

Learning to Optimize LSM-trees: Towards A Reinforcement Learning based Key-Value Store for Dynamic Workloads

Mo, Chen, Luo, Shan — SIGMOD 24'

Presenters: Alex Ott, James Chen

Why Use an LSM?



RocksDB



cassandra



Bigtable

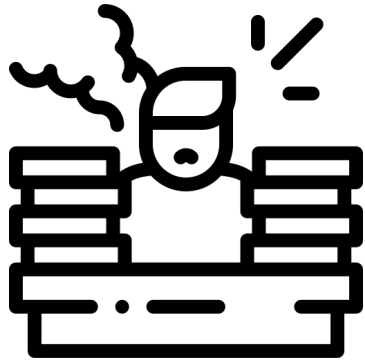


mongoDB®

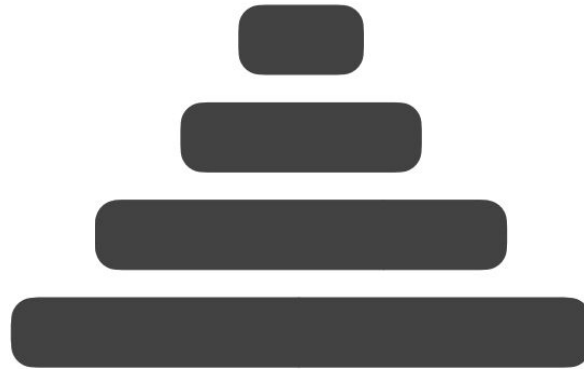


SCYLLA

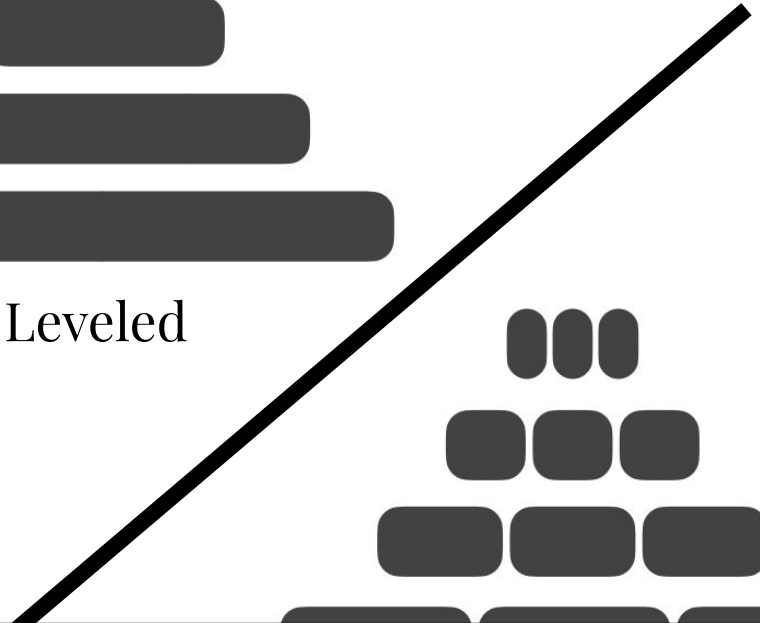
Tuning



Workload

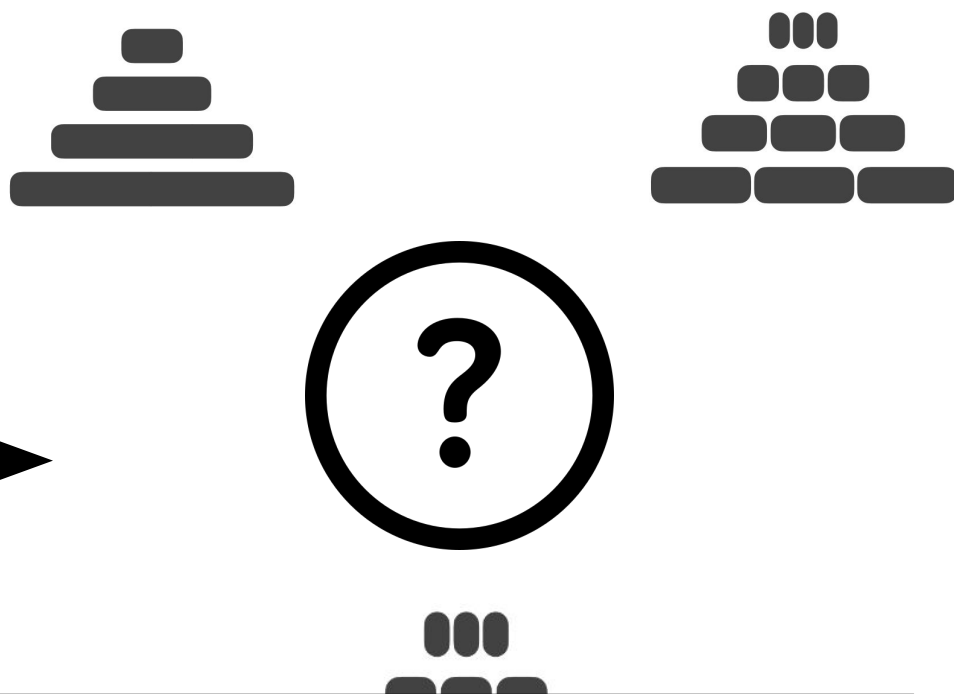
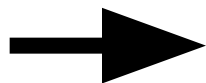
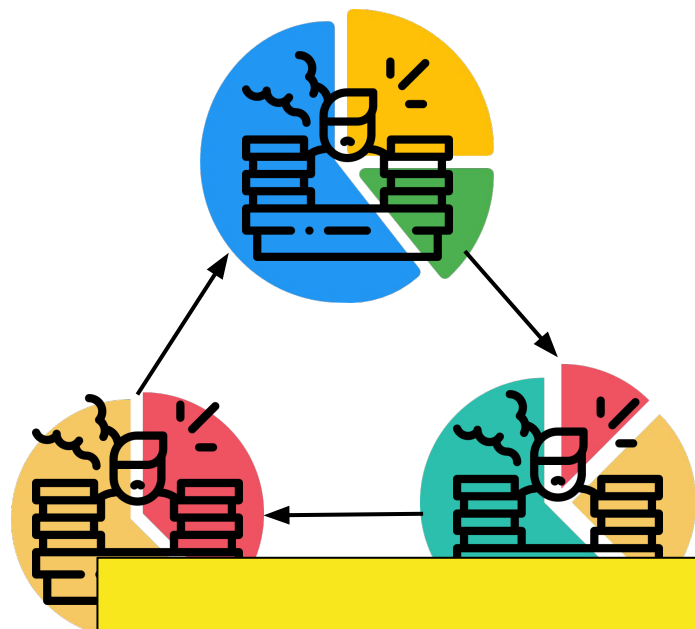


Leveled



What happens when the workload changes over time?

Dynamic Workload



How is robust tuning (Endure) different from dynamic tuning?

Motivation

Real world workloads are *dynamic*



Social Media



Data
Warehouse



Data Lake

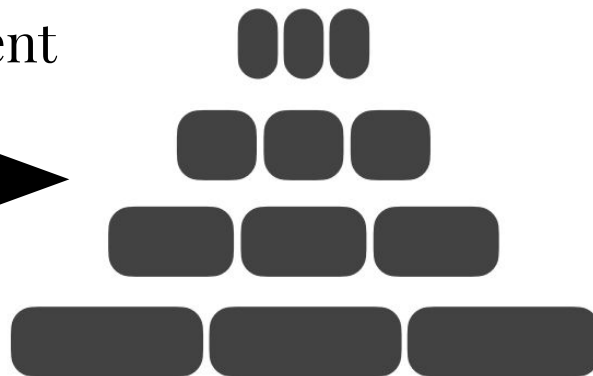
How to build a **workload-aware** LSM that adapts in **real-time**?

Intuition



Leveled

Data Movement
Policy



Tiered

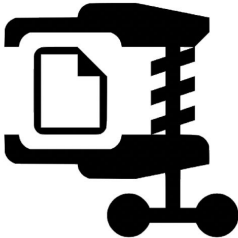
Background

Compactions Drive Data Layout

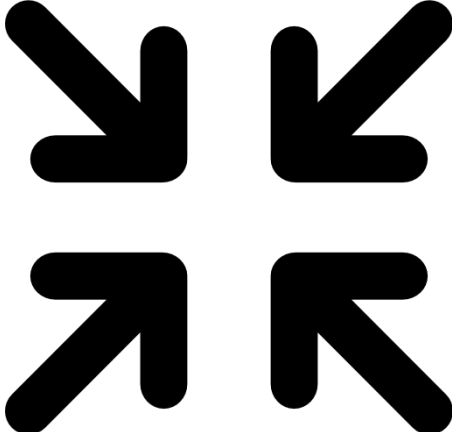
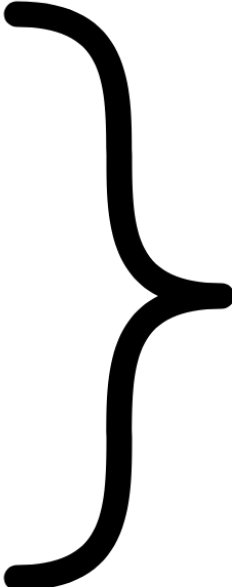
Fast Ingestion



Efficient Space Utilization



Competitive Reads



Compaction

What is Machine Learning?



how to learn, not
what to do

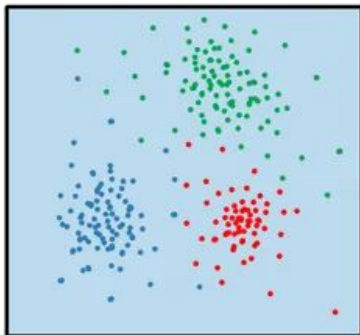


Generalizable

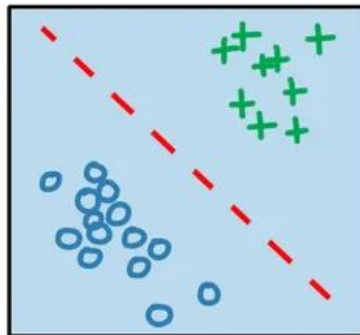
Data, Data, and Data

Too much data *machine learning*

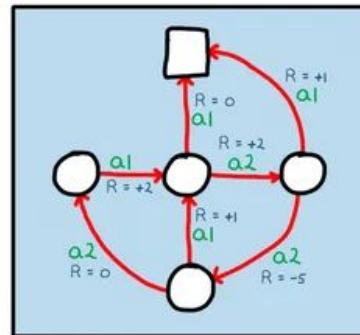
unsupervised learning



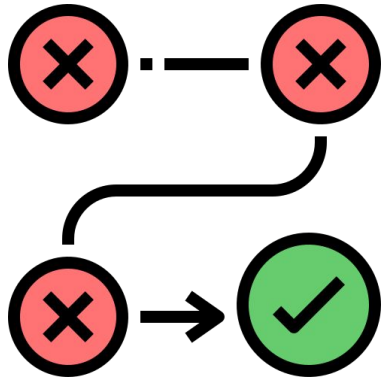
supervised learning



reinforcement learning

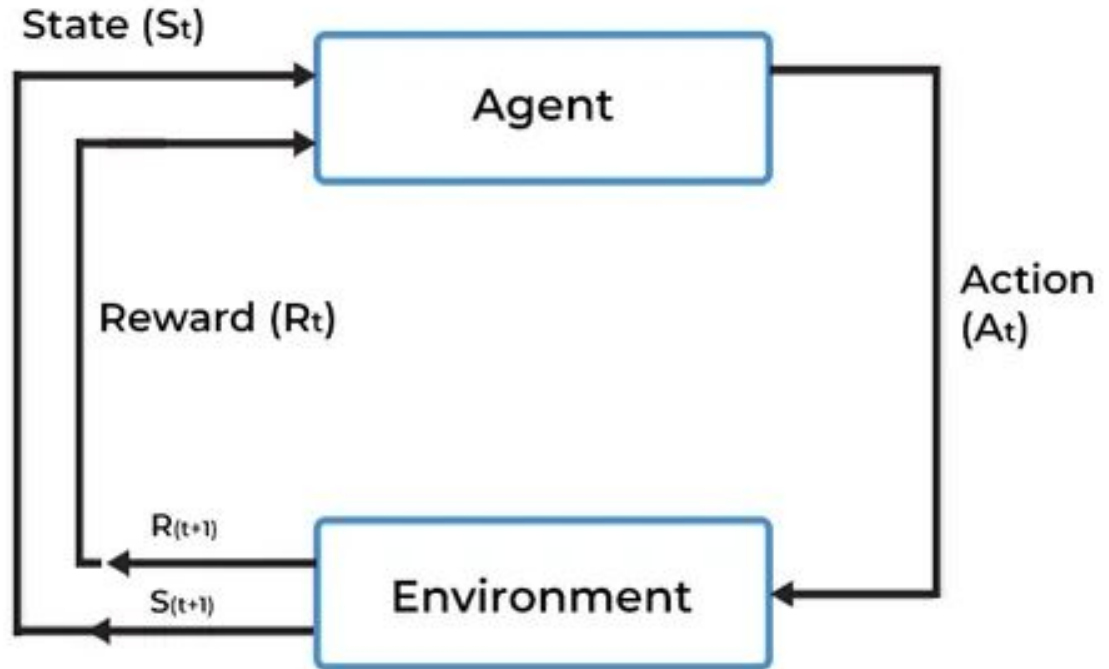


Reinforcement Learning (RL) Mental Model

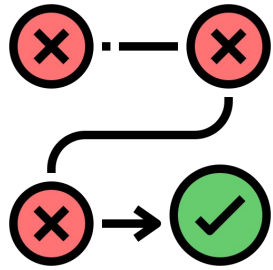


Trial and Error

REINFORCEMENT LEARNING MODEL

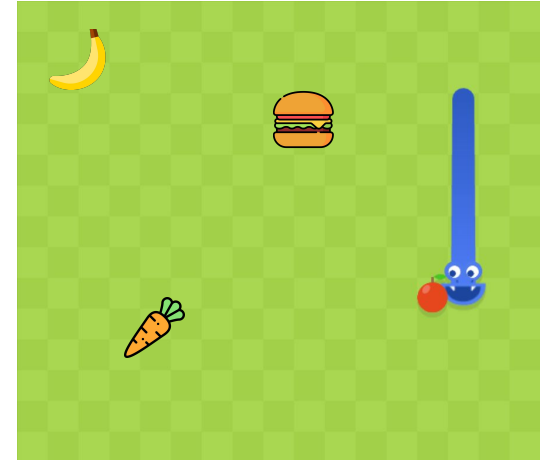
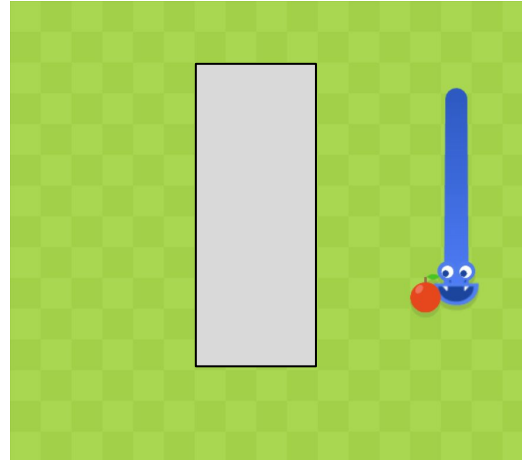
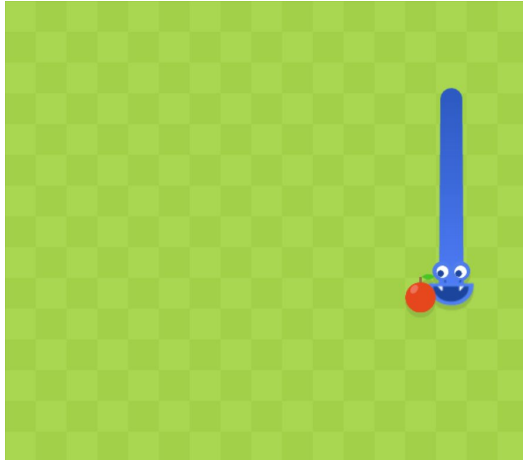
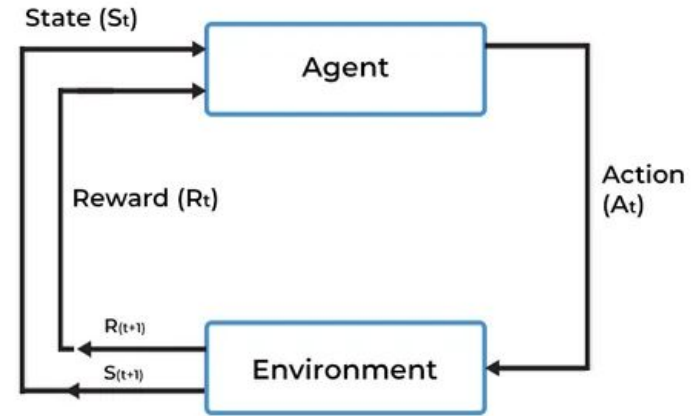


Reinforcement Learning (RL) Mental Model



Trial and Error

REINFORCEMENT LEARNING MODEL



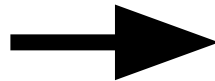
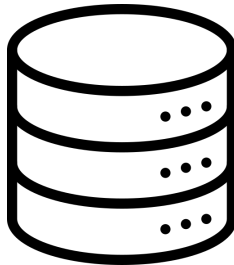
Advantages of RL?

No prior data








Challenges?

Need data to improve



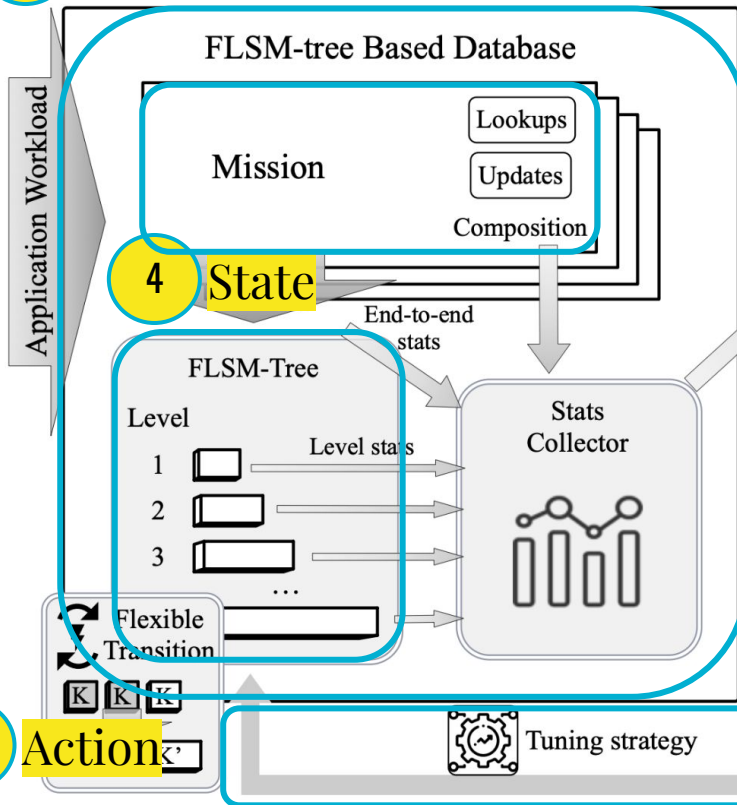
RusKey

RL in RusKey

	Agent	Environment	Actions	State	Reward
Snake Game	Snake character	2D grid	Turn left or right	Snake and apple position	+ food - Running into the wall
LSM					

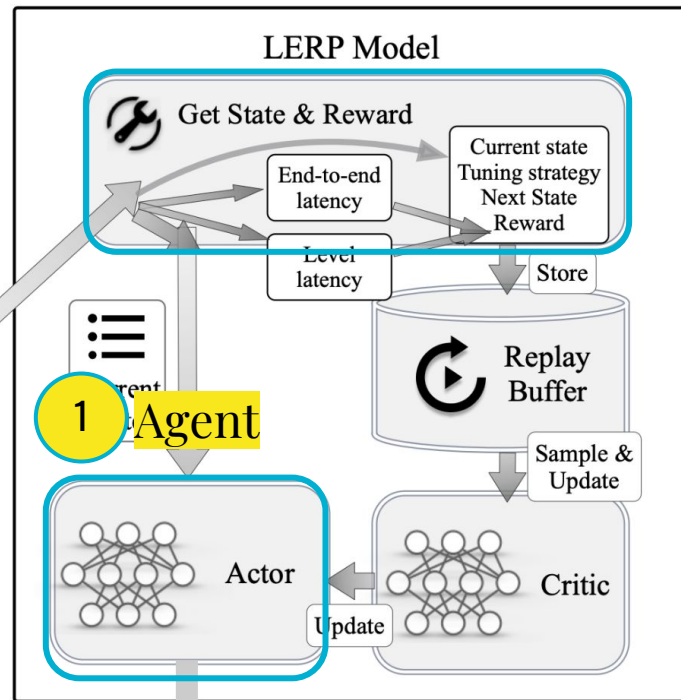
Key Components

2 Environment



4 State

5 Reward



1 Agent

3 Action

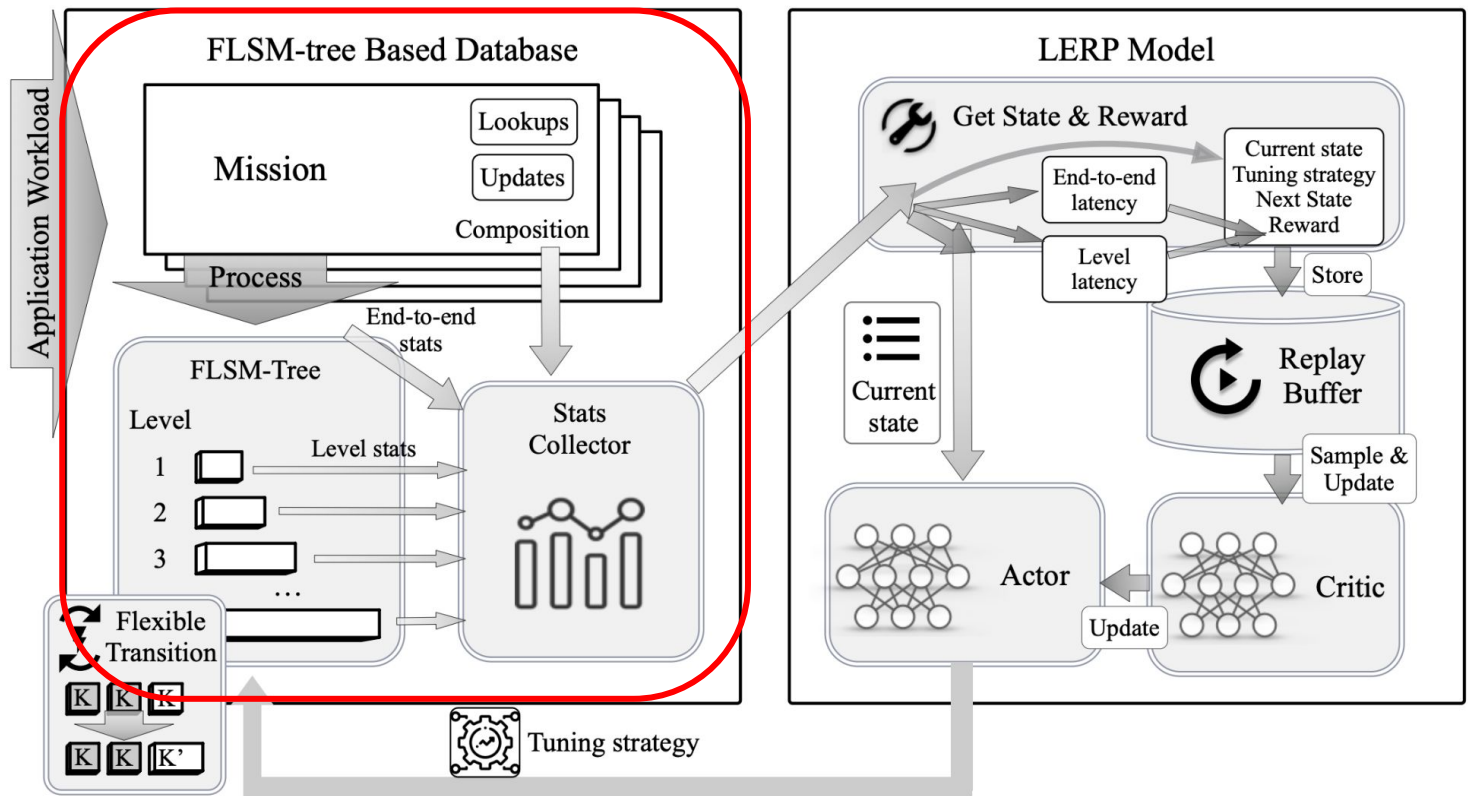
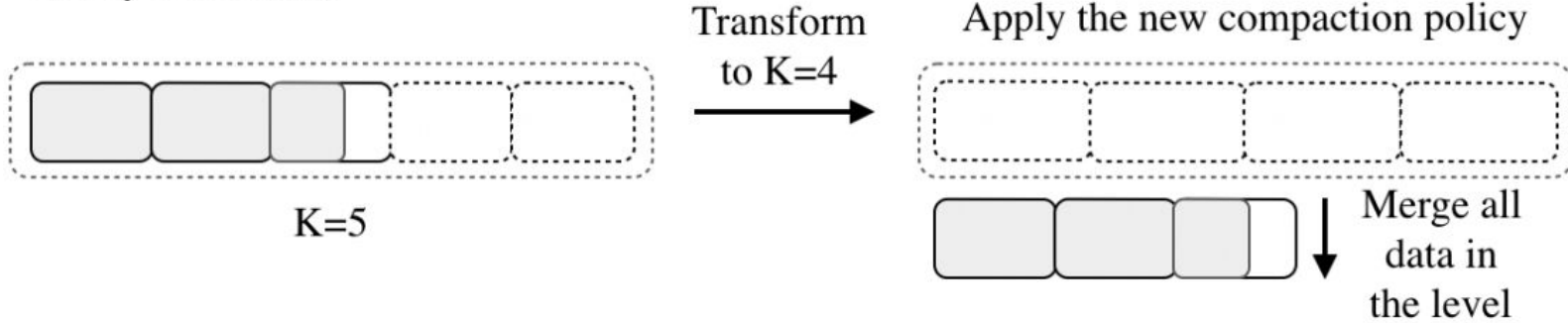


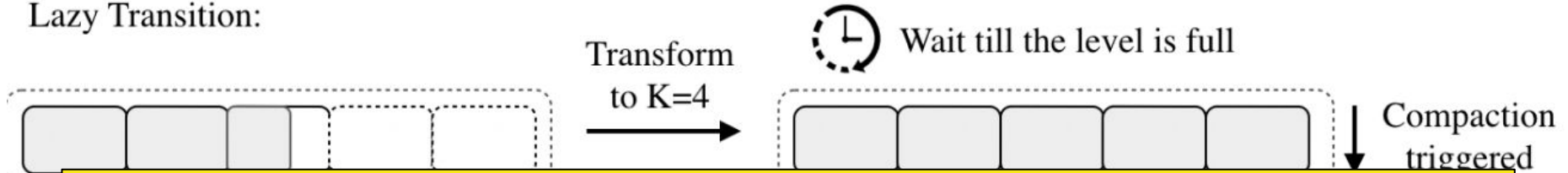
Fig. 1. The components and workflow of RusKEY.

Transitions: Greedy, Lazy, and Something in Between

Greedy Transition:



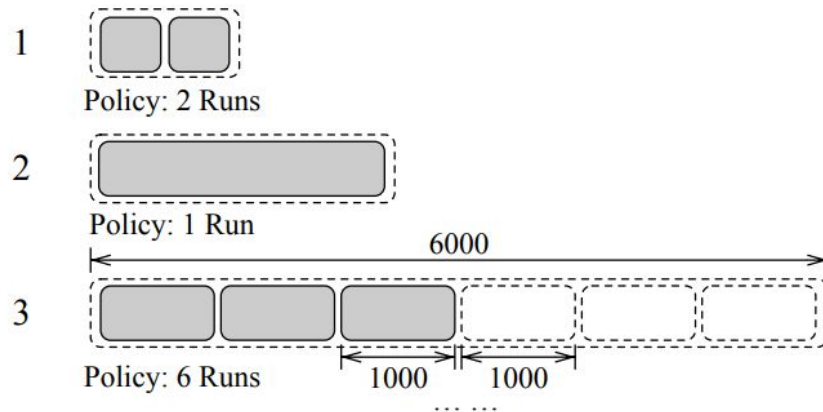
Lazy Transition:



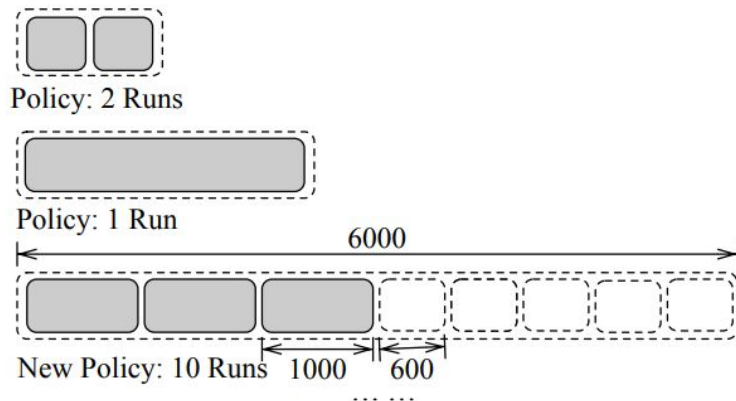
What is the problem with these schemes?

Apply the new compaction policy

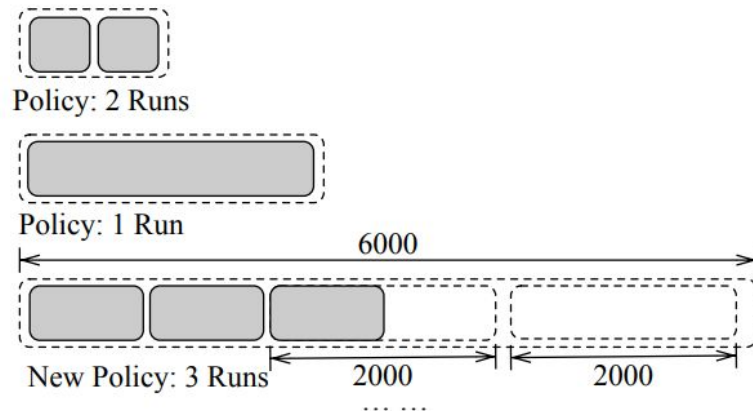
Level (A) Original LSM-tree



(B) Transform Policy at Level 3 to 10



(C) Transform Policy at Level 3 to 3



LERP Framework

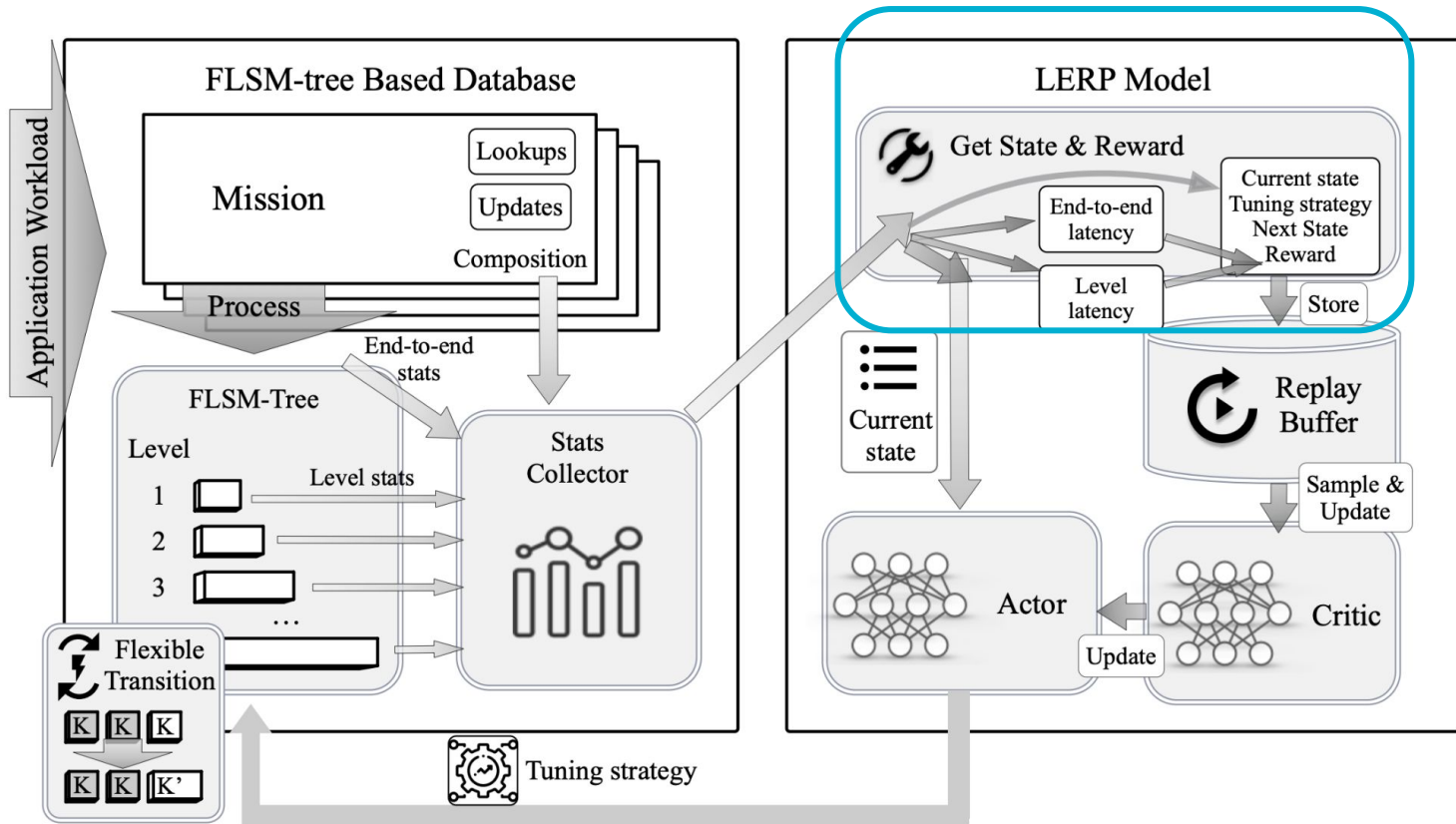


Fig. 1. The components and workflow of RusKEY.

LERP Framework

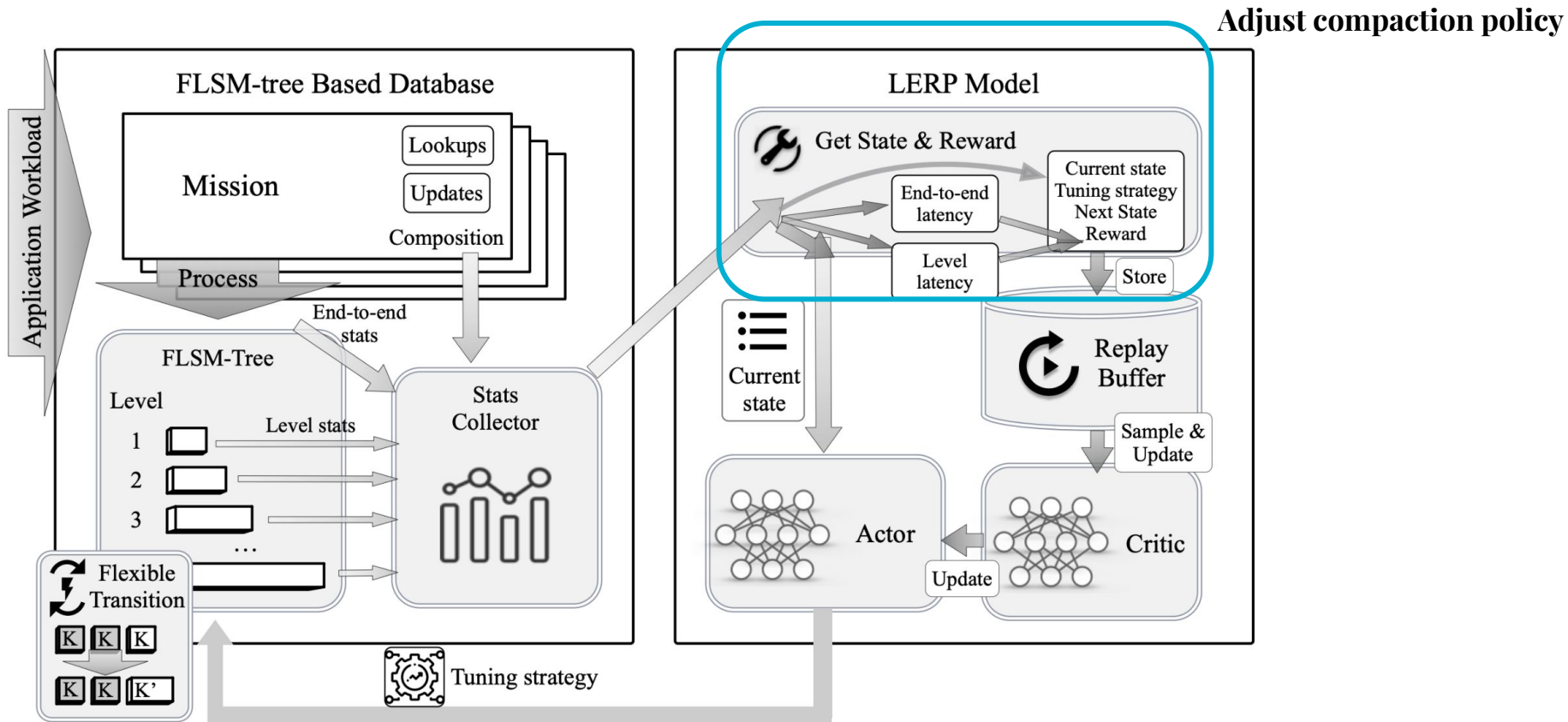


Fig. 1. The components and workflow of RusKEY.

LERP Framework

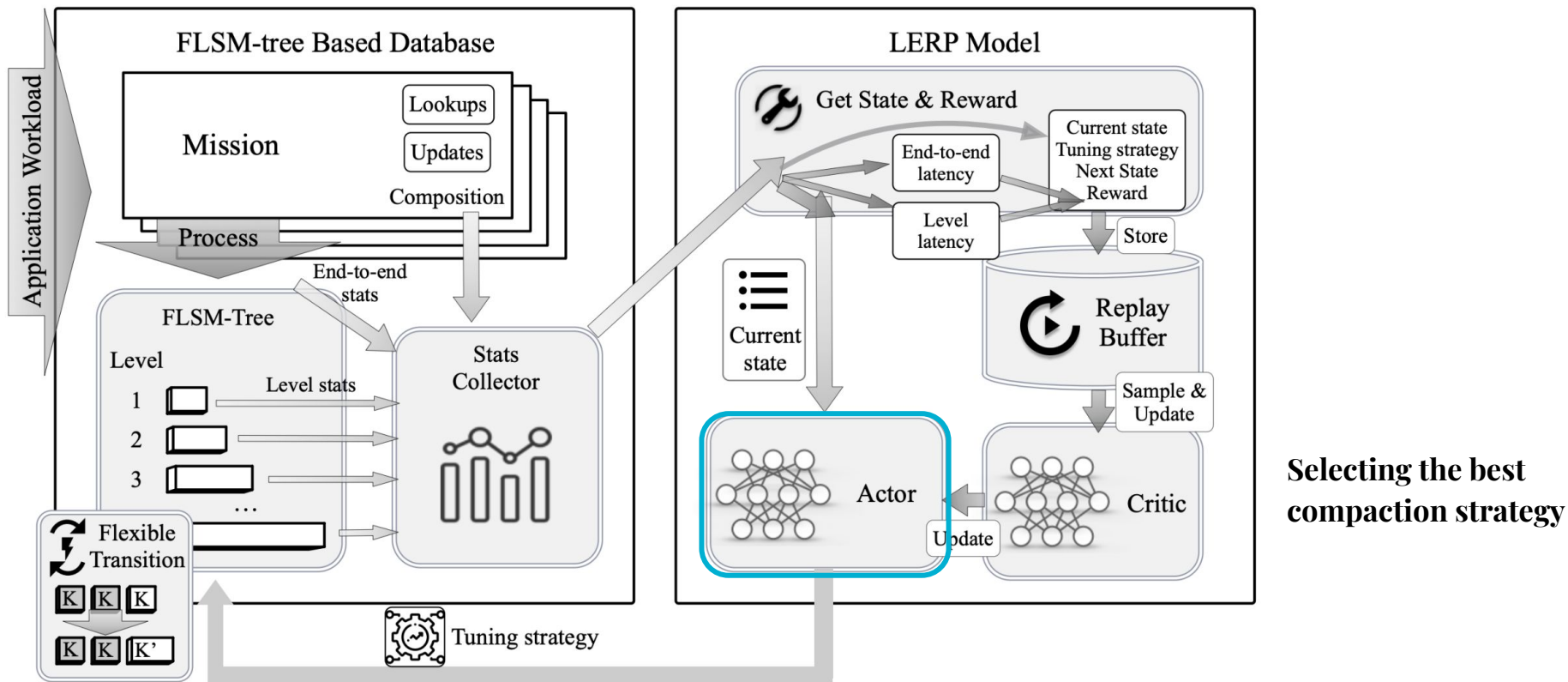


Fig. 1. The components and workflow of RusKEY.

LERP Framework

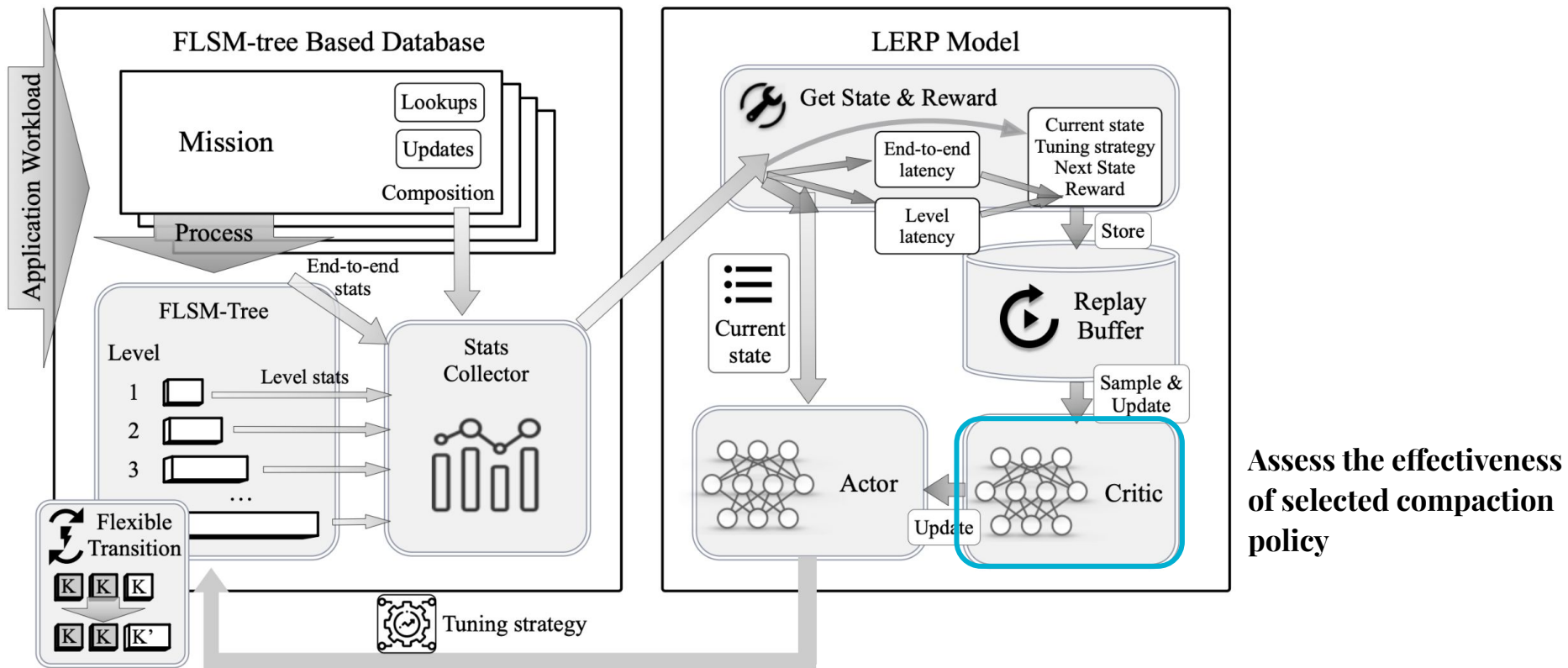


Fig. 1. The components and workflow of RusKEY.

LERP Framework

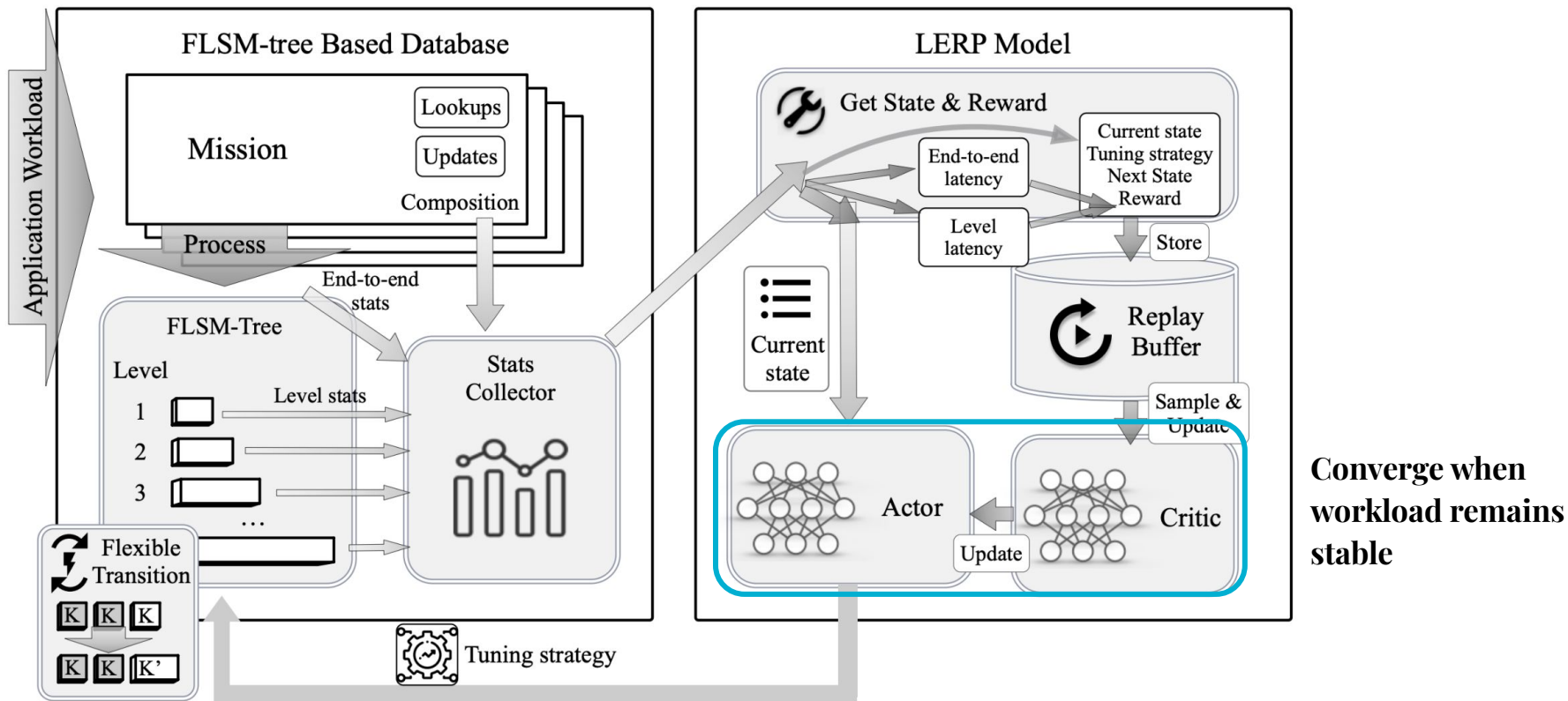
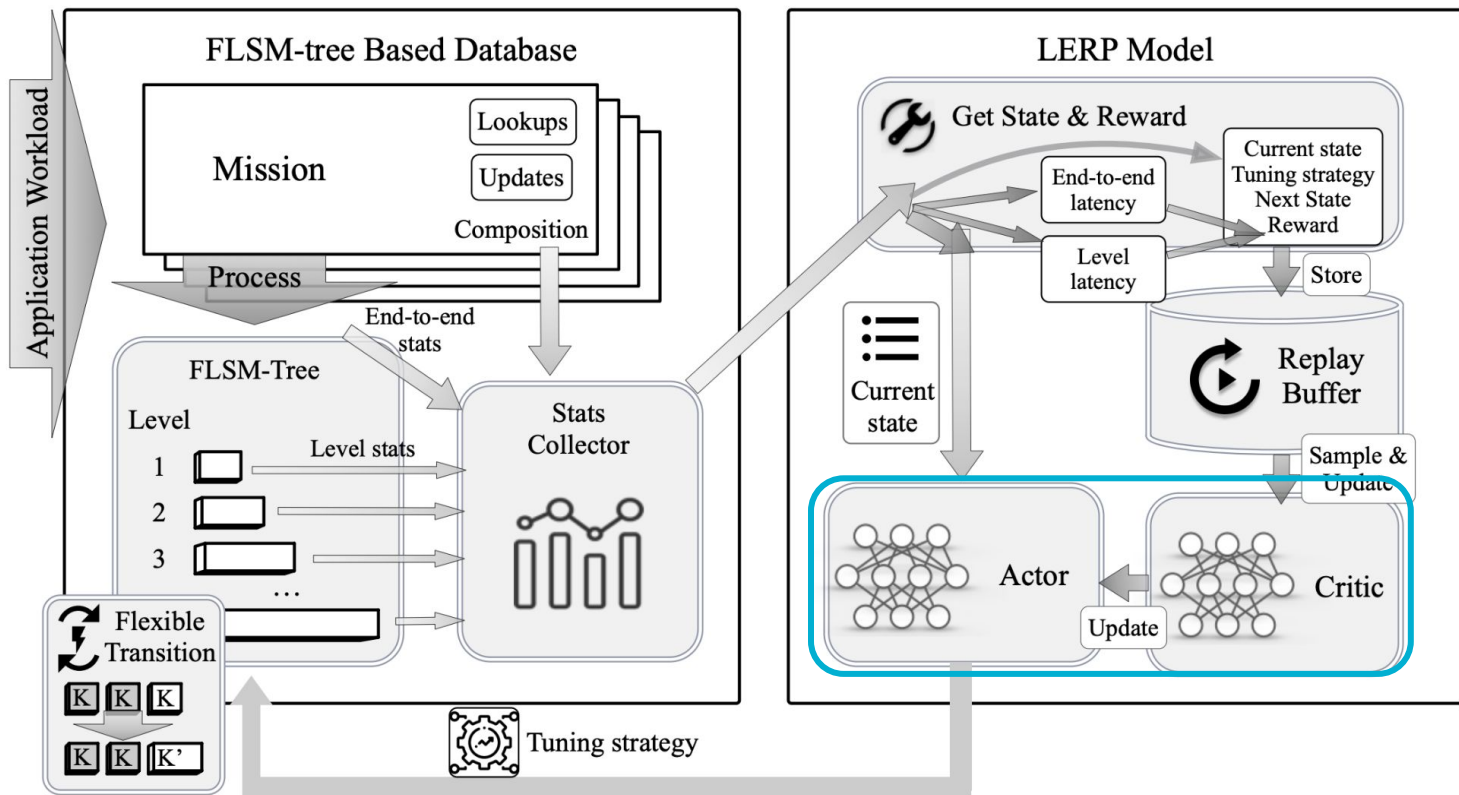


Fig. 1. The components and workflow of RusKEY.

LERP Framework



Exit convergence when workload changes.

LERP explores new compaction strategy adapted to new workload.

Fig. 1. The components and workflow of RusKEY.

Level Based RL Model



State

LSM-tree
configuration



Action Space

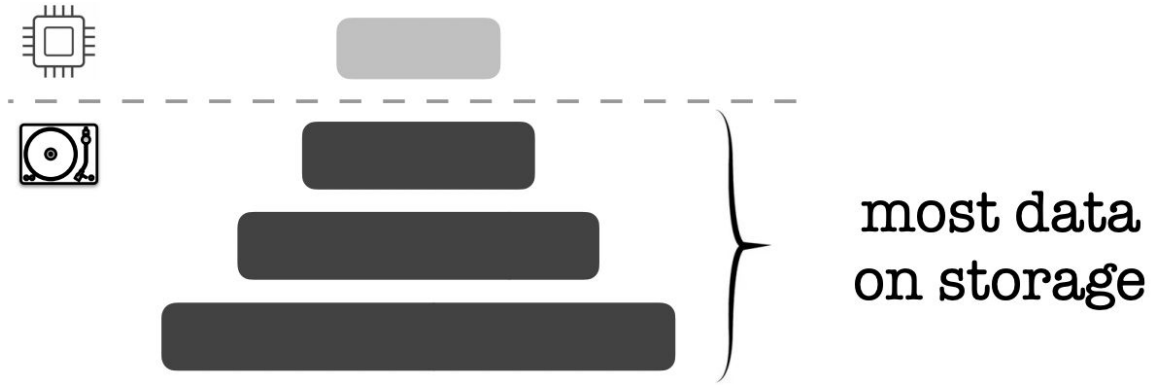
Increase or
decrease in
compaction
policies (level
specific)



Reward

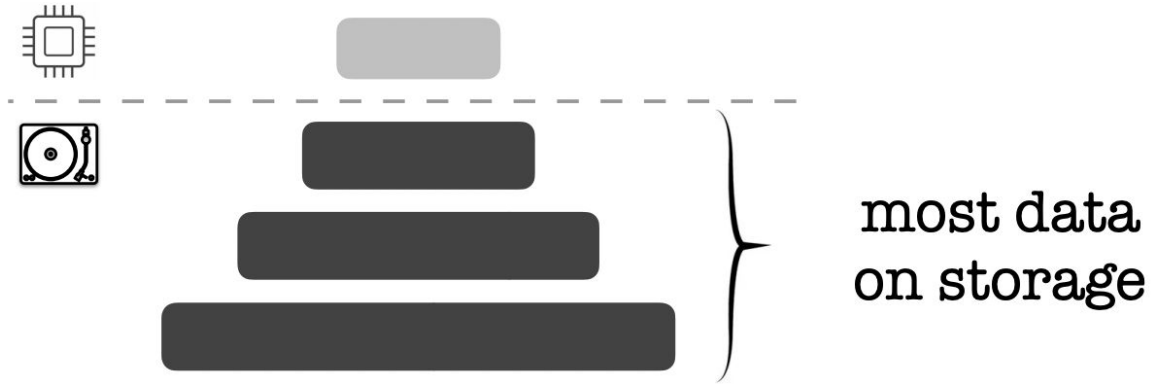
Minimize
latency

Why Policy Propagation to Optimize Training?



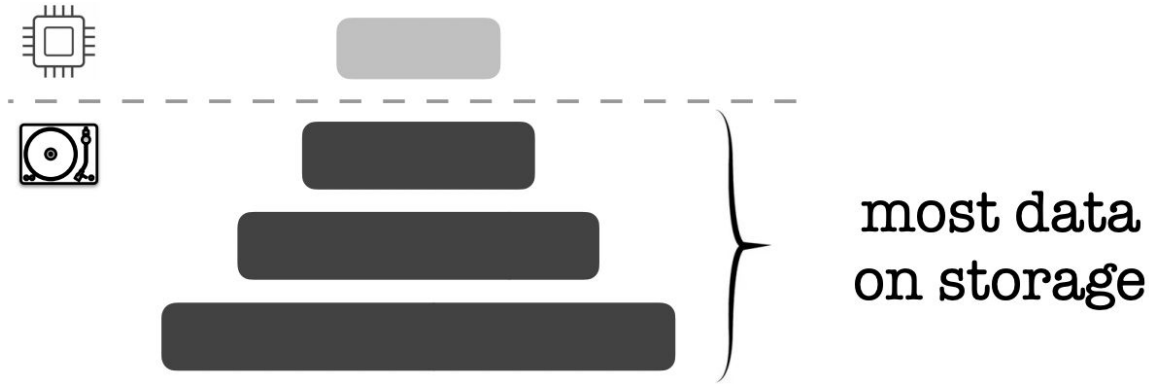
Deeper level size is exponentially larger

Why Policy Propagation to Optimize Training?



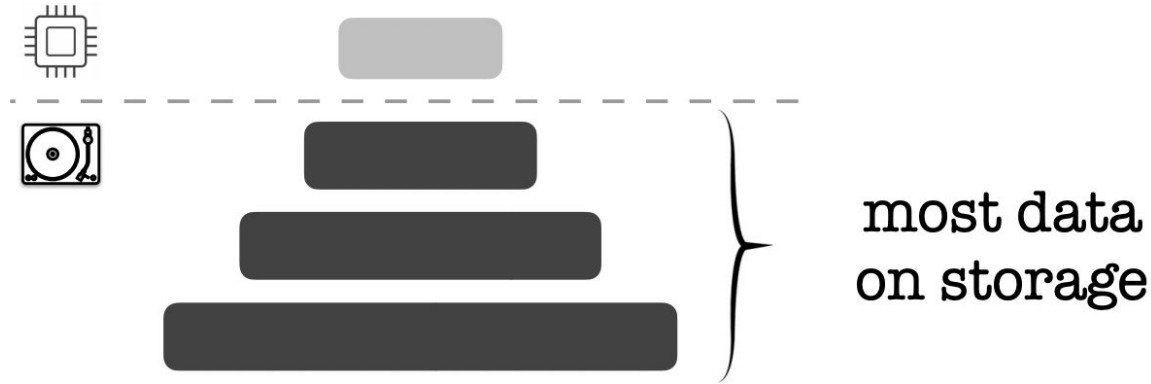
Deeper level size is exponentially larger
Compaction happens less frequently at deeper level

Why Policy Propagation to Optimize Training?



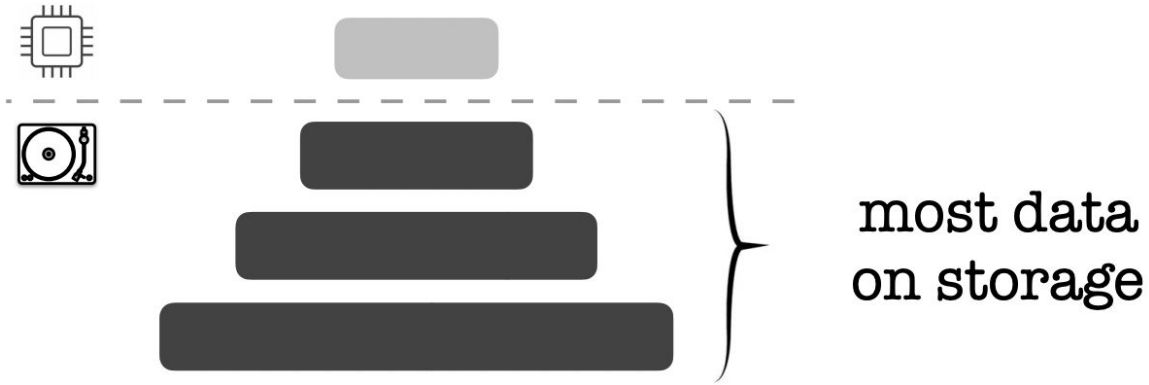
Deeper level size is exponentially larger
Compaction happen less frequently at deeper level
Less available training data

Why Policy Propagation to Optimize Training?



- Deeper level size is exponentially larger**
- Compaction happen less frequently at deeper level**
- Less available training data**
- Deeper level needs more training data**

Why Policy Propagation to Optimize Training?



How to learn the compaction policy at a deeper level?

Policy Propagation based on Cost Analysis



Uniform Bits-per-Key

Assigns the same bits-per-key to BF
(per level)

R/W cost ratio is similar across level



Monkey Allocation

Shallow level assigns more
bits-per-key than deeper
level

Policy Propagation based on Cost Analysis



Uniform Bits-per-Key

Assigns the same bits-per-key to BF
(per level)

R/W cost ratio is similar across level



Monkey Allocation

Shallow level assigns more
bits-per-key than deeper
level

Policy Propagation based on Cost Analysis



Uniform Bits-per-Key

Assigns the same bits-per-key to Bloom Filter (BF) (per level)

R/W cost ratio is similar across level
R/W amplification is same across level



Monkey Allocation

Shallow level assigns more bits-per-key than deeper level

Policy Propagation based on Cost Analysis



Uniform Bits-per-Key

Assigns the same bits-per-key to BF
(per level)

Use RL model to learn the policy of level 1
Propagate the policy to all levels



Monkey Allocation

Shallow level assigns more
bits-per-key than deeper
level

Policy Propagation based on Cost Analysis



Uniform Bits-per-Key

Assigns the same bits-per-key to BF
(per level)



Monkey Allocation

Shallow level assigns more
bits-per-key than deeper
level

R/W cost ratio is different

Policy Propagation based on Cost Analysis



Uniform Bits-per-Key

Assigns the same bits-per-key to BF
(per level)



Monkey Allocation

Shallow level assigns more
bits-per-key than deeper
level

R/W cost ratio is different
Policy can vary across level

Policy Propagation based on Cost Analysis



Uniform Bits-per-Key

Assigns the same bits-per-key to BF (per level)



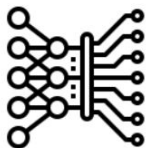
Monkey Allocation

Shallow level assigns more bits-per-key than deeper level

Infer remaining level policy based on the preceding two levels policy

Results

Experiment Design



Hardware:

Intel Xeon Gold
6326@2.9 GHz
CPU, NVMe
SSD, Ubuntu
22.04 OS.



Initial Data Load:

100 million kv
entries



Key Size: 128B

Value Size: 896B



Mission:

100 million operations
(lookup or updates)
divided into 2000 mission

Baseline Comparison

- **Aggressive Compaction ($K = 1$):**
 - low read cost
 - high write amplification
- **Lazy Compaction ($K = 10$):**
 - high read cost
 - low write amplification
- **Moderate Compaction ($K = 5$):**
 - Balance between aggressive and lazy
 - moderate read/write amplification

K

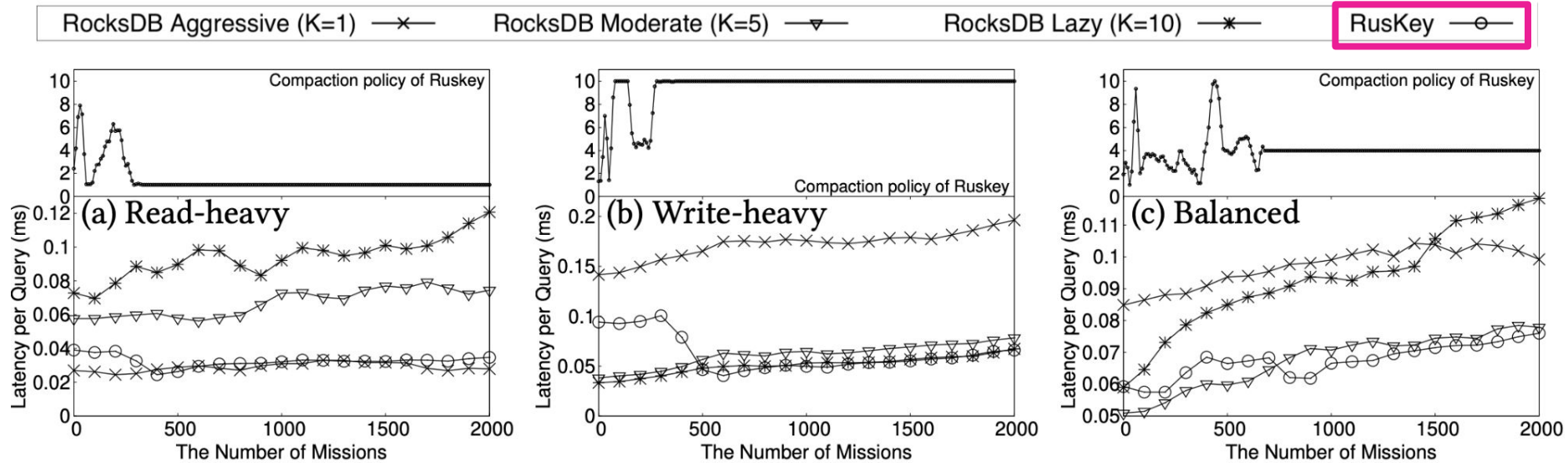
Max Number of
sorted runs per level

Workload

- **Workload Types:**
 - read-heavy (10% update)
 - balanced (50% update)
 - write-heavy (90% update),
 - write-inclined (70% update)
 - read-inclined (30% update)

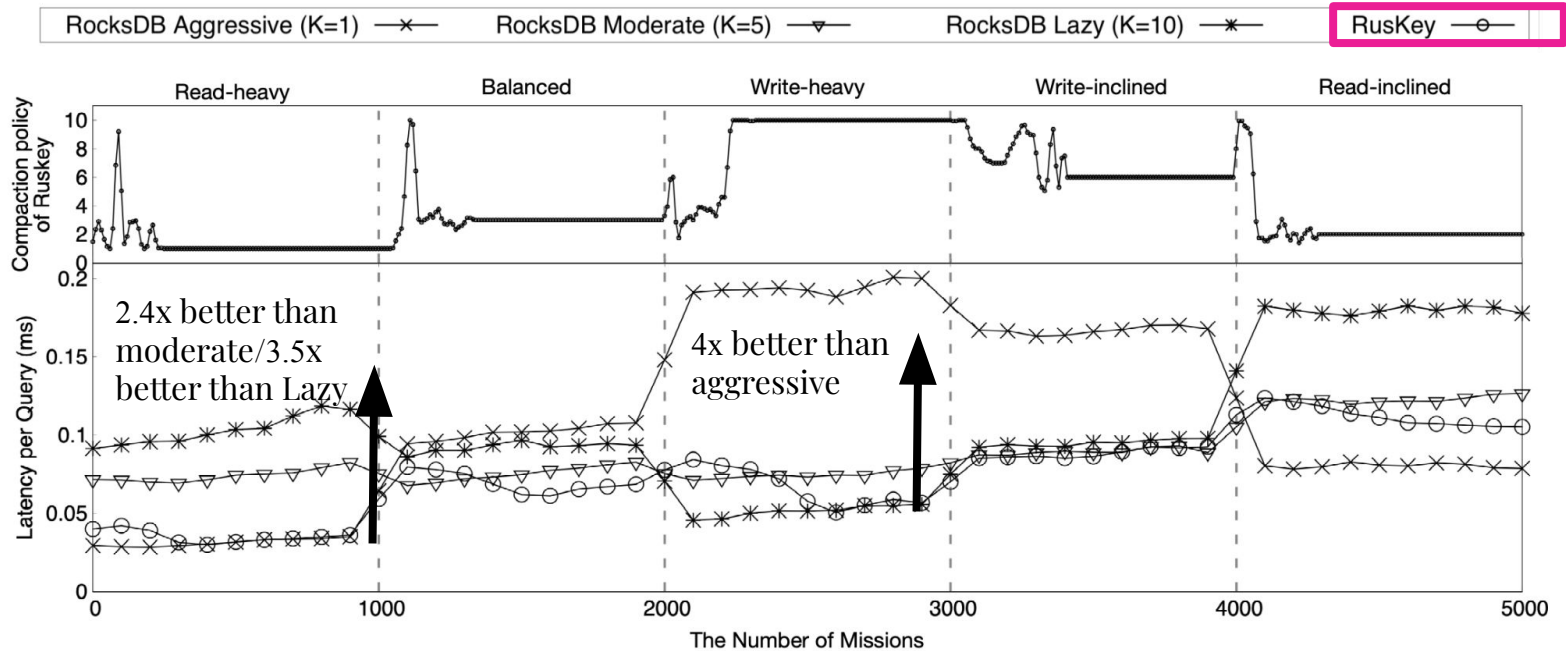
Each section includes 50 million operations which are divided into 1000 missions with 50000 operations for each.

Evaluation - Static Workload



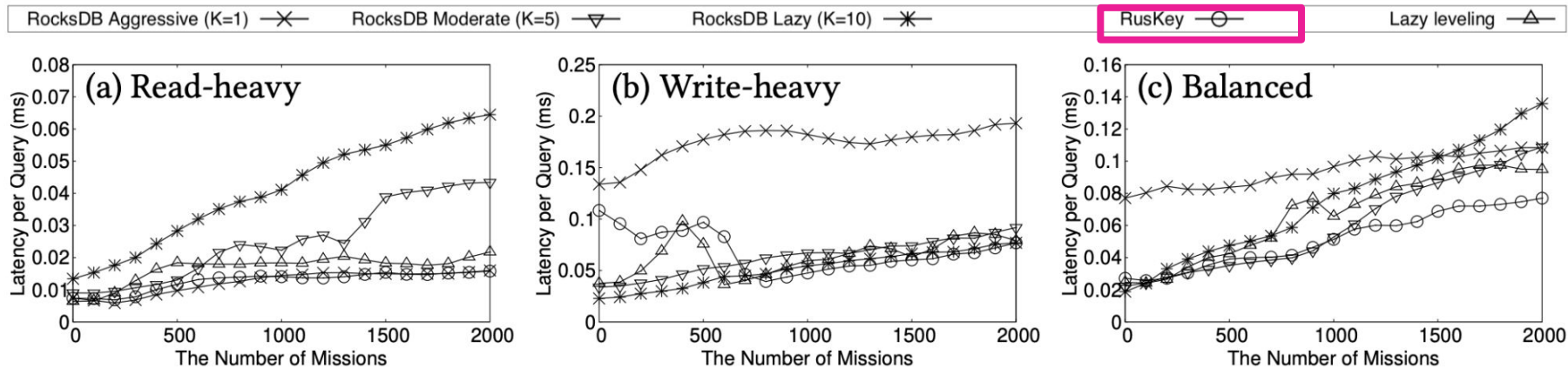
Ruskey achieves the lowest latency per query across all static workloads as missions increase

Evaluation - Dynamic Workload



RusKey maintains near-optimal latency across all dynamic workload sessions, while other baselines exhibit sub-optimal performance in at least one session.

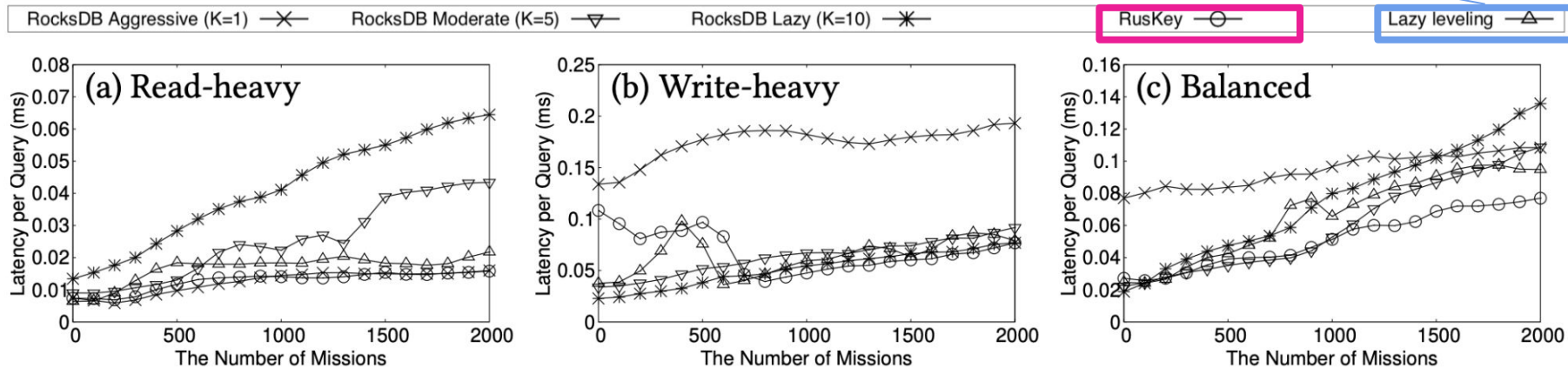
Ruskey Under Monkey Scheme



Compare RusKey with the baselines under the Monkey scheme with the same workload setting

Ruskey Under Monkey Scheme

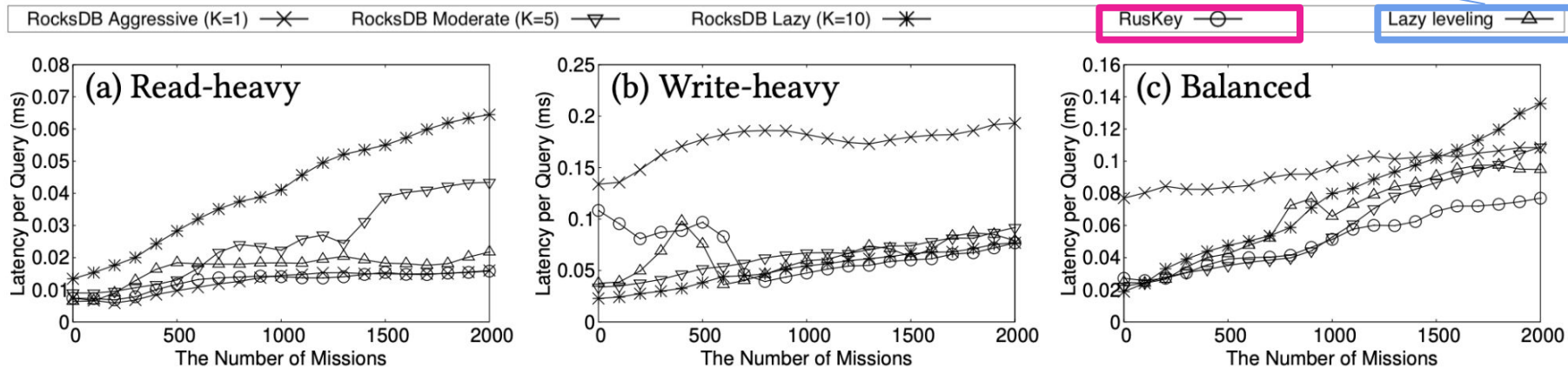
A state-of-the-art compaction policy



Compare RusKey with the baselines under the Monkey scheme with the same workload setting

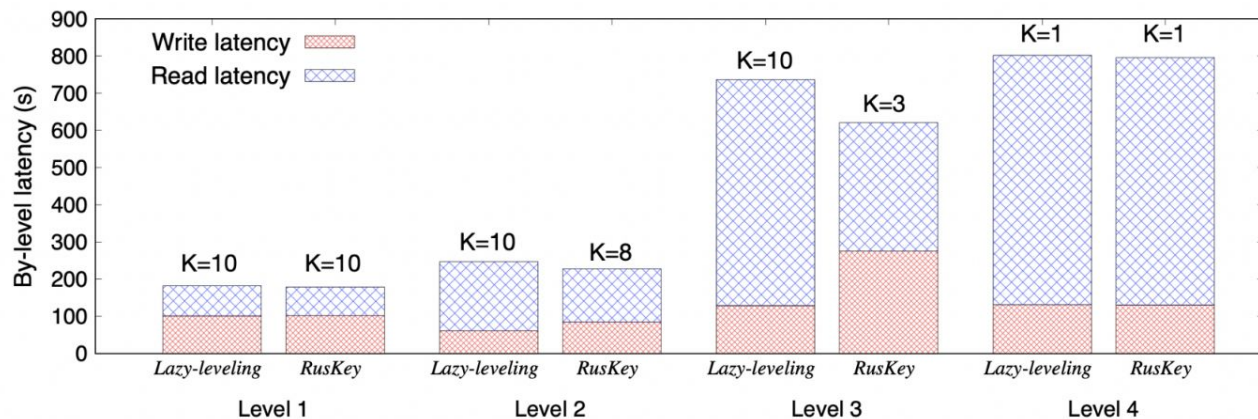
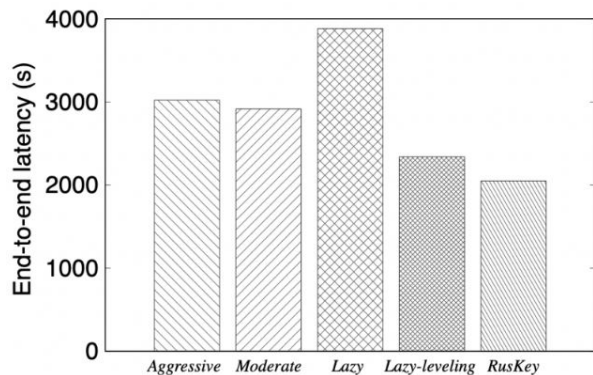
Ruskey Under Monkey Scheme

A state-of-the-art compaction policy



Ruskey and Lazy Leveling both achieve near-optimal performance, but Ruskey outperforms Lazy Leveling in every workload since it adopts novel policy setting through policy propagation

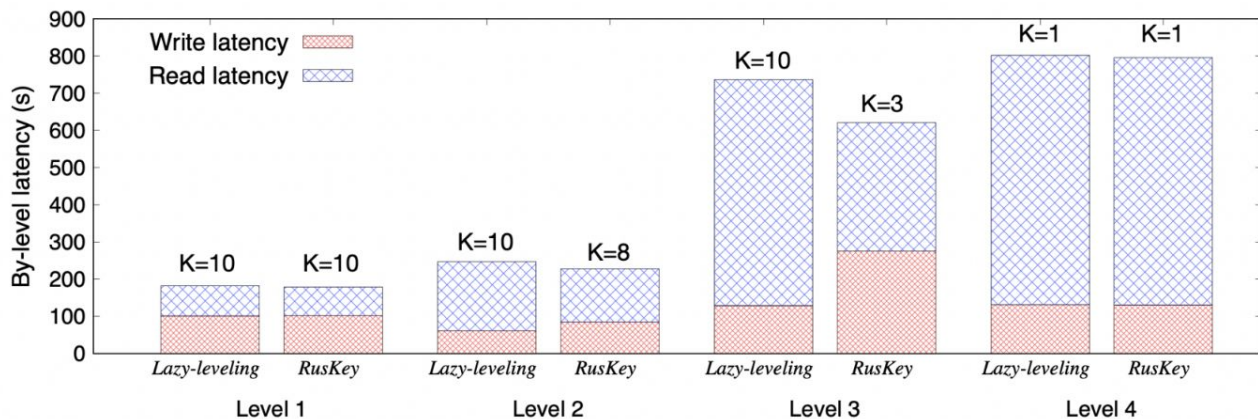
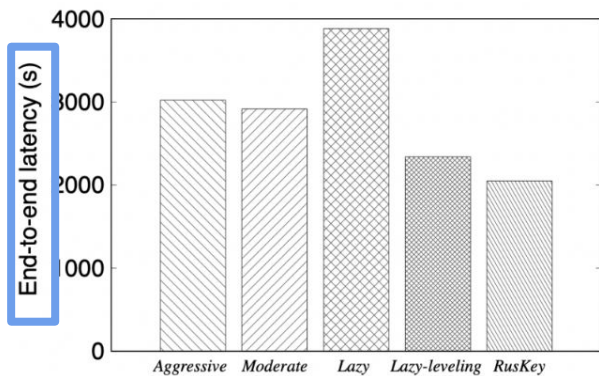
Compaction policy setting of RusKey



25 million operations of r/w balanced workload under the Monkey scheme

Compaction policy setting of RusKey

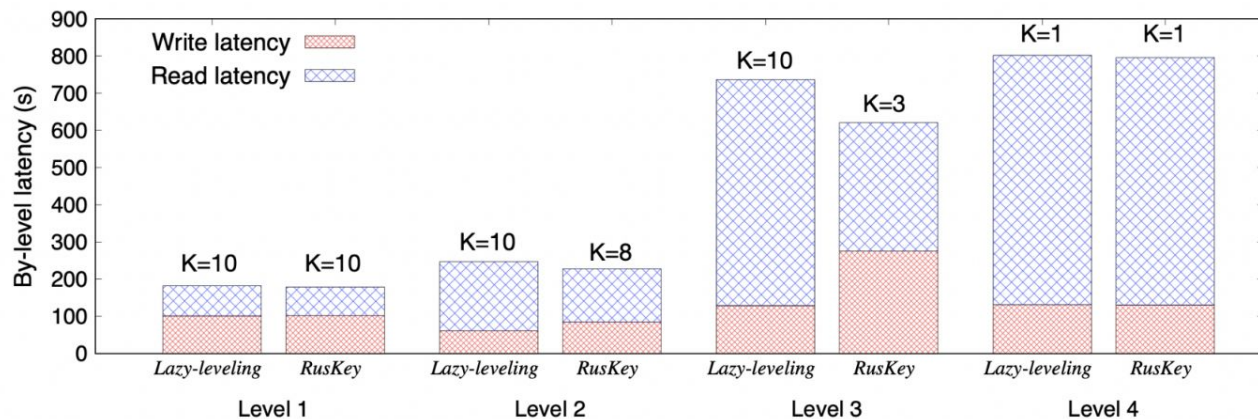
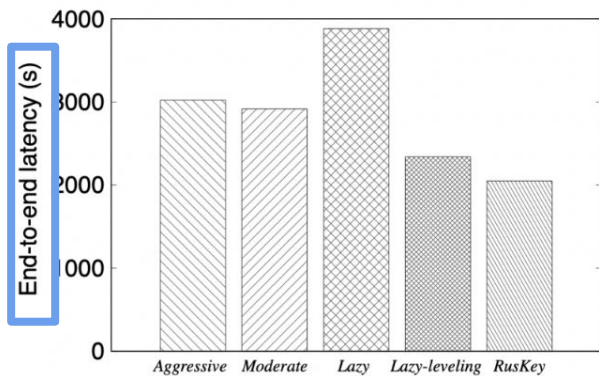
The running time of processing all the operation



25 million operations of r/w balanced workload under the Monkey scheme

Compaction policy setting of RusKey

The running time of processing all the operation

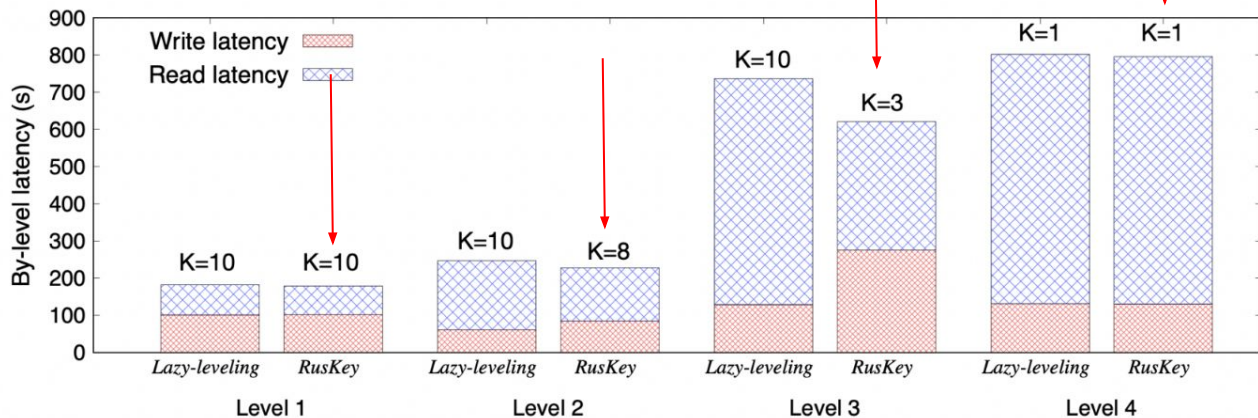
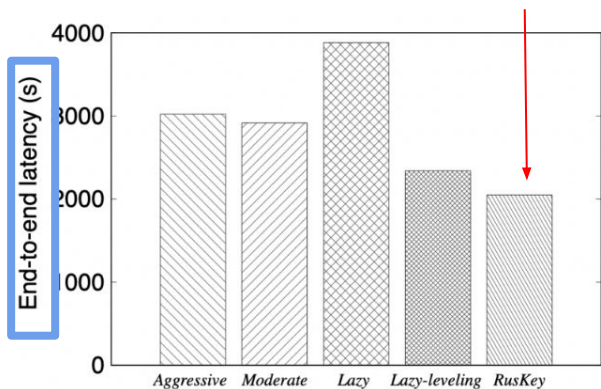


25 million operations of r/w balanced workload under the Monkey scheme

Different Compaction policy at each level (Ruskey)

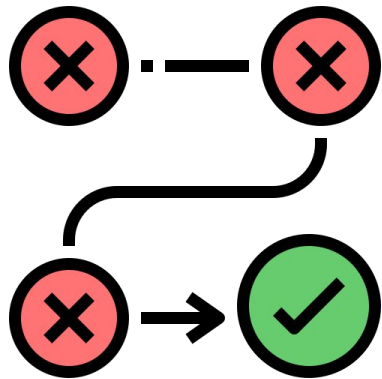
Compaction policy setting of RusKey

The running time of processing all the operation

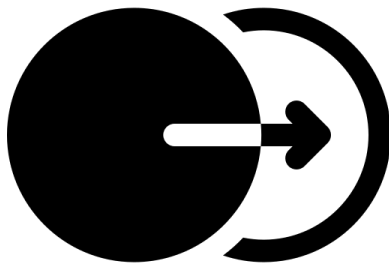


RusKey achieves optimal end-to-end and by-level latency by self-tuning its compaction policy under a balanced workload

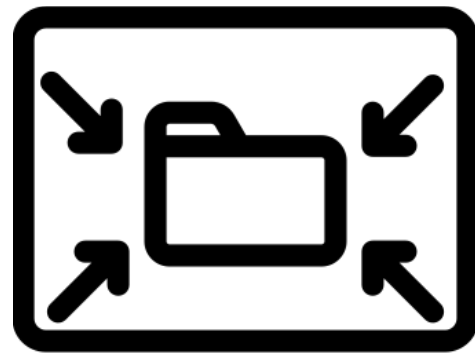
Summary



Reinforcement
Learning



Efficient
Transitions



Reduce
Required Data