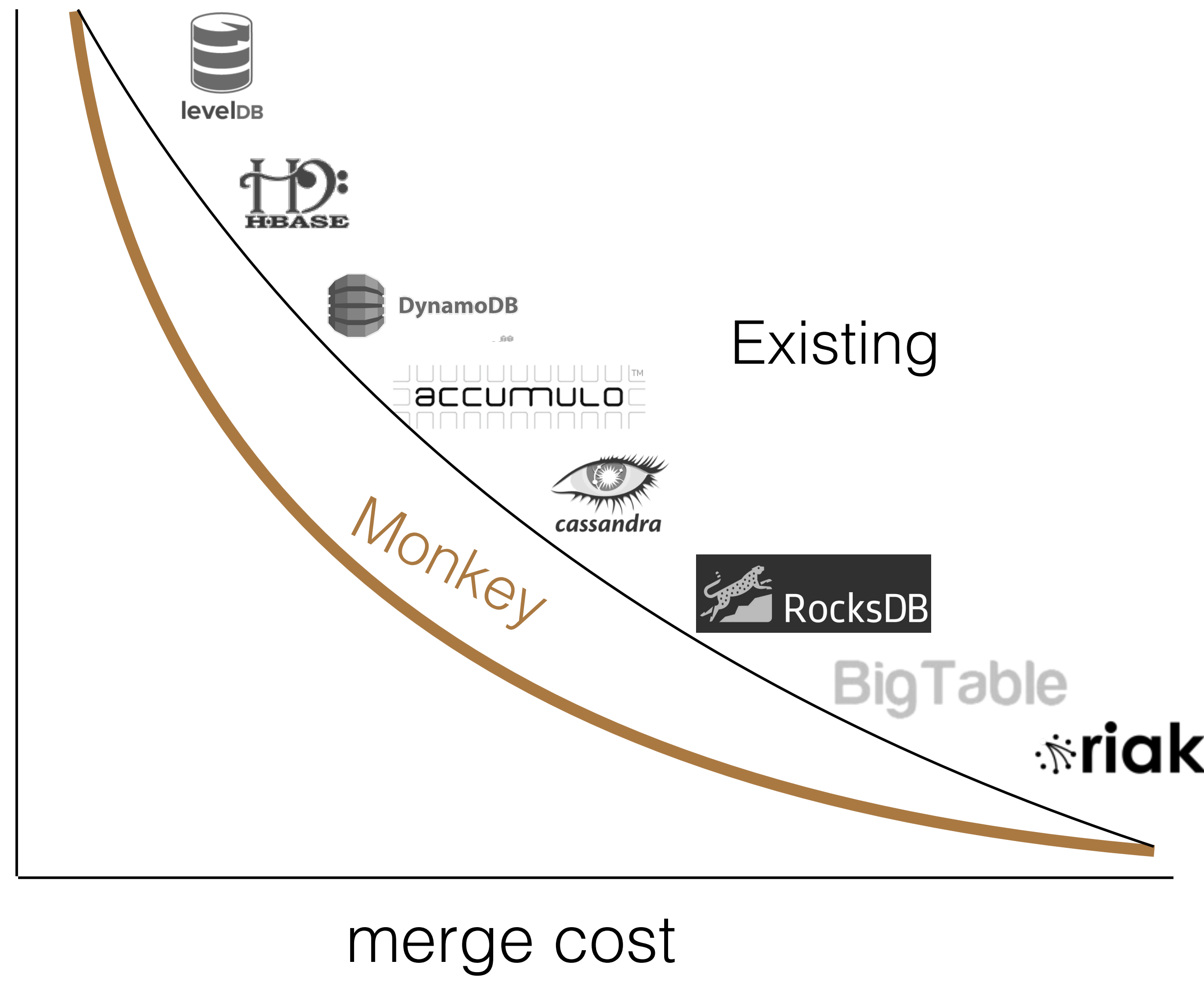


$$O(\mathbf{e^{-x}} \cdot \mathbf{\log_R(N)})$$

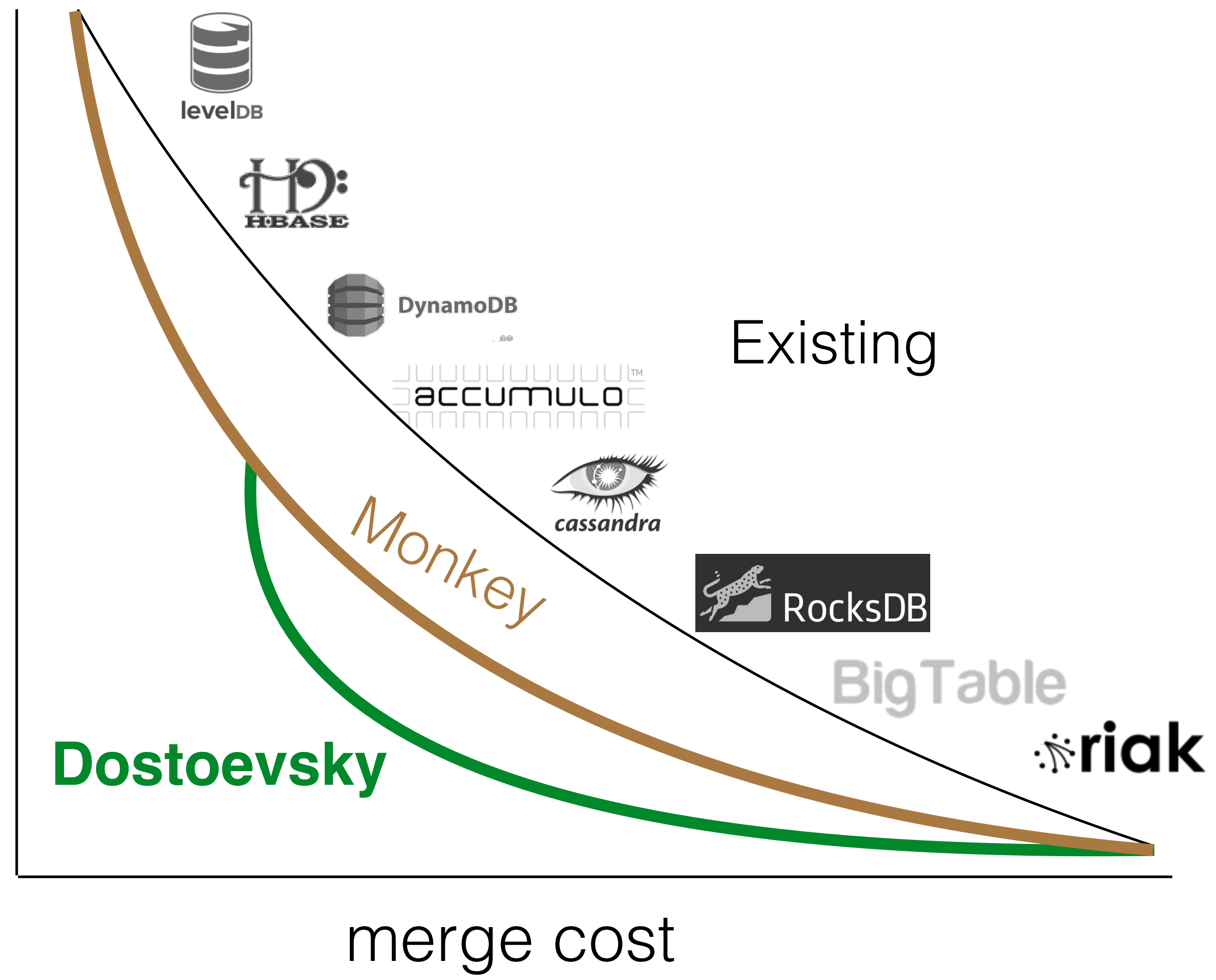
>

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

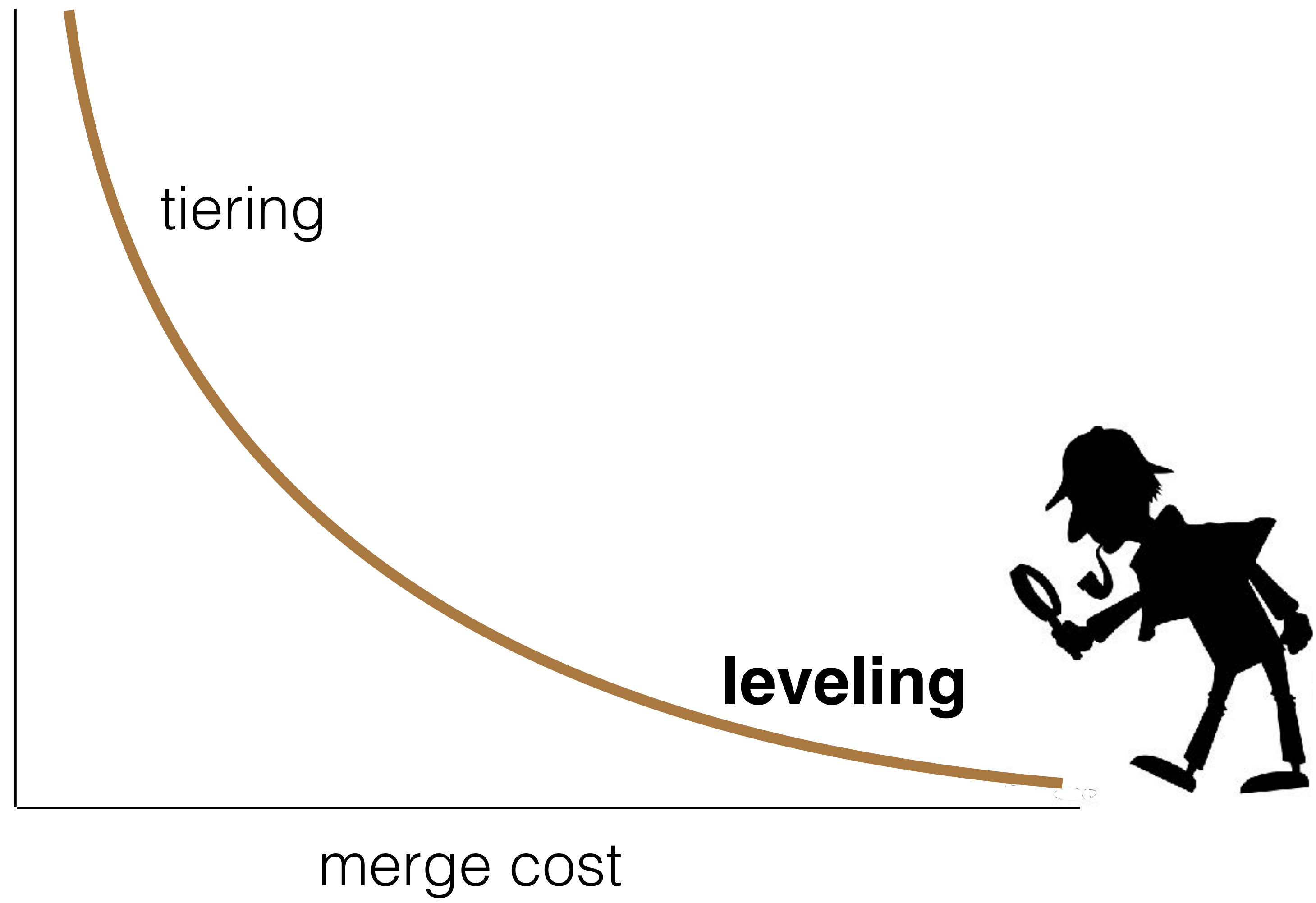
lookup cost



lookup cost



lookup cost



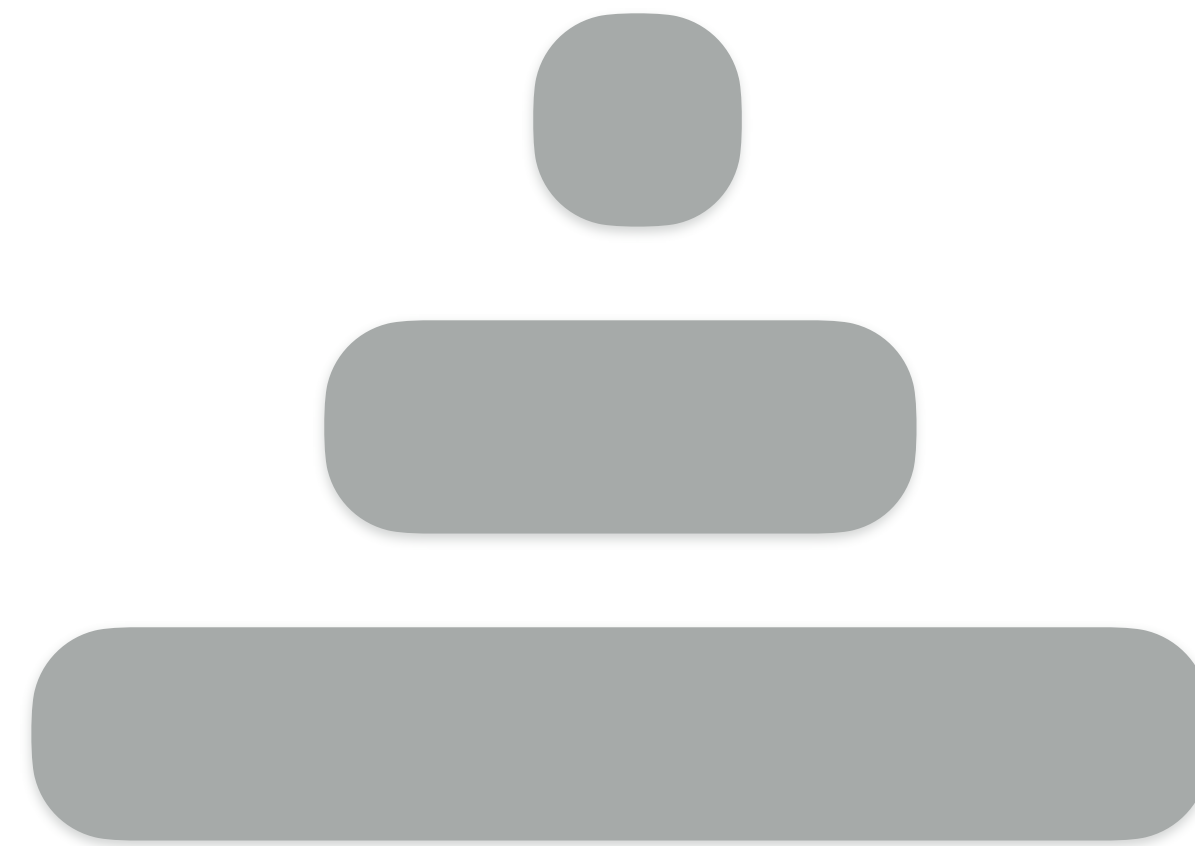
# I/O overheads with leveling

●  
point

→  
long range

→  
short range

↪  
merging

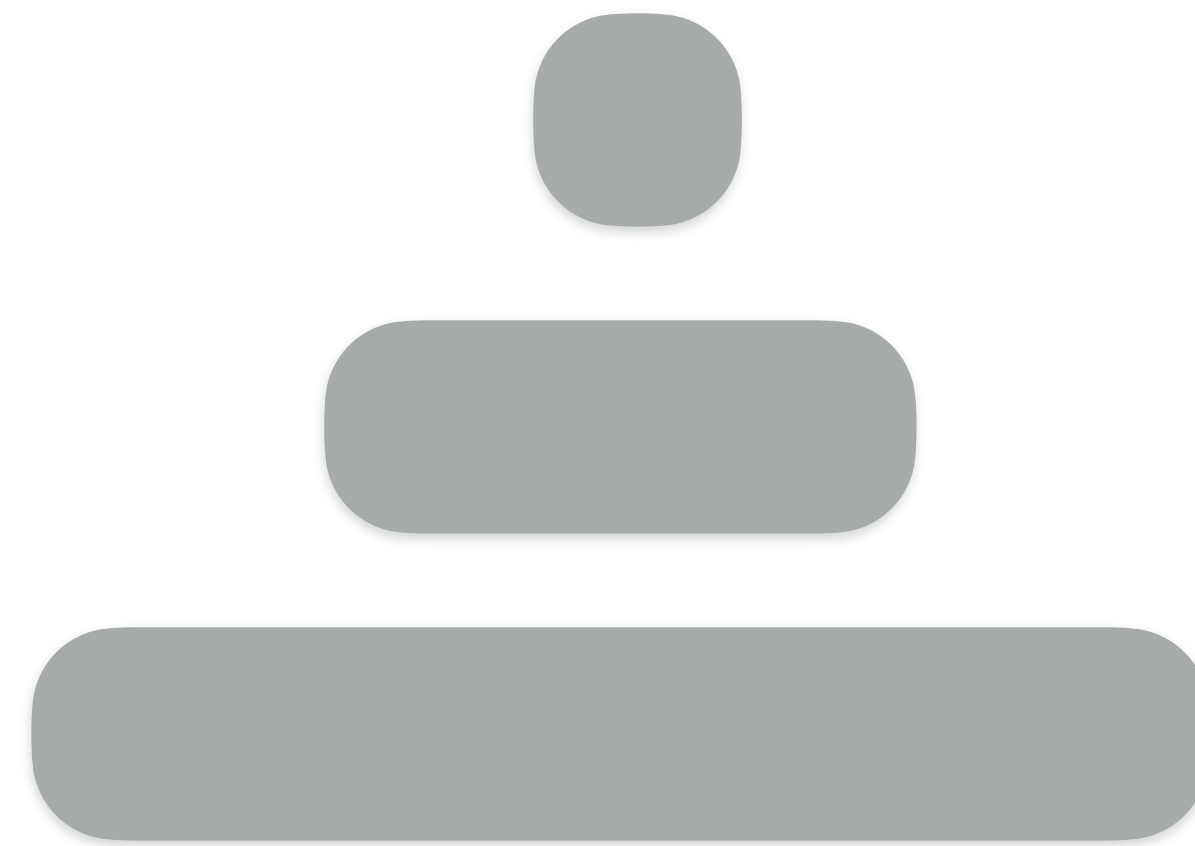


●  
point

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

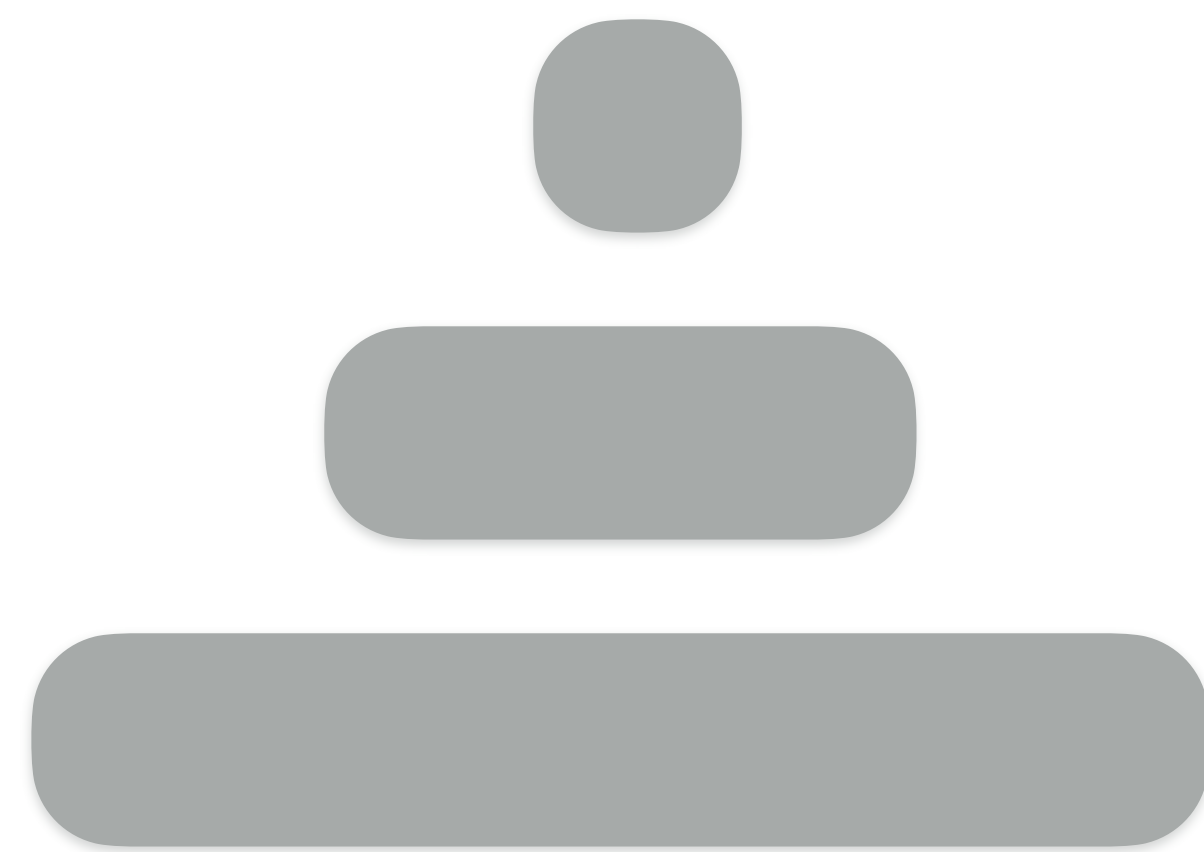
largest level

$O(e^{-x})$

  
**long range**

 short range

  
merging





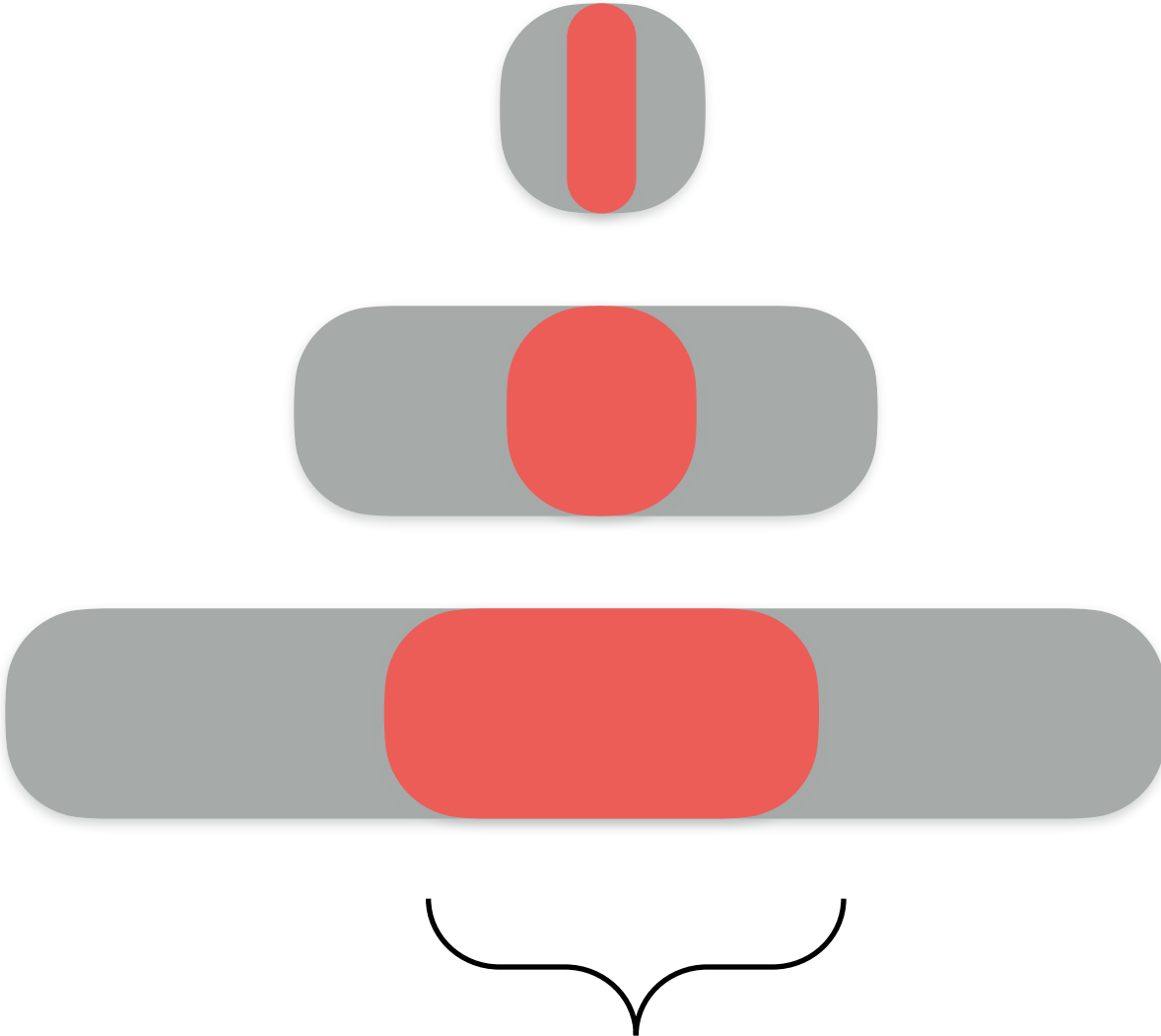
long range →

target range

$O(s/R^2)$

$O(s/R)$

$O(s)$



target key range

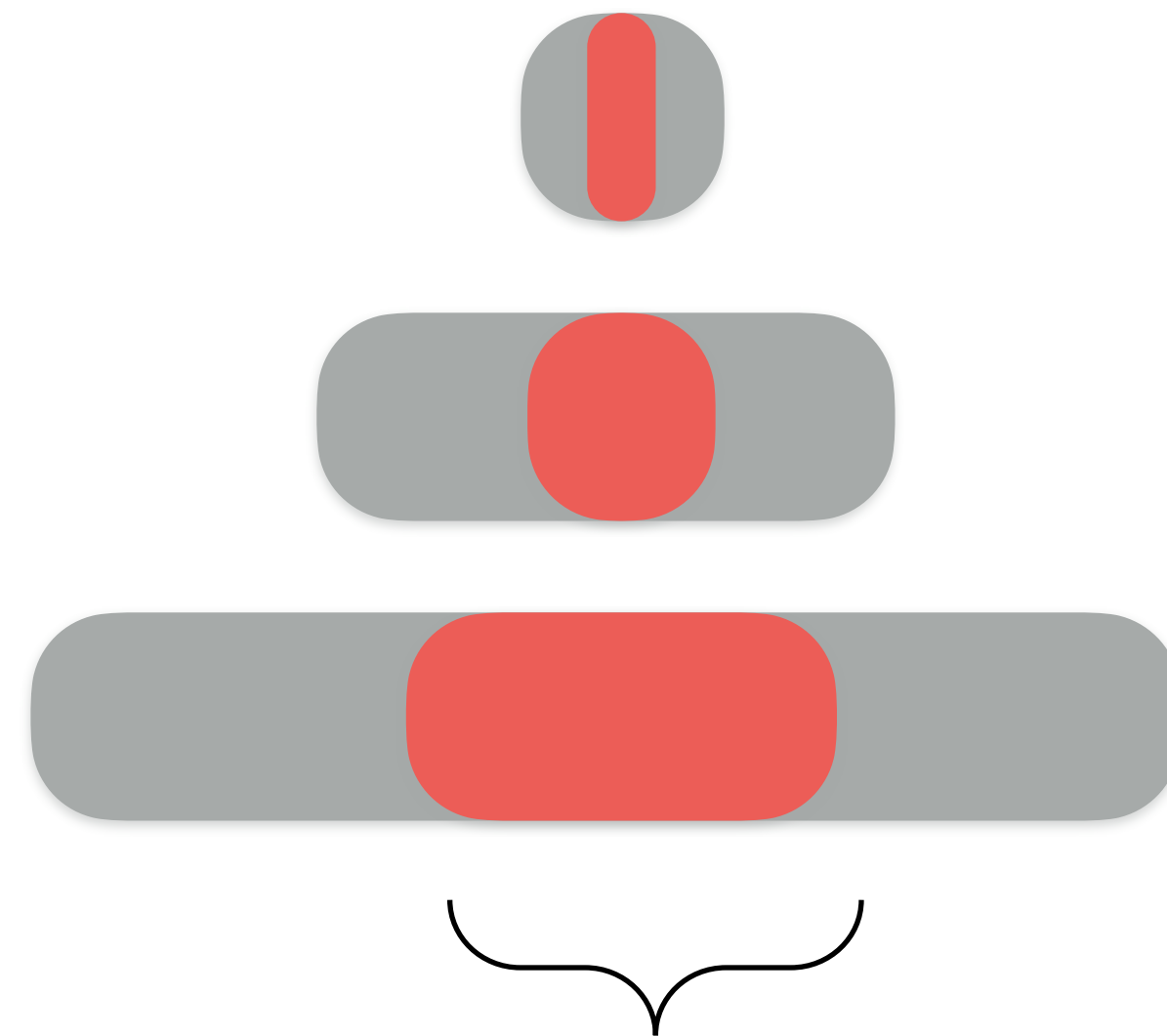
**long range** 

target range

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

$$O(s/R)$$

$$O(s)$$



**largest level**

**target key range**

●  
point

largest level

$O(e^{-x})$

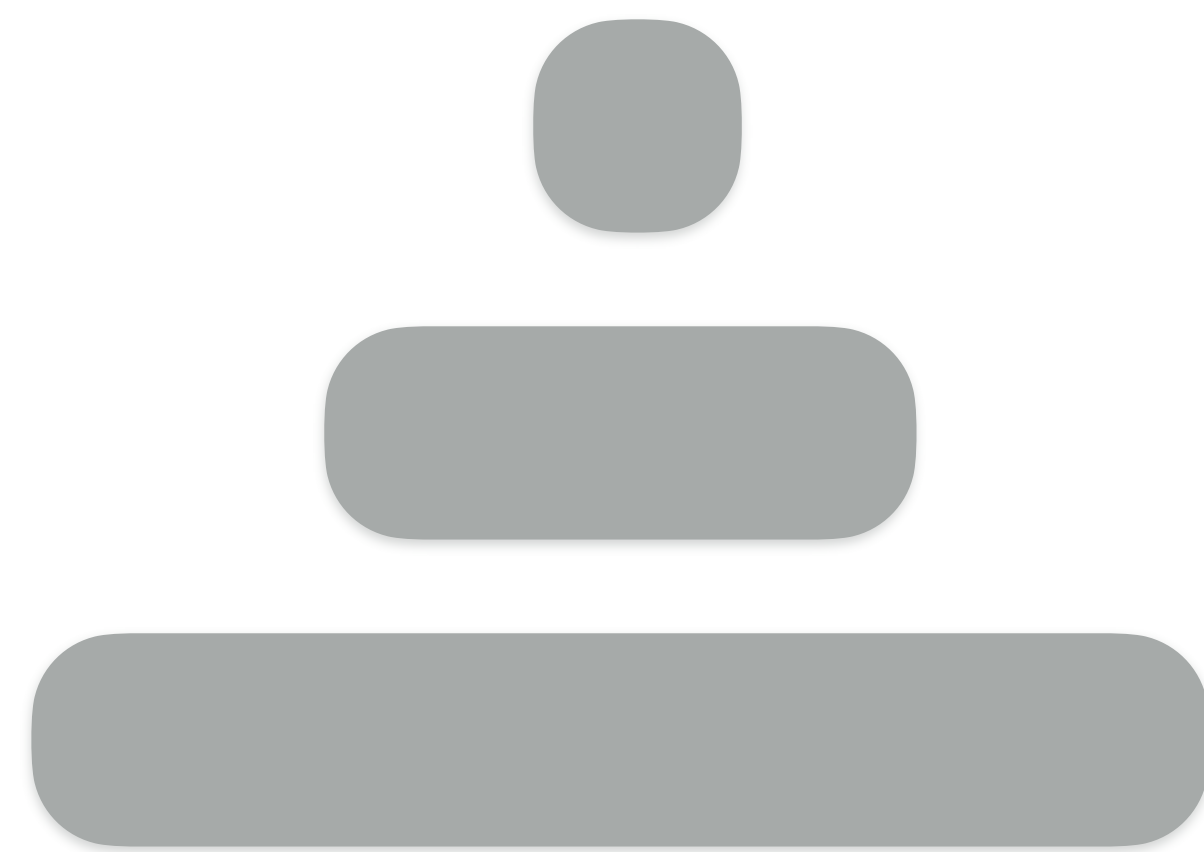
→  
long range

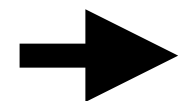
largest level

$O(s)$

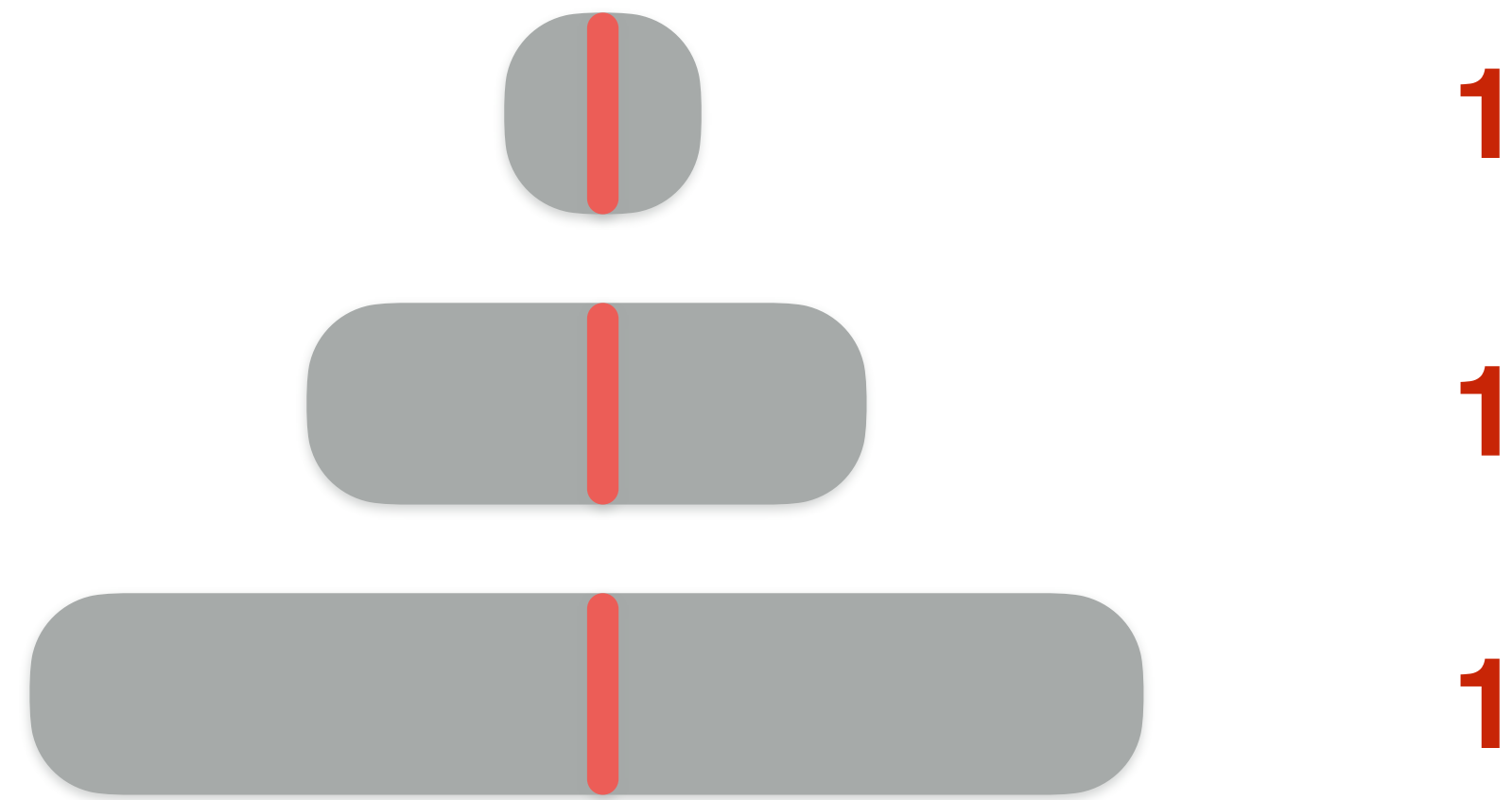
→  
**short range**

↶  
merging



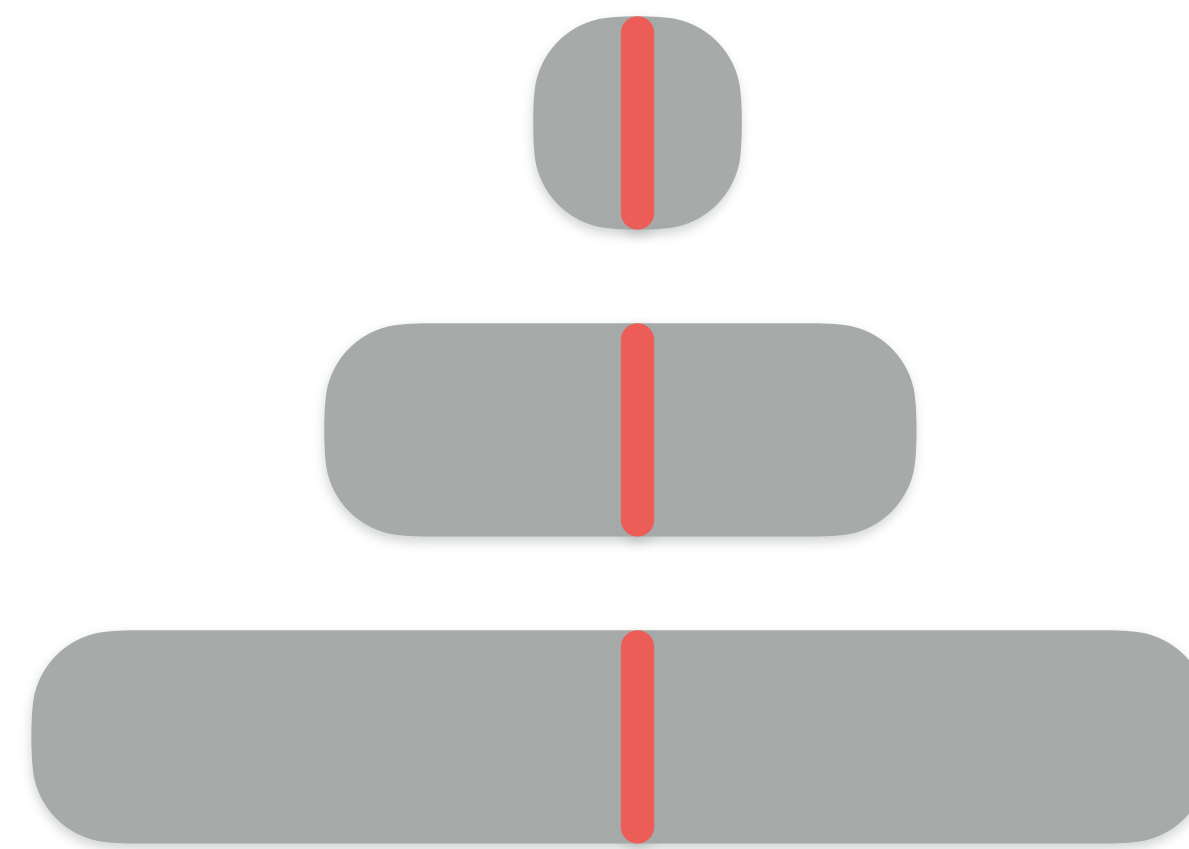
short  range

target range



short range →

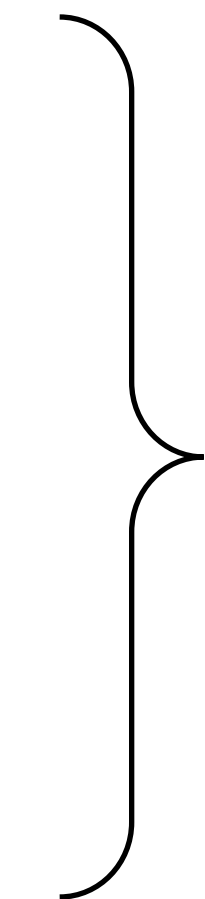
target range



1

1

1



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

●  
point

largest level

$O(e^{-x})$

→  
long range

largest level

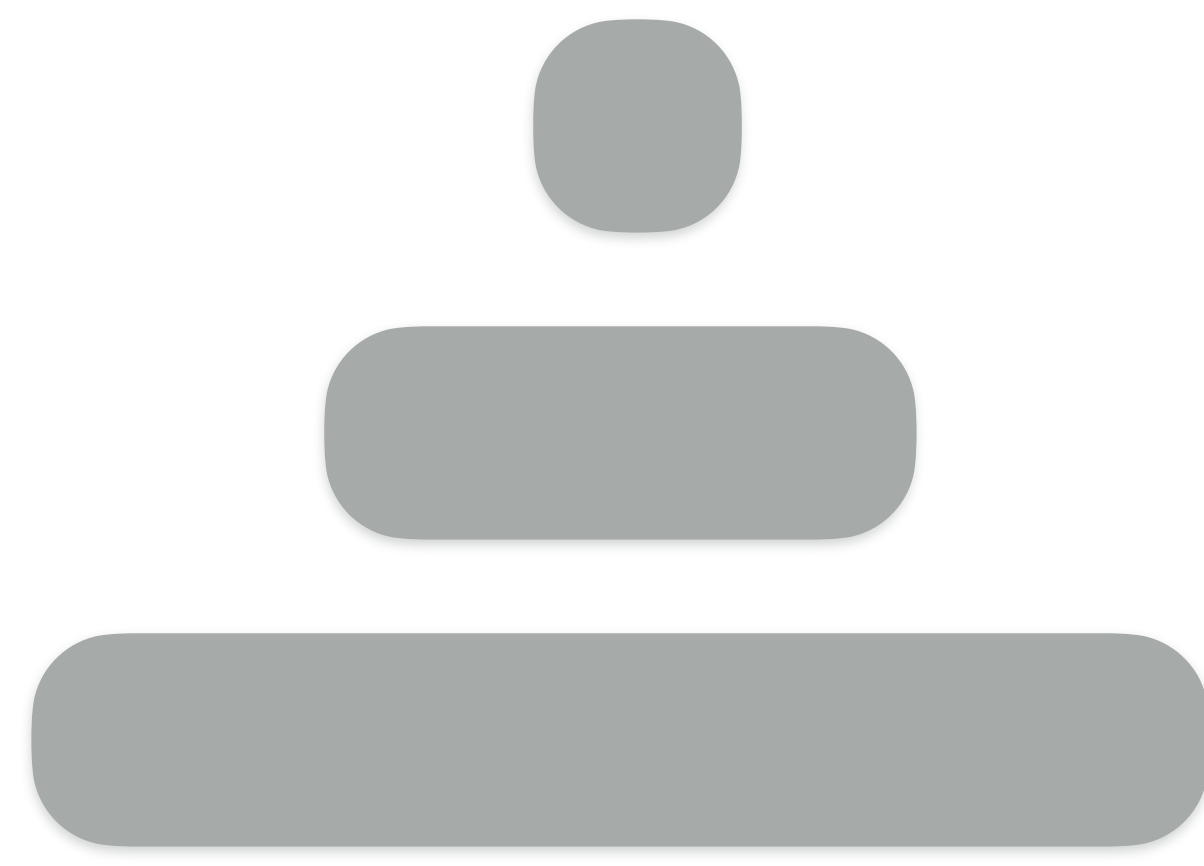
$O(s)$

→  
short range

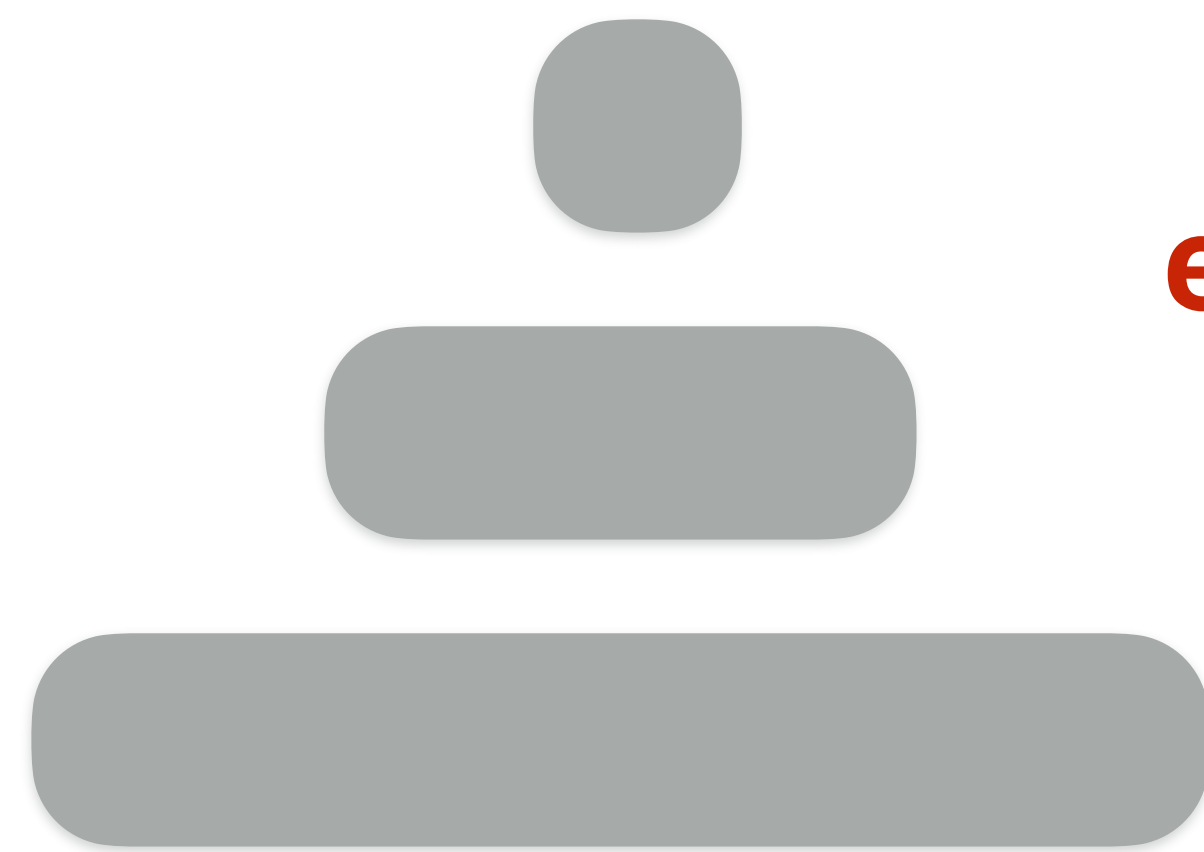
all levels

$O(\log_R(N))$

↶  
**merging**



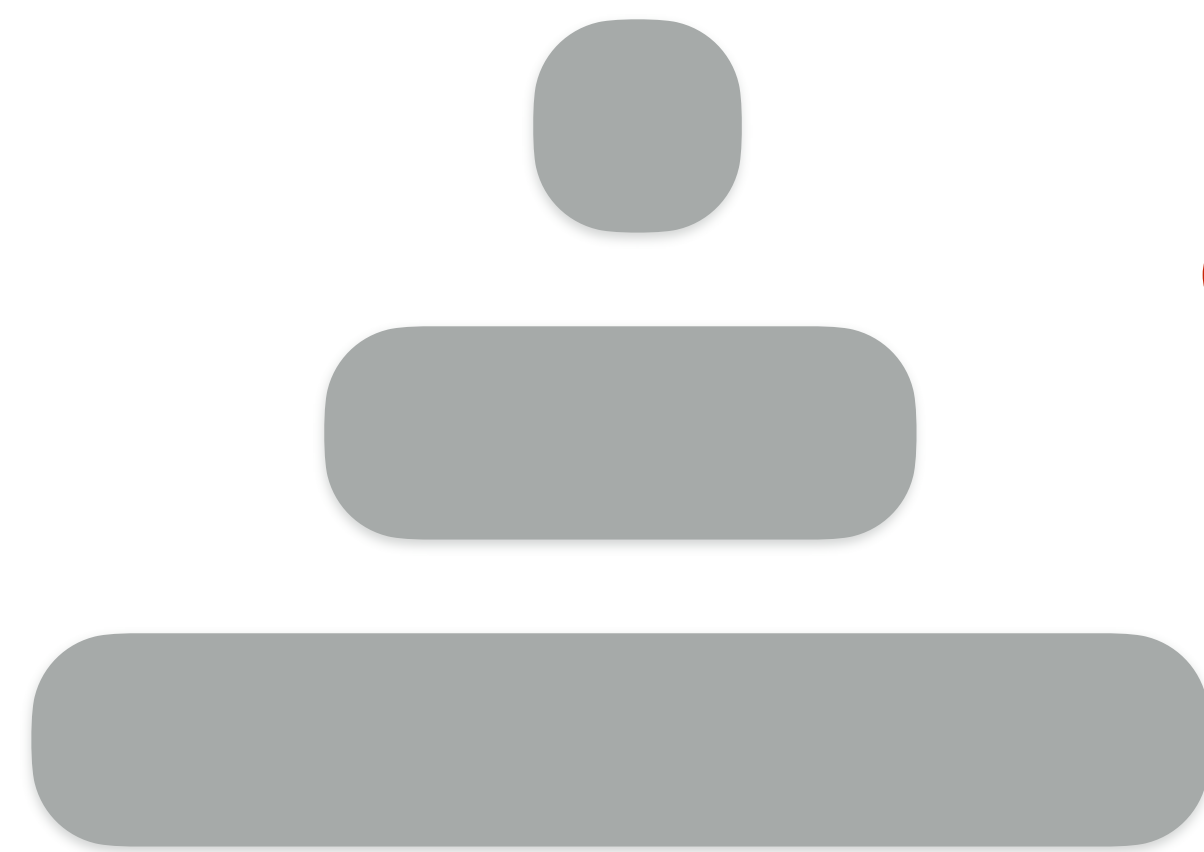
  
**merging**



**exponentially  
more work**



  
**merging**




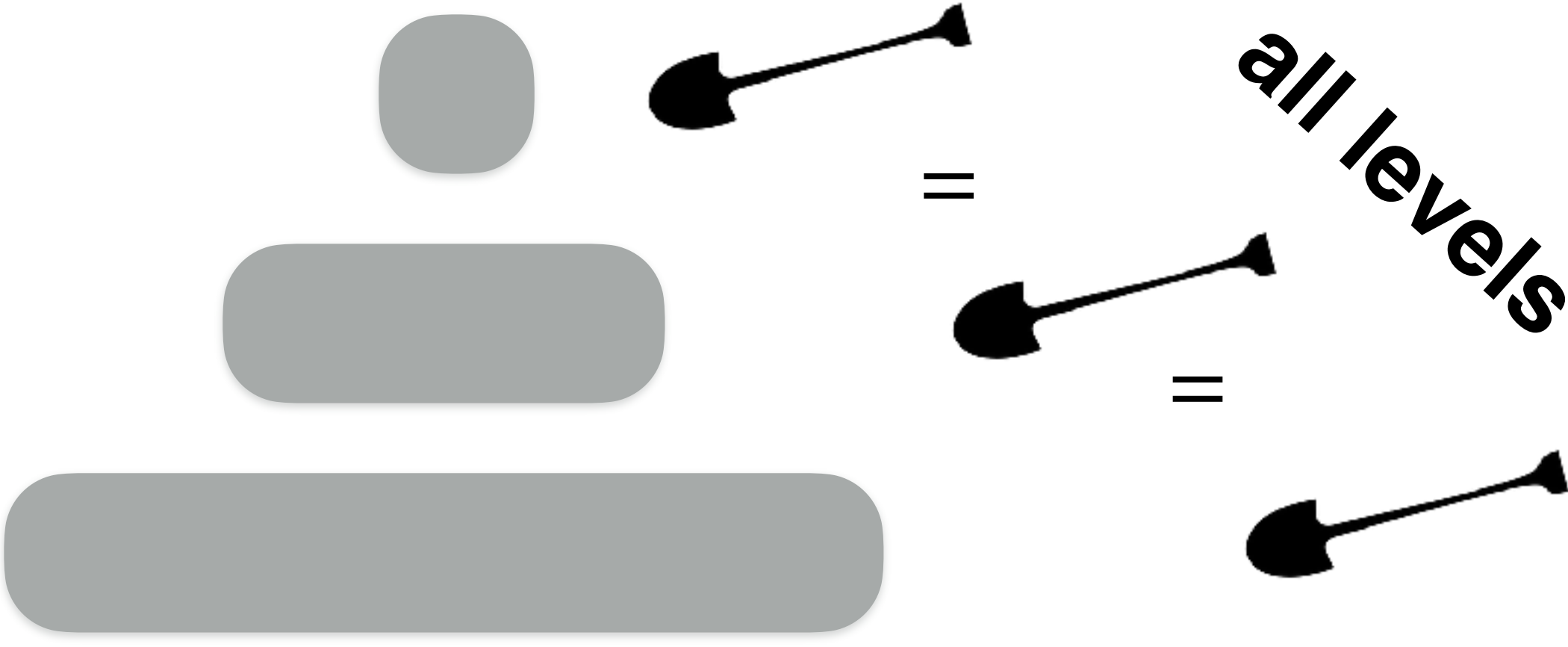
exponentially  
more work



**exponentially  
less frequent**



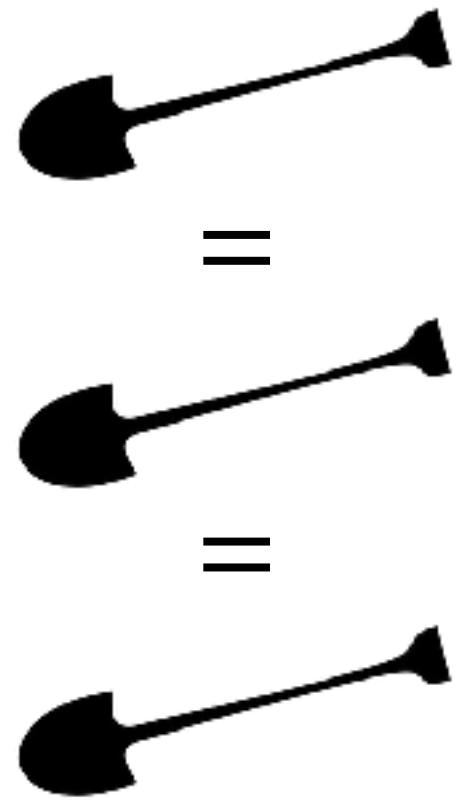
  
**merging**



**all levels**

**more work**  
↓  
**less frequent**

  
**merging**



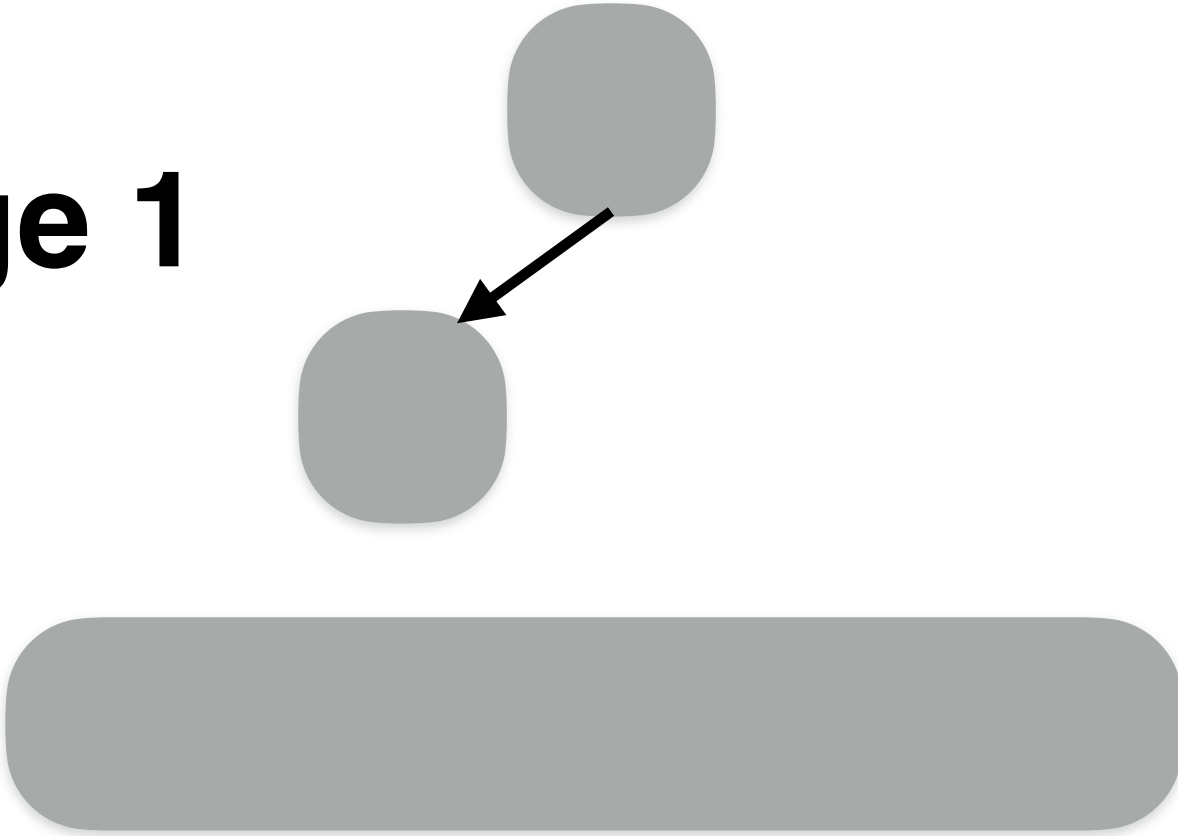
all levels



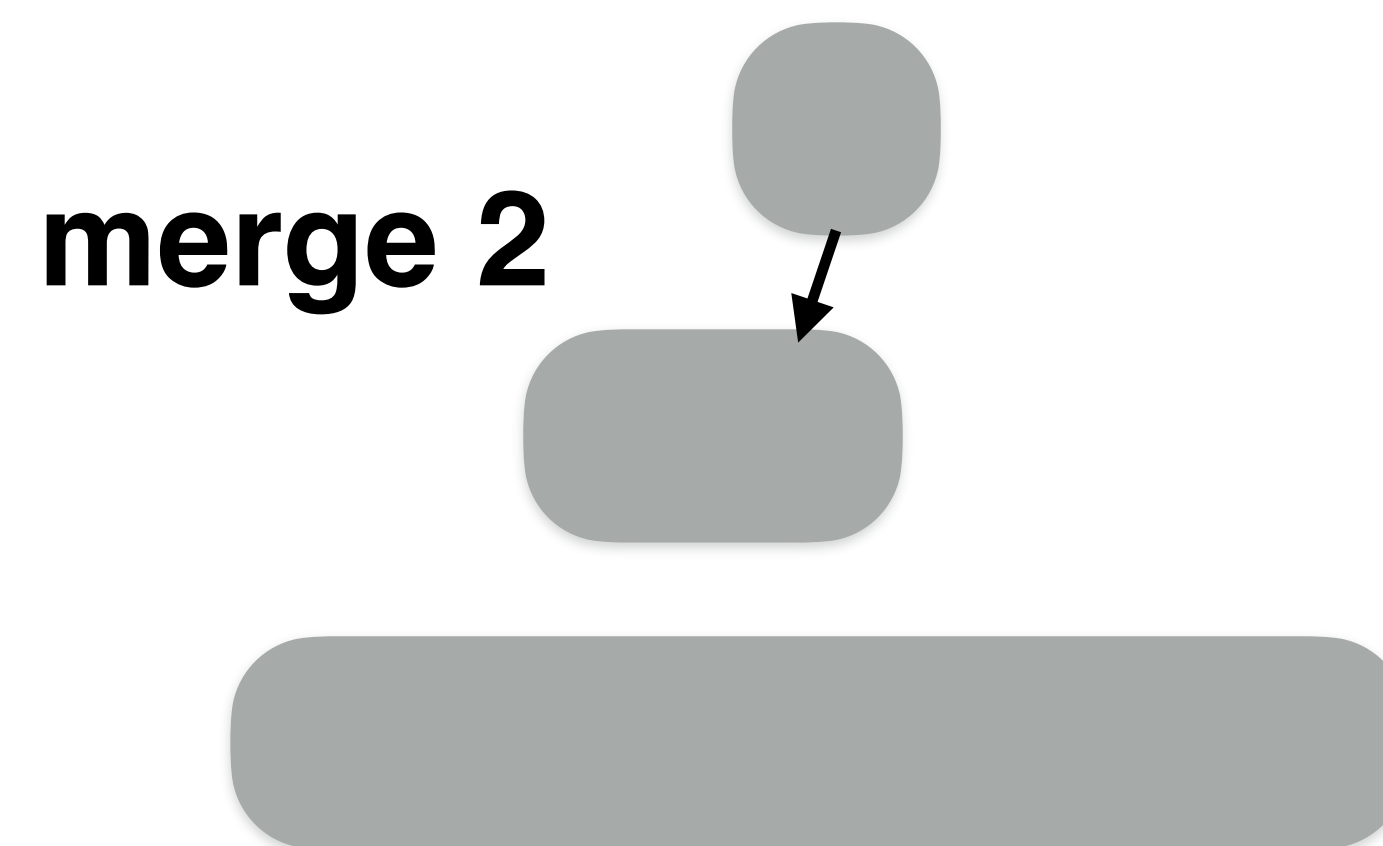


**merging**

**merge 1**

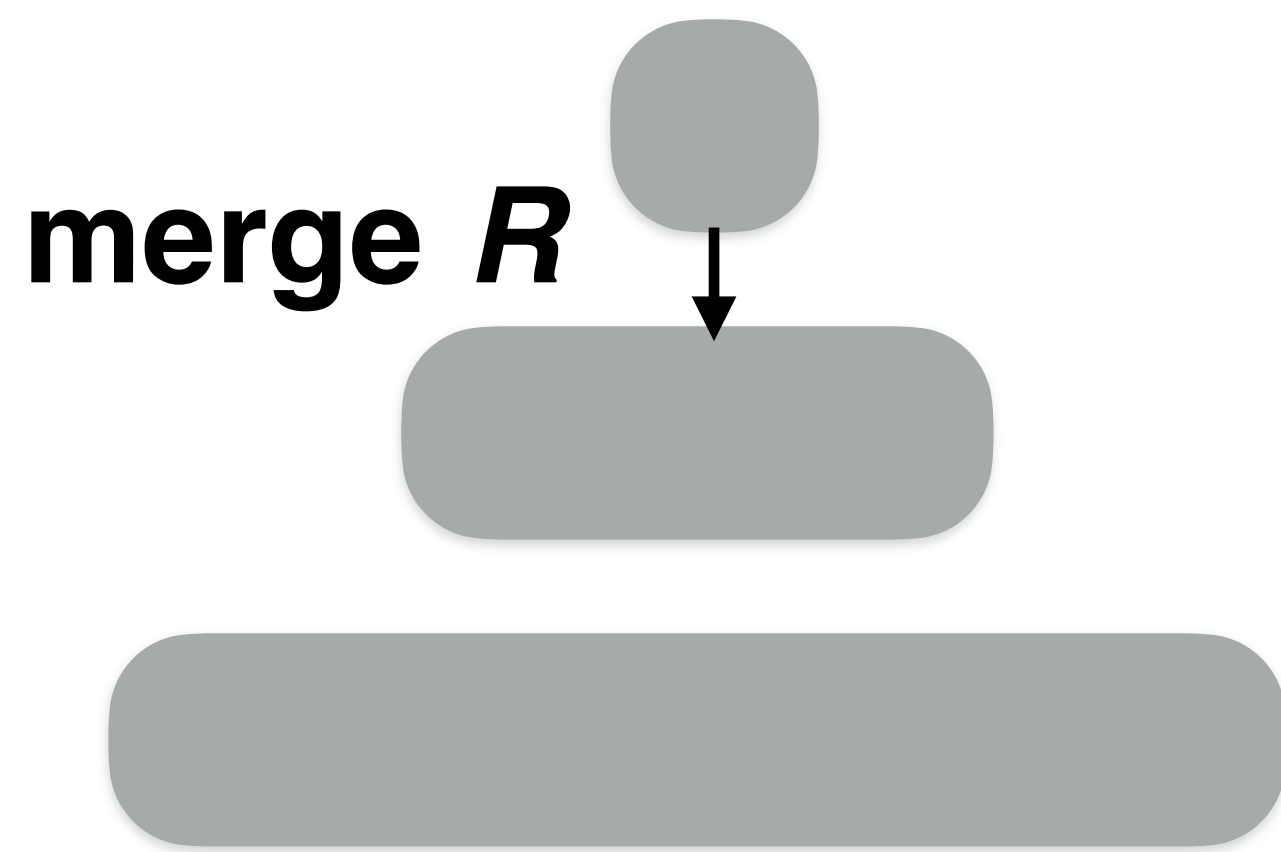


  
**merging**



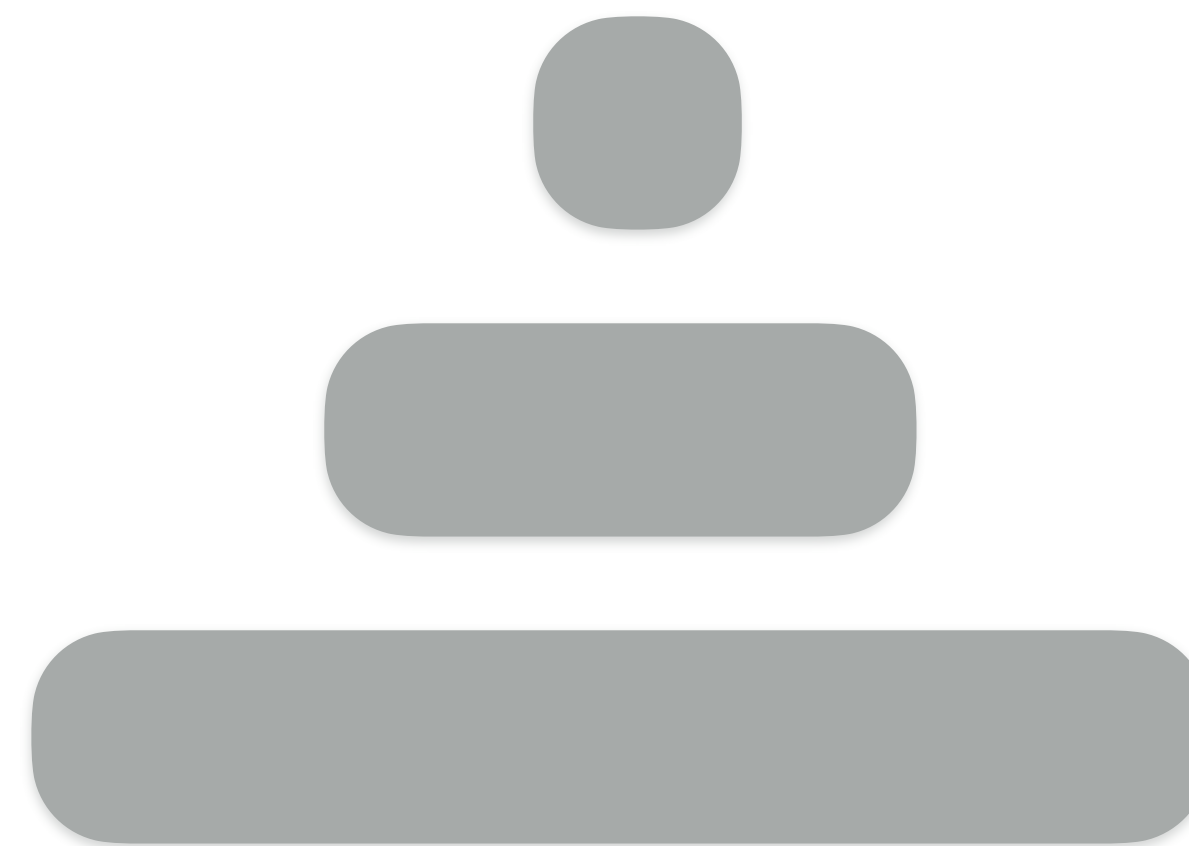


**merging**





write-amplification



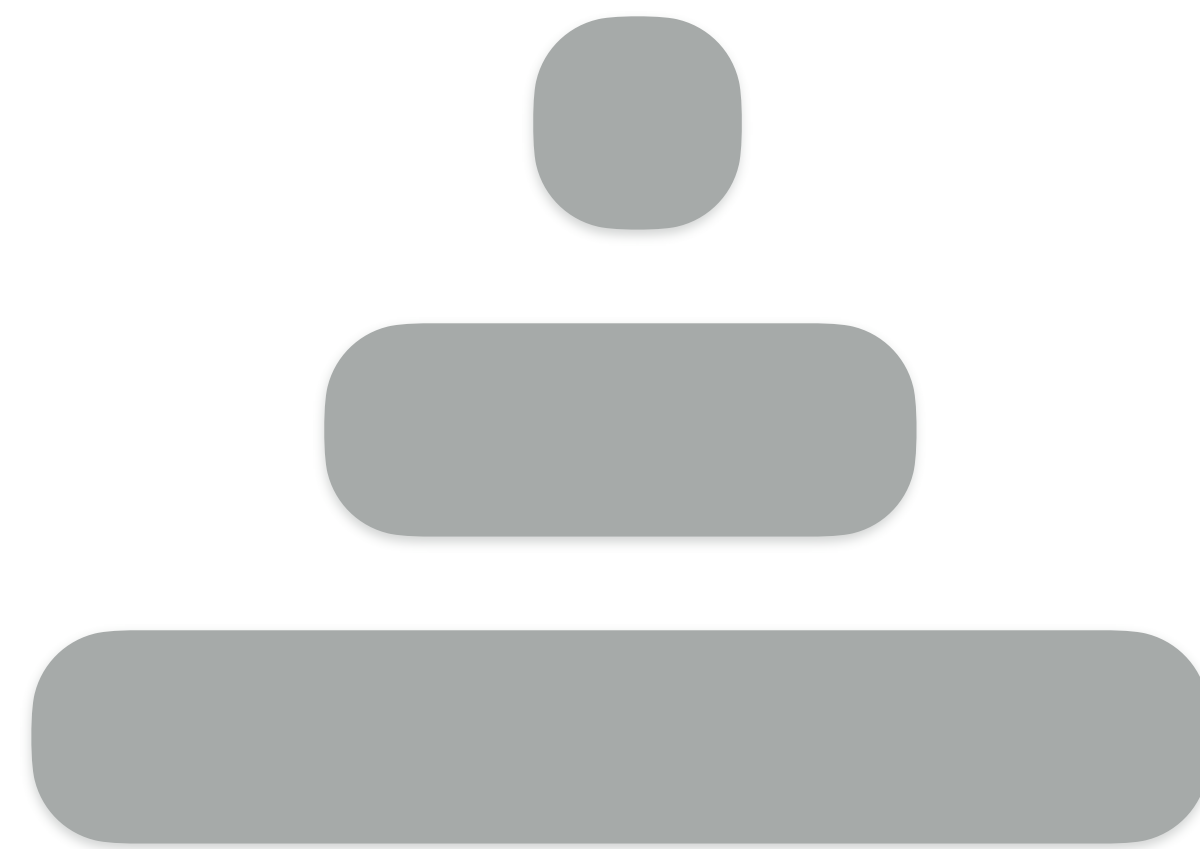
$O(R)$

$O(R)$

$O(R)$



**merging**



$O(R)$

$O(R)$

$O(R)$

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

●  
point

→  
long range

→  
short range

↵  
merging

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

$O(R)$

+

+

+

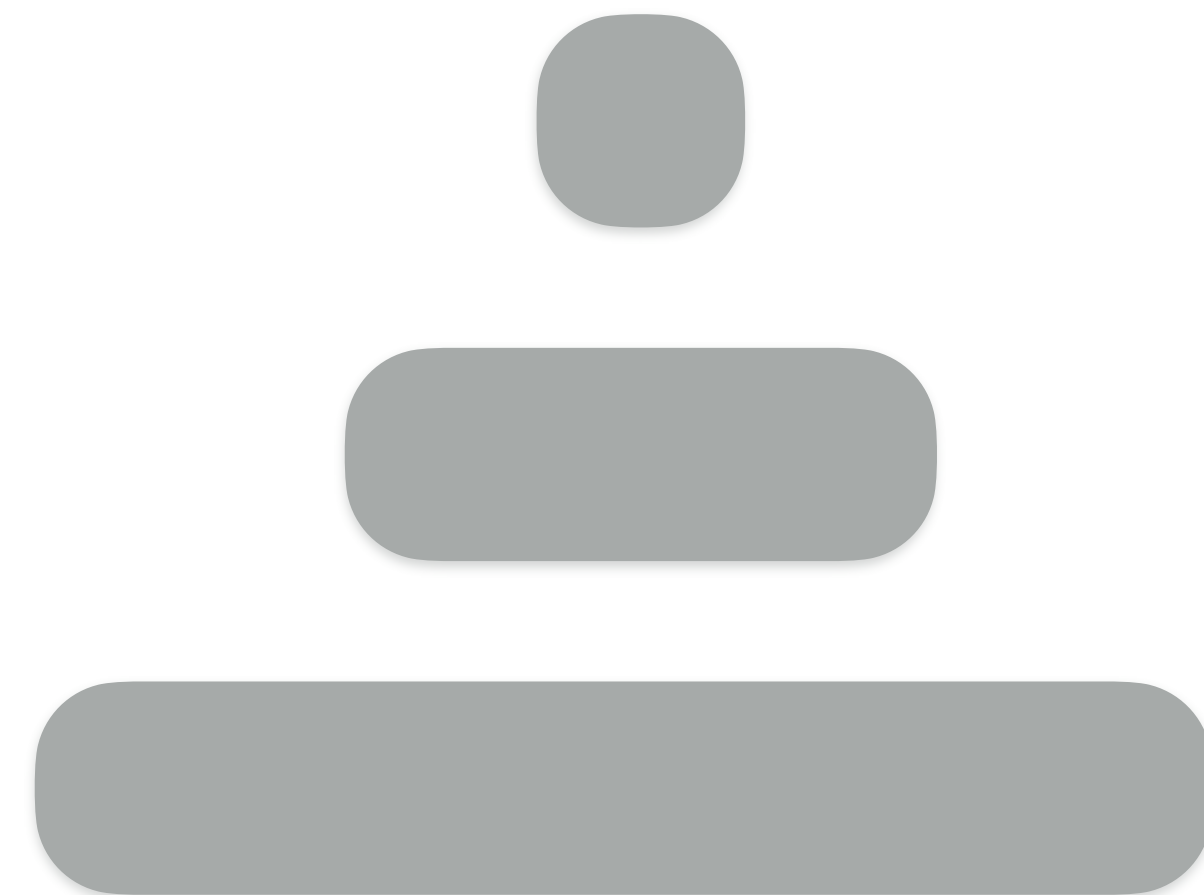
+

$O(e^{-x})$

$O(s)$

1

$O(R)$





●  
point

→  
long range

↵  
merging

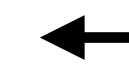
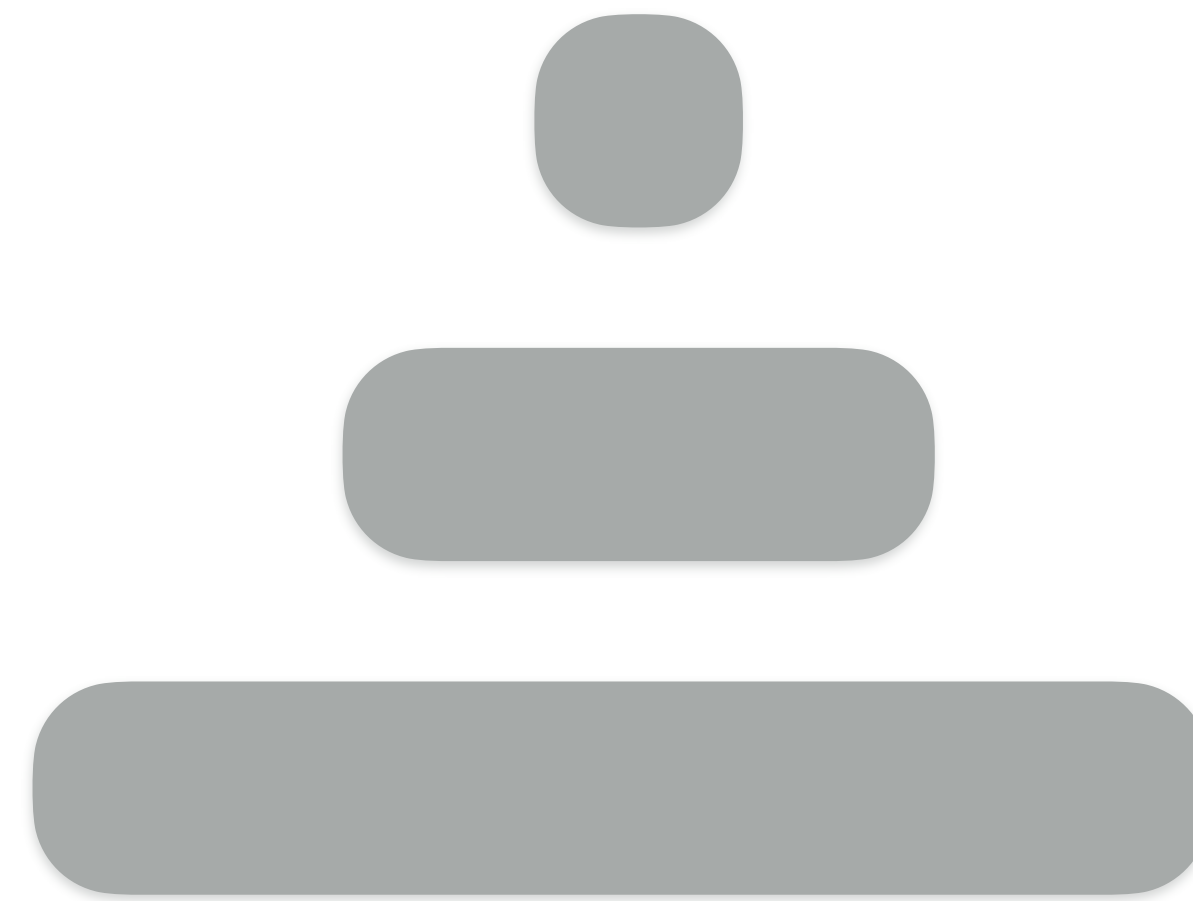
**largest level**

**largest level**

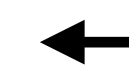
**all levels**

$O(e^{-x})$

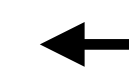
$O(s)$



$O(R)$



$O(R)$



$O(R)$

●  
point

→  
long range

↵  
merging

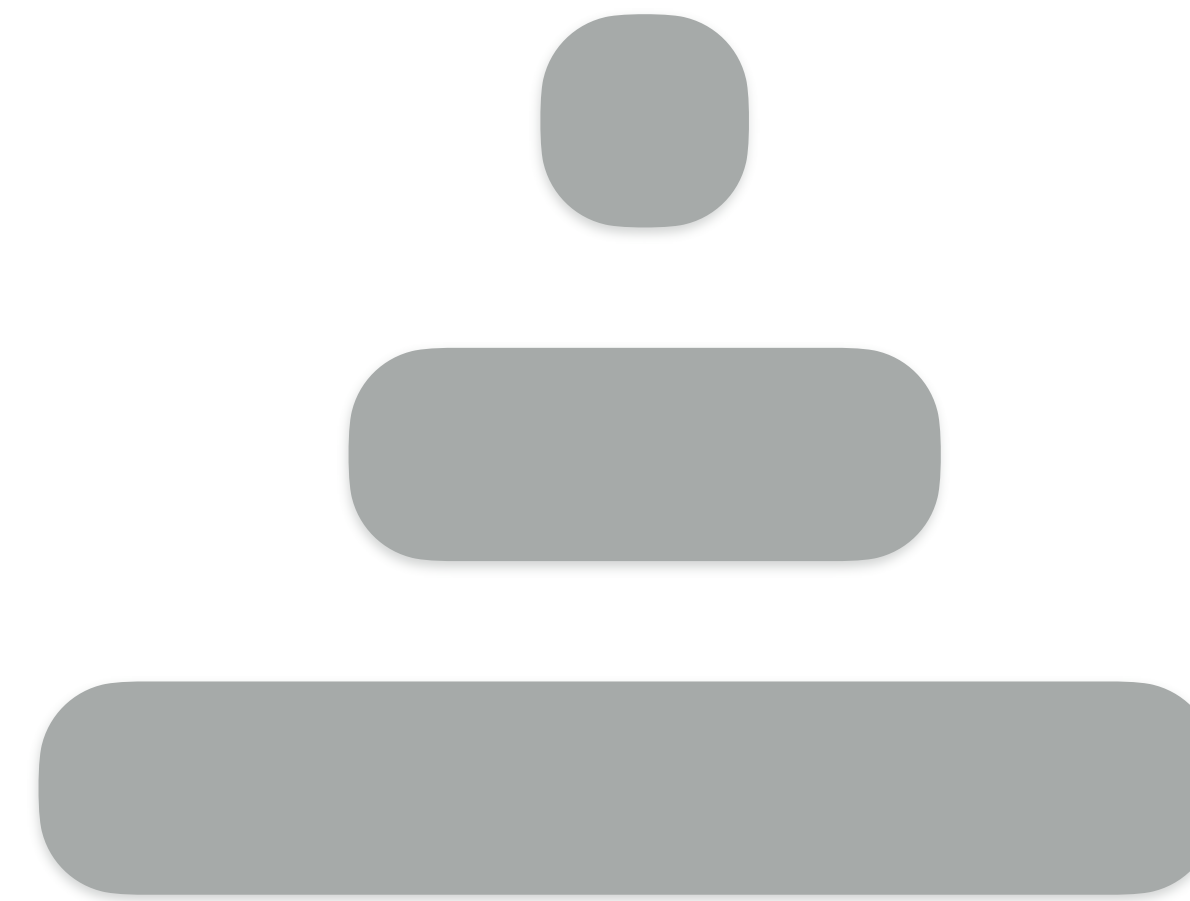
largest level

largest level

all levels

$O(e^{-x})$

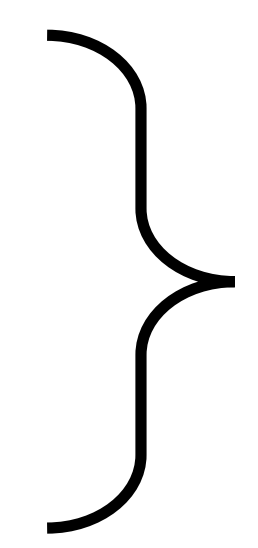
$O(s)$



←  $O(R)$

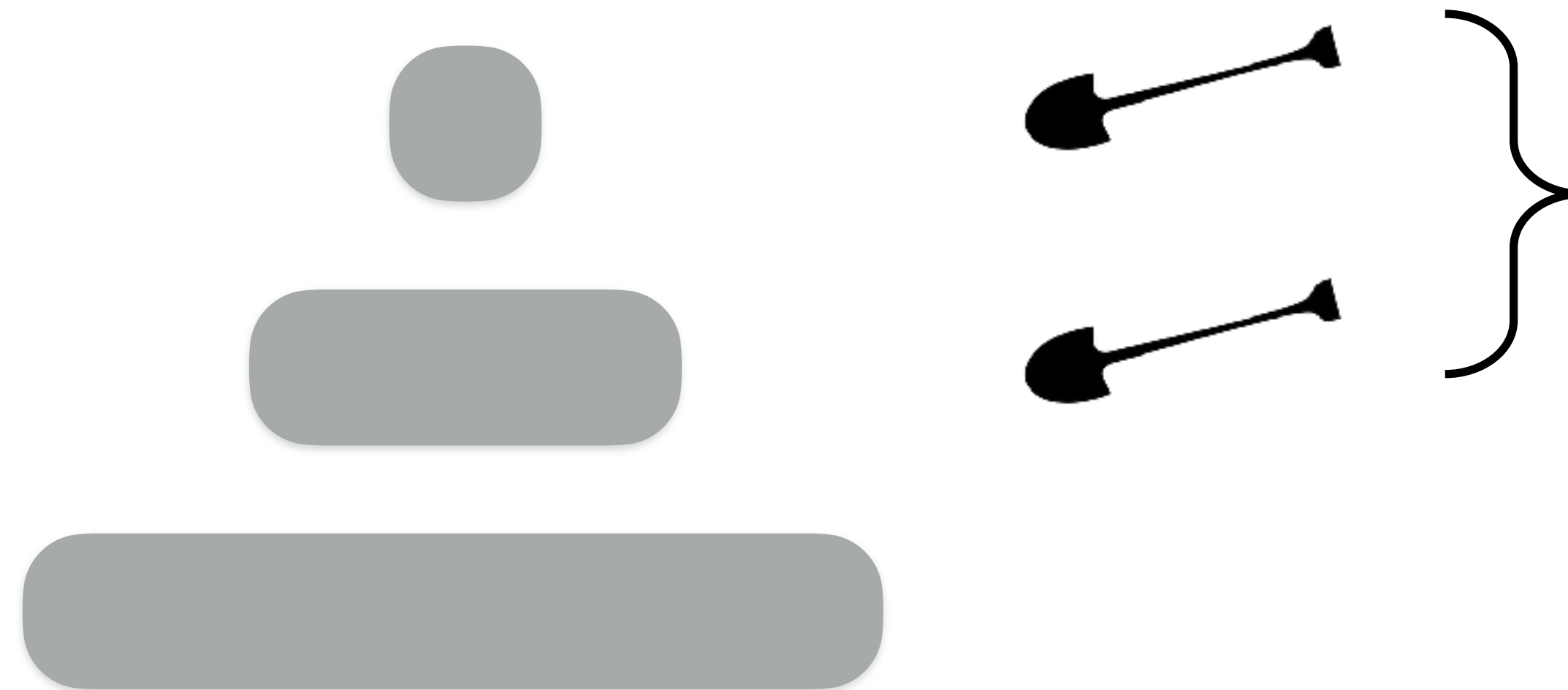
←  $O(R)$

←  $O(R)$



**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**



**D**ostoevsky: **S**pace-**T**ime **O**ptimized **E**volvab**l**e **S**calable **K**ey-Value Store



**very write-optimized**



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

**Tiering**  
merge-optimized

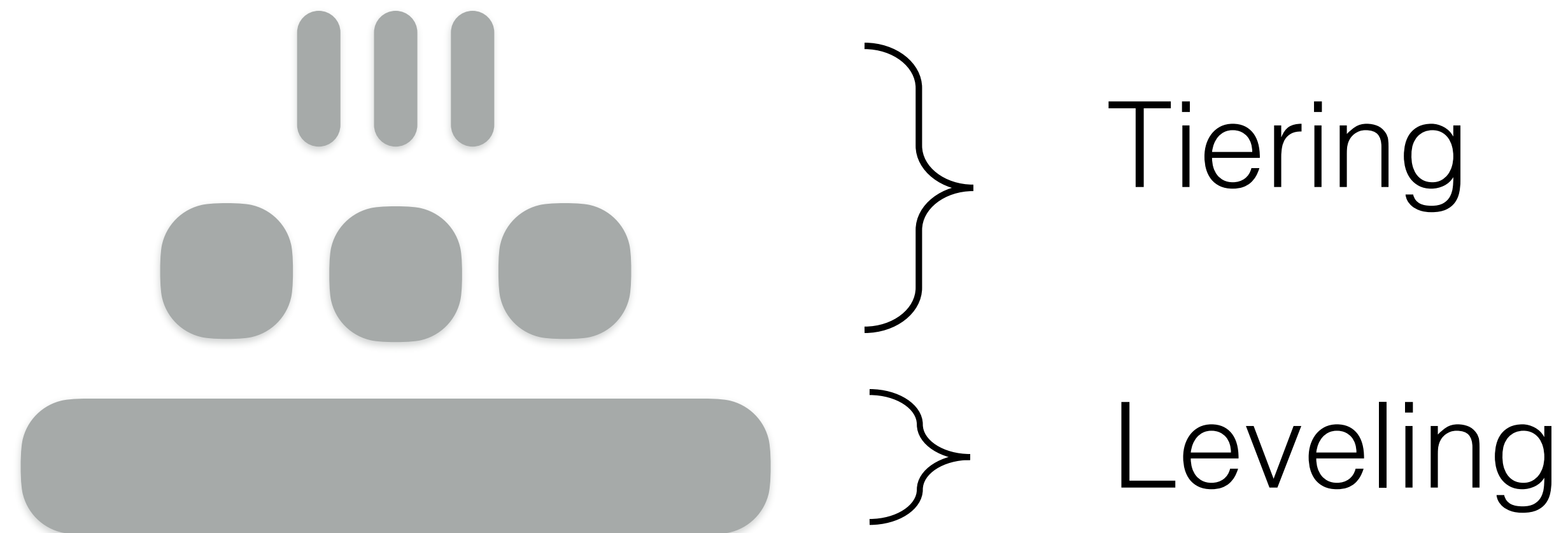
**Leveling**  
lookup-optimized

Tiering  
merge-optimized

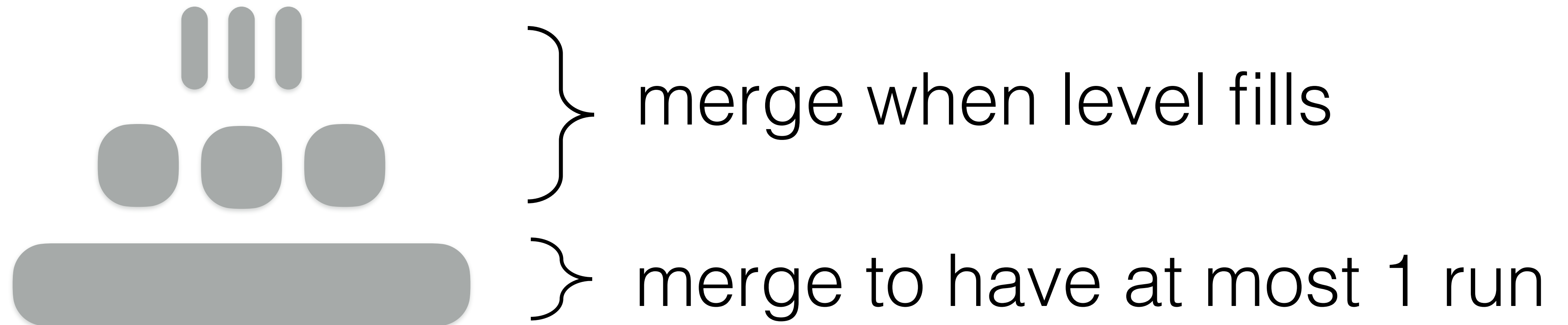
**Lazy Leveling**  
mixed-workloads

Leveling  
lookup-optimized

# Lazy Leveling



# Lazy Leveling



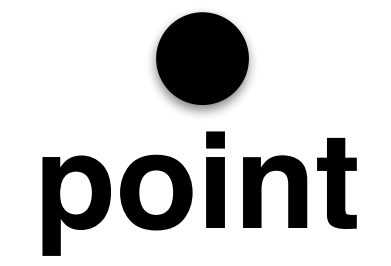
●  
point

→  
long range

→  
short range

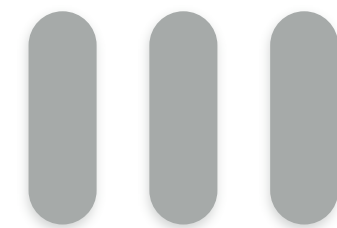
↶  
merging



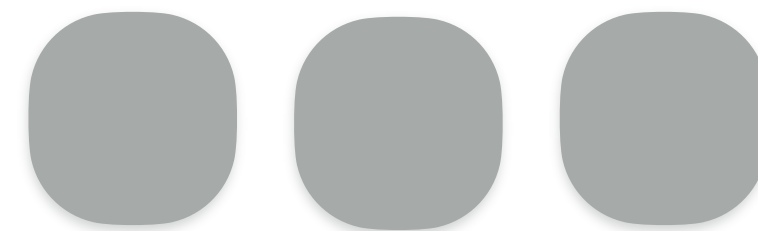


## false positive rates

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



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



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



●  
point

false positive rates

exponentially  
decreasing

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



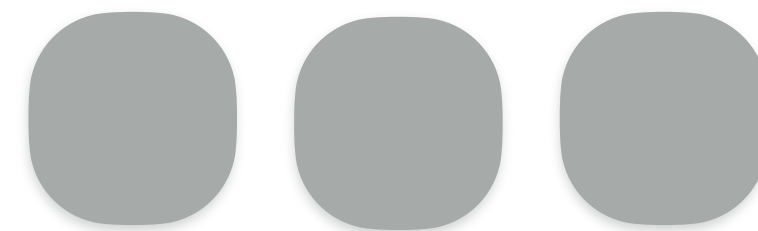
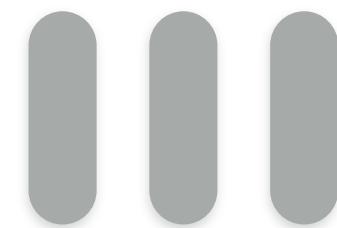


false positive rates

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

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

→  $O(e^{-x})$



**largest level**

●  
**point**

●  
point

$O(e^{-x})$

→  
**long range**

→  
short range

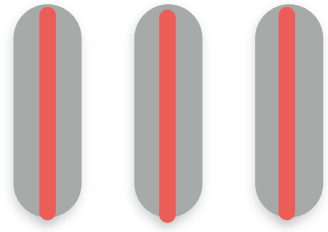
↪  
merging



→  
long range

target range

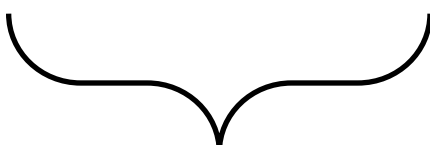
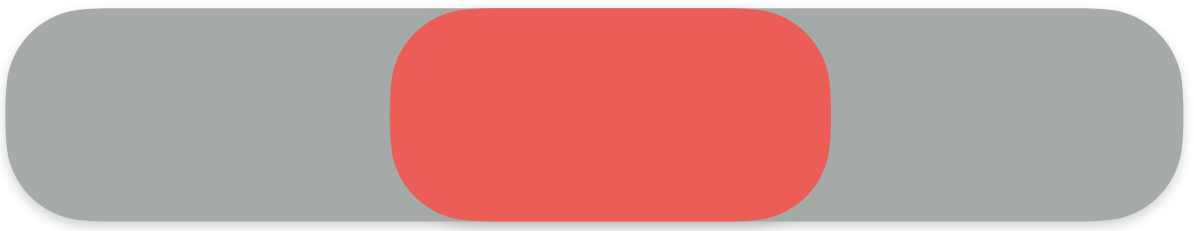
$O(s/R^2)$



$O(s/R)$



$O(s)$



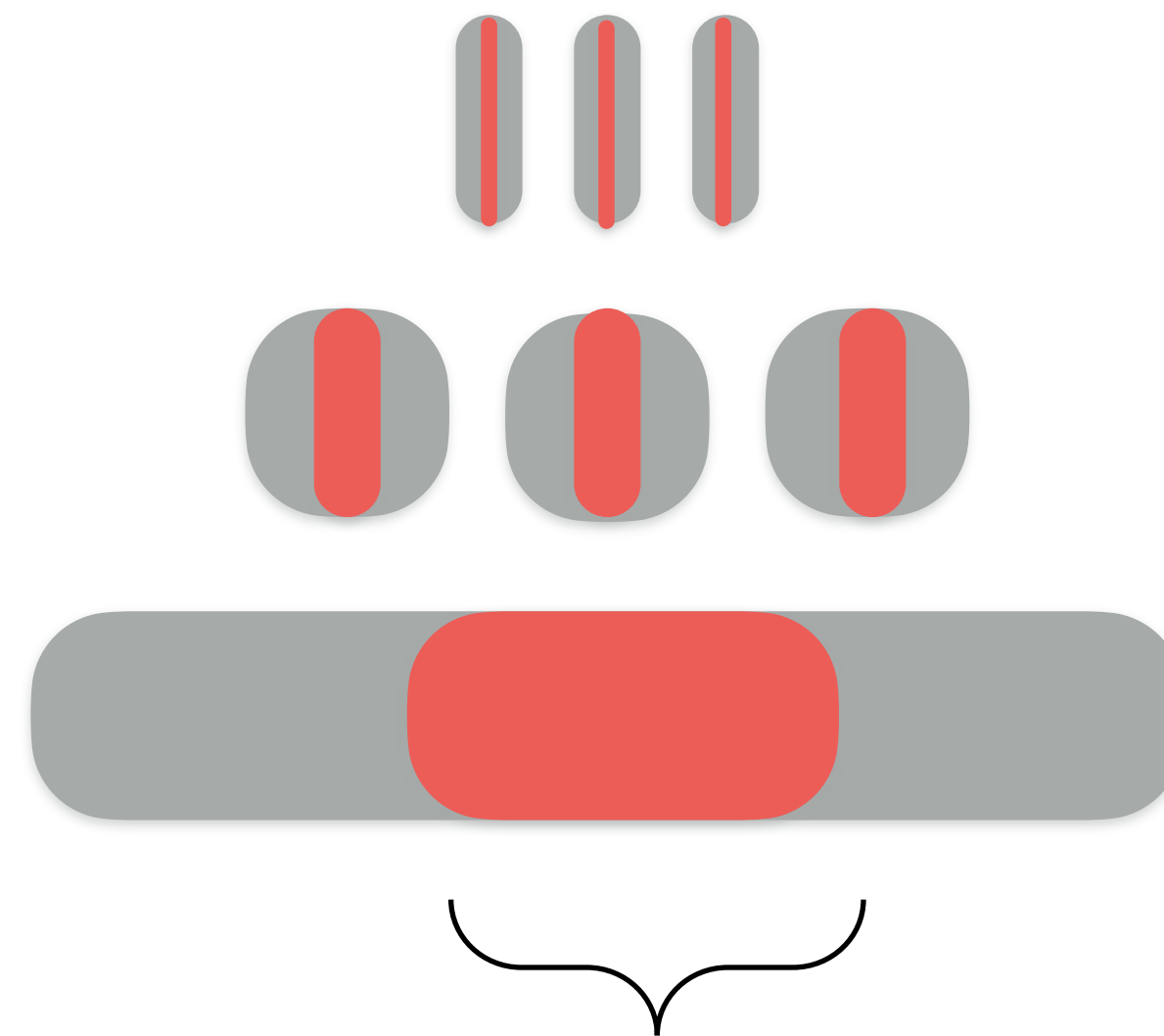
target key range

target range

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

$$O(s/R)$$

$$O(s)$$



largest level

  
long range

target key range

●  
point

$O(e^{-x})$

→  
long range

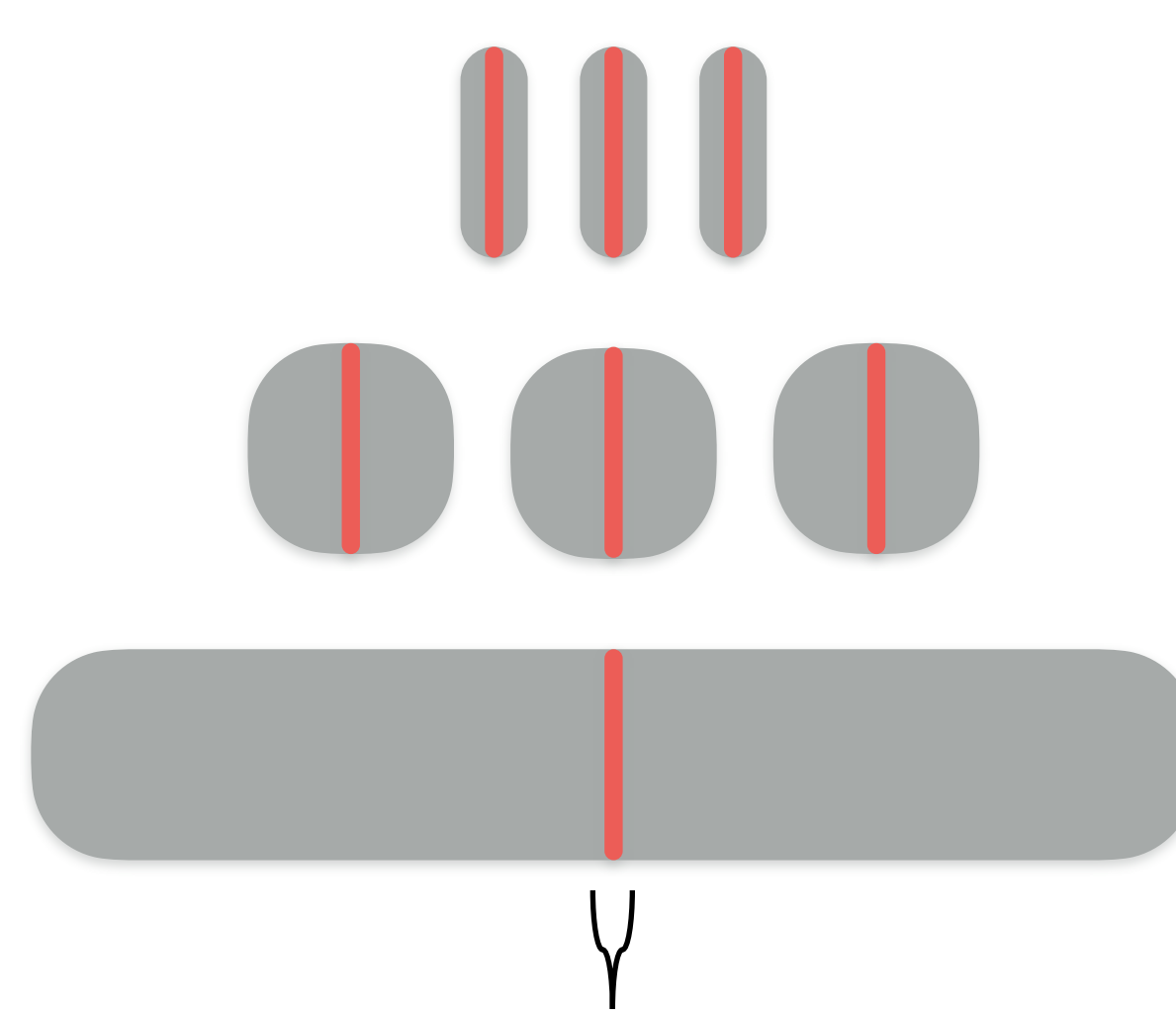
$O(s)$

→  
**short range**

↪  
merging



→  
**short range**



$O(R)$

$O(R)$

**1**

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

target key range

●  
point

$O(e^{-x})$

→  
long range


$O(s)$

→  
short range

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

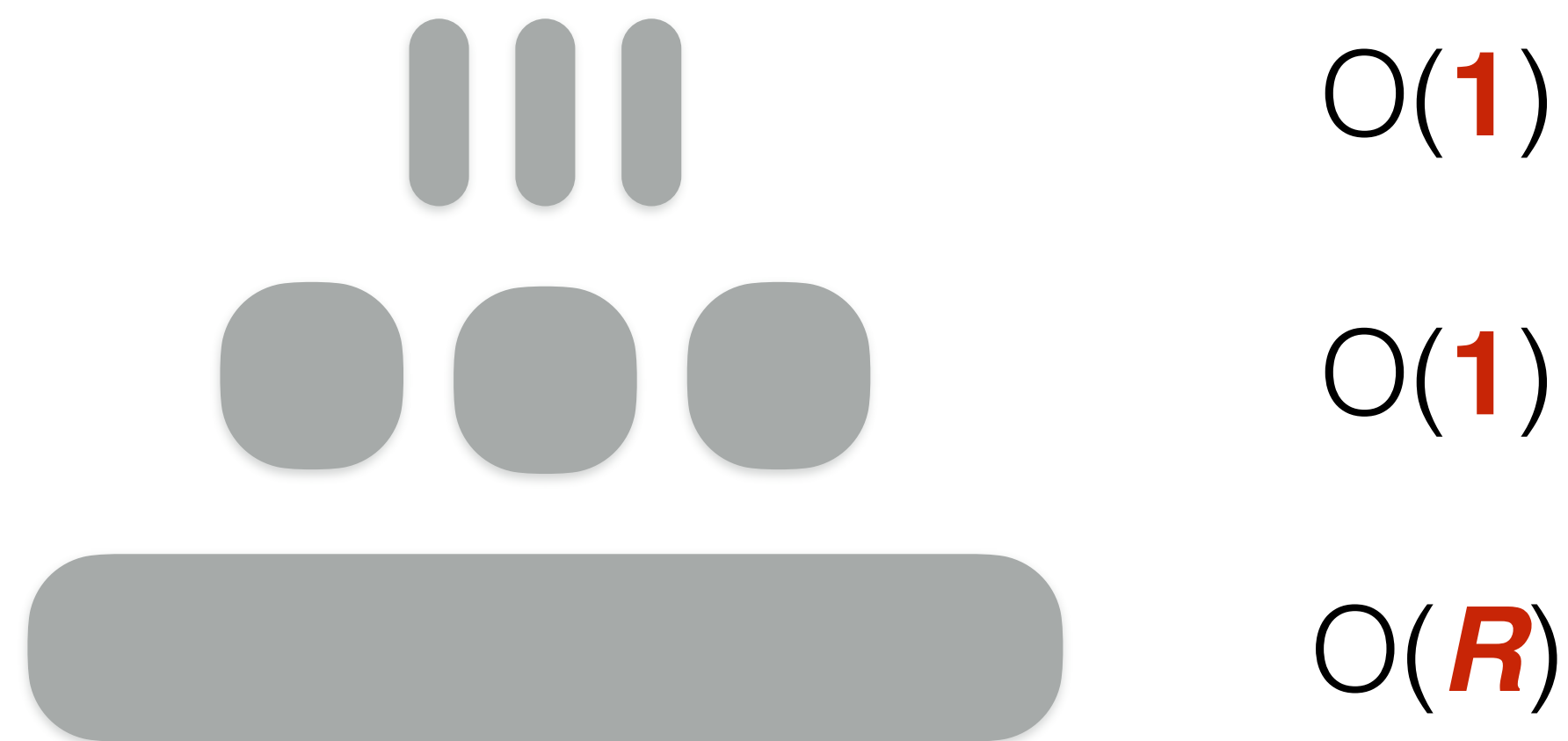
↪  
merging






**merging**

write-amplification





  
merging

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

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

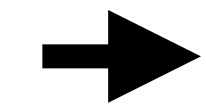
↪  
merging


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



  
point

  
long range

  
short range

  
merging

**Lazy Leveling**

$O(e^{-x})$

$O(s)$

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

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

=

=

V

Λ

**Leveling**

$O(e^{-x})$

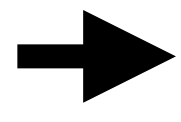
$O(s)$


$O(\log_R(N))$

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

  
point

  
long range

  
short range

  
merging

**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**

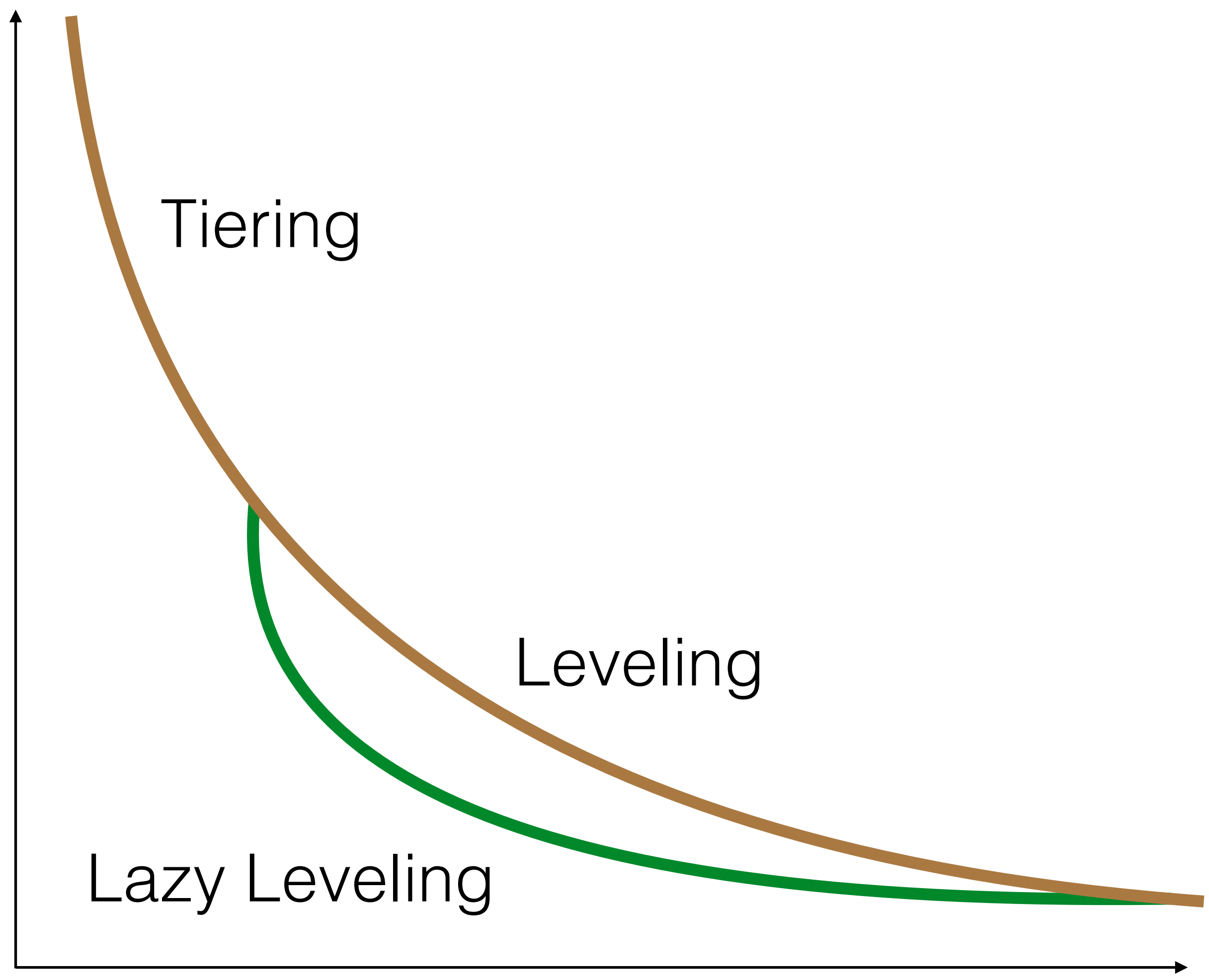
$O(e^{-x})$

$O(s)$

$O(\log_R(N))$

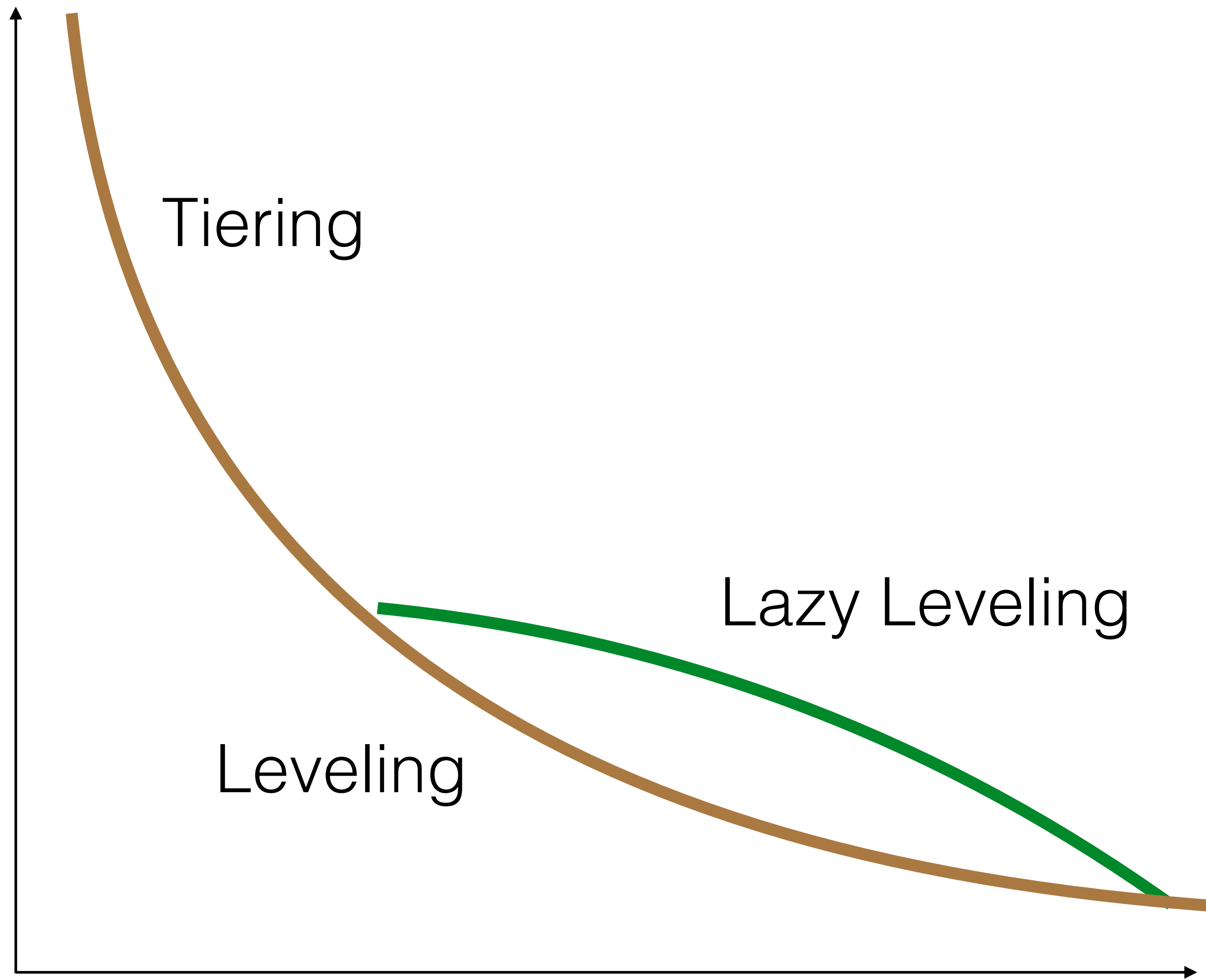
$O(R \cdot \log_R(N))$

●  
point

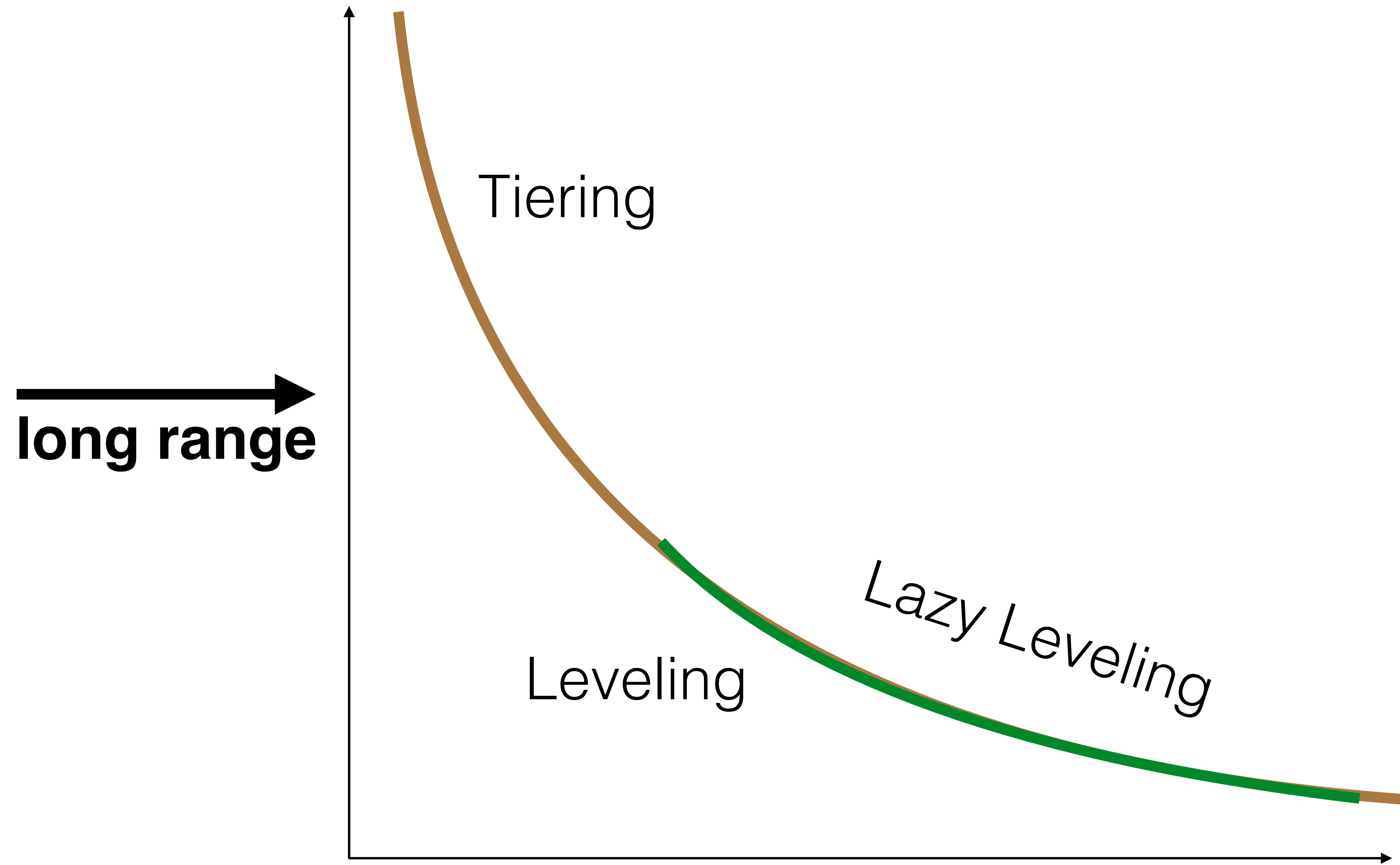


↖  
merging

**short range** →



merging



Tiering

Lazy Leveling

Leveling



Tiering  
updates 

Lazy Leveling

Leveling

Tiering  
updates 

Lazy Leveling

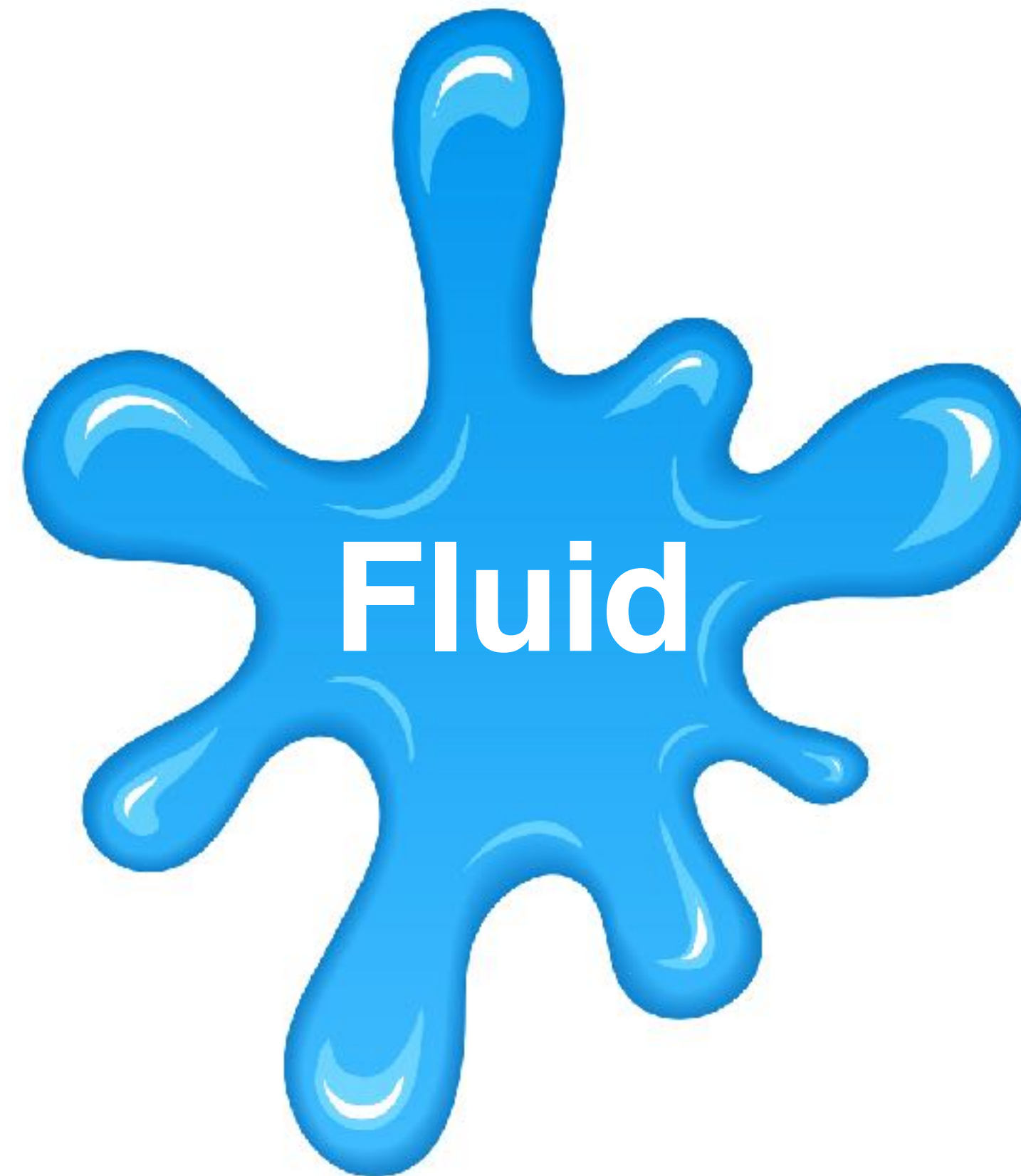
Leveling  
short range 

Tiering  
updates 

Lazy Leveling  
updates & point 

Leveling  
short range 

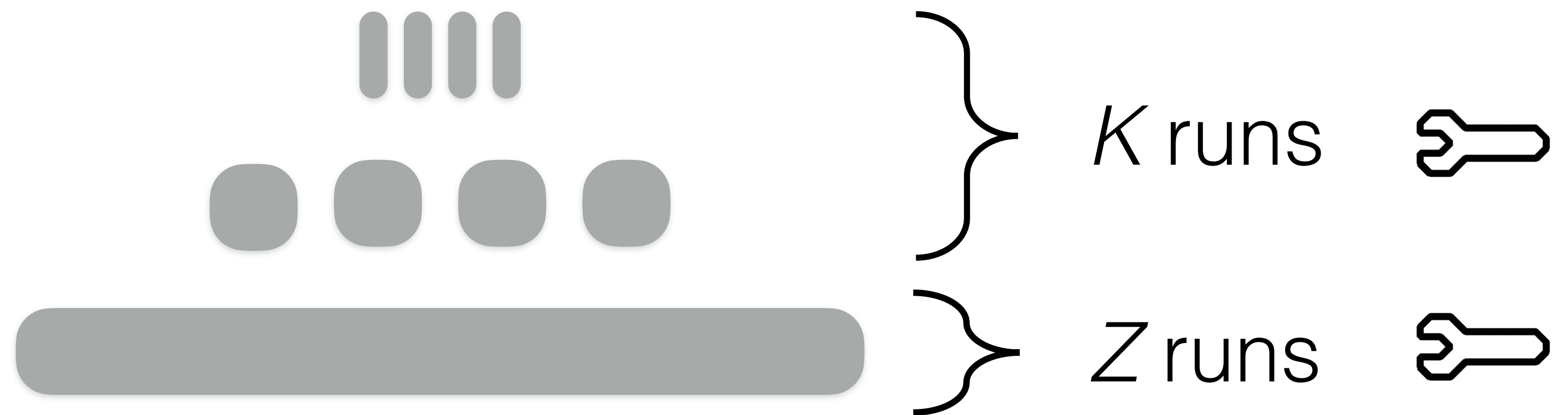
Lazy Leveling



Tiering

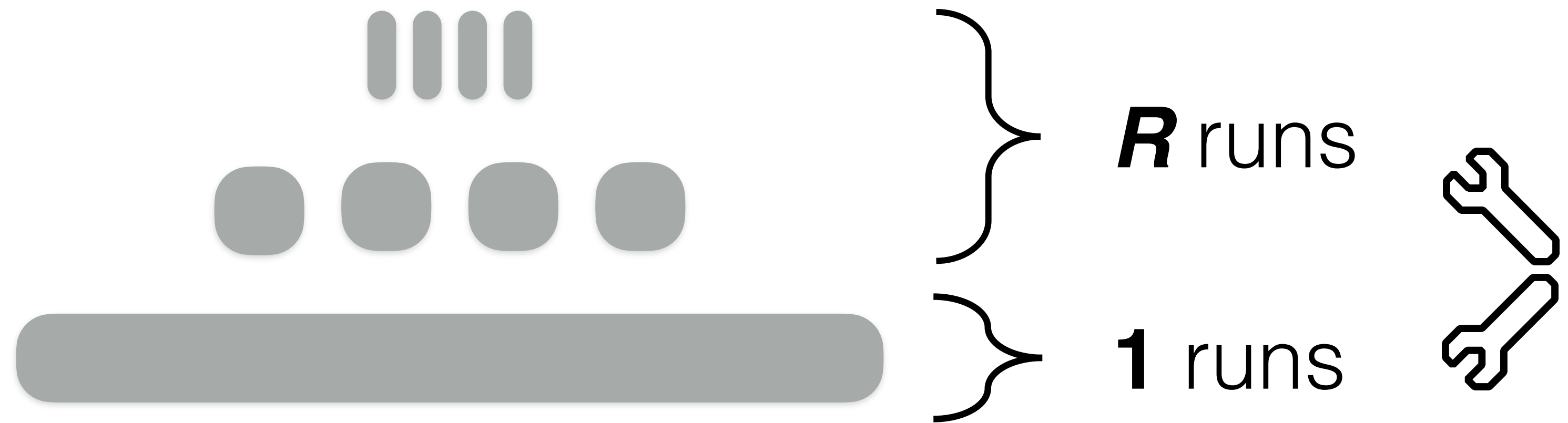
Leveling

# Fluid LSM-Tree



# Fluid LSM-Tree

Lazy Leveling



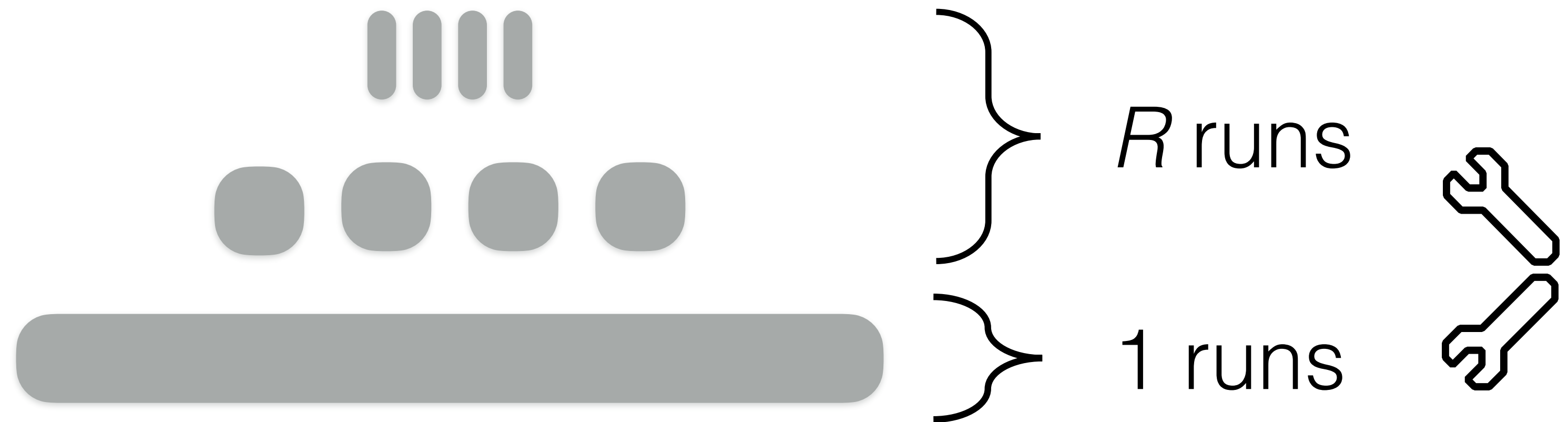
point

long range

short range

merging

Lazy Leveling



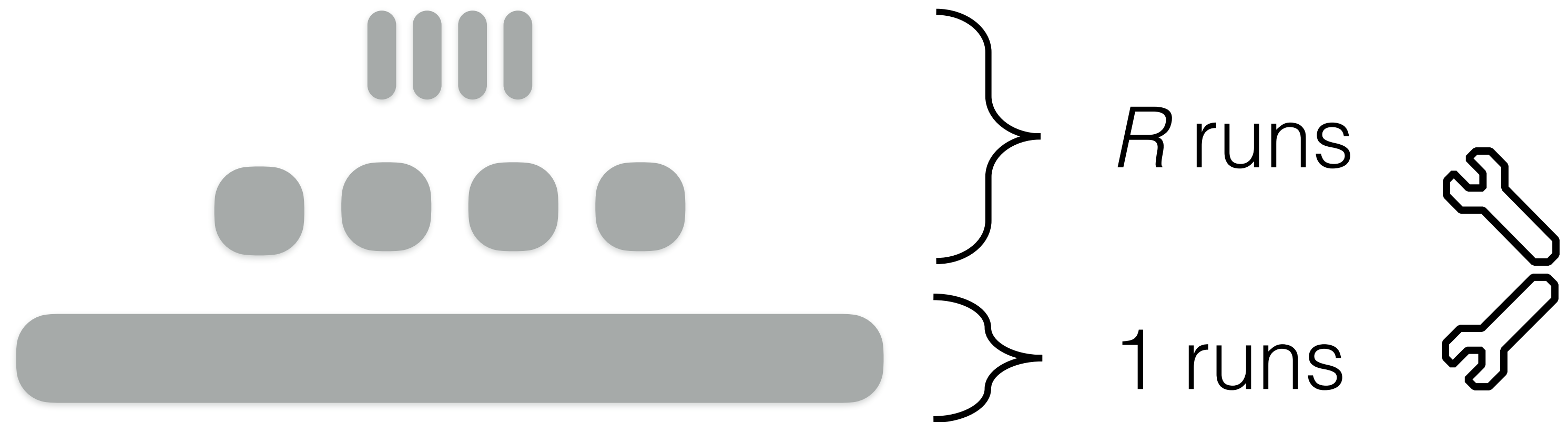
point

long range

optimize  
**short range**

merging

Lazy Leveling





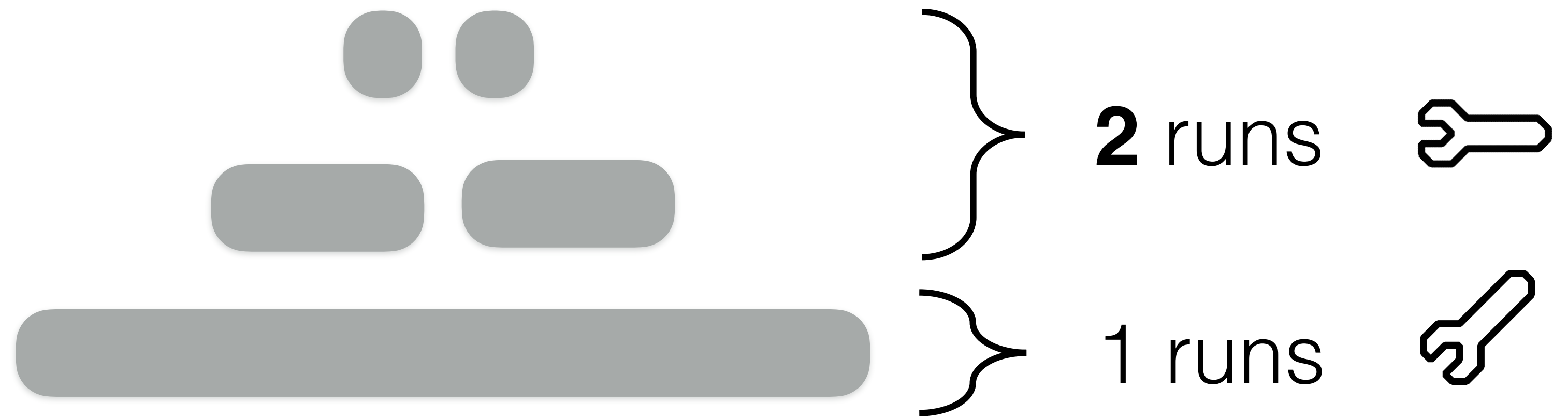
point

long range

optimize  
**short range**

**merging**

Lazy Leveling



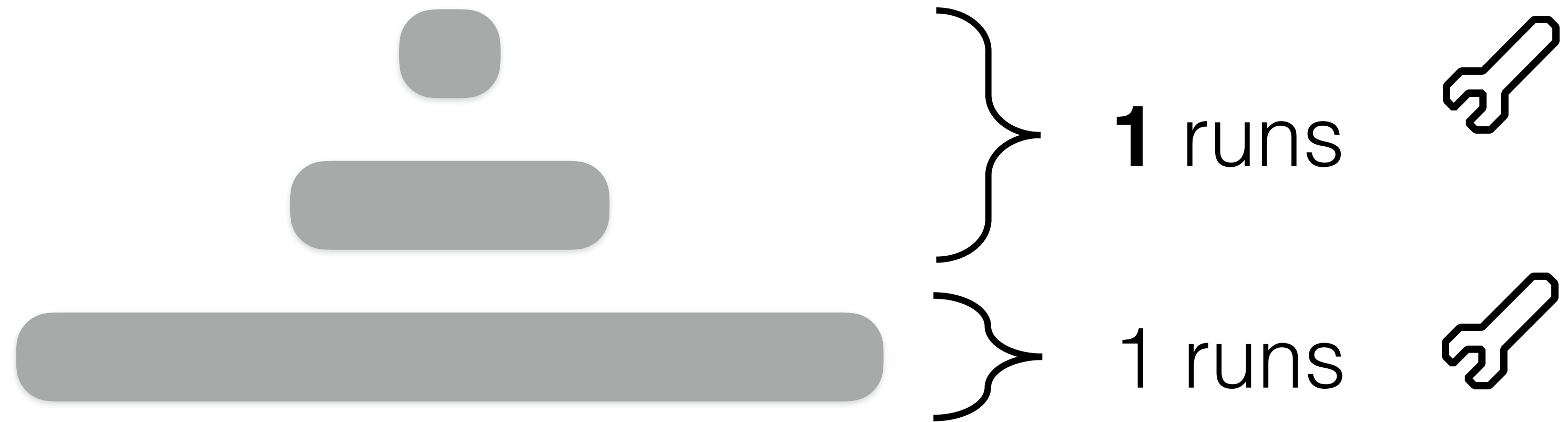
point

long range

optimize  
**short range**

**merging**

# Leveling



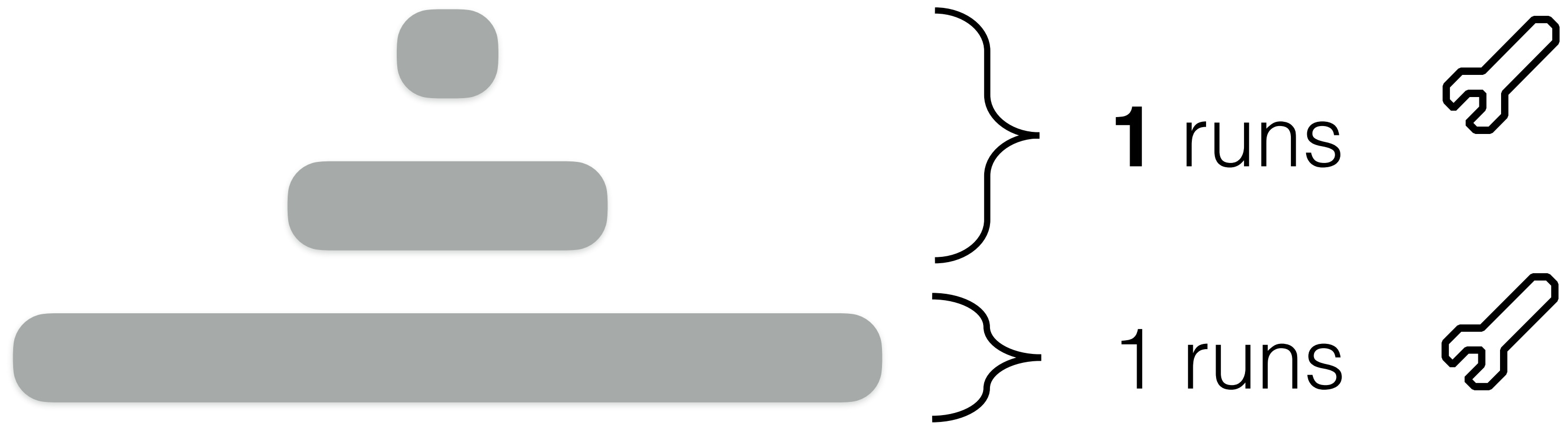
point

long range

short range

optimize  
**merging**

## Leveling



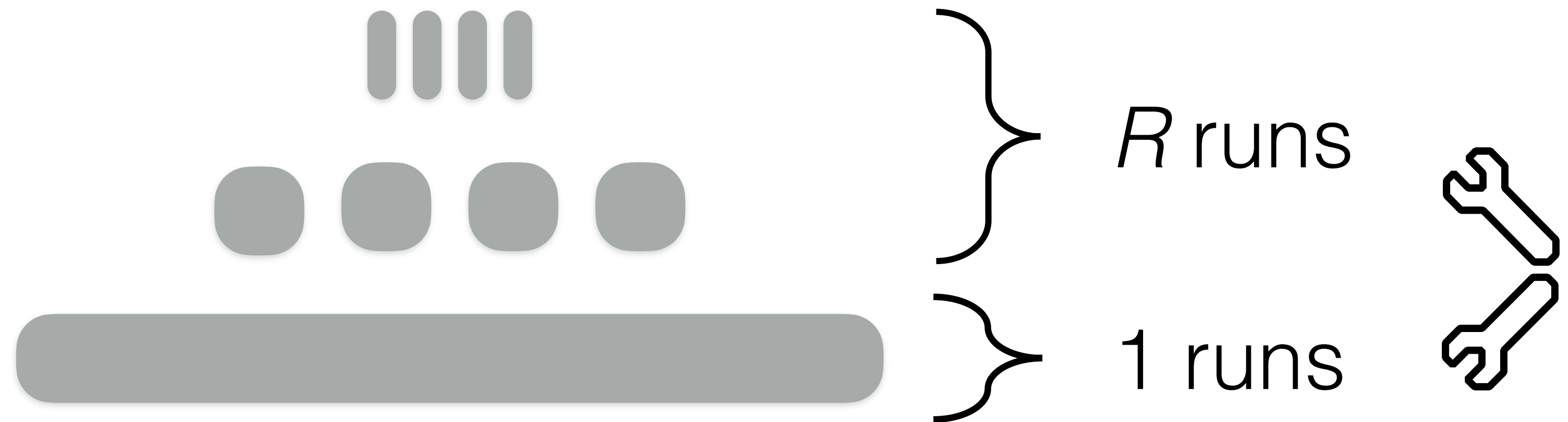
point

long range

**short range**

optimize  
**merging**

Lazy Leveling



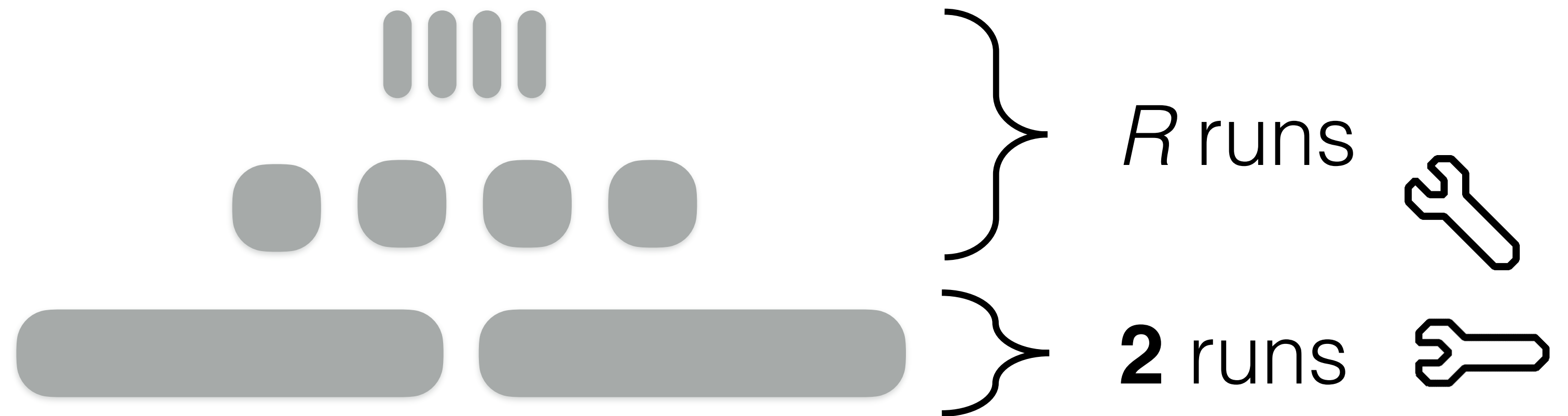
point

long range

short range

optimize  
**merging**

Lazy Leveling



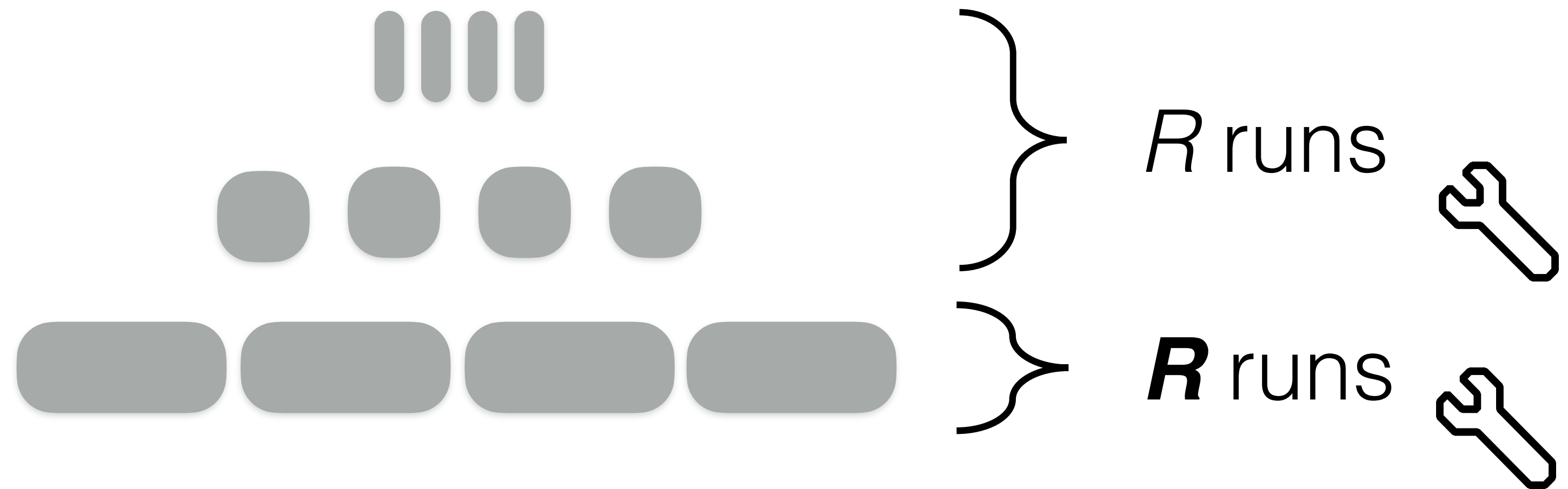
point

long range

short range

optimize  
**merging**

## Tiering



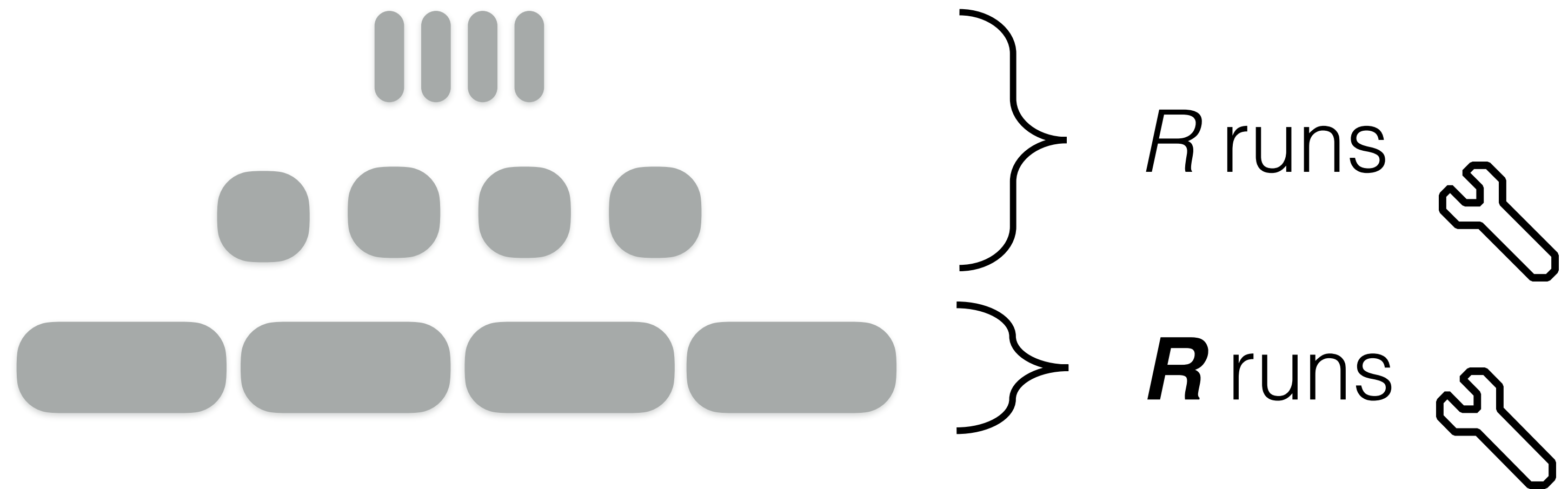
optimize  
**point**

long range

short range

merging

**Tiering**



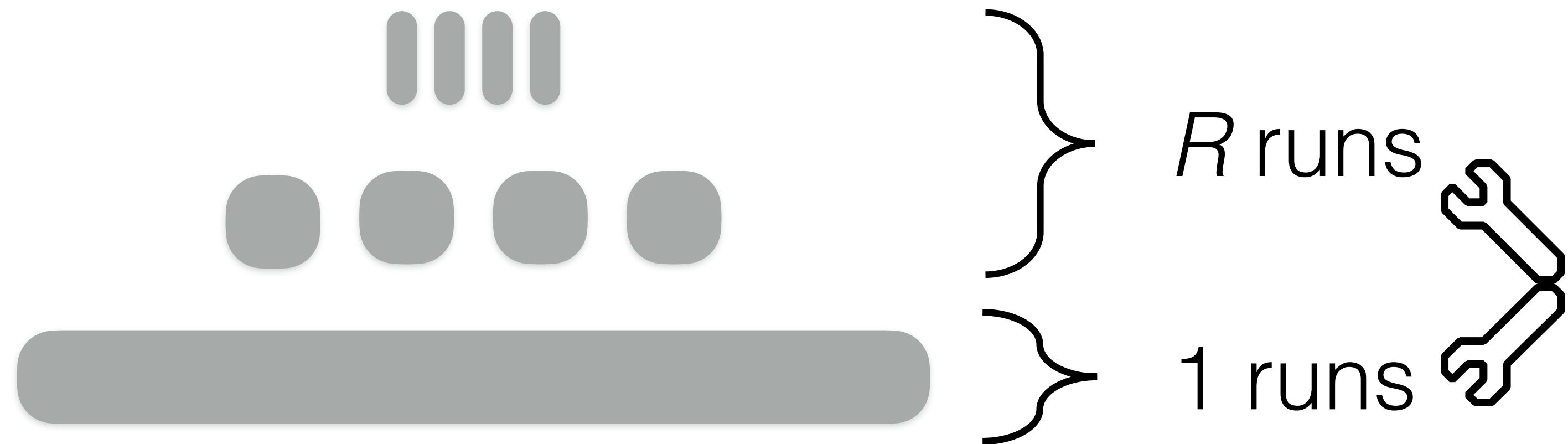
optimize  
**point**

long range

short range

merging

Lazy Leveling



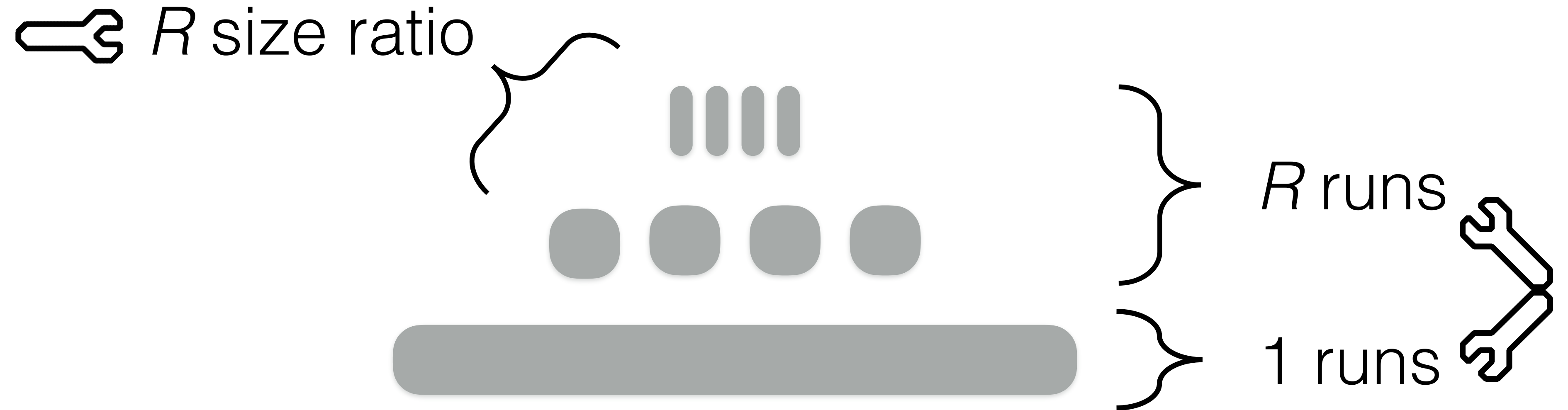


optimize  
**point**

long range

short range

merging



Lazy Leveling

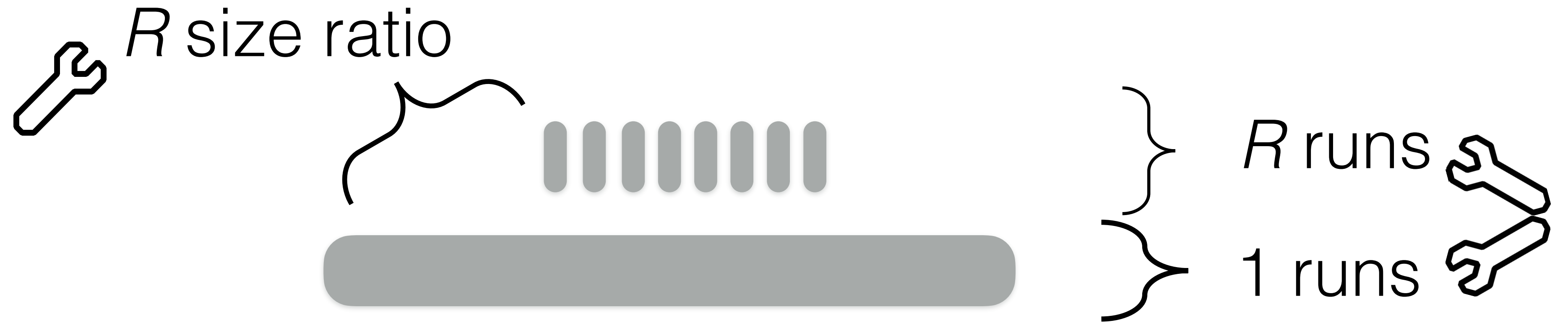
optimize  
**point**

long range

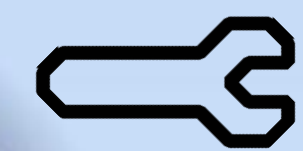
short range

merging

Lazy Leveling



## Fluid LSM-Tree



$R$  size ratio

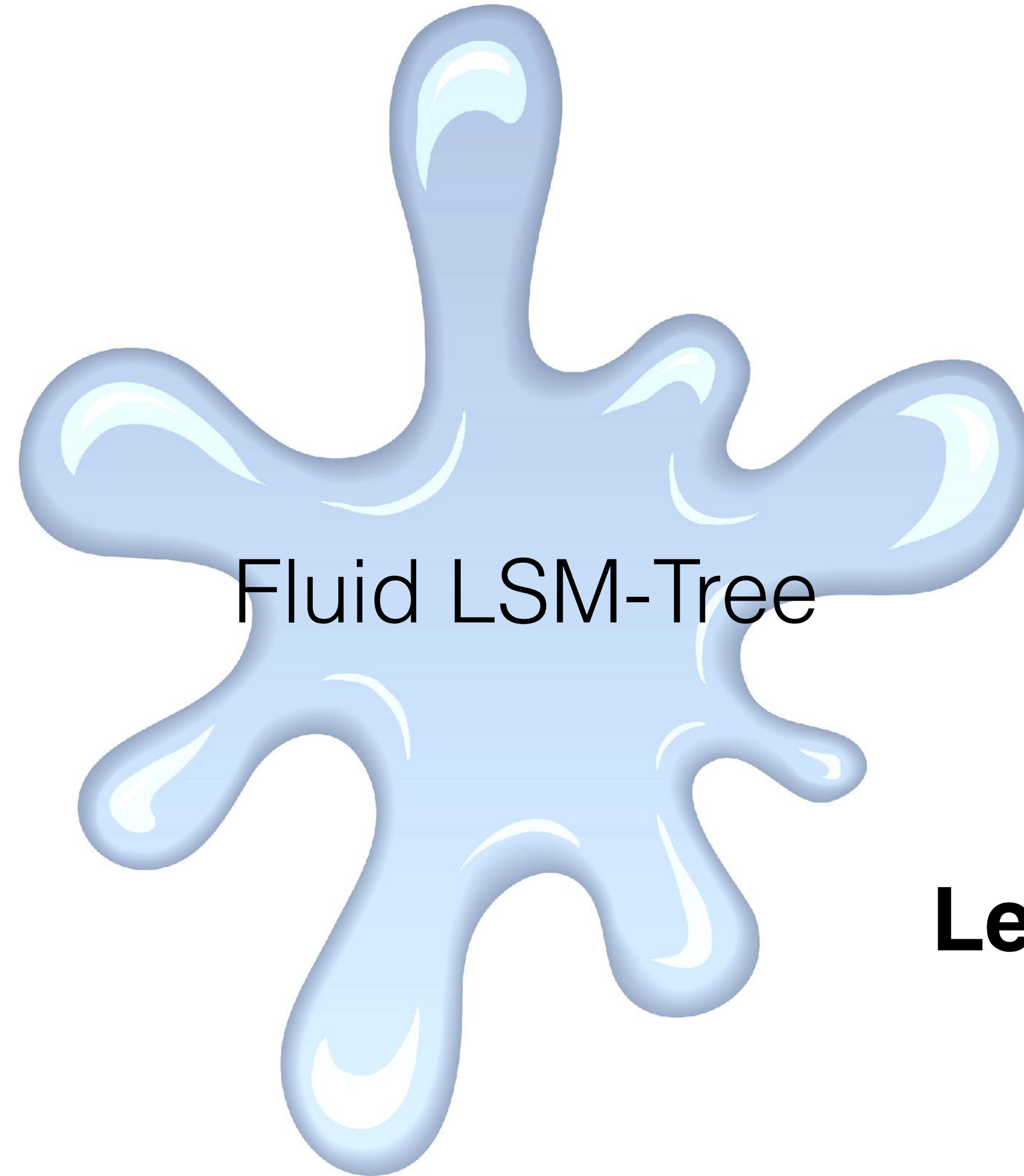


$K$  runs at smaller levels



$Z$  runs at largest level

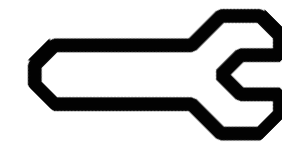
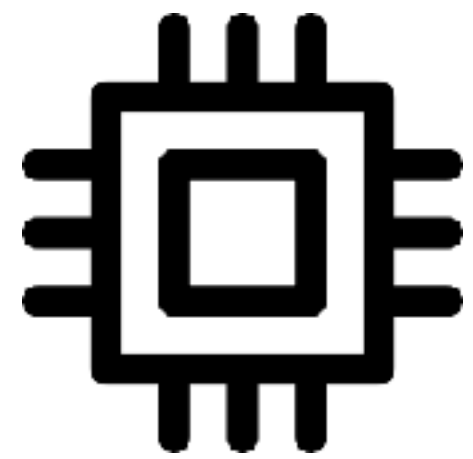
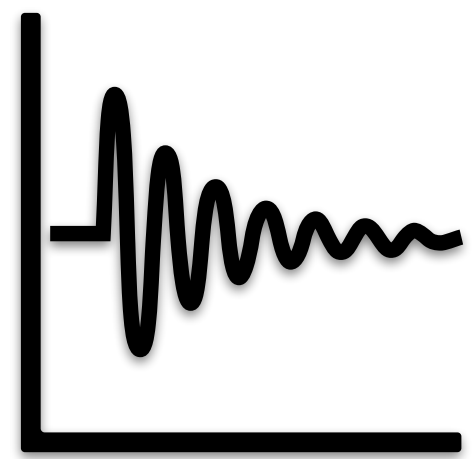
# Lazy Leveling

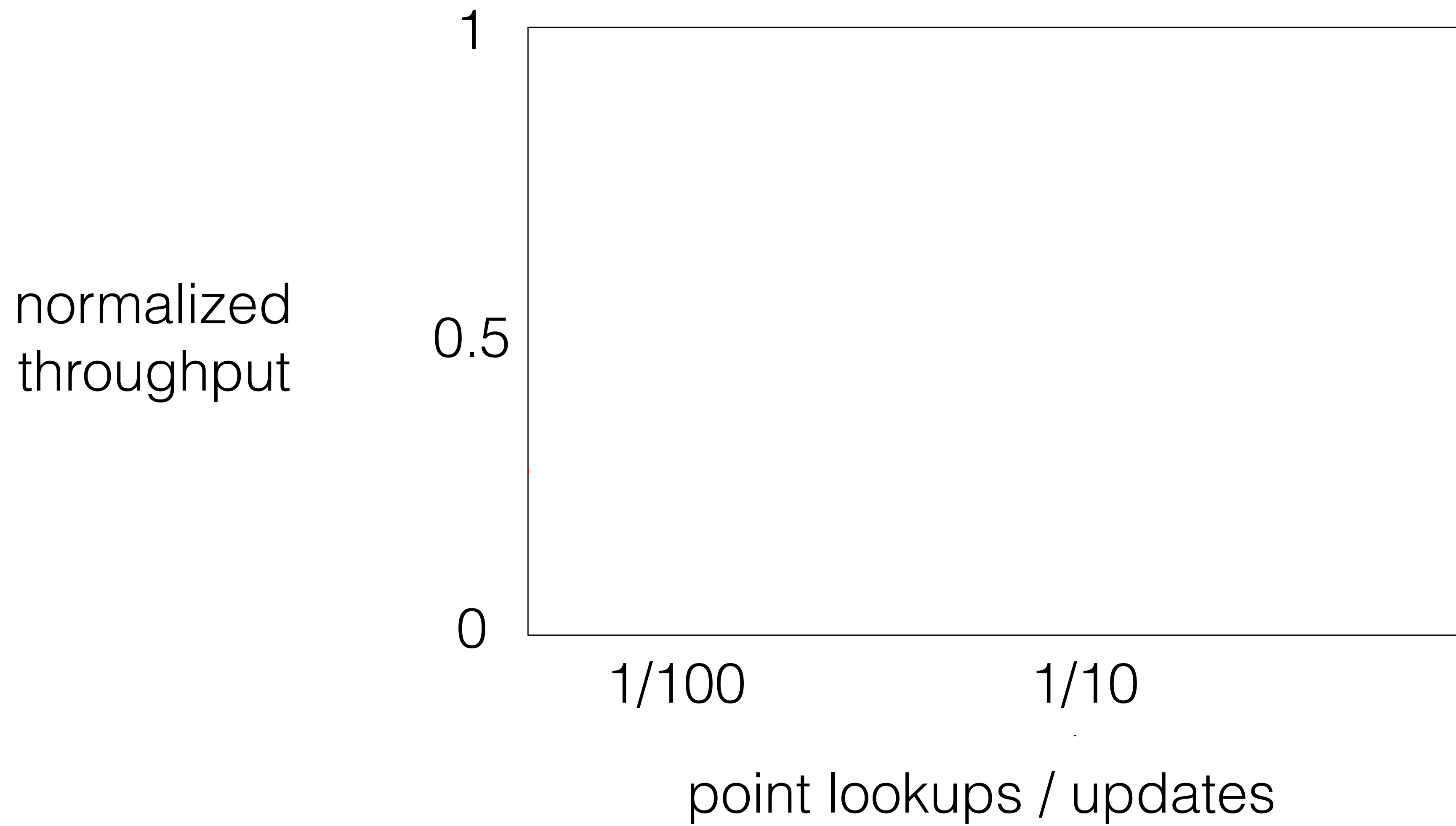


Fluid LSM-Tree

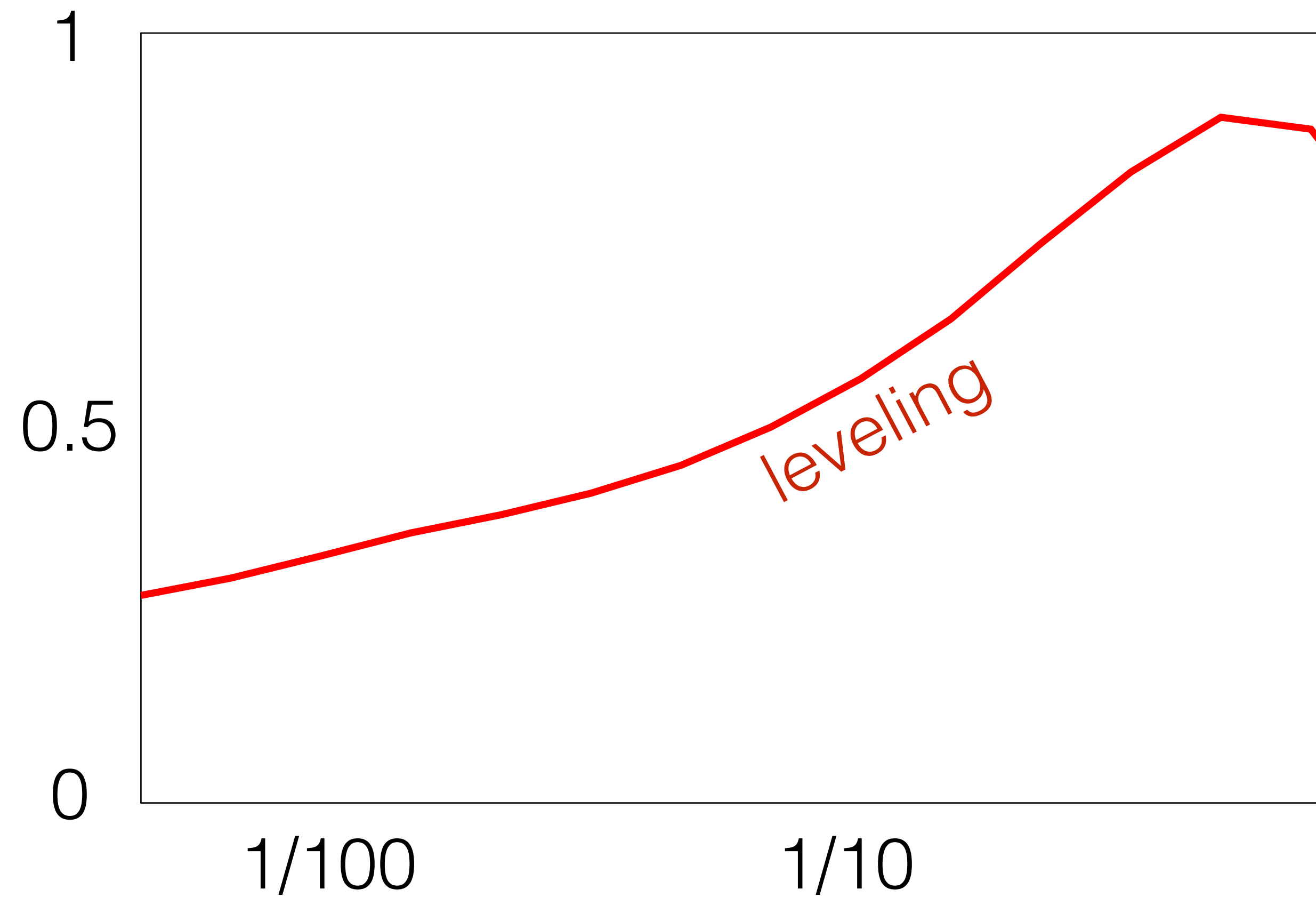
**Tiering**

**Leveling**



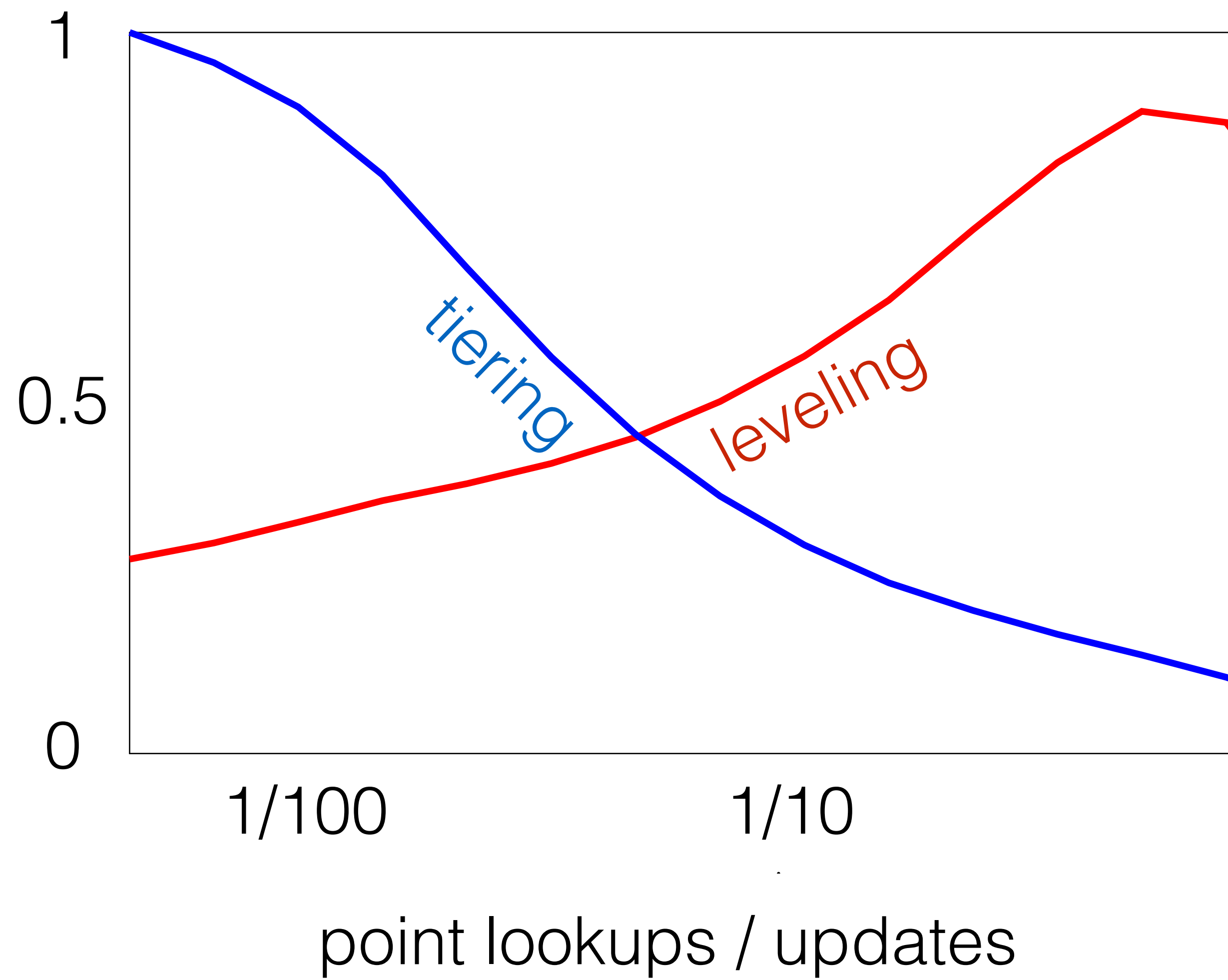


normalized  
throughput



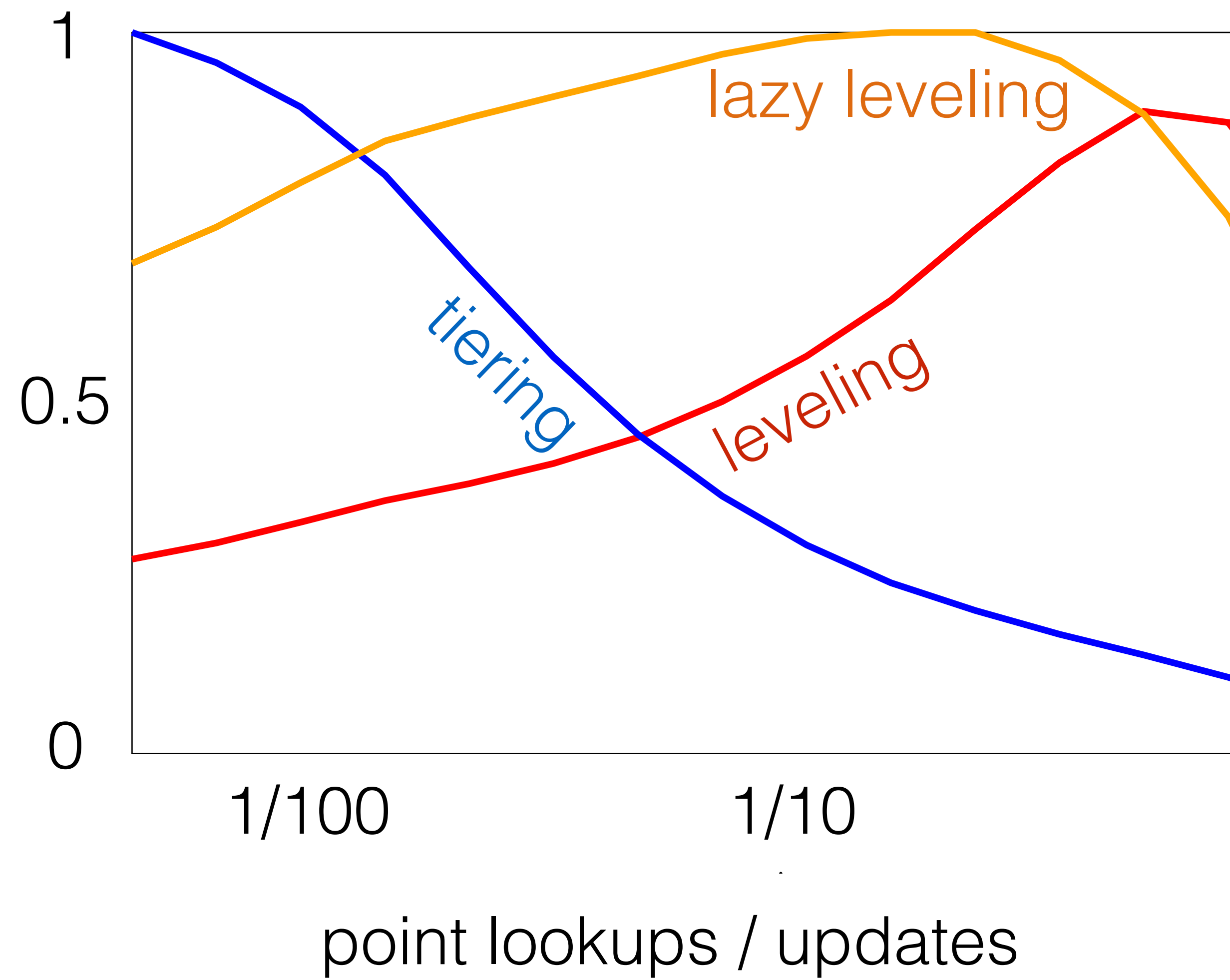
point lookups / updates

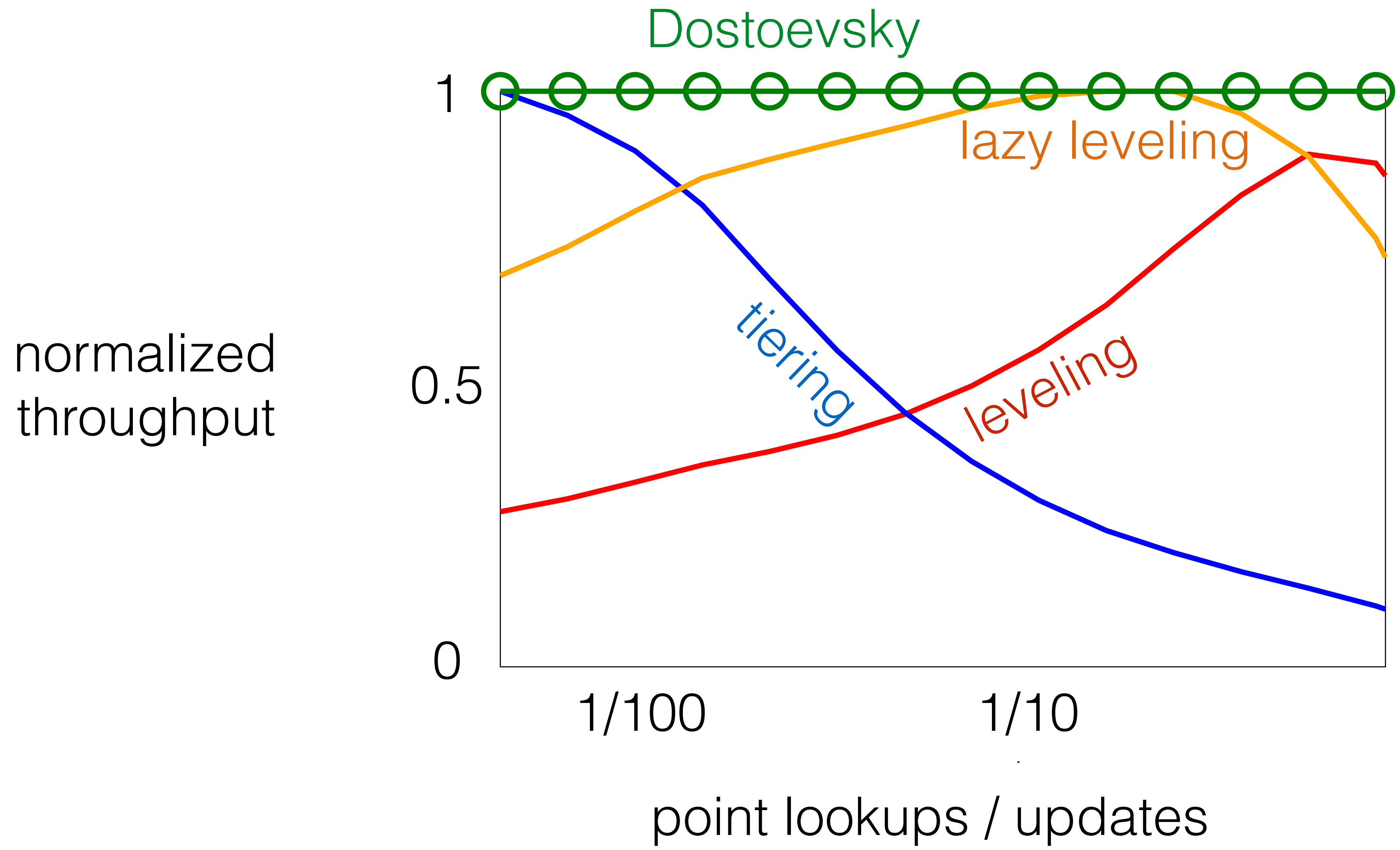
normalized  
throughput



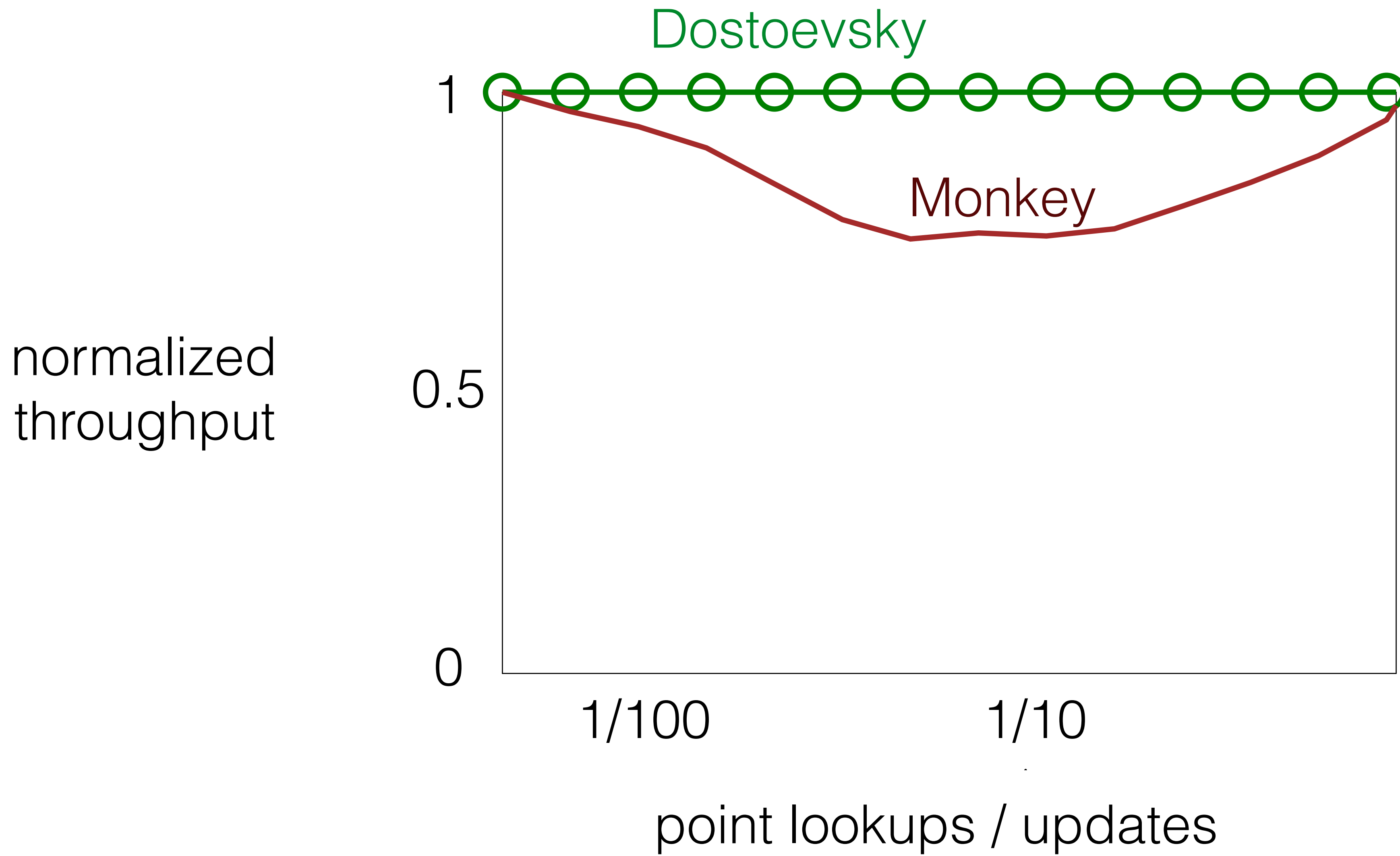


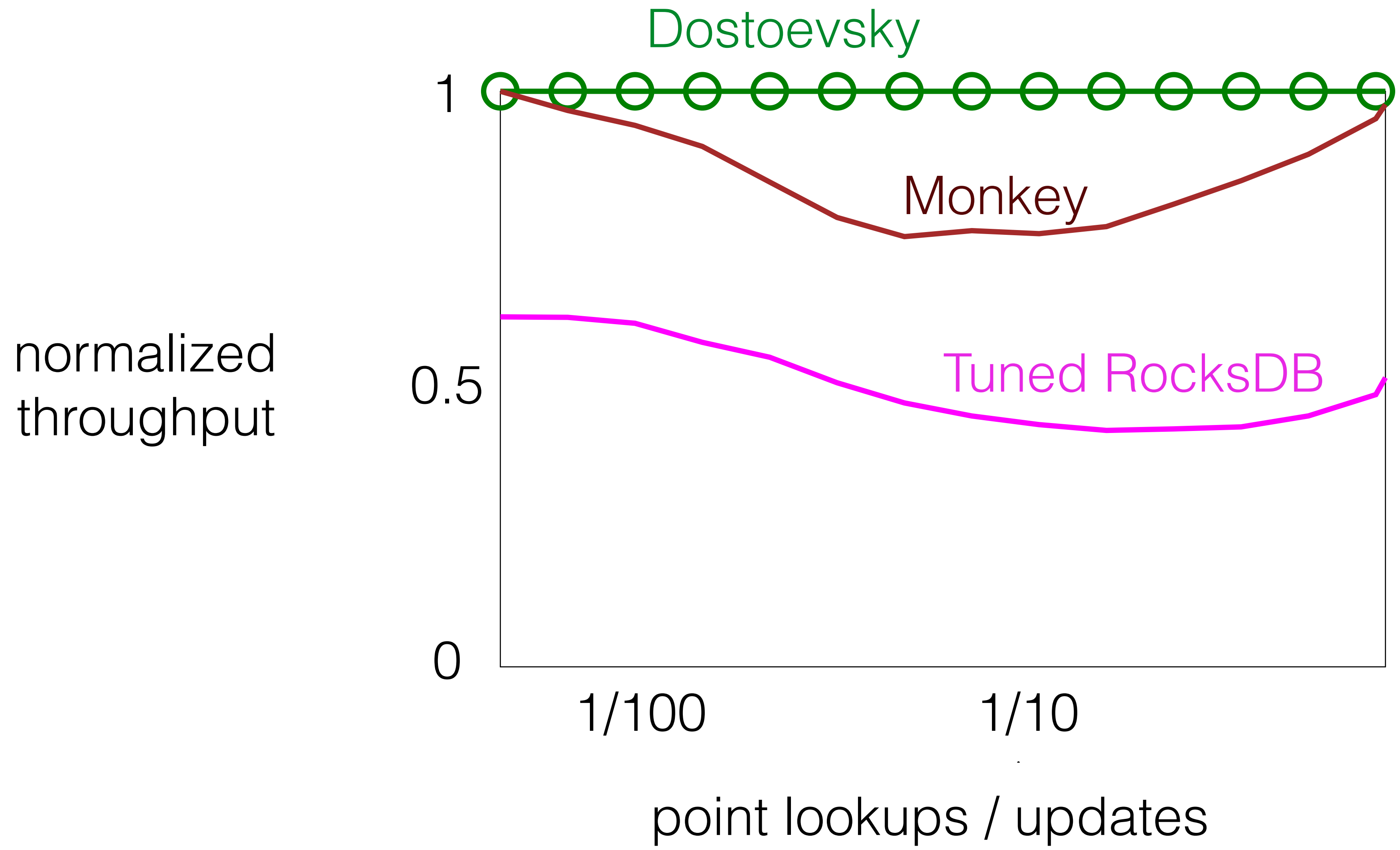
normalized  
throughput

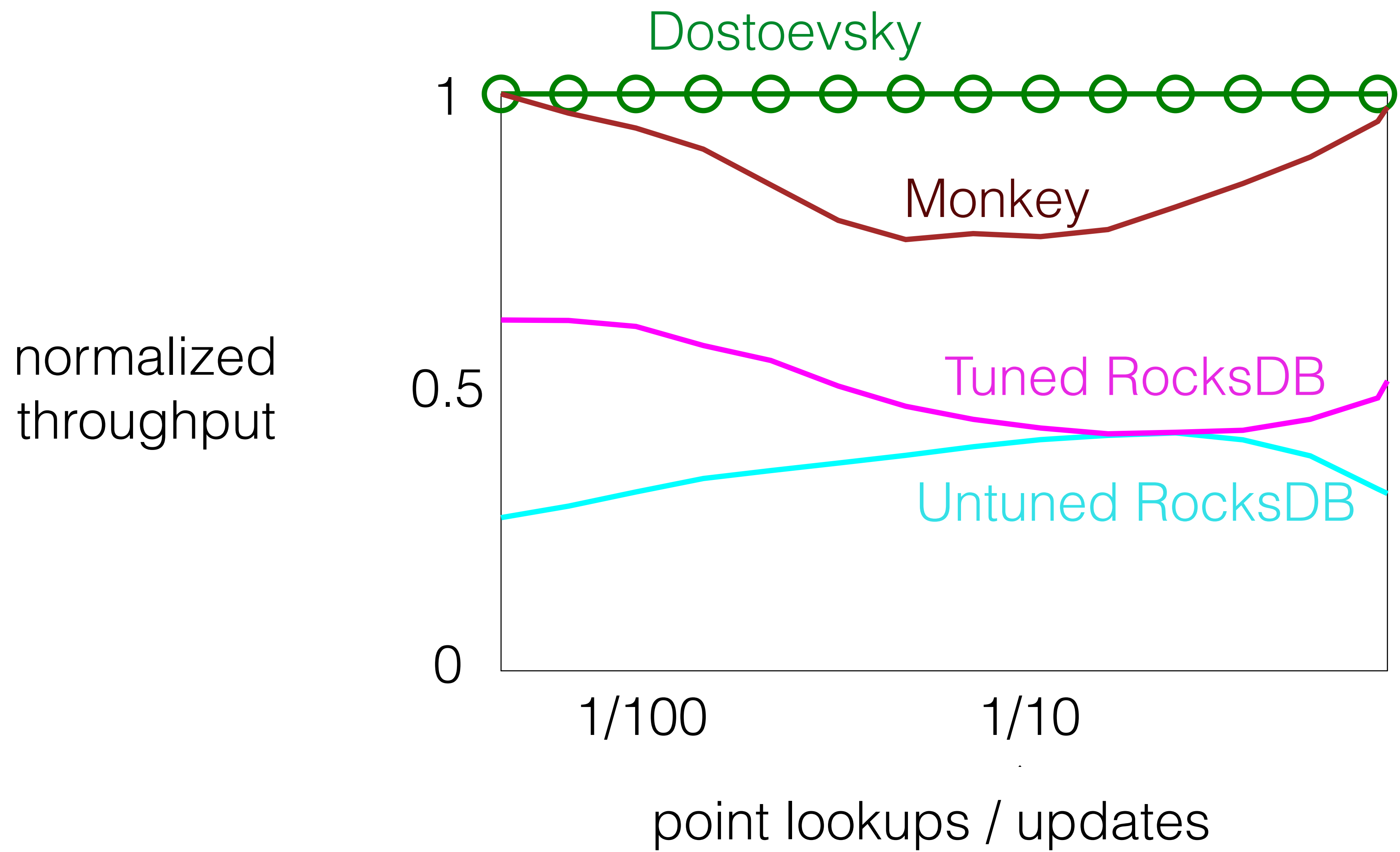








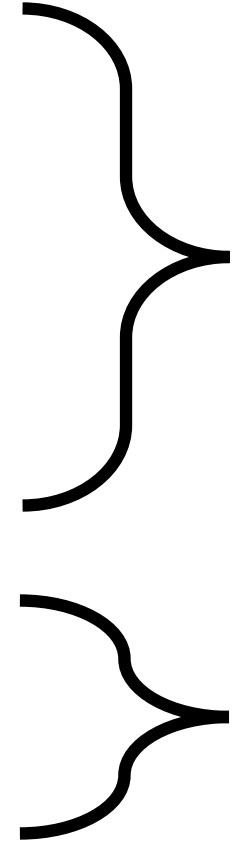




# Conclusion

# Conclusion

**Lazy Leveling**

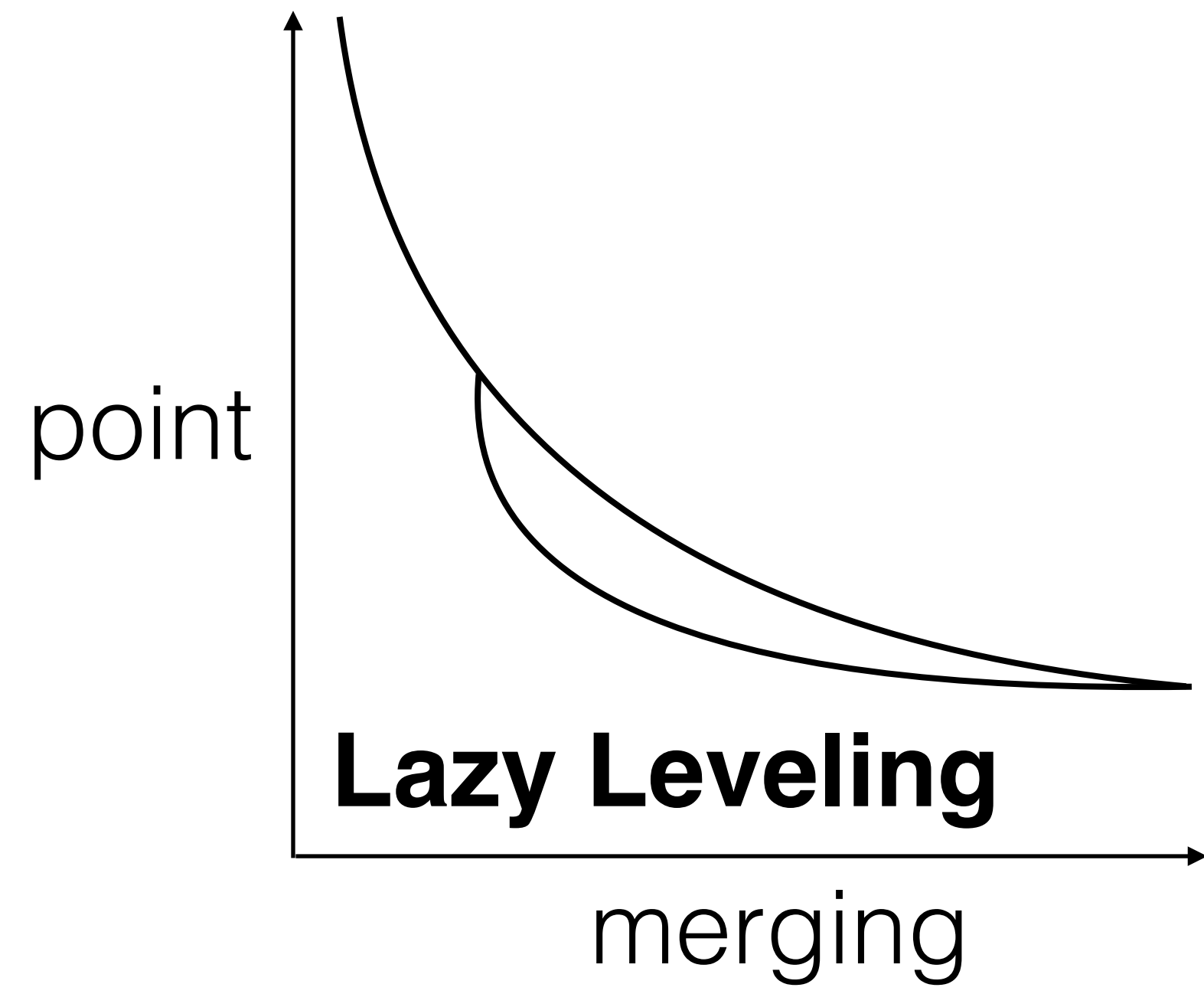


no merging

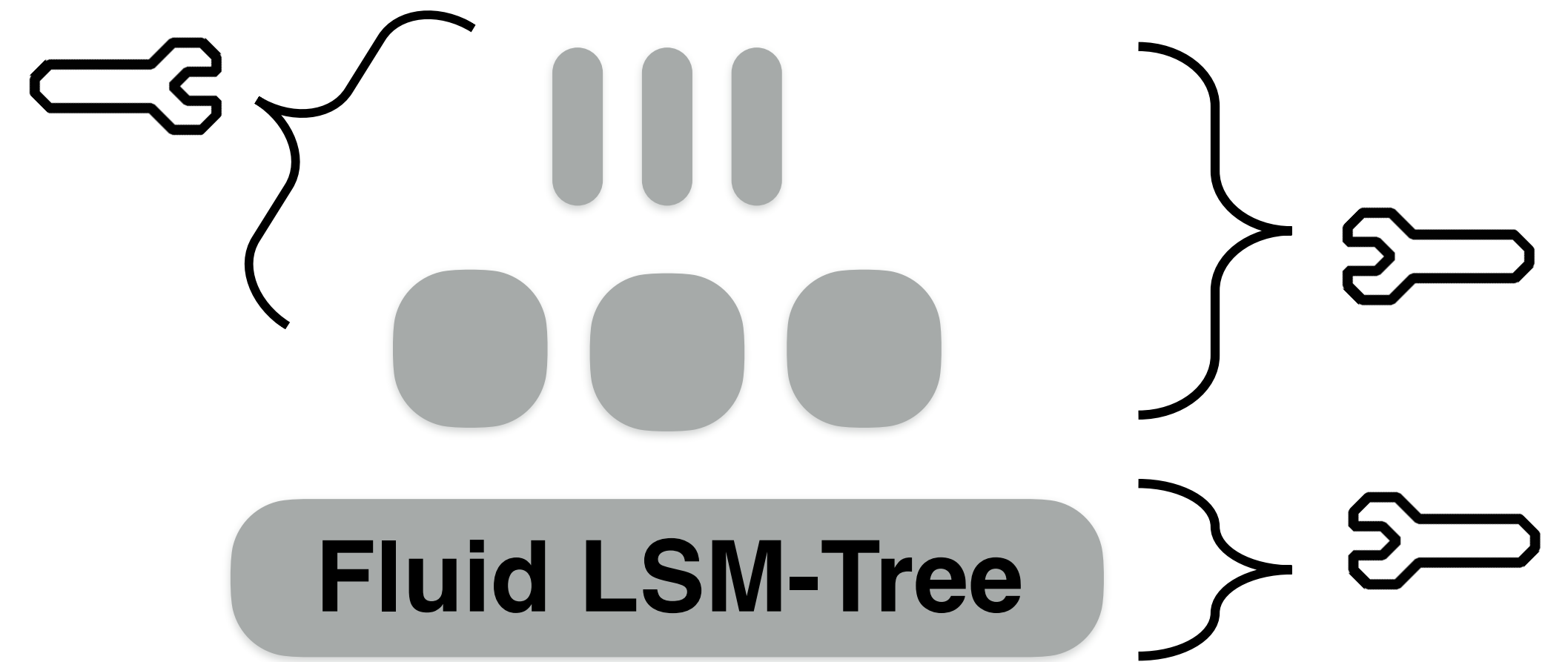
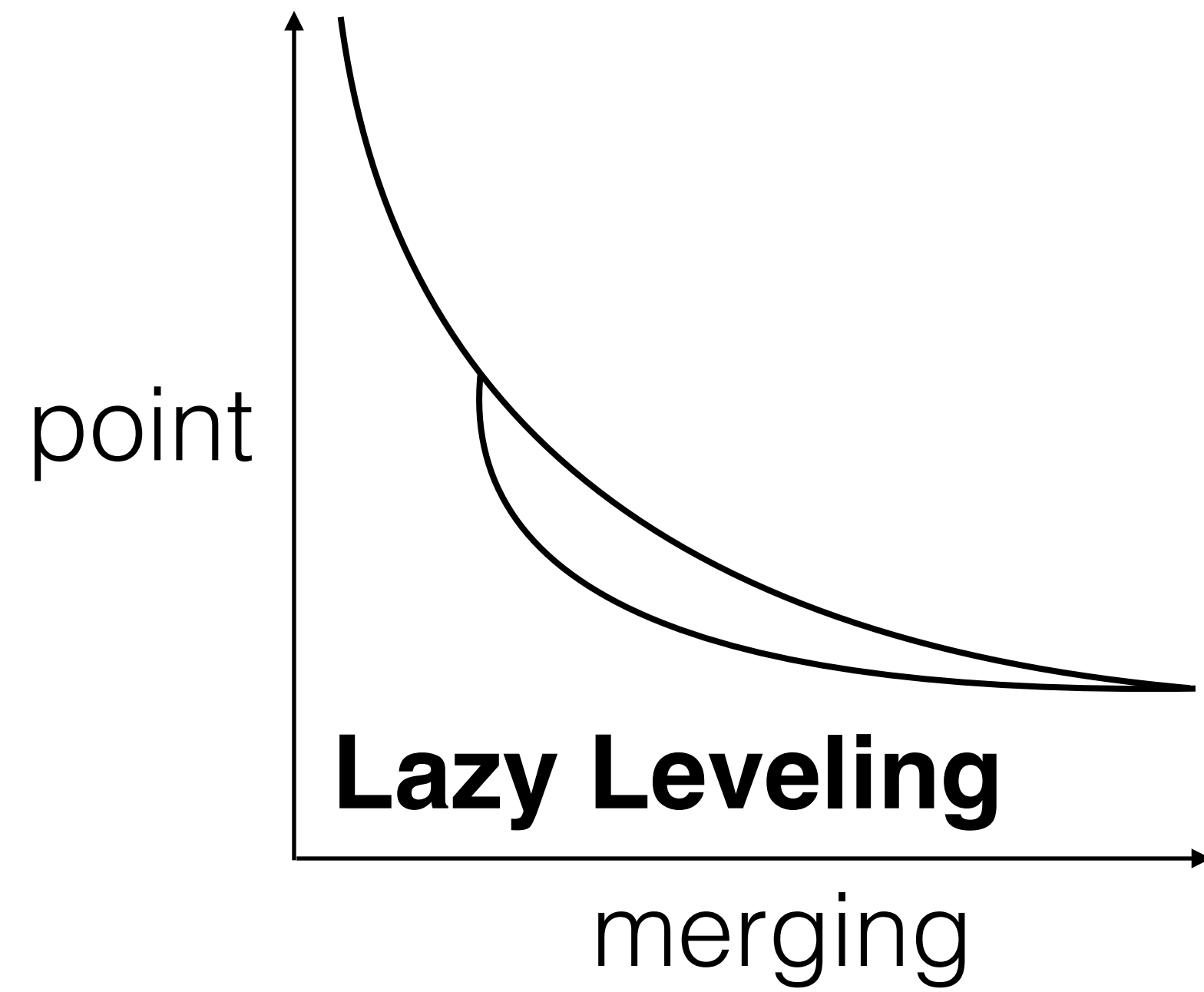
greedy merging



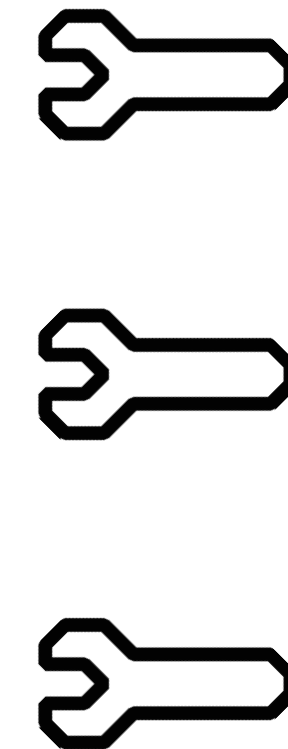
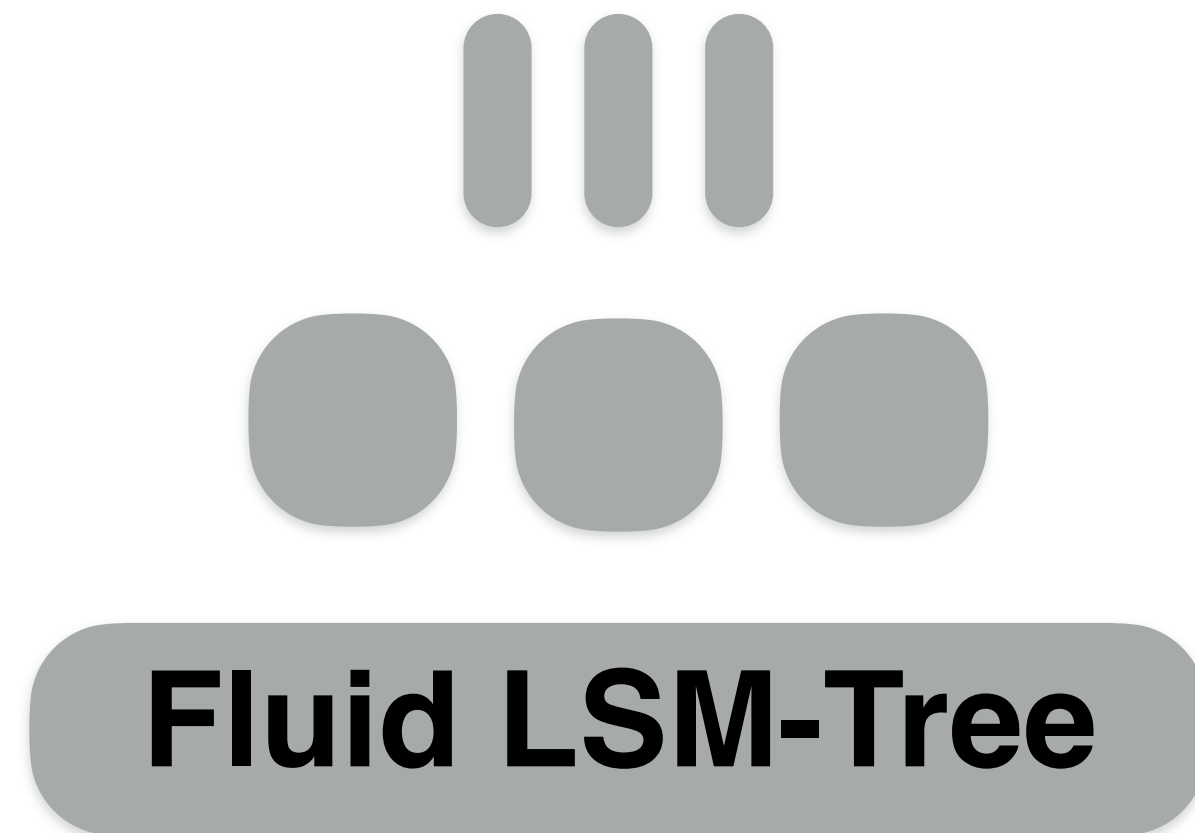
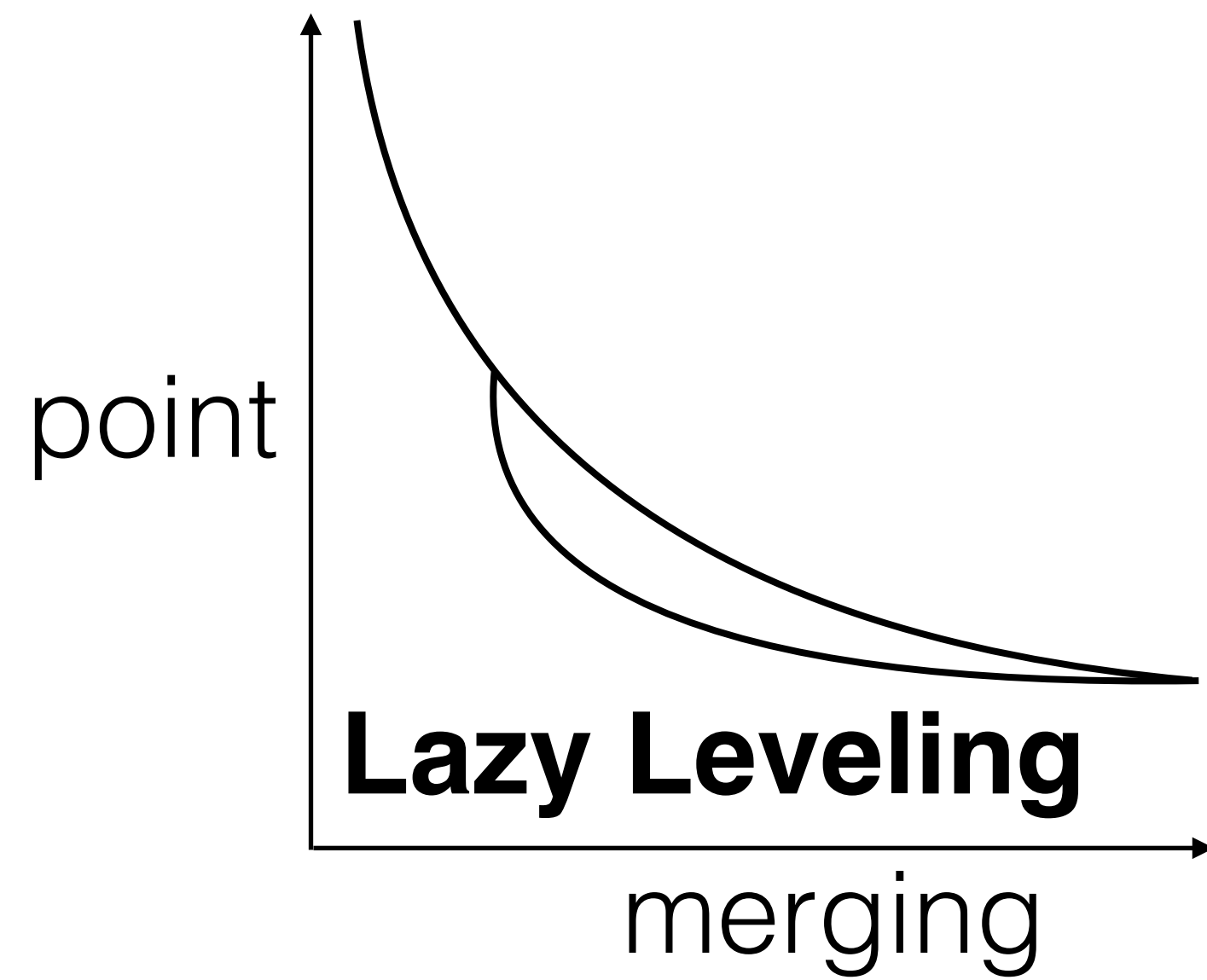
# Conclusion



# Conclusion



# Conclusion



**Dostoevsky**



**D**ostoevsky: **S**pace-**T**ime **O**ptimized **E**volvable **S**calable **K**ey-Value Store

VERY WRITE-OPTIMIZED

Thanks!



A self-designing key-value store

[www.crimsondb.org](http://www.crimsondb.org)