

os 265

Stratos Idreos

BIG DATA SYSTEMS

NoSQL | Neural Networks | Image AI | LLMs | Data Science

Paper/discussion phase to start after spring break.

- 1.FASTER: A Concurrent Key-Value Store with In-Place Updates.
- 2.The FastLanes Compression Layout: Decoding >100 Billion Integers per Second with Scalar Code.
- 3.PyTorch FSDP: Experiences on Scaling Fully Sharded Data Parallel
- 4.Megatron LM: Training Multi-Billion Parameter Language Models Using Model Parallelism.
- 5.RAG-design from a system design space
6. Overlap vector search time with LLM inference time
- 7.CROSSBOW: Scaling Deep Learning with Small Batch Sizes on Multi-GPU Servers
- 8.Nautilus: An Optimized System for Deep Transfer Learning over Evolving Training Datasets

Presentations

Task: This week:

Register for presentation. First come first serve. 4 students per paper max.

Link on Forum

Task: **Prep for presentation** — >Slides, Content, Delivery, Q&A

Read the paper several times.

Understand in-depth.

Use labs for questions.

After having a slides draft, meet with the TFs to get feedback.

TFs will schedule a meeting the week before class.

Preparing for presentations and reviews

review and slides should answer:

what is the problem

why is it important

why is it hard

why existing solutions do not work

what is the core intuition for the solution

solution step by step

does the paper prove its claims

exact setup of analysis/experiments

are there any gaps in the logic/proof

possible next steps

* follow a few citations to gain more background

And final question:

What can we use from this paper
technique or inspiration
towards self-designing systems?

how to prepare slides

no bullets 2 colors big text images animation for examples

how to prepare slides

no bullets 2 colors big text images animation for examples

story

one message per slide connection from slide to slide

how to prepare slides

no bullets 2 colors big text images animation for examples

story

one message per slide connection from slide to slide

all questions + discussion

add 3-4 interesting discussion points

Debug before presenting

1. Content: Did you answer all reviews questions?
2. Formatting: Did you apply all slide format rules?

Our goal:

Every presentation is better than the previous one

During discussion we will also touch on the presentation in every class

Grading will be relative to the presentation schedule

Reviews:

Submit in Canvas

One PDF page. 2-3 paragraphs max. Organize by Review questions.

The deadline is 30 minutes before class. Hard deadline.

Today: one more lecture on self-designing systems (Tuesday)

Then: 1 lecture on fast NN training through data movement optimization
(systems project)

One lecture on LLM+RAG

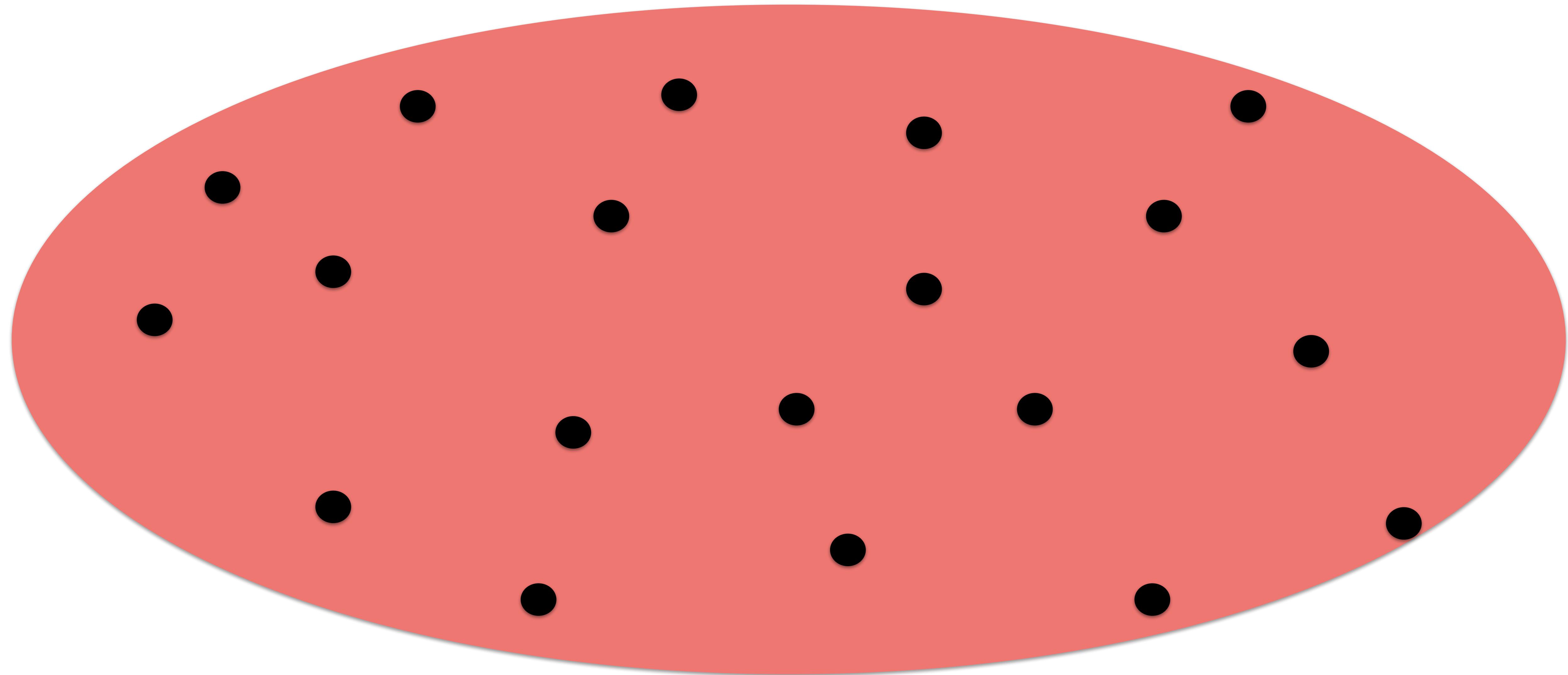
Then 3 lectures on end-to-end AI through reimagining storage

Then spring break

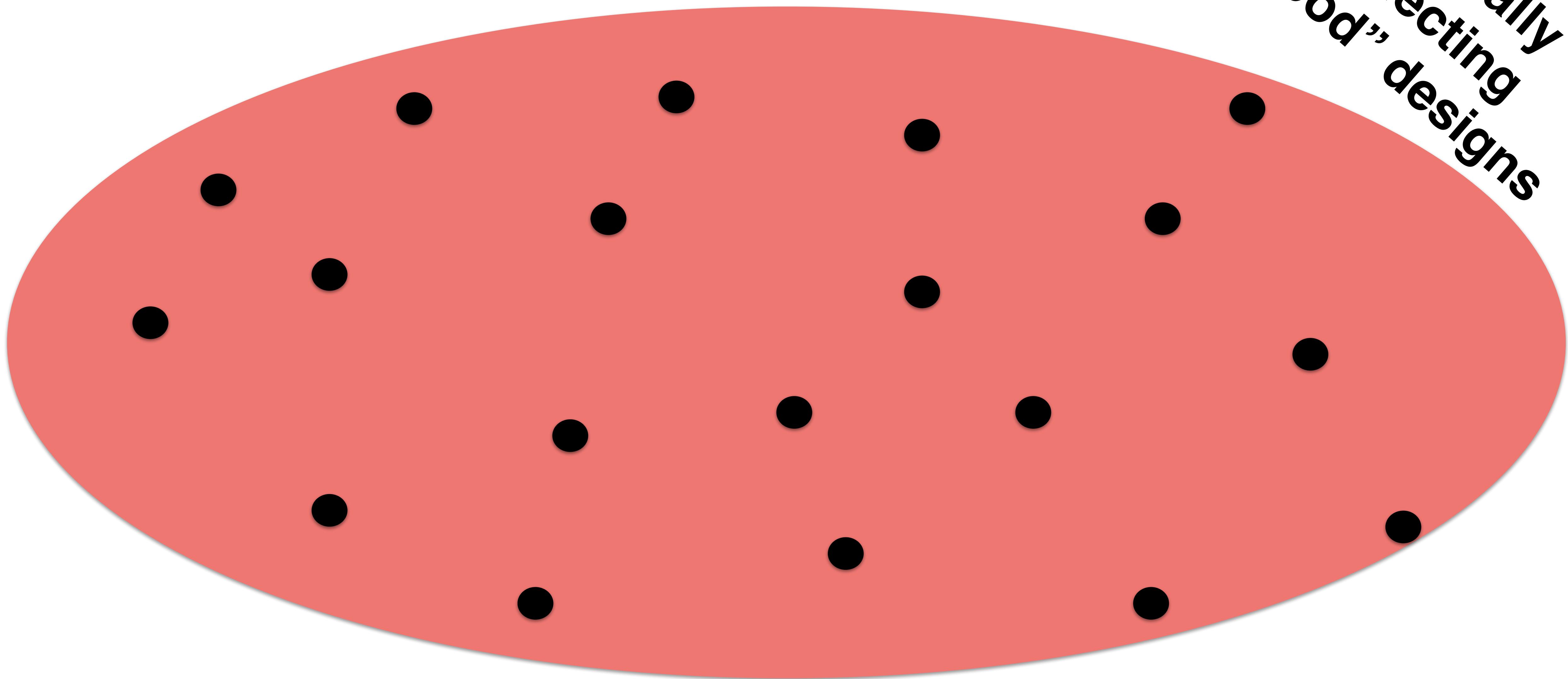
Then papers presentations

Systems projects should start coding this week
Research projects should know projects in next two weeks

The design space of systems is even larger



manually
selecting
“good” designs



*manually
selecting
“good”
designs*

**B-tree based
KV-system**

**LSM-tree based
KV-system**

**LSH based
KV-system**

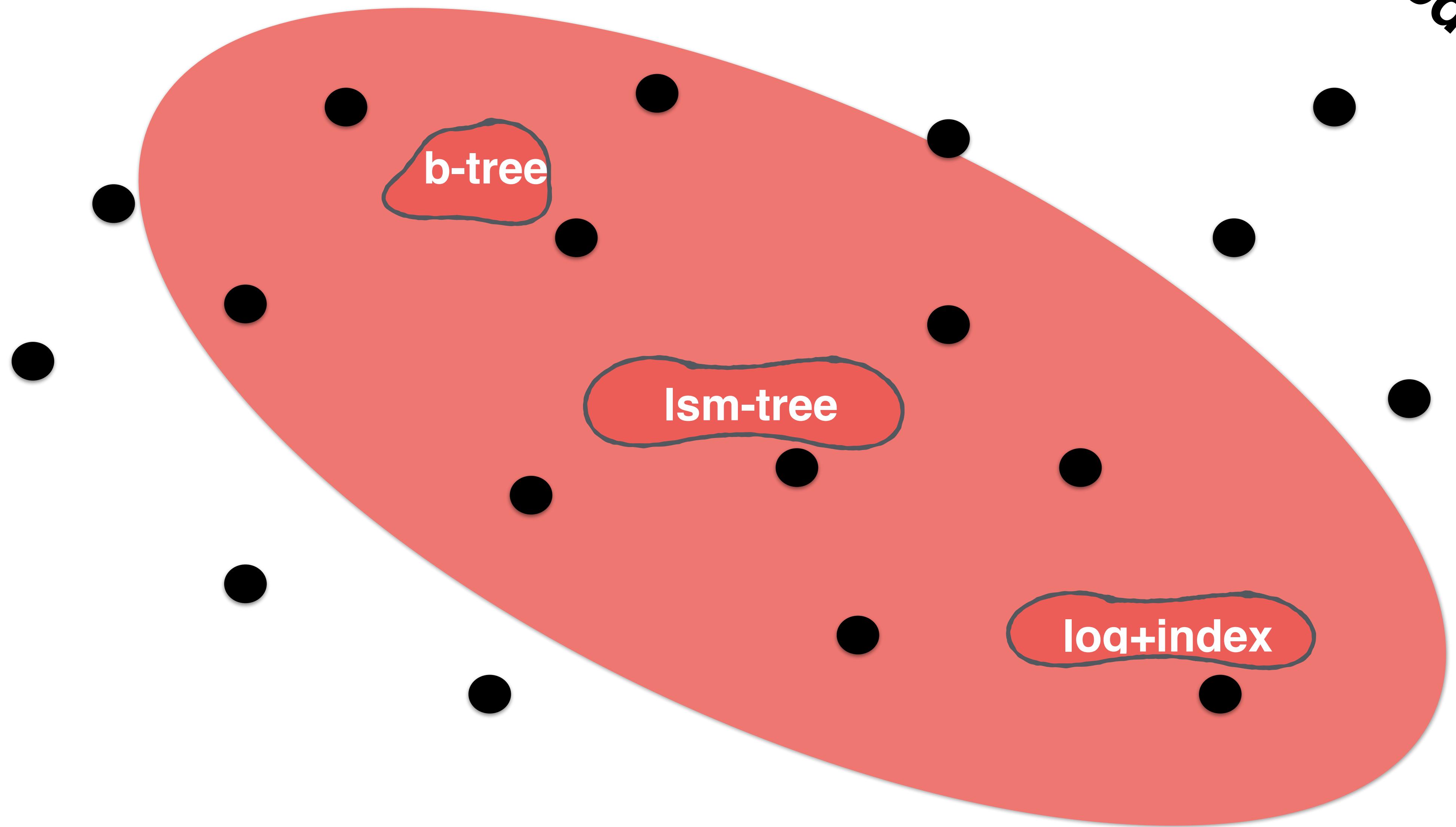
manually
selecting
“good” designs

b-tree

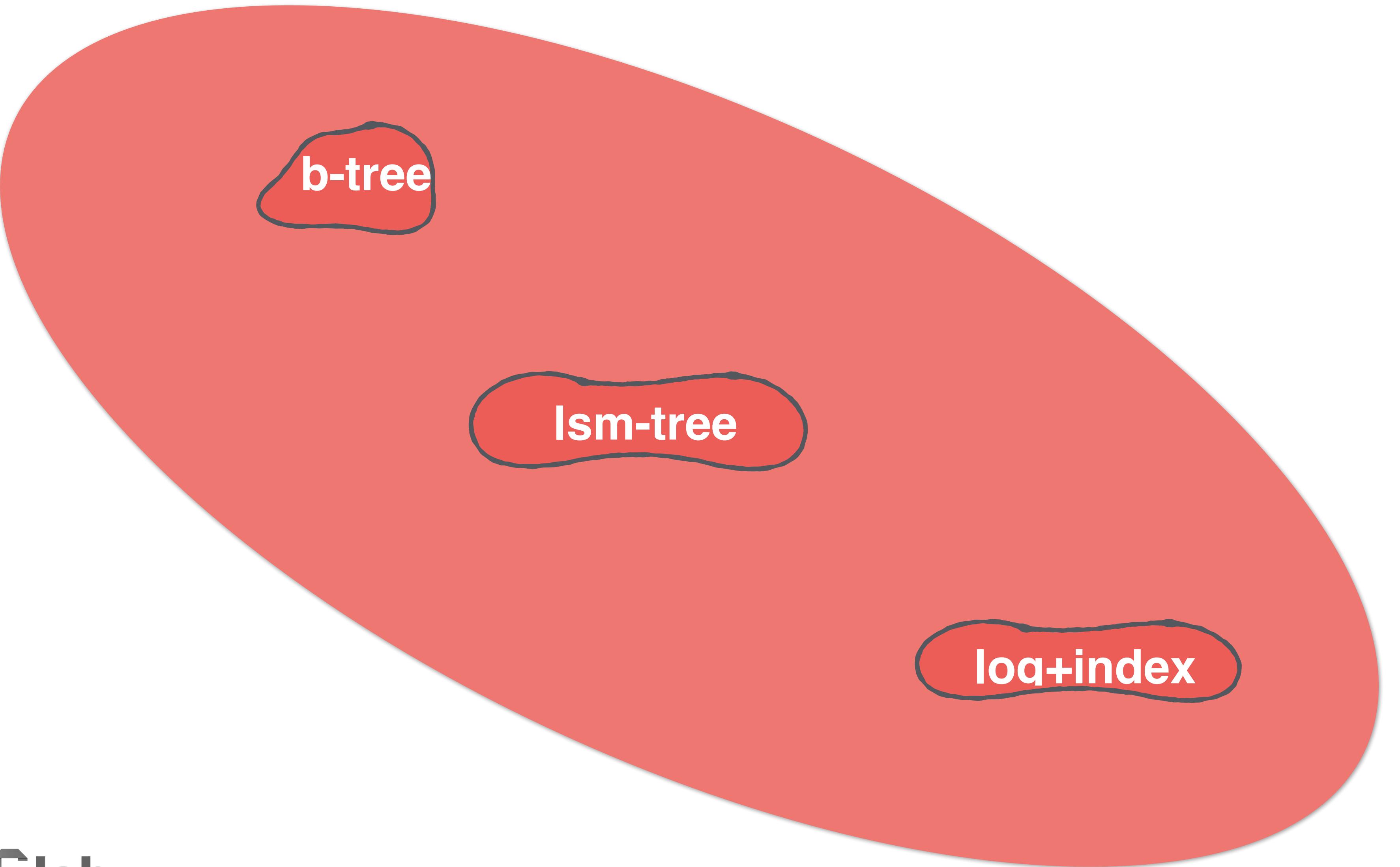
lsm-tree

log+index

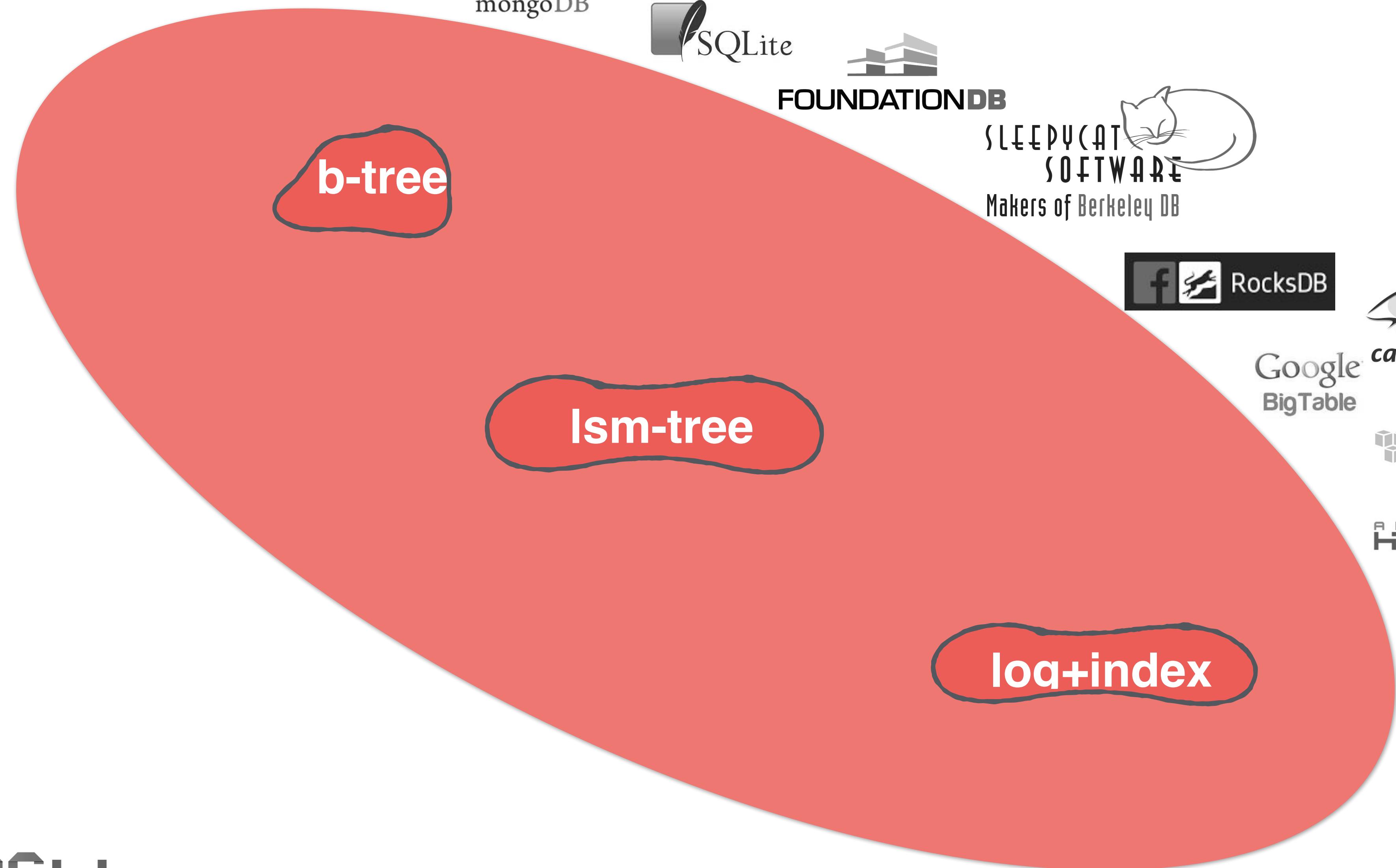
manually
selecting
“good” designs



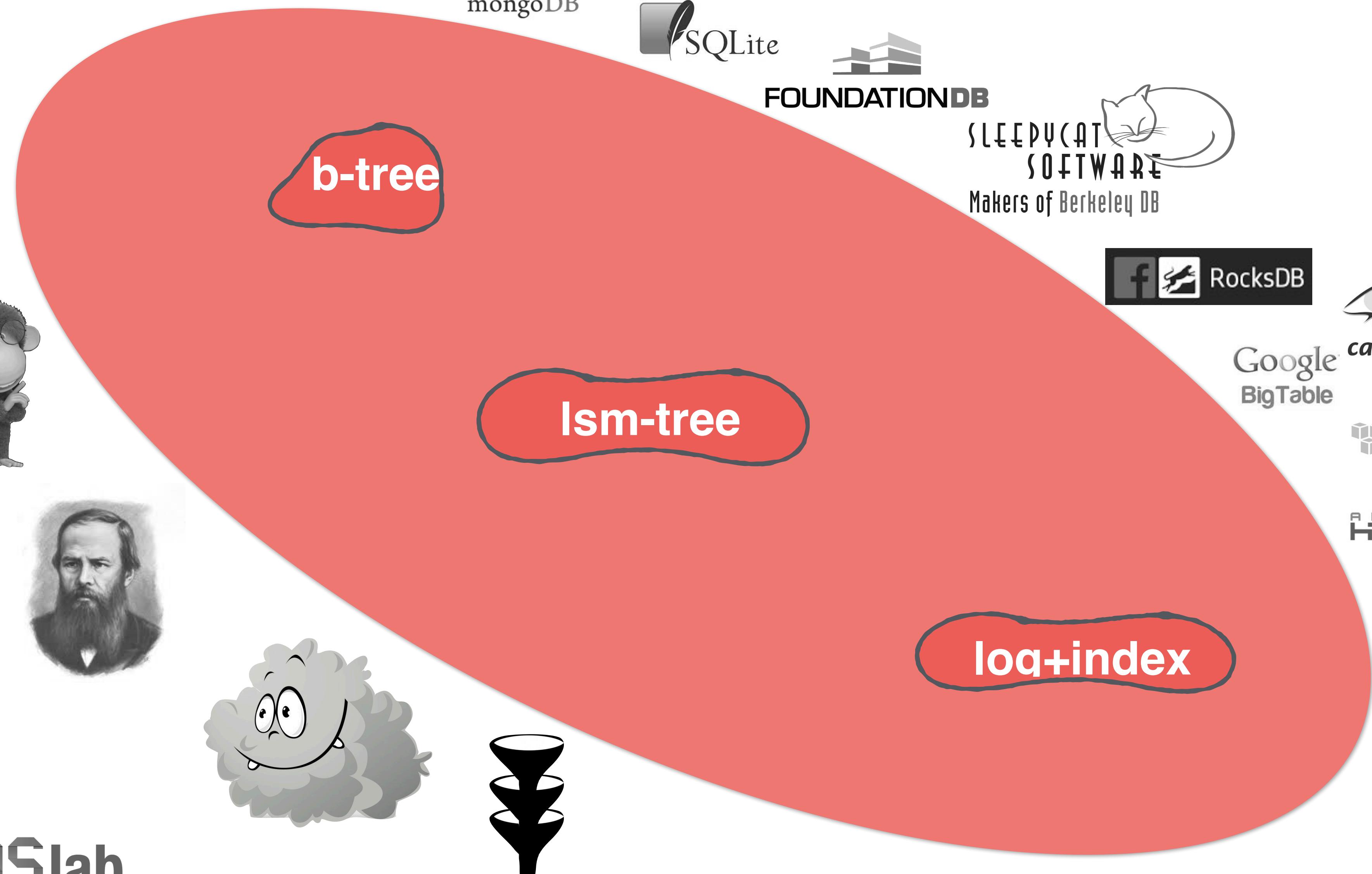
manually
selecting
“good” designs



manually
selecting
“good”
designs



manually
selecting
“good”
designs



WIREDTIGER



SQLite



FOUNDATIONDB

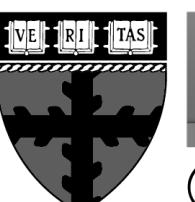


Google
BigTable



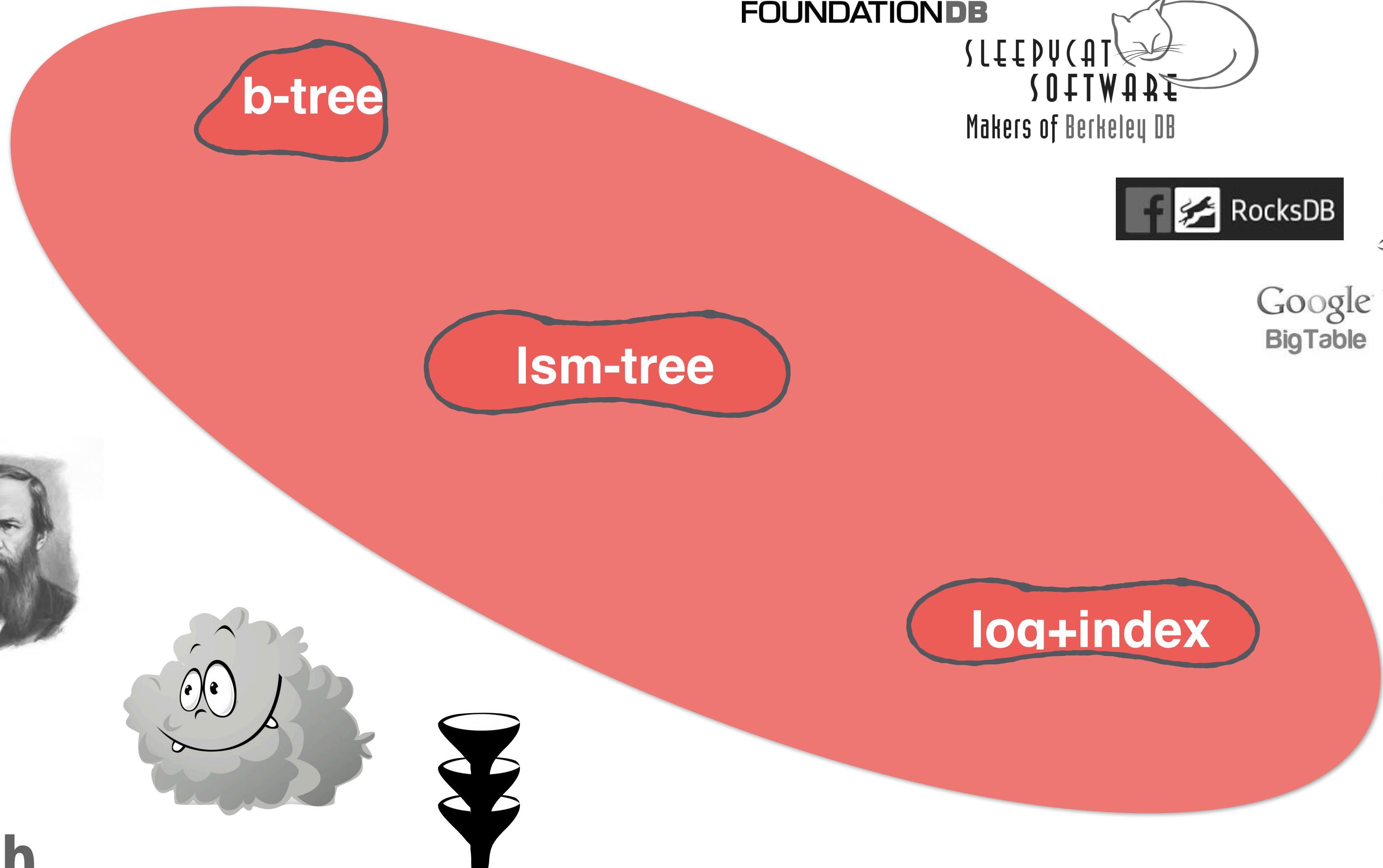
APACHE
HBASE

:riak



DASlab
@ Harvard SEAS

manually
selecting
“good”
designs



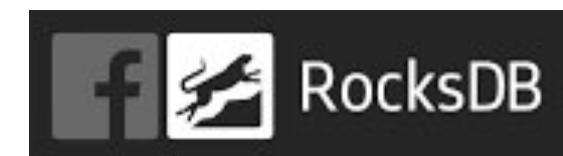
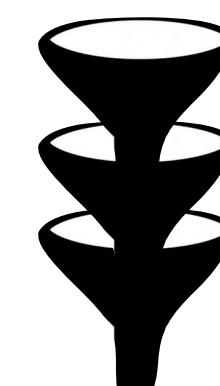
WIREDTIGER



FOUNDATIONDB



SLEEPYCAT
SOFTWARE
Makers of Berkeley DB



Google
BigTable

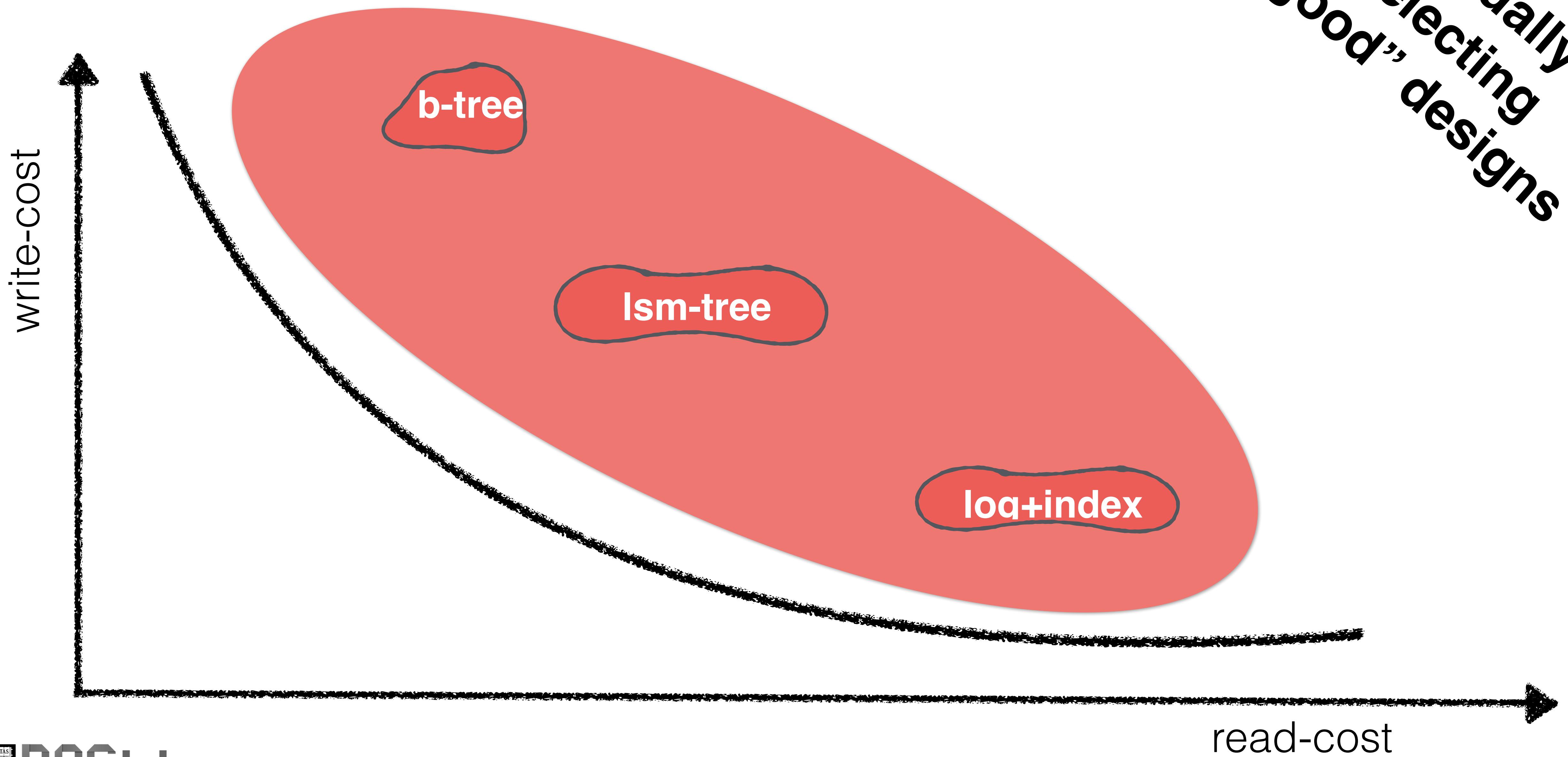


APACHE
HBASE

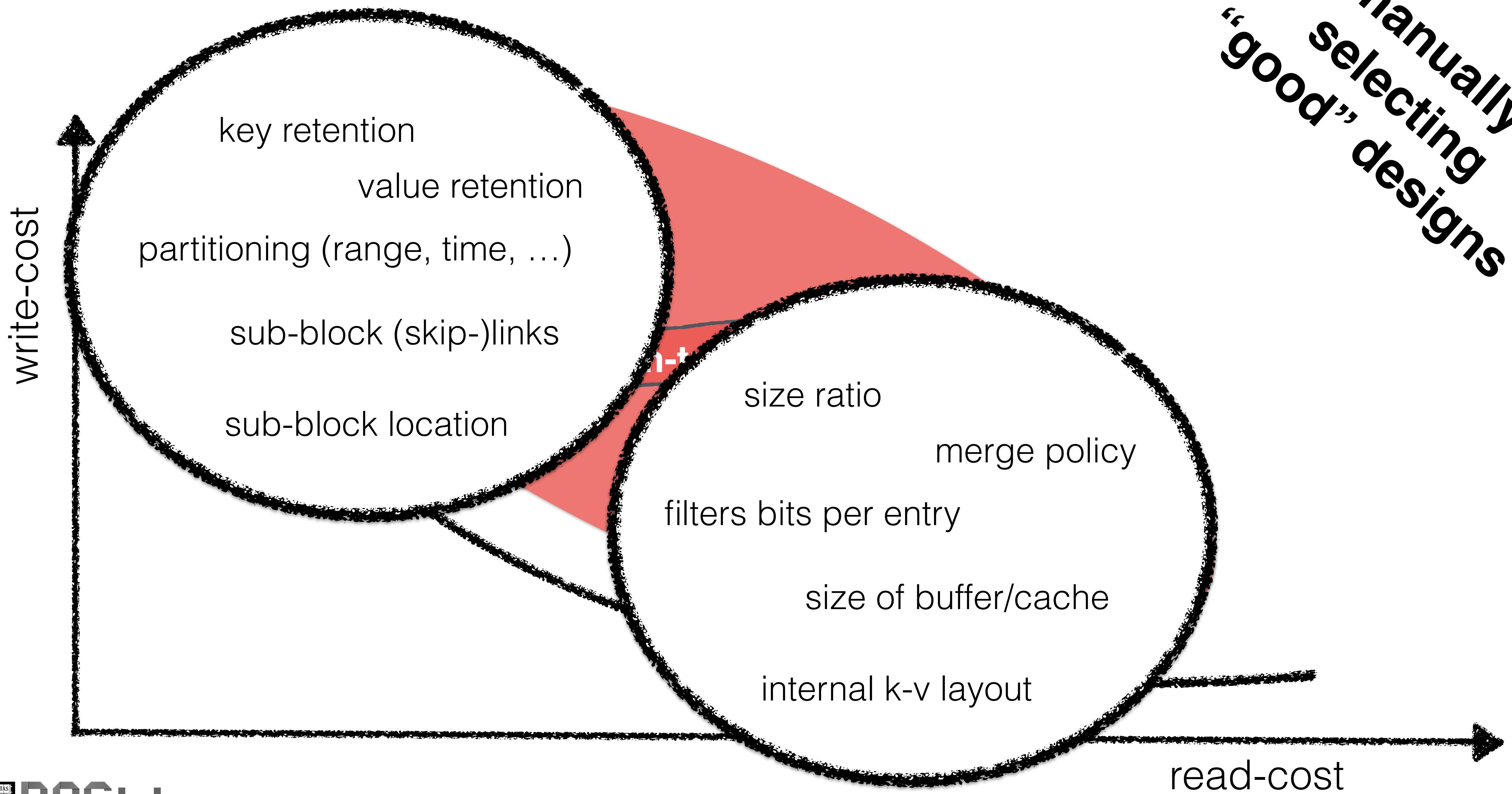
:riak

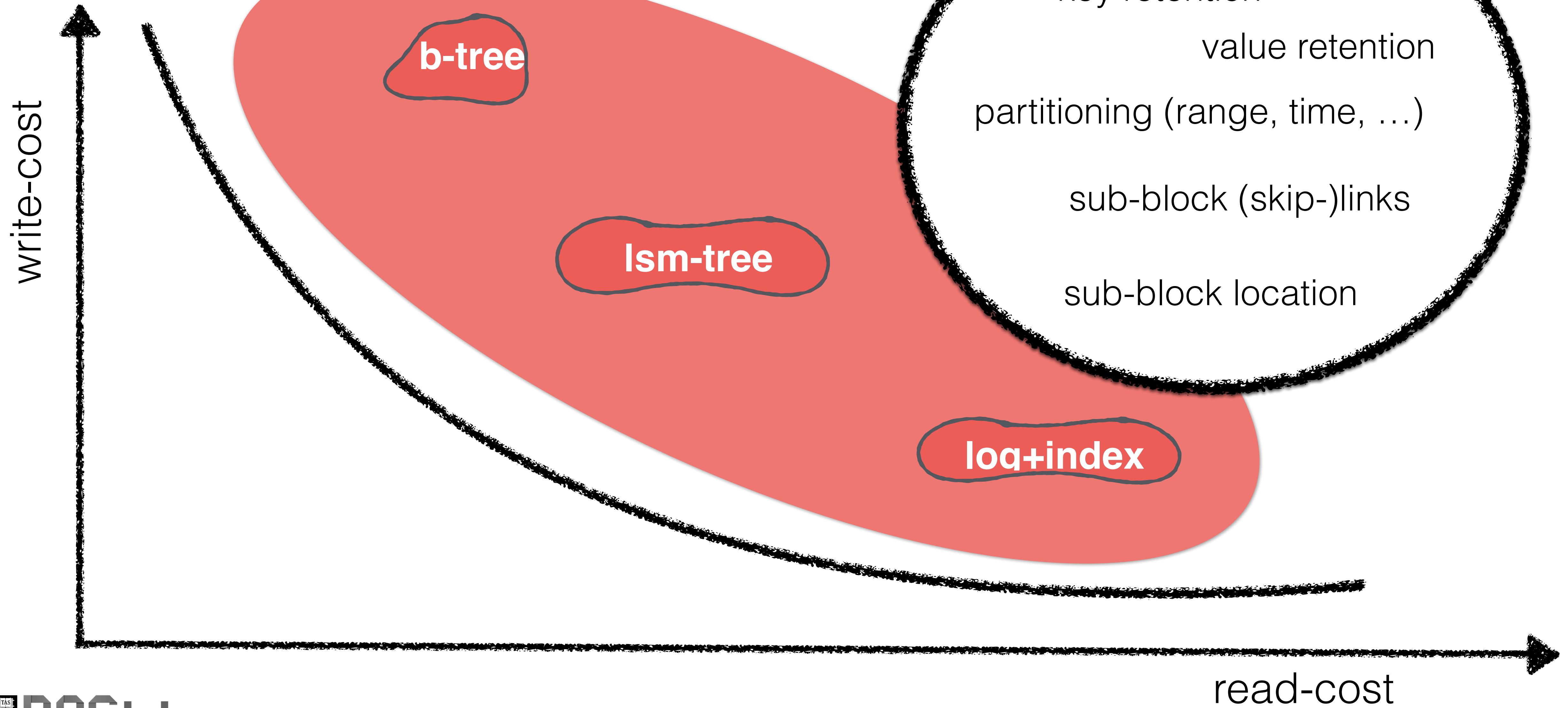
FASTER

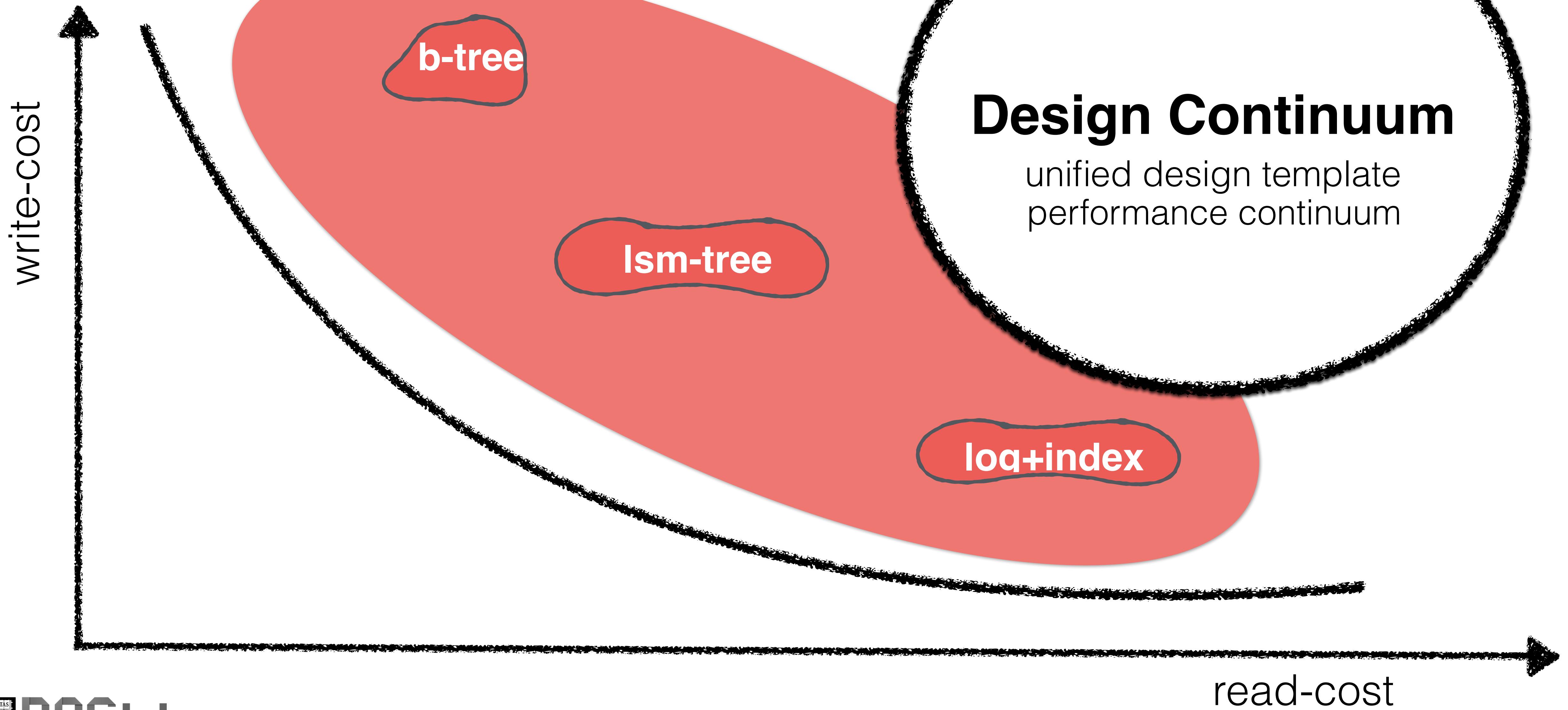
manually
selecting
“good” designs



*manually
selecting
“good”
designs*

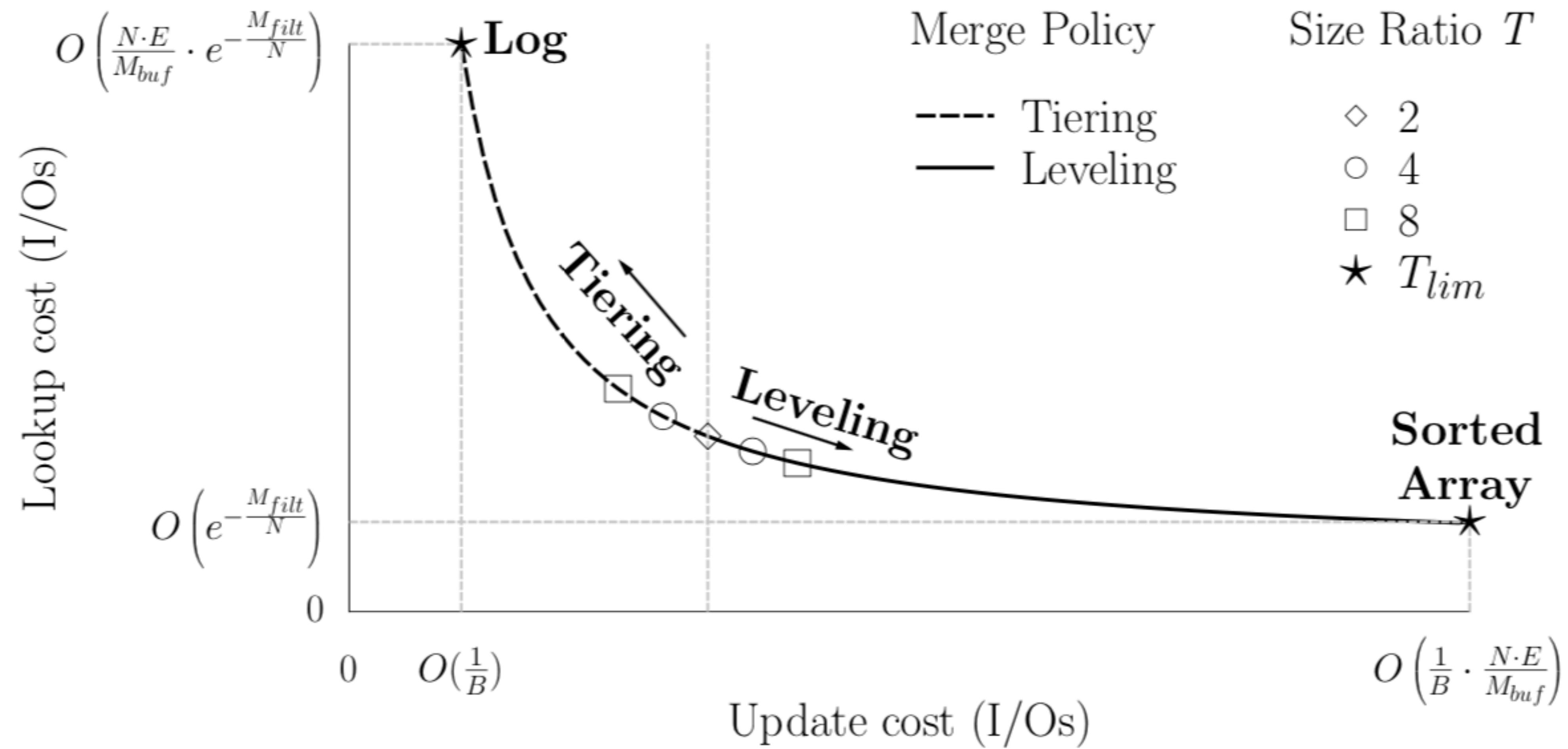






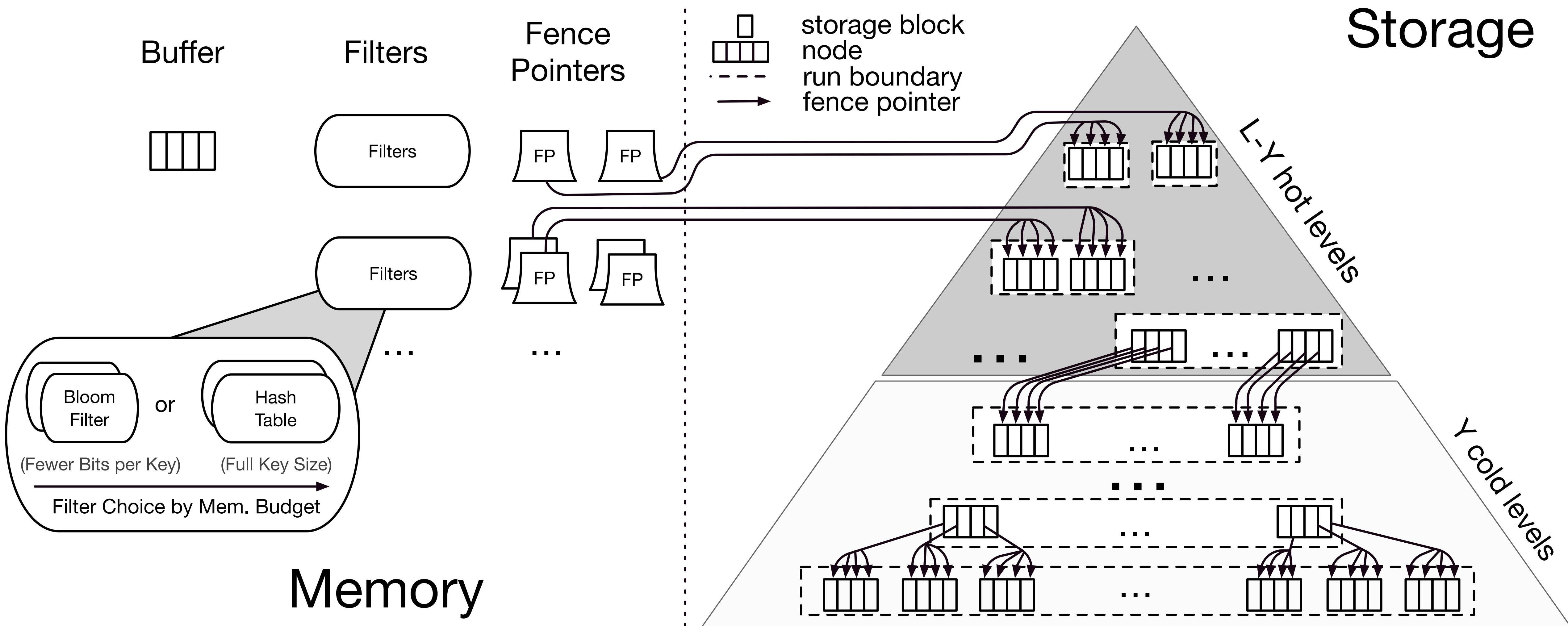
Design Continuum

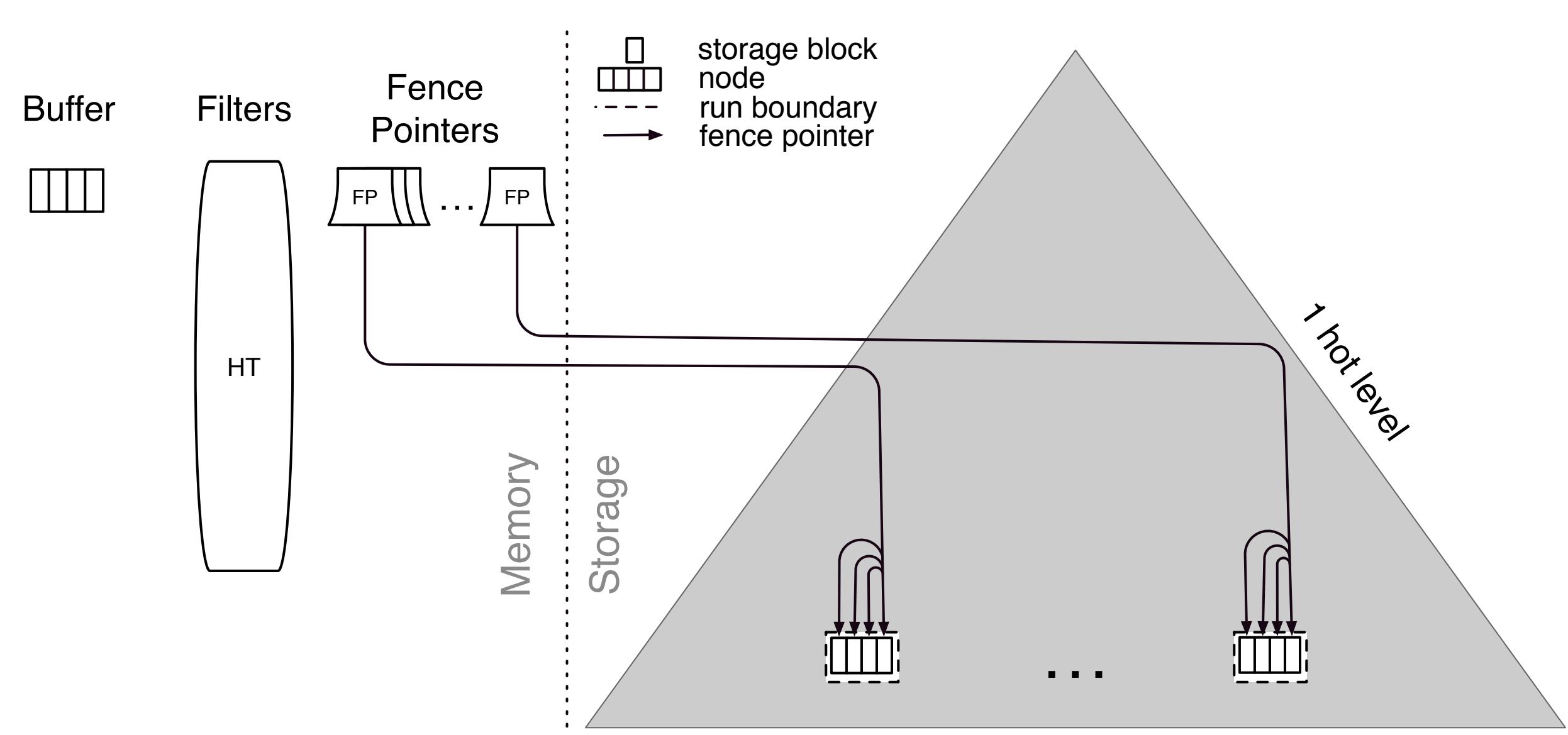
unified design template
performance continuum



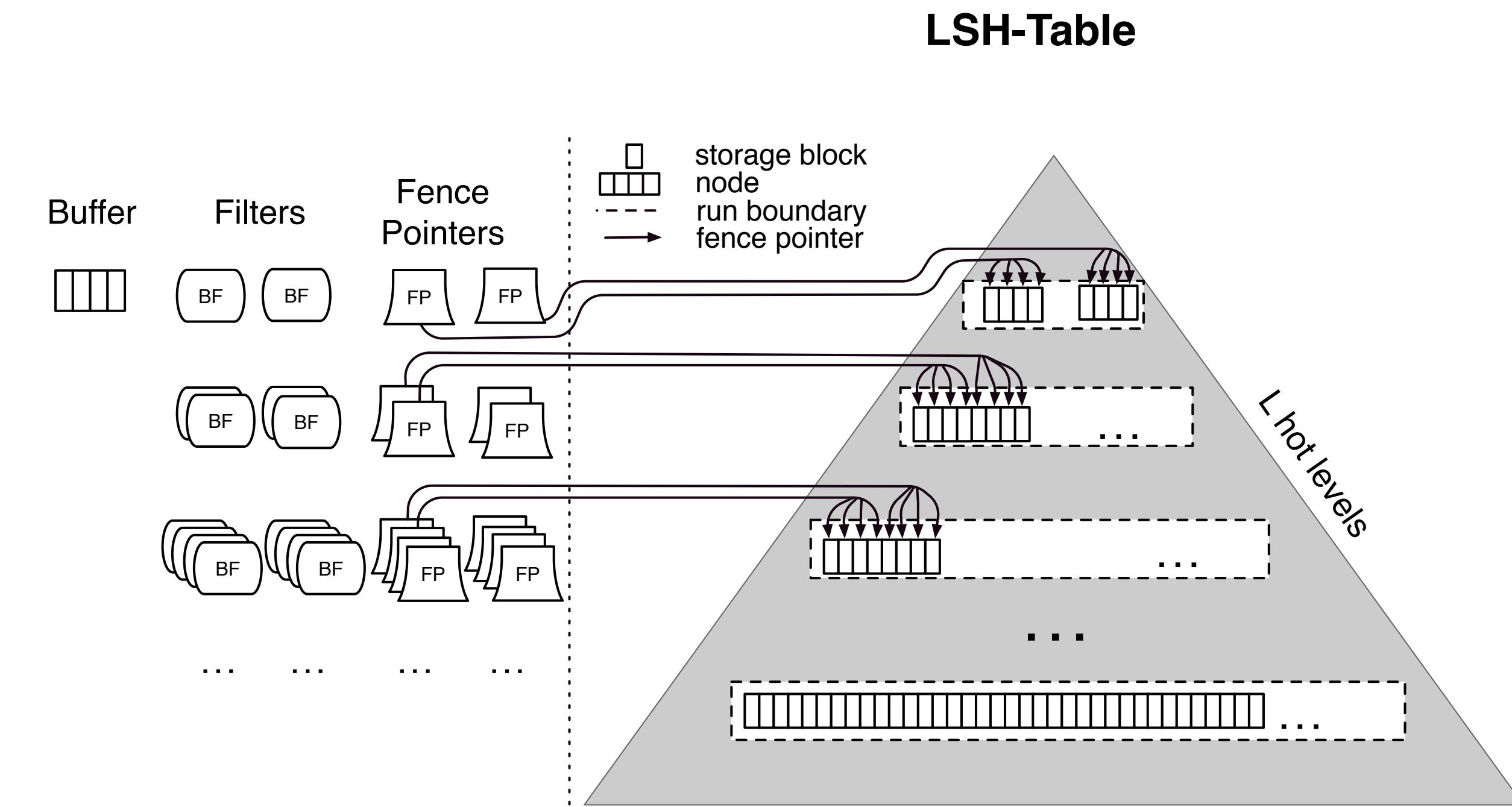
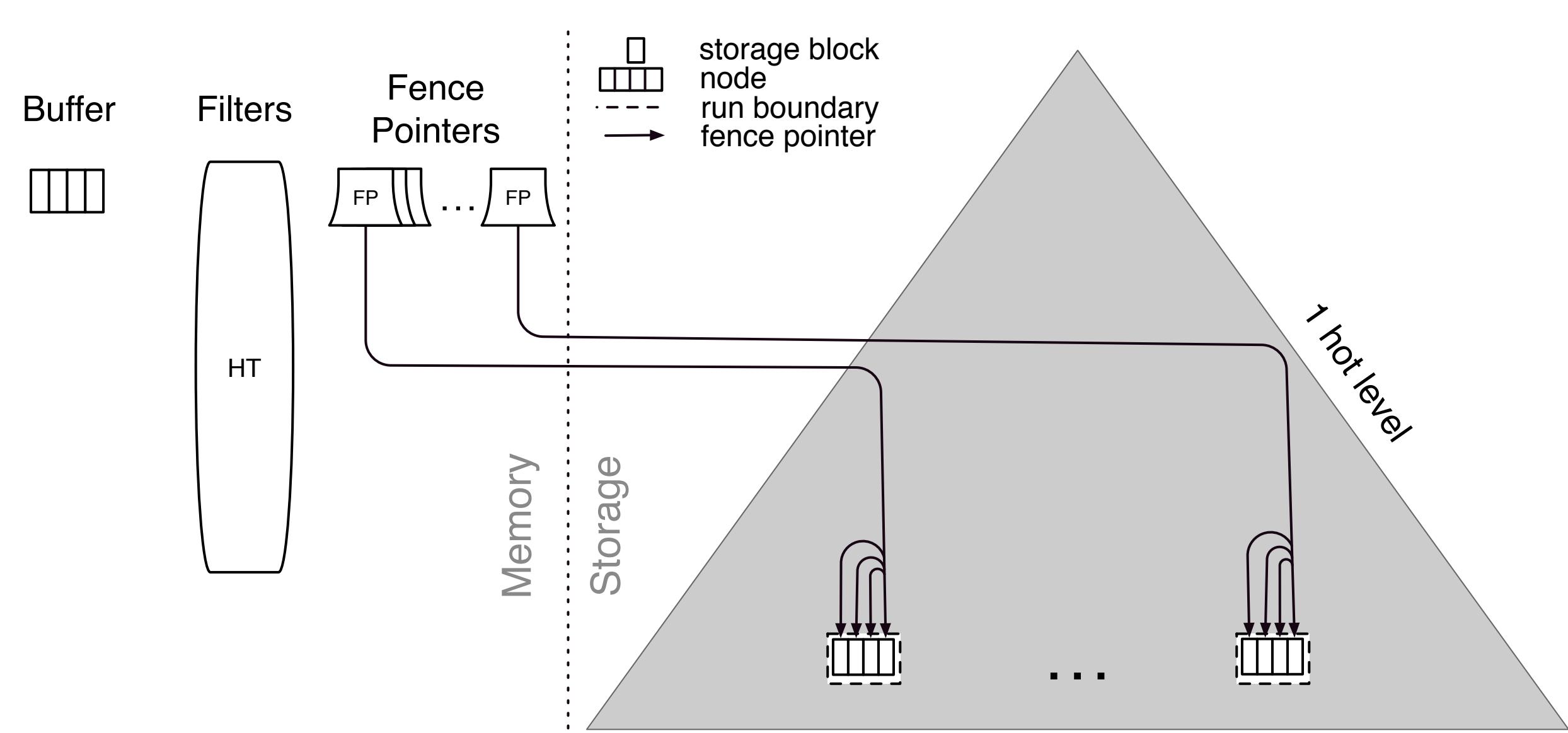
a unified design space

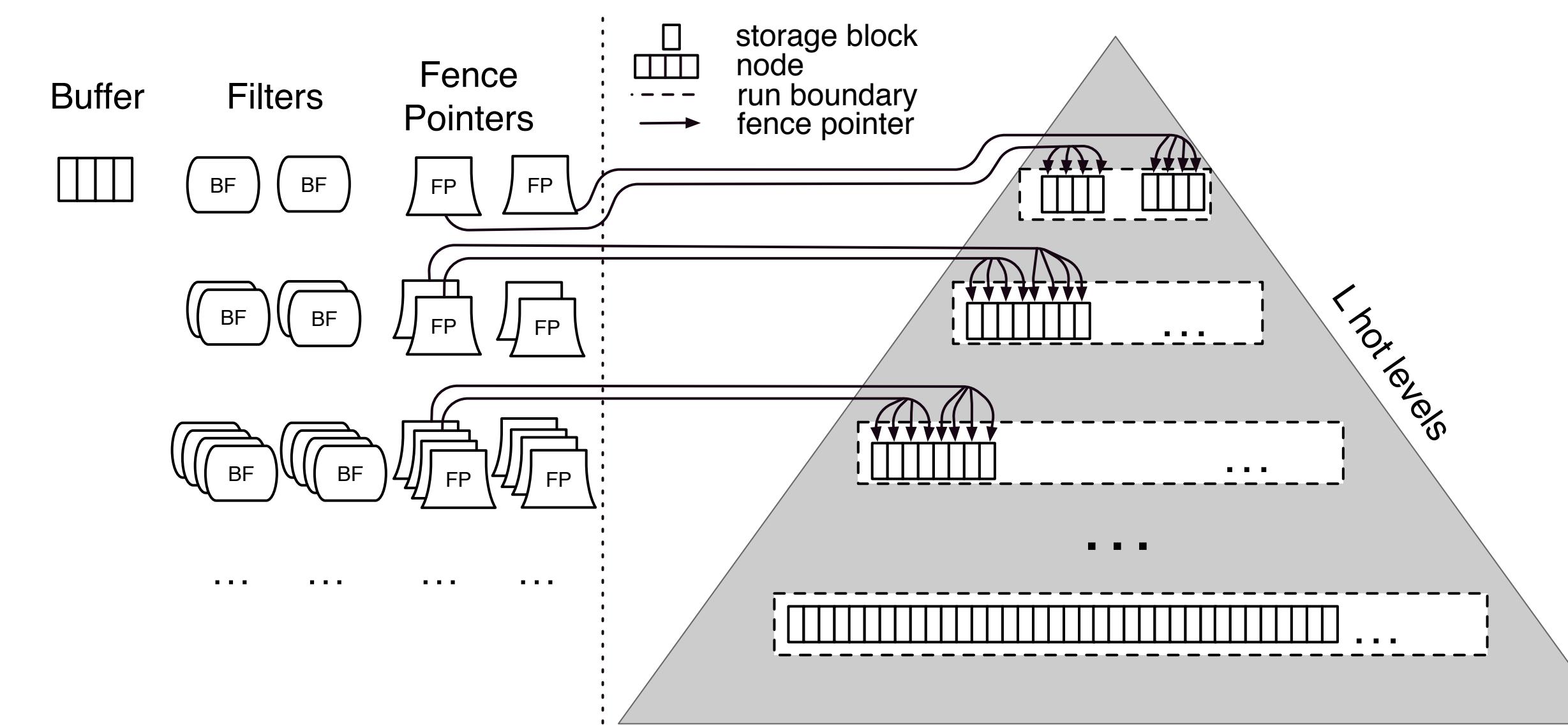
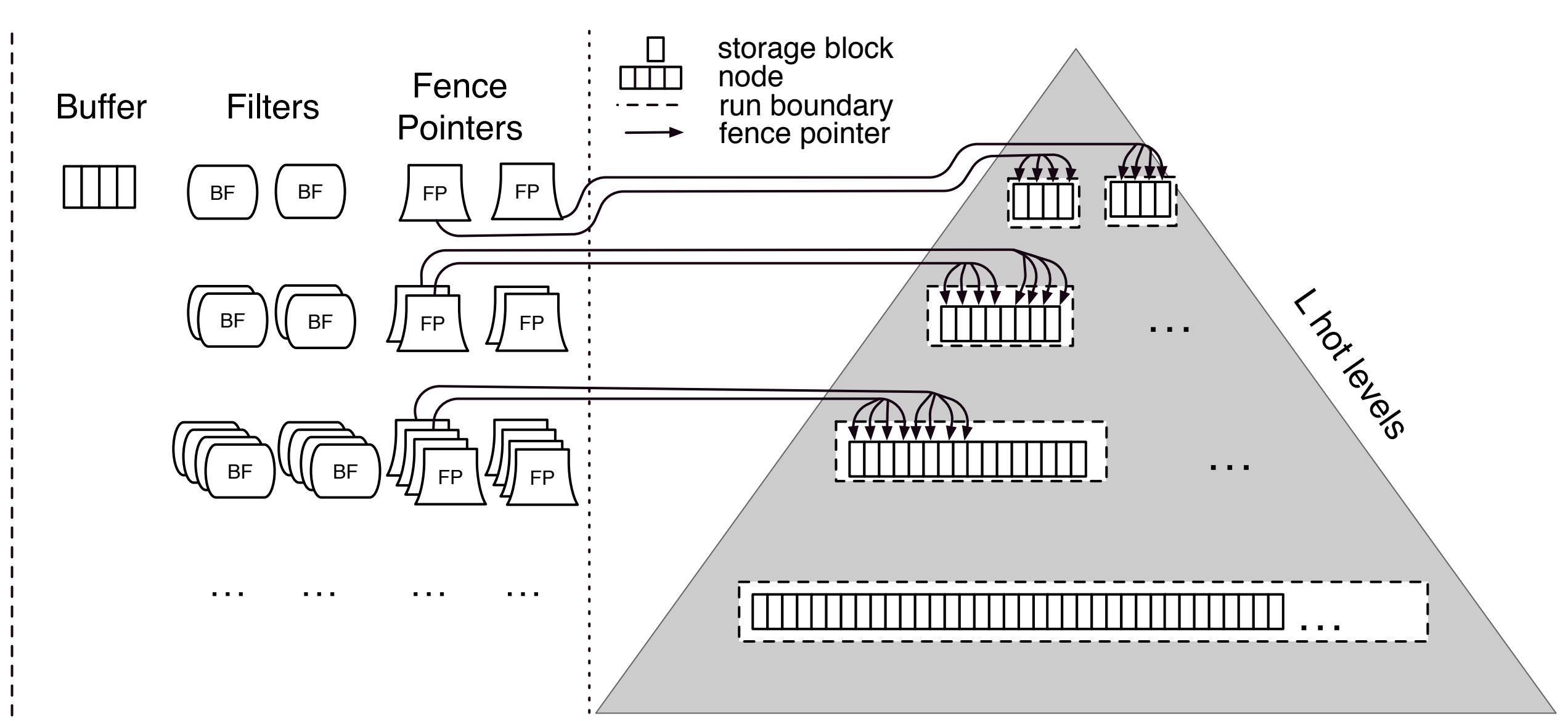
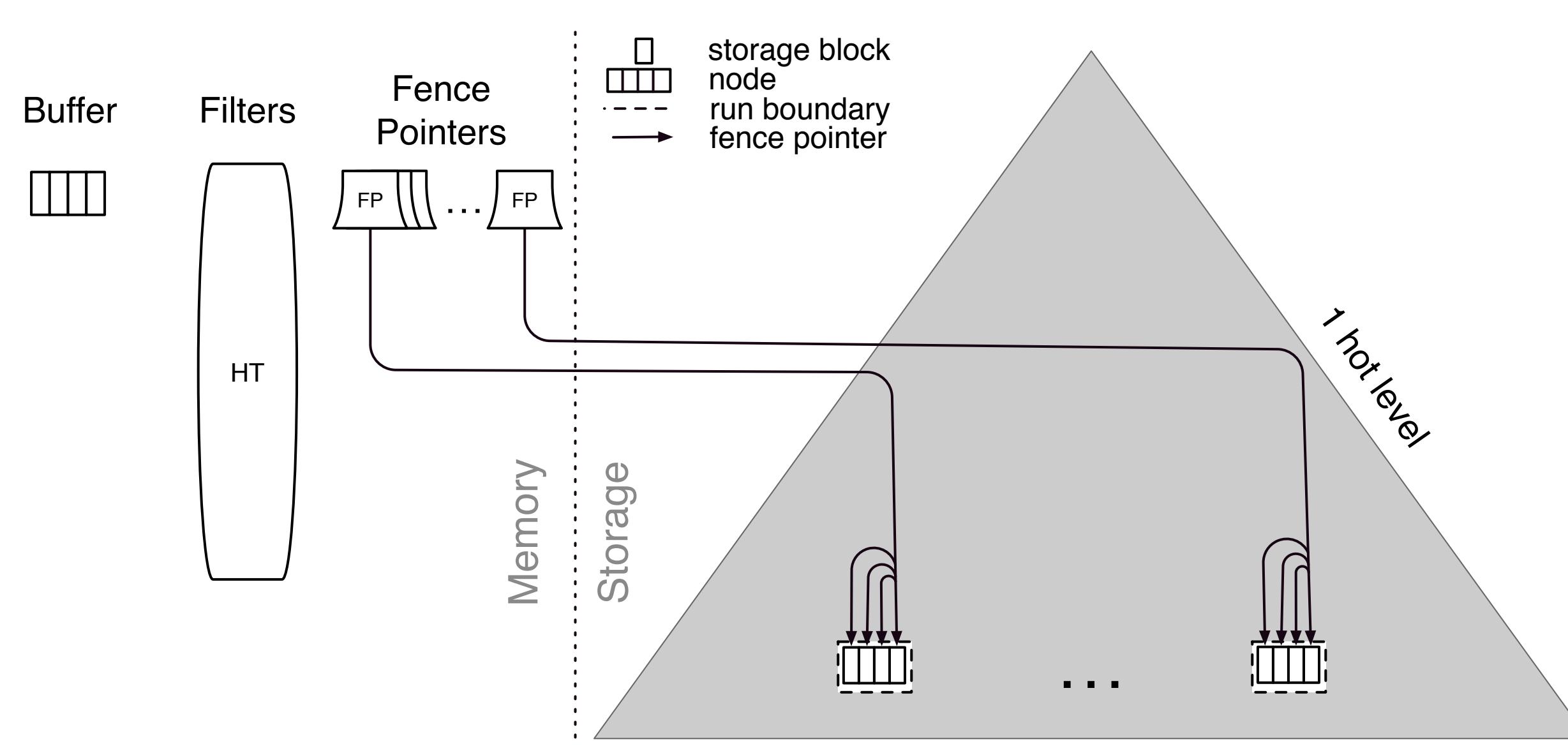
[log, log+hash, LSM-tree*, B^ETree, B-Tree, Sorted Array]

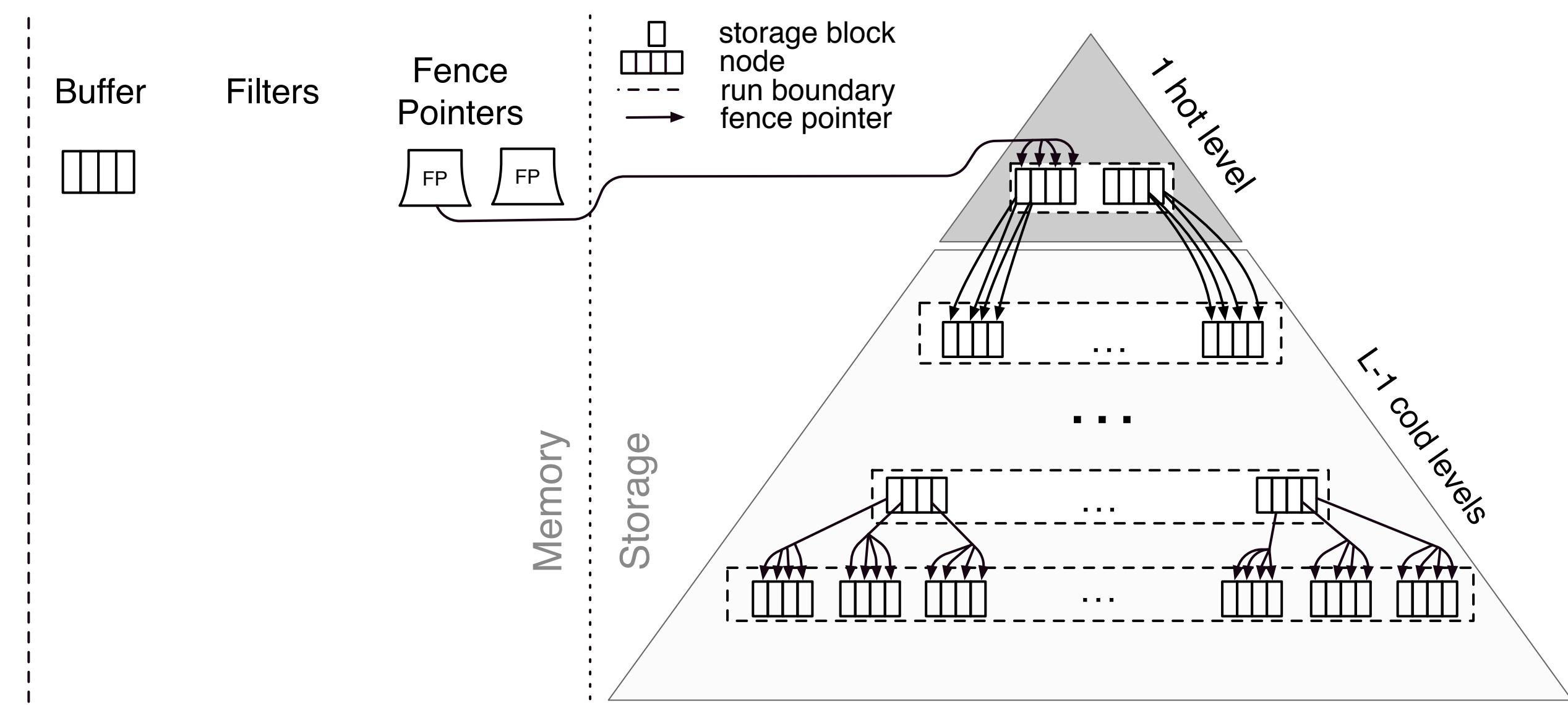
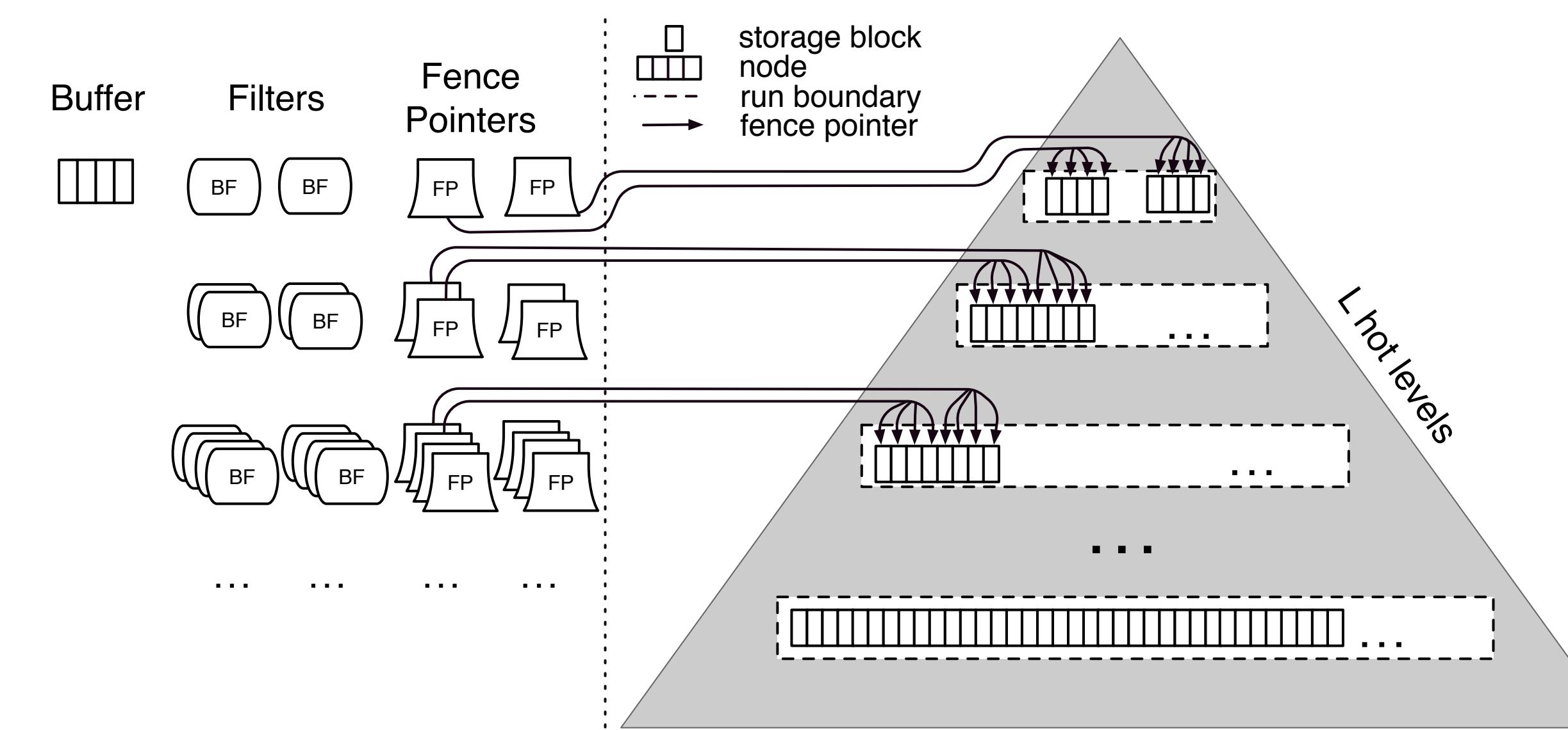
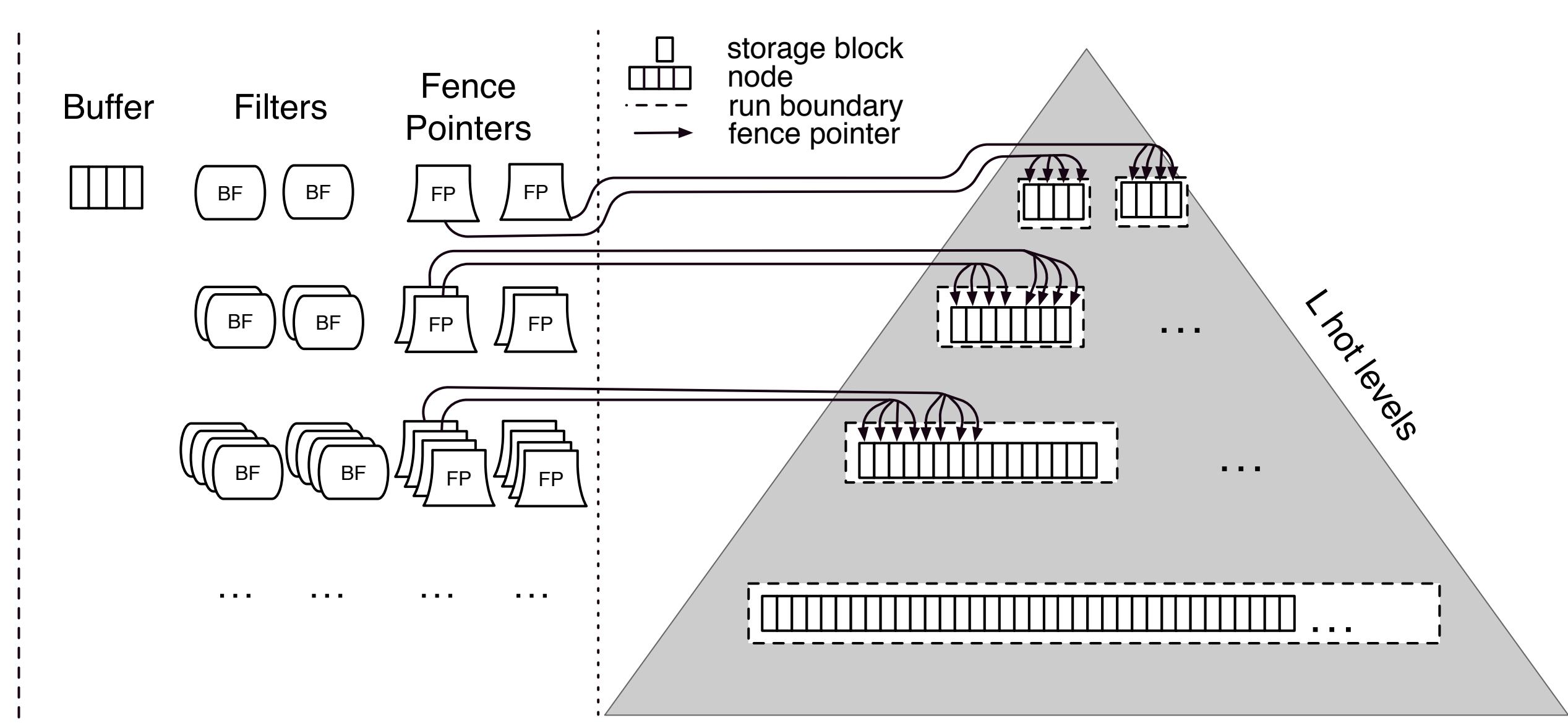
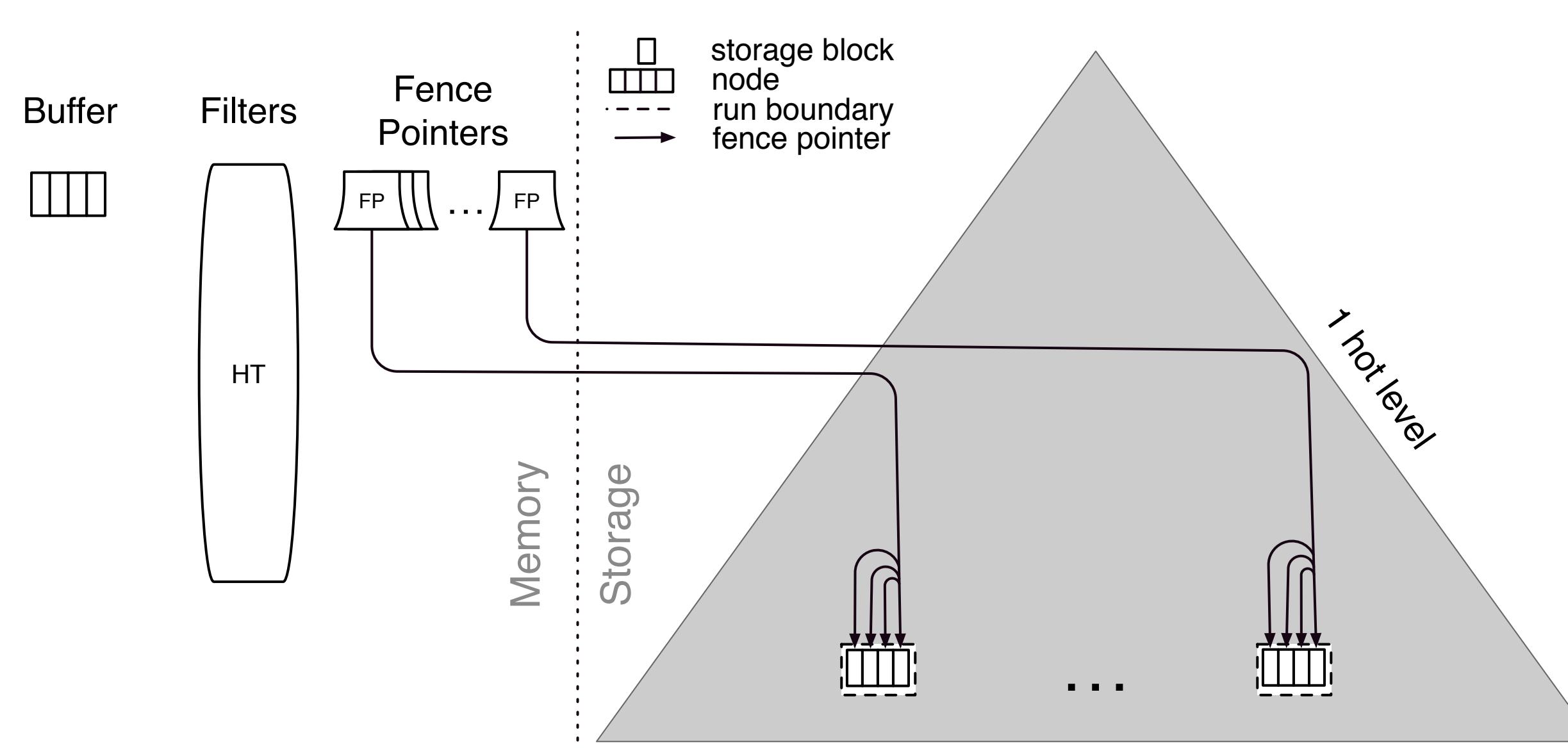




LSH-Table







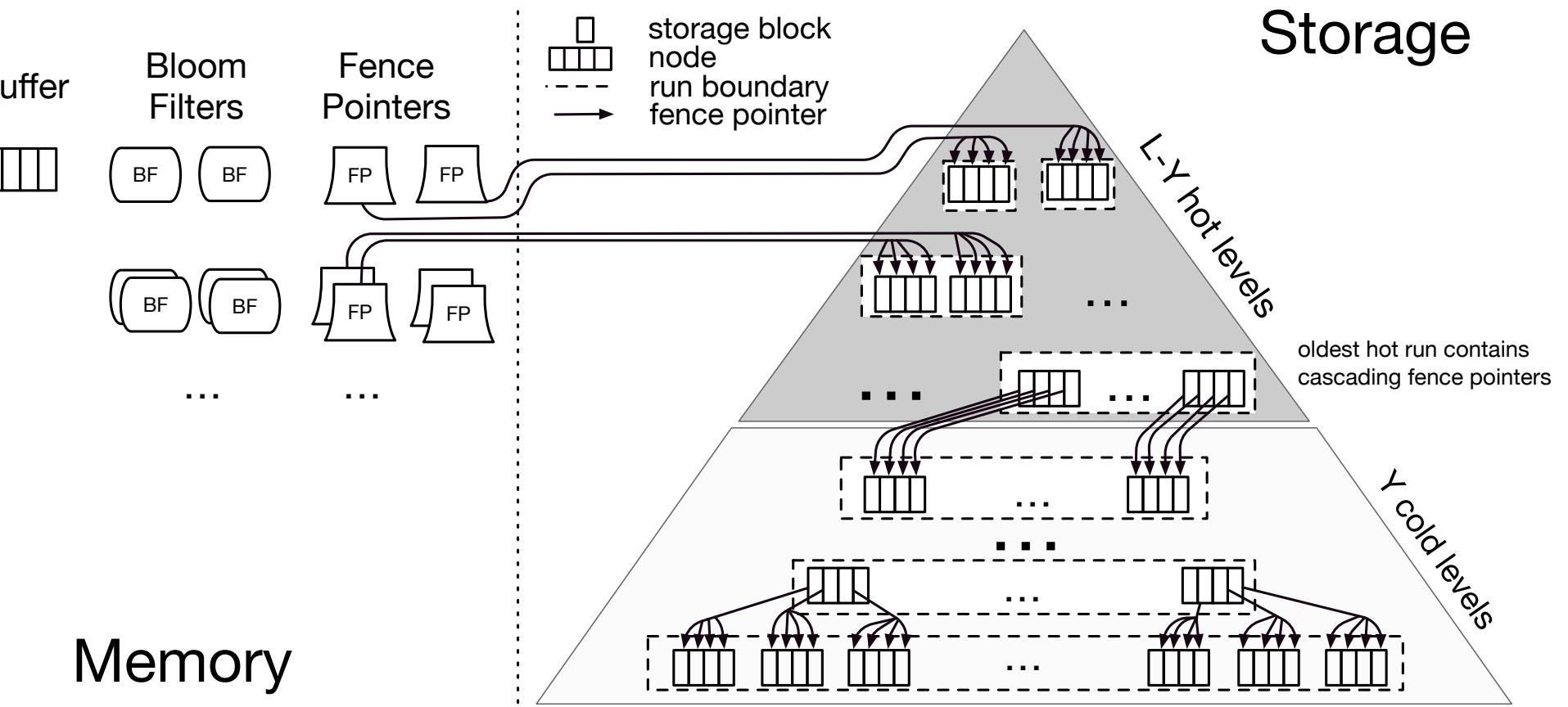
Components

Term	Name	Description	Min. Value	Max. Value	Units
B	Block Size	# data entries that fit in a storage block.			Entries
M	Memory	Total main memory budget.	$\frac{B \cdot E + F \cdot T \cdot M_B}{E \cdot B}$	$N \cdot E$	Bits
N	Dataset Size	# total data entries in the dataset.			Entries
E	Entry Size	Size of an entry.			Bits
F	Key Size	Size of a key, also used to approximate size of a fence (fence key and pointer).			Bits
s	Avg. Selectivity	Average selectivity of a long range query.			Entries
T	Growth Factor	Capacity ratio between adjacent levels.	2	B	Ratio
K	Hot Merge Threshold	Maximum # runs per hot level.	1	$T - 1$	Runs
Z	Cold Merge Threshold	Maximum # runs per cold level.	1	$T - 1$	Runs
D	Max. Node Size	Maximum size of a node; defines a contiguous data region.	1	$\frac{N}{B}$	Blocks
M_F	Fence & Filter Memory Budget	# bits of main memory budgeted to fence pointers and filters.	$\frac{F \cdot T \cdot M_B}{E \cdot B}$	M	Bits

Environment Parameters

Design Parameters

Super-Structure



Derived Term	Expression	Units
L (# total levels)	$\lceil \log_T \frac{N \cdot E}{M_B} \rceil$	Levels
X (Filters Memory Threshold)	$\frac{1}{\ln 2^2} \cdot (\frac{\ln T}{T-1} + \frac{\ln K - \ln Z}{T})$	Bits per Entry
M_{F_HI} (M_F Threshold: Hot Levels Saturation)	$N \cdot (\frac{X}{T} + \frac{F}{B})$	Bits
M_{F_LO} (M_F Threshold: Cold Levels Saturation)	$\frac{M_B \cdot F \cdot T}{E \cdot B}$	Bits
Y (# Cold Levels)	$\begin{cases} 0 & \text{if } M_F \geq M_{FHI} \\ \lfloor \log_T \frac{N}{M_F} \cdot (\frac{X}{T} + \frac{F}{B}) \rfloor & \text{if } M_{FLO} < M_F < M_{FHI} \\ L - 1 & \text{if } M_F = M_{FLO} \end{cases}$	Levels
M_{FP} (Fence Pointer Memory Budget)	$T^{L-Y+1} \cdot F \cdot \frac{M_b}{E \cdot B} \cdot \frac{T}{T-1}$	Bits
M_{BF} (Filter Memory Budget)	$M_F - M_{FP}$	Bits
M_B (Buffer Memory Budget)	$B \cdot E + (M - M_F)$	Bits
p_{sum} (Sum of BF False Positive Rates)	$e^{-\frac{M_{BE}}{N} \cdot \ln(2)^2 \cdot T^Y \cdot Z^{\frac{T-1}{T}} \cdot K^{\frac{1}{T}} \cdot \frac{T}{T-1}}$	
p_i (BF False Positive Rate at Level i)	$\begin{cases} 1 & \text{if } i > L - Y \\ \frac{p_{sum}}{Z} \cdot \frac{T-1}{T} & \text{if } i = L - Y \\ \frac{p_{sum}}{K} \cdot \frac{T-1}{T} \cdot \frac{1}{T^{L-Y-i}} & \text{if } i < L - Y \end{cases}$	Probability

Derived Design Rules

Operation	Cost Expression (I/O)
Update	$O(\frac{1}{B} \cdot (\frac{T}{K} \cdot (L - Y - 1) + \frac{T}{Z} \cdot (Y + 1)))$
Zero Result Lookup	$O(Z \cdot e^{-\frac{M_{BF}}{N} \cdot T^Y} + Y \cdot Z)$
Single Result Lookup	$O(1 + Z \cdot e^{-\frac{M_{BF}}{N} \cdot T^Y} + Y \cdot Z)$
Short Scan	$O(K \cdot (L - Y - 1) + Z \cdot (Y + 1))$
Long Scan	$O(\frac{s \cdot Z}{B})$

Components

Term	Name	Description	Min. Value	Max. Value	Units
B	Block Size	# data entries that fit in a storage block.			Entries
M	Memory	Total main memory budget.	$\frac{B \cdot E + F \cdot T \cdot M_B}{E \cdot B}$	$N \cdot E$	Bits
N	Dataset Size	# total data entries in the dataset			Entries
E	Entry Size	Size of an entry.			Bits
F	Key Size	Size of a key, also used to approximate size of a fence (fence key and pointer).			Bits
s	Avg. Selectivity	Average selectivity of a long range query.			Entries
T	Growth Factor	Capacity ratio between adjacent levels.	2	B	Ratio
K	Hot Merge Threshold	Maximum # runs per hot level.	1	$T - 1$	Runs
Z	Cold Merge Threshold	Maximum # runs per cold level.	1	$T - 1$	Runs
D	Max. Node Size	Maximum size of a node; defines a contiguous data region.	1	$\frac{N}{B}$	Blocks
M_F	Fence & Filter Memory Budget	# bits of main memory budgeted to fence pointers and filters.	$\frac{F \cdot T \cdot M_B}{E \cdot B}$	M	Bits

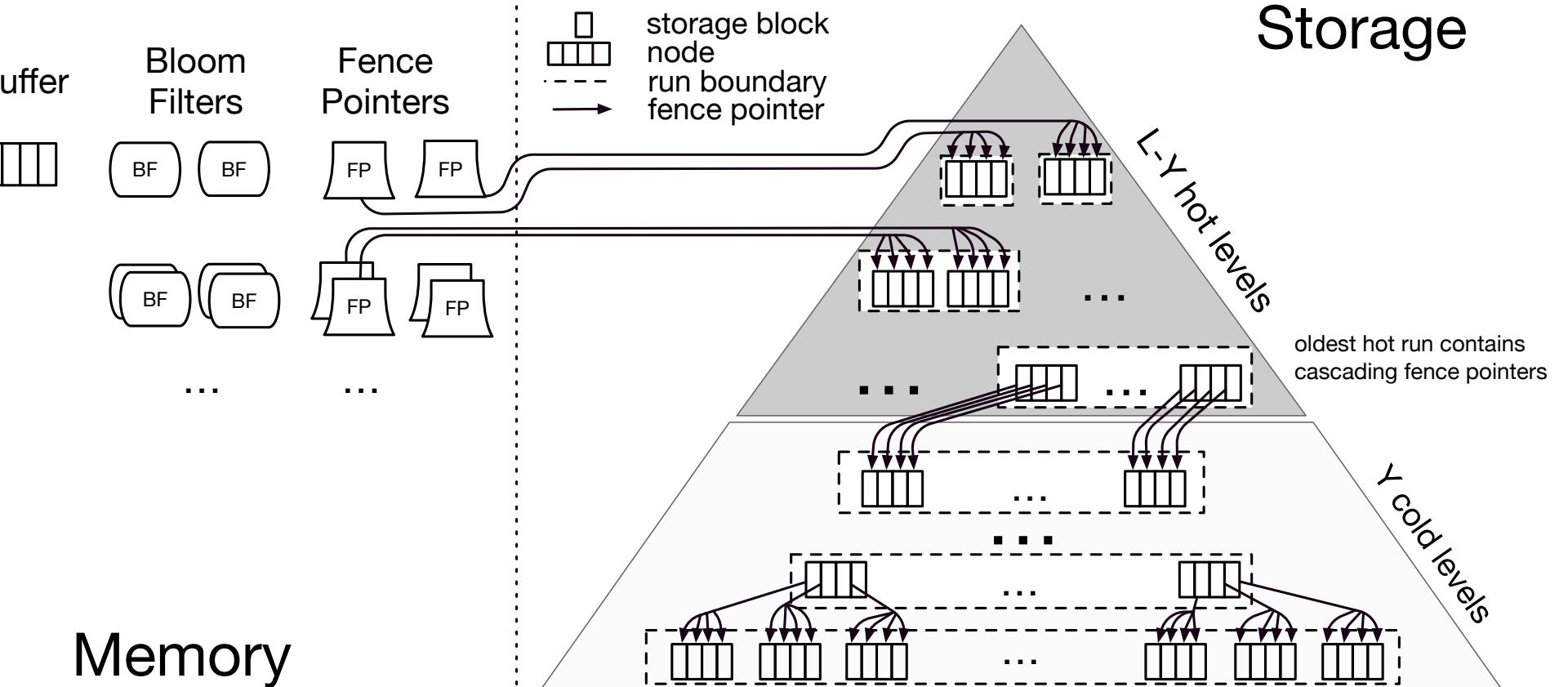
design

Environment Parameters

Design Parameters

Super-Structure

Memory



Derived Term	Expression	Units
L (# total levels)	$\lceil \log_T \frac{N \cdot E}{M_B} \rceil$	Levels
X (Filters Memory Threshold)	$\frac{1}{\ln 2^2} \cdot (\frac{\ln T}{T-1} + \frac{\ln K - \ln Z}{T})$	Bits per Entry
M_{FI} (M_F Threshold: Hot Levels Saturation)	$N \cdot (\frac{X}{T} + \frac{F}{B})$	Bits
M_{FO} (M_F Threshold: Cold Levels Saturation)	$\frac{M_B \cdot E \cdot T}{B}$	Bits
Y (# Cold Levels)	$\begin{cases} \lfloor \log_T \frac{N}{M_F} \cdot (\frac{X}{T} + \frac{F}{B}) \rfloor & \text{if } M_{FO} < M_F < M_{FI} \\ L - 1 & \text{if } M_F = M_{FO} \end{cases}$	Levels
M_{FP} (Fence Pointer Memory Budget)	$T^{L-Y+1} \cdot F \cdot \frac{M_b}{E \cdot B} \cdot \frac{T}{T-1}$	Bits
M_{BF} (Filter Memory Budget)	$M_F - M_{FP}$	Bits
M_B (Buffer Memory Budget)	$B \cdot E + (M - M_F)$	Bits
P_{sum} (Sum of BF False Positive Rates)	$e^{-\frac{M_{BE}}{N} \cdot \ln(2)^2 \cdot T^Y \cdot Z^{\frac{T-1}{T}} \cdot K^{\frac{1}{T}} \cdot T^{\frac{T}{T-1}}}$	
p_i (BF False Positive Rate at Level i)	$\begin{cases} \frac{1}{p_{sum}} & \text{if } i > L - Y \\ \frac{Z}{p_{sum}} \cdot \frac{T-1}{T} & \text{if } i = L - Y \\ \frac{1}{T^{L-Y-i}} & \text{if } i < L - Y \end{cases}$	Probability

Derived Design Rules

Operation	Cost Expression (T/C)
Update	$O(\frac{1}{B} \cdot (\frac{1}{K} \cdot (L - Y - 1) + \frac{T}{Z} \cdot (Y + 1)))$
Zero Result Lookup	$O(Z \cdot e^{-\frac{M_{BF}}{N} \cdot T^Y} + Y \cdot Z)$
Single Result Lookup	$O(1 + Z \cdot e^{-\frac{M_{BF}}{N} \cdot T^Y} + Y \cdot Z)$
Short Scan	$O(K \cdot L - Y) + Z \cdot (Y + 1))$
Long Scan	$O(Z \cdot \frac{1}{B})$

unified I/O models

Components

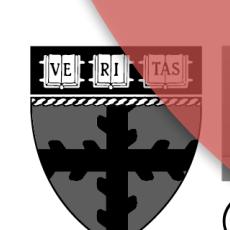
growth factor

hot merge threshold

cold merge threshold

fence/filter memory

max node size



DAS lab
@ Harvard SEAS

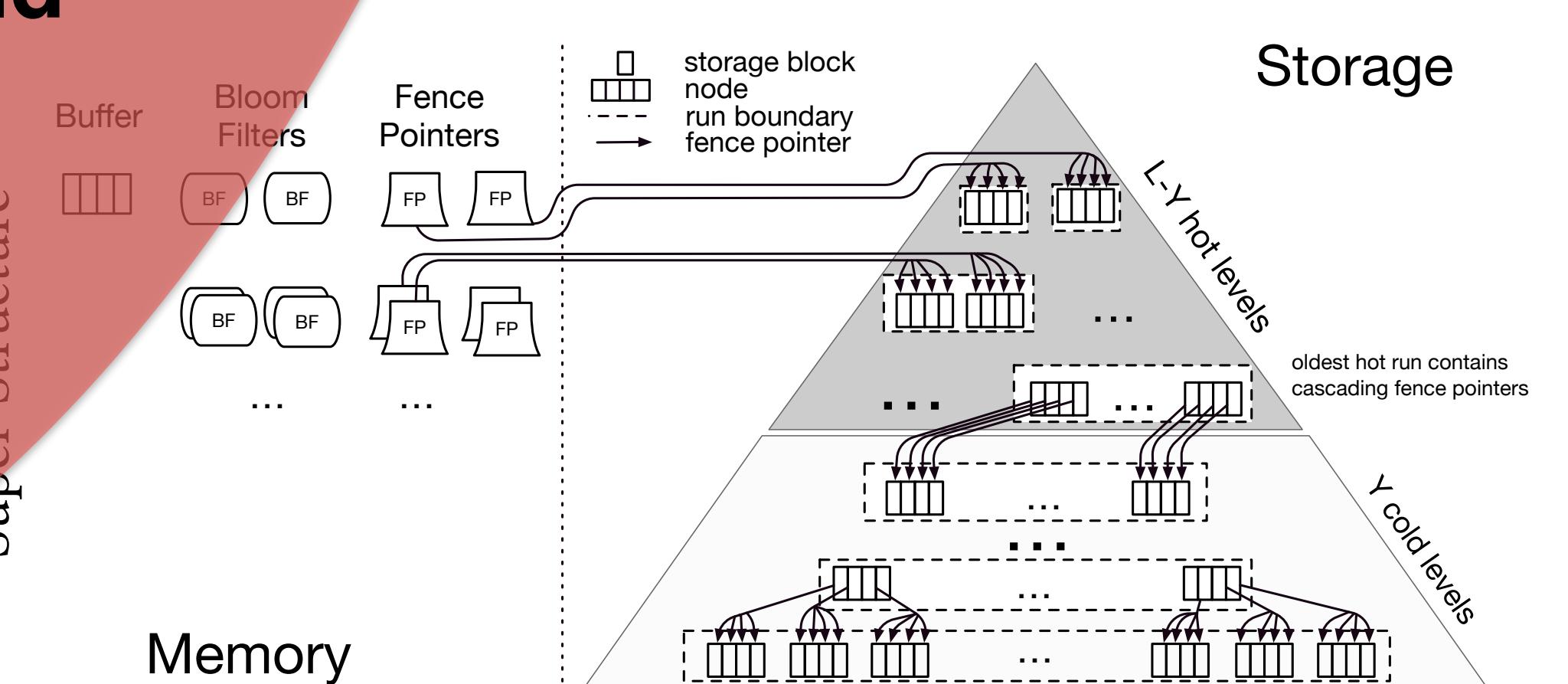
Term	Name	Description	Min. Value	Max. Value	Units
B	Block Size	# data entries that fit in a storage block.			Entries
M	Memory	Total main memory budget.	$\frac{B \cdot E + F \cdot T \cdot M_B}{E \cdot B}$	$N \cdot E$	Bits
N	Dataset Size	# total data entries in the dataset			Entries
E	Entry Size	Size of an entry.			Bits
F	Key Size	Size of a key, also used to approximate size of a fence (fence key and pointer).			Bits
s	Avg. Selectivity	Average selectivity of a long range query.			Entries
T	Growth Factor	Capacity ratio between adjacent levels.	2	B	Ratio
K	Hot Merge Threshold	Maximum # runs per hot level.	1	$T - 1$	Runs
Z	Cold Merge Threshold	Maximum # runs per cold level.	1	$T - 1$	Runs
D	Max. Node Size	Maximum size of a node; defines a contiguous data region.	1	$\frac{N}{B}$	Blocks
M_F	Fence & Filter Memory Budget	# bits of main memory budgeted to fence pointers and filters.	$\frac{F \cdot T \cdot M_B}{E \cdot B}$	M	Bits

Design Parameters

Environment Parameters

environment derived

design



Derived Term	Expression	Units
L (# total levels)	$\lceil \log_T \frac{N \cdot E}{M_B} \rceil$	Levels
X (Filters Memory Threshold)	$\frac{1}{\ln 2^2} \cdot (\frac{\ln T}{T-1} + \frac{\ln K - \ln Z}{T})$	Bits per Entry
M_{FH_I} (M_F Threshold: Hot Levels Saturation)	$N \cdot (\frac{X}{T} + \frac{F}{B})$	Bits
M_{FL_O} (M_F Threshold: Cold Levels Saturation)	$\frac{M_B \cdot P_T}{B}$	Bits
Y (# Cold Levels)	$\begin{cases} \lfloor \log_T \frac{N}{M_F} \cdot (\frac{X}{T} + \frac{F}{B}) \rfloor & \text{if } M_{FL_O} < M_F < M_{FH_I} \\ L - 1 & \text{if } M_F = M_{FH_I} \end{cases}$	Levels
M_{FP} (Fence Pointer Memory Budget)	$T^{L-Y+1} \cdot F \cdot \frac{M_b}{E \cdot B} \cdot \frac{T}{T-1}$	Bits
M_{BF} (Filter Memory Budget)	$M_F - M_{FP}$	Bits
M_B (Buffer Memory Budget)	$B \cdot E + (M - M_F)$	Bits
p_{sum} (Sum of BF False Positive Rates)	$e^{-\frac{M_{BF}}{N} \cdot \ln(2)^2 \cdot T^Y \cdot Z^{\frac{T-1}{T}} \cdot K^{\frac{1}{T}} \cdot \frac{T}{T-1}}$	
p_i (BF False Positive Rate at Level i)	$\begin{cases} \frac{1}{p_{sum}} & \text{if } i > L - Y \\ \frac{Z}{p_{sum}} \cdot \frac{T-1}{T} & \text{if } i = L - Y \\ \frac{1}{T^{L-Y-i}} & \text{if } i < L - Y \end{cases}$	Probability

Derived Design Rules

Operation	General Cost Model
Update	$O(\frac{1}{B} \cdot (\frac{1}{K} \cdot (L - Y - 1) + \frac{T}{Z} \cdot (Y + 1)))$
Zero Result Lookup	$O(Z \cdot e^{-\frac{M_{BF}}{N} \cdot T^Y} + Y \cdot Z)$
Single Result Lookup	$O(1 + Z \cdot e^{-\frac{M_{BF}}{N} \cdot T^Y} + Y \cdot Z)$
Short Scan	$O(K \cdot L - Y - 1) + Z \cdot (Y + 1))$
Long Scan	$O(Z \cdot e^{-\frac{M_{BF}}{N} \cdot T^Y} + Y \cdot Z)$

unified I/O models

from write to read optimized

Design Continuums @CDIR2019

Cosine

@PVLDB2022

Cloud-cost Optimized

Self Designing

Key-value Store

Design Abstractions of Template		Type/Domain	Example templates for diverse data structures			
			LSM variants	B-Tree variants	LSH variants	A new design
LAYOUT PRIMITIVES	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 <i>max size set to 1 GB</i>	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			

Storage engine template in Cosine and example initializations for diverse storage engine designs.

Cosine

@PVLDB2022

Cloud-cost Optimized

Self Designing

Key-value Store

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 <i>max size set to 1 GB</i>	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			

Storage engine template in Cosine and example initializations for diverse storage engine designs.

Cosine

@PVLDB2022

Cloud-cost Optimized

Self Designing

Key-value Store

Design Abstractions of Template		Type/Domain	Example templates for diverse data structures					
			LSM variants	B-Tree variants	LSH variants	A new design		
LAYOUT PRIMITIVES	ALGORITHMIC ABSTRACTIONS	<i>initialized by search through engine design space</i>	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 <i>max size set to 1 GB</i>	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, ..] (T is large)	[1000, 1001, ...]	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			

Storage engine template in Cosine and example initializations for diverse storage engine designs.

Cosine

@PVLDB2022

Cloud-cost Optimized

Self Designing

Key-value Store

Design Abstractions of Template		Type/Domain	Example templates for diverse data structures			
			LSM variants	B-Tree variants	LSH variants	A new design
LAYOUT PRIMITIVES	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 <i>max size set to 1 GB</i>	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			

Storage engine template in Cosine and example initializations for diverse storage engine designs.

Cosine

@PVLDB2022

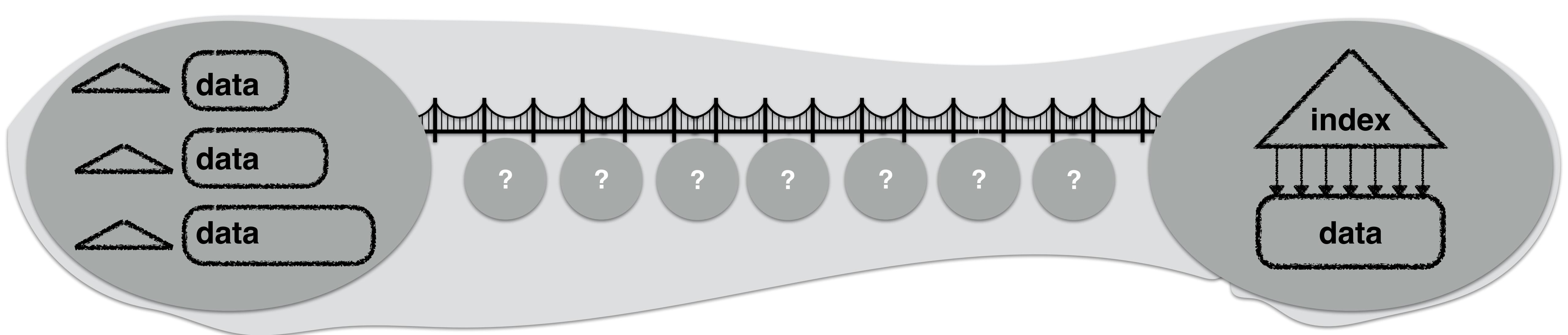
Cloud-cost Optimized

Self Designing

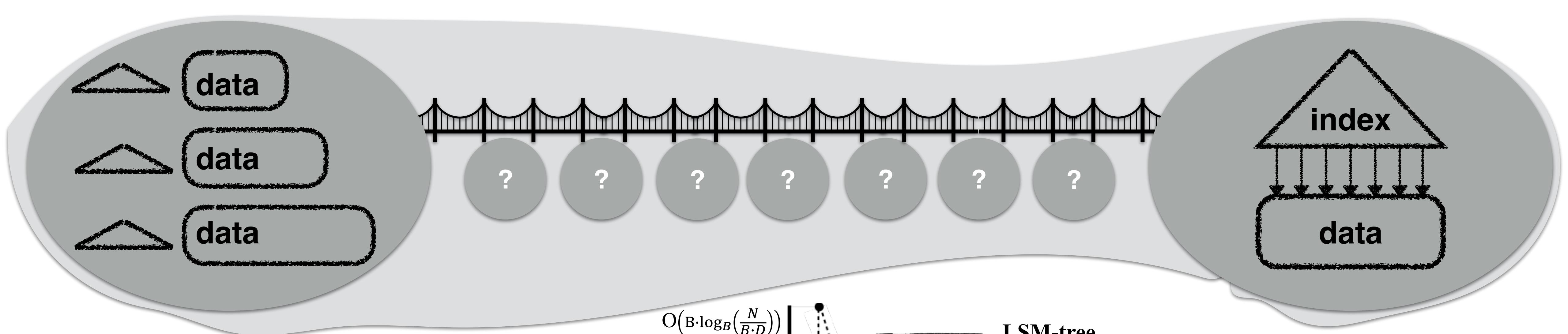
Key-value Store

Design Abstractions of Template		Type/Domain	Example templates for diverse data structures			
			LSM variants	B-Tree variants	LSH variants	A new design
LAYOUT PRIMITIVES	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 <i>max size set to 1 GB</i>	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			

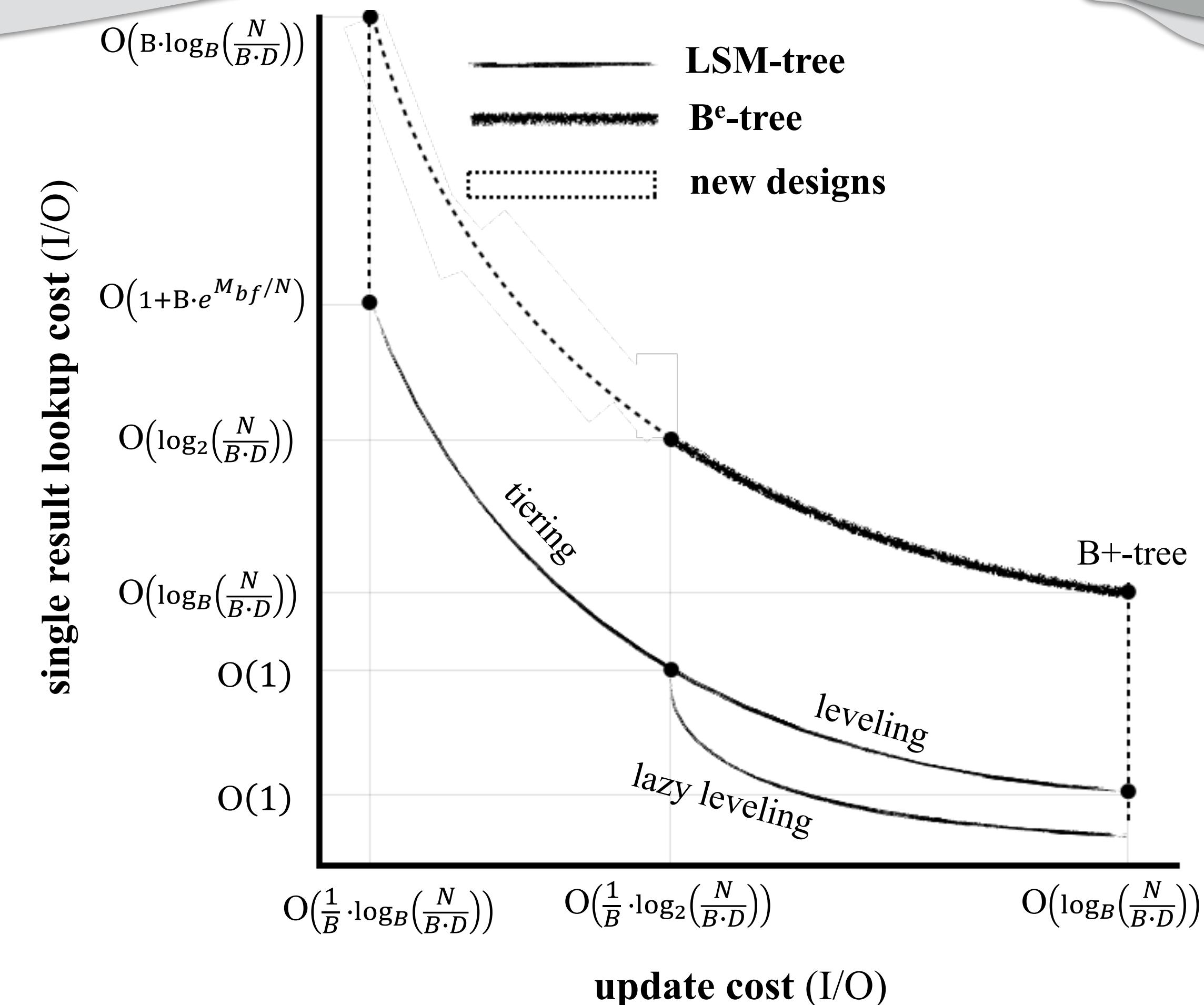
Storage engine template in Cosine and example initializations for diverse storage engine designs.



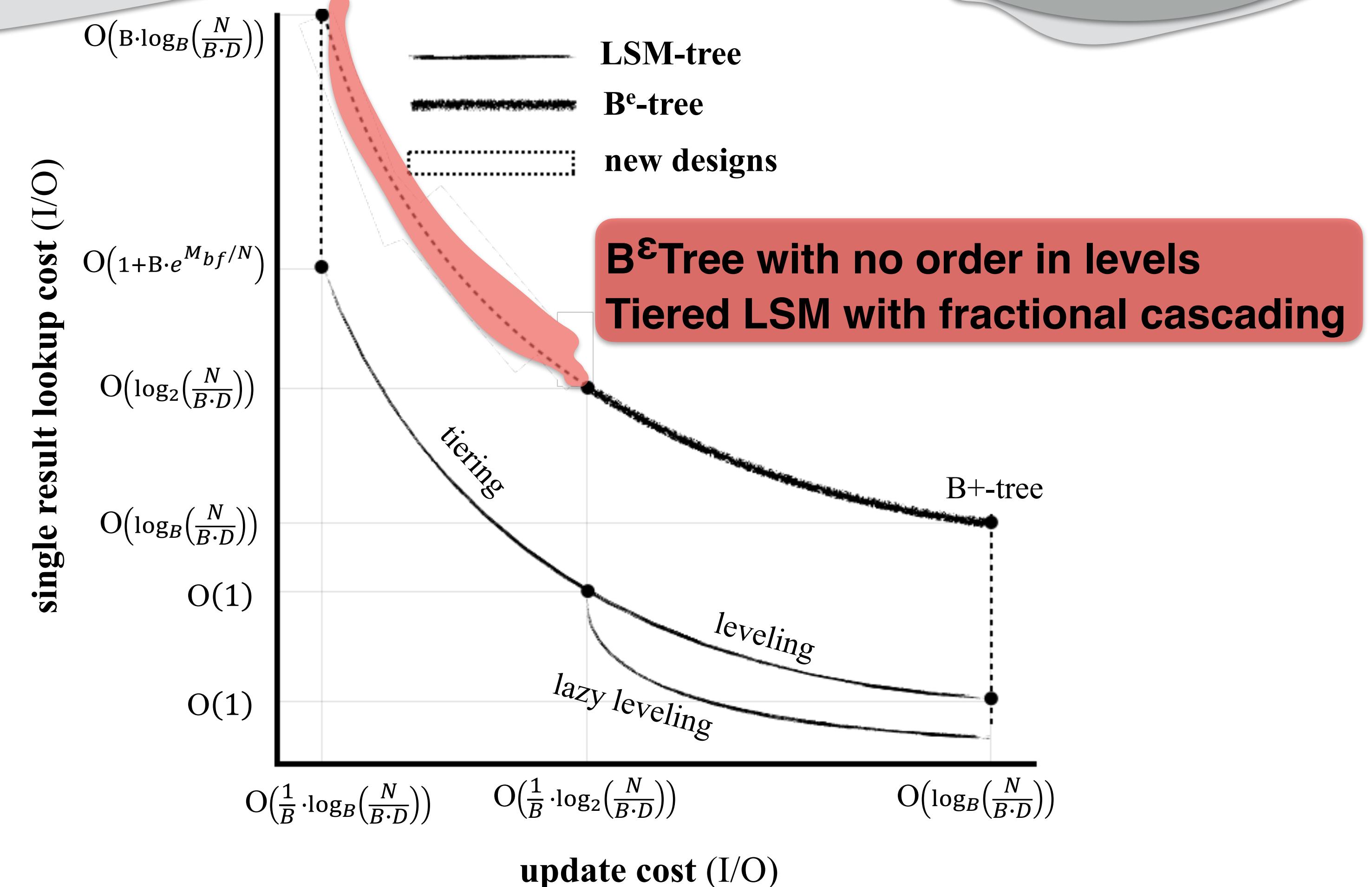
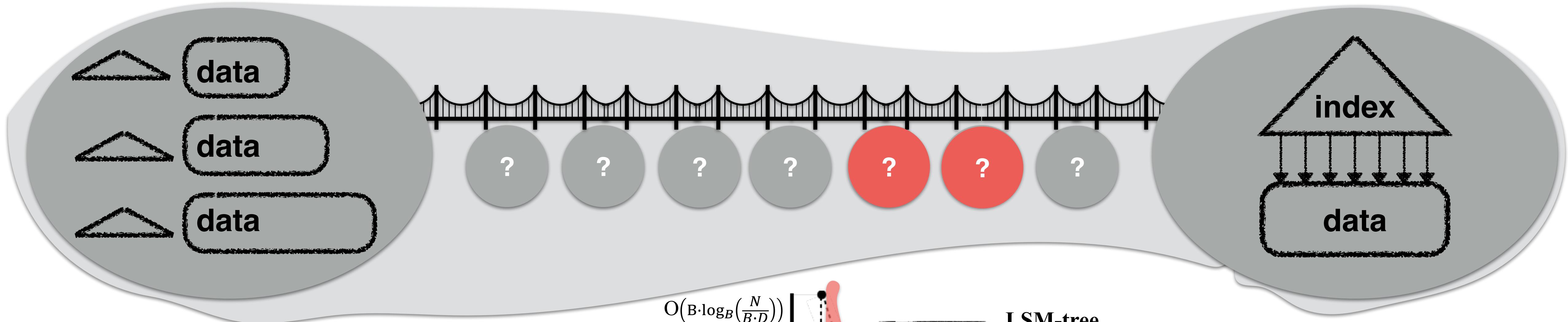
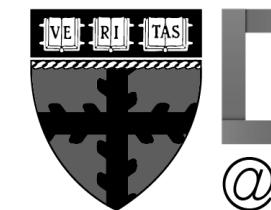
assisted discovery



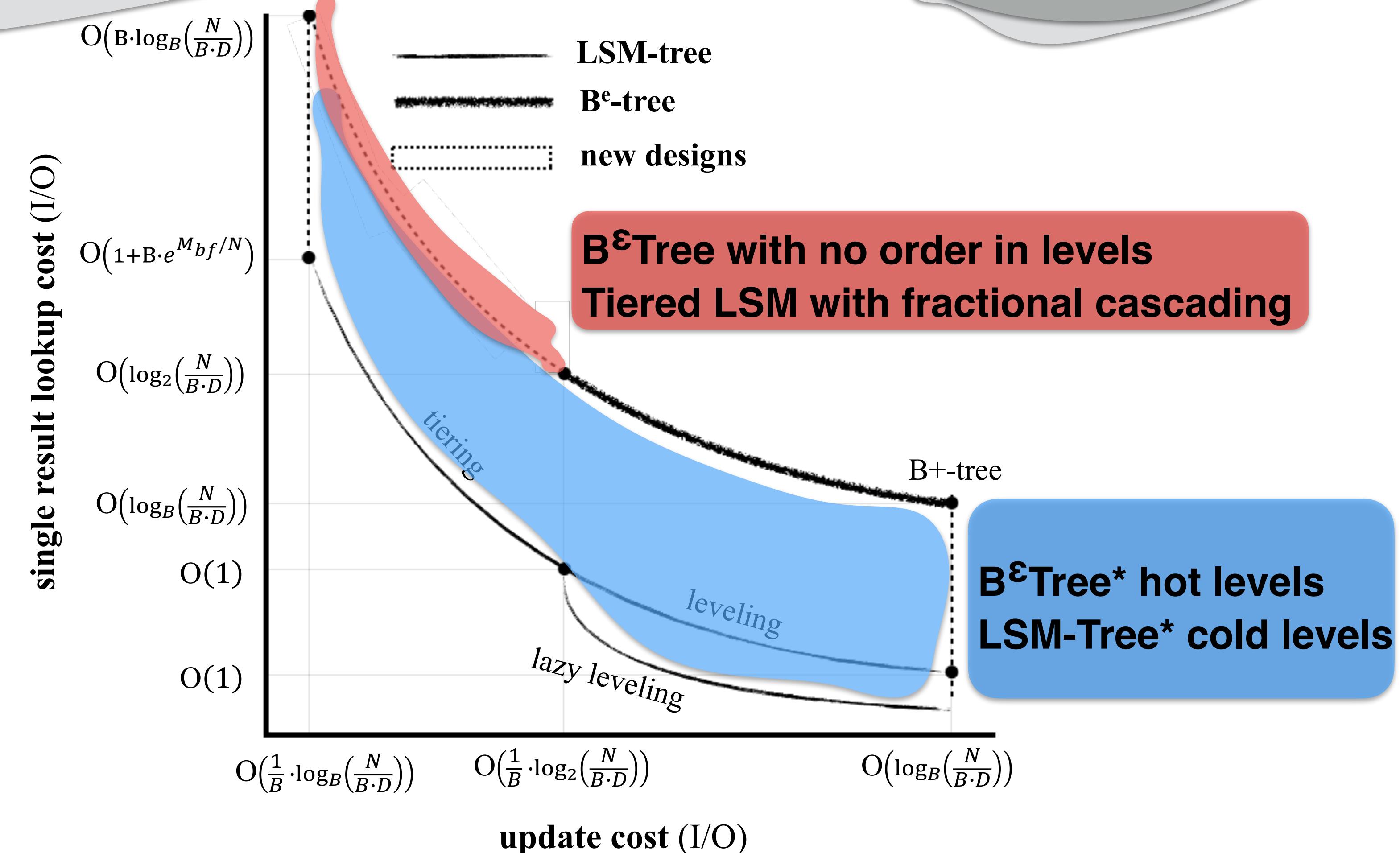
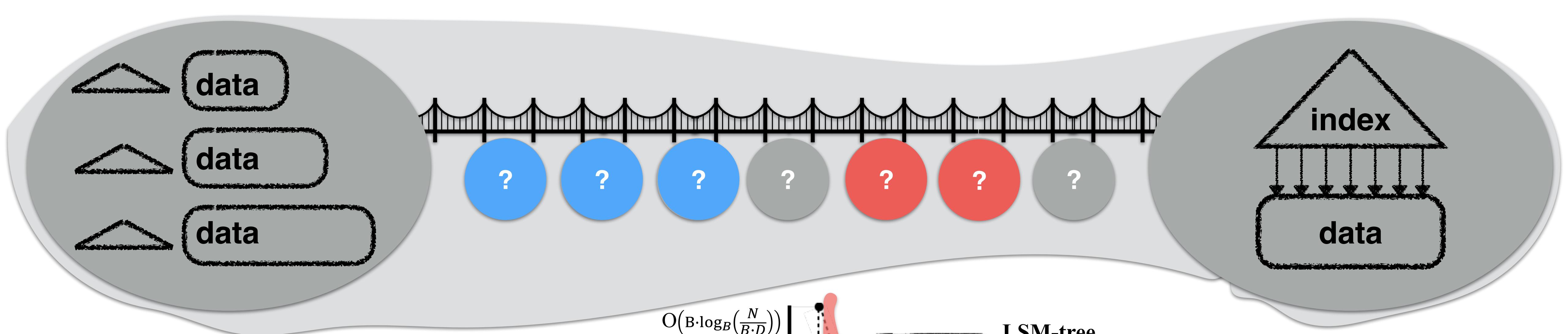
assisted discovery



assisted discovery



assisted discovery



Design Primitives

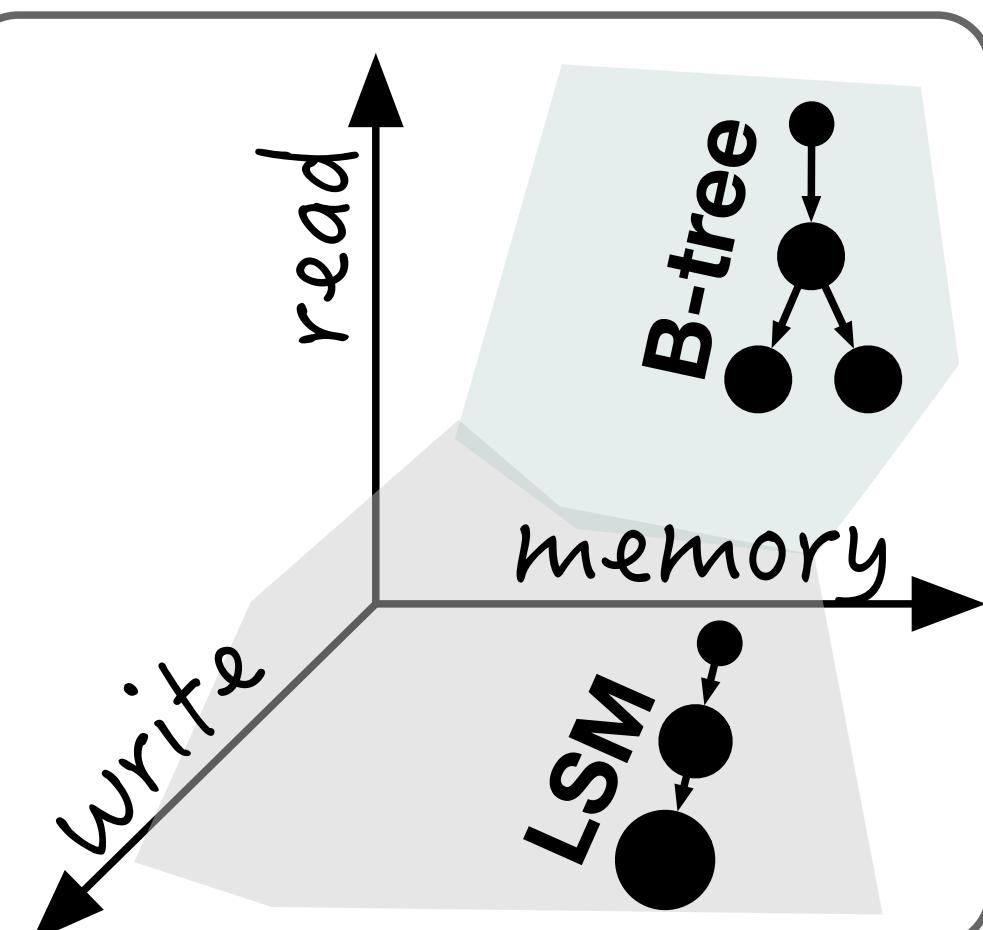
- fanout
- filter bits
- buffer size
- ...
- merge policy

data structures

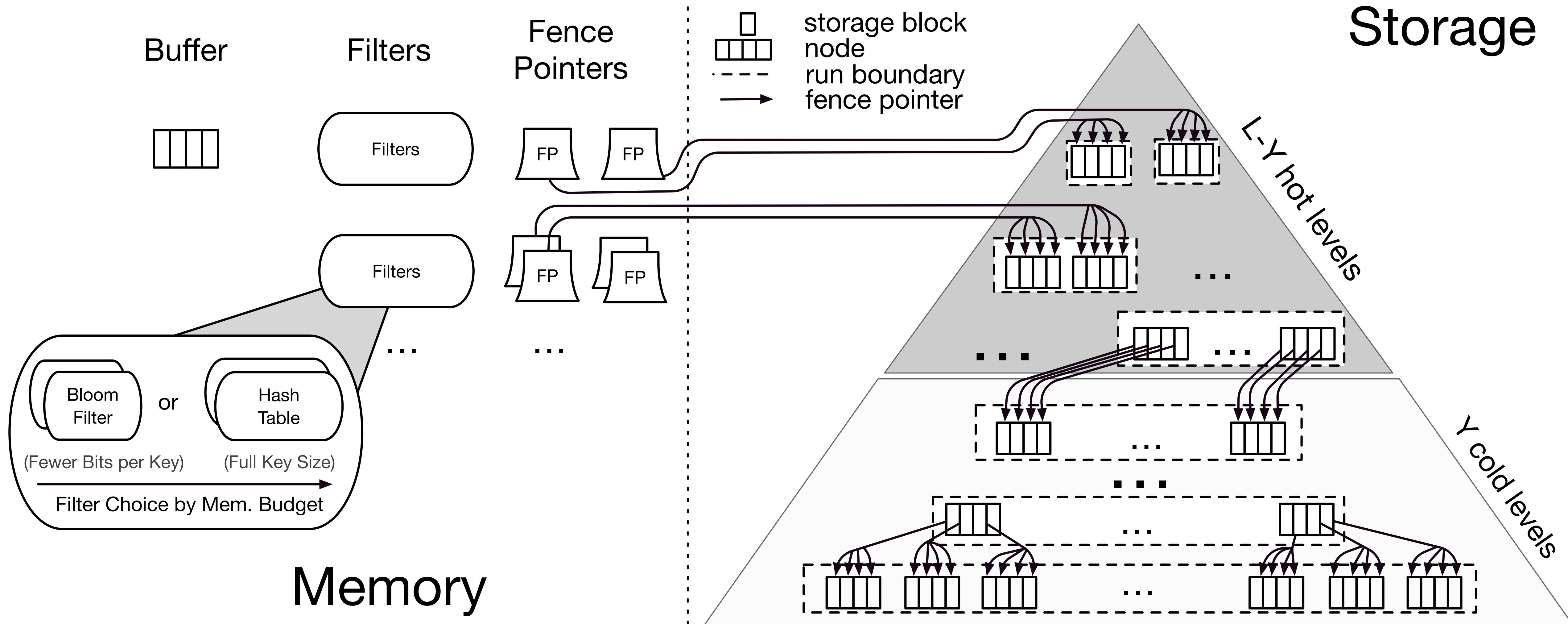
design space

Data structure designs are derived as combinations of fundamental design primitives

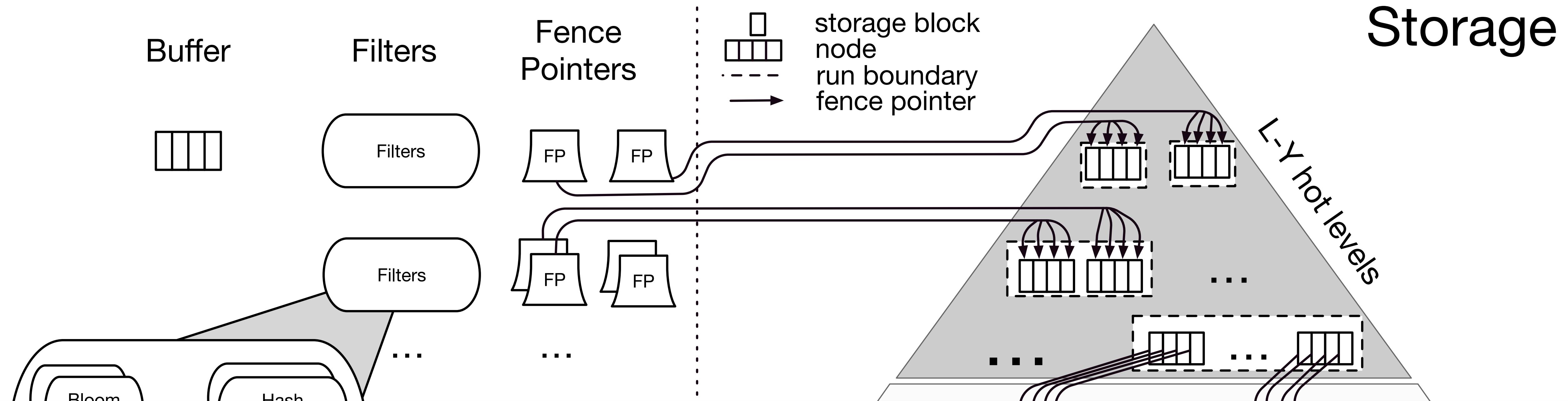
design continuum



unified design storage engine template



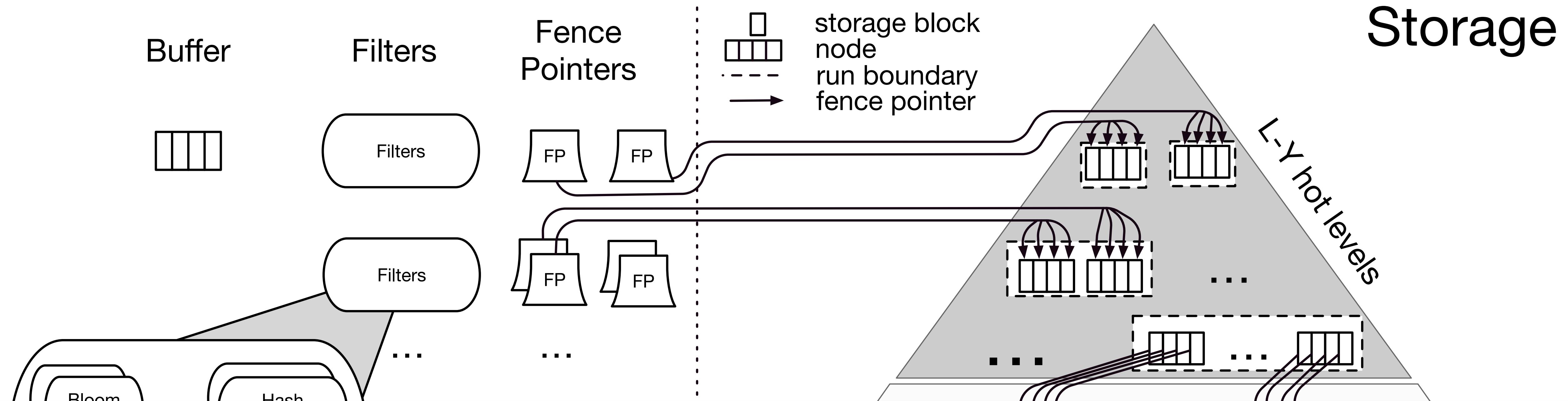
unified design storage engine template



analytical
I/O model

unified
closed form

unified design storage engine template



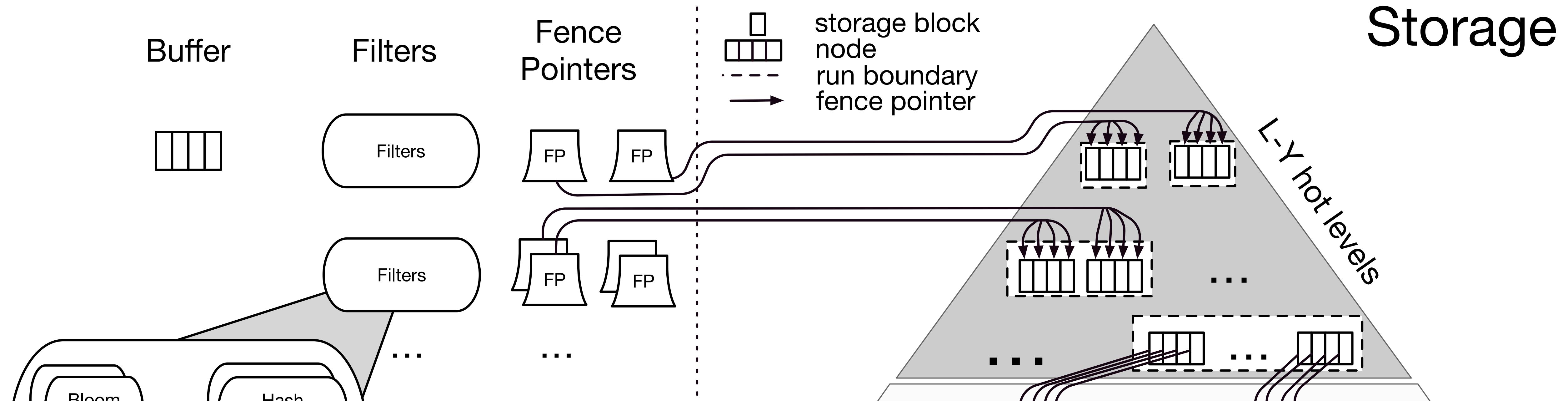
analytical
I/O model

unified
closed form

learned
CPU model

h/w,
parallelism

unified design storage engine template



analytical
I/O model

unified
closed form

learned
CPU model

h/w,
parallelism

cloud cost
mapping & SLA

AWS, Azure, Google



workload
budget
perf.

**analytical
I/O model**

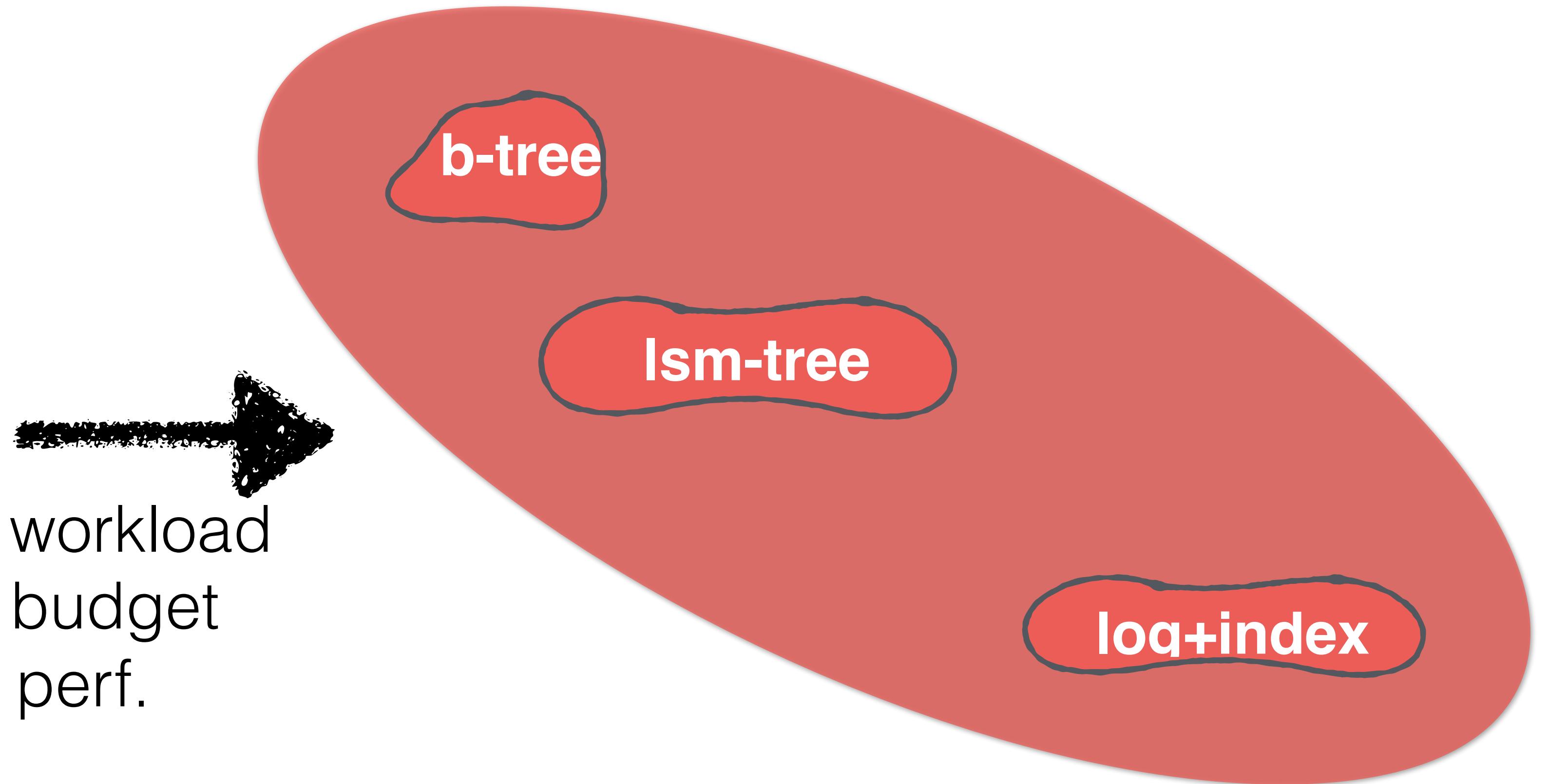
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

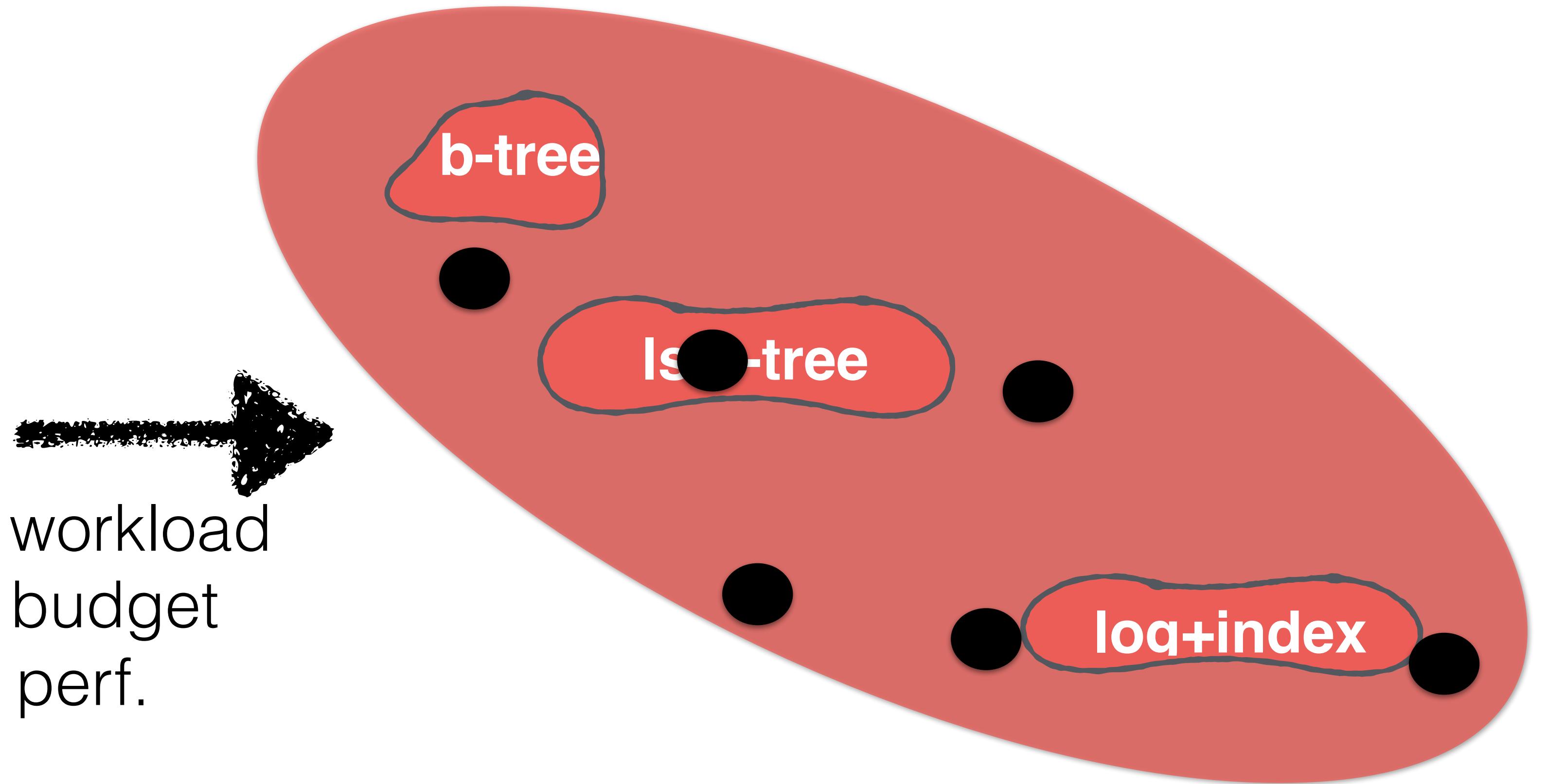
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

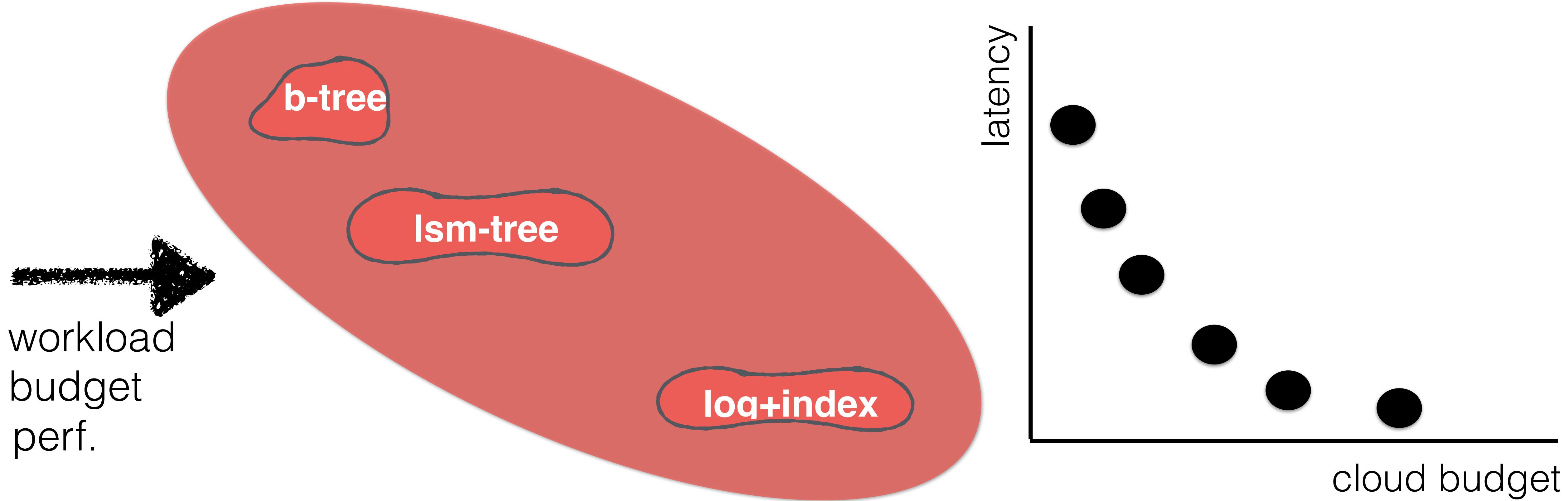
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

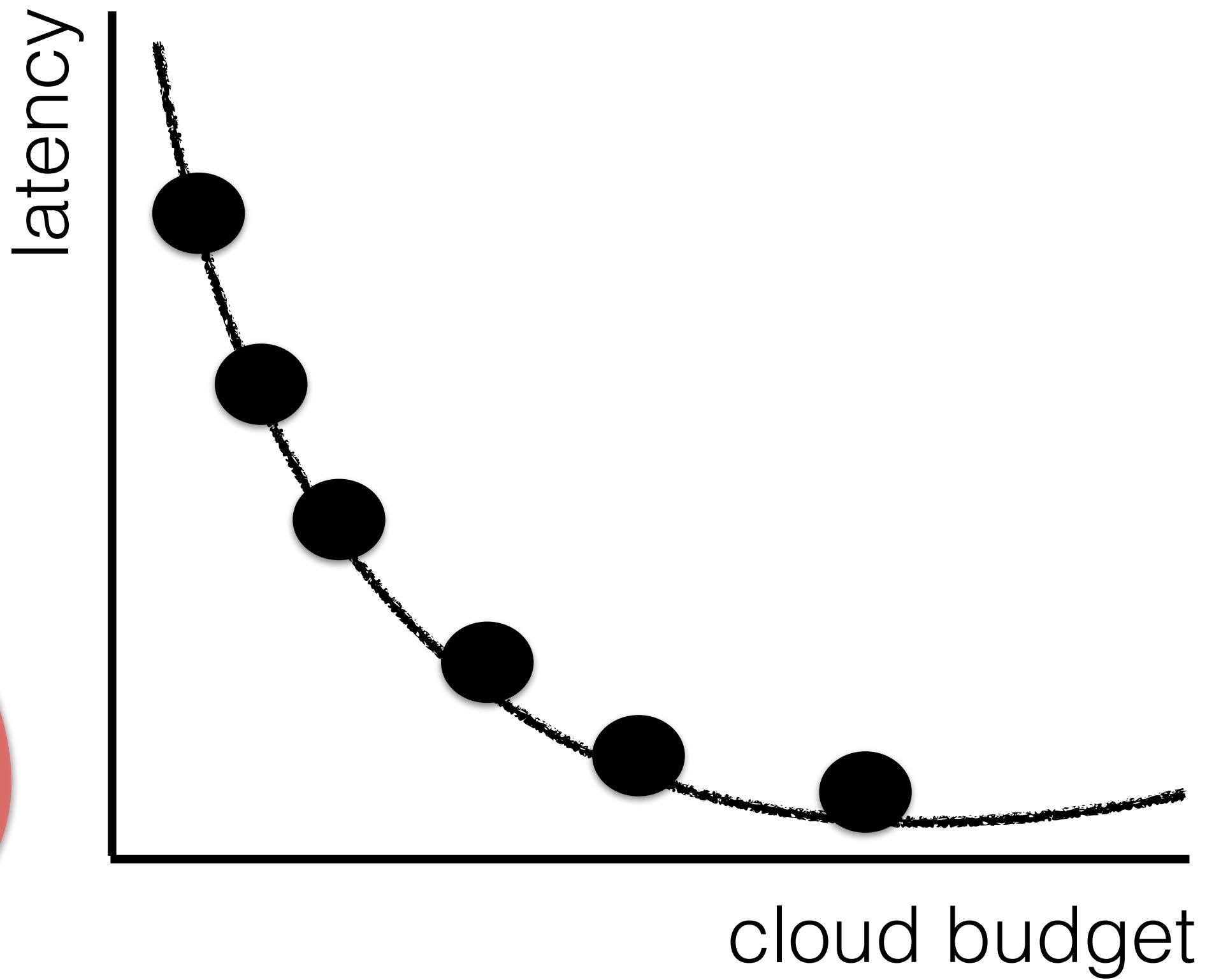
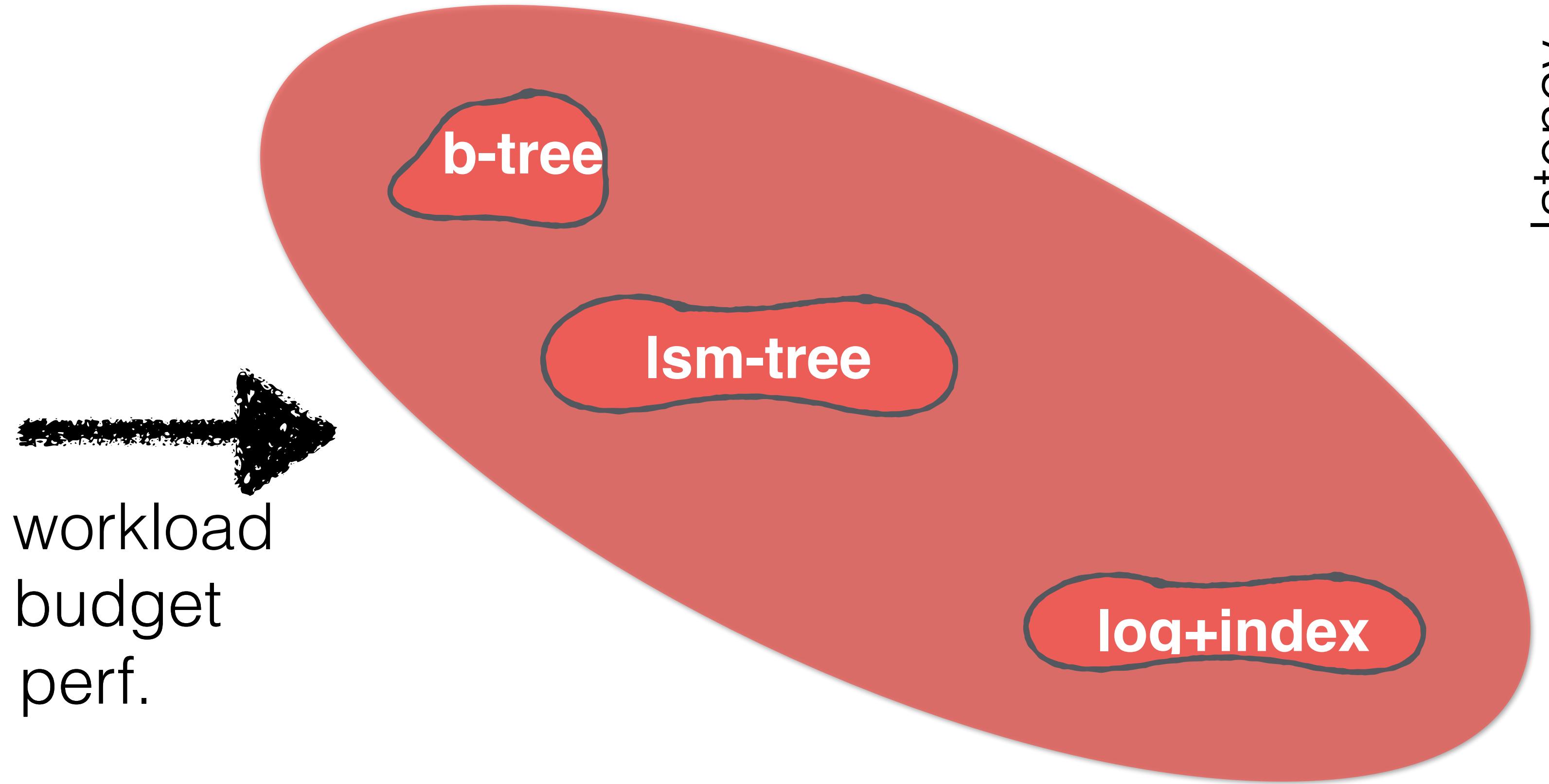
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

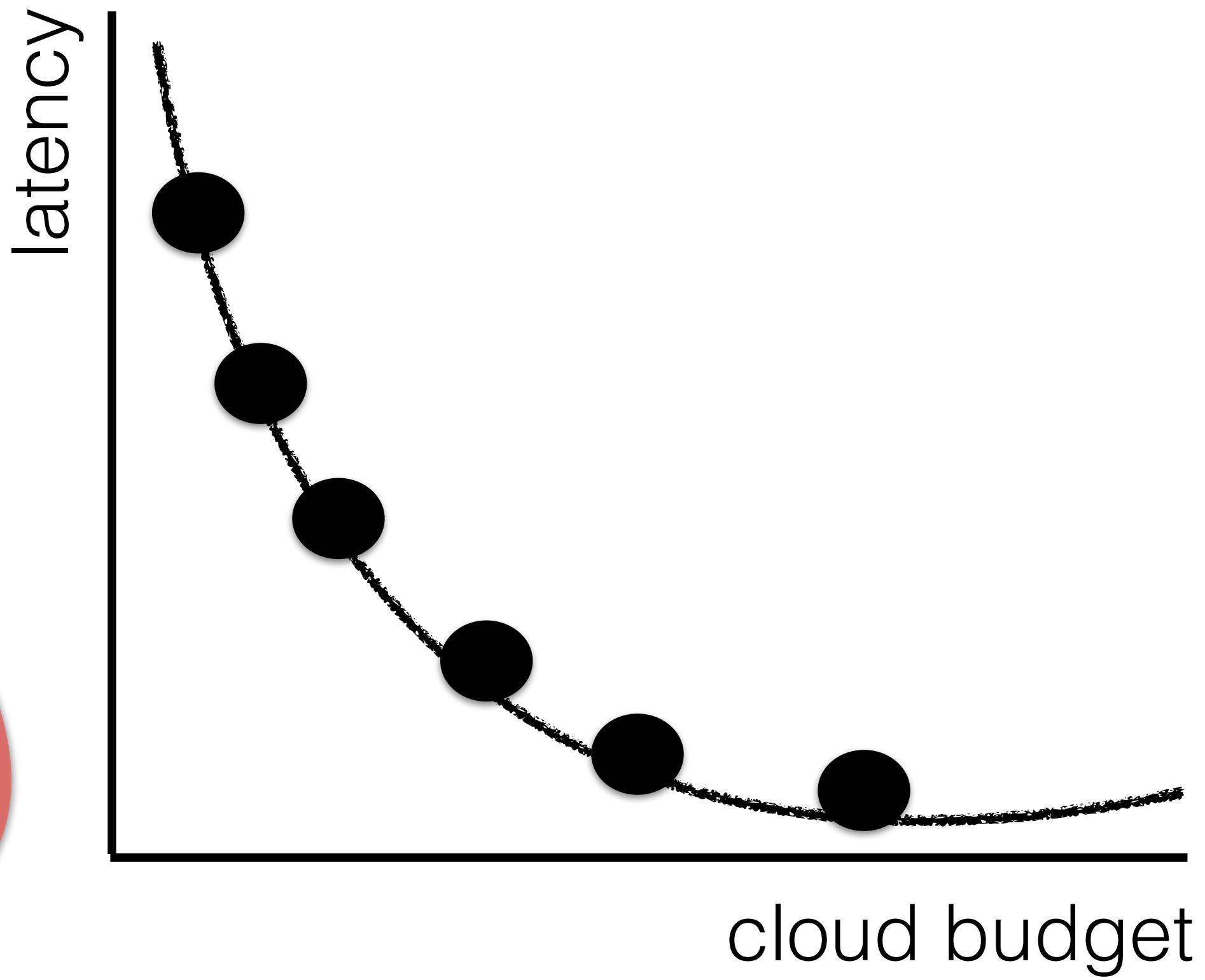
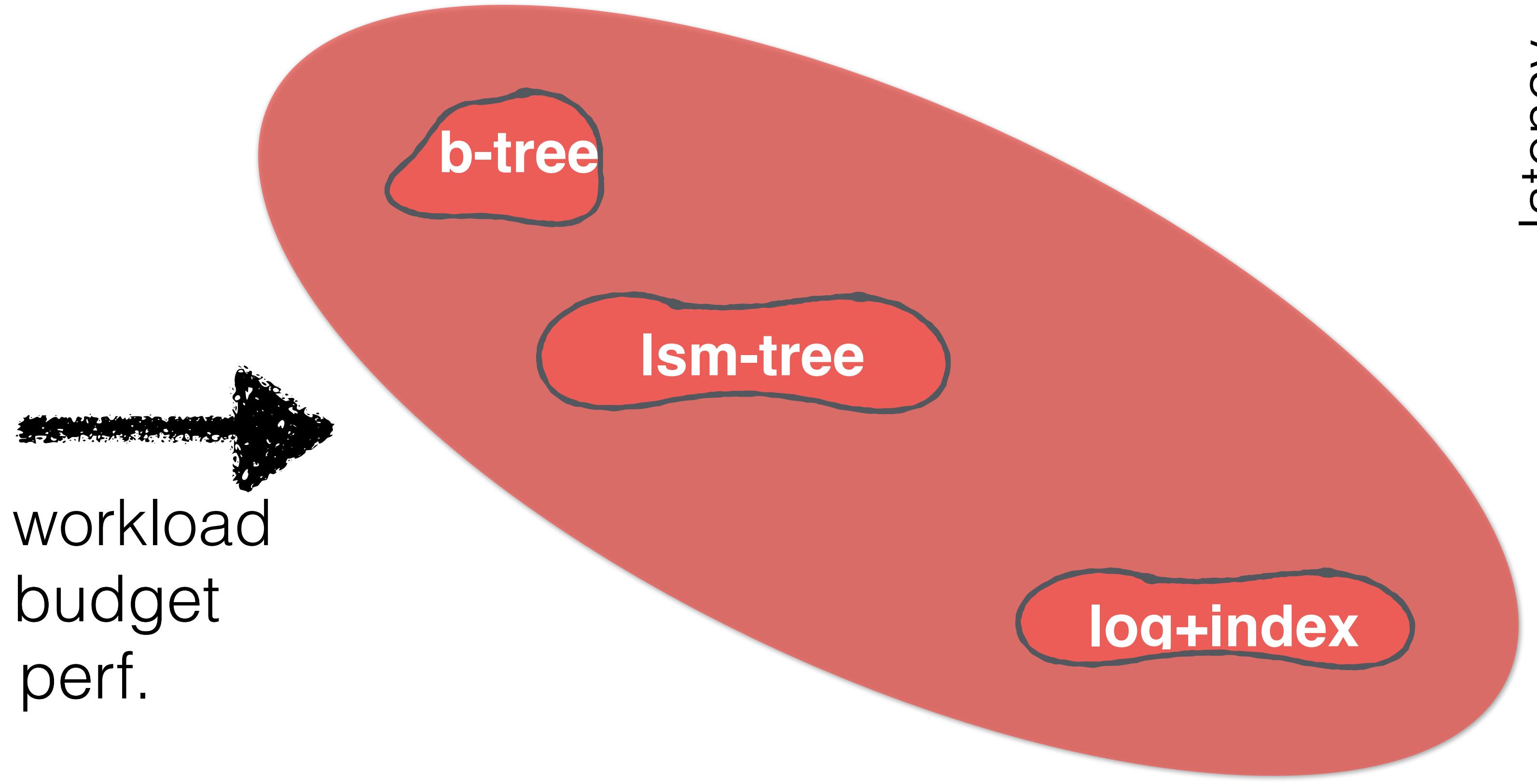
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

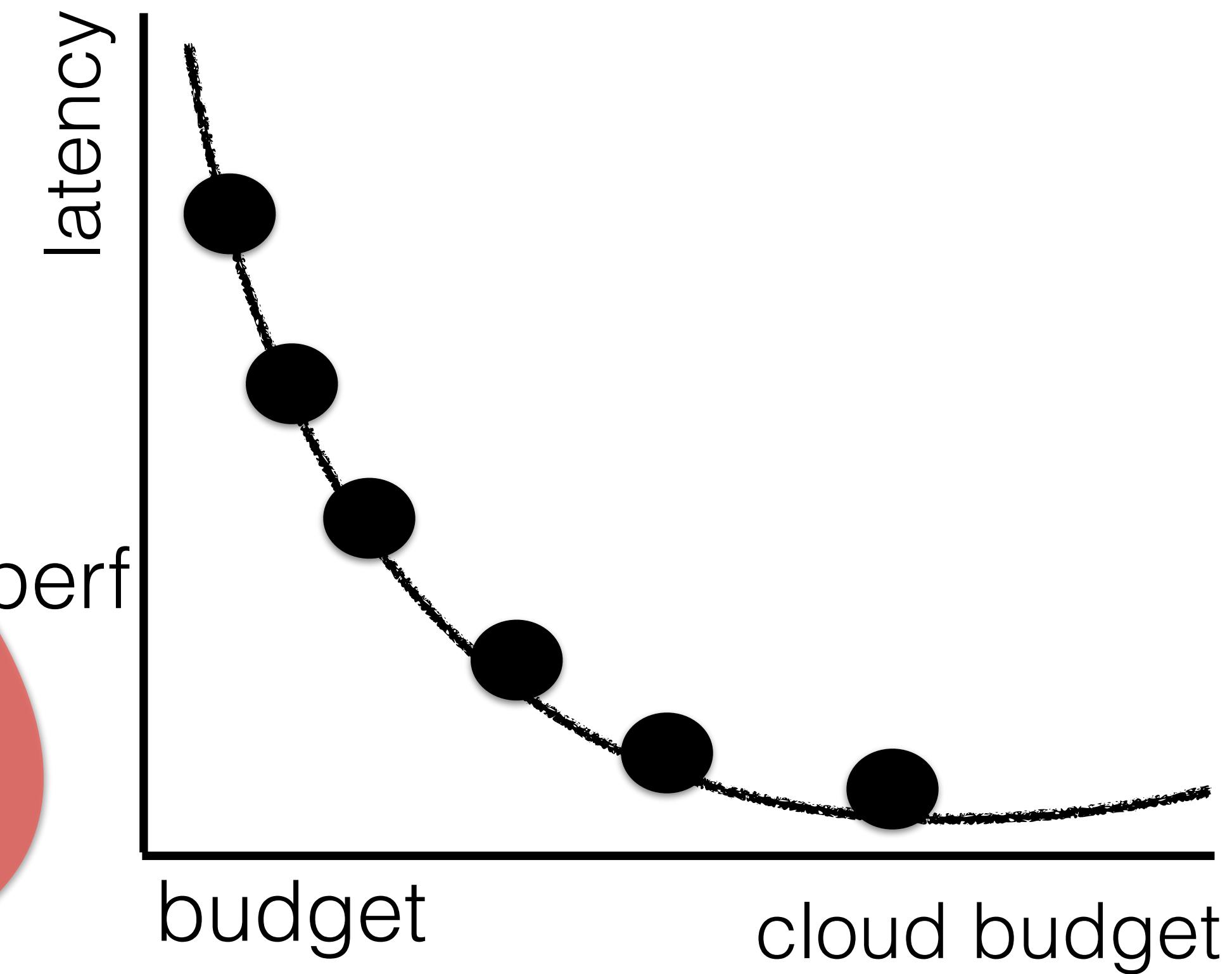
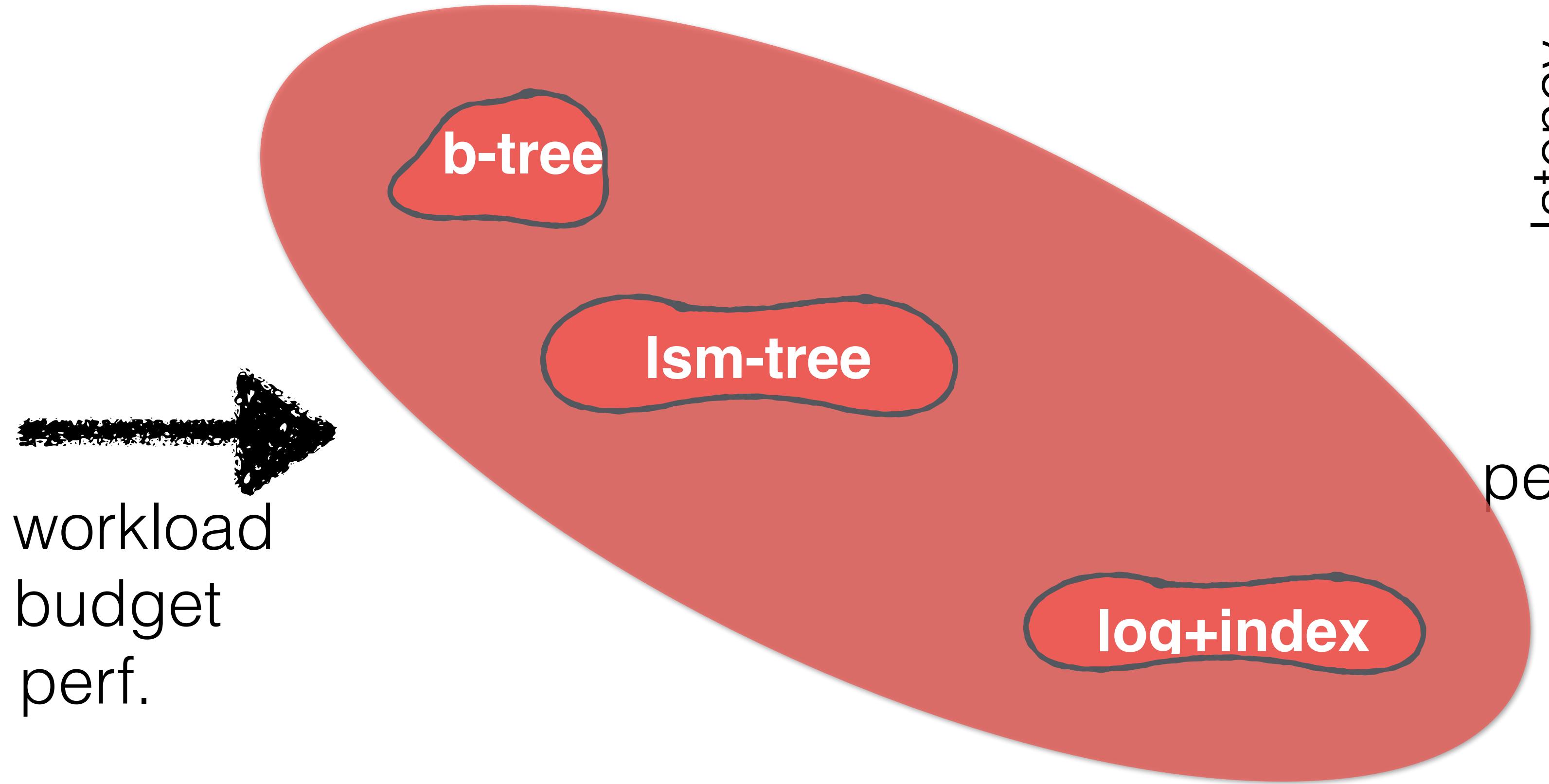
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

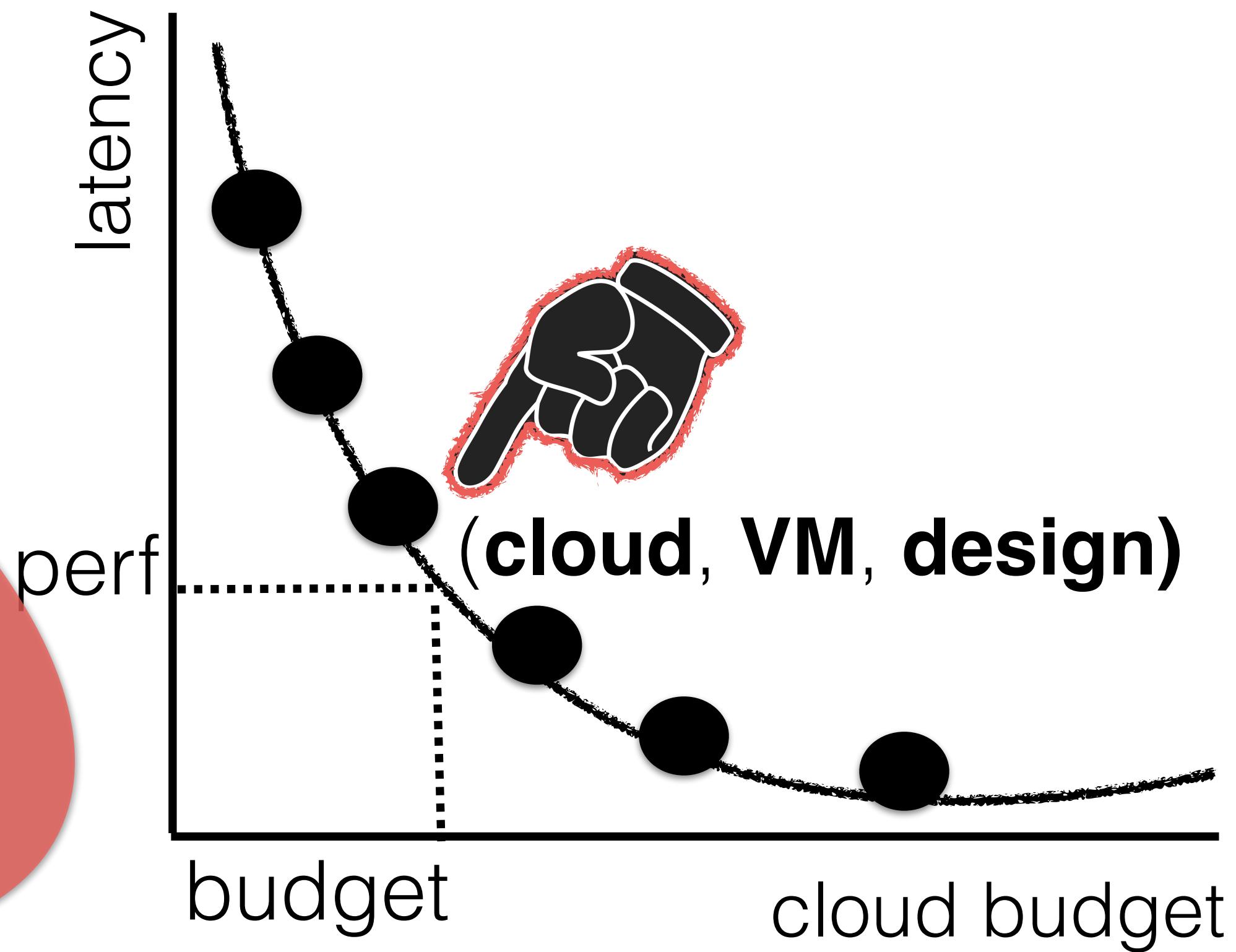
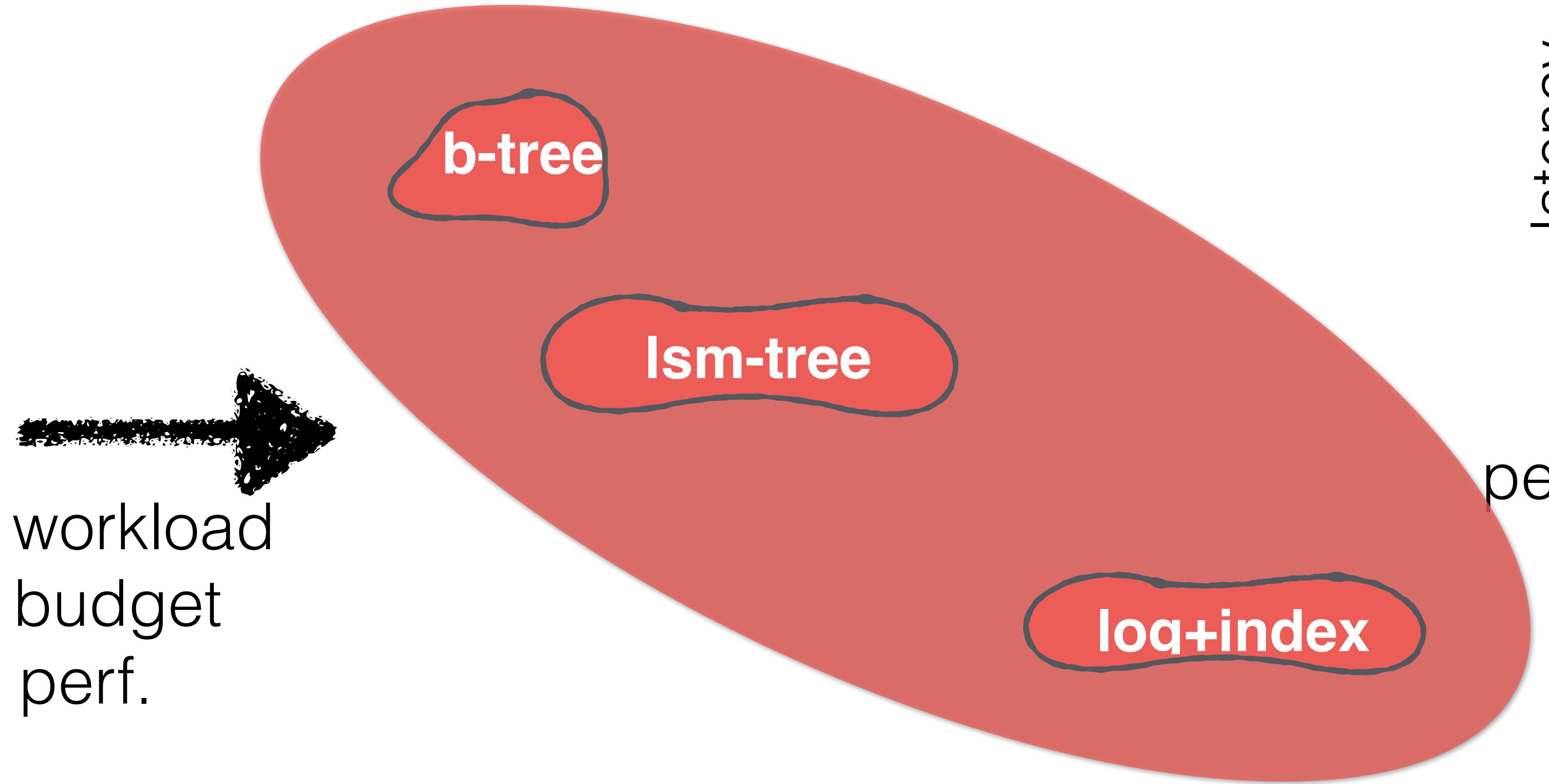
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

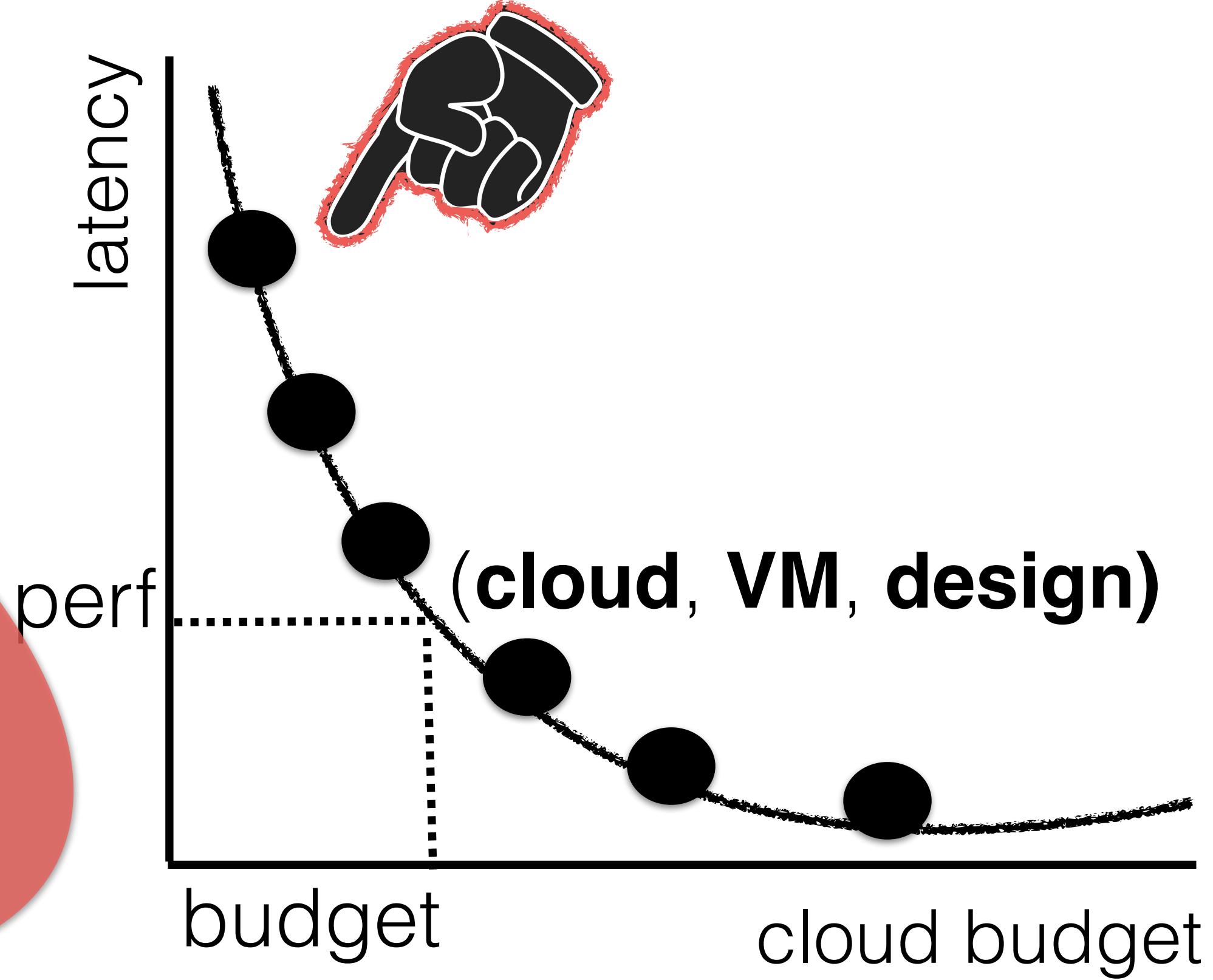
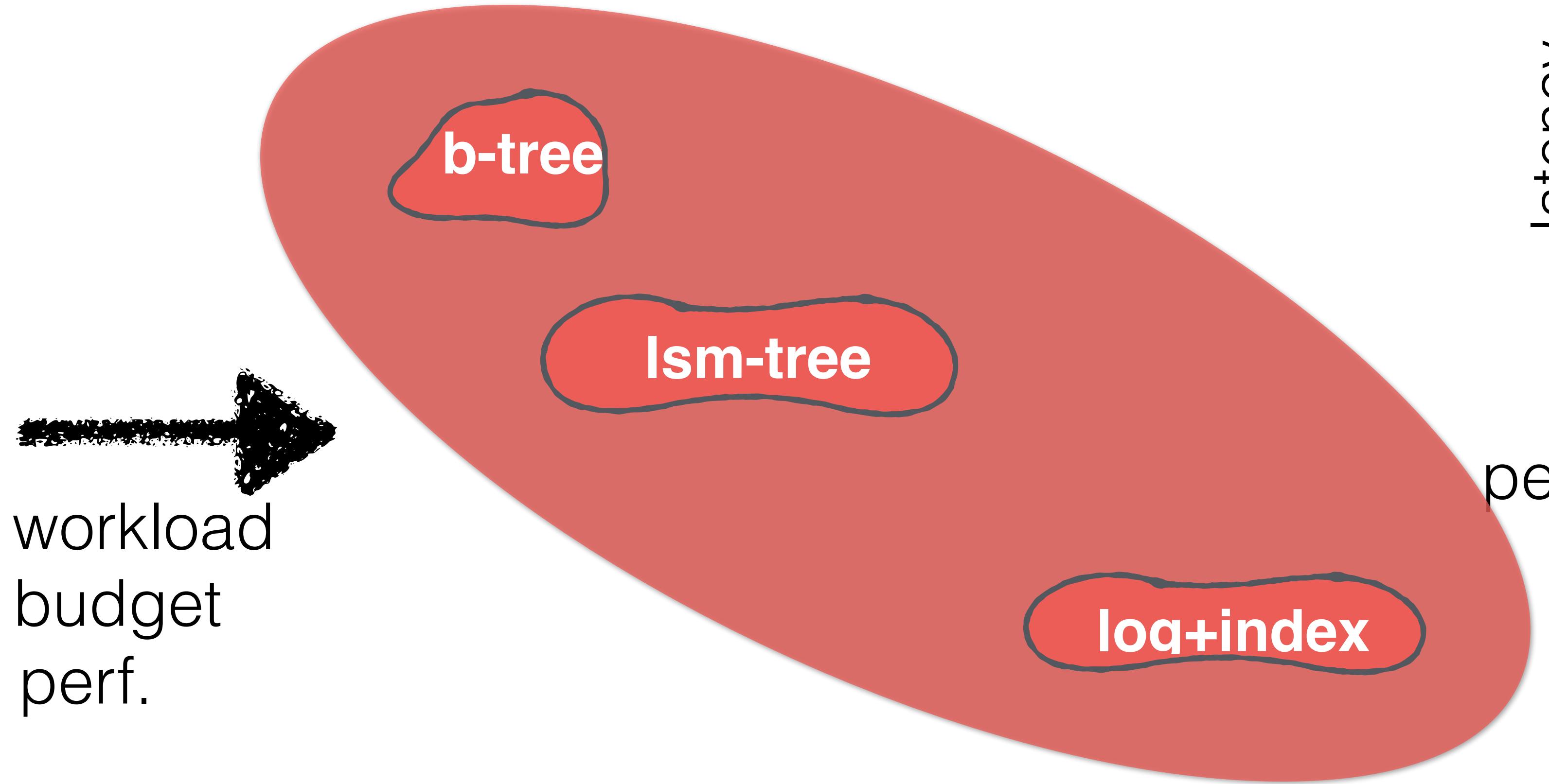
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

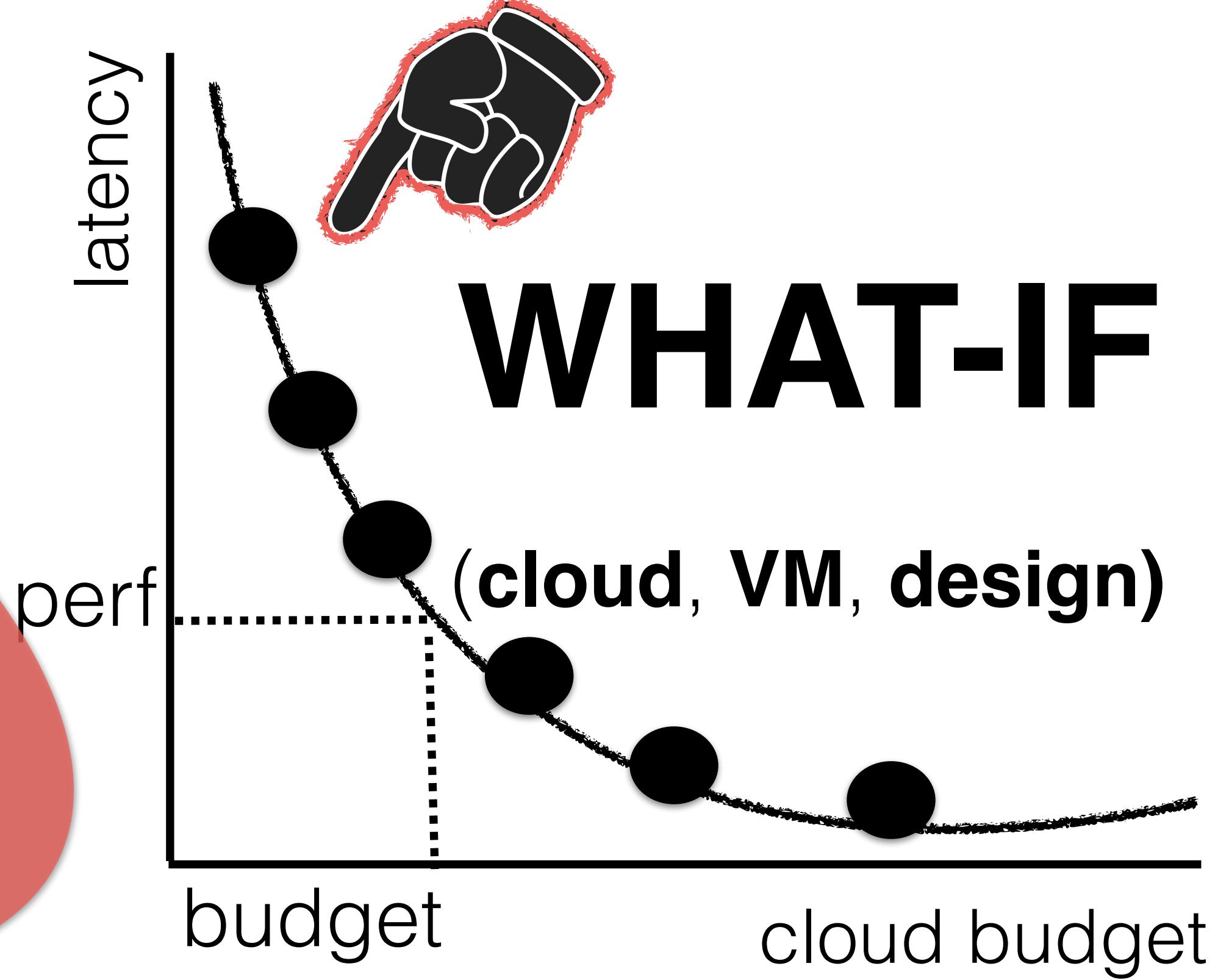
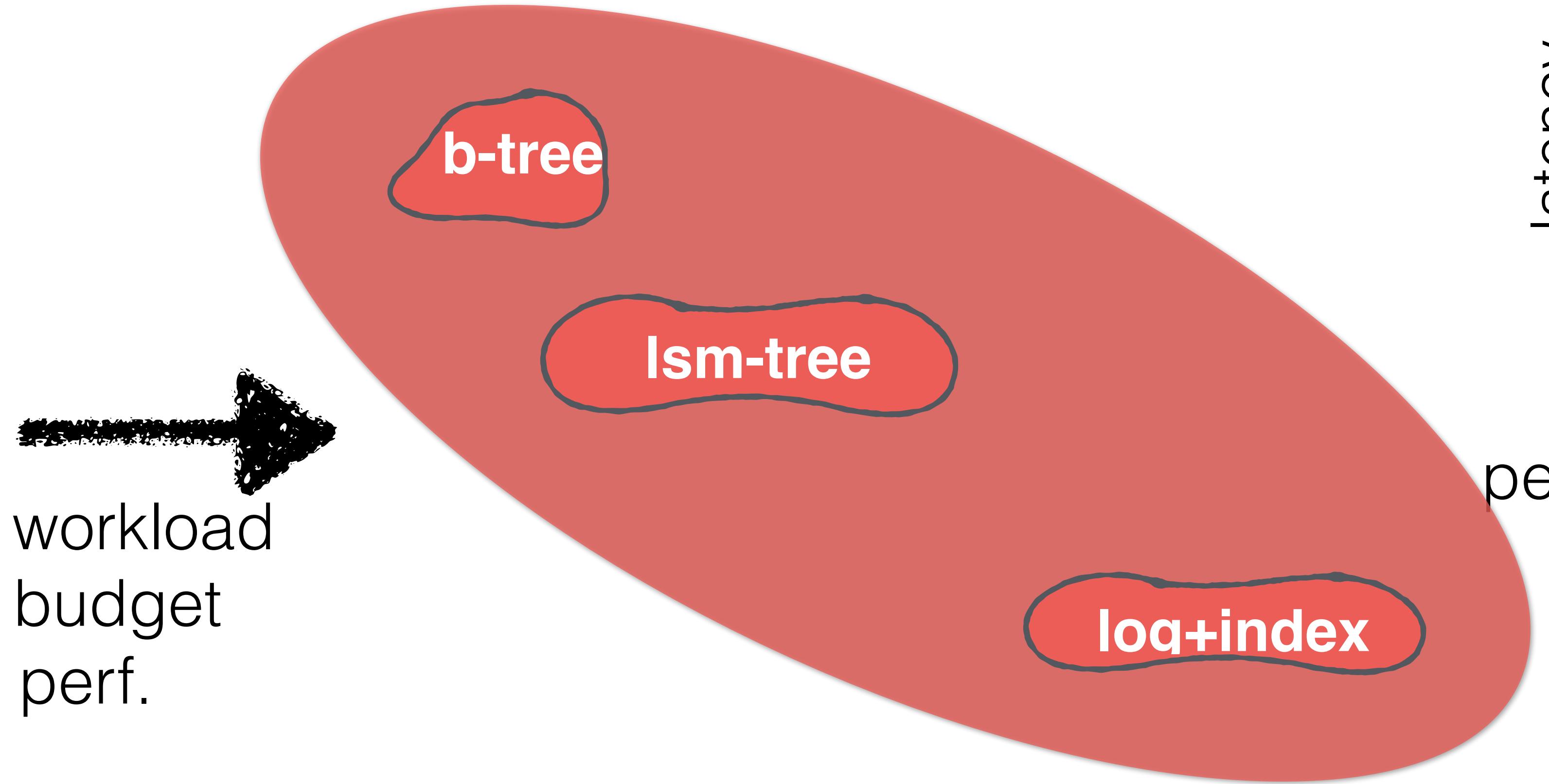
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

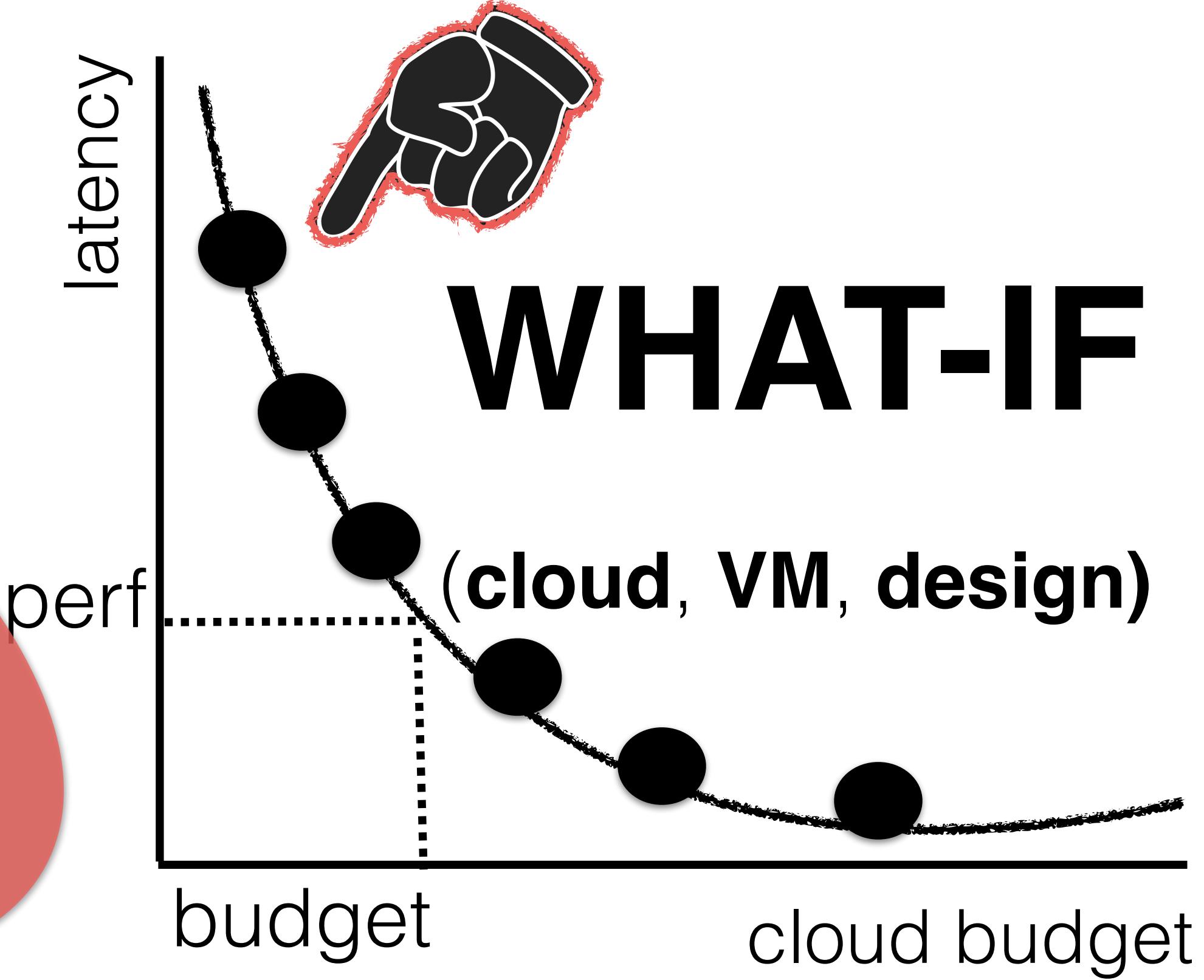
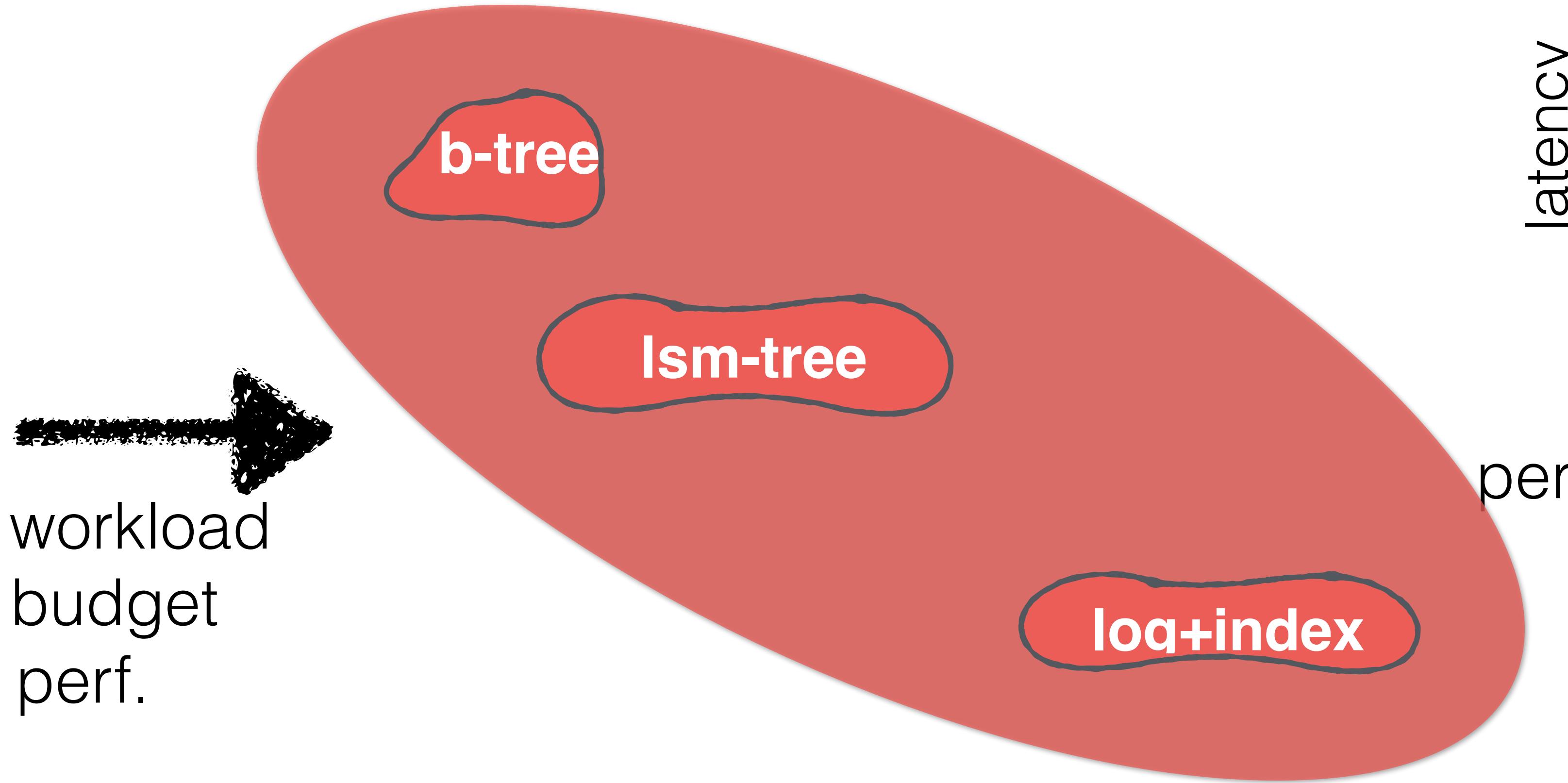
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



How to test?

workload/budget diversity

How to test?

workload/budget diversity

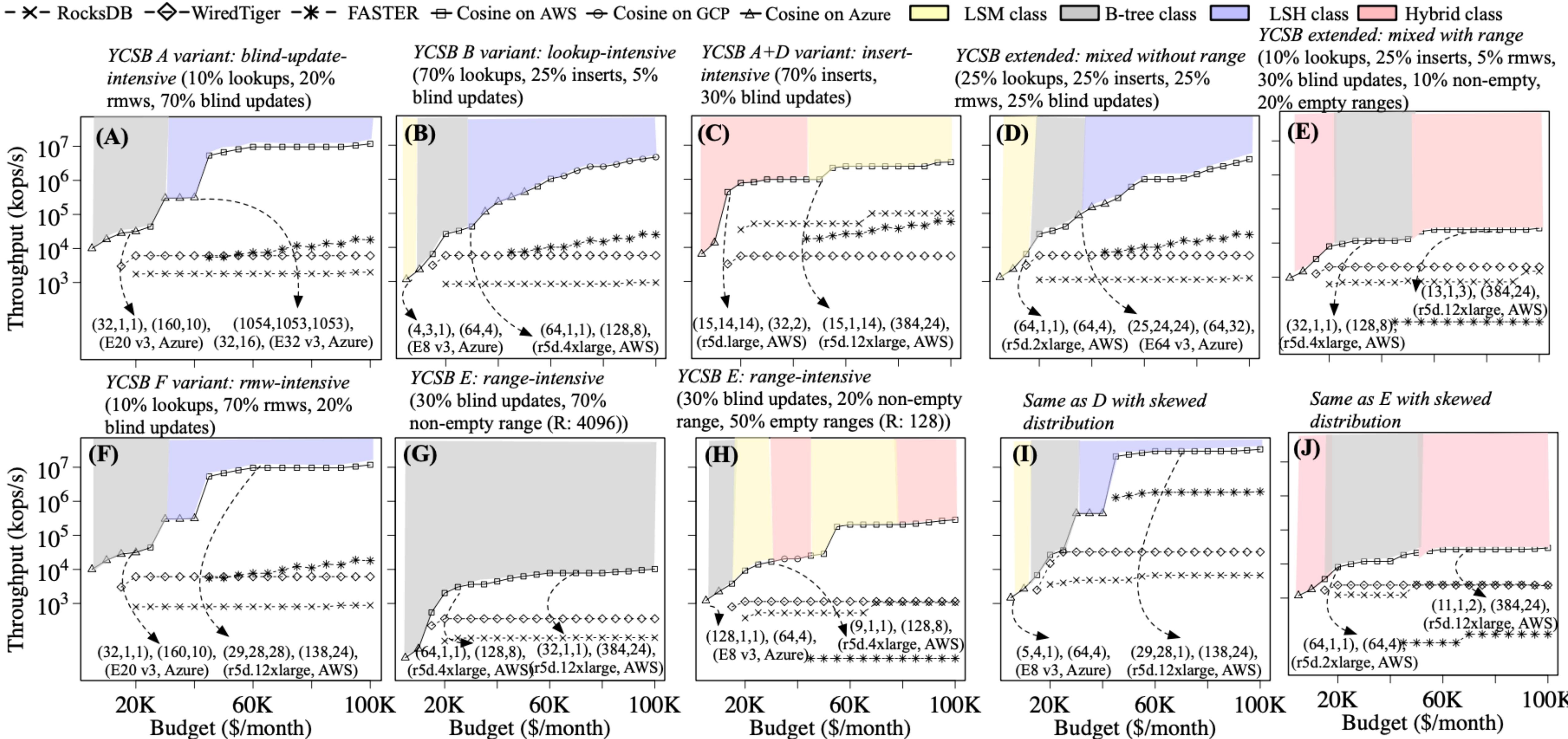
Can we beat the best system for every workload?

How to test?

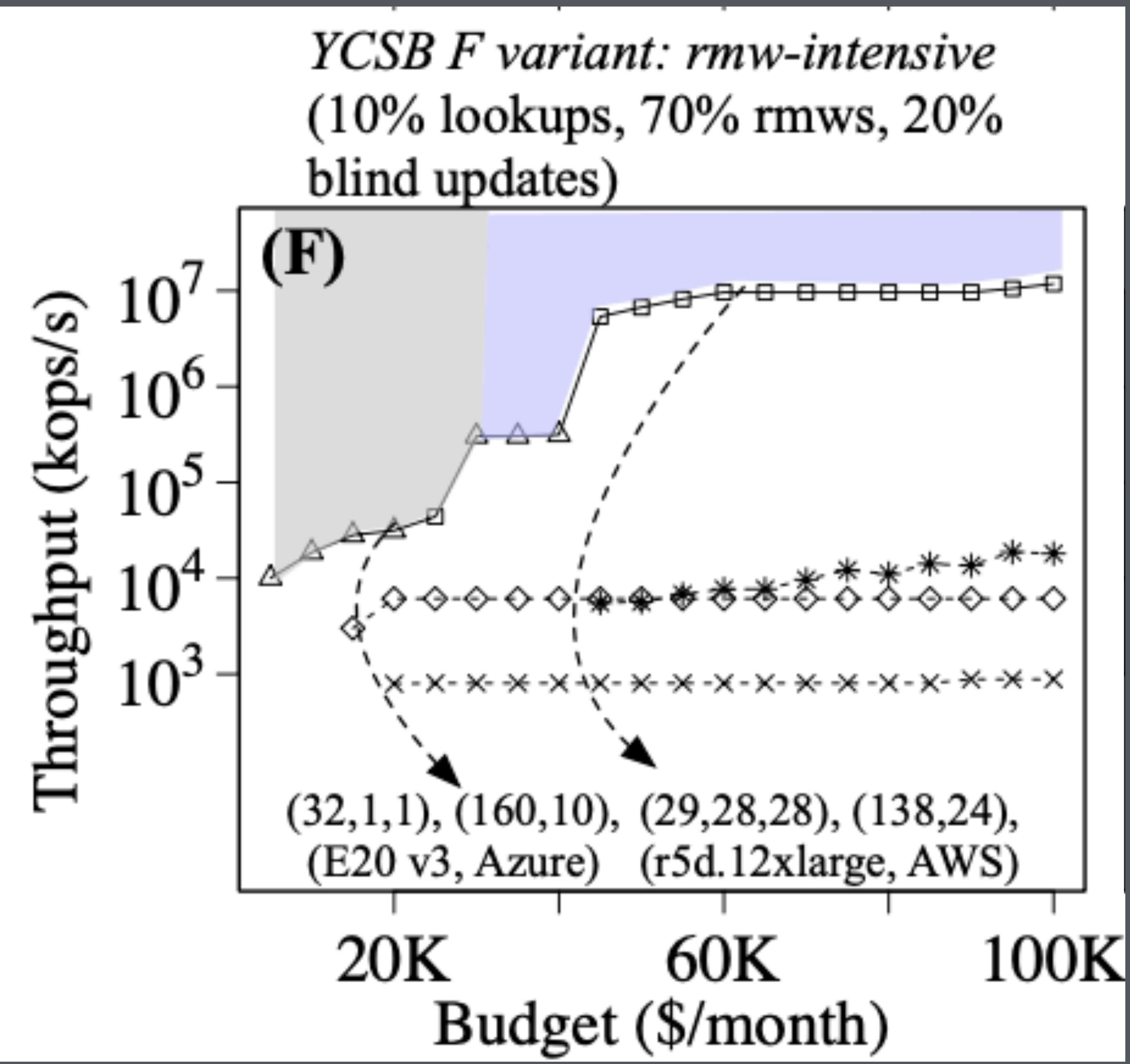
workload/budget diversity

Can we beat the best system for every workload?

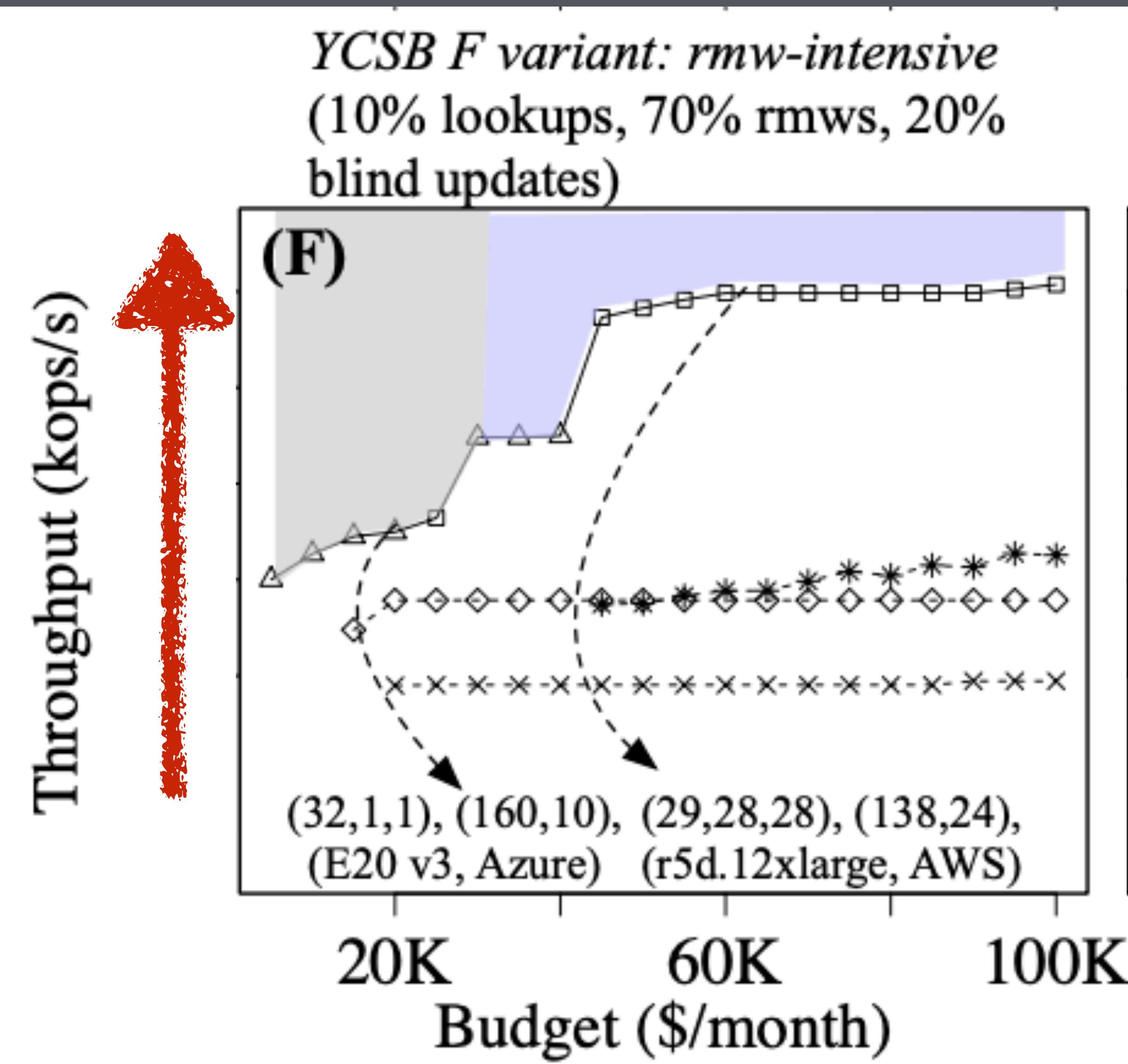




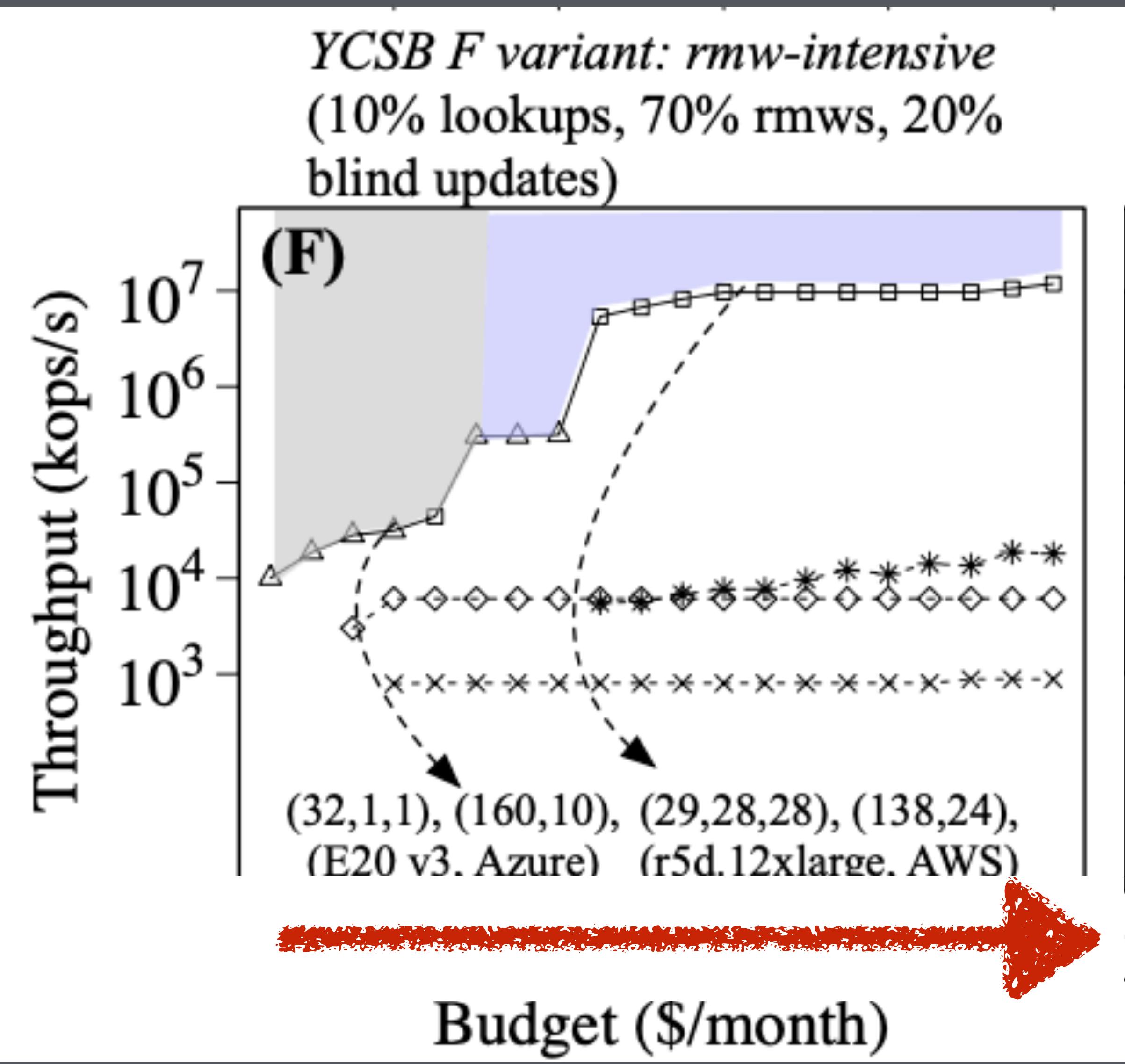
workload/budget diversity



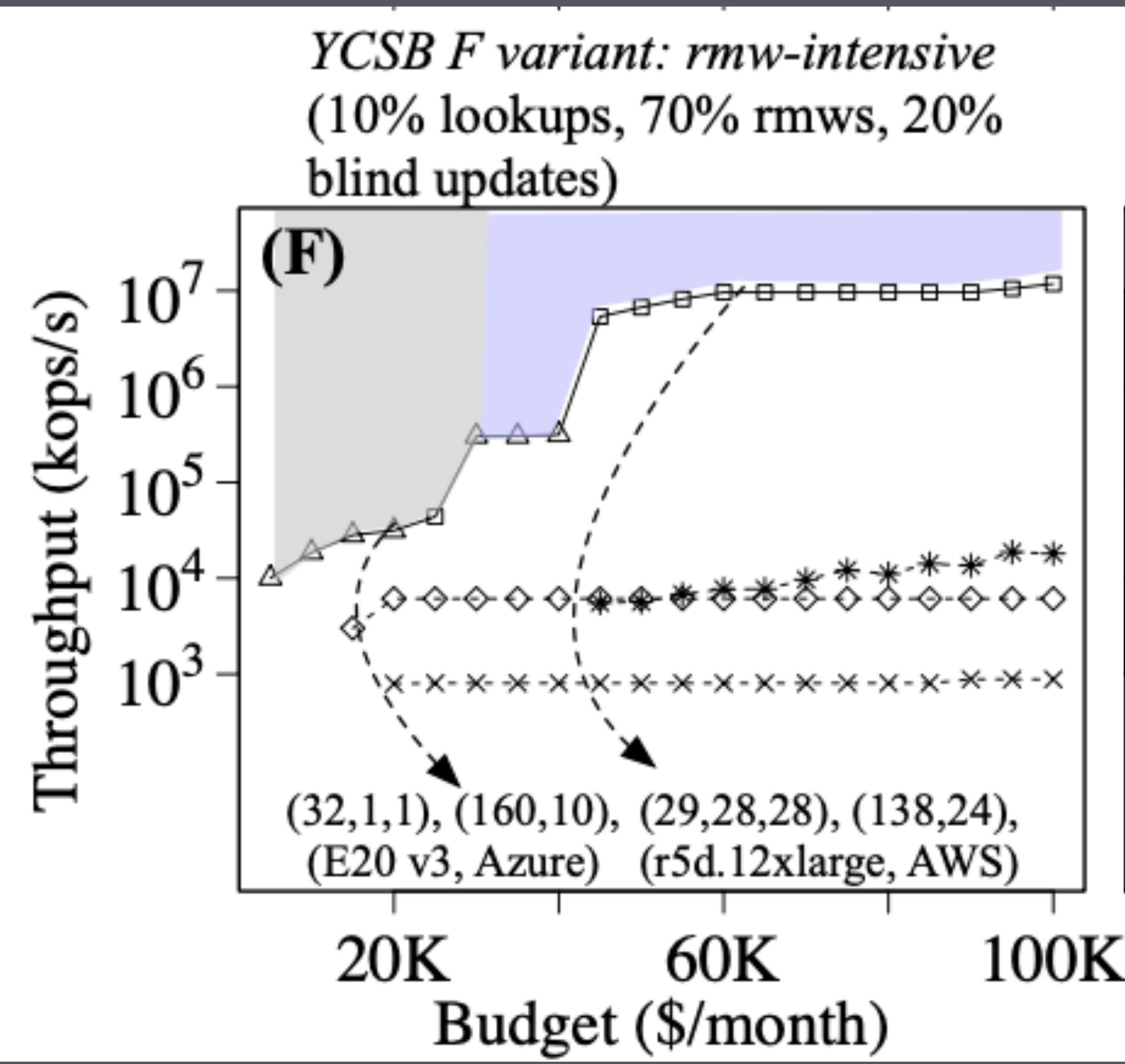
workload/budget diversity



workload/budget diversity

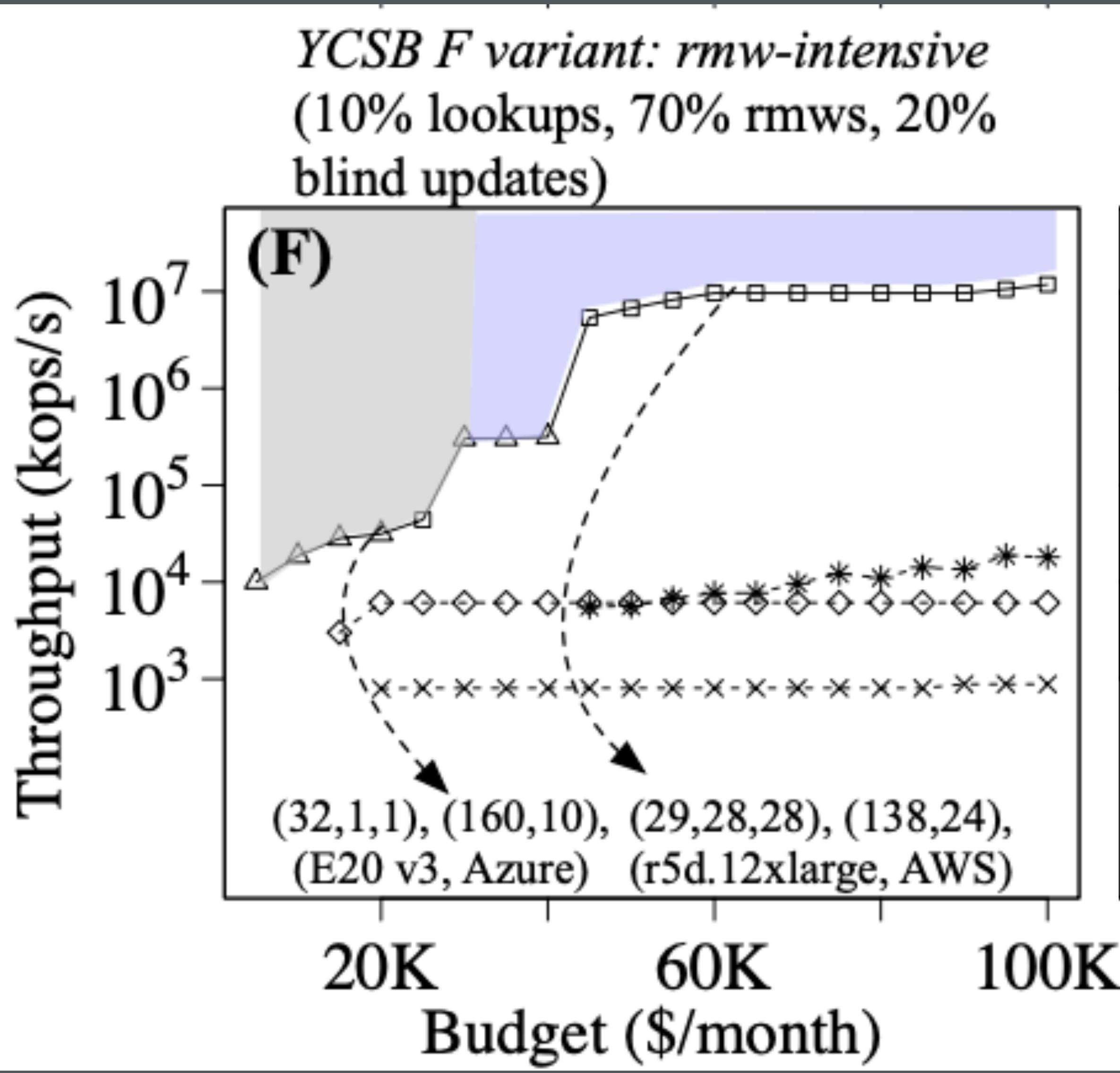


workload/budget diversity



• COSINE •

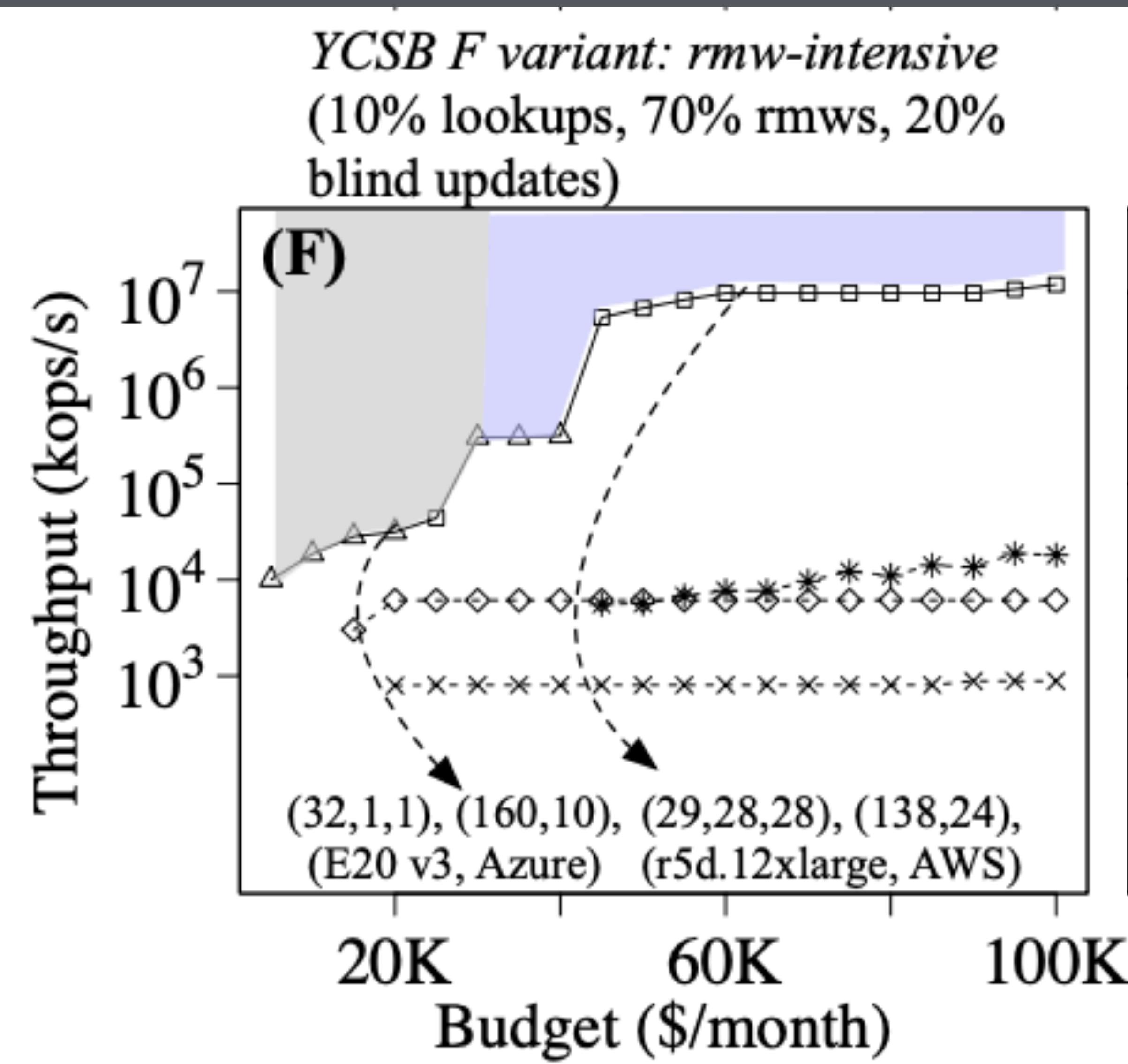
workload/budget diversity



⊕ COSINE ⊕

{ State of the Art
Meta, Microsoft, Mongo

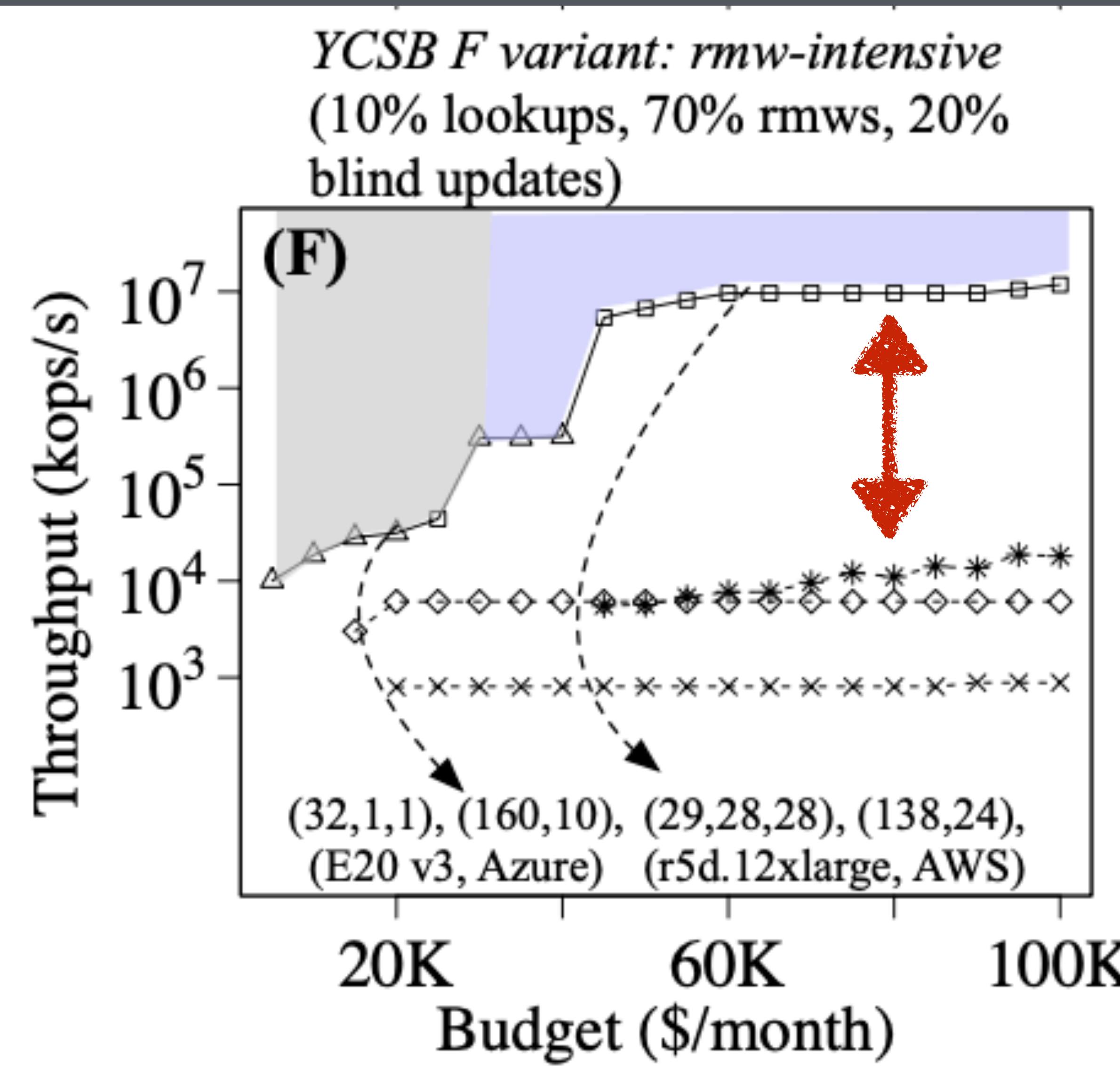
workload/budget diversity



• COSINE •

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity

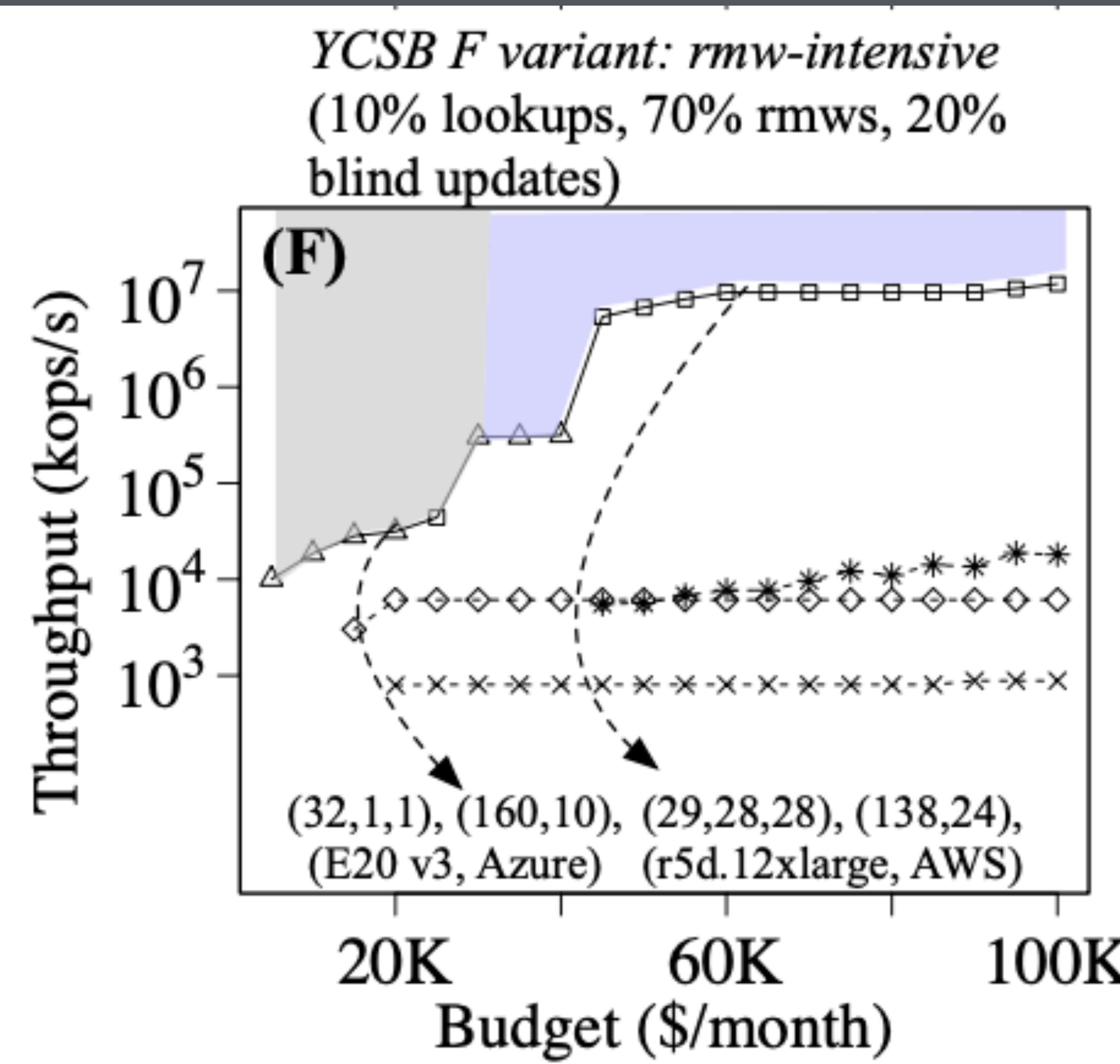


• COSINE •

Better throughput/cost

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity

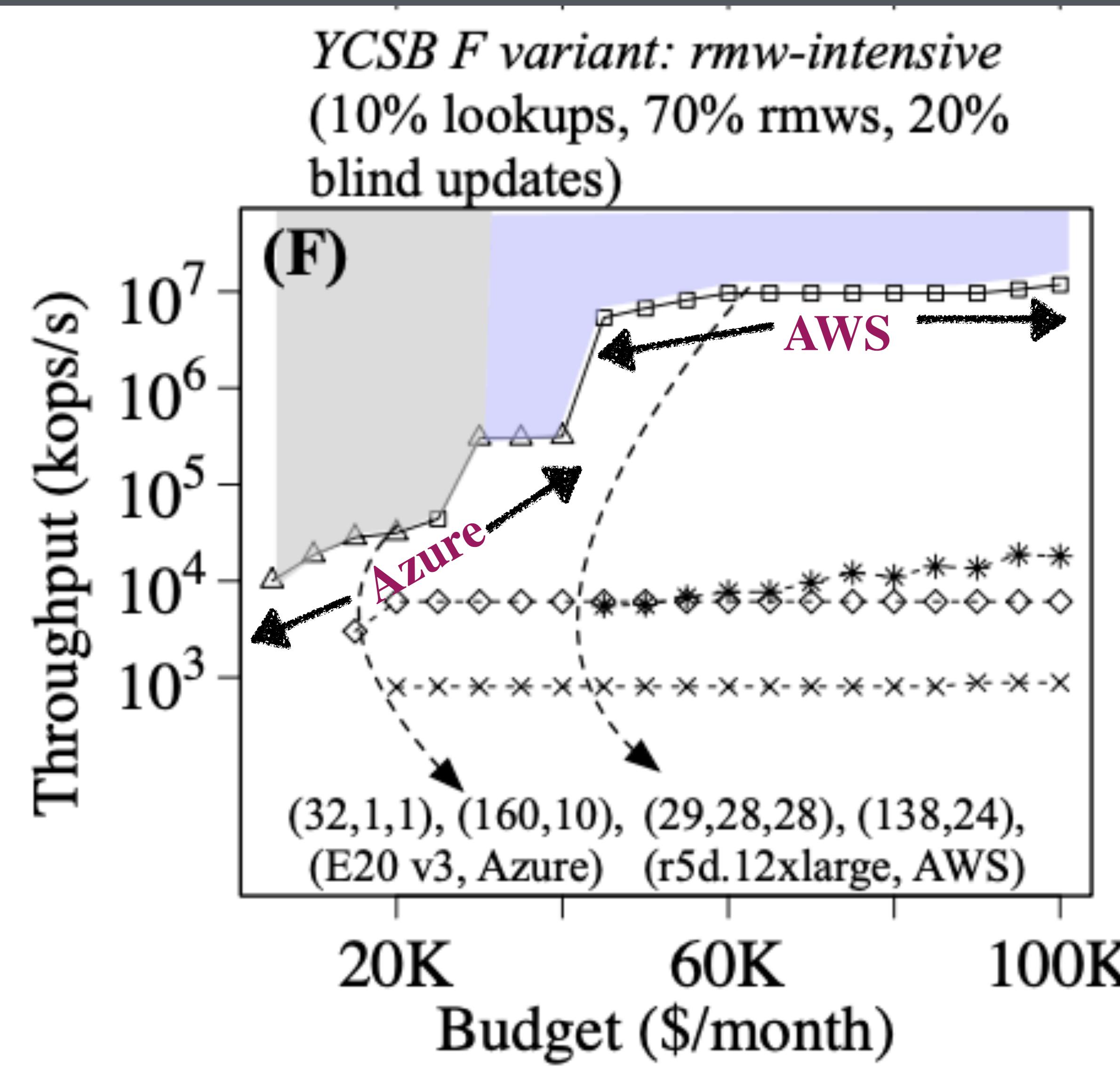


♣ COSINE ♣

Better throughput/cost
Self-designs

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity

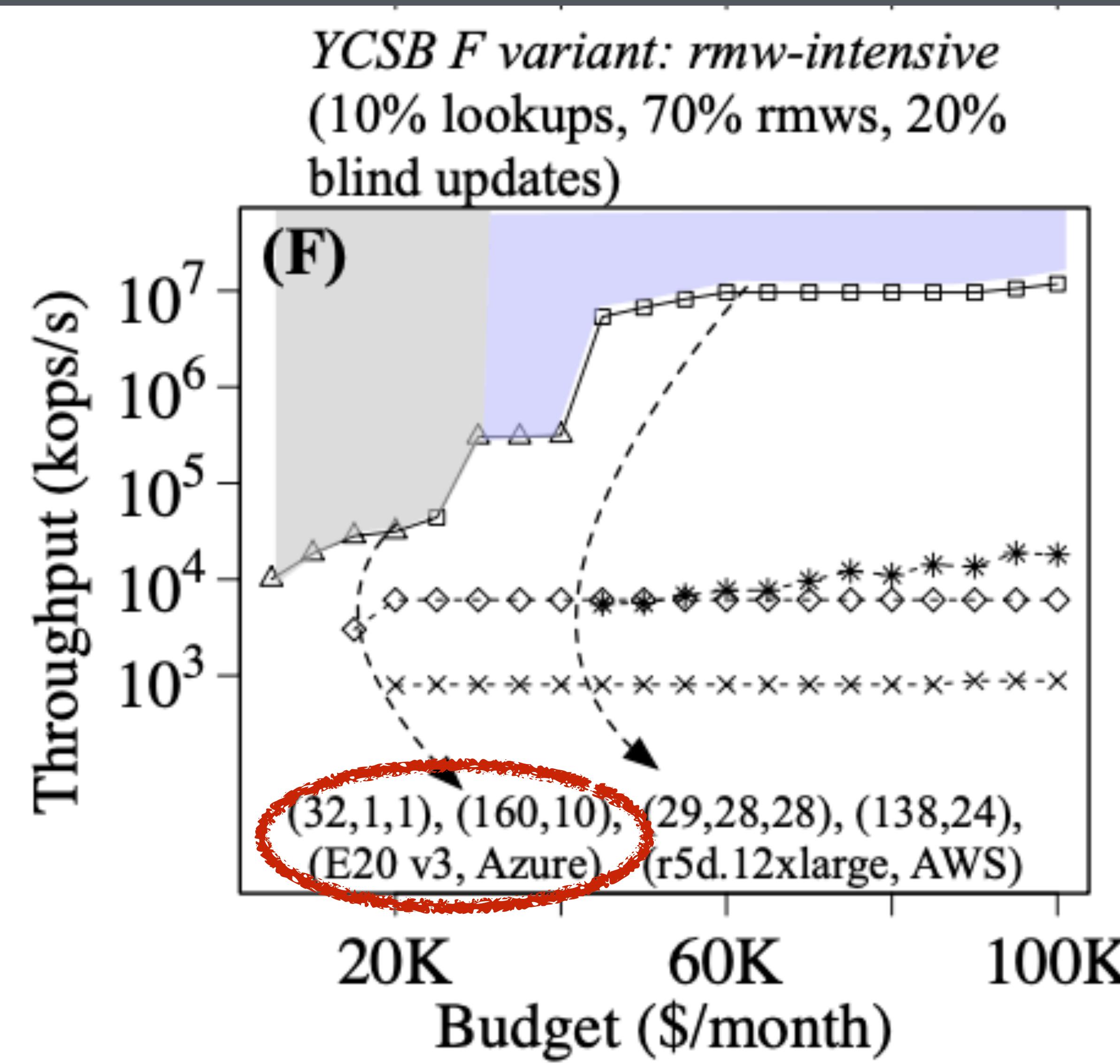


• COSINE •

Better throughput/cost
Self-designs

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity



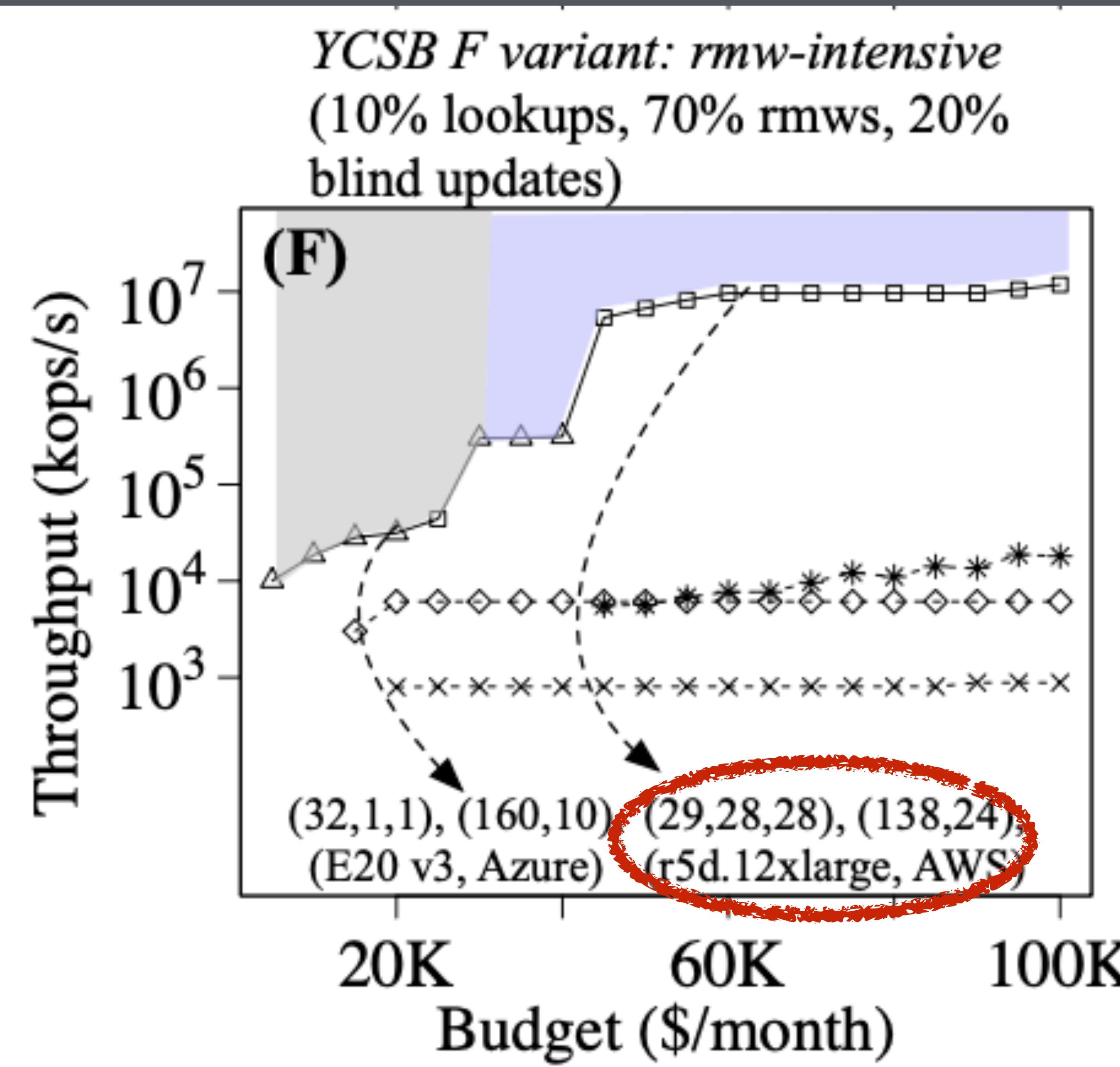
♣ COSINE ♣

Better throughput/cost

Self-designs

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity

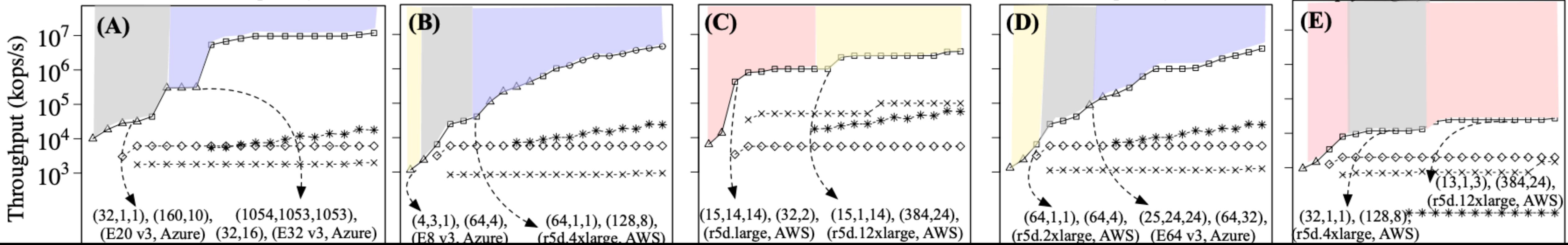
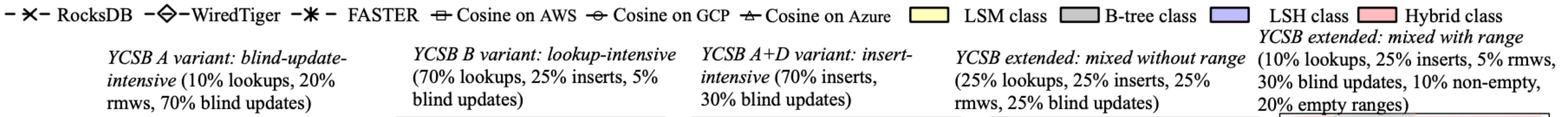


COSINE

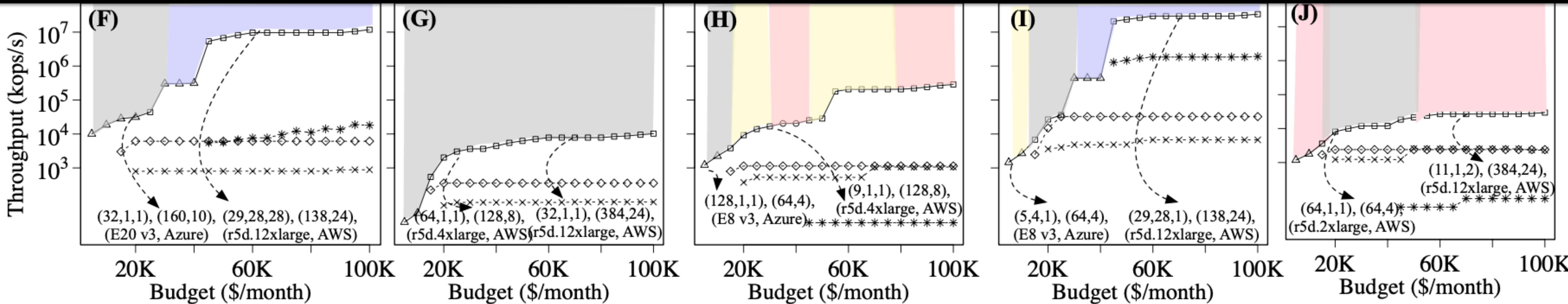
Better throughput/cost
Self-designs

{ State of the Art
Meta, Microsoft, Mongo

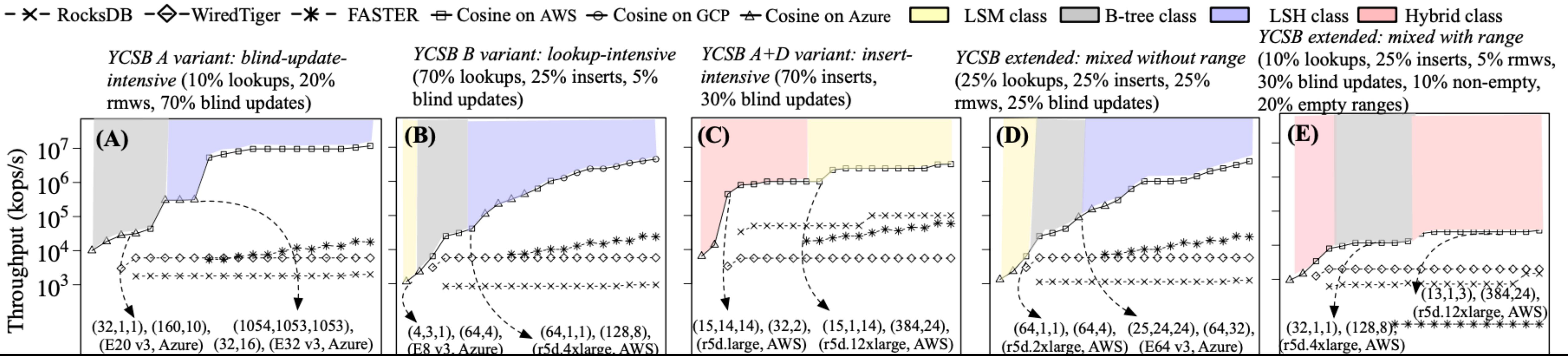
workload/budget diversity



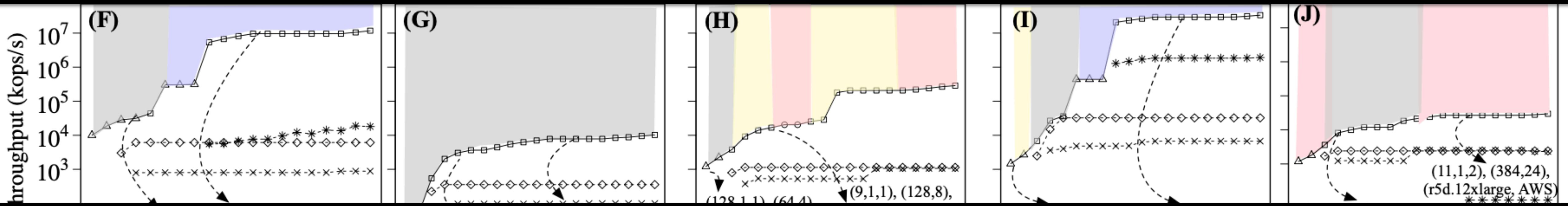
diversity beats top systems self-designs (provider, VM, new design)



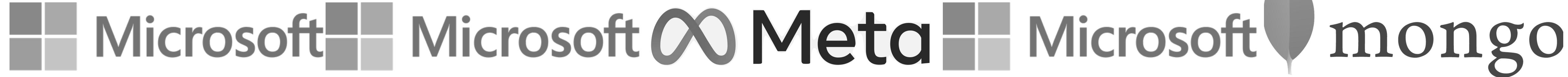
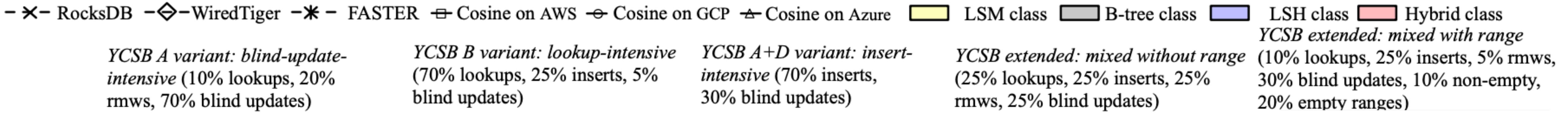
workload/budget diversity



diversity beats top systems self-designs (provider, VM, new design)



Cosine achieves the best perf. across all workloads by automatically designing a new system every time.



diversity beats top systems self-designs (provider, VM, new design)



Cosine achieves the best perf. across all workloads by automatically designing a new system every time.

We can automatically design 1000x faster new NoSQL systems

- 1) design space
- 2) navigation (math/ML)
- 3) code generation

Papers: **Cosine** PVLDB 2023, and new **Limousine** at SIGMOD 2024

We can automatically design 1000x faster new NoSQL systems

- 1) design space
- 2) navigation (math/ML)
- 3) code generation

Papers: **Cosine** PVLDB 2023, and new **Limousine** at SIGMOD 2024

How do these concepts translate to the other big data areas

neural networks, image AI, Blockchain, ...?

We can automatically design 1000x faster new NoSQL systems

- 1) design space
- 2) navigation (math/ML)
- 3) code generation

Papers: **Cosine** PVLDB 2023, and new **Limousine** at SIGMOD 2024

How do these concepts translate to the other big data areas

neural networks, image AI, Blockchain, ...?

again, it all starts from the storage design space



Limousine: Blending Learned and Classical Indexes to Self-Design Larger-than-Memory Cloud Storage Engines.

Subarna Chatterjee, Mark F. Pekala, Lev Kruglyak, and Stratos Idreos.
In Proceedings of the ACM Management of Data 2, 1, Article 47 (February 2024), (SIGMOD), 2024

Cosine: A Cloud-Cost Optimized Self-Designing Key-Value Storage Engine.

Subarna Chatterjee, Meena Jagadeesan, Wilson Qin, and Stratos Idreos.
In Proceedings of the Very Large Databases Endowment, (VLDB), 2022



The Image Calculator: 10x Faster Image-AI Inference by Replacing JPEG with Self-designing Storage Format.

Utku Sirin and Stratos Idreos.

In Proceedings of the ACM Management of Data 2, 1, Article 52 (February 2024), (SIGMOD), 2024

Project 1:

What are the boundaries between LSM-trees and B-trees?

Can we have a single strictly better design for big data NoSQL?

Project 2:

Can we design new systems using the Cosine/Limousine design space and cost/algo synthesis and an LLM?

os 265

Stratos Idreos

BIG DATA SYSTEMS

NoSQL | Neural Networks | Image AI | LLMs | Data Science