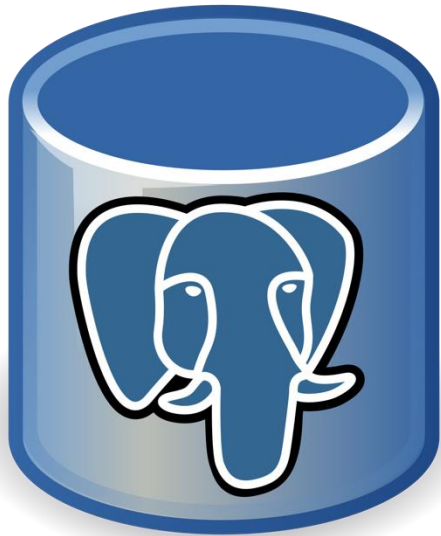




Robust Database Tuning with ENDURE

Andy Huynh, Harshal A. Chaudhari, Evimaria Terzi, Manos Athanassoulis

Databases Have Settings...



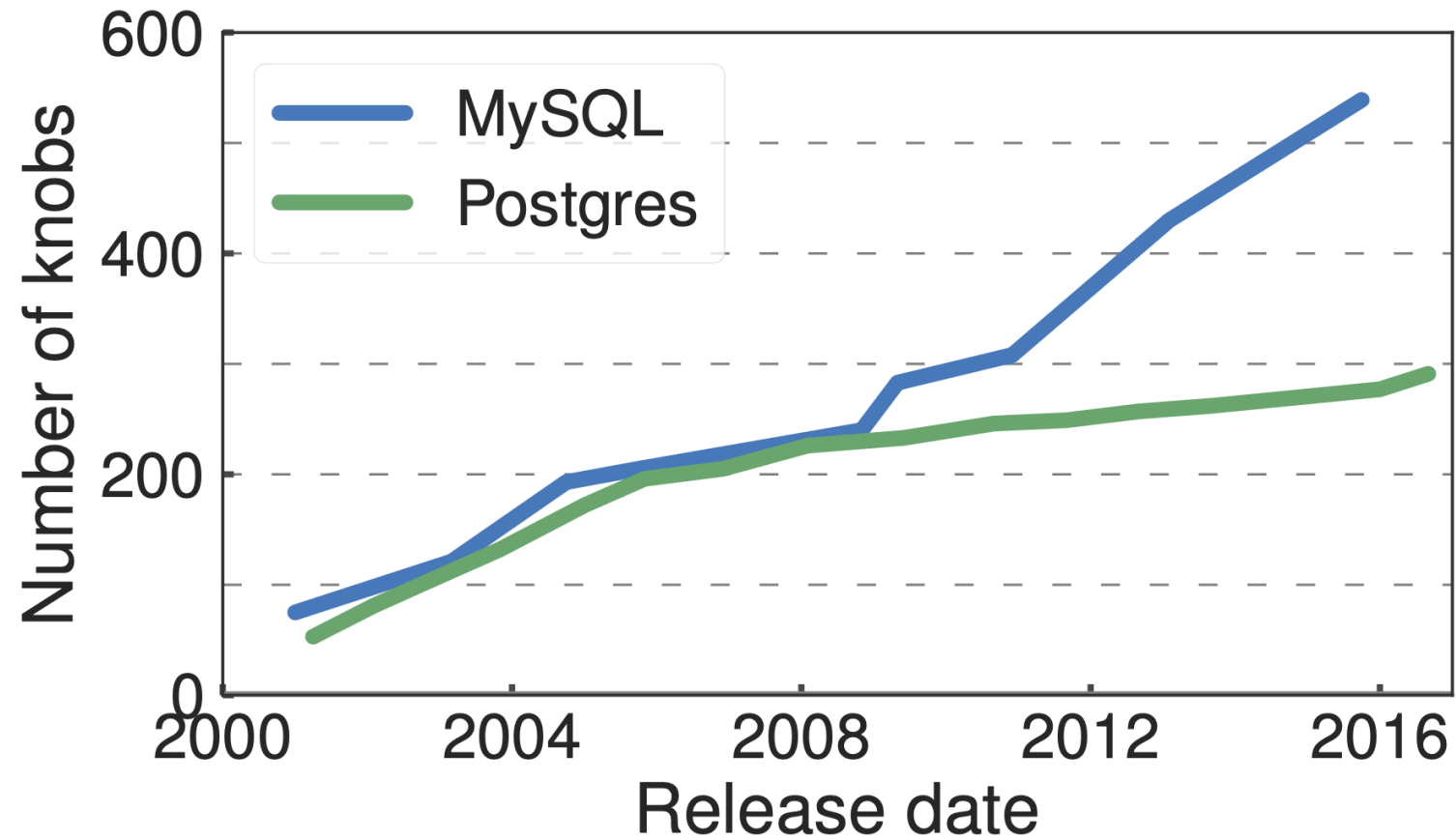
- ⚙ effective_cache_size
- ⚙ work_mem
- ⚙ wal_sync_method
- ⚙ max_prepared_transactions
- ⚙ random_page_cost
- ⚙ checkpoint_segments
- ⚙ maintenance_work_mem
- ...
- ⚙ shared_buffers

200+ settings



**Determines
performance**

Database Complexity



Tuning Makes a HUGE Difference

Databricks Sets Official Data Warehousing Performance Record



by Reynold Xin and Mostafa Mokhtar
Posted in COMPANY BLOG | November 15, 2021

AUTHOR



Benoit Dageville



Thierry Cruanes

Today, we are proud to announce that **Databricks SQL outperformed the previous 100TB TPC-DS**, the gold standard performance benchmark news, this result has been for

These results were corroborated by research which frequently runs TPC-DS on popular **benchmarked Databricks and Snowflake and 12x better in terms of price performance** warehouses such as Snowflake become production.

SHARE



SUBSCRIBE

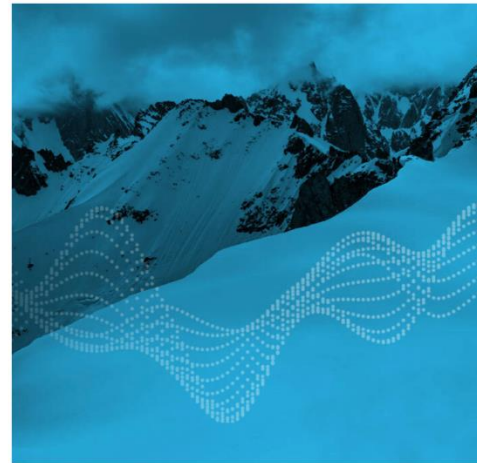
NOV 12, 2021

Industry Benchmarks and Competitors

Thought Leadership > Executive Platform



by Mostafa Mokhtar, Arsalan Tavakoli-Shiraji, Reynold Xin and Matei Zaharia
Posted in COMPANY BLOG | November 15, 2021

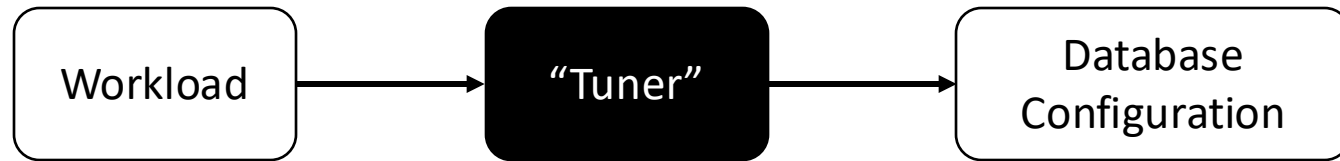


On Nov 2, 2021, we announced that **we set the official world record** for the fastest data warehouse with our Databricks SQL lakehouse platform. These results were audited and reported by the official Transaction Processing Performance Council (TPC) in a 37-page document **available online** at tpc.org. We also shared a third-party benchmark by the Barcelona Supercomputing Center (BSC) outlining that Databricks SQL is significantly faster and more cost effective than Snowflake.

A lot has happened since then: many congratulations, some questions, and some sour grapes. We take this opportunity to reiterate that **we stand by our blog post and the results: Databricks SQL provides superior performance and price performance over Snowflake, even on data warehousing workloads (TPC-DS).**

When we founded Snowflake, we set out to build an innovative platform. We had the opportunity to take into account what had worked well and what hadn't in prior architectures and implementations. We saw how we could leverage the cloud to rethink the limits of what was possible. We also focused on ease of use and building a system that "just worked." We knew there were many opportunities to improve upon prior implementations and innovate to lead on performance and scale, simplicity of administration, and data-driven collaboration.

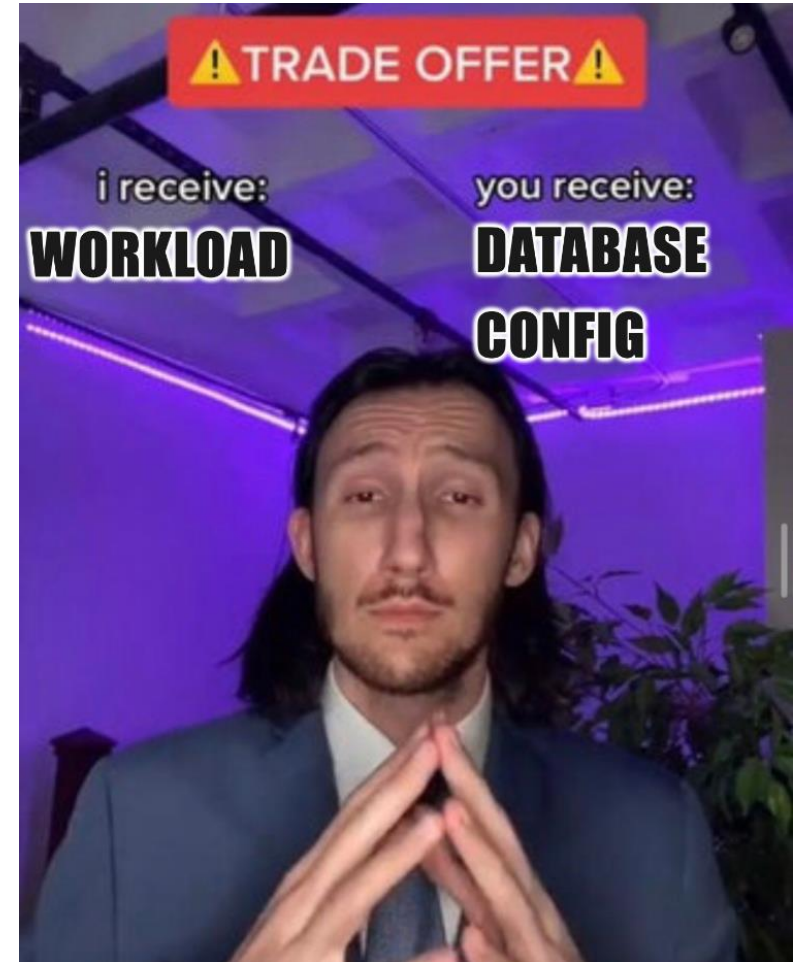
Tuning is Simple...?



What is a workload?

How do we tune?

What are you optimizing for?



Age of Log-Structured Merge-Trees



High impact tuning knobs



Compaction Buffer size Size ratio

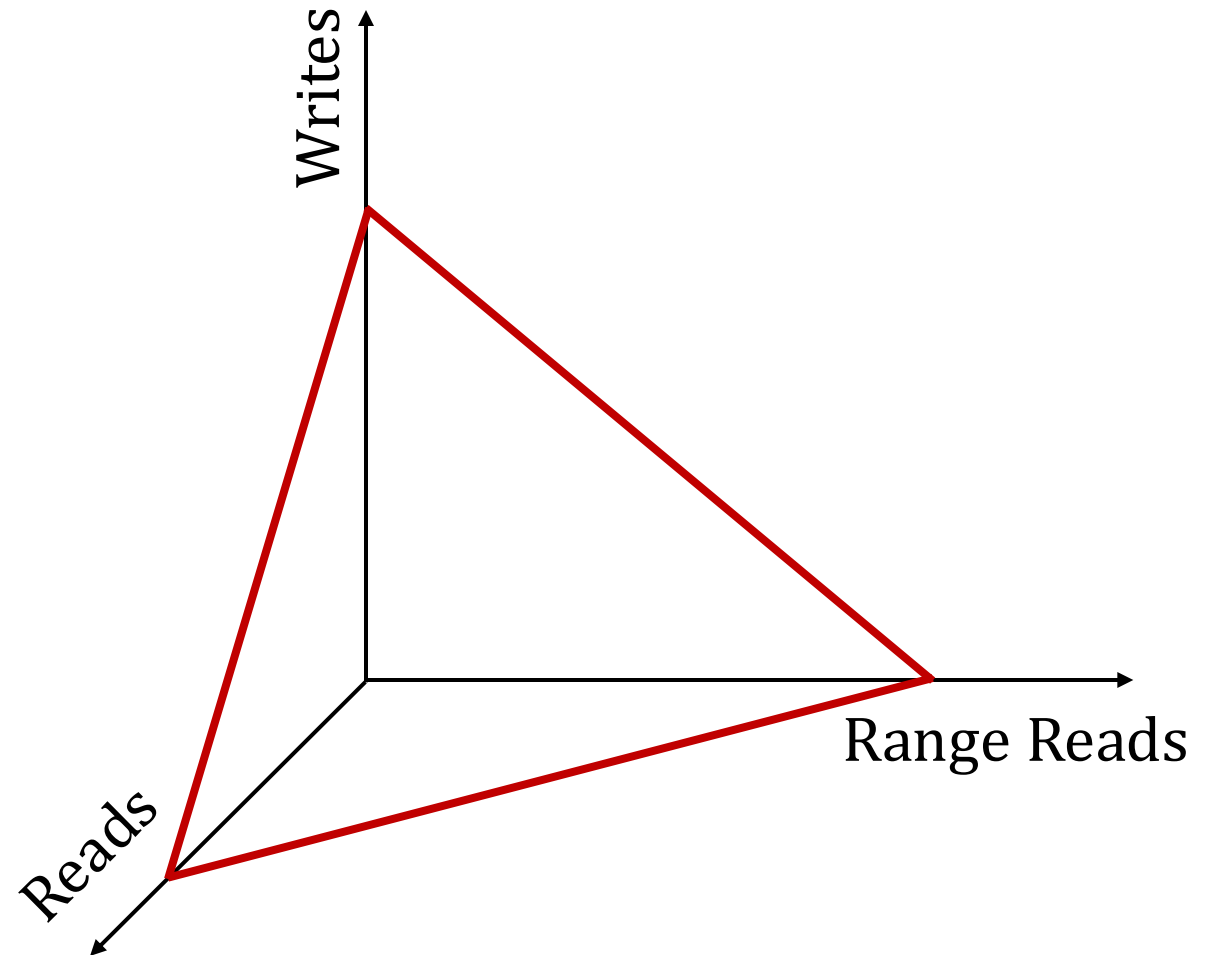
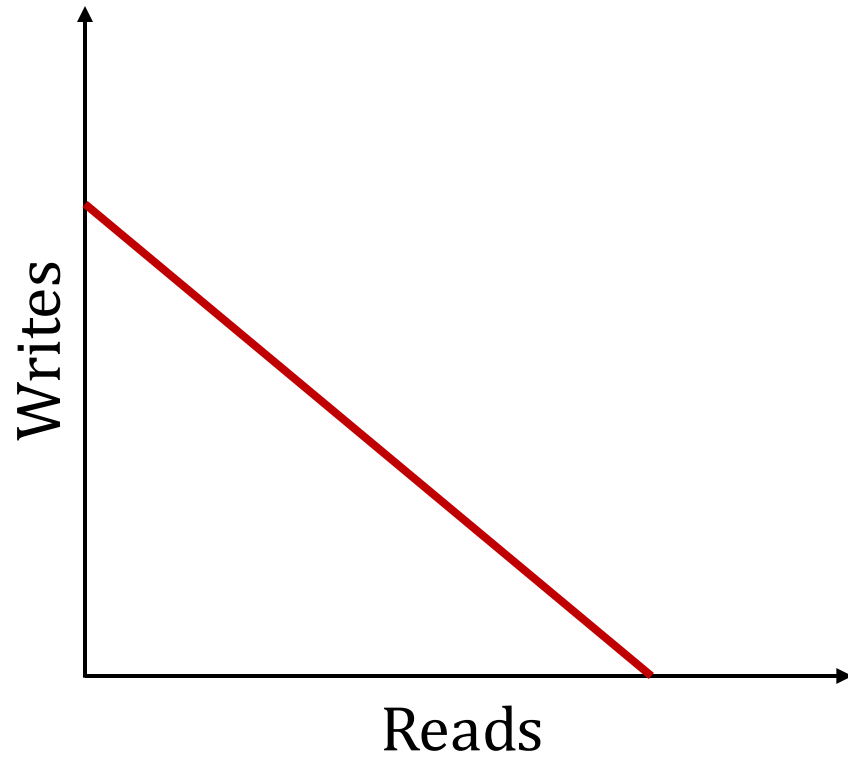


Dictates performance!



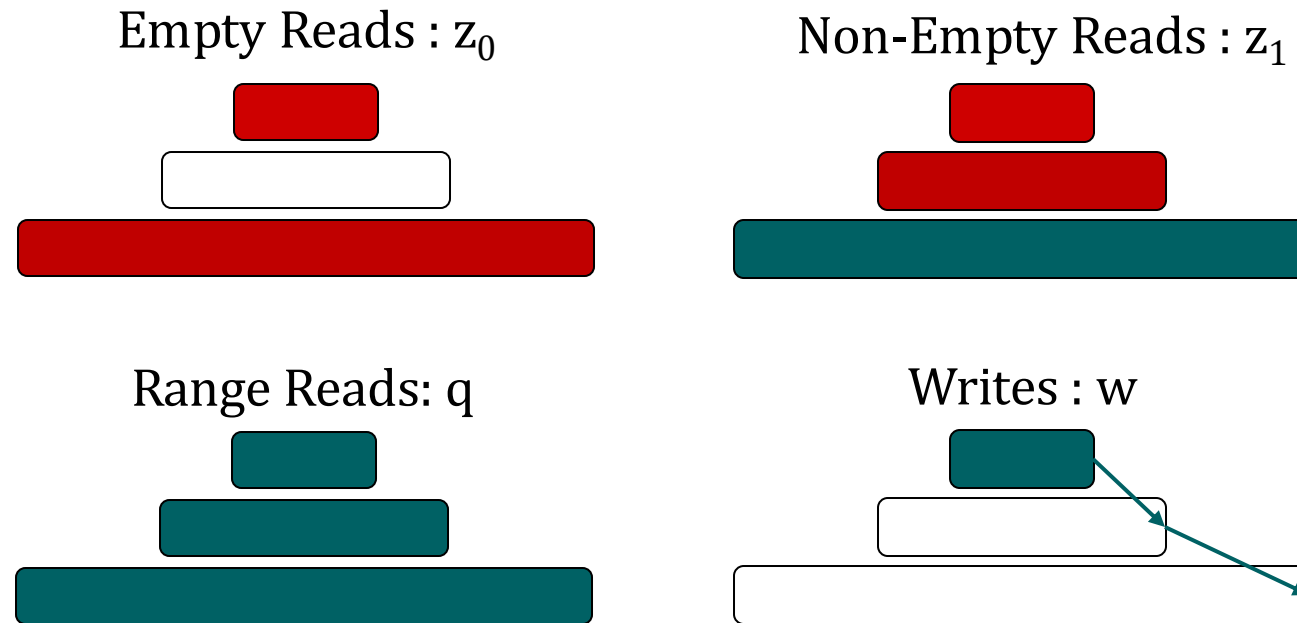
How do we go about tuning these knobs?

What is a Workload?



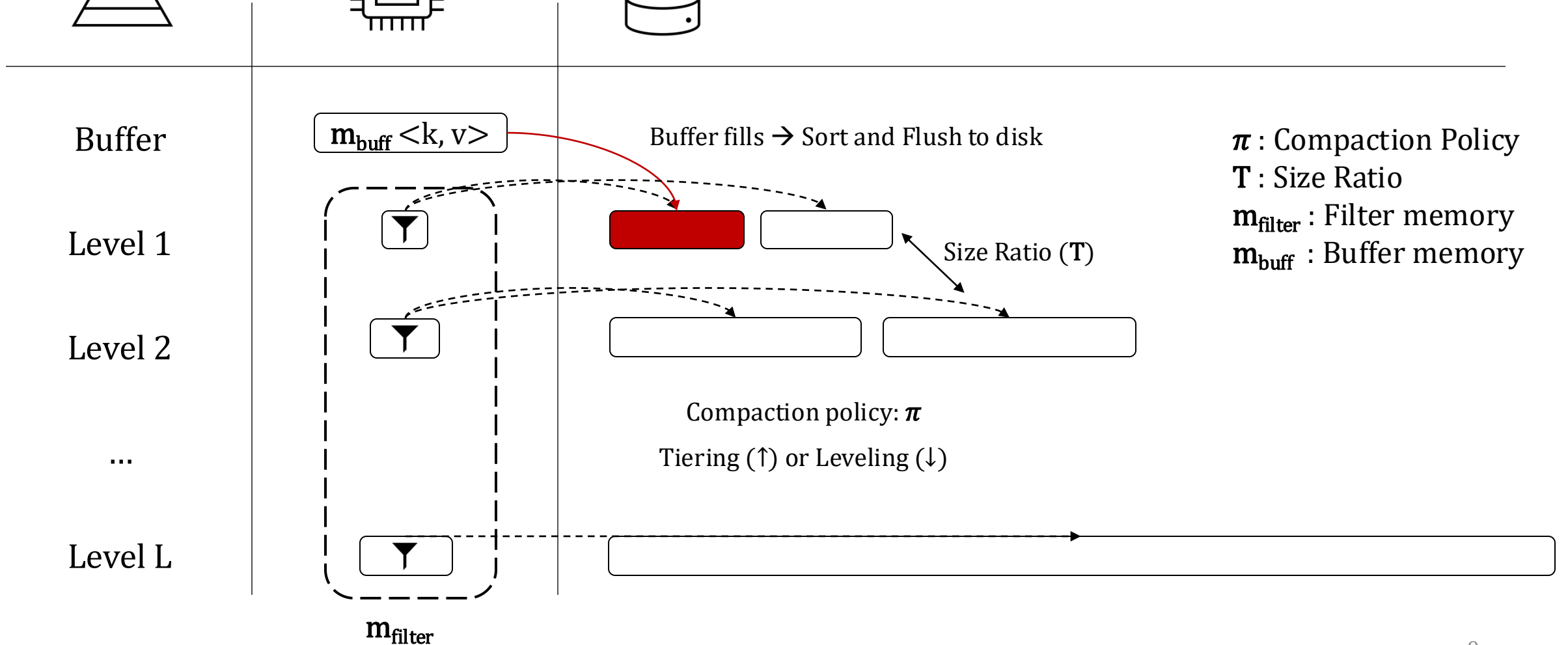
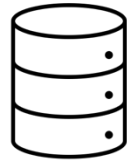
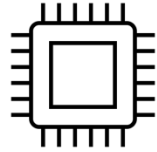
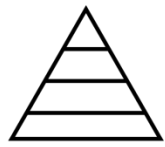
Query Types

Workload : (z_0, z_1, q, w)

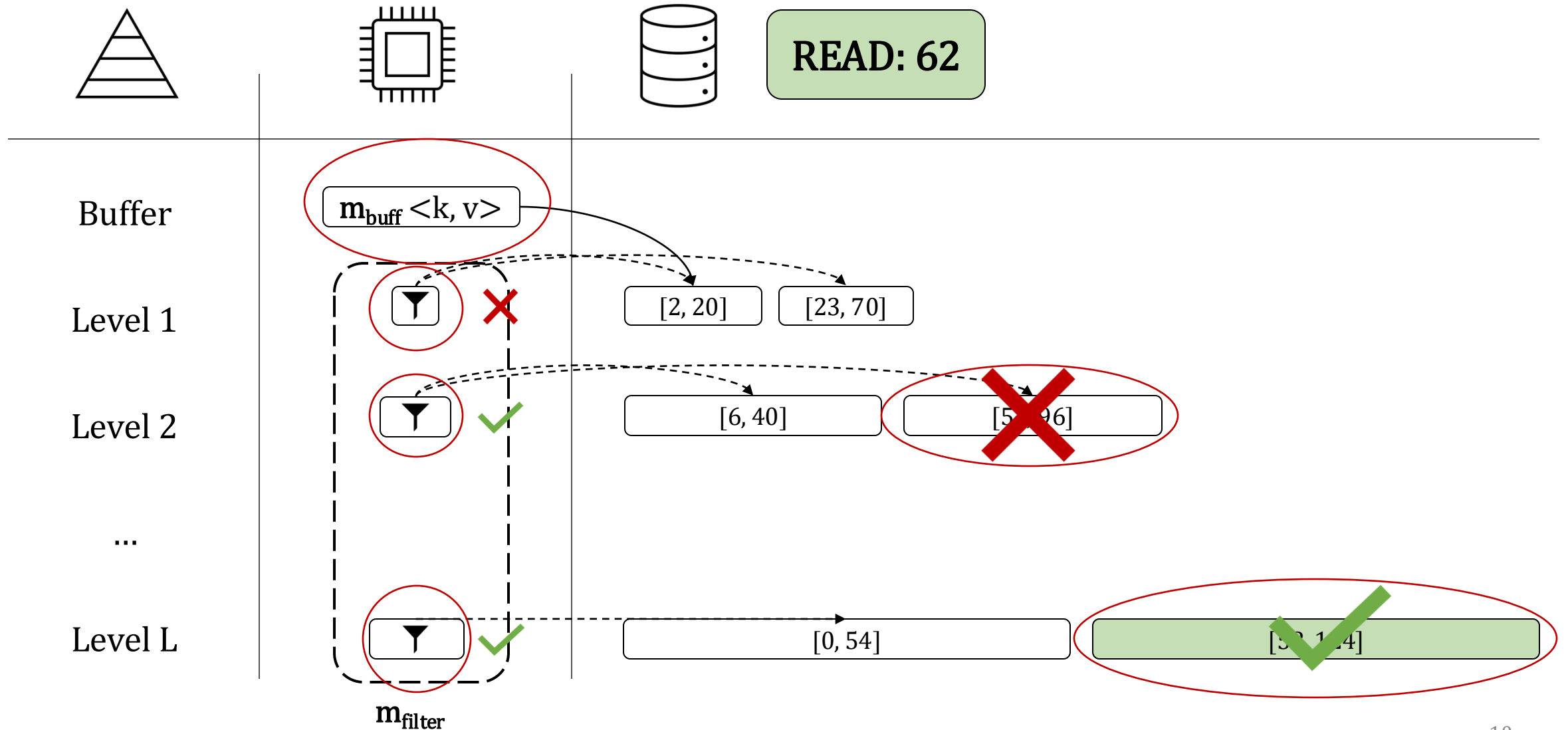


Cool! How do we go about tuning?

LSM Trees



LSM Trees – A Point Read



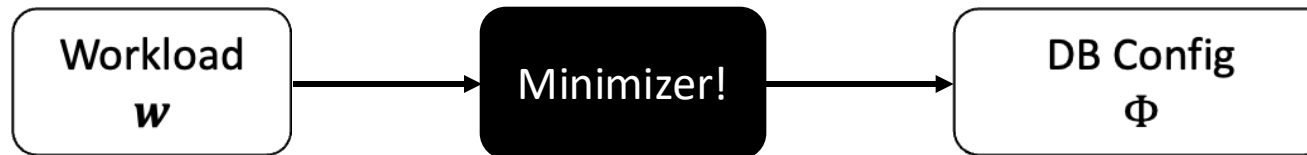
The LSM-Tuning Problem

\mathbf{w} : Workload (z_0, z_1, q, w)

Φ : LSM Tree Design $(m_{buffer}, m_{filter}, T, \pi)$

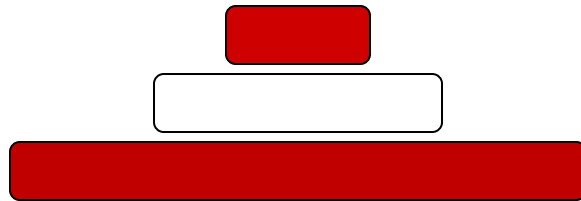
C : Cost

$$\Phi^* = \operatorname{argmin}_{\Phi} C(\mathbf{w}, \Phi)$$

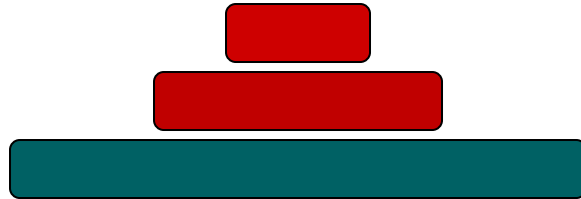


Point Reads

Empty Reads : z_0



Non-Empty Reads : z_1



Sum of
false
positives

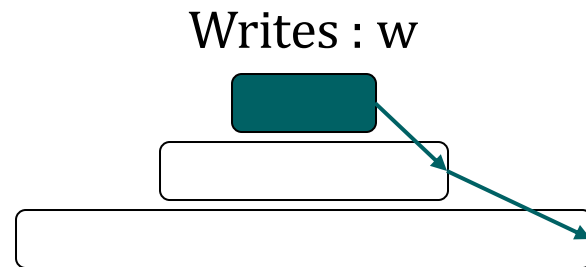
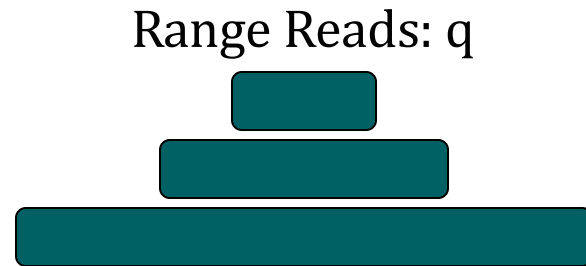
$$Z_0(\Phi) = \sum_{i=1}^L f_i$$

Probability query is
satisfied at level i

False positives
from levels above

$$Z_1(\Phi) = \sum_{i=1}^L \frac{T^{i-1} \cdot (T-1) \cdot m_{buf}}{N_f(T) \cdot E} \left(1 + \sum_{j=1}^{i-1} f_j \right)$$

Range-Reads and Writes



Sequential
read based on
selectivity

$$Q(\Phi) = S_{RQ} \cdot \frac{N}{B} + L \left. \vphantom{\frac{N}{B}} \right\} \begin{array}{l} 1 \text{ I/O per} \\ \text{Seek per} \\ \text{level} \end{array}$$

Average number of merges a write will participate in

$$W(\Phi) = \underbrace{\frac{L}{B} \cdot \frac{T-1}{2}}_{\text{Average number of merges}} \cdot (1 + A_{rw})$$

Writes only flush
once buffer is full

The LSM-Tuning Problem

\mathbf{w} : Workload (z_0, z_1, q, w)

Φ : LSM Tree Design $(m_{buffer}, m_{filter}, T, \pi)$

C : Cost (I/O)

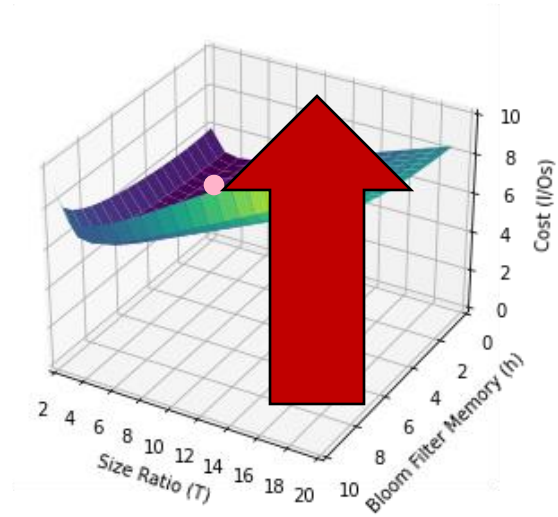
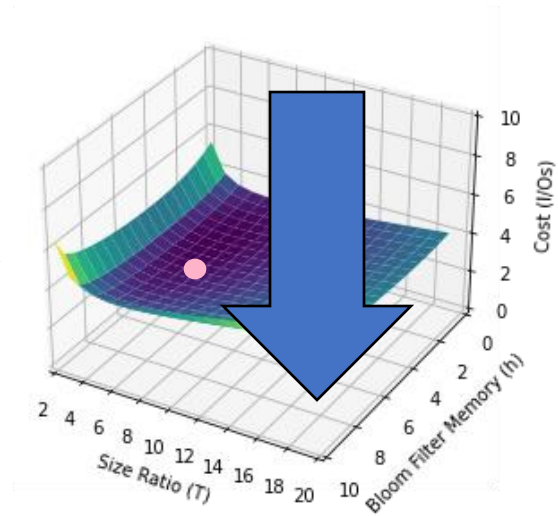
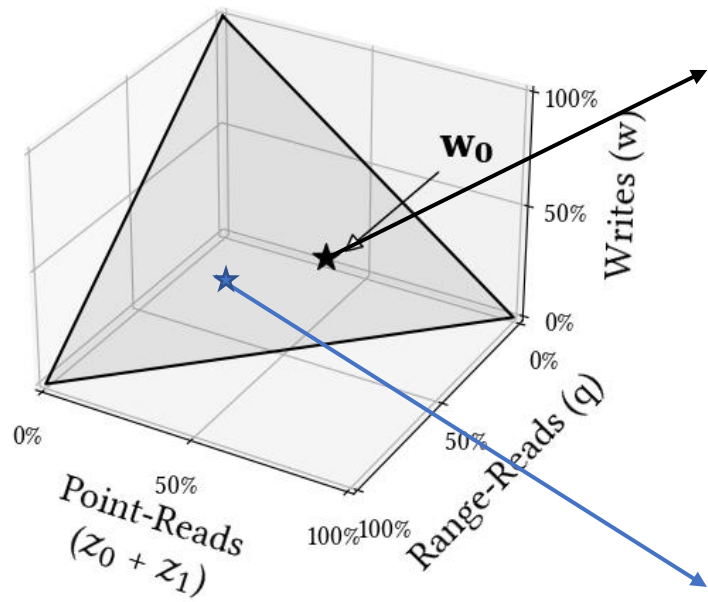
$$\Phi^* = \operatorname{argmin}_{\Phi} C(\mathbf{w}, \Phi)$$

Define our **cost function**

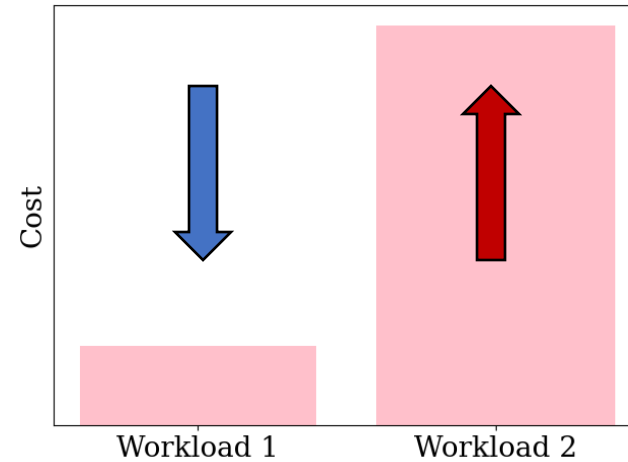
$$C(\hat{\mathbf{w}}, \Phi) = \hat{\mathbf{w}}^T \mathbf{c}(\Phi) = z_0 \cdot Z_0(\Phi) + z_1 \cdot Z_1(\Phi) + q \cdot Q(\Phi) + w \cdot W(\Phi)$$

Tuning Problems

w_0 : Workload (z_0, z_1, q, w)



● 'Best' Configuration



Optimal tuning depends on workload

Workload uncertainty leads to sub-optimal tuning

Outline

Introduction

LSM Trees Notation

Nominally Tuning LSM Trees

ENDURE: Robustly Tuning LSM Trees

The ENDURE Pipeline

ENDURE Evaluation

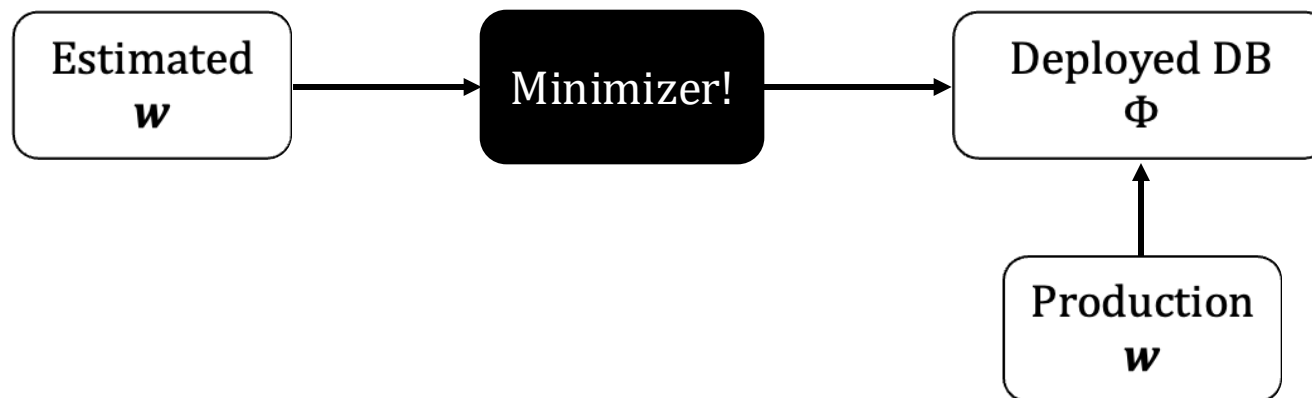
We've Got a Problem...

\mathbf{w} : Workload (z_0, z_1, q, w)

Φ : LSM Tree Design $(m_{buff}, m_{filter}, T, \pi)$

C : Cost (I/O)

$$\Phi^* = \operatorname{argmin}_{\Phi} C(\mathbf{w}, \Phi)$$



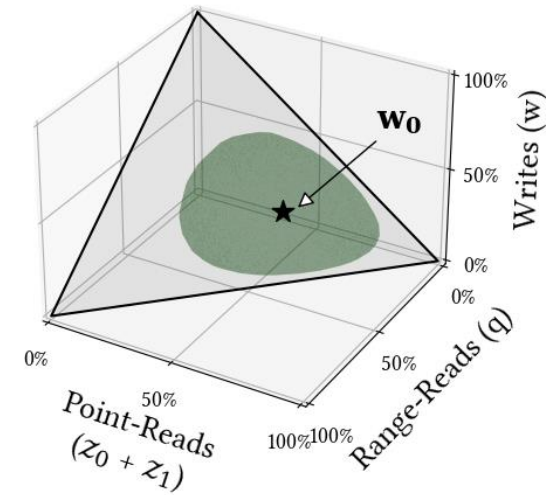
“It works on my machine!!!”

The LSM-Tuning Problem

\mathbf{w} : Workload (z_0, z_1, q, w)

Φ : LSM Tree Design $(m_{buff}, m_{filter}, T, \pi)$

C : Cost (I/O)



$$\Phi^* = \operatorname{argmin}_{\Phi} C(\mathbf{w}, \Phi)$$

Nominal

$U_{\mathbf{w}}^{\rho}$: Uncertainty Neighborhood of Workloads

ρ : Size of this neighborhood

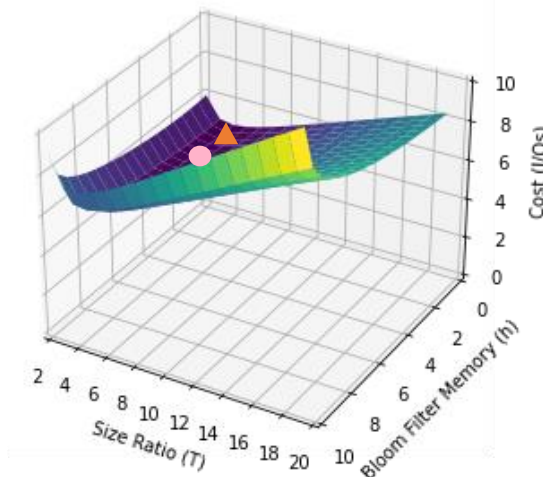
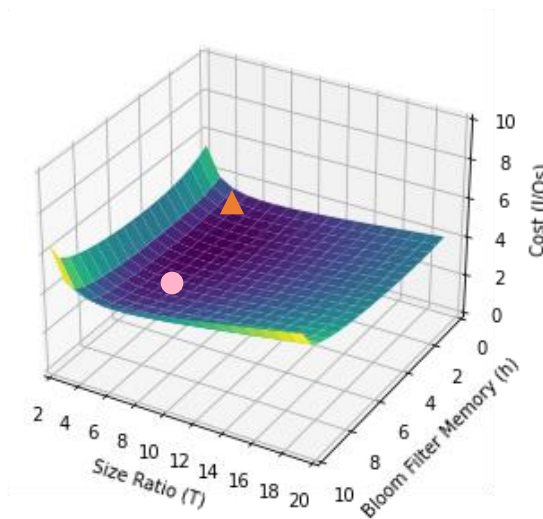
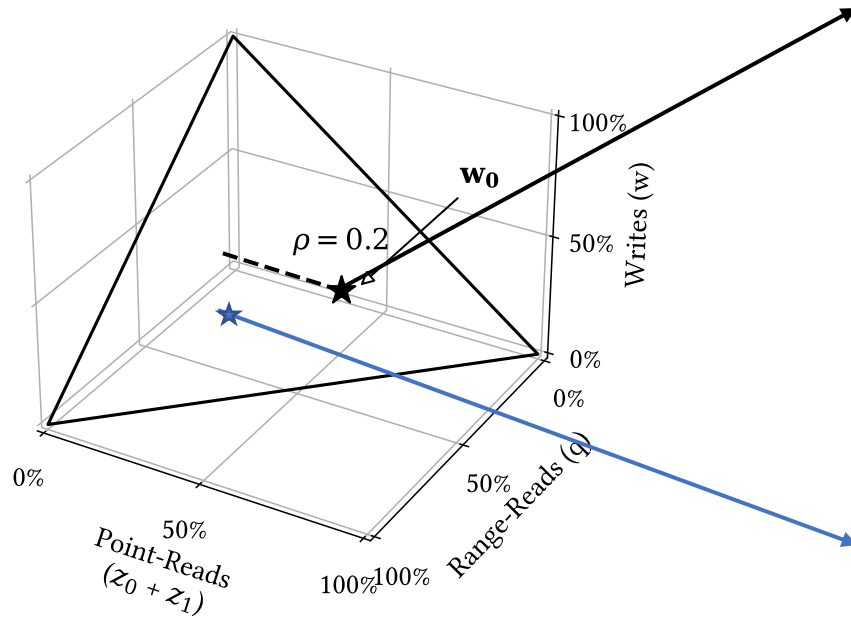
Robust

$$\Phi^* = \operatorname{argmin}_{\Phi} C(\hat{\mathbf{w}}, \Phi)$$

$$s. t., \quad \hat{\mathbf{w}} \in U_{\mathbf{w}}^{\rho}$$

Robust Tuning

w_0 : Workload (z_0, z_1, q, w)

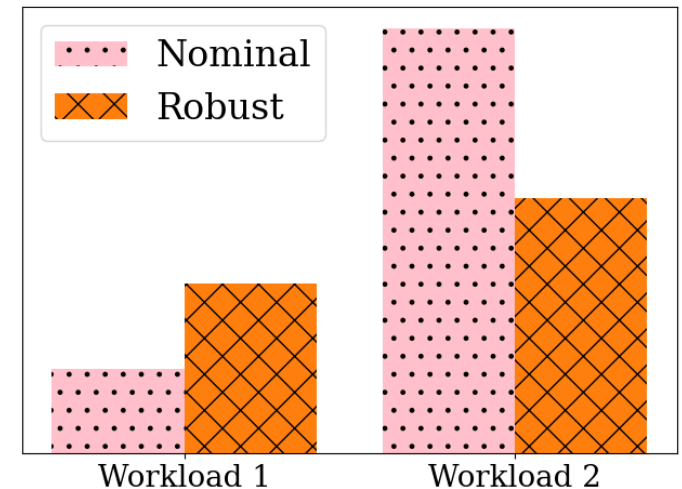


$$\Phi^* = \operatorname{argmin}_{\Phi} C(\hat{w}, \Phi)$$

$$s. t., \quad \hat{w} \in U_w^\rho$$

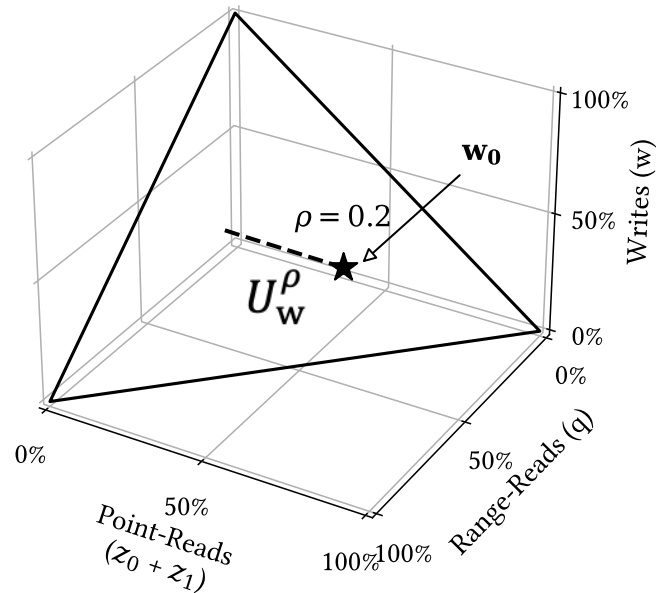
● Optimal configuration for expected workload

▲ Robust configuration for the workload neighborhood



Uncertainty Neighborhood

Workload Characteristic



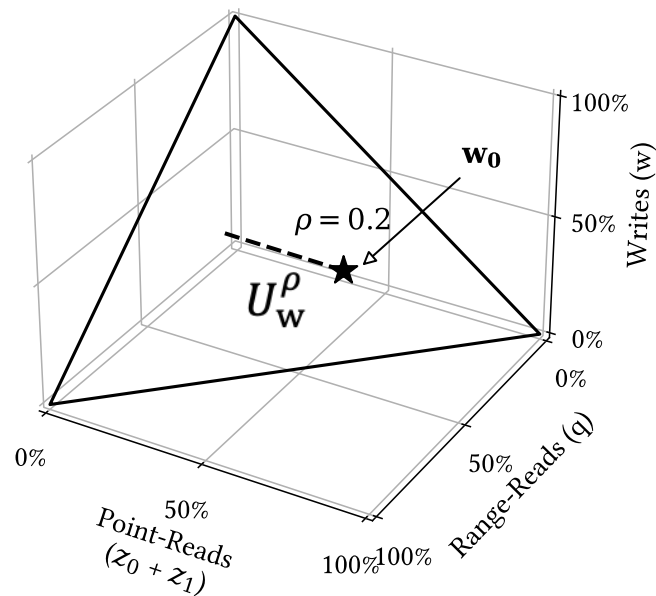
Neighborhood of workloads (ρ) via the KL-divergence

$$I_{KL}(\hat{w}, w) = \sum_{i=1}^m \hat{w}_i \cdot \log\left(\frac{\hat{w}_i}{w_i}\right)$$

U_w^ρ : Uncertainty Neighborhood of Workloads
 ρ : Size of this neighborhood

Calculating Neighborhood Size

Workload Characteristic



Historical workloads
maximum/average uncertainty among workload pairings

User provided workload uncertainty

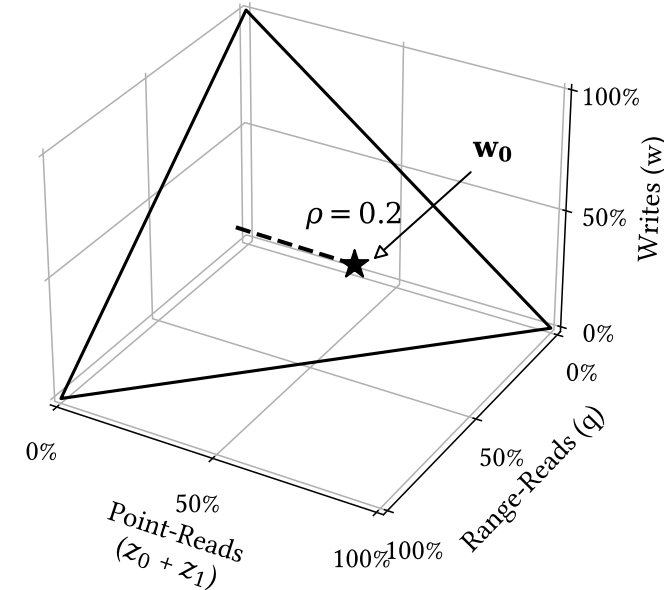
U_w^ρ : Uncertainty Neighborhood of Workloads
 ρ : Size of this neighborhood

Solving Robust Problem

Iterating over every possible workload is expensive

$$\Phi_R = \arg \min_{\Phi} \hat{\mathbf{w}}^\top \mathbf{c}(\Phi)$$

$$\text{s.t. } \hat{\mathbf{w}} \in \mathcal{U}_{\mathbf{w}}^{\rho}$$



Solving Robust Problem

Iterating over every possible workload is expensive

Rewrite as a min-max

Find the dual of the maximization problem to reduce to a feasible problem [2]

$$\Phi_R = \arg \min_{\Phi} \hat{\mathbf{w}}^\top \mathbf{c}(\Phi)$$

$$\text{s.t. } \hat{\mathbf{w}} \in \mathcal{U}_{\mathbf{w}}^{\rho}$$



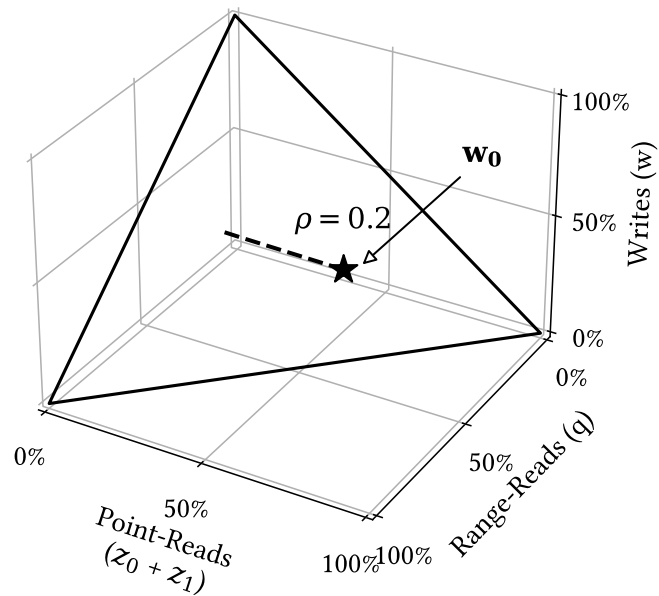
$$\min_{\Phi} \max_{\hat{\mathbf{w}} \in \mathcal{U}_{\mathbf{w}}^{\rho}} \hat{\mathbf{w}}^\top \mathbf{c}(\Phi)$$



$$\min_{\Phi, \lambda \geq 0, \eta} \left\{ \eta + \rho\lambda + \lambda \sum_{i=1}^m w_i \phi_{KL}^* \left(\frac{c_i(\Phi) - \eta}{\lambda} \right) \right\}$$

ENDURE Pipeline

Workload Characteristic

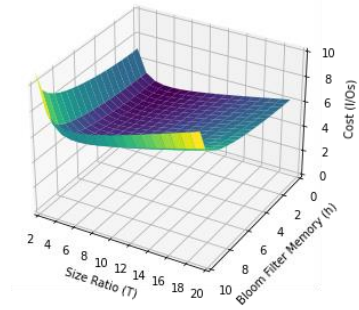


System Information

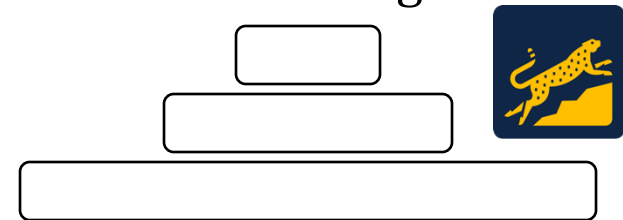
Page Size
Memory Budget

ENDURE
Solves the
Robust Problem

Expected performance



RocksDB Configuration



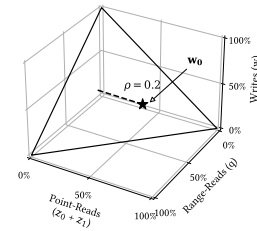
Testing Suite



ENDURE in Python, implemented in tandem with RocksDB

Uncertainty benchmark

- 15 expected workloads
- 10K randomly sampled workloads as a test-set



Normalized delta throughput

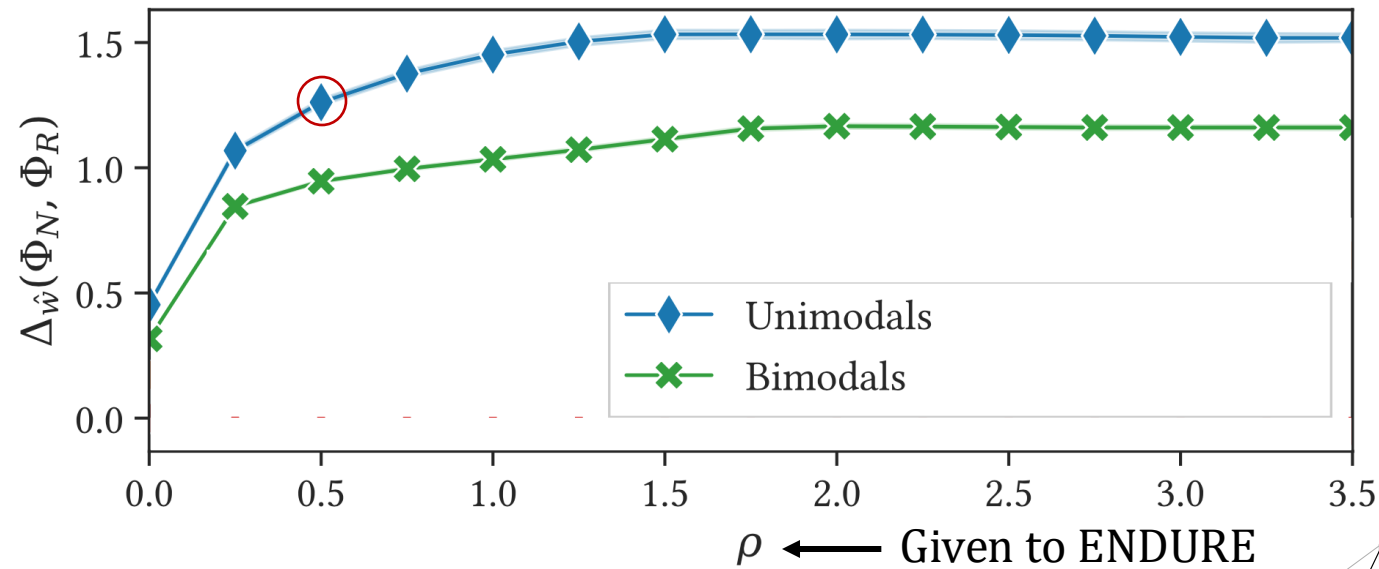
$$\Delta_{\mathbf{w}}(\Phi_1, \Phi_2) = \frac{1/C(\mathbf{w}, \Phi_2) - 1/C(\mathbf{w}, \Phi_1)}{1/C(\mathbf{w}, \Phi_1)}$$

Nominal vs Robust: > 0 is better

1 means 2x speedup

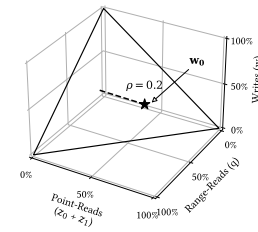
Index	(z_0, z_1, q, w)				Type
0	25%	25%	25%	25%	Uniform
1	97%	1%	1%	1%	Unimodal
2	1%	97%	1%	1%	
3	1%	1%	97%	1%	
4	1%	1%	1%	97%	
5	49%	49%	1%	1%	Bimodal
6	49%	1%	49%	1%	
7	49%	1%	1%	49%	
8	1%	49%	49%	1%	
9	1%	49%	1%	49%	
10	1%	1%	49%	49%	
11	33%	33%	33%	1%	Trimodal
12	33%	33%	1%	33%	
13	33%	1%	33%	33%	
14	1%	33%	33%	33%	

Impact of Workload Type

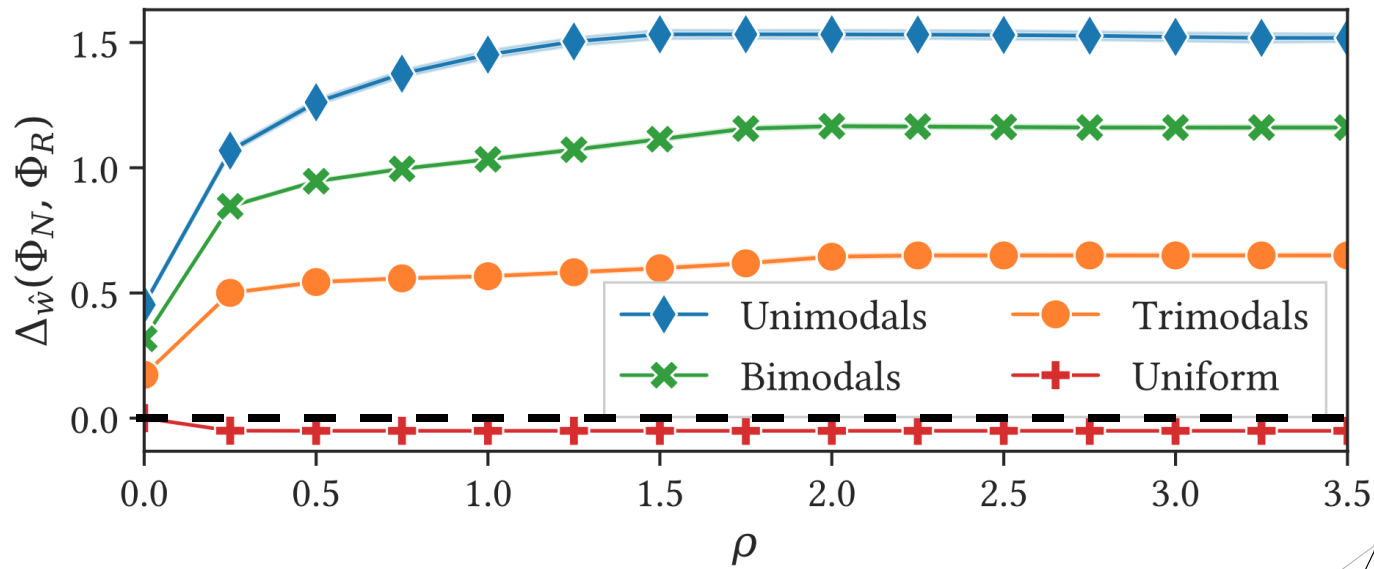


Index	(z_0, z_1, q, w)				Type
0	25%	25%	25%	25%	Uniform
1	97%	1%	1%	1%	Unimodal
2	1%	97%	1%	1%	
3	1%	1%	97%	1%	
4	1%	1%	1%	97%	
5	49%	49%	1%	1%	Bimodal
6	49%	1%	49%	1%	
7	49%	1%	1%	49%	
8	1%	49%	49%	1%	
9	1%	49%	1%	49%	
10	1%	1%	49%	49%	
11	33%	33%	33%	1%	Trimodal
12	33%	33%	1%	33%	
13	33%	1%	33%	33%	
14	1%	33%	33%	33%	

Unbalanced workloads result in overfitted nominal tunings



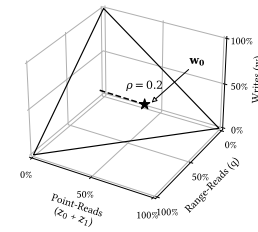
Impact of Workload Type



Index	(z_0, z_1, q, w)				Type
0	25%	25%	25%	25%	Uniform
1	97%	1%	1%	1%	Unimodal
2	1%	97%	1%	1%	
3	1%	1%	97%	1%	
4	1%	1%	1%	97%	
5	49%	49%	1%	1%	Bimodal
6	49%	1%	49%	1%	
7	49%	1%	1%	49%	
8	1%	49%	49%	1%	
9	1%	49%	1%	49%	
10	1%	1%	49%	49%	
11	33%	33%	33%	1%	Trimodal
12	33%	33%	1%	33%	
13	33%	1%	33%	33%	
14	1%	33%	33%	33%	

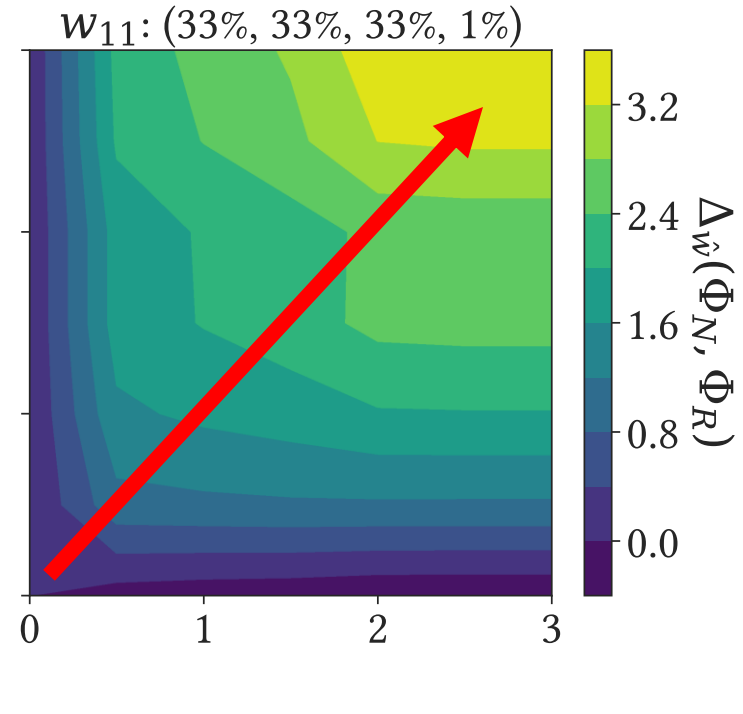
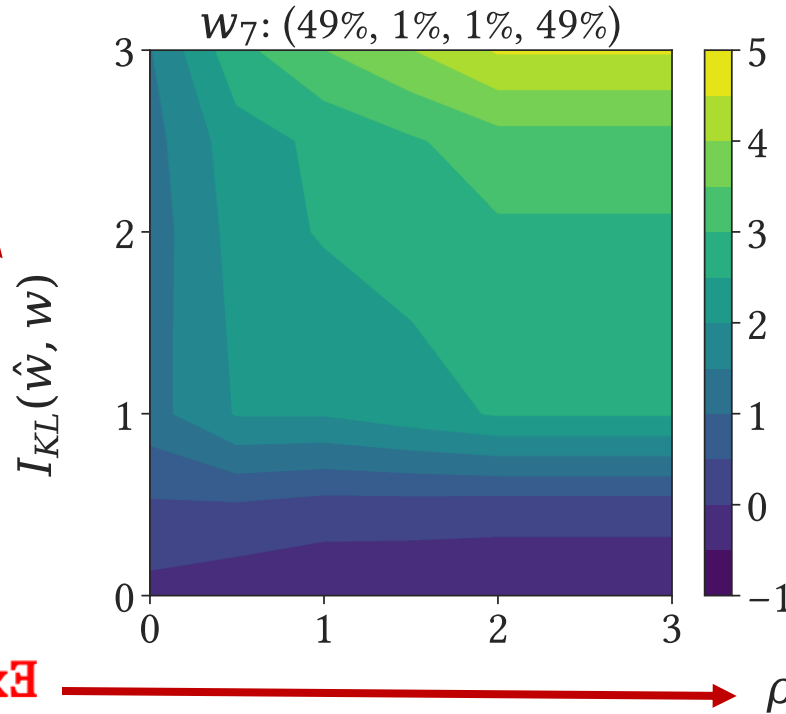
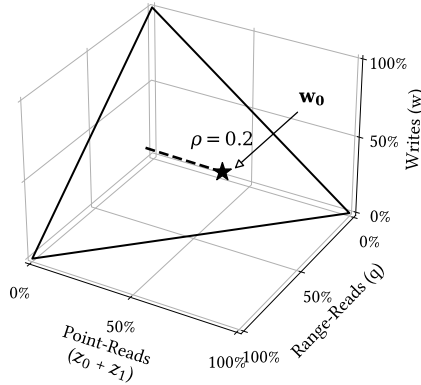
Unbalanced workloads result in overfitted nominal tunings

Tuning with uncertainty ($\rho > 0.5$) provides benefits



Relationship of Expected and Observed ρ

Observed ρ : distance from executed workload to expected workload

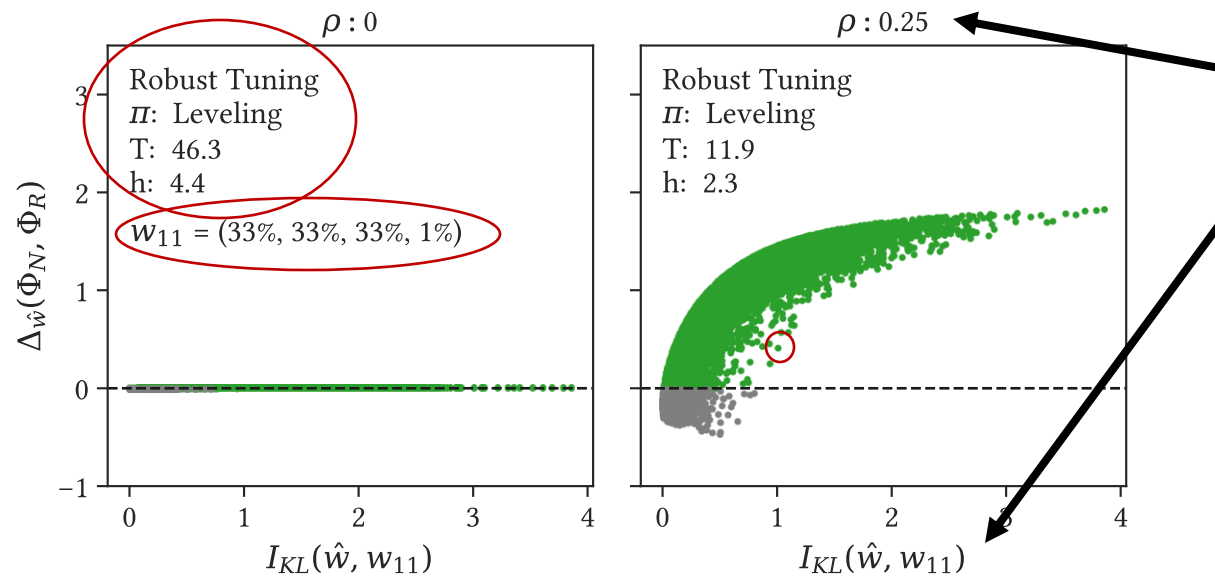


Expected ρ : workload given to tuner

Highest throughput when observed and expected ρ match

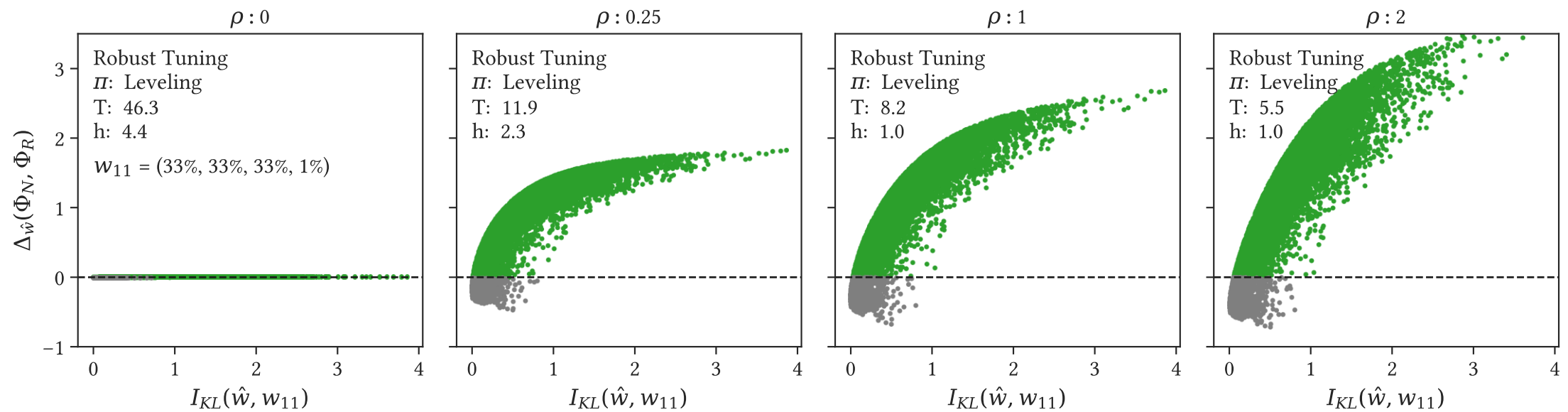
Lowest throughput when ρ is mismatched

Impact of Observed vs Expected ρ



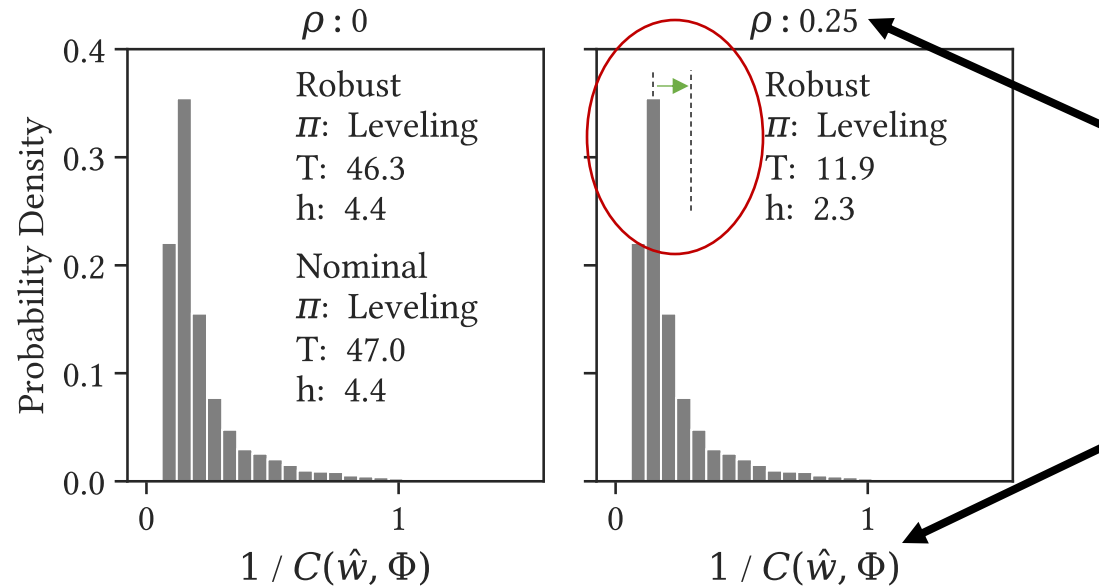
Expected ρ : given to tuner
 Observed ρ : distance from expected workload to expected workload

Impact of Observed vs Expected ρ



- Higher expected ρ accounts for more uncertainty
- Potential speed up of Δ
- Higher expected ρ \leftarrow anticipates writes \leftarrow shallow tree

ρ and Performance Gain Distribution

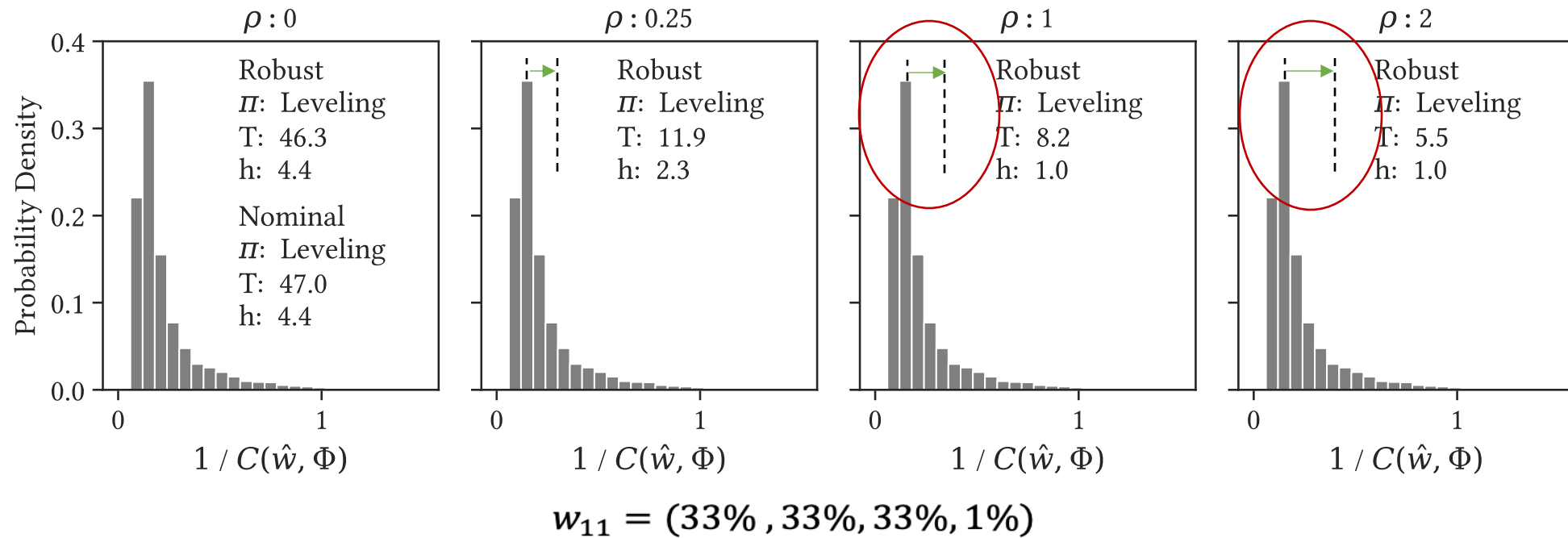


Expected ρ : given to tuner

Throughput

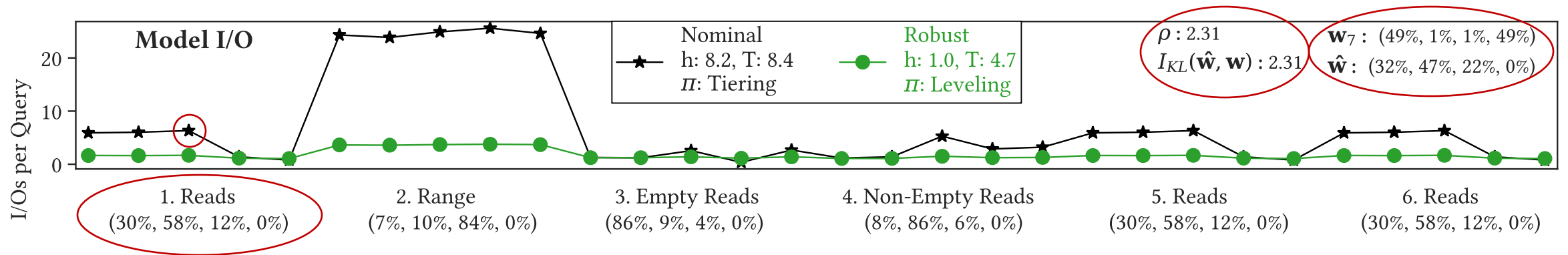
$$w_{11} = (33\%, 33\%, 33\%, 1\%)$$

ρ and Performance Gain Distribution



Peak of the distribution moves towards higher throughput as we consider higher uncertainty

Workload Sequence on RocksDB

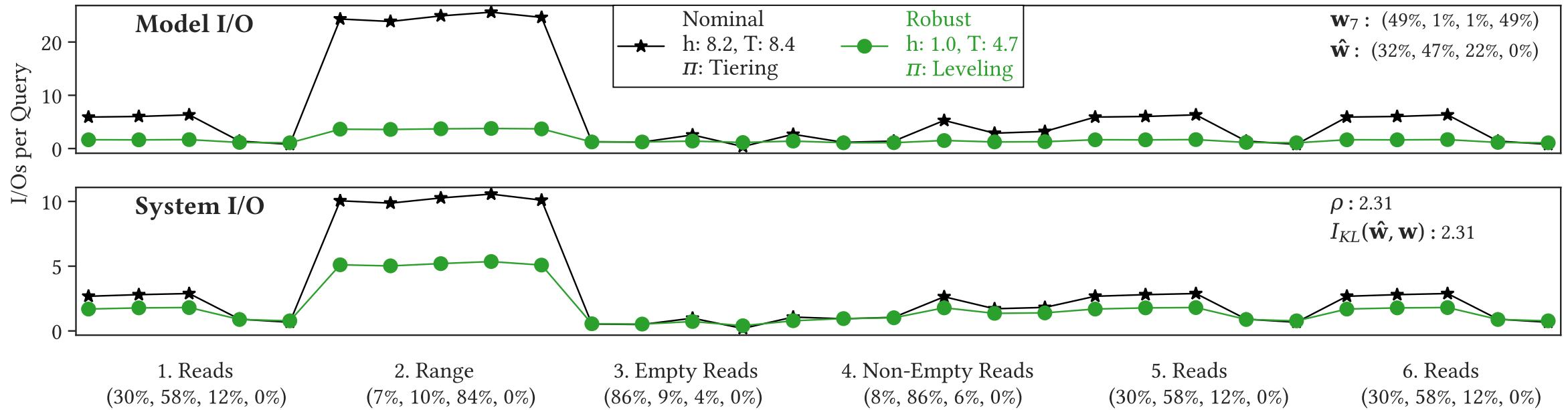


RocksDB instance setup with 10 million unique key-value pairs of size 1KB

Each observation period is 200K queries, with 5 observations per session 6 million queries to the DB

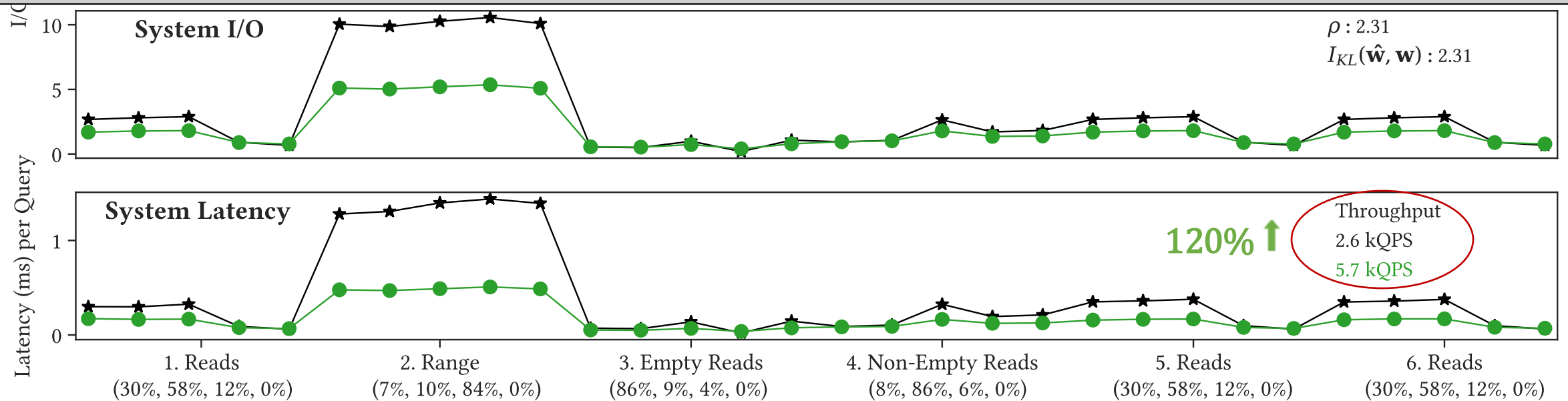
Writes are unique, range queries average 1-2 pages per level

Workload Sequence



Workload Sequence

Small subset of results! Take a look at the paper for a more detailed analysis



Thanks!

Workload uncertainty creates suboptimal tunings

ENDURE: robust tuning using neighborhood of workloads

Deployed ENDURE on RocksDB

disc.bu.edu/

www.ndhuynh.com/

@nd_huynh

