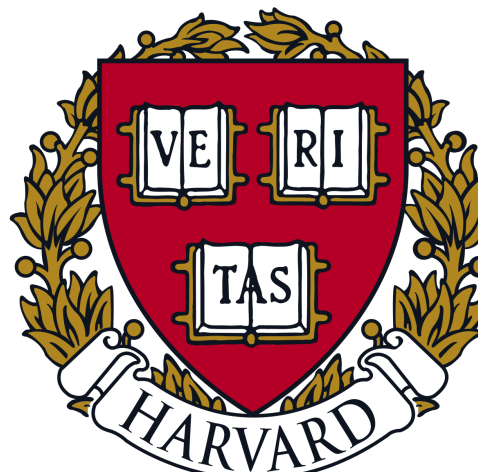# Cosine: A Cloud-Cost Optimized Self-Designing Key-Value Storage Engine
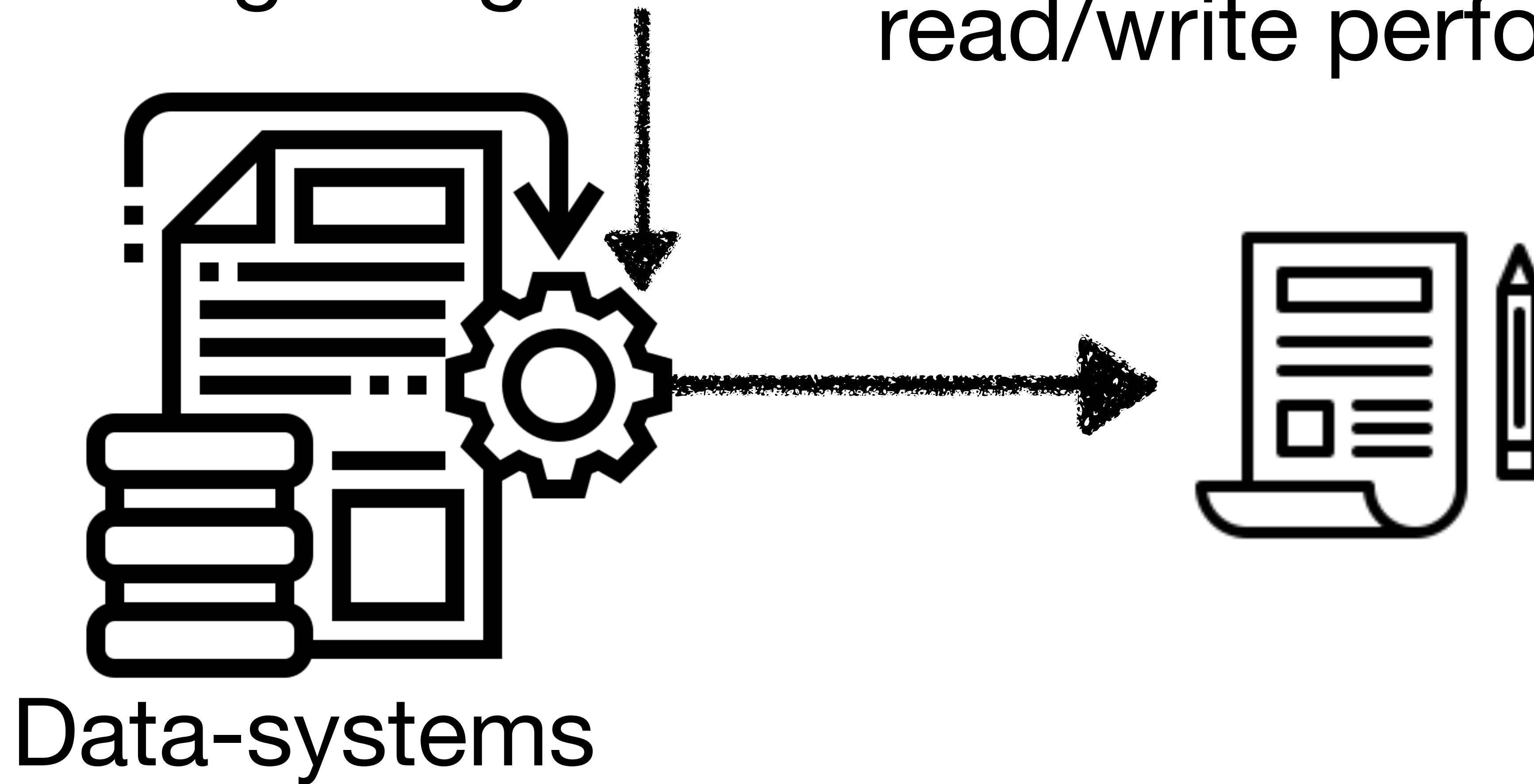
Subarna Chatterjee, Meena Jagadeesan,
Wilson Qin, Stratos Idreos
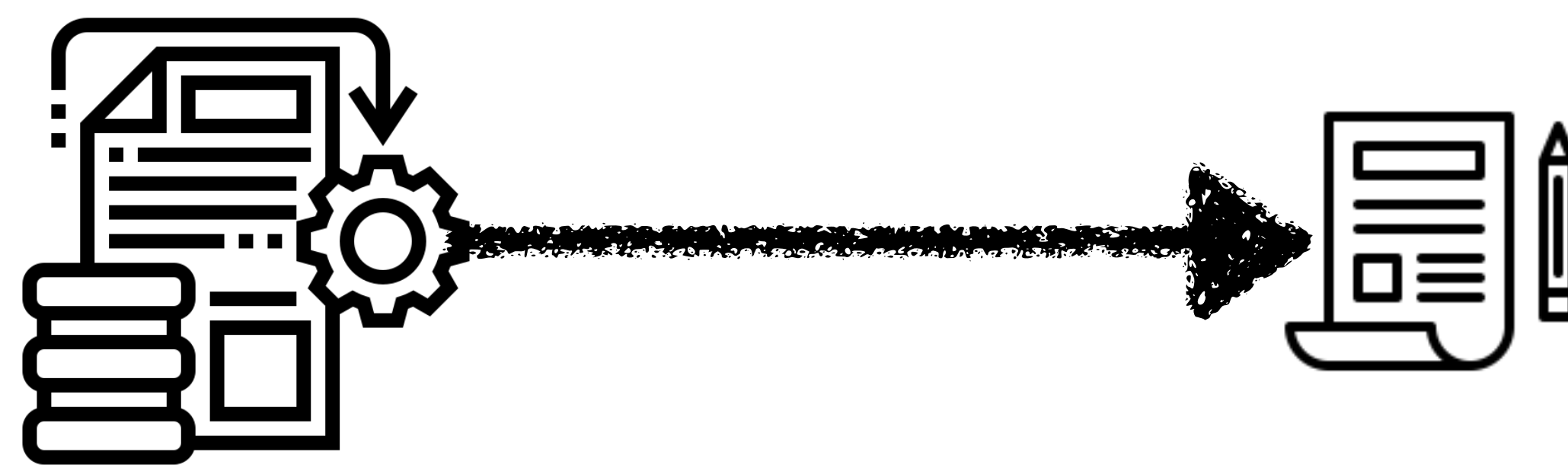
Data Systems Laboratory (DASLab)
Harvard University

storage-engines

read/write performance

Data-systems

# CONTEXT

data

hardware

applications
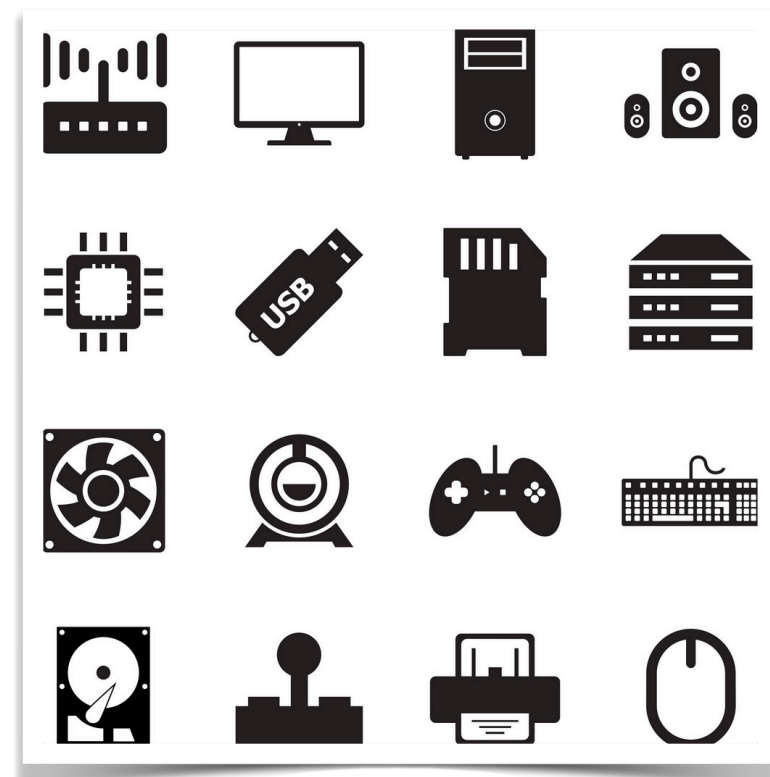
performance goals

# The **CONTEXT** keeps changing …
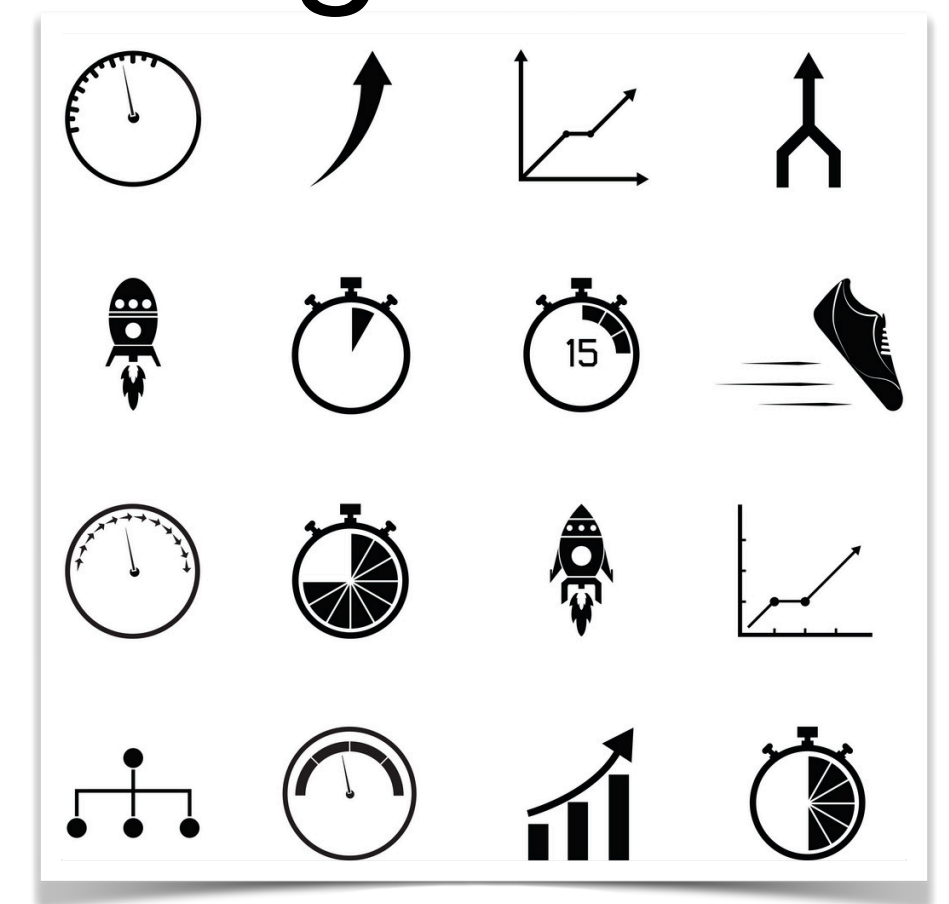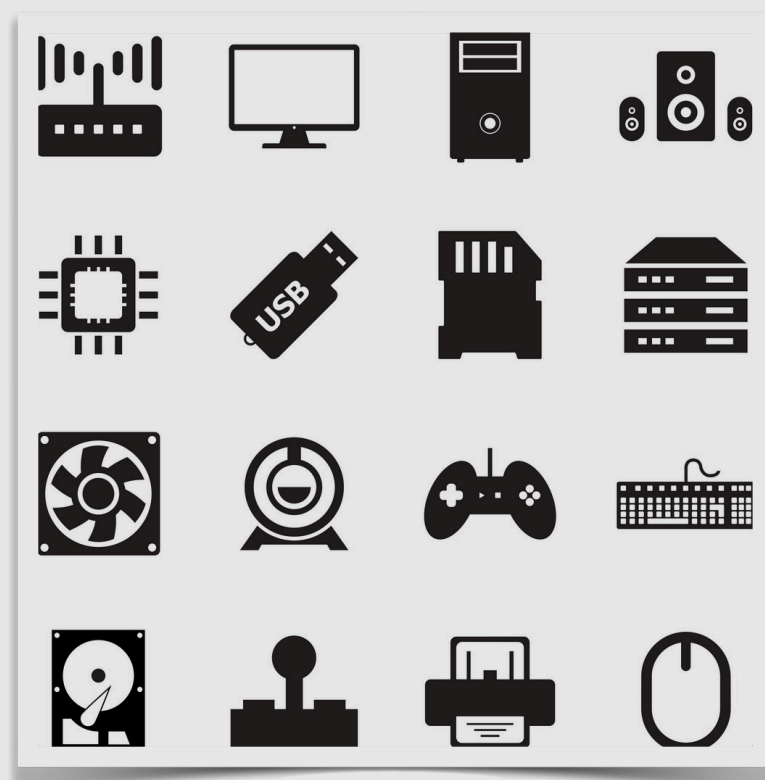


performance goals

hardware

applications
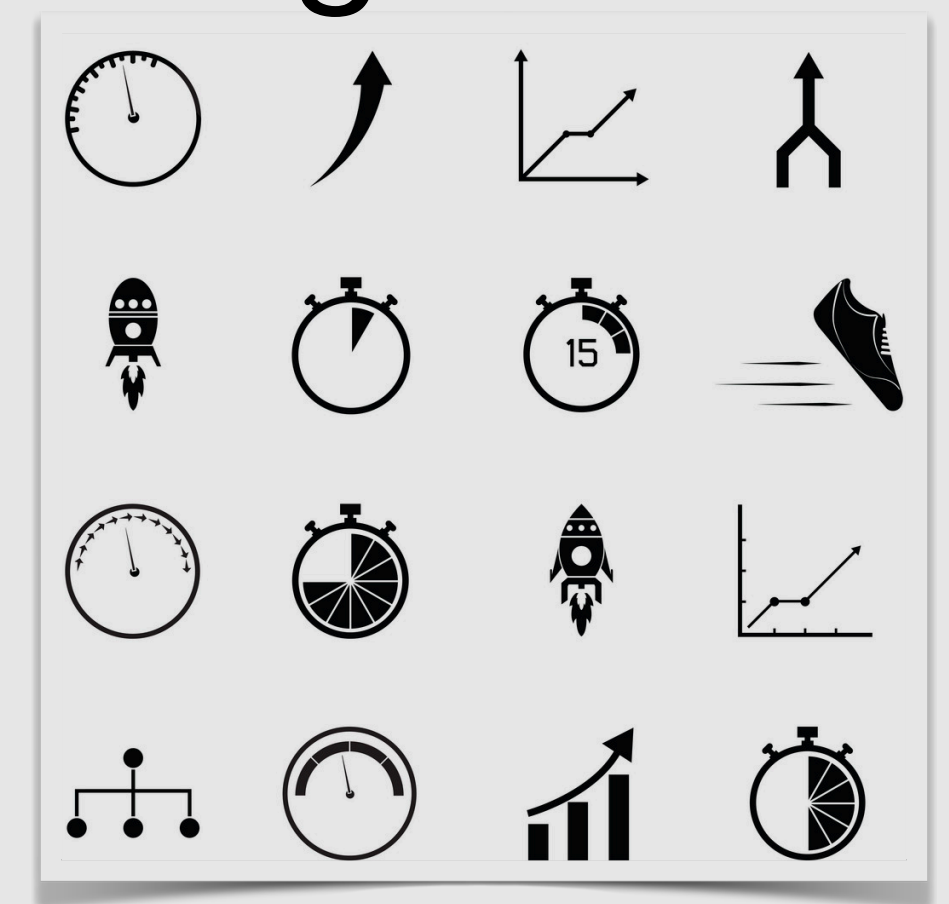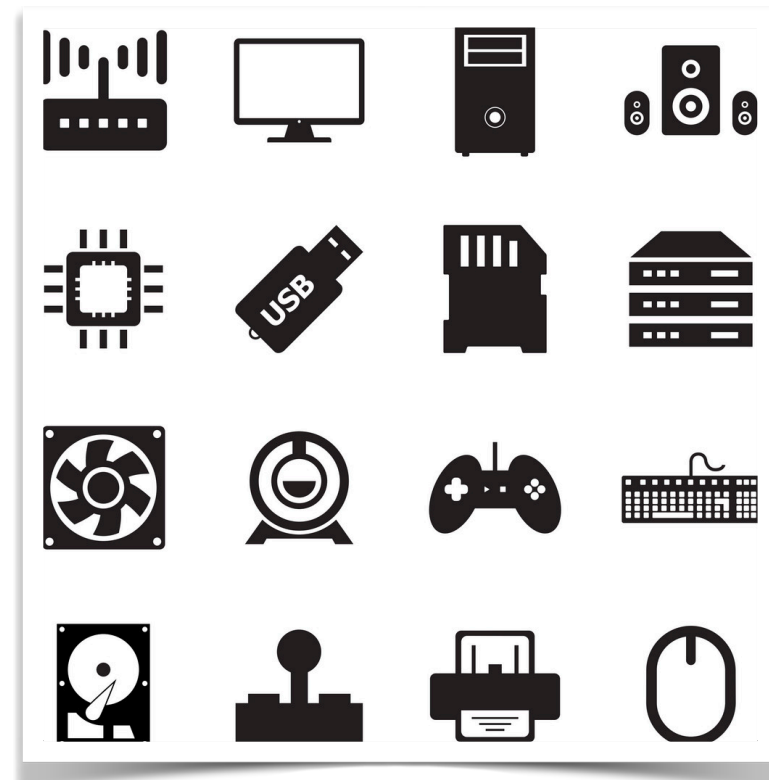
data

# The **CONTEXT** keeps changing …



data

hardware

applications

performance goals

# The **CONTEXT** keeps changing …



data

hardware

applications

performance goals

Bottleneck: Sub-Optimal Data Systems

# Bottleneck: Sub-Optimal Data Systems



**PERFORMANCE**

**COST**

Ryan Booth @that1guy_15 · Mar 4
Oh fun! Another **high cloud bill**. This is exactly how I wanted to spend the rest of my week...

Matt Getty @aspen · Feb 16
The **Cloud** is as **high** in the sky as the **bill** from AWS*

* for a tremendous amount of workloads

CiscoEvents @CiscoEvents · Jan 28
Is your **cloud bill** too **high**? Learn how you can control **cloud** costs with CloudCenter Suite.

Translucent Computing @translucentcomp · Jan 20
#Kubernetes **cloud** costs are getting out of hand. Working with many clients and different **cloud** service providers, in more than 70% of cases, we see the VM cluster nodes are underutilized, which leads to a **high cloud bill**.

DASlab
@ Harvard SEAS

**PERFORMANCE**

**COST**

Cloud-provider and
hardware space

VMs of different
capacity

different storage type
(SSD, HDD, EBS)

amazon
web services™

Microsoft
Azure

Google Cloud

DASlab
@ Harvard SEAS

**PERFORMANCE**

**COST**

LSM

write

read

memory

LSH

B-Trees

Storage-engine design space

Cloud-provider and hardware space

amazon web services™

Microsoft Azure

Google Cloud

VMs of different capacity

different storage type (SSD, HDD, EBS)

DASlab
@ Harvard SEAS

PERFORMANCE

COST

Performance space

write-optimized

FASTER
RocksDB
massive number of unexplored designs
WiredTiger

read-optimized

Storage-engine design space

LSM
B-Trees

write

read

memory

LSH

Cloud-provider and hardware space

amazon web services™

Microsoft Azure

Google Cloud

VMs of different capacity

different storage type (SSD, HDD, EBS)

DASlab
@ Harvard SEAS

PERFORMANCE

COST

Workload space

*blind update + insert*

*rmw + get*

*range + blind update + get*

*non-empty range + get*

*empty range + insert*

Performance space

massive number of unexplored designs

WiredTiger

RocksDB

FASTER

*write-optimized*

*read-optimized*

Storage-engine design space

LSM

B-Trees

LSH

write

read

memory

Cloud-provider and hardware space

amazon web services™

Microsoft Azure

Google Cloud

VMs of different capacity

different storage type (SSD, HDD, EBS)

DASlab
@ Harvard SEAS

Workload space

*empty range + insert*

*rmw + get*

*range + blind update + get*

*non-empty range + get*

*blind update + insert*

Performance space

massive number of unexplored designs

WiredTiger

FASTER

RocksDB

*write-optimized*

*read-optimized*

Storage-engine design space

LSM

B-Trees

write

read

memory

LSH

Cloud-provider and hardware space

amazon web services™

Microsoft Azure

Google Cloud

different storage type (SSD, HDD, EBS)

VMs of different capacity

Exhaustive Decision Space

Complex

Vast
(10^35 possibilities)

Manual
(Limited exploration
of systems)

DASlab
@ Harvard SEAS

Workload space

empty range + insert

rmw + get

range + blind update + get

non-empty range + get

blind update + insert

Complex

Vast (10^35 possibilities)

Manual (Limited exploration

GOAL: To create the perfect data system tailored for each context

Cloud-provider and hardware space

memory

amazon web services™

Microsoft Azure

Google Cloud different storage type (SSD, HDD, EBS)

VMs of different capacity

DASlab @ Harvard SEAS

Gilad David Maayan
Posted on Sep 25, 2020 · Updated on Dec 10, 2020

AWS to Azure: Making the Move

#azure   #aws

Both Amazon Web Services and Microsoft Azure are considered top cloud

# Key Intuition

Storage engine = amalgamation of data structures

# Key Intuition

Storage engine = amalgamation of data structures



buffer · filter · cache · index

# Key Intuition

Storage engine = amalgamation of data structures

# Key Intuition

Every-growing space
of data structures

Massive design space
of key-value stores



*and more ...*

What if we could **reason** about the **massive design space** of key-value storage engines?

# Key Intuition

Every-growing space of data structures → Massive design space of key-value stores



*and more ...*

What if we could **reason** about the **massive design space** of key-value storage engines?

# Key Intuition

few existing designs

Workload

Cloud budget

What if we could **reason** about the **massive design space** of key-value storage engines?

workload

$ budget

Cosine

optimal configuration

SE design    h/w    cloud provider

# Storage Engine Template

## Layout Primitives

buffer?   filters?

cache?   indexes?

hot-cold partition?   growth factor?

level size?   greediness of merge?

## Algorithmic Abstractions

restructuring strategy

filter design

*full*

*partial*

*block*
*file*
*run*

*hybrid*

indexing algorithm

*fence pointer*
*hash*

file picking strategy

*oldest merged*
*oldest flushed*
*oldest*

**Storage Engine Template**

Layout Primitives

buffer? filters?
cache? indexes?
hot-cold growth
partition? factor?
level size? greediness of merge?

Algorithmic Abstractions

restructuring strategy
filter indexing
design algorithm
file picking strategy

| MEMORY | DISK |
|---|---|
| | storage pattern?<br>{flat logs, hierarchical}<br><br>growth factor?<br>[1, .., multiples of block size]<br><br>level size?←■→<br>[1, .., L]<br><br>greediness of merge?<br>[1 (high), .., T (low)]<br><br>file size?<br>MB ... GB<br><br>hot-cold partition? |

# MEMORY

# DISK

buffer?

[1..M]

filters?

Bloom? Cuckoo?

...

[1..M]

indexes?

hash table?
zone map?

...

cache?

[1..M]

storage pattern?

{flat logs, hierarchical}

growth factor?

[1, .., multiples of block size]

level size?

[1, .., L]

greediness of merge?

[1 (high), .., T (low)]

file size?

MB ... GB

hot-cold partition?

# Superstructure



MEMORY

DISK

buffer?

storage pattern?

growth factor?

hot

filters?

level size?

greediness of merge?

indexes?

hot-cold partition?

cold

file size?

cache?

DASlab
@ Harvard SEAS

MEMORY

buffer?

filters?

indexes?

cache?

DISK

hot

cold

storage pattern?

growth factor?

level size?

greediness of merge?

hot-cold partition?

file size?

Leveled LSM

Tiered LSM

BTree

New layout

# Storage Engine Template

**Layout Primitives**

buffer?  filters?

cache?  indexes?

hot-cold partition?  growth factor?

level size?  greediness of merge?

# Storage Engine Template

**Layout Primitives**

buffer?     filters?

cache?     indexes?

hot-cold partition?     growth factor?

level size?     greediness of merge?

**Algorithmic Abstractions**

restructuring strategy

filter design

*block*   *file*   *run*

*full*   *partial*   *hybrid*

indexing algorithm

*fence pointer*   *hash*

file picking strategy

*oldest merged*   *oldest flushed*   *oldest*

DASlab
@ Harvard SEAS

# Storage Engine Template

**Layout Primitives**

buffer?    filters?

cache?    indexes?

hot-cold partition?    growth factor?

level size?    greediness of merge?

**Algorithmic Abstractions**

restructuring strategy

filter design

*full*

*partial*

*hybrid*

*block*

*file*

*run*

indexing algorithm

*fence pointer*

*hash*

file picking strategy

*oldest flushed*

*oldest merged*

DASlab
@ Harvard SEAS

## Storage Engine Template

Layout Primitives

buffer? filters? cache? indexes? hot-cold partition? growth factor? level size? greediness of merge?

Algorithmic Abstractions

restructuring strategy — full, partial, hybrid (block, file, run) · filter design · indexing algorithm · fence pointer, hash · file picking strategy · oldest merged, oldest flushed, oldest

LAYOUT PRIMITIVES ← → ALGORITHMIC ABSTRACTIONS

| | | # | Design Abstractions of Template | Type/Domain | Example templates for diverse data structures | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | RocksDB variants | WiredTiger variants | FASTER variants | A new design |
| Design and hardware specification — *initialized by search through engine design space* | | 1. | **Key size:** Denotes the size of keys in the workload. | unsigned int | auto-configured from the sample workload | | | |
| | | 2. | **Value size:** Denotes the size of values in the workload. All values are accepted as variable-length strings. | string/slice *max size set to 1 GB* | auto-configured from the sample workload | | | |
| | | 3. | **Size ratio (T):** The maximum number of entries in a block (e.g. growth factor in LSM trees or fanout of B-trees. | unsigned integer \| function (func) | [2,.. 32] | [32, 64, 128, 256, ..] | [1000, 1001, …] (T is large) | 2 |
| | | 4. | **Runs per hot level (K):** At what capacity hot levels are compacted. Rule: should be less than size ratio. | unsigned int | [1.. T] | | [T-1] | 7 |
| | | 5. | **Runs per cold level (Z):** At what capacity cold levels are compacted. Rule: should be less than size ratio. | unsigned int | [1.. T] | [1] | | 32 |
| | | 6. | **Logical block size (B):** Number of consecutive disk blocks. | unsigned int | [2048, 4096, …] | | | |
| | | 7. | **Buffer capacity** $(M_B)$**:** Denotes the amount of memory allocated to in-memory buffer/memtables. Configurable w.r.t file size. | 64-bit floating point \| function (func) | [64 MB, 128 MB, …] | [1 MB, 2 MB, …] | [64 MB, 128 MB, …] | h/w dependent |
| | | 8. | **Indexes** $(M_{FP})$**:** Amount of memory allocated to indexes (fence pointers/hashtables). | 64-bit floating point \| function (func) | memory to cover L | memory for first level | memory for hash table | h/w dependent |
| | | 9. | **Bloom filter memory** $(M_{BF})$**:** Denotes the bits/entry assigned to Bloom filters. | 64-bit float \| func(FPR) | 10 bits/key | | | func(FPR) |
| Data access — *derived with empirically verified rules* | | 10. | **Bloom filter design:** Denotes the granularity of Bloom filters, e.g., one Bloom filter instance per block or per file or per run. The default is file. | block \| file \| run | file | | | file |
| | | 11. | **Compaction/Restructuring algorithm:** Full does level-to-level compaction; partial is file-to-file; and hybrid uses both full and partial at separate levels. | partial \| full \| hybrid | full, partial | partial | partial | hybrid |
| | | 12. | **Run strategy:** Denotes which run to be picked for compaction (only for partial/hybrid compaction). | first \| last_full \| fullest | first, fullest, last_full | | first | fullest |
| | | 13. | **File picking strategy:** Denotes which file to be picked for compaction (for partial/hybrid compaction). For LSM-trees we set default to dense_fp as it empirically works the best. B-trees pick the first file found to be full. LSH-table restructures at the granularity of runs. | oldest_merged \| oldest_flushed \| dense_fp \| sparse_fp \| choose_first | dense_fp | choose_first | | dense_fp (hot), choose_first (cold) |
| Parallelism | | 14. | **Merge threshold:** If a level is more than x% full, a compaction is triggered. | 64-bit floating point | [0.7..1] | 0.5 | | 0.75 |
| | | 15. | **Full compaction levels:** Denotes how many levels will have full compaction (only for hybrid compaction). The default is set to 2. | unsigned integer \| function (func) | [1..L] | | | L-Y (from optimal config) |
| | | 16. | **No. of CPUs:** Number of available cores to use in a VM. | unsigned int | Use all available cores | | | |
| | | 17. | **No of threads:** Denotes how many threads are used to process the workload. | unsigned int | Use 1 thread per CPU core | | | |

DASlab @ Harvard SEAS

**Storage Engine Template**

Layout Primitives

buffer?　filters?
cache?　indexes?
hot-cold partition?　growth factor?
level size?　greediness of merge?

Algorithmic Abstractions

filter design
restructuring strategy
full
partial
hybrid
block
file
run
indexing algorithm
fence pointer
hash
oldest merged
oldest flushed
oldest
file picking strategy

**Distribution-Aware I/O Model**

Shape of SE

Per operation IO cost

probabilistic existence in disk

key1
key2
key3

L1　L2　L3　L4　L5

**Concurrency-Aware CPU Model**

Empirically learn proportion of parallelizable component

Get speedup coefficient

Get actual end-to-end performance

Amdahl's Law

DASlab
@ Harvard SEAS

**Storage Engine Template**

Layout Primitives
- buffer?
- filters?
- cache?
- indexes?
- hot-cold partition?
- growth factor?
- level size?
- greediness of merge?

Algorithmic Abstractions
- filter design
- restructuring strategy
  - *full*
  - *partial*
  - *hybrid*
  - *block*
  - *file*
  - *run*
- indexing algorithm
  - *fence pointer*
  - *hash*
- file picking strategy
  - *oldest merged*
  - *oldest flushed*
  - *oldest*

**Distribution-Aware I/O Model**

Shape of SE | Per operation IO cost

probabilistic existence in disk

key1, key2, key3

L1 L2 L3 L4 L5

**Concurrency-Aware CPU Model**

Empirically learn proportion of parallelizable component → Amdahl's Law → Get speedup coefficient → Get actual end-to-end performance

**Search Algorithms**

Search space of configurations

Cloud Provider and Hardware Space | Storage Engine Design Space

SSD, aws

Cost-Performance Continuum

*Pareto frontier*

Latency vs Cost

Navigation through search space

Repeat for all configurations

Co-optimize hardware and SLA

✓*reliability*
✓*fault tolerance* ✓*backup*
✓*DB migration*

Evaluate and rank configurations → Pick optimal

*optimized configuration for a single provider*

Latency vs Cost

DASlab @ Harvard SEAS

## Storage Engine Template

### Layout Primitives

buffer? filters?

cache? indexes?

hot-cold partition? growth factor?

level size? greediness of merge?

### Algorithmic Abstractions

filter design

restructuring strategy

*full*

*partial*

*hybrid*

*block*

*file*

*run*

indexing algorithm

*fence pointer*

*hash*

*oldest merged*

*oldest flushed*

*oldest*

file picking strategy

## Distribution-Aware I/O Model

Shape of SE

Per operation IO cost

probabilistic existence in disk

key1

key2

key3

L1 L2 L3 L4 L5

## Concurrency-Aware CPU Model

Empirically learn proportion of parallelizable component

Amdahl's Law

Get speedup coefficient

Get actual end-to-end performance

## Search Algorithms

### Search space of configurations

Cloud Provider and Hardware Space

Storage Engine Design Space

SSD

aws

### Cost-Performance Continuum

Latency

*Pareto frontier*

Cost

### Navigation through search space

Repeat for all configurations

Co-optimize hardware and SLA

✓*reliability*

✓*fault tolerance* ✓*backup*

✓*DB migration*

*optimized configuration for a single provider*

Latency

Cost

Evaluate and rank configurations

Pick optimal

## Cosine Engine

### Memory (layout primitive)

Filters Fences Buffer

### On-Disk Run (layout primitive)

LSM Run ✅

B-tree Run

Fanout Cascading Fence Pointers

### Restructuring Module (algorithm)

full

partial ✅

hybrid

### KV OPERATIONS
(get, put, range, update, rmw)

B+-Tree

LSM-Tree

LSHTable

New Design

Workload + Budget + Target Perf. + SLA Specs.

DASlab @ Harvard SEAS

MODEL

(latency)

MODEL

(latency)

I/O          concurrency

DASlab
@ Harvard SEAS

MODEL

(latency)

concurrency

CPU

abstraction

Storage

I/O

DASlab
@ Harvard SEAS

(latency)

analytical ? learned ?

**10^35 possibilities**

DASlab
@ Harvard SEAS

I/O
analytical

concurrency
learned

$L$

$\dfrac{L}{k}$

(latency)

**1** design morphology

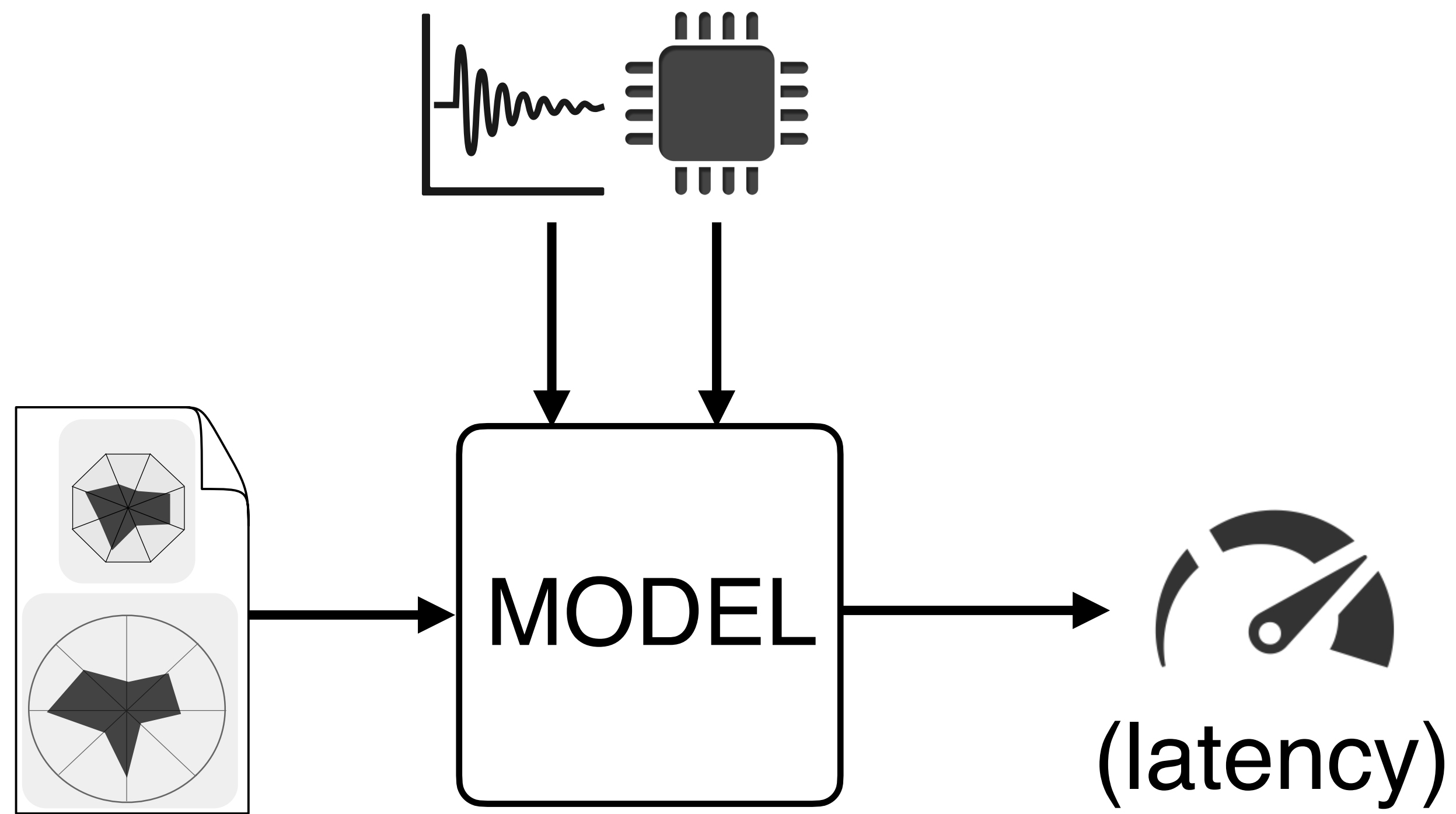| Design Abstractions of Template | Type/Domain | Example templates for diverse data structures | | | |
|---|---|---|---|---|---|
| | | LSM variants | B-Tree variants | LSH variants | A new design |
| 1. **Key size:** Denotes the size of keys in the workload. | unsigned int | auto-configured from the sample workload | | | |
| 2. **Value size:** Denotes the size of values in the workload. All values are accepted as variable-length strings. | string/slice *max size set to 1 GB* | auto-configured from the sample workload | | | |
| 3. **Size ratio (T):** The maximum number of entries in a block (e.g. growth factor in LSM trees or fanout of B-trees. | unsigned integer \| function (func) | [2,.. 32] | [32, 64, 128, 256, ..] | [1000, 1001, …] (T is large) | 2 |
| 4. **Runs per hot level (K):** At what capacity hot levels are compacted. Rule: should be less than size ratio. | unsigned int | [1.. T] | | [T-1] | 7 |
| 5. **Runs per cold level (Z):** At what capacity cold levels are compacted. Rule: should be less than size ratio. | unsigned int | [1.. T] | [1] | | 32 |
| 6. **Logical block size (B):** Number of consecutive disk blocks. | unsigned int | [2048, 4096, …] | | | |
| 7. **Buffer capacity** $(M_B)$**:** Denotes the amount of memory allocated to in-memory buffer/memtables. Configurable w.r.t file size. | 64-bit floating point \| function (func) | [64 MB, 128 MB, …] | [1 MB, 2 MB, …] | [64 MB, 128 MB, …] | h/w dependent |
| 8. **Indexes** $(M_{FP})$**:** Amount of memory allocated to indexes (fence pointers/hashtables). | 64-bit floating point \| function (func) | memory to cover L | memory for first level | memory for hash table | h/w dependent |
| 9. **Bloom filter memory** $(M_{BF})$**:** Denotes the bits/entry assigned to Bloom filters. | 64-bit float \| func(FPR) | 10 bits/key | | | func(FPR) |
| 10. **Bloom filter design:** Denotes the granularity of Bloom filters, e.g., one Bloom filter instance per block or per file or per run. The default is file. | block \| file \| run | file | | | file |
| 11. **Compaction/Restructuring algorithm:** Full does level-to-level compaction; partial is file-to-file; and hybrid uses both full and partial at separate levels. | partial \| full \| hybrid | full, partial | partial | partial | hybrid |
| 12. **Run strategy:** Denotes which run to be picked for compaction (only for partial/ hybrid compaction). | first \| last_full \| fullest | first, fullest, last_full | | first | fullest |
| 13. **File picking strategy:** Denotes which file to be picked for compaction (for partial/ hybrid compaction). For LSM-trees we set default to dense_fp as it empirically works the best. B-trees pick the first file found to be full. LSH-table restructures at the granularity of runs. | oldest_merged \| oldest_flushed \| dense_fp \| sparse_fp \| choose_first | dense_fp | choose_first | | dense_fp (hot), choose_first (cold) |
| 14. **Merge threshold:** If a level is more than x% full, a compaction is triggered. | 64-bit floating point | [0.7..1] | 0.5 | | 0.75 |
| 15. **Full compaction levels:** Denotes how many levels will have full compaction (only for hybrid compaction). The default is set to 2. | unsigned integer \| function (func) | [1..L] | | | L-Y (from optimal config) |
| 16. **No. of CPUs:** Number of available cores to use in a VM. | unsigned int | Use all available cores | | | |
| 17. **No of threads:** Denotes how many threads are used to process the workload. | unsigned int | Use 1 thread per CPU core | | | |

LAYOUT PRIMITIVES — ALGORITHMIC ABSTRACTIONS

Design and hardware specification (initialized by search through engine design space)

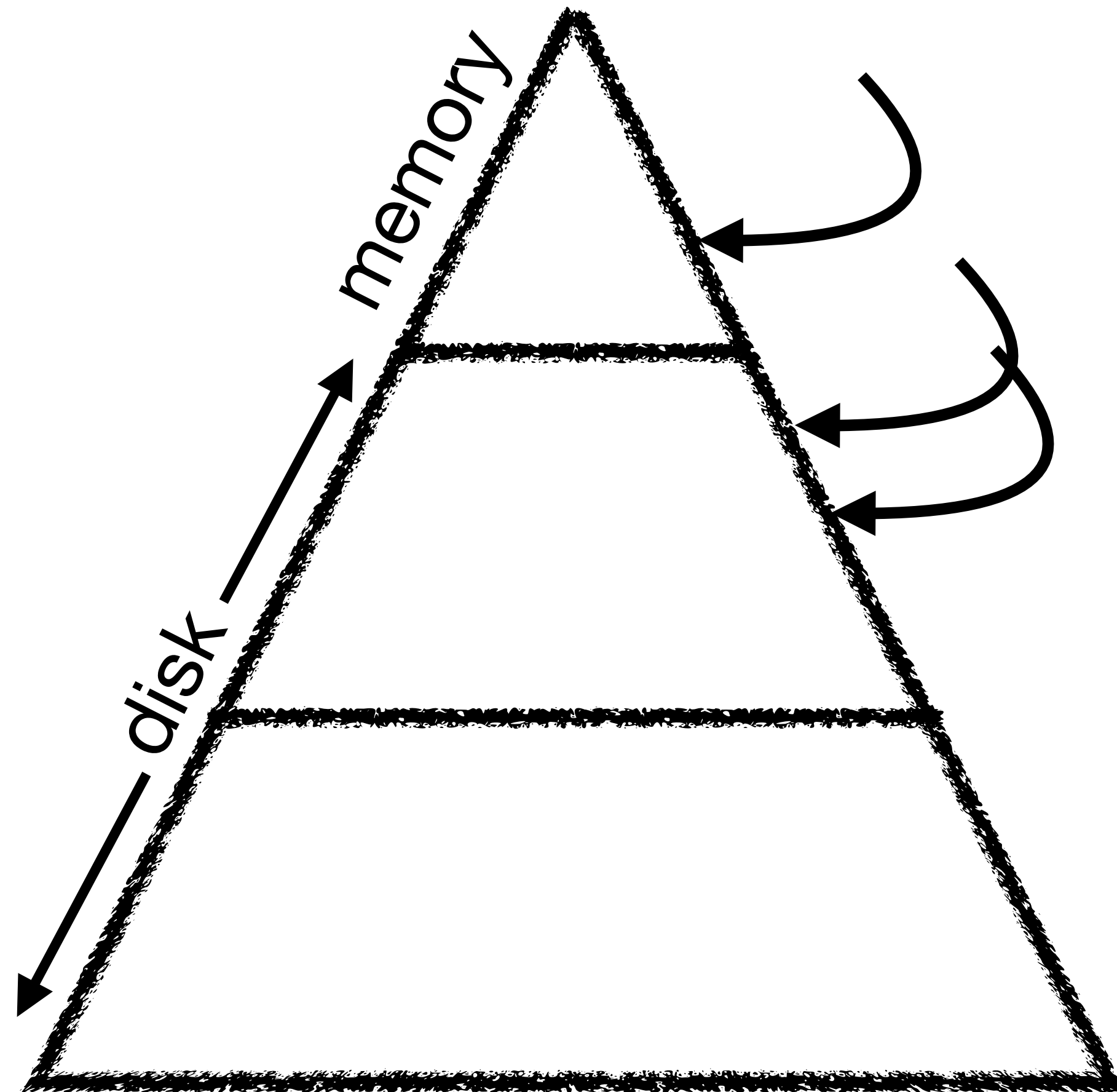Data access / Parallelism (derived with empirically verified rules)

DASlab @ Harvard SEAS

# ① design morphology

| | | | Design Abstractions of Template | Type/Domain | Example templates for diverse data structures | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | **LSM variants** | **B-Tree variants** | **LSH variants** | **A new design** |
| LAYOUT PRIMITIVES | Design and hardware specification / initialized by search through engine design space | 1. | **Key size:** Denotes the size of keys in the workload. | unsigned int | auto-configured from the sample workload | | | |
| | | 2. | **Value size:** Denotes the size of values in the workload. All values are accepted as variable-length strings. | string/slice *max size set to 1 GB* | auto-configured from the sample workload | | | |
| | | 3. | **Size ratio (T):** The maximum number of entries in a block (e.g. growth factor in LSM trees or fanout of B-trees). | unsigned integer \| function (func) | [2,.. 32] | [32, 64, 128, 256, ..] | [1000, 1001, …] (T is large) | 2 |
| | | 4. | **Runs per hot level (K):** At what capacity hot levels are compacted. Rule: should be less than size ratio. | unsigned int | [1.. T] | | [T-1] | 7 |
| | | 5. | **Runs per cold level (Z):** At what capacity cold levels are compacted. Rule: should be less than size ratio. | unsigned int | [1.. T] | [1] | | 32 |
| | | 6. | **Logical block size (B):** Number of consecutive disk blocks. | unsigned int | [2048, 4096, …] | | | |
| | | 7. | **Buffer capacity** $(M_B)$**:** Denotes the amount of memory allocated to in-memory buffer/memtables. Configurable w.r.t file size. | 64-bit floating point \| function (func) | [64 MB, 128 MB, …] | [1 MB, 2 MB, …] | [64 MB, 128 MB, …] | h/w dependent |
| | | 8. | **Indexes** $(M_{FP})$**:** Amount of memory allocated to indexes (fence pointers/hashtables). | 64-bit floating point \| function (func) | memory to cover L | memory for first level | memory for hash table | h/w dependent |
| | | 9. | **Bloom filter memory** $(M_{BF})$**:** Denotes the bits/entry assigned to Bloom filters. | 64-bit float \| func(FPR) | 10 bits/key | | | func(FPR) |
| ALGORITHMIC ABSTRACTIONS | Data access / derived with empirically verified rules | 10. | **Bloom filter design:** Denotes the granularity of Bloom filters, e.g., one Bloom filter instance per block or per file or per run. The default is file. | block \| file \| run | file | | | file |
| | | 11. | **Compaction/Restructuring algorithm:** Full does level-to-level compaction; partial is file-to-file; and hybrid uses both full and partial at separate levels. | partial \| full \| hybrid | full, partial | partial | partial | hybrid |
| | | 12. | **Run strategy:** Denotes which run to be picked for compaction (only for partial/ hybrid compaction). | first \| last_full \| fullest | first, fullest, last_full | | first | fullest |
| | | 13. | **File picking strategy:** Denotes which file to be picked for compaction (for partial/ hybrid compaction). For LSM-trees we set default to dense_fp as it empirically works the best. B-trees pick the first file found to be full. LSH-table restructures at the granularity of runs. | oldest_merged \| oldest_flushed \| dense_fp \| sparse_fp \| choose_first | dense_fp | choose_first | | dense_fp (hot), choose_first (cold) |
| | | 14. | **Merge threshold:** If a level is more than x% full, a compaction is triggered. | 64-bit floating point | [0.7..1] | 0.5 | | 0.75 |
| | | 15. | **Full compaction levels:** Denotes how many levels will have full compaction (only for hybrid compaction). The default is set to 2. | unsigned integer \| function (func) | [1..L] | | | L-Y (from optimal config) |
| | Parallelism | 16. | **No. of CPUs:** Number of available cores to use in a VM. | unsigned int | Use all available cores | | | |
| | | 17. | **No of threads:** Denotes how many threads are used to process the workload. | unsigned int | Use 1 thread per CPU core | | | |

bulkwrites · read · universe size · writes · uniform · size ratio · write skew proportion · cold merge threshold · read skew proportion · hot merge threshold

**1** design morphology

**2** per operation I/O cost

**1** design morphology  **2** per operation I/O cost
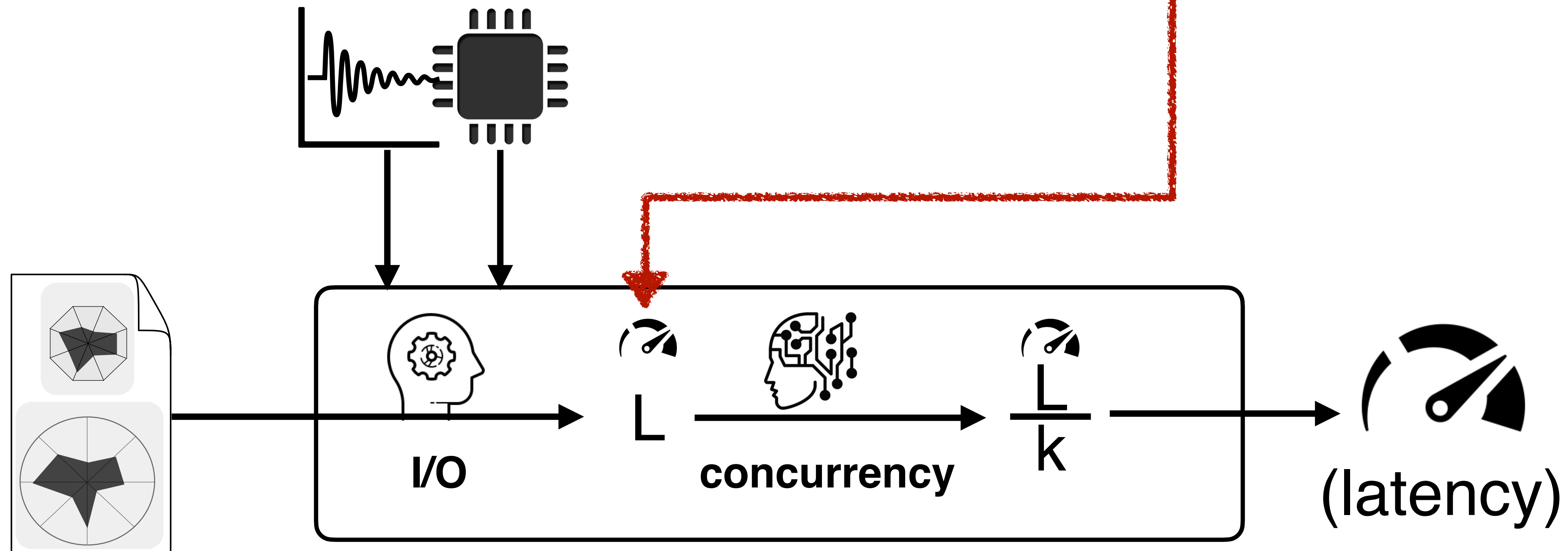


early-stopping

infrequent merging

**1** design morphology　**2** per operation I/O cost　**3** total cost

$$\text{latency} = \frac{\text{Total I/O}_{\text{workload}}}{\text{IOPS}} = L$$

# ① design morphology  ② per operation I/O cost  ③ total cost

$$\text{latency} = \frac{\text{Total I/O}_{\text{workload}}}{\text{IOPS}} = L$$



I/O

L    concurrency    $\frac{L}{k}$

(latency)

# Amdahl's Law (1967)



sequential (1-f)    parallel (f)    1 core

T

2 core

4 core

8 core

theoretical speedup k = $\dfrac{T}{T - fT + fT/n}$ = $\dfrac{1}{1 - f(1 - 1/n)}$

DASlab
@ Harvard SEAS

| Design class | | | | |
|---|---|---|---|---|
| LSM | Hybrid-1 | Hybrid-2 | B-tree | LSH |
| r1 | r2 | r3 | r4 | r5 |
| w1 | w2 | w3 | w4 | w5 |

H/W 1: D16s v3 (16 vcpus, 64 GB RAM), 25600 IOPS, premium SSD
H/W 2: D8s v3 (8 vcpus, 32 GB RAM), 12800 IOPS, premium SSD
H/W 3: D16a v4 (16 vcpus, 64 GB RAM), 32x500 IOPS, non-premium SSD

| (Base data, #reads) | H/W 1 | H/W 2 | H/W 3 |
|---|---|---|---|
| 10M, 1M | 0.965 | 0.953 | 0.95 |
| 50M, 5M | 0.94 | 0.94 | 0.92 |
| 100M, 10M | 0.942 | 0.952 | 0.923 |

| Ops | Design class | | | | |
|---|---|---|---|---|---|
| | LSM | Hybrid-1 | Hybrid-2 | B-tree | LSH |
| 📖 | r1 | r2 | r3 | r4 | r5 |
| 🖋 | w1 | w2 | w3 | w4 | w5 |

DASlab
@ Harvard SEAS

| Ops | Design class | | | | |
|---|---|---|---|---|---|
| | LSM | Hybrid-1 | Hybrid-2 | B-tree | LSH |
| 📖 | 0.91 | 0.93 | 0.92 | 0.94 | 0.97 |
| 🖊 | 0.71 | 0.79 | 0.74 | 0.65 | 0.8 |

20% reads +
80% writes

proportion of parallelizable code = 0.2*0.91 + 0.8*0.71 = 0.75

DASlab
@ Harvard SEAS

SEARCH ALGORITHM

MODEL

Workload space

empty range + insert

rmw + get

range + blind update + get

non-empty range + get

blind update + insert

Performance space

massive number of unexplored designs

WiredTiger

RocksDB

FASTER

write-optimized

read-optimized

Storage-engine design space

LSM

B-Trees

write

read

memory

LSH

Cloud-provider and hardware space

amazon web services™

Microsoft Azure

Google Cloud

VMs of different capacity

different storage type (SSD, HDD, EBS)

YCSB E variant (30% blind update, 20% non-empty range, 50% empty range)

Throughput (kops/s)

$10^7$
$10^6$
$10^5$
$10^4$
$10^3$

20K    60K    100K

Budget ($/month)

YCSB E variant (30% blind update, 20% non-empty range, 50% empty range)

YCSB E variant (30% blind update, 20% non-empty range, 50% empty range)

YCSB E variant (30% blind update, 20% non-empty range, 50% empty range)

**Legend:** RocksDB · WiredTiger · FASTER · Cosine on AWS · Cosine on GCP · Cosine on Azure · LSM class · B-tree class · LSH class · Hybrid class

*YCSB A variant: blind-update-intensive (10% lookups, 20% rmws, 70% blind updates)*

**(A)** (32,1,1), (160,10), (E20 v3, Azure) (32,16) (1054,1053,1053), (E32 v3, Azure)

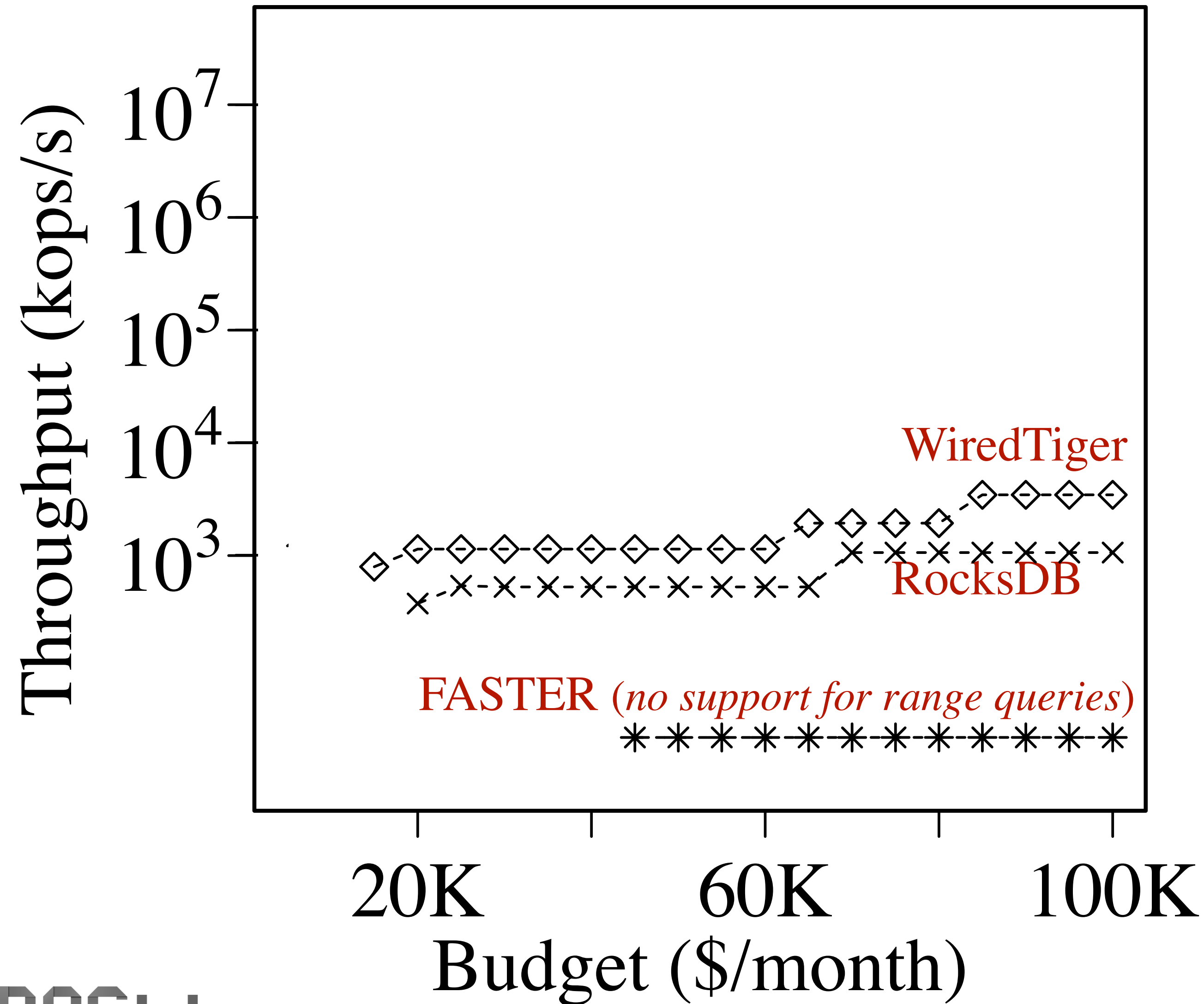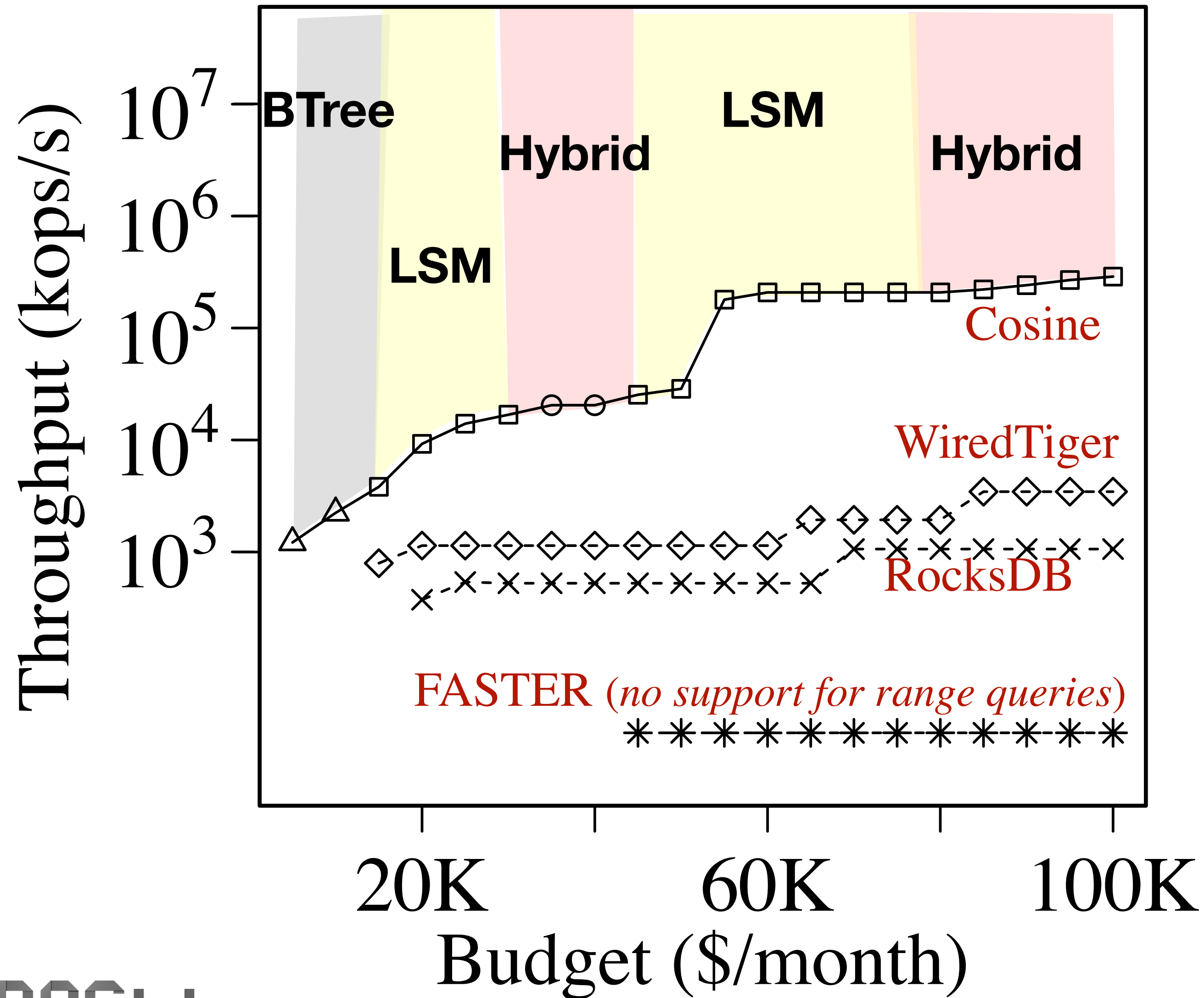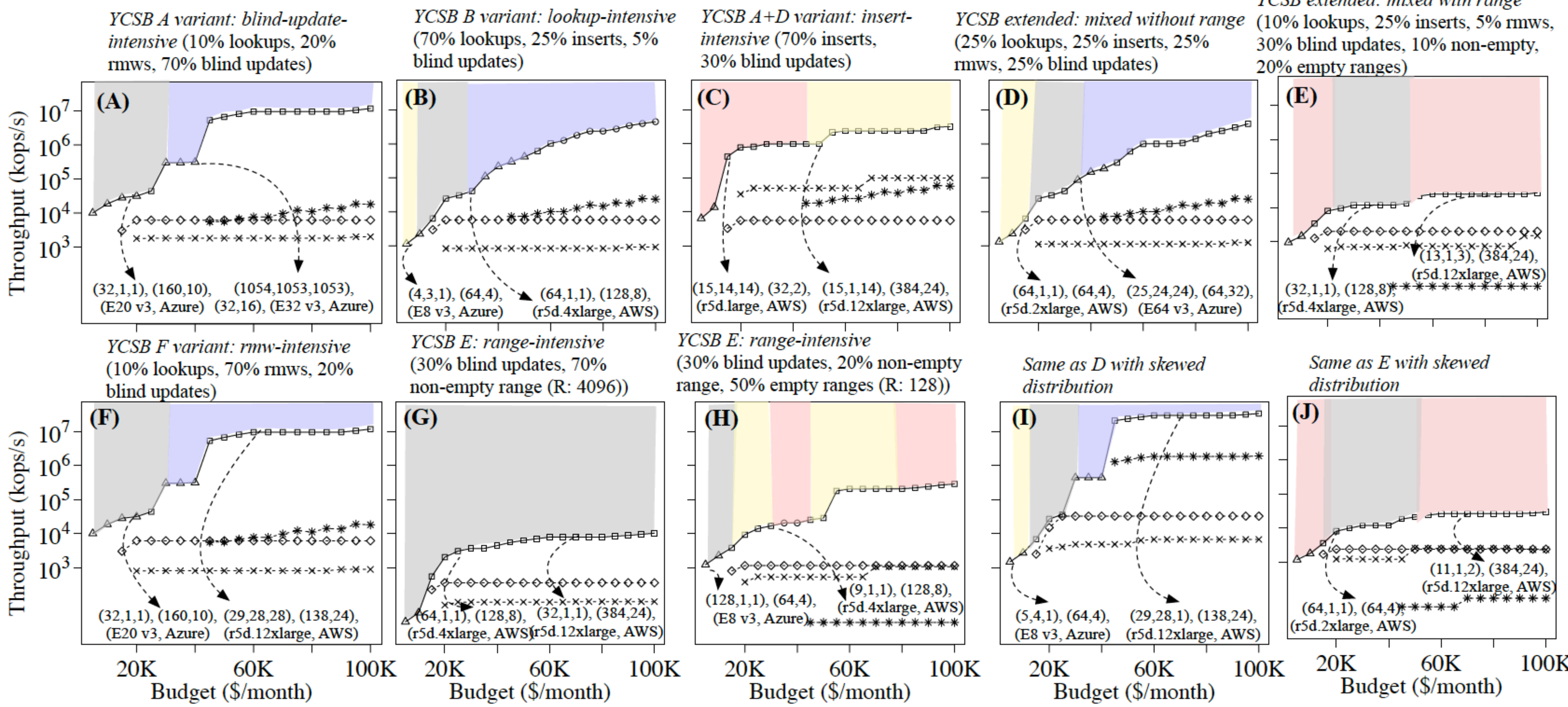*YCSB B variant: lookup-intensive (70% lookups, 25% inserts, 5% blind updates)*

**(B)** (4,3,1), (64,4), (E8 v3, Azure) (64,1,1), (128,8), (r5d.4xlarge, AWS)

*YCSB A+D variant: insert-intensive (70% inserts, 30% blind updates)*

**(C)** (15,14,14), (32,2), (r5d.large, AWS) (15,1,14), (384,24), (r5d.12xlarge, AWS)

*YCSB extended: mixed without range (25% lookups, 25% inserts, 25% rmws, 25% blind updates)*

**(D)** (64,1,1), (64,4), (r5d.2xlarge, AWS) (25,24,24), (64,32), (E64 v3, Azure)

*YCSB extended: mixed with range (10% lookups, 25% inserts, 5% rmws, 30% blind updates, 10% non-empty, 20% empty ranges)*

**(E)** (32,1,1), (128,8), (r5d.4xlarge, AWS) (13,1,3), (384,24), (r5d.12xlarge, AWS)

*YCSB F variant: rmw-intensive (10% lookups, 70% rmws, 20% blind updates)*

**(F)** (32,1,1), (160,10), (E20 v3, Azure) (29,28,28), (138,24), (r5d.12xlarge, AWS)

*YCSB E: range-intensive (30% blind updates, 70% non-empty range (R: 4096))*

**(G)** (64,1,1), (128,8), (r5d.4xlarge, AWS) (32,1,1), (384,24), (r5d.12xlarge, AWS)

*YCSB E: range-intensive (30% blind updates, 20% non-empty range, 50% empty ranges (R: 128))*

**(H)** (128,1,1), (64,4), (E8 v3, Azure) (9,1,1), (128,8), (r5d.4xlarge, AWS)

*Same as D with skewed distribution*

**(I)** (5,4,1), (64,4), (E8 v3, Azure) (29,28,1), (138,24), (r5d.12xlarge, AWS)

*Same as E with skewed distribution*

**(J)** (64,1,1), (64,4), (r5d.2xlarge, AWS) (11,1,2), (384,24), (r5d.12xlarge, AWS)