

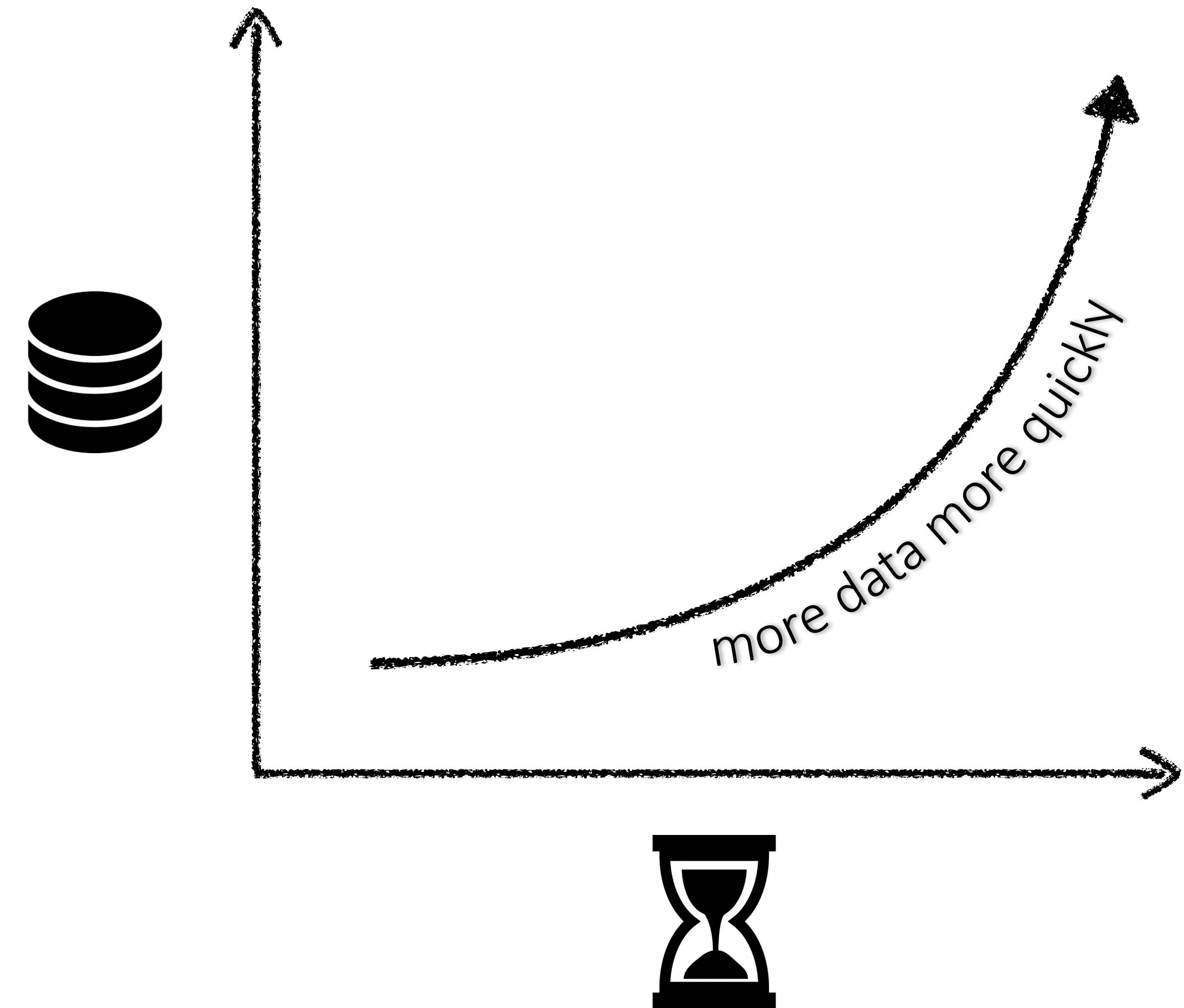


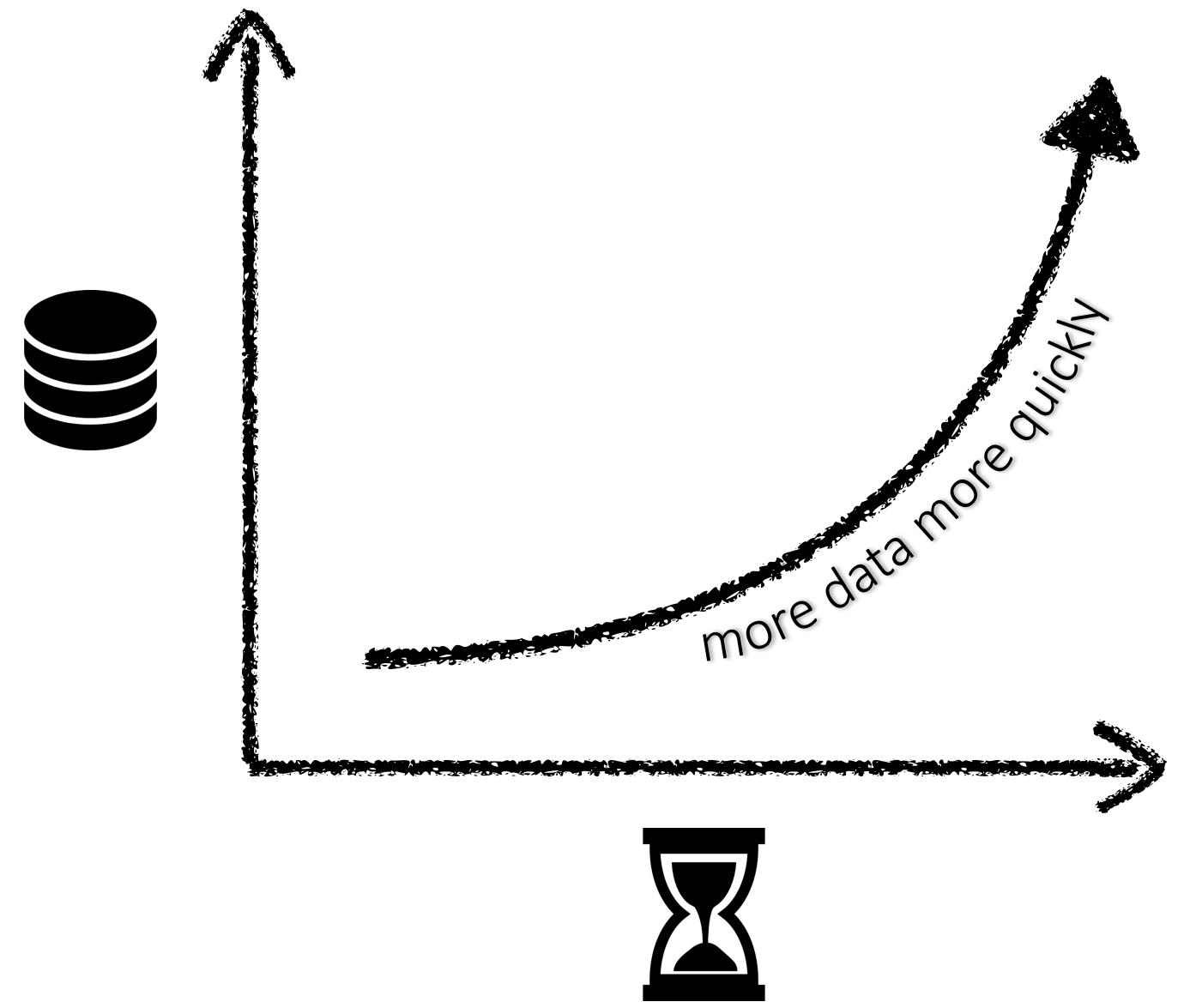
The Log-Structured Merge-Bush & the Wacky Continuum

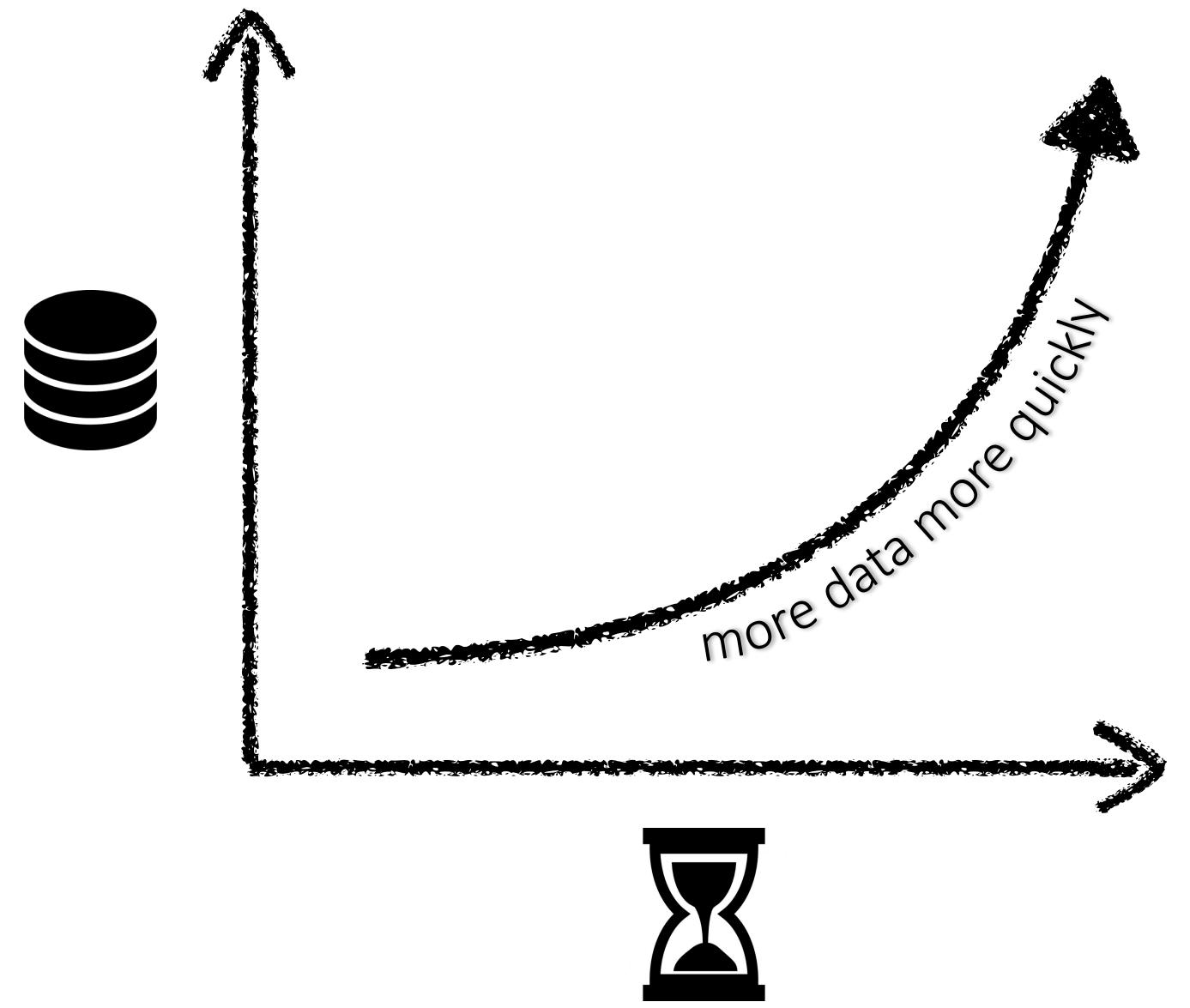


DASlab
@ Harvard SEAS

Niv Dayan &
Stratos Idreos





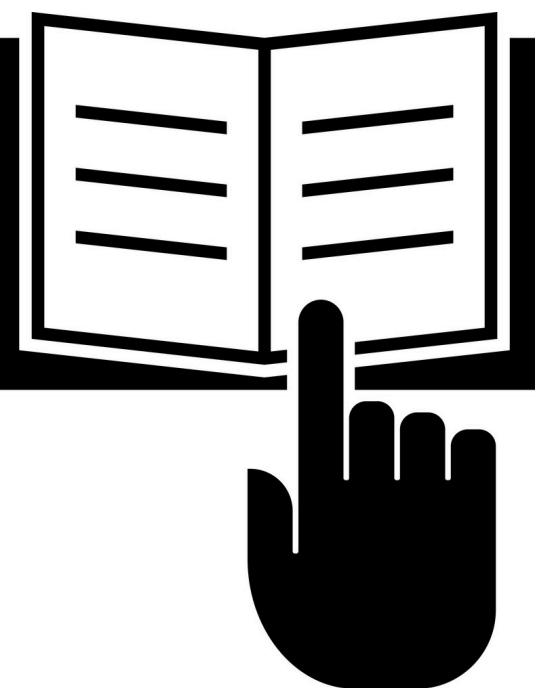




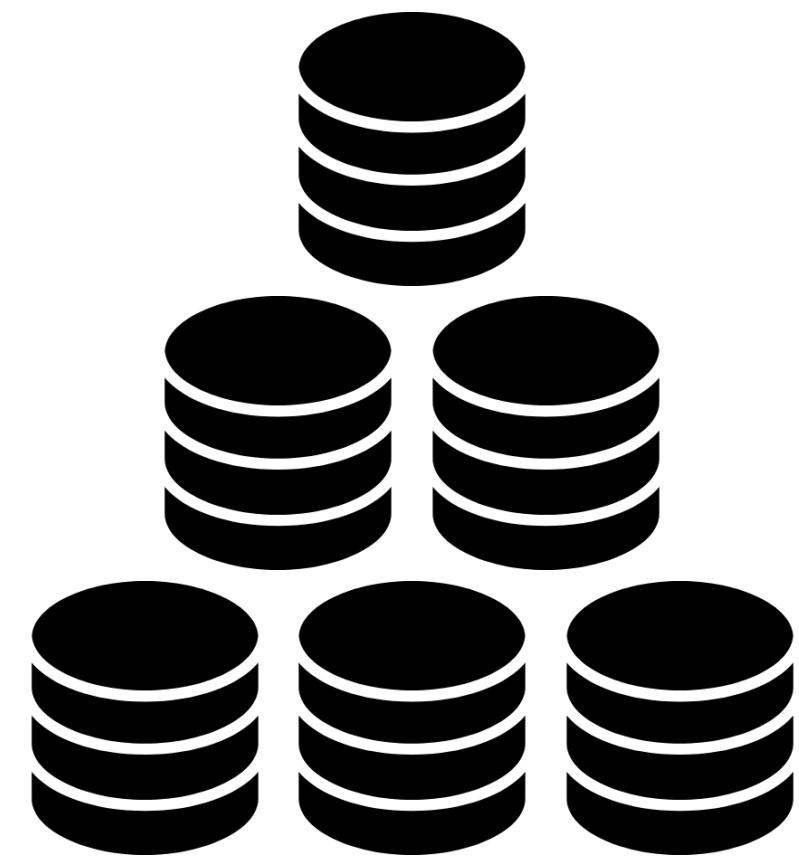
fast writes



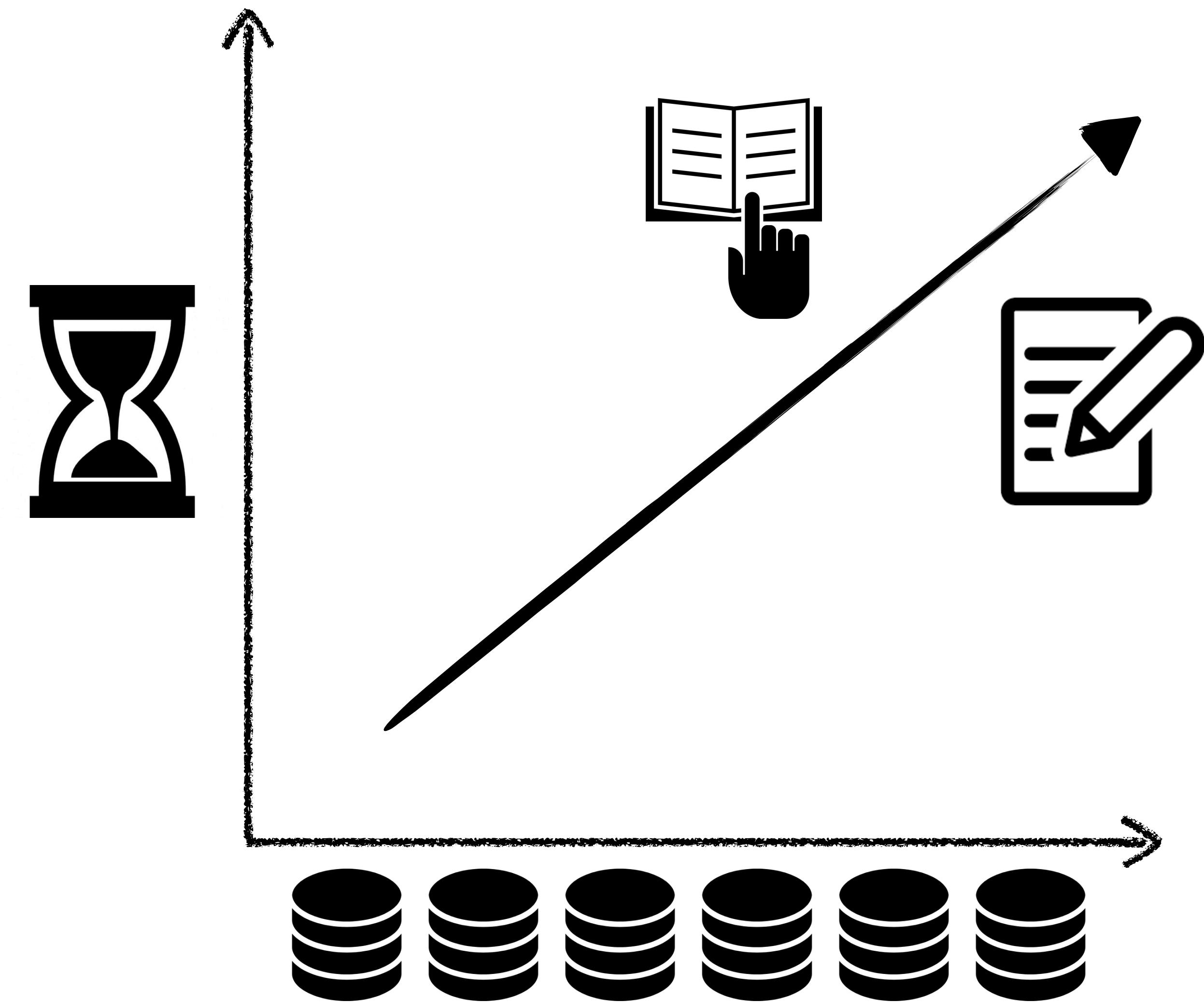
fast writes

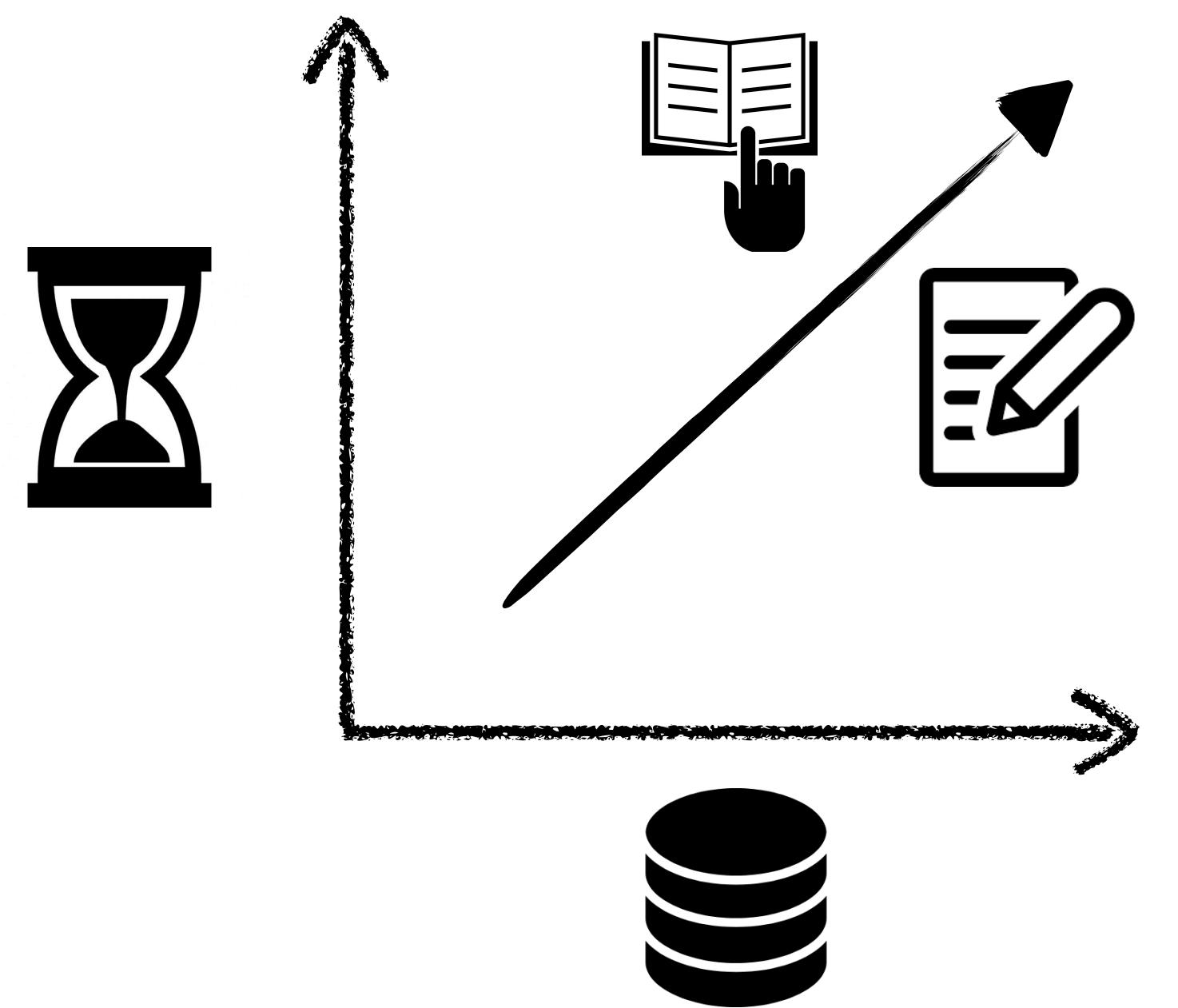


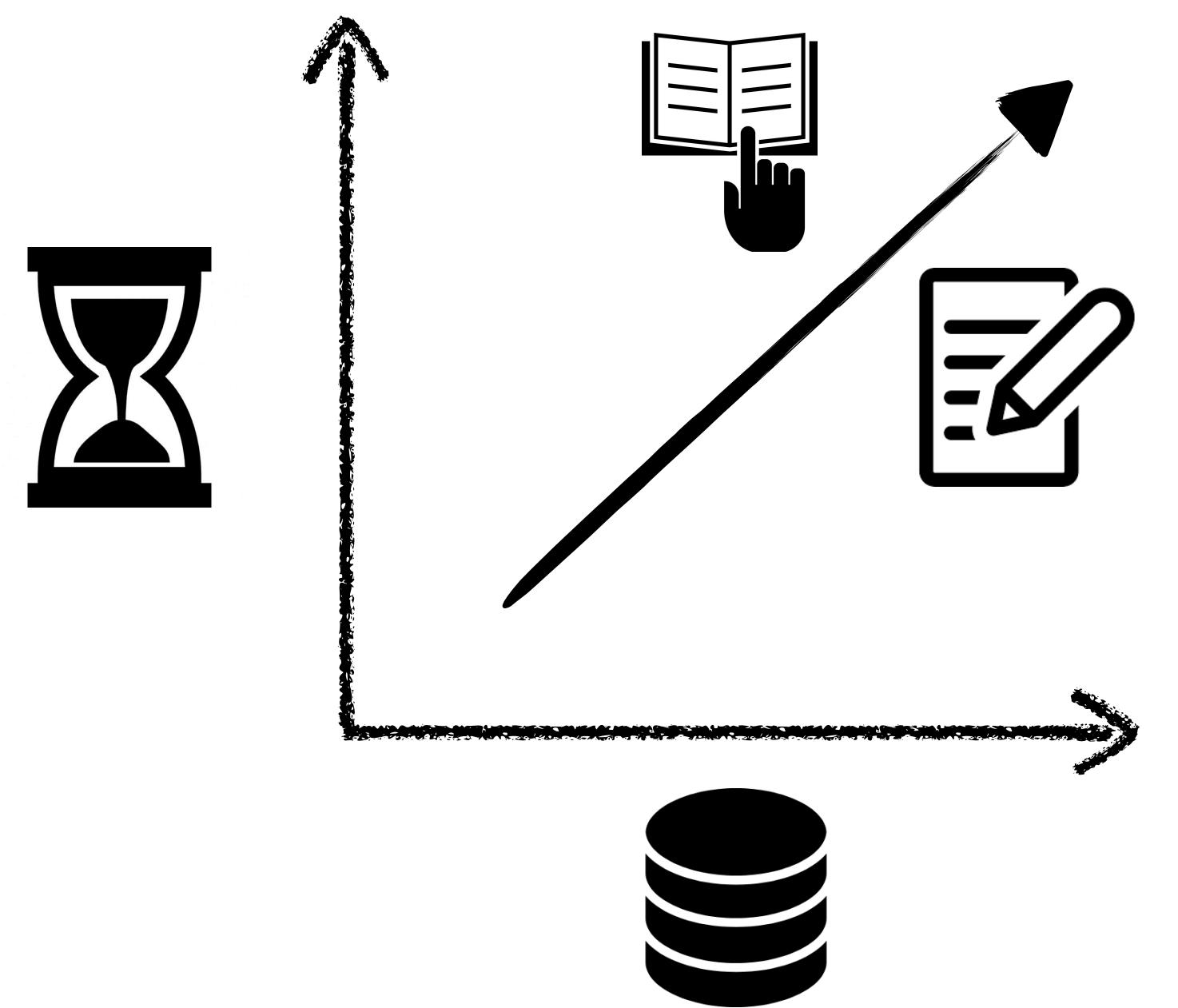
fast reads



massive data







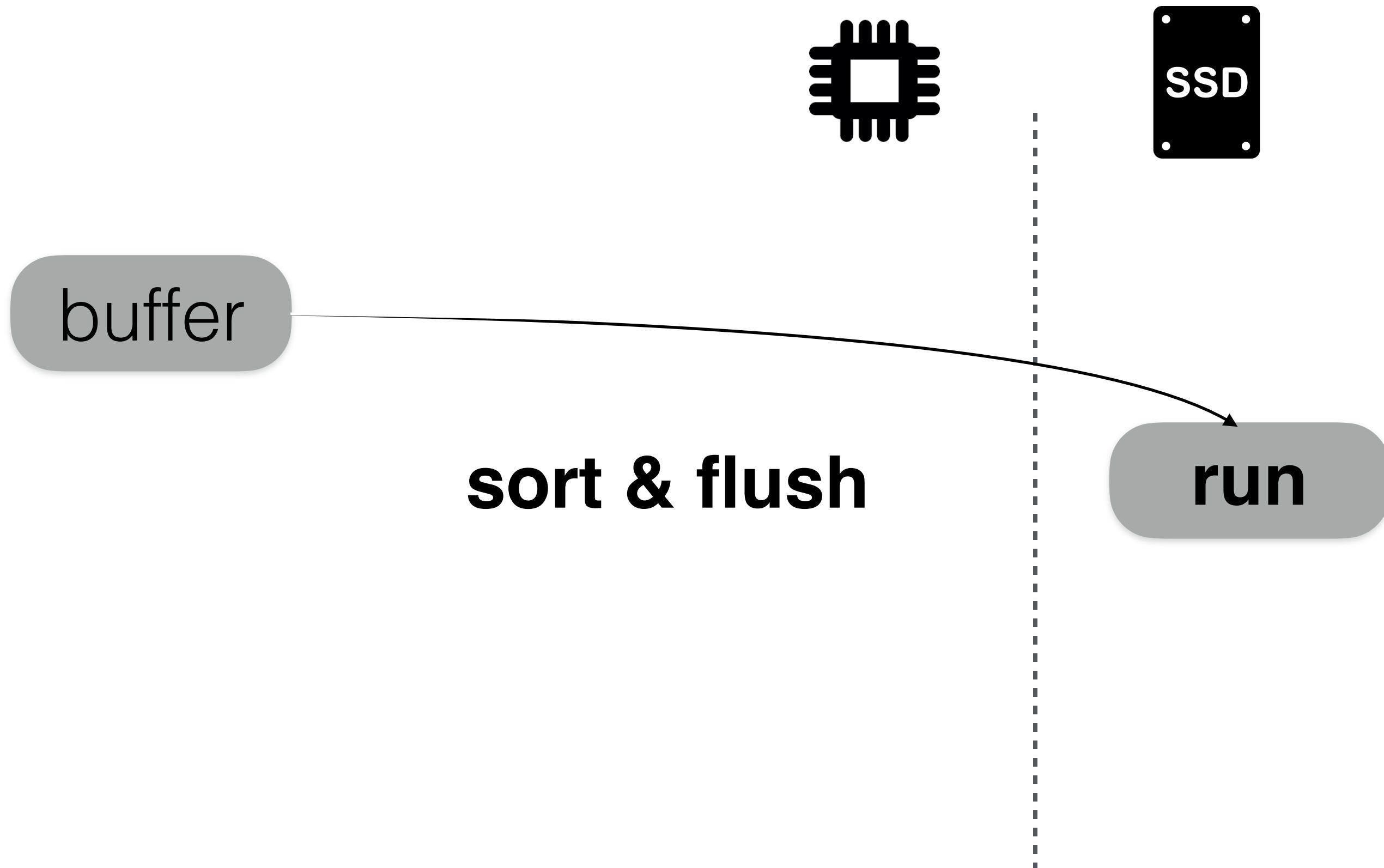
Background

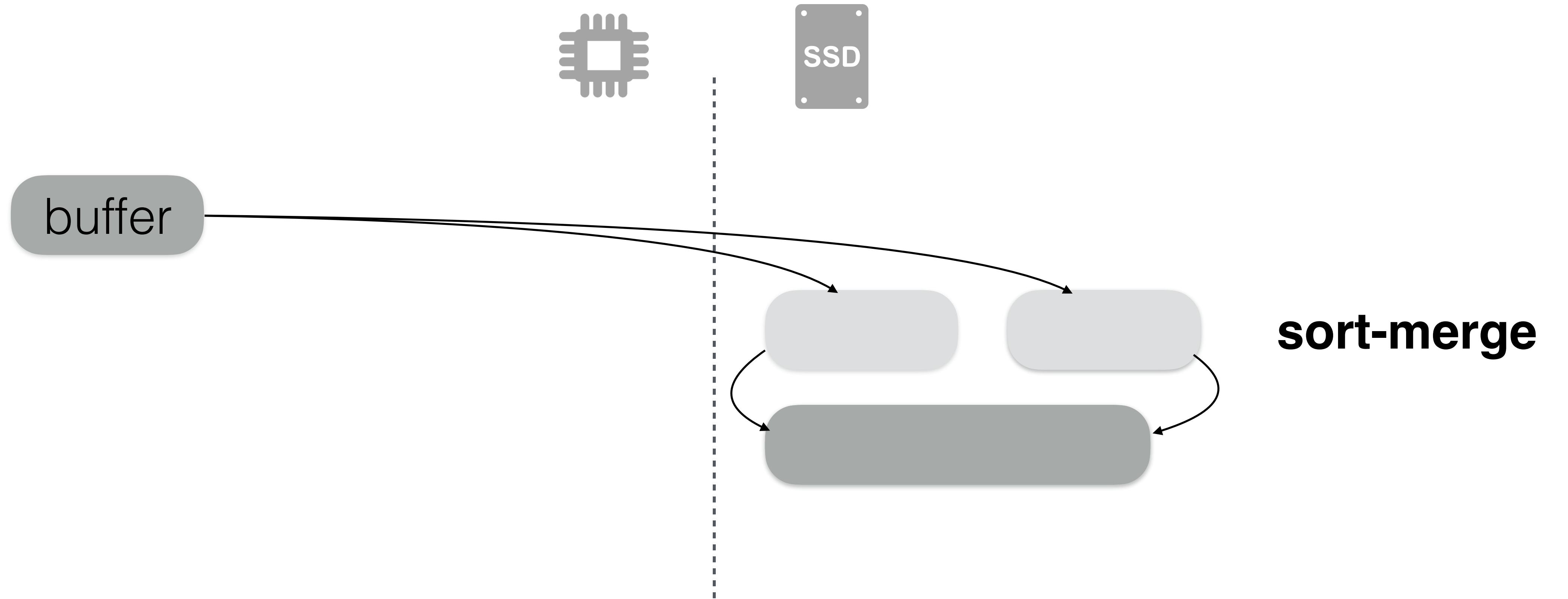
writes

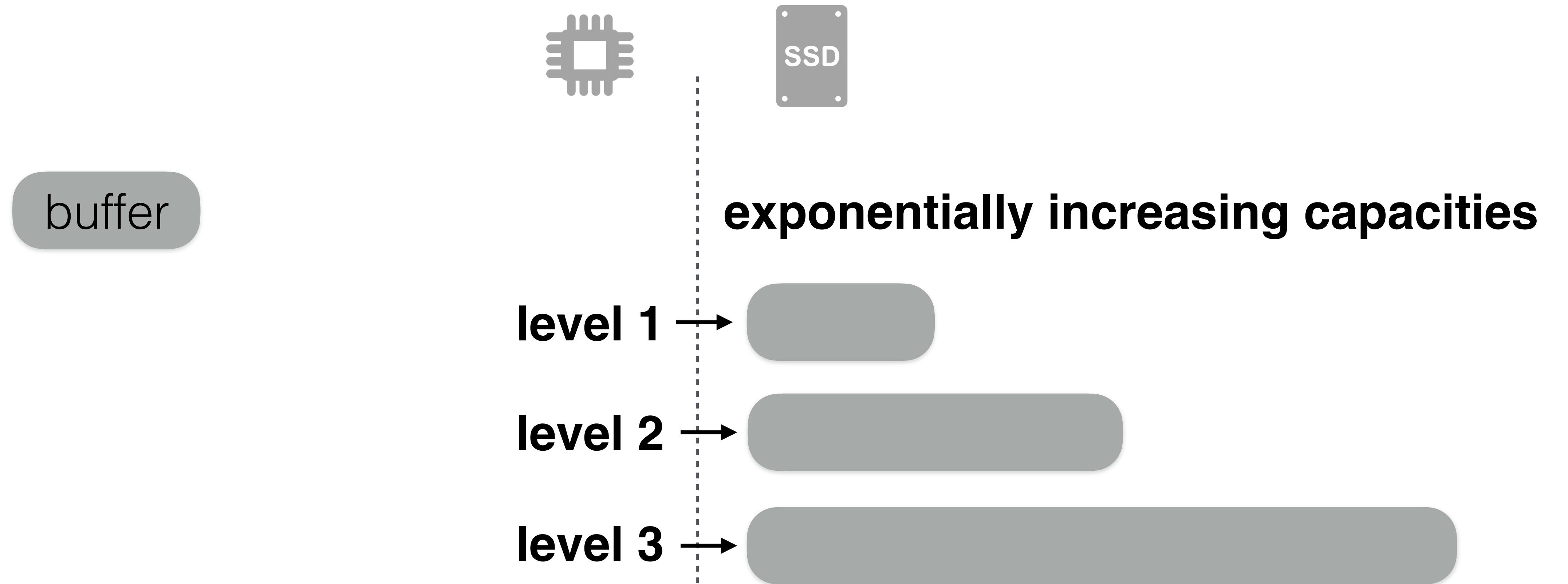


buffer

Background

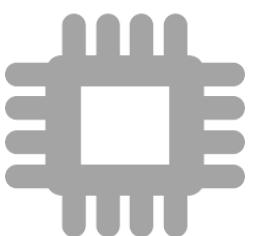






read
key X

buffer

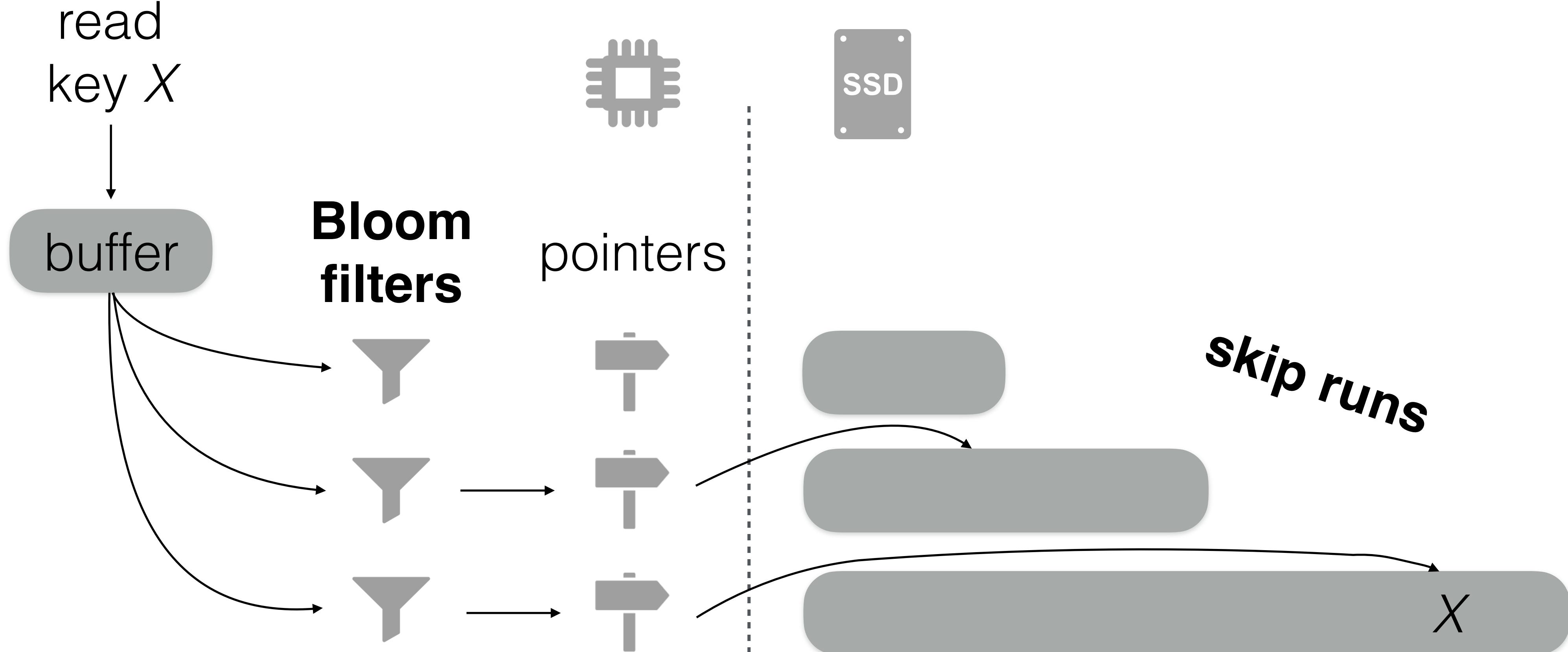


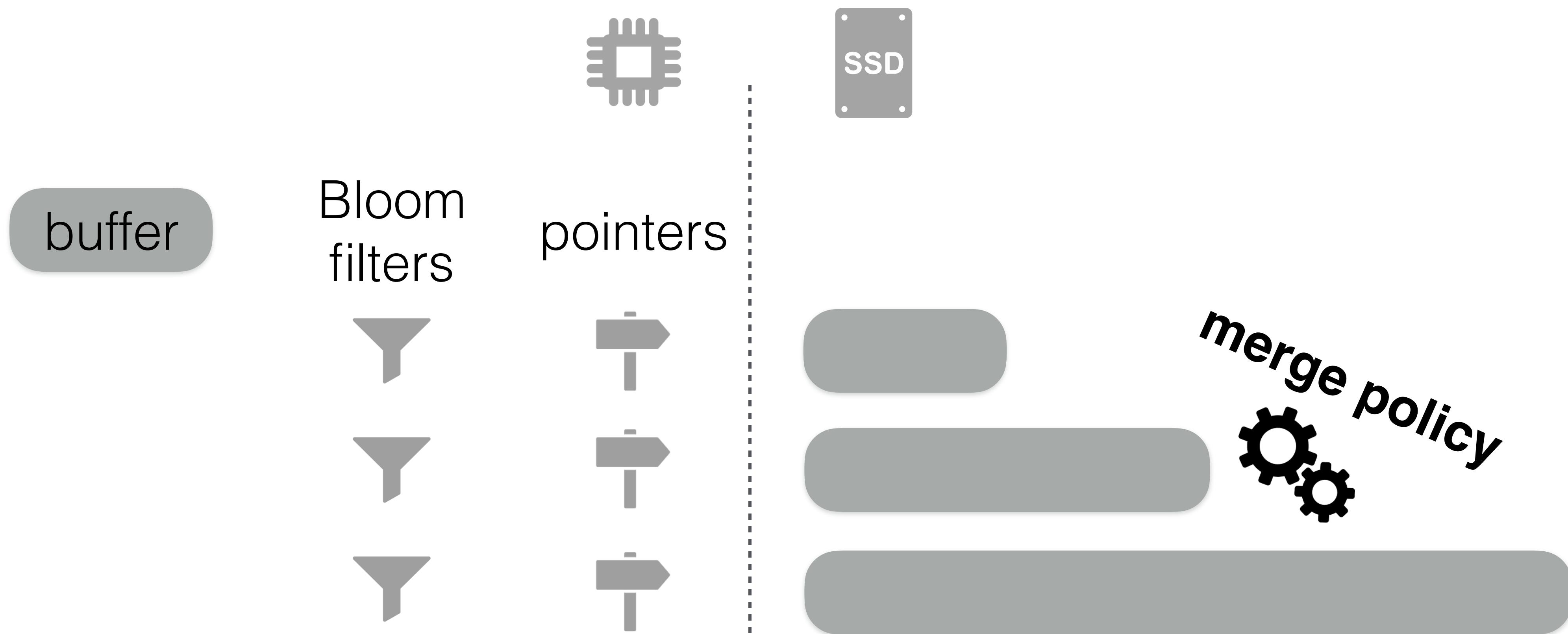
pointers



one I/O per run

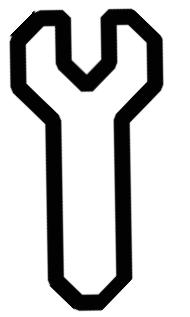
X

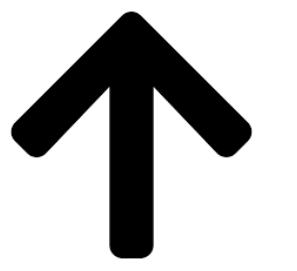




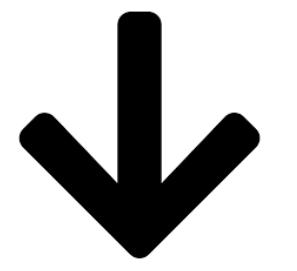


merge policy





merge policy



two merge policies



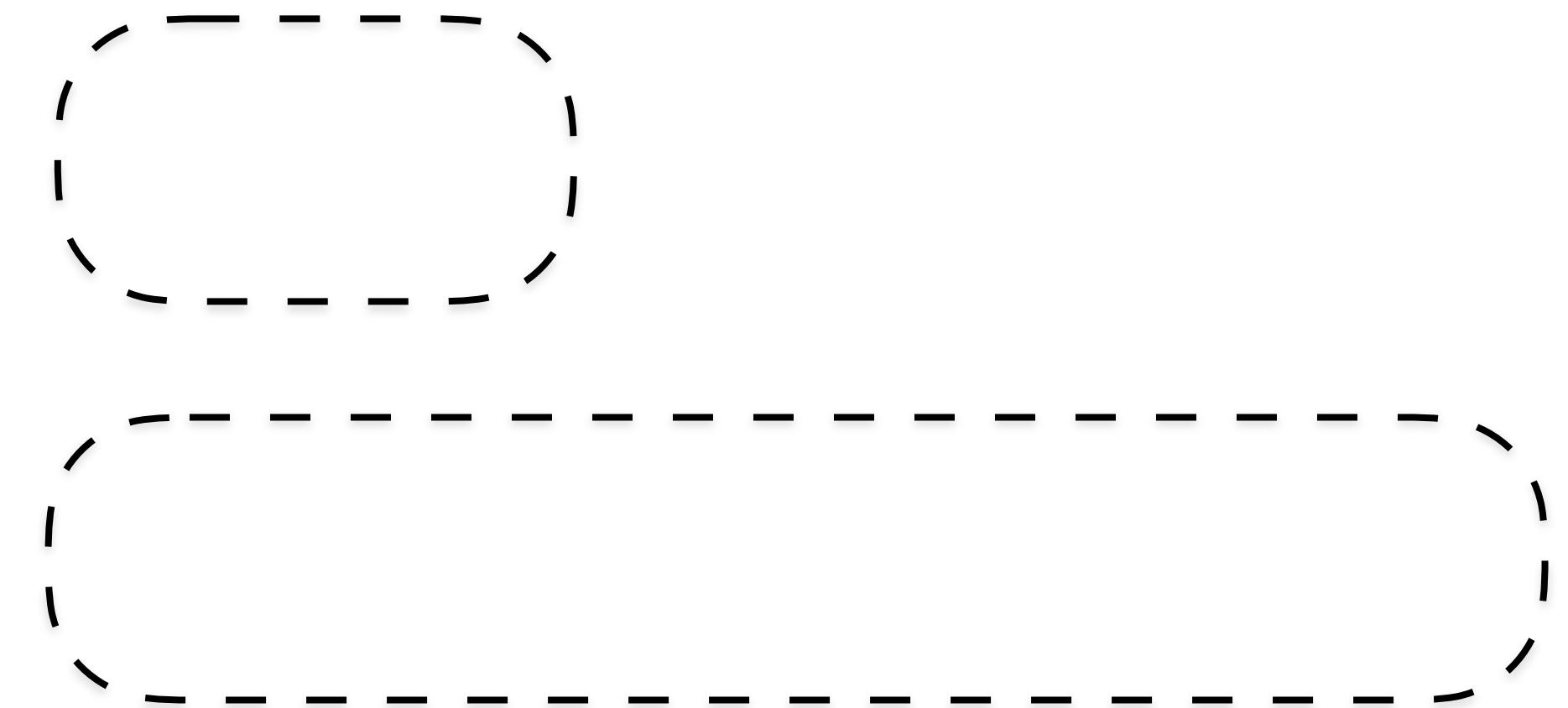
Tiering

Leveling

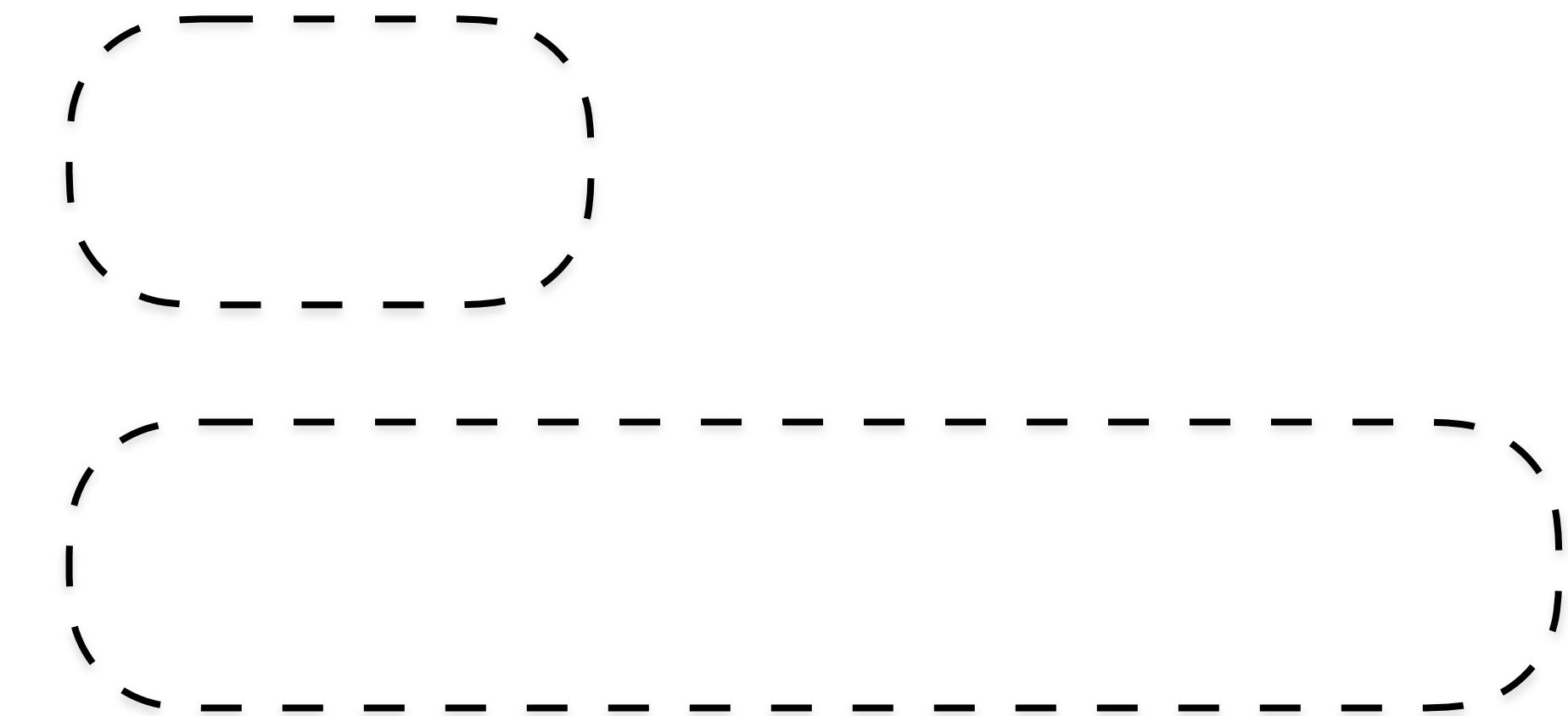
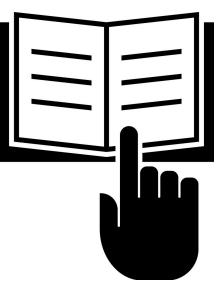




Tiering

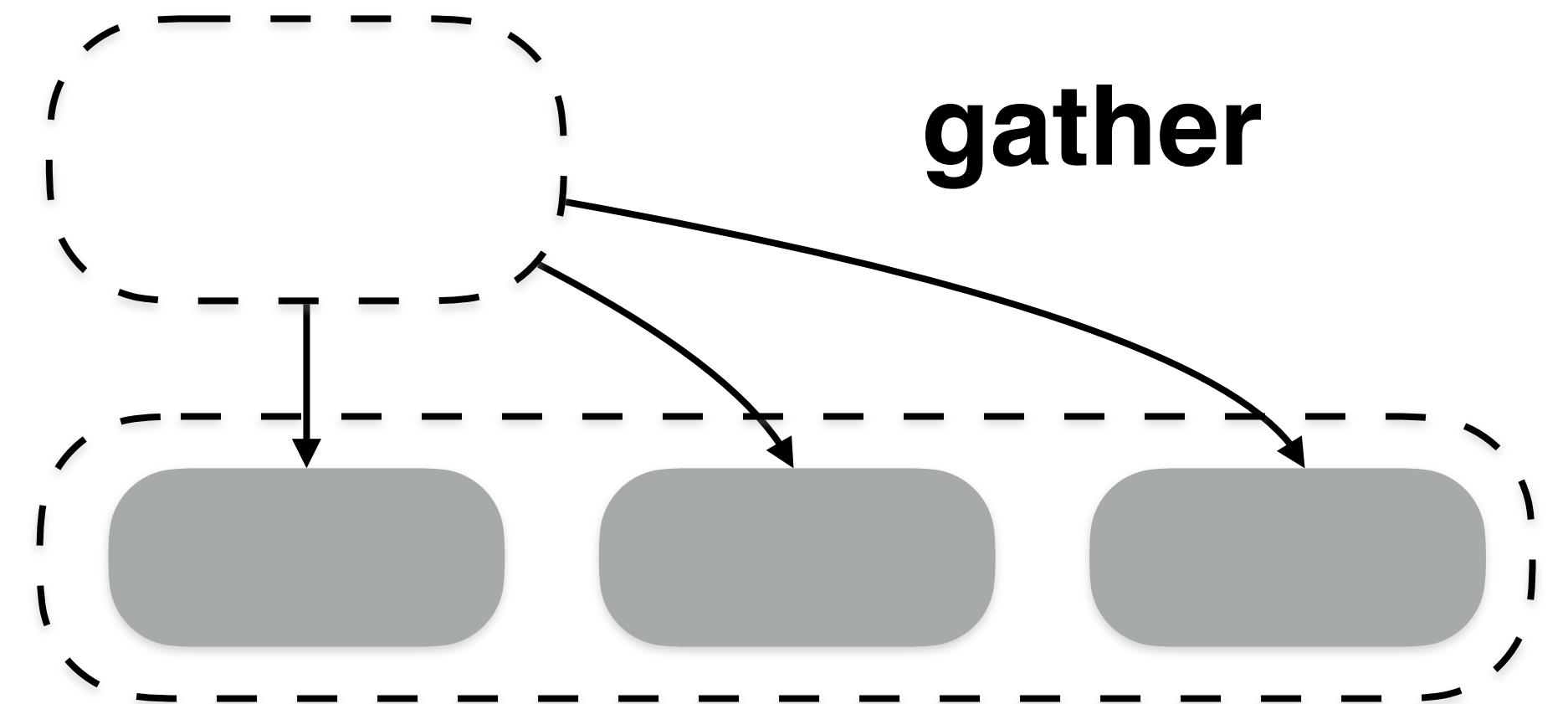


Leveling

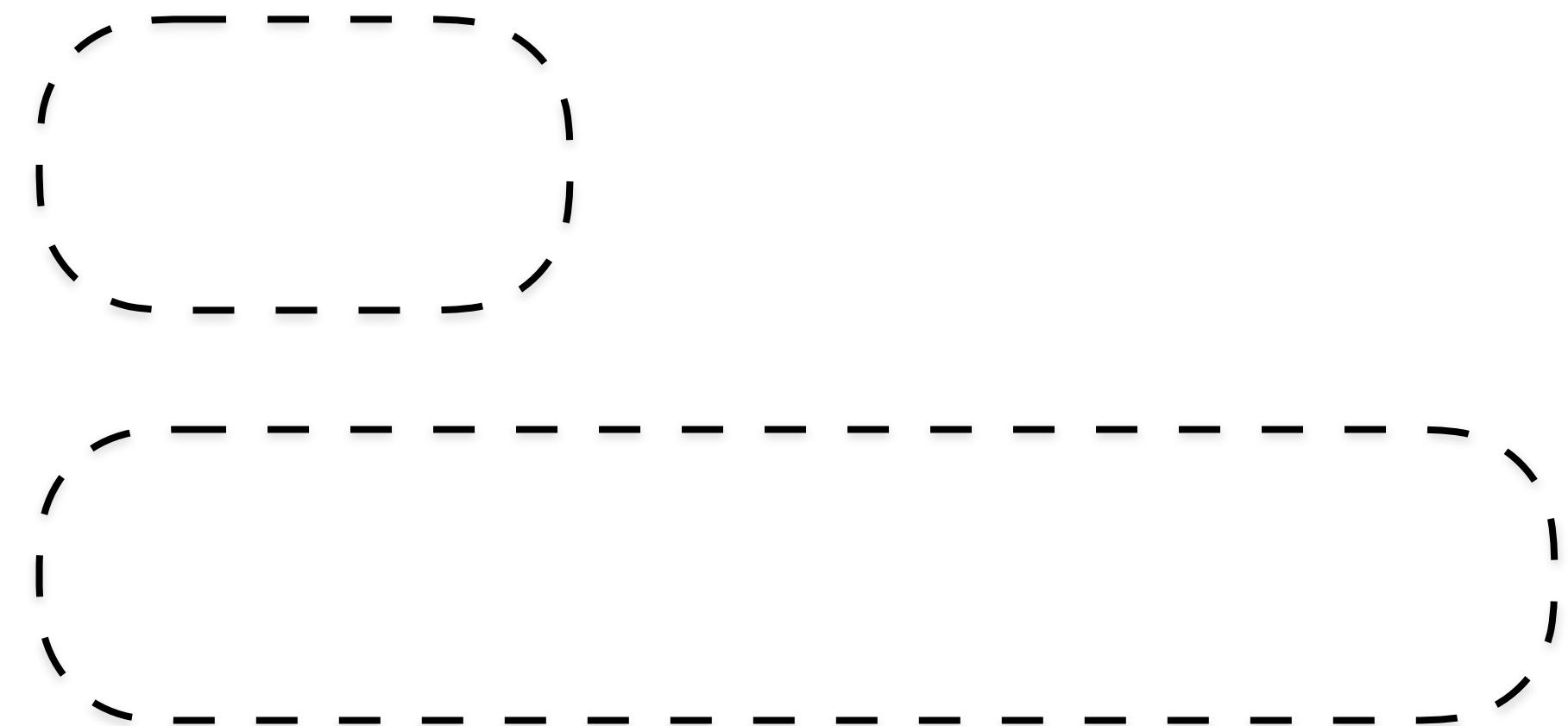
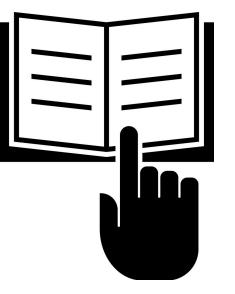




Tiering

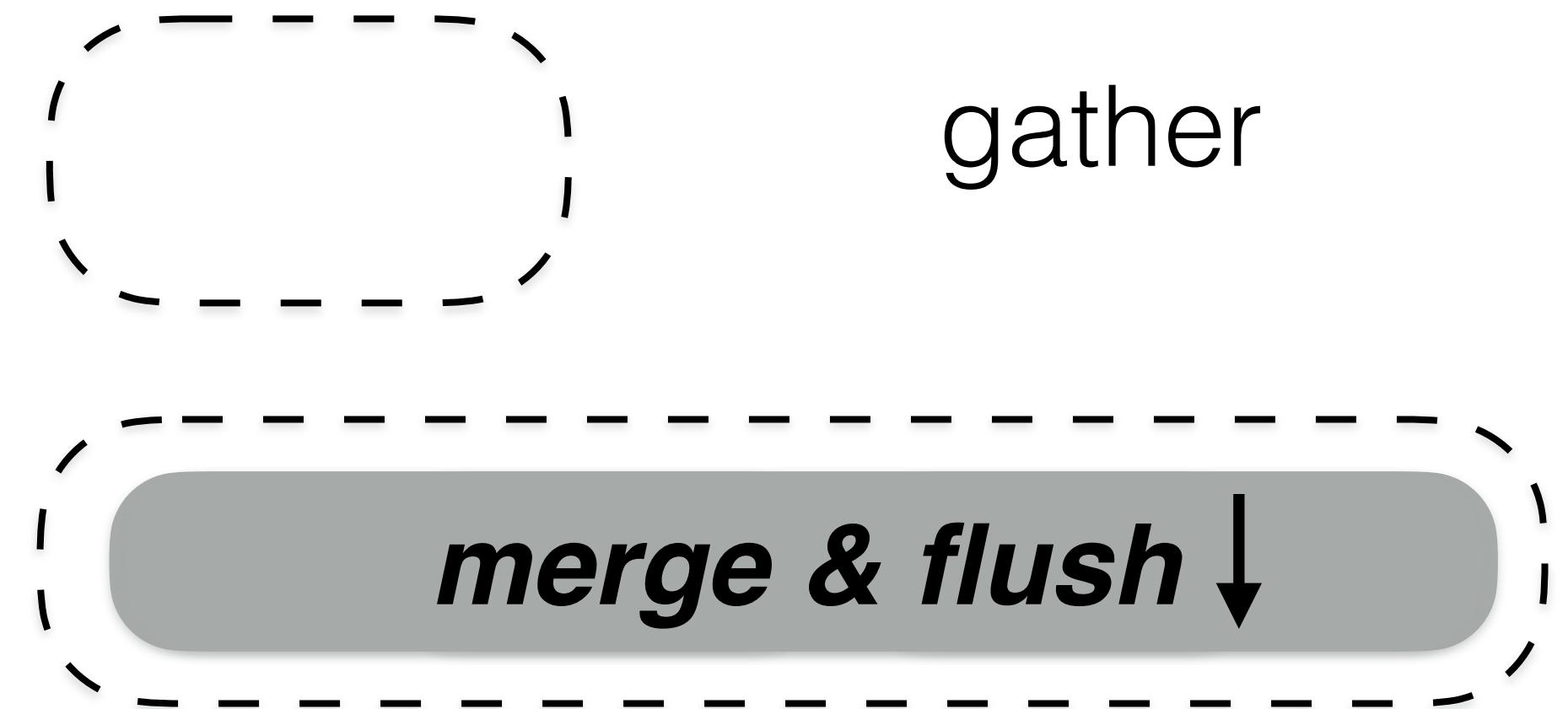


Leveling

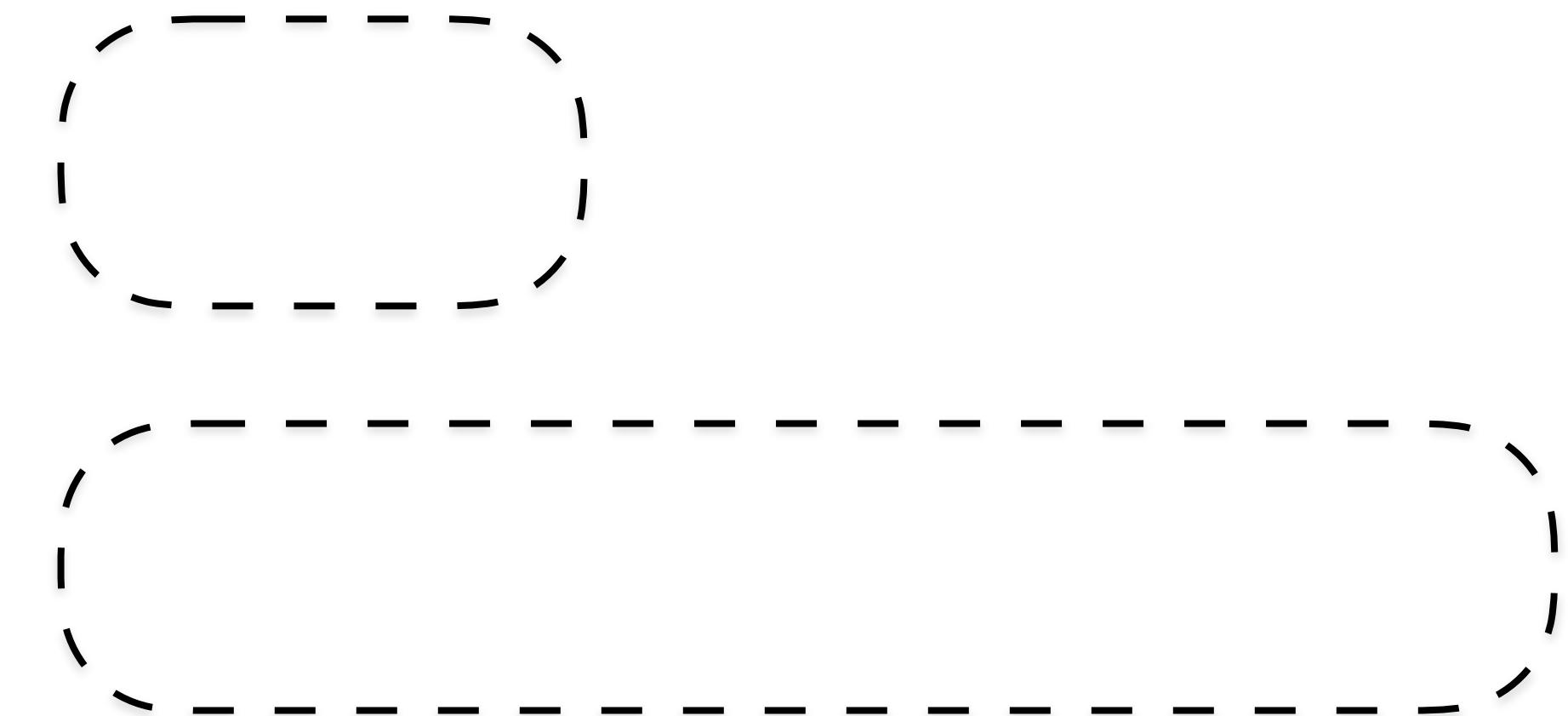
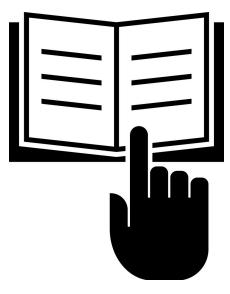




Tiering

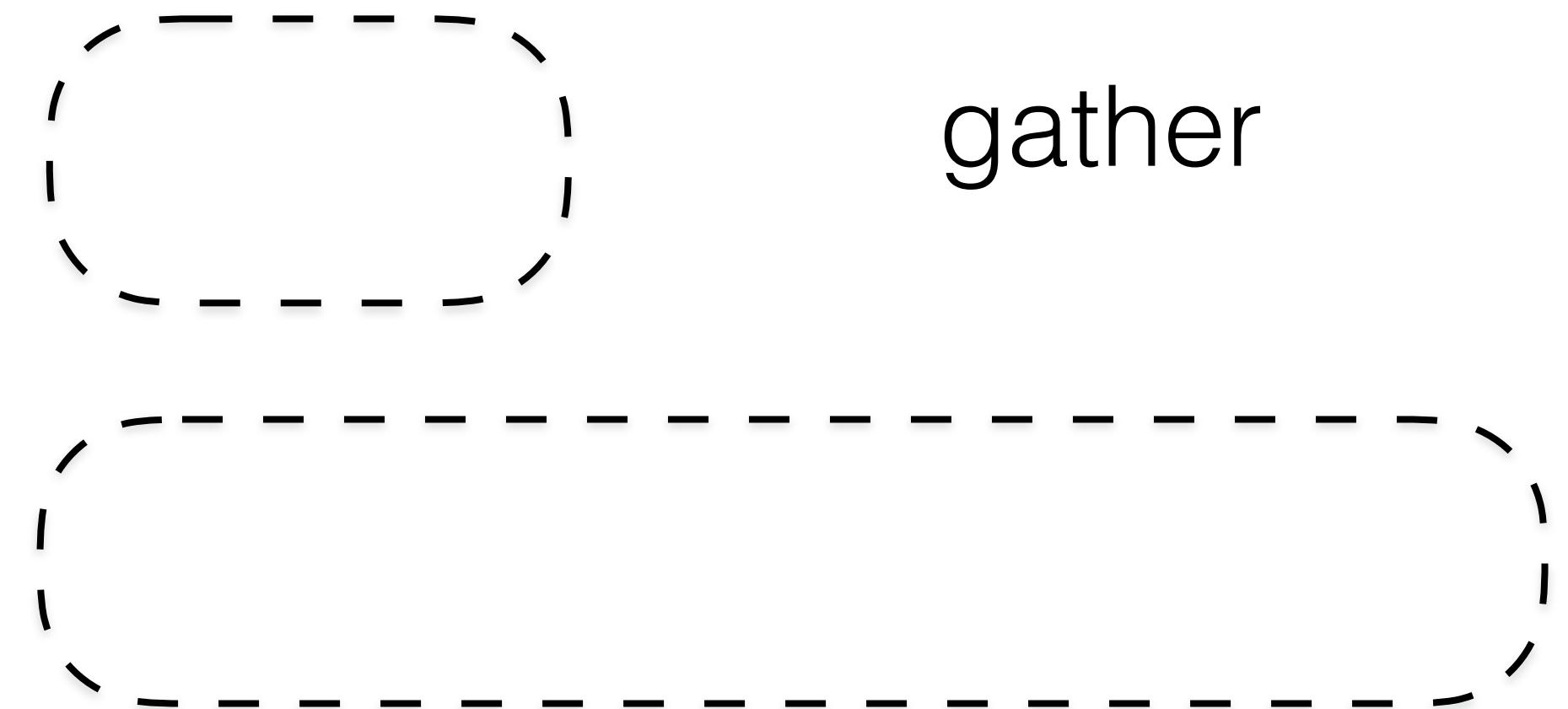


Leveling

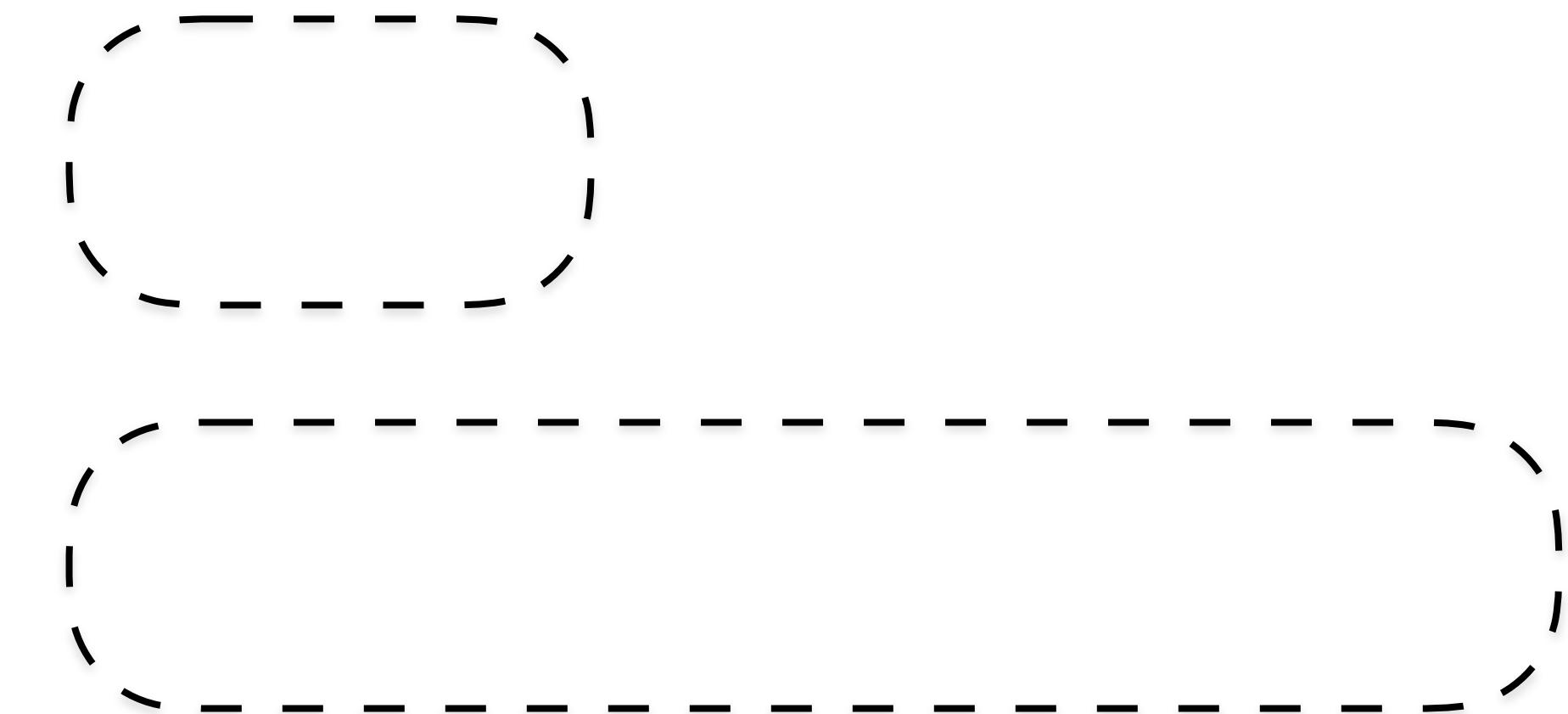




Tiering

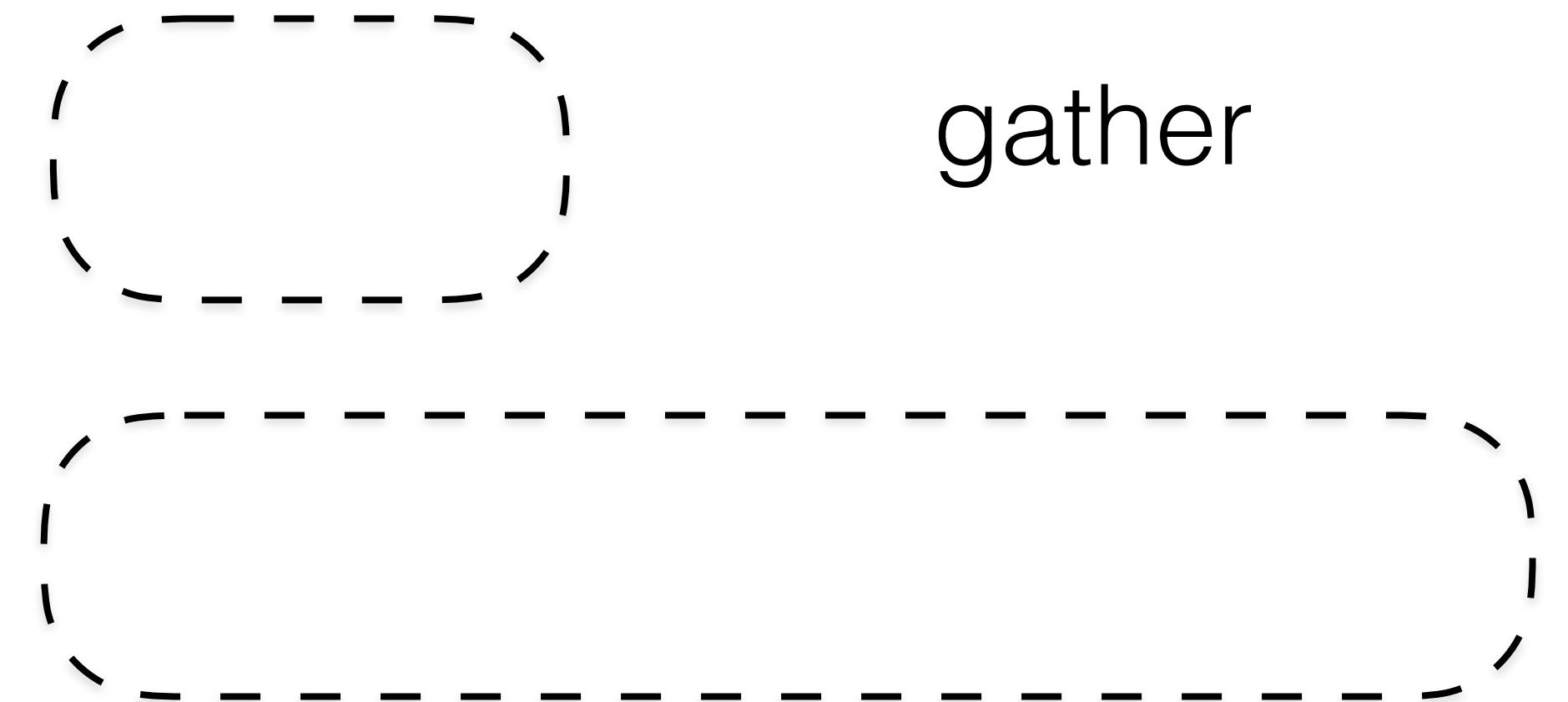


Leveling

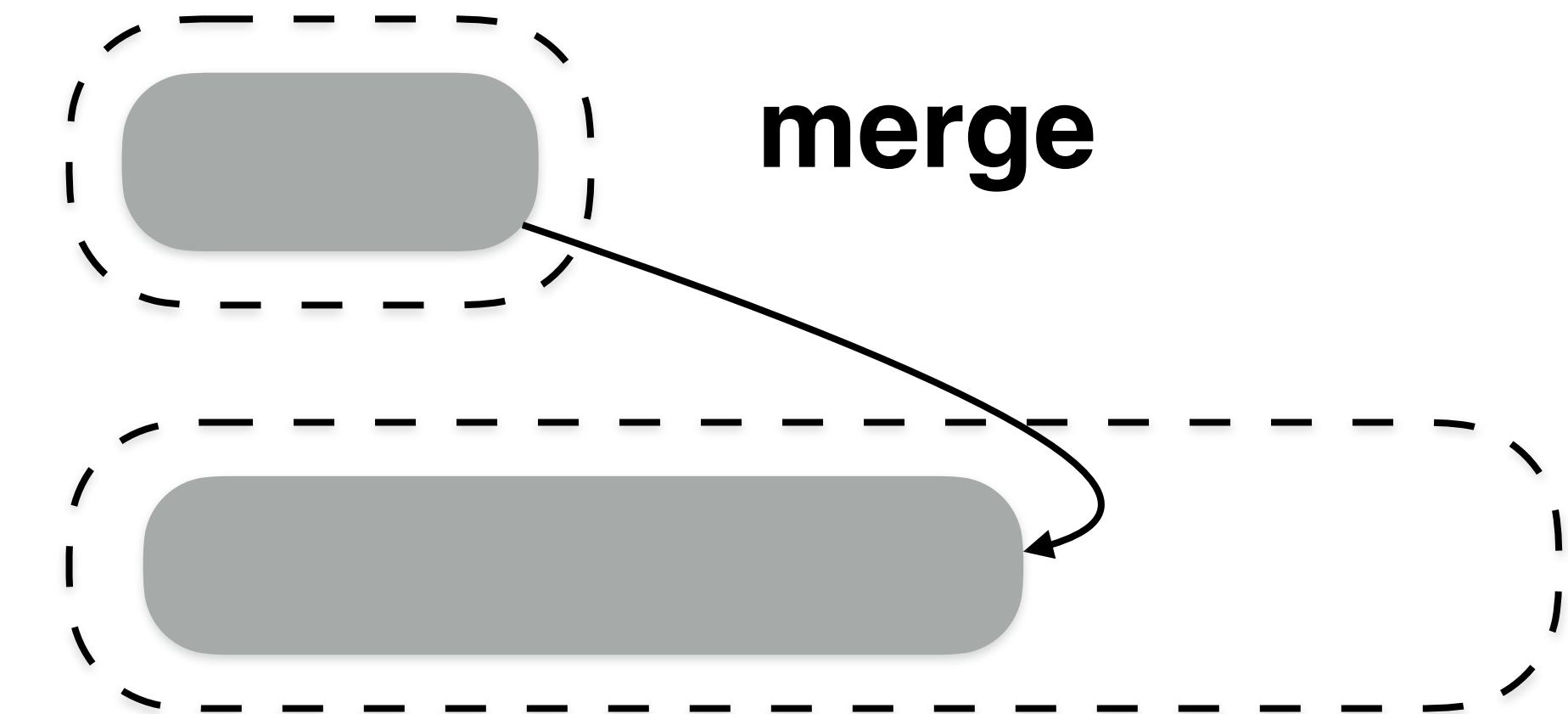
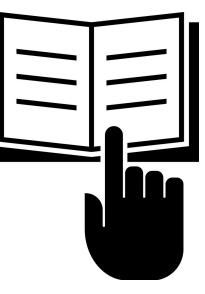




Tiering

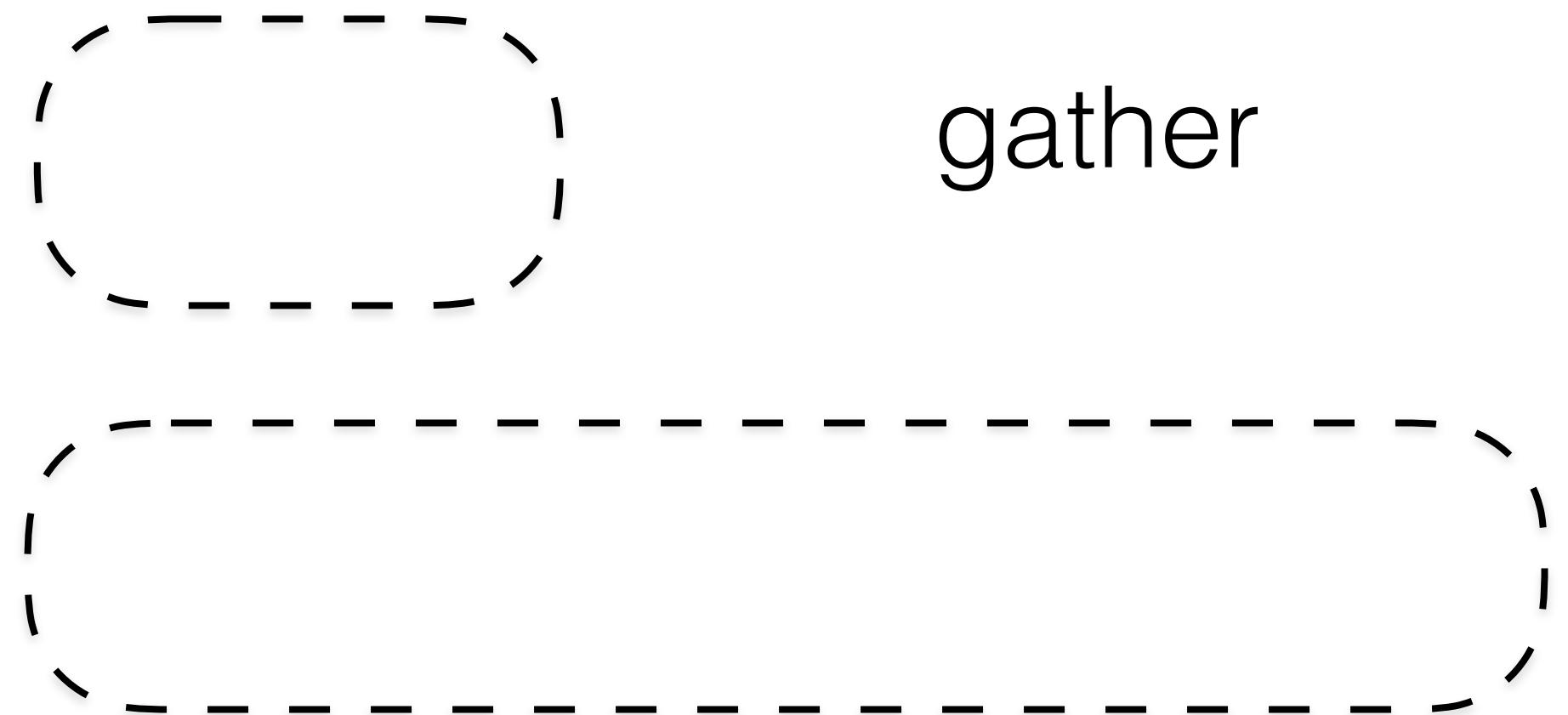


Leveling

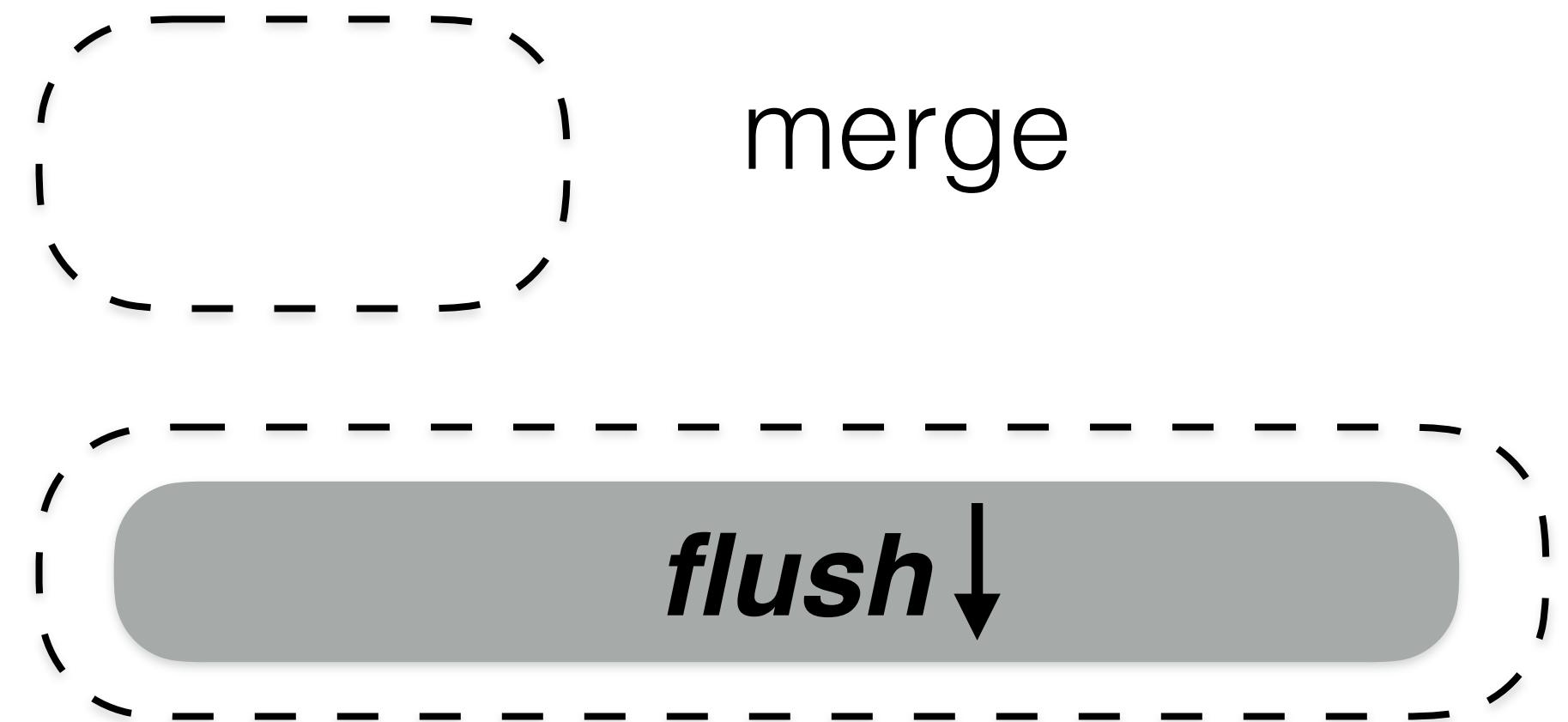




Tiering

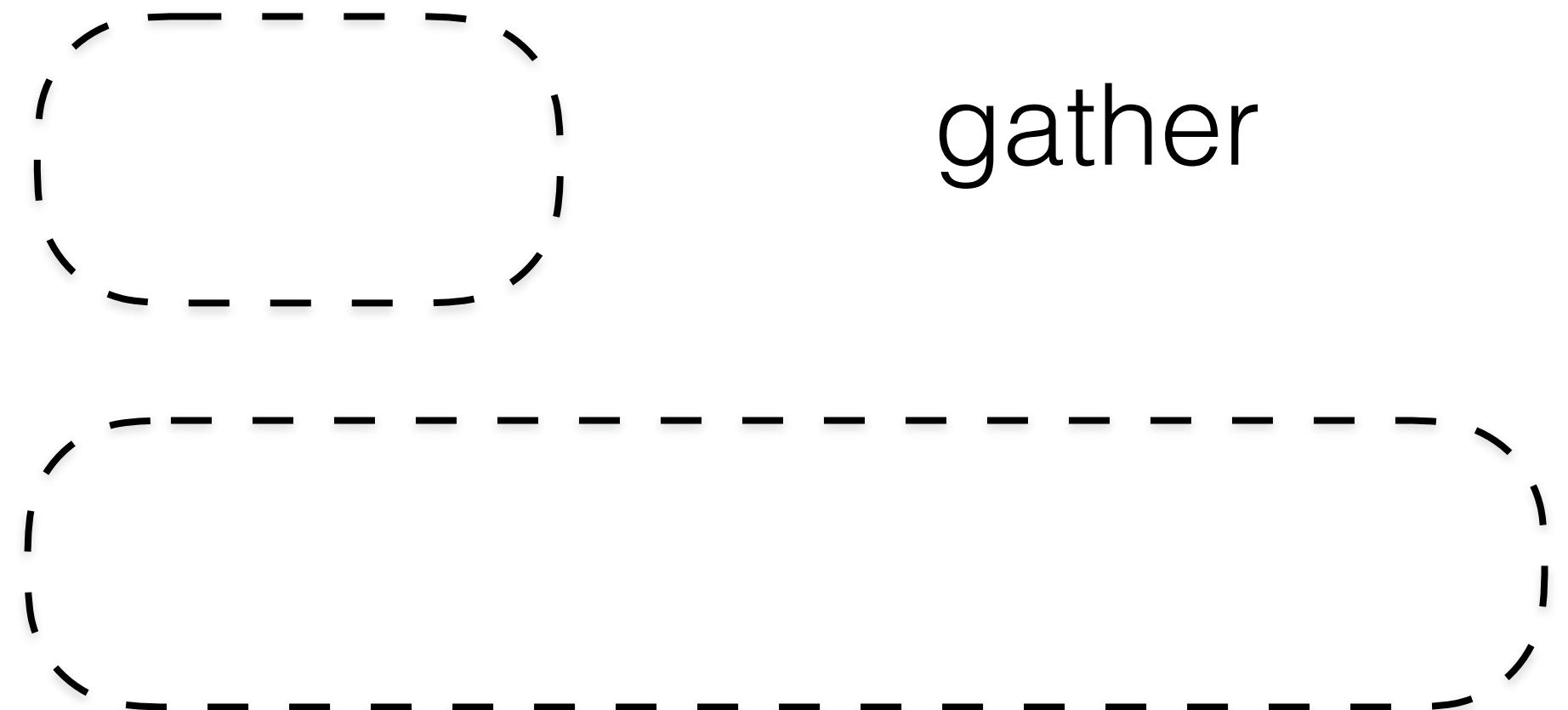


Leveling

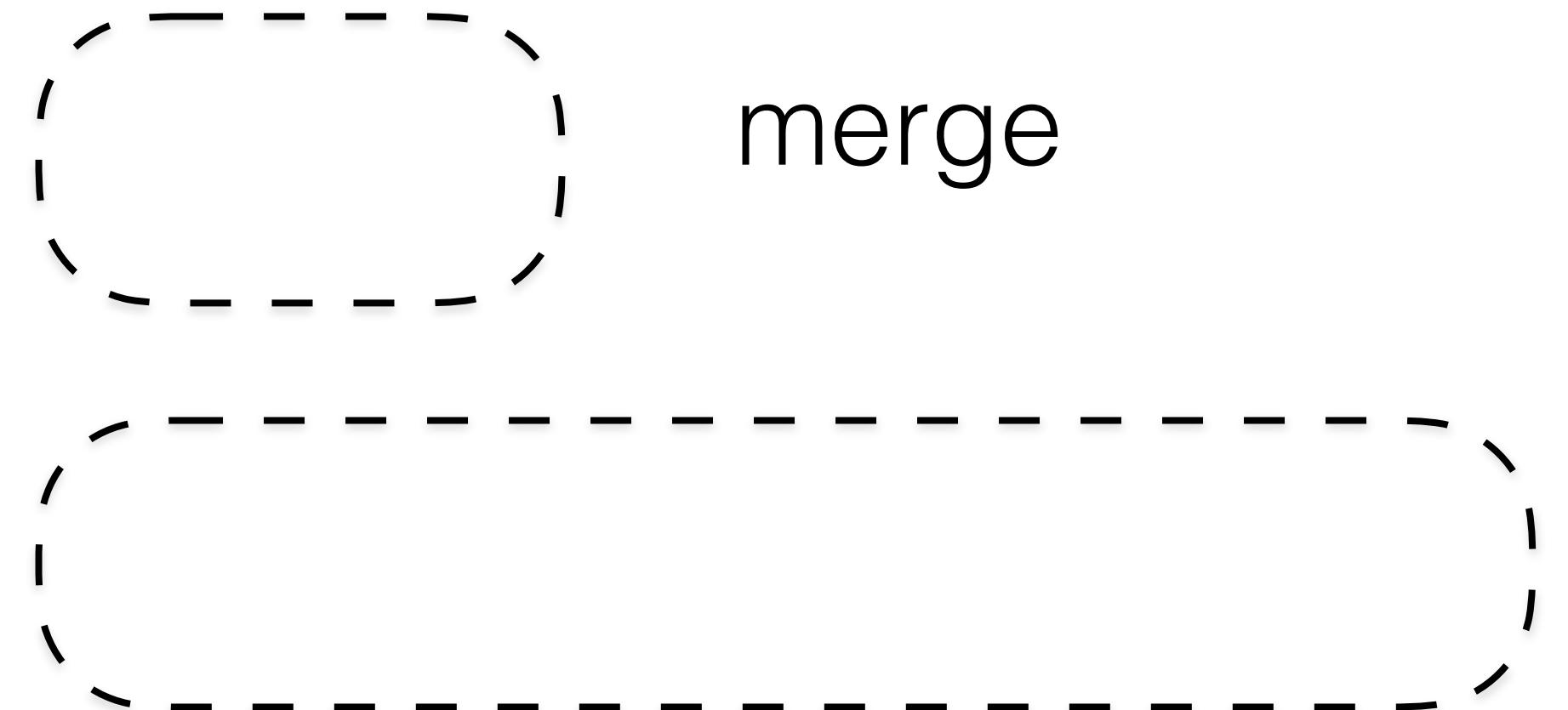




Tiering



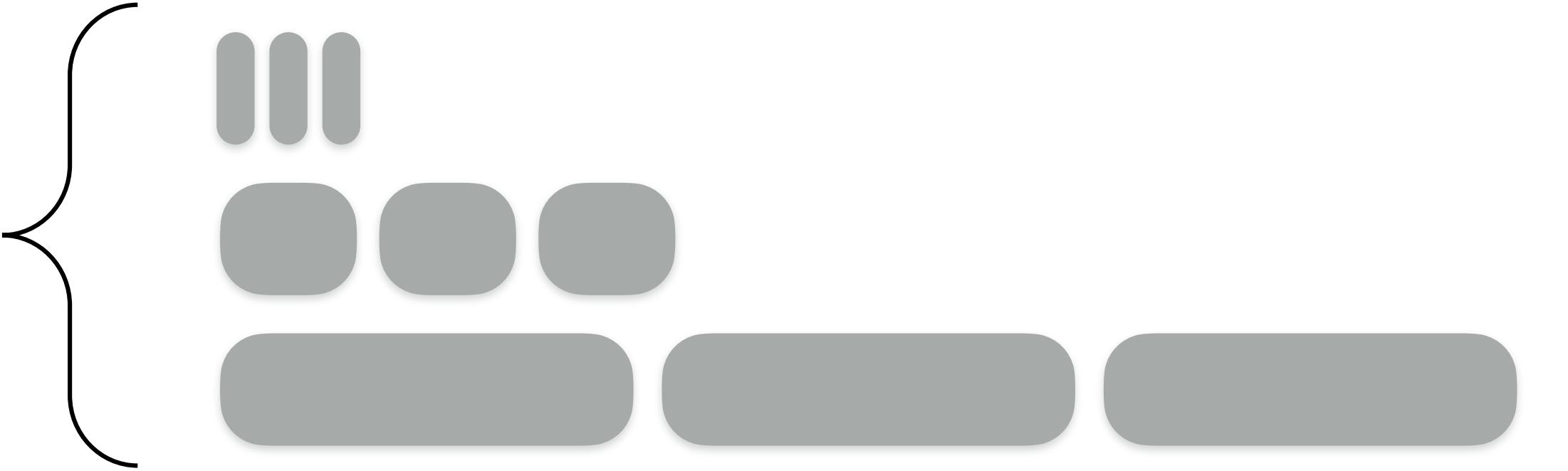
Leveling





Tiering

$\log_R(N)$



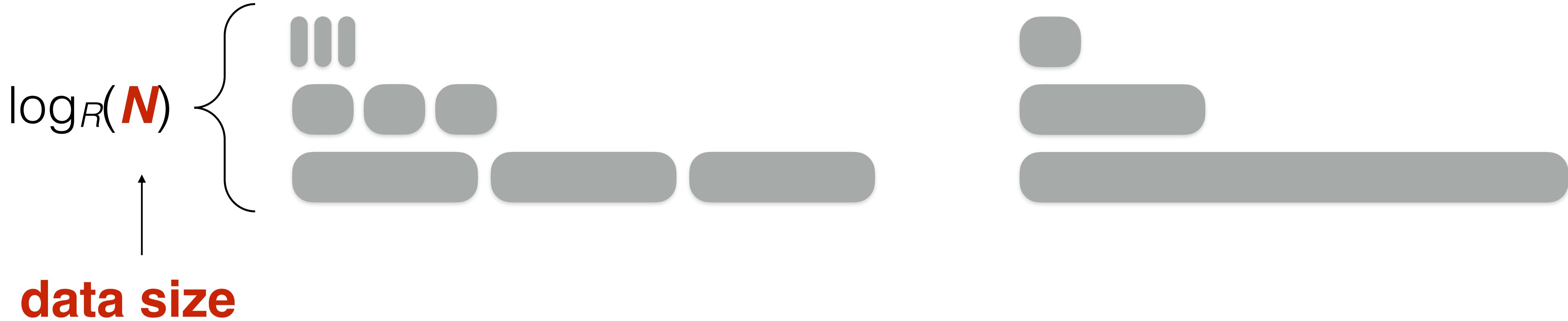
Leveling





Tiering

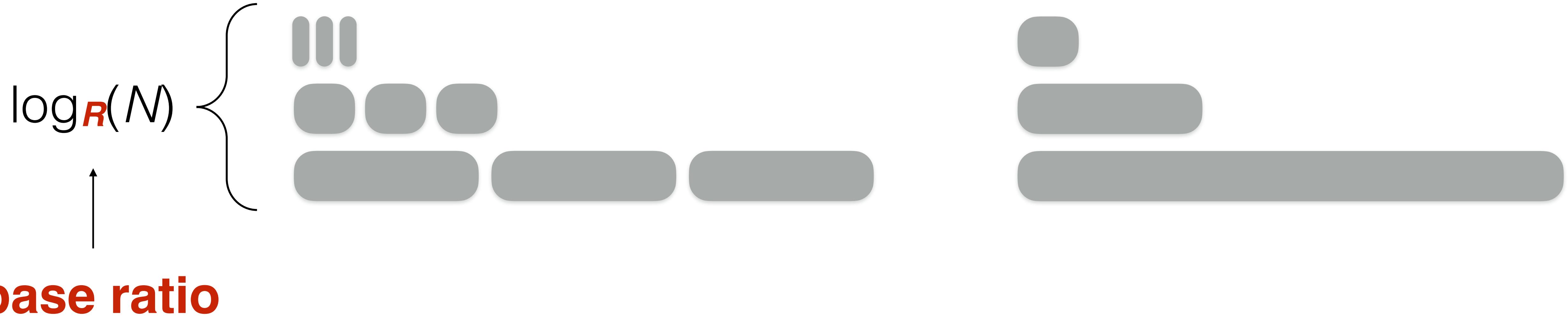
Leveling





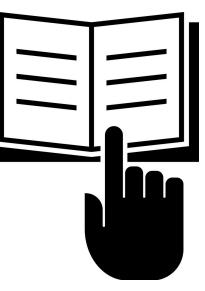
Tiering

Leveling





Tiering

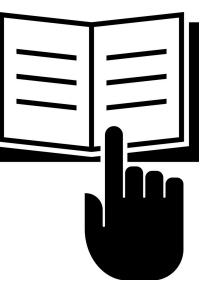


Leveling

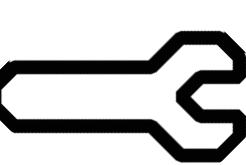




Tiering



Leveling

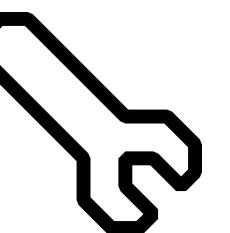
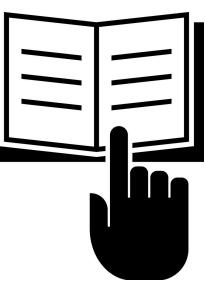


base ratio R



Tiering

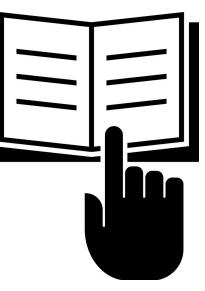
Leveling



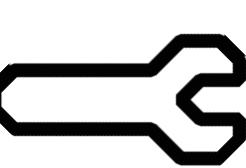
base ratio $R \searrow$



Tiering



Leveling



base ratio R

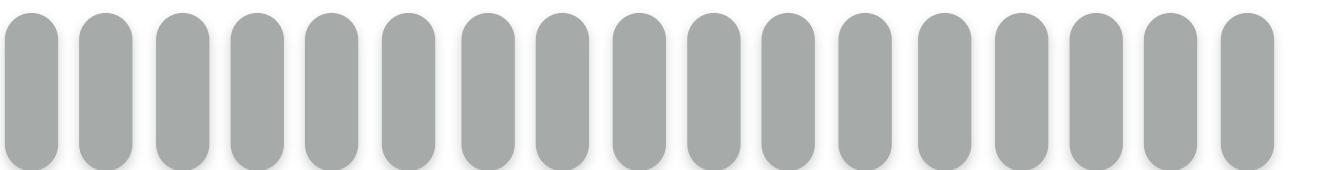


Tiering

Leveling



$O(N)$ runs per level



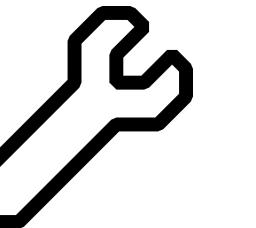
log

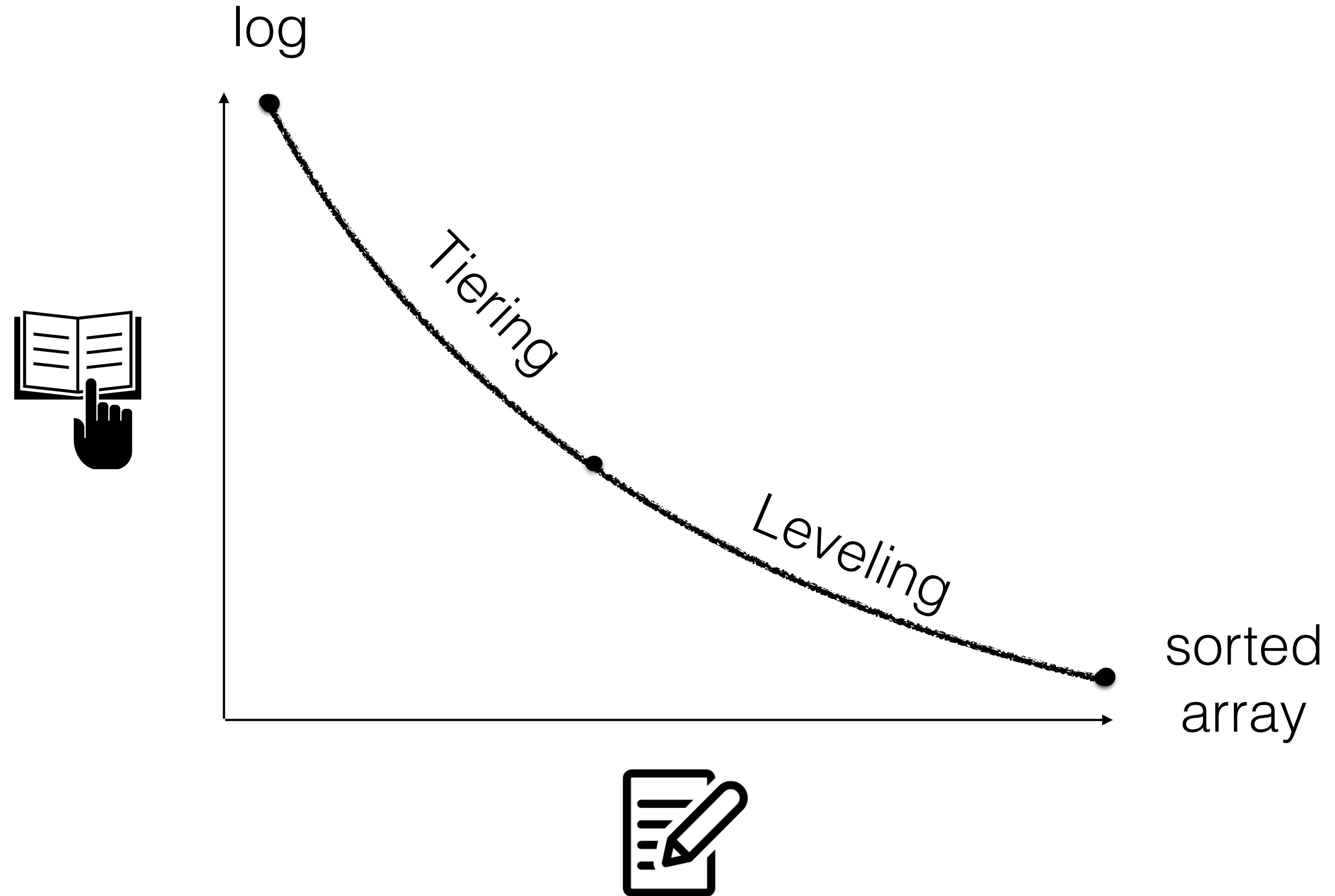
1 run per level

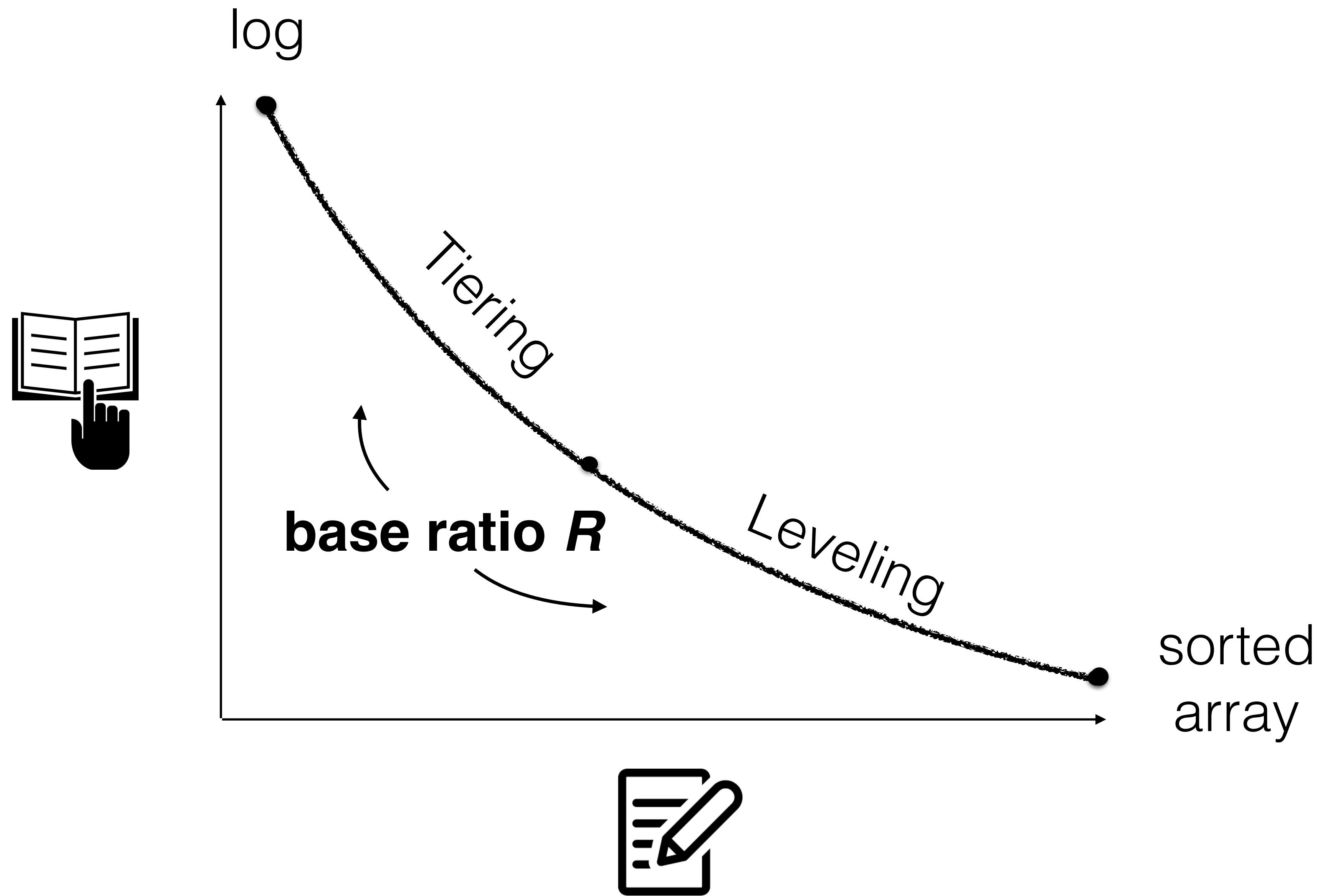


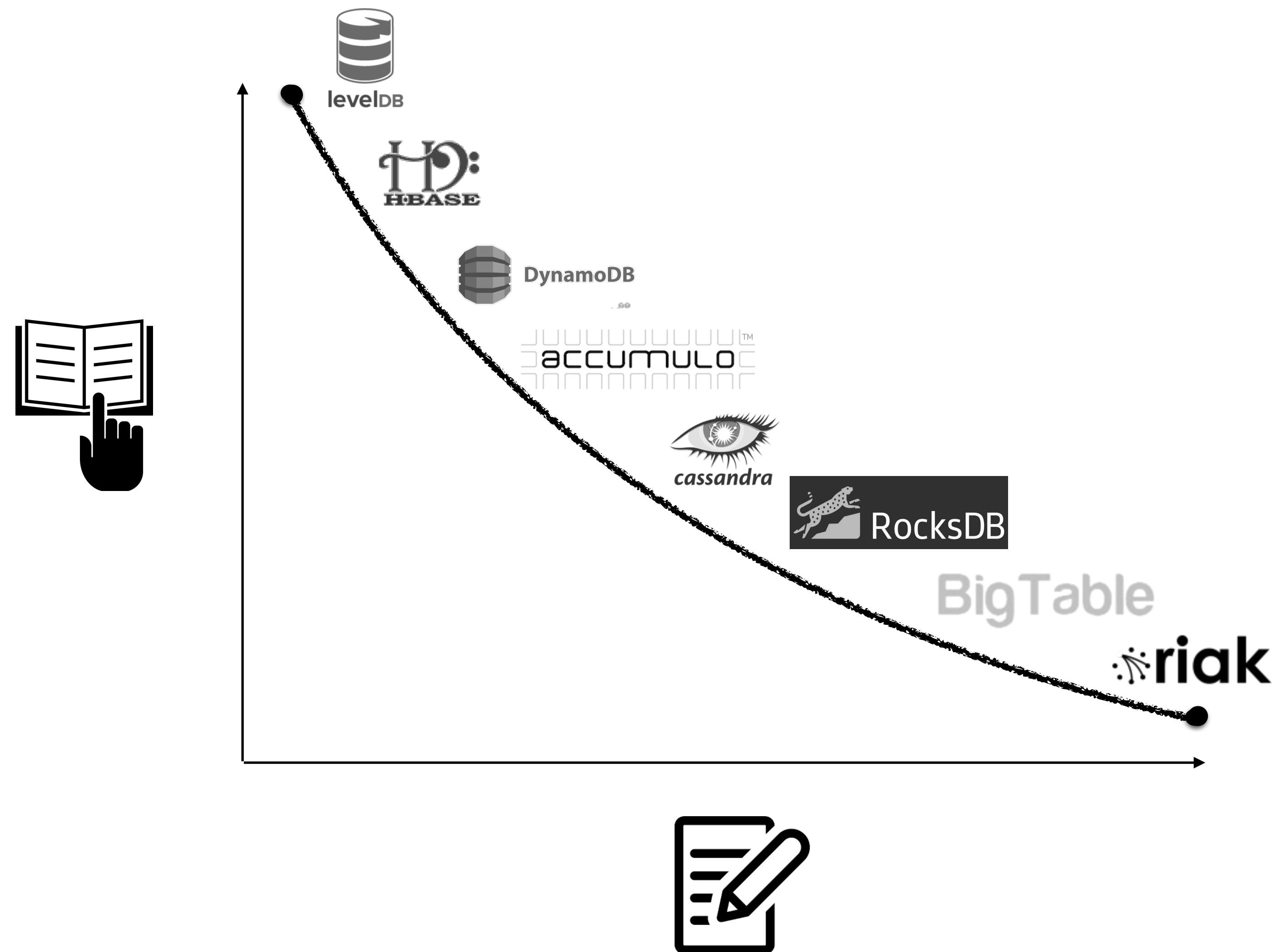
sorted
array

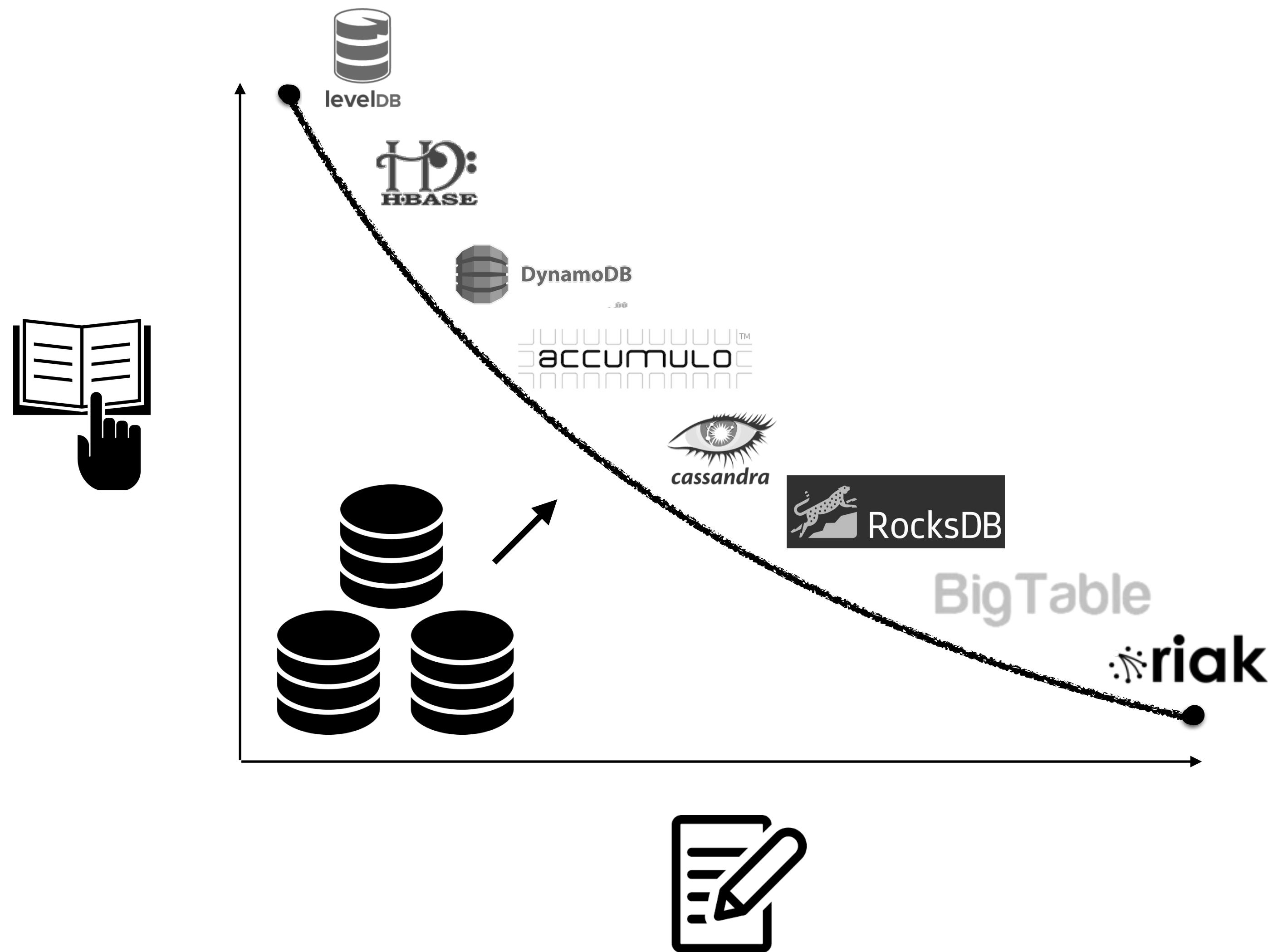
base ratio $R \nearrow$

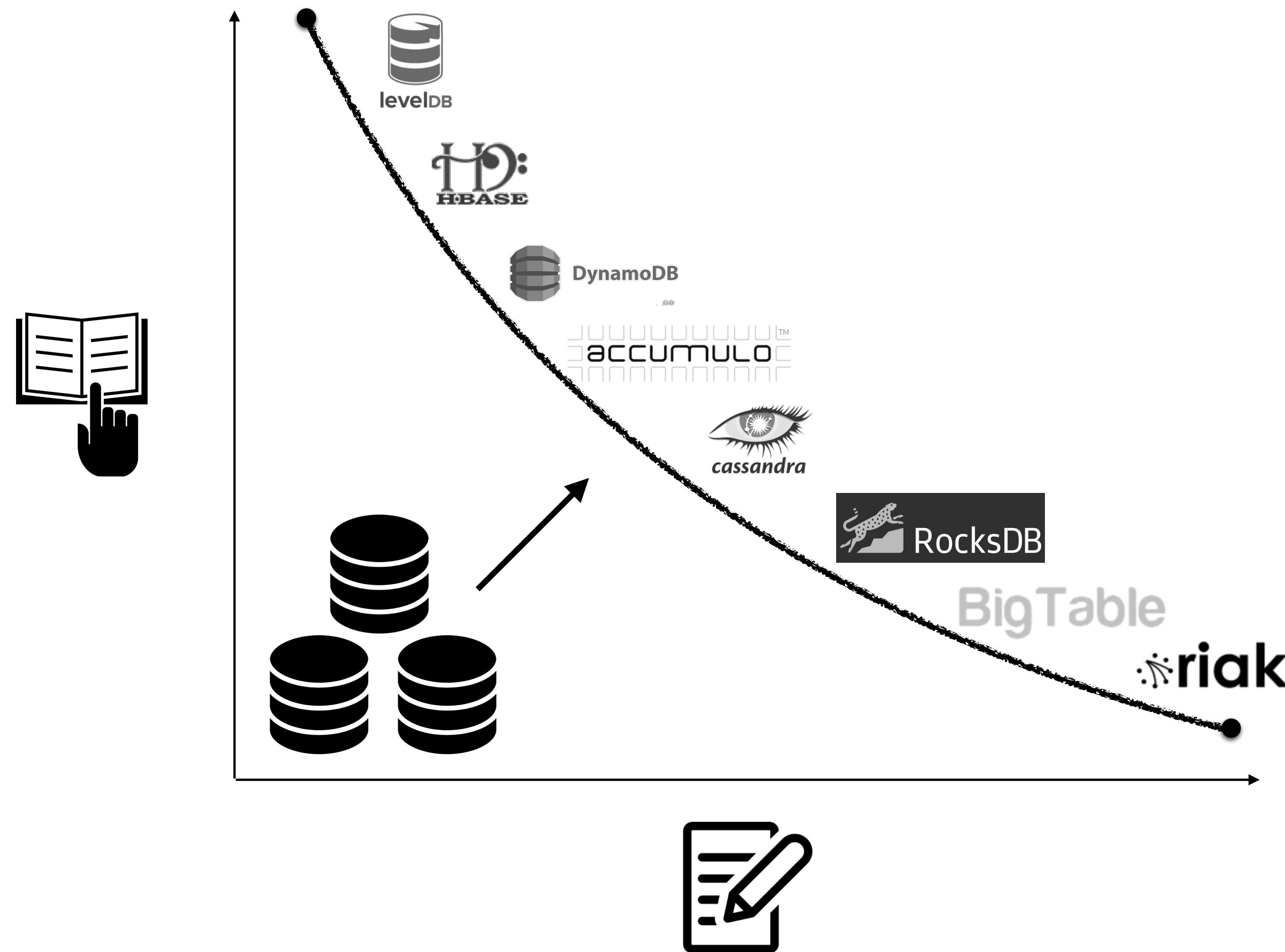


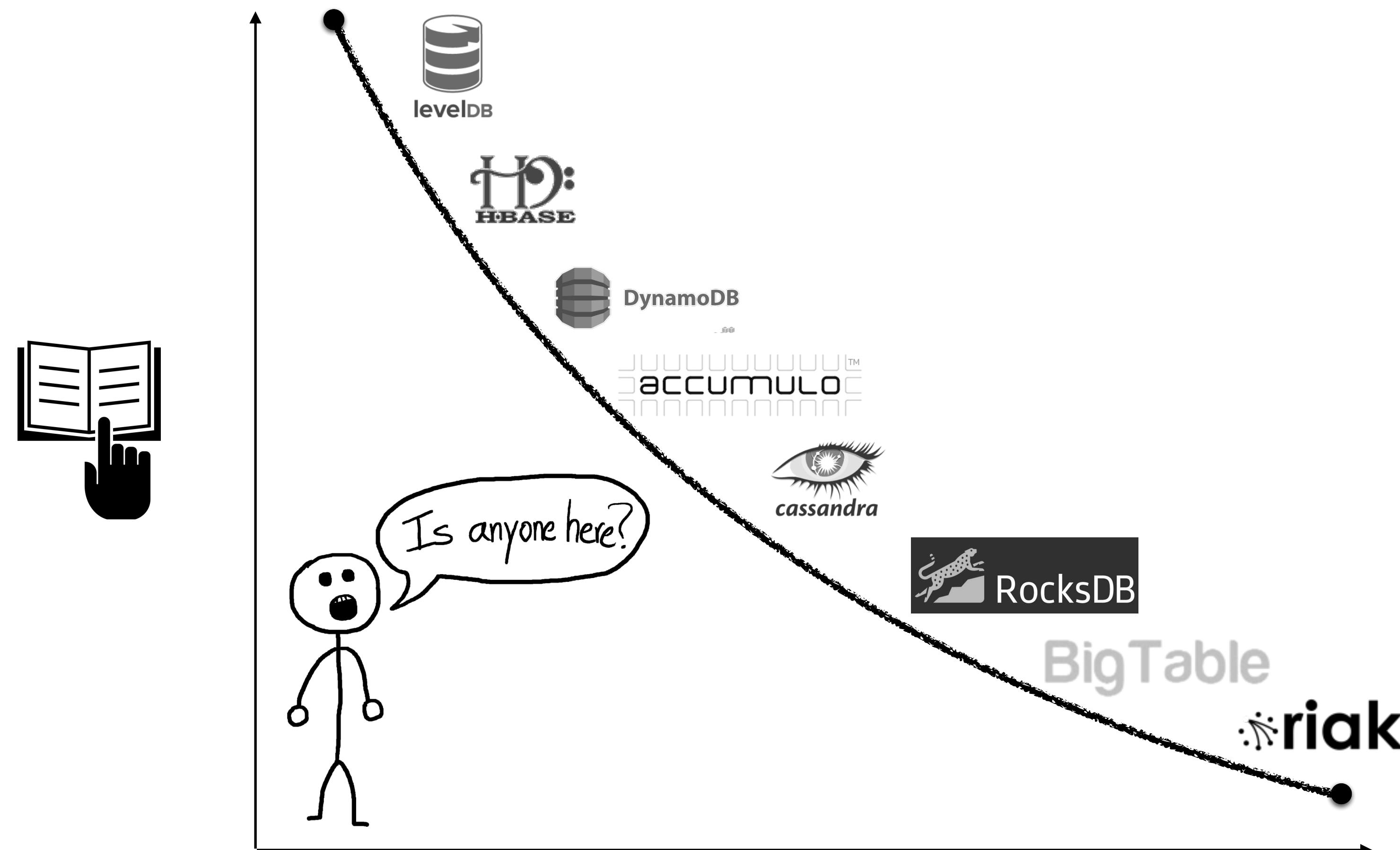


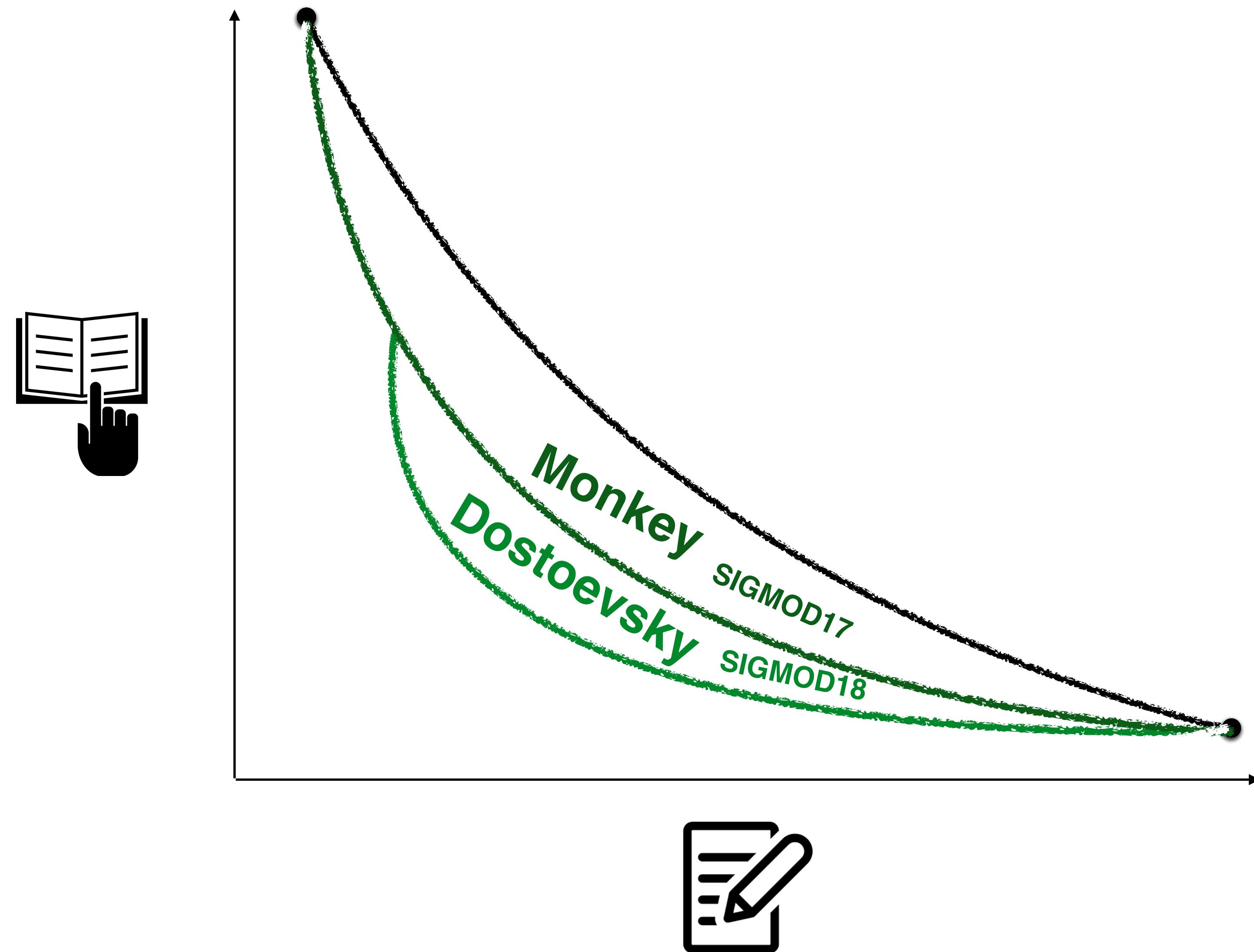


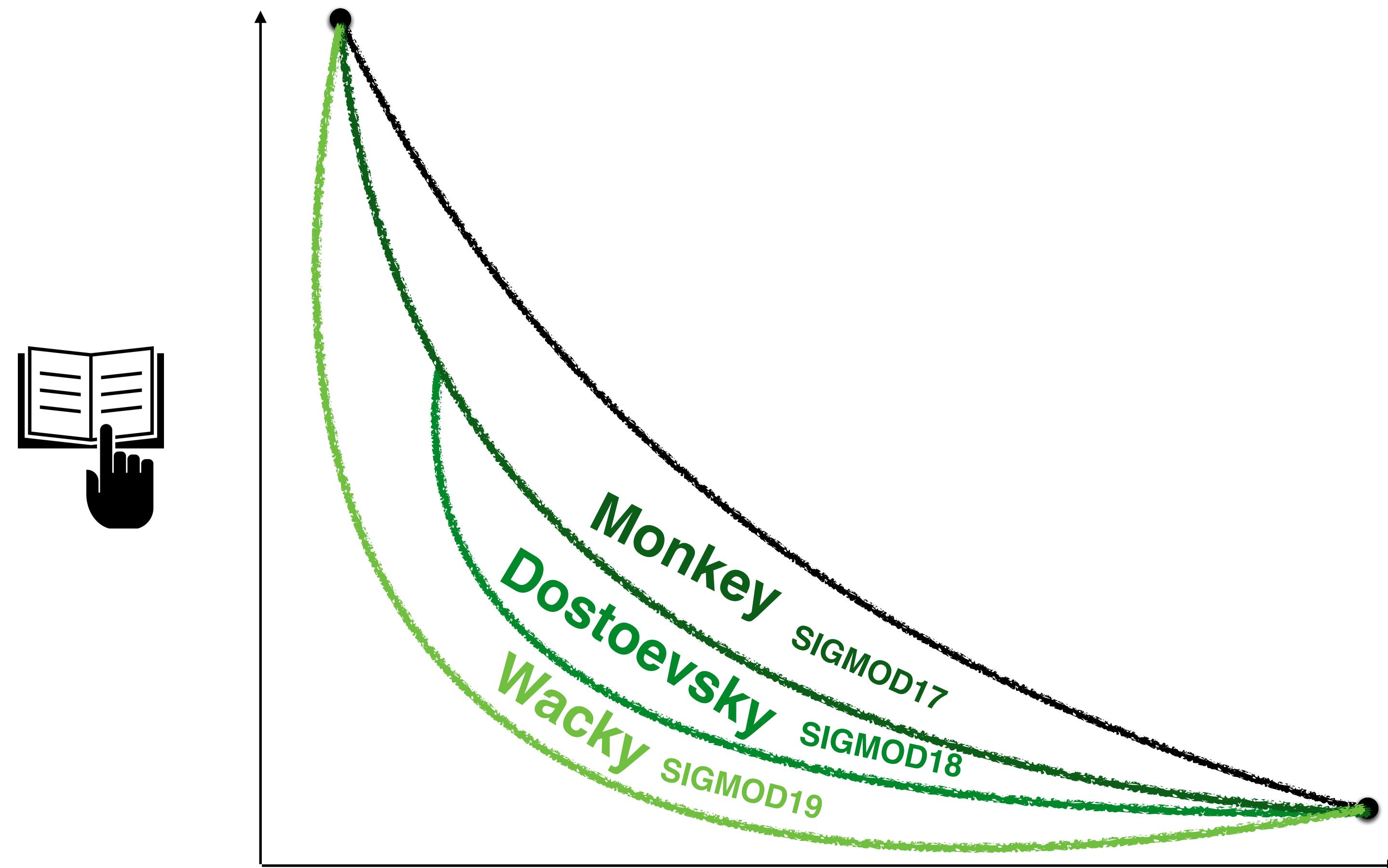












Monkey: Optimal Navigable **Key**-Value Store

SIGMOD17



Monkey: Optimal Navigable Key-Value Store

SIGMOD17

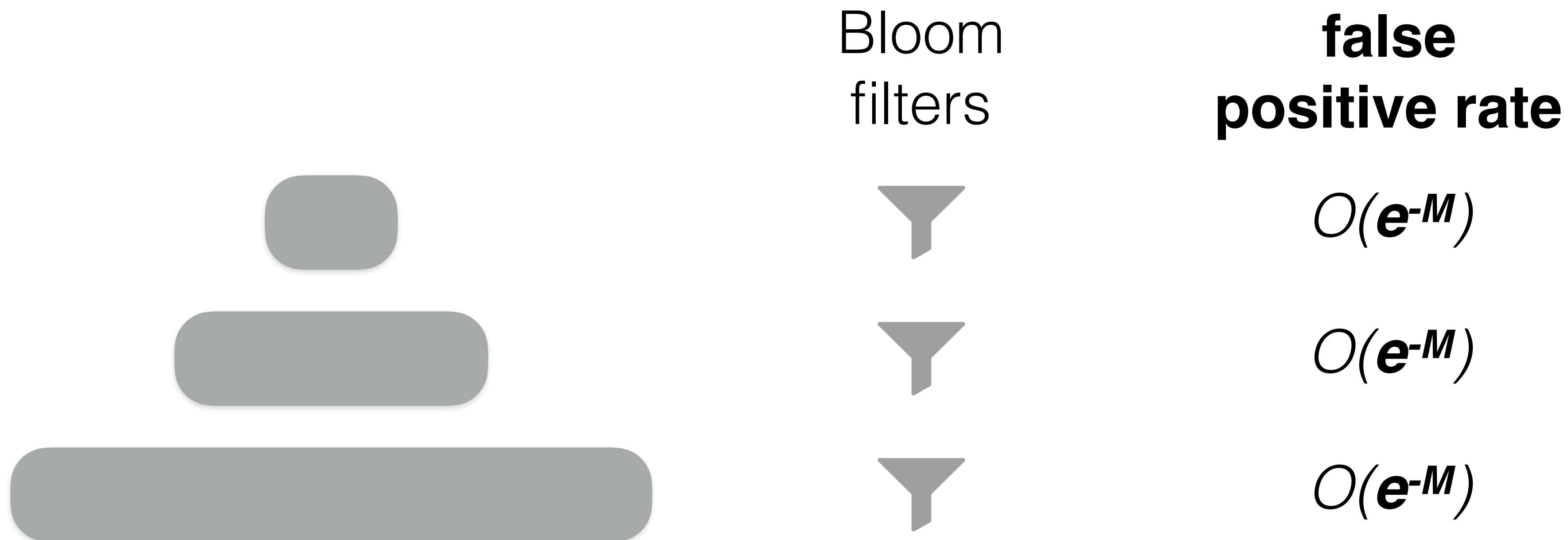


Bloom
filters



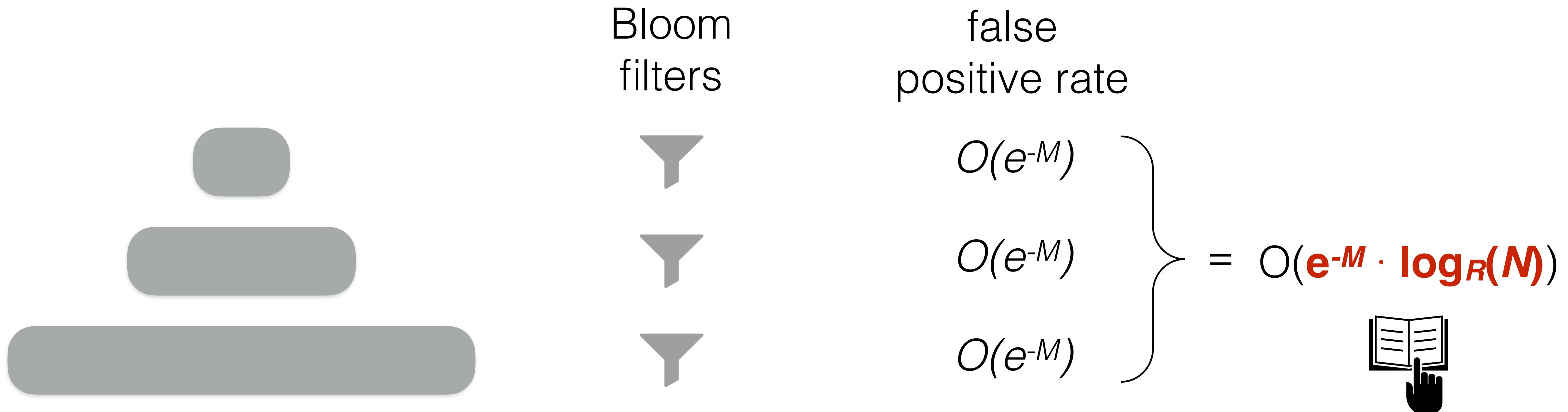
Monkey: Optimal Navigable Key-Value Store

SIGMOD17



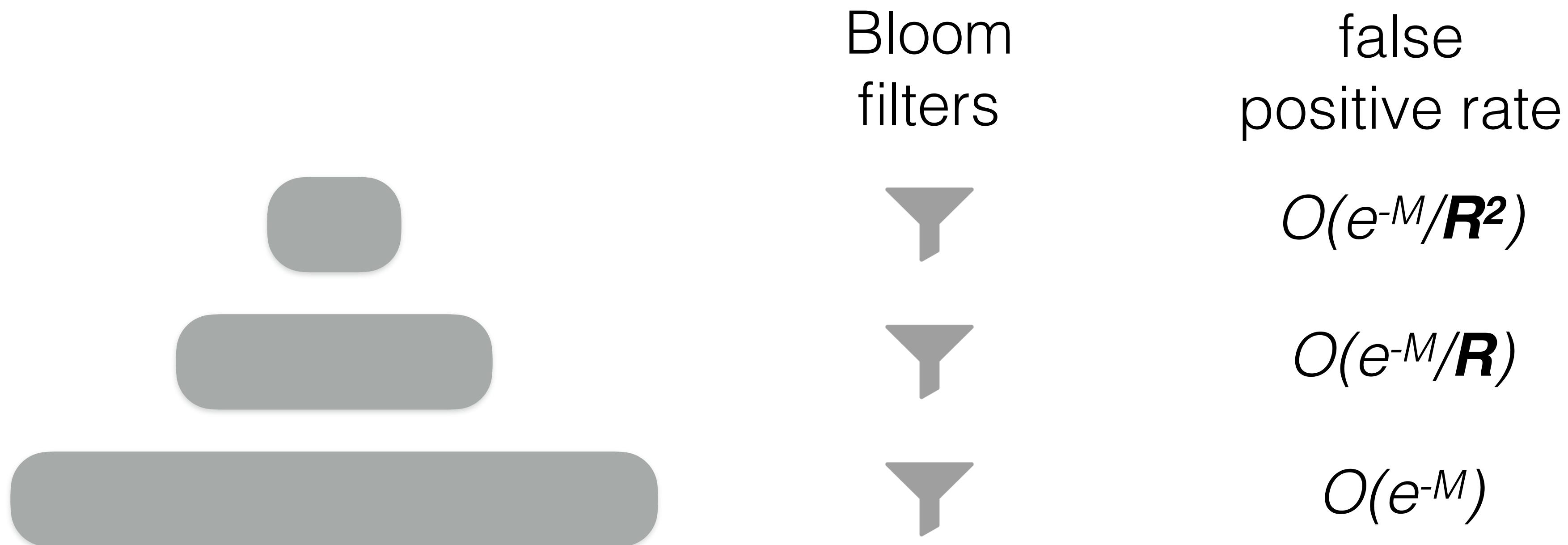
Monkey: Optimal Navigable Key-Value Store

SIGMOD17



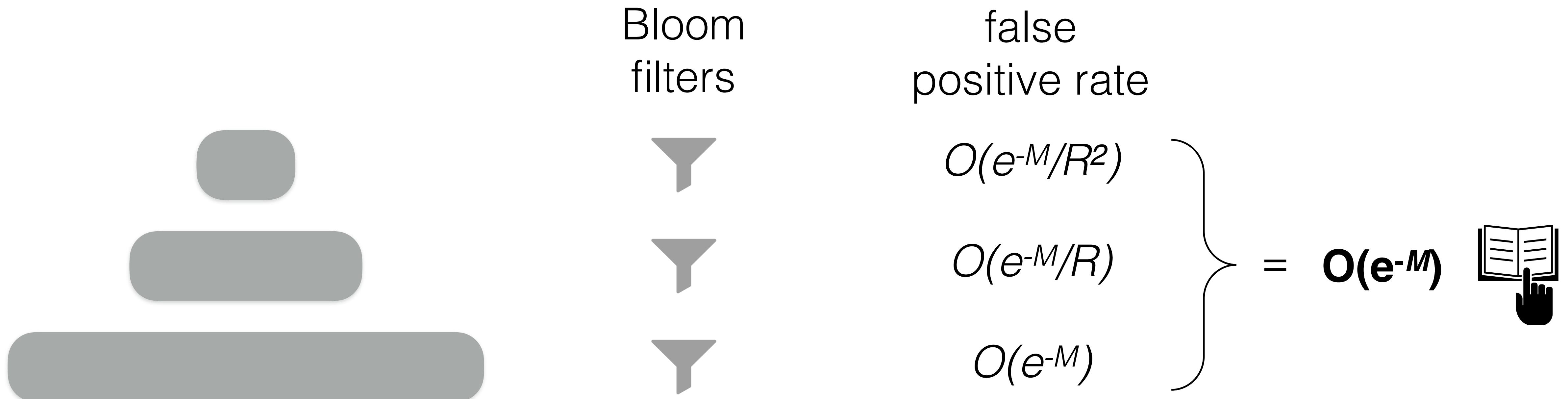
Monkey: Optimal Navigable Key-Value Store

SIGMOD17



Monkey: Optimal Navigable Key-Value Store

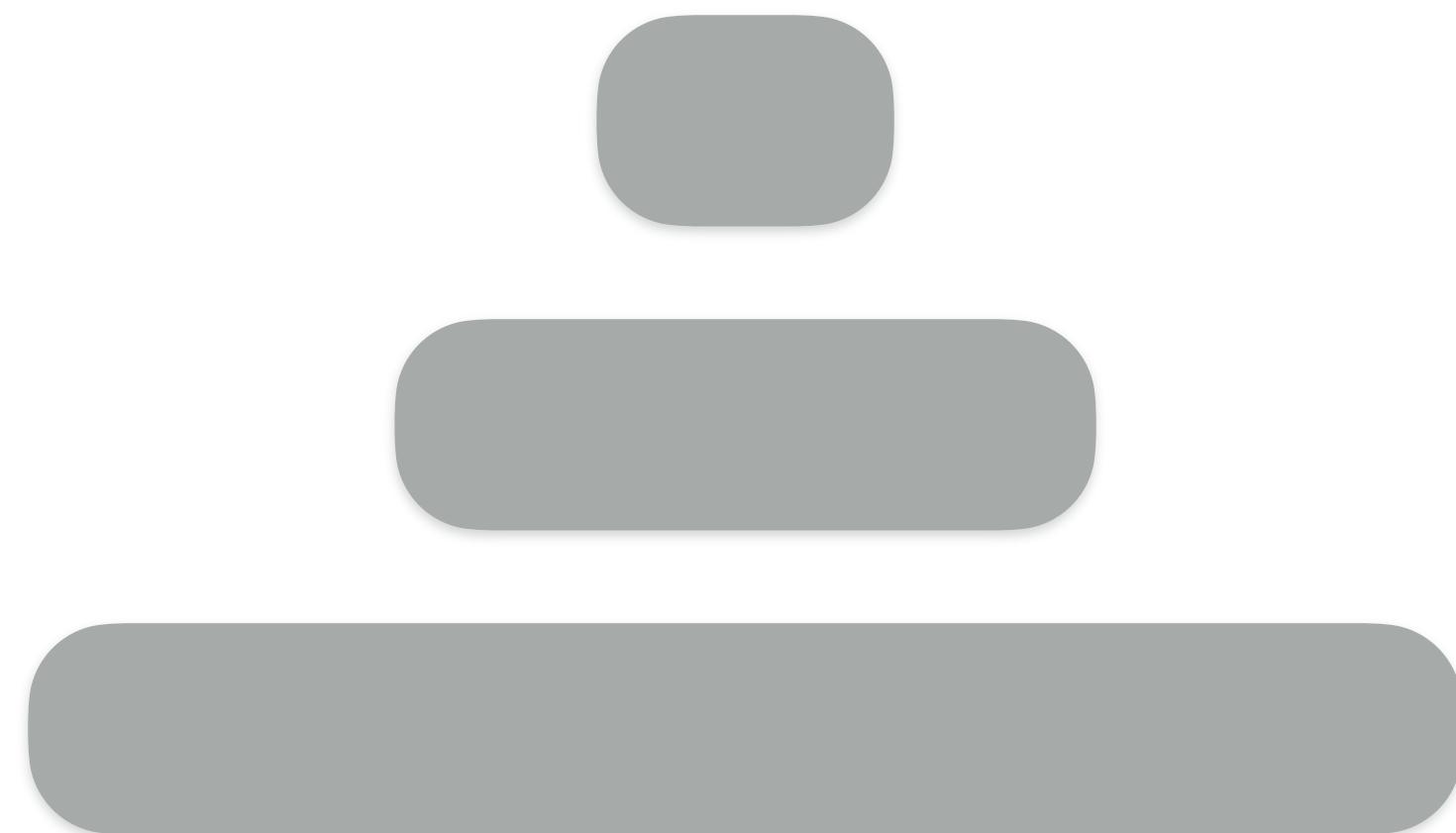
SIGMOD17



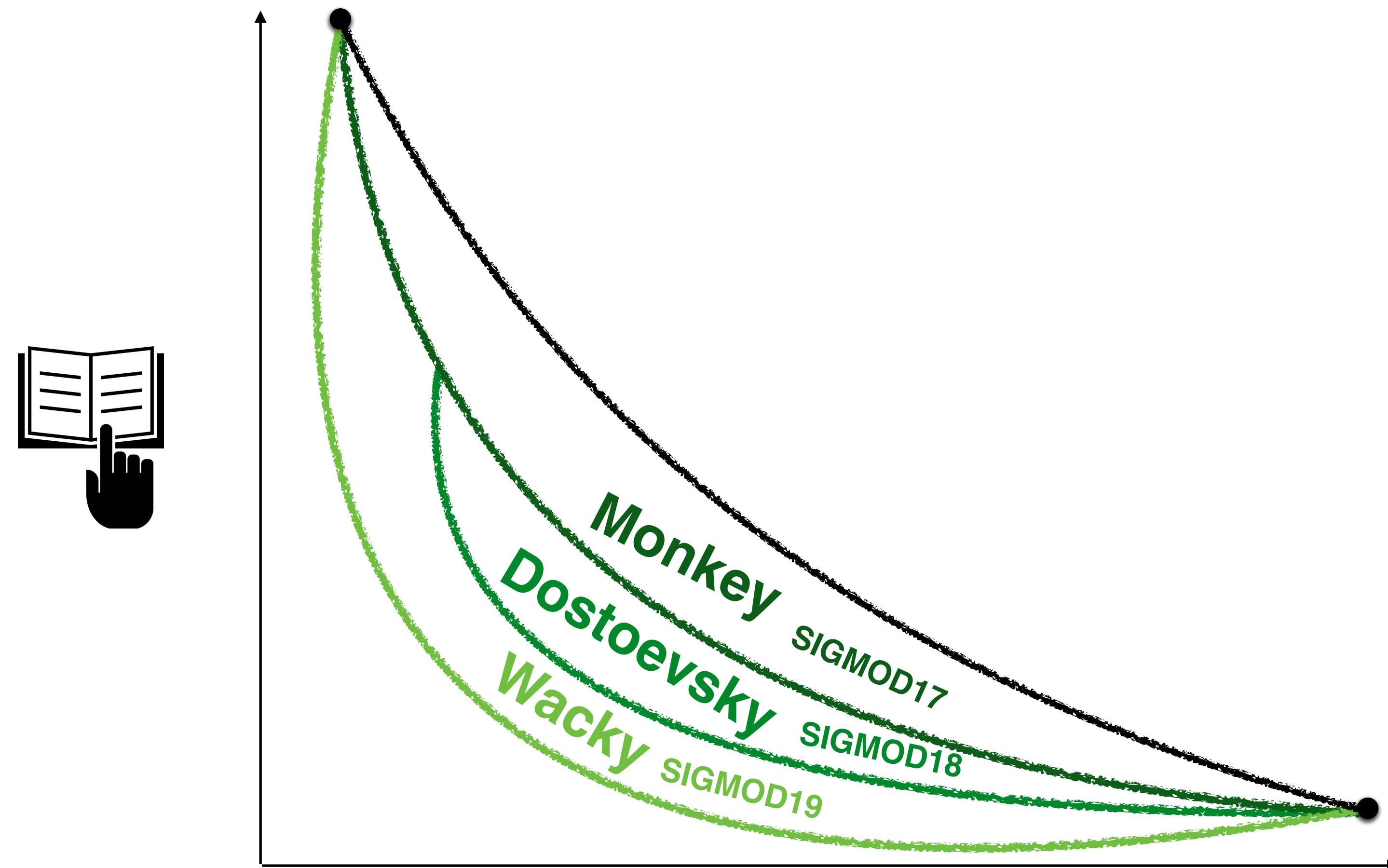
Monkey: Optimal Navigable Key-Value Store

SIGMOD17

Bloom
filters



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



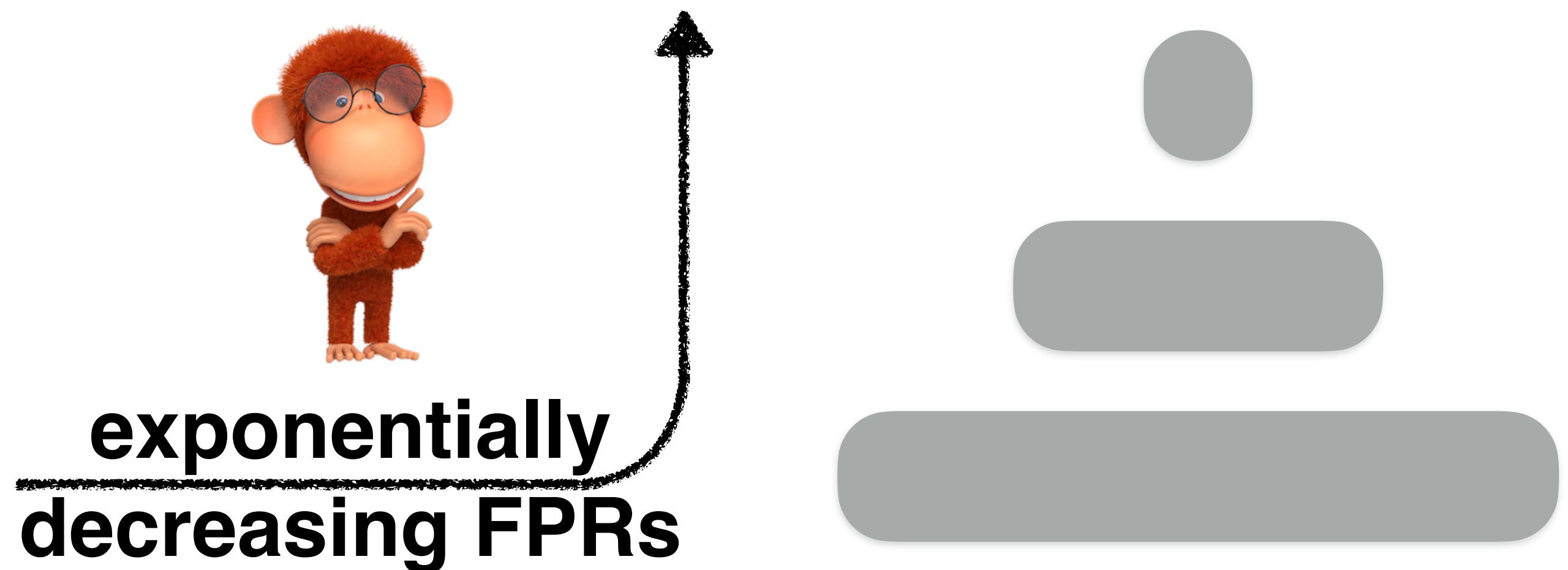
Dostoevsky: Space-Time Optimized Evolvable Scalable Key-Value Store

SIGMOD18



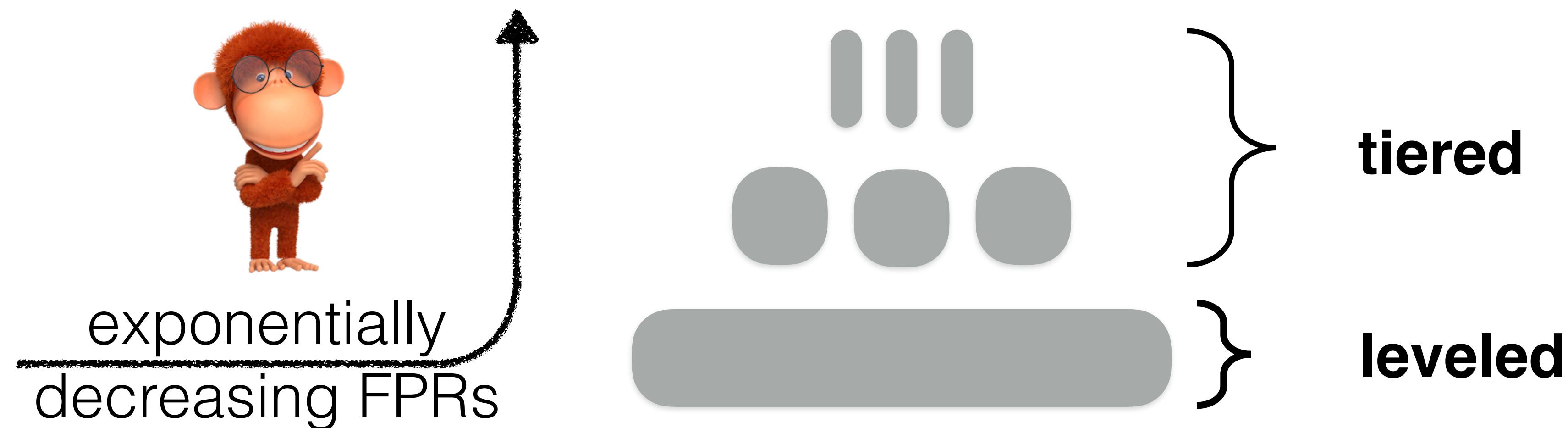
Dostoevsky: Space-Time Optimized Evolvable Scalable Key-Value Store

SIGMOD18



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

SIGMOD18



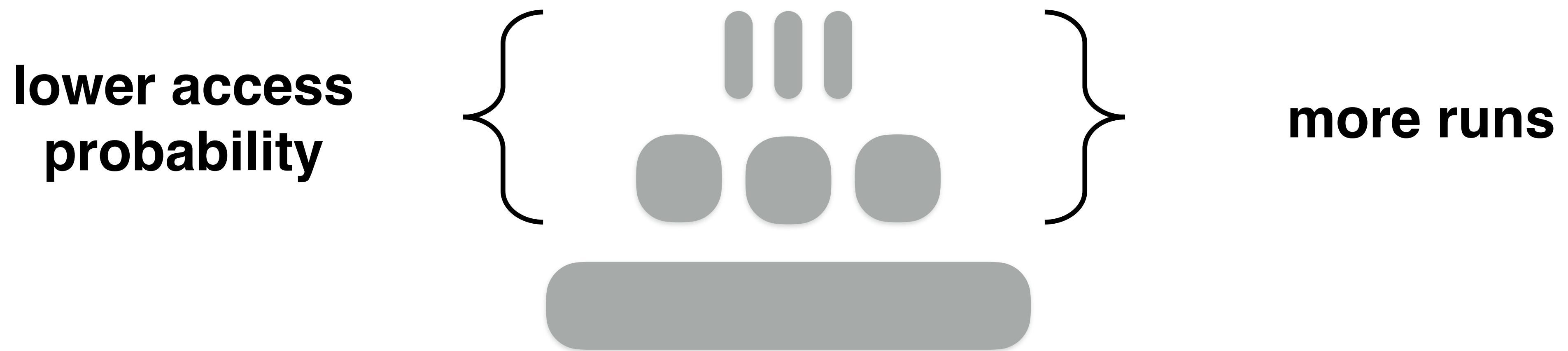
Dostoevsky: Space-Time Optimized Evolvable Scalable Key-Value Store

SIGMOD18



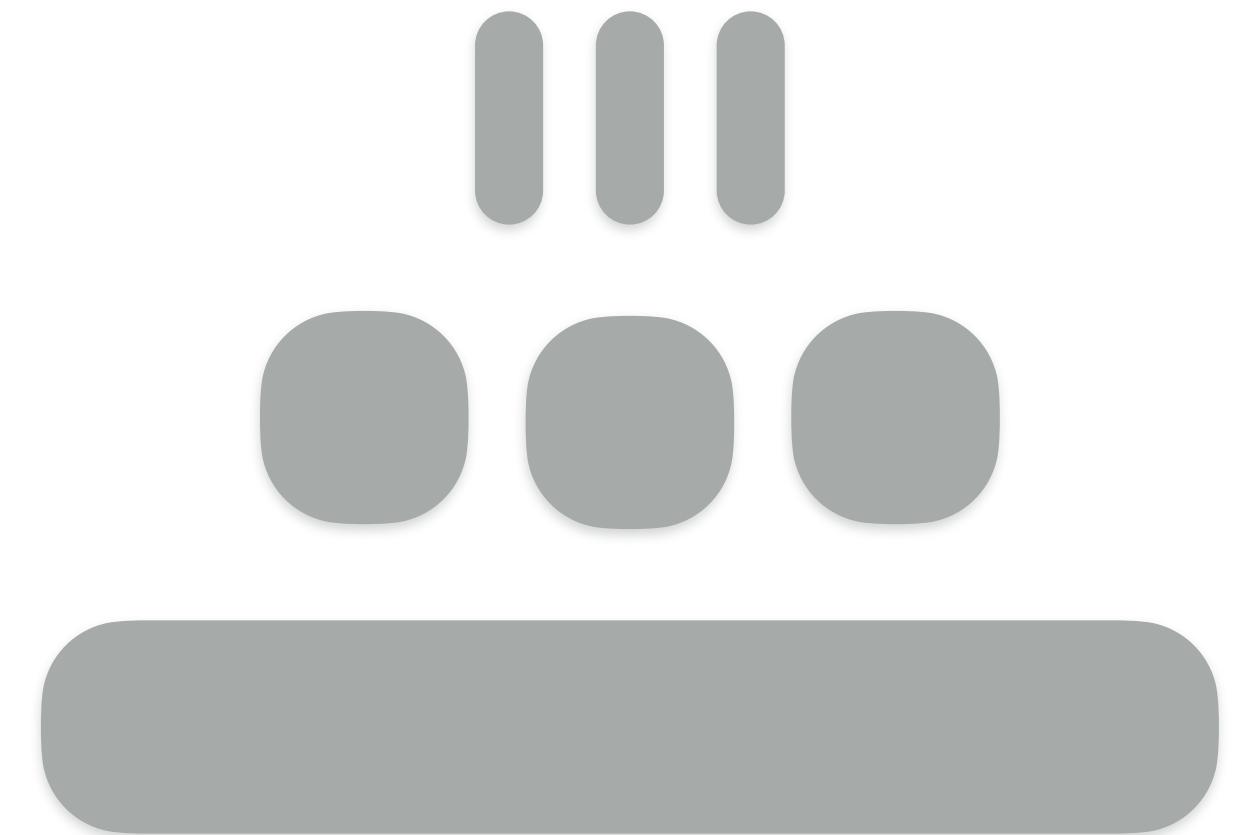
Dostoevsky: Space-Time Optimized Evolvable Scalable Key-Value Store

SIGMOD18

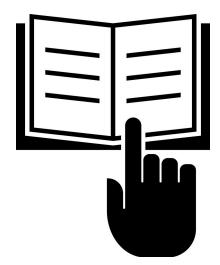


Dostoevsky: Space-Time Optimized Evolvable Scalable Key-Value Store

SIGMOD18



$O(e^{-M})$



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



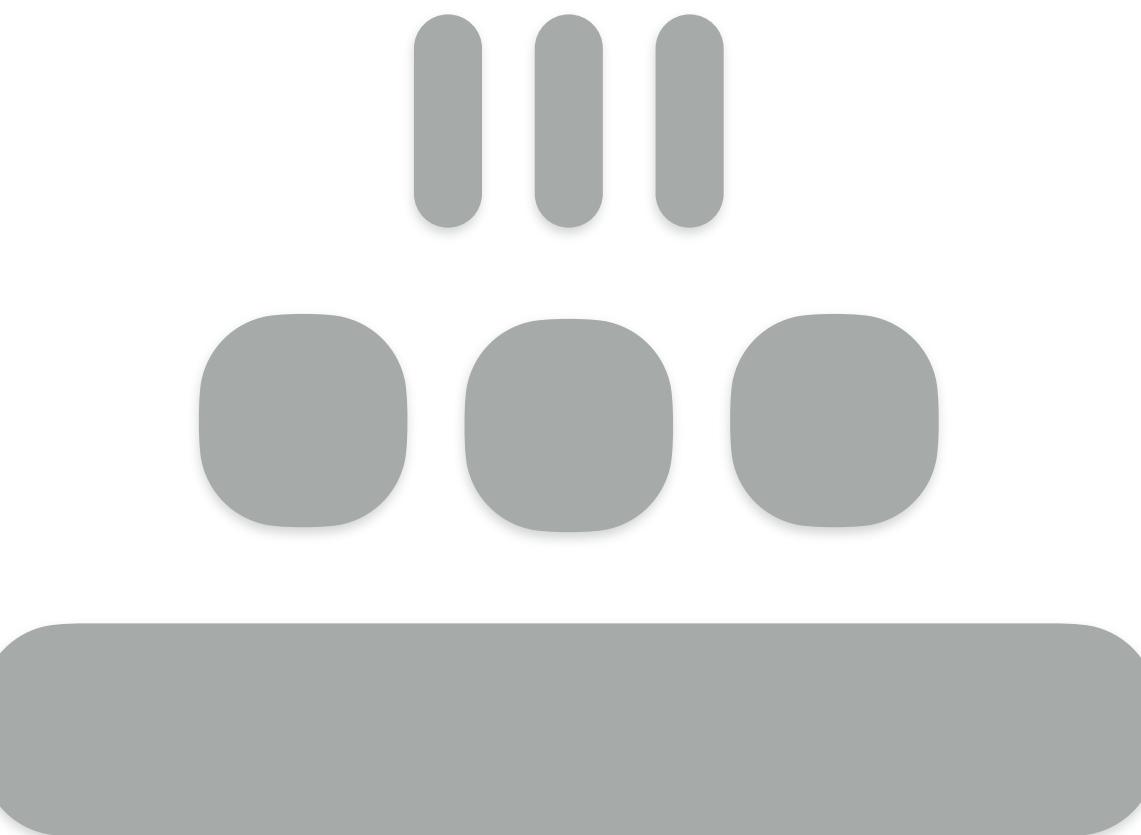
Dostoevsky: Space-Time Optimized Evolvable Scalable Key-Value Store

SIGMOD18

Leveling



**Lazy
Leveling**

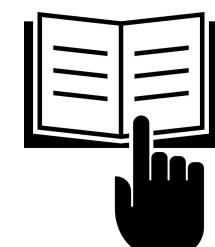


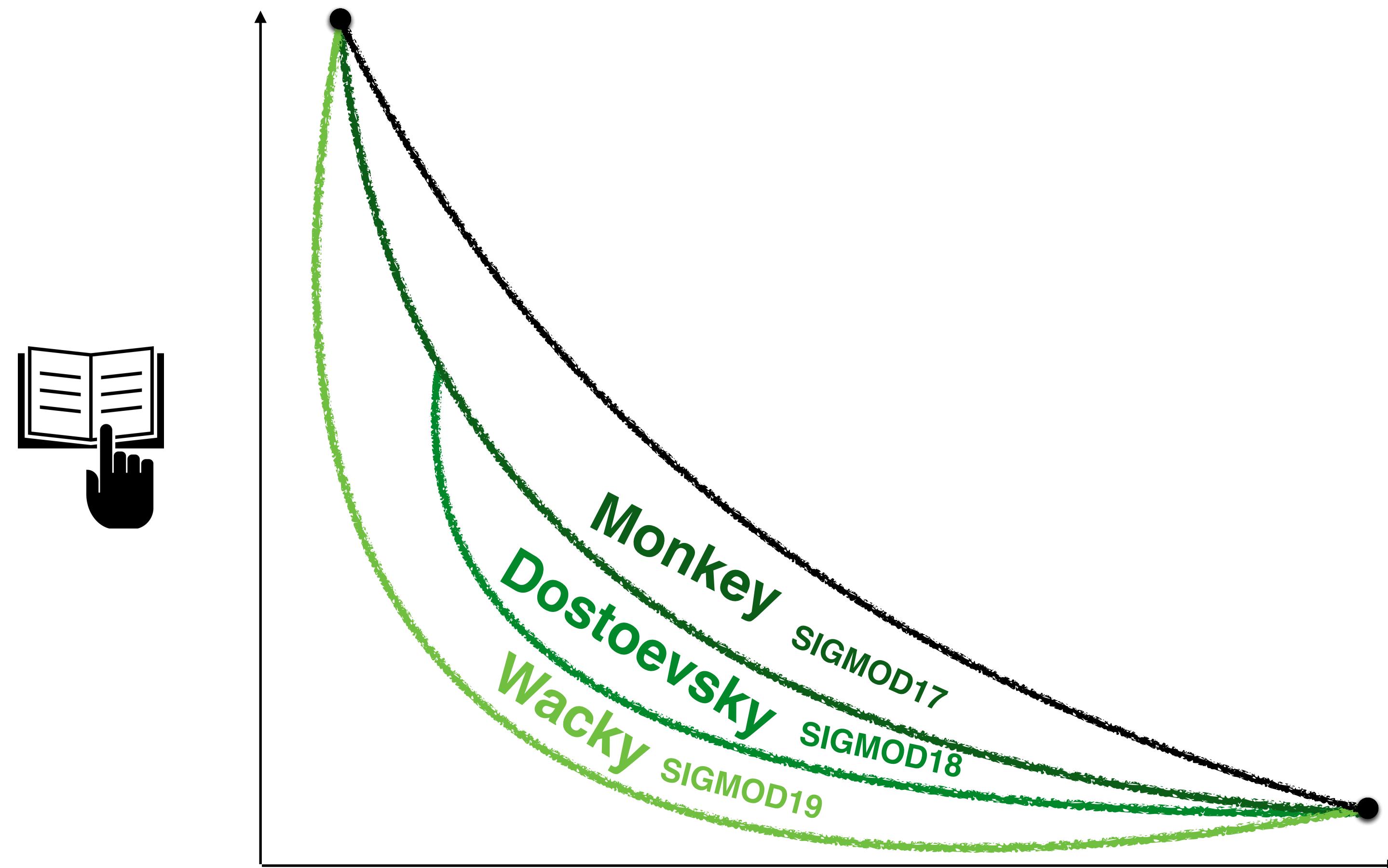
$$O(\mathbf{e}^{-M}) =$$

$$O(\mathbf{R} \cdot \log_{\mathbf{R}}(\mathbf{N})) >$$

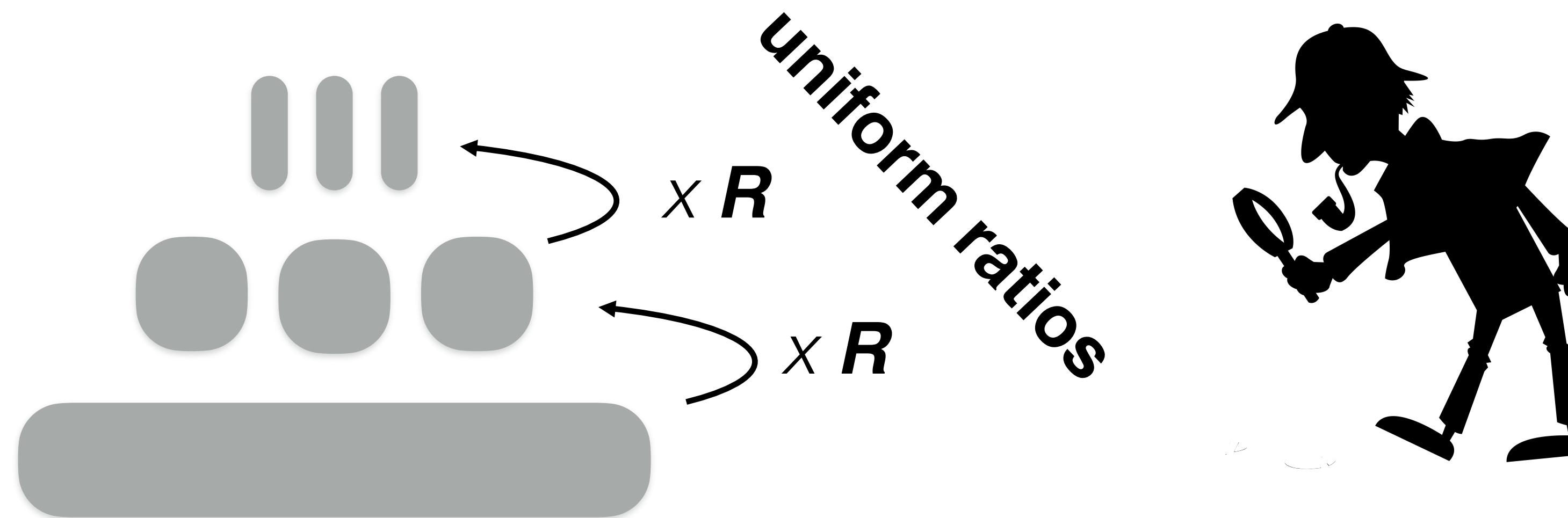
$$O(\mathbf{e}^{-M})$$

$$O(\mathbf{R} + \log_{\mathbf{R}}(\mathbf{N}))$$

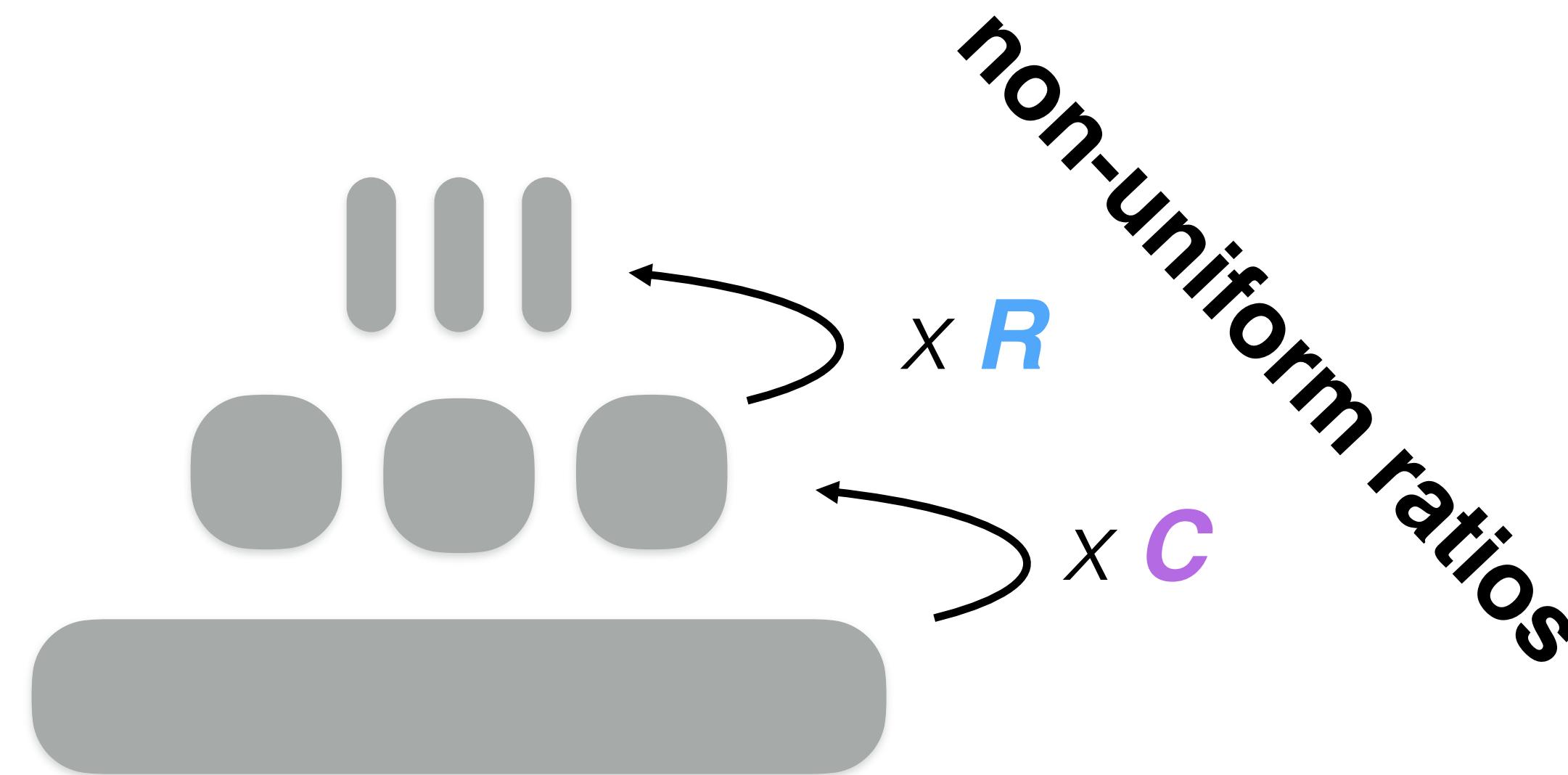




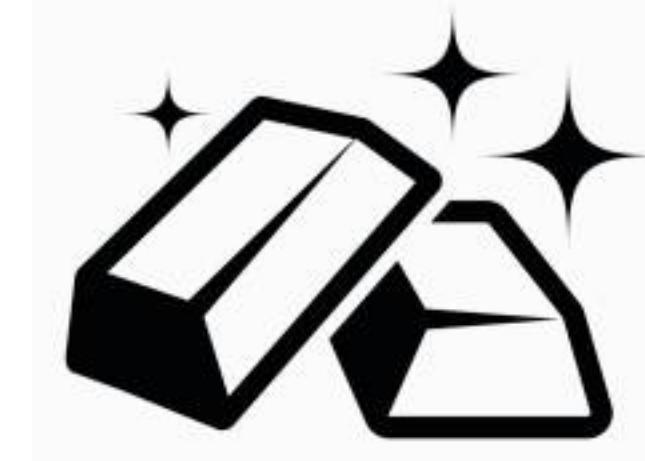
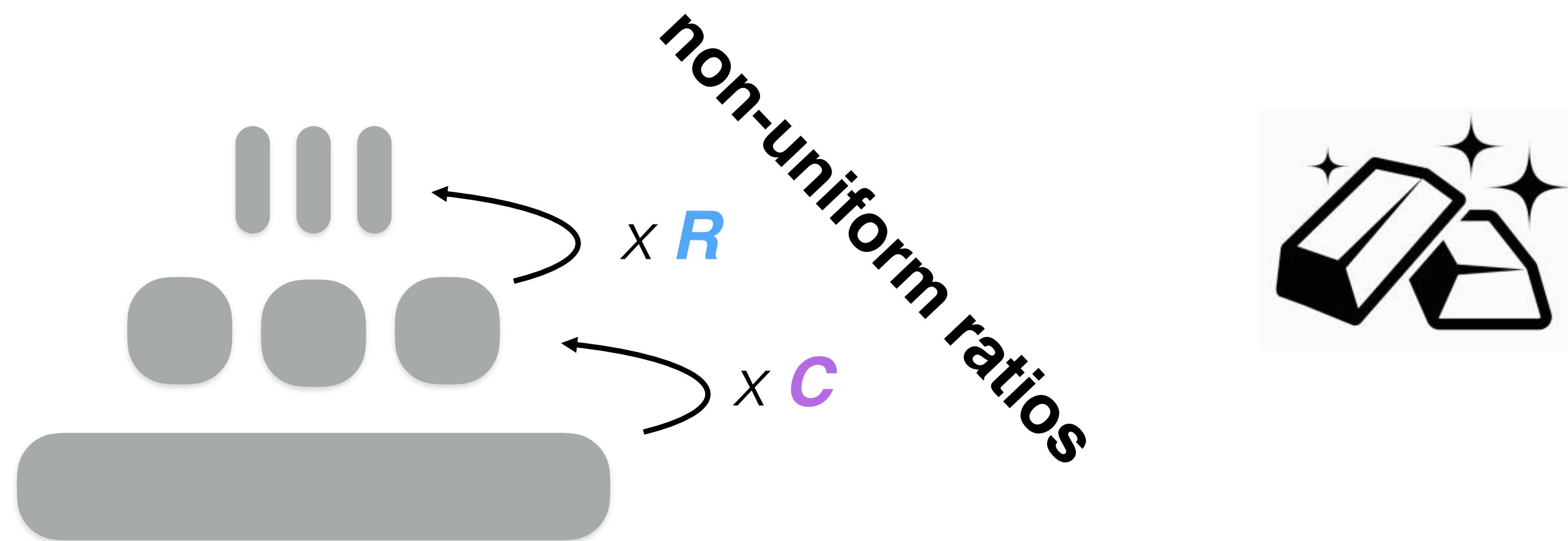
Wacky SIGMOD19



Wacky SIGMOD19



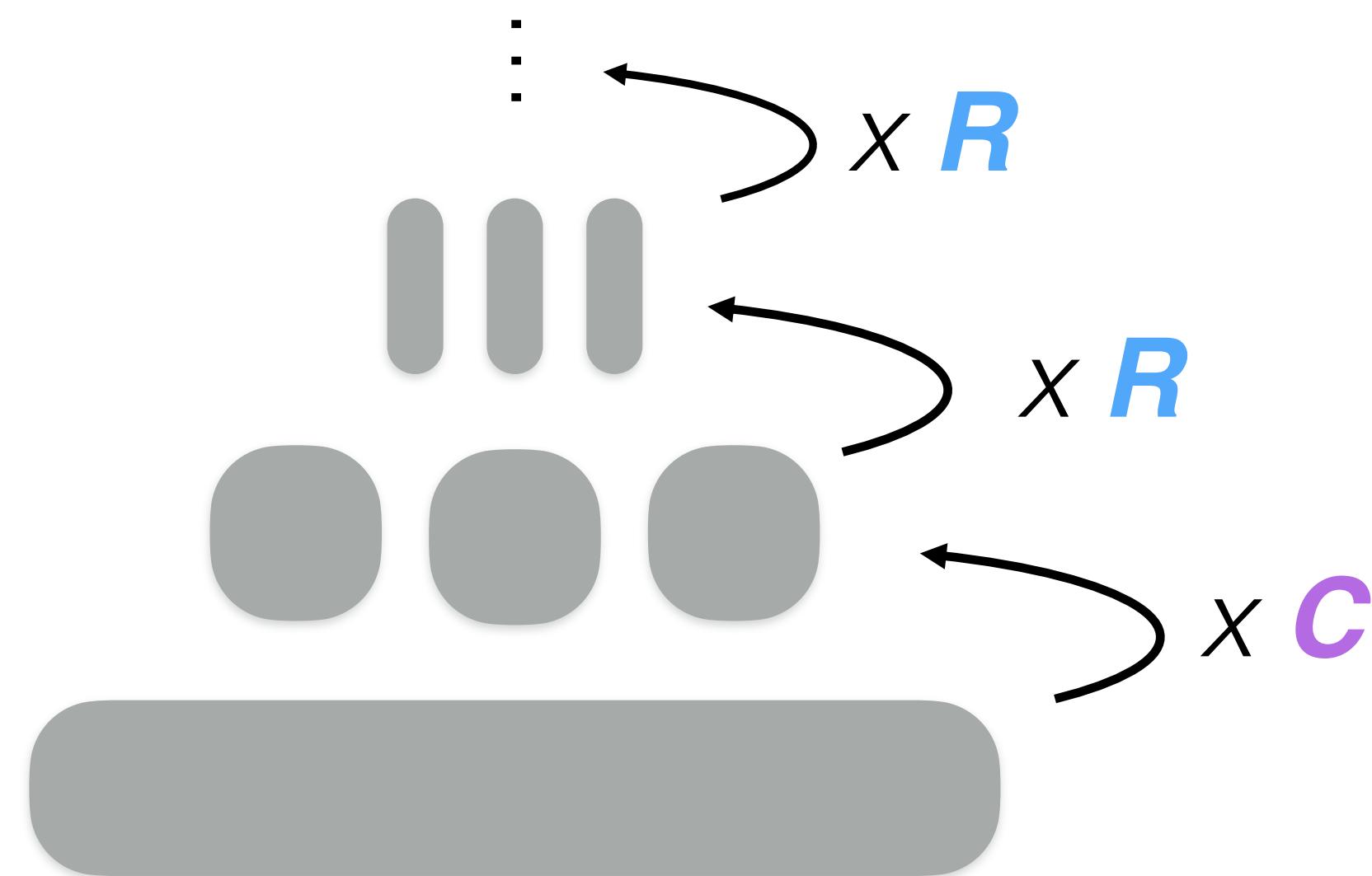
Wacky SIGMOD19



Wacky

SIGMOD19

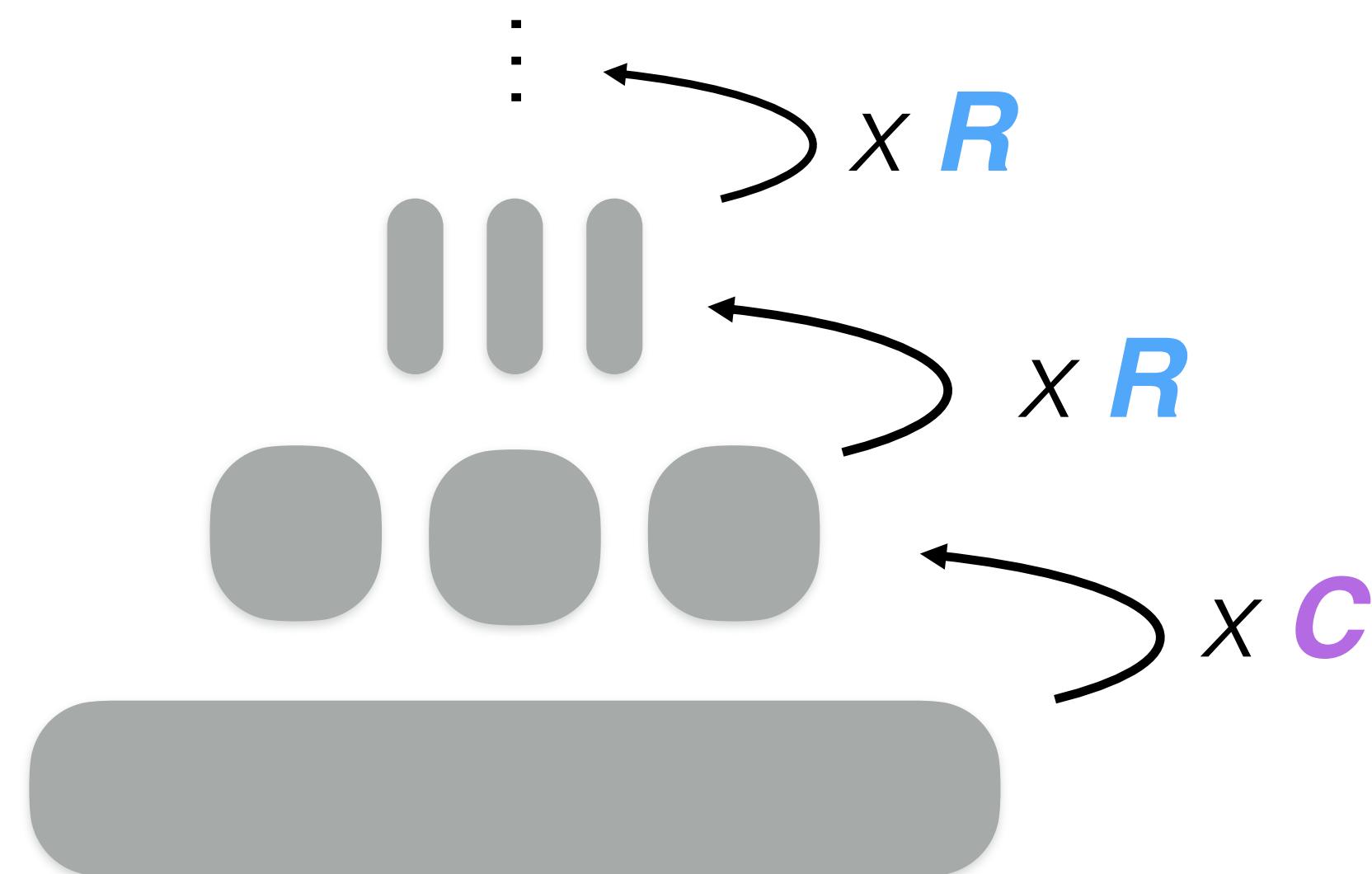
1. Squared Capped Lazy Leveling



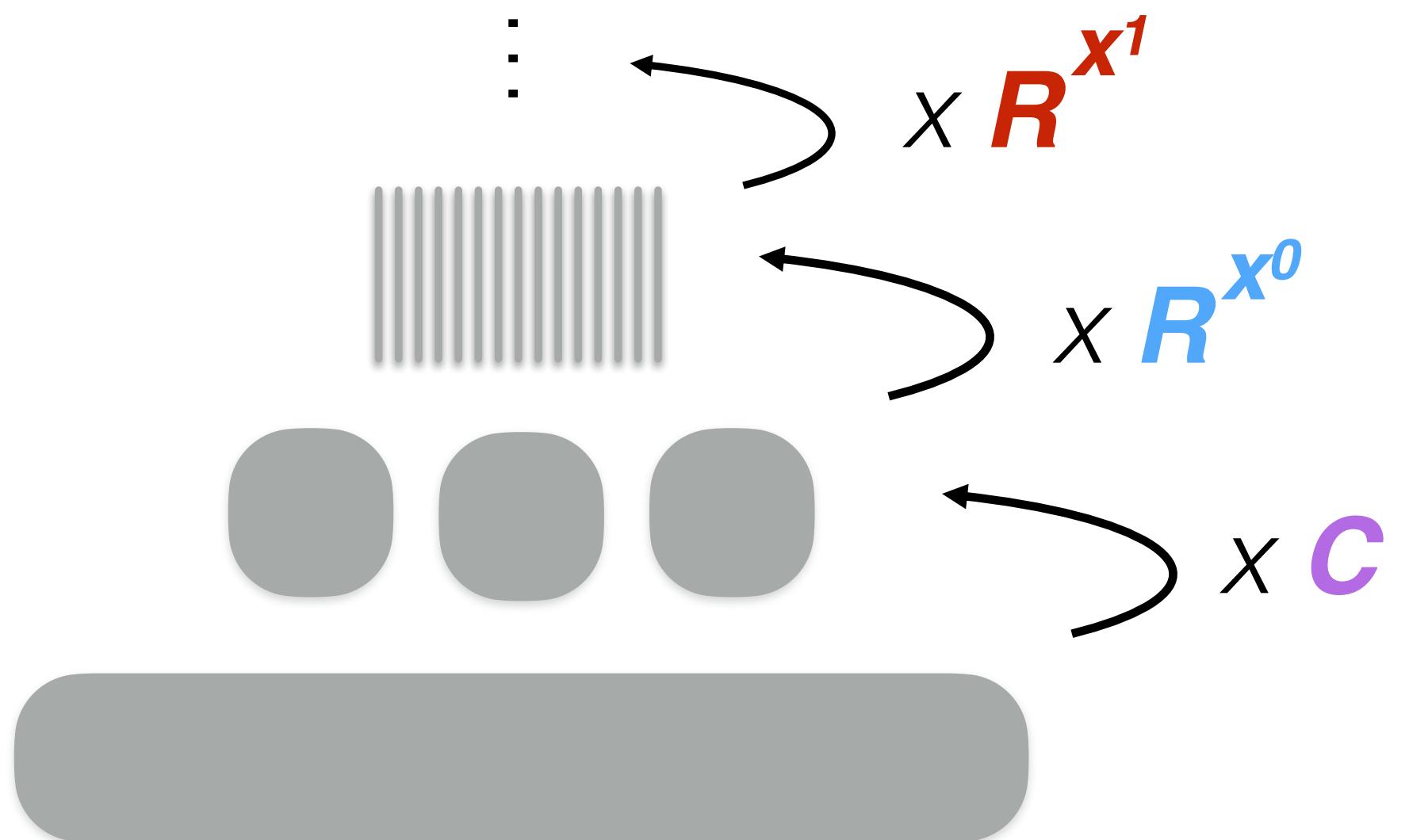
Wacky

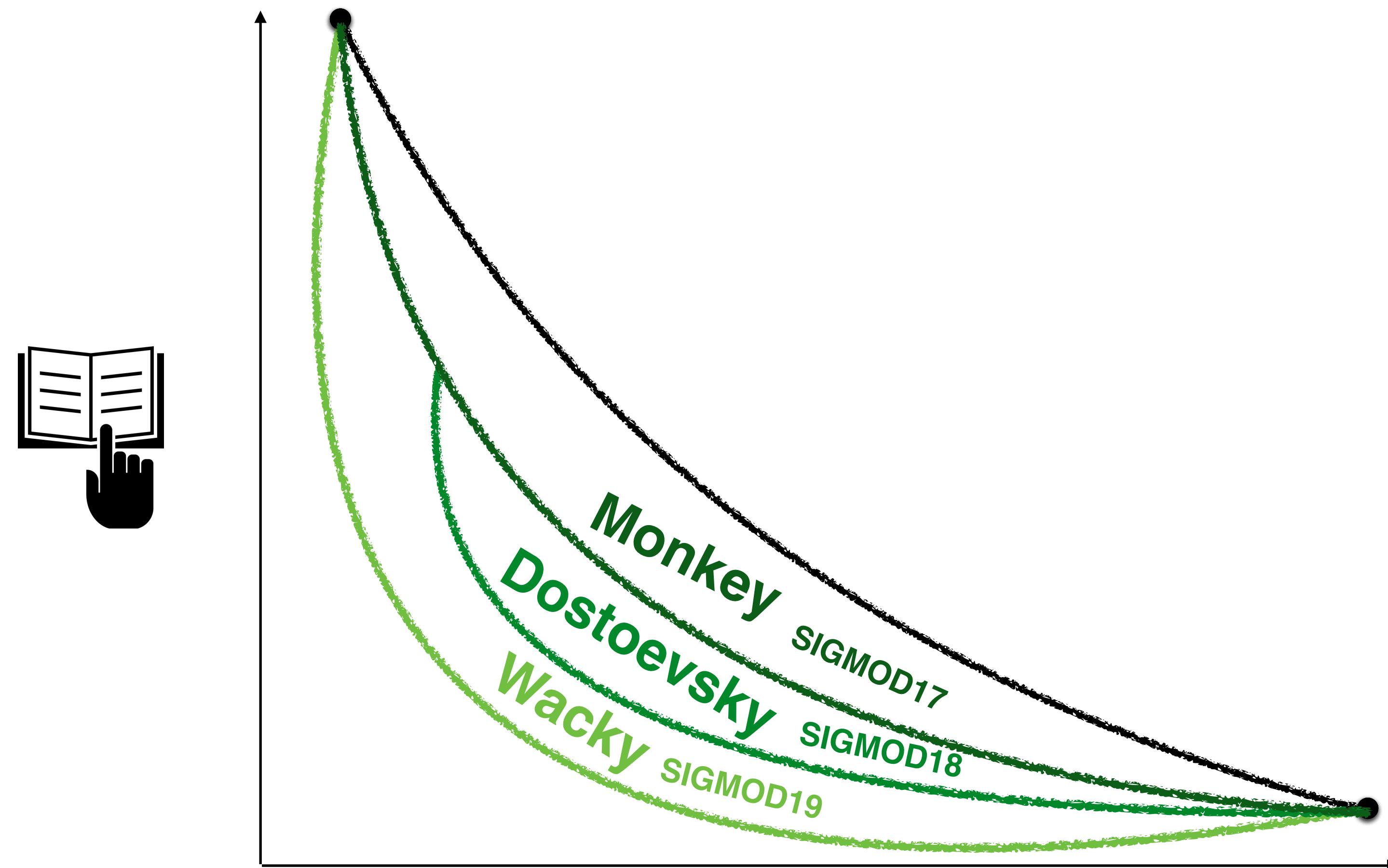
SIGMOD19

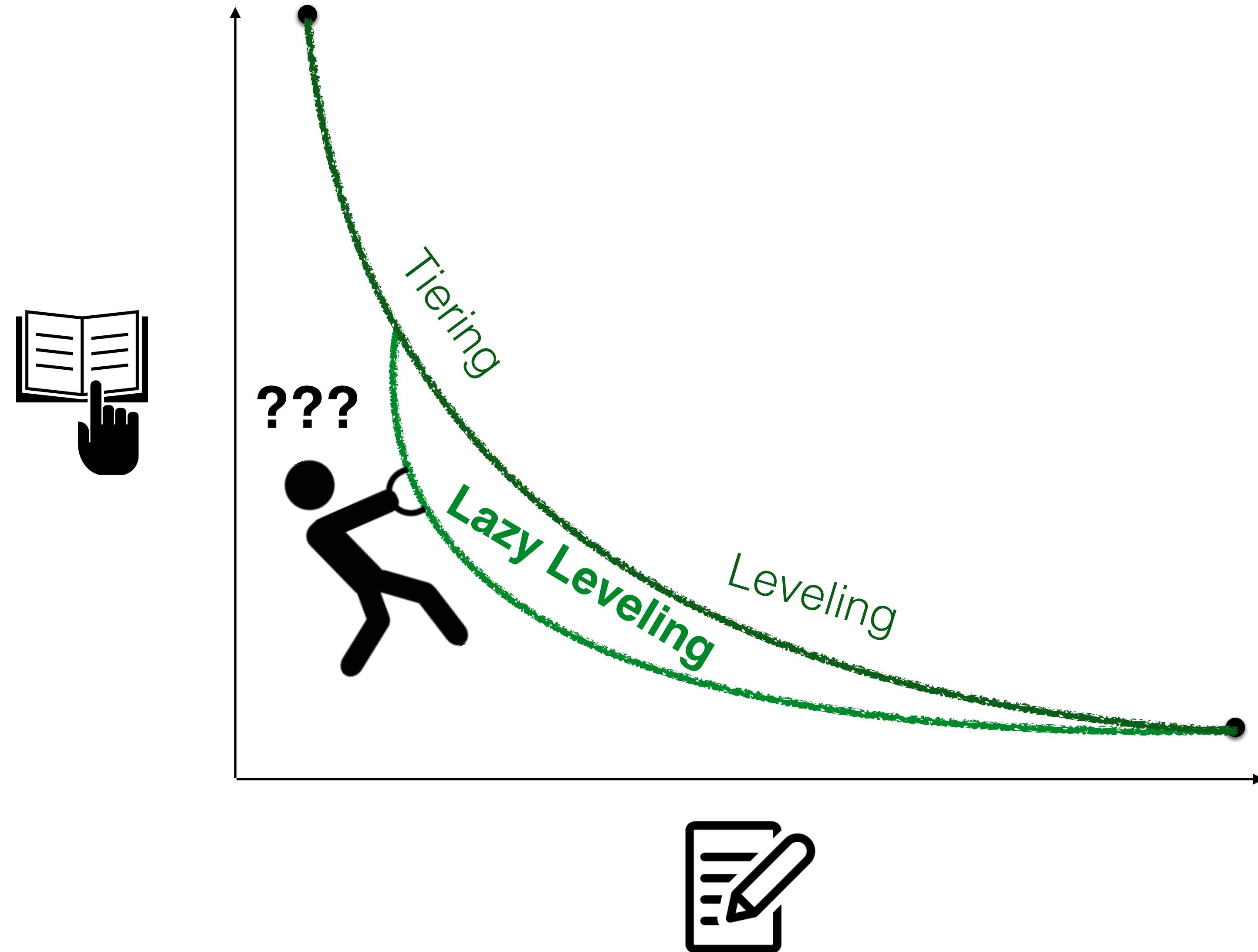
1. Squared Capped Lazy Leveling



2. LSM-bush



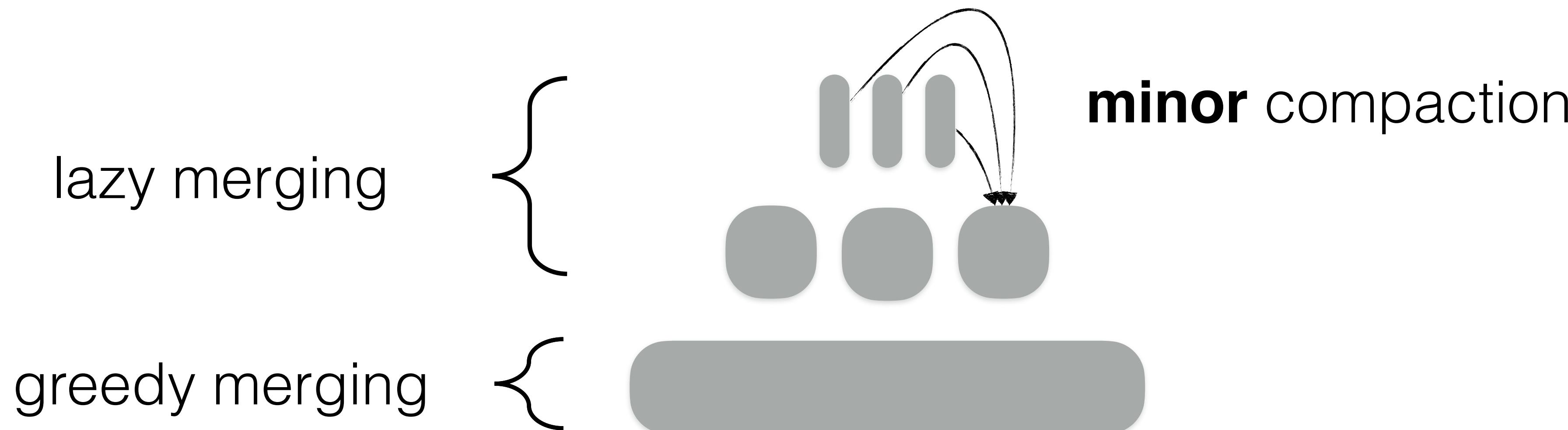




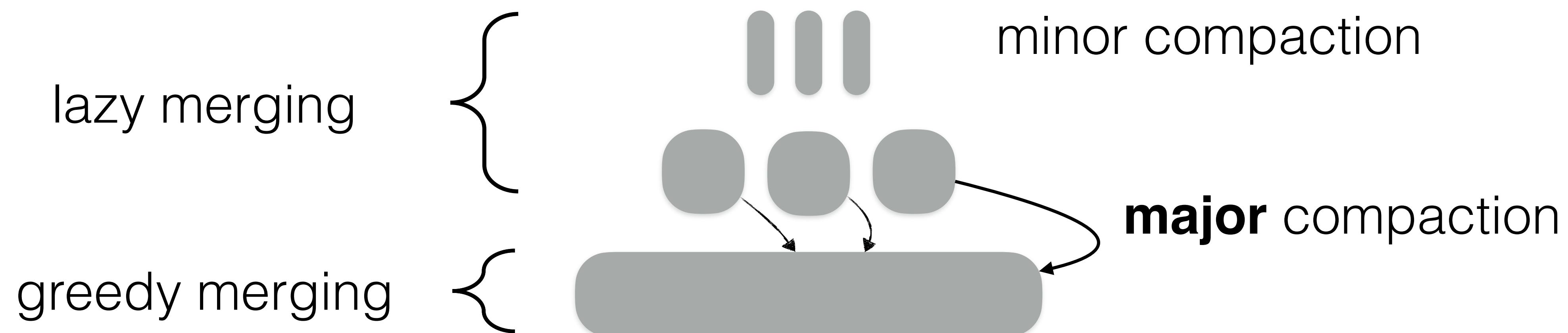
Problem Analysis: Lazy Leveling



Problem Analysis: Lazy Leveling



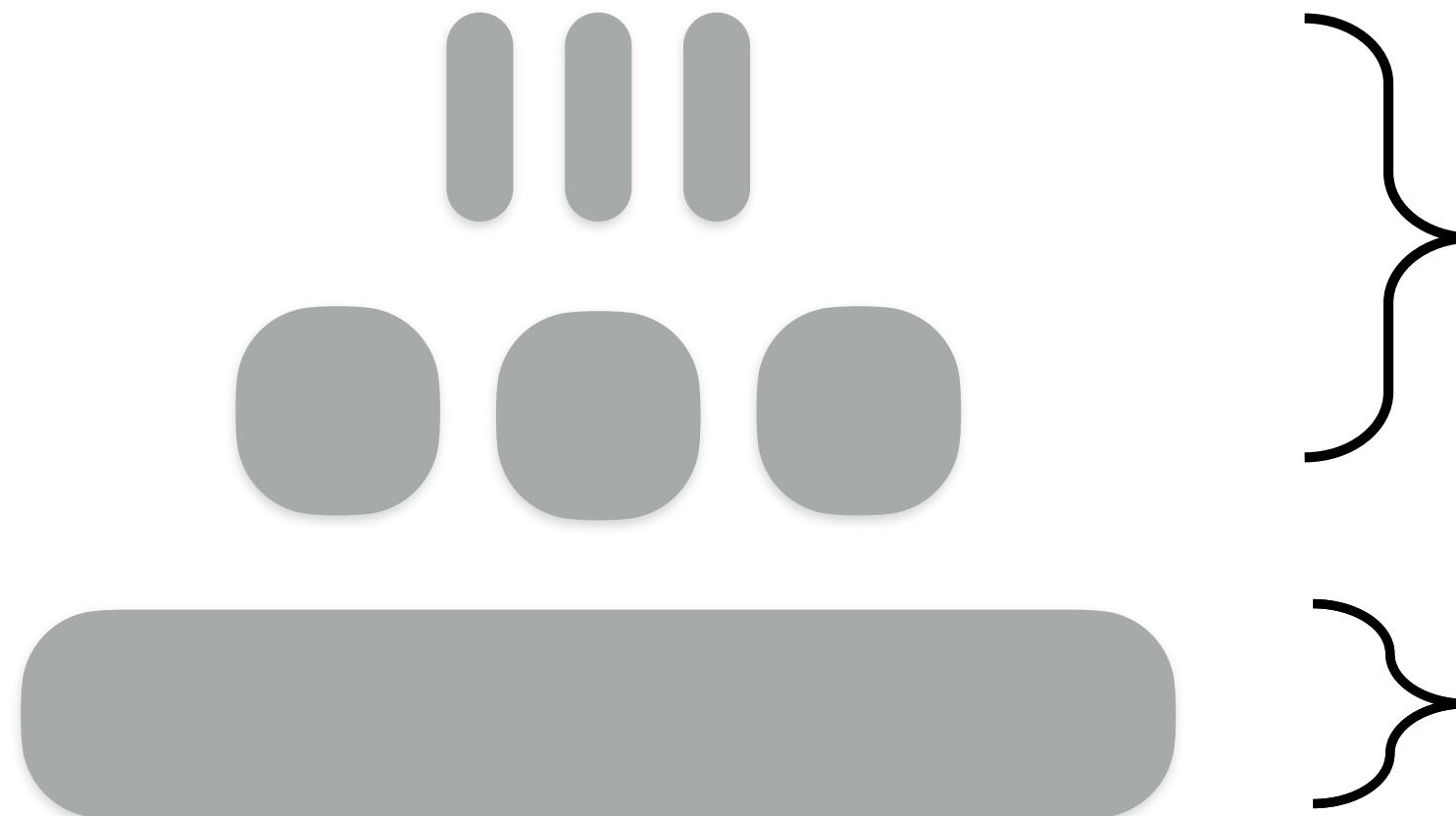
Problem Analysis: Lazy Leveling



Problem Analysis: Lazy Leveling



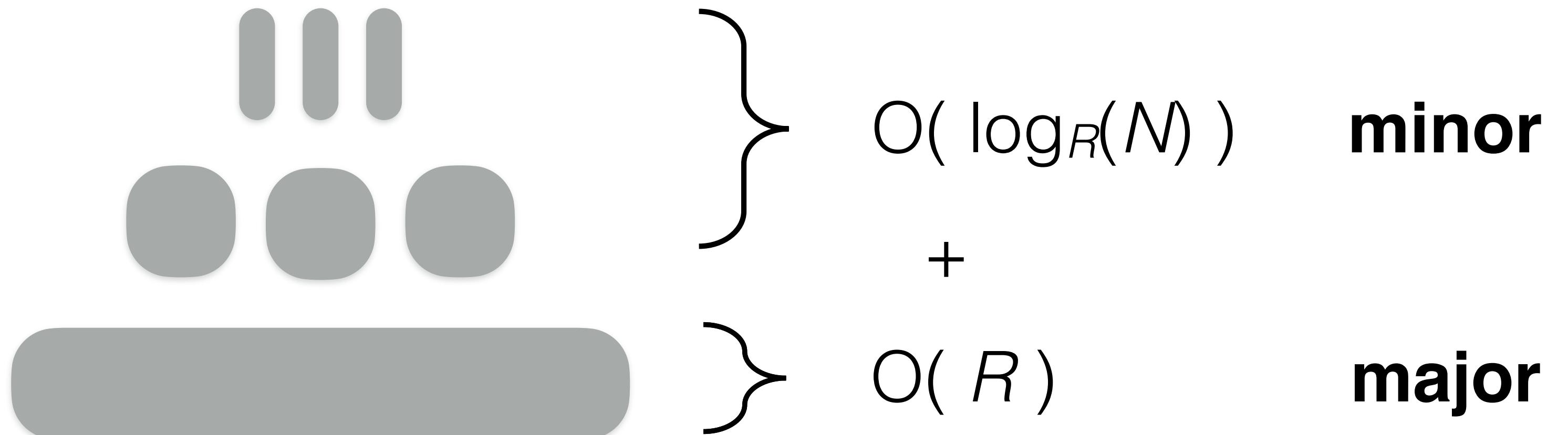
$O(R + \log_R(N))$



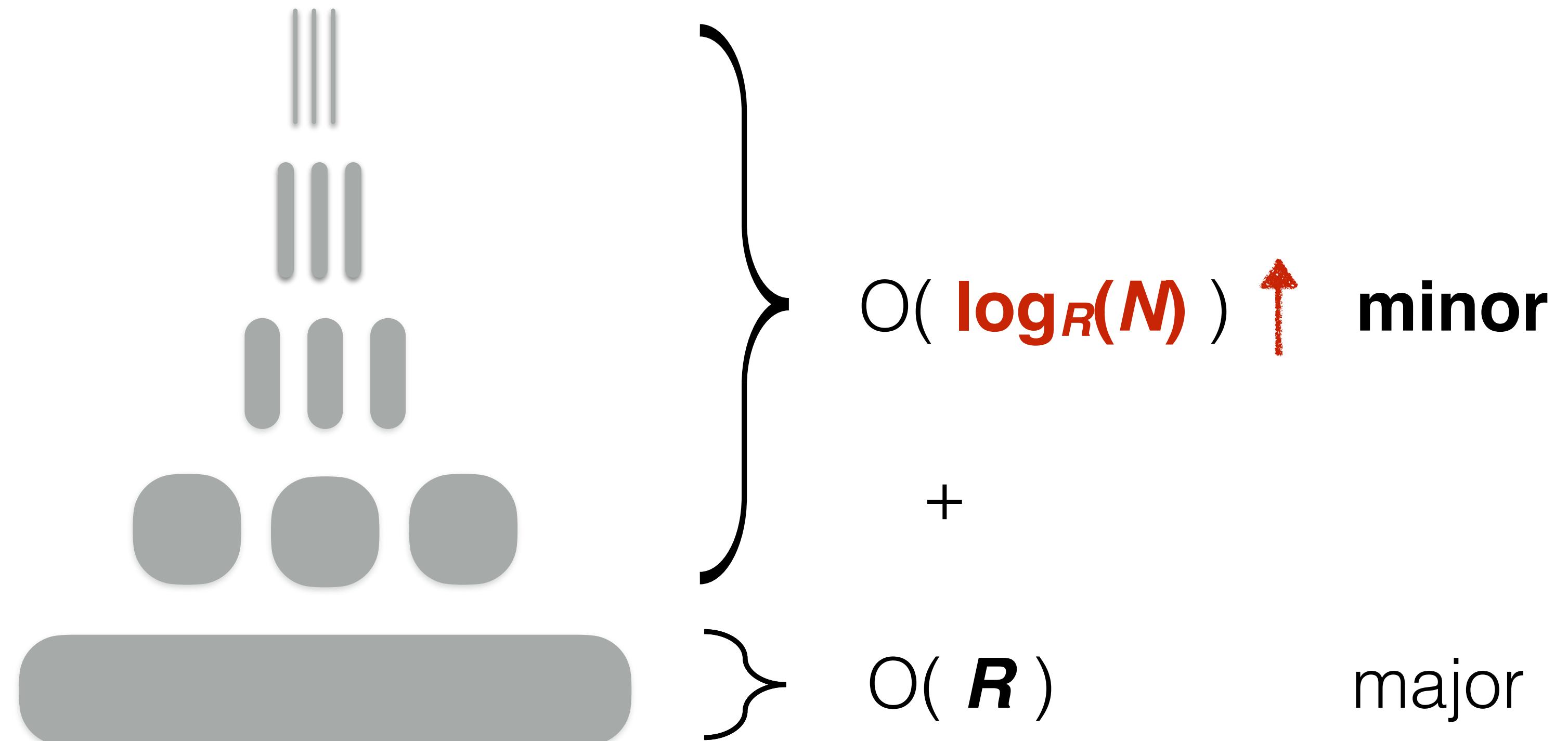
minor compactions

major compaction

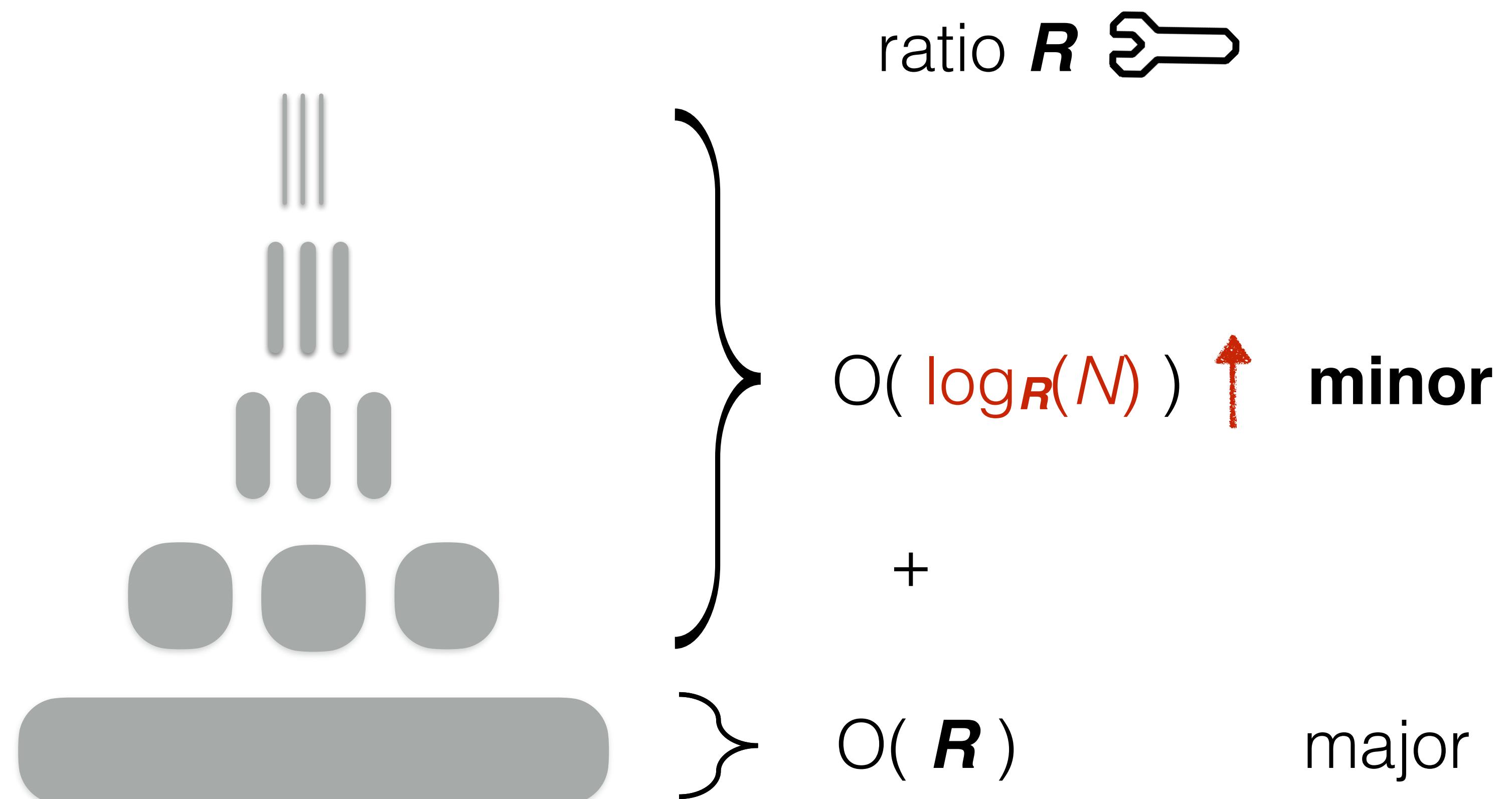
Problem Analysis: Lazy Leveling



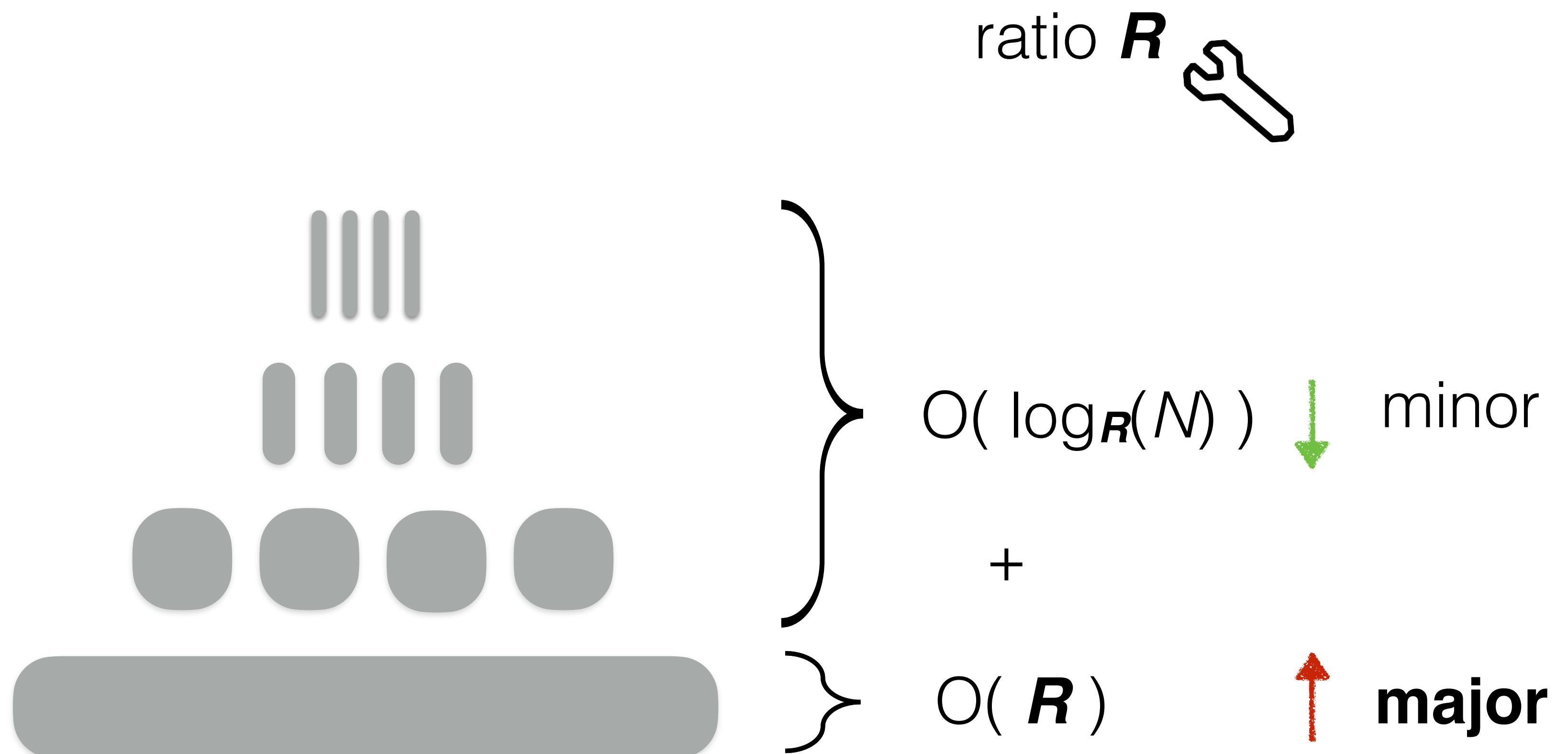
Problem Analysis: Lazy Leveling



Problem Analysis: Lazy Leveling



Problem Analysis: Lazy Leveling



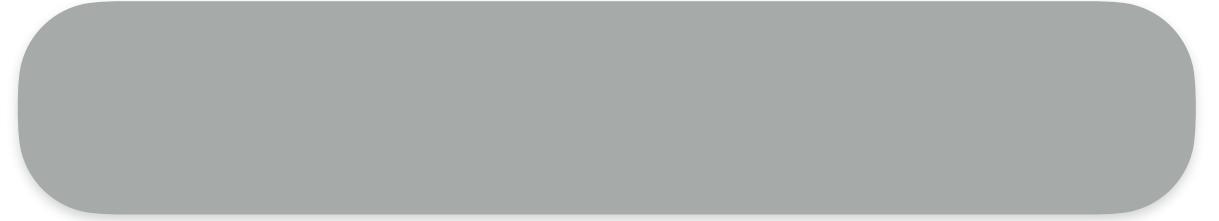
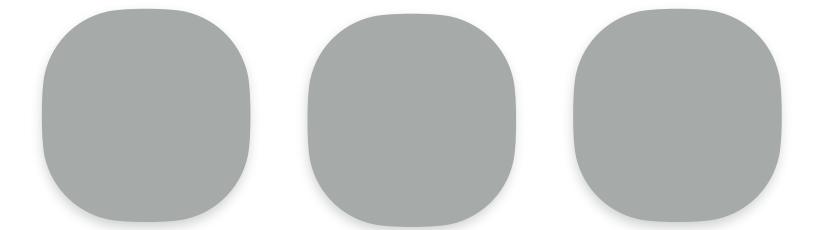
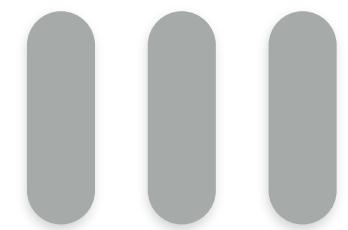
$$O(R + \log_R(N))$$

major



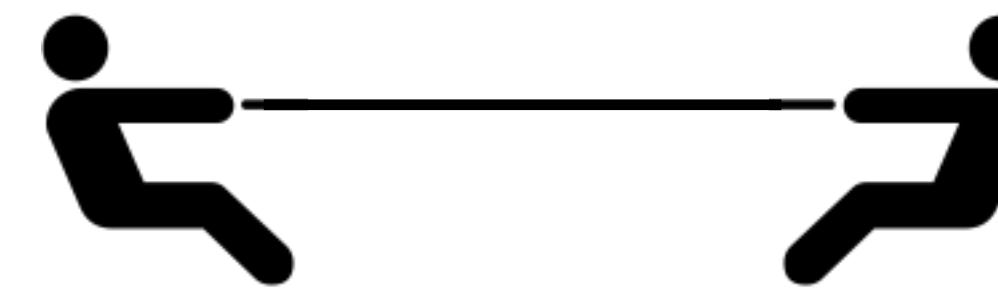
minor

:

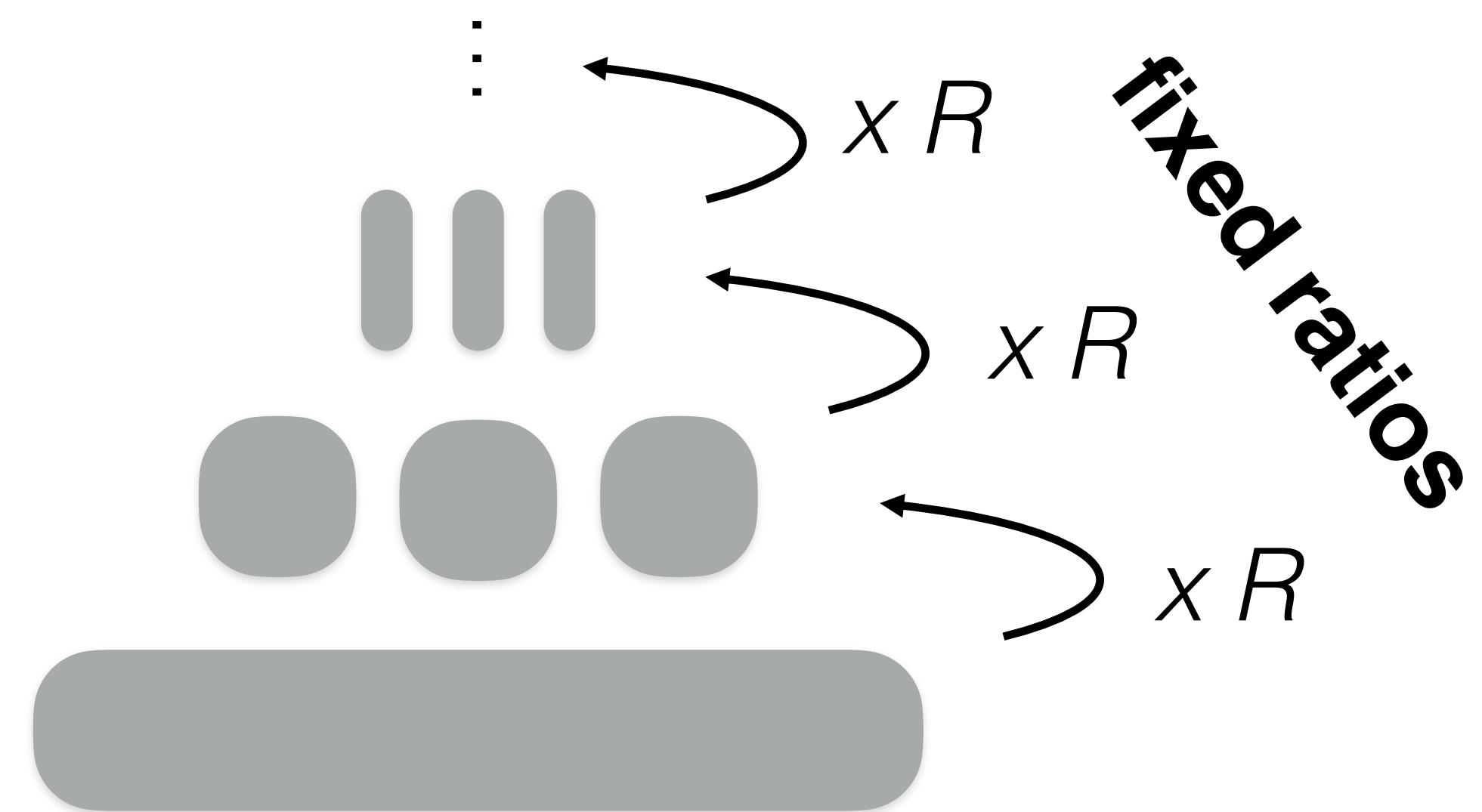


$O(R + \log_R(N))$

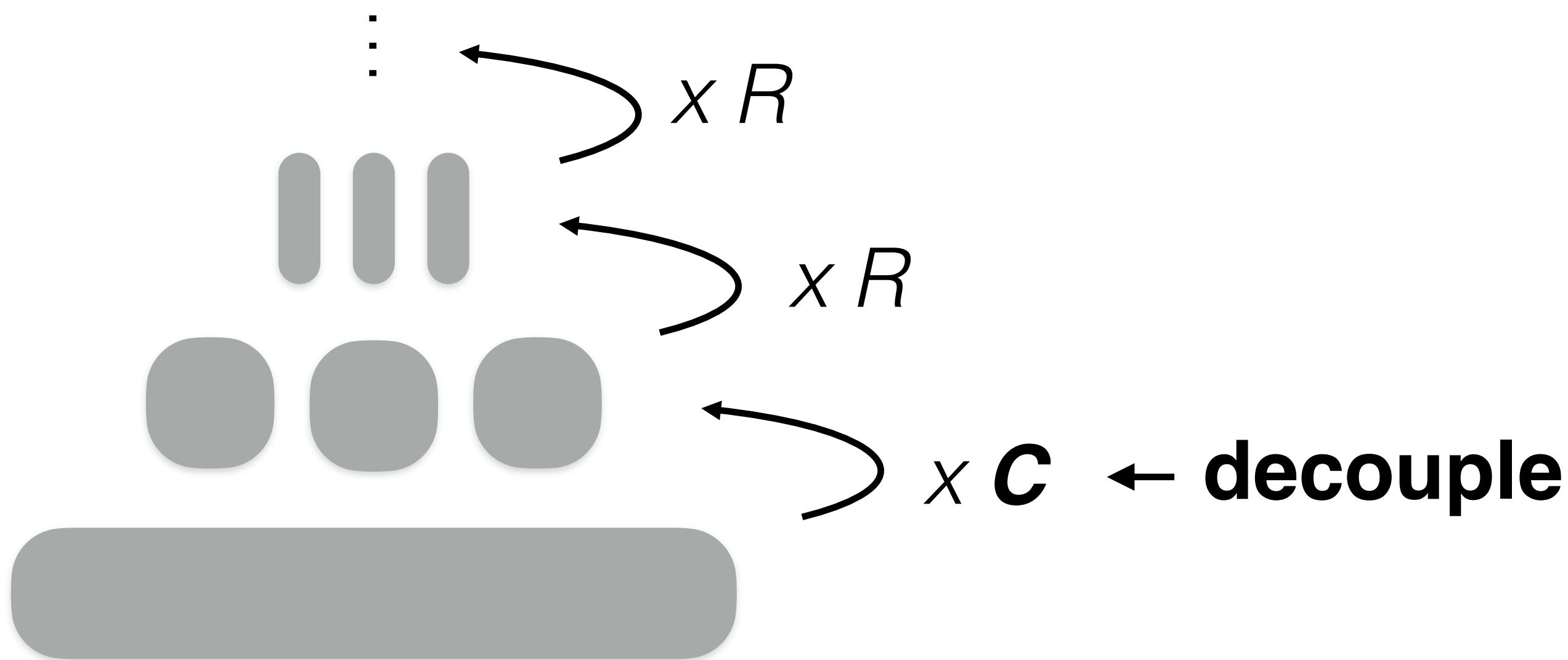
major



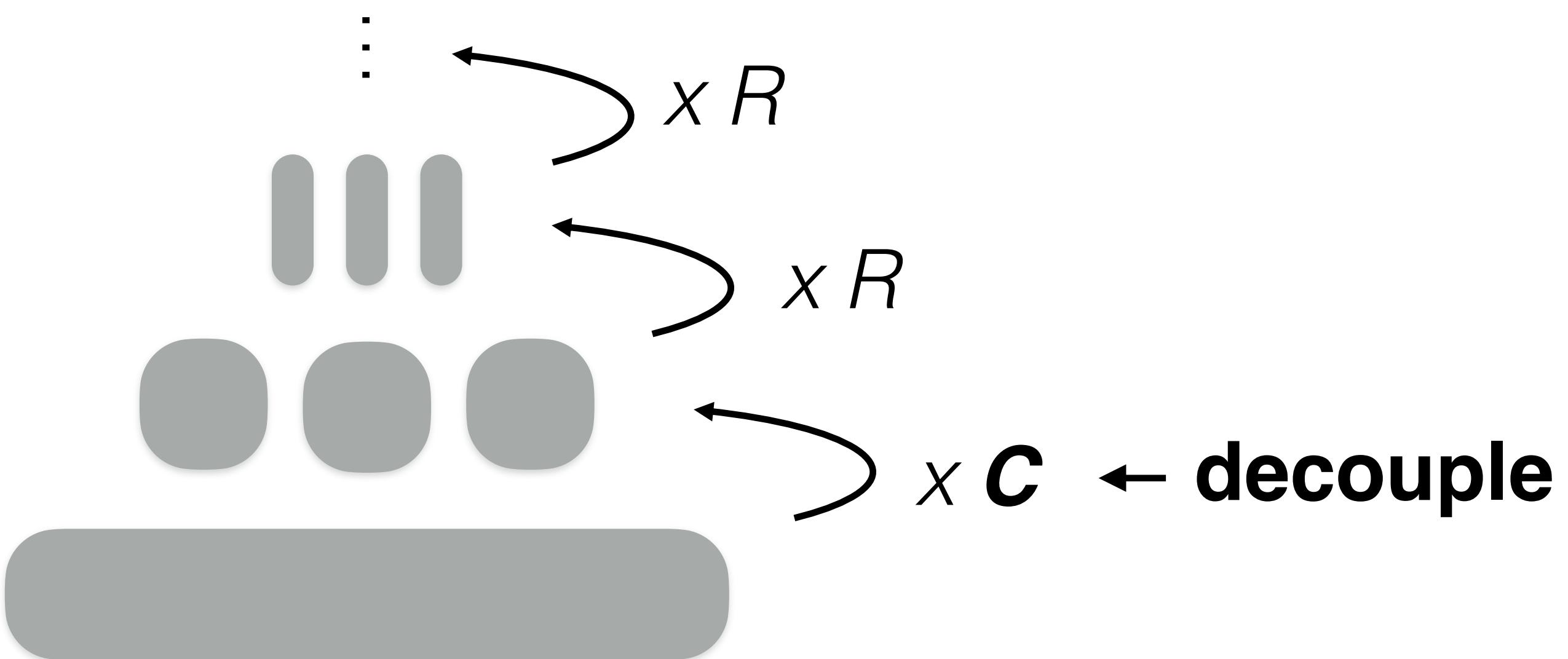
minor



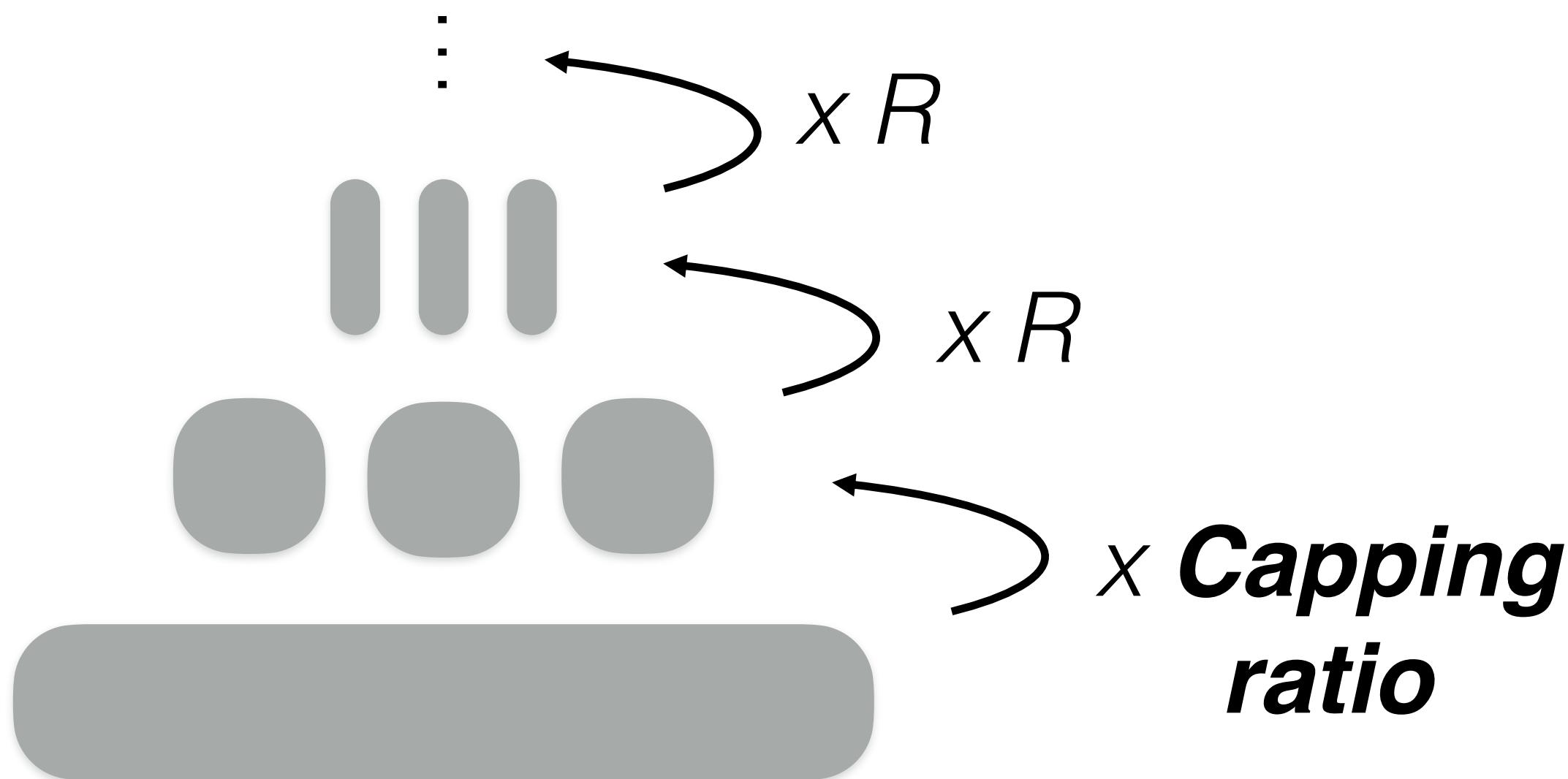
Squared Capped Lazy Leveling



SCLL

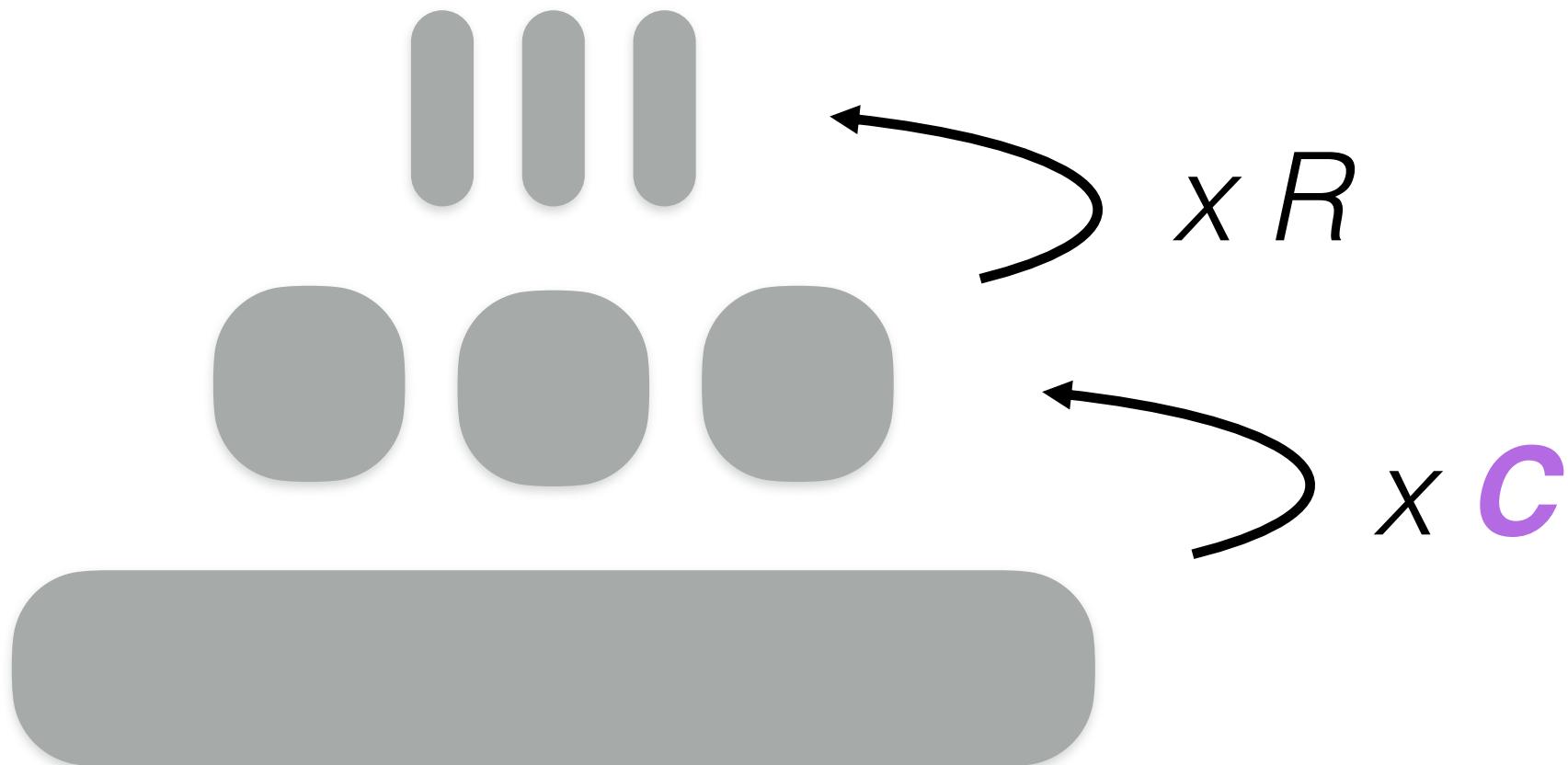


SCLL



SCLL

model:



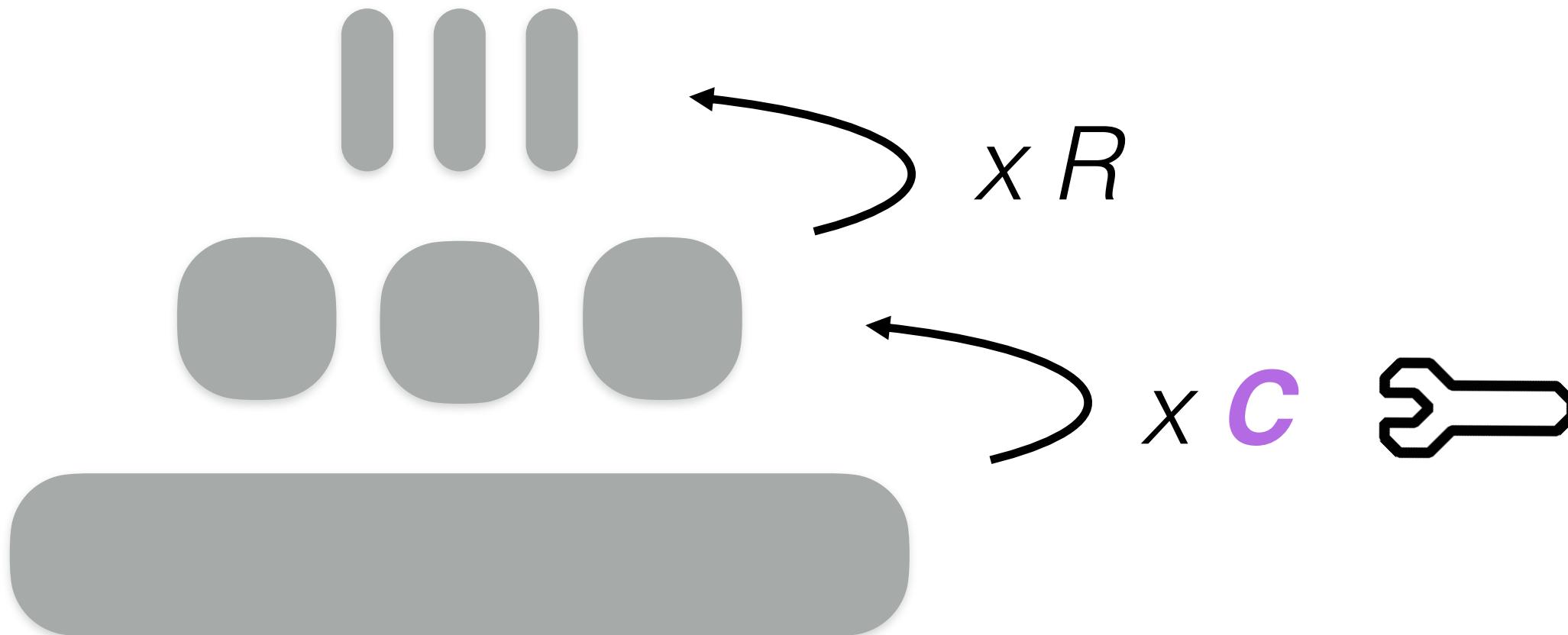
SCLL



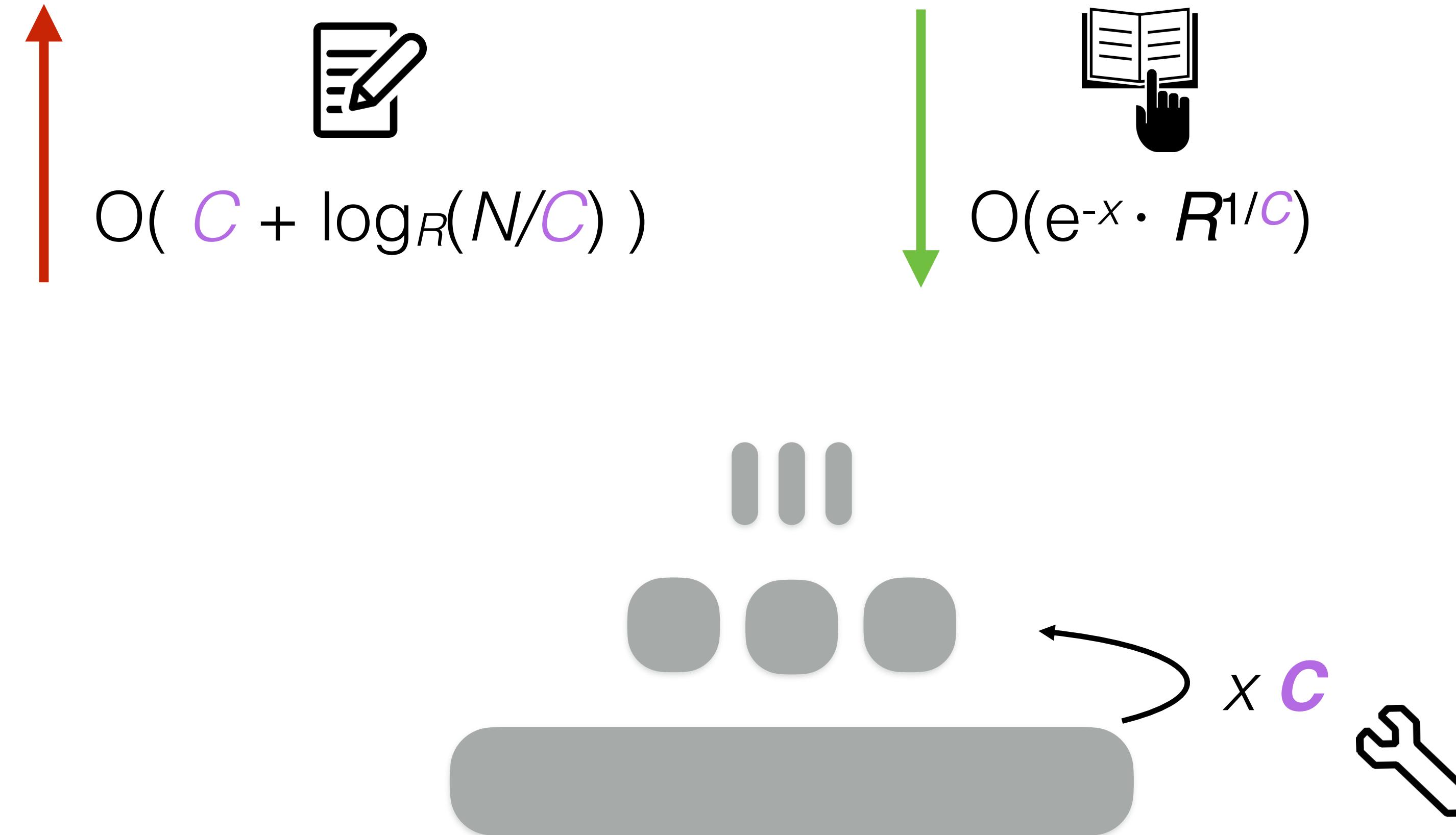
$$O(\textcolor{violet}{C} + \log_R(N/\textcolor{violet}{C}))$$



$$O(e^{-x} \cdot R^{1/\textcolor{violet}{C}})$$



SCLL



SCLL

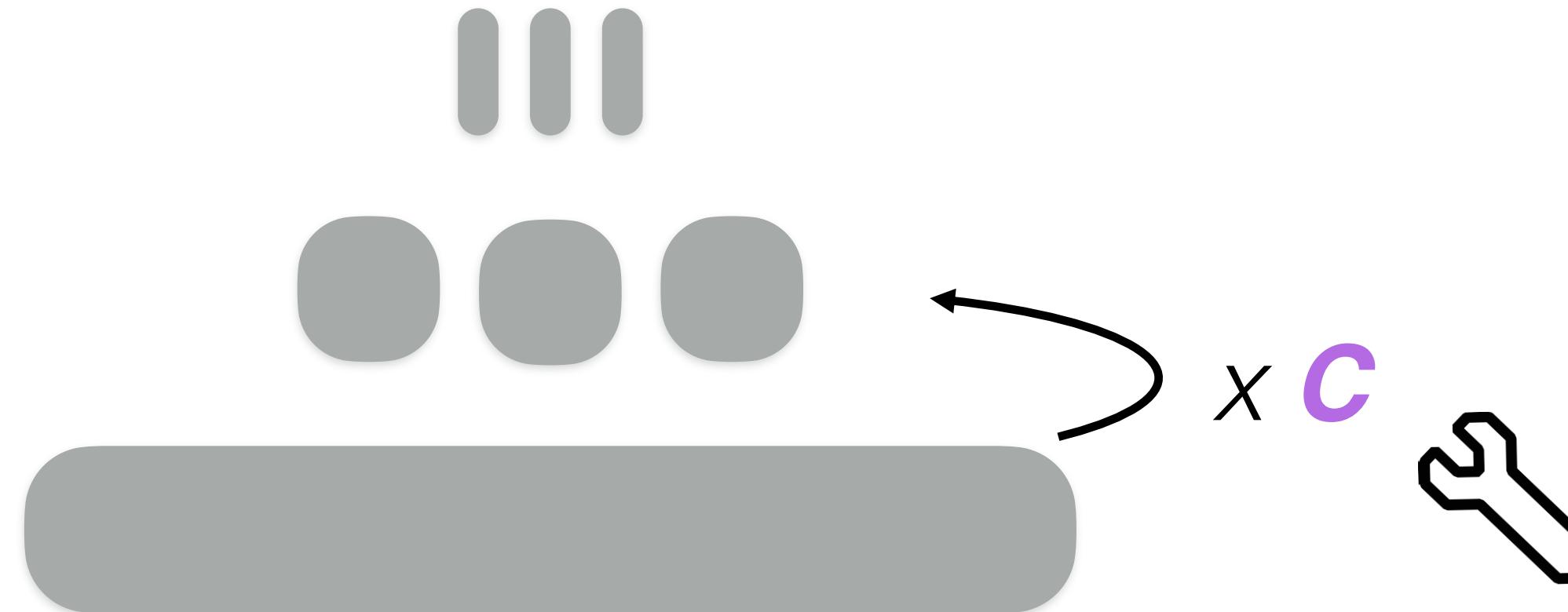


$$O(C + \log_R(N/C))$$



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

**costlier major
compactions**



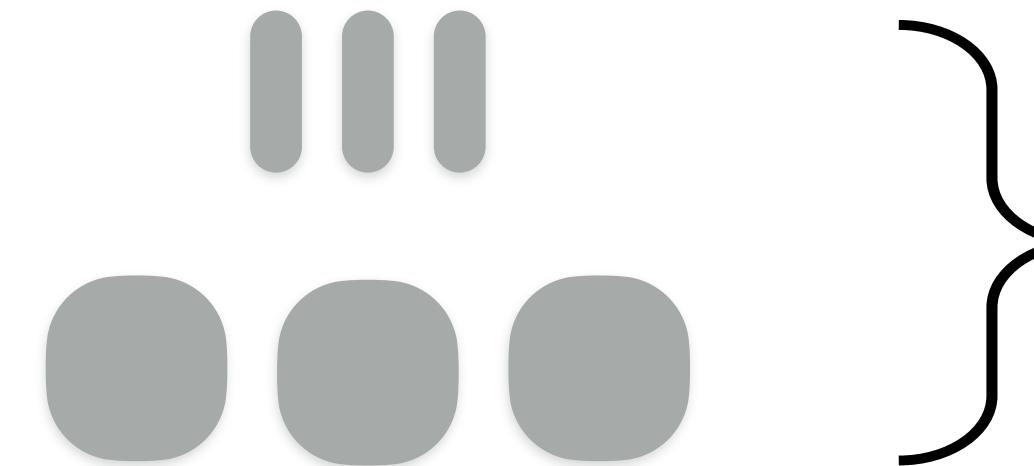
SCLL



$O(C + \log_R(N/C))$



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



**more effective
Bloom filters**

costlier major
compactions

SCLL



$$O(\ C + \log_R(N/C) \)$$

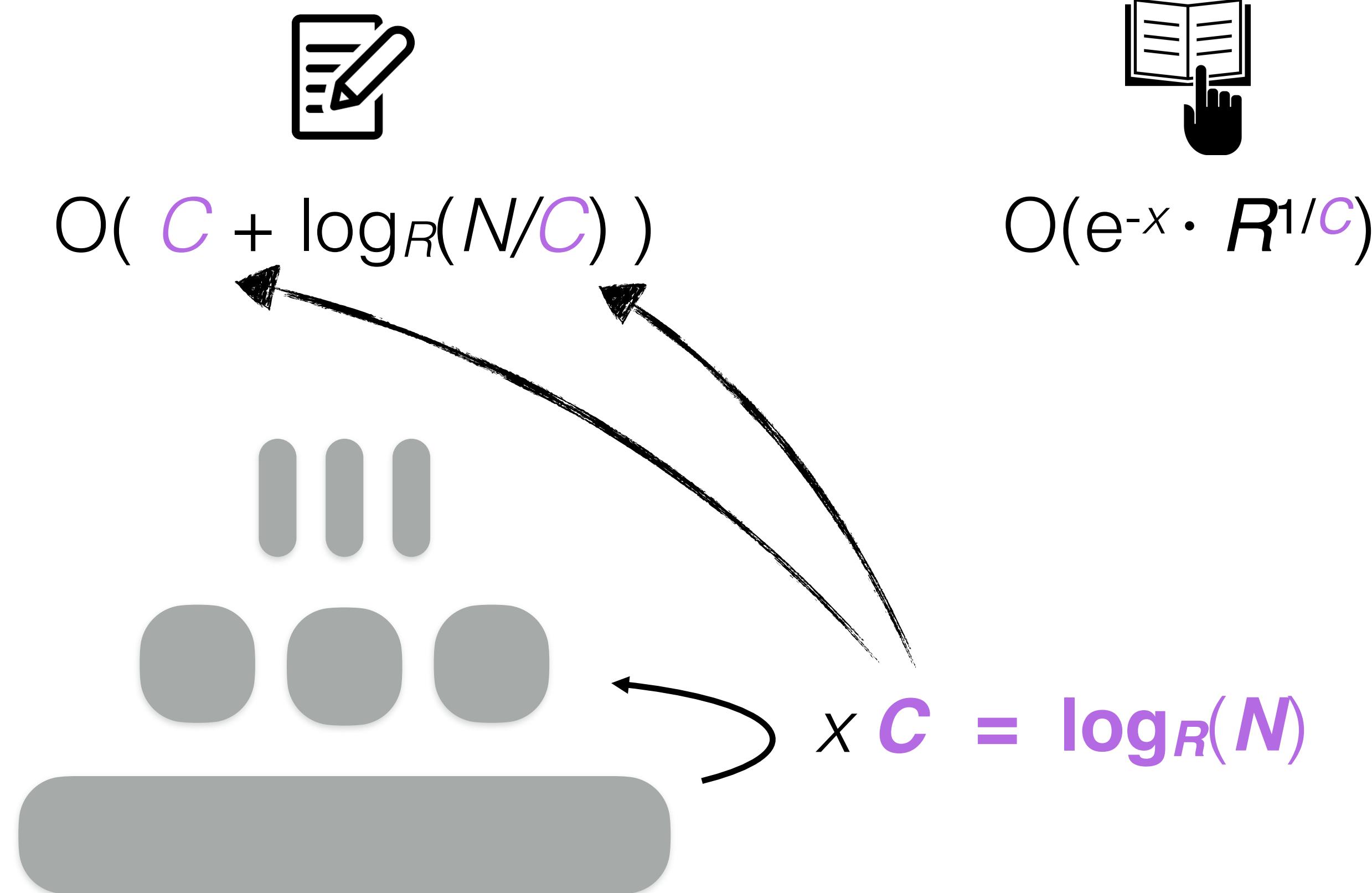


$$O(e^{-x} \cdot R^{1/C})$$



$$\times C = \log_R(N)$$

SCLL



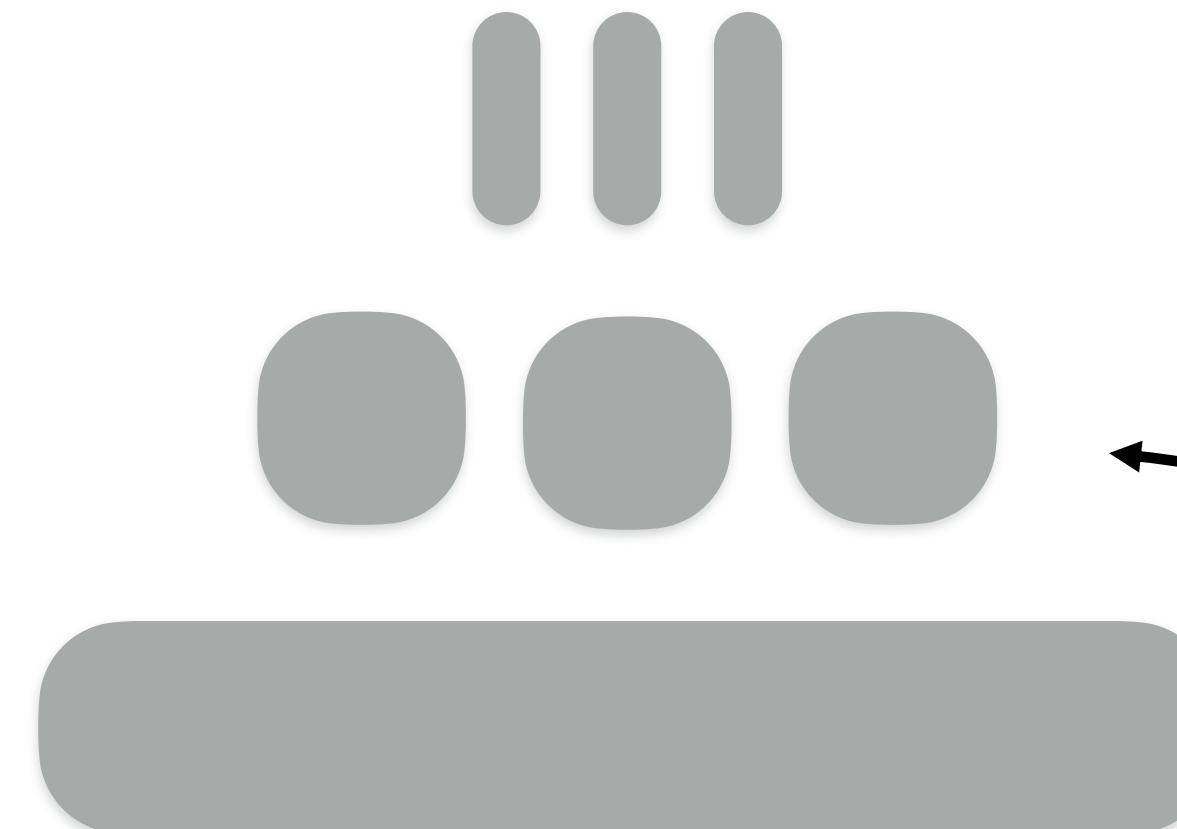
SCLL



$O(\log_R(N))$



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

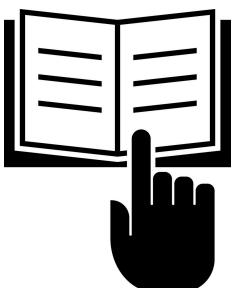


$$x C = \log_R(N)$$

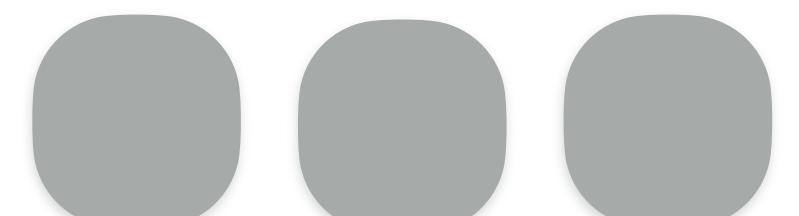
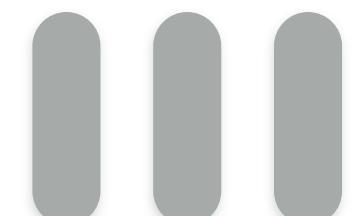
SCLL



$O(\log_R(N))$

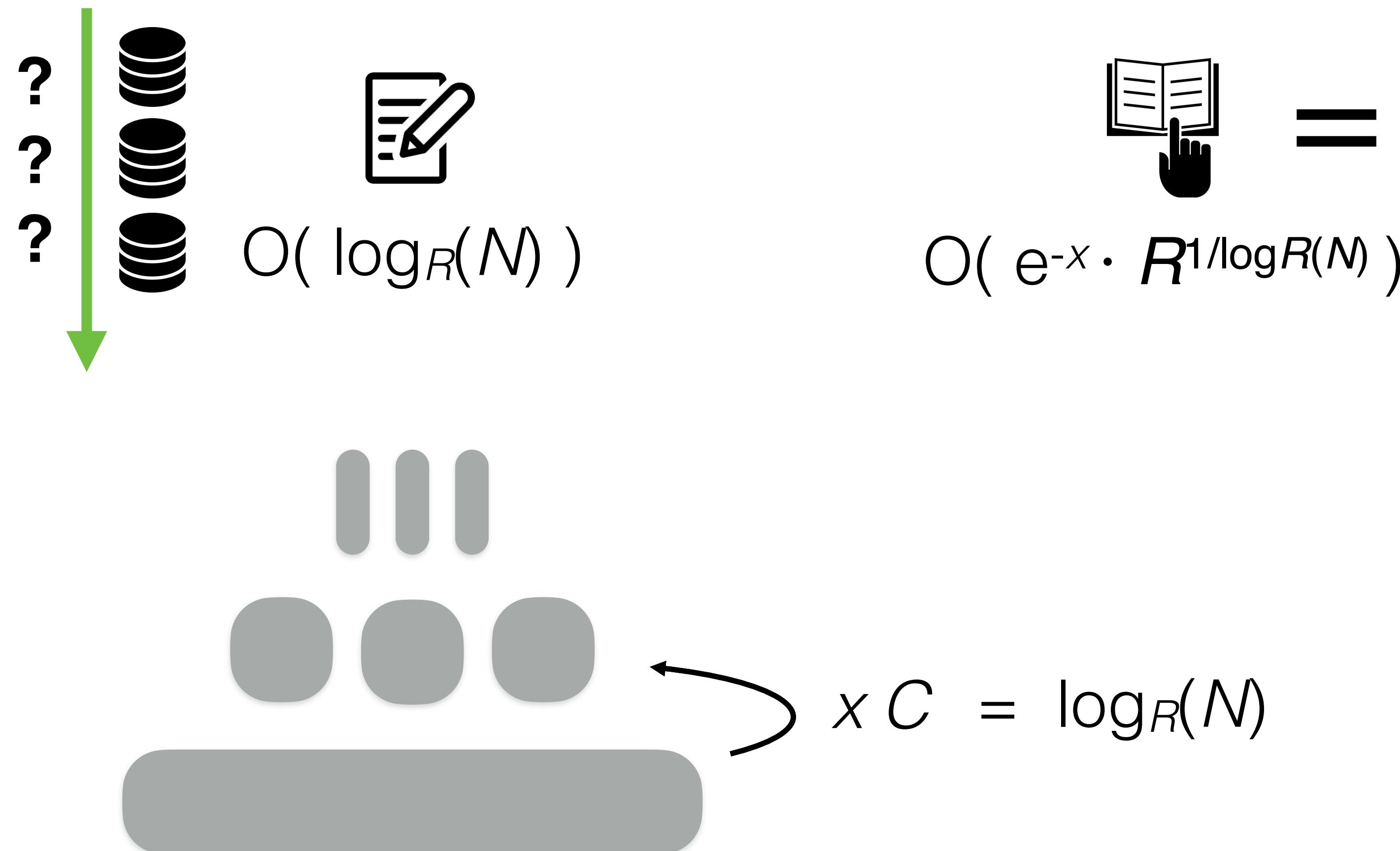


$O(e^{-x} \cdot R^{1/\log R(N)})$

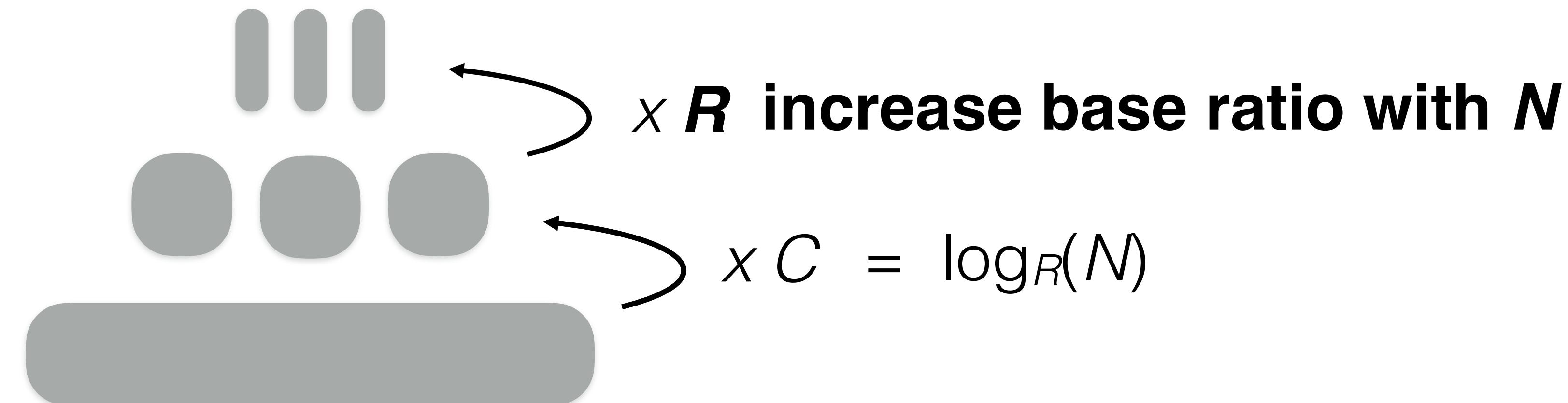
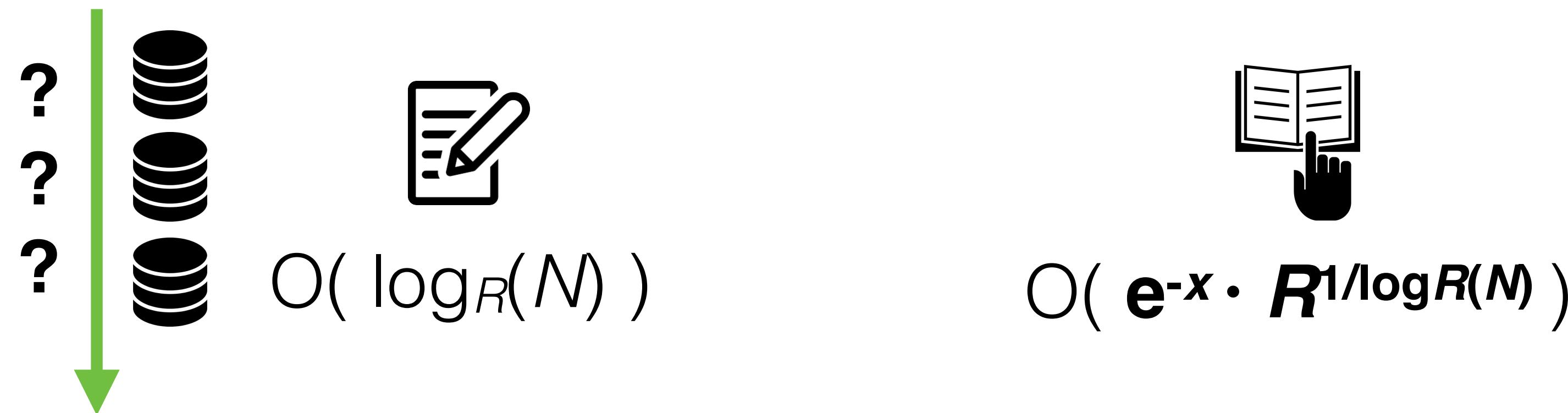


$$x C = \log_R(N)$$

SCLL



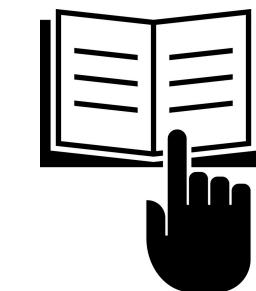
SCLL



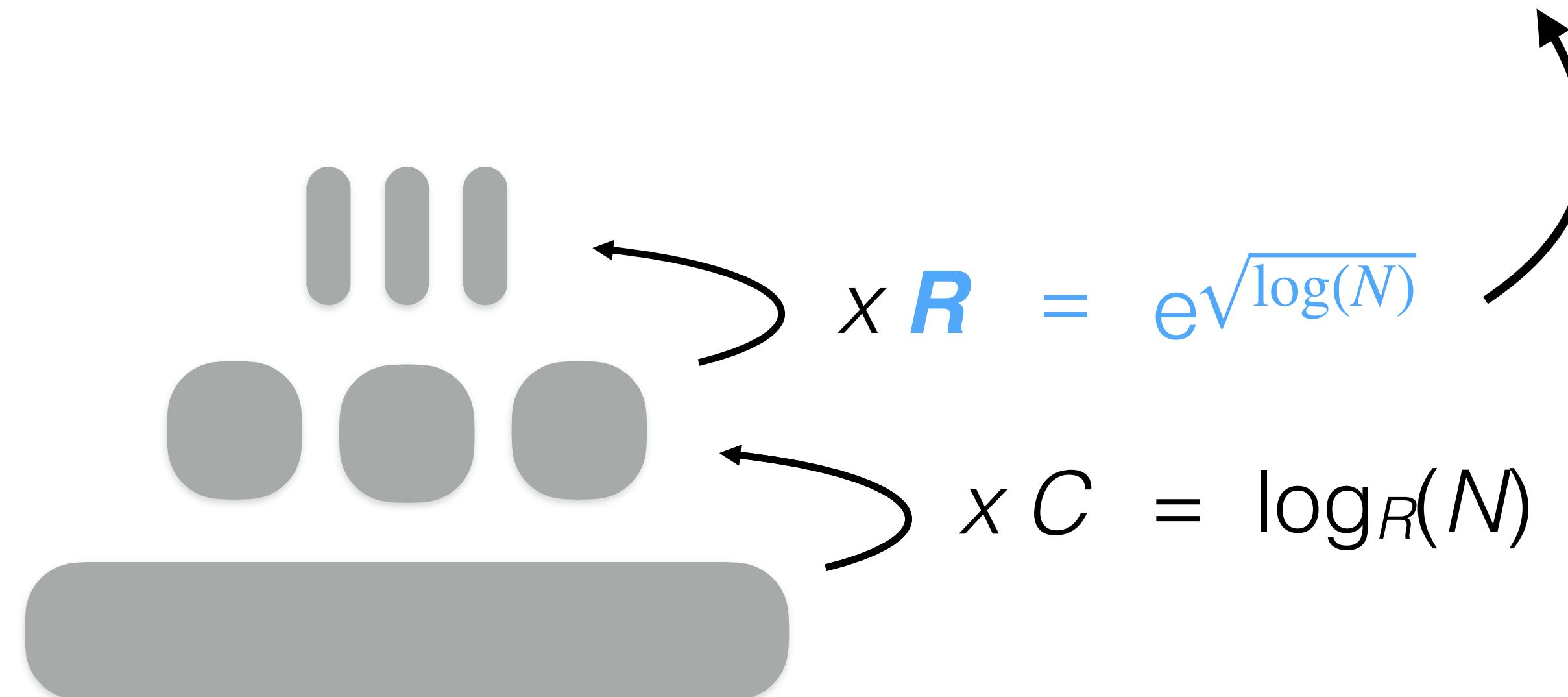
SCLL



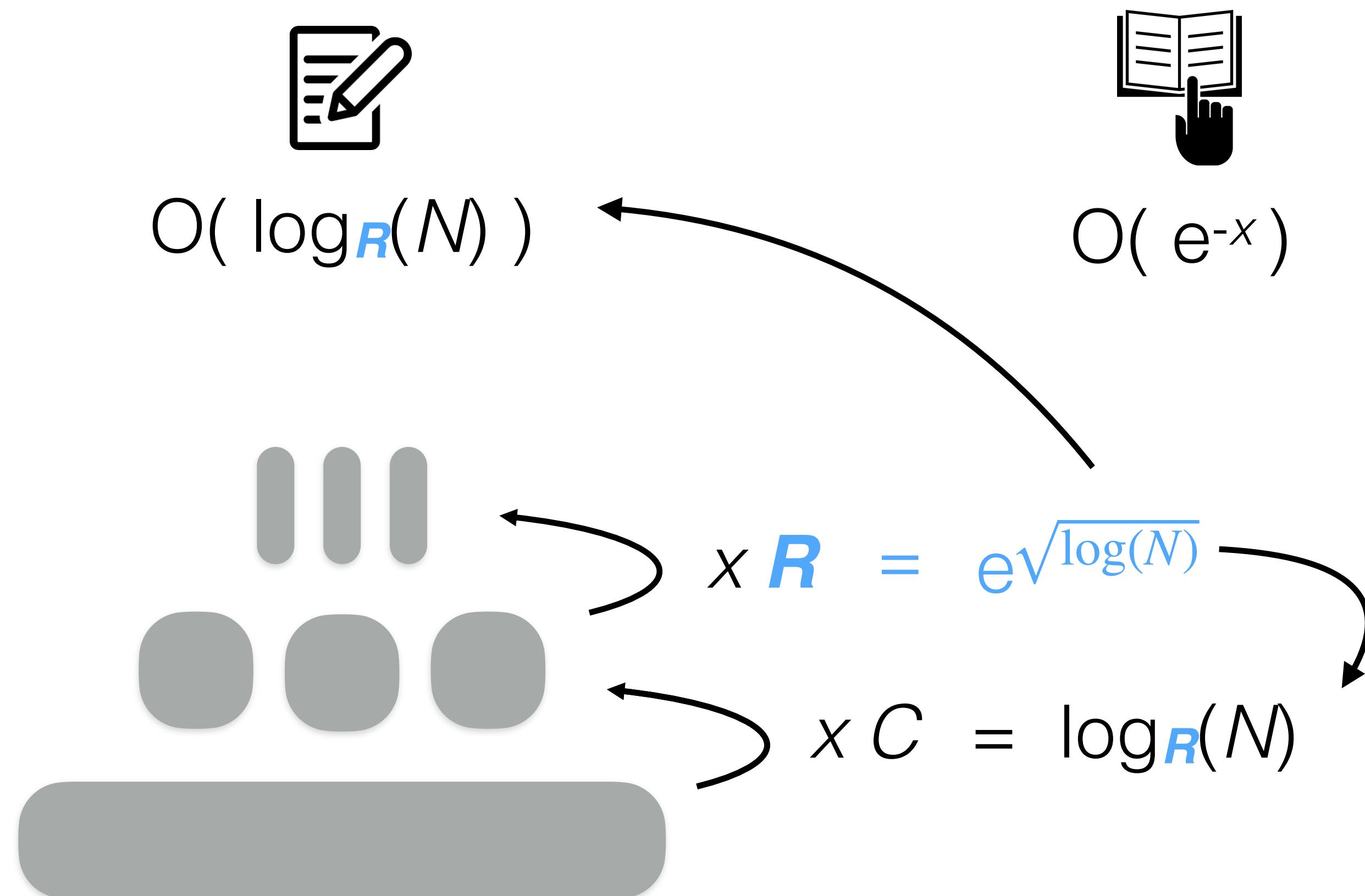
$O(\log_R(N))$



$O(e^{-x} \cdot R^{1/\log R(N)})$



SCLL



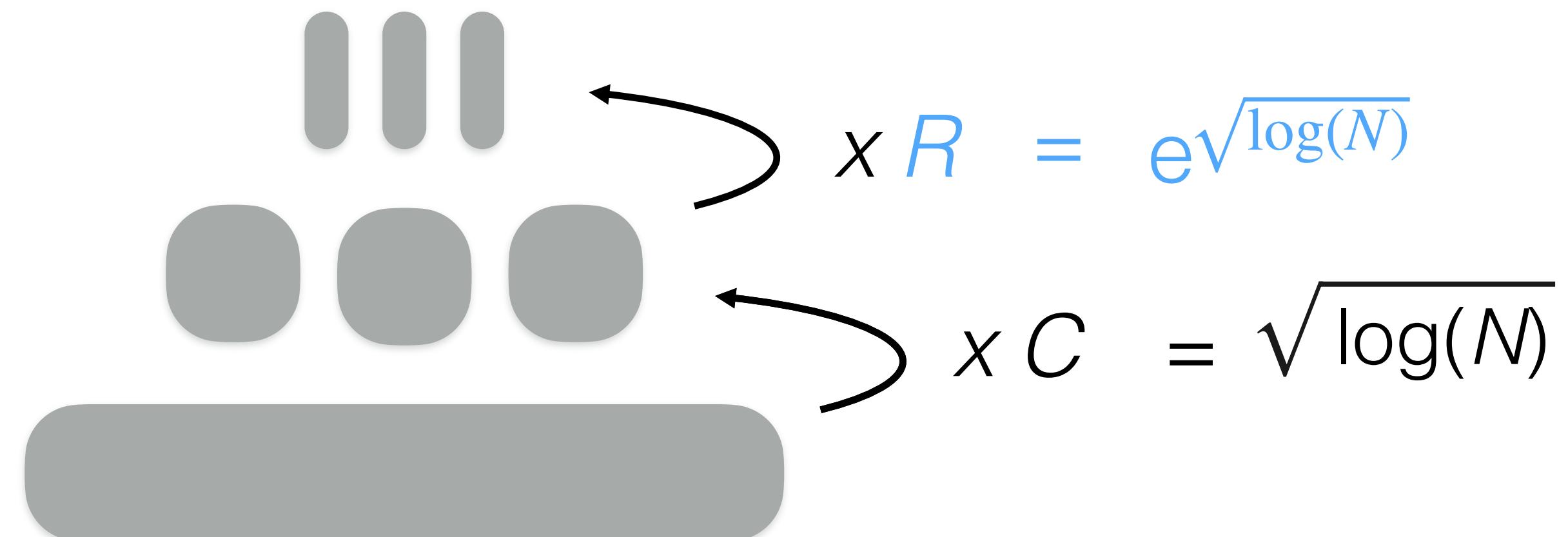
SCLL



$O(\sqrt{\log(N)})$



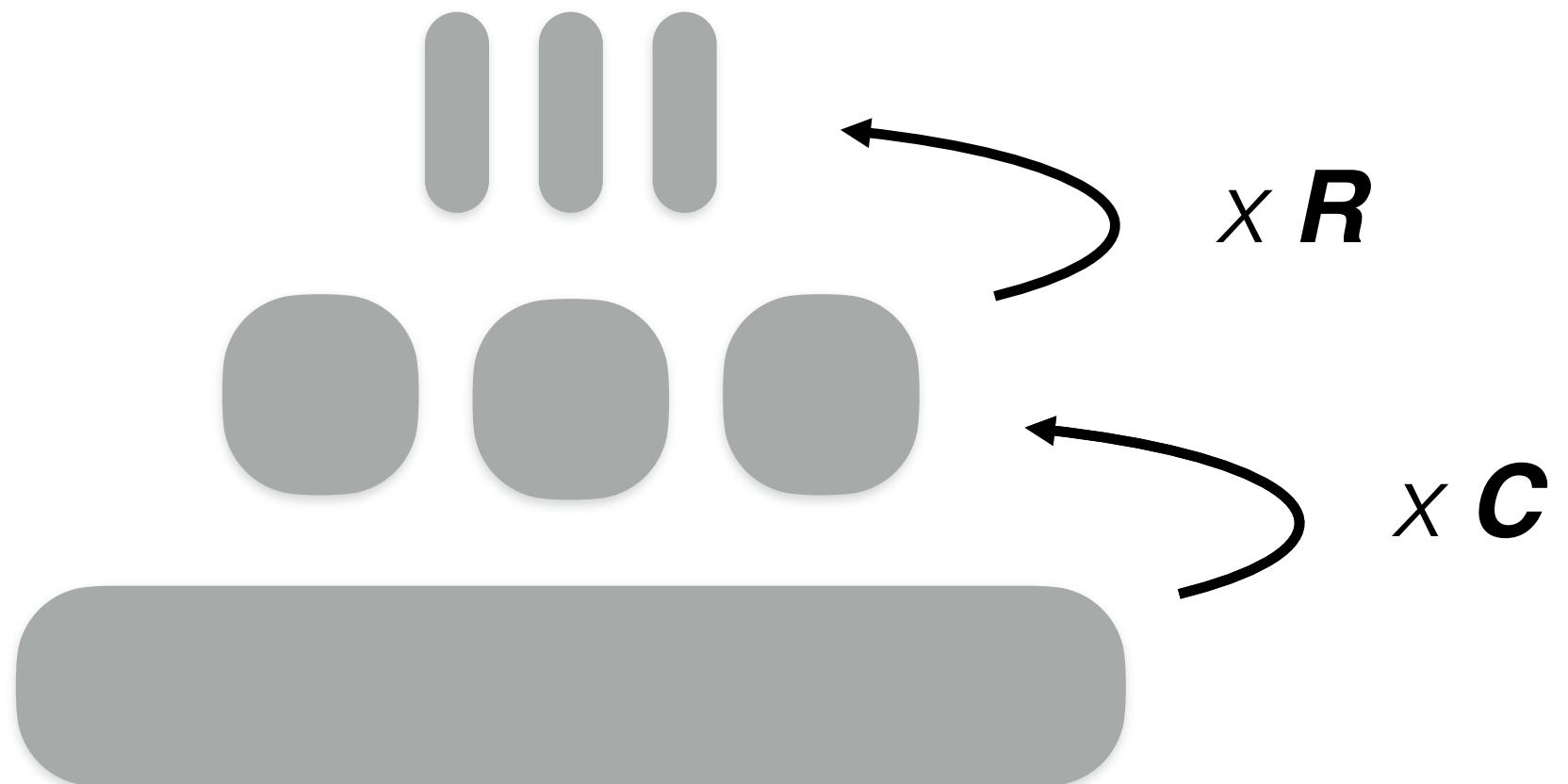
$O(e^{-x})$



SCLL

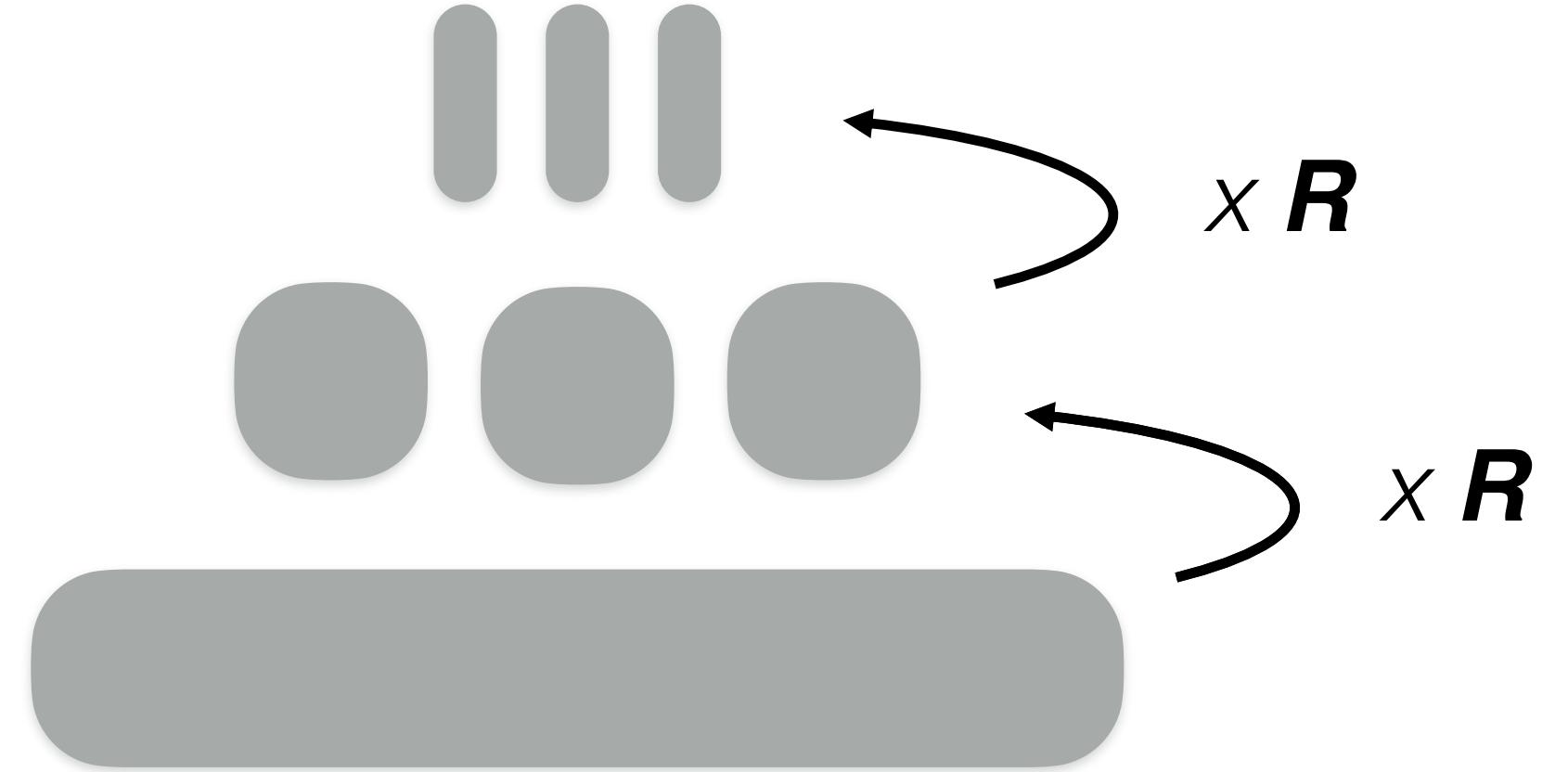


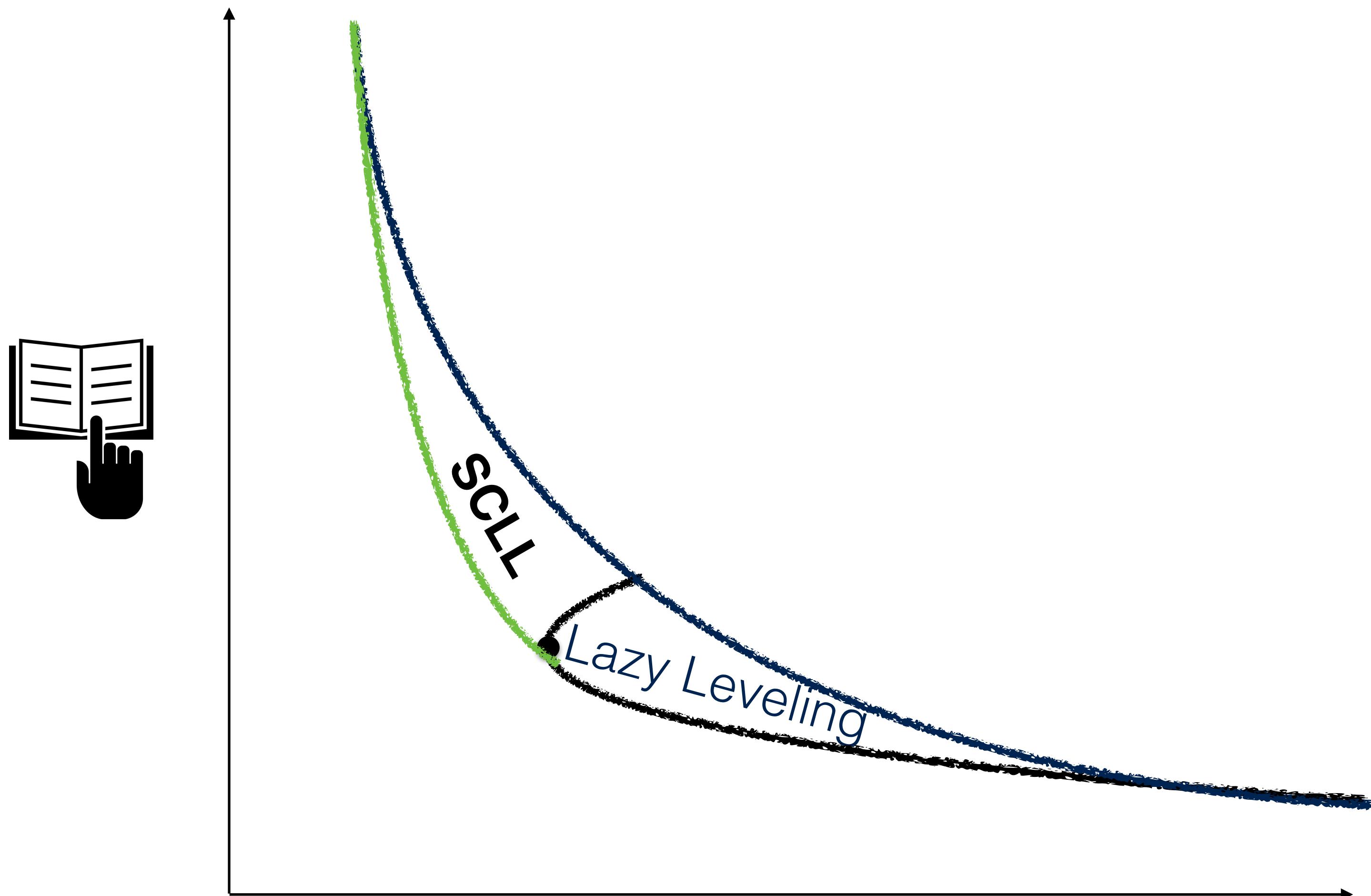
$O(\sqrt{\log(N)})$

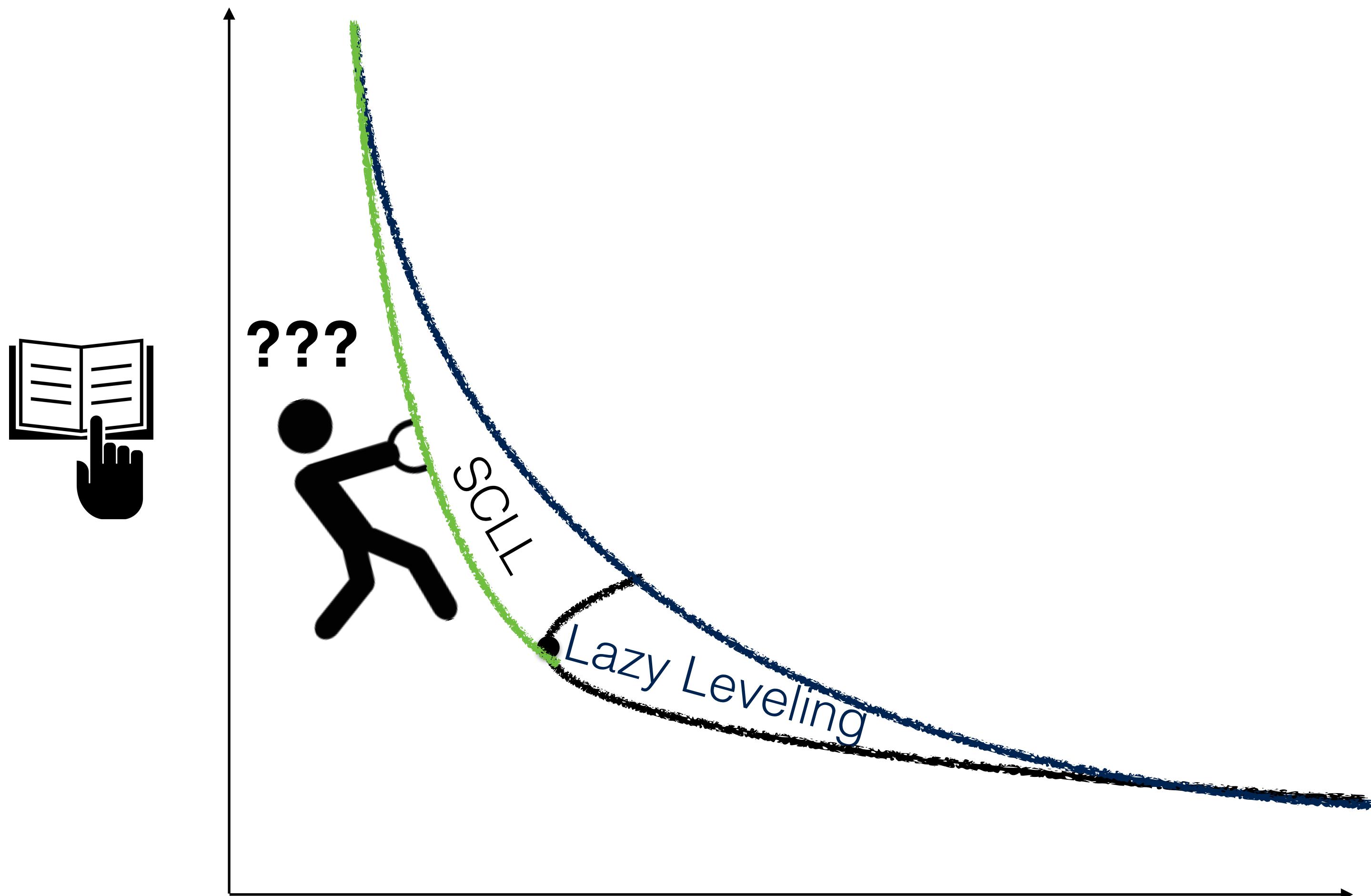


Traditional LSM

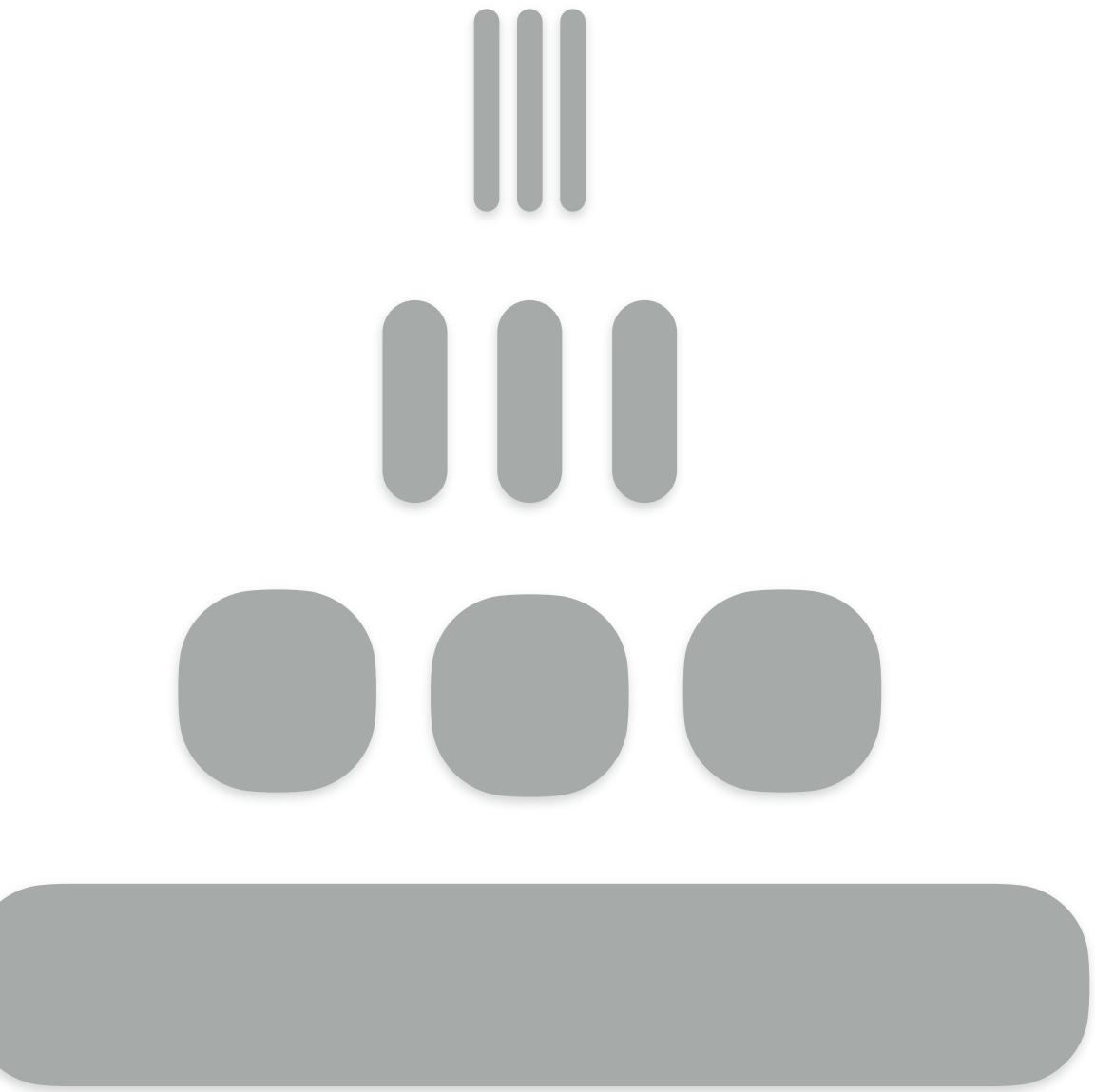
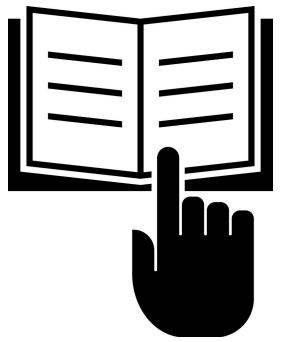
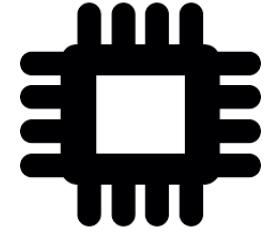
$O(\log(N))$



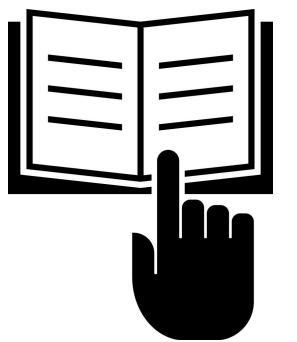
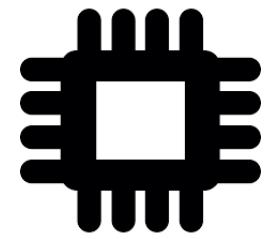




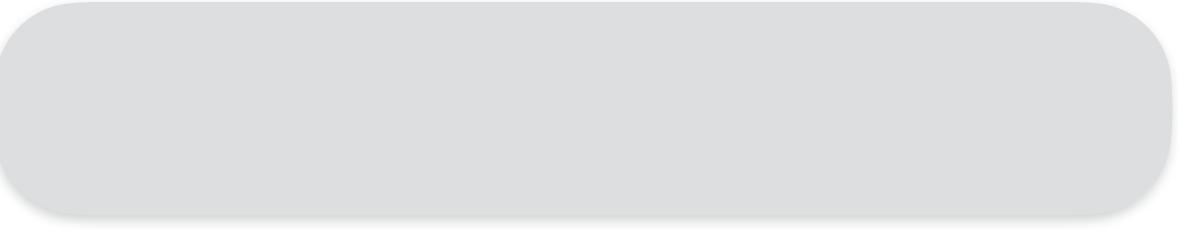
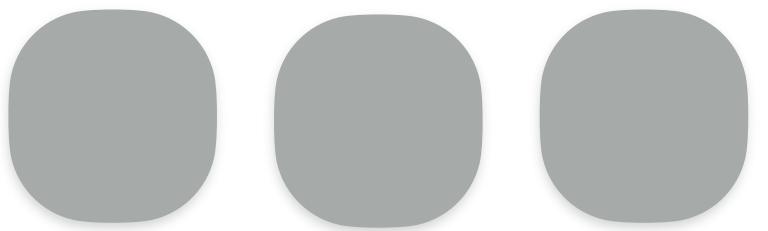
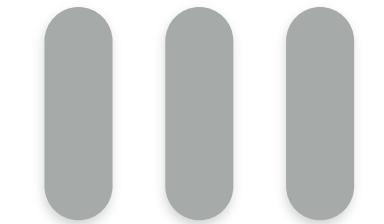
SCLL & Lazy Leveling **Cost Breakdown**



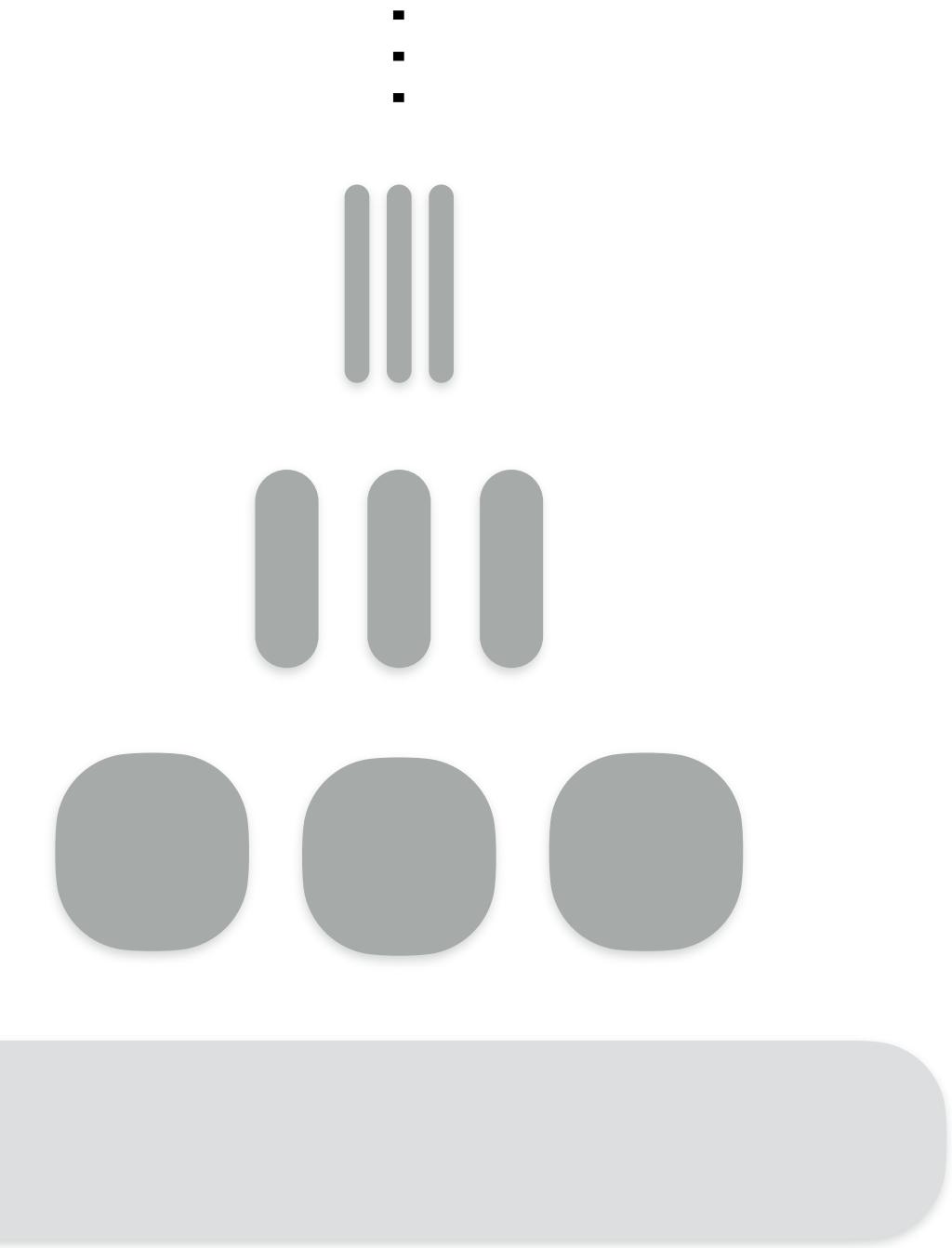
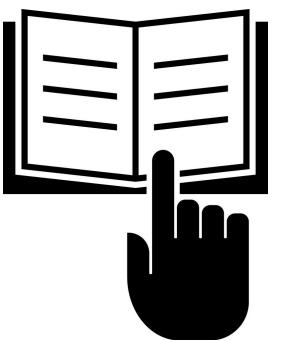
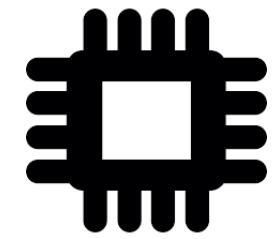
SCLL & Lazy Leveling Cost Breakdown



:



SCLL & Lazy Leveling Cost Breakdown



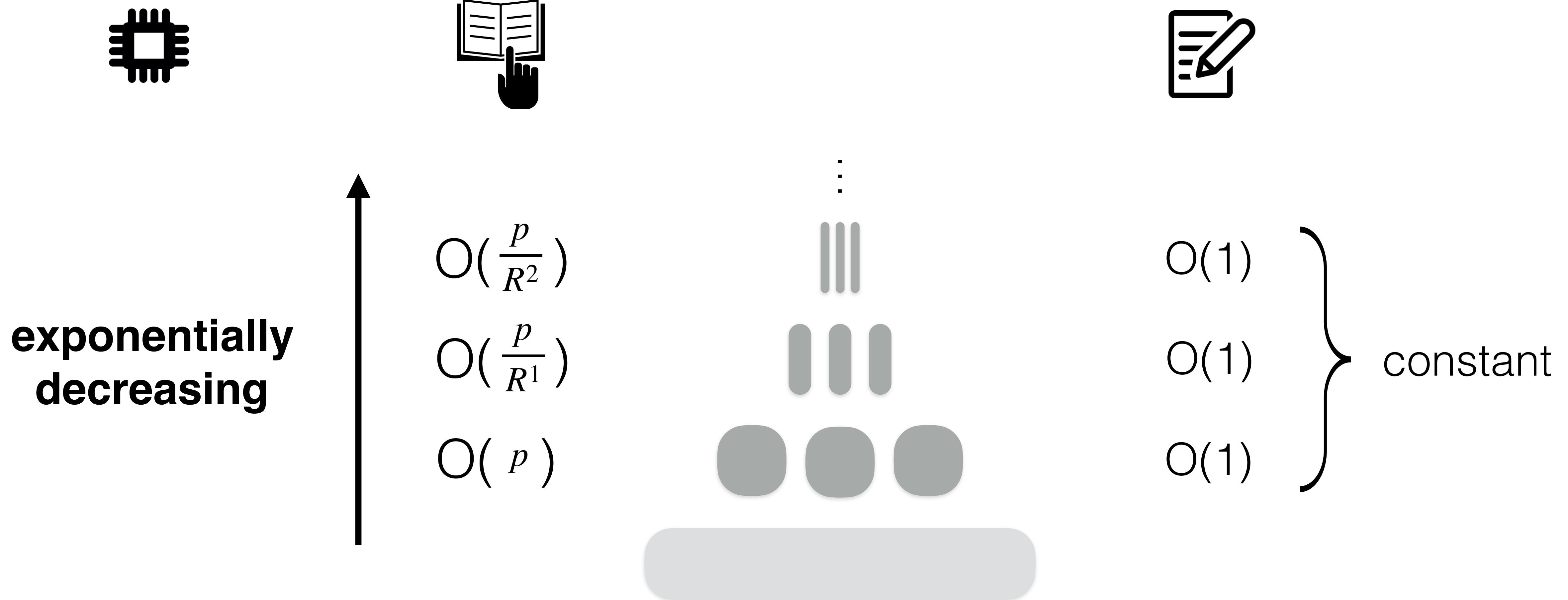
$O(1)$

$O(1)$

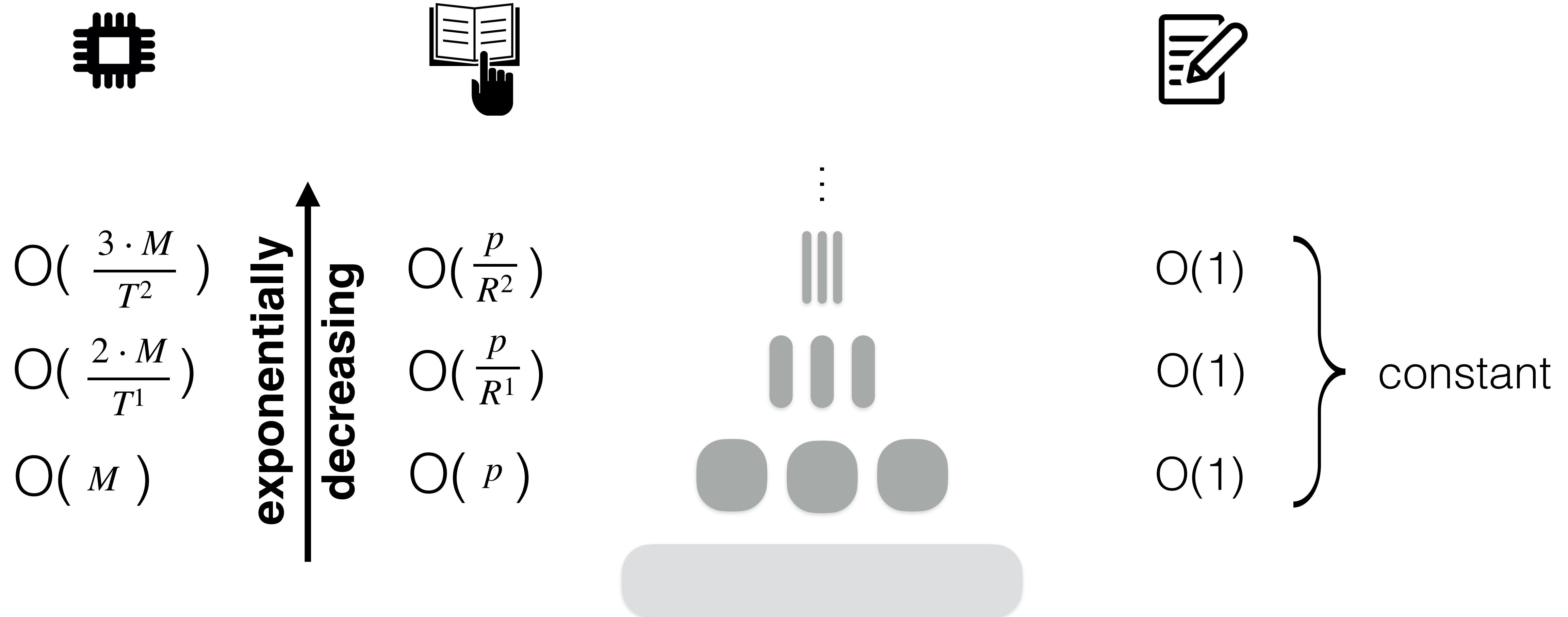
$O(1)$

constant

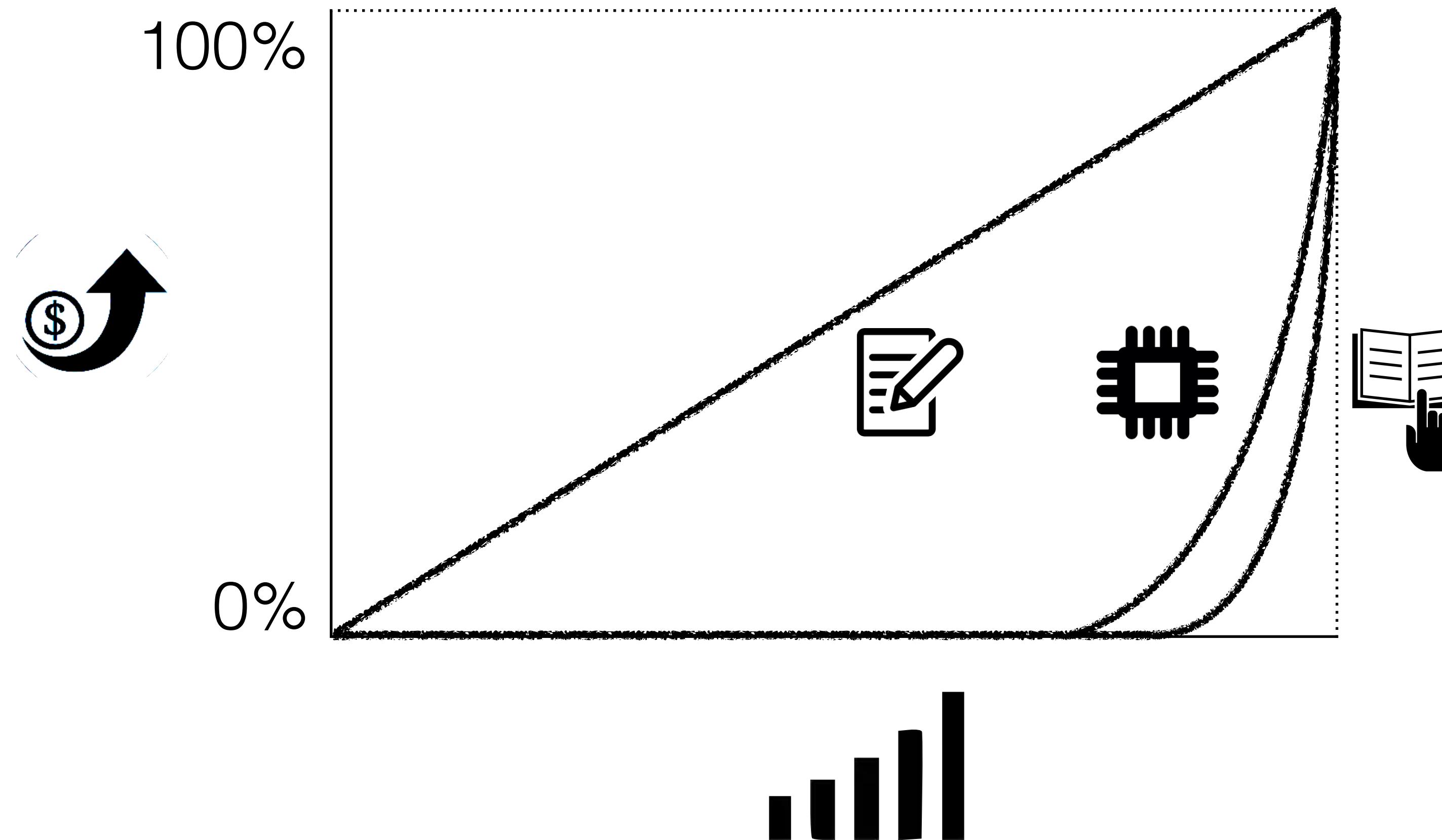
SCLL & Lazy Leveling Cost Breakdown



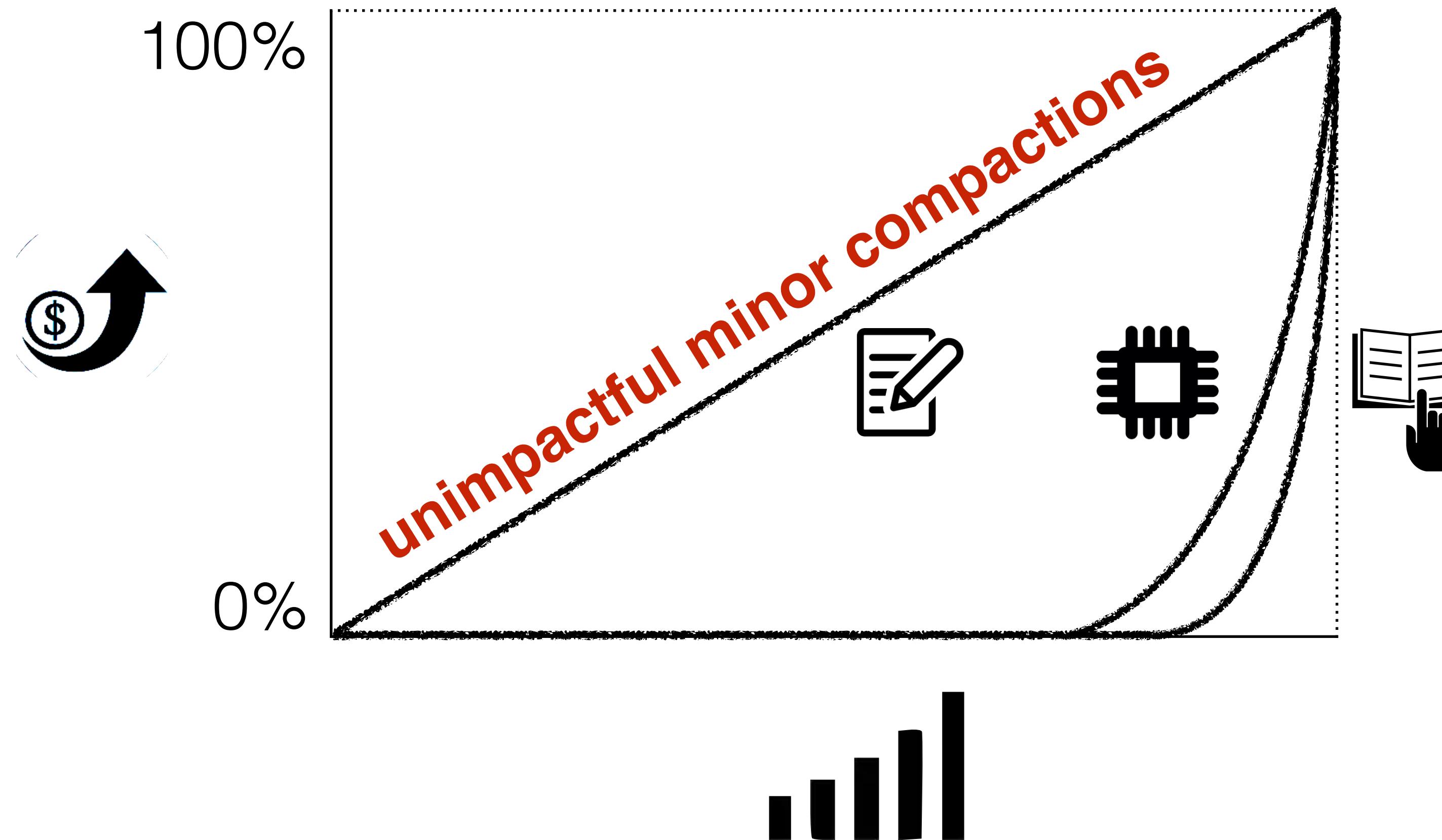
SCLL & Lazy Leveling Cost Breakdown



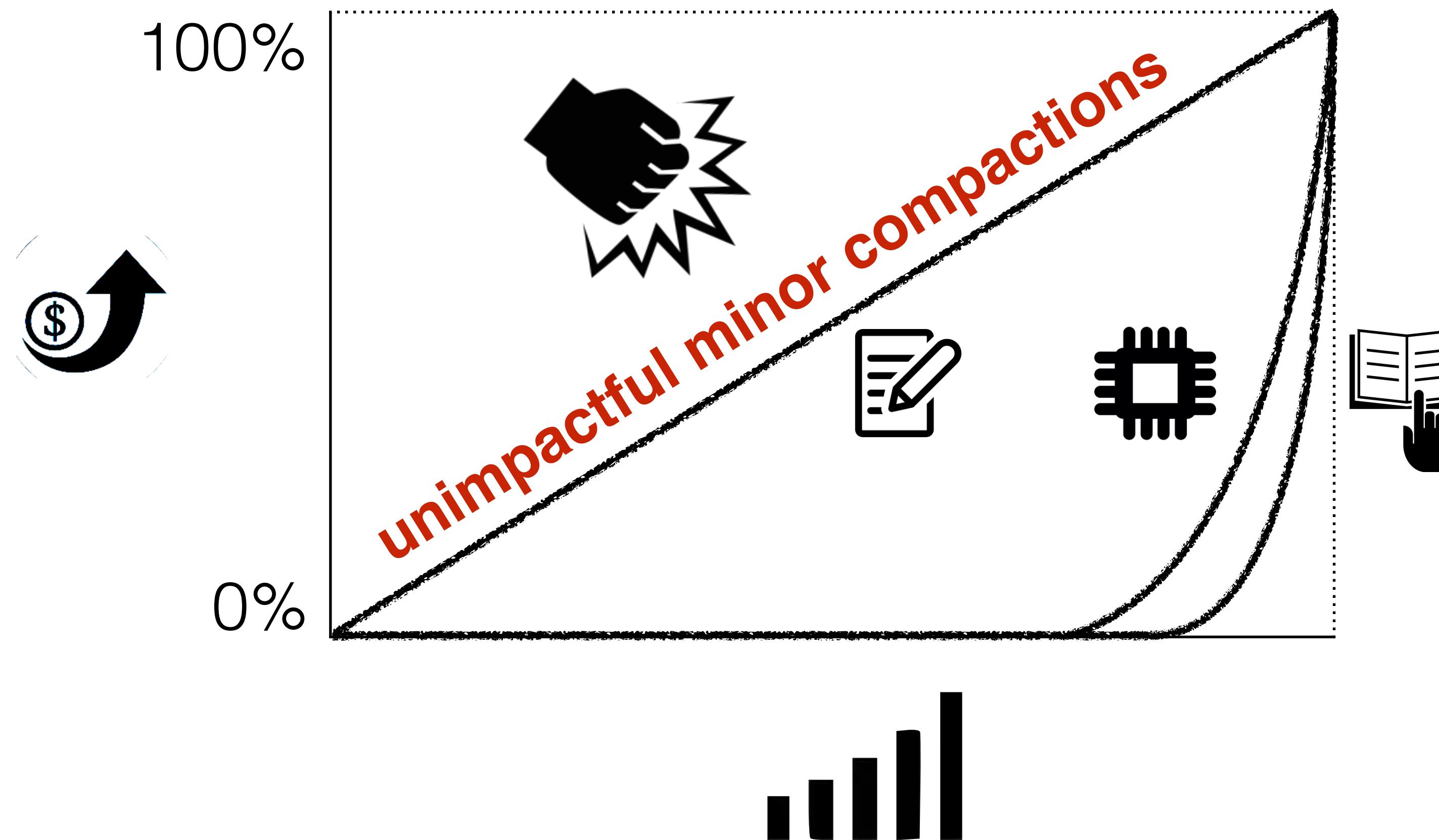
SCLL & Lazy Leveling **Cost Breakdown**



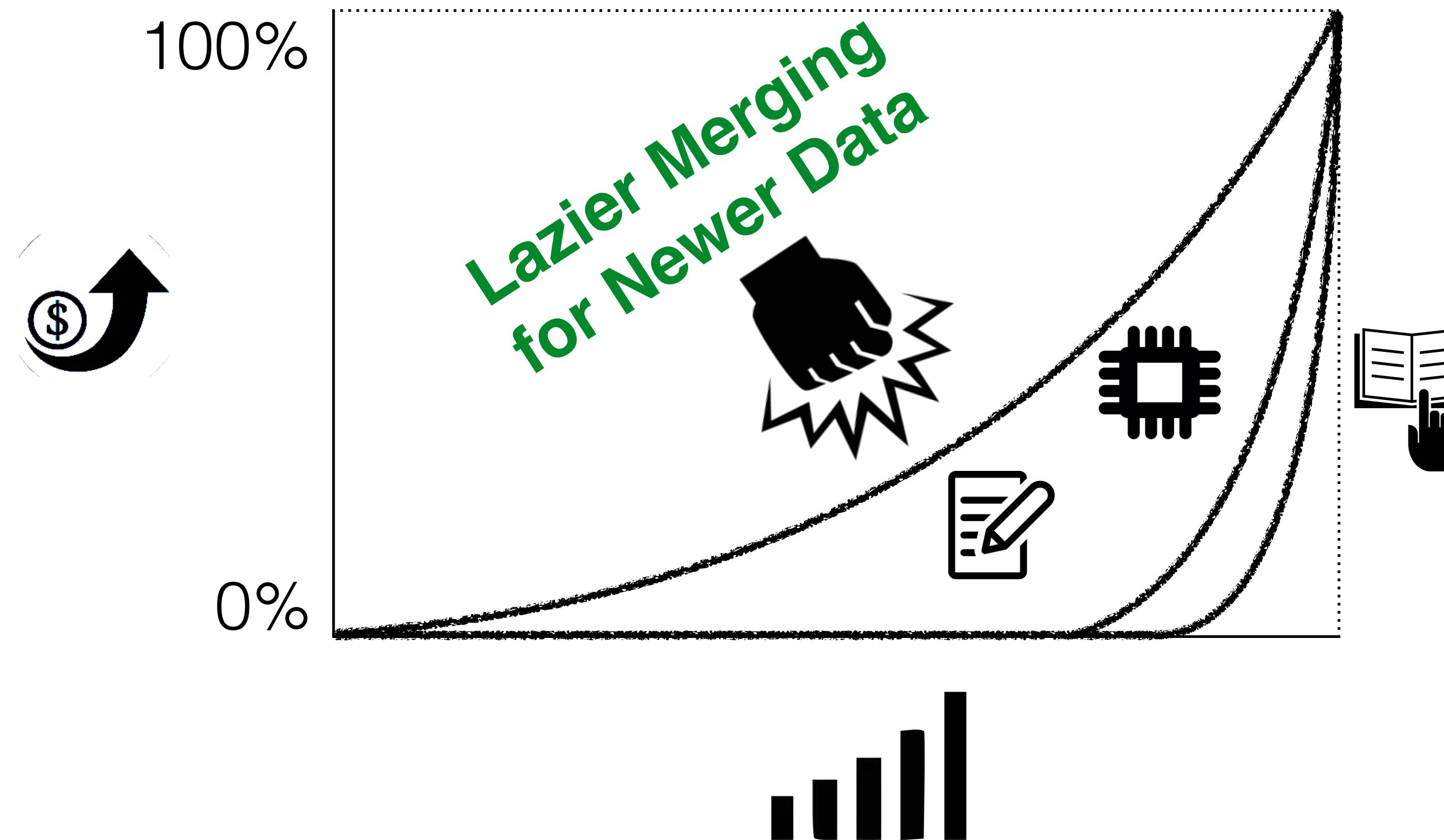
SCLL & Lazy Leveling Cost Breakdown



SCLL & Lazy Leveling Cost Breakdown

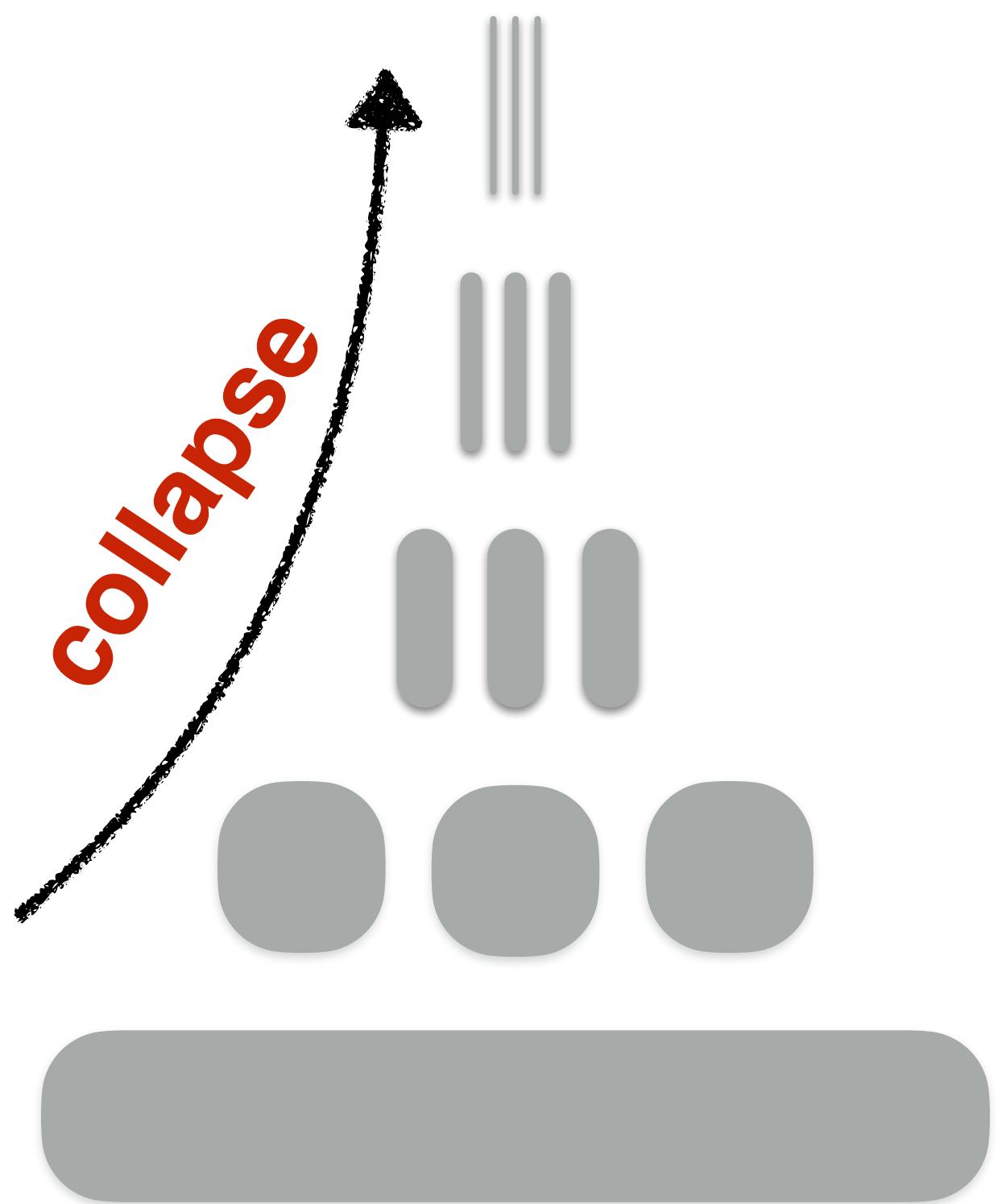


SCLL & Lazy Leveling Cost Breakdown

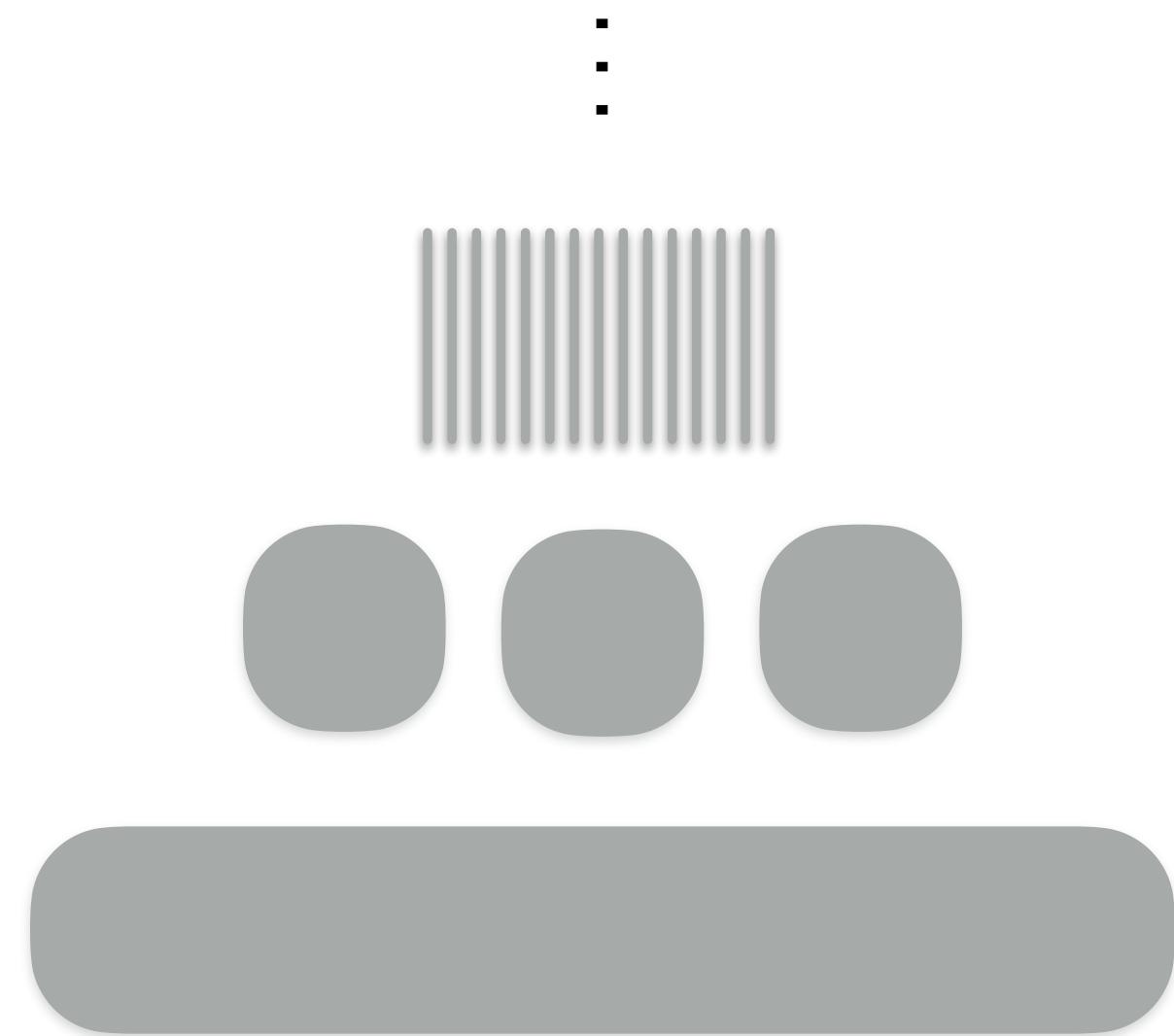


LSM-bush - Lazier Merging for Newer Data

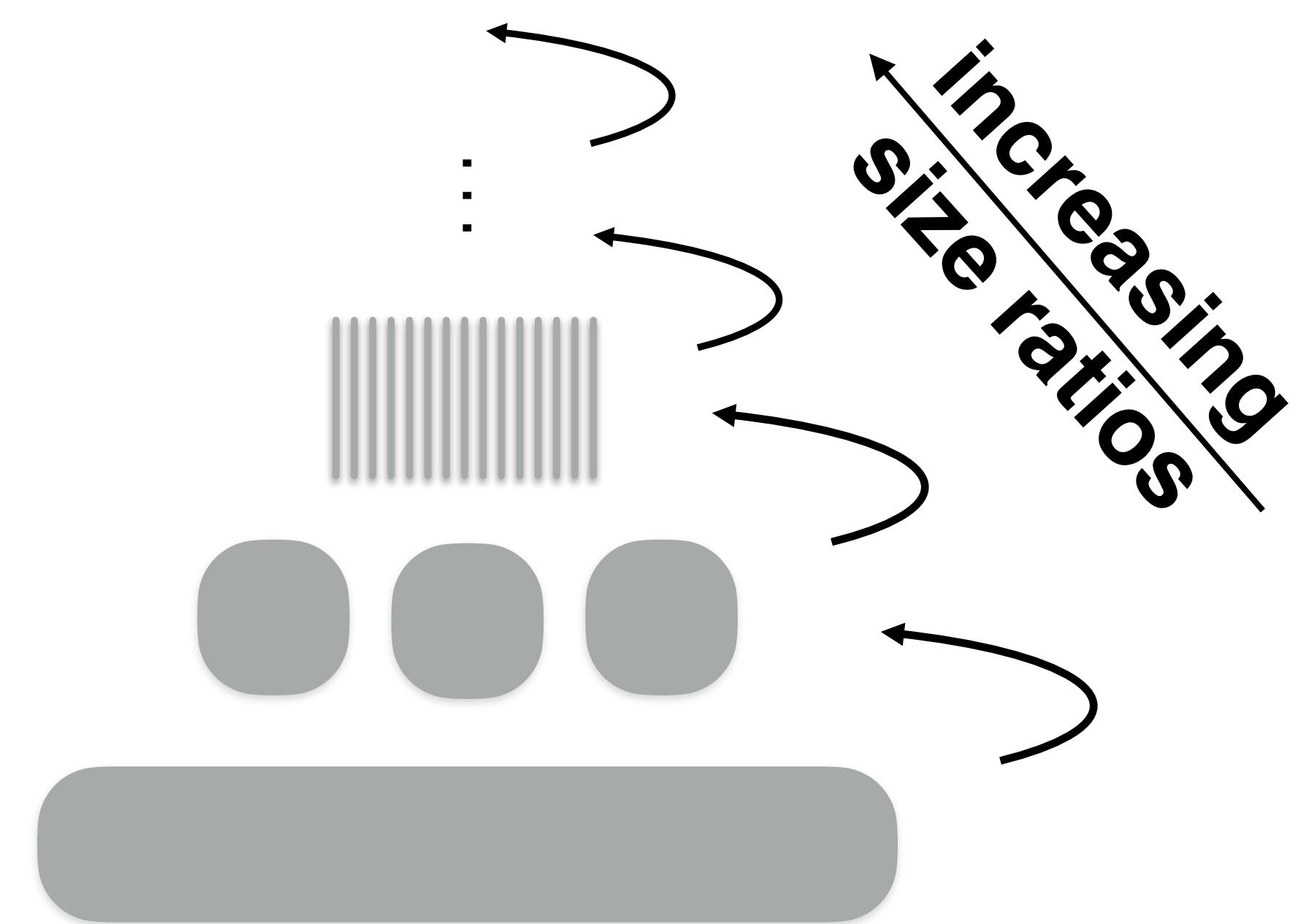
LSM-bush - Lazier Merging for Newer Data



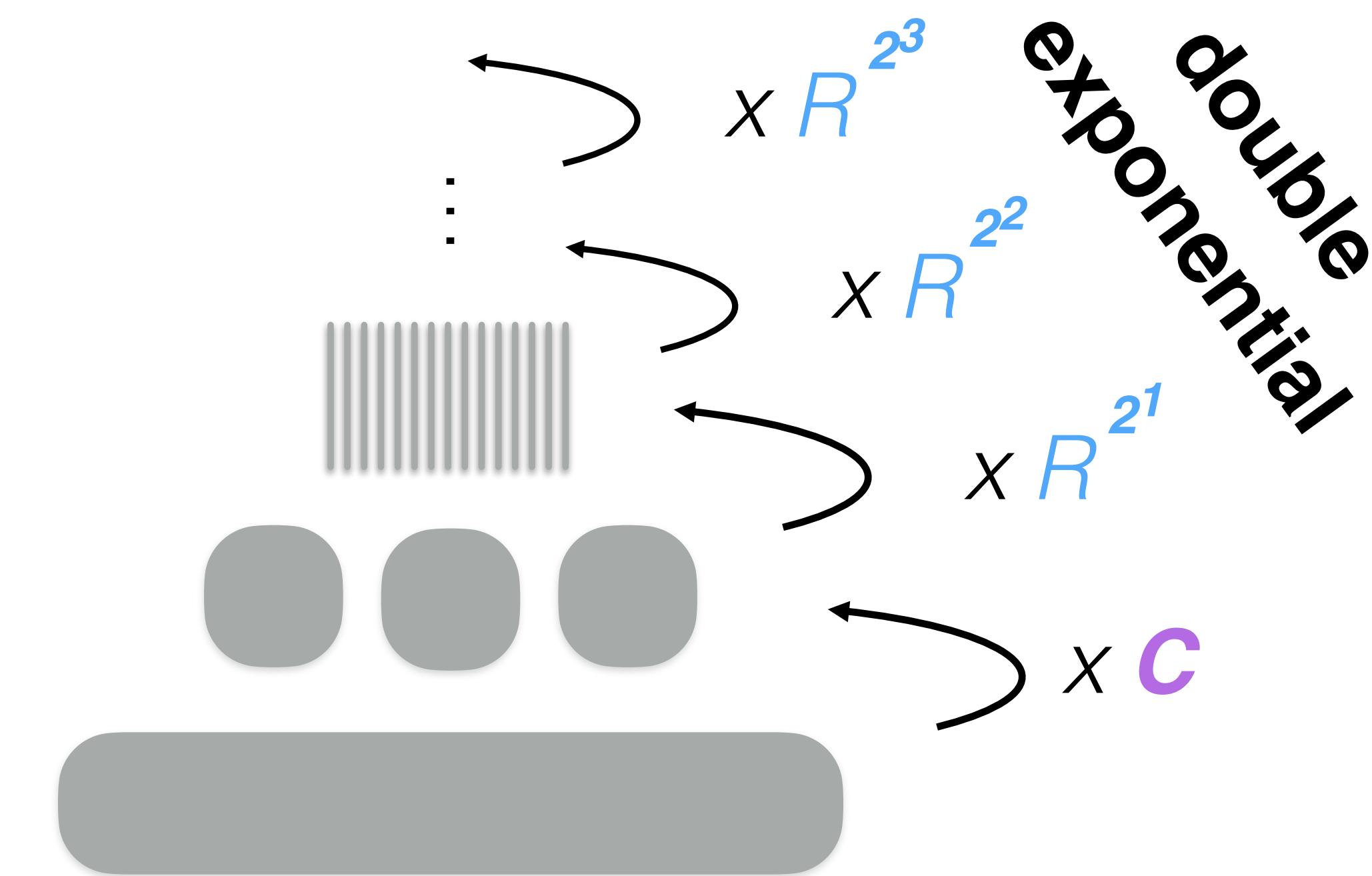
LSM-bush - Lazier Merging for Newer Data



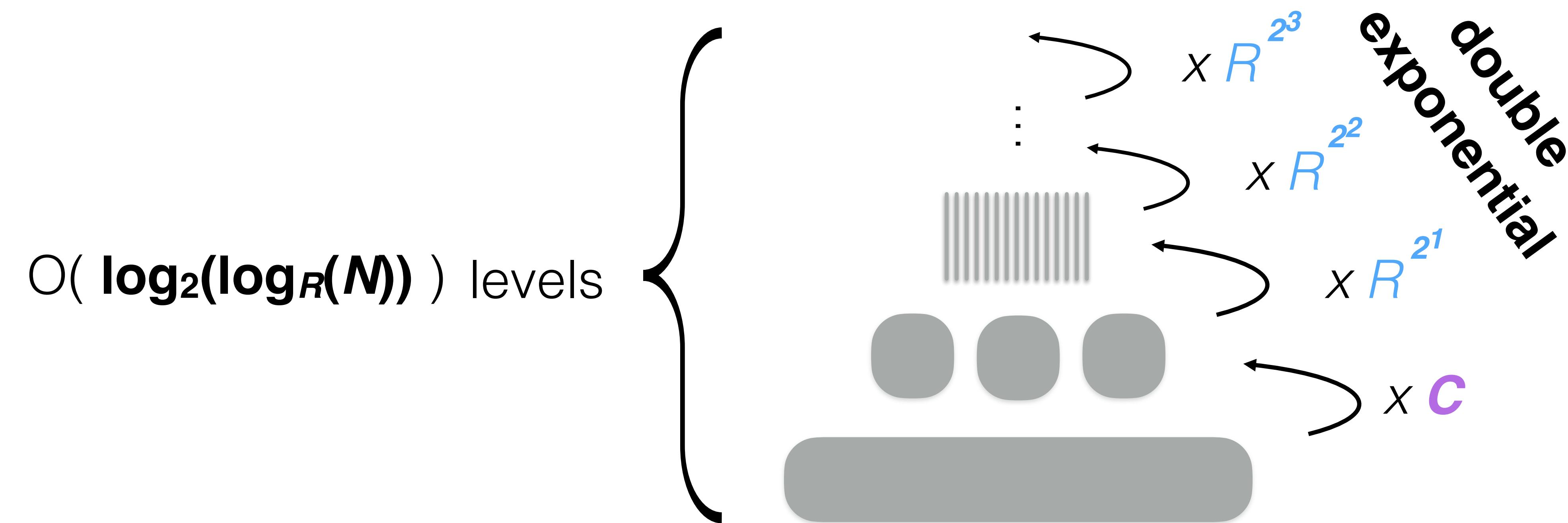
LSM-bush - Lazier Merging for Newer Data



LSM-bush - Lazier Merging for Newer Data



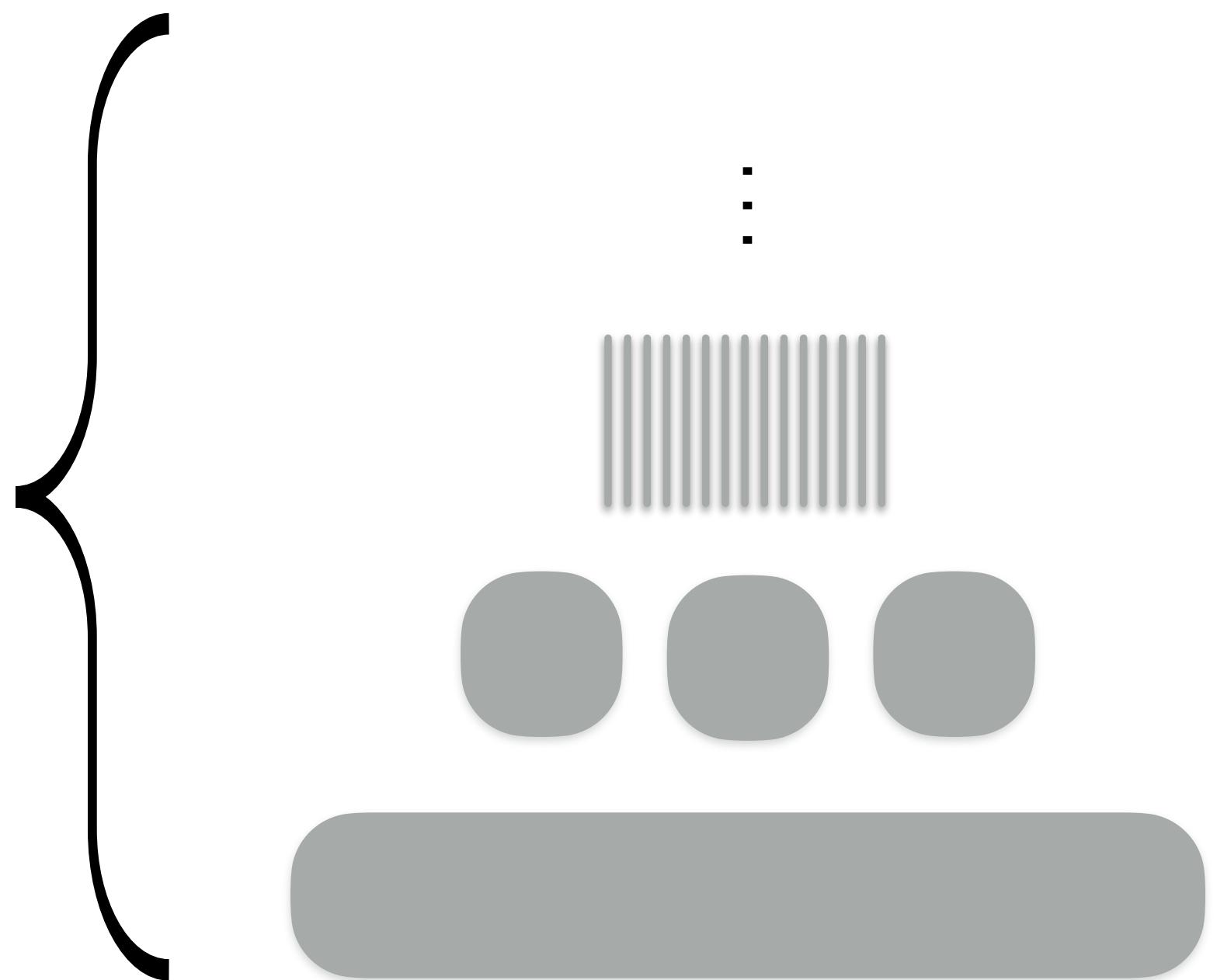
LSM-bush - Lazier Merging for Newer Data



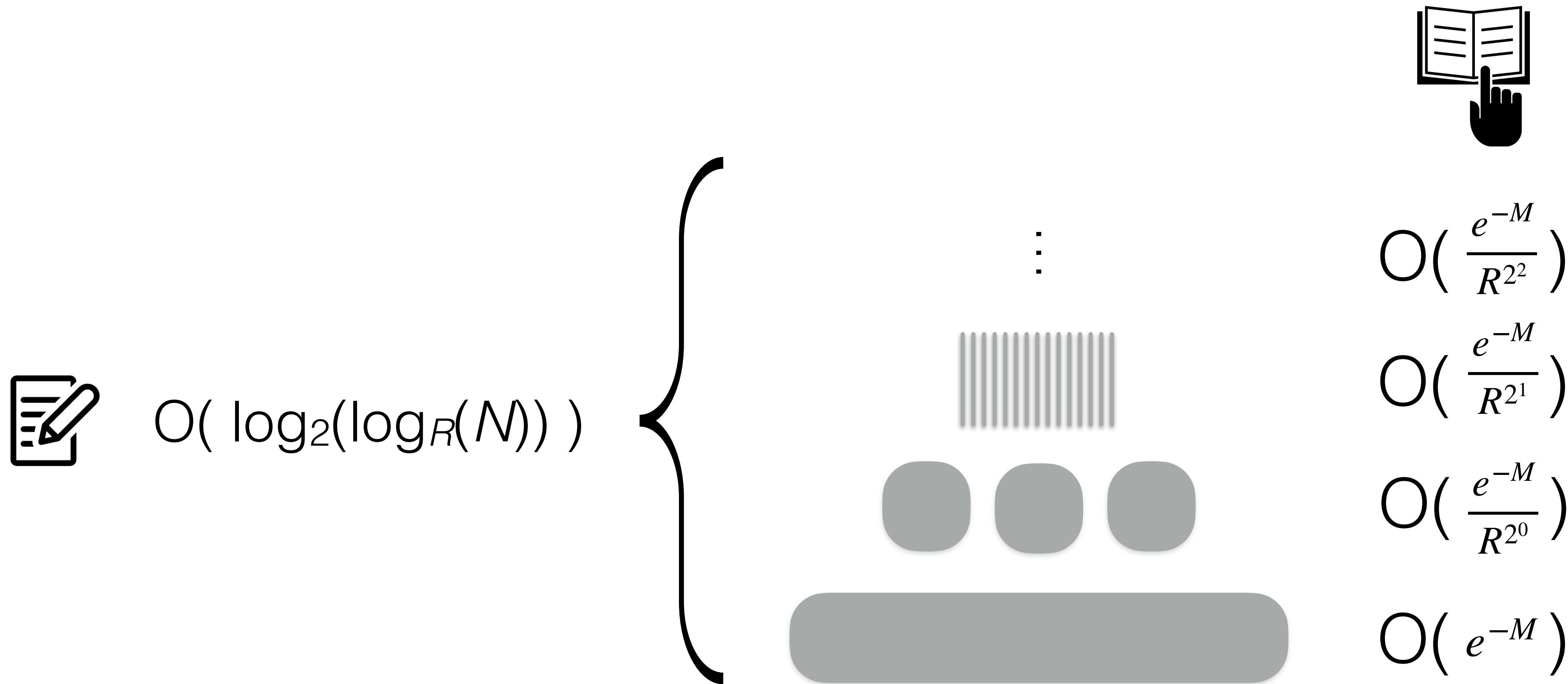
LSM-bush - Lazier Merging for Newer Data



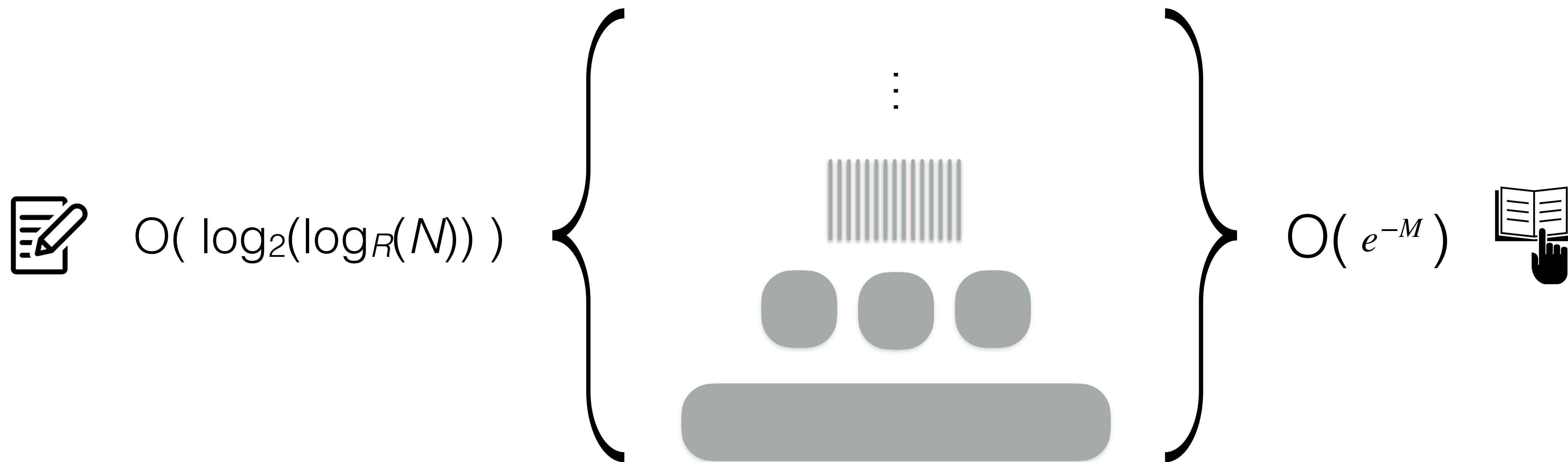
$\mathcal{O}(\log_2(\log_R(N)))$



LSM-bush - Lazier Merging for Newer Data



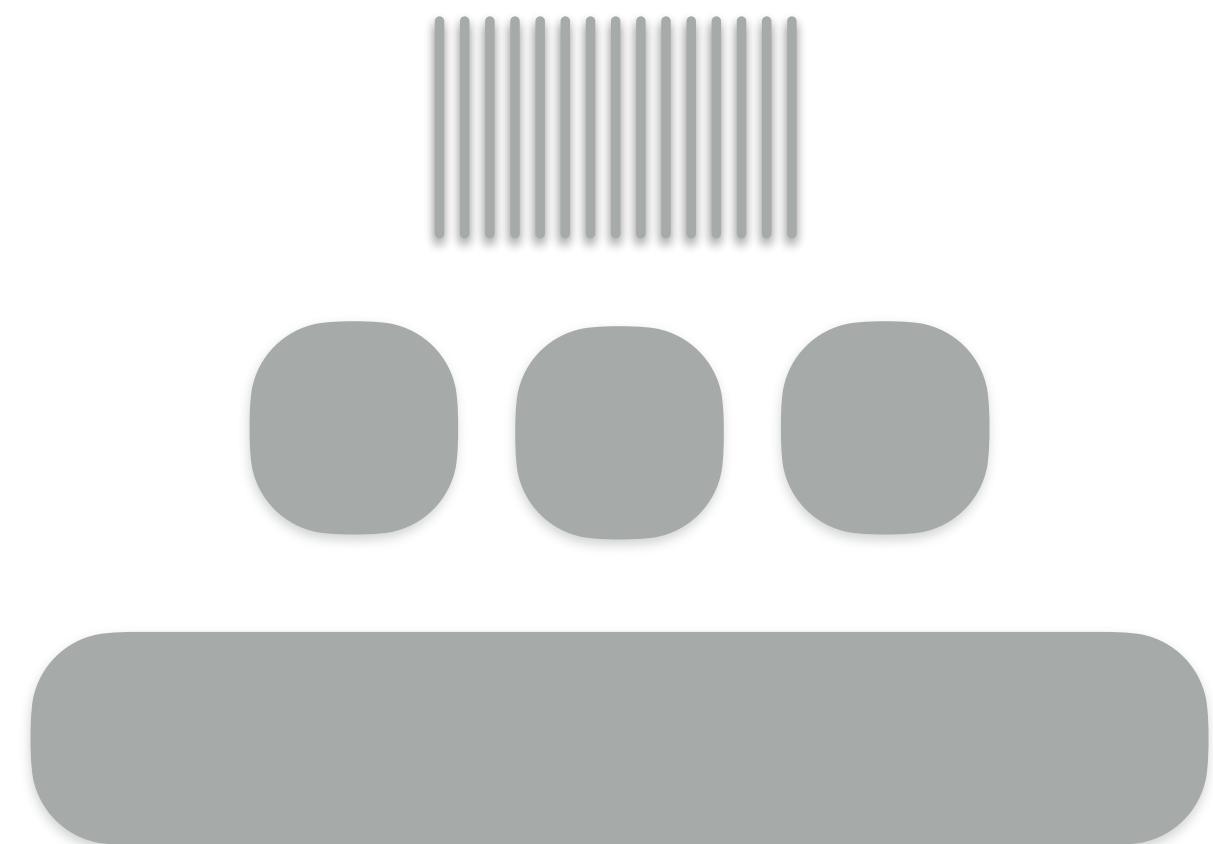
LSM-bush - Lazier Merging for Newer Data



LSM-bush



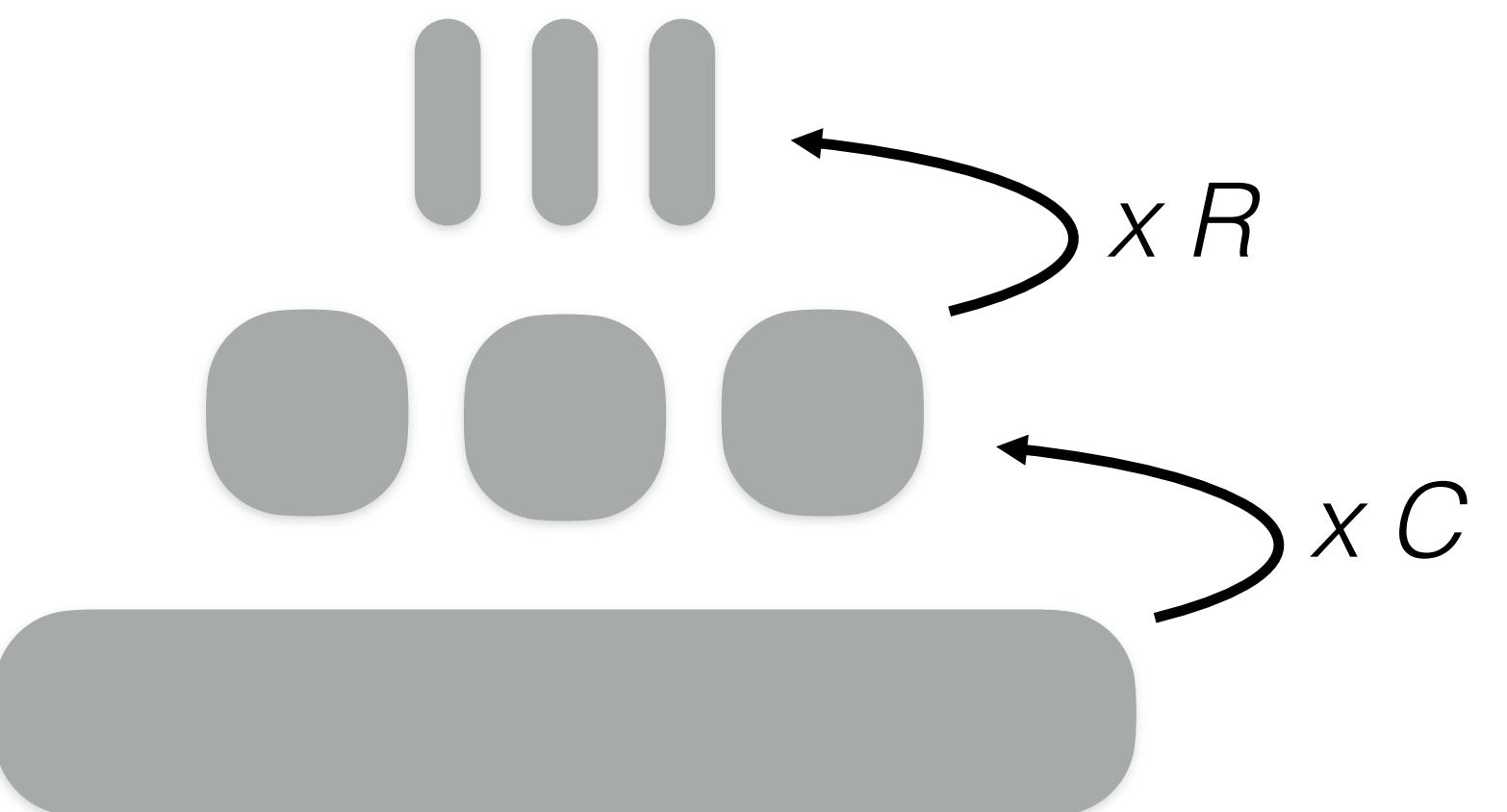
$O(\log(\log(N)))$



SCLL

<

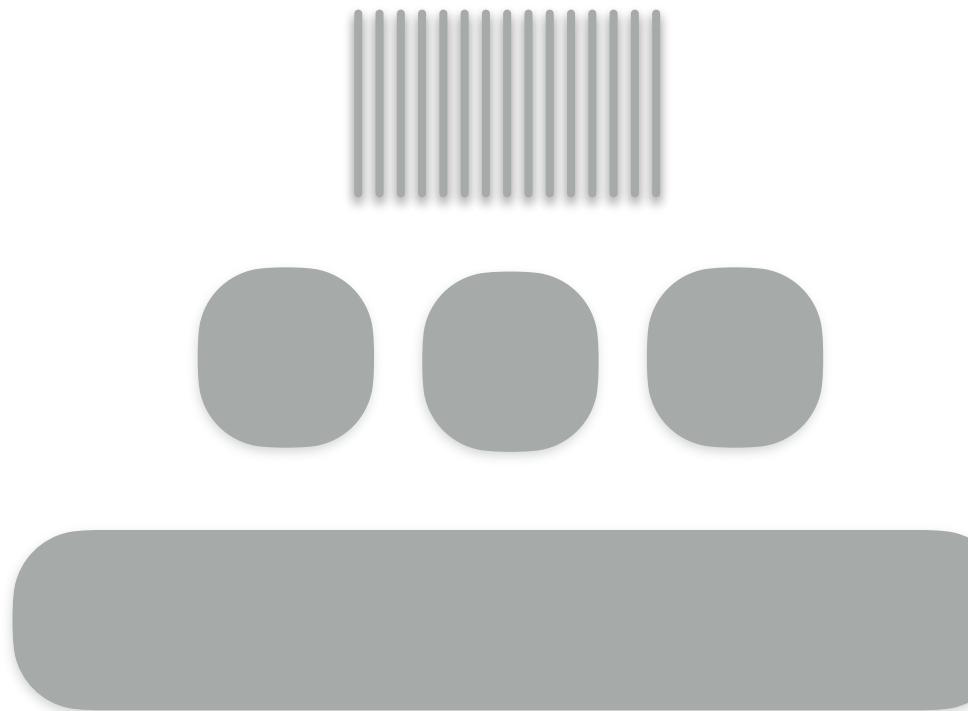
$O(\sqrt{\log(N)})$



LSM-bush

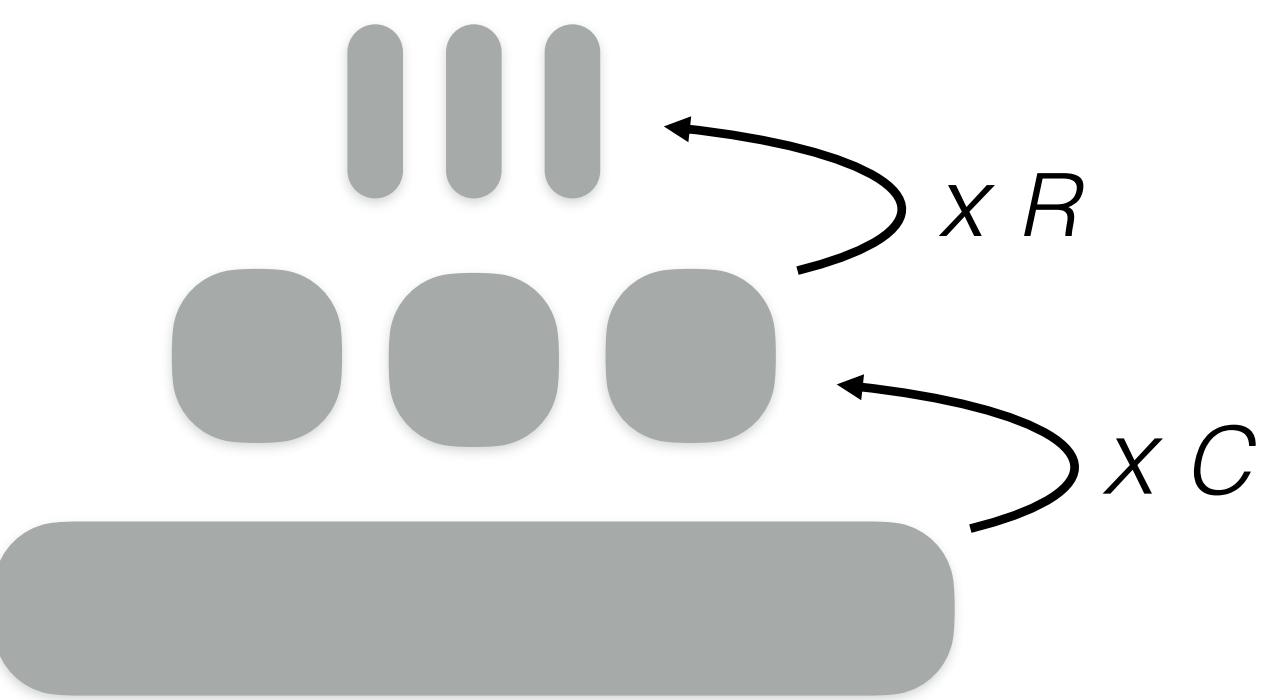


$O(\log(\log(M)))$



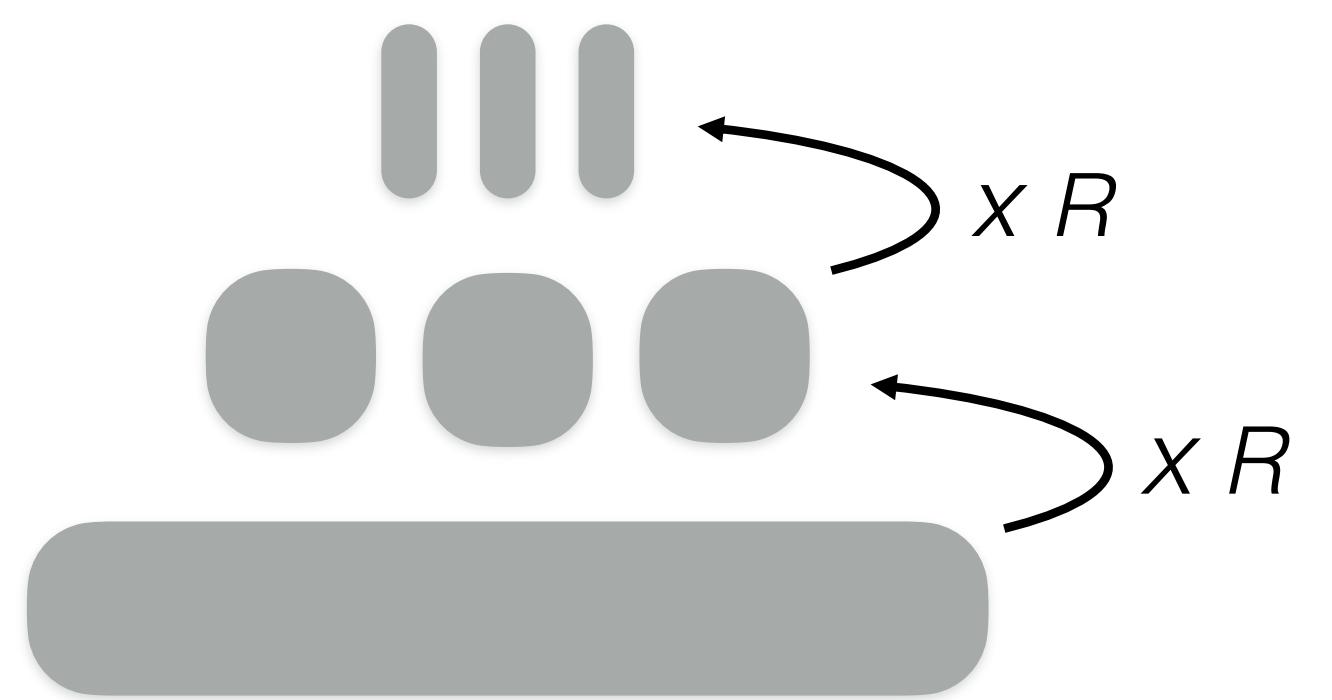
SCLL

$O(\sqrt{\log(M)})$

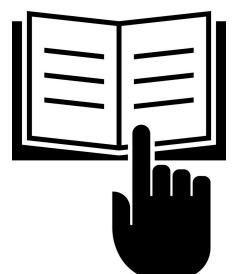


LSM-tree

$O(\log(M))$



LSM-bush



$O(e^{-M})$

=

$O(e^{-M})$

=

$O(e^{-M})$



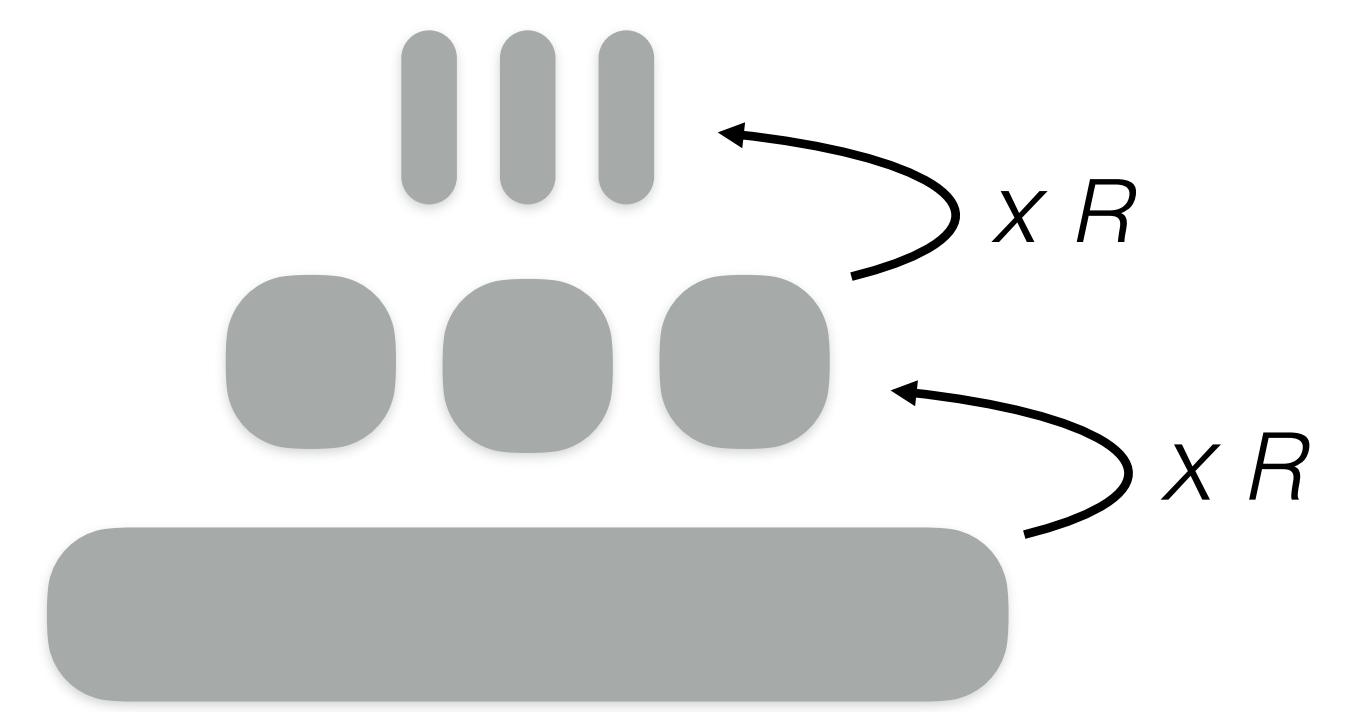
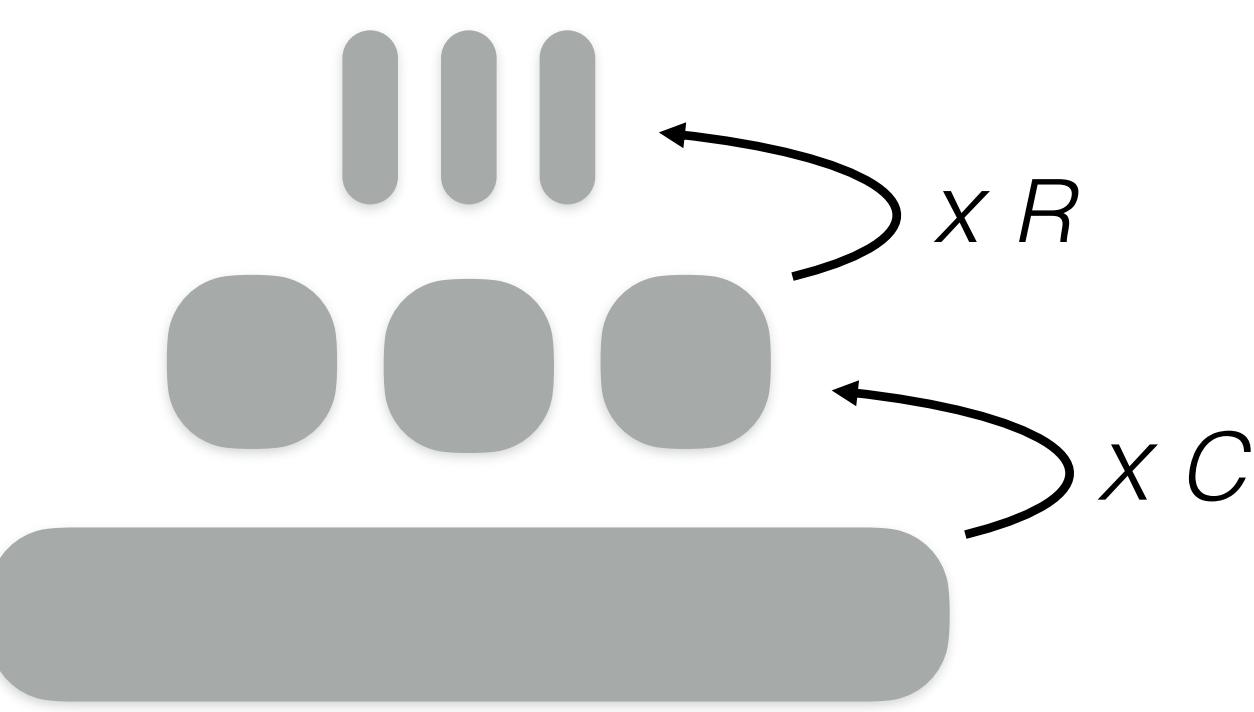
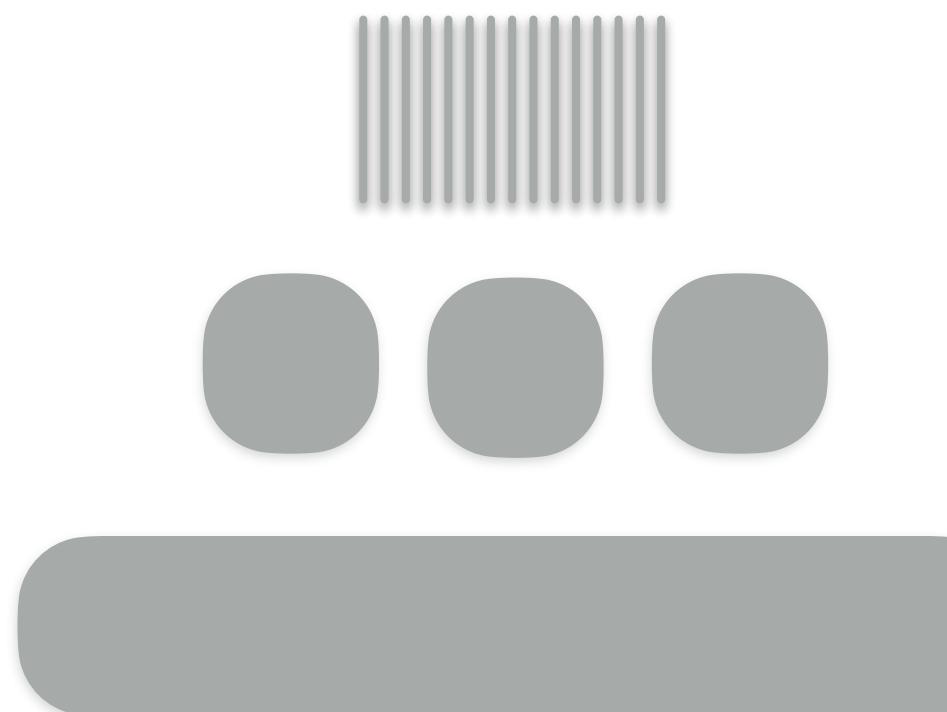
$O(\log(\log(N)))$

<

$O(\sqrt{\log(N)})$

<

$O(\log(N))$

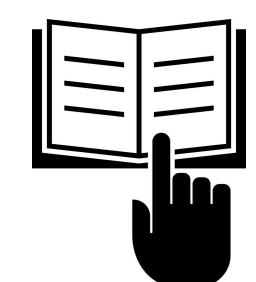


log

LSM-bush

SCLL

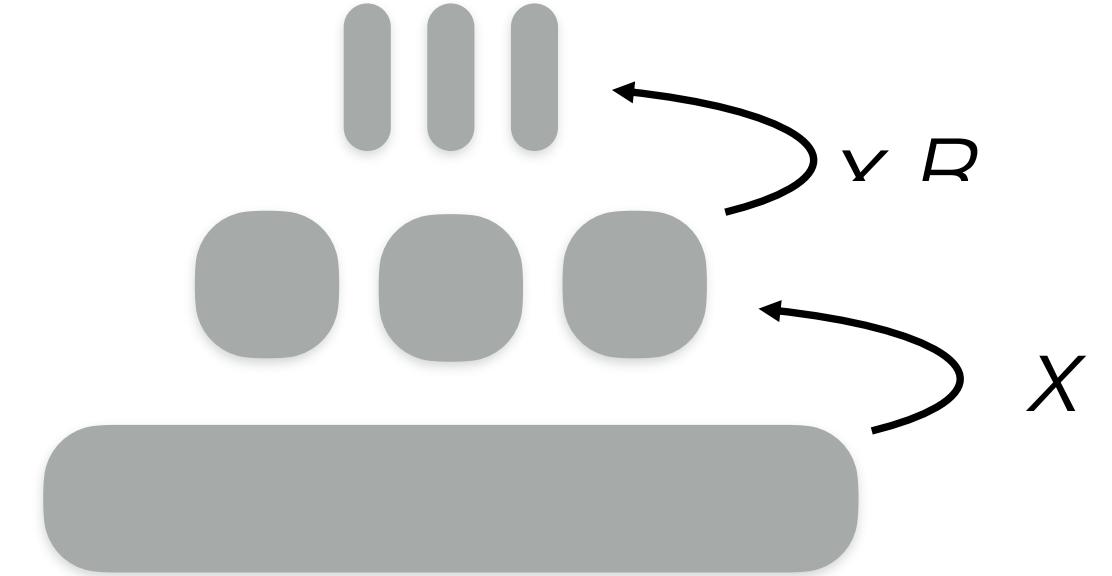
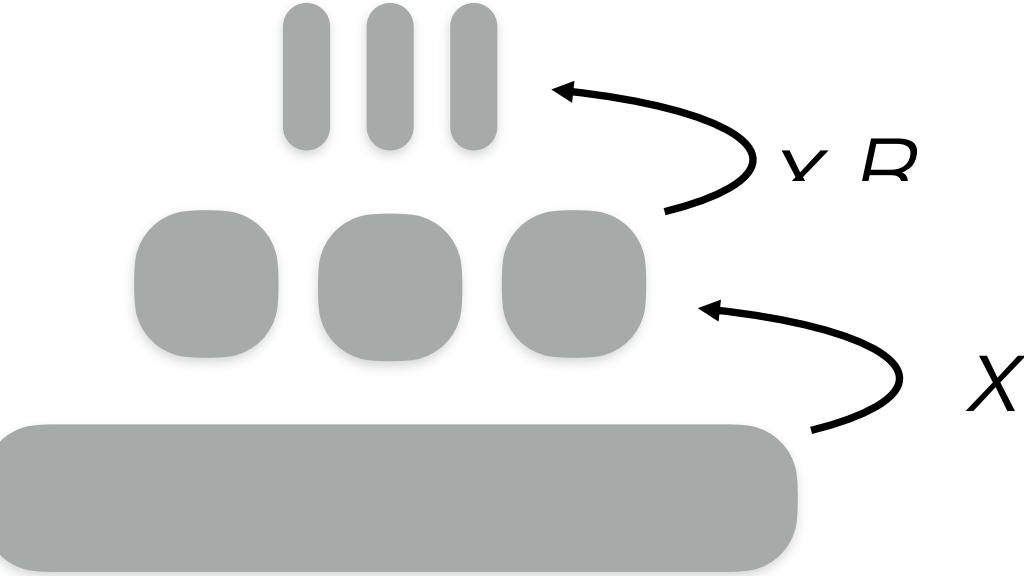
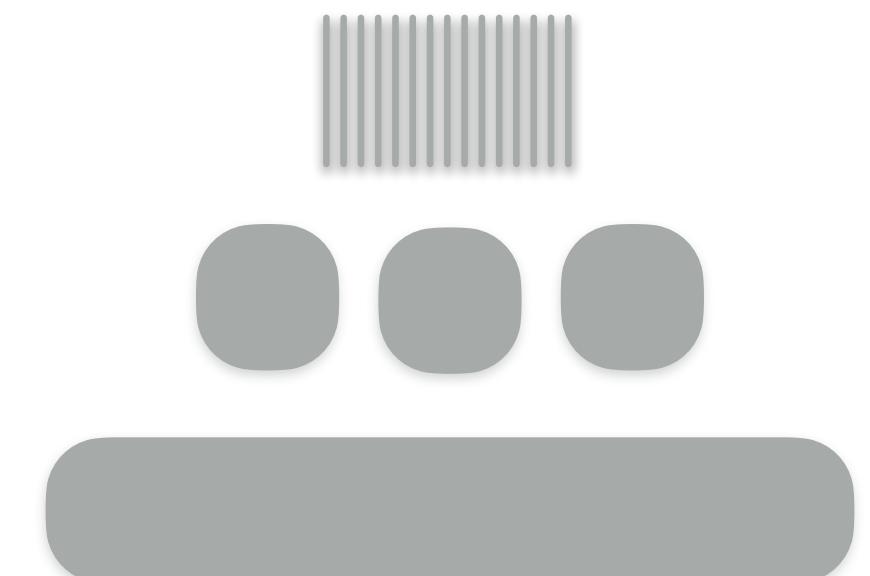
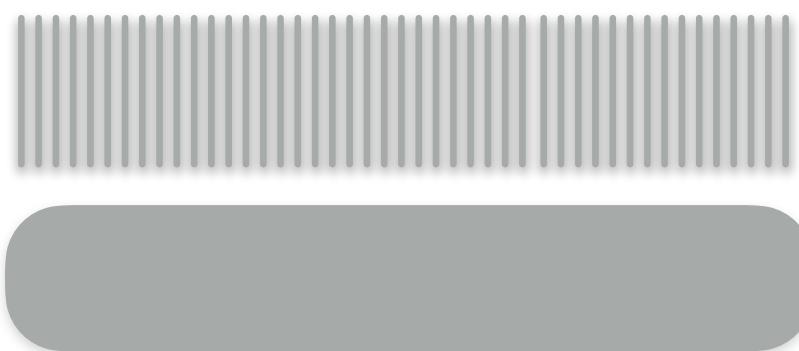
LSM-tree

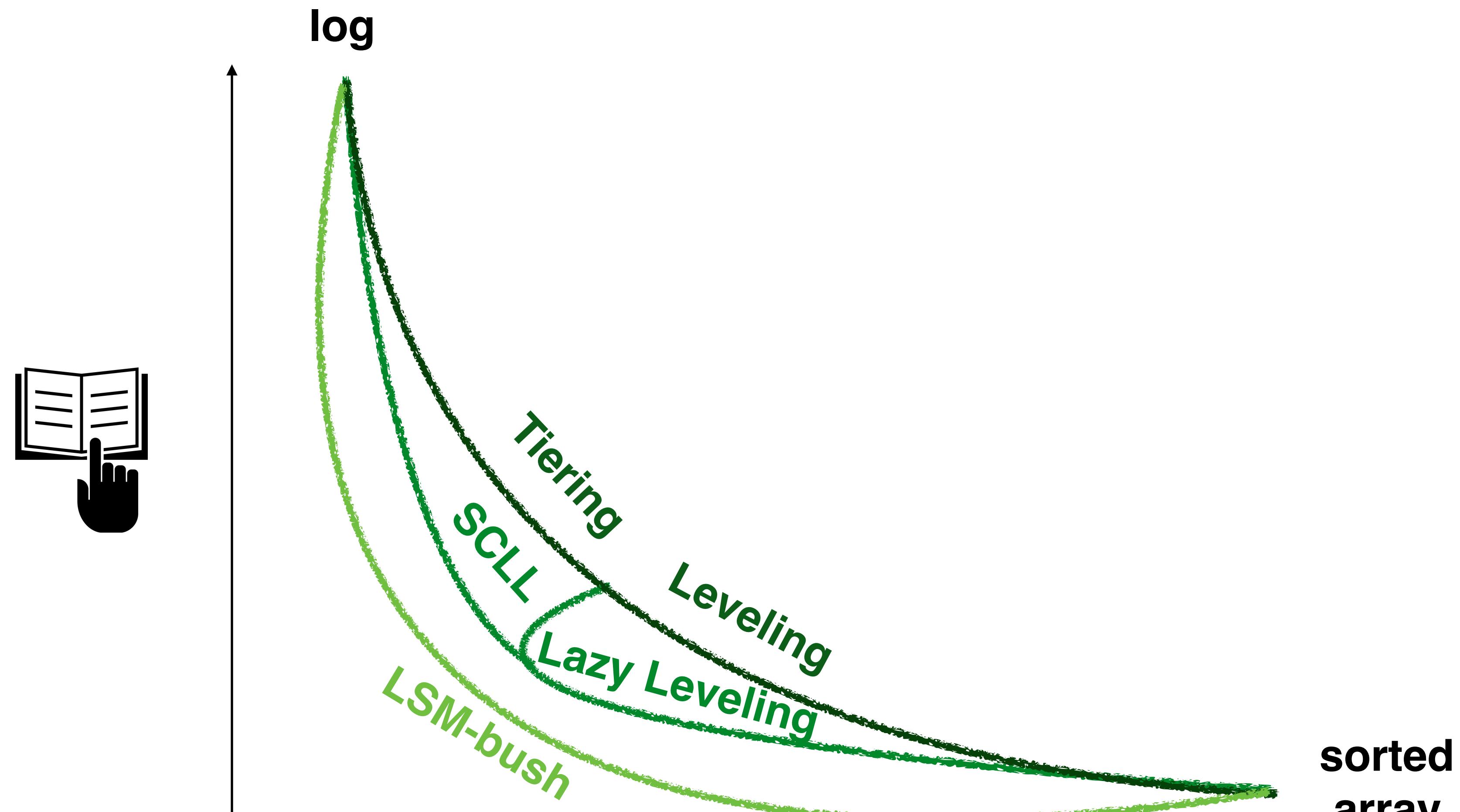


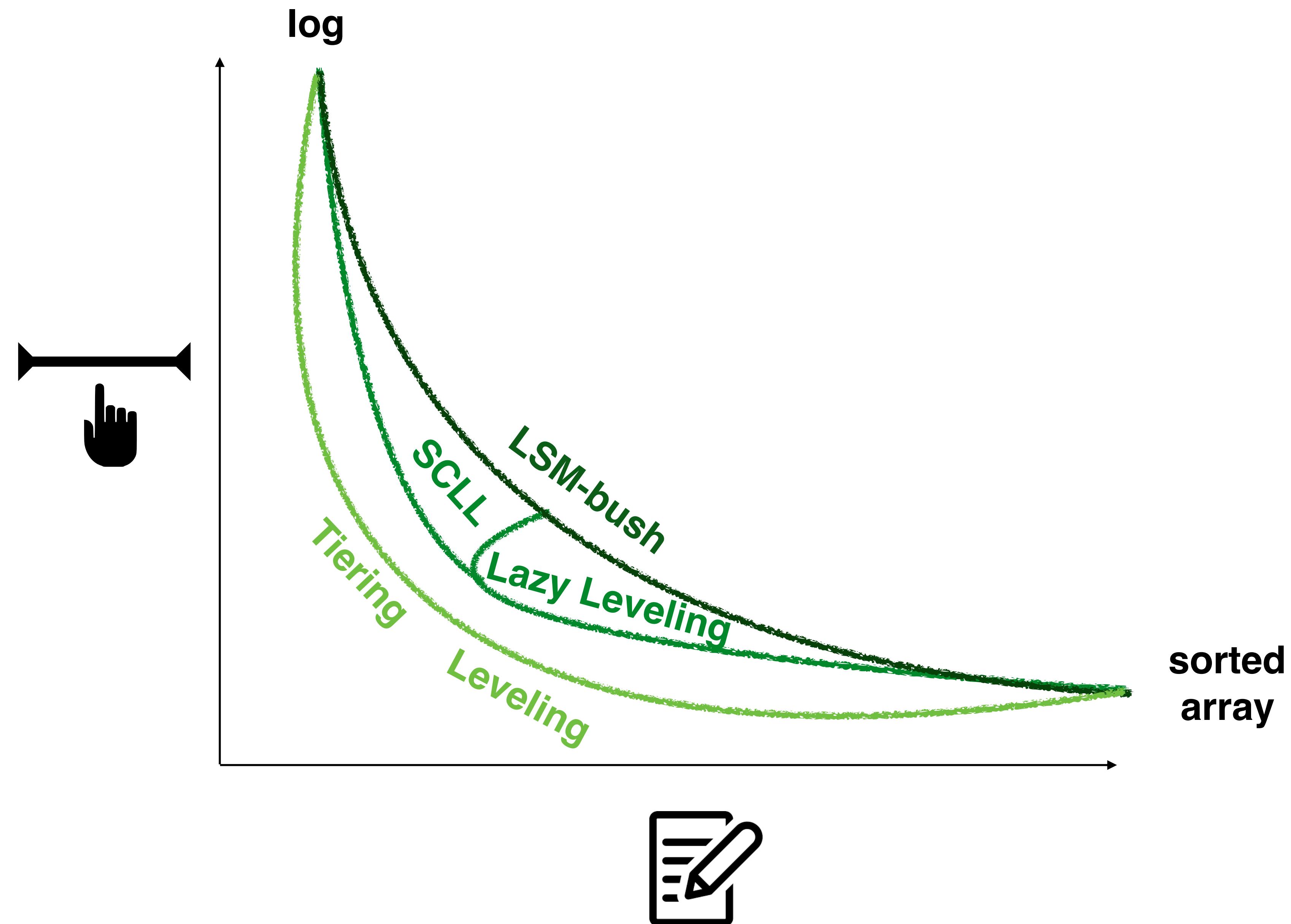
$$\text{O}(N \cdot e^{-M}) > \text{O}(e^{-M}) = \text{O}(e^{-M}) = \text{O}(e^{-M})$$

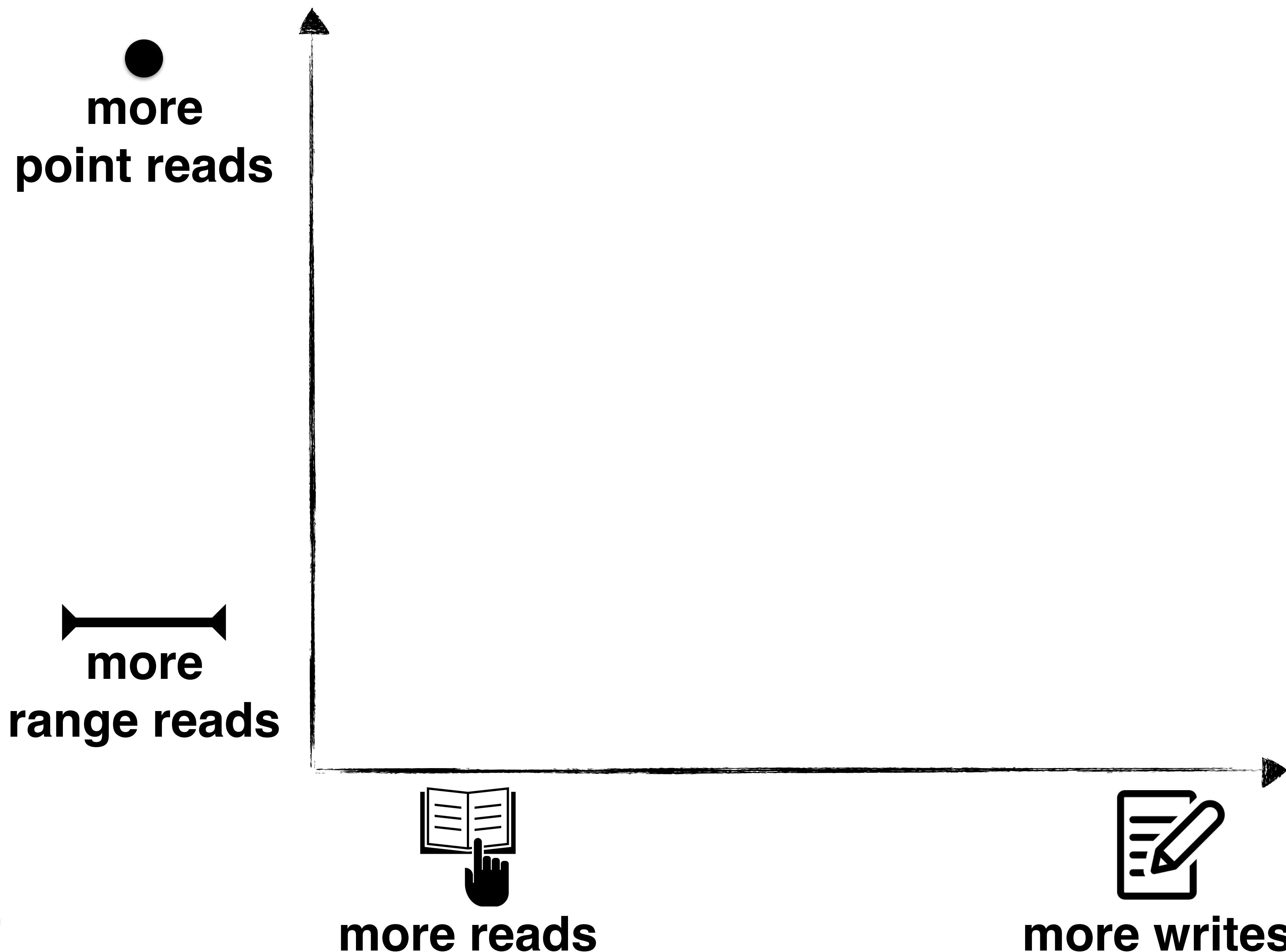


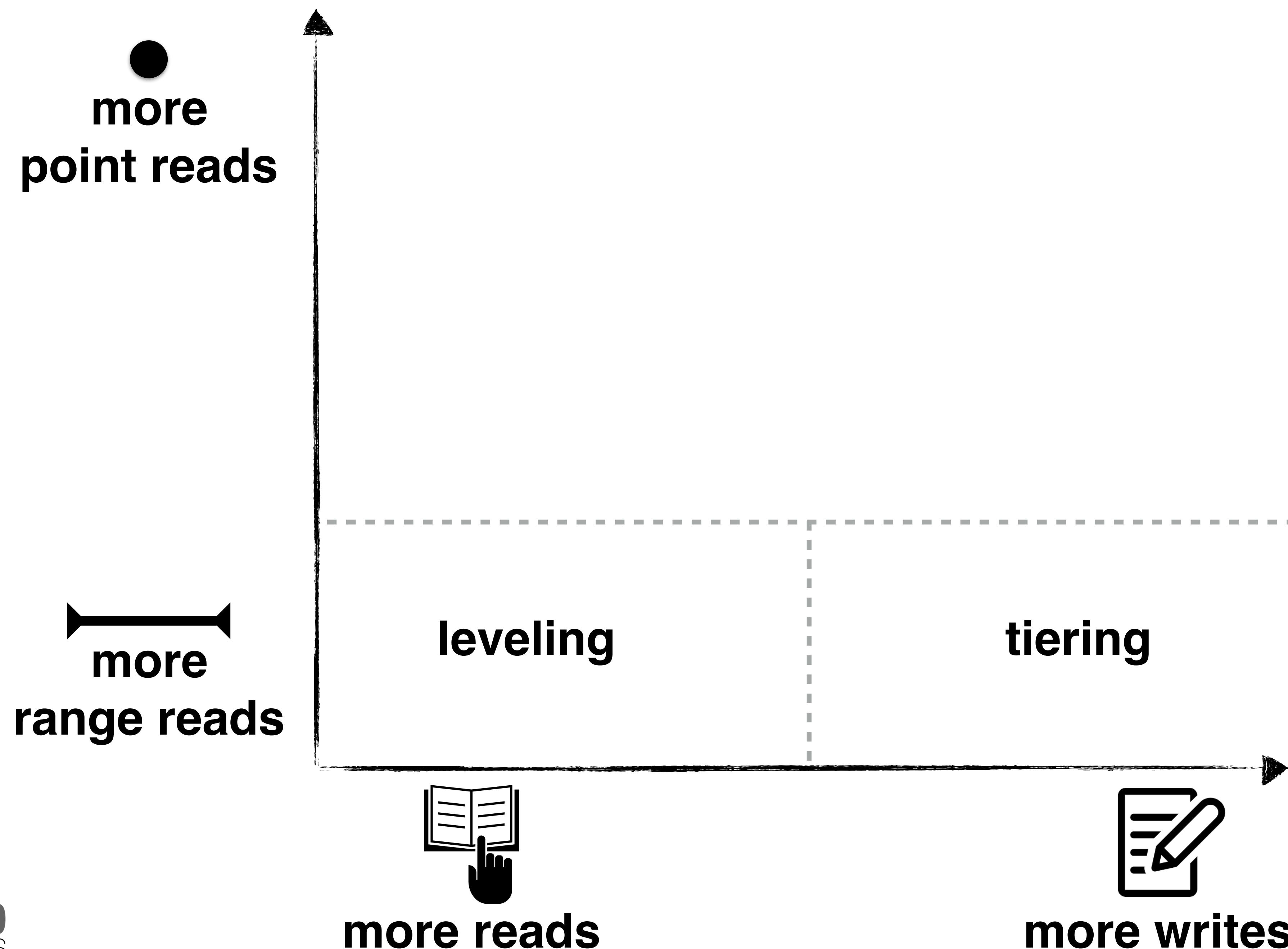
$$\text{O}(1) < \text{O}(\log(\log(N))) < \text{O}(\sqrt{\log(N)}) < \text{O}(\log(N))$$

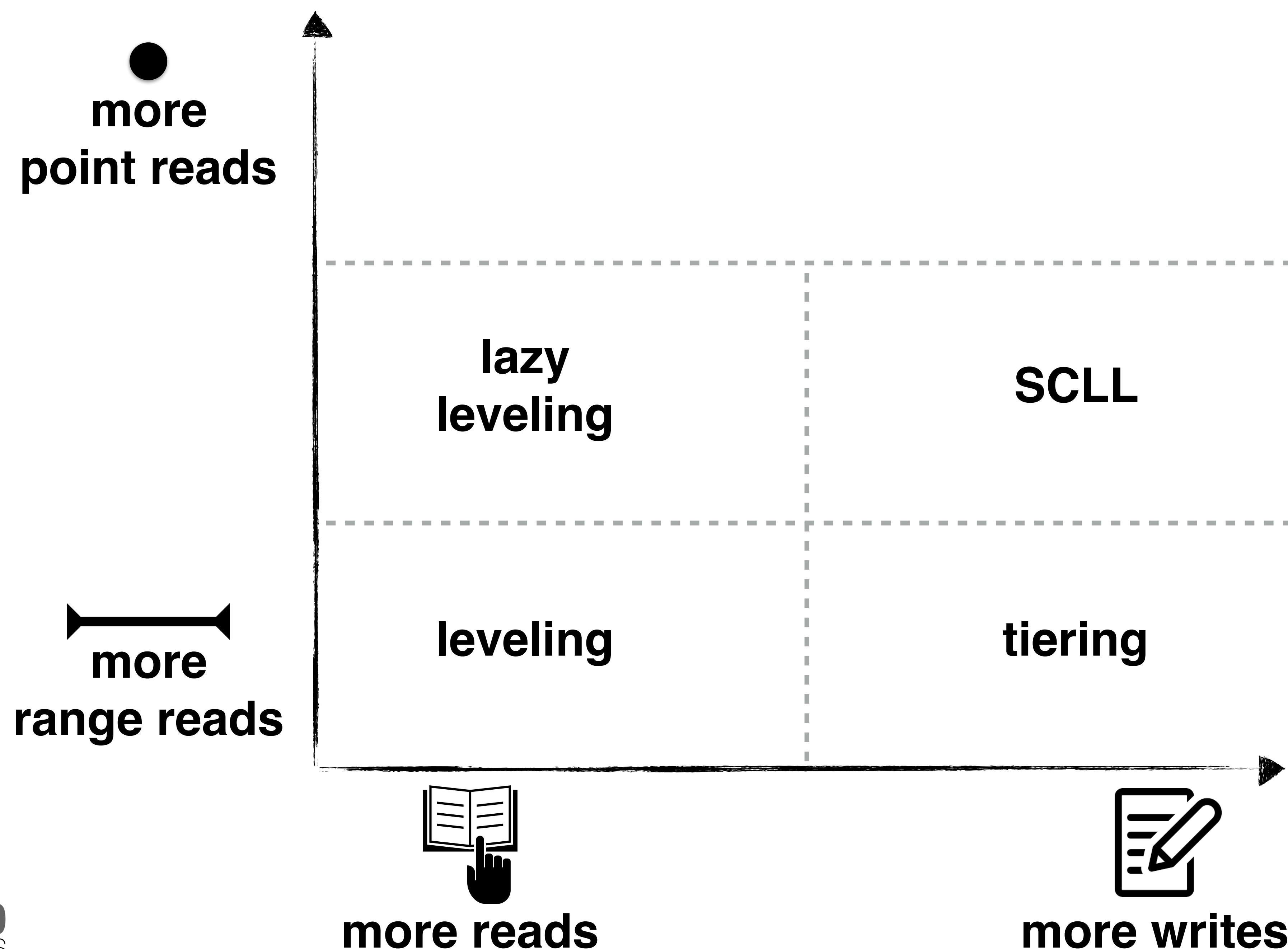


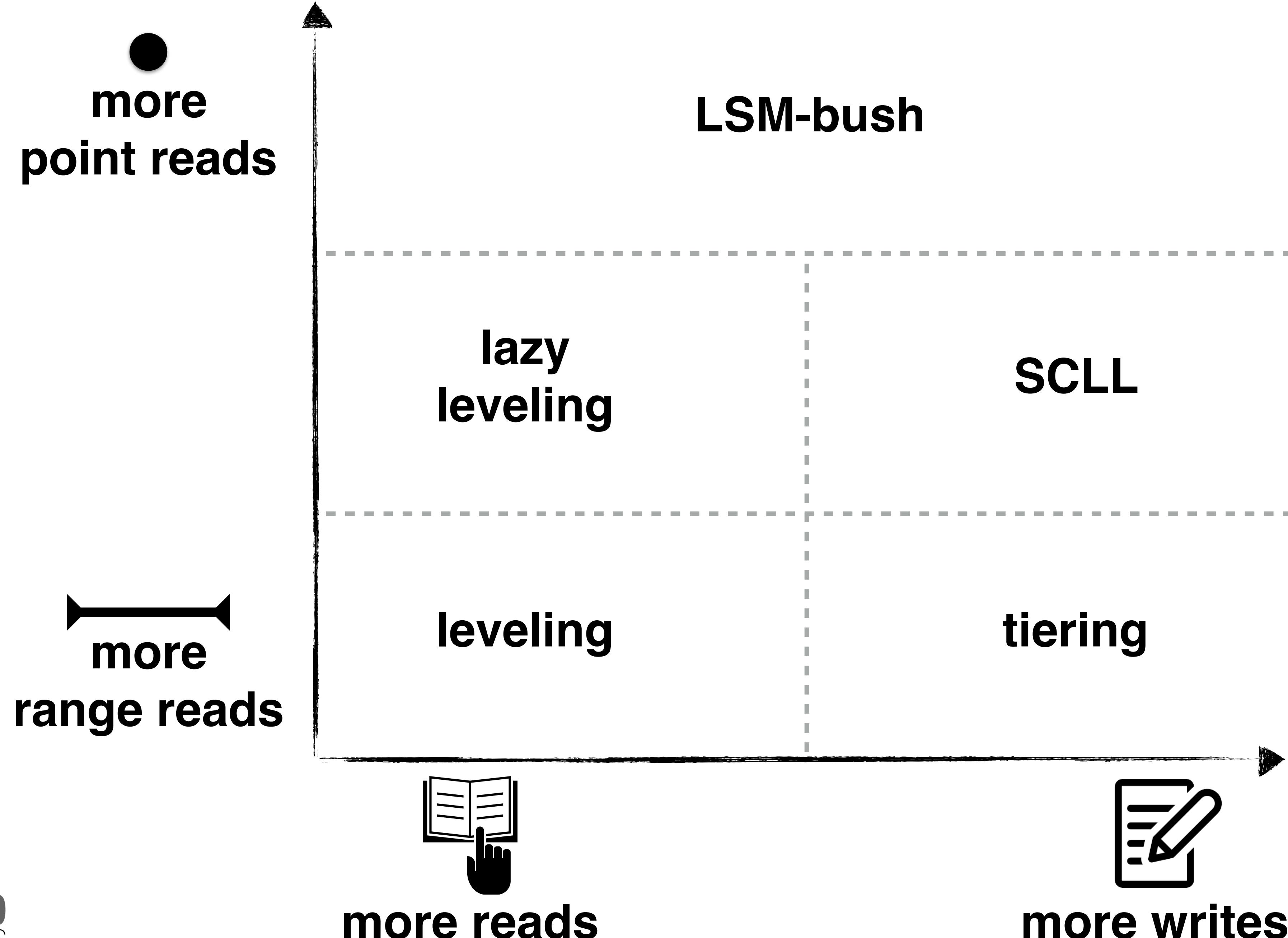












Wacky: Amorphous Calculable Key-Value Store

LSM-bush

**lazy
leveling**

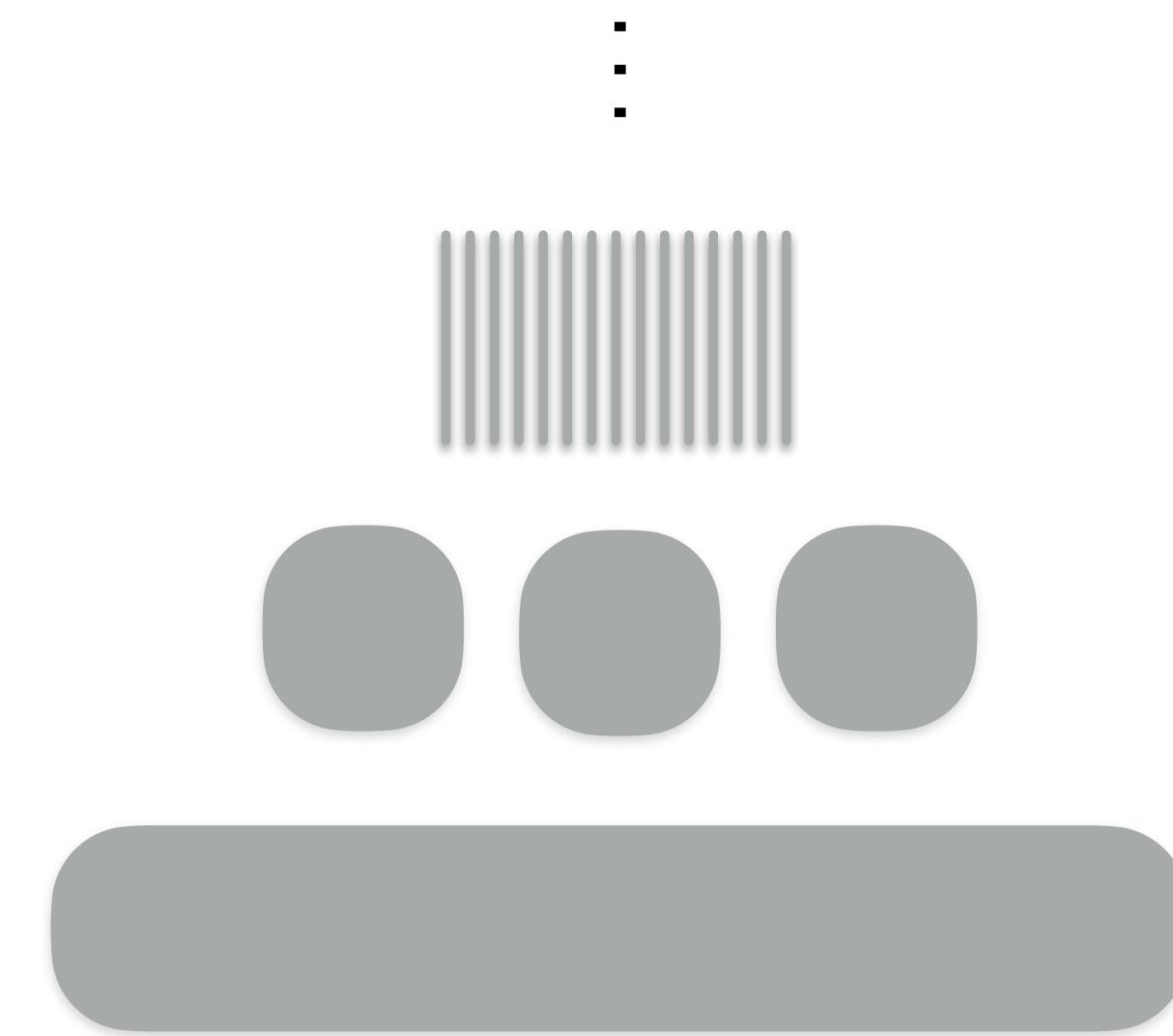


SCLL

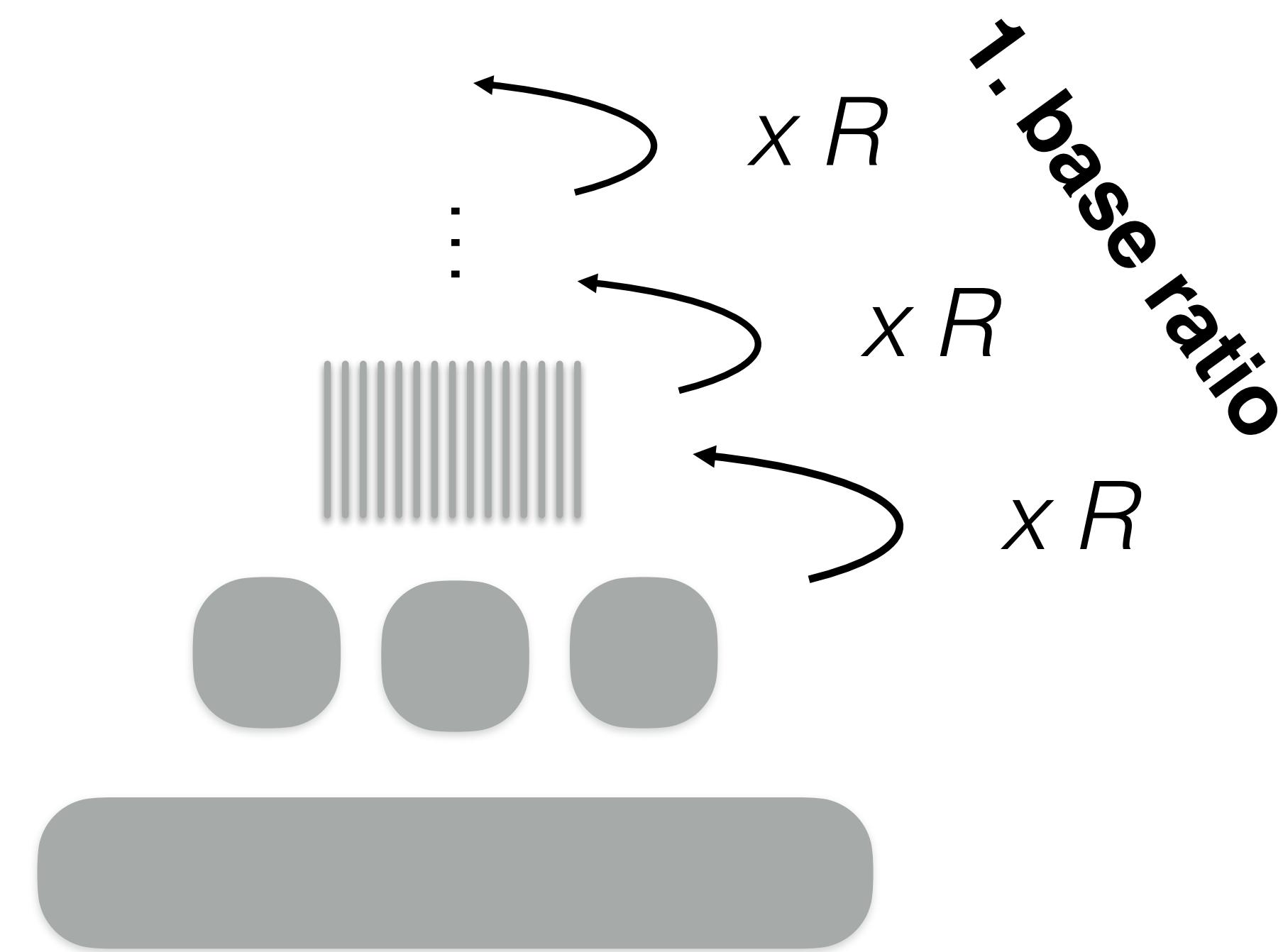
tiering

leveling

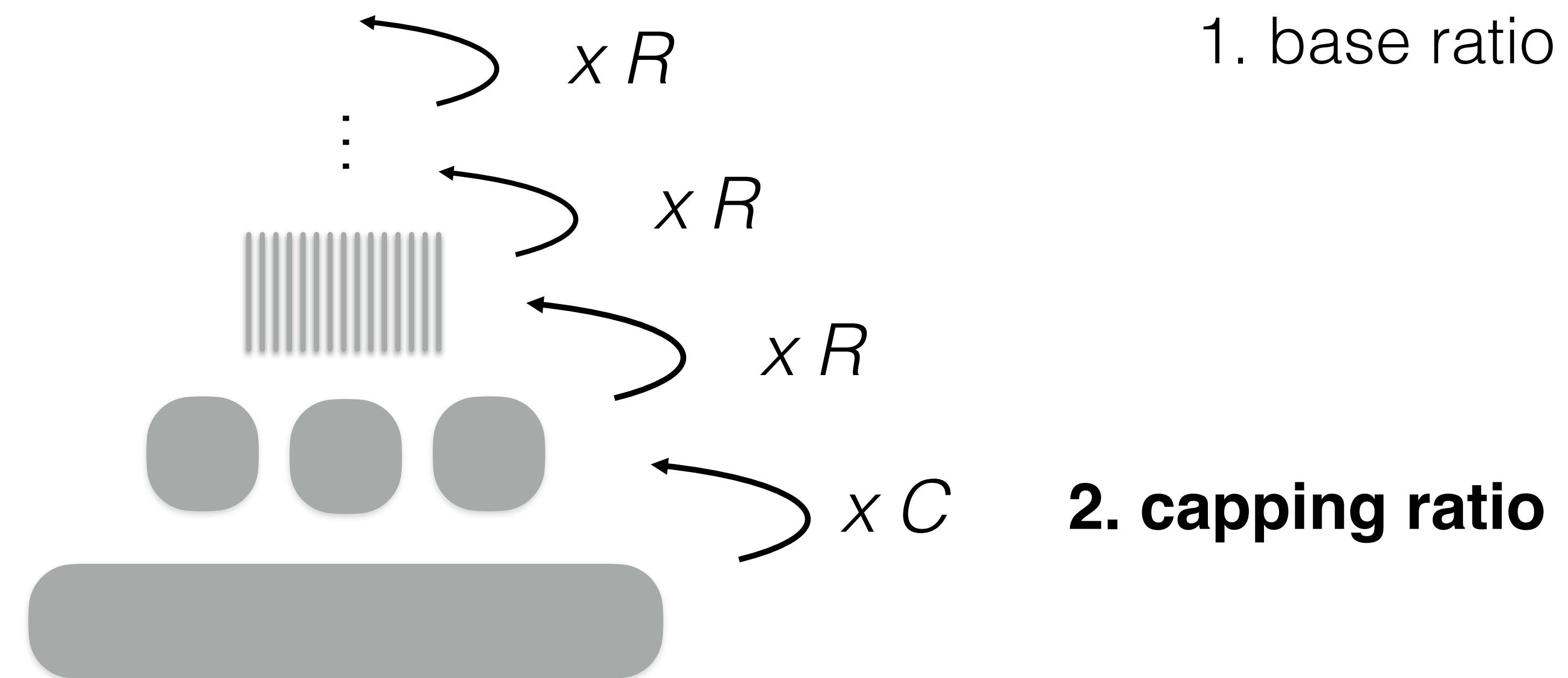
Wacky: Amorphous Calculable Key-Value Store



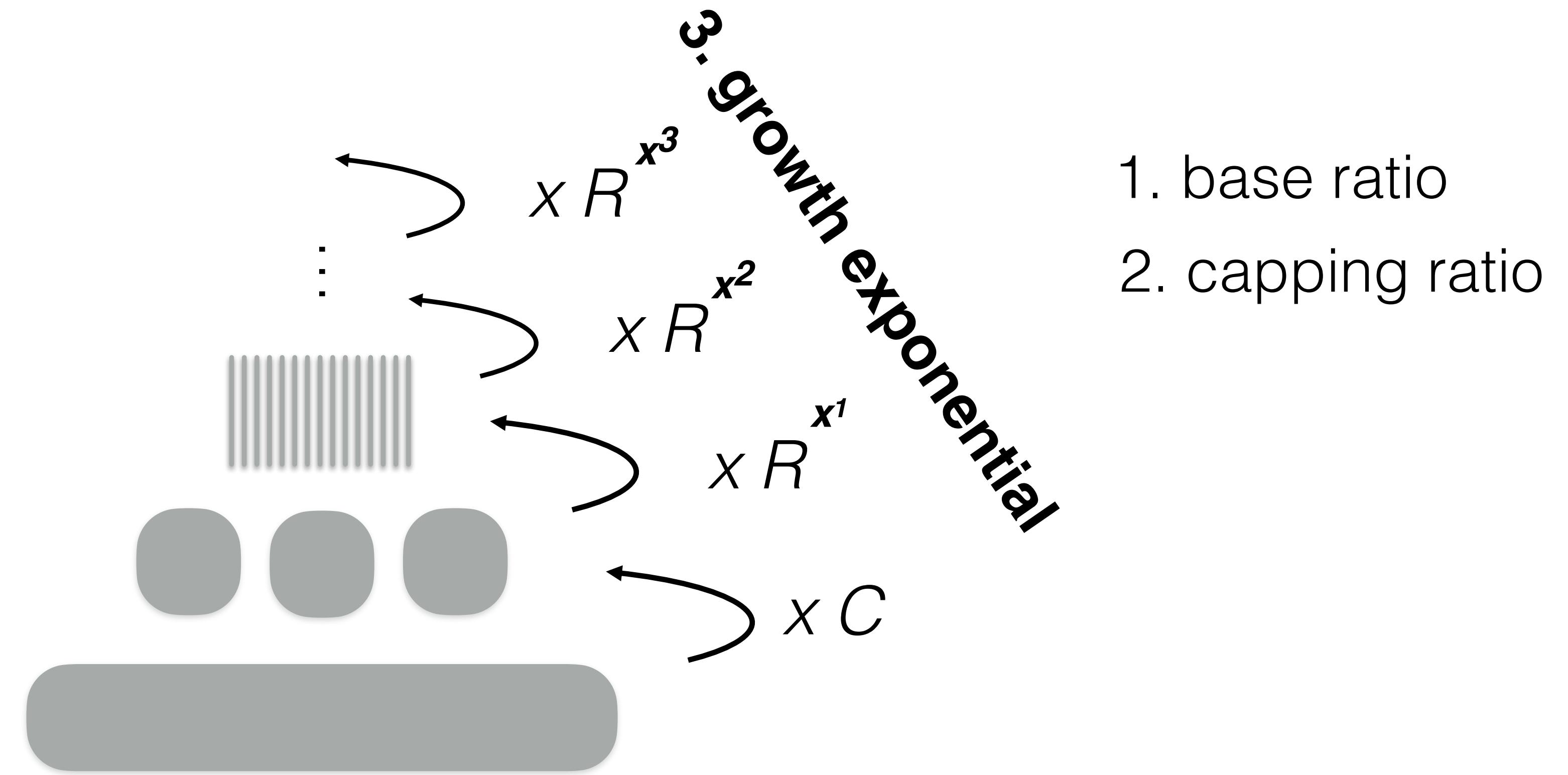
Wacky: Amorphous Calculable Key-Value Store



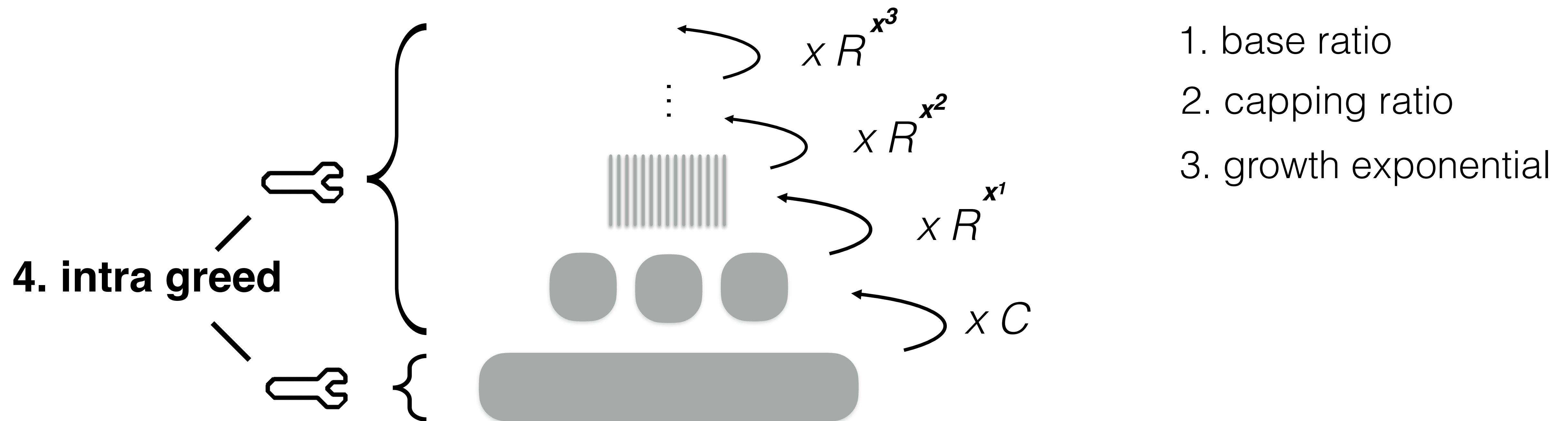
Wacky: Amorphous Calculable Key-Value Store

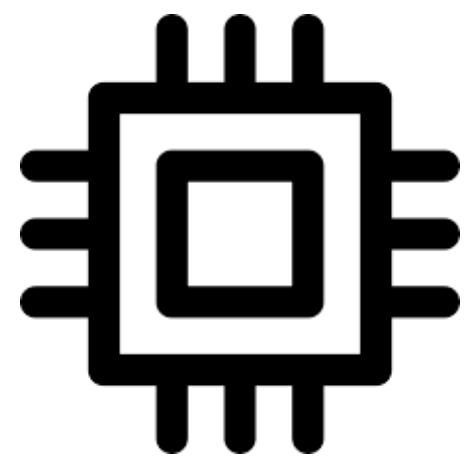
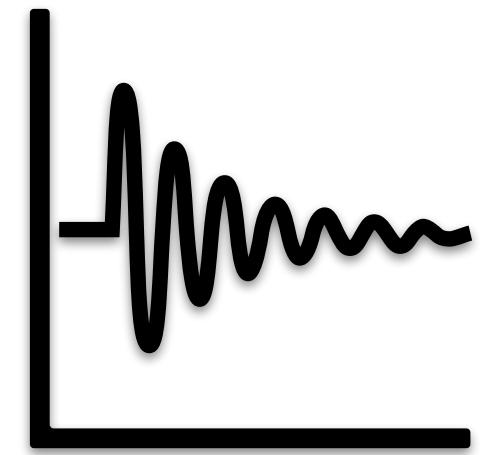


Wacky: Amorphous Calculable Key-Value Store



Wacky: Amorphous Calculable Key-Value Store





1. base ratio
2. capping ratio
3. growth exponential
4. intra greed

normalized
throughput

1

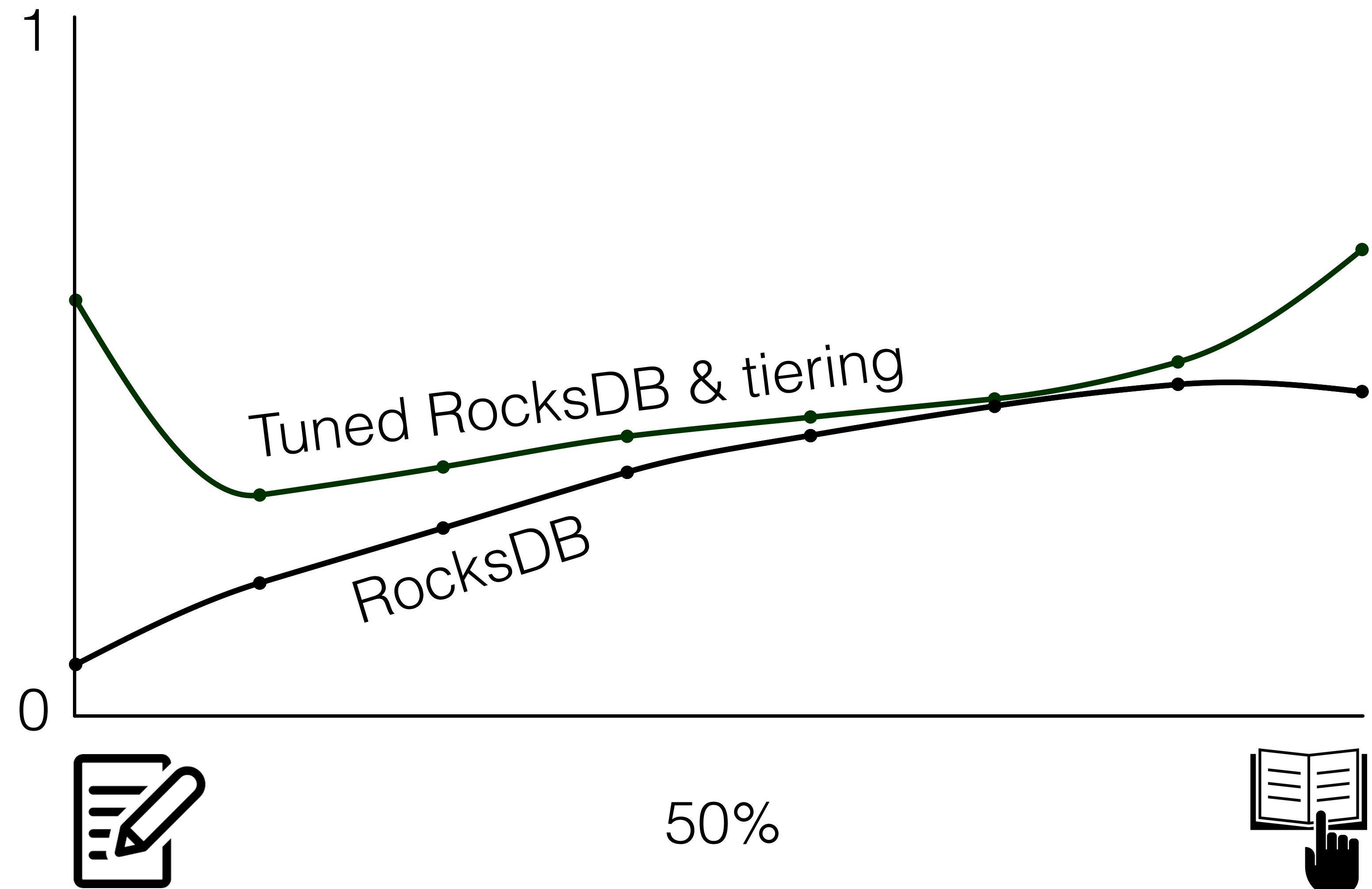
0

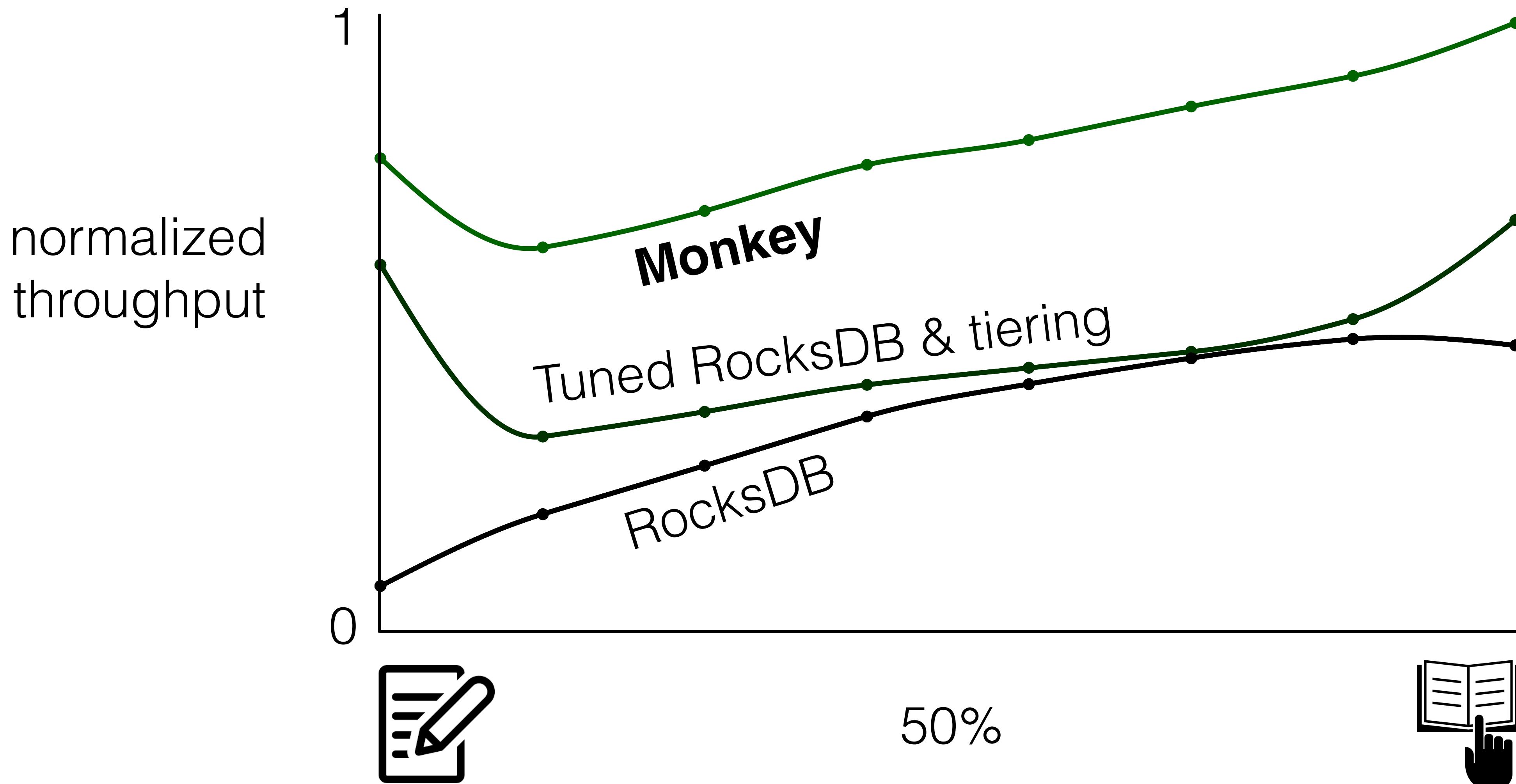


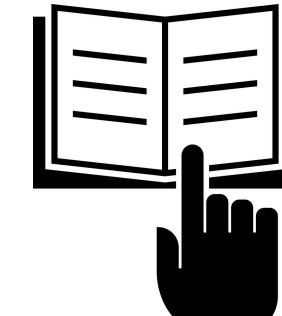
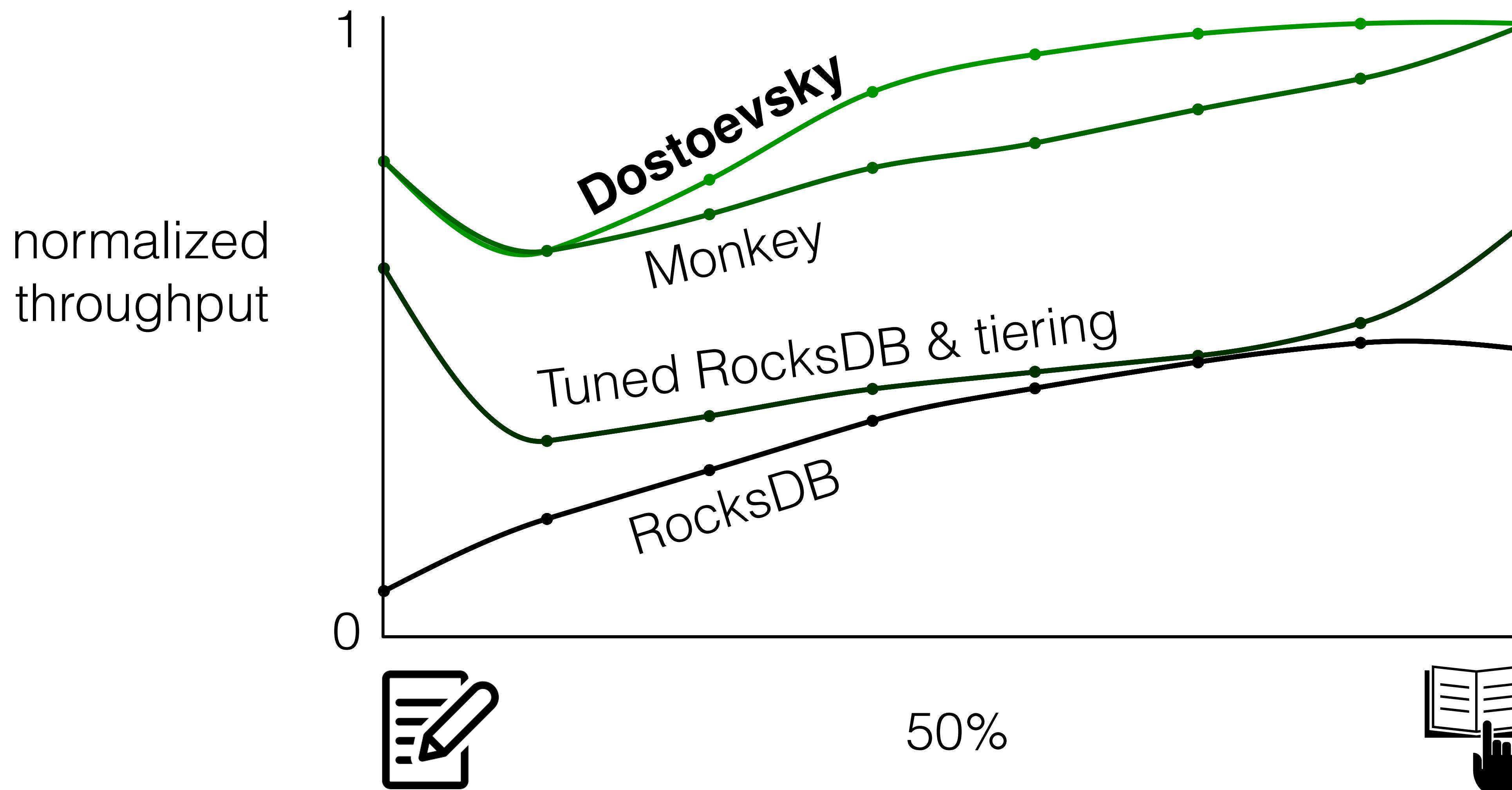
50%

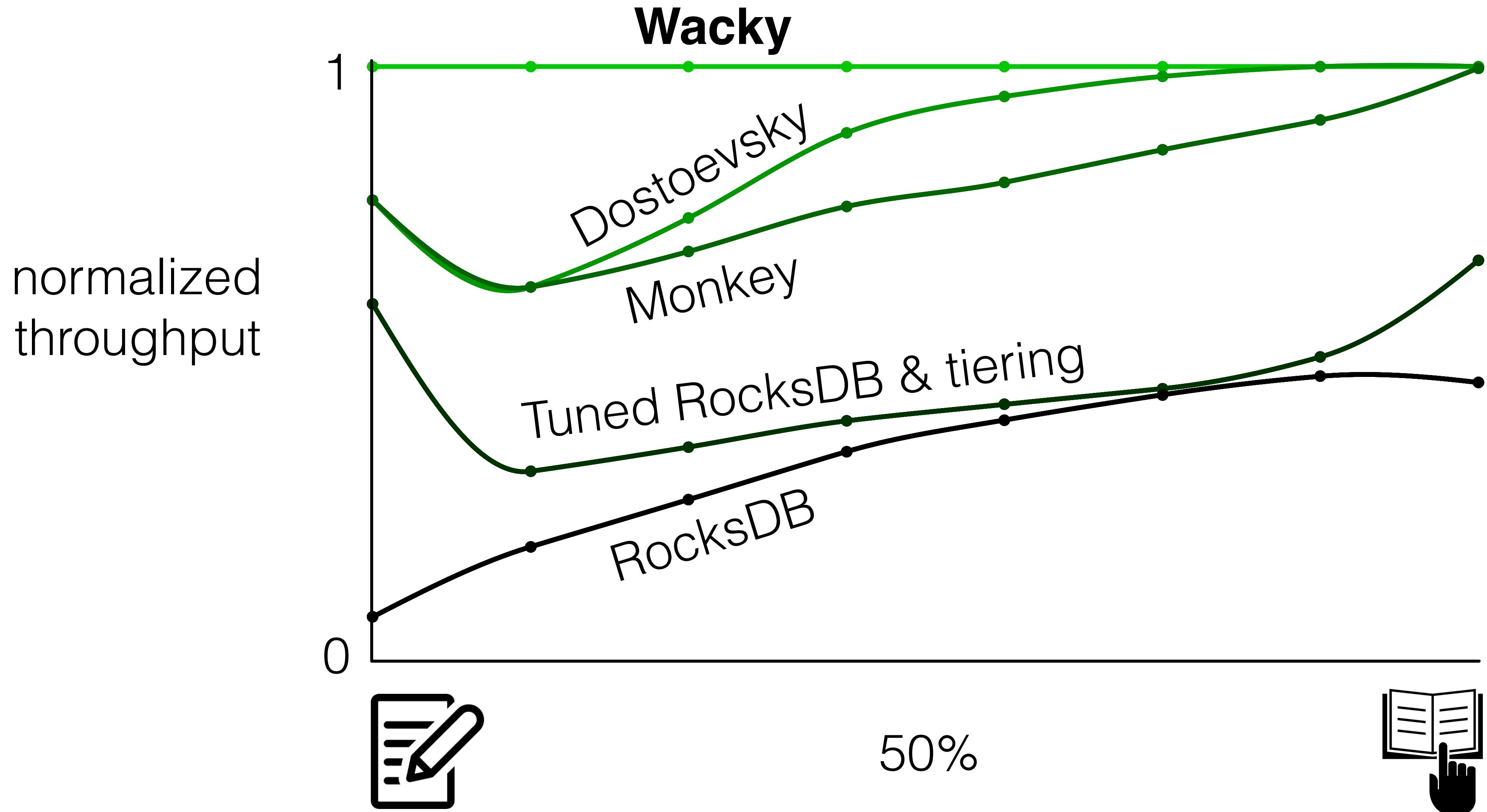


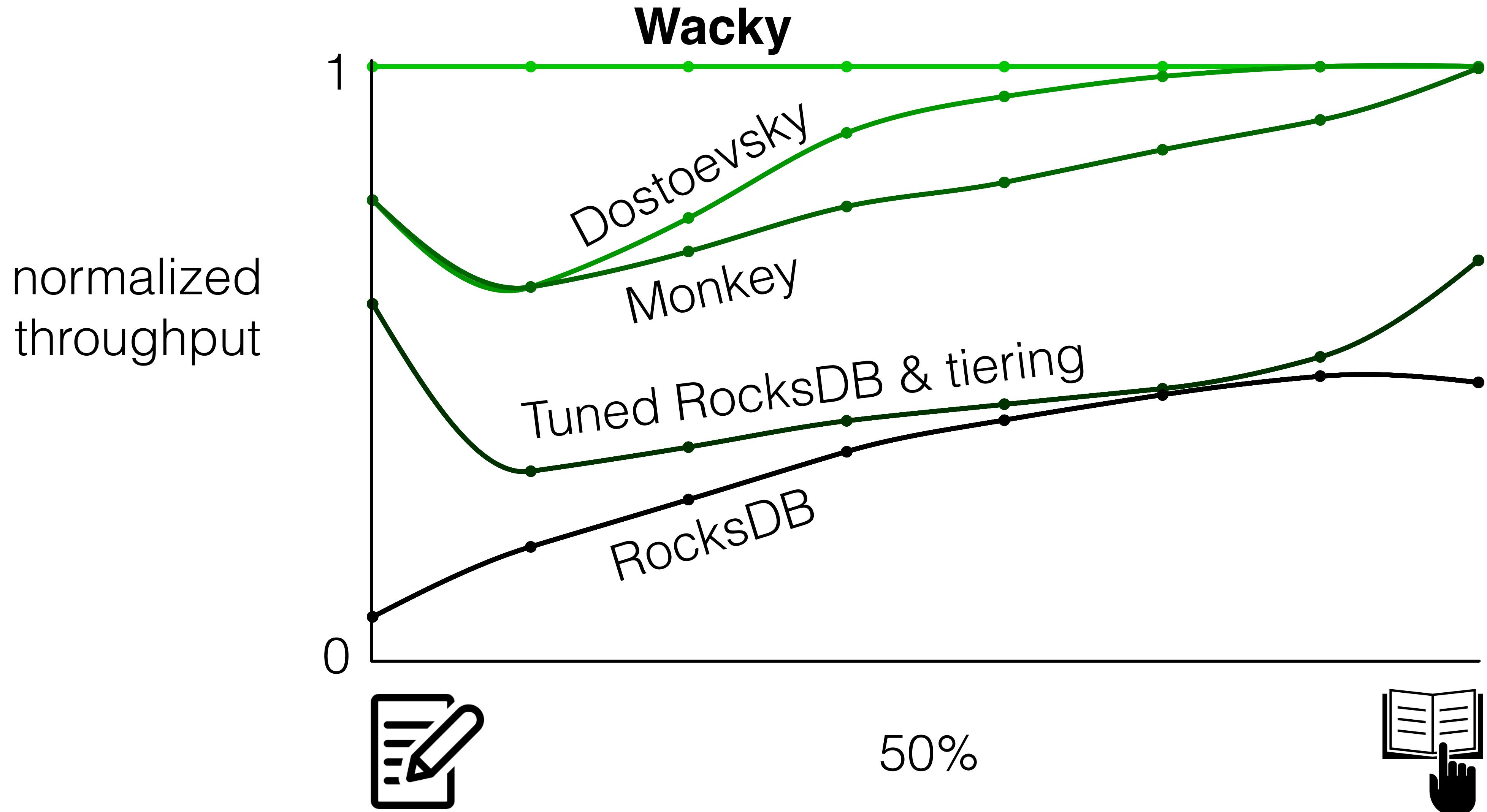
normalized
throughput











thanks!



Monkey



Dostoevsky



Wacky