



Constructing and Analyzing **the LSM Compaction Design Space**

Subhadeep Sarkar

Zichen Zhu

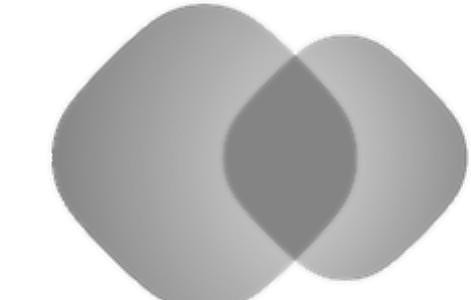
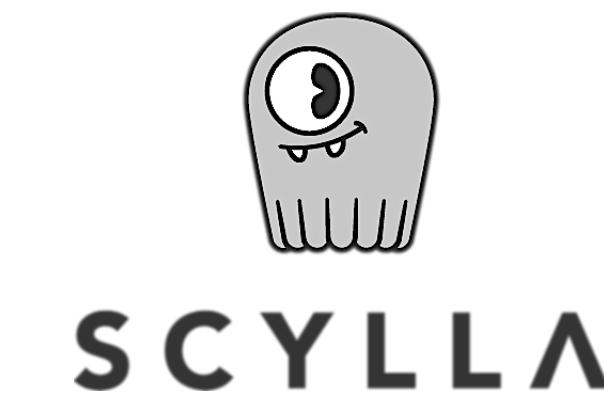
Dimitris Staratzis

Manos Athanassoulis

Log-Structured Merge-tree

LSM-tree

LSM-tree



tarantool



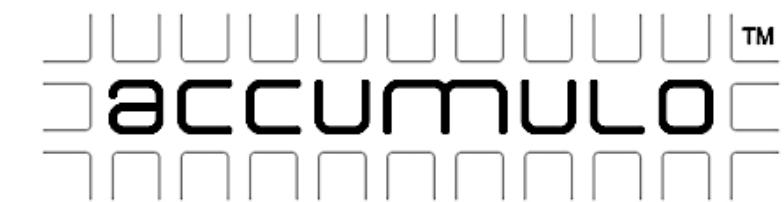
Bigtable



SQLite

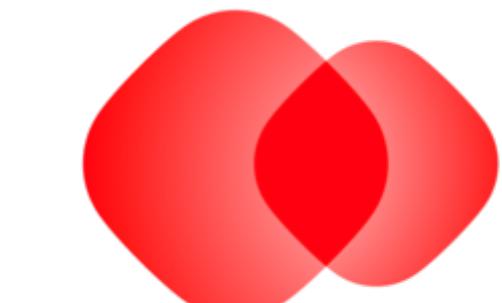
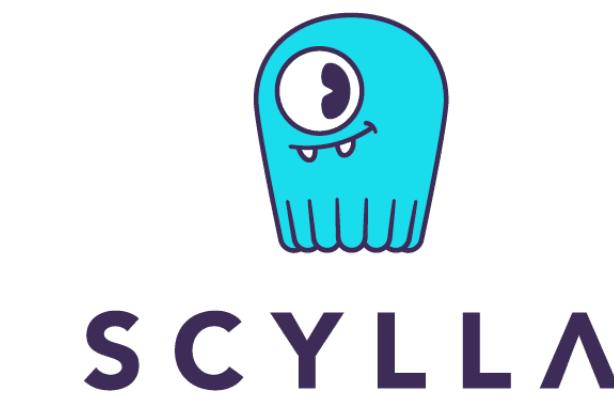


QuasarDB



2021

LSM-tree



tarantool



Bigtable

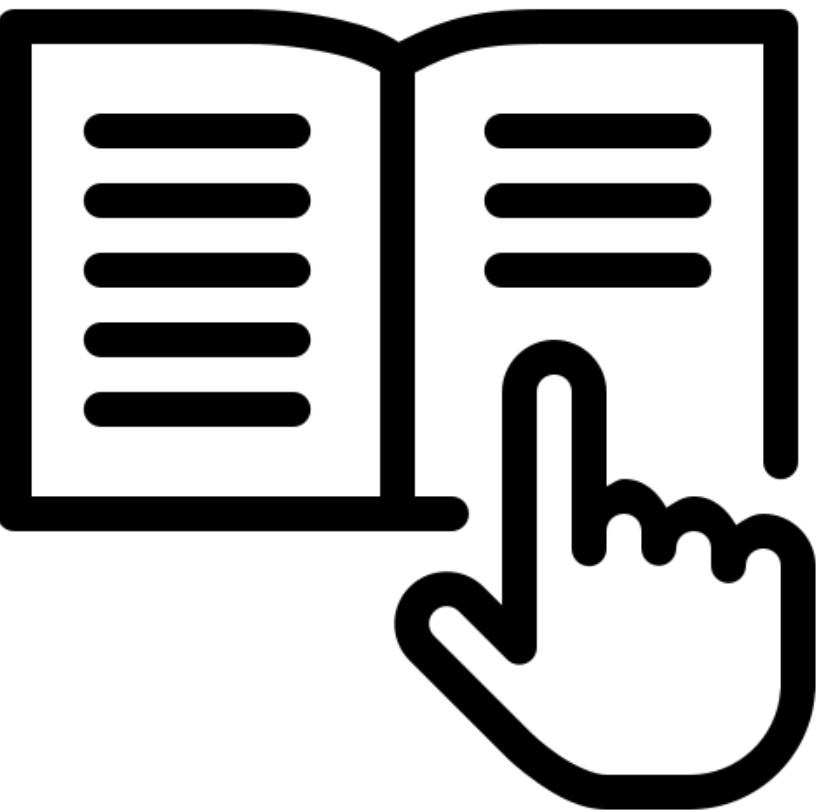


Why **LSM** ?

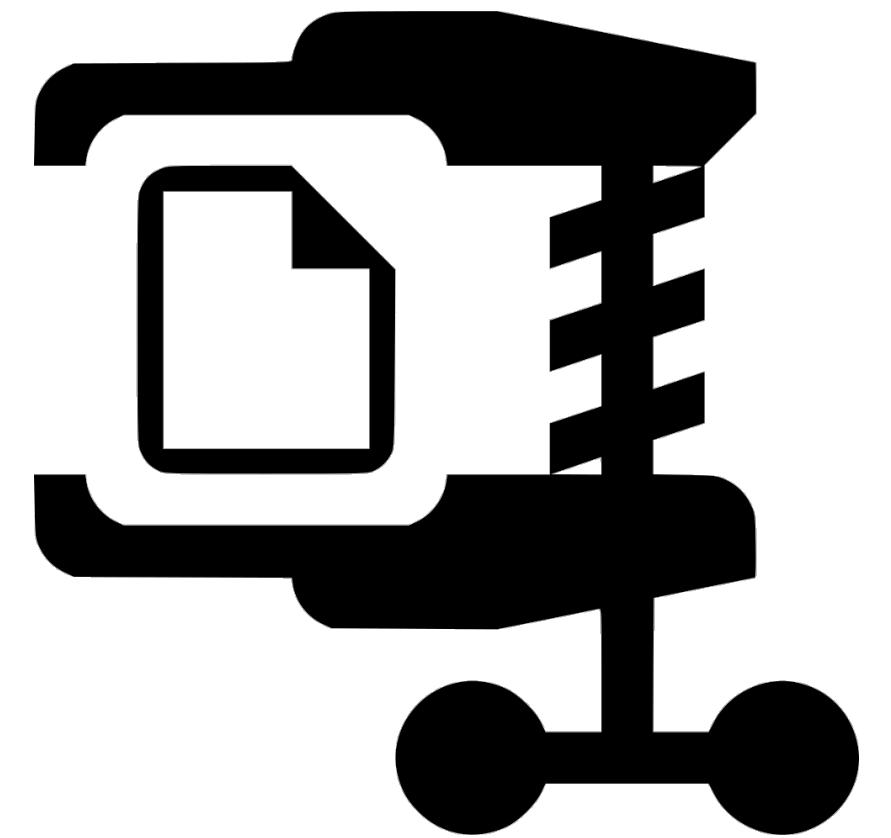
Why **LSM** ?



fast writes



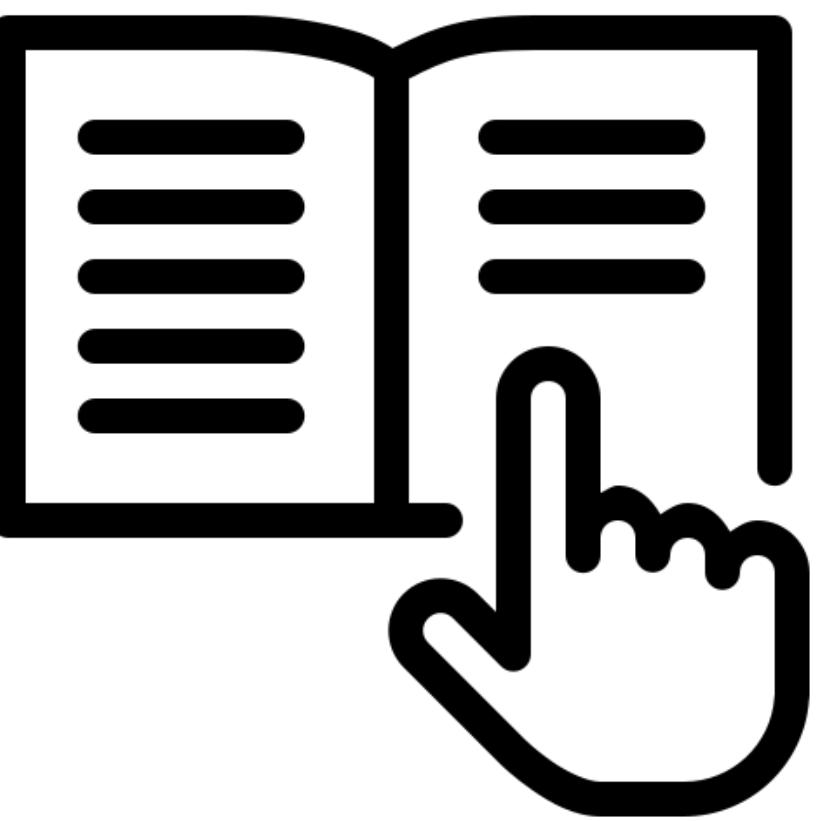
competitive reads



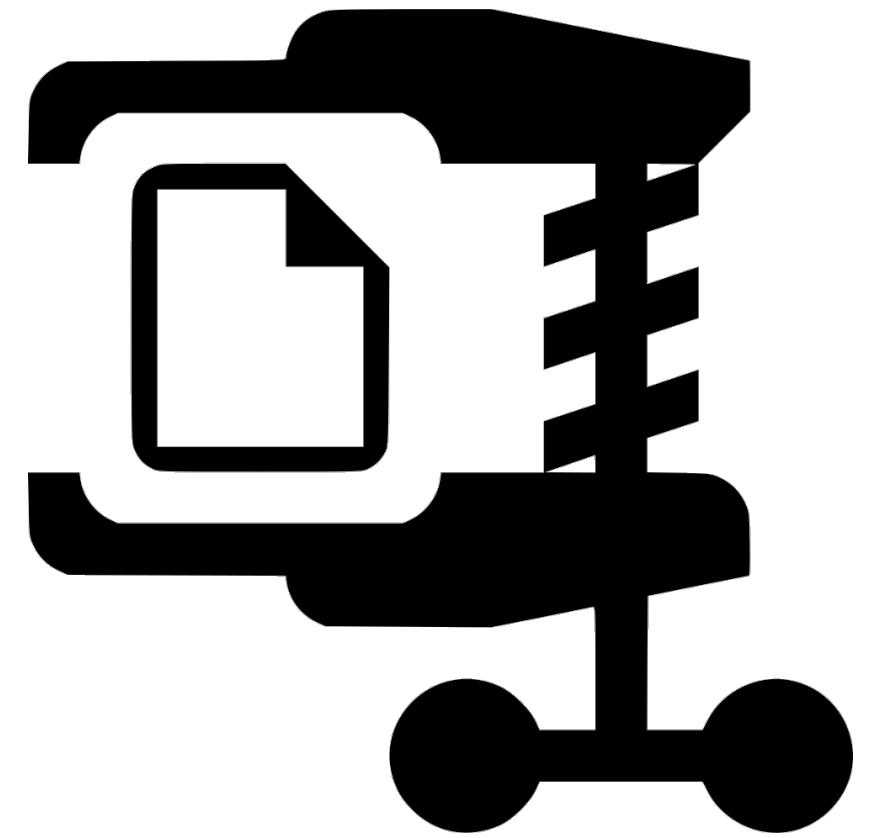
good space
utilization



fast writes



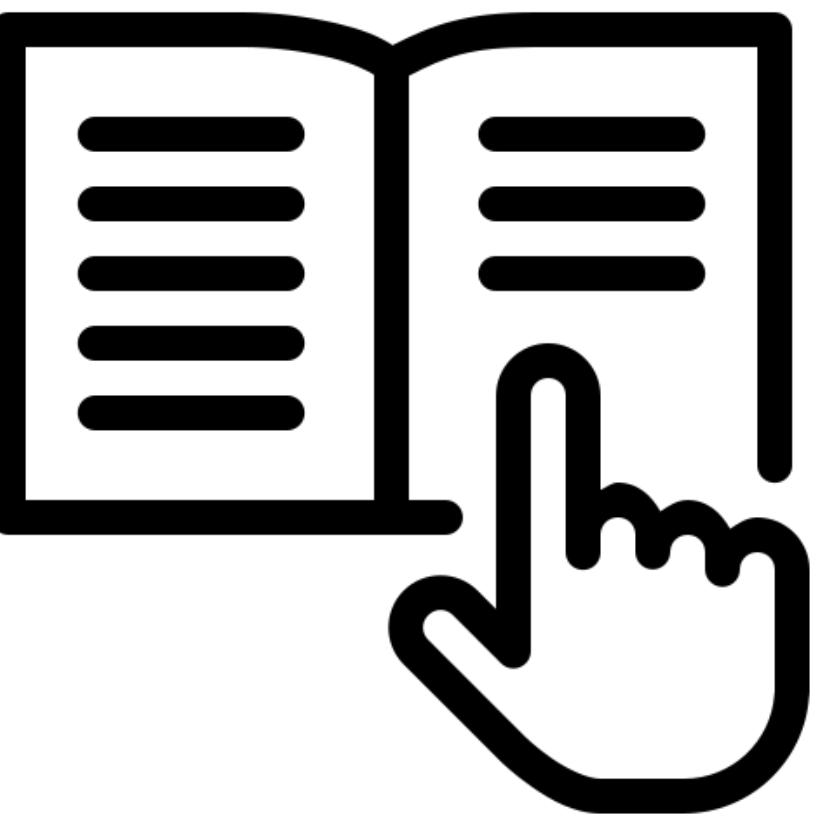
competitive reads



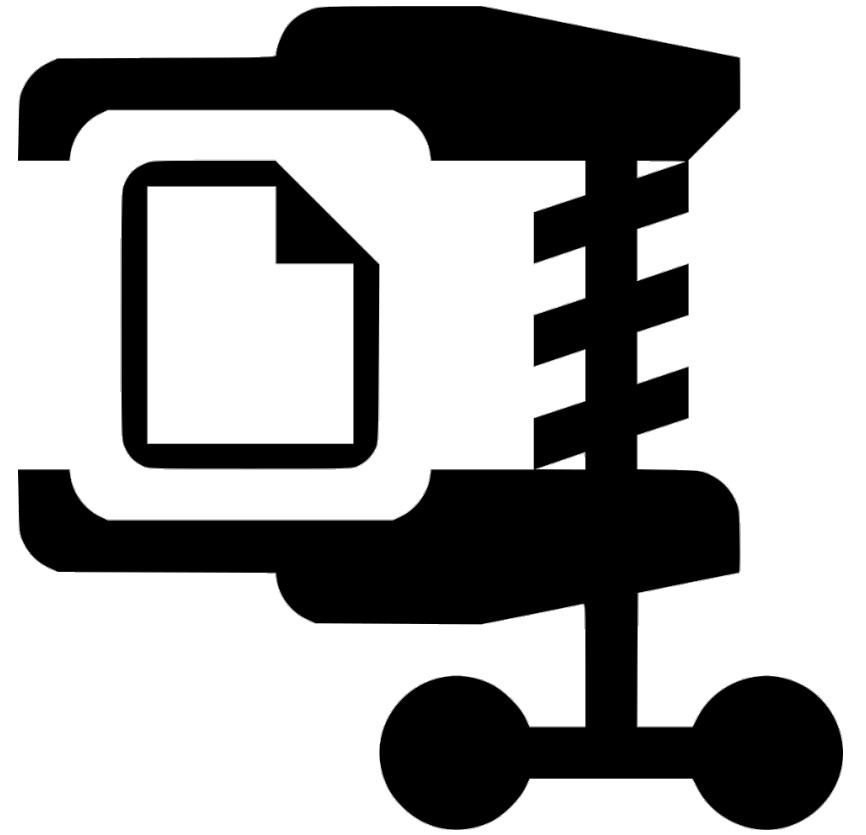
good space
utilization



fast writes



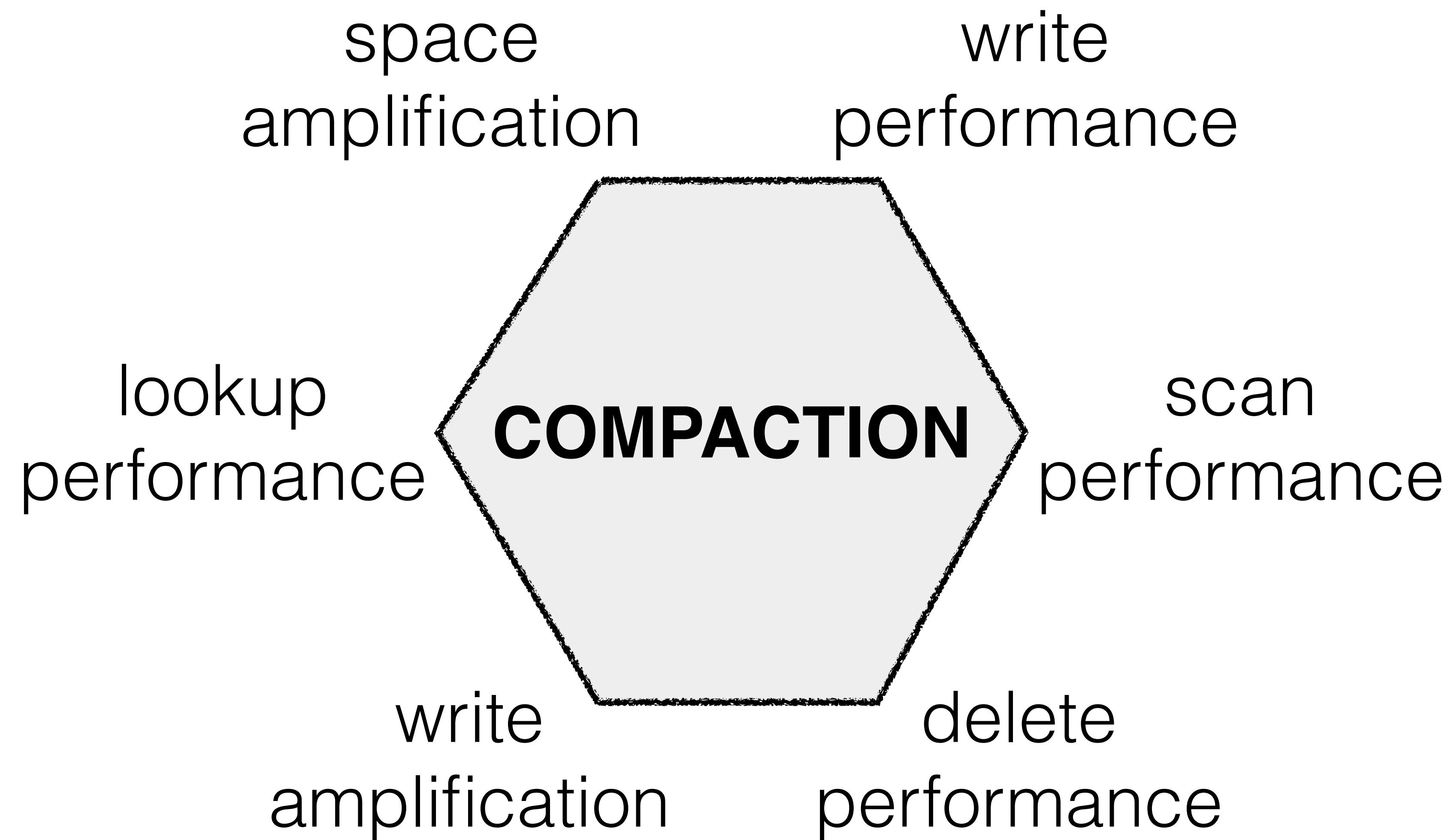
competitive reads



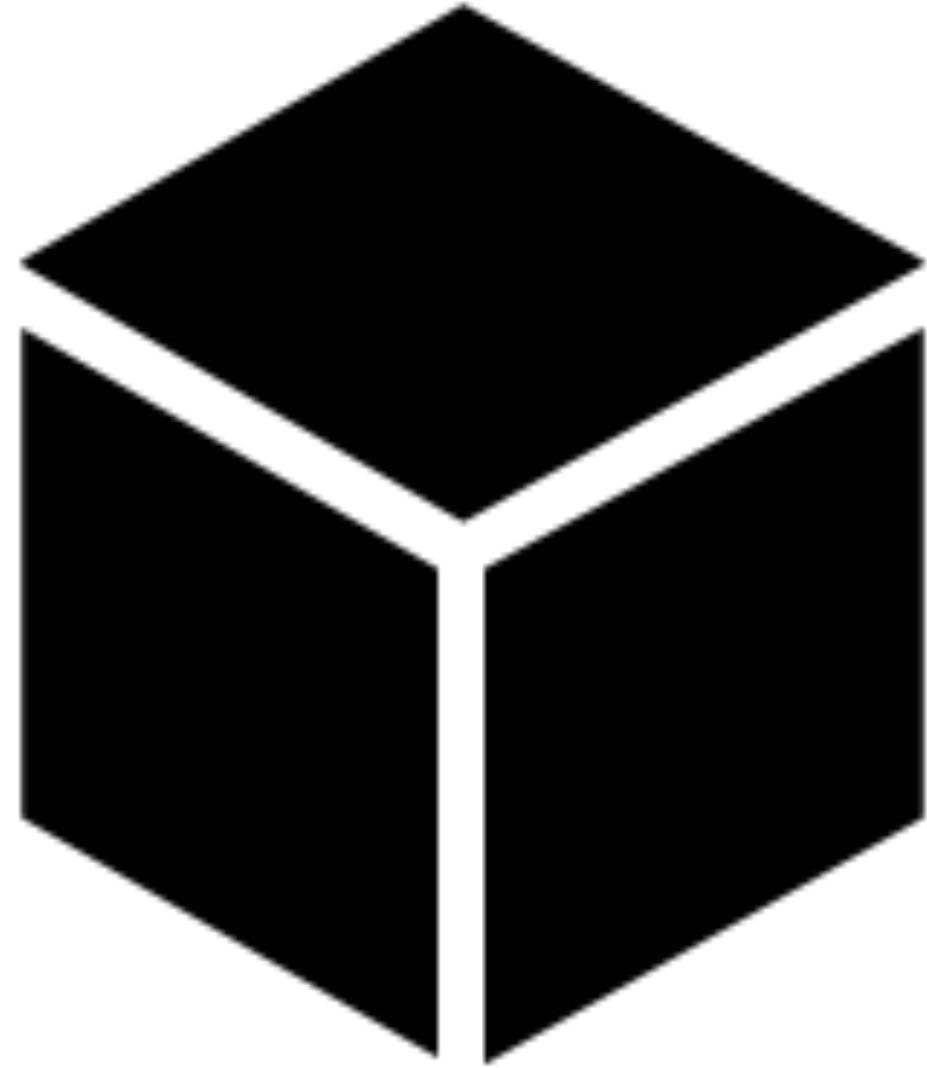
good space
utilization

COMPACTI^ON

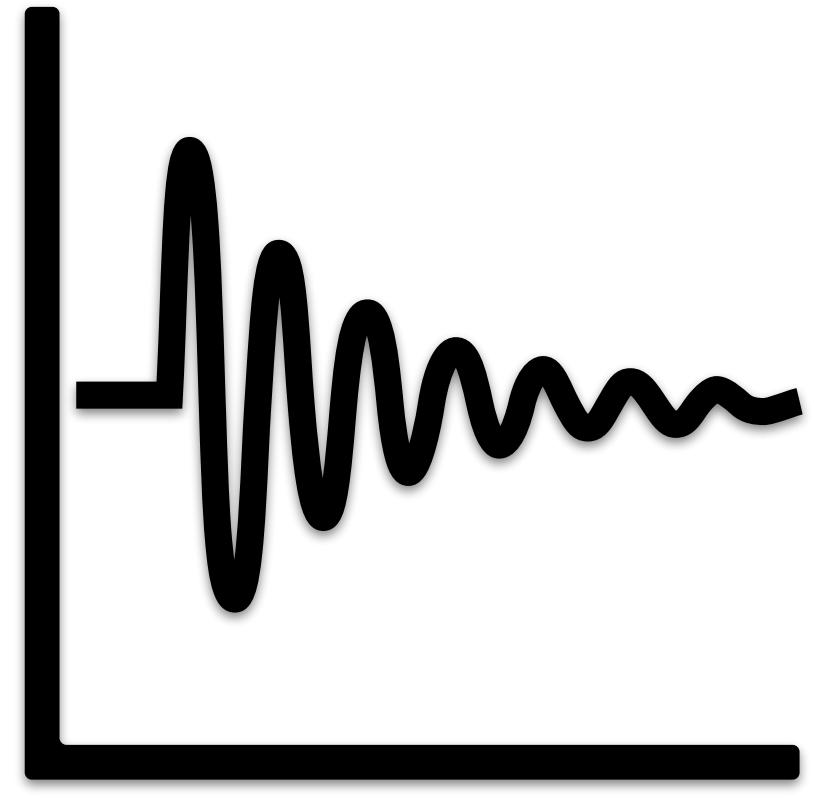
COMPACTION



COMPACTION



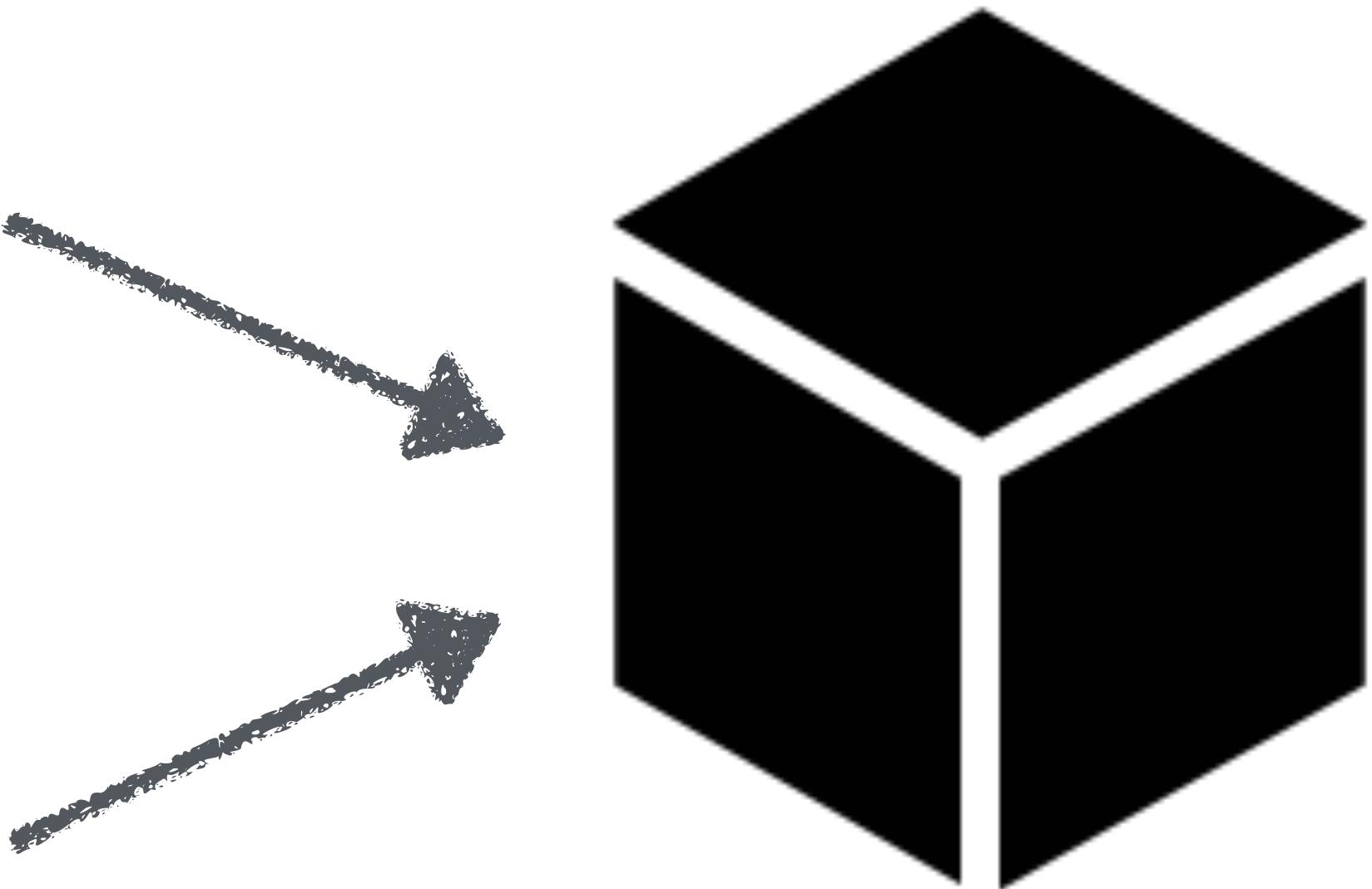
COMPACTION



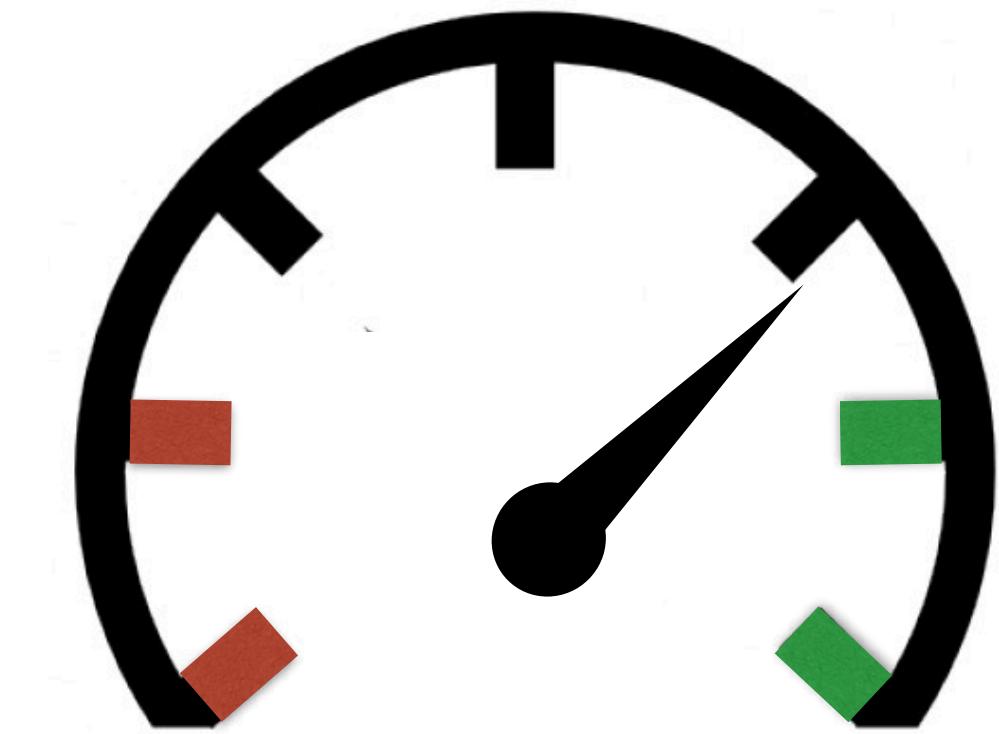
workload



LSM tuning



COMPACTOR



performance

Our Goal

Our Goal

1



Roadmap to pick
compactions

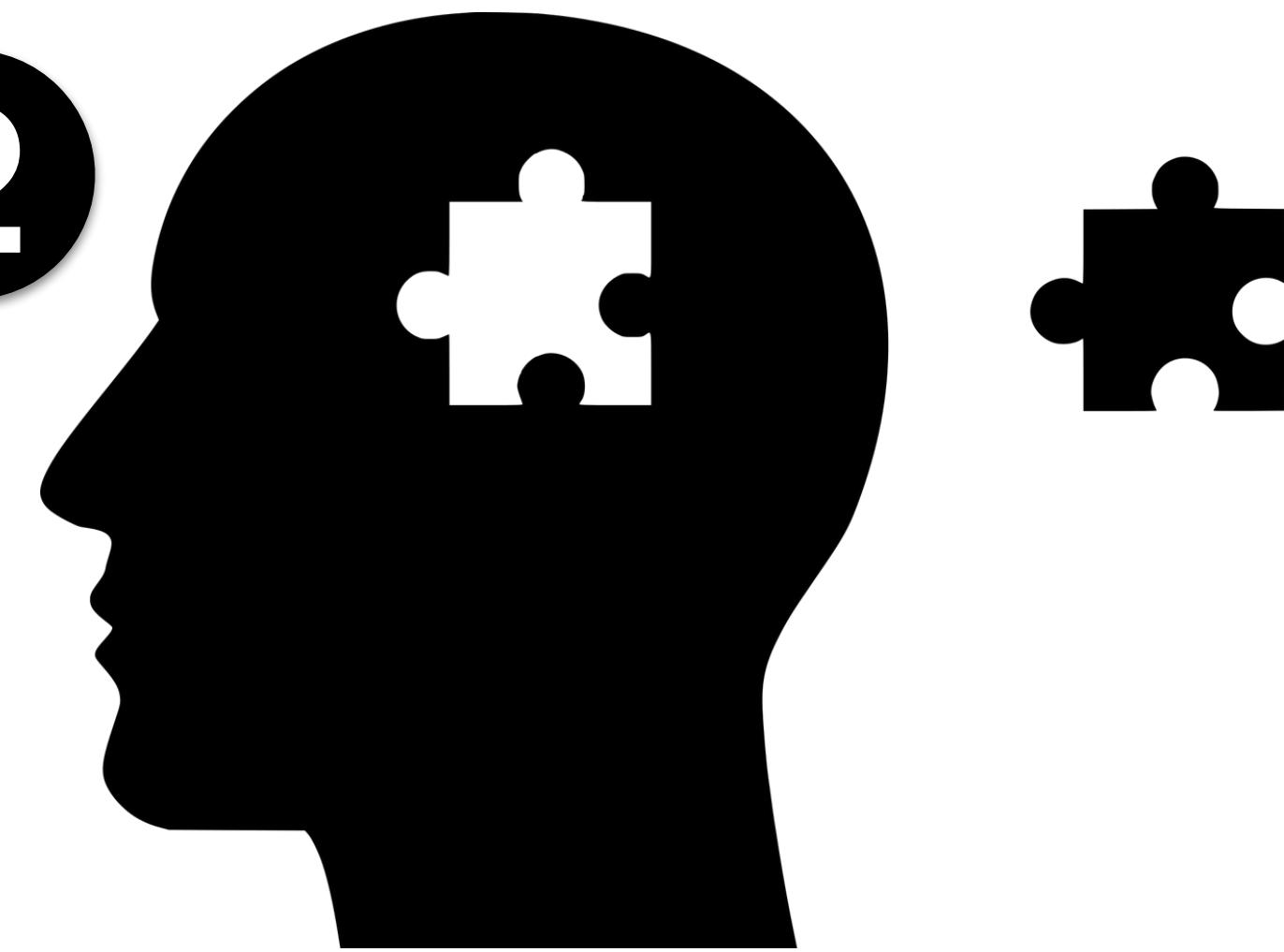
Our Goal

1



Roadmap to pick
compactions

2



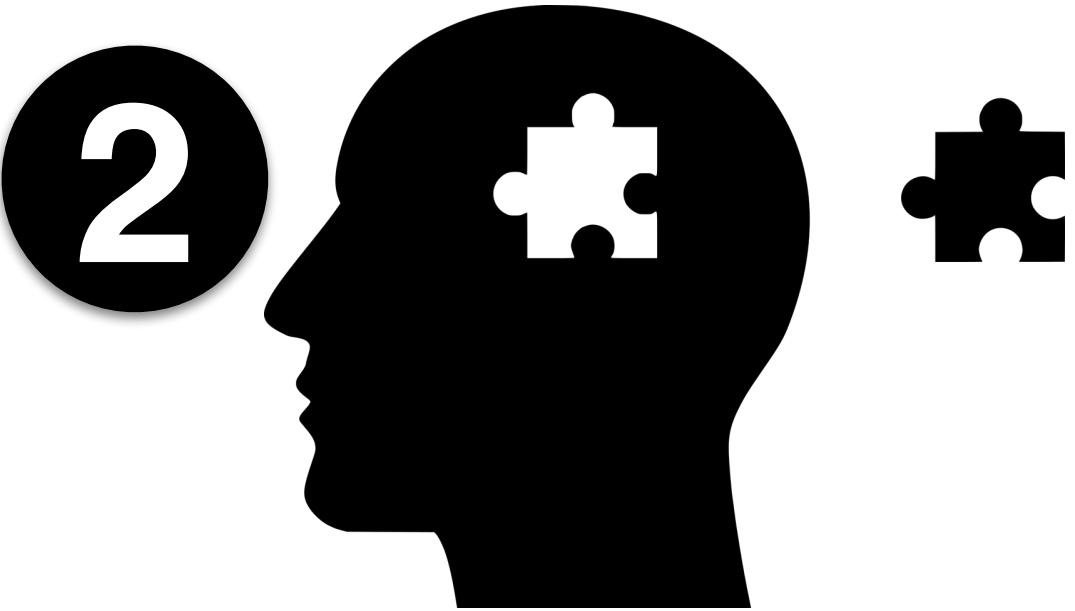
Answer to complex
design questions

Our Goal

1



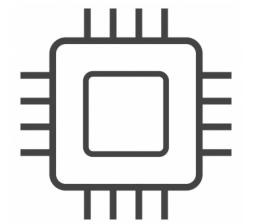
2



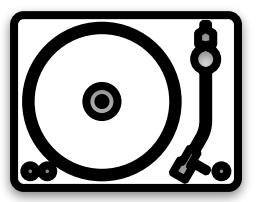
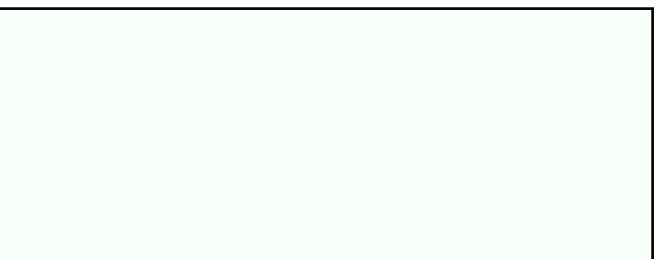
break the black box

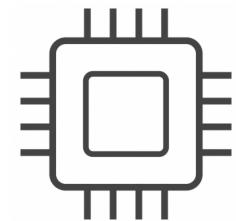


learn from 2000+ experiments



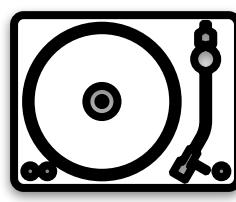
buffer

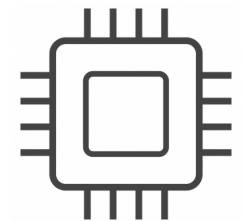




buffer

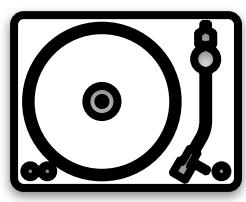
2	6	1	4
---	---	---	---

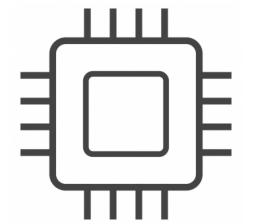




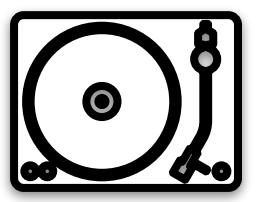
buffer

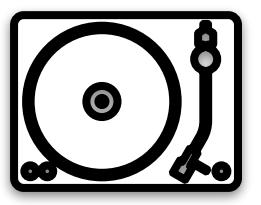
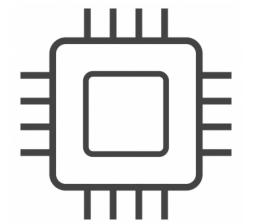
1	2	4	6
---	---	---	---



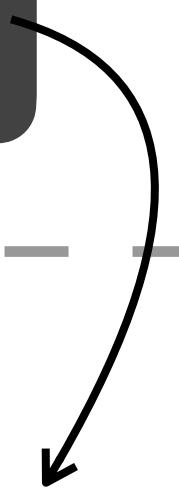


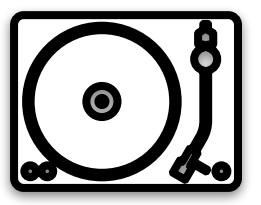
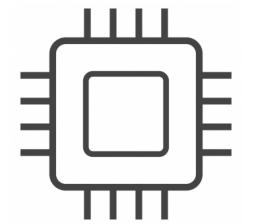
buffer



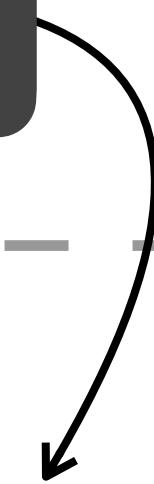


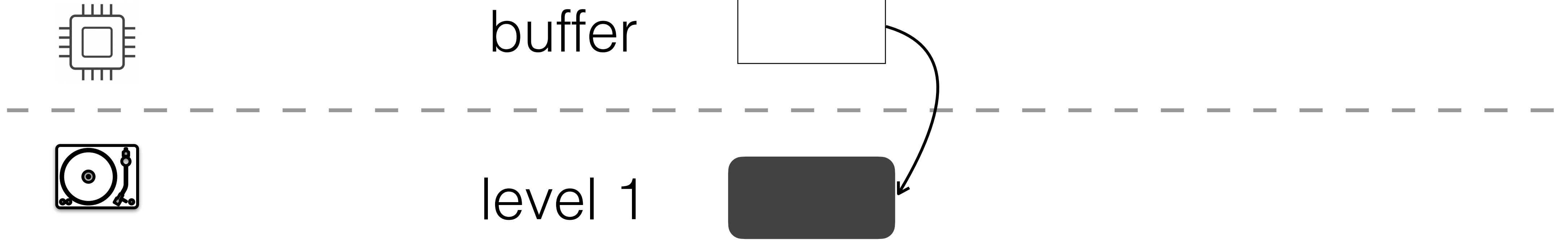
buffer

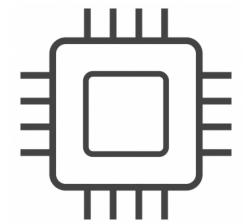




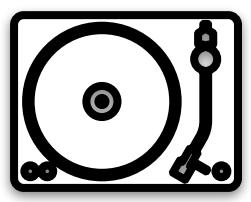
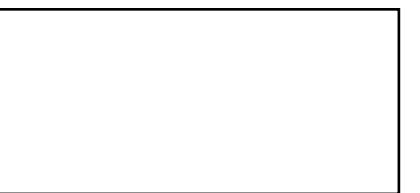
buffer





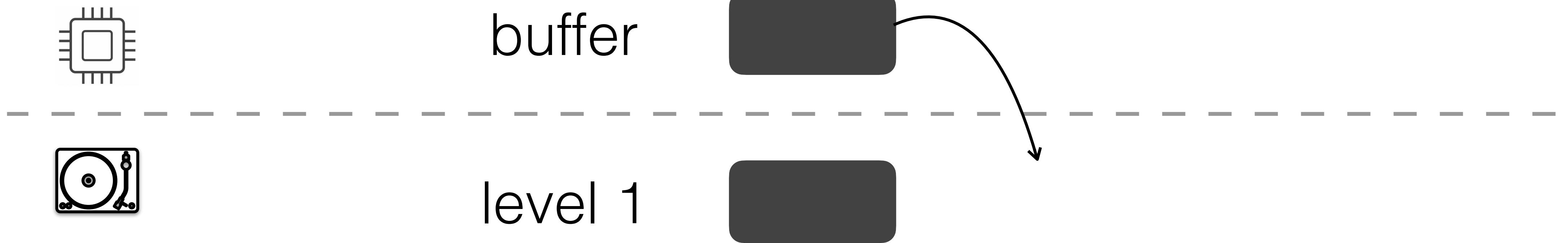


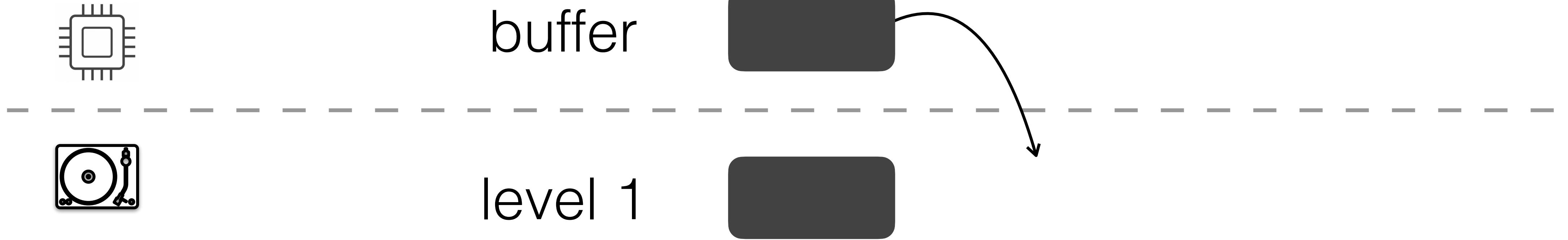
buffer

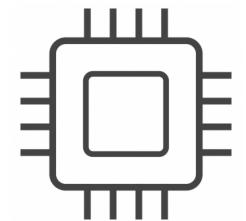


level 1

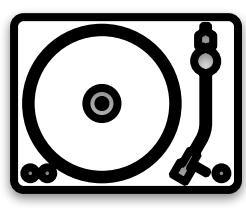




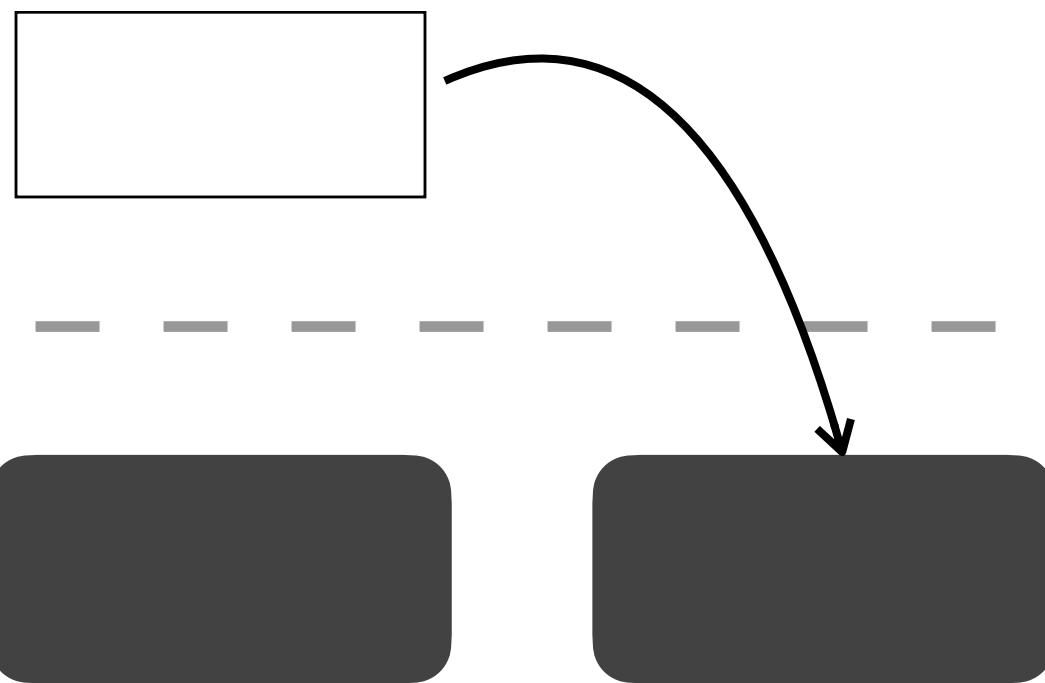


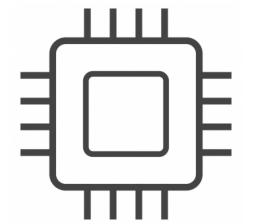


buffer

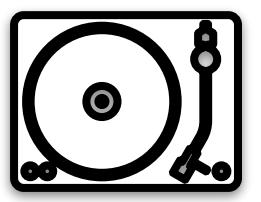
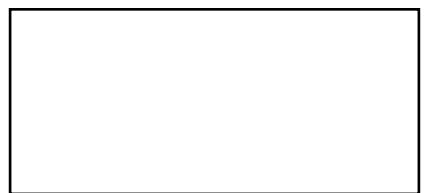


level 1



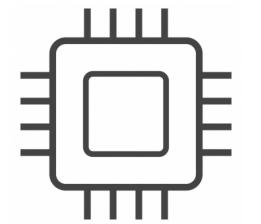


buffer

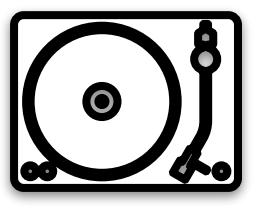
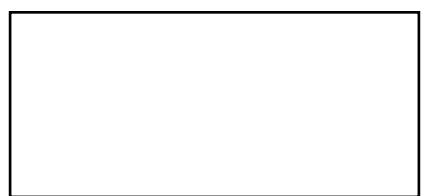


level 1



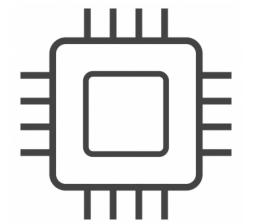


buffer

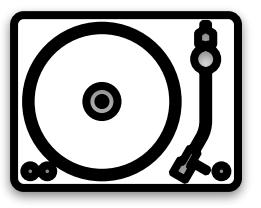


level 1





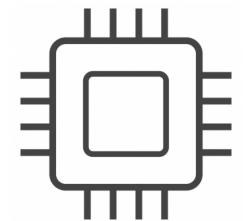
buffer



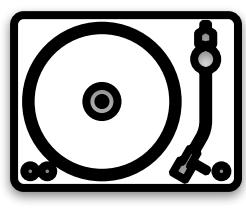
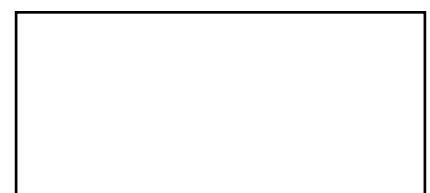
level 1

level 2





buffer

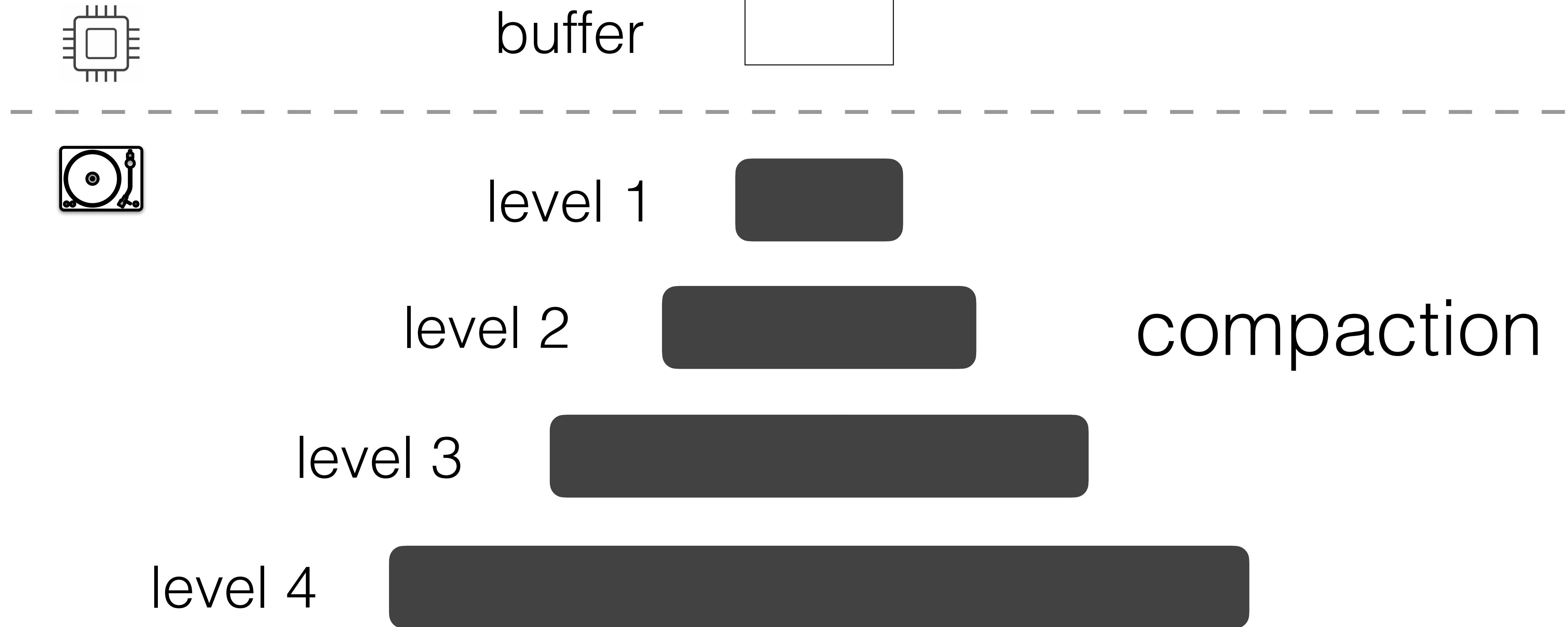


level 1

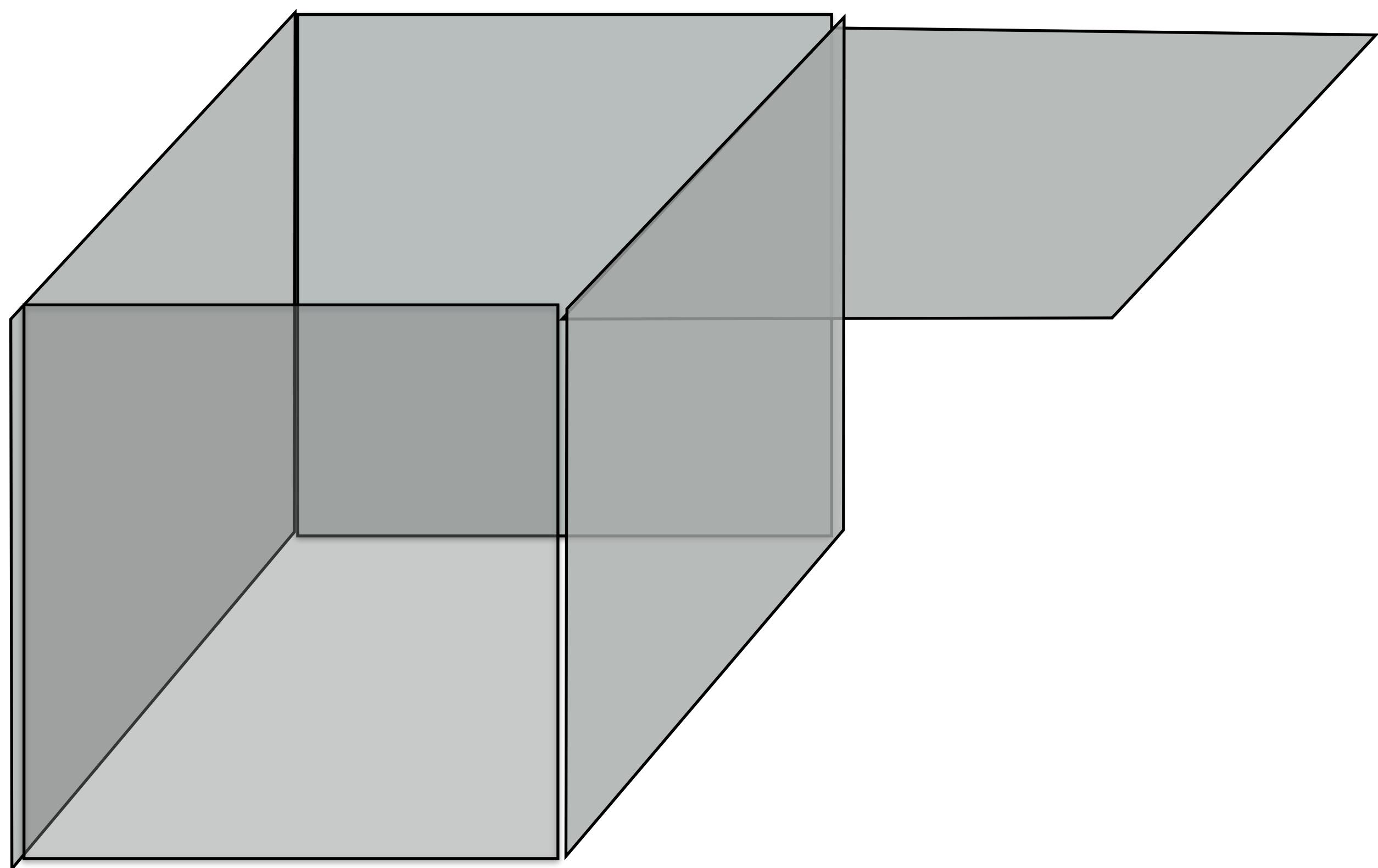
level 2

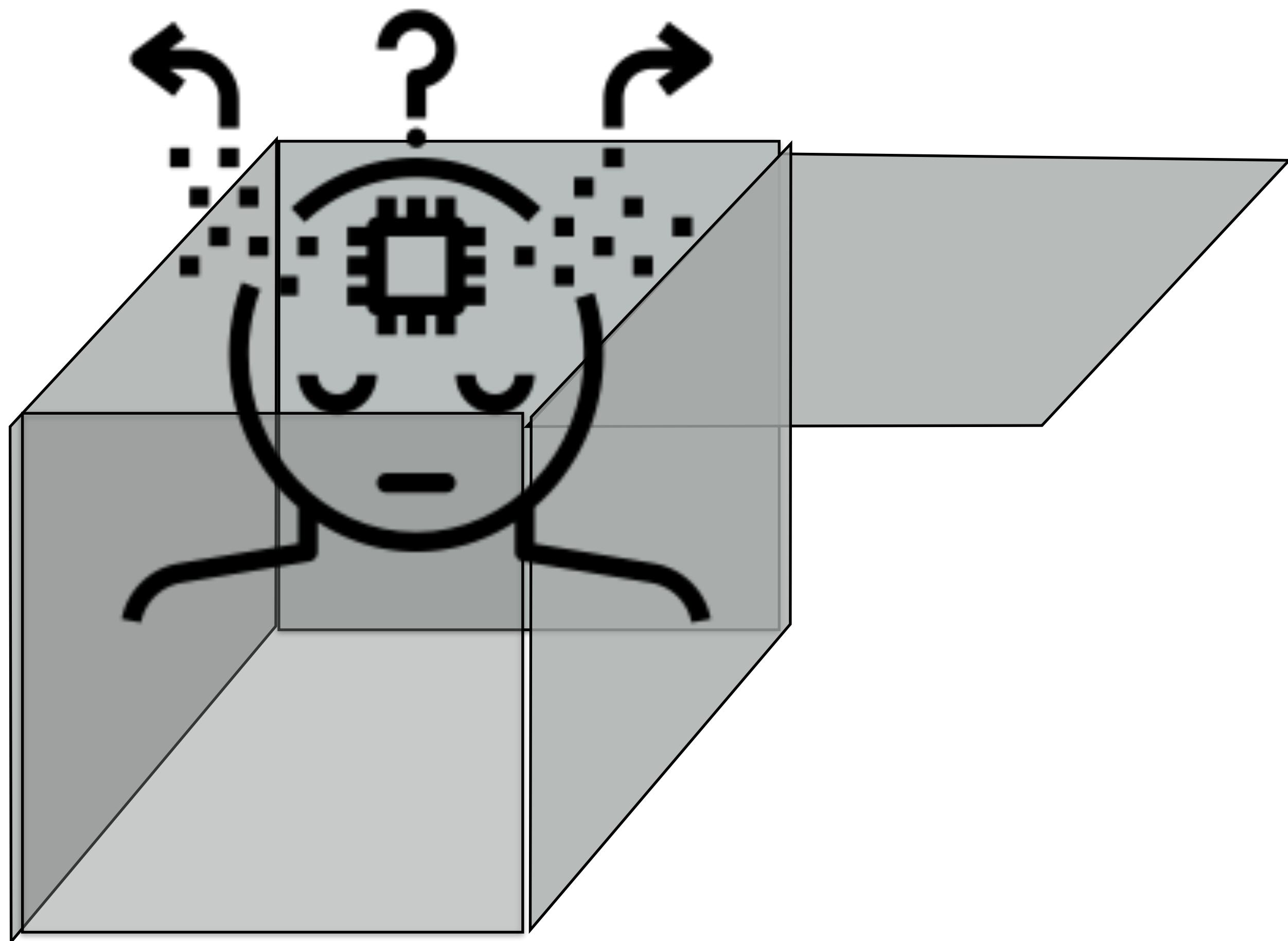


compaction

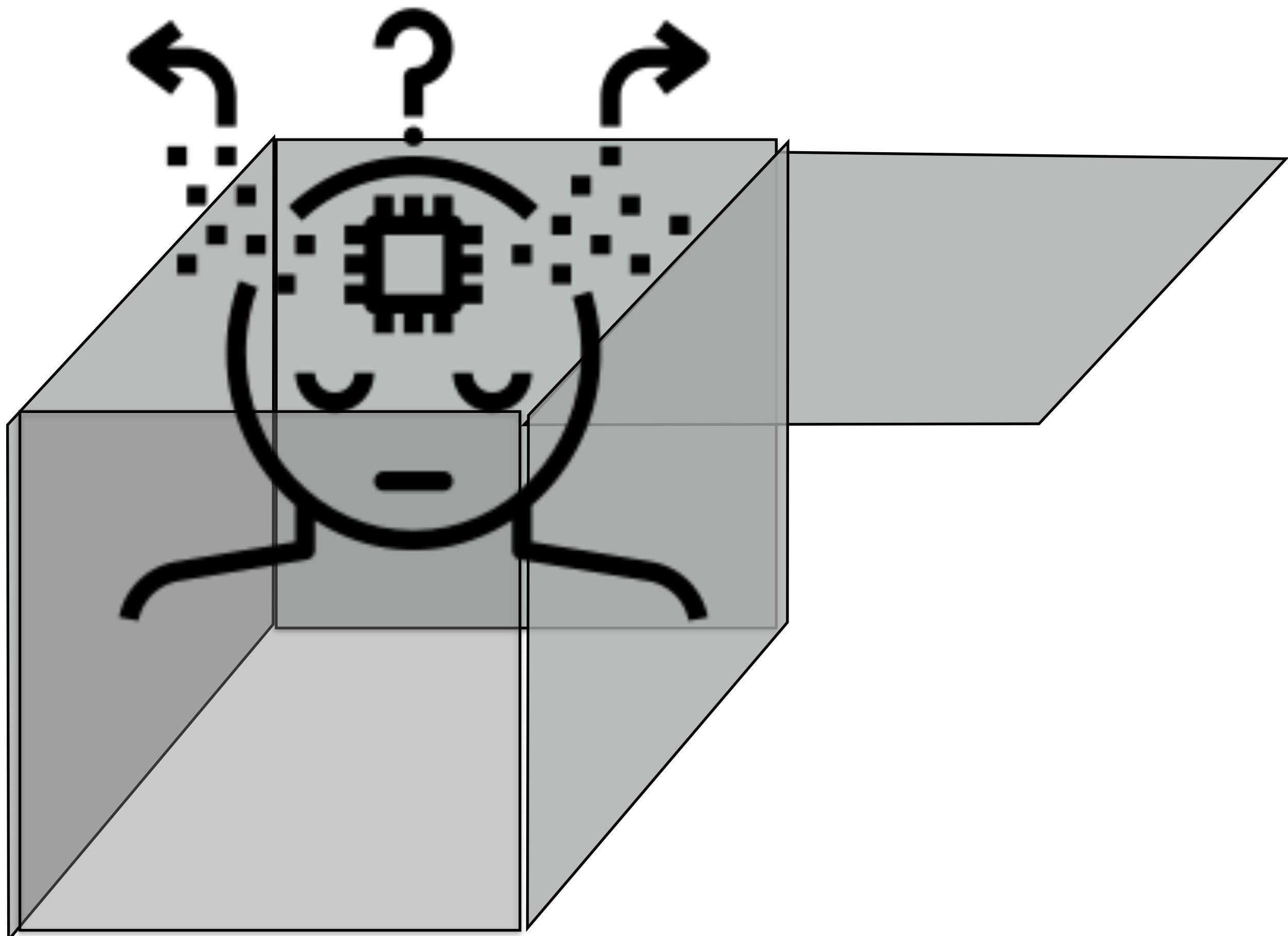


compaction?

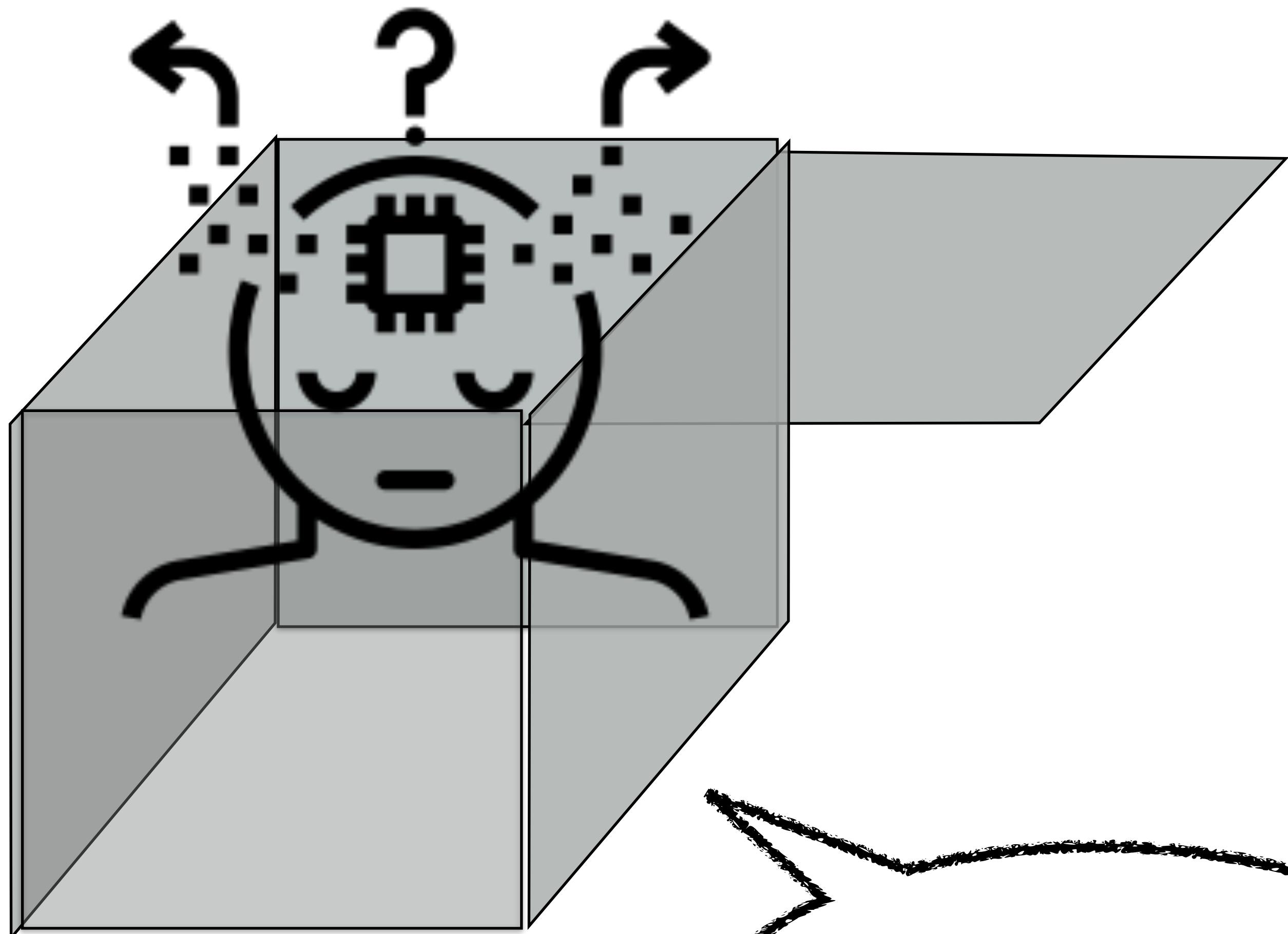




How many runs
per level?

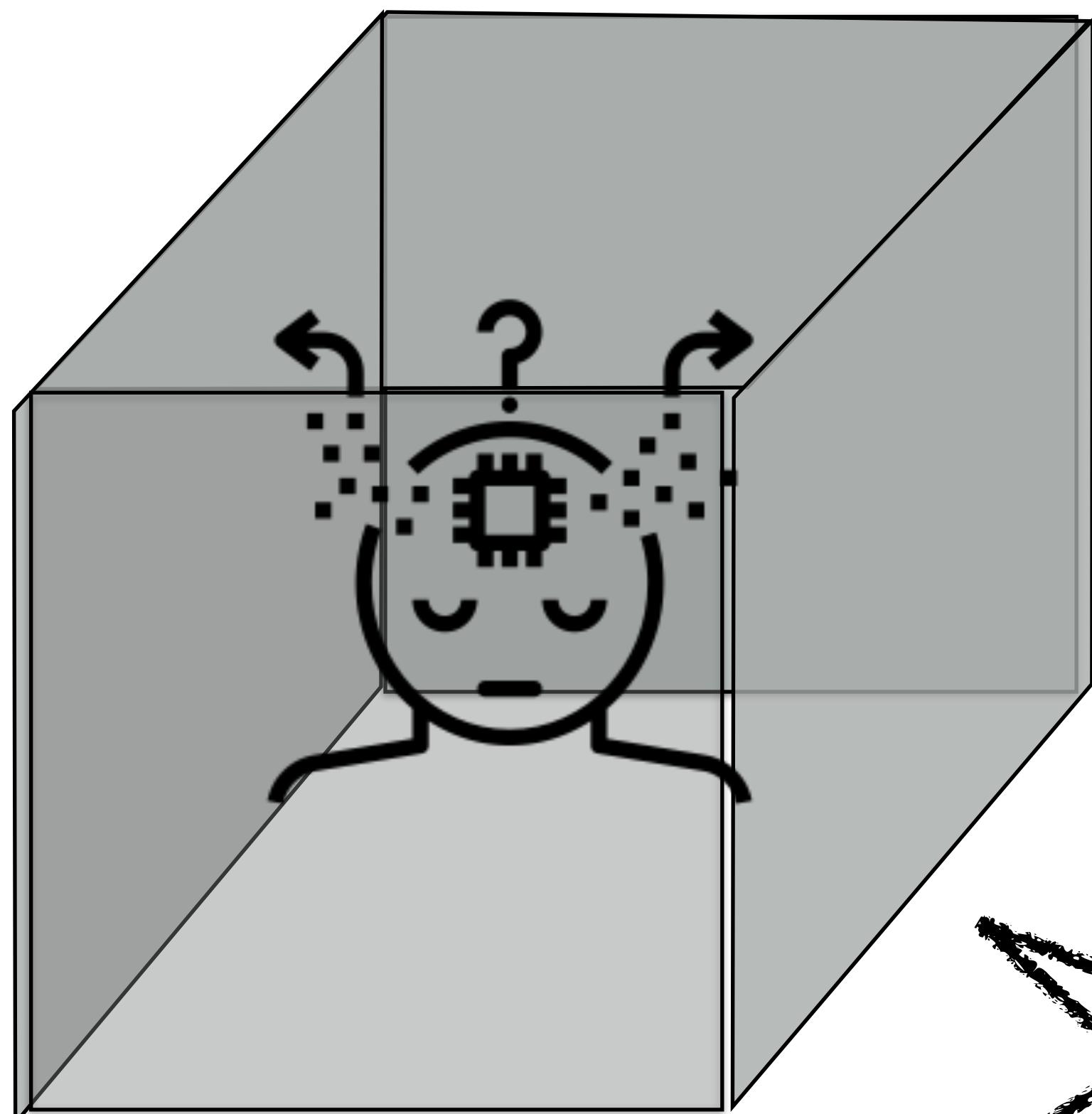


How many runs
per level?

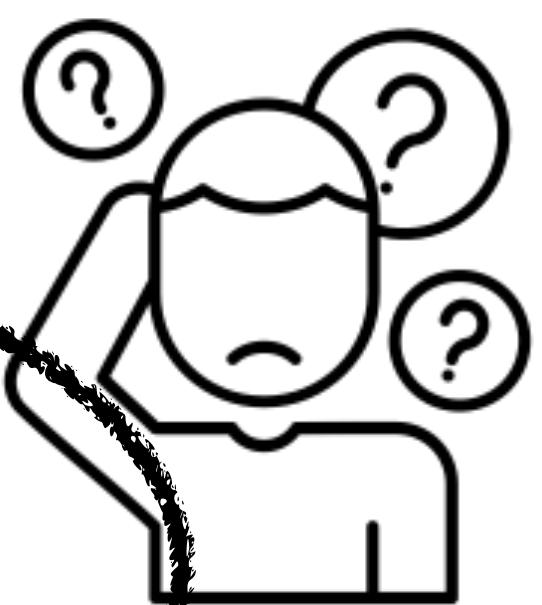


How much data to
compact at once?

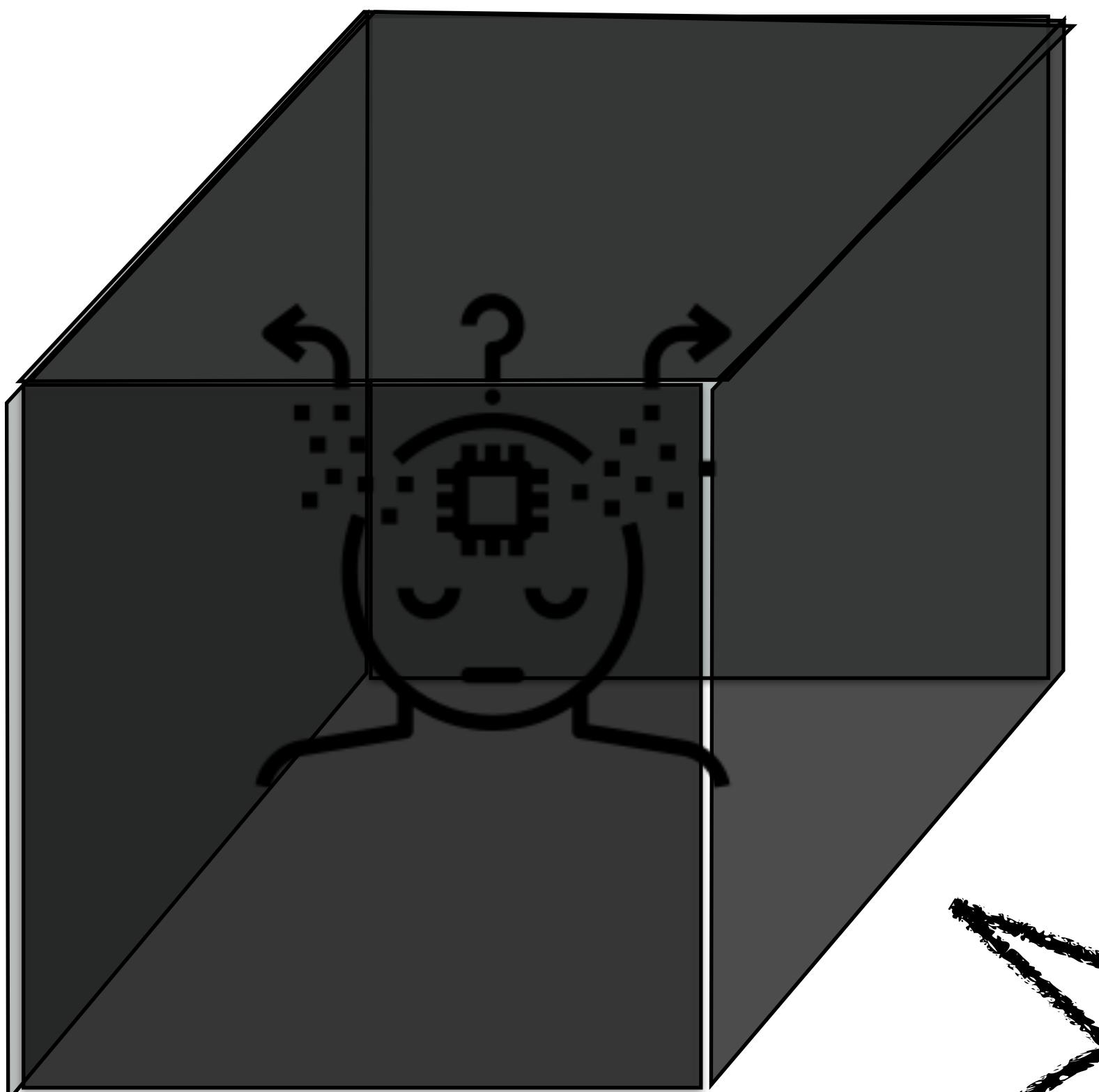
How many runs
per level?



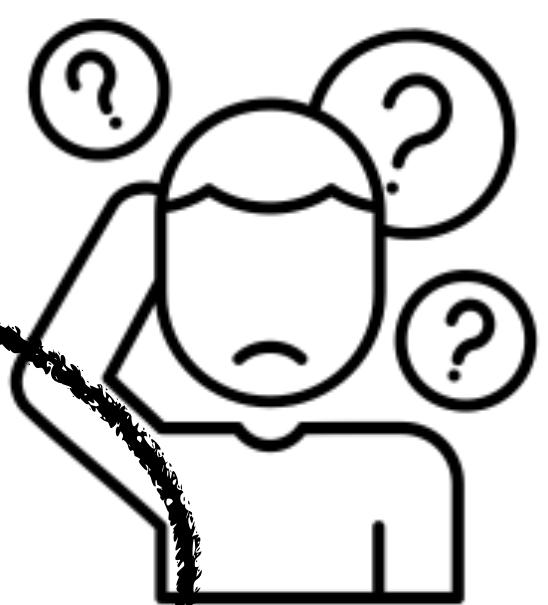
How much data to
compact at once?



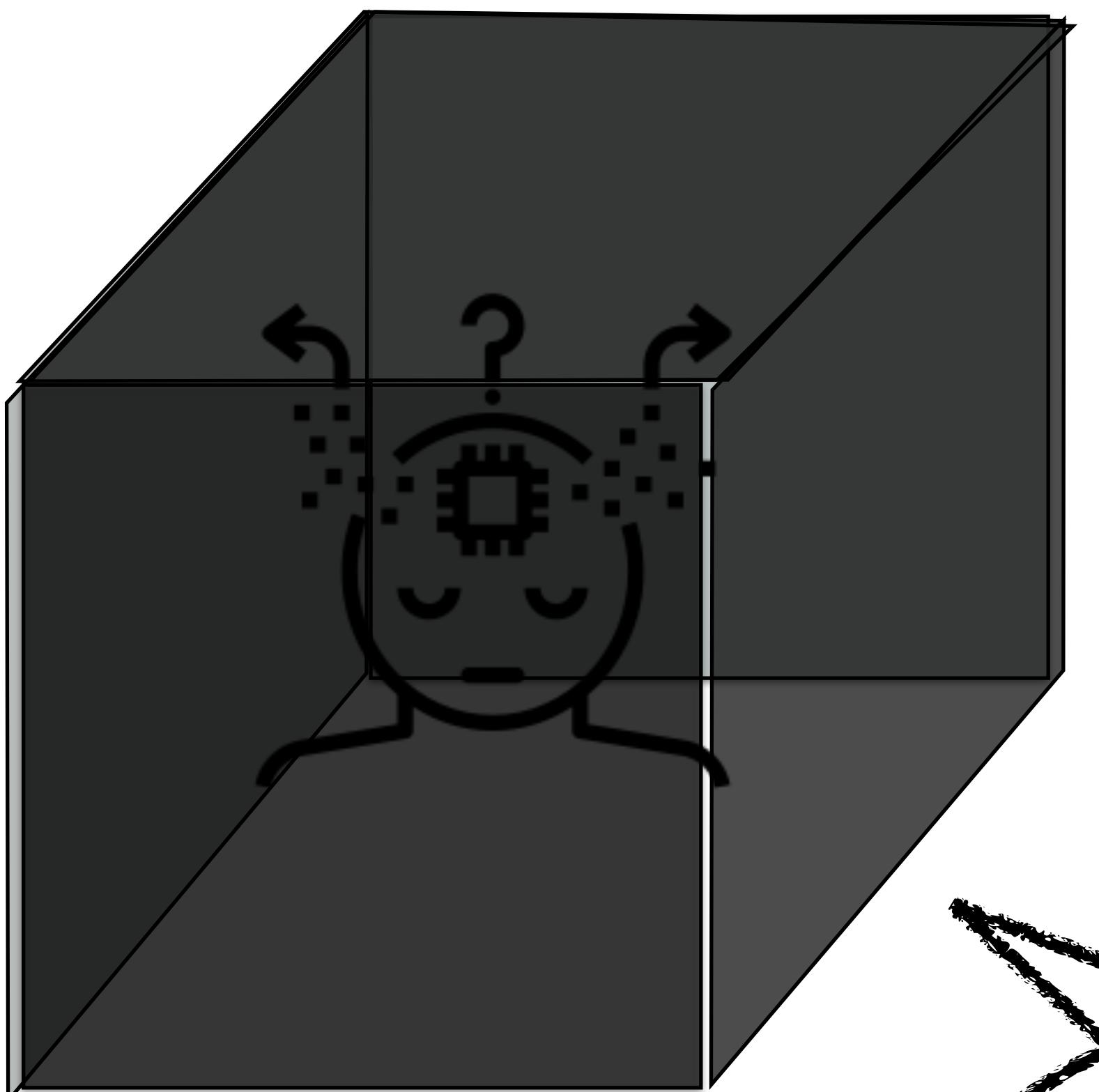
How many runs
per level?



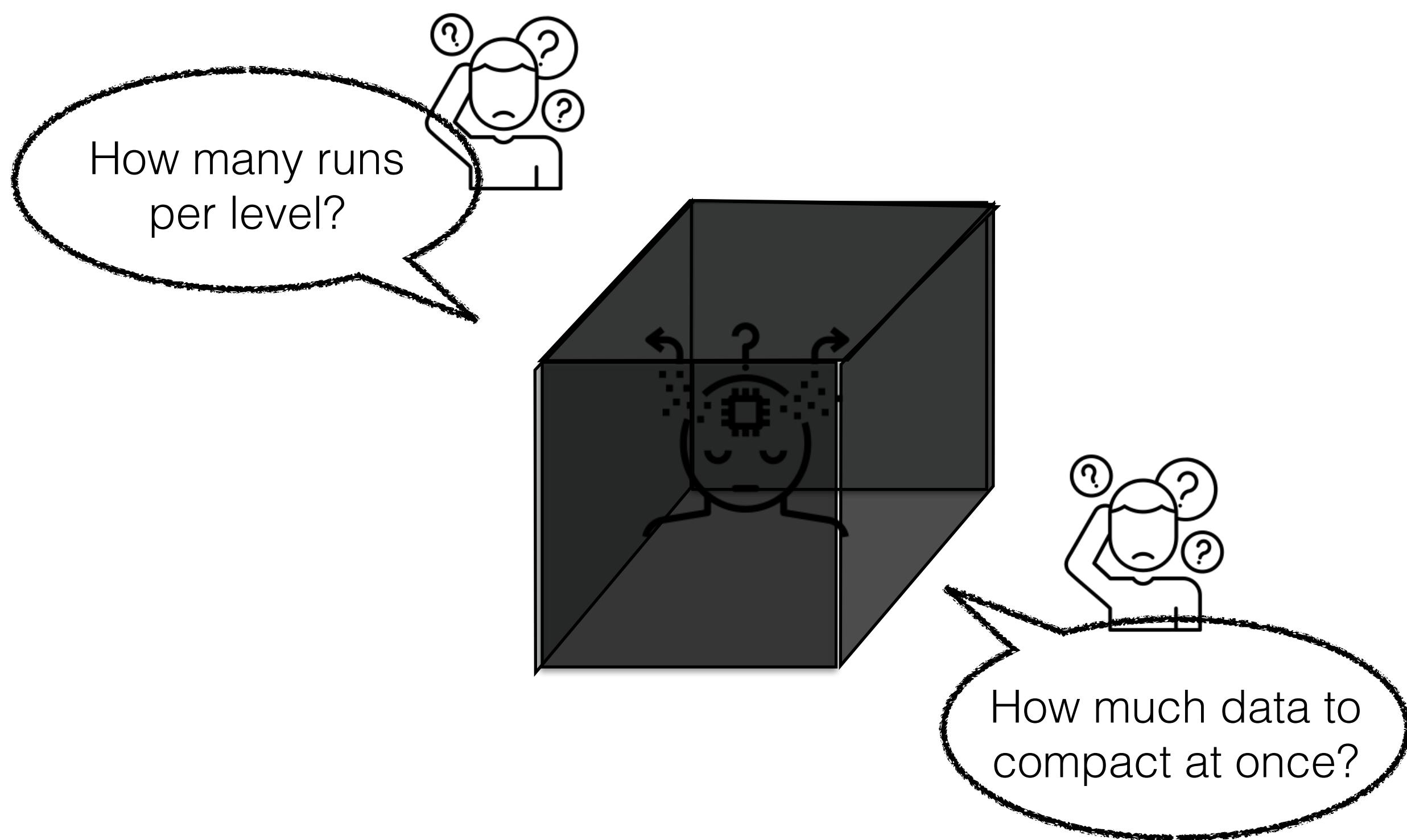
How much data to
compact at once?



How many runs
per level?

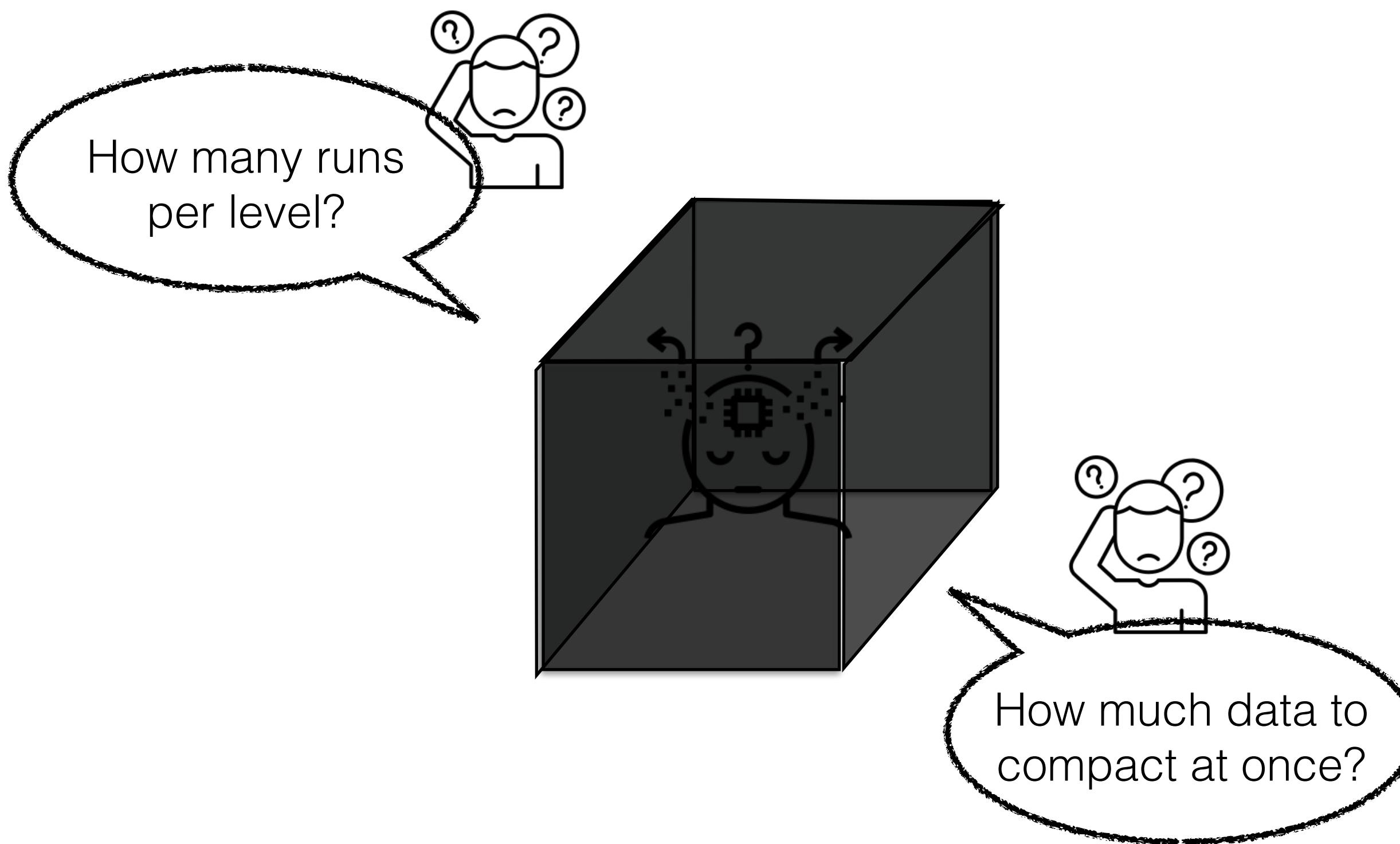


How much data to
compact at once?



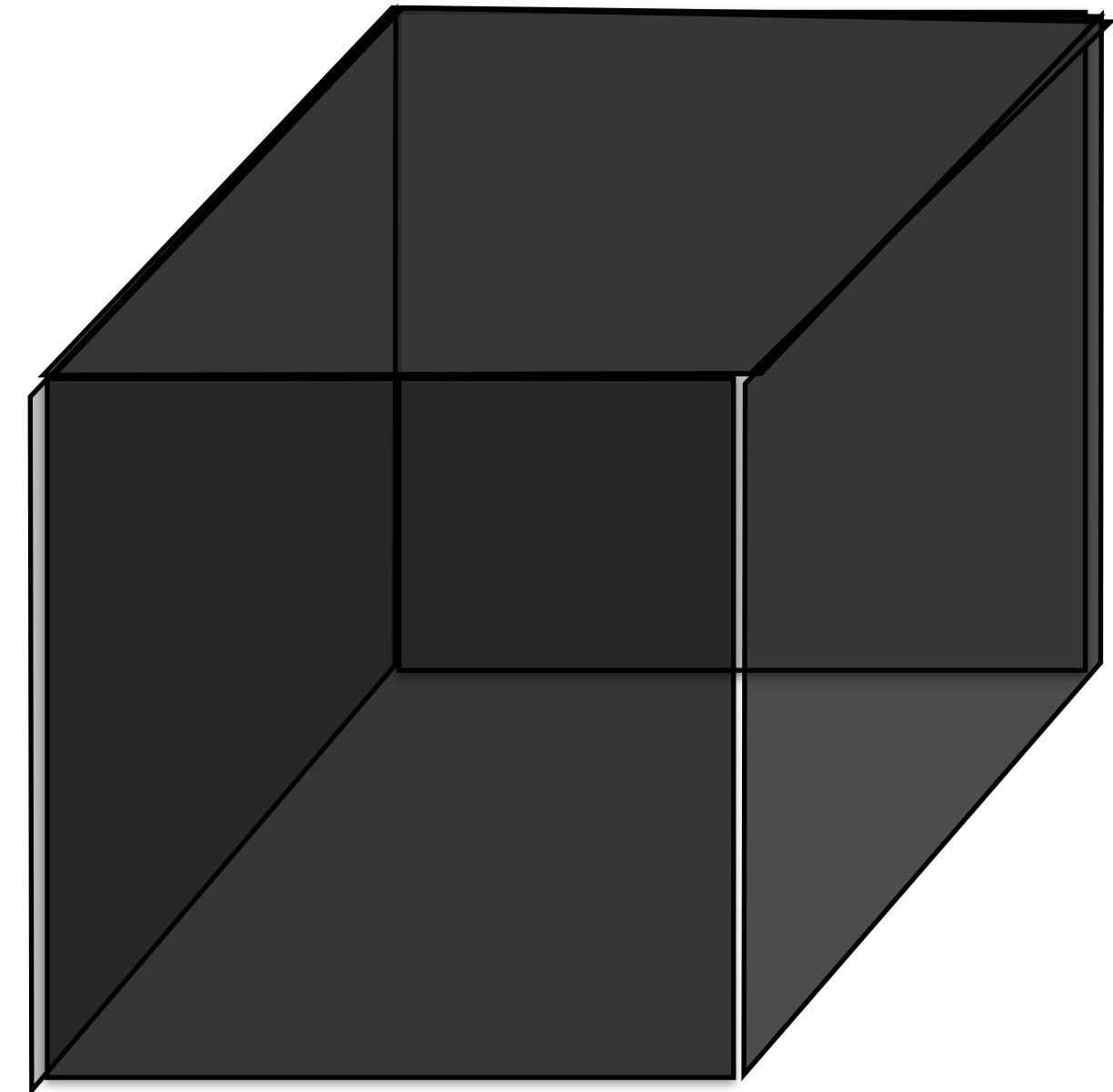
What are the **design choices**?

How does a choice **affect performance**?



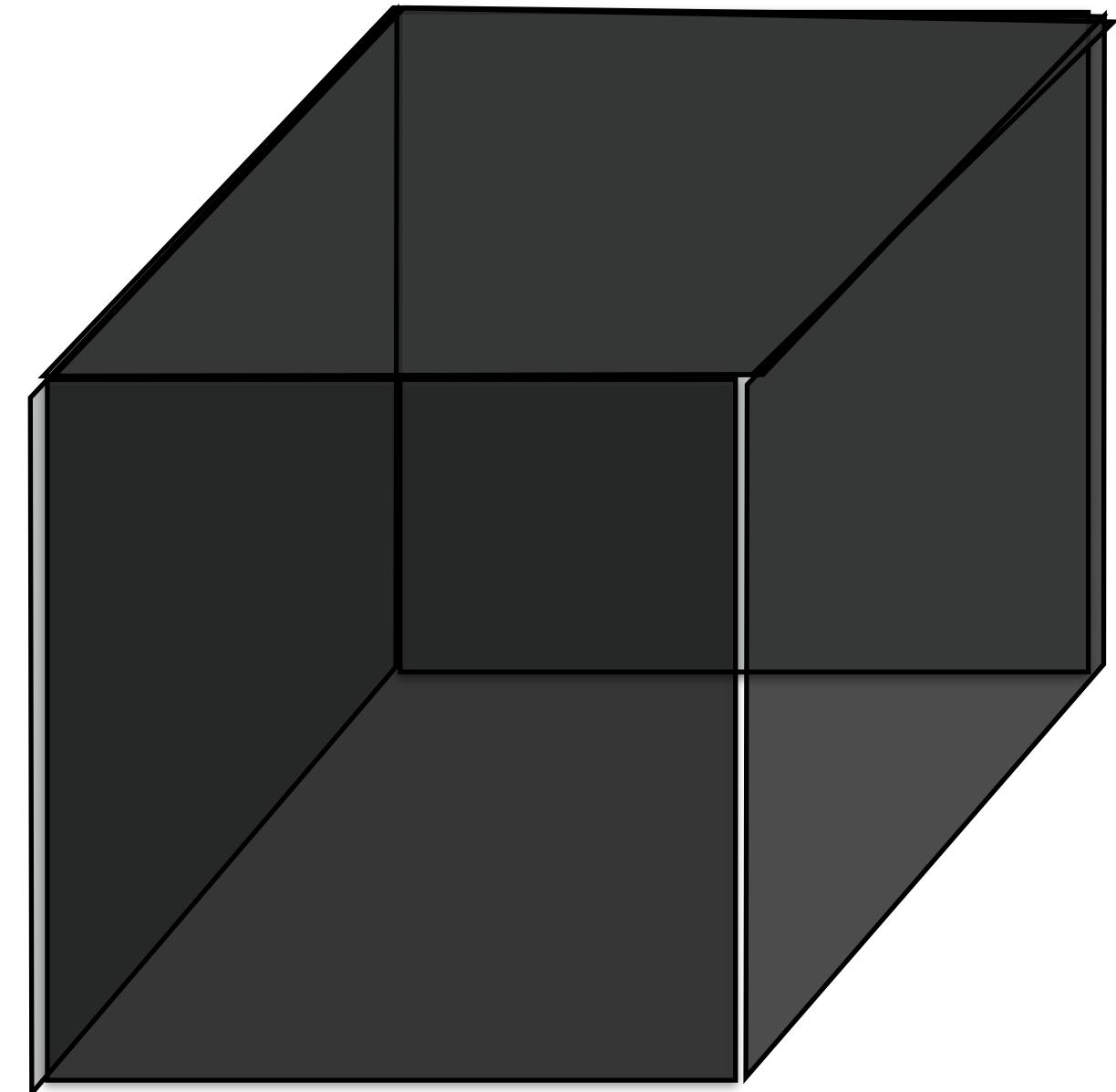
What are the **design choices**?

How does a choice
affect performance?



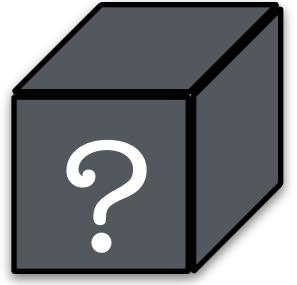
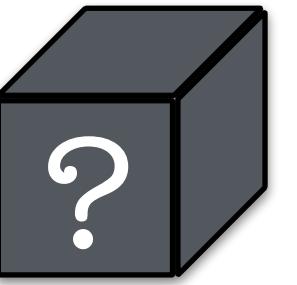
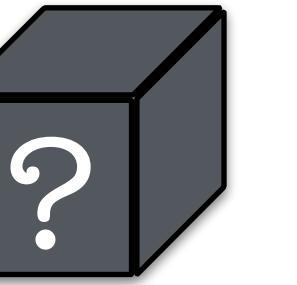
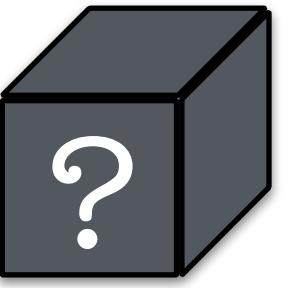
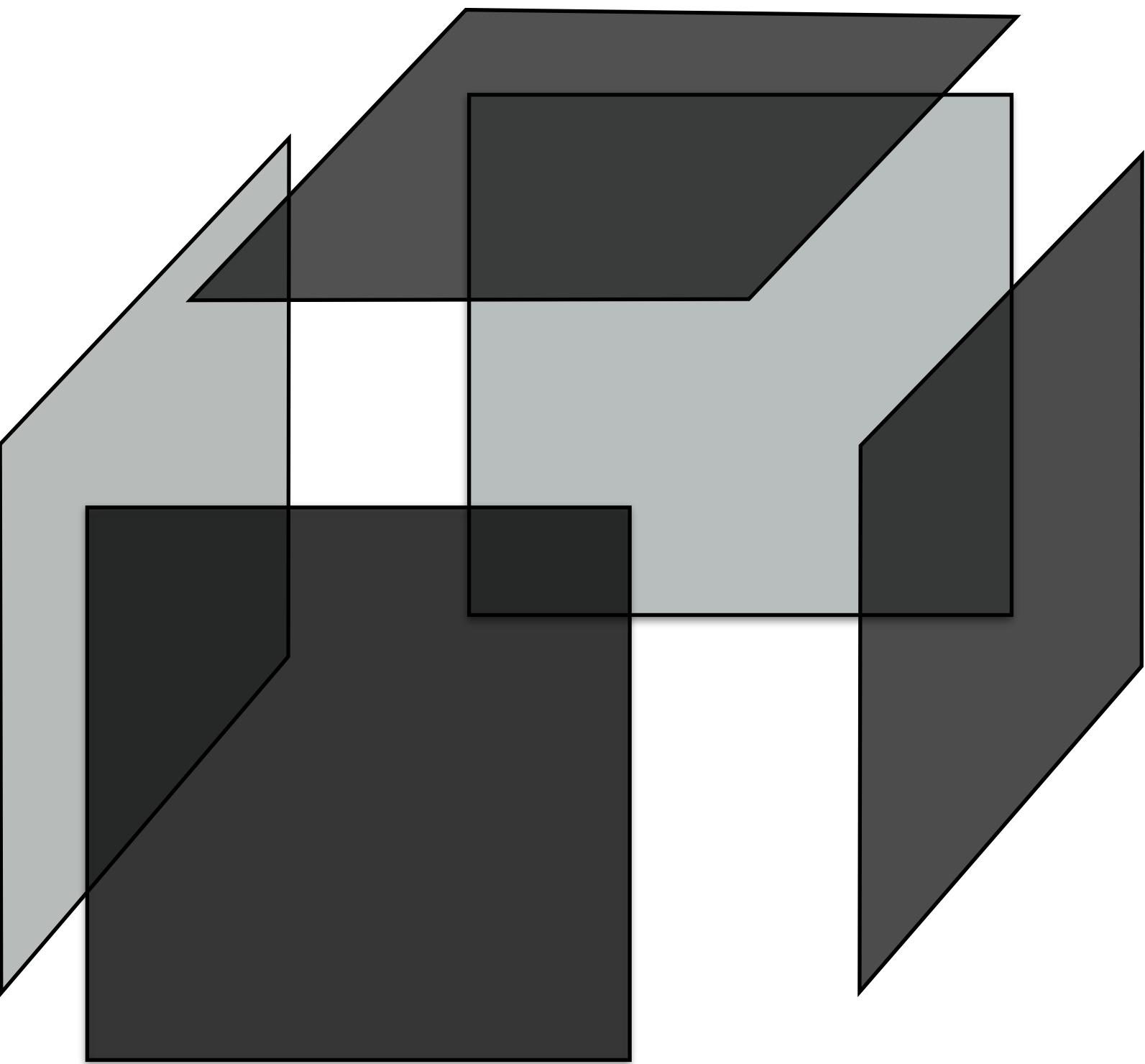
What are the **design choices**?

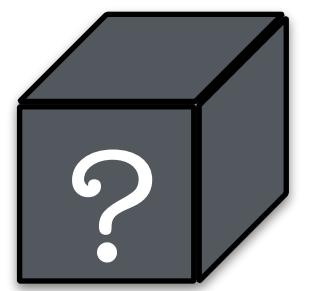
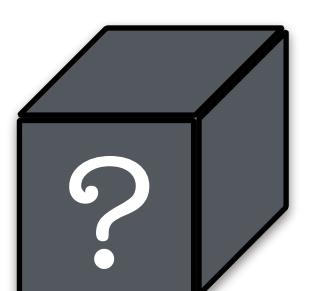
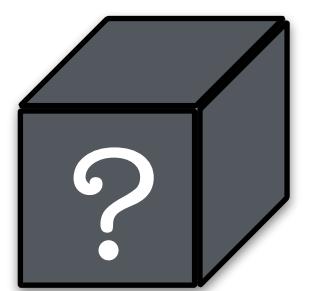
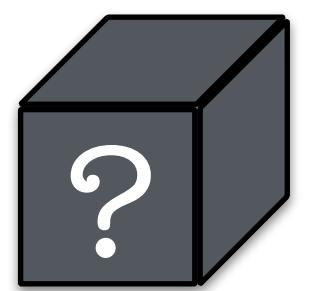
How does a choice
affect performance?

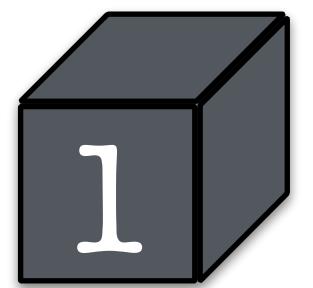


What are the **design choices**?

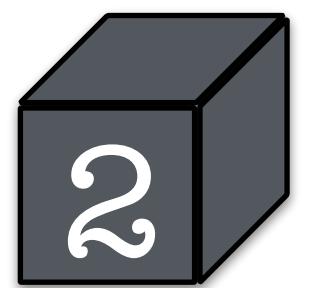
How does a choice
affect performance?



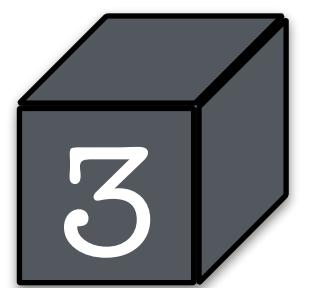




How to organize the data on device?



How much data to move at-a-time?



Which block of data to be moved?



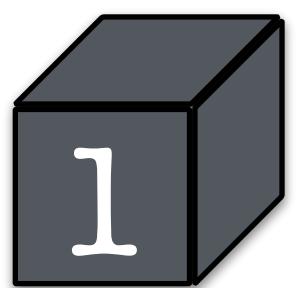
When to re-organize the data layout?

Data Layout

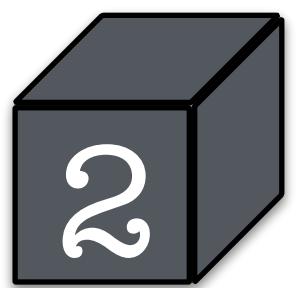
Compaction
granularity

Data movement
policy

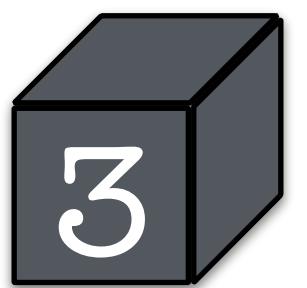
Trigger



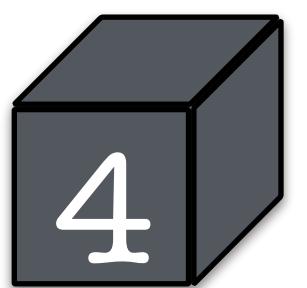
How to organize the data on device?



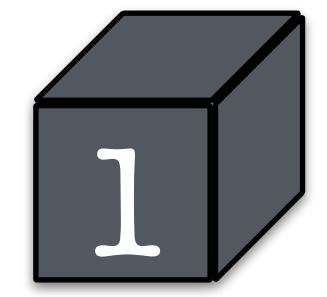
How much data to move at-a-time?



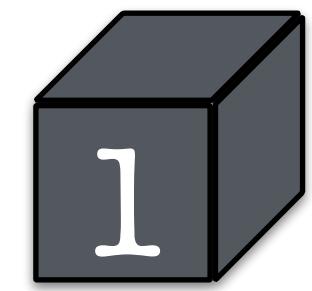
Which block of data to be moved?



When to re-organize the data layout?

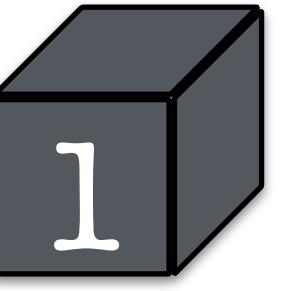


Data Layout



Data Layout

number of runs per level

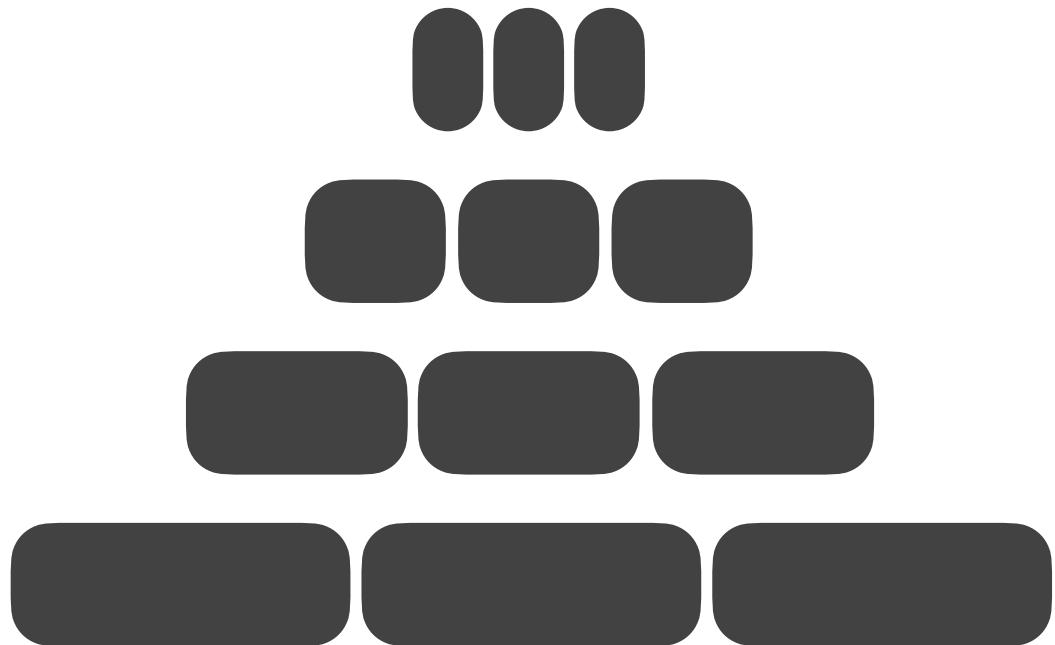
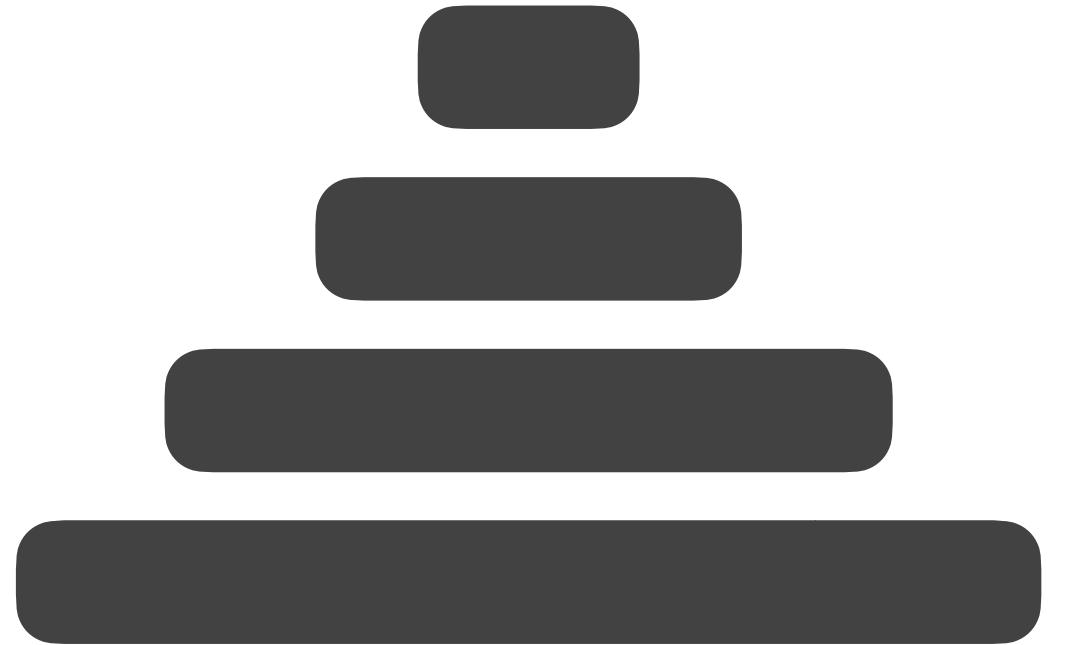


Data Layout

number of runs per level

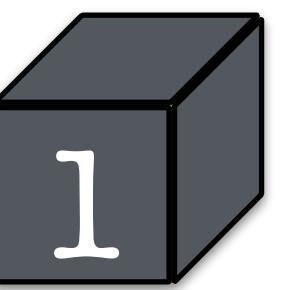
leveling

[eager]



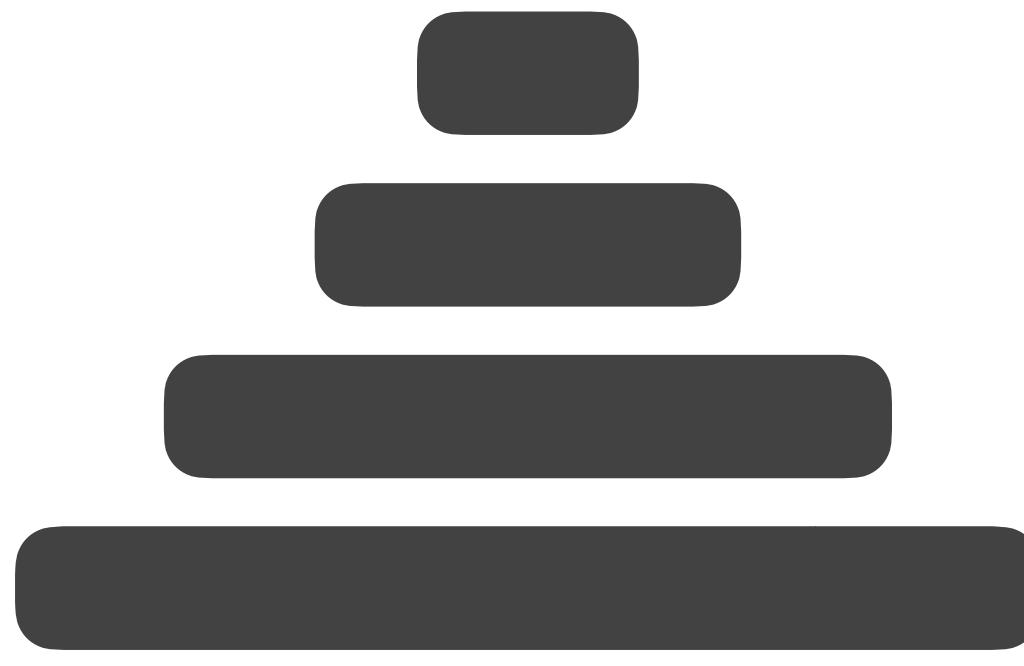
tiering

[lazy]

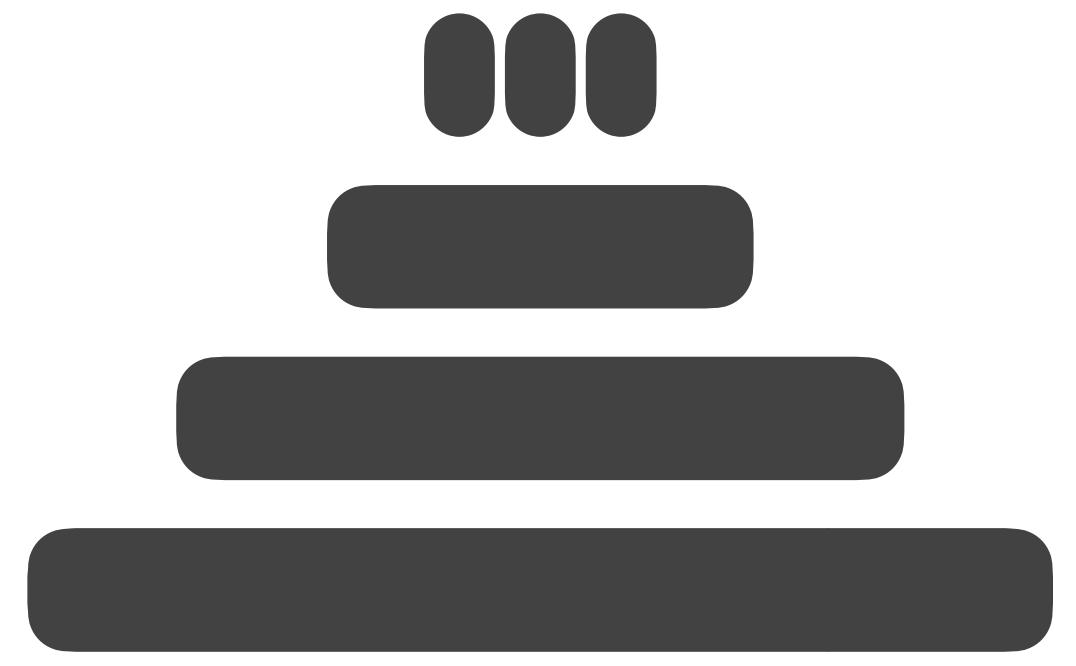


Data Layout

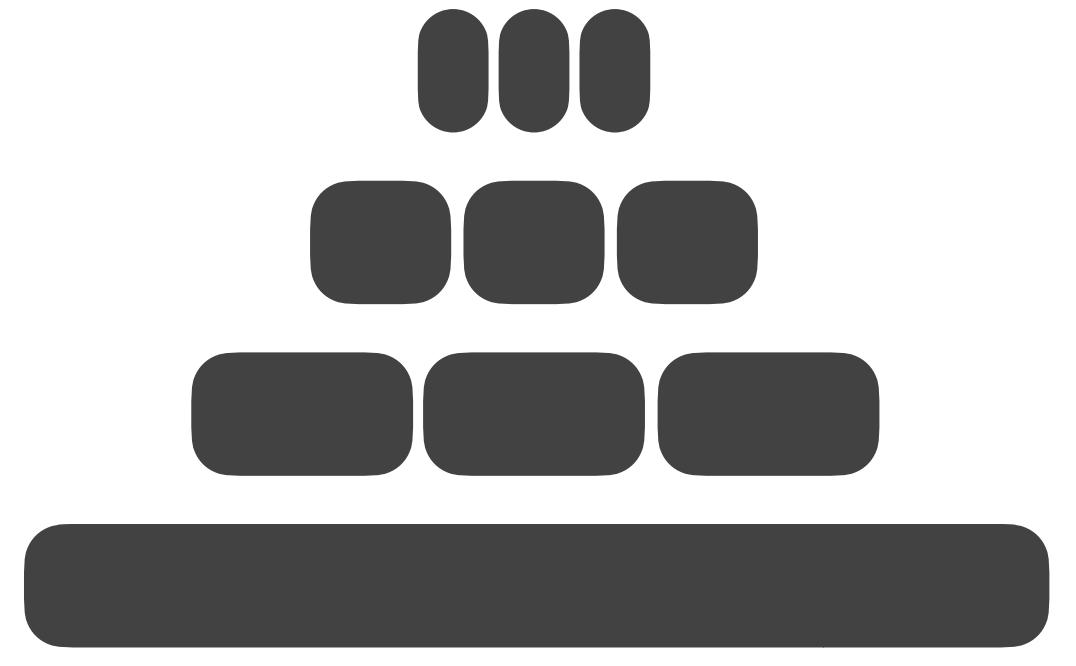
number of runs per level



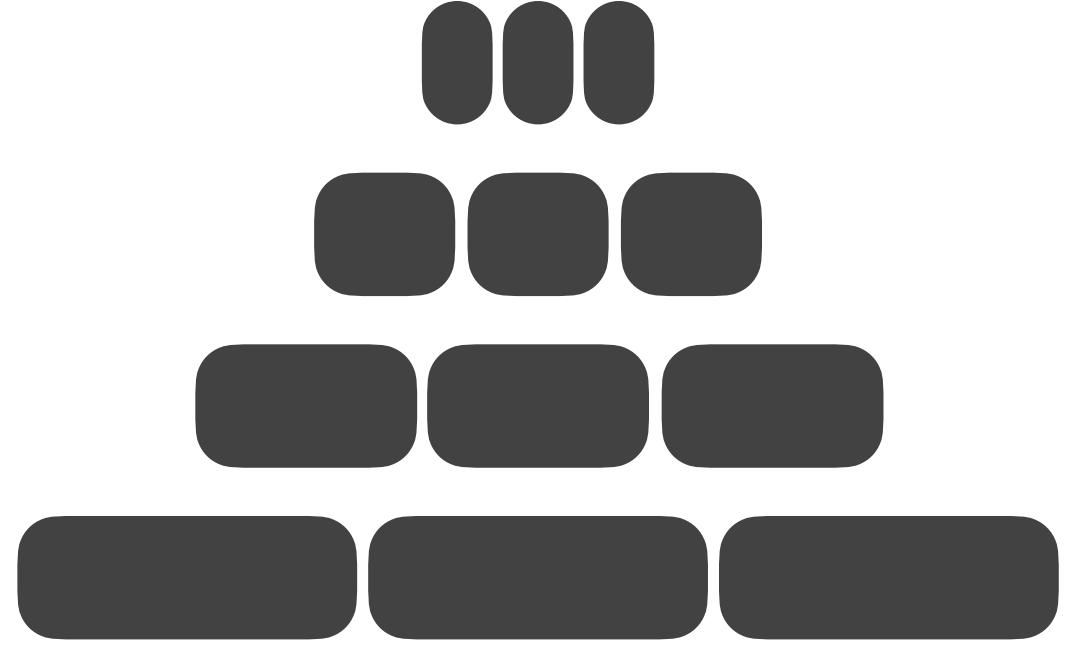
leveling



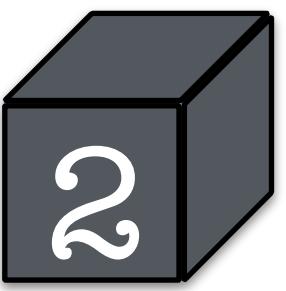
1-leveling



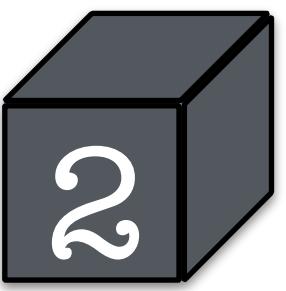
L-leveling



tiering

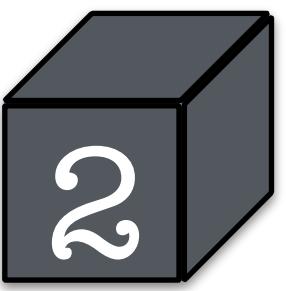


Compaction Granularity



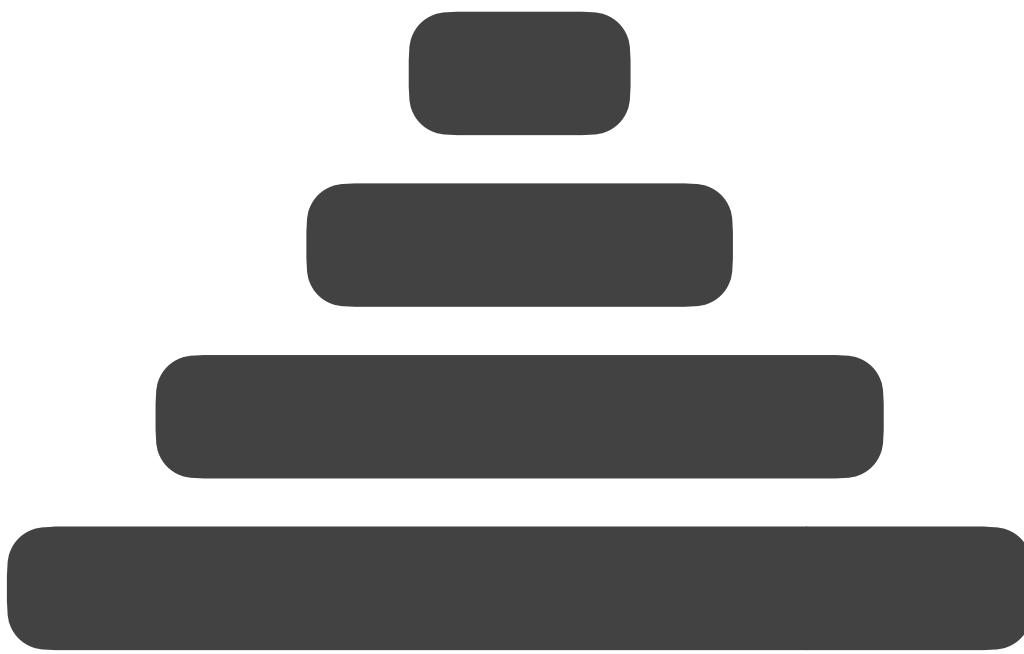
Compaction Granularity

data moved per compaction

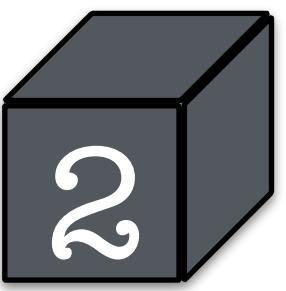


Compaction Granularity

data moved per compaction

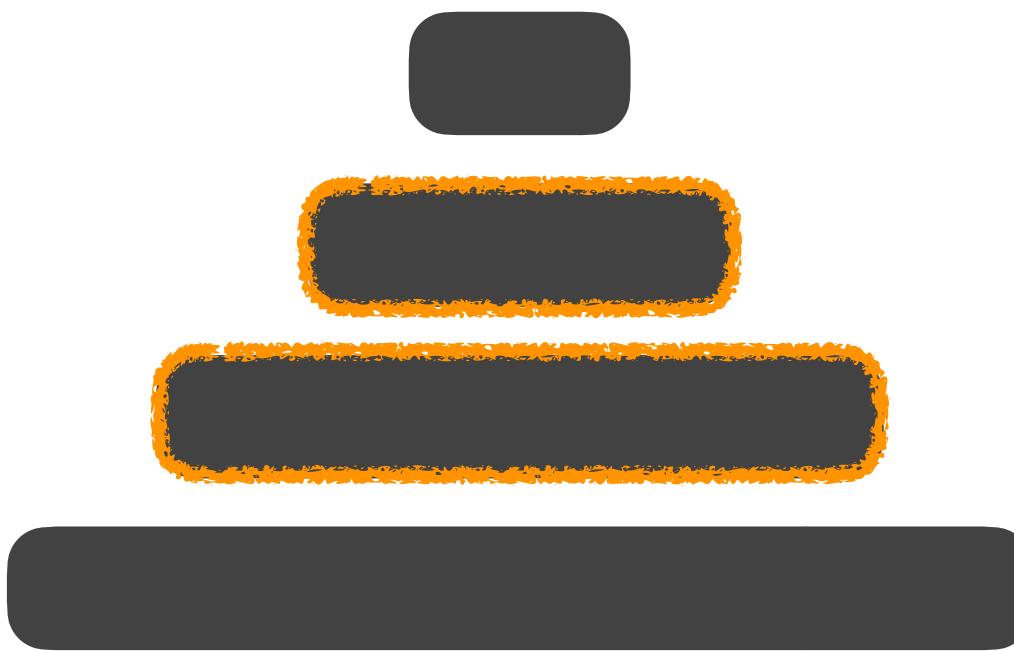


levels

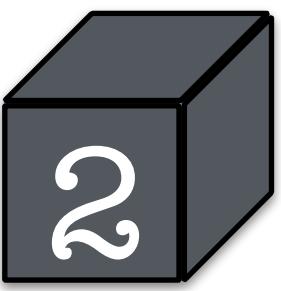


Compaction Granularity

data moved per compaction

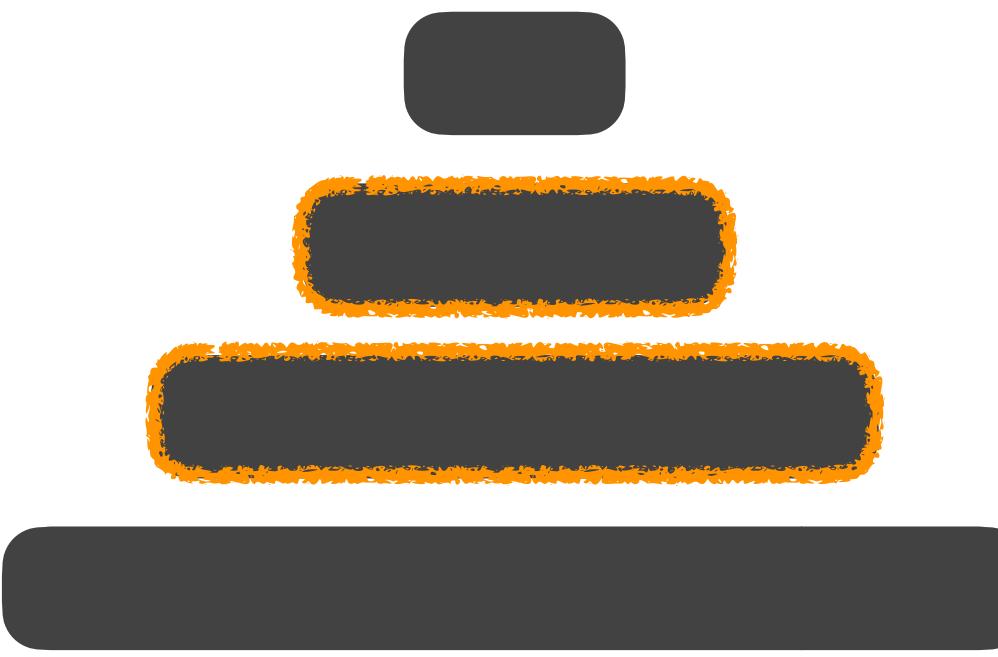


levels

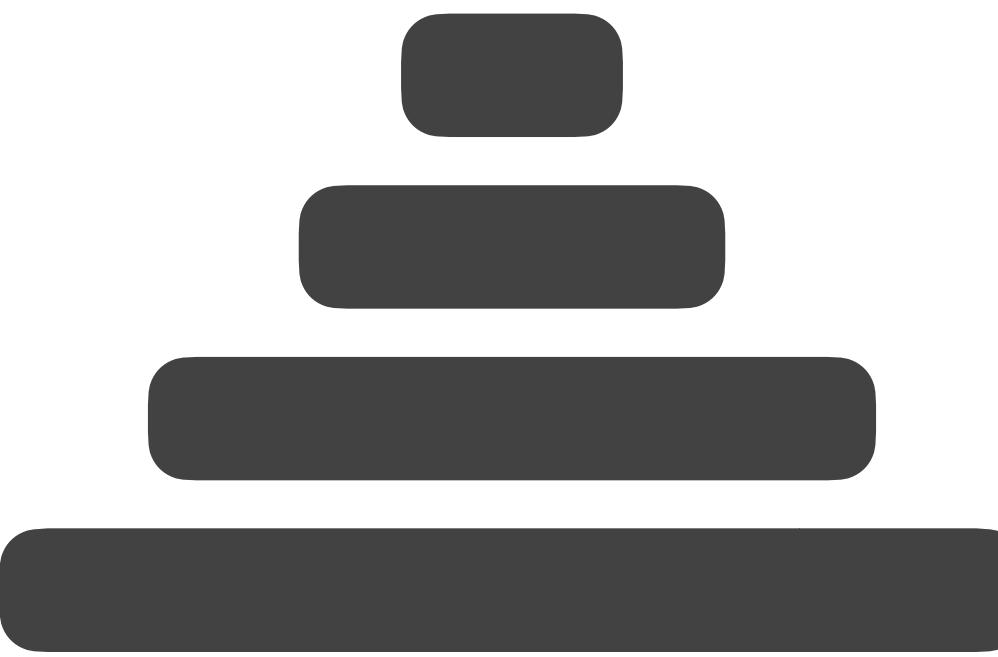


Compaction Granularity

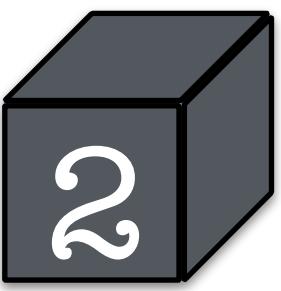
data moved per compaction



levels

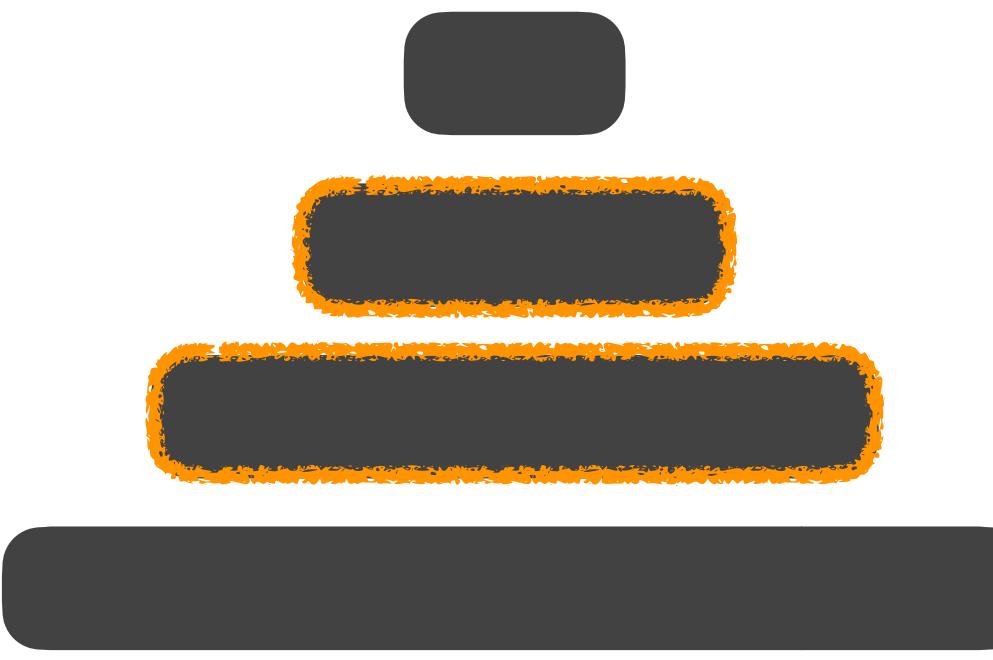


files

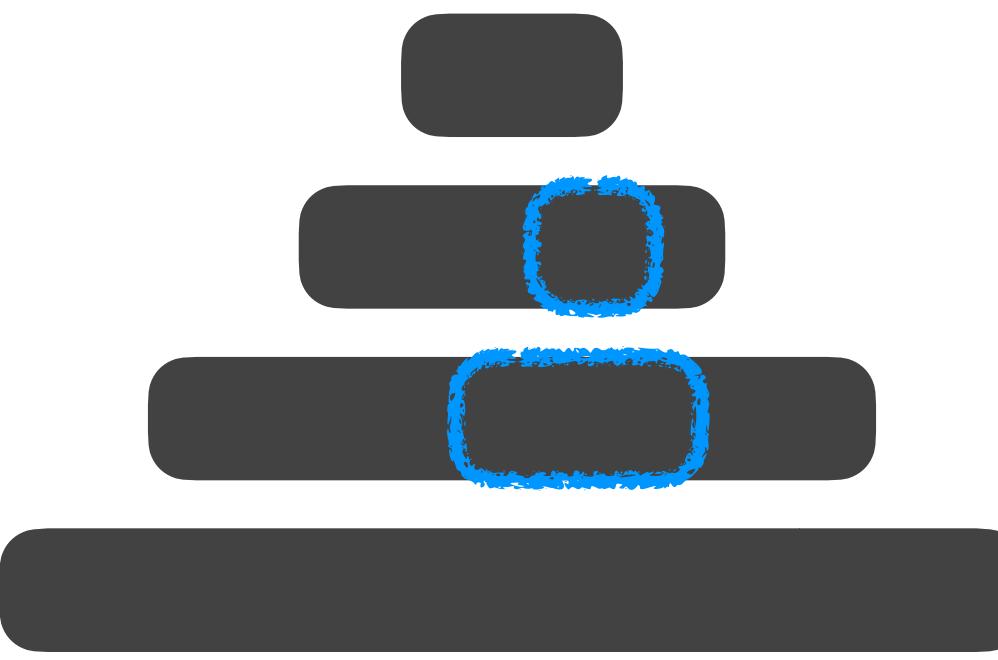


Compaction Granularity

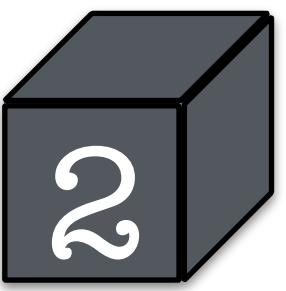
data moved per compaction



levels

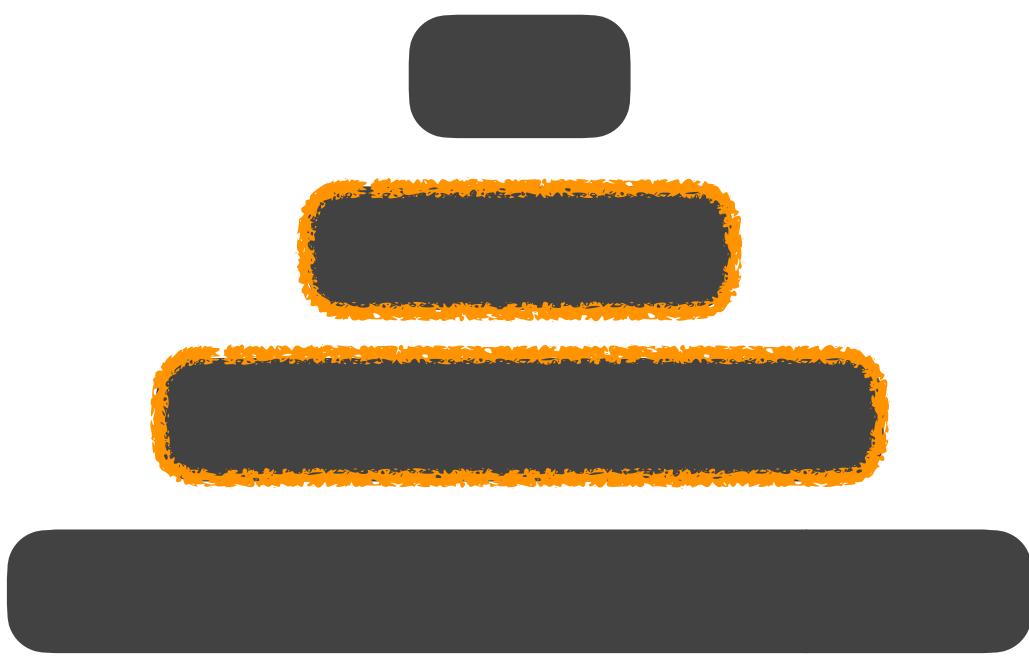


files

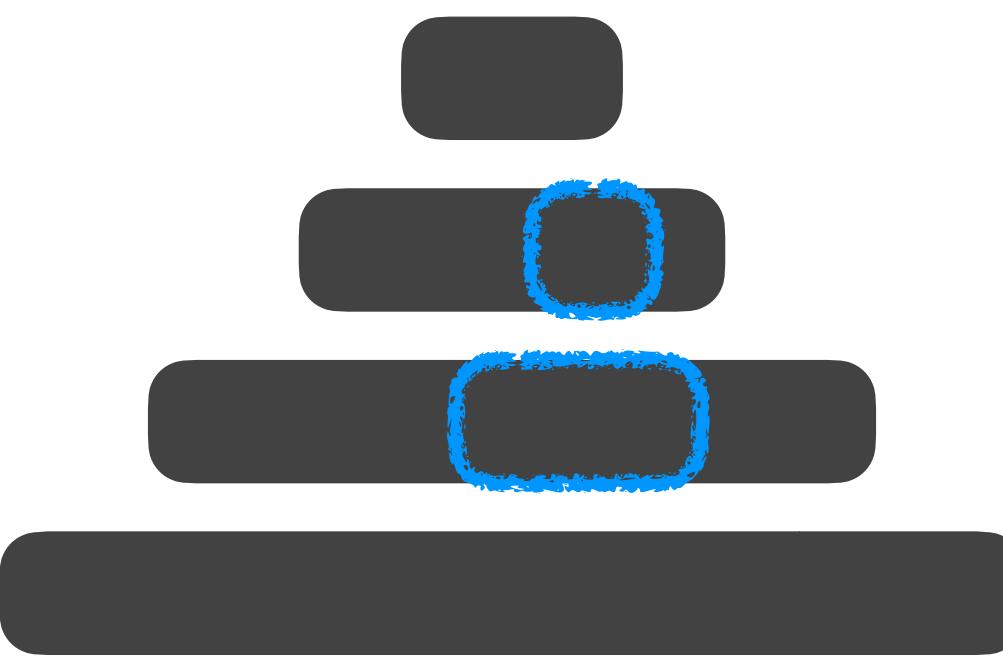


Compaction Granularity

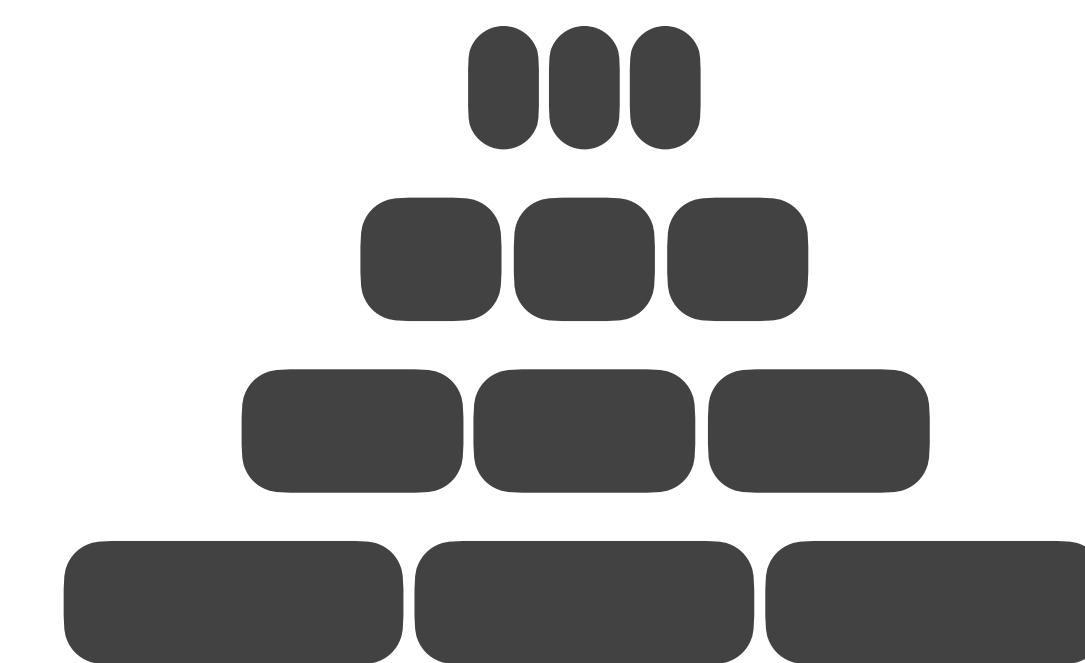
data moved per compaction



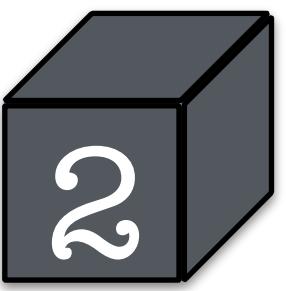
levels



files

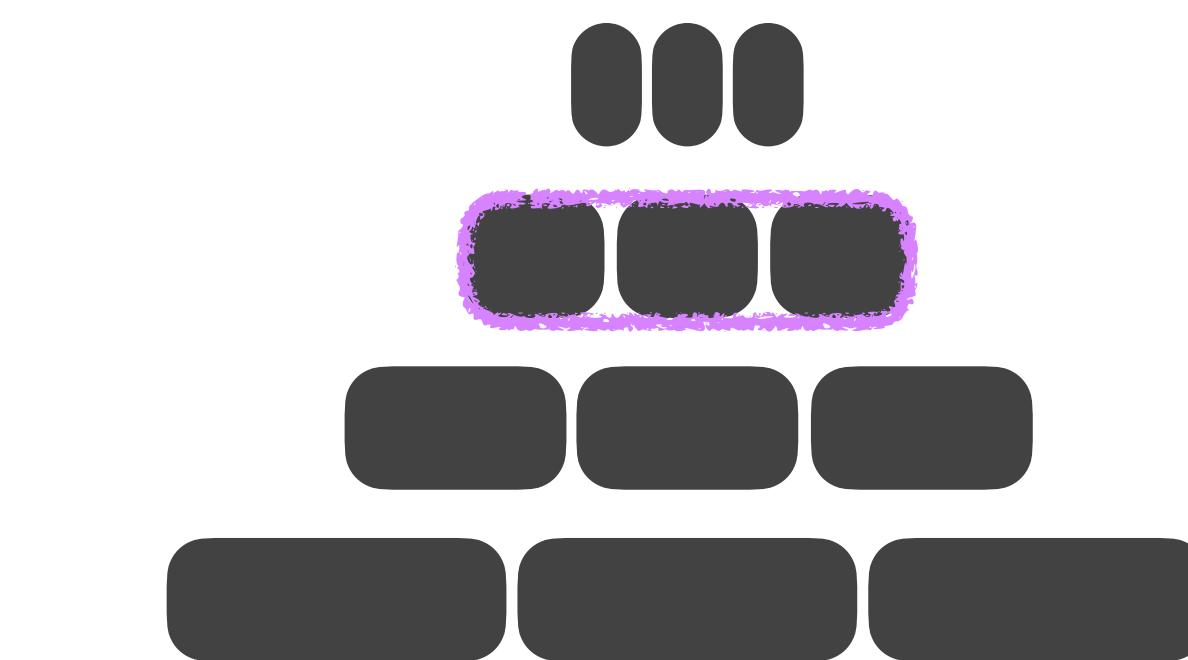
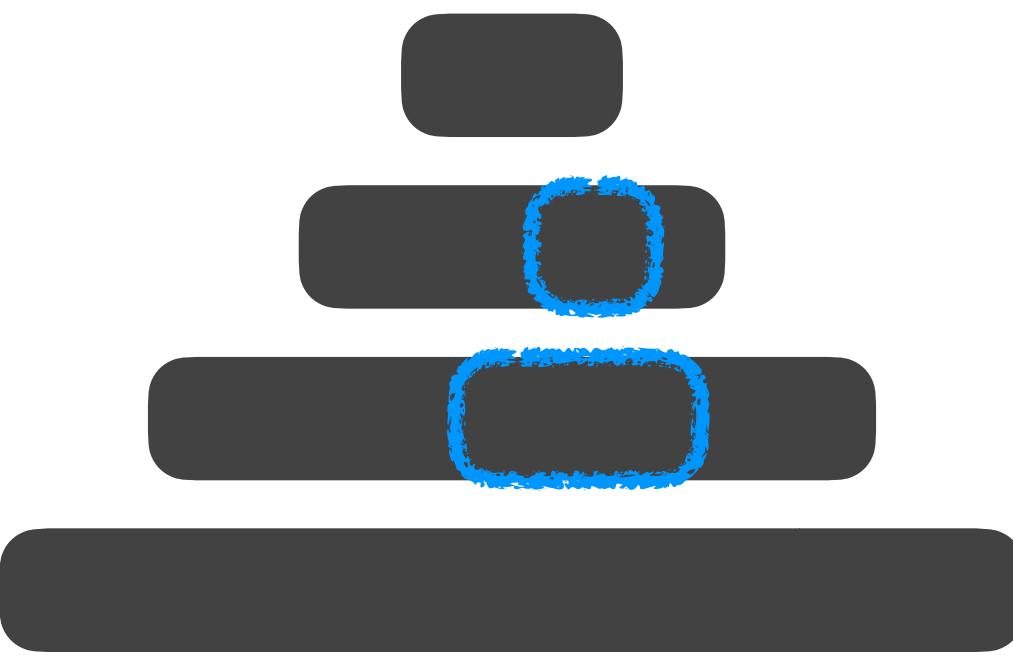
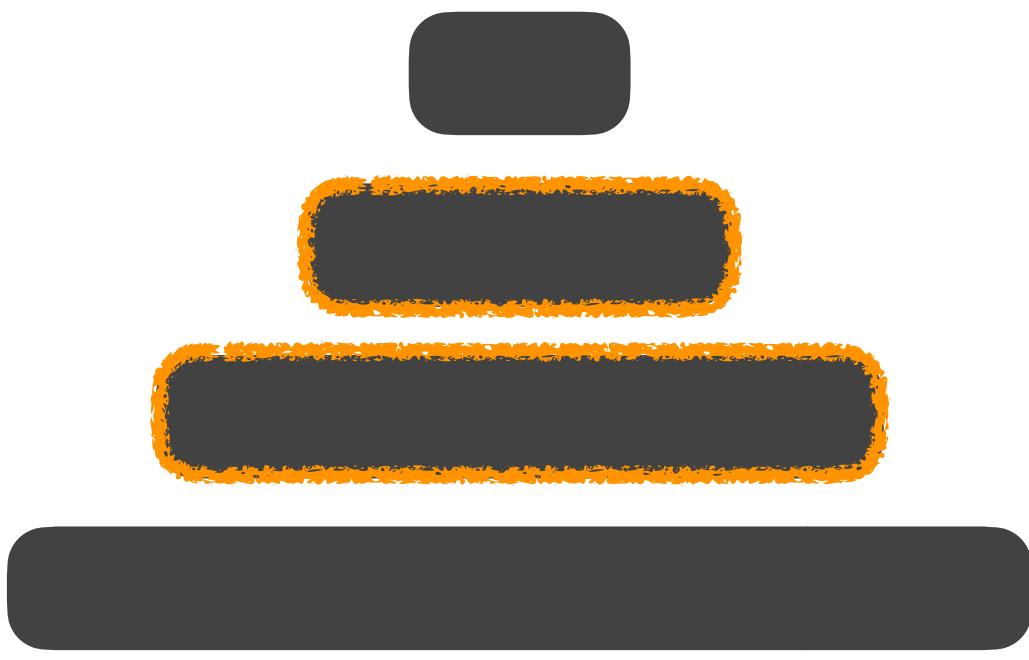


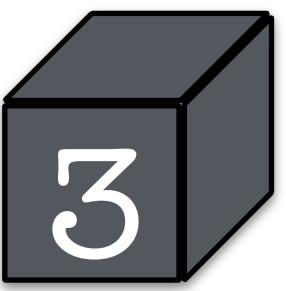
sorted runs in a level



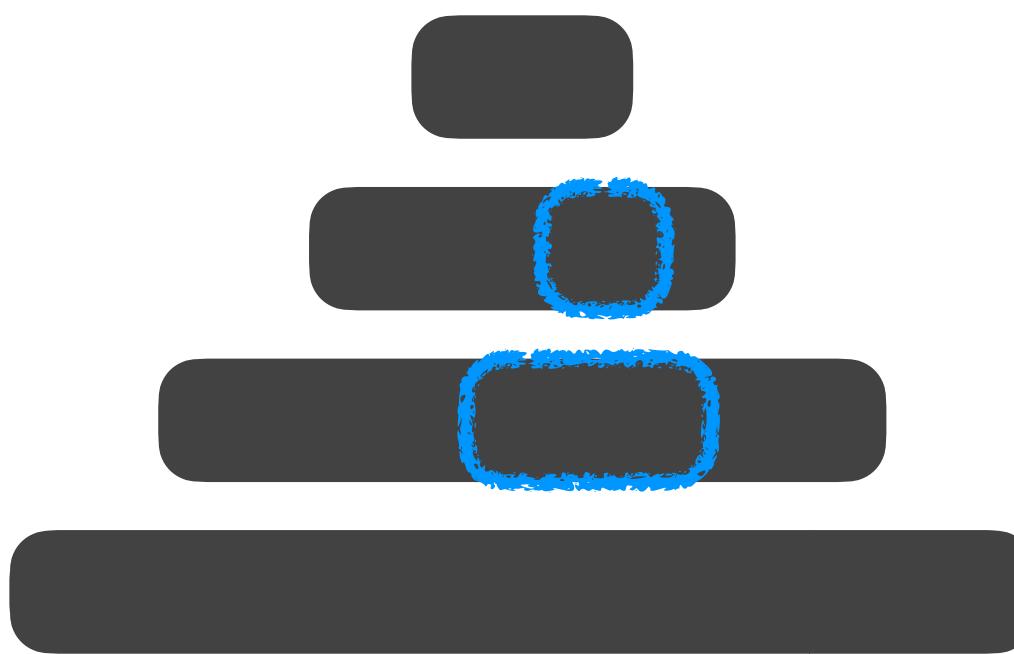
Compaction Granularity

data moved per compaction

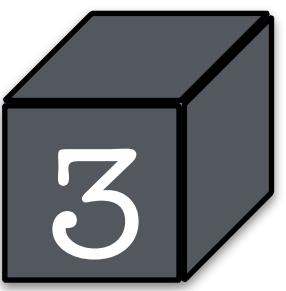




Data Movement Policy

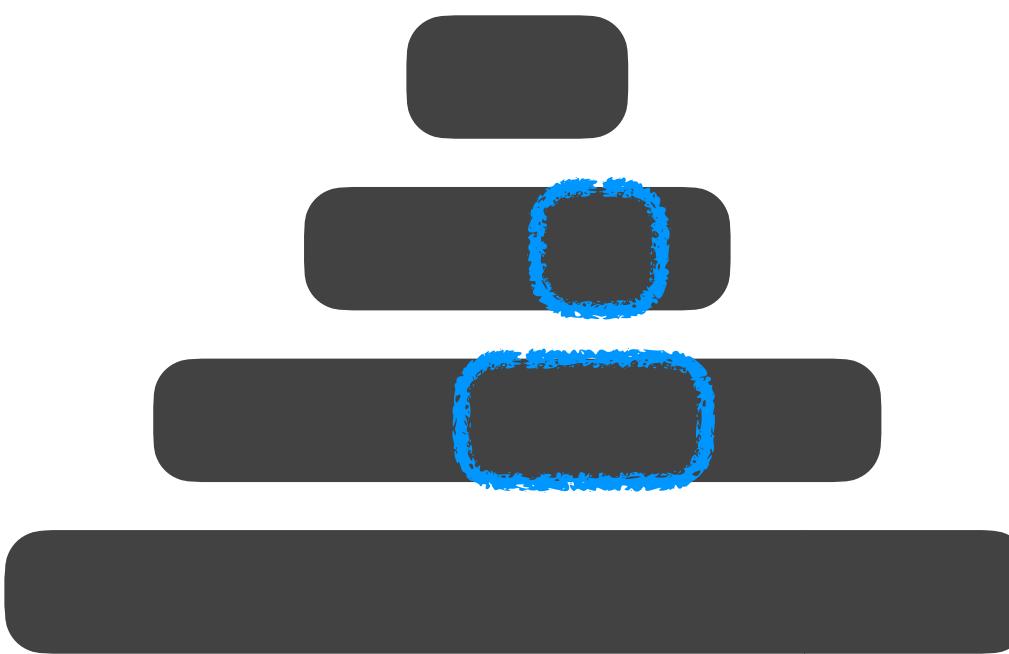


files

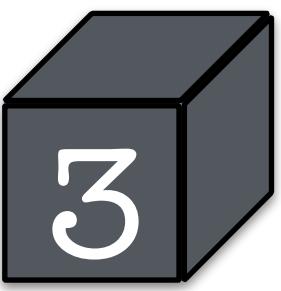


Data Movement Policy

which data to compact



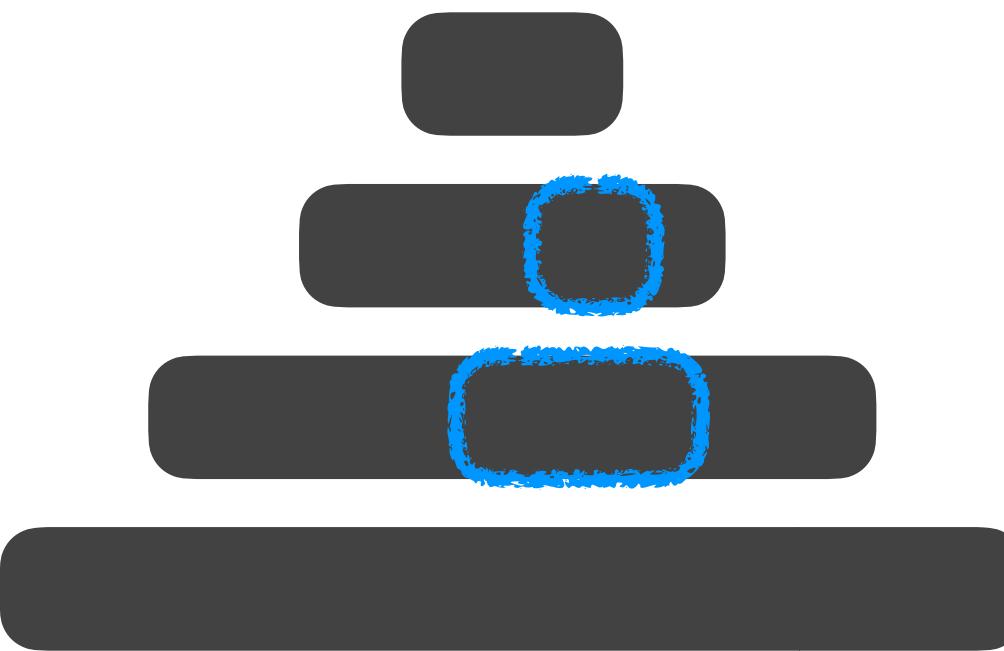
files



Data Movement Policy

which data to compact

round-robin



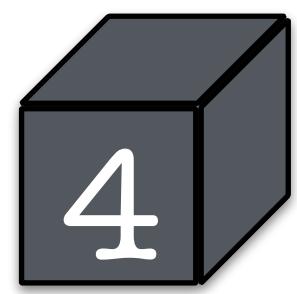
minimum **overlap with parent** level

file with most **tombstones**

coldest file



Compaction Trigger



Compaction Trigger

invoking the compaction routine



Compaction Trigger

invoking the compaction routine

level **saturation**





Compaction Trigger

invoking the compaction routine

level **saturation**





Compaction Trigger

invoking the compaction routine

level **saturation**

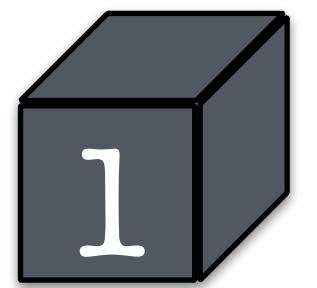


number of **sorted runs**

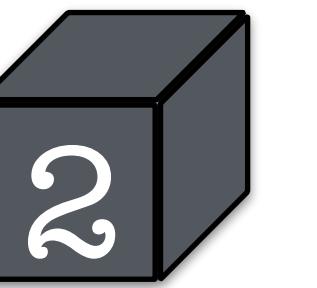
age of a file

space amplification

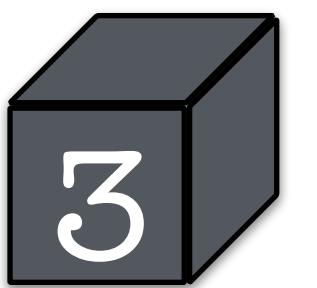
BOSTON
UNIVERSITY



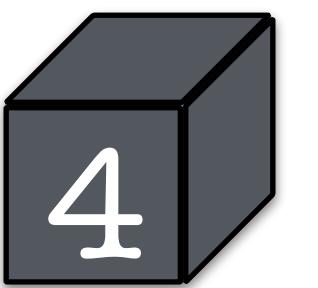
Data Layout



Compaction
Granularity



Data Movement
Policy



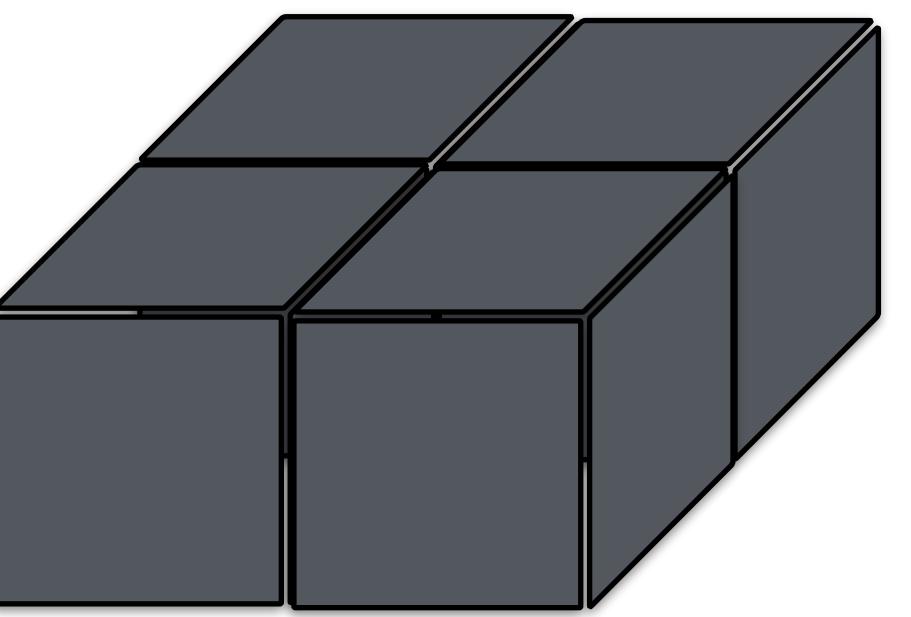
Compaction
Trigger

Data Layout

Compaction
Granularity

Data Movement
Policy

Compaction
Trigger



Any Compaction Algorithm

Database	Data layout	Compaction Trigger				Compaction Granularity		Data Movement Policy								
		Level saturation	#Sorted runs	File staleness	Space amp.	Tombstone-TTL	Level	Sorted run	File (single)	File (multiple)	Round-robin	Least overlap (+1)	Least overlap (+2)	Coldest file	Oldest file	Tombstone density
RocksDB [30], Monkey [22]	Leveling / 1-Leveling	✓	✓					✓	✓		✓	✓	✓	✓	✓	
	Tiering		✓	✓	✓		✓								✓	
LevelDB [32], Monkey (J.) [21]	Leveling	✓						✓			✓	✓	✓			
SlimDB [47]	Tiering	✓						✓	✓							✓
Dostoevsky [23]	L-leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T			✓ ^L				✓ ^T	
LSM-Bush [24]	Hybrid leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T			✓ ^L				✓ ^T	
Lethe [51]	Leveling	✓		✓				✓	✓		✓				✓	
Silk [11], Silk+ [12]	Leveling	✓						✓	✓		✓					
HyperLevelDB [35]	Leveling	✓						✓			✓	✓	✓			
PebblesDB [46]	Hybrid leveling	✓						✓	✓						✓	
Cassandra [8]	Tiering		✓	✓	✓	✓		✓							✓	
	Leveling	✓		✓				✓	✓		✓			✓	✓	
WiredTiger [62]	Leveling	✓				✓									✓	
X-Engine [34], Leaper [63]	Hybrid leveling	✓						✓	✓		✓				✓	
HBase [7]	Tiering		✓				✓								✓	
AsterixDB [3]	Leveling	✓			✓										✓	
	Tiering	✓			✓										✓	

Blueprint for Experiments

Blueprint for **Experiments**

10 compaction strategies

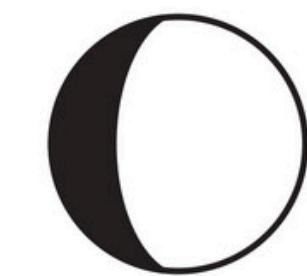
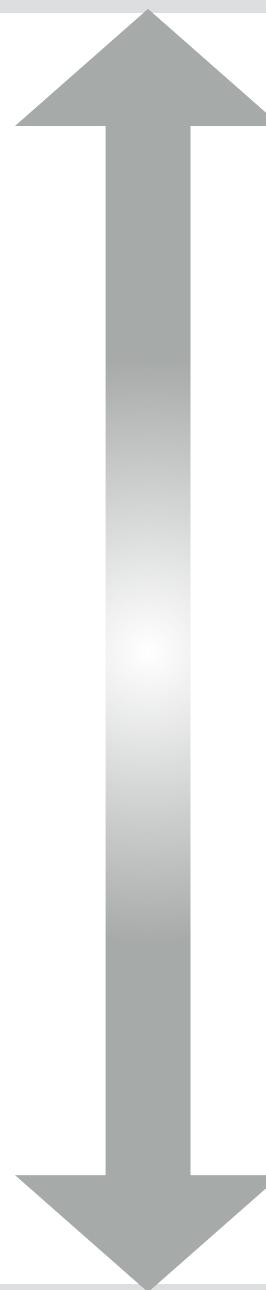
Blueprint for **Experiments**

10 compaction strategies

612 metrics

Blueprint for **Experiments**

10 compaction strategies

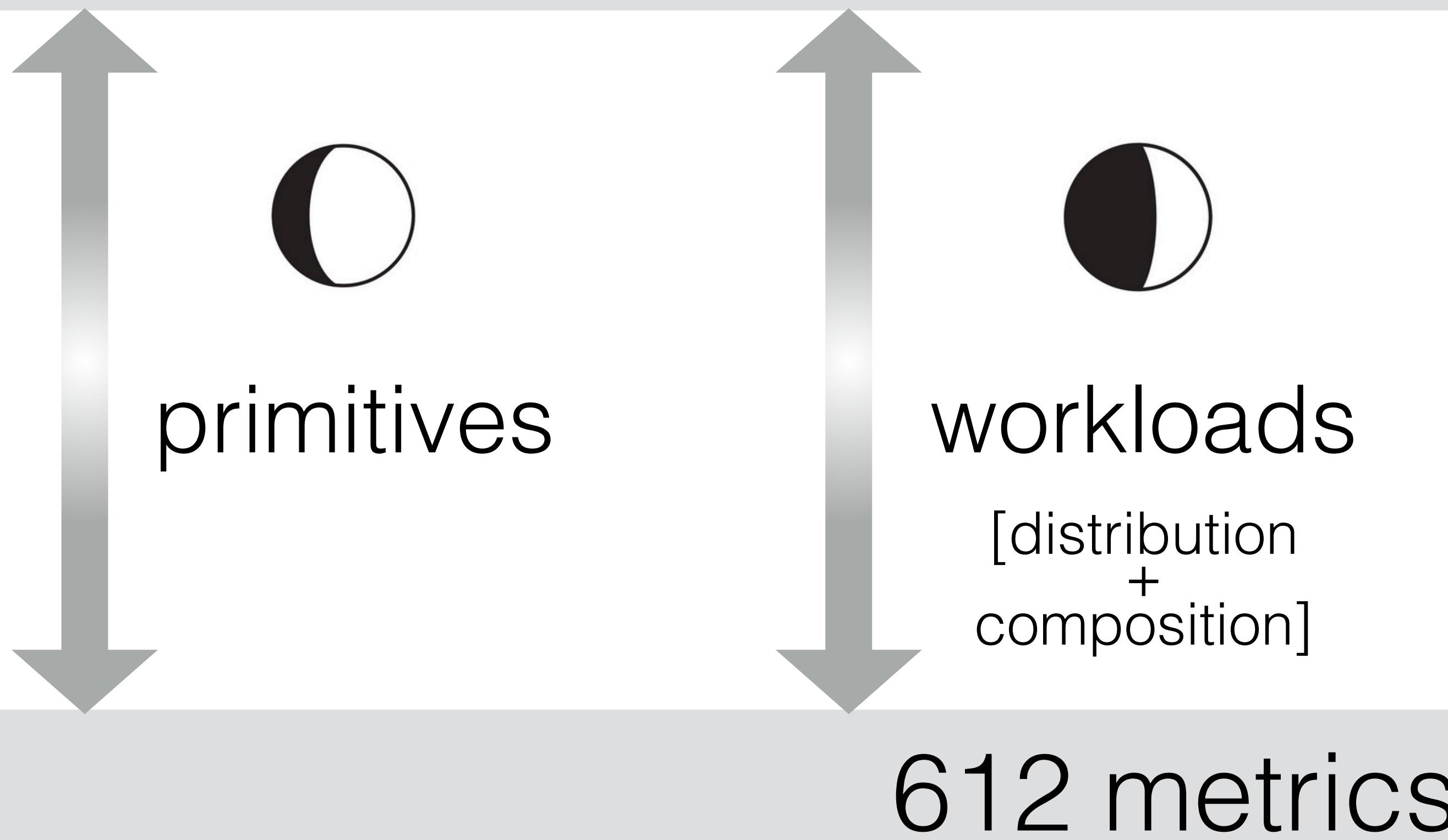


primitives

612 metrics

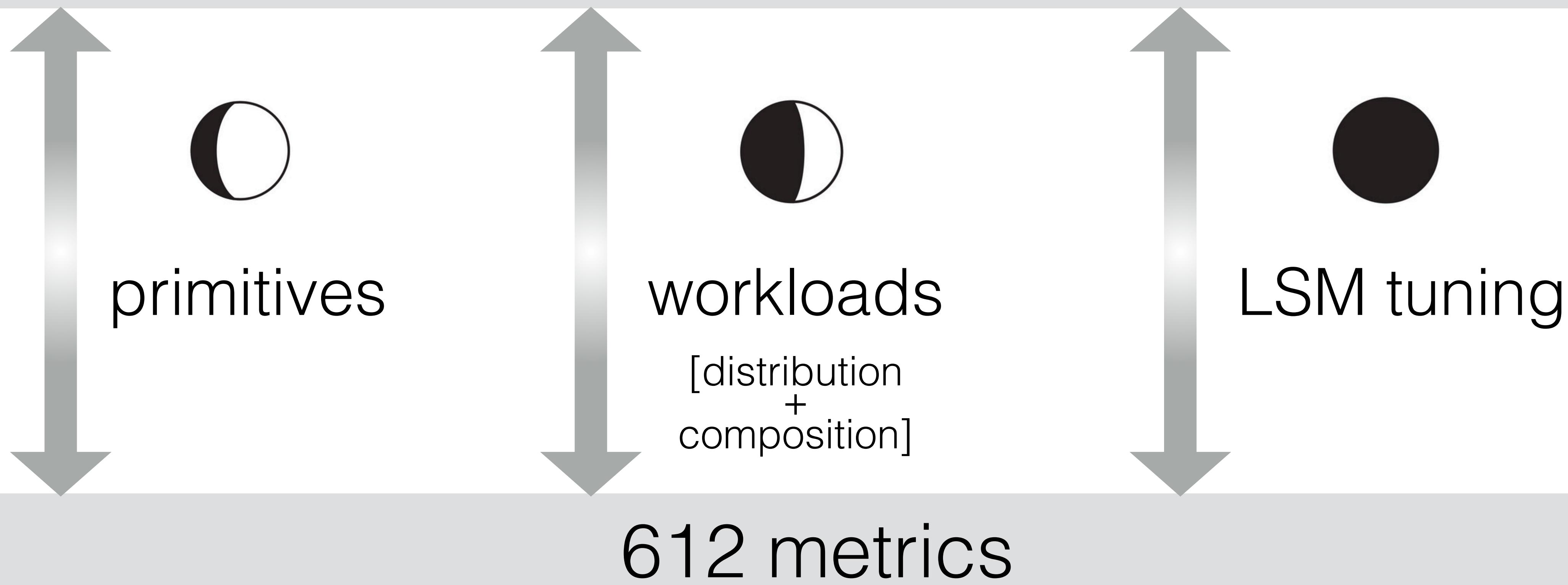
Blueprint for Experiments

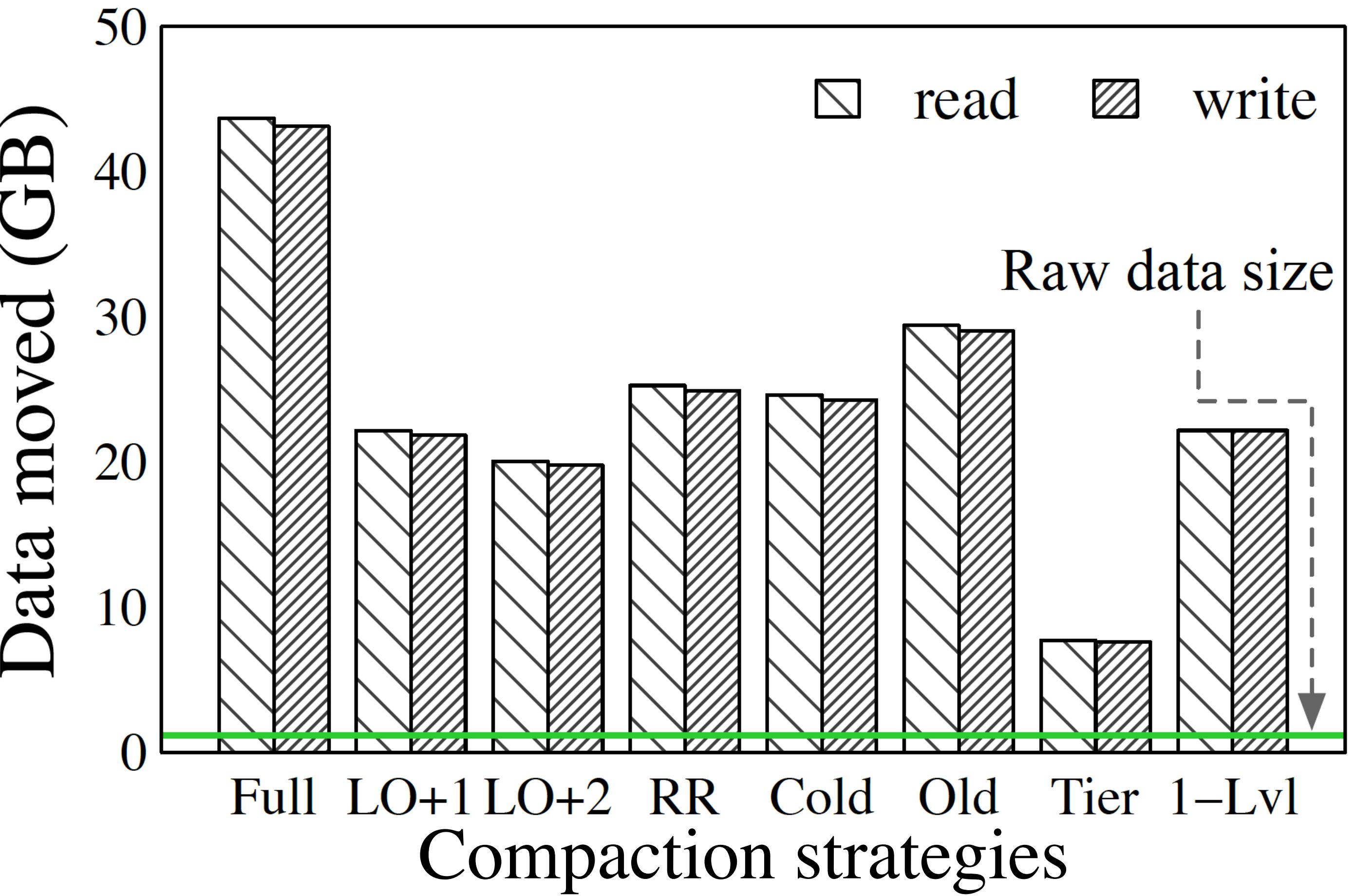
10 compaction strategies



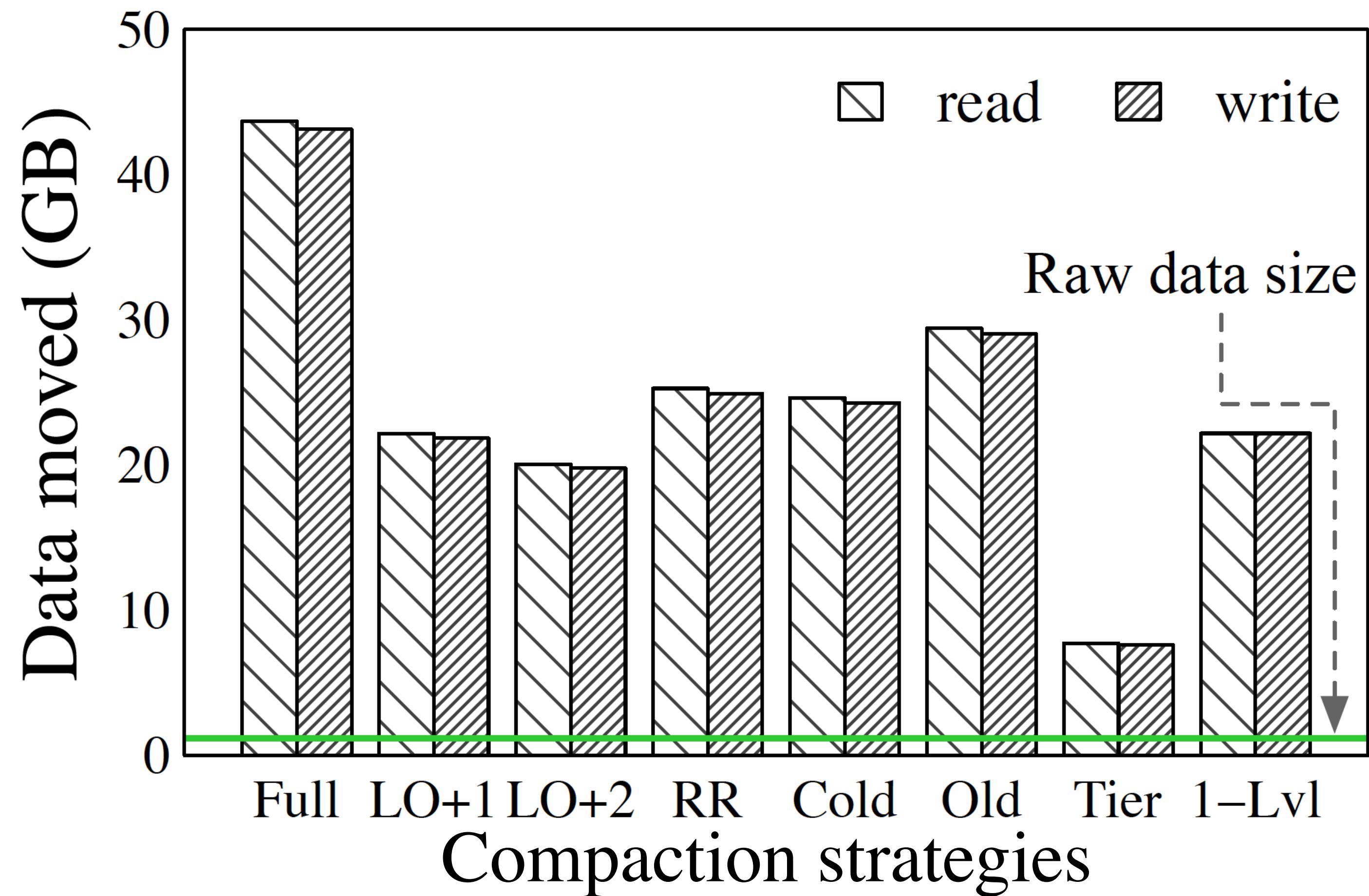
Blueprint for Experiments

10 compaction strategies

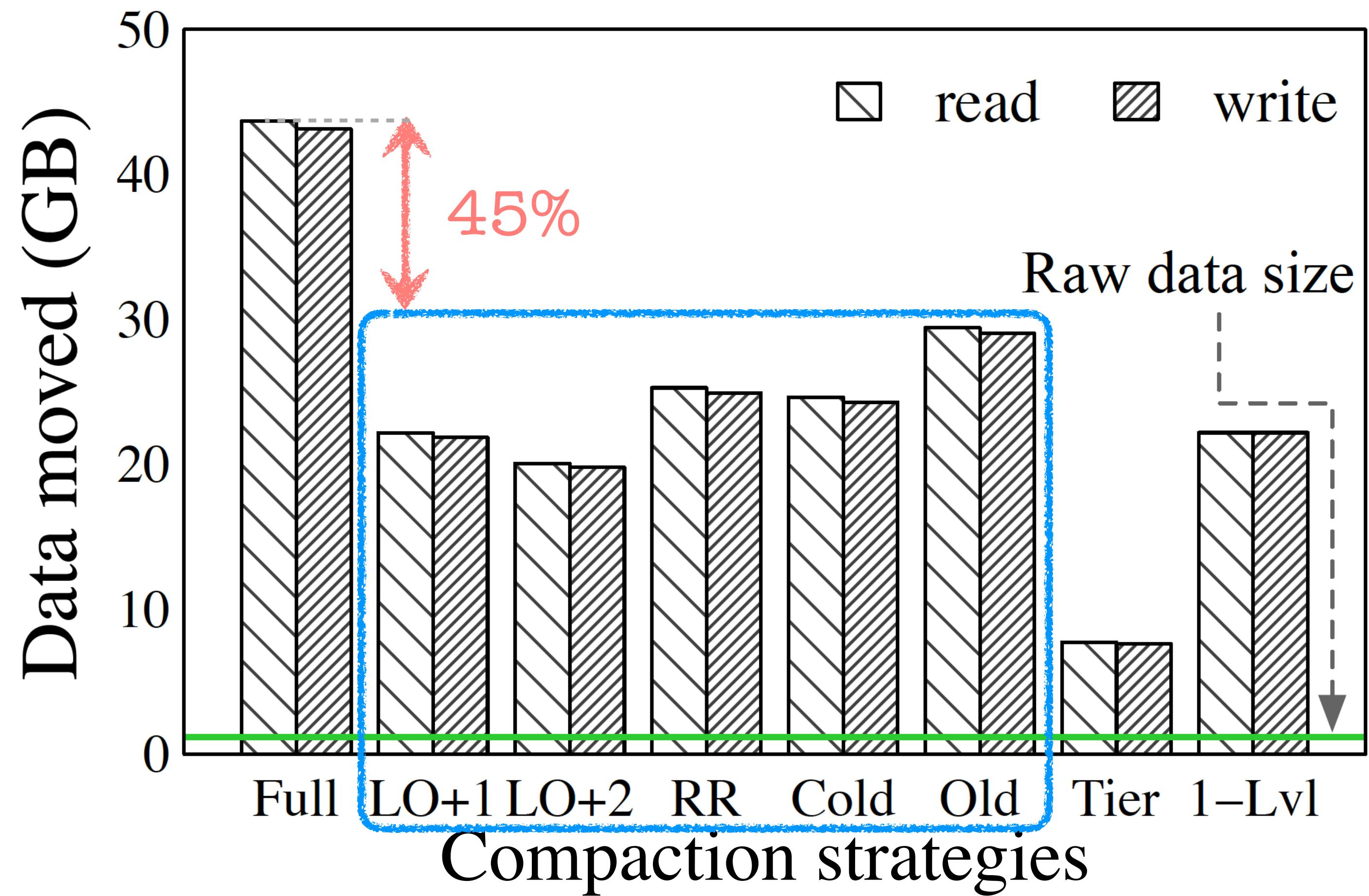




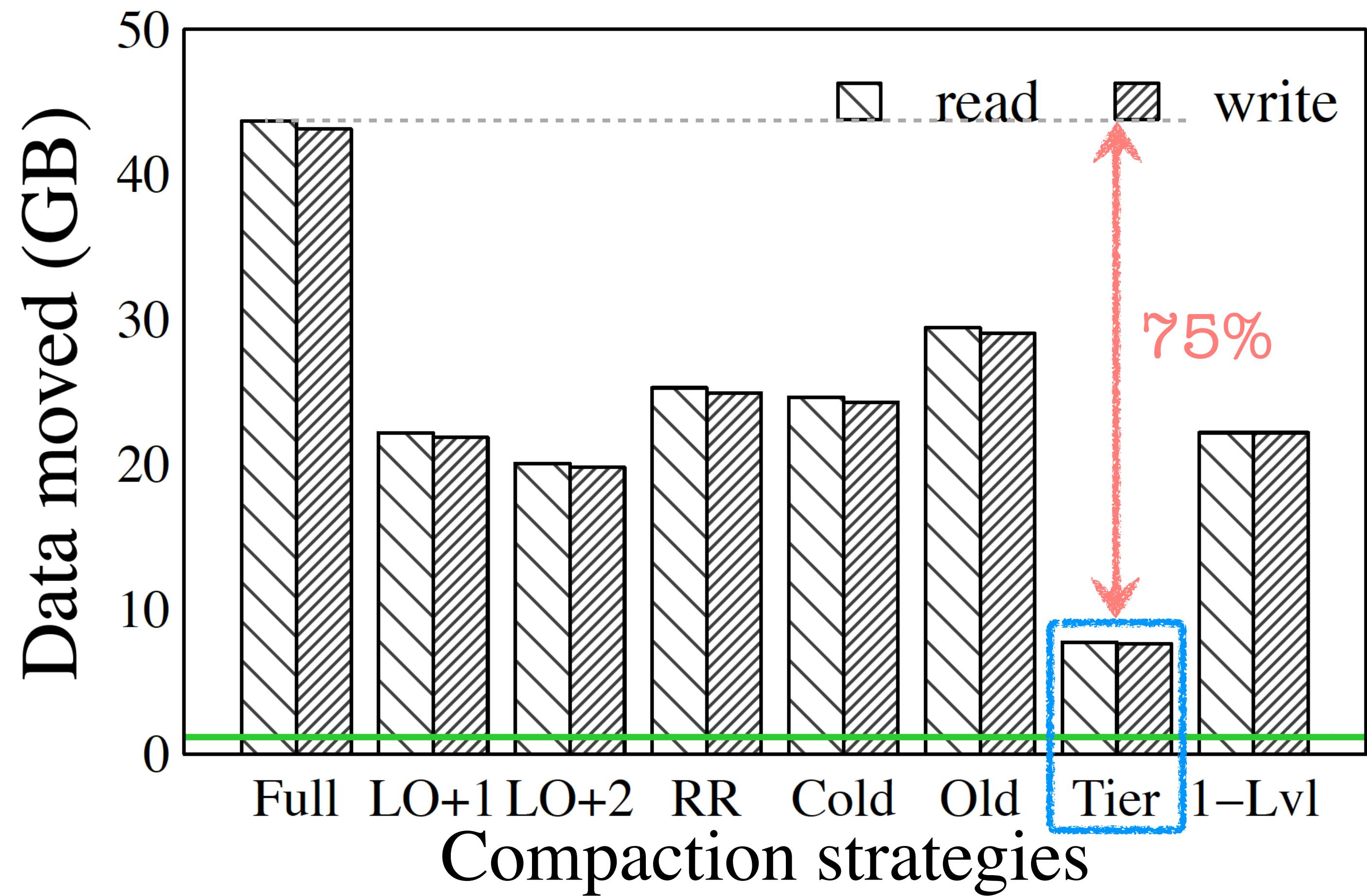
Compacting data at smaller granularity reduces data movement.



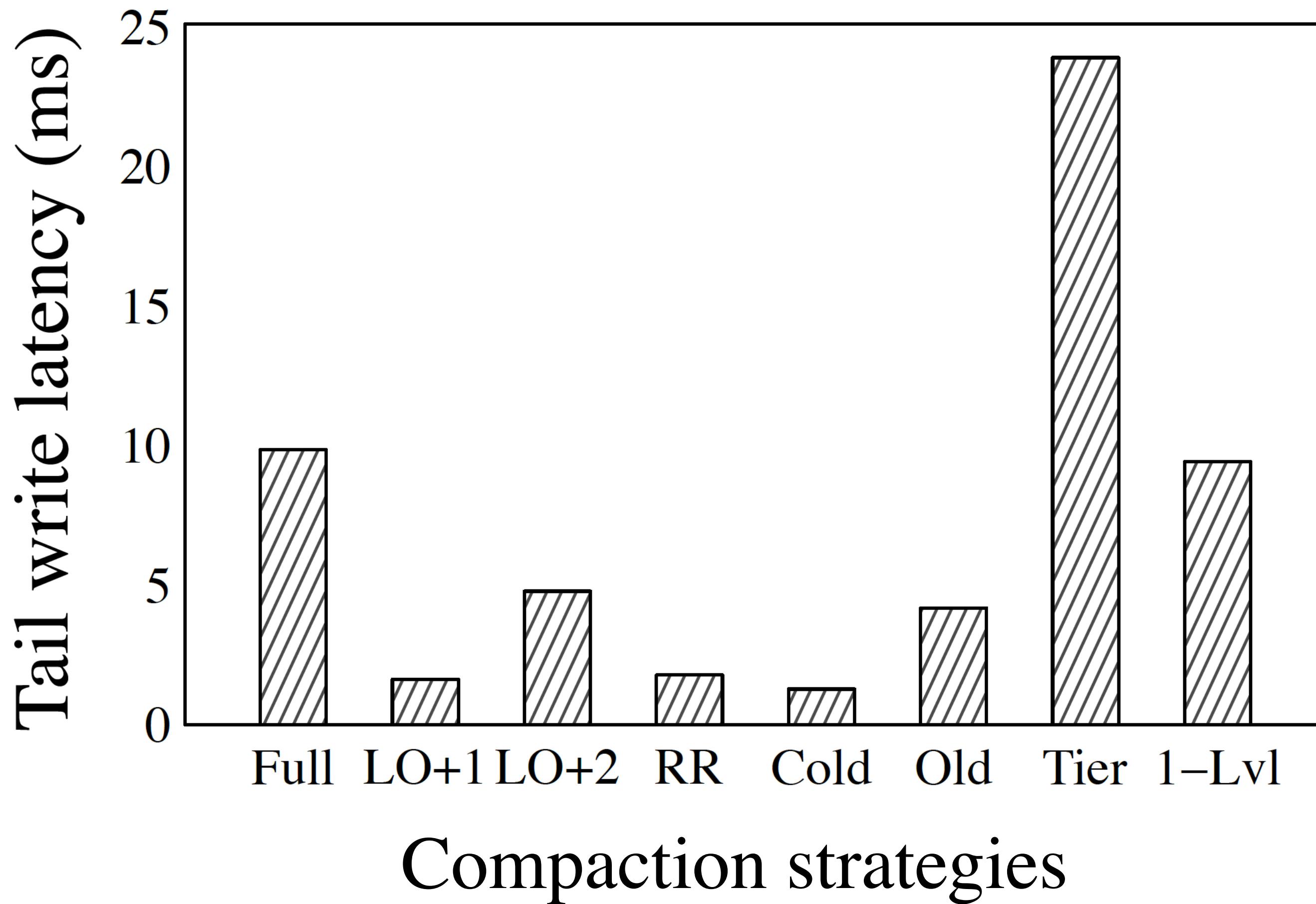
Compacting data at smaller granularity reduces data movement.



Compacting data at smaller granularity reduces data movement.



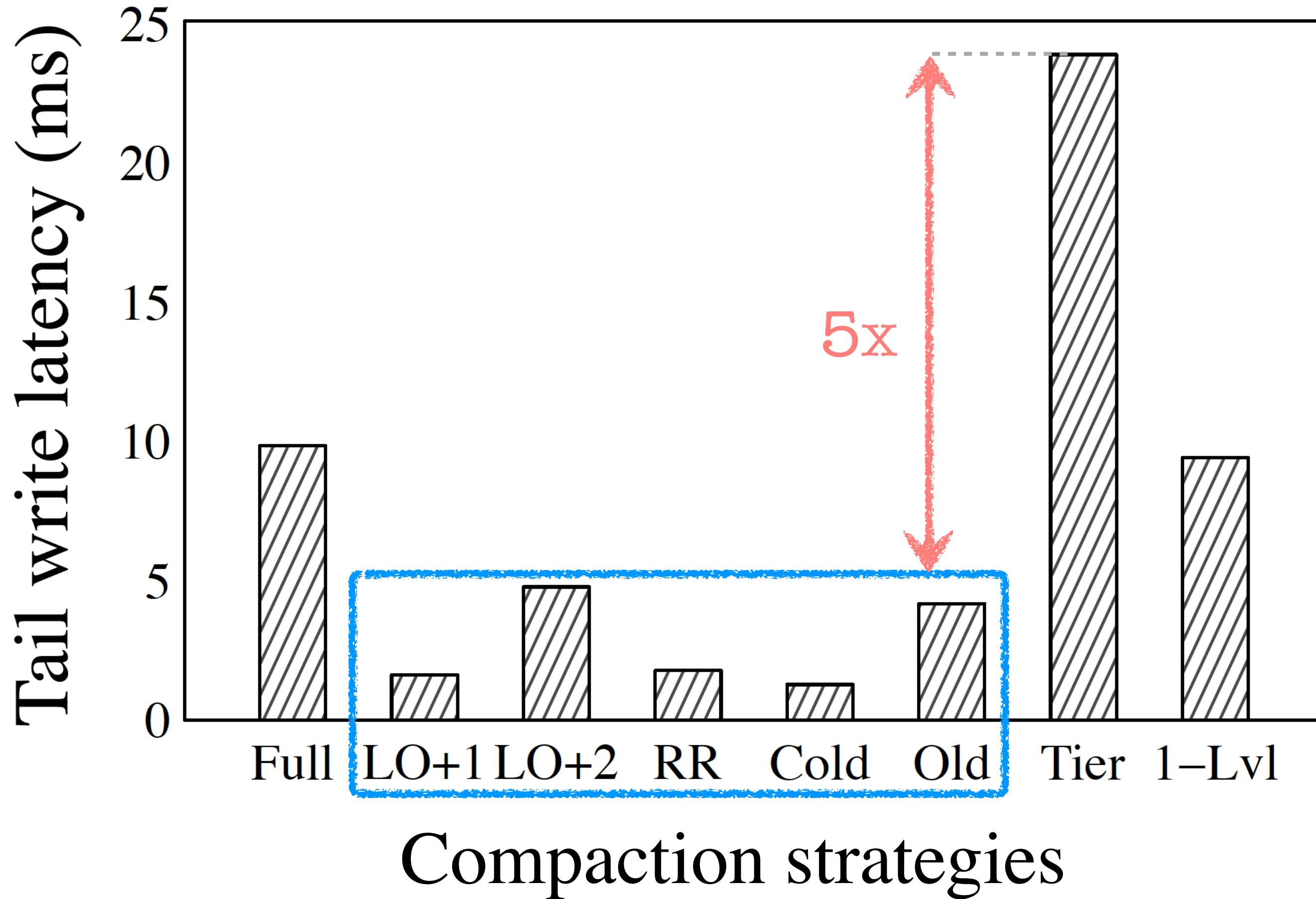
Compacting data at smaller granularity reduces data movement.



Compacting data at smaller granularity reduces data movement.



Tiered data layout has the highest write throughput but also the highest tail write latency.



Compacting data at smaller granularity reduces data movement.



Tiered data layout has the highest write throughput but also the highest tail write latency.



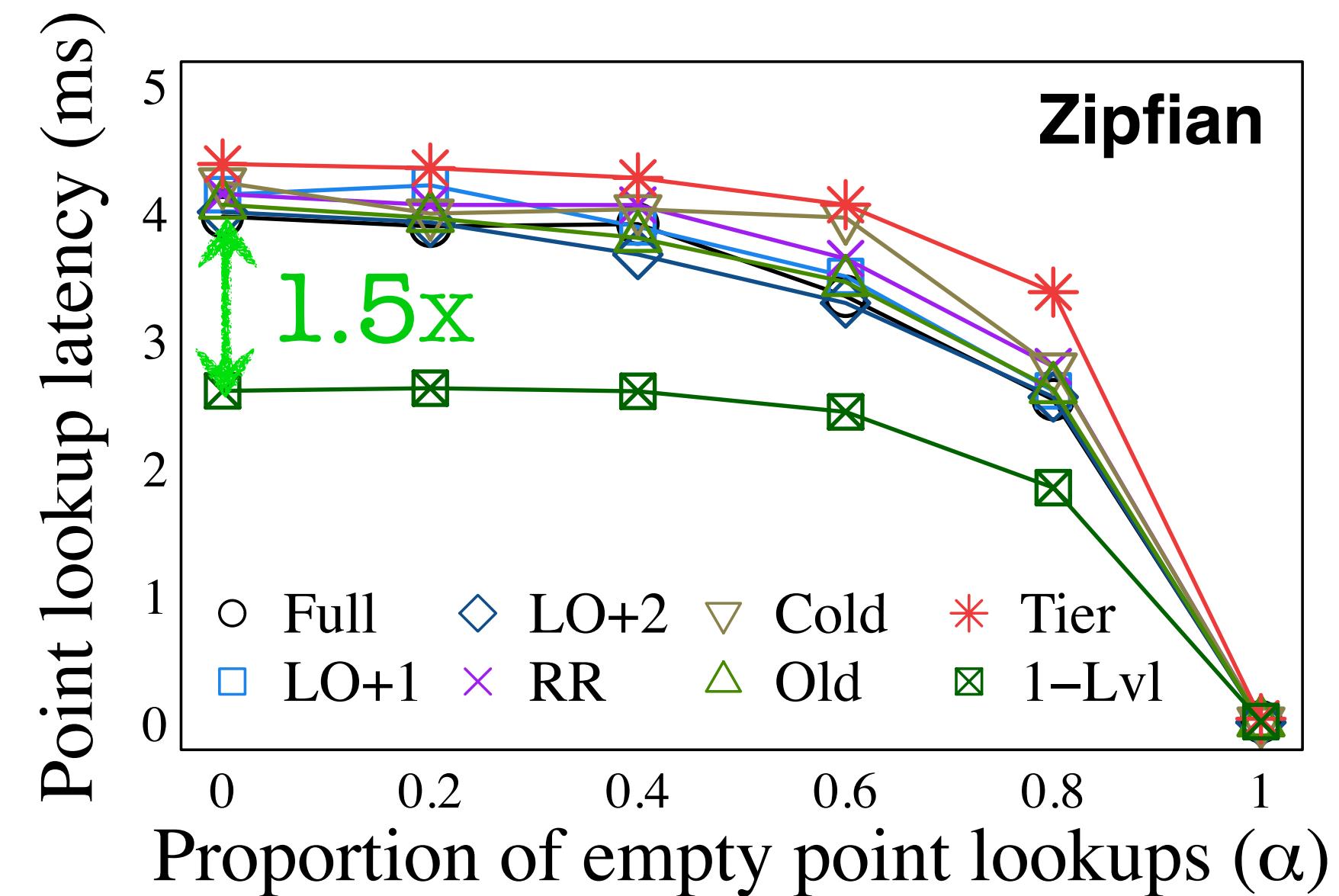
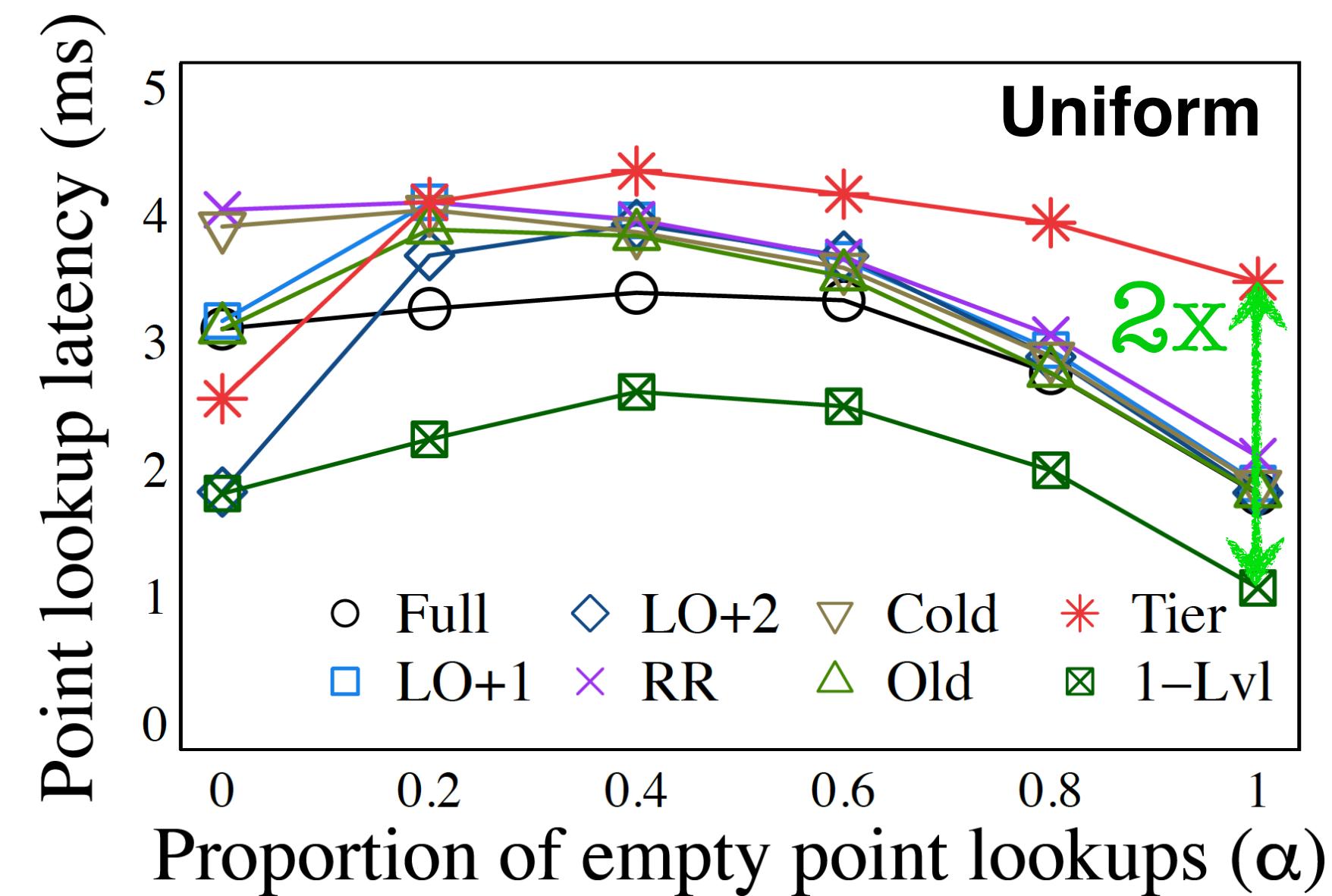
Compacting data at smaller granularity reduces data movement.



Tiered data layout has the highest write throughput but also the highest tail write latency.



Hybrid data layouts dominate point lookup performance.



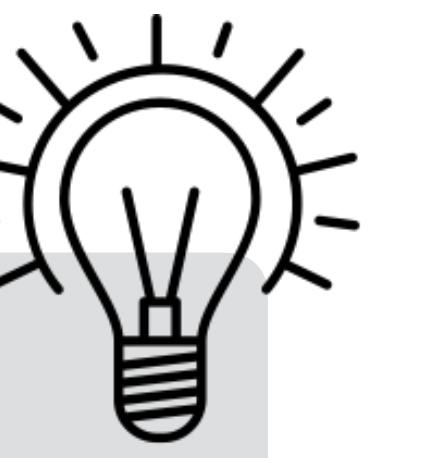
Compacting data at smaller granularity reduces data movement.



Tiered data layout has the highest write throughput but also the highest tail write latency.



Hybrid data layouts dominate point lookup performance.



Compacting data at smaller granularity reduces data movement.



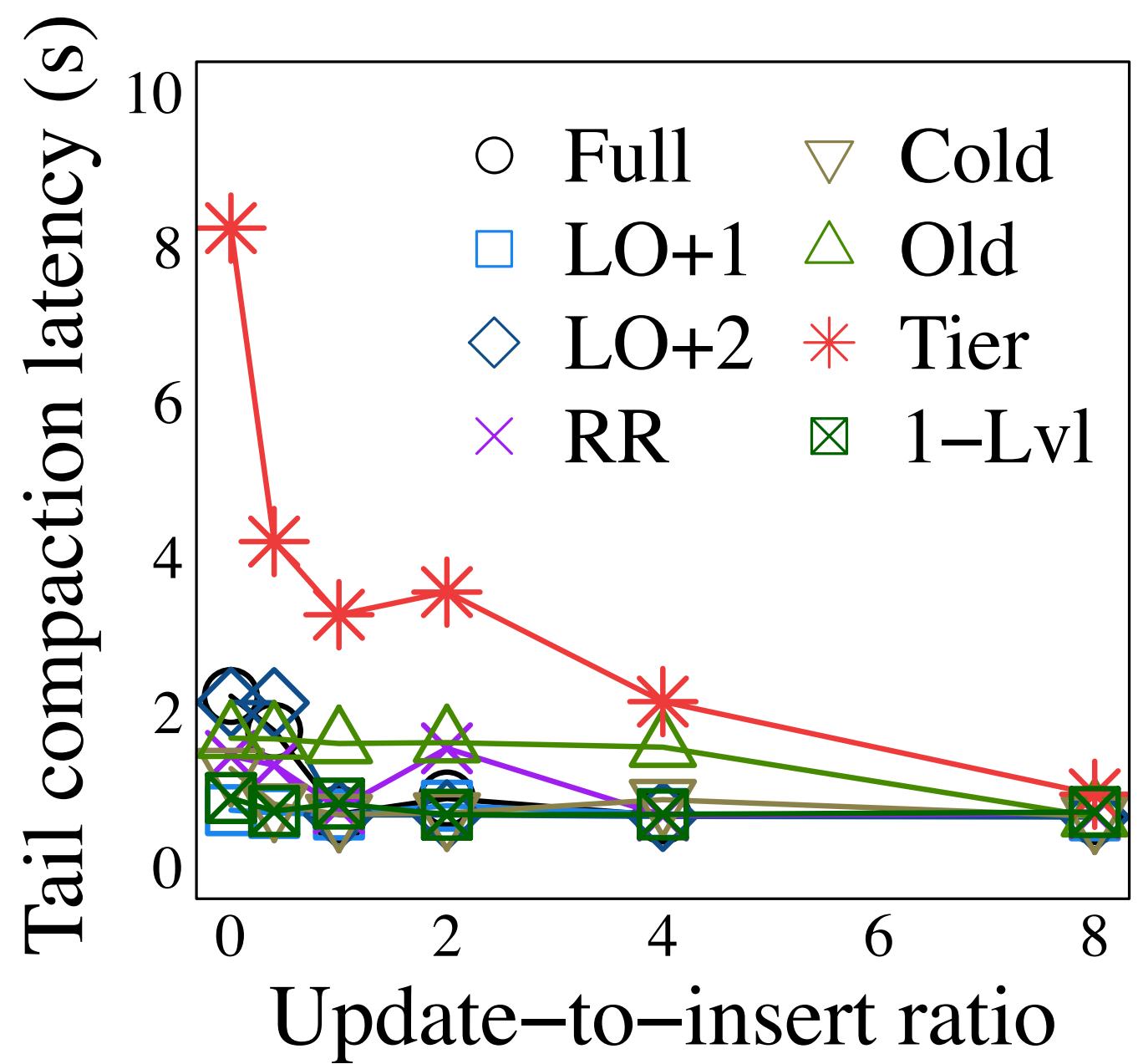
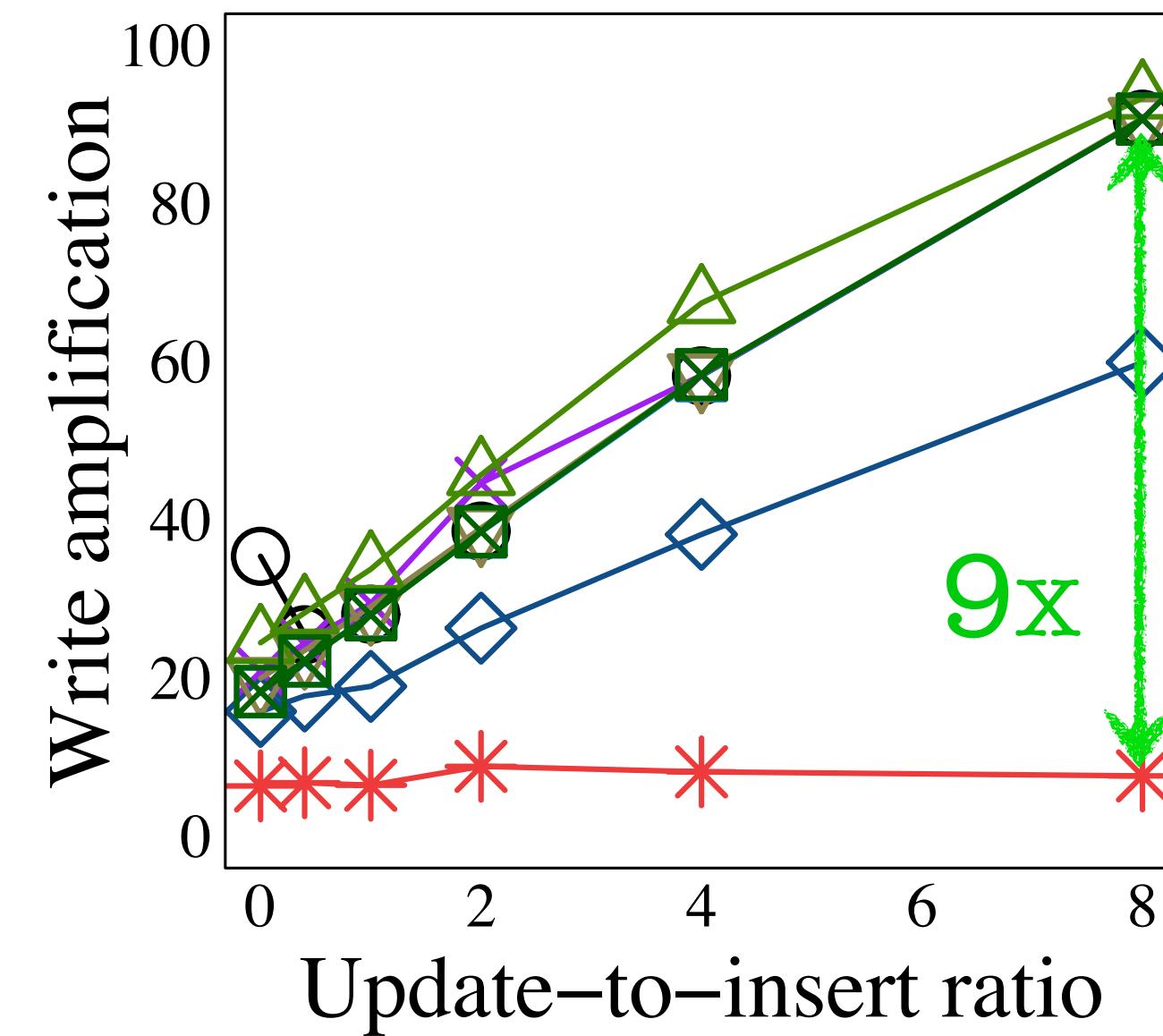
For update-intensive workloads, tiering dominates the performance space.



Tiered data layout has the highest write throughput but also the highest tail write latency.



Hybrid data layouts dominate point lookup performance.





Compacting data at smaller granularity reduces data movement.



For update-intensive workloads, tiering dominates the performance space.



Tiered data layout has the highest write throughput but also the highest tail write latency.



Hybrid data layouts dominate point lookup performance.



Compacting data at smaller granularity reduces data movement.



For update-intensive workloads, tiering dominates the performance space.



Tiered data layout has the highest write throughput but also the highest tail write latency.



The relative benefits of compaction strategies are marginally affected by LSM-tuning.



Hybrid data layouts dominate point lookup performance.

Summary

Summary

Compaction is **key to LSM-performance**.

Summary

Compaction is **key to LSM-performance**.

Compaction as first-order **design primitives**.

Summary

Compaction is **key to LSM-performance**.

Compaction as first-order **design primitives**.

Guidelines to design and tuning through experiments.

Summary

Compaction is **key to LSM-performance**.

Compaction as first-order **design primitives**.

Guidelines to design and tuning through experiments.

Thank You!