



Tuning Log-Structured Merge Trees

Andy Huynh

February 25th 2025

Serving Diverse Data Domains



Distributed IOT

Low Latency



HIPAA Compliant



NoSQL Distributed Database



LLM Vector Databases



Time Series



Messaging*



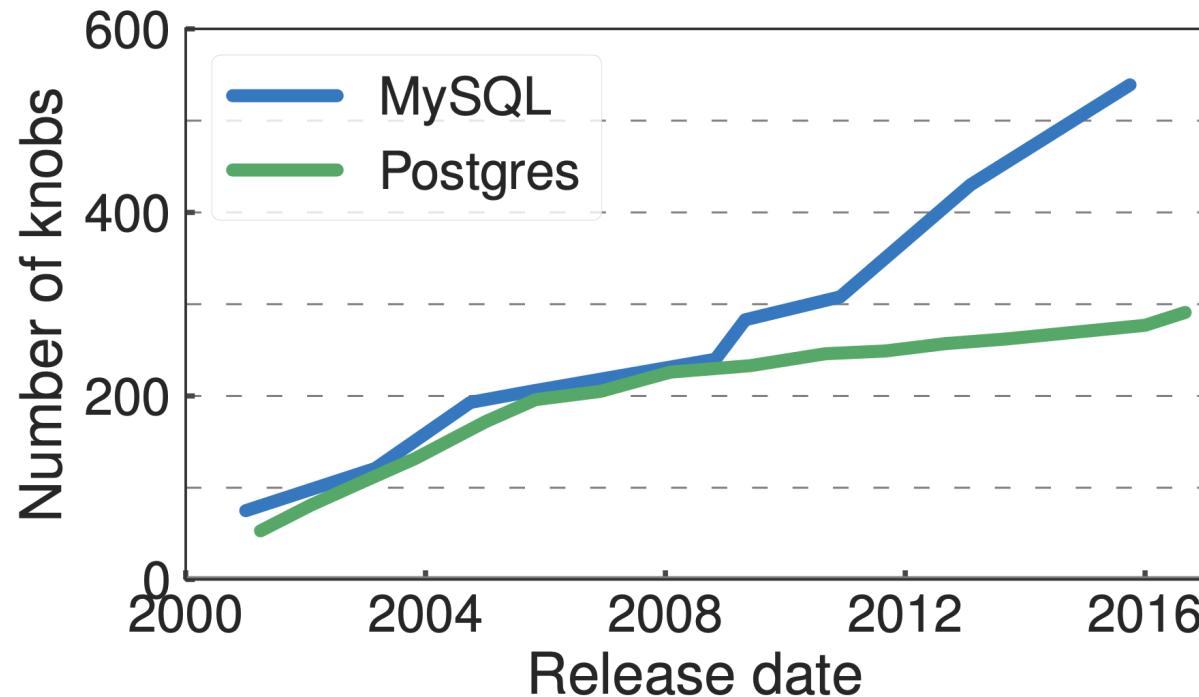
Flexibility Due To Configuration

NoSQL Distributed Database



- num_tokens
- max_hints_delivery_threads
- batchlog_endpoint_strategy
- credentials_validity
- row_cache_size
- column_index_size
- column_index_cache_size
- ...
- memtable_flush_writers

Database Complexity



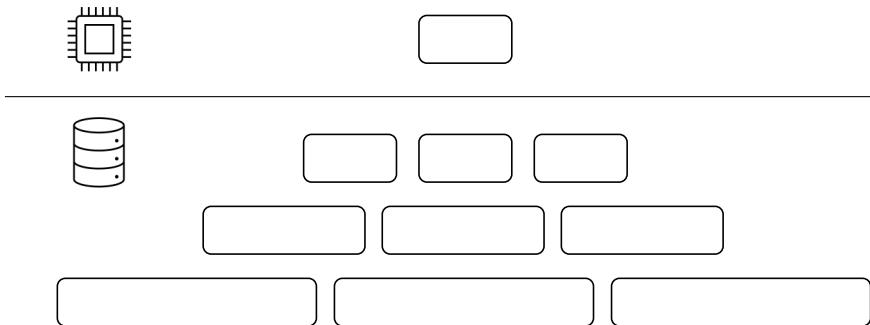
Van Aken D. et. al., "Automatic Database Management System Tuning Through Large-scale Machine Learning". SIGMOD 2017

LSM Trees Are Everywhere!

NoSQL Distributed Database



Log-Structured Merge-Trees



Outline

A Primer on LSM Trees

Flexibility in Compaction Design [VLDB-J 24]

Modeling LSM Trees

Tuning LSM Trees

ENDURE: Finding Robust Tunings for LSM Trees [VLDB 22]

AXE: Learning to Tune LSM Trees By Task Decomposition [UR*]

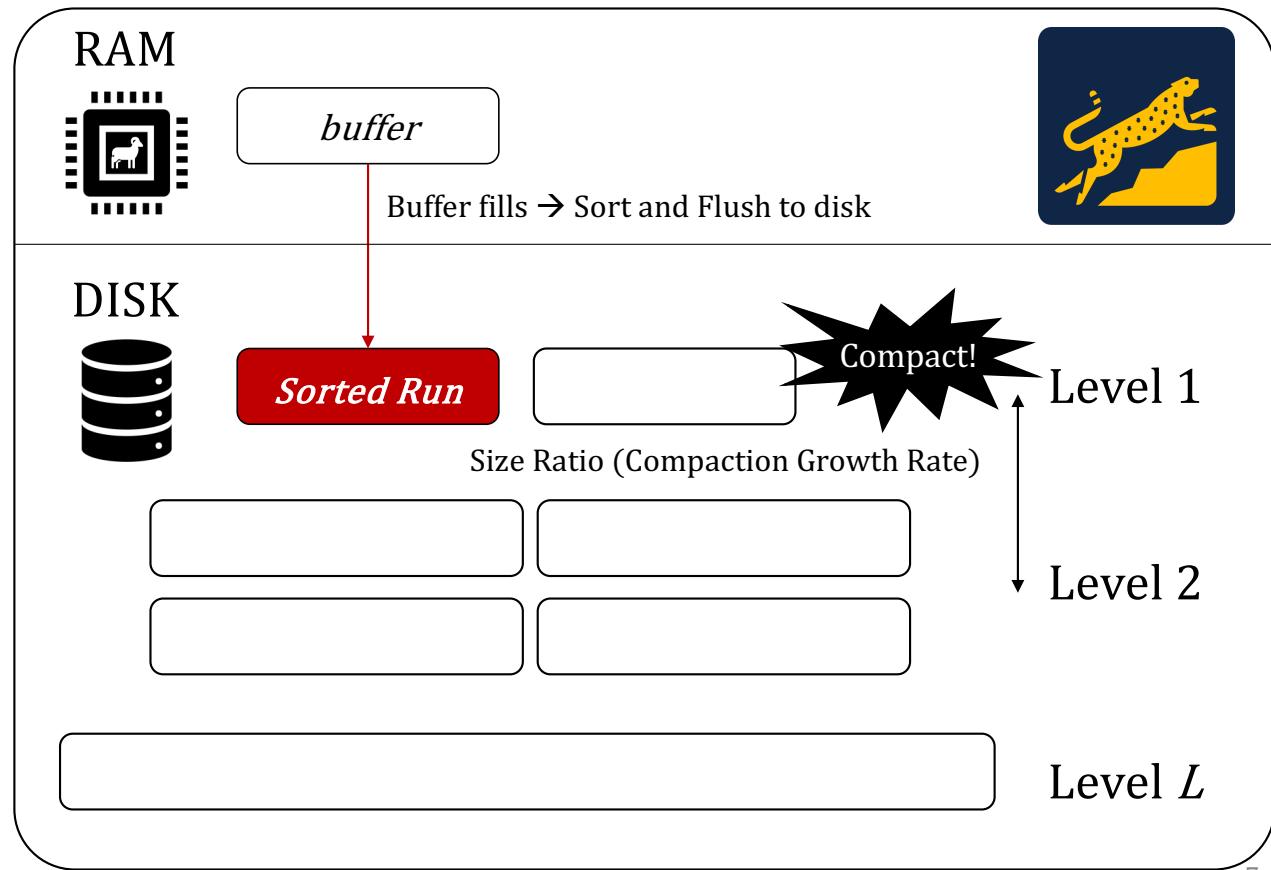
Primer: Log-Structured Merge-Trees

Application

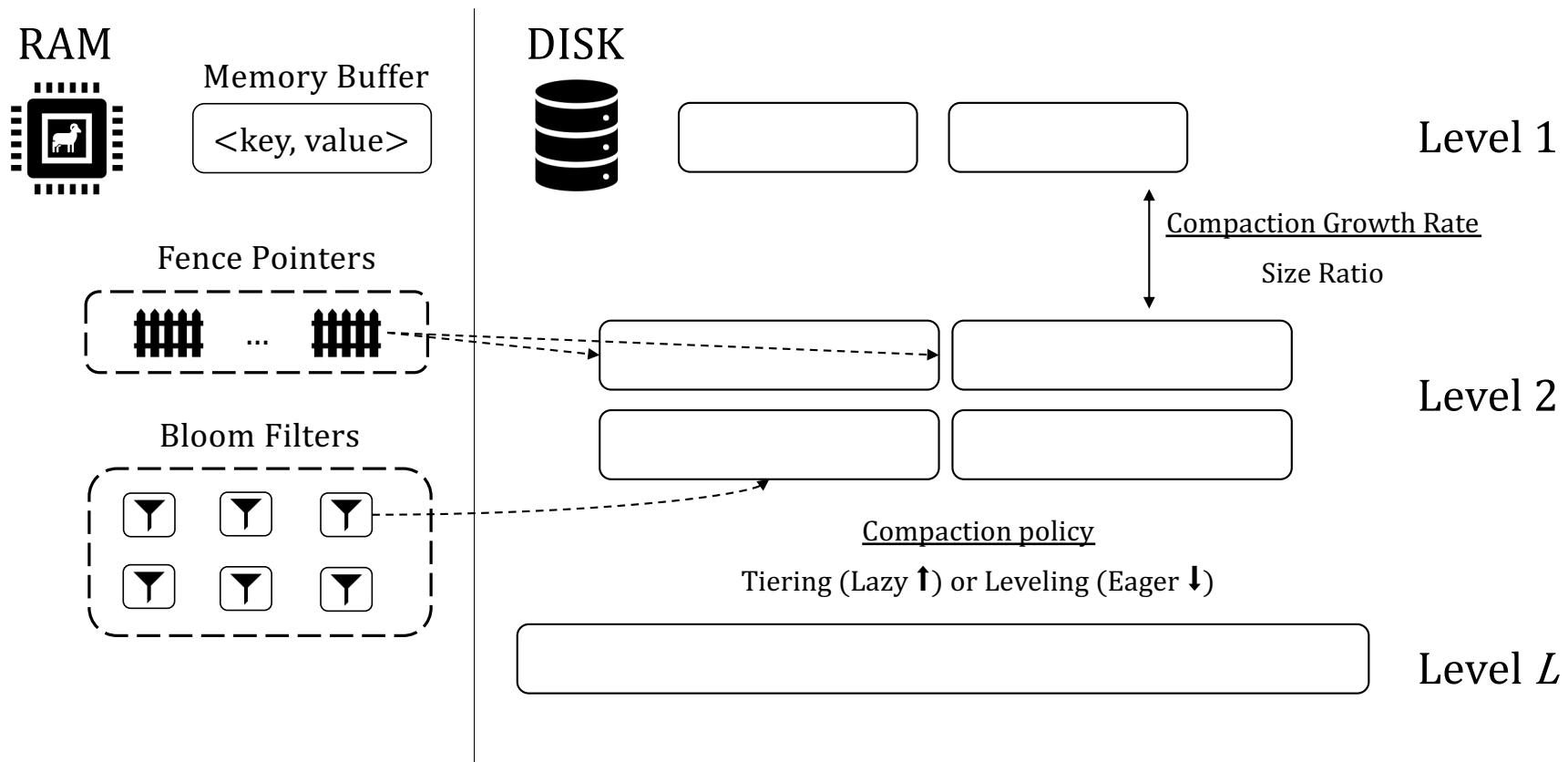


<key, value>
↔

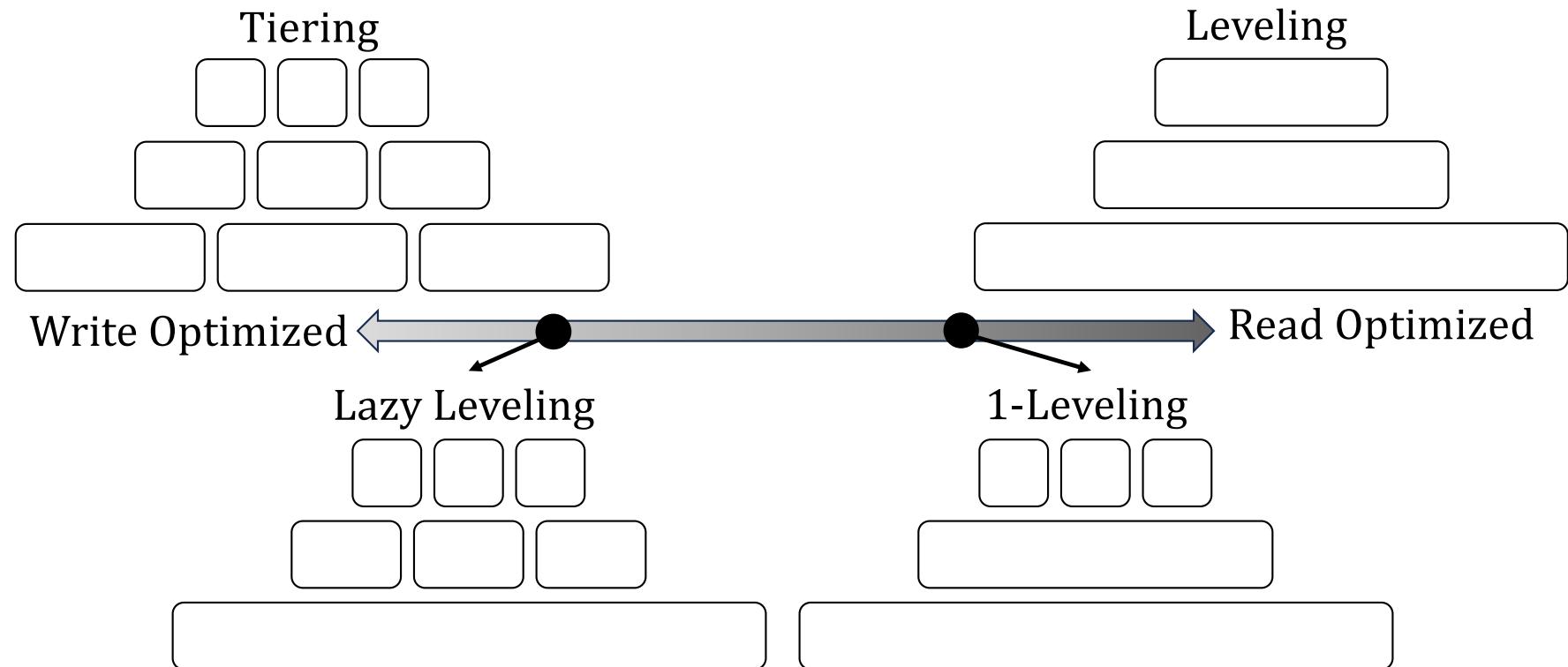
GET 25
PUT 87, DATA
SCAN 10 – 275
PUT 86, 
GET 99
GET 47
...
PUT 25



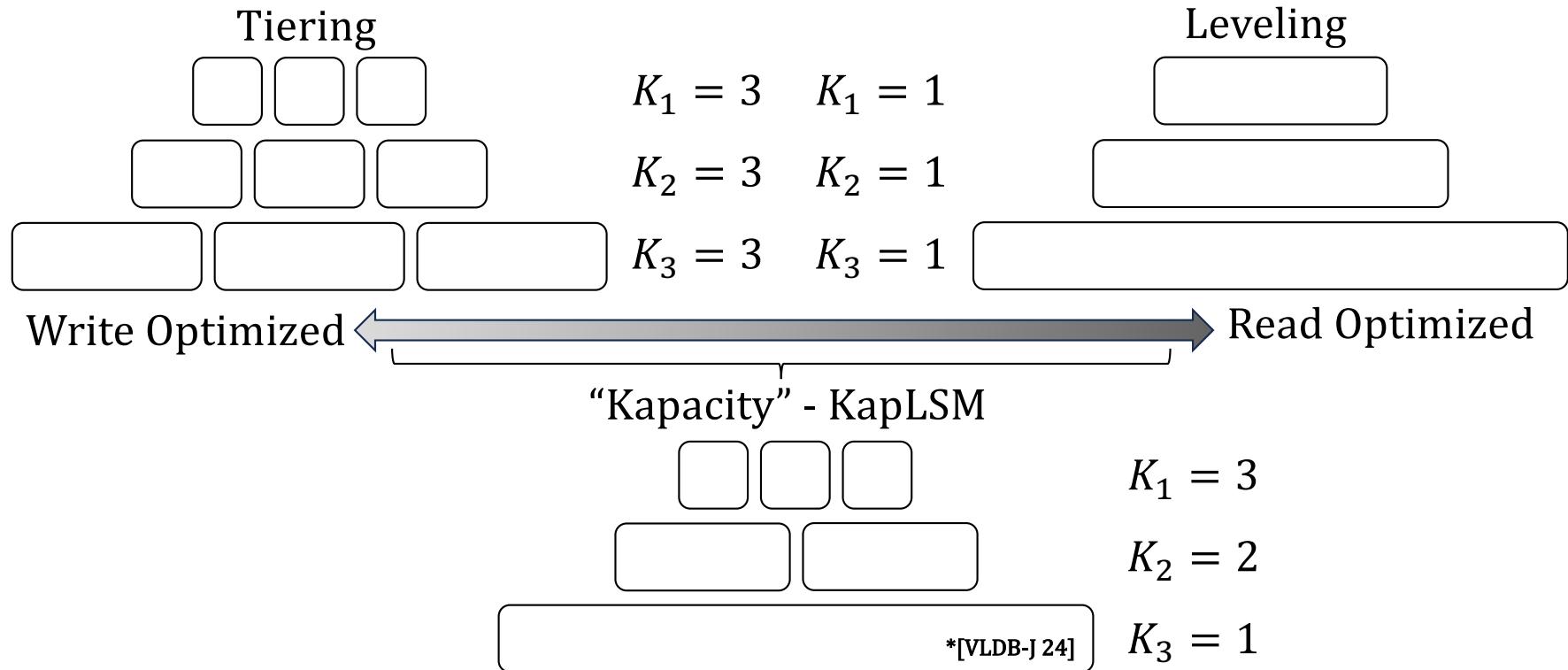
Primer: Log-Structured Merge-Trees



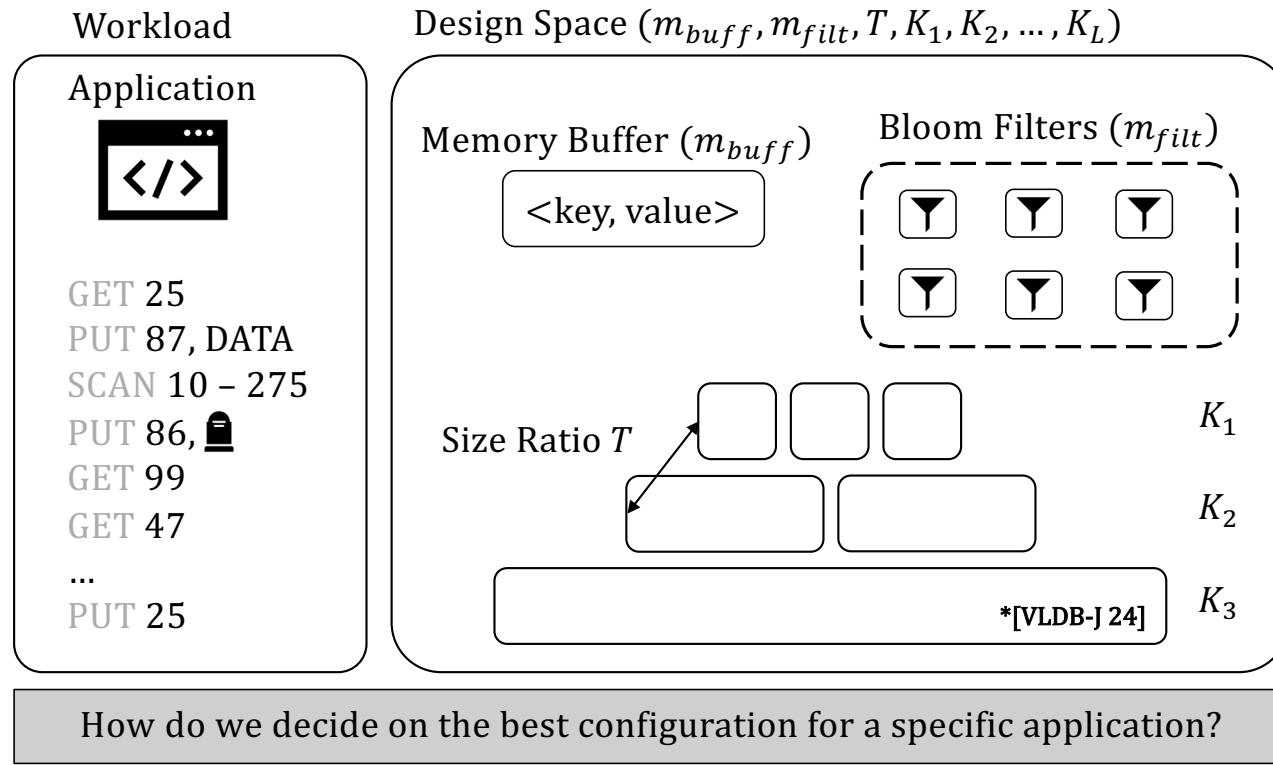
Compaction Design



Unifying Compaction Design Language



Navigating Flexibility



The Tuning Problem

w : Workload

φ : Design $(m_{buff}, m_{filter}, T, K_1, \dots, K_L)$

C : Cost

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

Automatic Database Management System Tuning Through
Large-scale Machine Learning

CAAMA: Optimizing LSM-trees via Active Learning
Rover: An Online Spark SQL Tuning Service via Generalized

Yu Shen*
School of CS, Peking University
ByteDance Inc.
Beijing, China
shenyu@pku.edu.cn

Huaixun Jiang
Center for Data Science, Peking
University
ByteDance Inc.
Beijing, China
jianghuaixun@pku.edu.cn

Yang Li
School of CS, Peking University
Beijing, China
liyang.cs@pku.edu.cn

CAIHUA SHAN, Microsoft Research Asia, China

Transfer Learning

**Efficient Deep Learning Pipelines for Accurate Cost Estimations
Over Large Scale Query Workload**

Johan Kok Zhi Kang
johan.kok@u.nus.edu
National University of Singapore
Singapore

Feng Cheng
feng.cheng@grab.com
GrabTaxi Holdings
Singapore

wentao.zhang@mila.quebec

Gaurav
gaurav@grab.com
GrabTaxi Holdings
Singapore

Shixuan Sun
sunxs@comp.nus.edu.sg
National University of Singapore
Singapore

INSTITUTE OF COMPUTATIONAL SCIENCE

Science, Peking University (Qingdao)

Beijing, China

bin.cui@pku.edu.cn

Monkey: Optimal Navigable Key-Value Store

Sien Yi Tan
sienyi.tan@grab.com
GrabTaxi Holdings
Singapore

Bingsheng He
hebs@comp.nus.edu.sg
National University of Singapore
Singapore

Wilson Qin
Harvard University
wilson@seas.harvard.edu

infeasible Stratos Idreos

UI
OS]

Surajit Chaudhuri

Vivek Narasayya
Microsoft Research, One Microsoft Way, Redmond, WA, 98052.
{surajitc, viveknar}@microsoft.com

QTune: A Query-Aware Database Tuning System with Deep

An Efficient, Cost-Driven Index Selection Tool for Microsoft SQL Server

Dana Van Aken
Carnegie Mellon University
dvanaken@cs.cmu.edu

Andrew Pavlo
Carnegie Mellon University
pavlo@cs.cmu.edu

Geoffrey J. Gordon
Carnegie Mellon University
ggordon@cs.cmu.edu

Bohan Zhang
Peking University
bohan@pku.edu.cn

LlamaTune: Sample-Efficient DBMS Configuration Tuning

Meena Jagadeesan
Harvard University
mjagadeesan@seas.harvard.edu

Stratos Idreos
Harvard University
stratos@seas.harvard.edu

Andreas Müller
Microsoft Gray Systems Lab
amueller@microsoft.com

Konstantinos Kanellis
University of Wisconsin-Madison
kkanelis@cs.wisc.edu

Cong Ding
University of Wisconsin-Madison
cding@cs.wisc.edu

Carlo Curino
Microsoft Gray Systems Lab
ccurino@microsoft.com

Company
wei.com

Brian Kroth
Microsoft Gray Systems Lab
bpkroth@microsoft.com

Shivaram Venkataraman
University of Wisconsin-Madison
shivaram@cs.wisc.edu

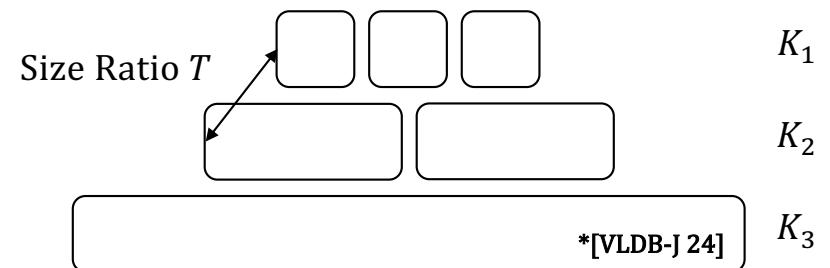
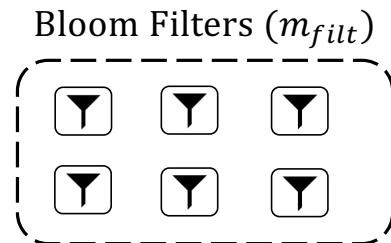
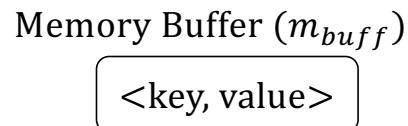
The LSM Tuning Problem

w : Workload

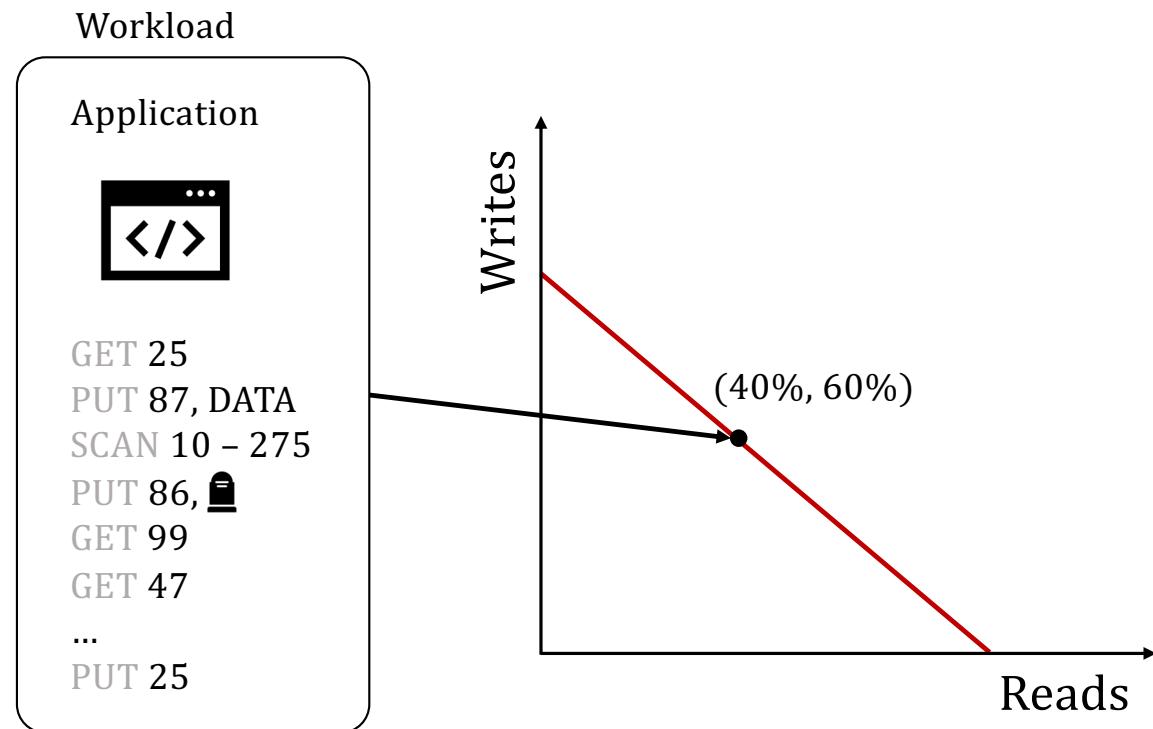
φ : Design $(m_{buff}, m_{filter}, T, K_1, \dots, K_L)$

C : Cost

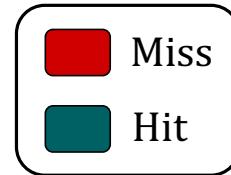
$$\varphi^* = \operatorname{argmin}_{\varphi} C(w, \varphi)$$



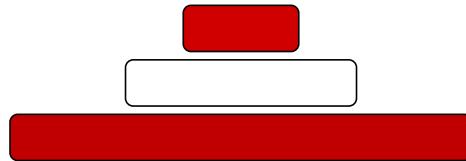
Workloads



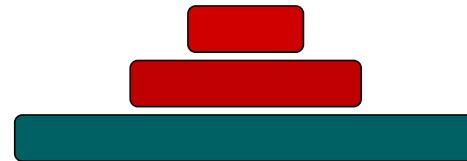
Query Types



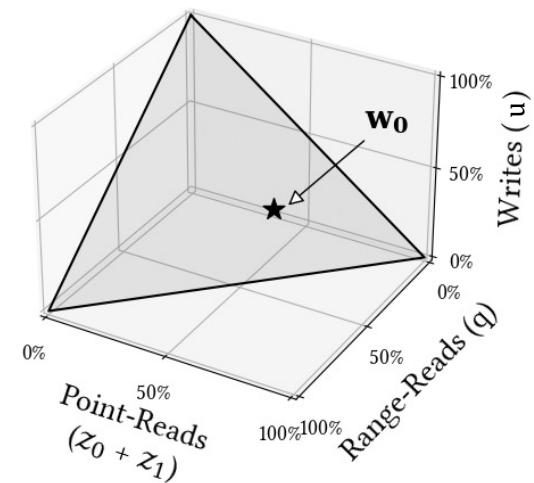
Empty Reads : z_0



Non-Empty Reads : z_1



Workload : (z_0, z_1, q, u)



Range Reads: q

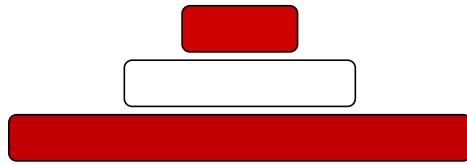


Writes : u



Point Reads

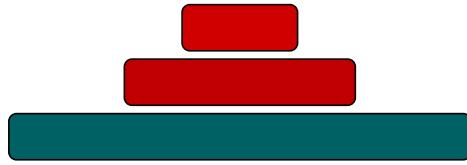
Empty Reads : z_0



Sum of false positives

$$Z_0(\varphi) = \sum_{i=1}^{L(T)} K_i \cdot f_i(T)$$

Non-Empty Reads : z_1



Probability query is satisfied at Level i

$$Z_1(\varphi) = \sum_{i=1}^{L(T)} \frac{(T-1) \cdot T^{i-1}}{N_f(T)} \cdot \frac{m_{buf}}{E} \left(1 + \sum_{j=1}^{i-1} K_j \cdot f_j(T) + \frac{K_{i-1}}{2} \cdot f_i(T) \right)$$

False positives from levels above

Additional Notation	
E	Num entries in 1 page
$L(T)$	Total levels
$f_i(T)$	BF false positive rate at level i
$N_f(T)$	Total keys if tree was full

Range-Reads and Writes

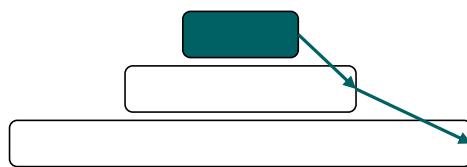
Range Reads: q



Sequential read based on selectivity

$$Q(\varphi) = \left[f_{seq} \cdot \overbrace{S_{RQ}}^{\text{Range query selectivity}} \cdot \frac{N}{B} + \sum_{i=1}^{L(T)} K_i \right] \text{ 1 I/O per Seek per level}$$

Writes : u



Average number of merges a write will participate in

$$U(\varphi) = \left[f_{seq} \cdot \overbrace{\frac{1+f_a}{B}}^{\text{Average number of merges}} \cdot \sum_{i=1}^{L(T)} \overbrace{\frac{T-1+K_i}{2K_i}}^{\text{Average number of merges per level}} \right]$$

Writes only flush once buffer is full

Additional Notation	
B	Num. entries in buffer
$L(T)$	Total levels
S_{RQ}	Range query selectivity
f_{seq}	Sequential access cost

The LSM Tuning Problem

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

φ : Design $(m_{buff}, m_{filter}, T, K_1, K_2, \dots, K_L)$

C : Cost (I/O)

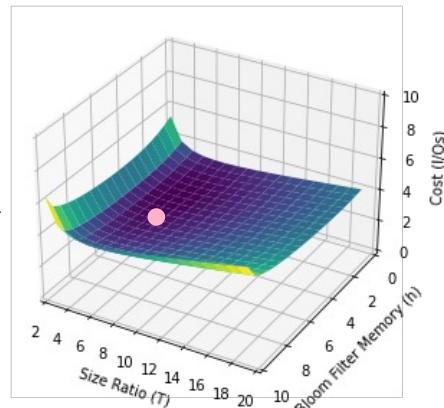
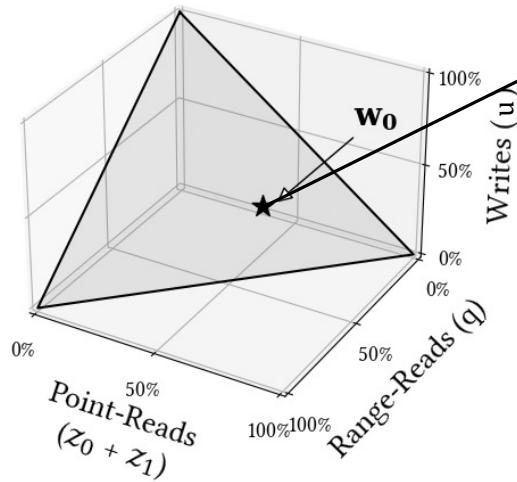
$$\varphi^* = \operatorname{argmin}_{\varphi} C(w, \varphi)$$

Cost function is a weighted sum over each operation type

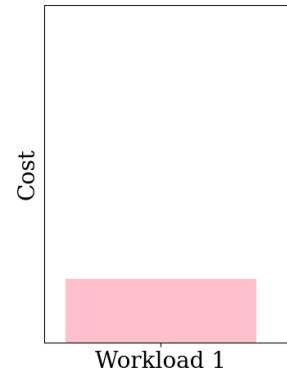
$$C_S(w, \varphi) = z_0 \cdot Z_0(\varphi) + z_1 \cdot Z_1(\varphi) + q \cdot Q(\varphi) + u \cdot U(\varphi)$$

Tuning Problems

\mathbf{w}_0 : Workload (z_0, z_1, q, u)



● Best Configuration



Outline

A Primer on LSM Trees

Flexibility in Compaction Design [VLDB-J 24]

Modeling LSM Trees

Tuning LSM Trees

ENDURE: Finding Robust Tunings for LSM Trees [VLDB 22]

AXE: Learning to Tune LSM Trees By Task Decomposition

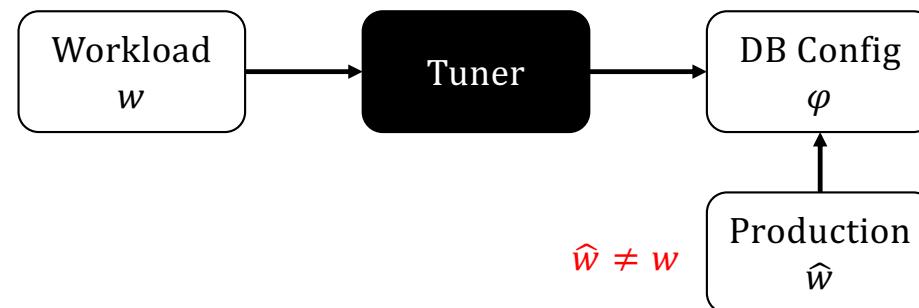
We Are Not ORACLEs

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

φ : LSM Tree Design ($m_{buf}, m_{filter}, T, K_1, K_2, \dots, K_L$)

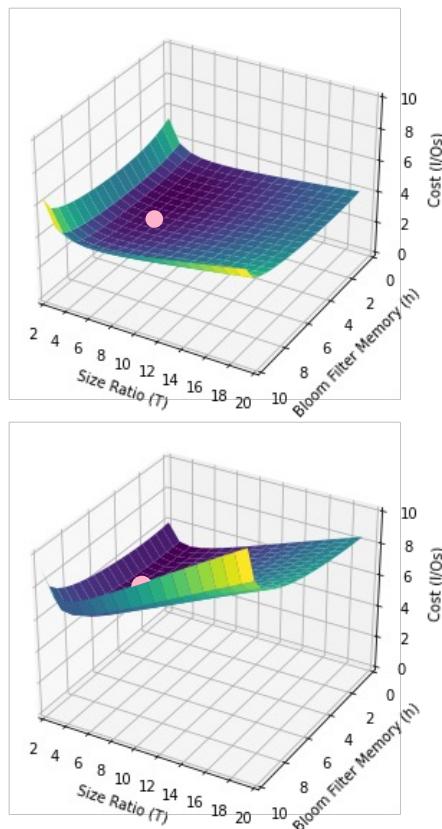
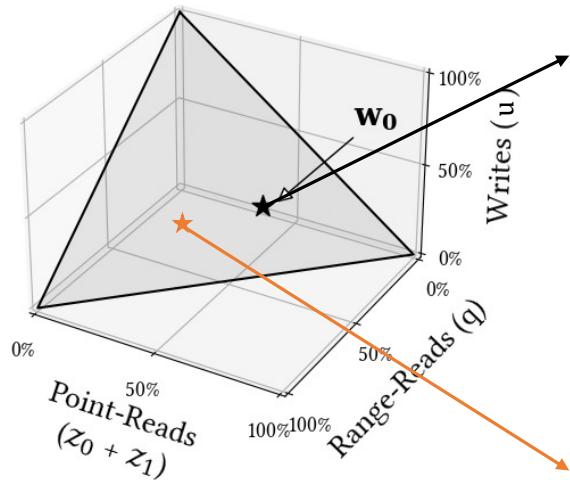
C : Cost (I/O)

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

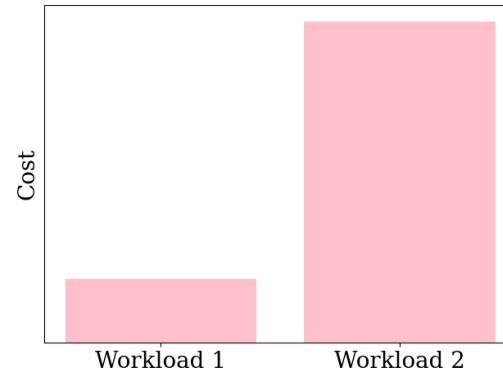


Tuning Problems

\mathbf{w}_0 : Workload (z_0, z_1, q, u)



● 'Best' Configuration



Optimal tuning depends on workload

Workload uncertainty leads to
sub-optimal tuning

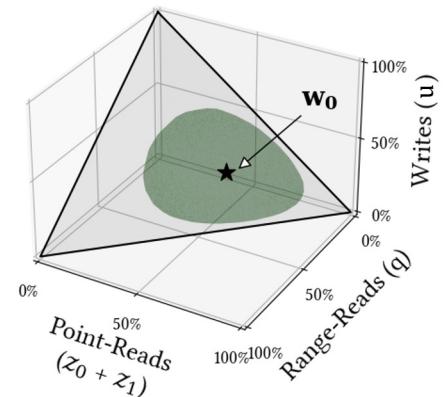
The LSM-Tuning Problem

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

φ : Design ($m_{buff}, m_{filter}, T, K_1, K_2, \dots, K_L$)

C : Cost (I/O)

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



Nominal

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

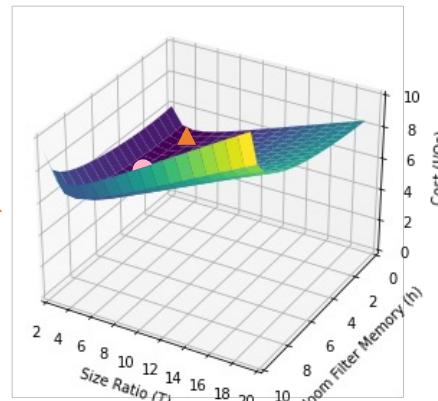
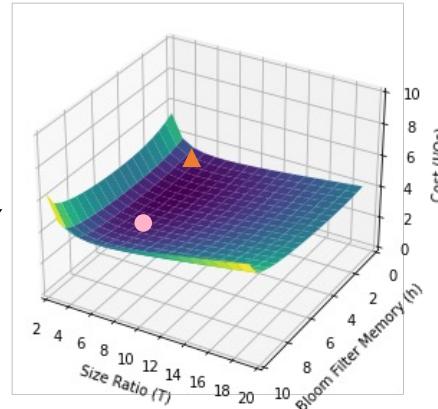
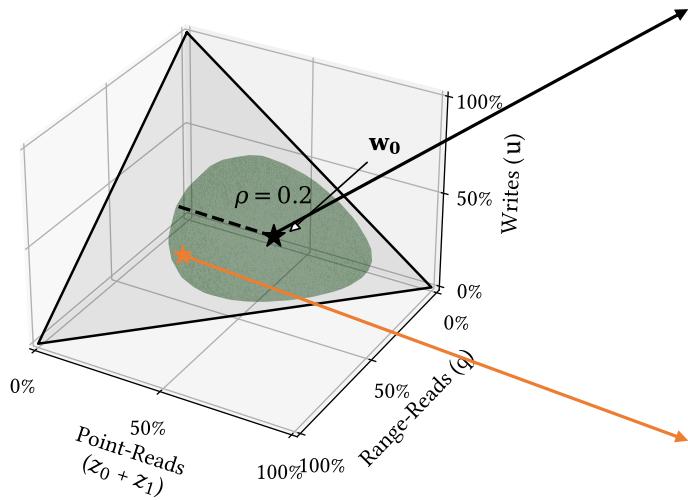
Robust

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

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

Robust Tuning

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

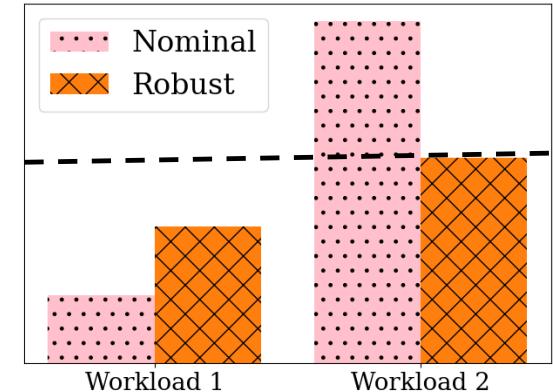


$$\varphi^* = \operatorname{argmin}_\varphi C(\hat{w}, \varphi)$$

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

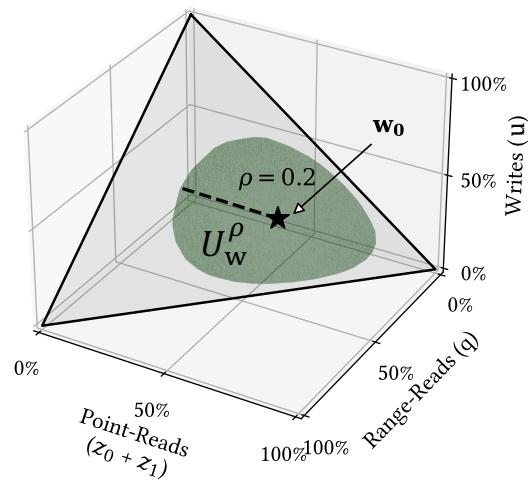
● Optimal configuration for expected workload

▲ Robust configuration for the workload neighborhood



Uncertainty Neighborhood

Workload Characteristic



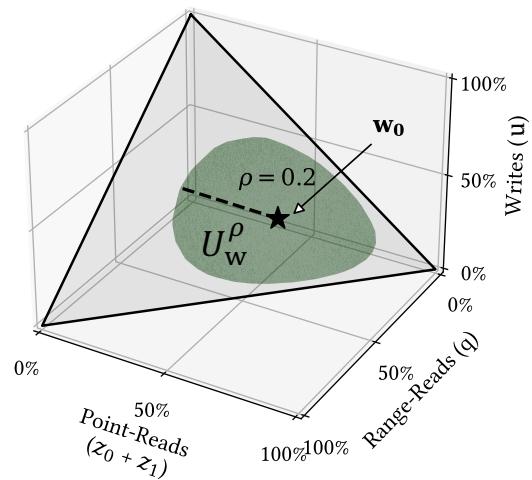
Neighborhood of workloads (ρ) via the KL-divergence

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

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

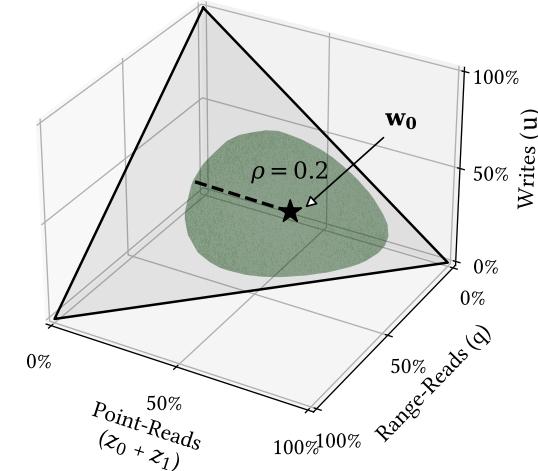
User provided workload uncertainty

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

Solving Robust Problem

Iterating over every possible workload is expensive

$$\begin{aligned}\varphi_R = \arg \min_{\varphi} C_S(\hat{w}, \varphi) \\ s.t., \quad \hat{w} \in \mathcal{U}_w\end{aligned}$$



Solving Robust Problem

Iterating over every possible workload is expensive

$$\begin{aligned}\varphi_R &= \arg \min_{\varphi} C_S(\hat{w}, \varphi) \\ s.t., \quad \hat{w} &\in \mathcal{U}_w\end{aligned}$$

Formulation of the dual problem to reduce to a feasible problem[‡]

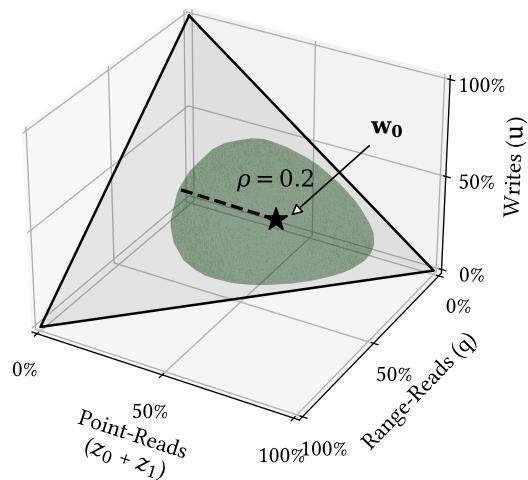
$$\min_{\varphi} \max_{\hat{w} \in \mathcal{U}_w} C_S(\hat{w}, \varphi)$$

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

[‡]Aharon Ben-Tal, Dick den Hertog, Anja De Waegenaere, Bertrand Melenberg, and Gijs Rennen.
2013. Robust Solutions of Optimization Problems Affected by Uncertain Probabilities.

ENDURE Pipeline

Workload Characteristic



System Information

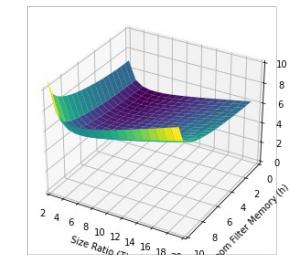
Page Size
Memory Budget



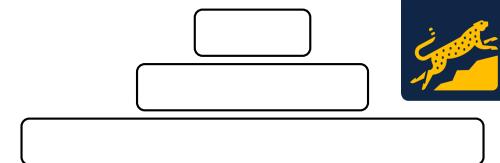
ENDURE
Solves the
Robust Problem



Expected performance



RocksDB Configuration



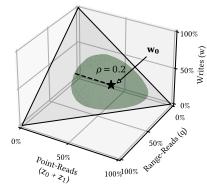
Testing Suite



ENDURE in Python, implemented in tandem with RocksDB

Uncertainty benchmark

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



Normalized delta throughput

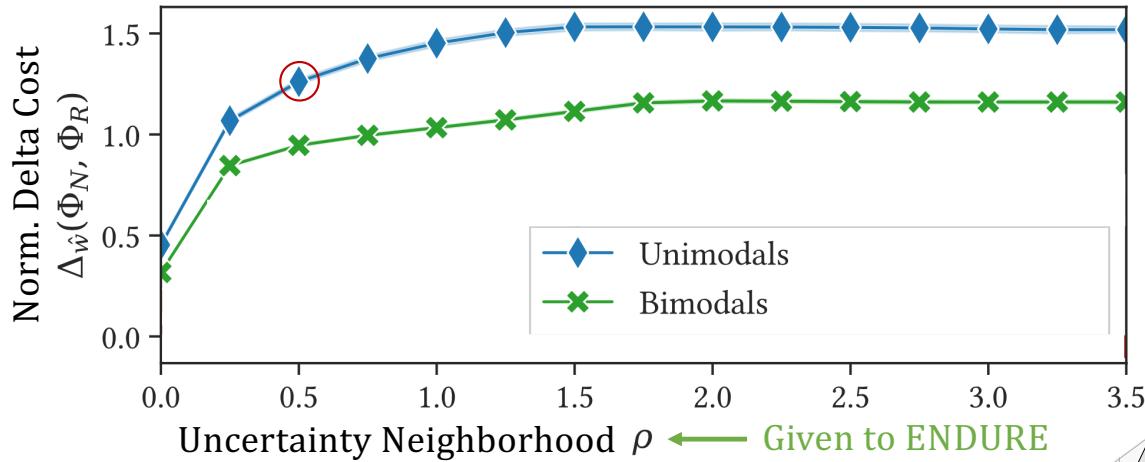
$$\Delta_w(\Phi_1, \Phi_2) = \frac{1/C(w, \Phi_2) - 1/C(w, \Phi_1)}{1/C(w, \Phi_1)}$$

Nominal vs Robust: > 0 is better

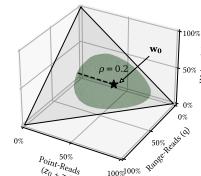
1 means 2x speedup

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

Impact of Workload Type

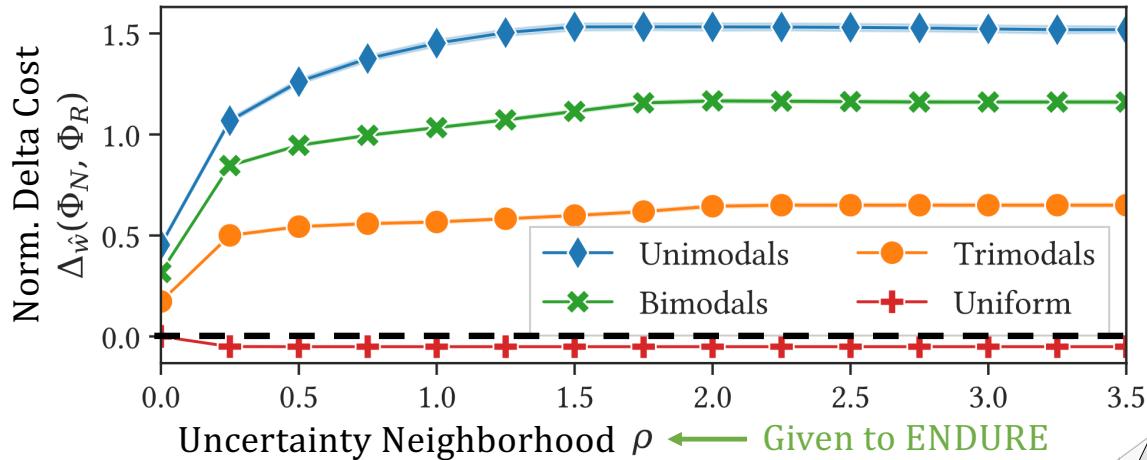


Unbalanced workloads result in overfitted nominal tunings



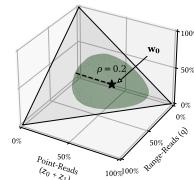
Index	(z_0, z_1, q, w)					Type
	25%	25%	25%	25%	Uniform	
0	97%	1%	1%	1%	1%	Unimodal
1	1%	97%	1%	1%	1%	Unimodal
2	1%	1%	97%	1%	1%	Unimodal
3	1%	1%	1%	97%	1%	Unimodal
4	1%	1%	1%	1%	97%	Unimodal
5	49%	49%	1%	1%	1%	Bimodal
6	49%	1%	49%	1%	1%	Bimodal
7	49%	1%	1%	49%	1%	Bimodal
8	1%	49%	49%	1%	1%	Bimodal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

Impact of Workload Type



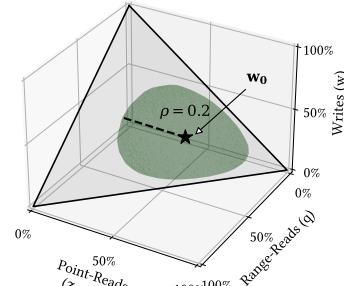
Unbalanced workloads result in overfitted nominal tunings

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

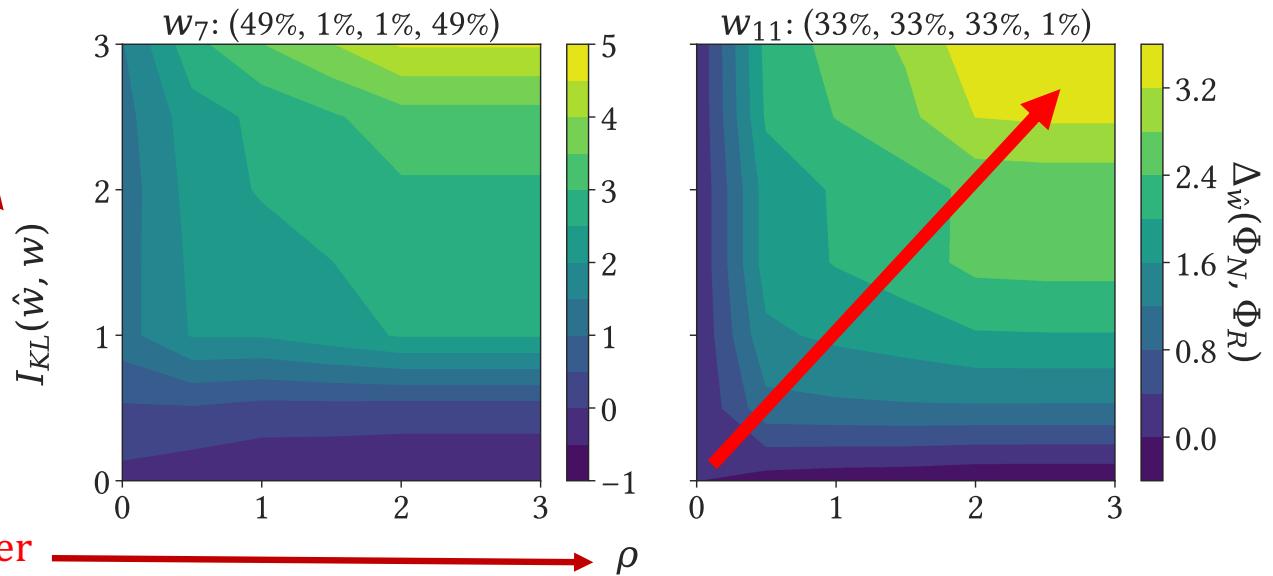


Relationship of Expected and Observed ρ

Observed ρ : distance from executed workload to expected workload



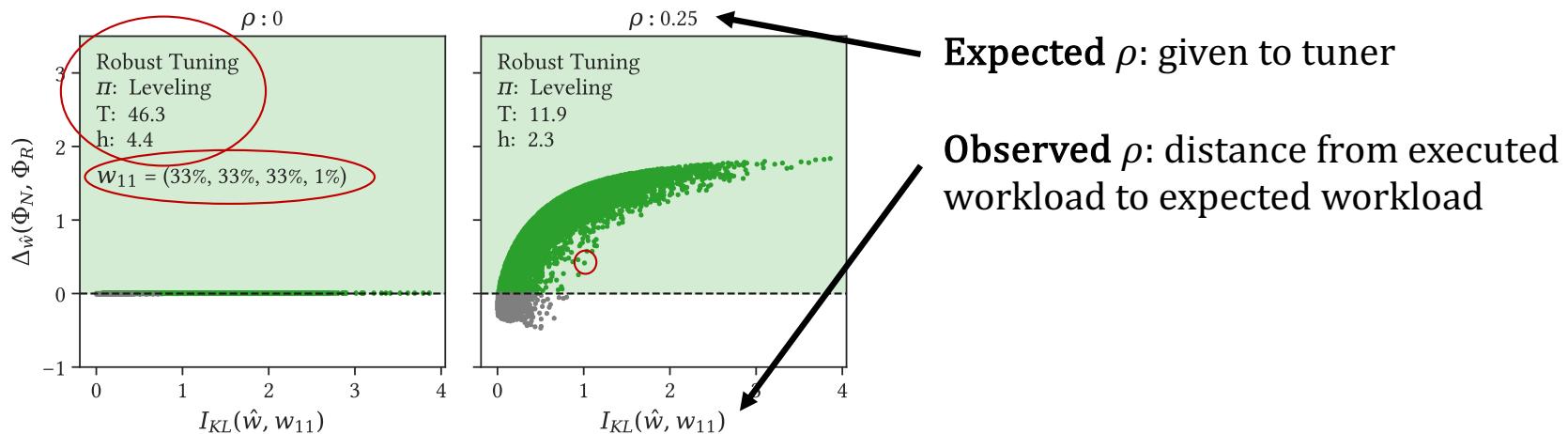
Expected ρ : workload given to tuner



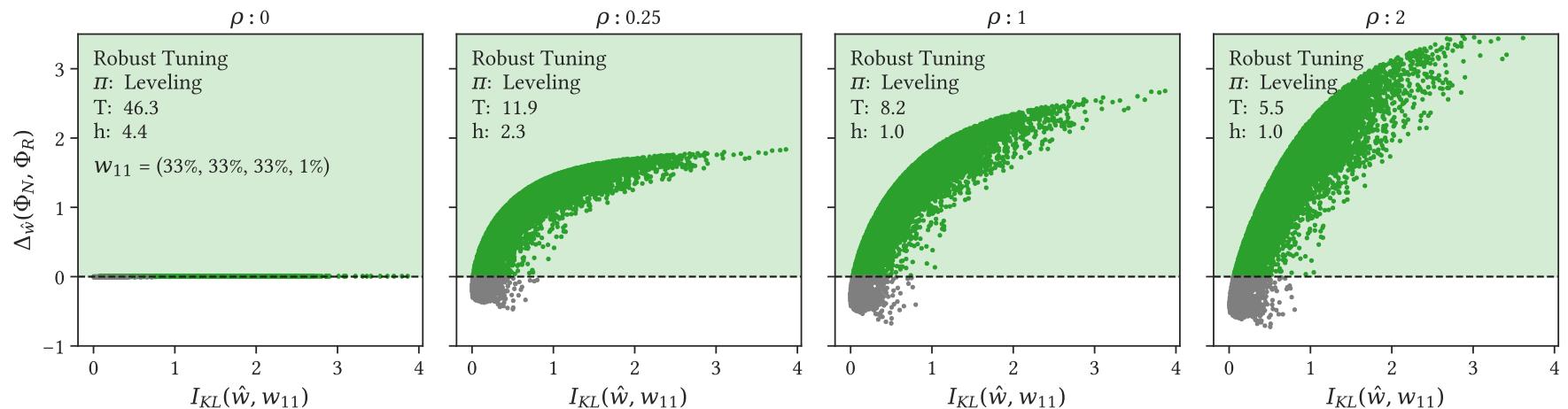
Highest throughput when observed and expected ρ match

Lowest throughput when ρ is mismatched

Impact of Observed vs Expected ρ

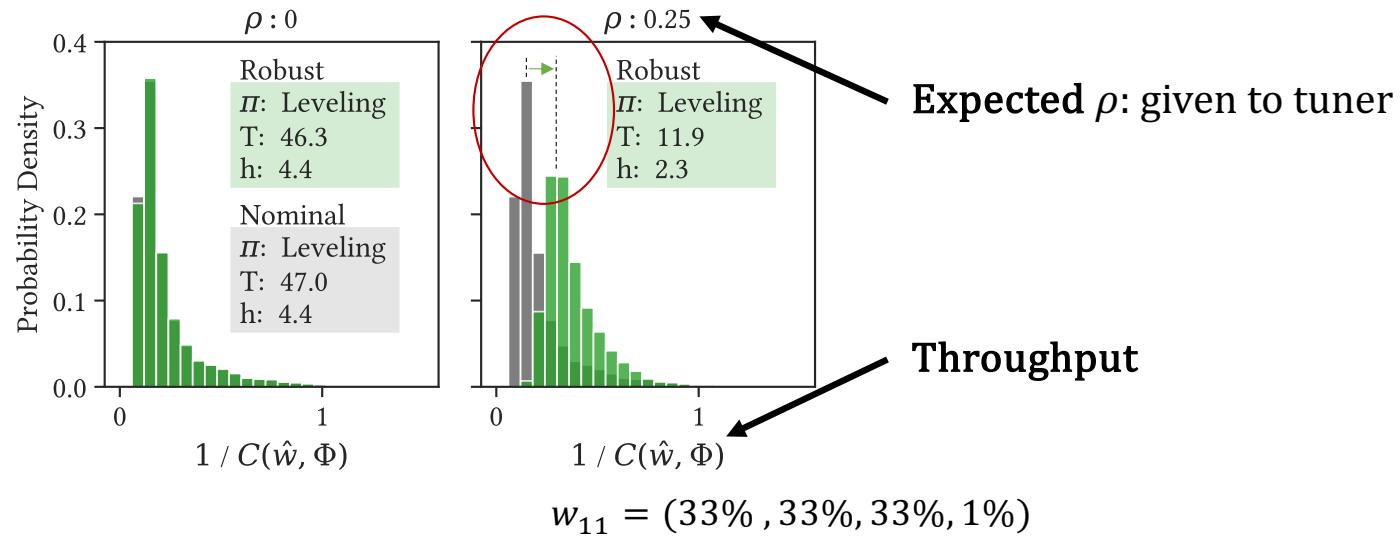


Impact of Observed vs Expected ρ

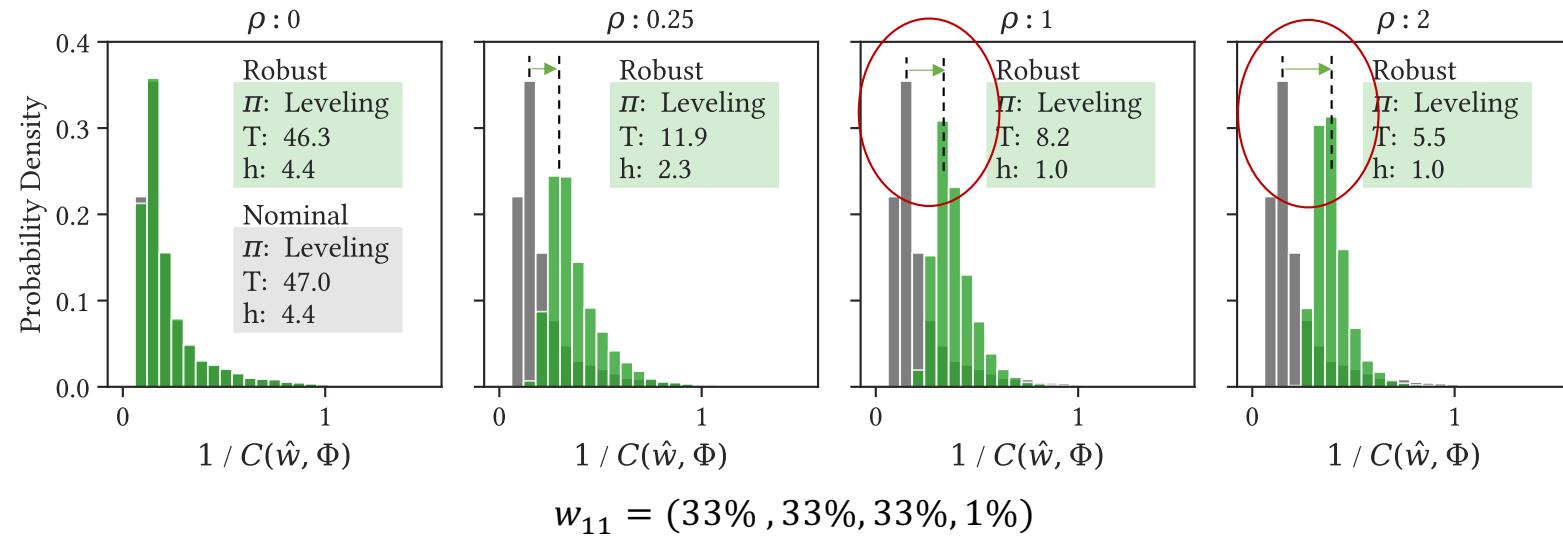


- Higher expected ρ accounts for more uncertainty,
- Potential speed up of 4x
- Higher expected $\rho \rightarrow$ anticipates writes \rightarrow shallow tree

ρ and Performance Gain Distribution

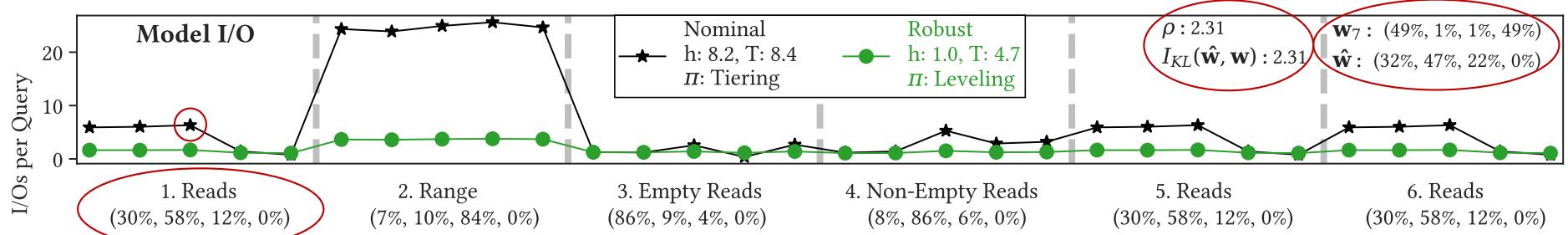


ρ and Performance Gain Distribution



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

Workload Sequence on RocksDB

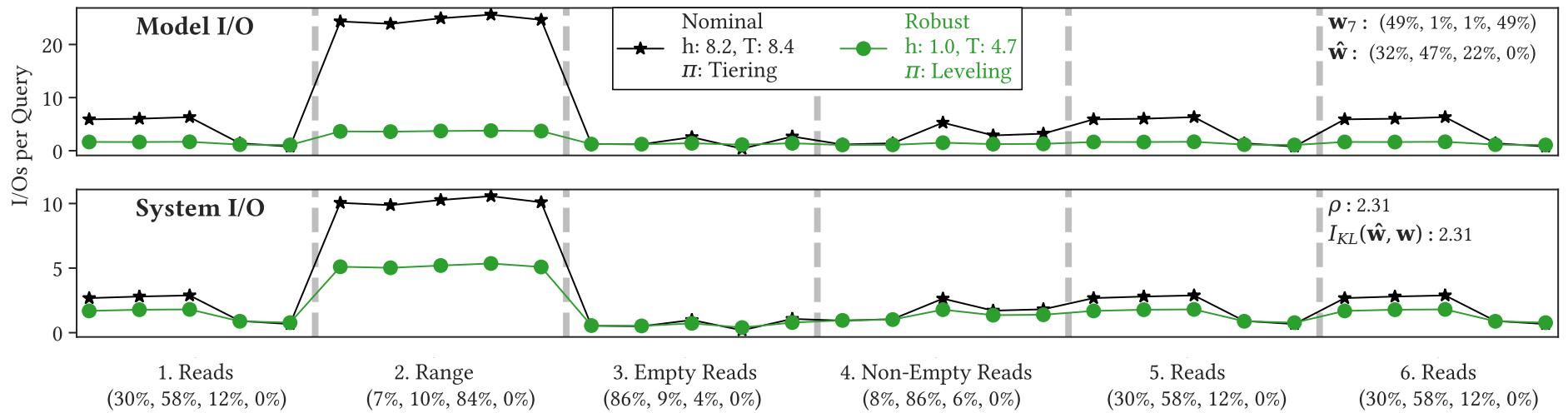


RocksDB instance setup with 10 million unique key-value pairs of size 1KB

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

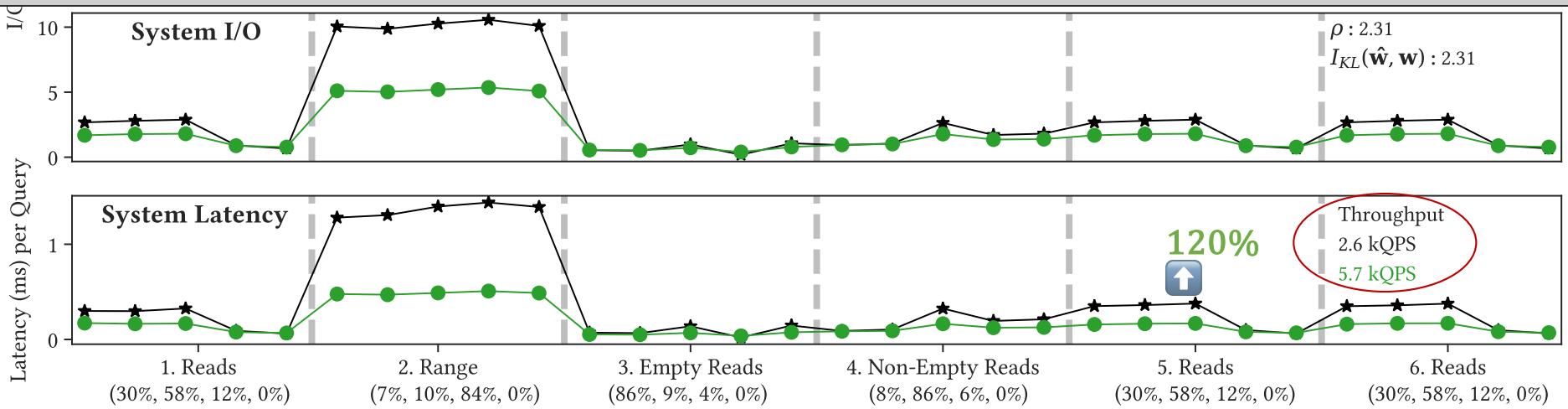
Writes are unique, range queries average 1-2 pages per level

Workload Sequence



Workload Sequence

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



Outline

A Primer on LSM Trees

Flexibility in Compaction Design [VLDB-J 24]

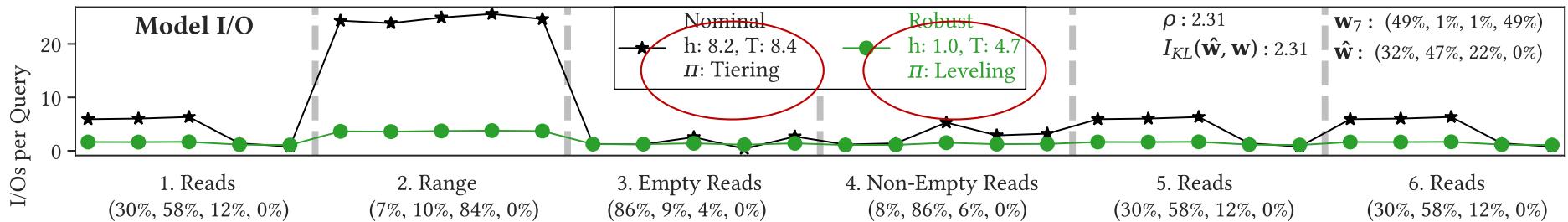
Modeling LSM Trees

Tuning LSM Trees

ENDURE: Finding Robust Tunings for LSM Trees [VLDB 22]

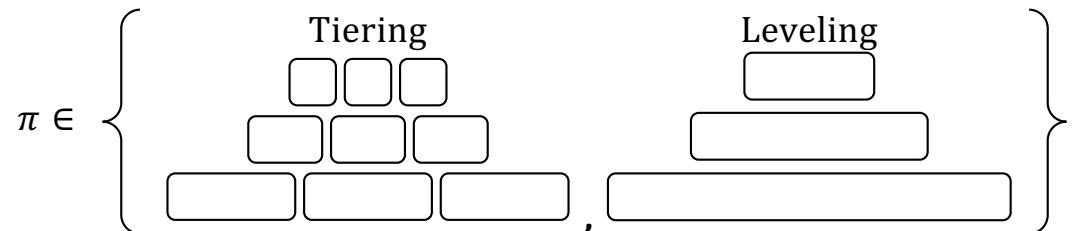
AXE: Learning to Tune LSM Trees By Task Decomposition

Limitations in ENDURE



Optimizers tends to scale poorly as the decision space grows, hence, we limited the decision to Tiering or Leveling

Cost model took a lot of man hours



Model

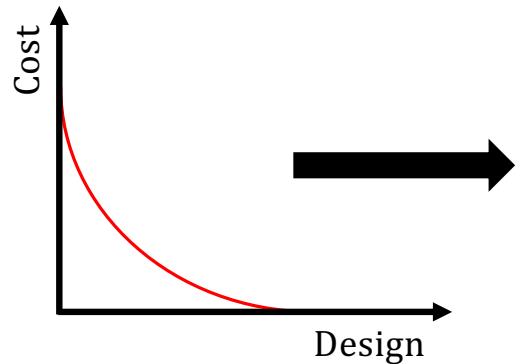
$$Z_0(\varphi) = \sum_{i=1}^{L(T)} K_i \cdot f_i(T)$$

$$U(\varphi) = f_{seq} \cdot \frac{1+f_a}{B} \cdot \sum_{i=1}^{L(T)} \frac{T-1+K_i}{2K_i}$$

Optimizer

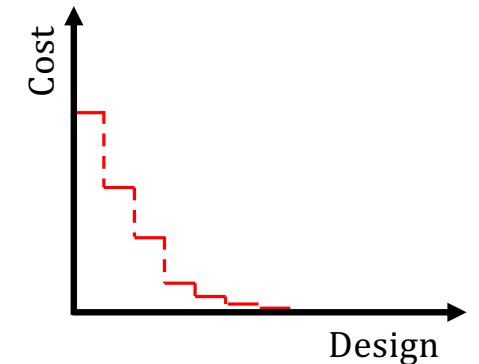
Complex Decision Space

Knob	Value
K_1	1.0
K_2	4.0
K_3	2.0
K_4	3.0
T	6.0



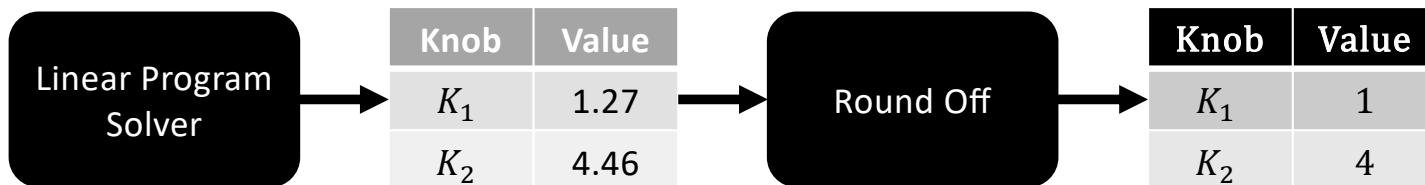
We like to think of knobs as continuous values

Knob	1	2	3	4	5	6
K_1	1	0	0	0	0	0
K_2	0	0	0	1	0	0
K_3	0	1	0	0	0	0
K_4	0	0	1	0	0	0
T	0	0	0	0	0	1



Reality is they take on a set number of values

Solution: Pretend we're in a continuous space then round off

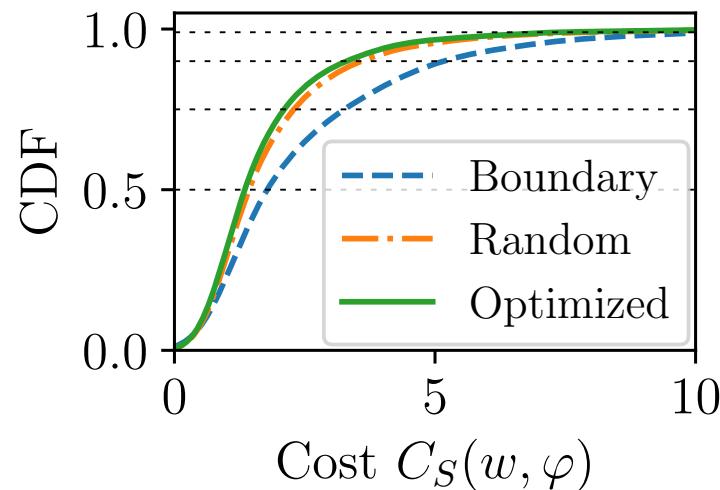


Problem: Errors propagate fast

Sensitive Optimizers

1. Build optimizer with different initial starting conditions
2. For each optimizer, test recommended designs for various workloads

No workload drift. ($w_0 = \hat{w}$)



Median Slowdown from Optimal

Percentile	Boundary	Random
P50%	1.0x	1.0x
P75%	1.7x	1.1x
P90%	2.4x	1.2x
P99%	4.1x	1.7x
P99.9%	6.2x	3.3x

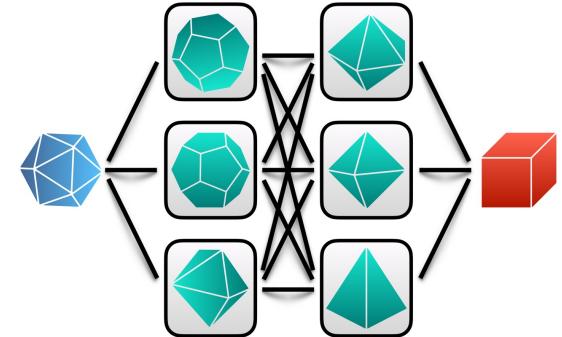
Tackling Discrete Optimization

Machine learning strategies proven to work well in discrete optimization tasks

Adapt techniques for systems tuning



Machine Learning for
Combinatorial Optimization
—COMPETITION 2021—



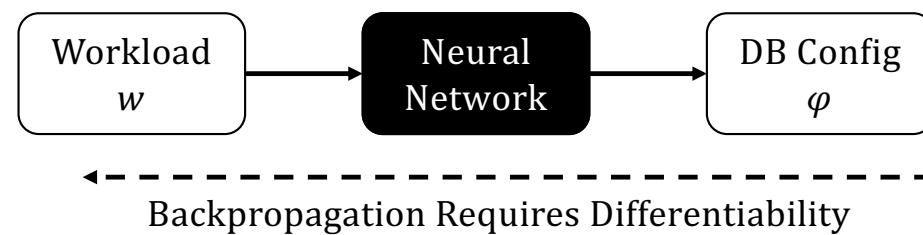
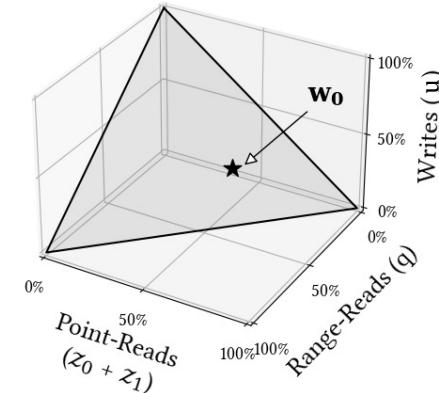
The LSM Tuning Problem

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

φ : Design $(m_{buff}, m_{filter}, T, K_1, K_2, \dots, K_L)$

C : Cost (I/O)

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

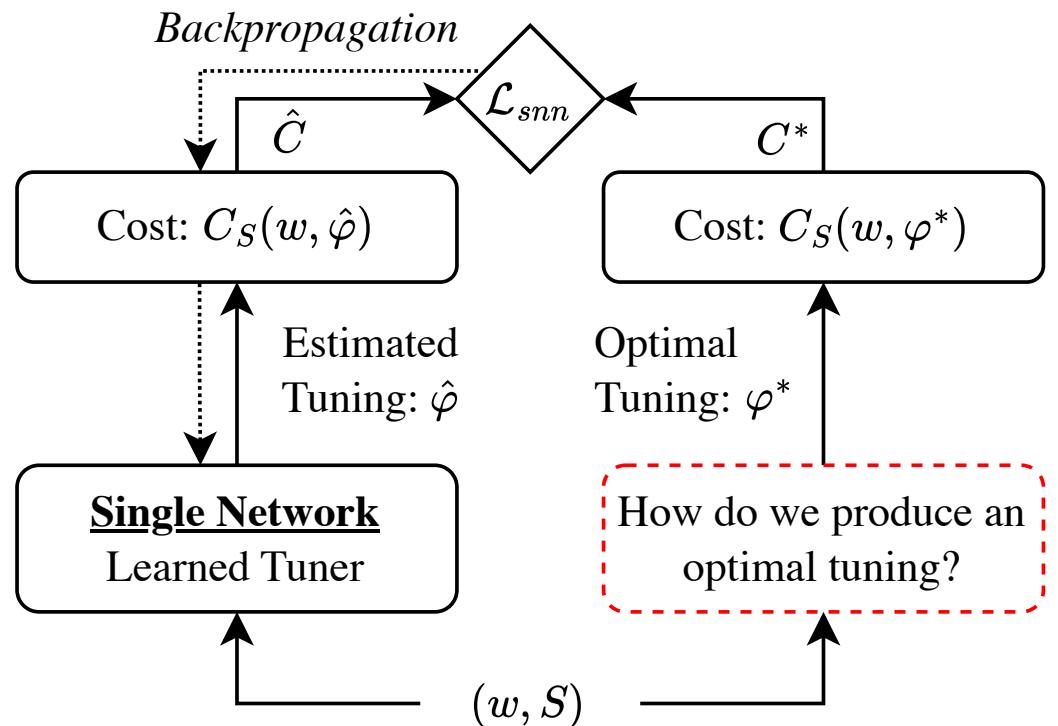


Training a Tuner

How would we train a tuner?

Single neural network for our task is infeasible

$$\mathcal{L}_{snn} = \mathbb{E}_{(w,S) \sim D} [C_S(w, \hat{\varphi}) - C^*]$$

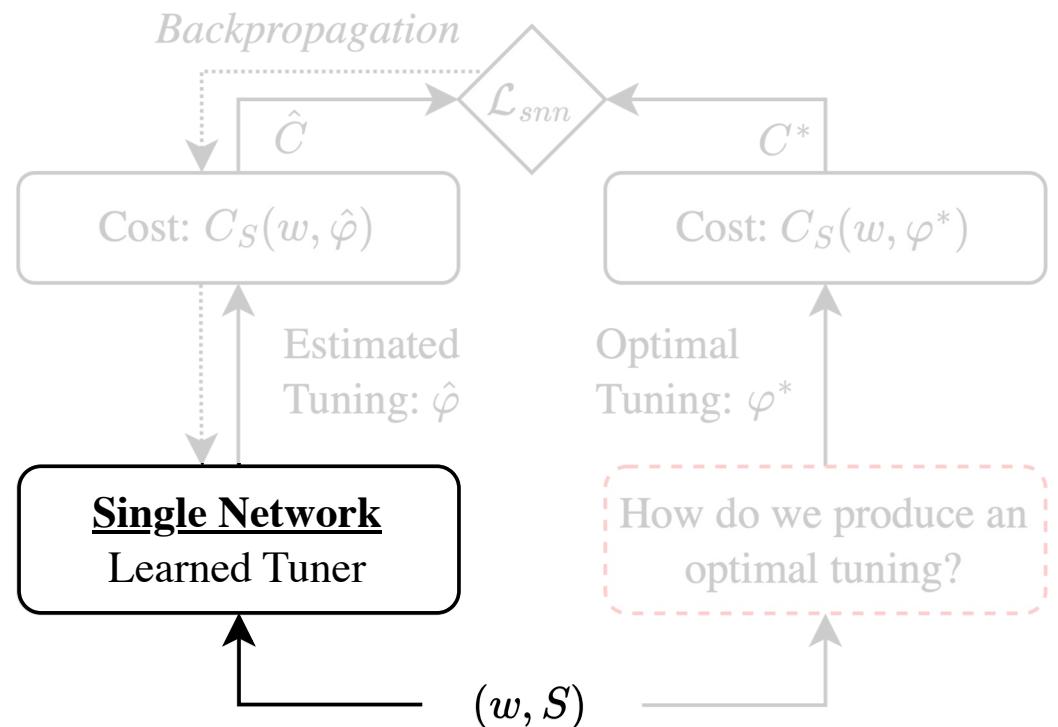


Training a Tuner

How would we train a tuner?

Single neural network for our task is infeasible

$$\mathcal{L}_{snn} = \mathbb{E}_{(w,S) \sim D} [C_S(w, \hat{\varphi}) - C^*]$$

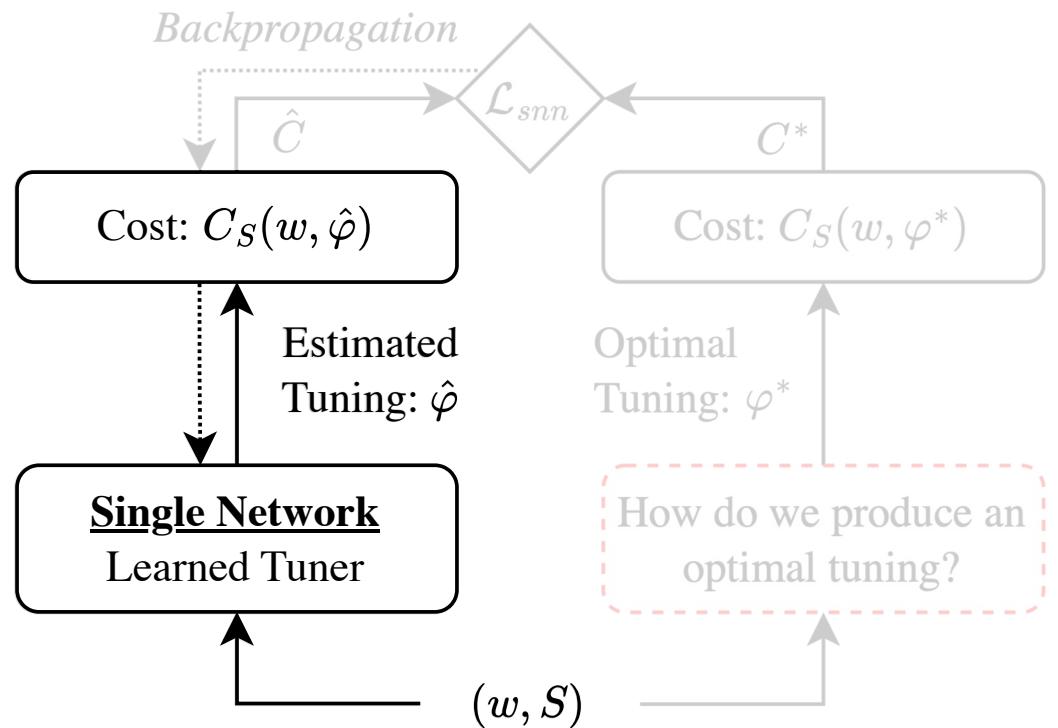


Training a Tuner

How would we train a tuner?

Single neural network for our task is infeasible

$$\mathcal{L}_{snn} = \mathbb{E}_{(w,S) \sim D} [C_S(w, \hat{\varphi}) - C^*]$$

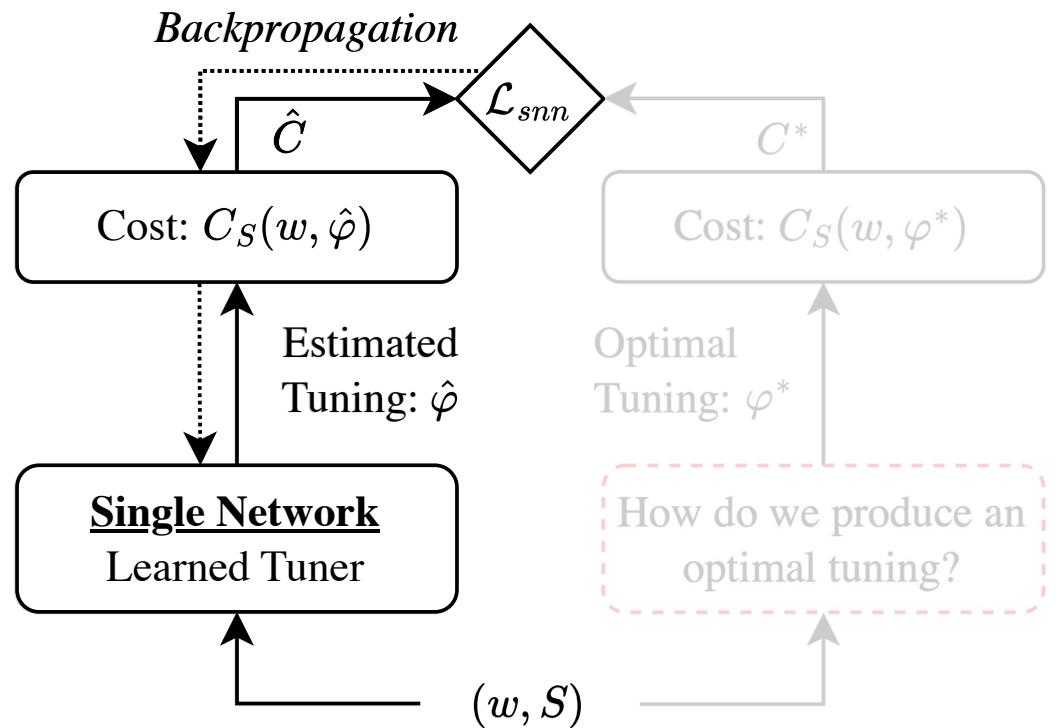


Training a Tuner

How would we train a tuner?

Single neural network for our task is infeasible

$$\mathcal{L}_{snn} = \mathbb{E}_{(w,S) \sim D} [C_S(w, \hat{\varphi}) - C^*]$$

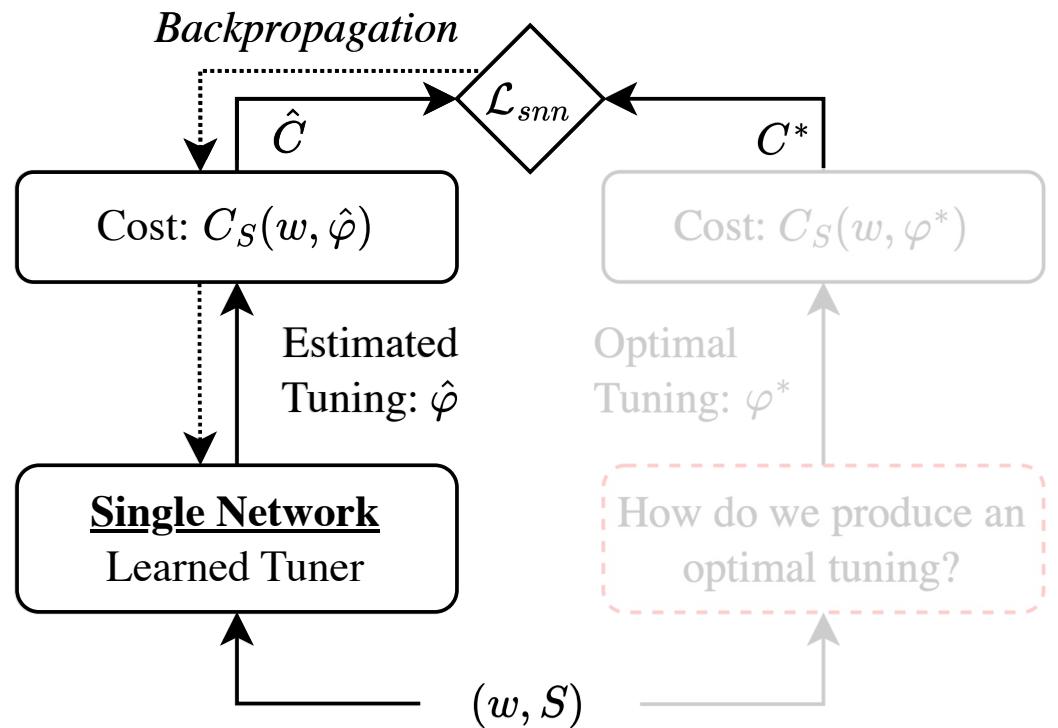


Training a Tuner

How would we train a tuner?

Single neural network for our task is infeasible

$$\mathcal{L}_{snn} = \mathbb{E}_{(w,S) \sim D} [C_S(w, \hat{\varphi}) - C^*]$$

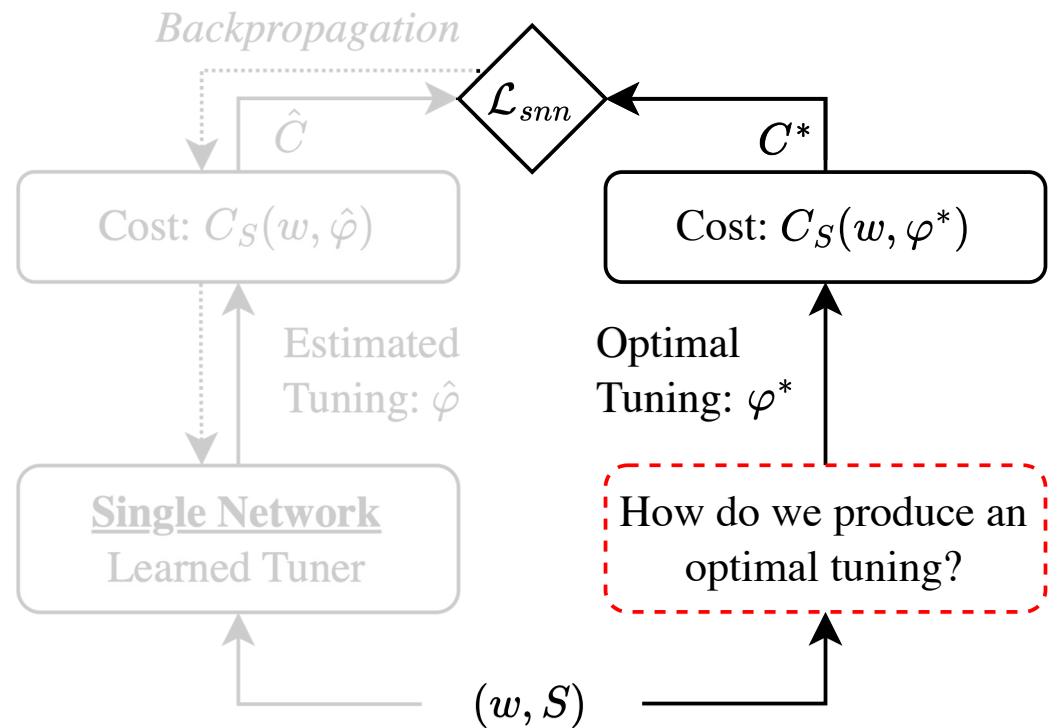


Training a Tuner

How would we train a tuner?

Single neural network for our task is infeasible

$$\mathcal{L}_{snn} = \mathbb{E}_{(w,S) \sim D} [C_S(w, \hat{\varphi}) - C^*]$$



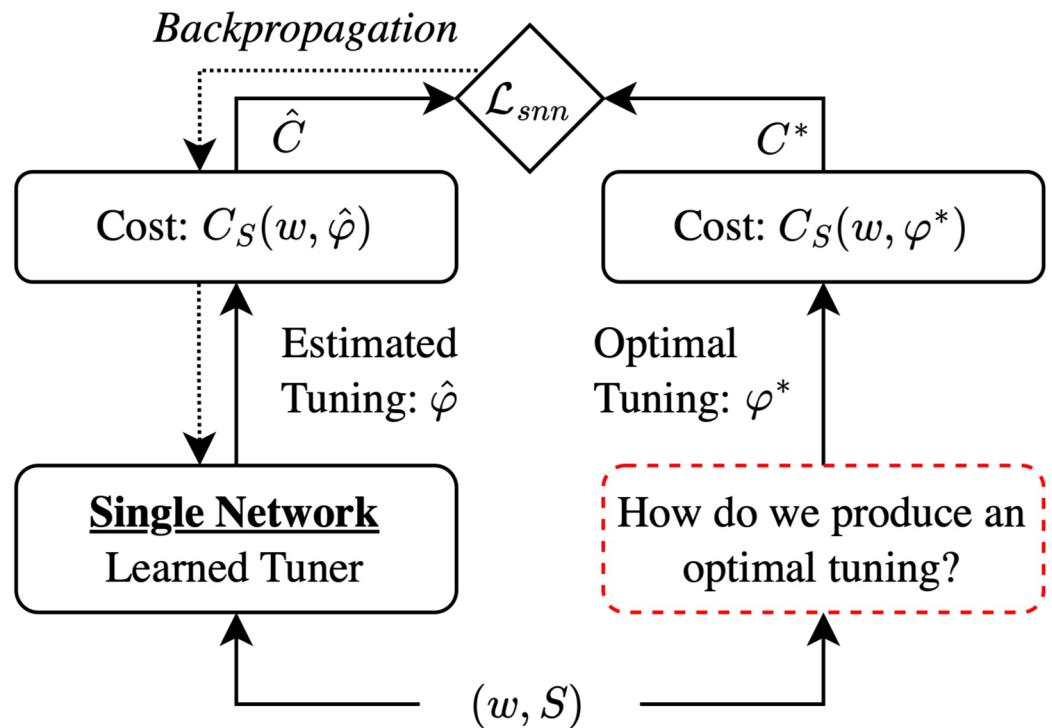
Training a Tuner

How would we train a tuner?

Single neural network for our task is infeasible

We would need to have already solved the problem to obtain the correct data

$$\mathcal{L}_{snn} = \mathbb{E}_{(w,S) \sim D} [C_S(w, \hat{\varphi}) - C^*]$$



The LSM Tuning Problem

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

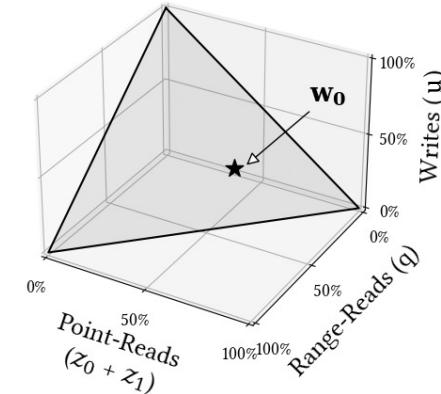
φ : Design $(m_{buff}, m_{filter}, T, K_1, K_2, \dots, K_L)$

C : Cost (I/O)

$$\varphi^* = \underset{\varphi}{\operatorname{argmin}} C(w, \varphi)$$

What is the **optimal tuning** for a specific **workload**?

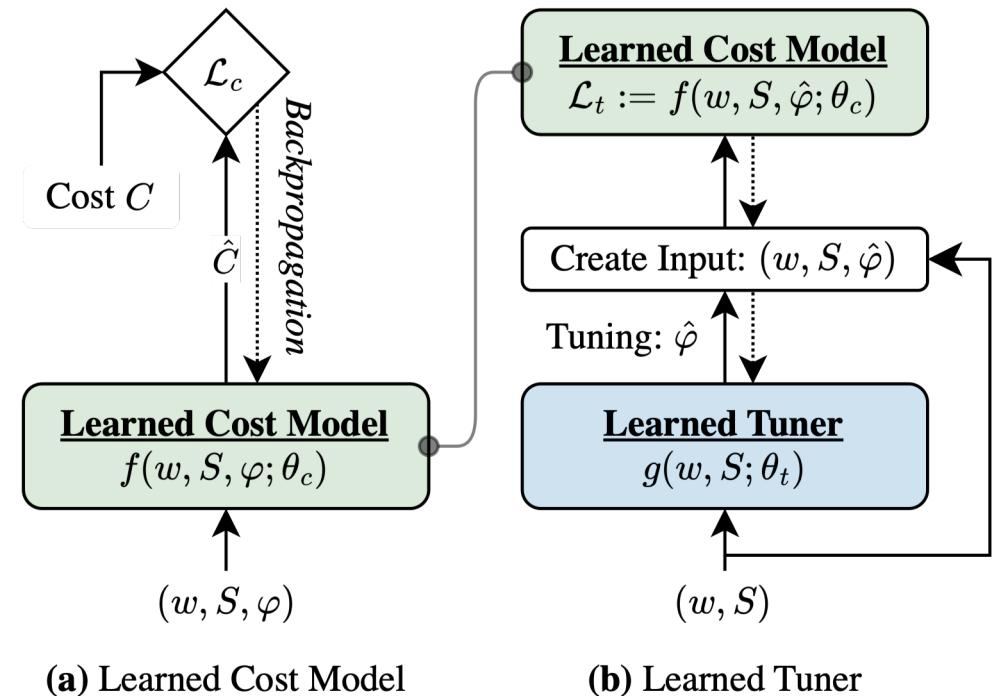
What is the **cost** of a **workload** executed on a **tuning**?



Task Decomposition of the Tuning Problem

We approach each subtask with a neural network

$$\mathcal{L}_c = \mathbb{E}_{(w, S, \varphi, C) \sim D} \left[\|\hat{C} - C\| \right]$$

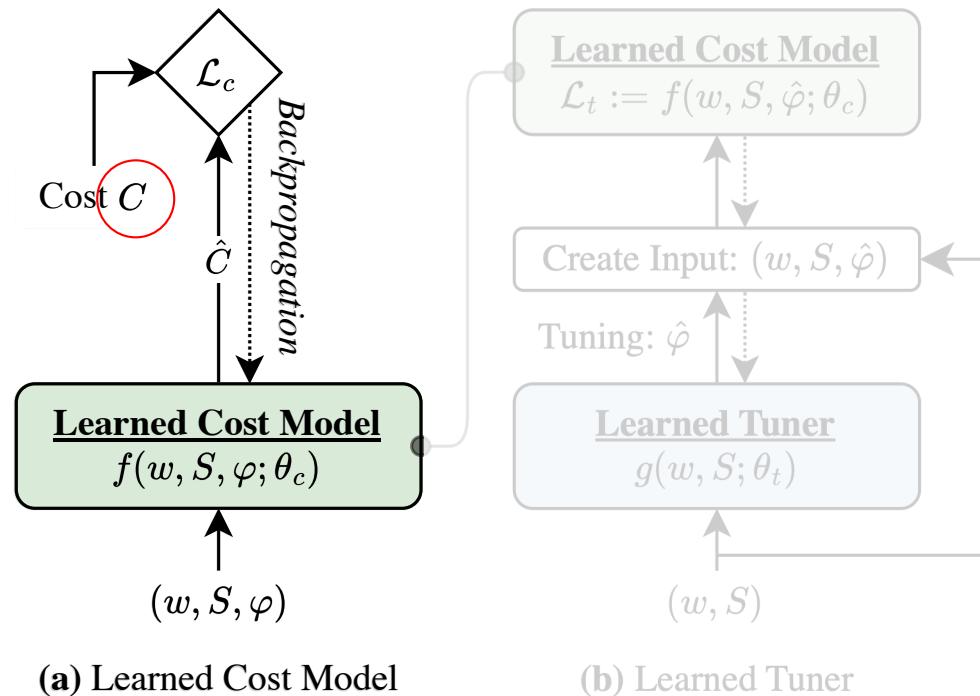


Task Decomposition of the Tuning Problem

We approach each subtask with a neural network

$$\mathcal{L}_c = \mathbb{E}_{(w, S, \varphi, C) \sim D} \left[||\hat{C} - C|| \right]$$

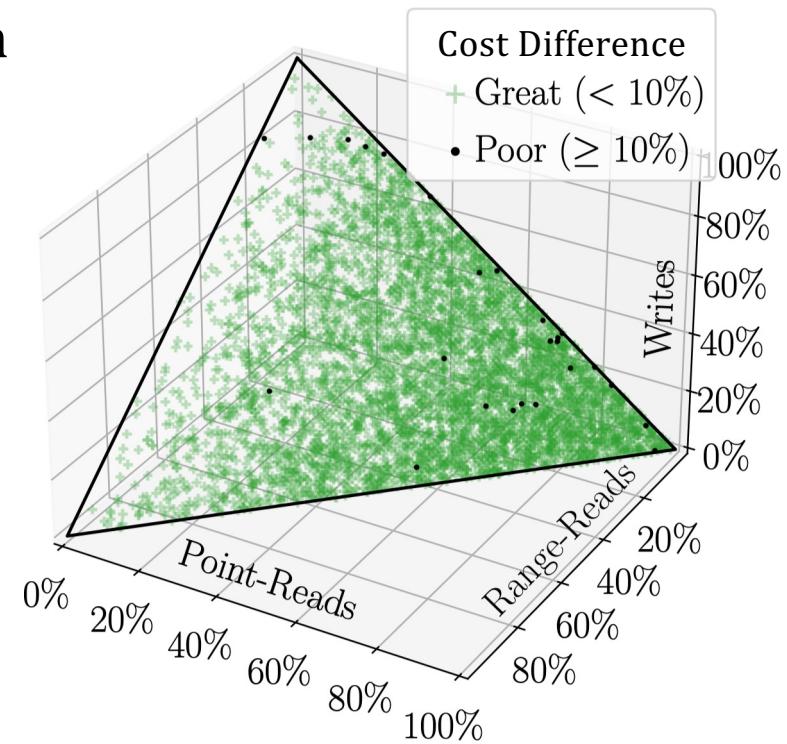
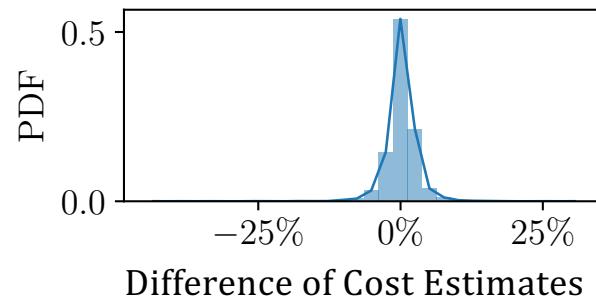
What is the **cost** of a **workload** executed on a **tuning**?



Learned Cost Model Performance

Randomly sampled workloads with a fixed environment and tuning

99% of points evaluated are within 10% of the true cost



Task Decomposition of the Tuning Problem

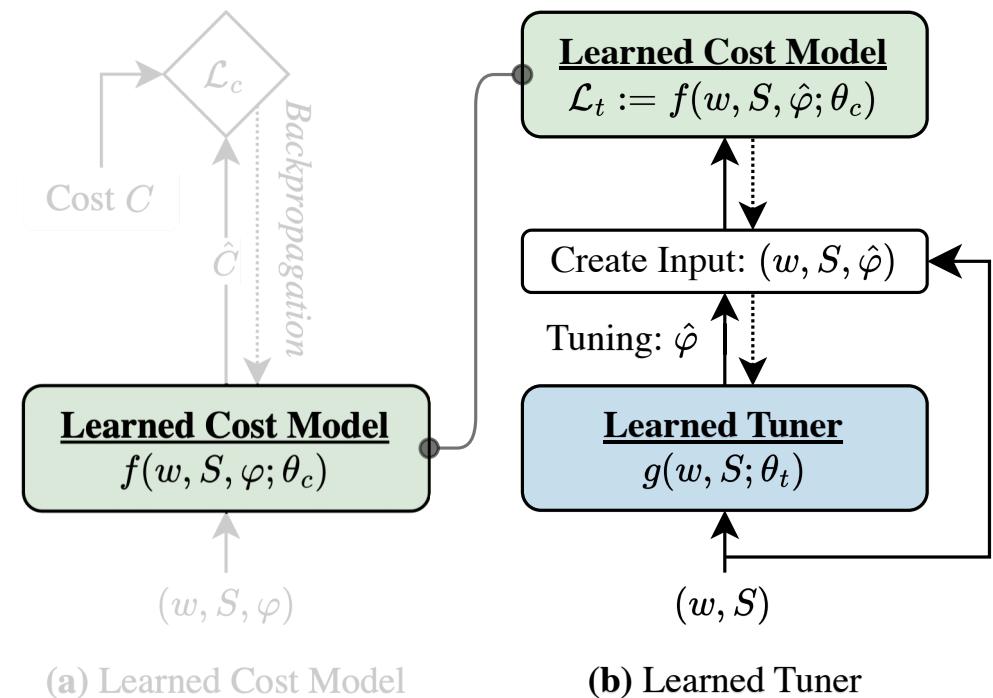
What is the **optimal tuning** for a specific **workload**?

$$\varphi^* = \arg \min_{\varphi} C_S(w, \varphi)$$

$$\mathcal{L}_t = \mathbb{E}_{(w, S) \sim \mathcal{D}} [C_S(w, \hat{\varphi})]$$

$$\mathcal{L}_t = \mathbb{E}_{(w, S) \sim \mathcal{D}} [f(w, S, \hat{\varphi})]$$

Learned cost model is used as the loss for the tuner



Task Decomposition of the Tuning Problem

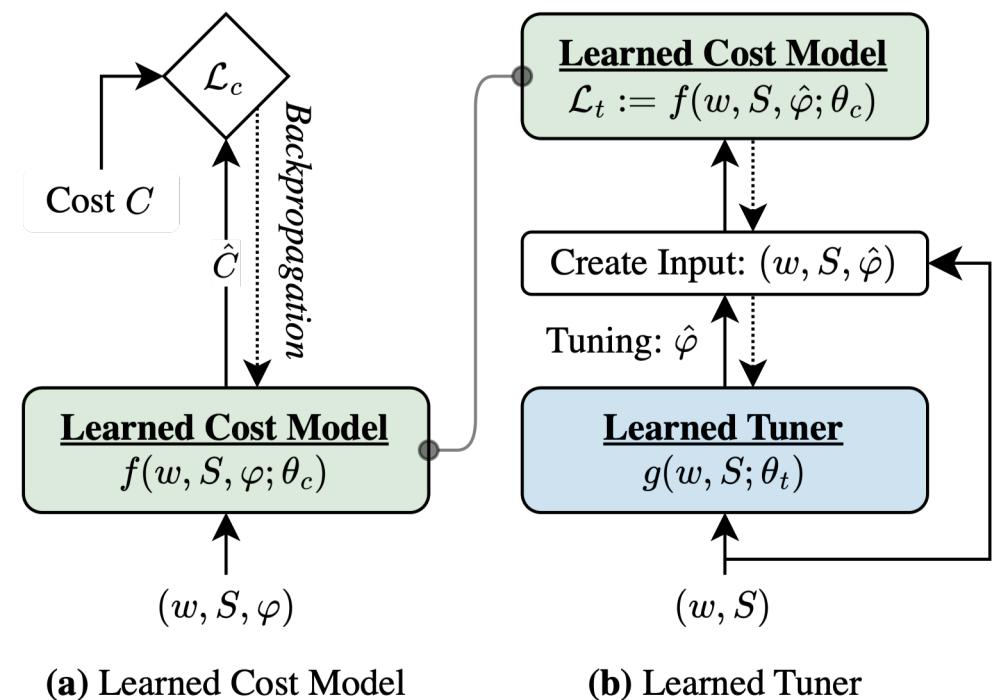
What is the **optimal tuning** for a specific **workload**?

$$\varphi^* = \arg \min_{\varphi} C_S(w, \varphi)$$

$$\mathcal{L}_t = \mathbb{E}_{(w, S) \sim \mathcal{D}} [C_S(w, \hat{\varphi})]$$

$$\mathcal{L}_t = \mathbb{E}_{(w, S) \sim \mathcal{D}} [f(w, S, \hat{\varphi})]$$

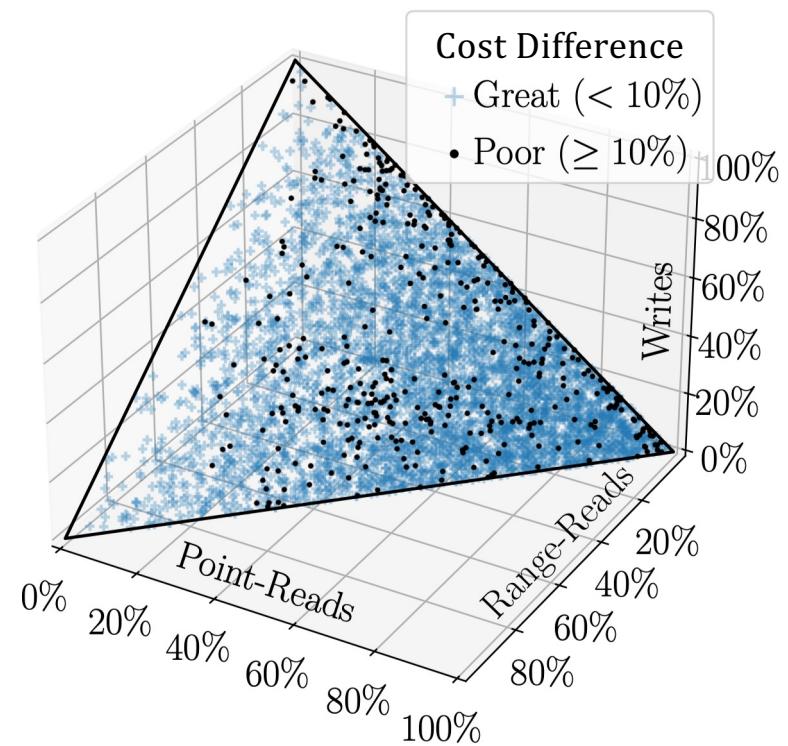
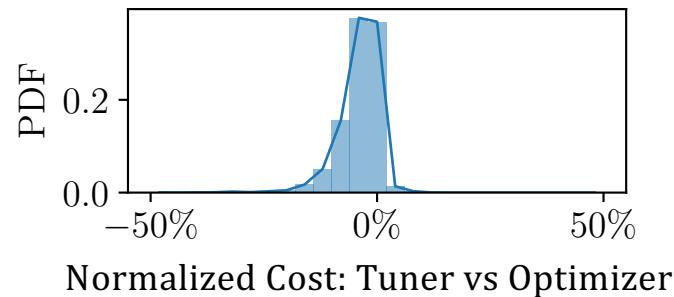
Learned cost model is used as the loss for the tuner



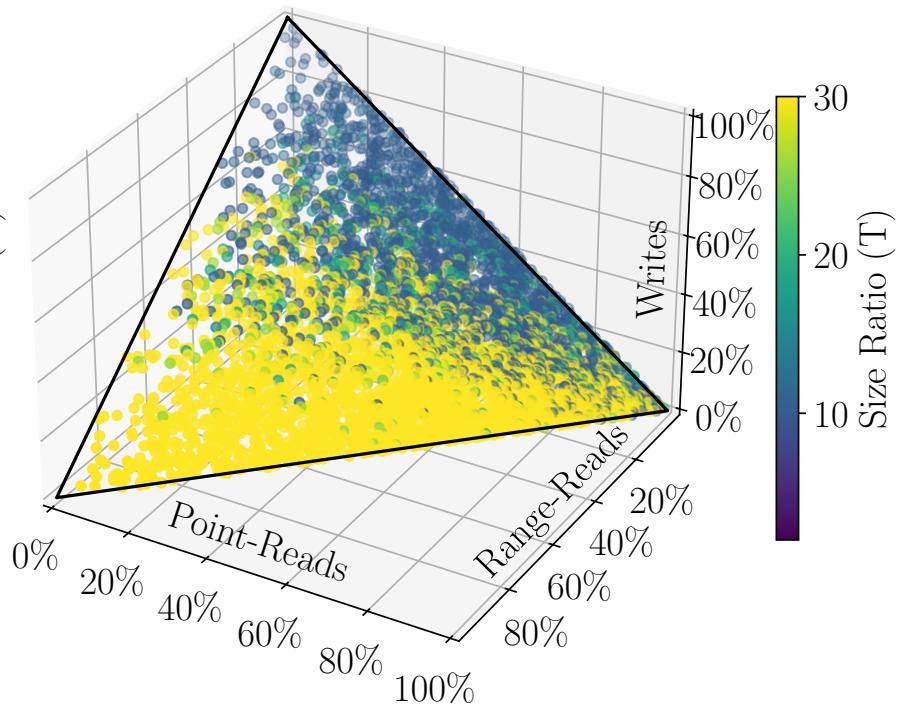
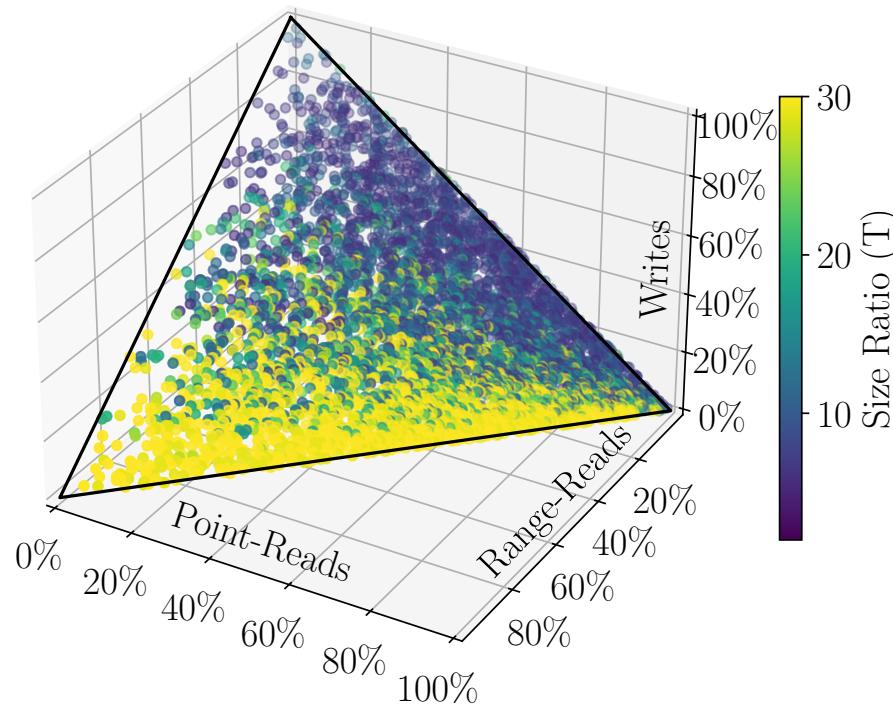
Learned Tuner Performance

Comparison to an expert
configured analytical optimizer

88% of tunings are within 10%
from the optimal

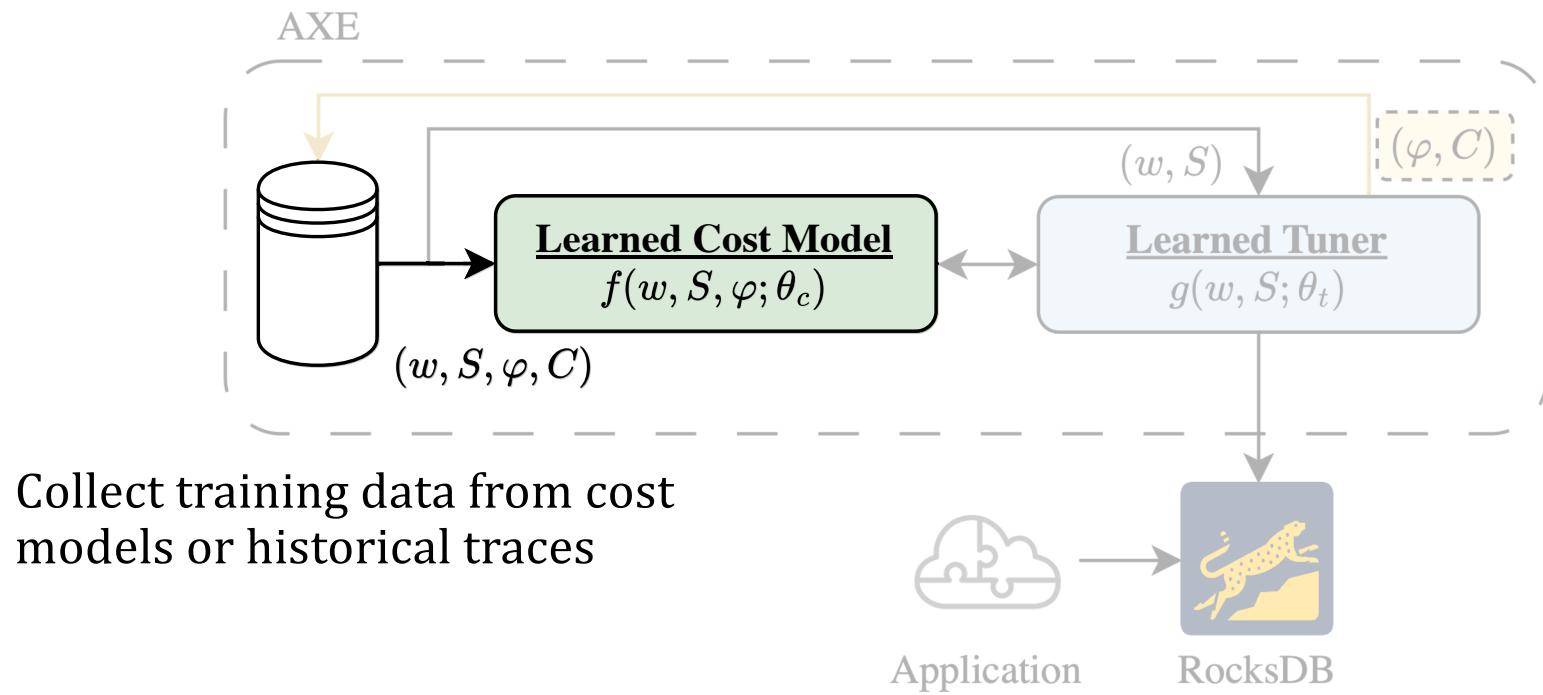


Size Ratio Recommendations

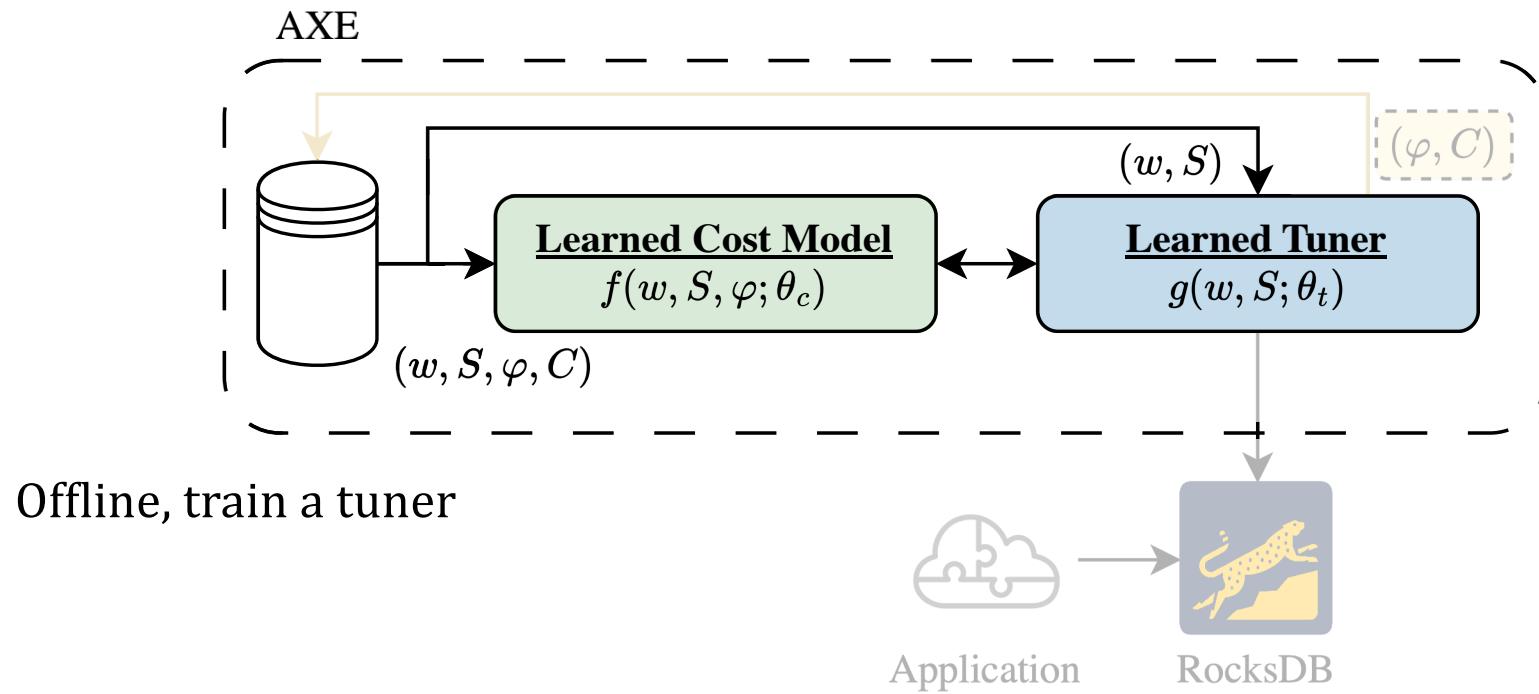


Analytical (L) and Learned Tuner (R) effectively navigates optimal size ratios

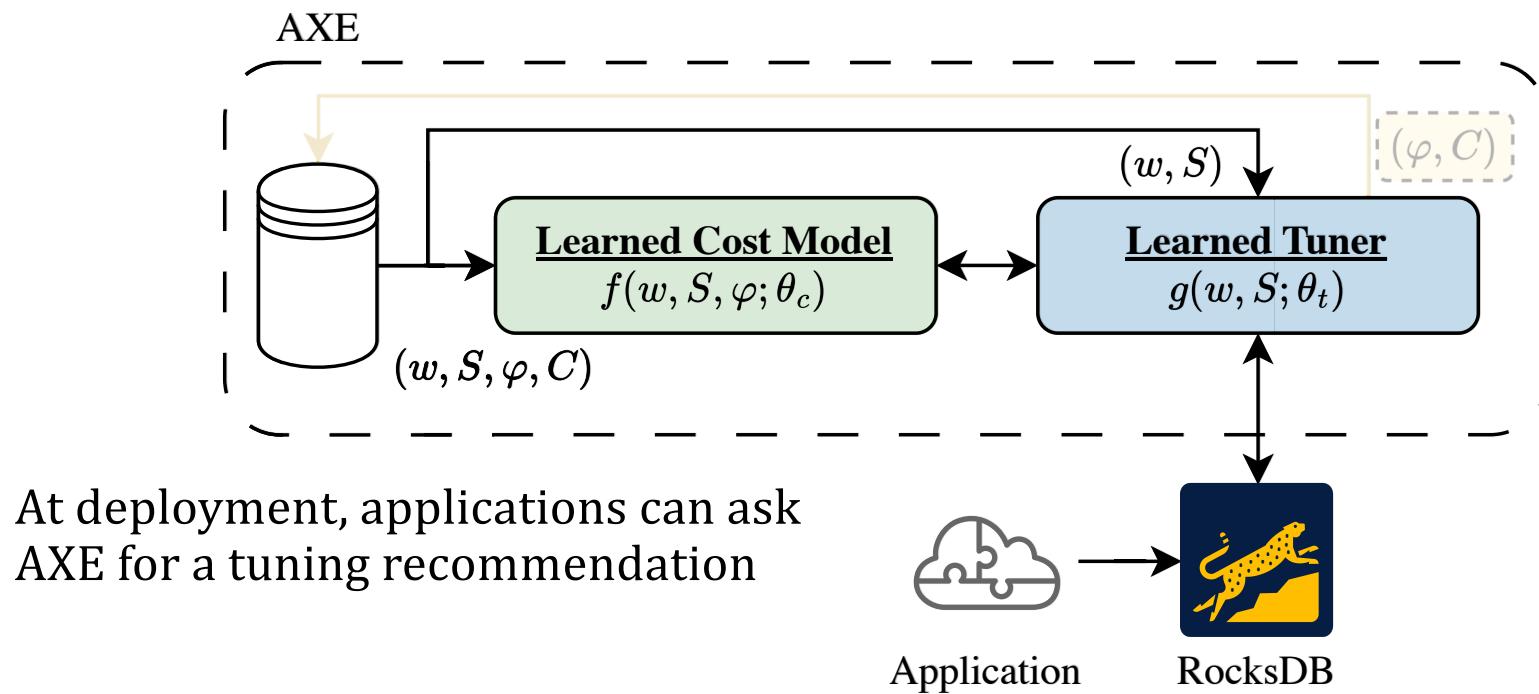
Building AXE



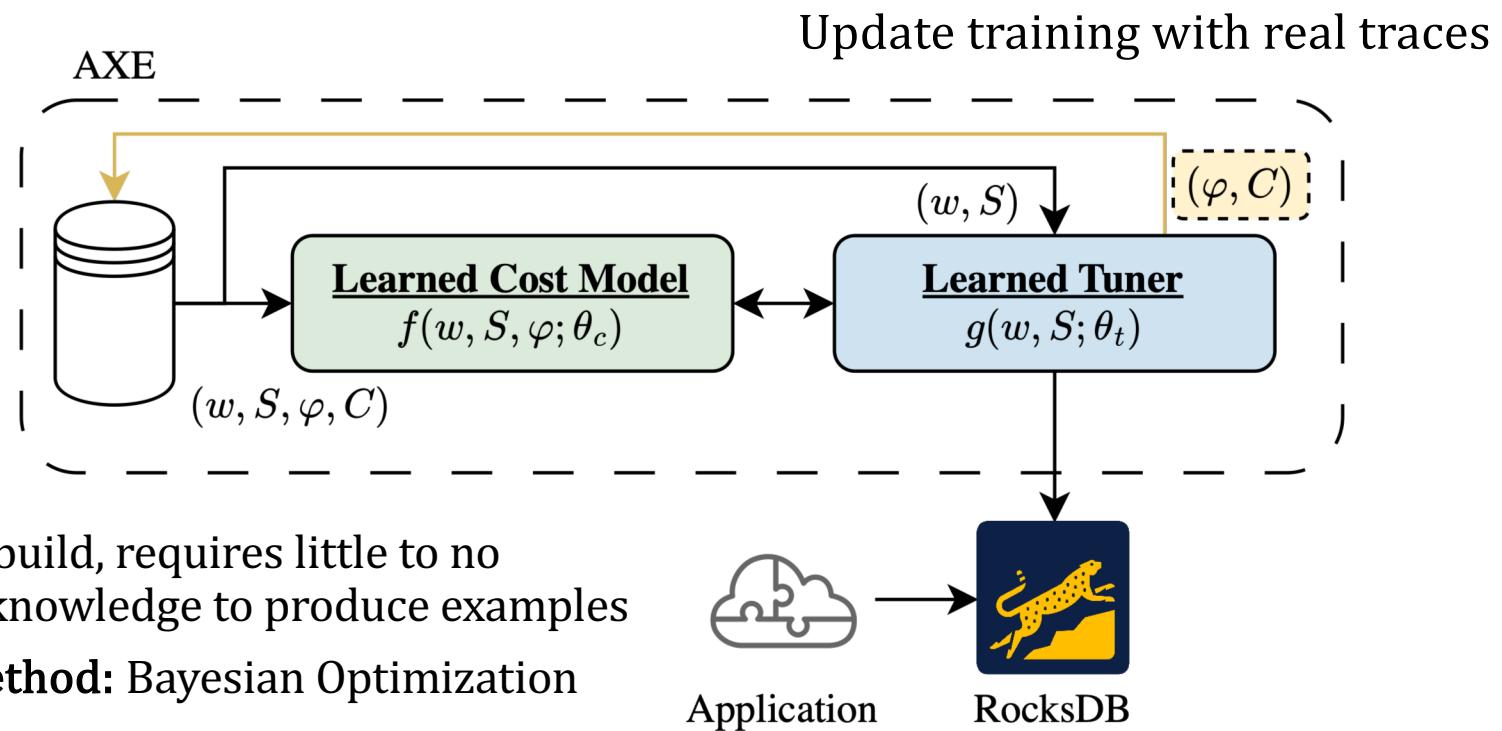
Building AXE



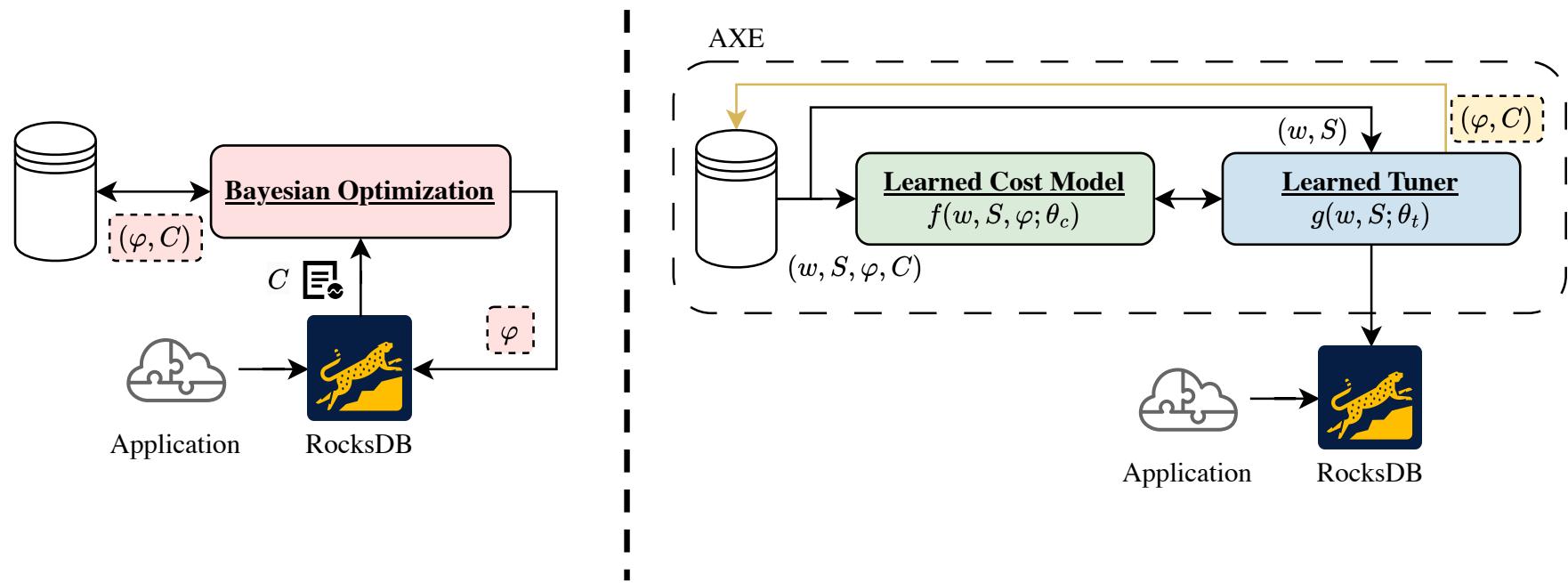
Building AXE



Building AXE



AXE Compared To BO

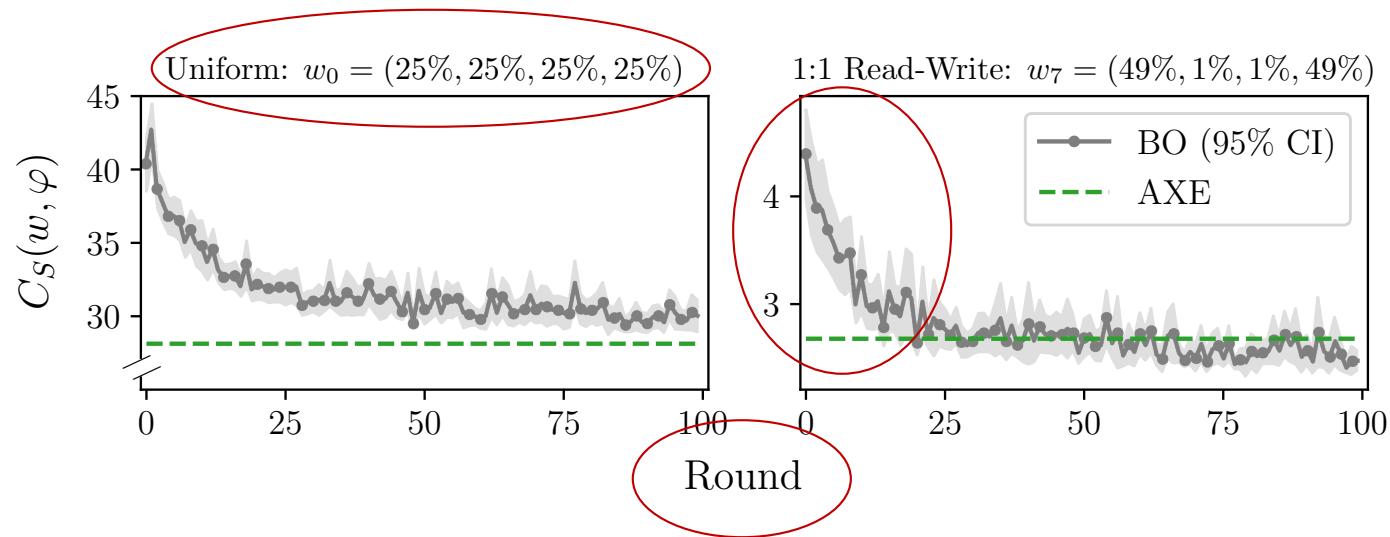


BO adapts to new workloads

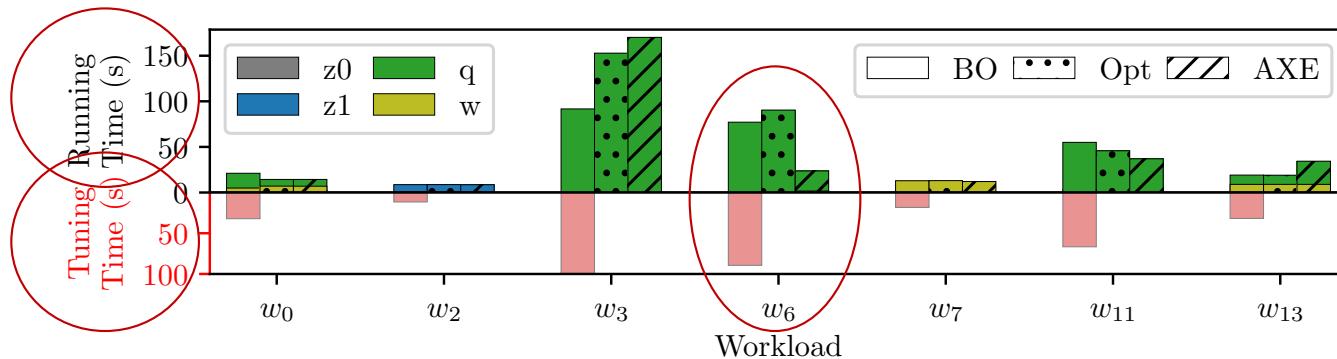
Mechanism for feedback is expensive

The Cost of AXE is Marginal

Pays a **hefty** upfront cost compared to AXE.



Tunings Deployed on RocksDB



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

AXE on RocksDB creates competitive tunings

Less domain expertise needed compared to an optimizer

Less upfront cost compared to iterative tuning

Thanks!

Learn more at

disc.bu.edu/



www.ndhuynh.com/

