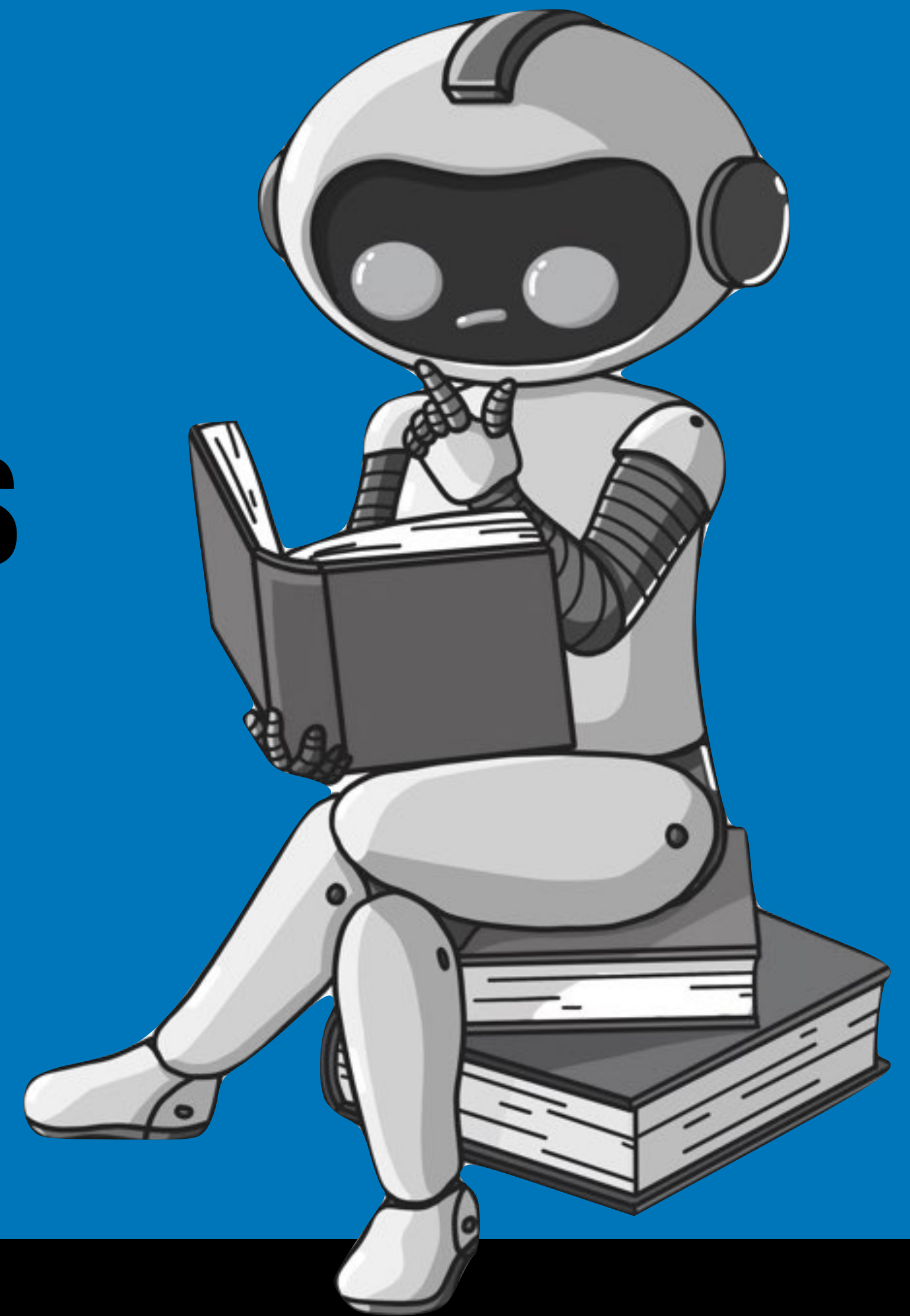


ALPHABETS, GRAMMARS, CALCULATORS, AND THE END OF HAND-CRAFTED SYSTEMS

Stratos
Idreos



DASlab
@ Harvard SEAS

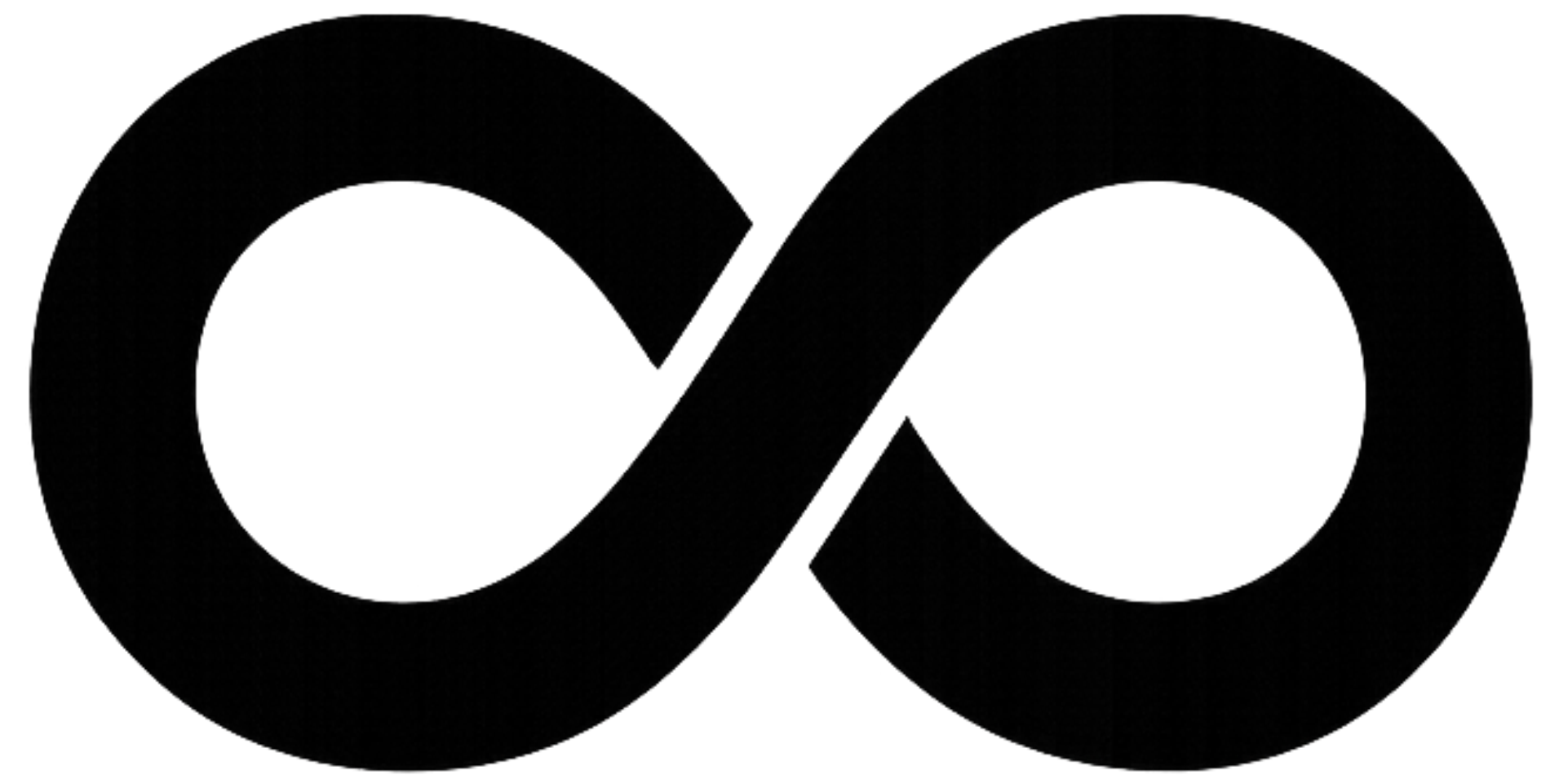


What if we can reason about systems design?

IS THIS GOING TO
REPLACE RESEARCHERS?

NOPE!

Infinite
Research
Topics



POSSIBLE
SYSTEM
DESIGNS

WHAT WE HAVE
INVENTED SINCE
THE DAWN OF
COMPUTER
SCIENCE

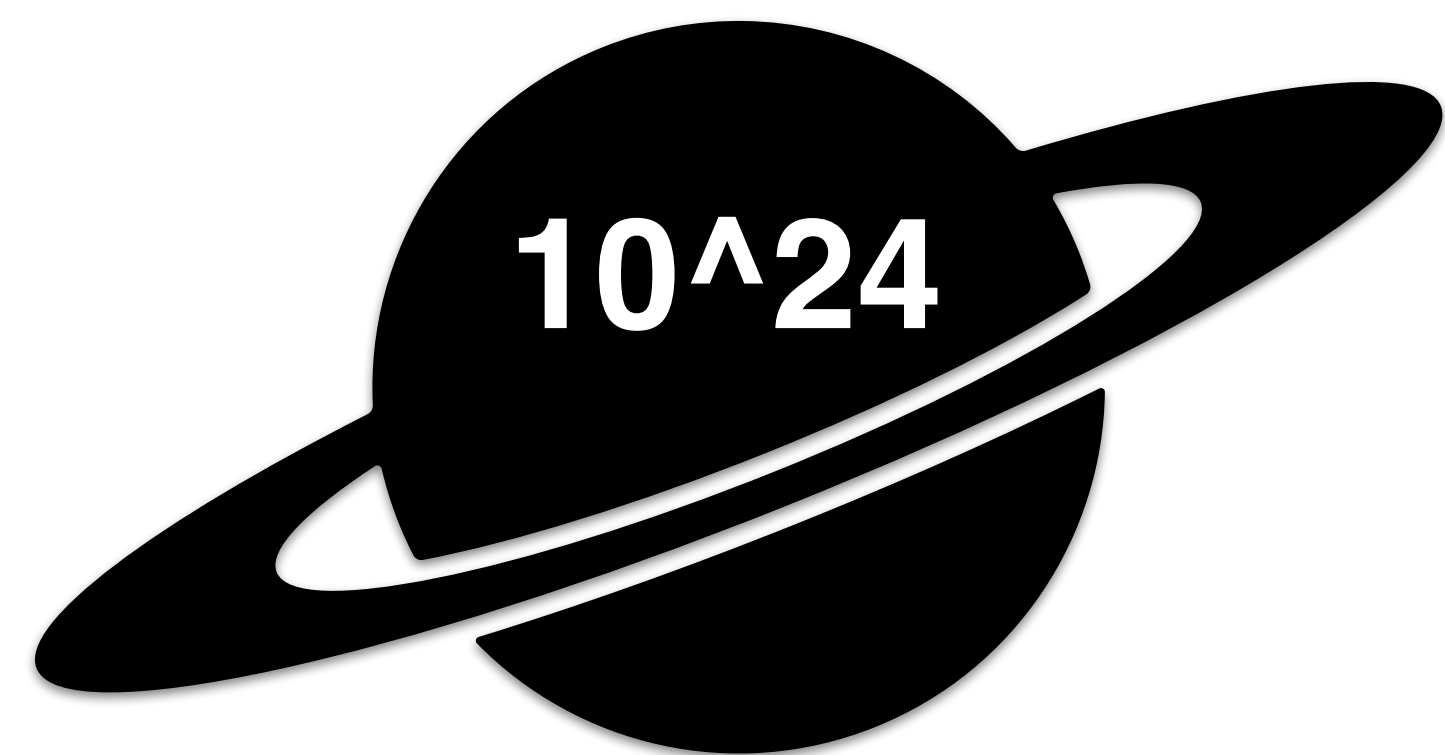


Possible
Data
Structures = **10^{48}**

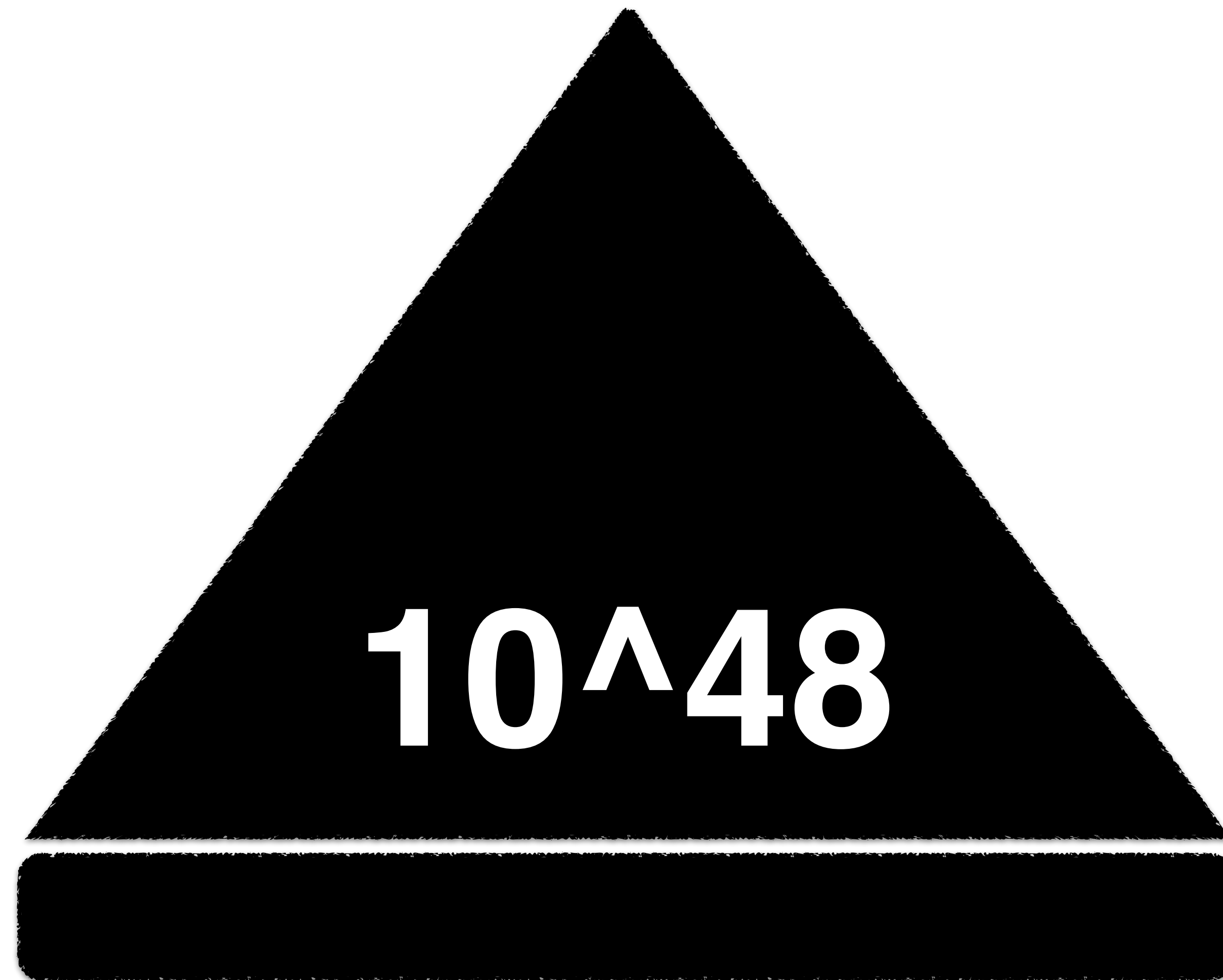
POSSIBLE
SYSTEM
DESIGNS



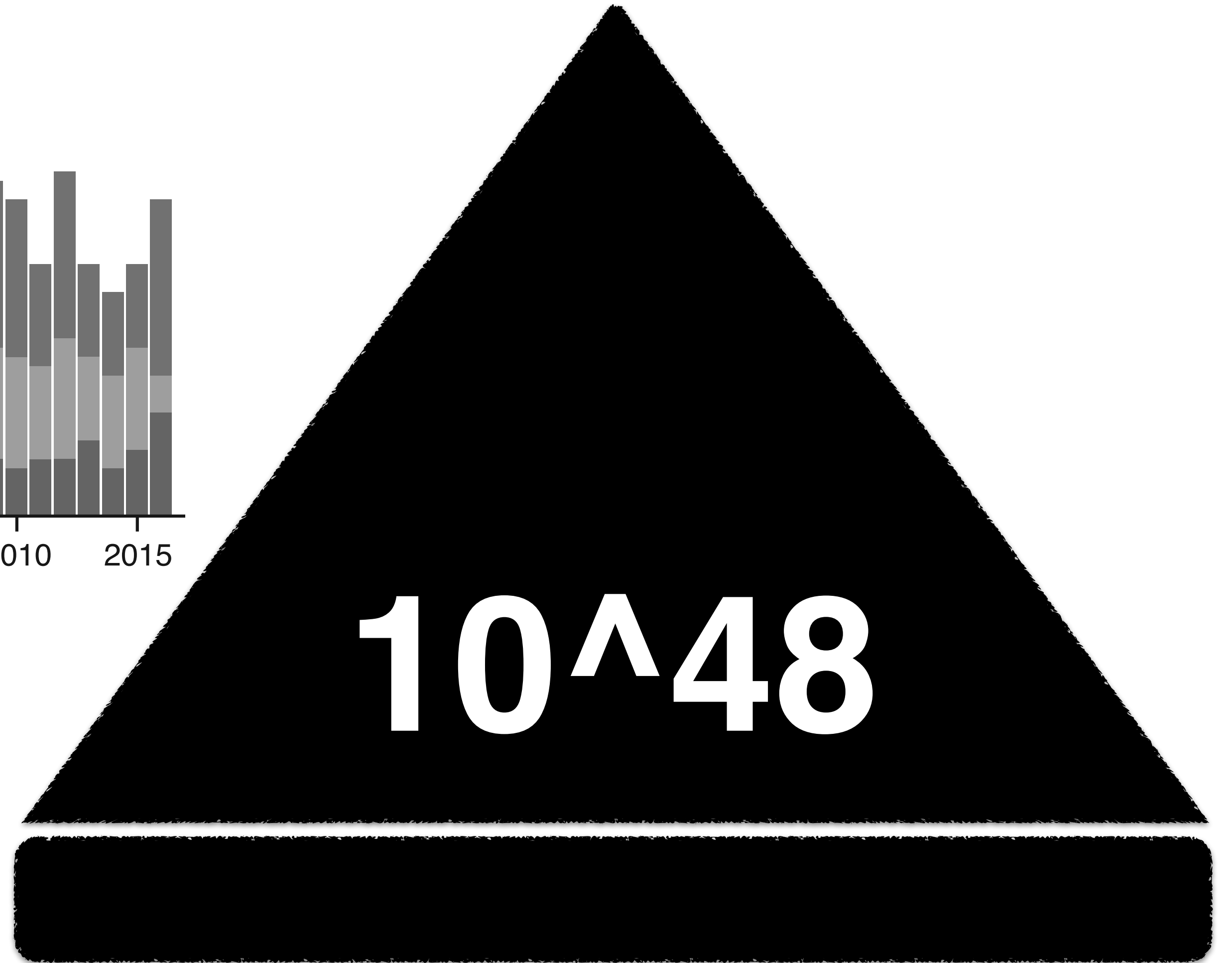
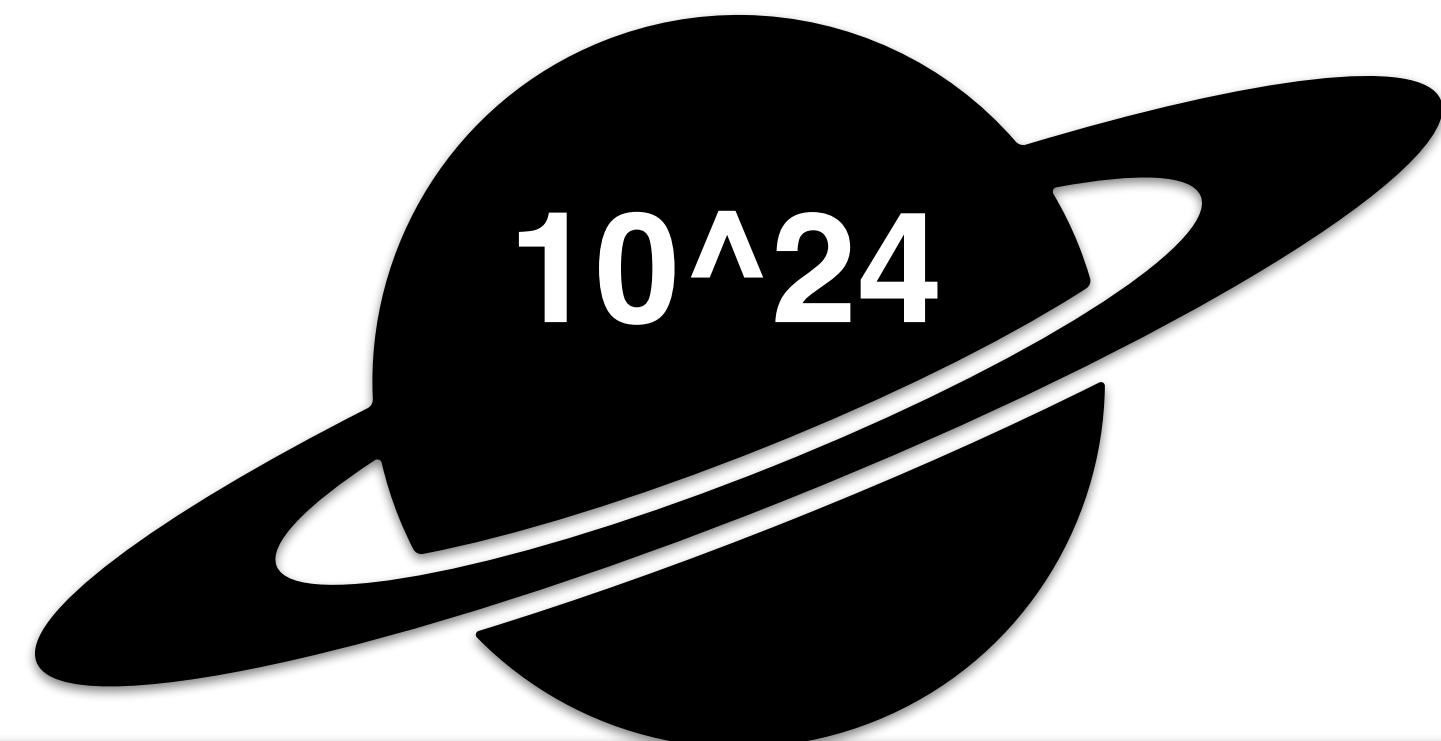
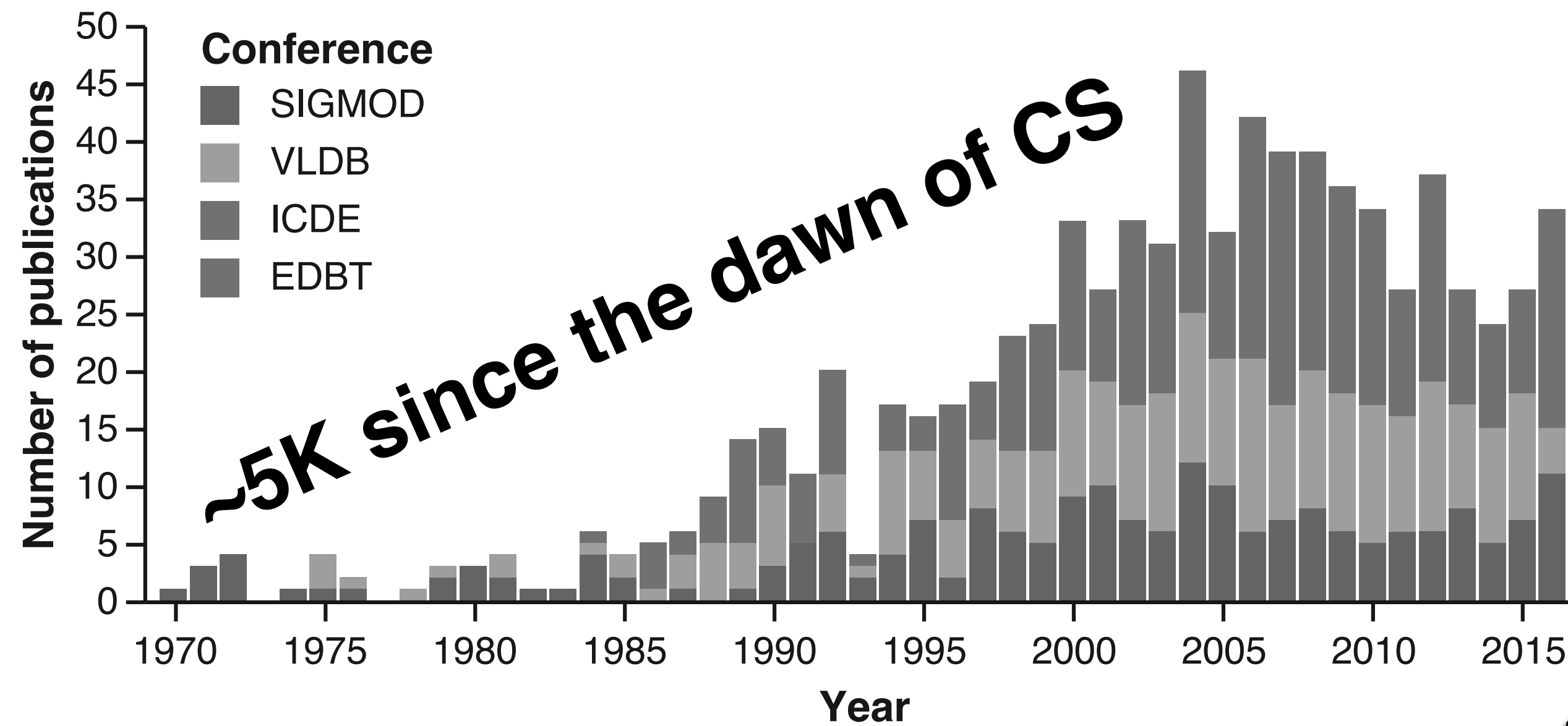
WHAT WE HAVE
INVENTED SINCE
THE DAWN OF
COMPUTER
SCIENCE



STARS ON THE SKY



POSSIBLE DATA STRUCTURES



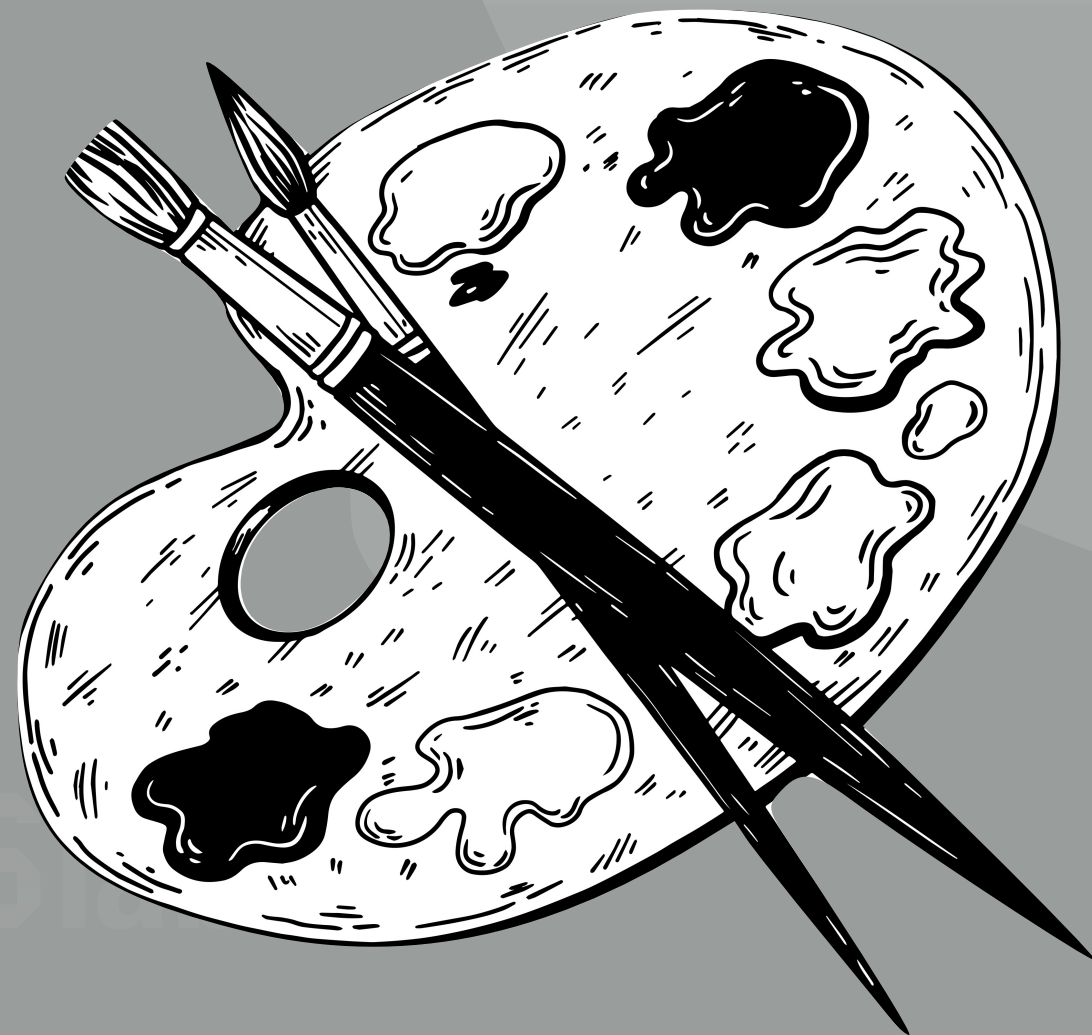
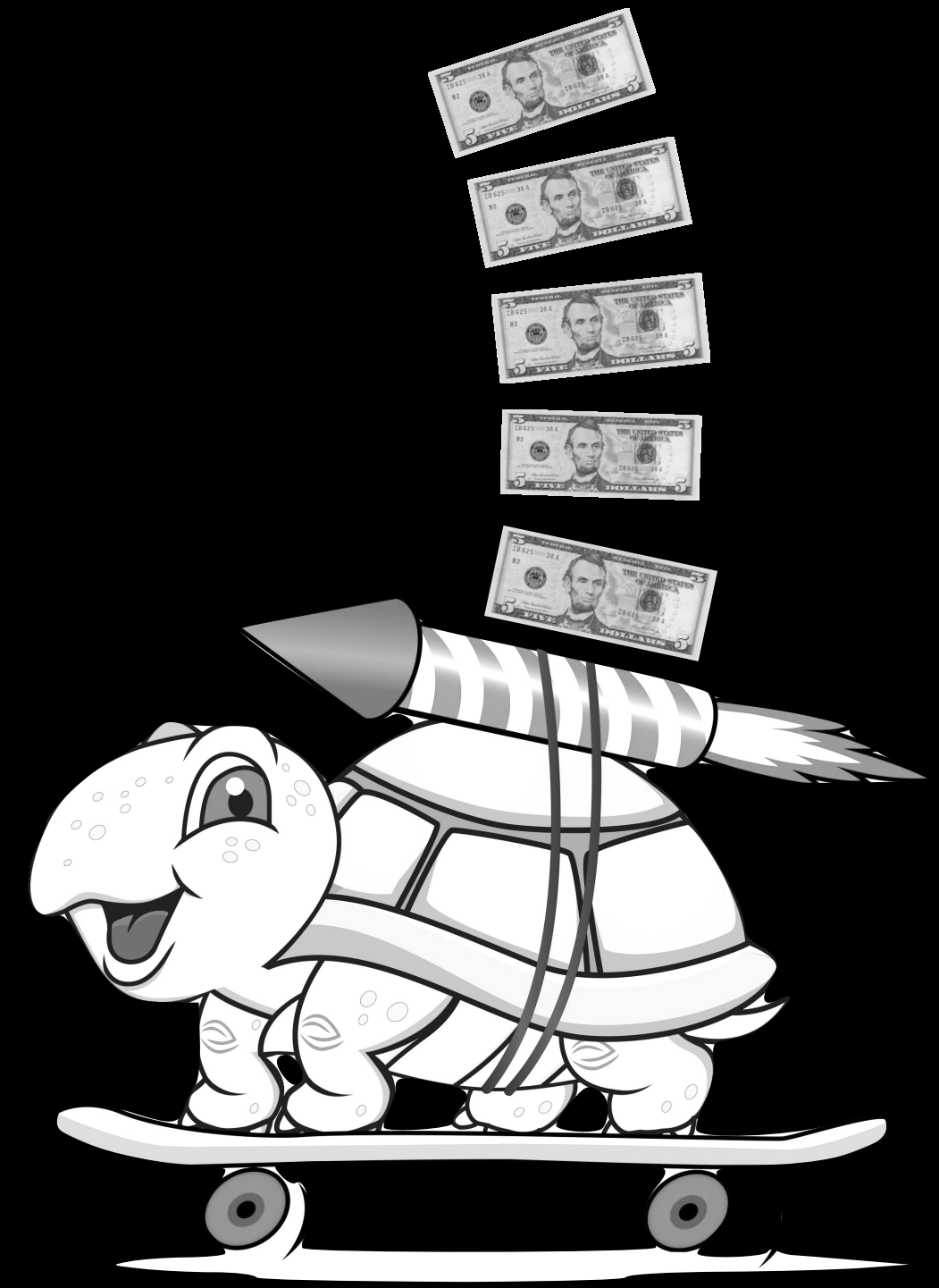
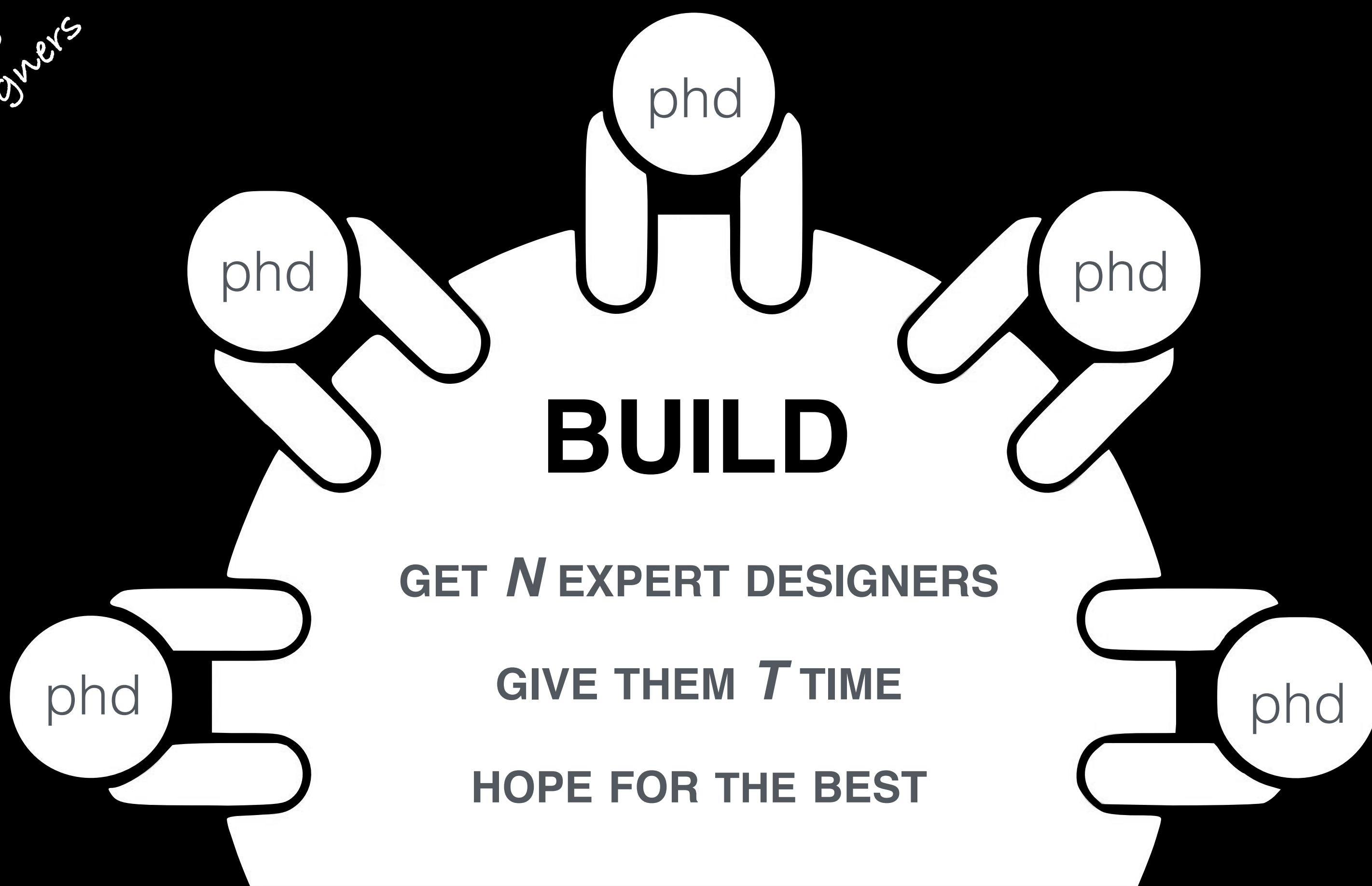
$$10^{48} - 5 \times 10^3 = 10^{48} \quad \text{zero progress}$$



Key-value Stores/Databases/ML Systems

$>10^{100}$

the dining
systems designers



design is an art

Humans do not lack
imagination ...
but manual design
exploration is too slow



**EVERY PAPER
(1 DESIGN)
IS 1-2 YEARS**

**EVERY FULL
SYSTEM IS
7-10 YEARS**

How can we navigate the space of all possible system designs?

Available designs?
Useful designs?
In what context?

...

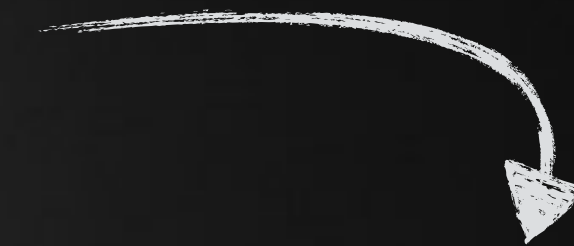


Humans do not lack
imagination ...
but manual design
exploration is too slow

**EVERY PAPER
(1 DESIGN)
IS 1-2 YEARS**

**EVERY FULL
SYSTEM IS
7-10 YEARS**

Input (Requirements):
Data, Queries, H/W,
Cost, Speed, Accuracy, ...



Output:
Tailored System

WHY DO WE NEED NEW SYSTEMS?

But first: *“What is a data system design?”*

A TYPICAL BIG DATA TASK

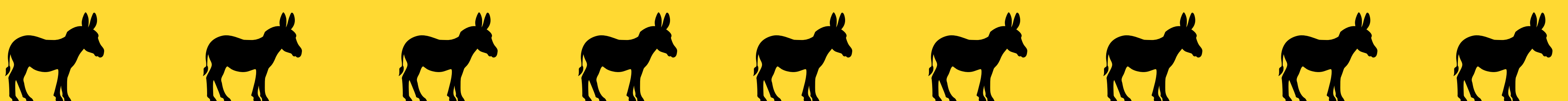
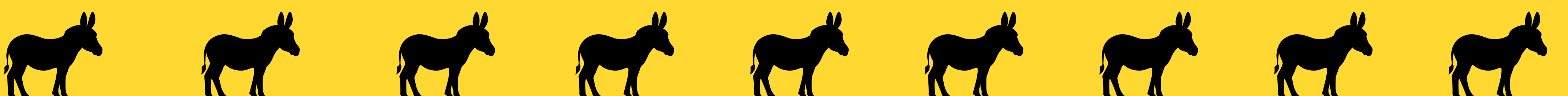
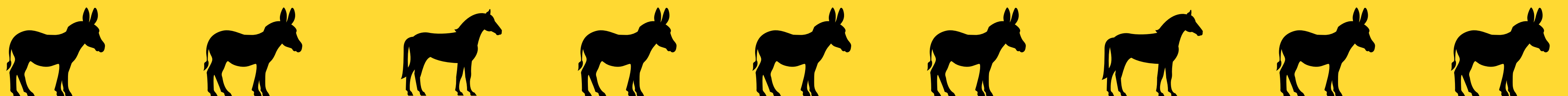
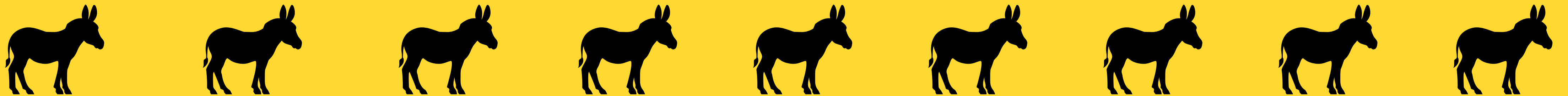
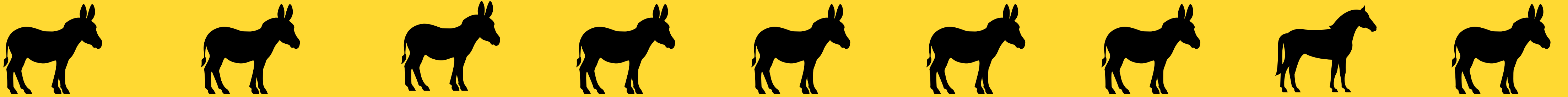
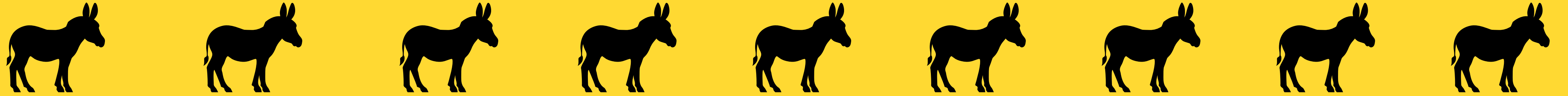
image analysis: e.g., detect the number of horses

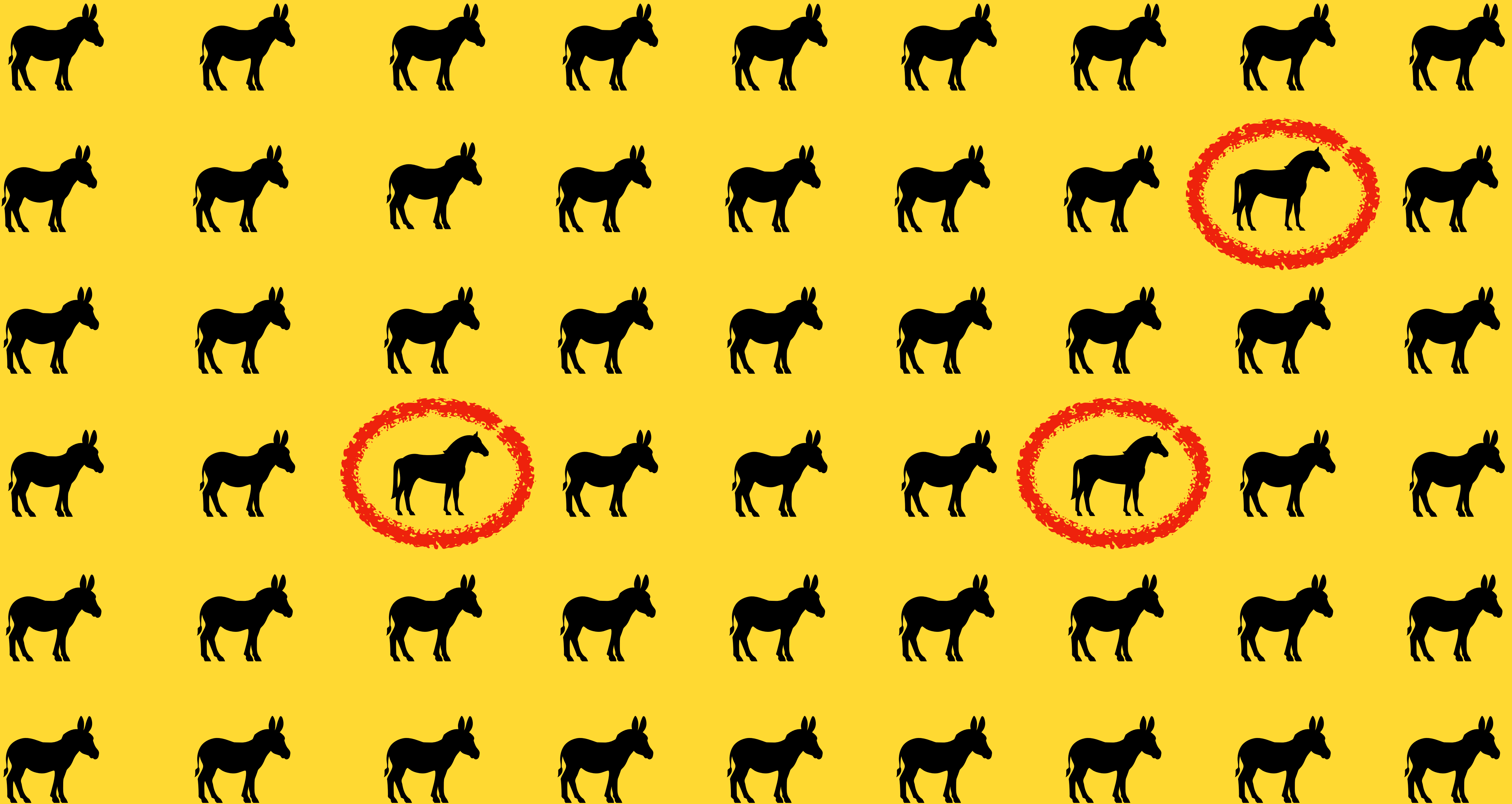


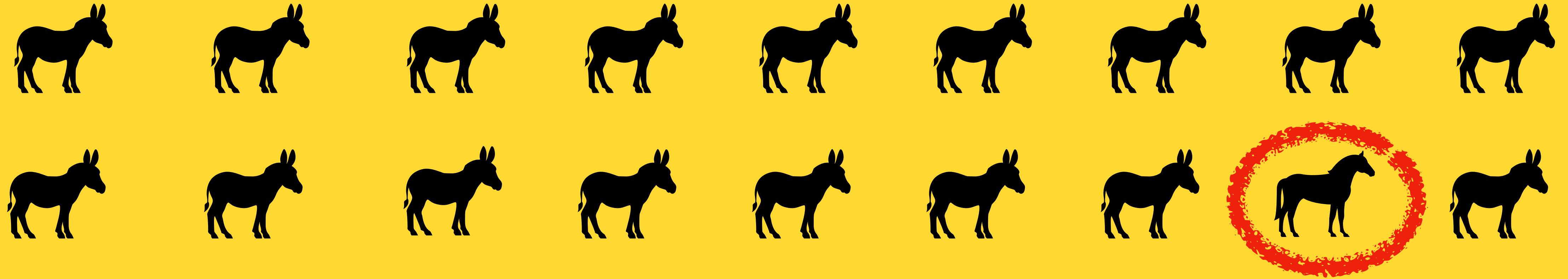
A TYPICAL BIG DATA TASK

image analysis: e.g., detect the number of horses

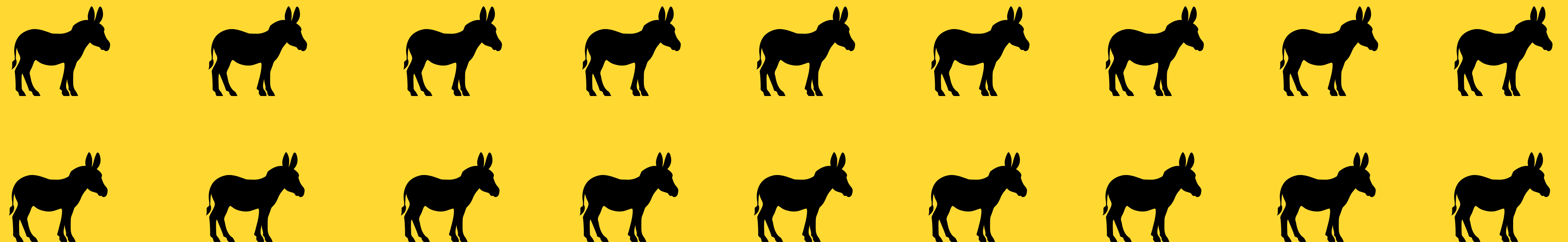




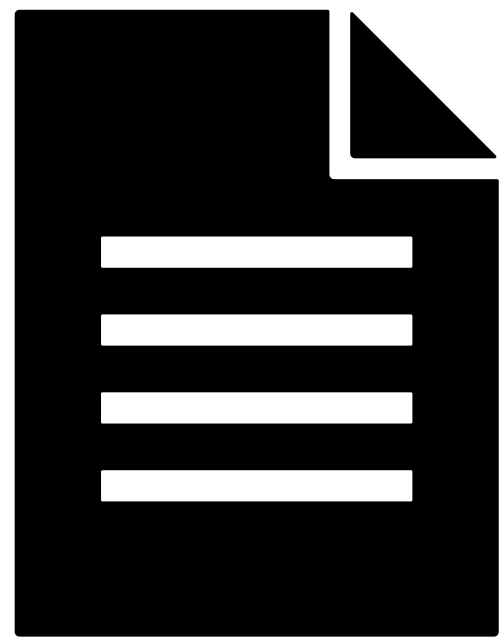




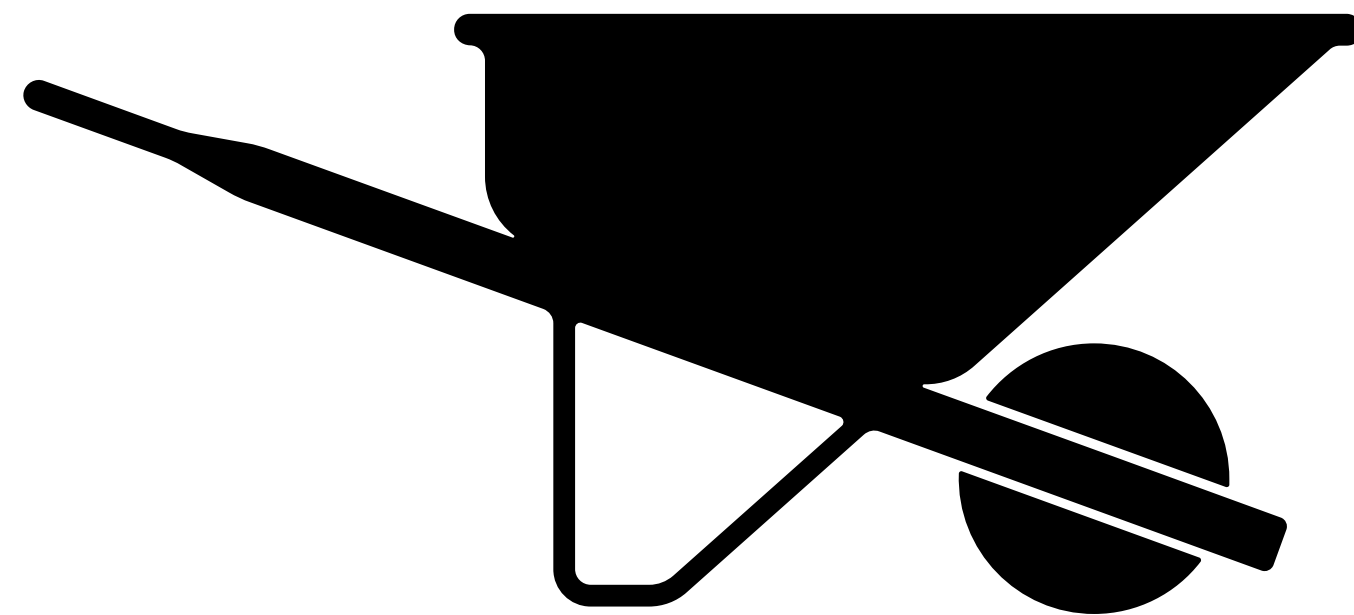
**The size and organization of the data
define how we can process requests**



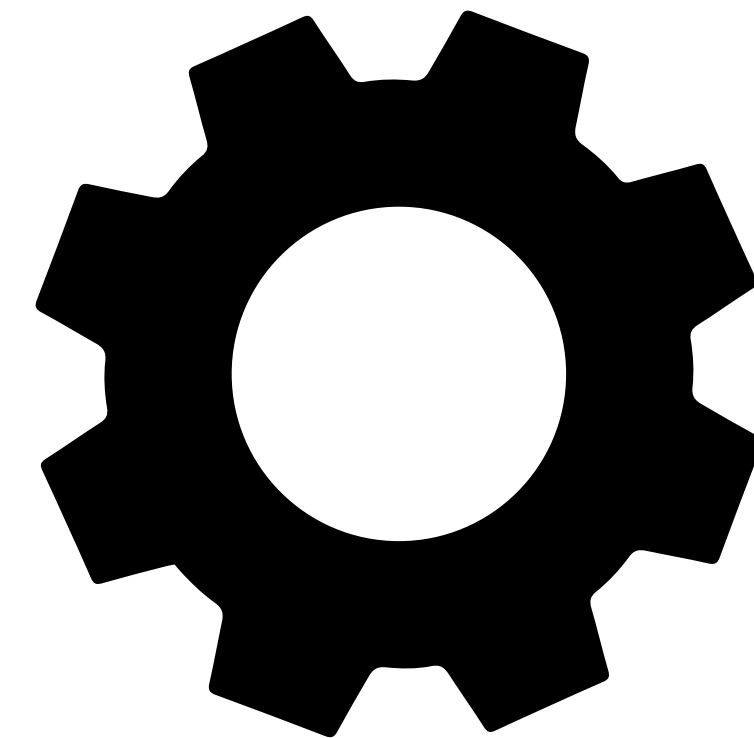
Three components in e2e system design



STORE



MOVE



PROCESS

NEW PERFORMANCE/FEATURES

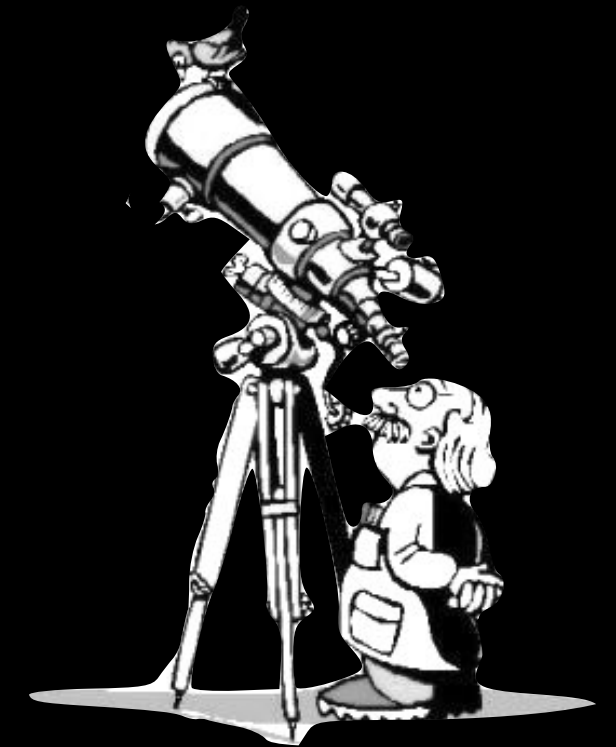
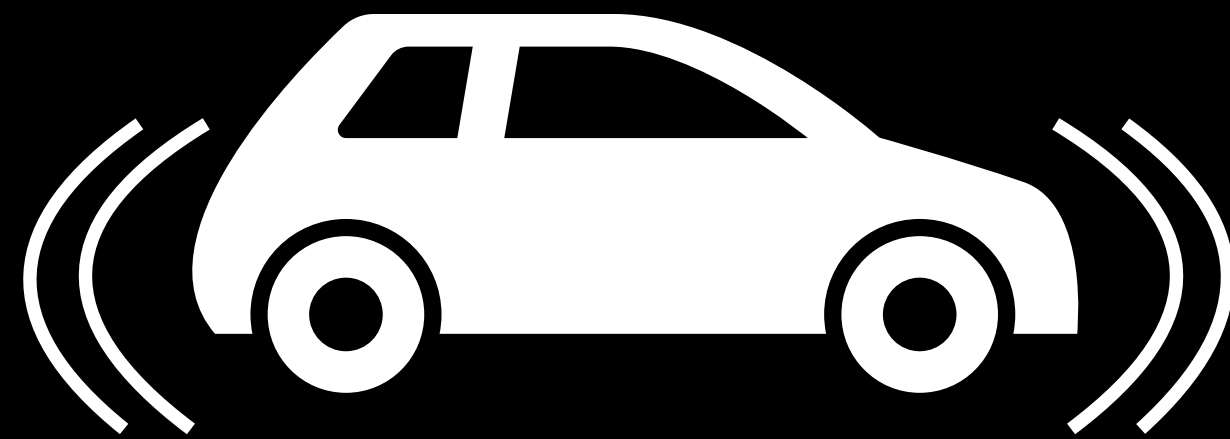
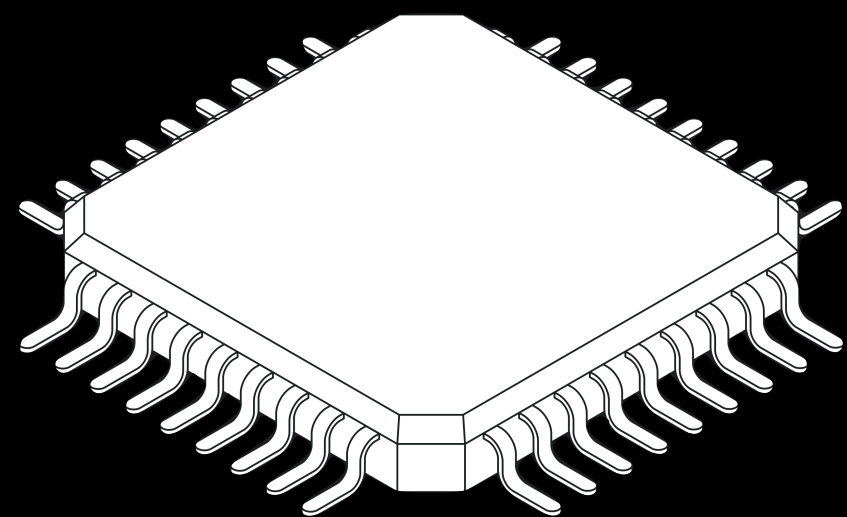
REDESIGN FROM STORAGE



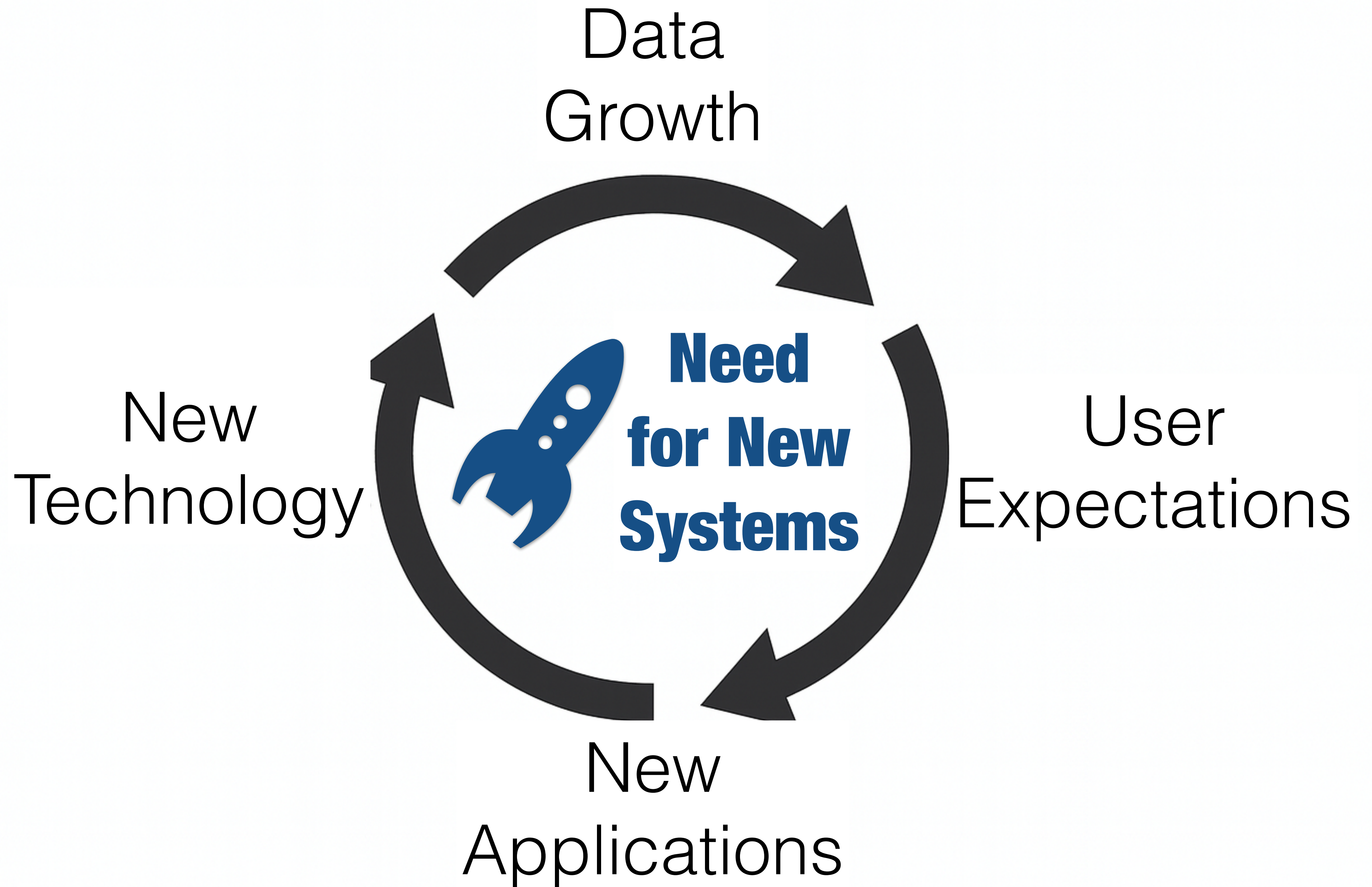
BIG DATA/AI



NEW APPLICATIONS/PERFORMANCE



NEW SYSTEMS



“Getting a new data structure into
production takes years.
And by the time it’s ready,
your assumptions are already wrong.”

Mark Callaghan



“Getting a new data structure into
production takes years.
And by the time it’s ready,
your assumptions are already wrong.”

Mark Callaghan



“We likely have the data to cure cancer.
We just do not know what query to ask.”

Martin Kersten

How do we design a data system that is **X times faster for a workload W?** 

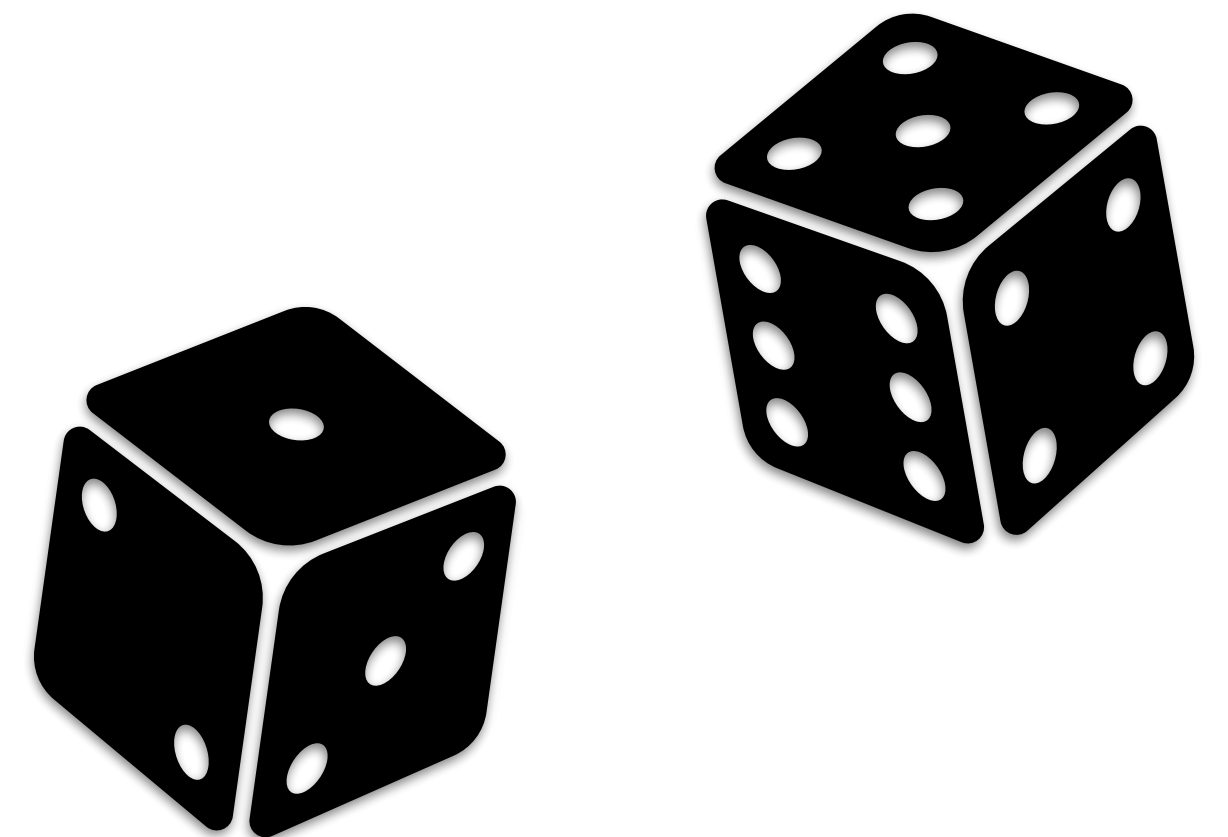


How do we design a data system that allows for control of **cloud cost**?

What happens if we introduce **new application feature Y?**

Should we **upgrade** to new version Z?

What will **break** our system?



Beyond Performance: WHAT IF design



HOW ABOUT ADAPTIVE SYSTEMS?

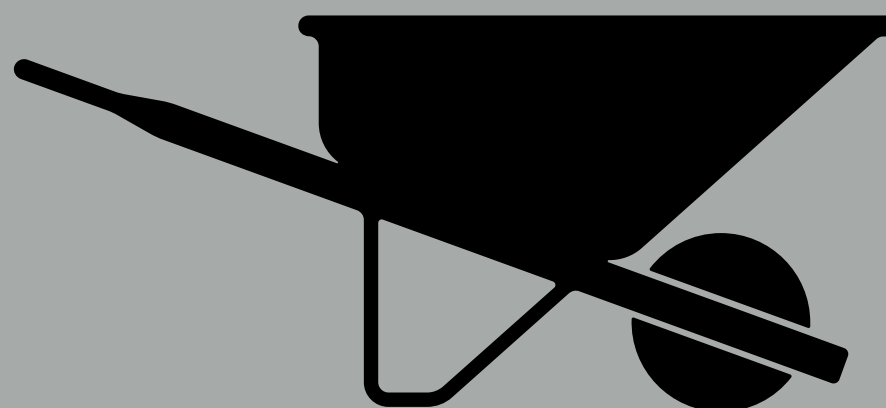
Databases= Adaptive (Optimizers, Auto-tuning)?
Cracking/Adaptive Indexing/Learned Components?...

 **small variations**

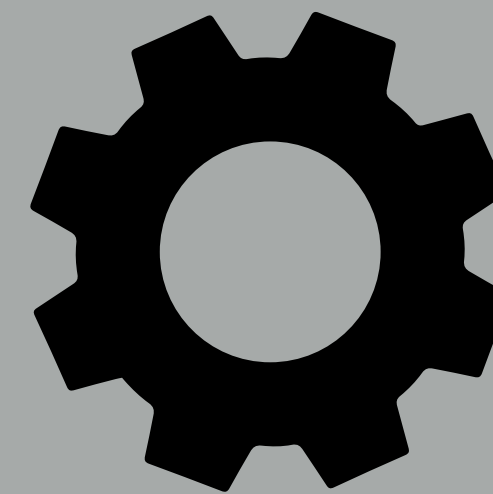
 **small-big variations**



STORE

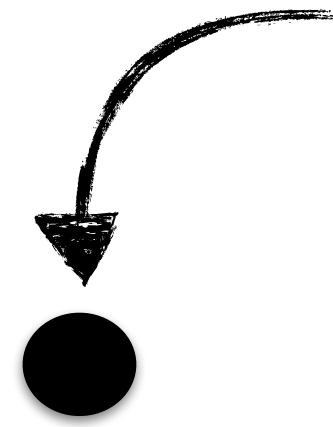


MOVE



PROCESS

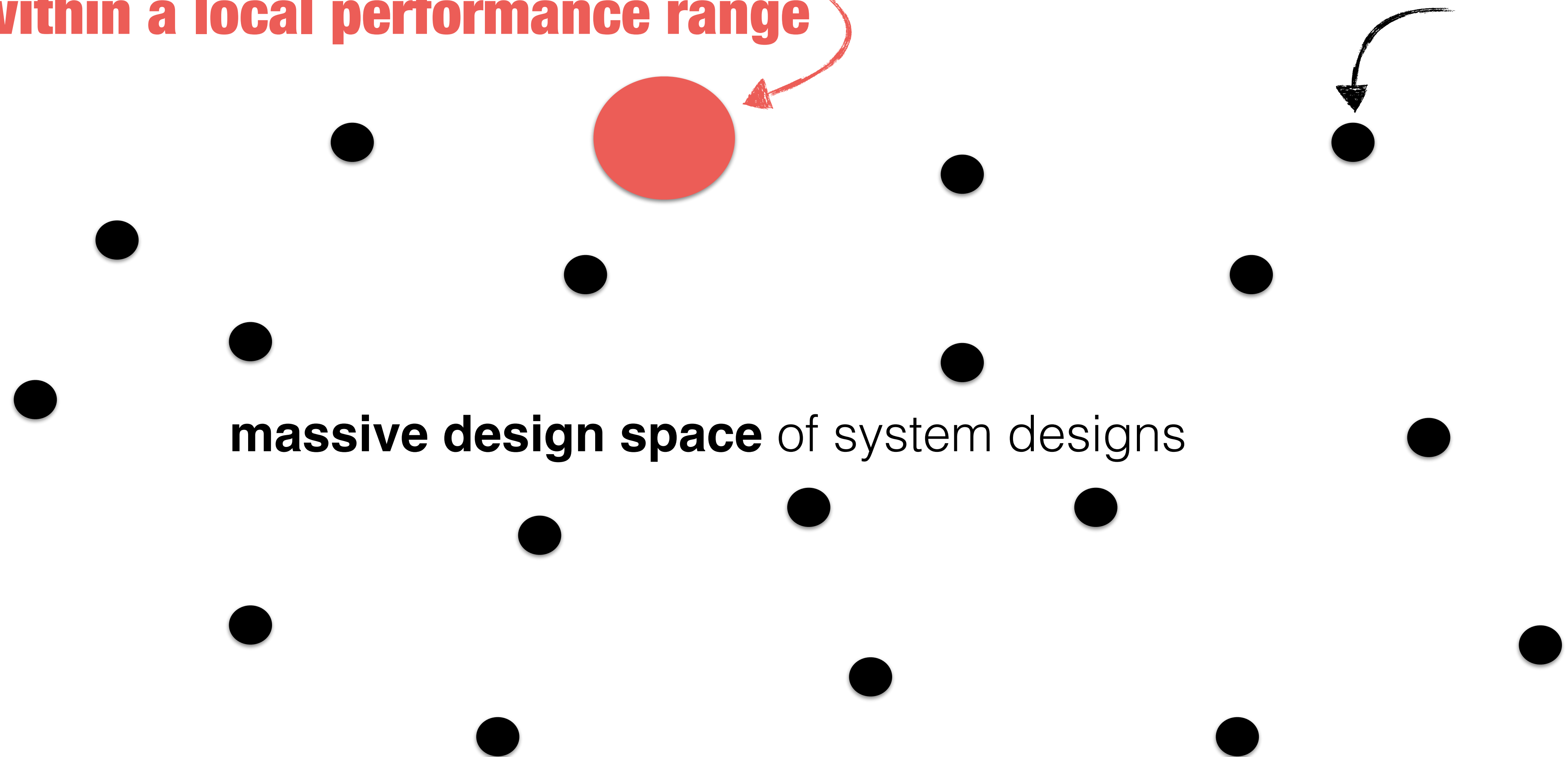
A unique full system design

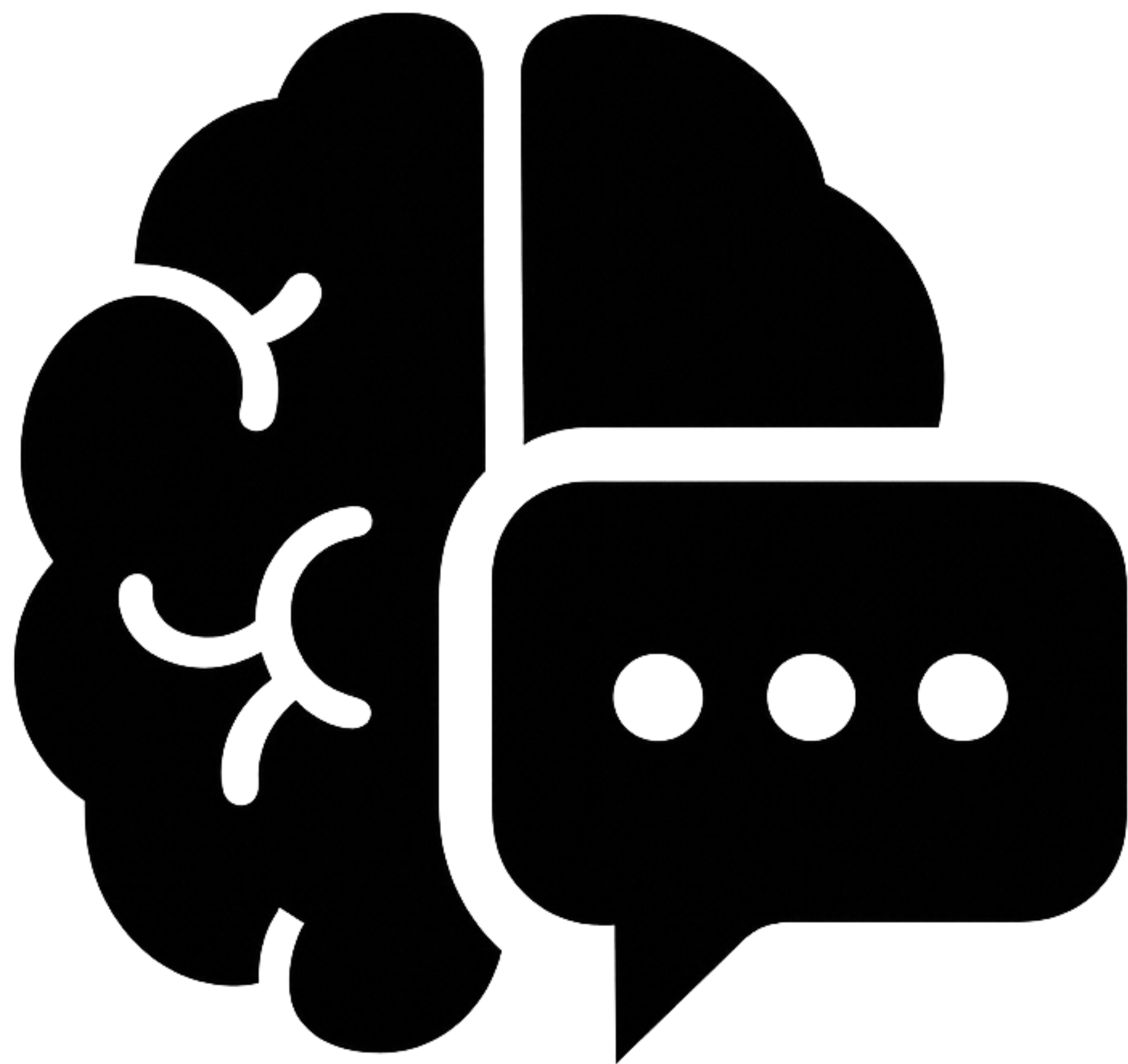


massive design space of system designs

**adaptive designs help fine-tune designs
within a local performance range**

A unique full system design

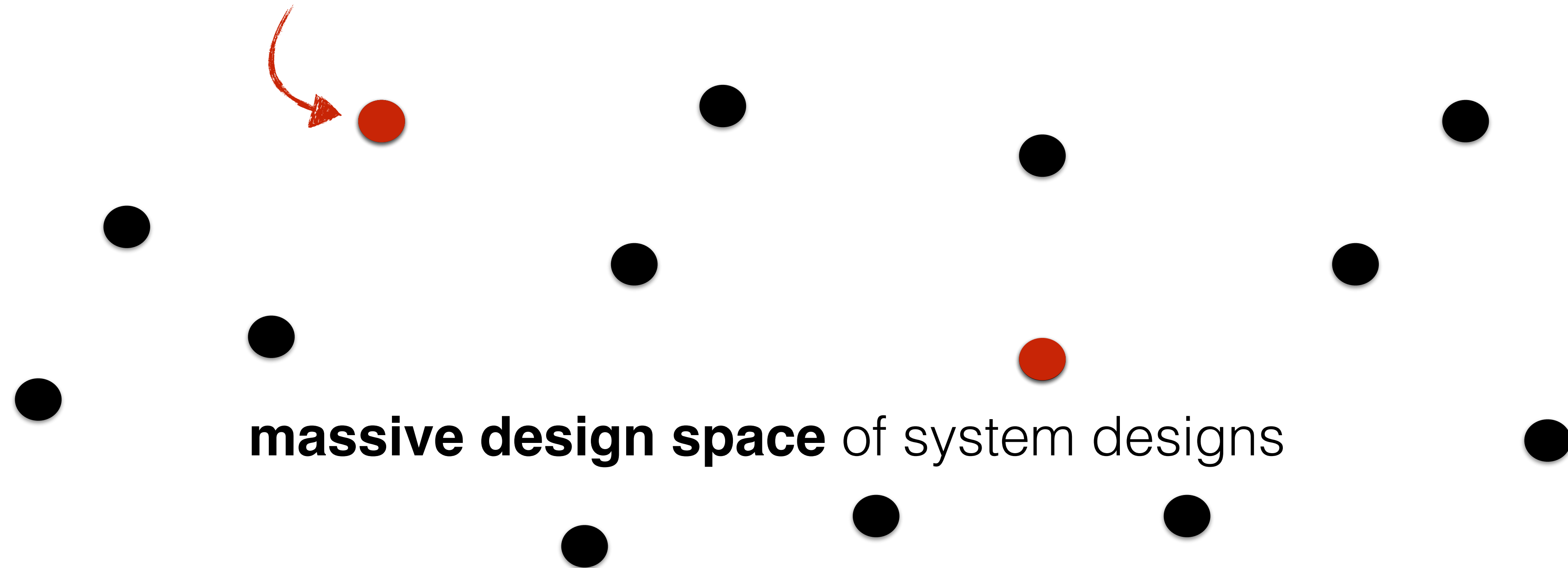




CAN LLMs OR NEW LARGE MODELS MAKE A DIFFERENCE?

Yes, amazing reasoning, but
need tons of data, money, time

few existing designs and mostly complex closed-source code



NO WEBSCALE DATA FOR SYSTEM DESIGN

Adaptive designs & large models are part of the solution, but we also need something else to:

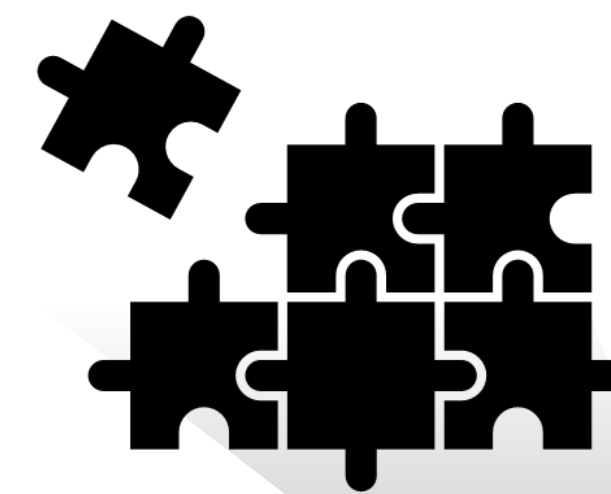


**FIND FAST THE BEST
POSSIBLE DESIGN**

SELF-DESIGNING SYSTEMS

automatically invent & build the perfect system for any new application

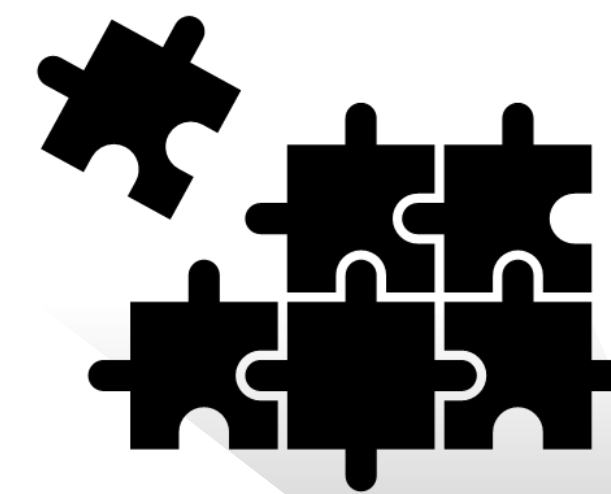
few existing designs



system design=
a set of low-level
design decisions

massive design space of system designs

few existing designs

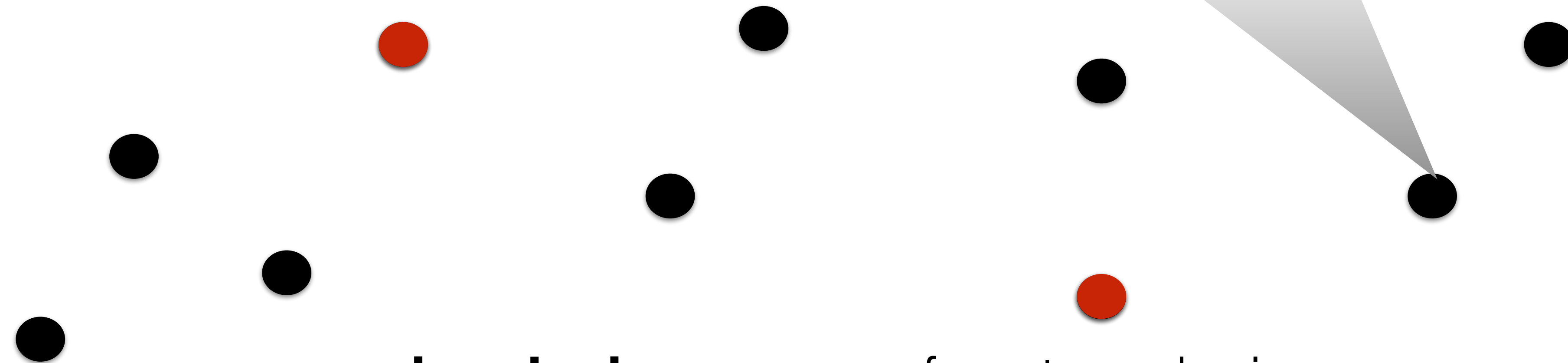
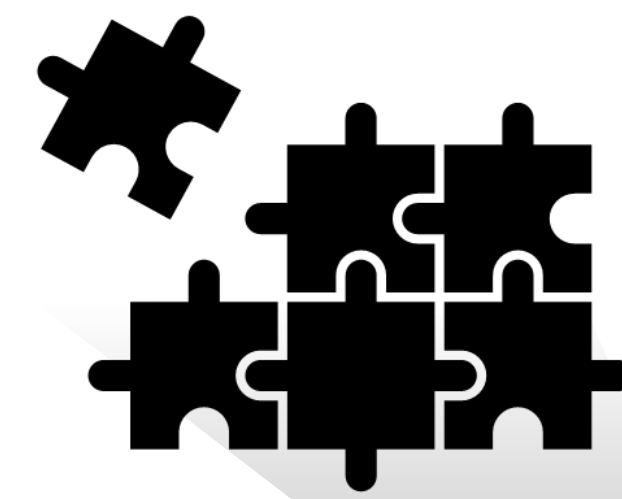


system design=
a set of low-level
design decisions

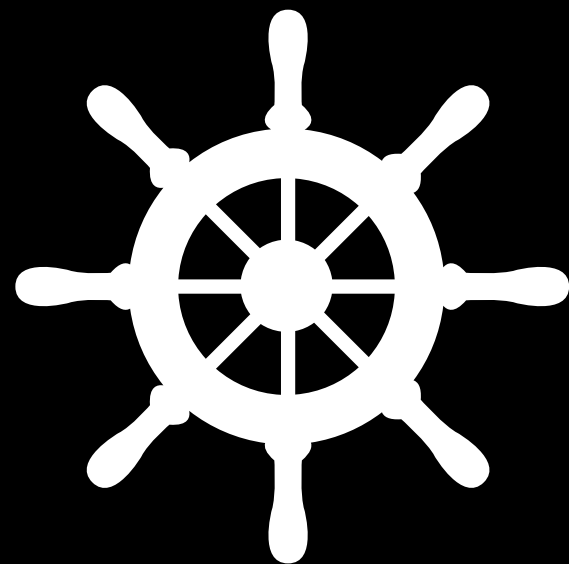
massive design space of system designs

workload
cloud
budget





massive design space of system designs



reasoning: understand all the design decisions & their impact



**AUTO DESIGN IS AS OLD
AS COMPUTER SCIENCE**



Rob Tarjan, Turing Award 1986

“IS THERE A CALCULUS OF DATA STRUCTURES

by which one can choose the appropriate representation
and techniques for a given problem?” (SIAM, 1978)

[P vs NP , average case, constant factors vs asymptotic, low bounds]



IS THERE A CALCULUS OF DATA SYSTEMS?



Rob Tarjan, Turing Award 1986

“IS THERE A CALCULUS OF DATA STRUCTURES

by which one can choose the appropriate representation
and techniques for a given problem?” (SIAM, 1978)

[P vs NP, average case, constant factors vs asymptotic, low bounds]

the **grammar** of data systems design



the **grammar** of data systems design

*action is for nothing
hope the most holy
fear free form of
am ultimate I theory*

Nikos Kazantzakis, philosopher

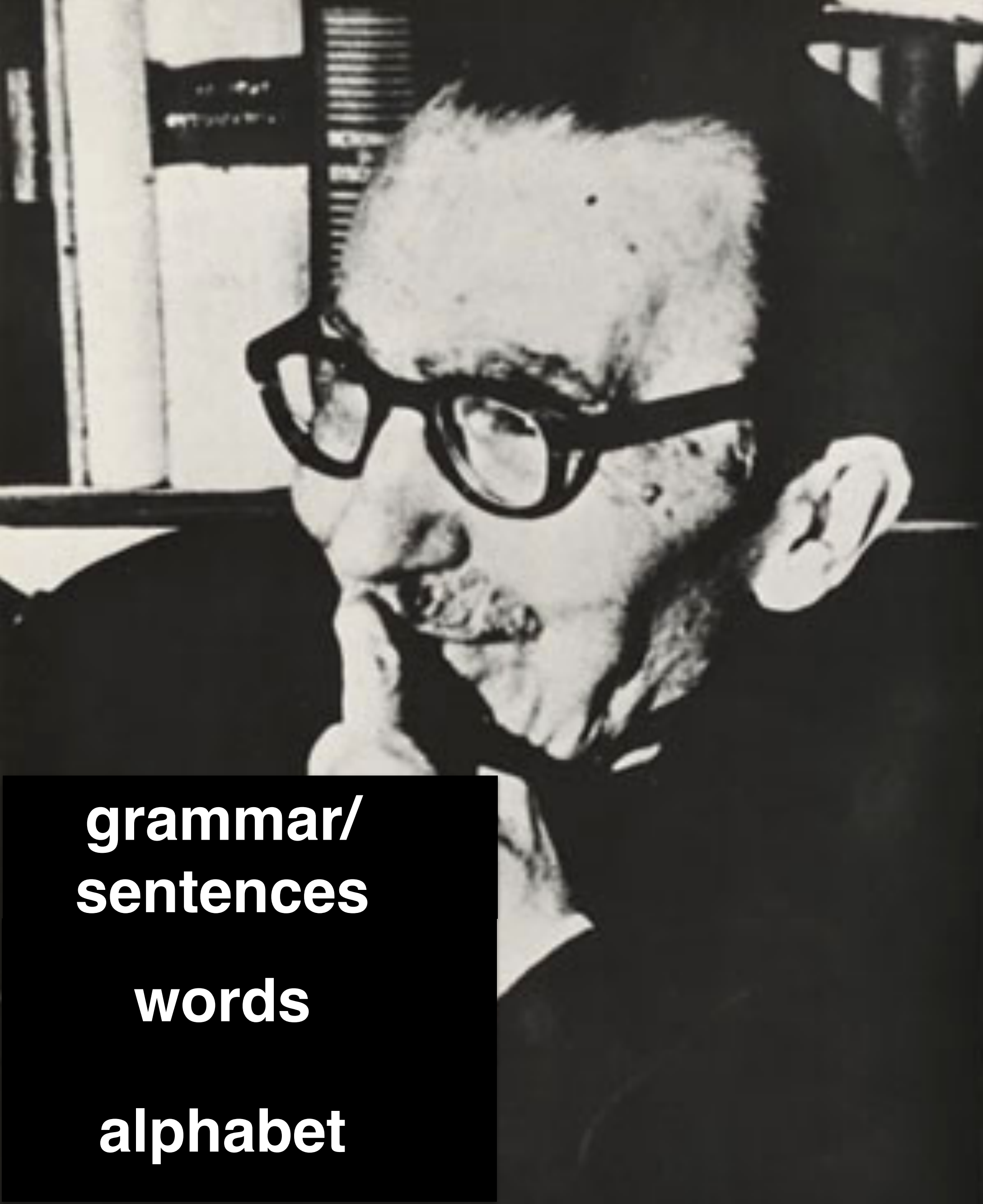


the **grammar** of data systems design

*action is
the most holy
ultimate form
theory*

*I hope for nothing
I fear nothing
I am free*

Nikos Kazantzakis, philosopher



**grammar/
sentences**

words

alphabet

Nikos Kazantzakis, philosopher

the **grammar** of data systems design

*action is
the most holy
ultimate form
theory*

*I hope for nothing
I fear nothing
I am free*



**grammar/
sentences**

words

alphabet

interactions

data structures

principles

Nikos Kazantzakis, philosopher

the **grammar** of data systems design

*action is
the most holy
ultimate form
theory*

*I hope for nothing
I fear nothing
I am free*



**grammar/
sentences**

words

alphabet

interactions

data structures

principles

Nikos Kazantzakis, philosopher

the **grammar** of data systems design

*action is
the most holy
ultimate form
theory*

NEW

*I hope for nothing
I fear nothing
I am free*



**grammar/
sentences**

words

alphabet

interactions

data structures

principles

Nikos Kazantzakis, philosopher

the **grammar** of data systems design

action is

the

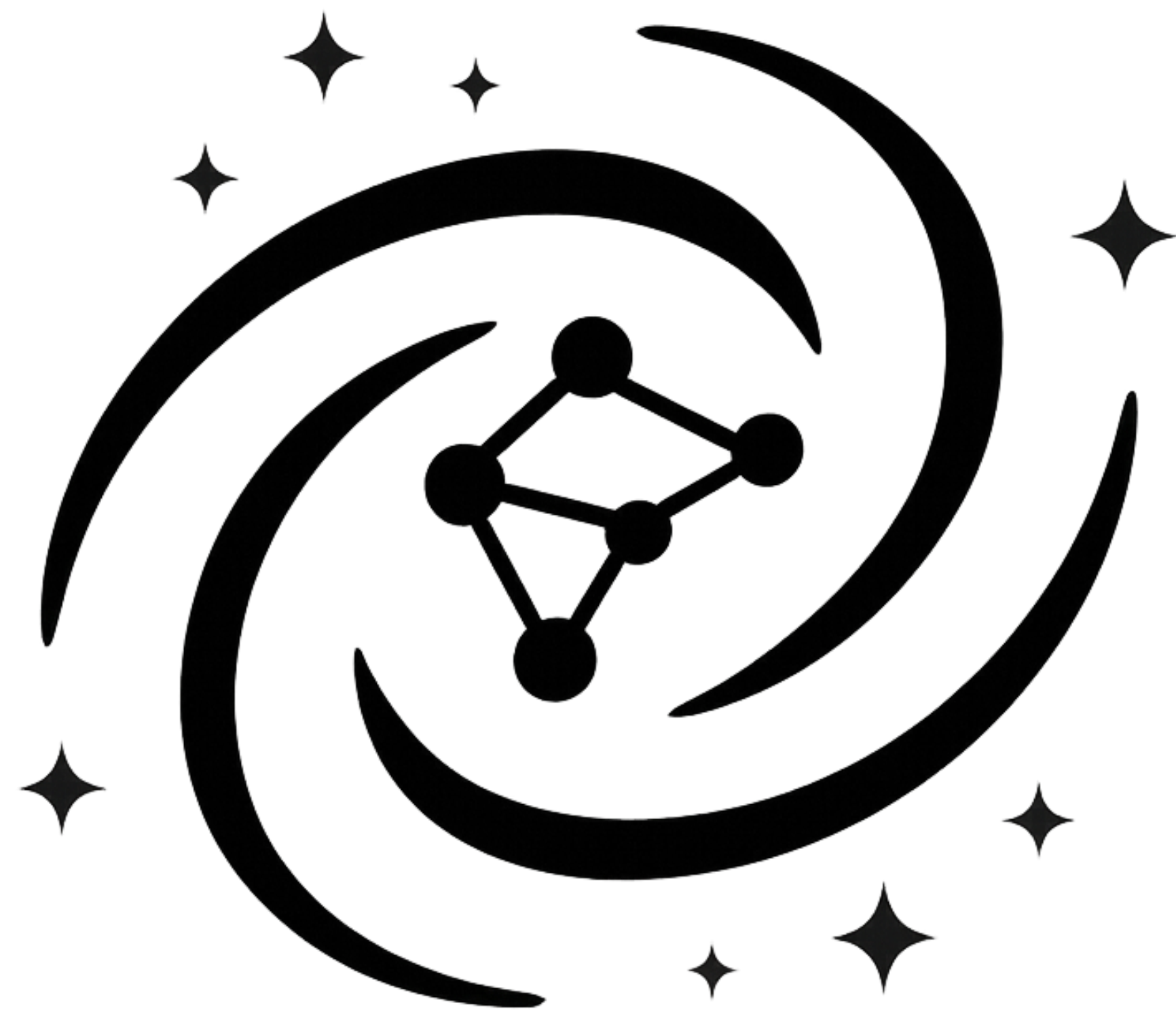
*most holy
of
form
theory*

which are “all”
possible *data systems*
we may ever invent?

I hope for nothing

I fear nothing

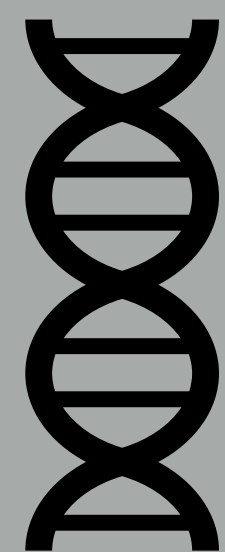
I am free



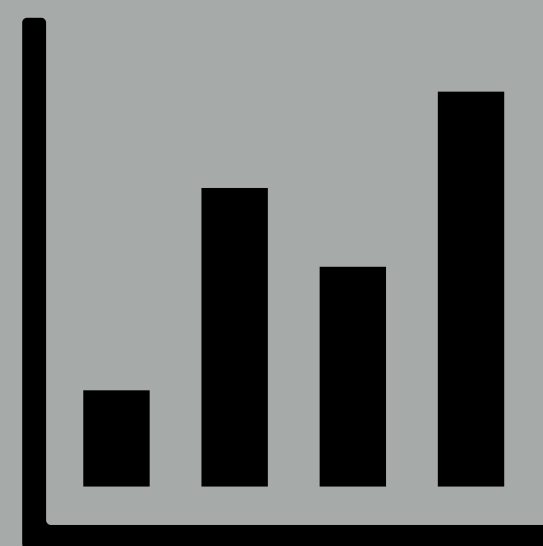
There exist trillions of possible designs for:

- 1) **Data Structures**
- 2) **Key-value Stores**
- 3) **Storage for Image AI**
- 4) **Large Model Training Algorithms**

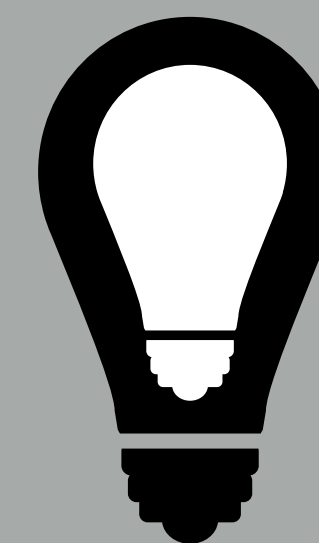
...and we can navigate their design space



DESIGN SPACE

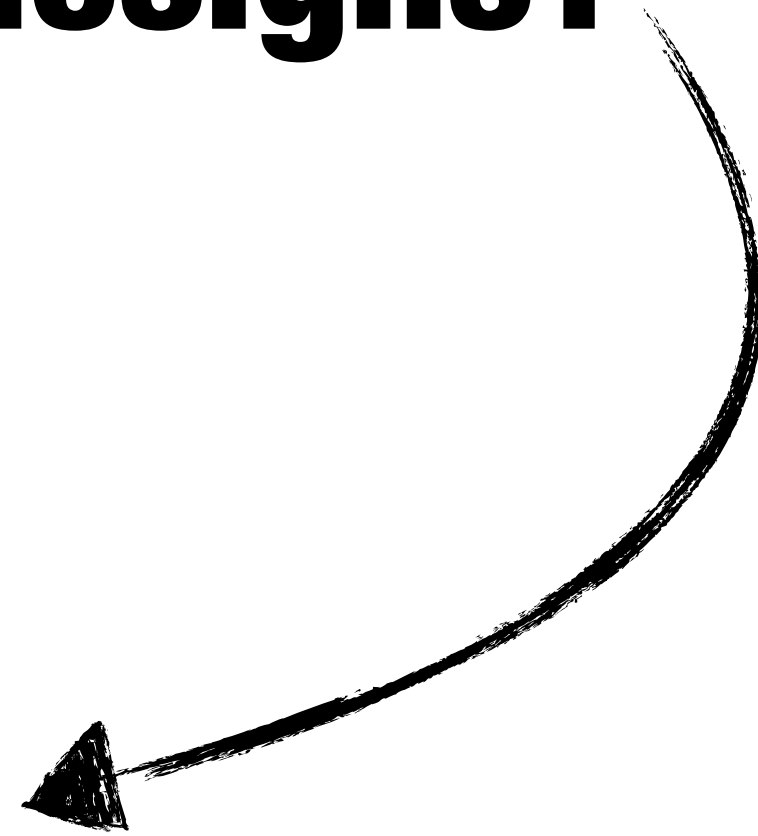


COST ESTIMATION

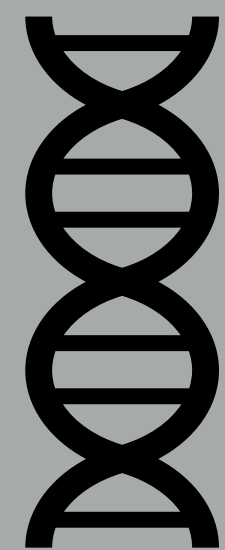


SEARCH

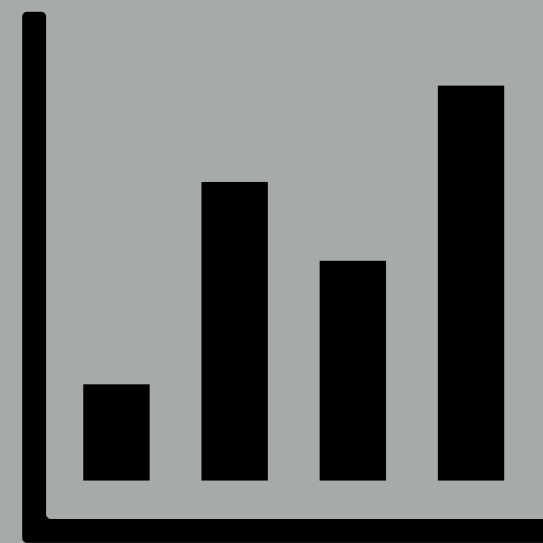
Which are all design principles & how they “connect” to synthesize all possible designs?



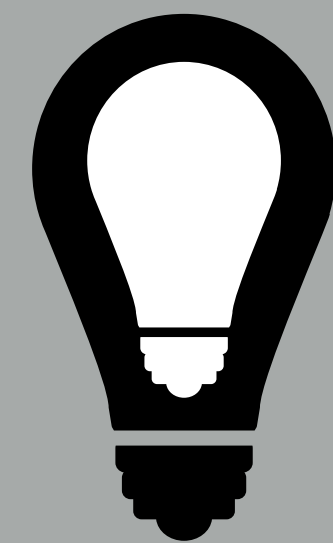
If we implemented in the best possible way 2 designs, how would they behave on data X & hardware Y?



DESIGN SPACE

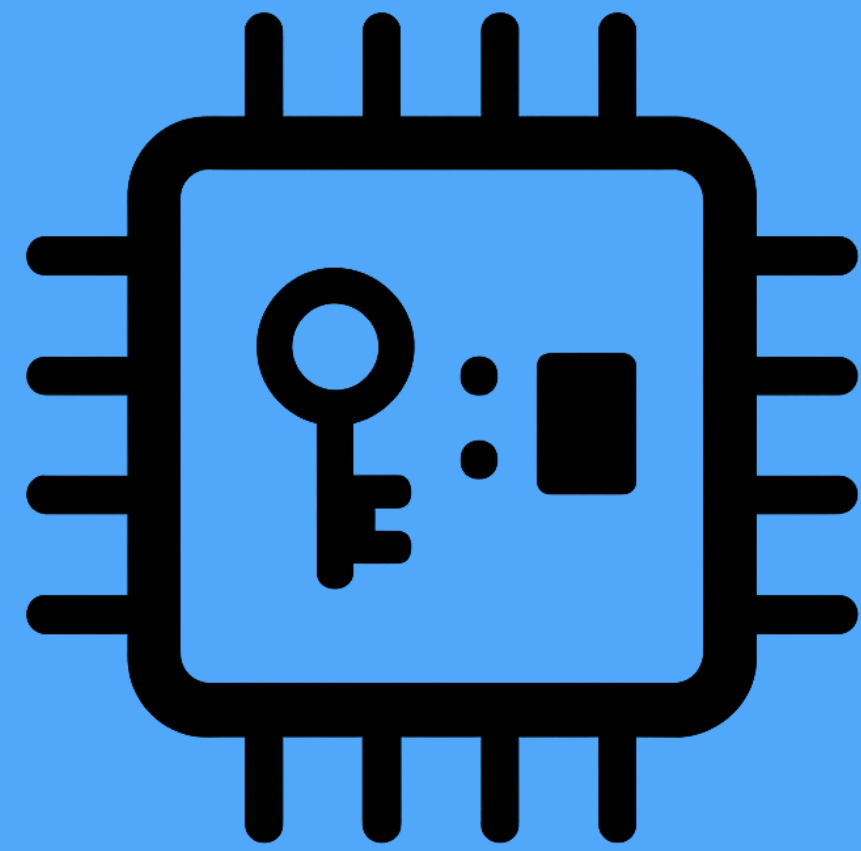


COST ESTIMATION



SEARCH

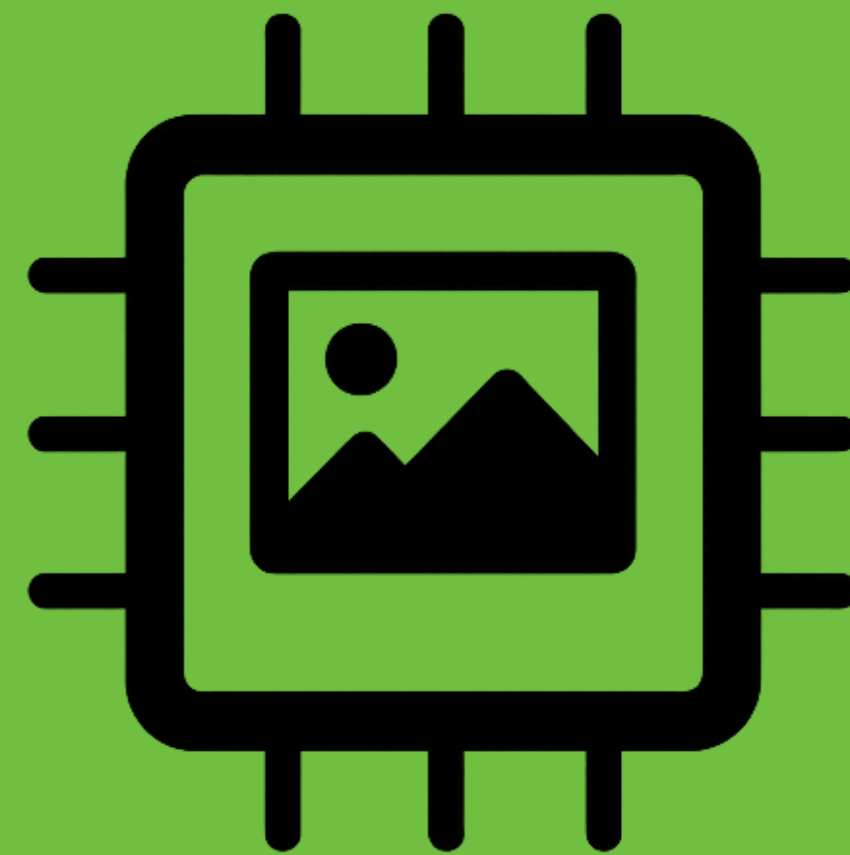
Cosine



1000x faster
key-value stores

SIGMOD'18/24, VLDB'22

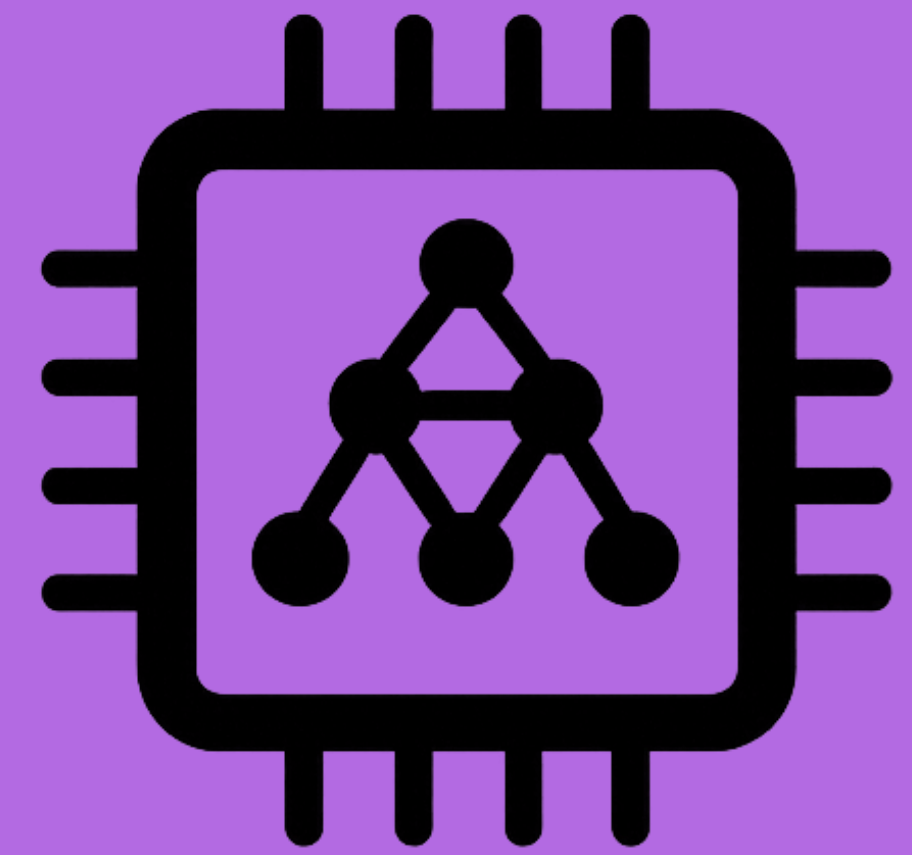
Image Calculator



10x faster
image AI inference

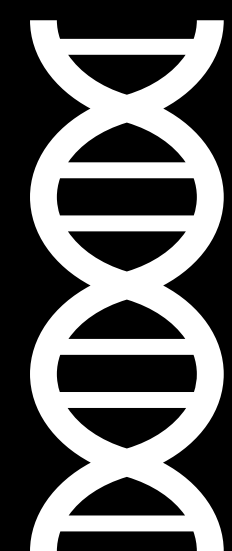
SIGMOD'24

TorchTitan

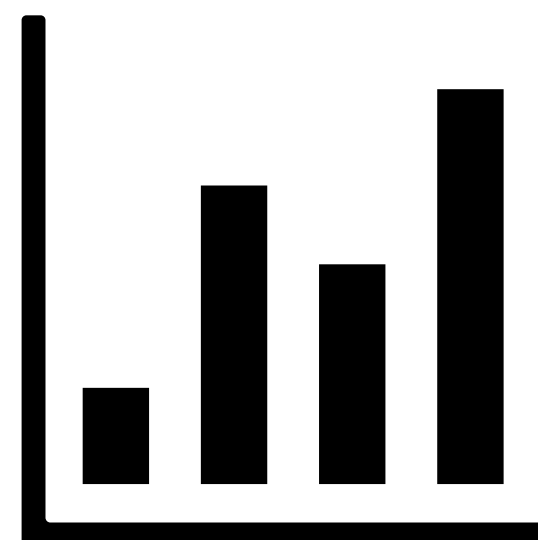


3x faster
large model training

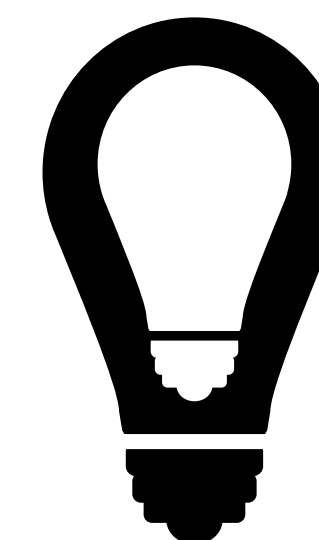
MLsys 2023, ICLR'25



DESIGN SPACE

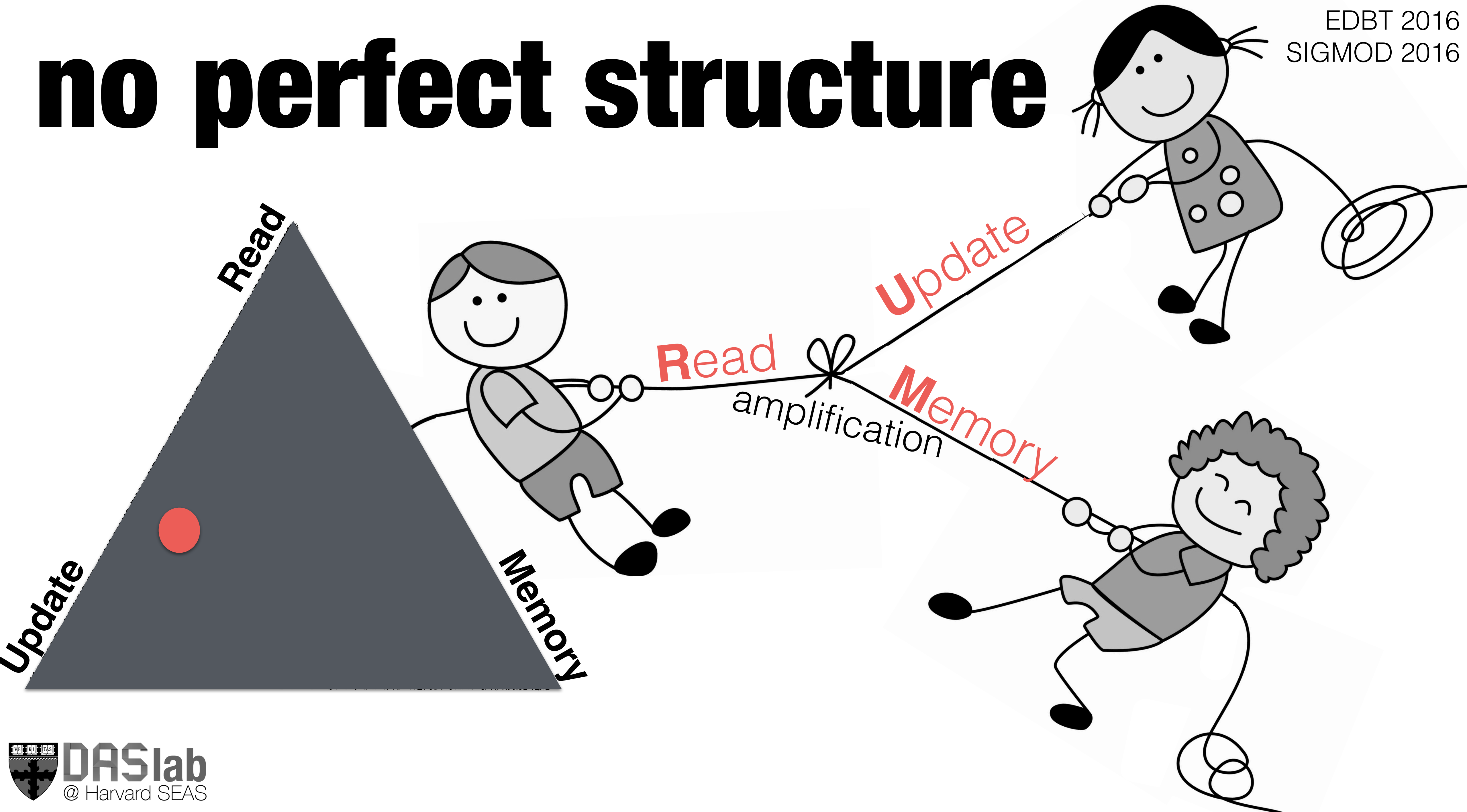


COST ESTIMATION

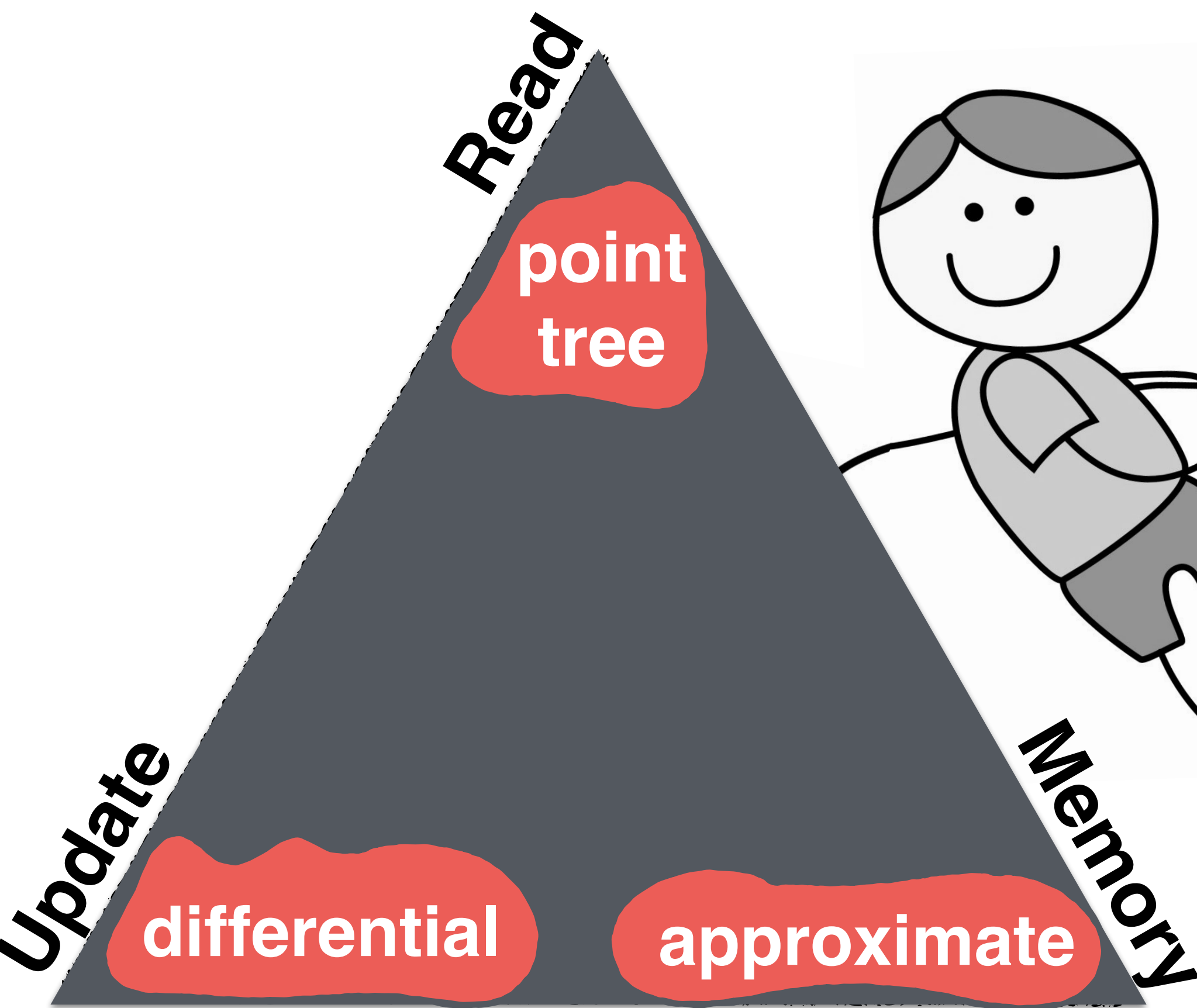


SEARCH

no perfect structure



no perfect structure



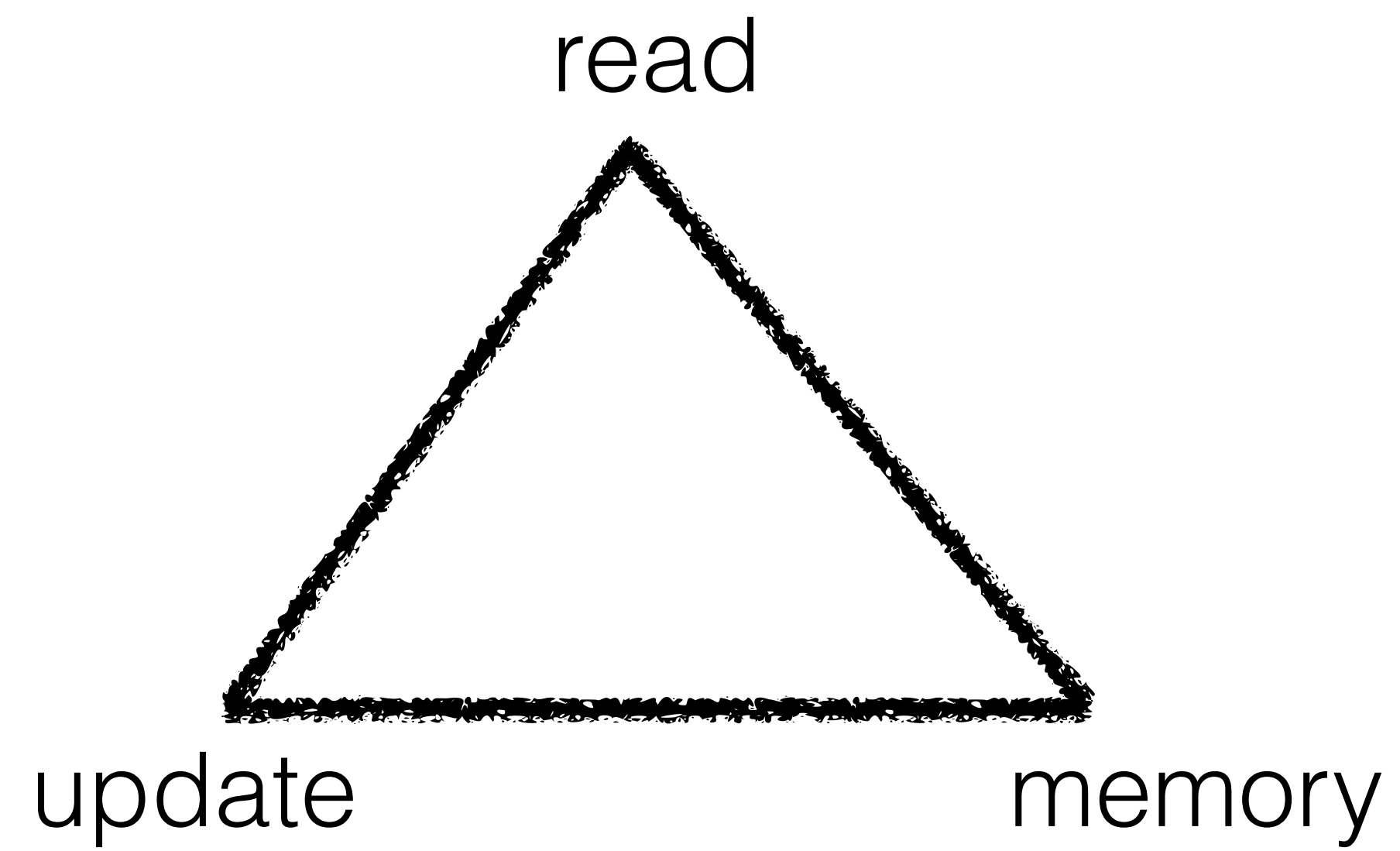
Read

amplification

Update

Memory





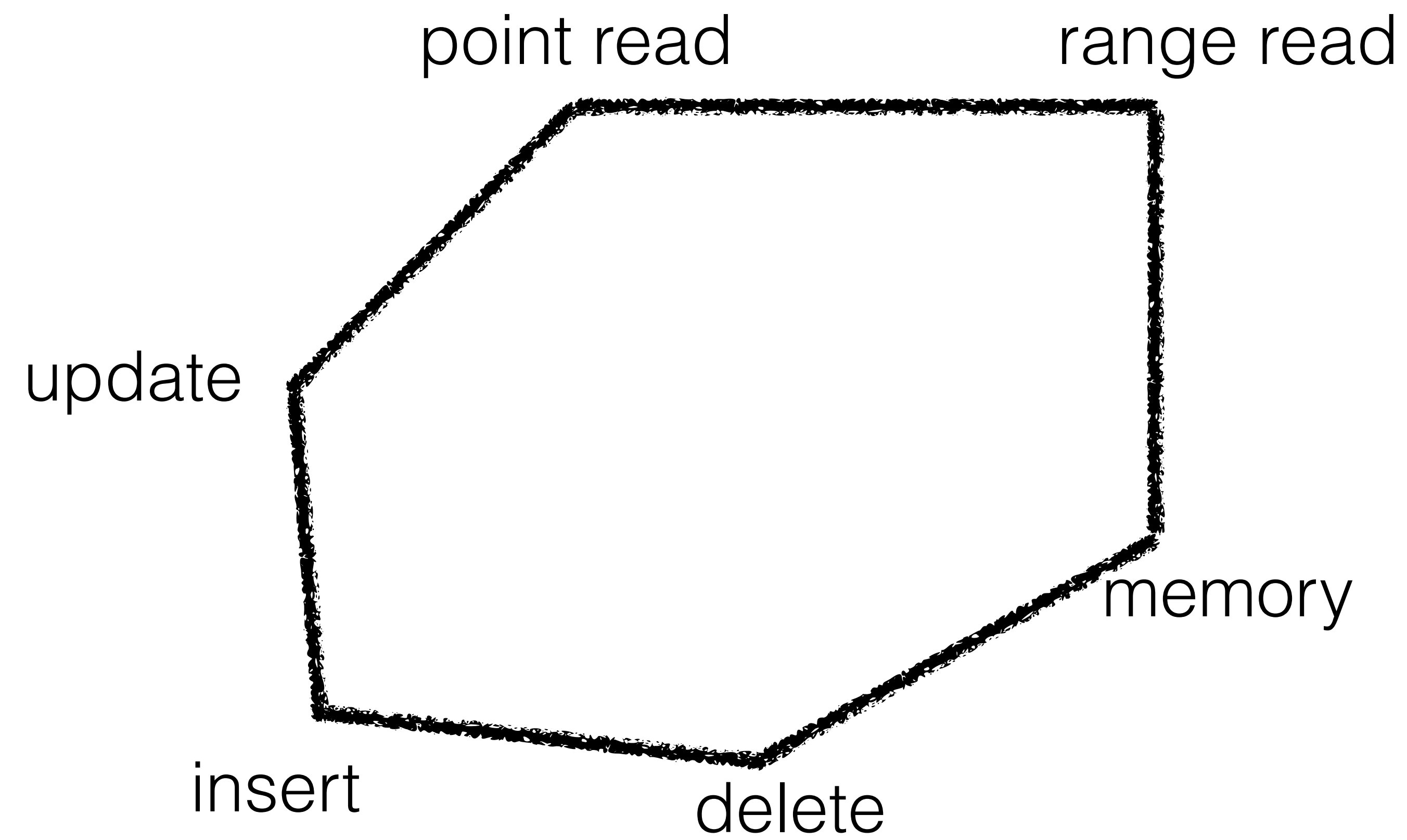
point read

range read



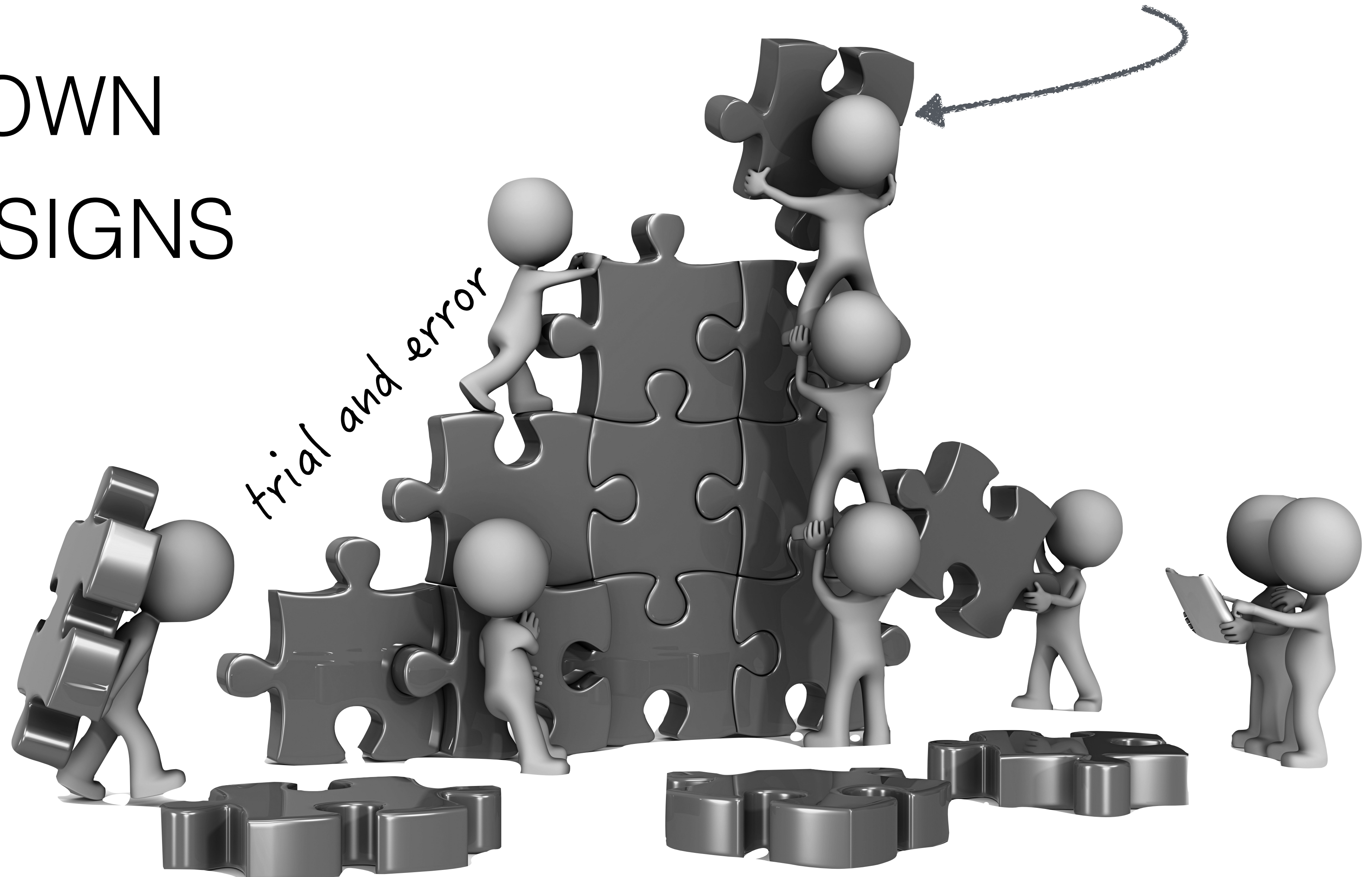
update

memory



FIRST PRINCIPLE: DESIGN CONCEPT THAT IS NOT POSSIBLE OR MEANINGFUL TO BREAK FURTHER

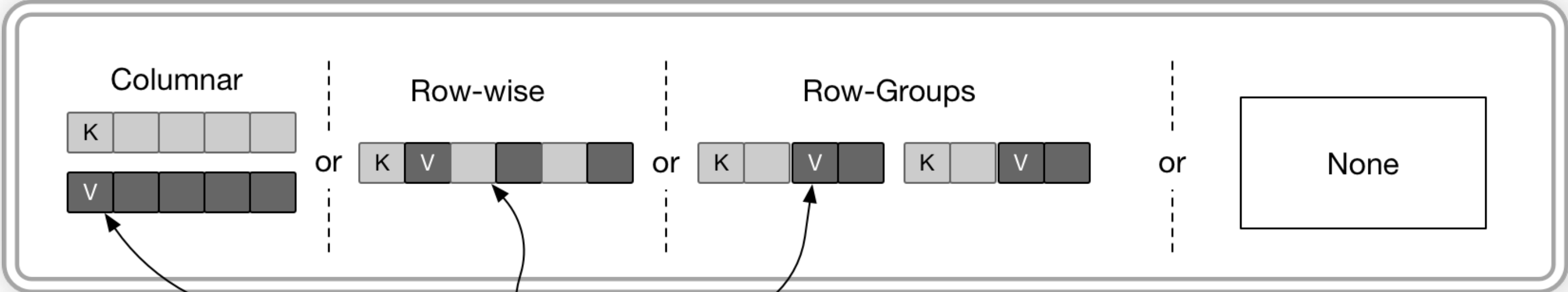
BREAK DOWN
KNOWN DESIGNS



Are keys retained? (yes, no, function)

Are values retained?

Utilization? (e.g., >50%)



Key and value layout

Fanout (fixed/functional | unlimited | terminal |)

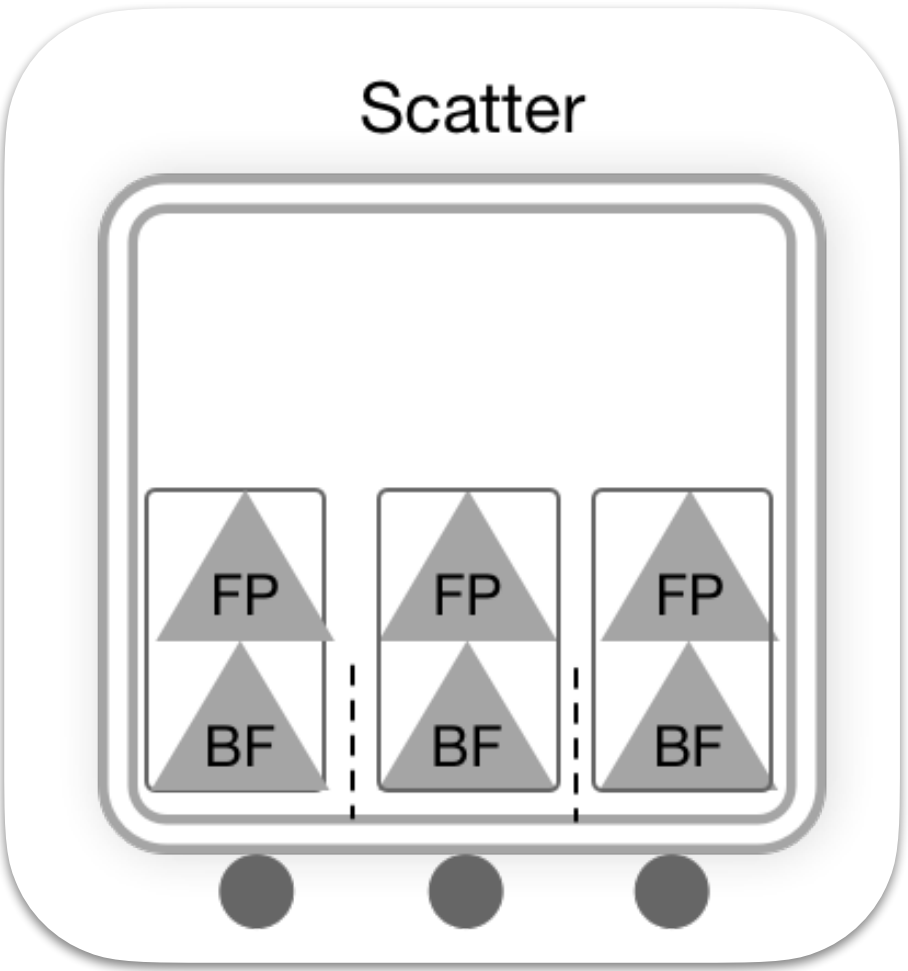
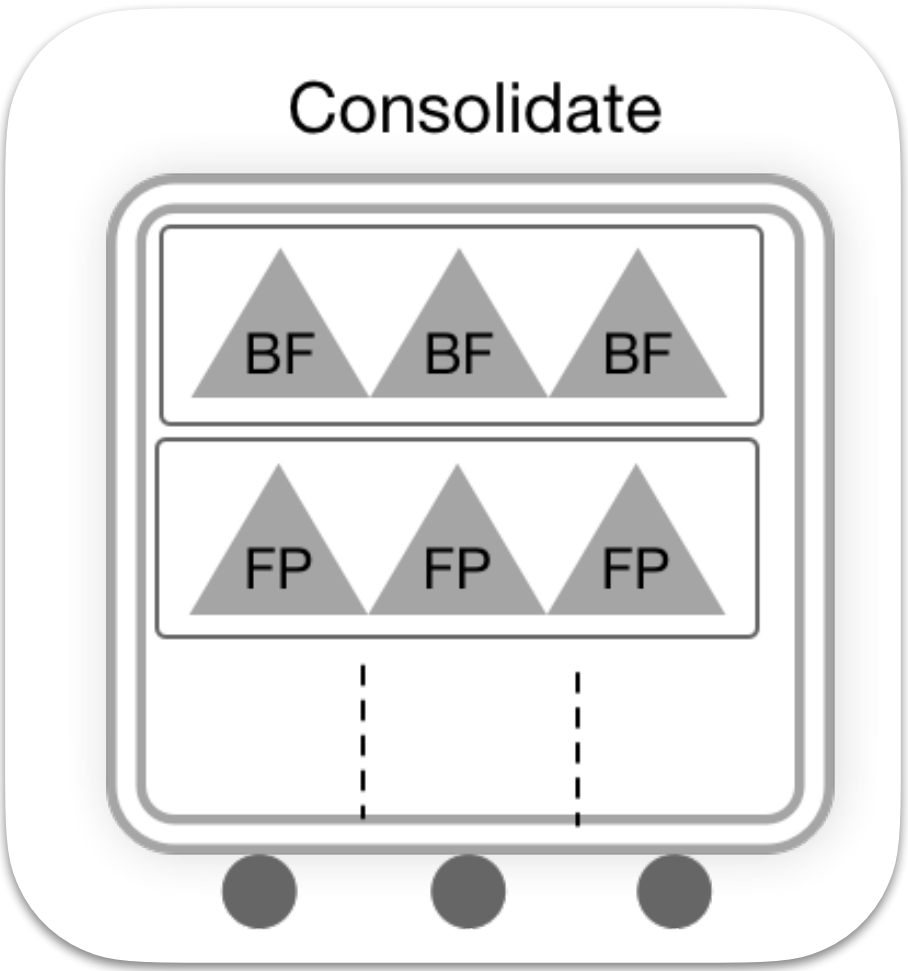
Key partitioning (none(fw-append | bw-append) | sorted | range() | radix() | function (func) | temporal(...))

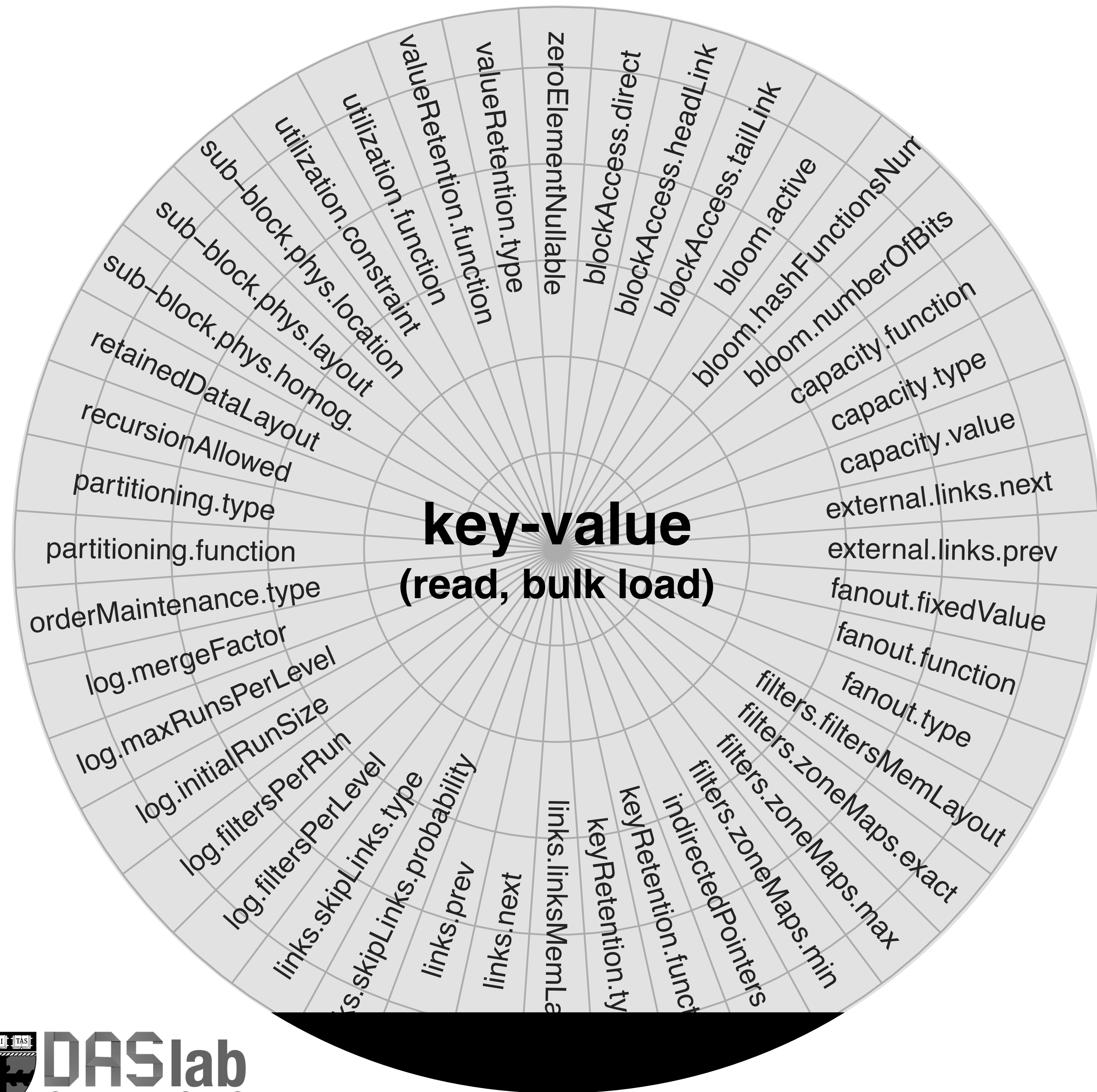
Zone Maps (min | max | both | exact | off)

Bloom filters (off | on(num_hashes: int, num_bits: int))

Filters layout (consolidate | scatter)

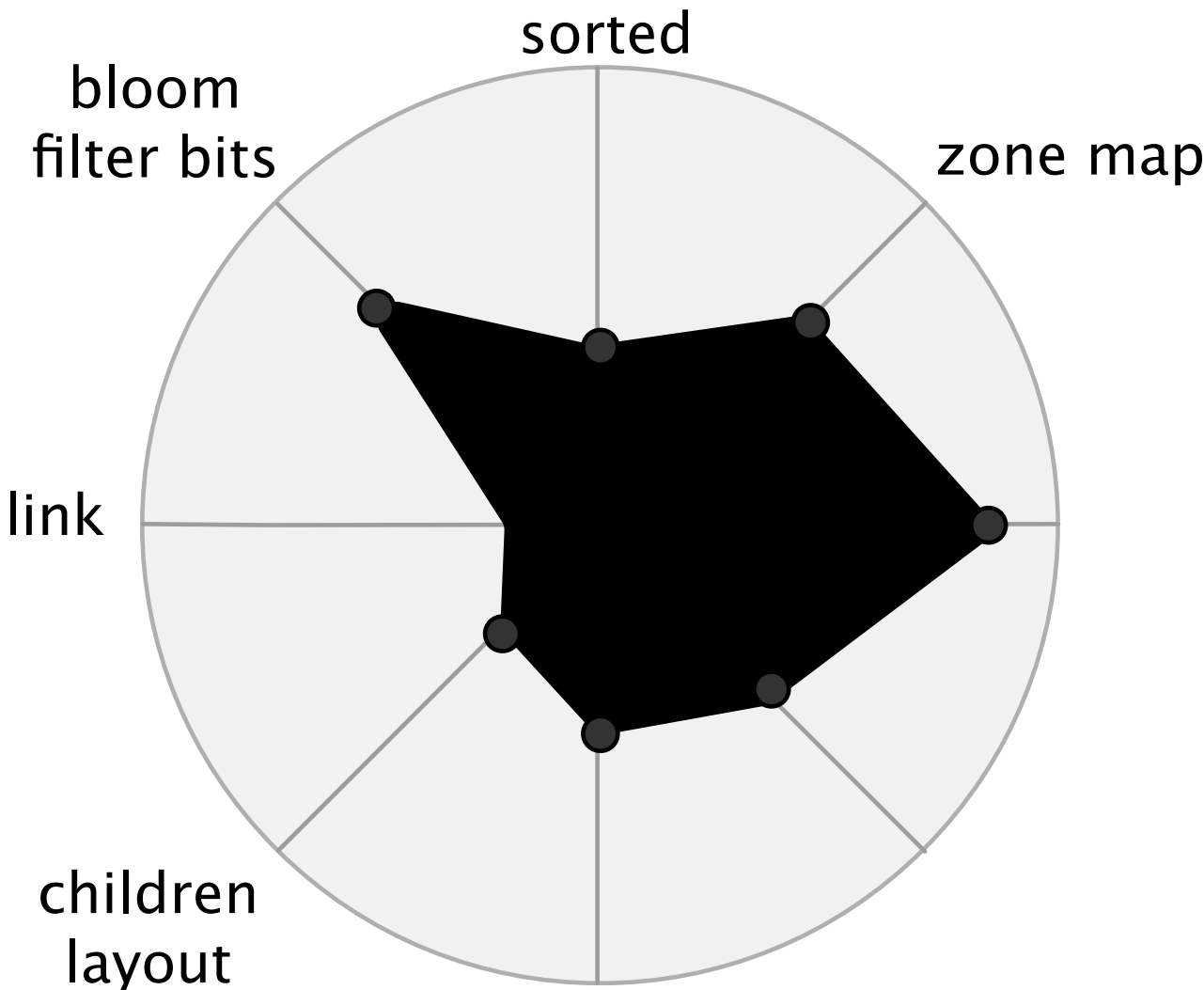
Links layout (consolidate | scatter)



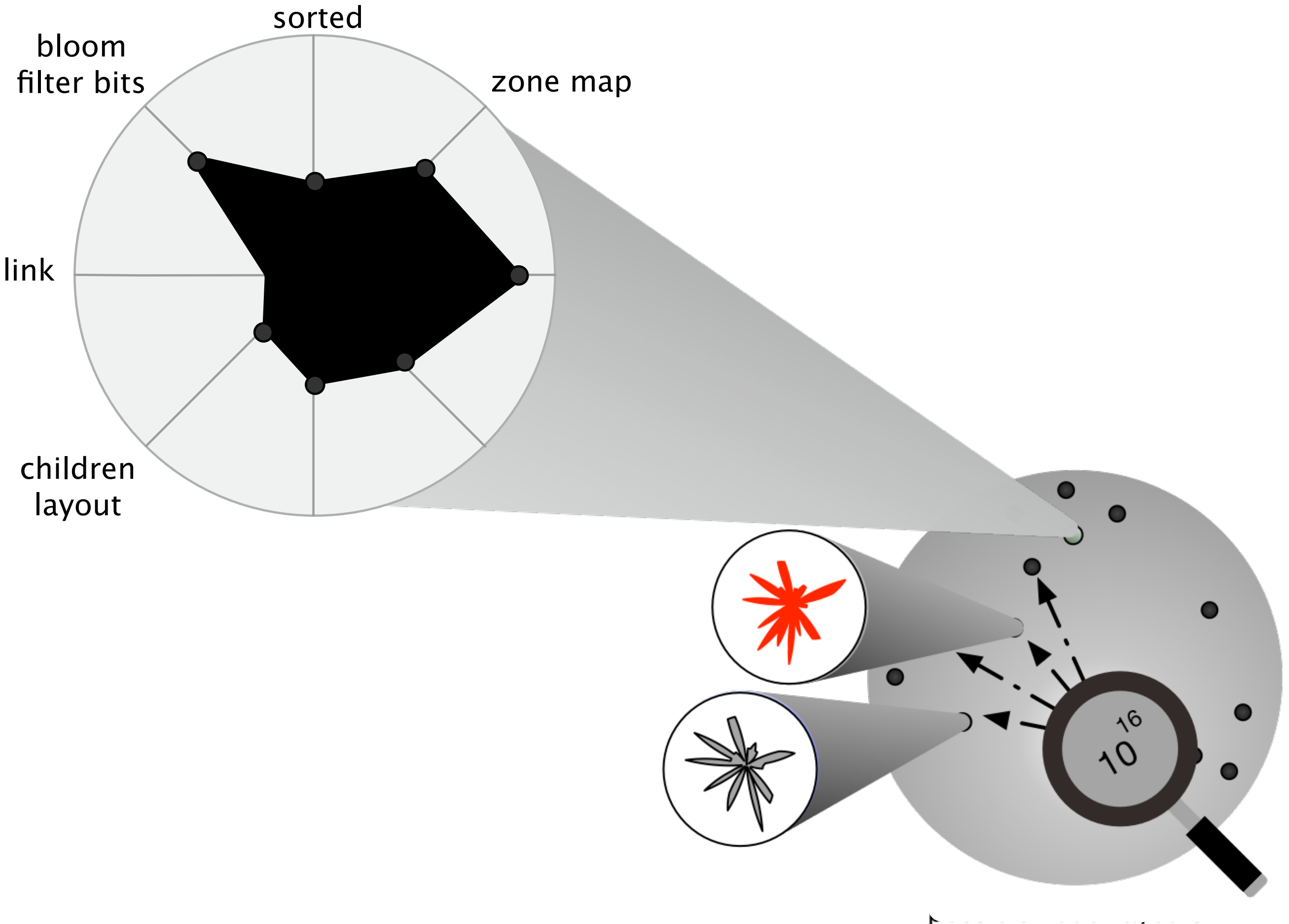


DESCRIBE ONE DATA BLOCK AT A TIME
AS A SET OF CONCEPTS
physical layout and domain partitioning

SETS OF CONCEPTS



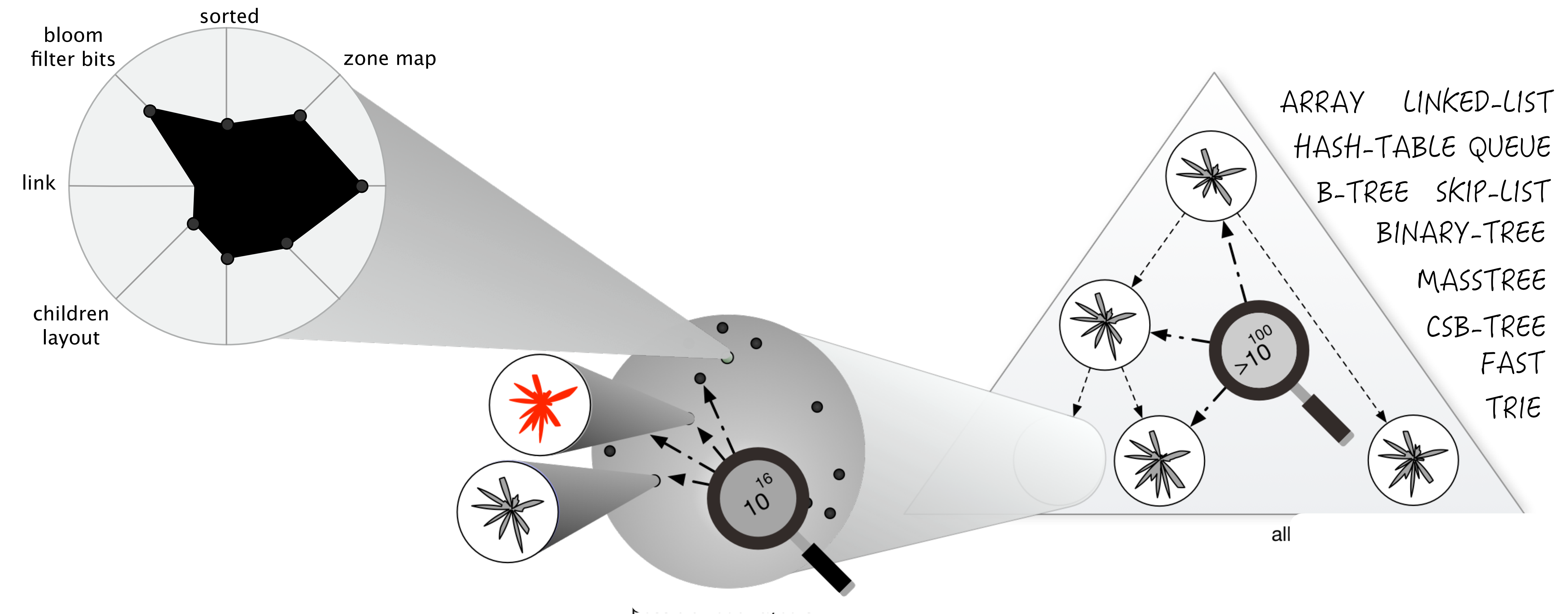
SETS OF CONCEPTS POSSIBLE NODE DESIGNS



SETS OF CONCEPTS

POSSIBLE NODE DESIGNS

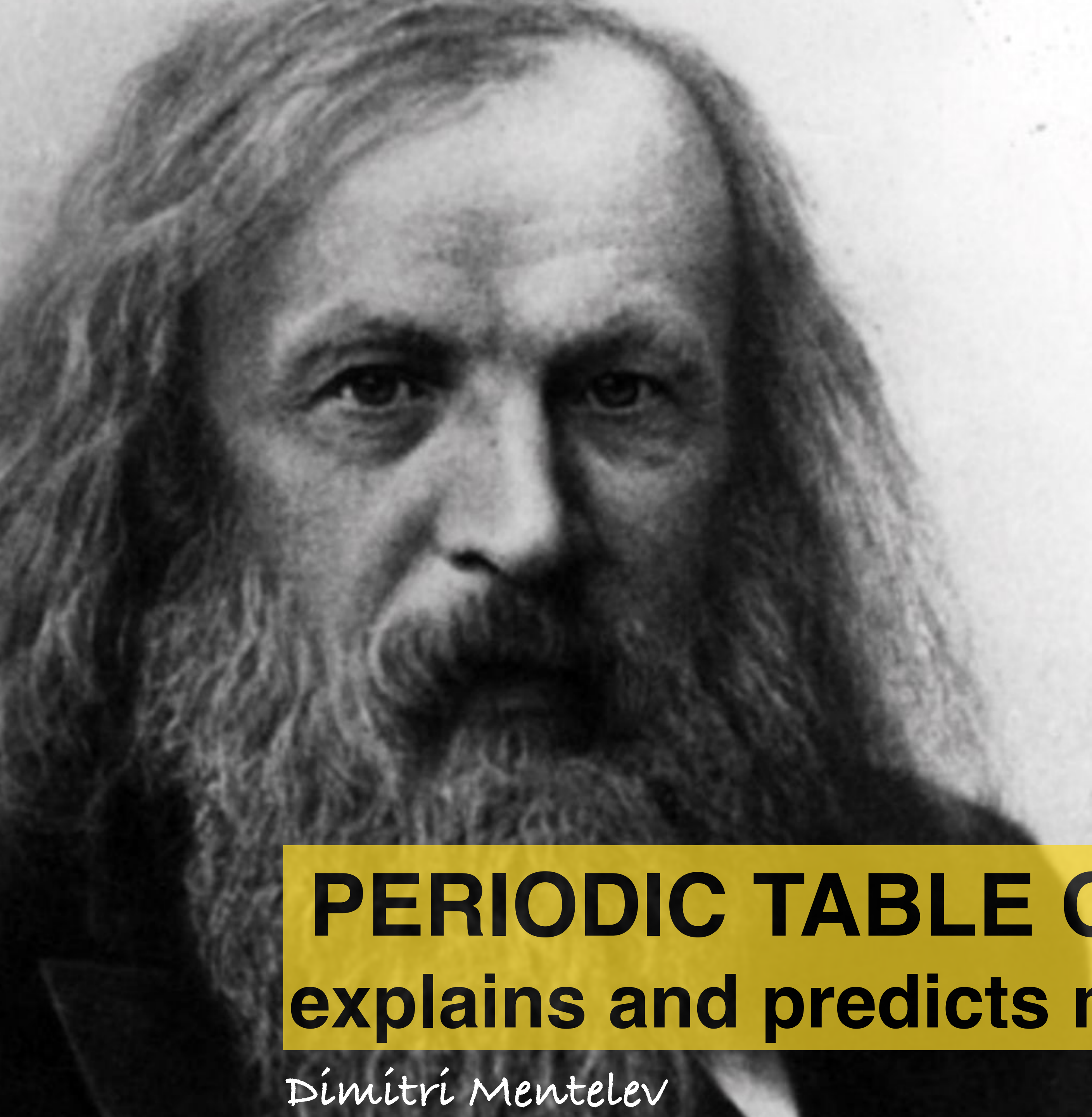
POSSIBLE STRUCTURES



Design Primitives to Auto Generate Trillions of Data Structures

				Unless otherwise specified, we use a reduced default values domain of 100 values for integers, 10 values for doubles, and 1 value for functions.		Hash Table			B+Tree/CSB+Tree/FAST				
						LPL							
				Primitive	Domain	size	H	LL	UDP	B+	CSB+	FAST	ODP
Node organization	1	Key retention. <u>No</u> : node contains no real key data, e.g., intermediate nodes of b+trees and linked lists. <u>Yes</u> : contains complete key data, e.g., nodes of b-trees, and arrays. <u>Function</u> : contains only a subset of the key, i.e., as in tries.		yes no function(func)	3		no	no	yes	no	no	no	yes
	2	Value retention. <u>No</u> : node contains no real value data, e.g., intermediate nodes of b+trees, and linked lists. <u>Yes</u> : contains complete value data, e.g., nodes of b-trees, and arrays. <u>Function</u> : contains only a subset of the values.		yes no function(func)	3		no	no	yes	no	no	no	yes
	3	Key order. Determines the order of keys in a node or the order of fences if real keys are not retained.		none sorted k-ary (k: int)	12		none	none	none	sorted	sorted	4-ary	sorted
	4	Key-value layout. Determines the physical layout of key-value pairs.		row-wise columnar col-row-groups(size: int)	12								
		<u>Rules</u> : requires key retention != no or value retention != no.						col.			col.		
	5	Intra-node access. Determines how sub-blocks (one or more keys of this node) can be addressed and retrieved within a node, e.g., with direct links, a link only to the first or last block, etc.		direct head_link tail_link link_function(func)	4		direct	head	direct	direct	direct	direct	direct
6	Utilization. Utilization constraints in regards to capacity. For example, >= 50% denotes that utilization has to be greater than or equal to half the capacity.		= (X%) function(func) none <i>(we currently only consider X=50)</i>	3		none	none	none	>= 50%	>= 50%	>= 50%	none	
Node filters	7	Bloom filters. A node's sub-block can be filtered using bloom filters. Bloom filters get as parameters the number of hash functions and number of bits.		off on(num_hashes: int, num_bits: int) <i>(up to 10 num_hashes considered)</i>	1001		off	off	off	off	off	off	off
	8	Zone map filters. A node's sub-block can be filitered using zone maps, e.g., they can filter based on mix/max keys in each sub-block.		min max both exact off	5		off	off	off	min	min	min	off
	9	Filters memory layout. Filters are stored contiguously in a single area of the node or scattered across the sub-blocks.		consolidate scatter	2					scatter	scatter	scatter	
<u>Rules</u> : requires bloom filter != off or zone map filters != off.													
	10	Fanout/Radix. Fanout of current node in terms of sub-blocks. This can either be unlimited (i.e., no restriction on the number of sub-blocks), fixed to a number, decided by a function or the node is terminal and thus has a fixed capacity.		fixed(value: int) function(func) mited terminal(cap: int) <i>(up to 10 different capacities and up to 10 fixed fanout values are considered)</i>	22		fixed(100)	unlimited	term(256)	fixed(20)	fixed(20)	fixed(16)	term(256)
	11	Key partitioning. Set if there is a pre-defined key partitioning imposed, e.g. the)						

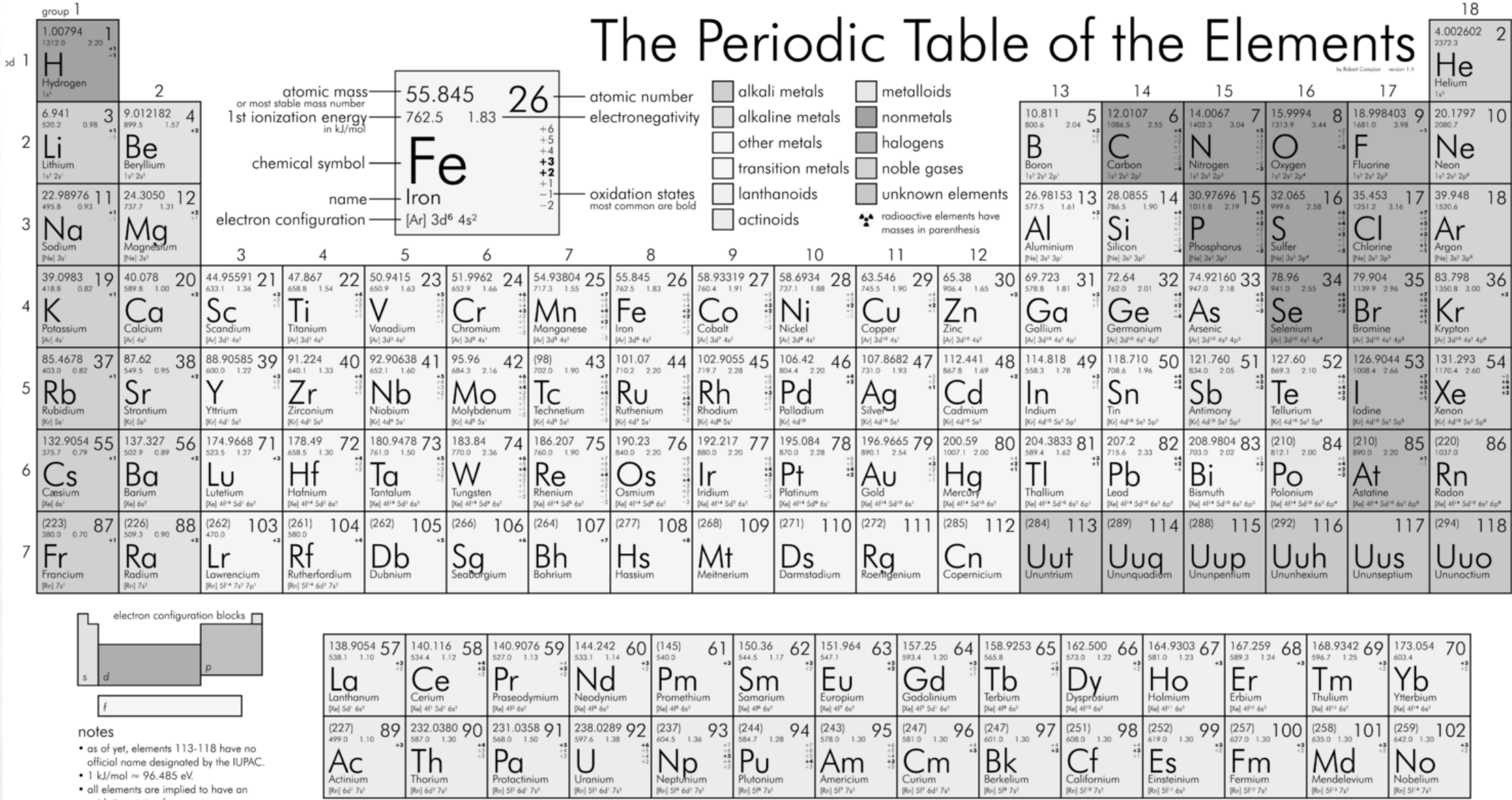
Partitioning		or balanced (i.e., all sub-blocks have the same size), unrestricted or functional.	function(func) parameter restricted function(func) (up to 10 different fixed capacity values are considered)	13	unrestricted	fixed(25)		balanced	balanced	balanced	
		Rules: requires key partitioning != none.									
	13	Immediate node links. Whether and how sub-blocks are connected.	next previous both none	4	none	next	none	none	none	none	none
	14	Skip node links. Each sub-block can be connected to another sub-block (not only the next or previous) with skip-links. They can be perfect, randomized or custom.	perfect randomized(prob: double) function(func) none	13	none	none	none	none	none	none	none
	15	Area-links. Each sub-tree can be connected with another sub-tree at the leaf level throu area links. Examples include the linked leaves of a B+Tree.	forward backward both none	4	none	none	forw.	none	none	none	none
Children layout	16	Sub-block physical location. This represents the physical location of the sub-blocks. Pointed: in heap, Inline: block physically contained in parent. Double-pointed: in heap but with pointers back to the parent.	inline pointed double-pointed	3	pointed	inline		pointed	pointed	pointed	
		Rules: requires fanout/radix != terminal.									
	17	Sub-block physical layout. This represents the physical layout of sub-blocks. Scatter: random placement in memory. BFS: laid out in a breadth-first layout. BFS layer list: hierarchical level nesting of BFS layouts.	BFS BFS layer(level-grouping: int) scatter (up to 3 different values for layer-grouping are considered)	5	scatter	scatter		scatter	BFS	BFS-LL	
		Rules: requires fanout/radix != terminal.									
	18	Sub-blocks homogeneous. Set to true if all sub-blocks are of the same type.	boolean	2	true	true		true	true	true	
		Rules: requires fanout/radix != terminal.									
	19	Sub-block consolidation. Single children are merged with their parents.	boolean	2	false	false		false	false	false	
		Rules: requires fanout/radix != terminal.									
	20	Sub-block instantiation. If it is set to eager, all sub-blocks are initialized, otherwise they are initialized only when data are available (lazy).	lazy eager	2	lazy	lazy		lazy	lazy	lazy	
		Rules: requires fanout/radix != terminal.									
21	Sub-block links layout. If there exist links, are they all stored in a single array (consolidate) or spread at a per partition level (scatter).	consolidate scatter	2		scatter						
	Rules: requires immediate node links != none or skip links != none.										
Recursion	22	Recursion allowed. If set to yes, sub-blocks will be subsequently inserted into a node of the same type until a maximum depth (expressed as a function) is reached. Then the terminal node type of this data structure will be used.	yes(func) no	3				yes(logn)	yes(logn)	yes(logn)	
		Rules: requires fanout/radix != terminal.			no	no					



PERIODIC TABLE OF ELEMENTS

explains and predicts missing elements

Dimitrí Mendeleev



4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

原子番号 (陽子数)	
1	元素名 (和名)
水素	元素記号
H	

					2 ヘリウム He
5 リチウム	6 ベリリウム	7 硼	8 炭素	9 窒素	10 酸素
B	C	N	O	F	Ne
13 アルミニウム	14 ケイ素	15 リン	16 硫黄	17 塩素	18 アルゴン
Al	Si	P	S	Cl	Ar
31 ガリウム	32 ゲルマニウム	33 ヒ素	34 セレン	35 臭素	36 クリプトン
n Ga	Ge	As	Se	Br	Kr
49 インジウム	50 スズ	51 アンチモン	52 テルル	53 ヨウ素	54 キセノン
d In	Sn	Sb	Te	I	Xe
81 タリウム	82 鉛	83 ビスマス	84 ポロニウム	85 アスタチン	86 ラドン
g Tl	Pb	Bi	Po	At	Rn
113	114 フルロビウム	115	116 リホベリウム	117	118
113		115	Lv	117	118

・超重元素

7	58	59	60	61	62	63	64	65	66	67	68	69	70	71
ランタニウム	セリウム	プラセオジム	ネオジム	プロメチウム	サマリウム	ユウロピウム	ガドリニウム	テルビウム	ジスプロシウム	ホウミウム	エルビウム	イットリウム	ランタニウム	ルテチウム
La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Tb	Dy				Yb	Lu
89	90	91	92	93	94	95	96	97	98	99	100	101	102	103
アクチニウム	トランシウム	アクチニウム	ウラン	ネプツウム	プルトニウム	アメリシウム	キュリウム	ベルグショウム	カリフォルニウム	エーレンバークニウム	ヘーグニウム	コペルニウム	ドブニウム	ホーグニウム
Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Mn	Uu	Uub

nihonium



理化学

Kosuke Morita

the periodic table of data structures

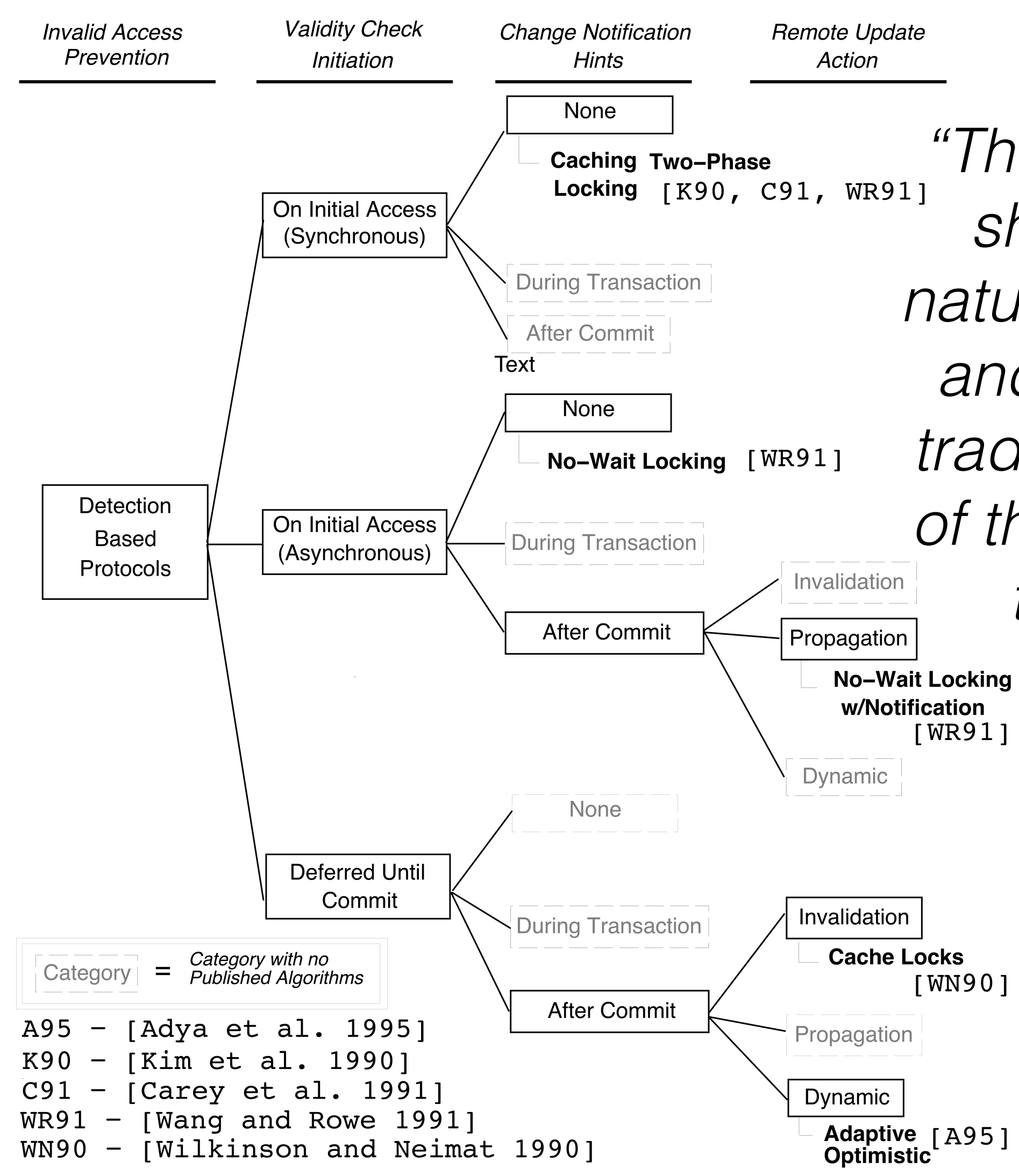
		the periodic table of data structures									
classes of primitives	classes of designs										
	B-trees & Variants	Tries & Variants	LSM-Trees & Variants	Differential Files	Membership Tests	Zone maps & Variants	Bitmaps & Variants	Hashing	Base Data & Columns		
	Partitioning	DONE	DONE	DONE					DONE	DONE	↓↑↑ RUM
	Logarithmic Design	DONE	DONE	DONE							↓↓↑ RUM
	Fractional Cascading	DONE		DONE	DONE						↓↑↑ RUM
	Log- Structured	DONE		DONE	DONE						↑↓↑ RUM
	Buffering	DONE			DONE				DONE		↓♦↑ RUM
	Differential Updates	DONE			DONE						↑↓↓ RUM
	Sparse Indexing	DONE				DONE	DONE				↓♦↑ RUM
Adaptivity	DONE								DONE		



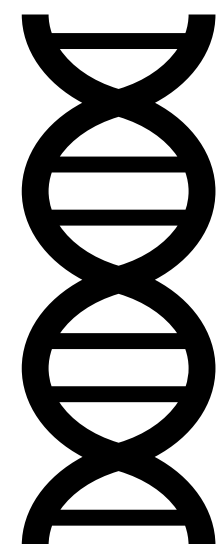
TAXONOMY OF COMPLEX ALGORITHMS

transactional cache consistency maintenance

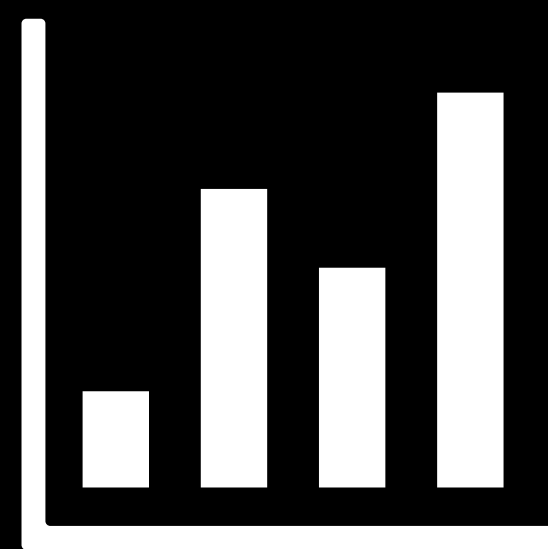
Mike Franklin



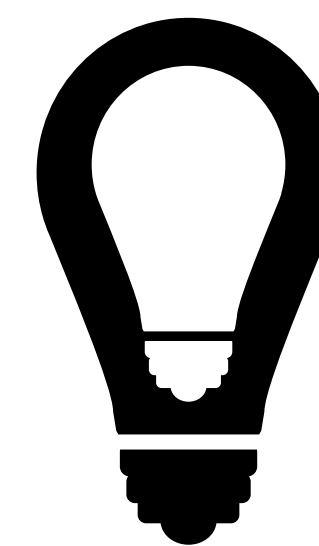
“The taxonomy is used to shed light both on the nature of the design space and on the performance tradeoffs implied by many of the choices that exist in the design space.”



DESIGN SPACE



COST ESTIMATION



SEARCH

HOW TO JUDGE A DESIGN?



1

**COMPLEXITY
ANALYSIS**



2

**IMPLEMENTATION
& TESTING**

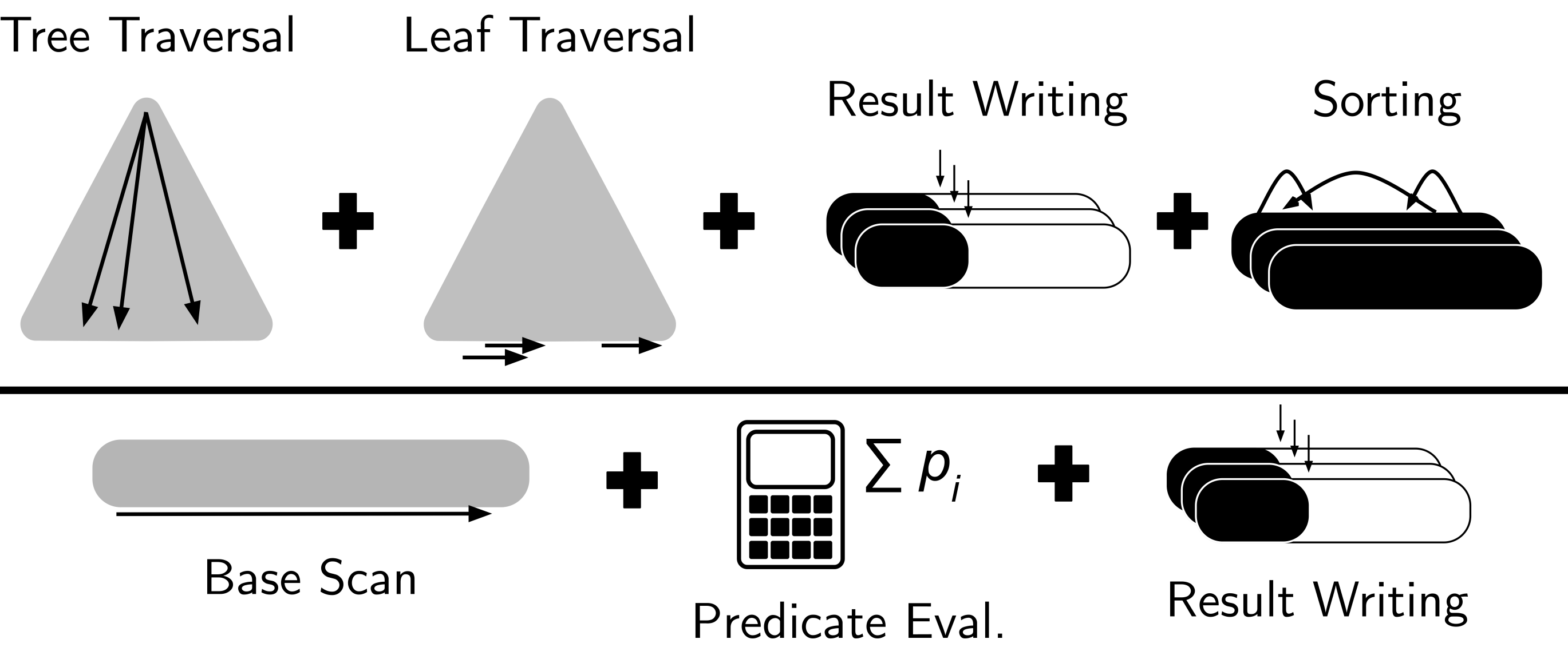


3

**GENERALIZED
MODELS**

HARD & SLOW

Access path selection @SIGMOD2017

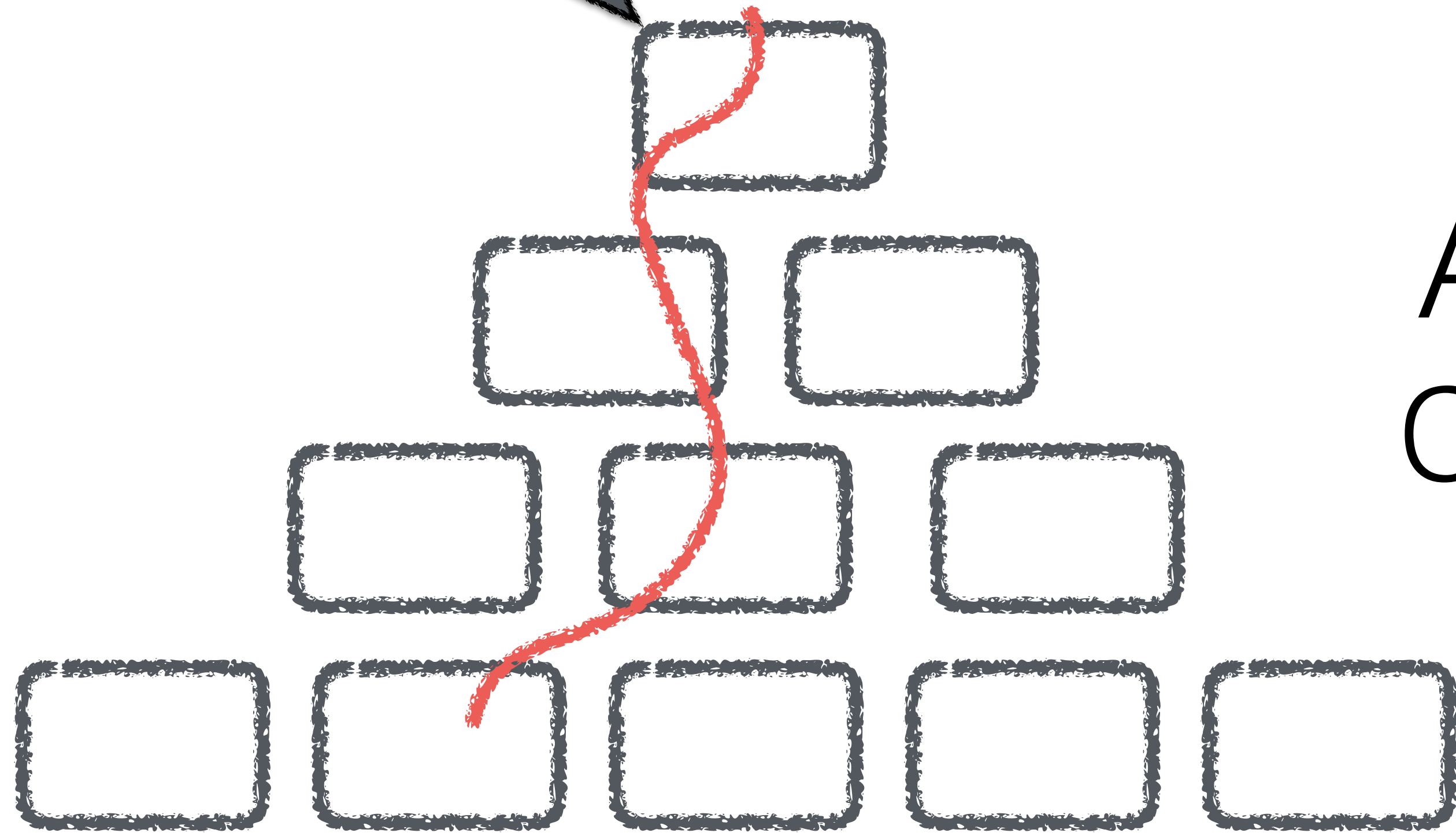


$$\begin{aligned}
 APS(q, S_{tot}) = & \frac{q \cdot \frac{1 + \lceil \log_b(N) \rceil}{N} \cdot \left(BW_S \cdot C_M + \frac{b \cdot BW_S \cdot C_A}{2} + \frac{b \cdot BW_S \cdot f_p \cdot p}{2} \right)}{\max(ts, 2 \cdot f_p \cdot p \cdot q \cdot BW_S) + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}} \\
 & + \frac{S_{tot} \left(\frac{BW_S \cdot C_M}{b} + (aw + ow) \cdot \frac{BW_S}{BW_I} + rw \cdot \frac{BW_S}{BW_R} \right)}{\max(ts, 2 \cdot f_p \cdot p \cdot q \cdot BW_S) + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}} \\
 & + \frac{S_{tot} \cdot \log_2(S_{tot} \cdot N) \cdot BW_S \cdot C_A}{\max(ts, 2 \cdot f_p \cdot p \cdot q \cdot BW_S) + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}}
 \end{aligned}$$

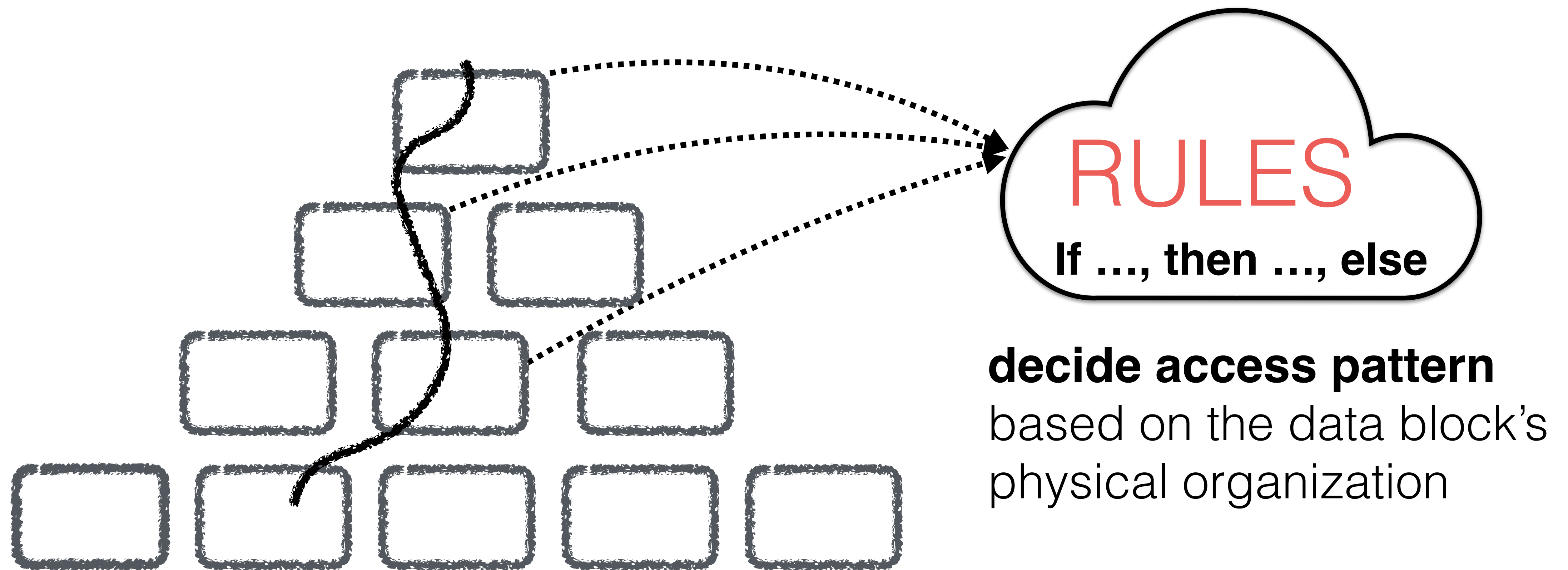
Workload	q s_i S_{tot}	number of queries selectivity of query i total selectivity of the workload
Dataset	N ts	data size (tuples per column) tuple size (bytes per tuple)
Hardware	C_A C_M BW_S BW_R BW_I p f_p	L1 cache access (sec) LLC miss: memory access (sec) scanning bandwidth (GB/s) result writing bandwidth (GB/s) leaf traversal bandwidth (GB/s) The inverse of CPU frequency Factor accounting for pipelining
Scan & Index	rw b aw ow	result width (bytes per output tuple) tree fanout attribute width (bytes of the indexed column) offset width (bytes of the index column offset)

DESIGN
SPACE
OF POSSIBLE
STORAGE LAYOUTS

operation



ALGORITHM &
COST SYNTHESIS



**sorted keys
columnar layout**

RULES



**sorted
search**

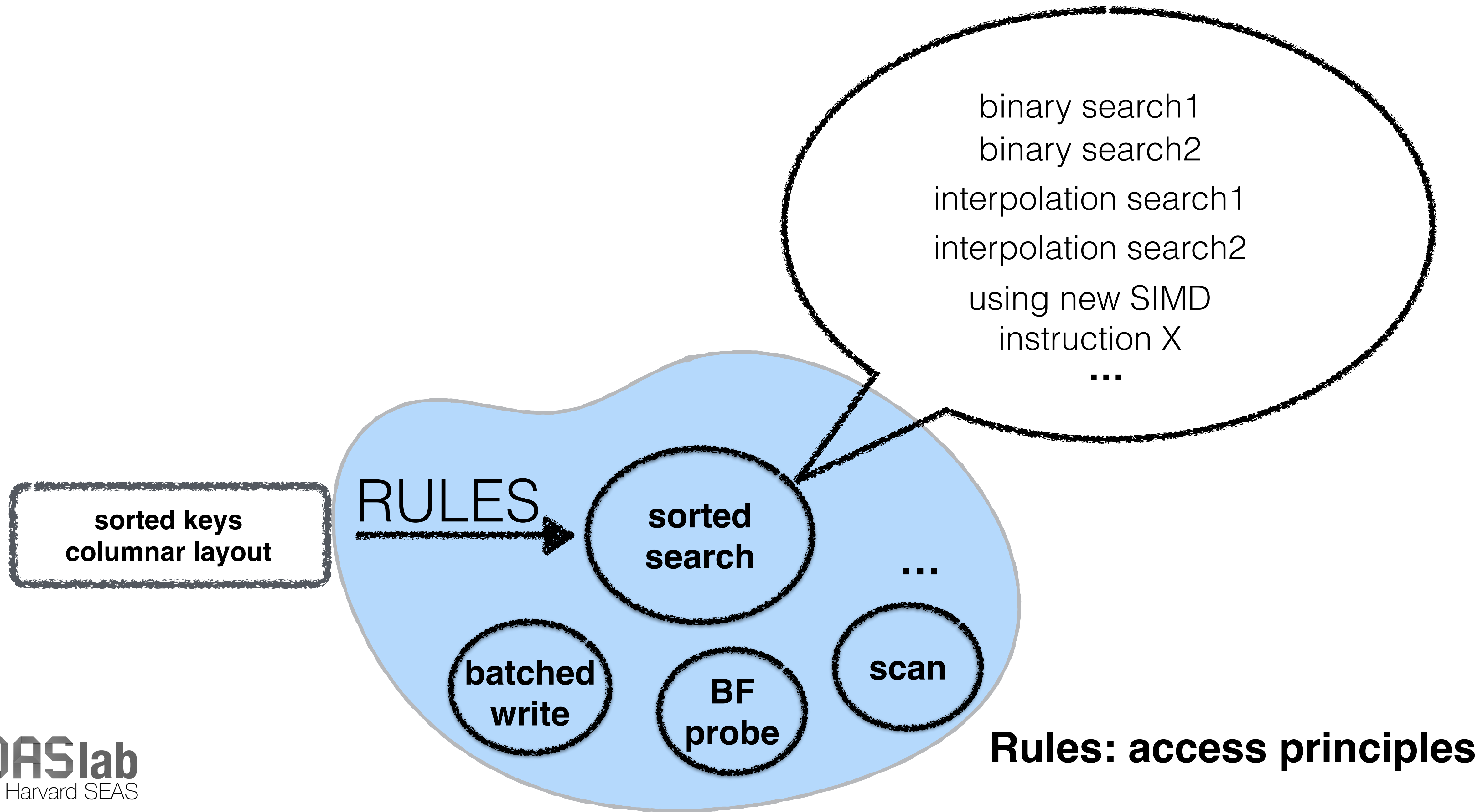
DEPENDS ON **HARDWARE ENGINEERING**

sorted keys
columnar layout

RULES →

sorted
search

binary search1
binary search2
interpolation search1
interpolation search2
using new SIMD
instruction X
...



Rules: access principles

sorted keys
columnar layout

RULES

sorted
search

batched
write

BF
probe

scan

...

binary search1

binary search2

interpolation search1

interpolation search2

using new SIMD

instruction X

...

code,
model

code,
model

code,
model

**Learning of fine-grained
access patterns**

Rules: access principles

SYNTHESIS FROM LEARNED MODELS

coding, modeling, generalized models, and a touch of ML



1. MINIMAL CODE

e.g., binary search

```
if (data[middle] < search_val) {  
    low = middle + 1;  
} else {  
    high = middle;  
}  
middle = (low + high)/2;
```

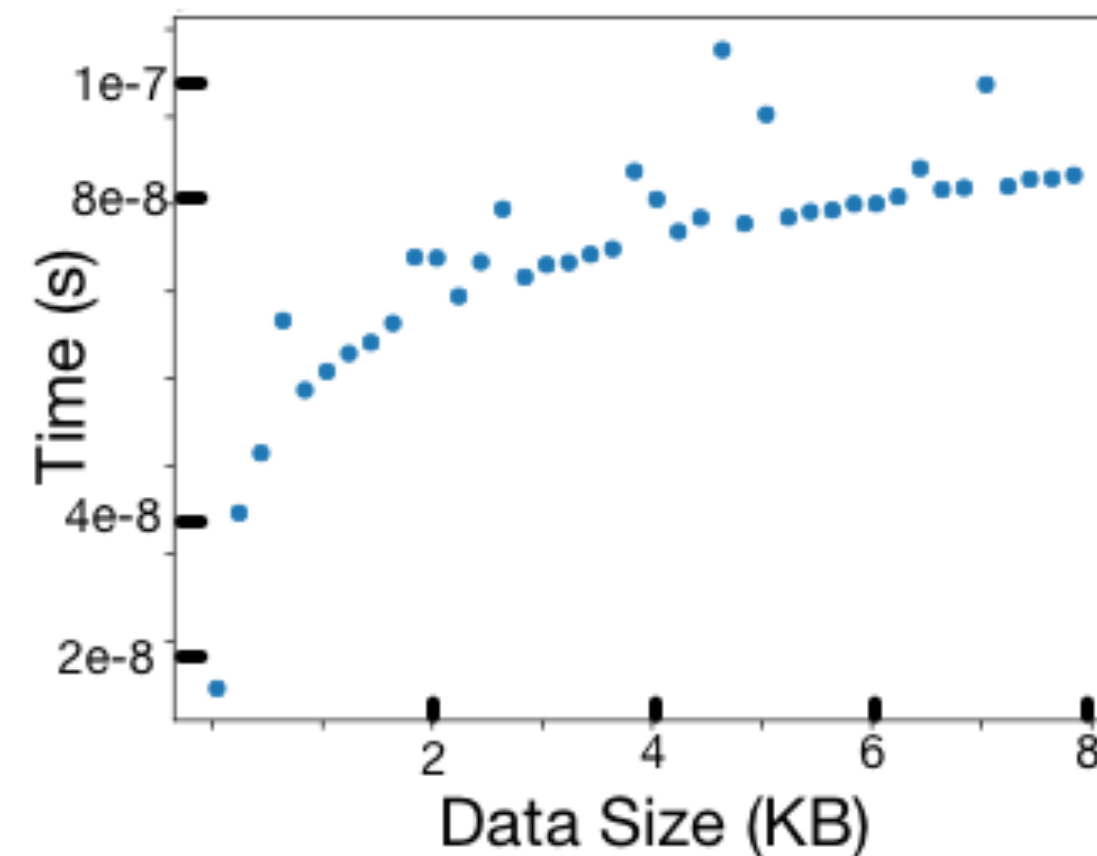
C++

1	11	17	37	51	66	80	94
---	----	----	----	----	----	----	----



Run

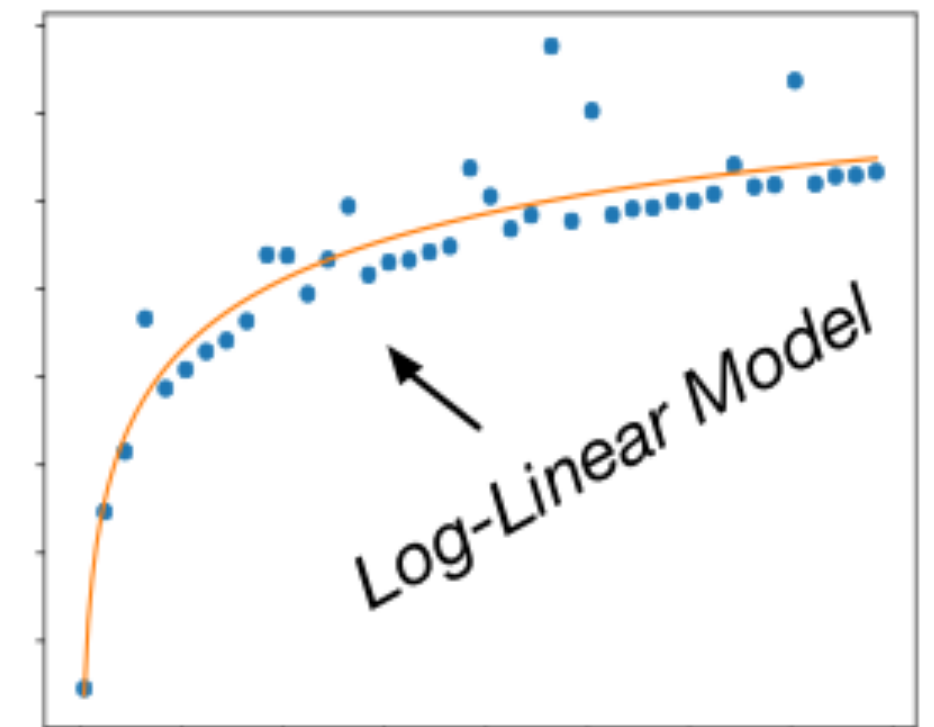
2. BENCHMARK



$f(x)$

Train

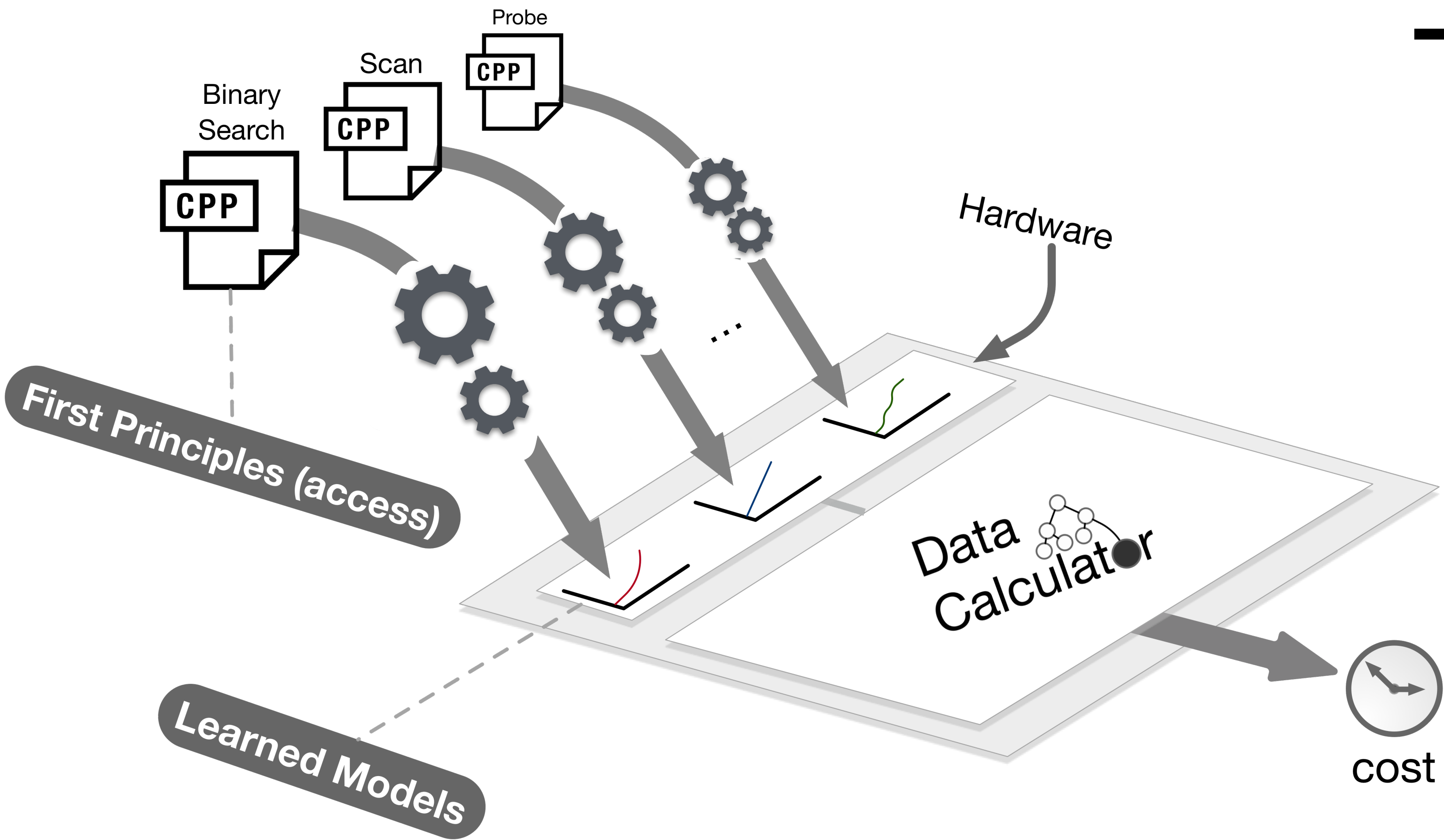
3. FIT MODEL



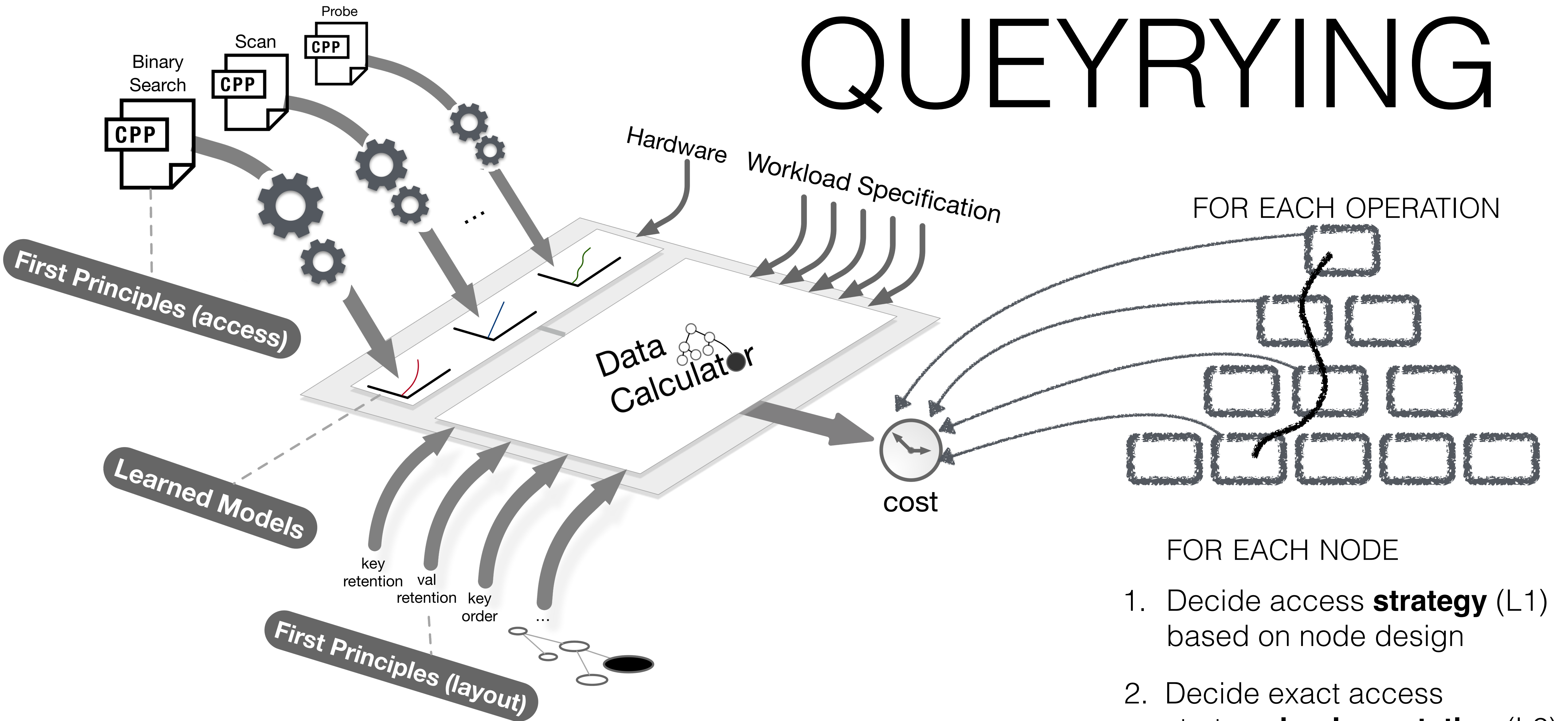
$$f(x) = ax + b \log x + c$$

FOLDING ALGORITHMIC, ENGINEERING, AND H/W, PROPERTIES INTO THE COEFFICIENTS

TRAINING

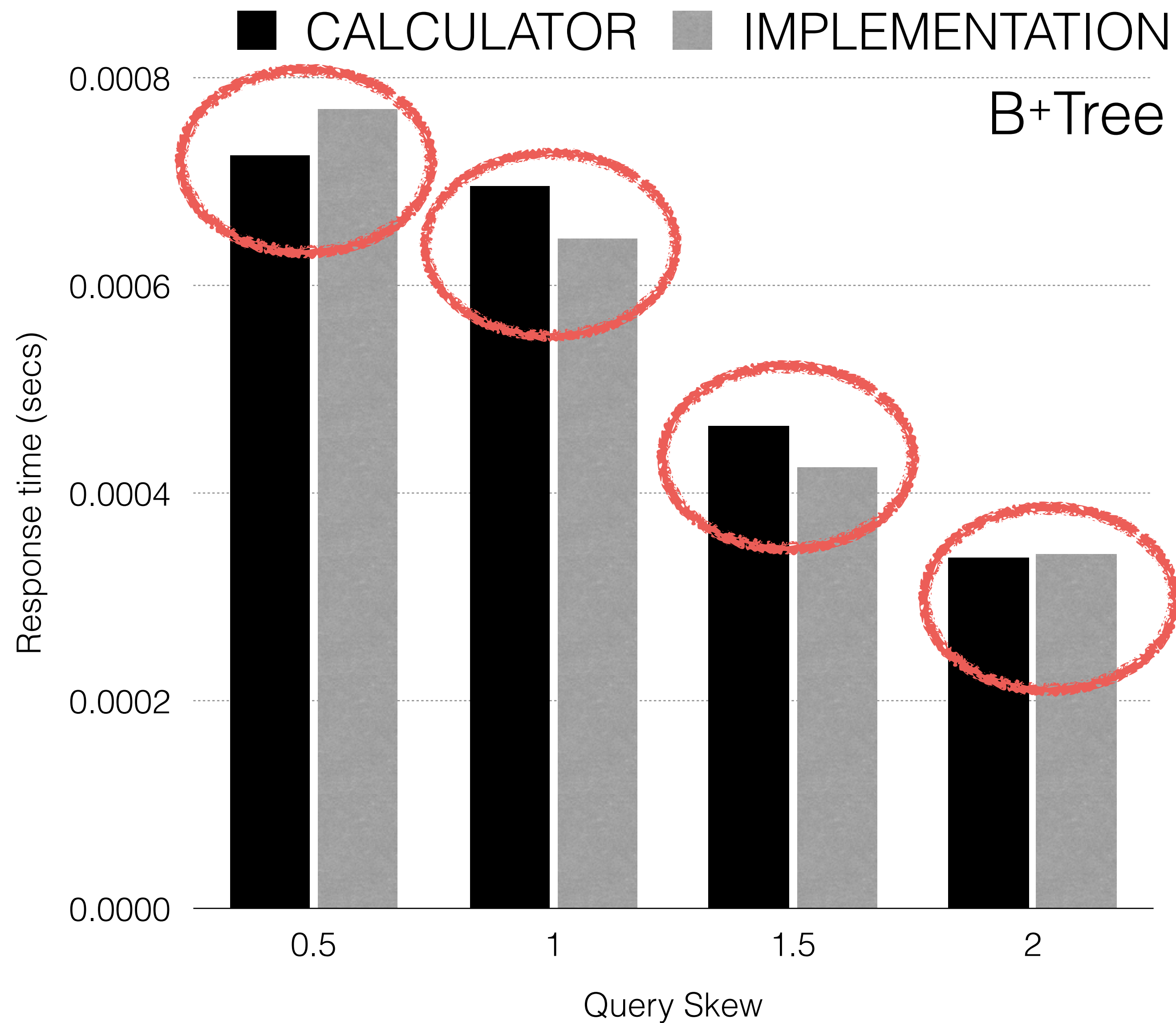


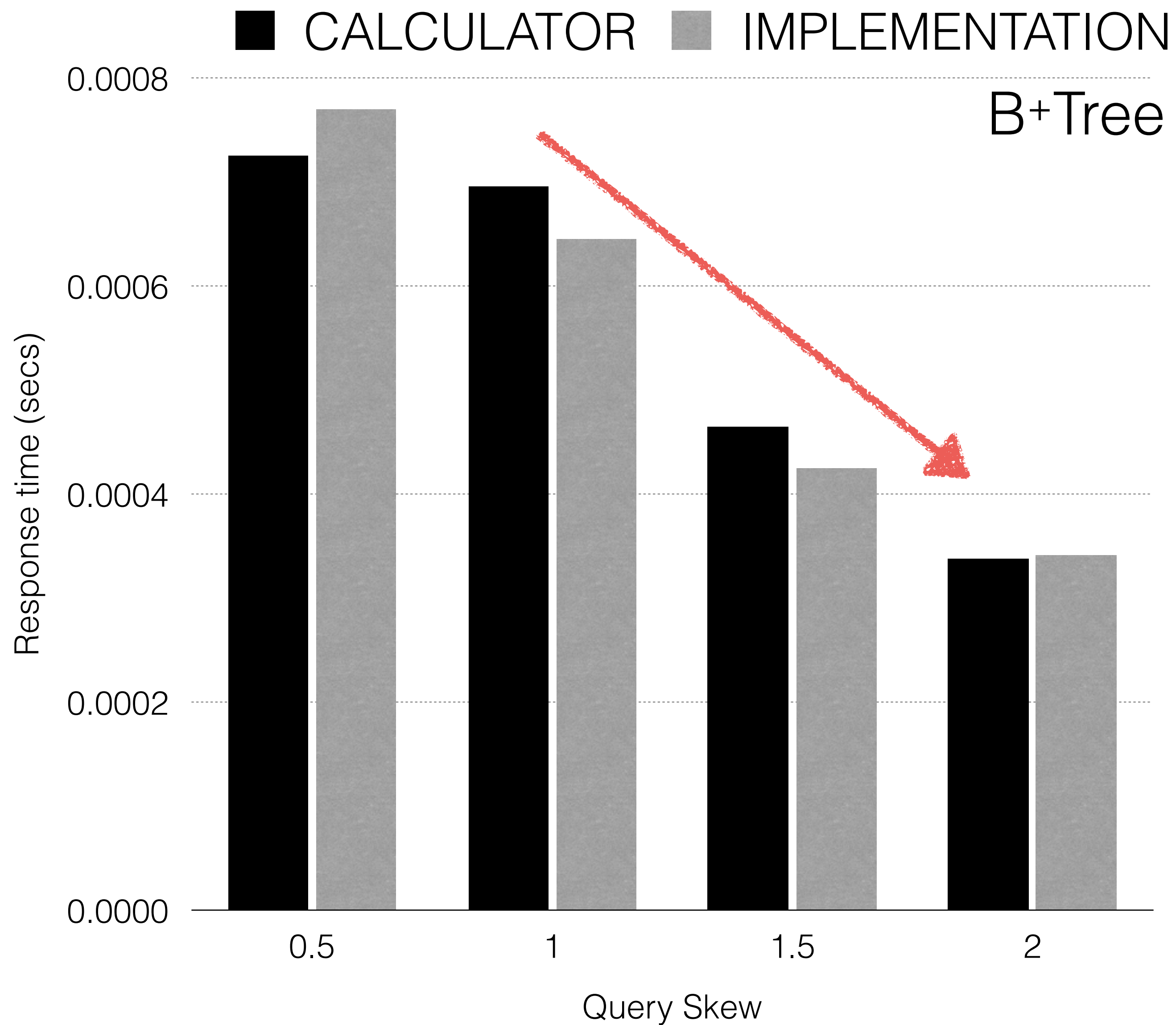
QUERYRYING

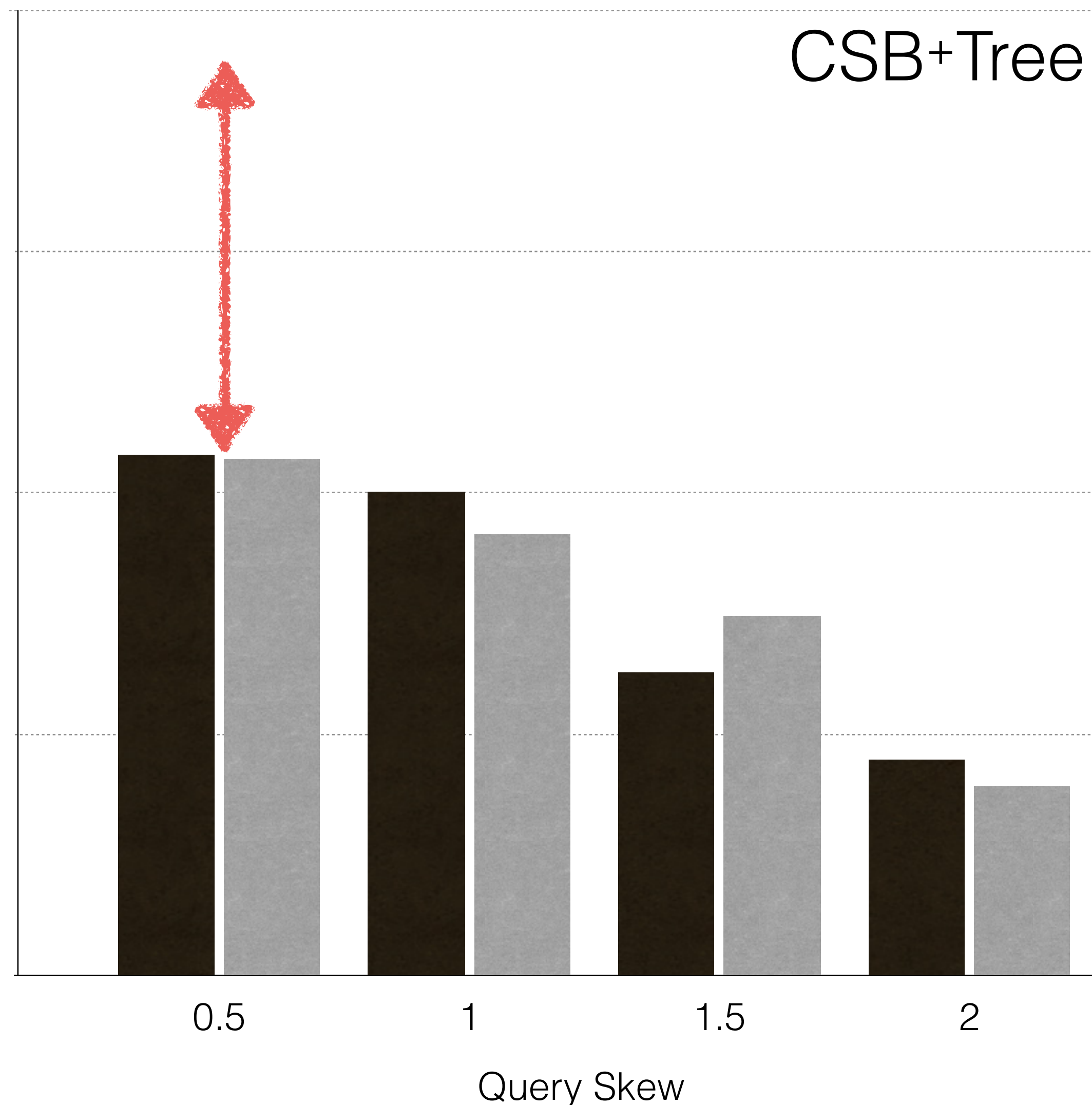
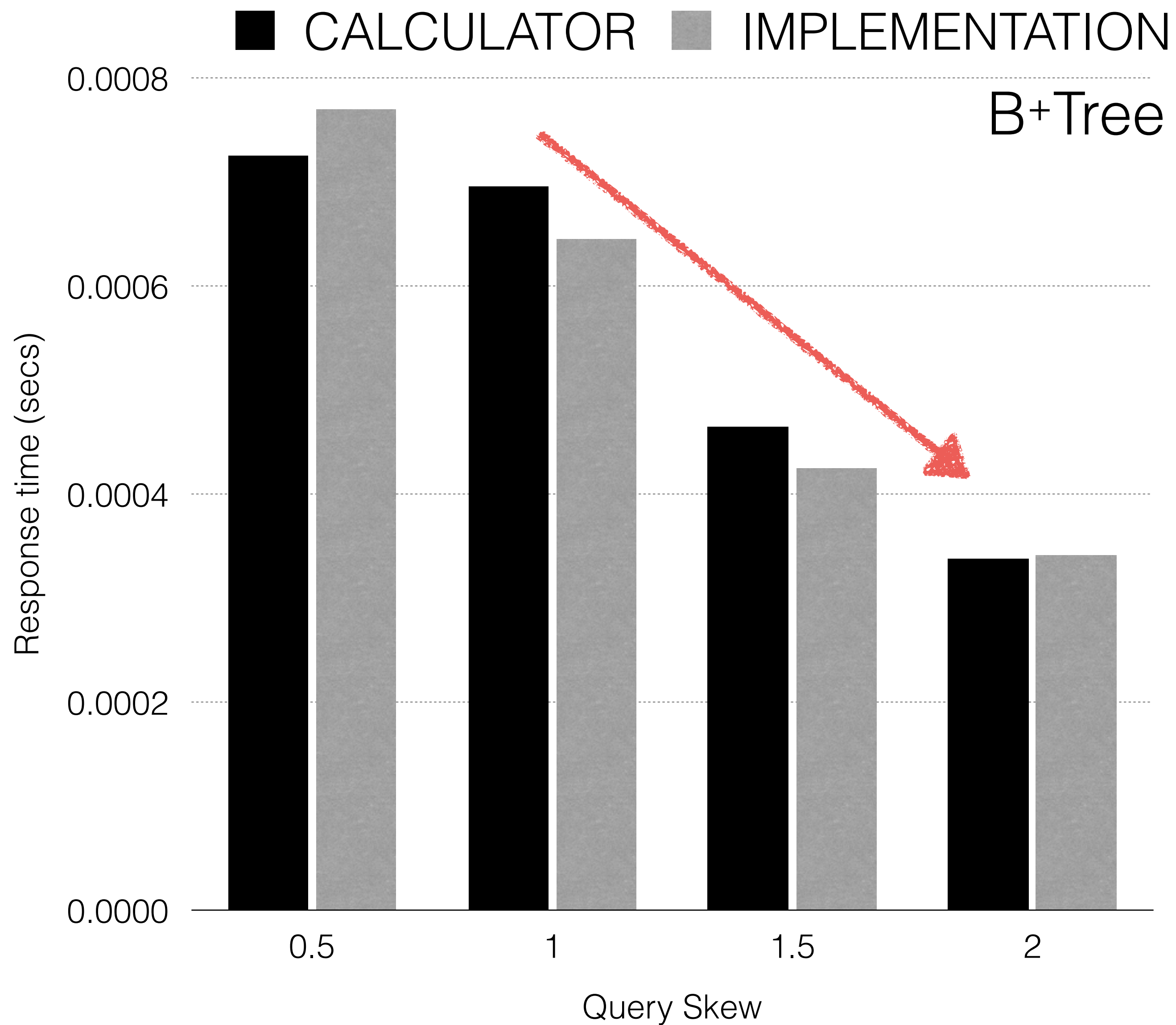


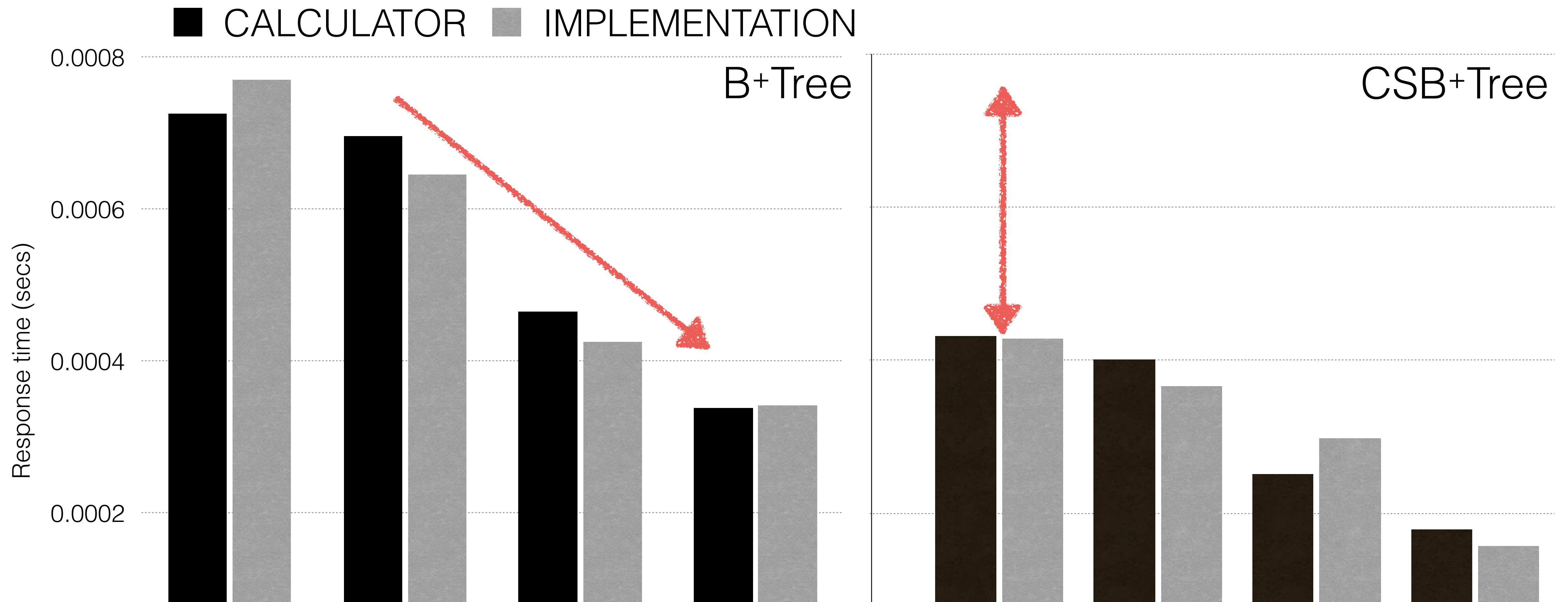
CAN WE COMPUTE PERFORMANCE ACCURATELY?



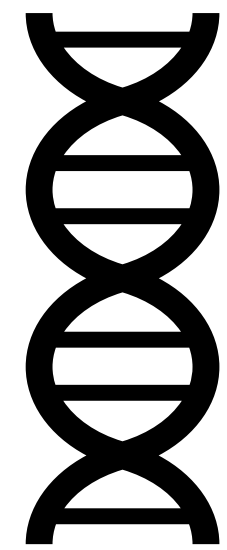




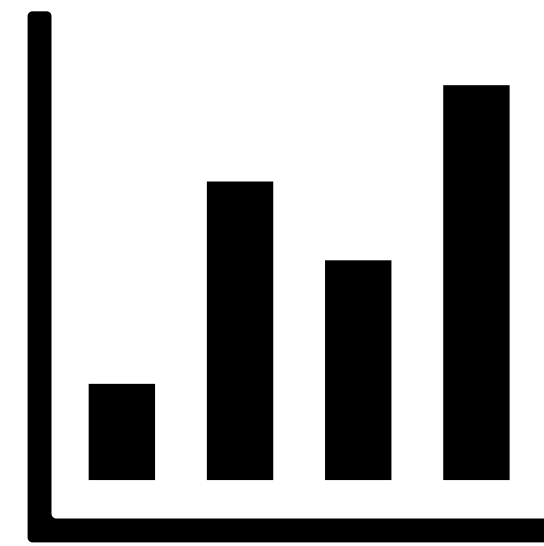




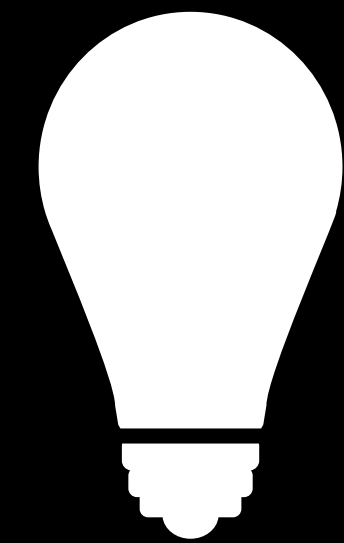
**IF IT WORKS FOR >>1 KNOWN DESIGNS
WE CAN TRUST IT FOR NEW ONES**



DESIGN SPACE



COST ESTIMATION



SEARCH



Key-value Stores/Databases/ML Systems

$>10^{100}$

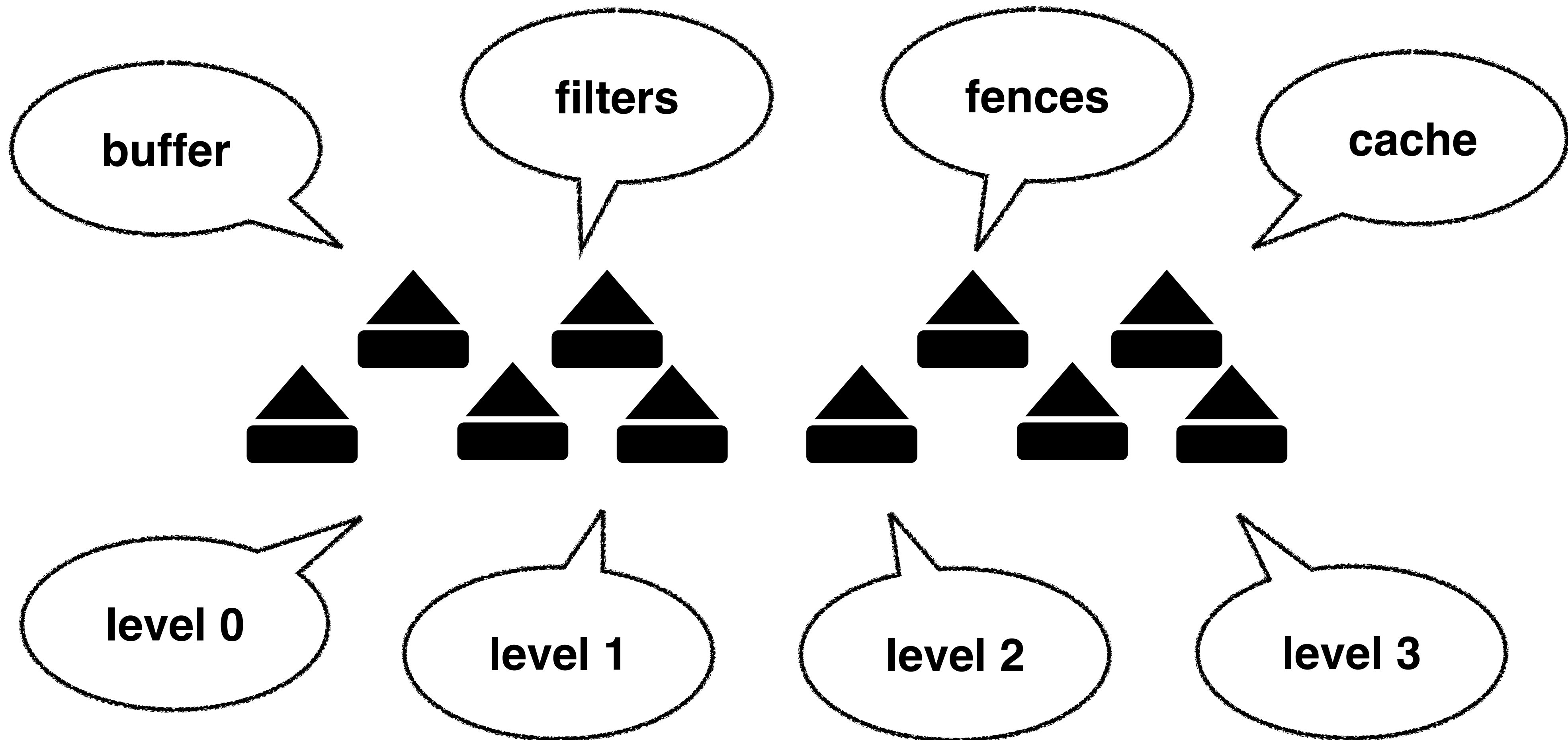
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



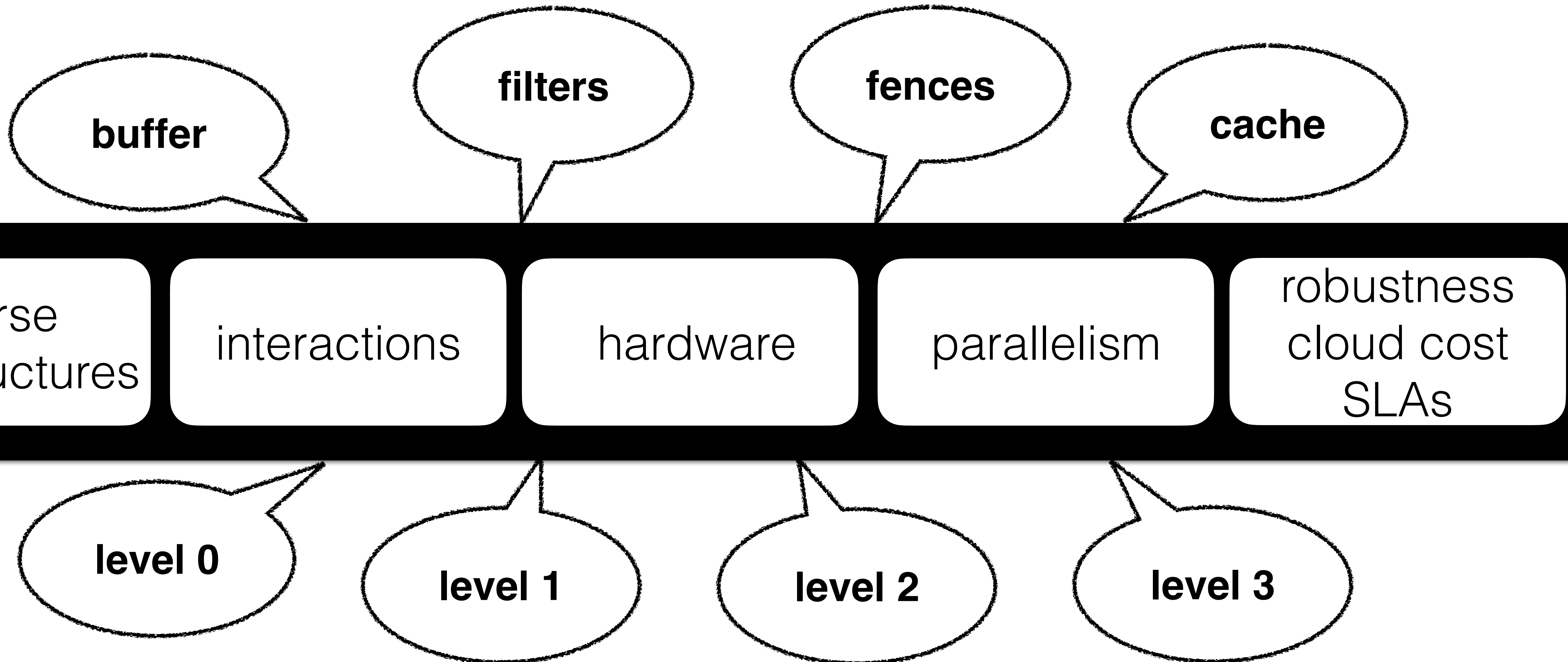
NoSQL systems are the backbone of the BigData and AI era

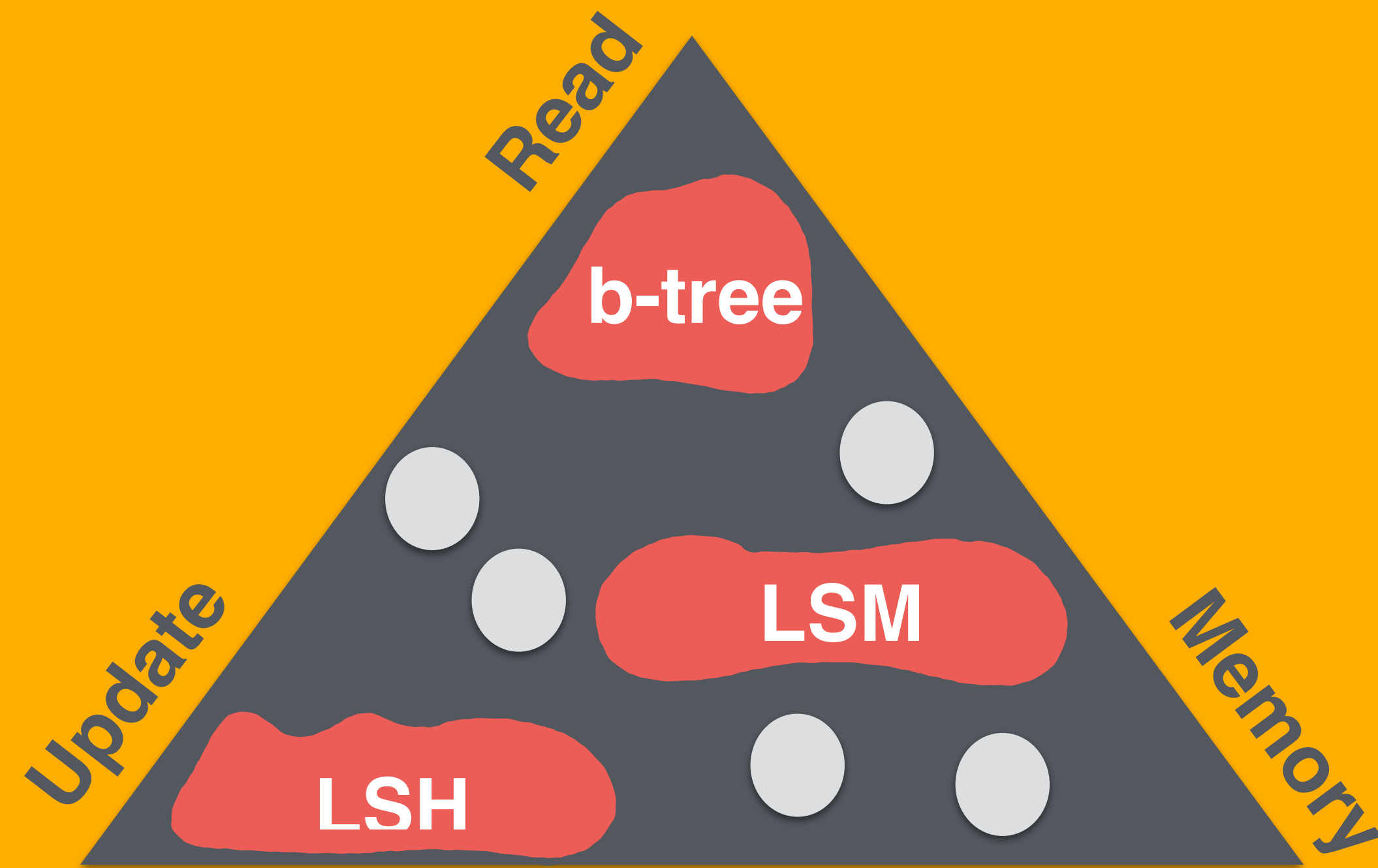
LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE





diverse
data structures

interactions

hardware

parallelism

robustness
cloud cost
SLAs

There exist three core variations of NoSQL KV-stores
But there is a massive possible set of designs

diverse
data structures

interactions

hardware

parallelism

robustness
cloud cost
SLAs

Requirements/Goals

Context

data & queries

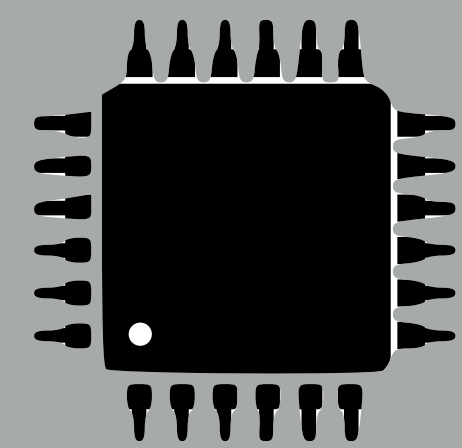
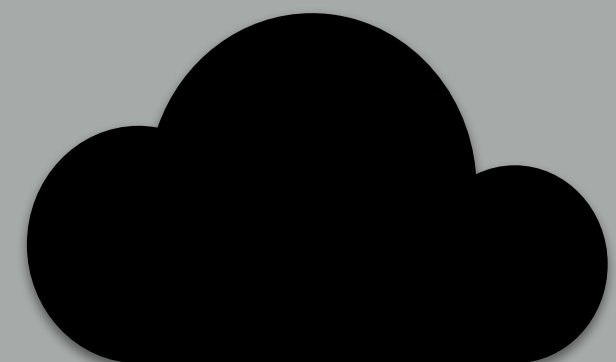


performance

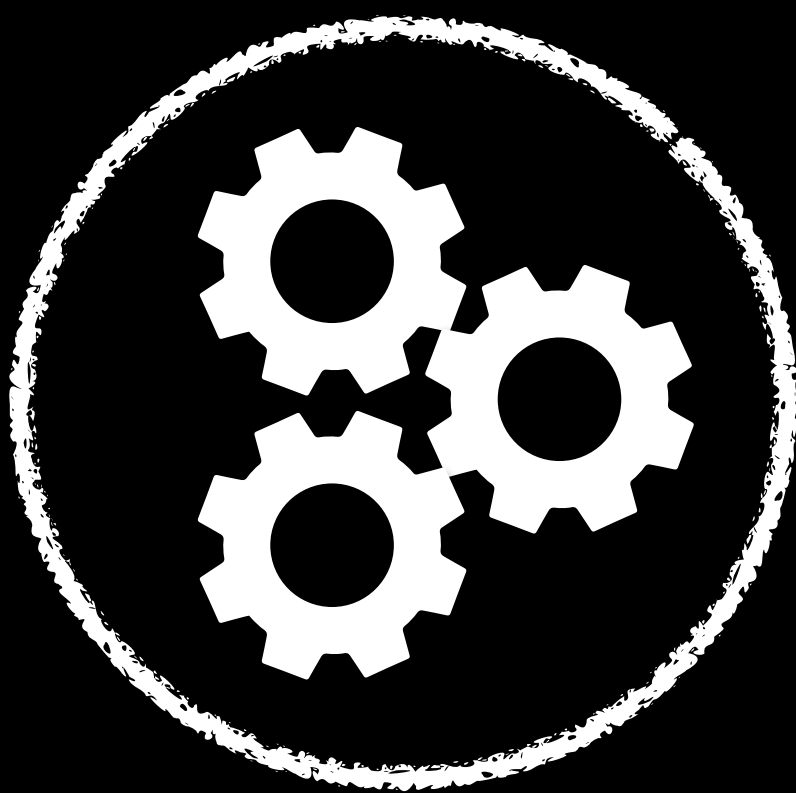
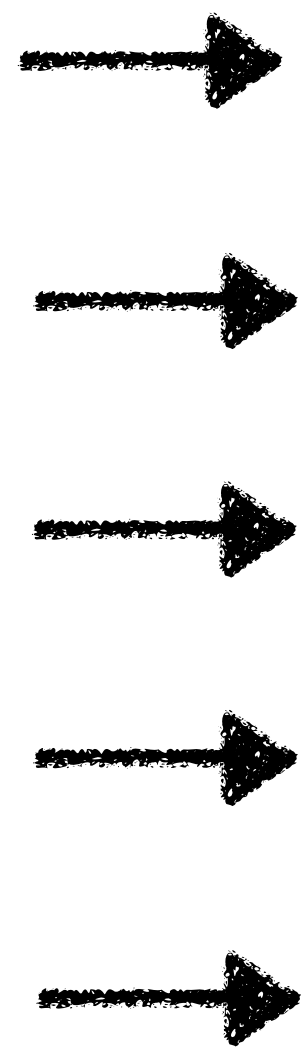


budget

\$\$\$



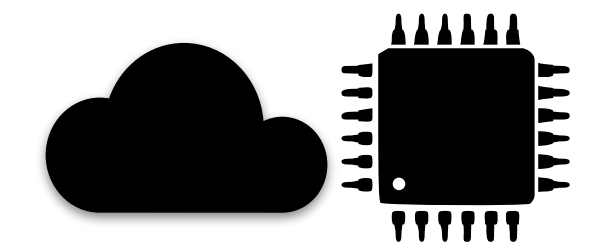
SLA

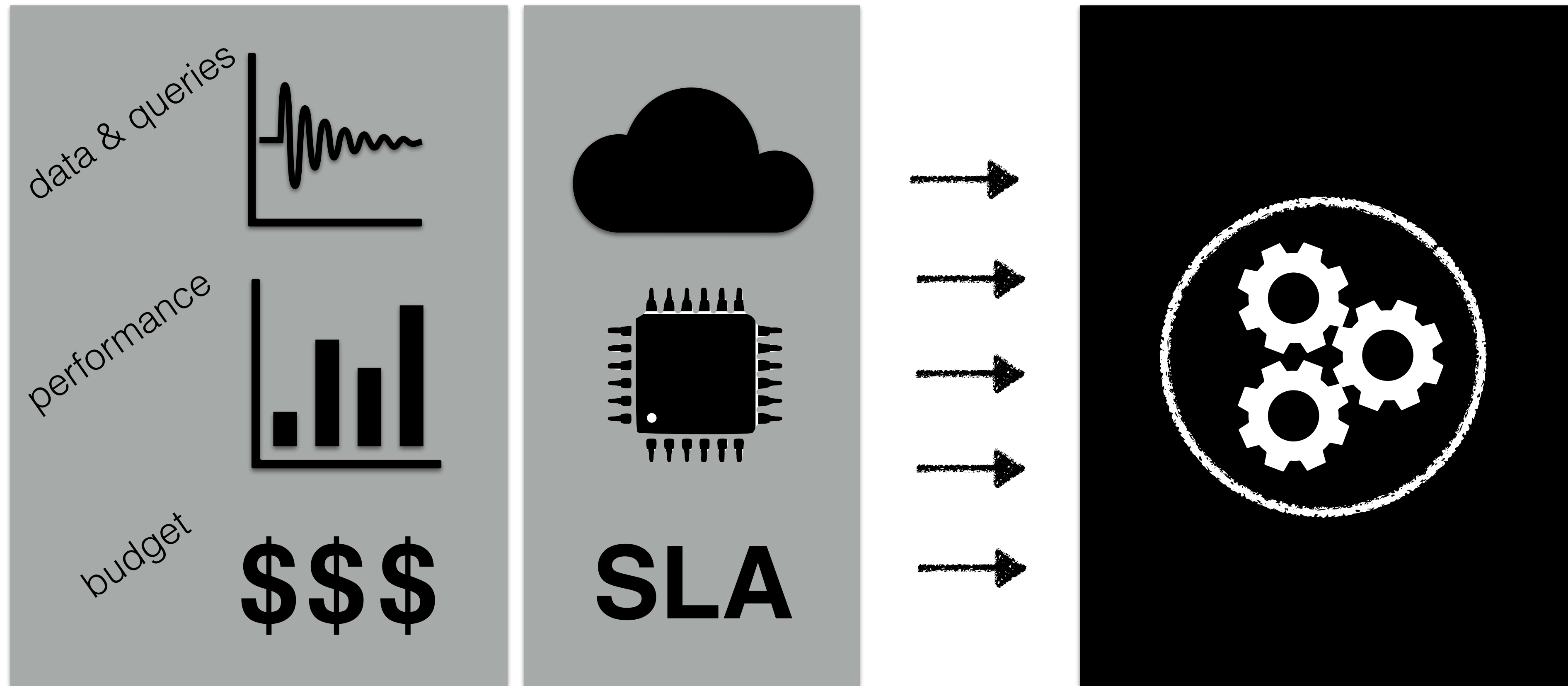


best

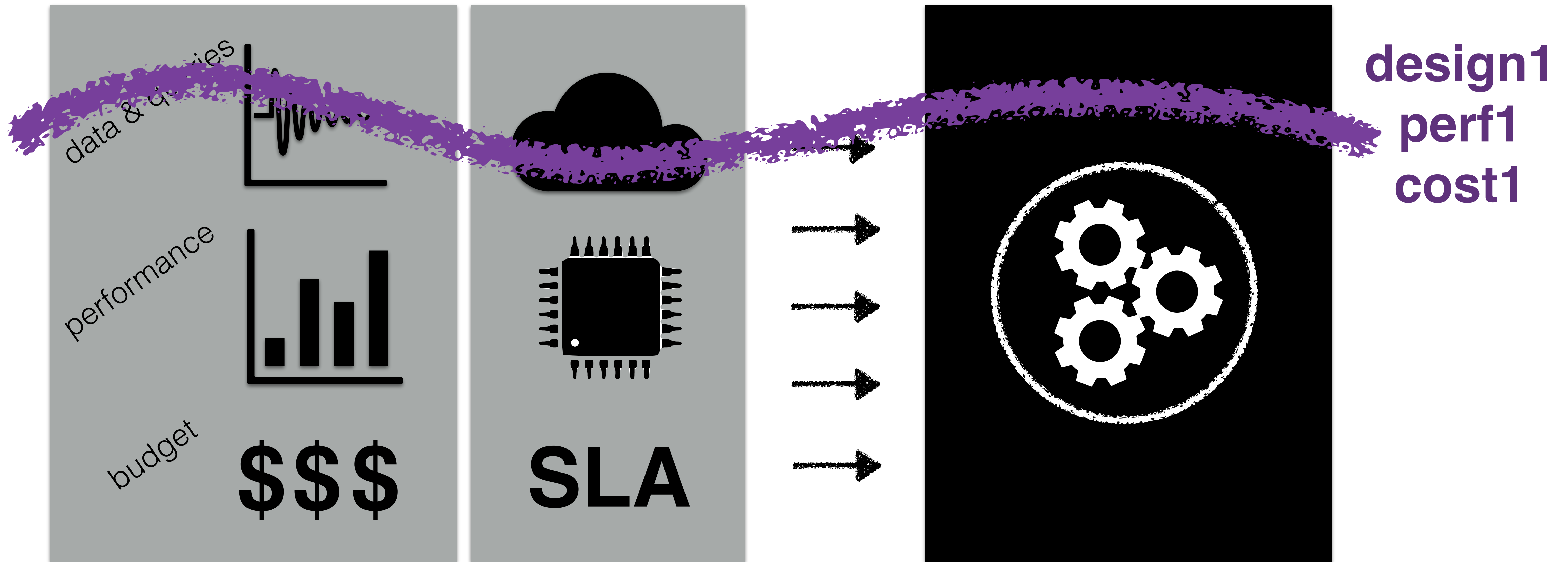


DATA
SYSTEM
design & code

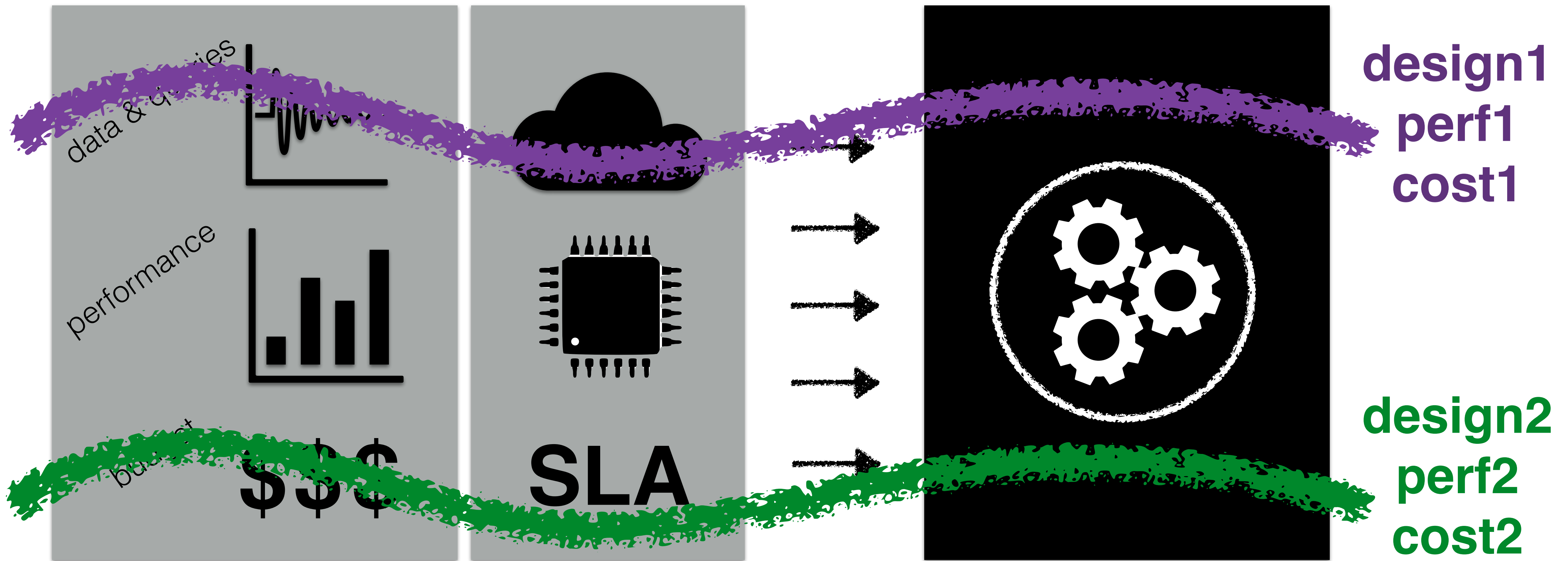




 **what-if reasoning**

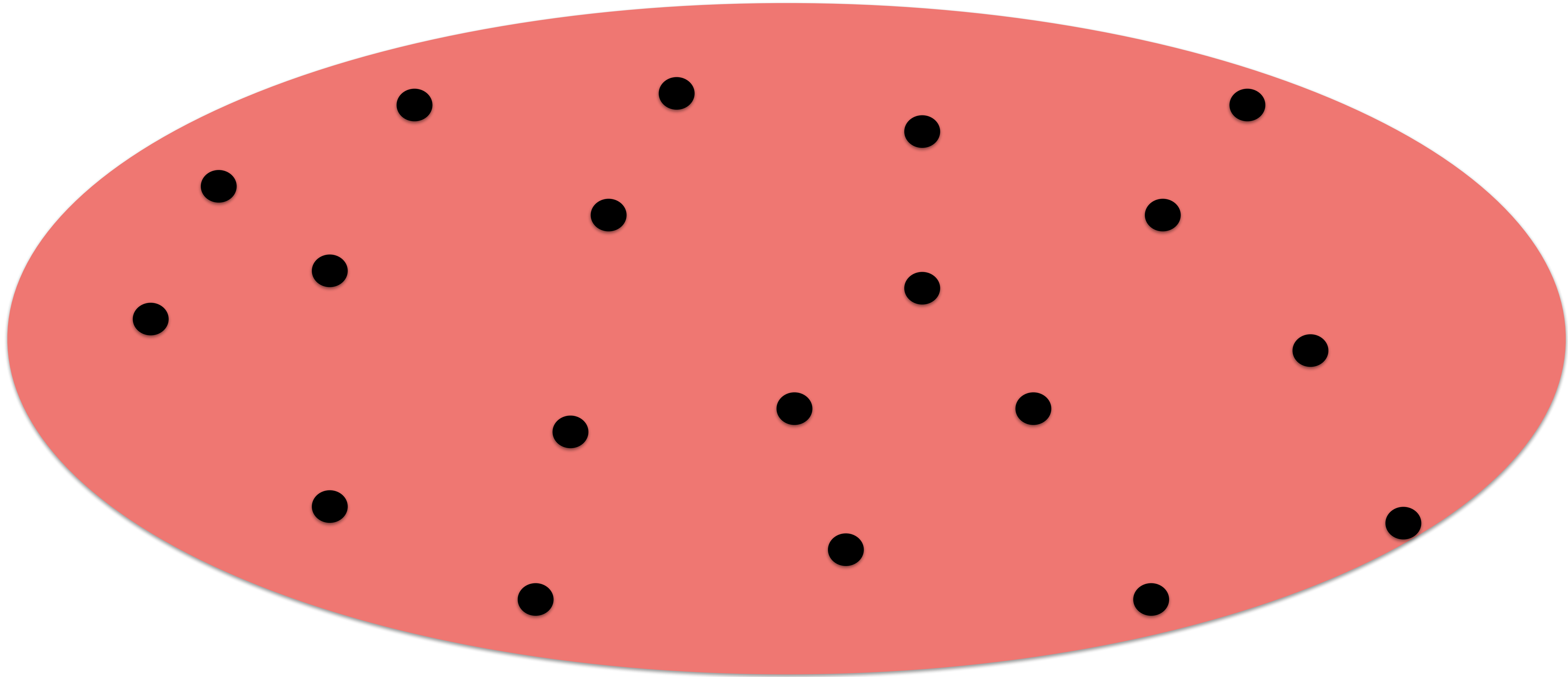


 **what-if reasoning**



 **what-if reasoning**

shrinking the massive design space



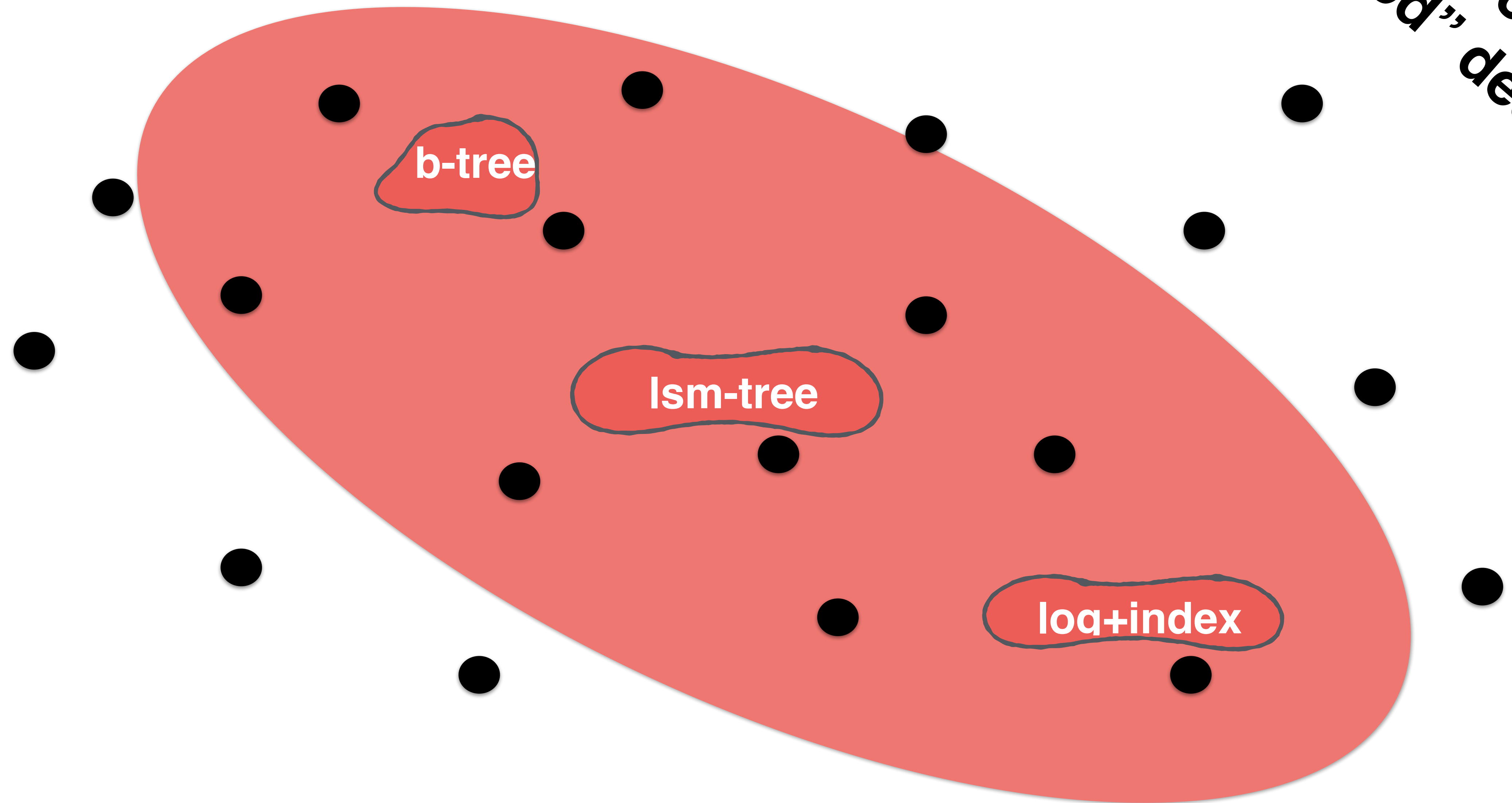
manually
selecting
“good” designs

b-tree

lsm-tree

log+index

manually
selecting
“good” designs



manually
selecting
“good” designs

b-tree

lsm-tree

log+index

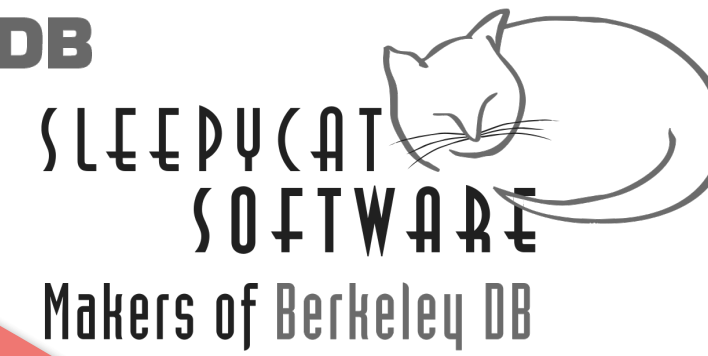
manually
selecting
“good” designs



WIREDTIGER



FOUNDATIONDB



Google
BigTable



APACHE
HBASE

riak



FASTER

b-tree

lsm-tree

log+index

manually
selecting
“good” designs



WIREDTIGER



FOUNDATIONDB



Google
BigTable



cassandra



APACHE
HBASE

riak



b-tree

lsm-tree

log+index



Monkey
SIGMOD 2017



Dostoevsky
SIGMOD 2018



Wacky
SIGMOD 2019



Stacked Filters
SIGMOD 2020

manually
selecting
“good” designs



WIREDTIGER



FOUNDATIONDB



Google
BigTable



cassandra



APACHE
HBASE

riak



b-tree

lsm-tree

log+index



Monkey
SIGMOD 2017



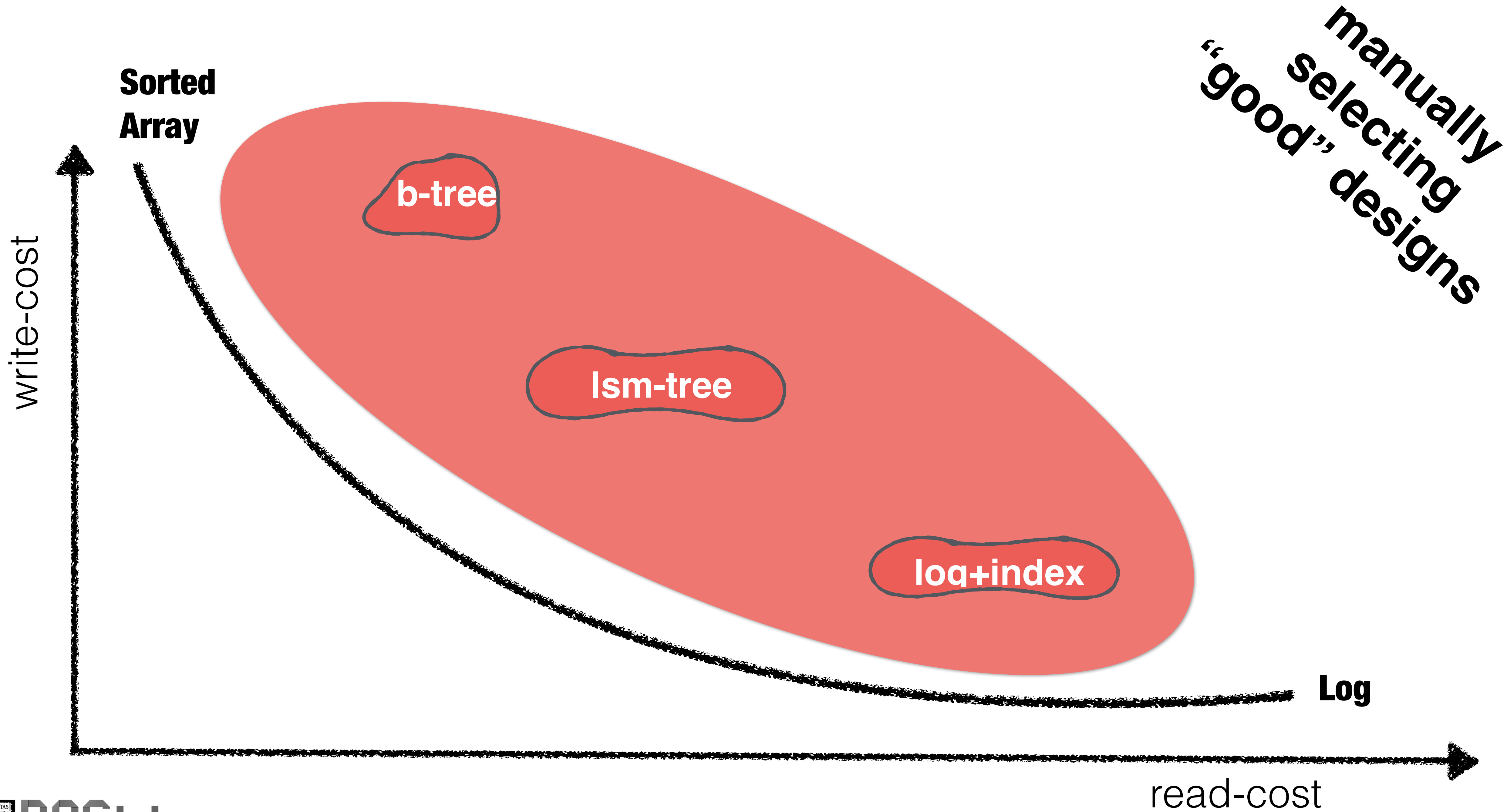
Dostoevsky
SIGMOD 2018

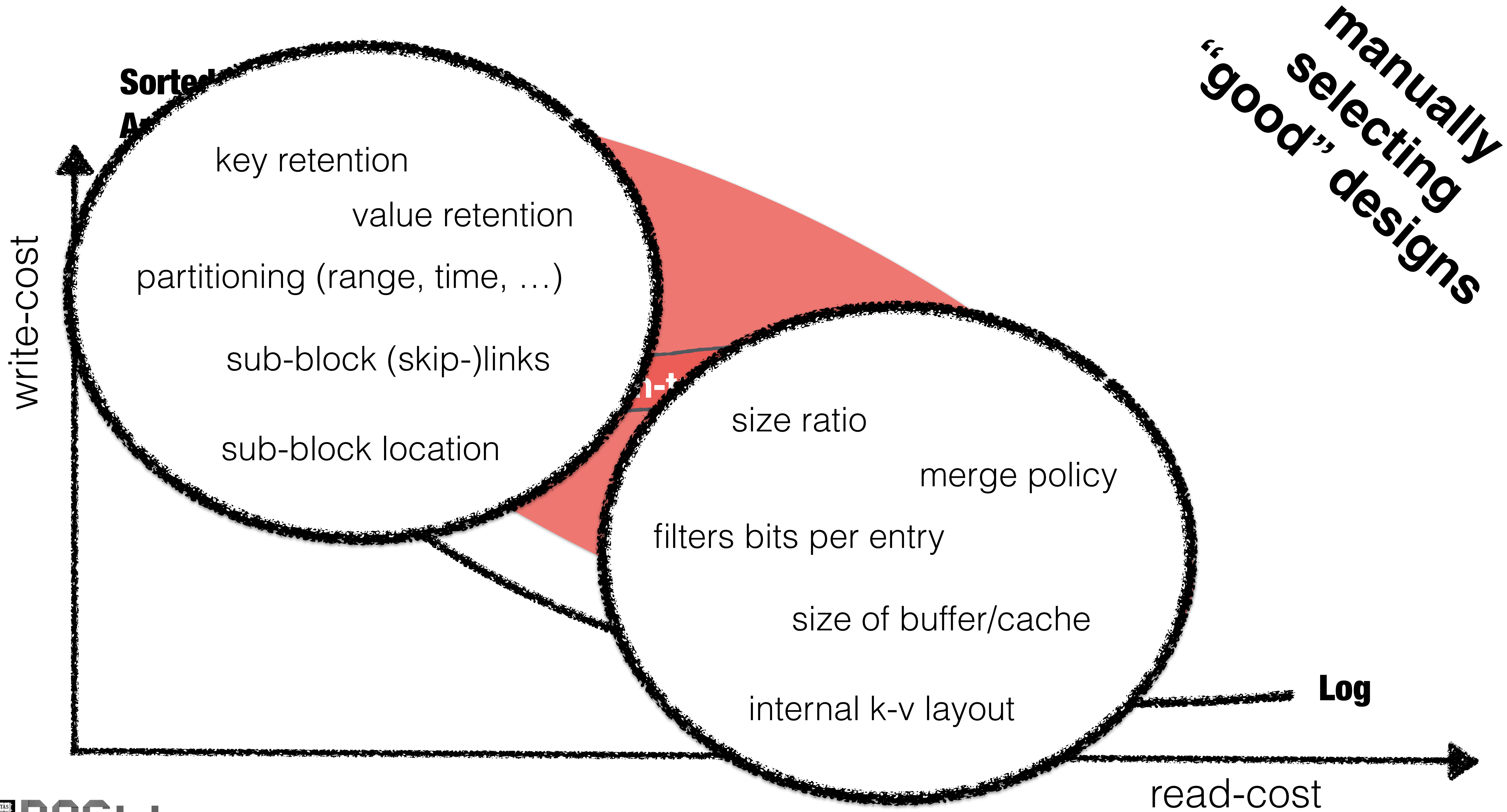


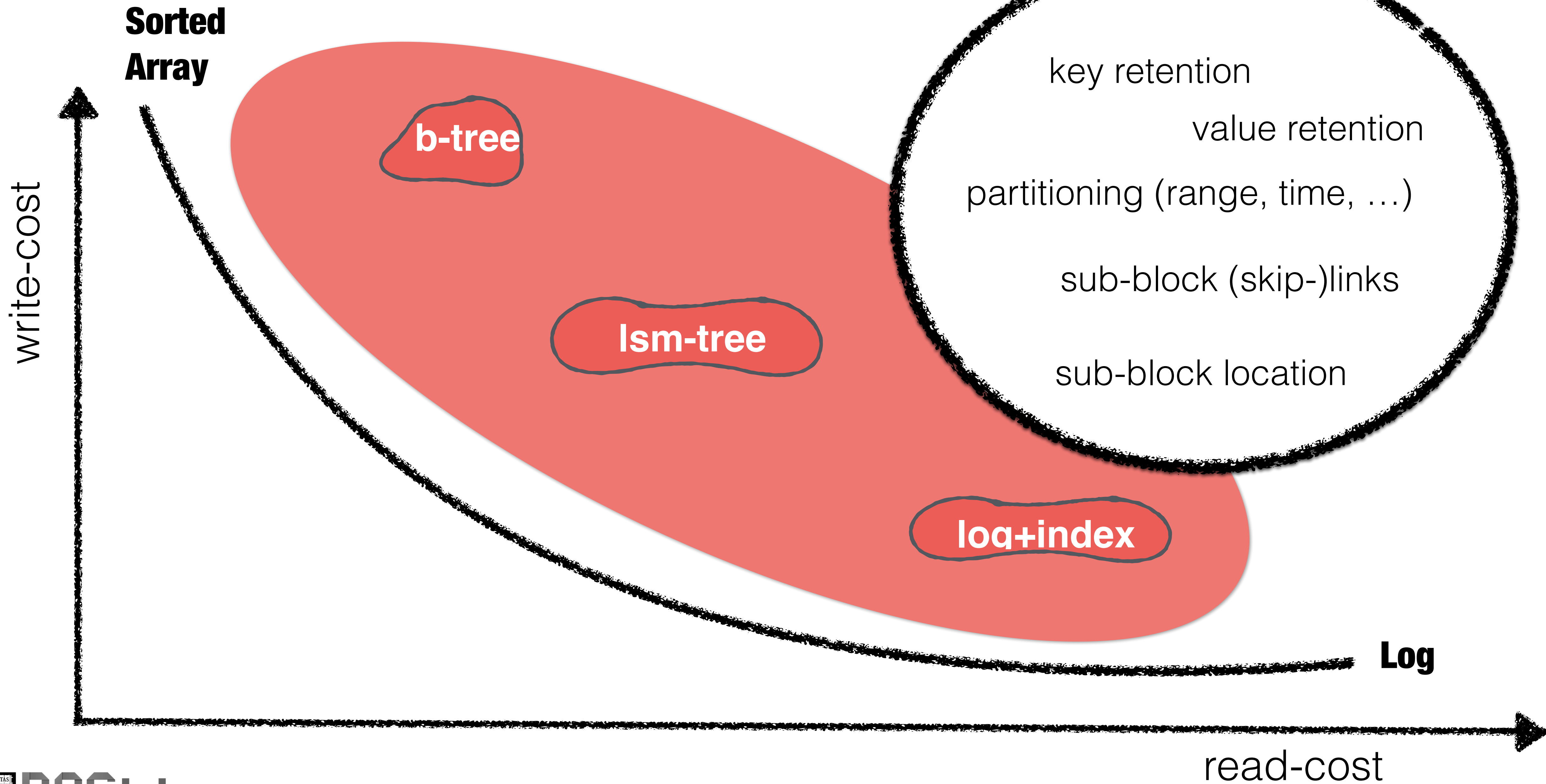
Wacky
SIGMOD 2019

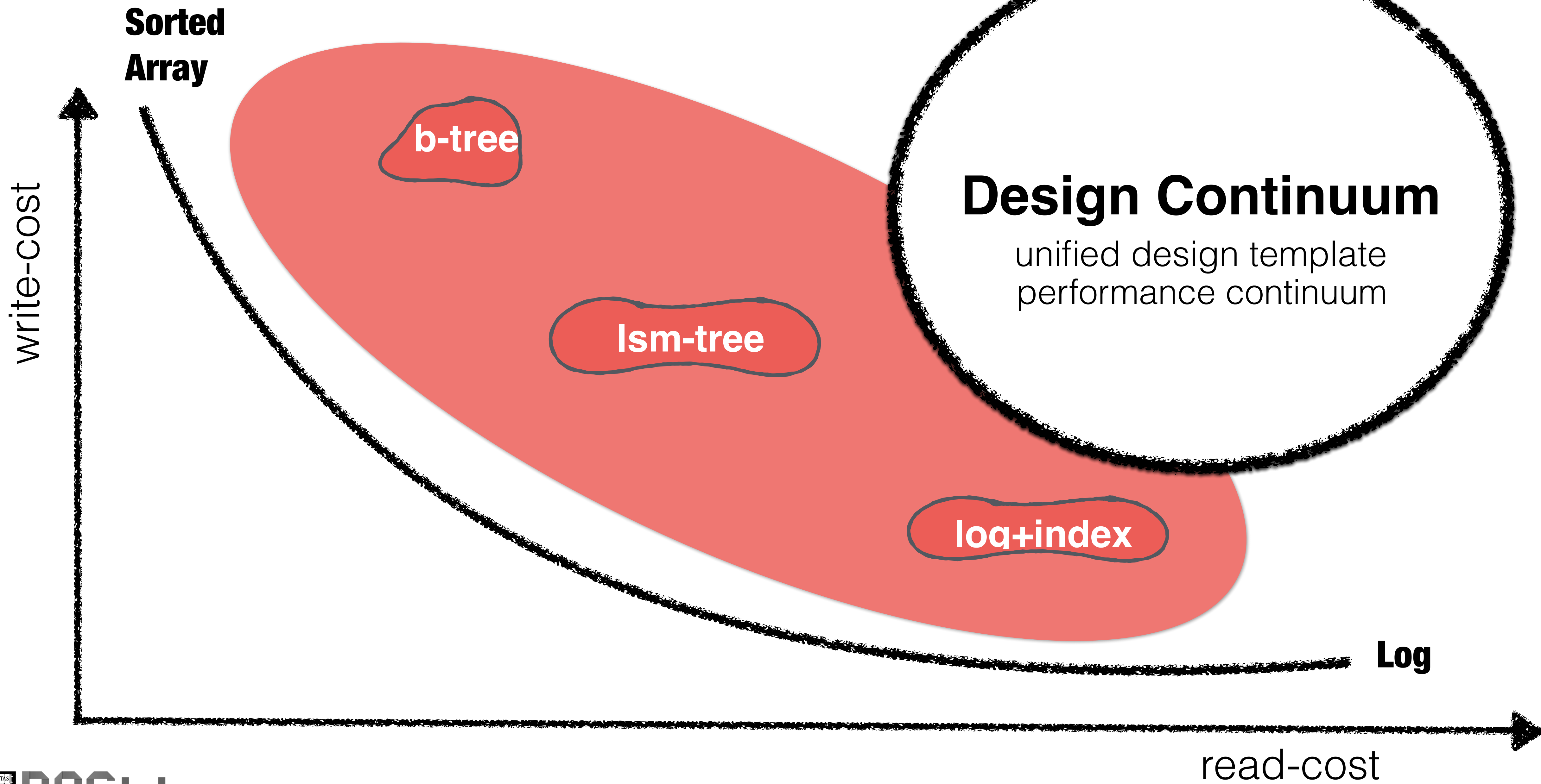


Stacked Filters
SIGMOD 2020









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
ALGORITHMIC ABSTRACTIONS	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 <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)
	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)
		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			
	Parallelism						

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
ALGORITHMIC ABSTRACTIONS	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 <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)
	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)
		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		
	Parallelism						

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
ALGORITHMIC ABSTRACTIONS	Design and hardware specification <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, ..]	[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 <i>derived with empirically verified rules</i>	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		
	Parallelism						

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

Cosine
@PVLDB2022

Cloud-cost
Optimized

Self
Designing

Key-value
Store

ALGORITHMIC ABSTRACTIONS

LAYOUT PRIMITIVES

		Design Abstractions of Template	Type/Domain	Example templates for diverse data structures			
				LSM variants	B-Tree variants	LSH 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 <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)
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)
		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			
Parallelism							

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

Cosine
@PVLDB2022

Cloud-cost
Optimized

Self
Designing

Key-value
Store

ALGORITHMIC ABSTRACTIONS ↔ LAYOUT PRIMITIVES

		Design Abstractions of Template	Type/Domain	Example templates for diverse data structures			
				LSM variants	B-Tree variants	LSH 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 <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)
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)
		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			
Parallelism							

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

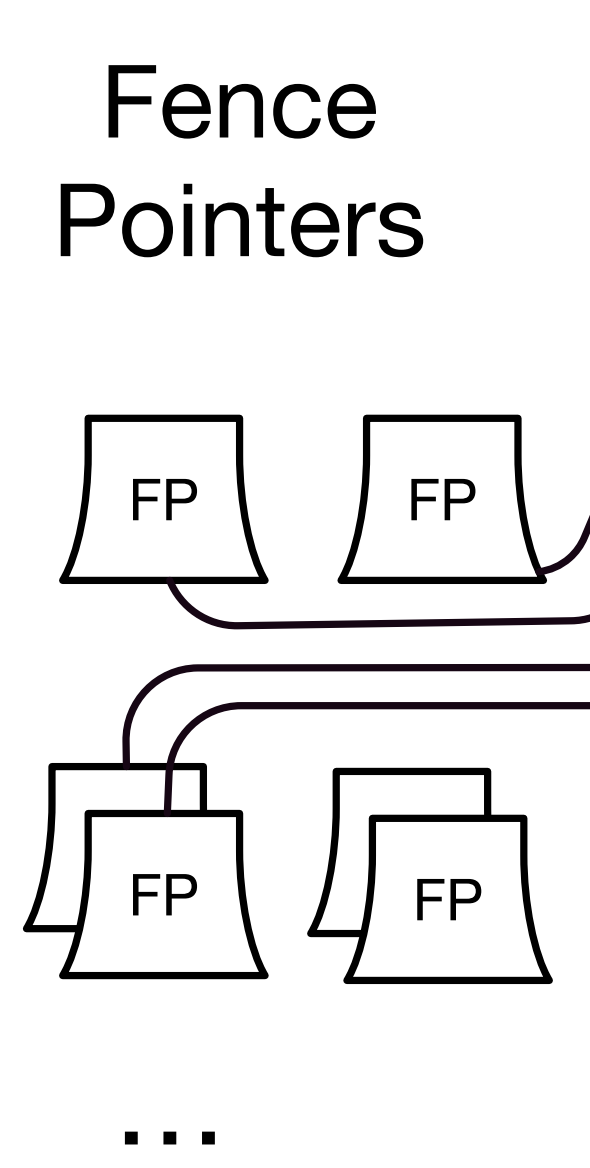
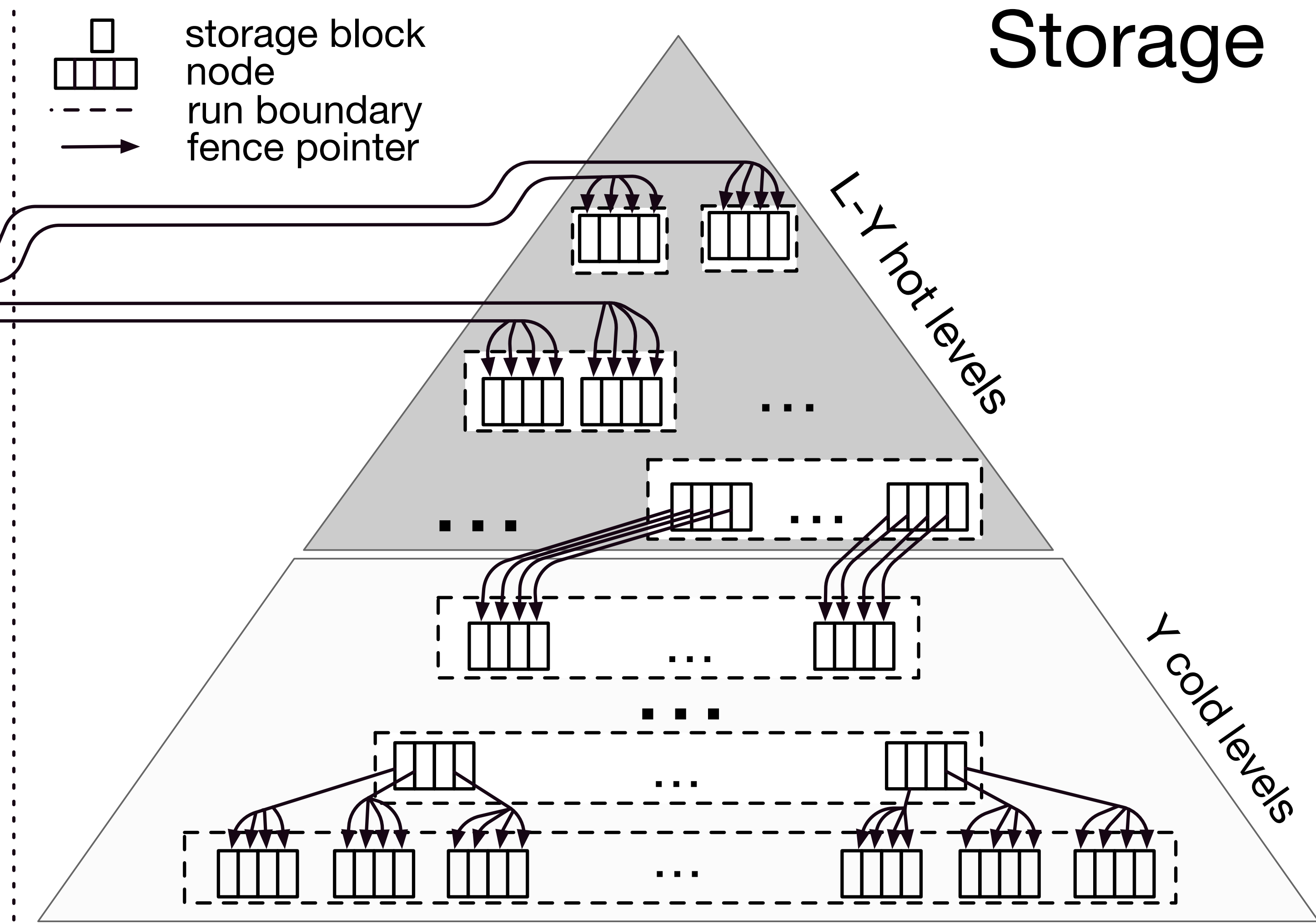
Design Continuums @CDIR2019

Design Continuums @CDIR2019

[illegible]

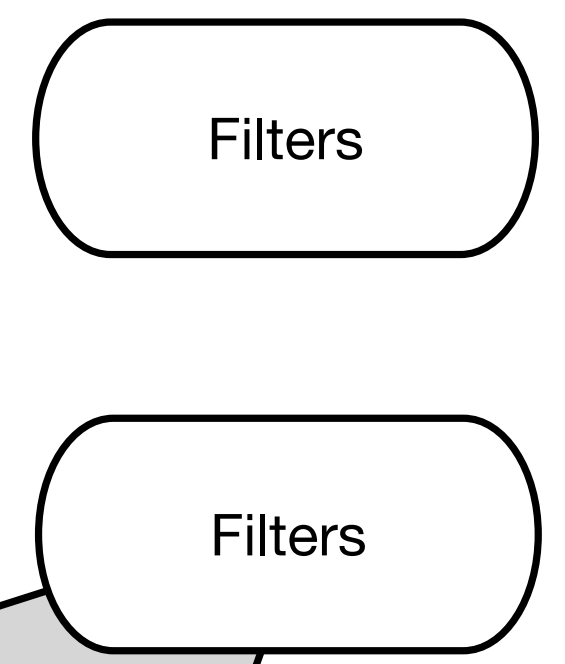
unified design storage engine template

Storage

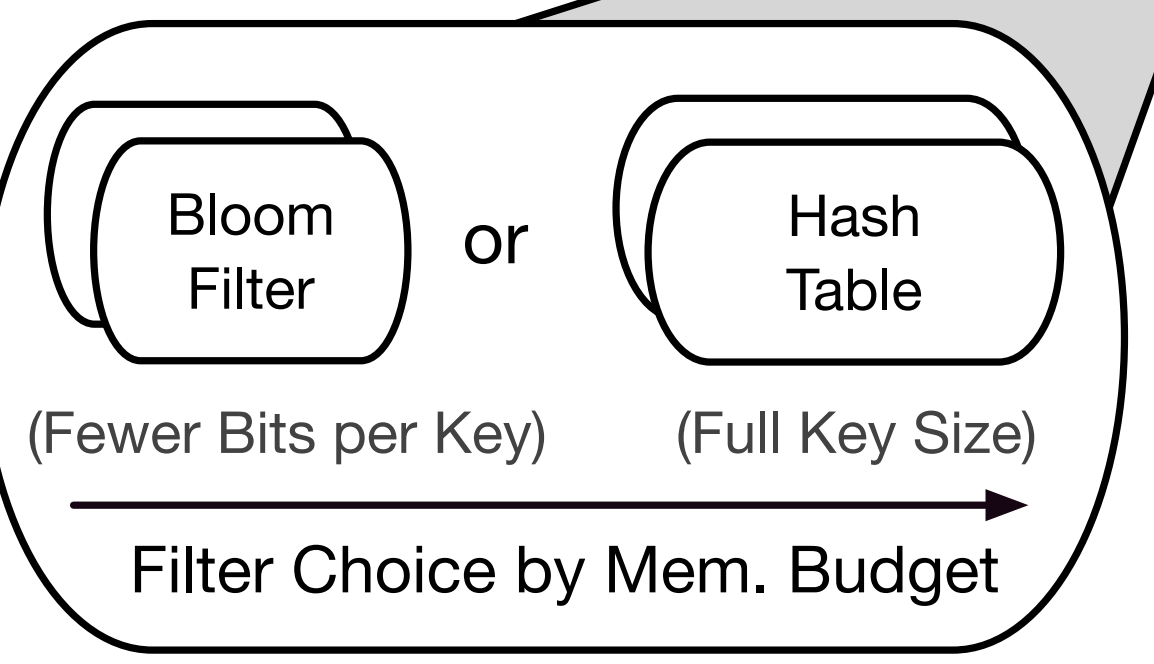
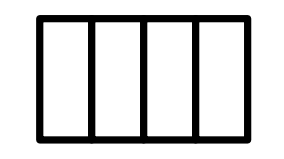


storage block
node
run boundary
fence pointer

Filters



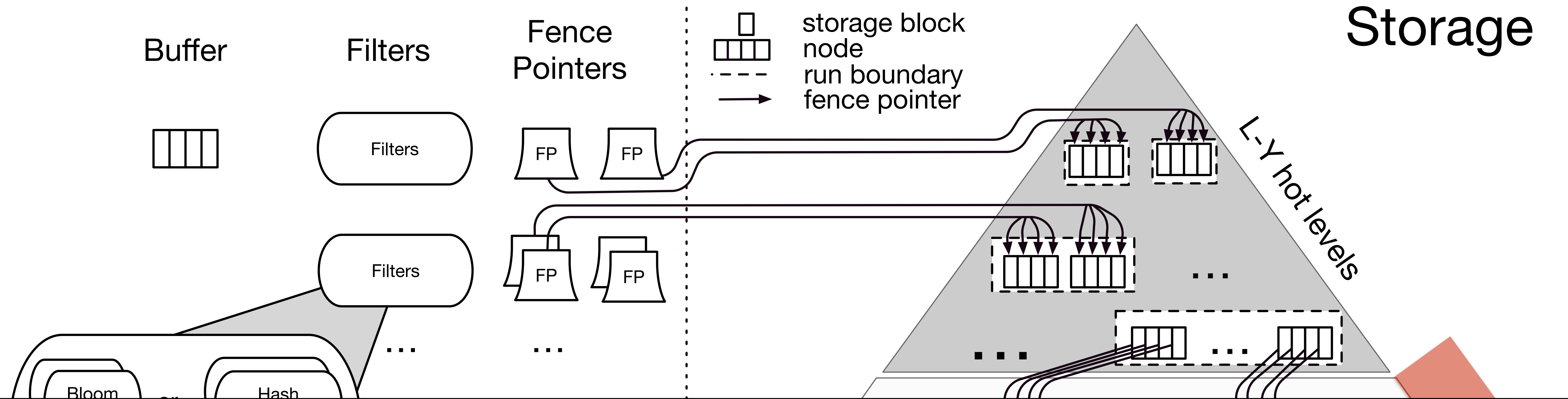
Buffer



Memory

unified design
storage engine template

Storage



**analytical
I/O model**

**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



workload
budget
perf.

**analytical
I/O model**

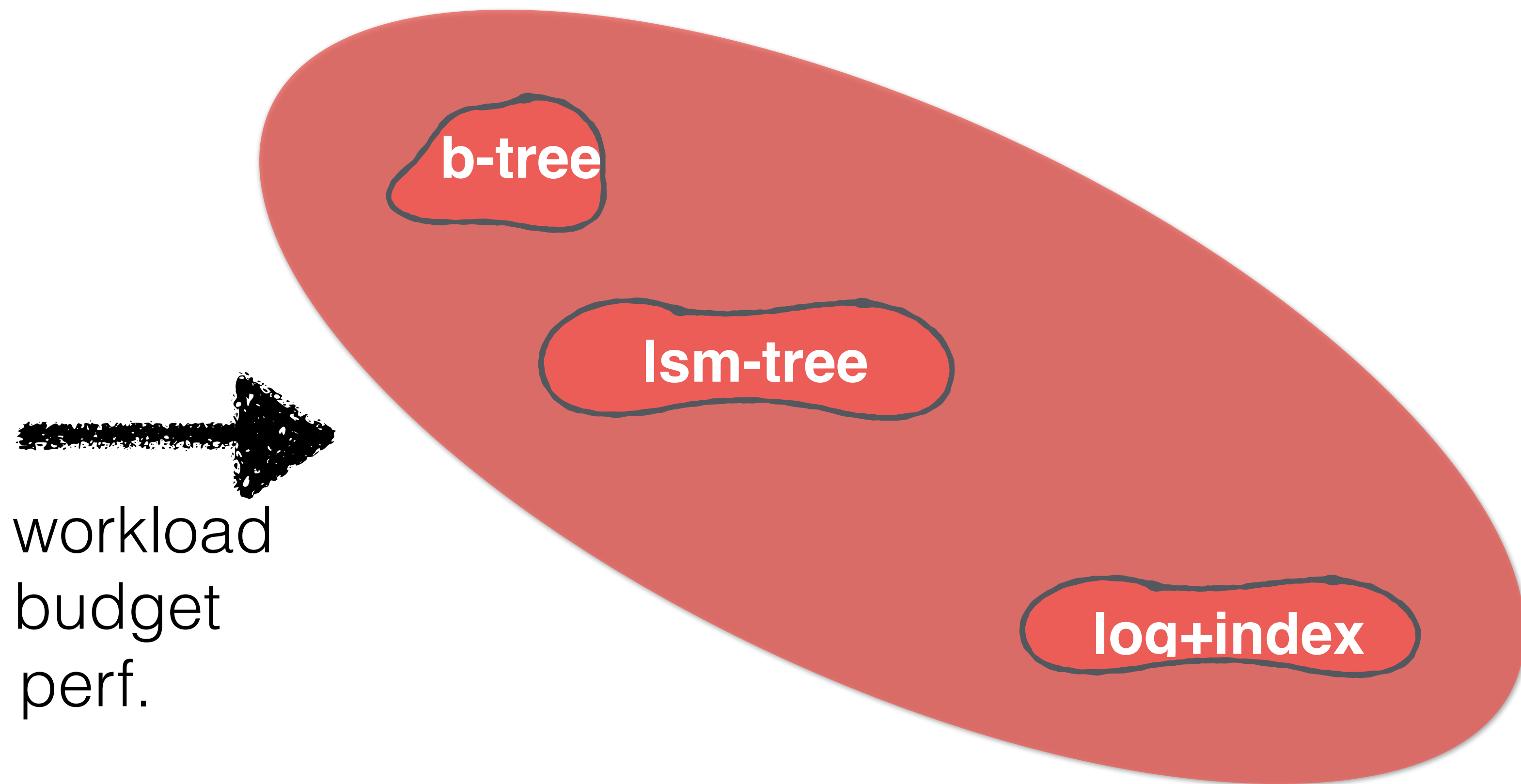
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

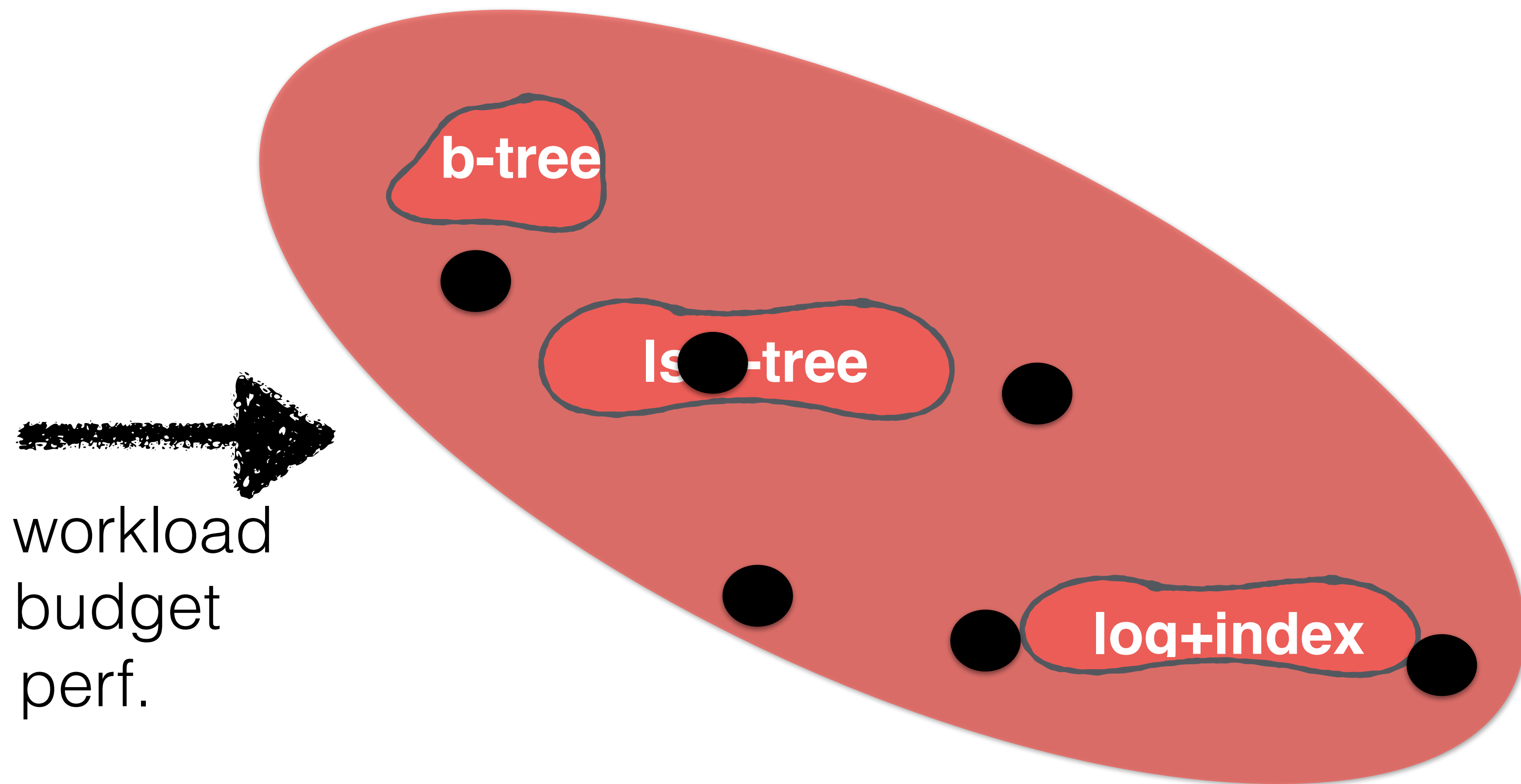
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

unified
closed form

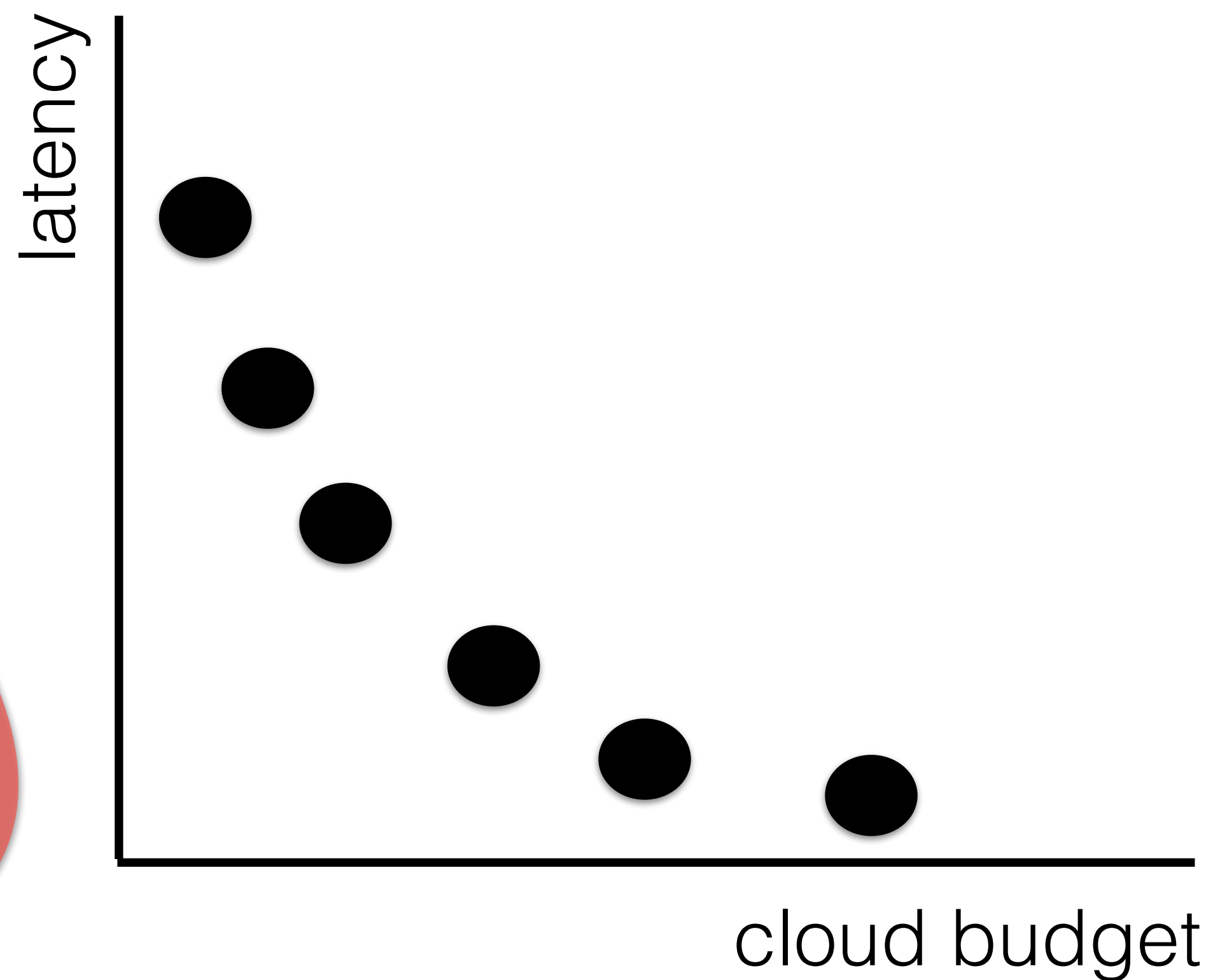
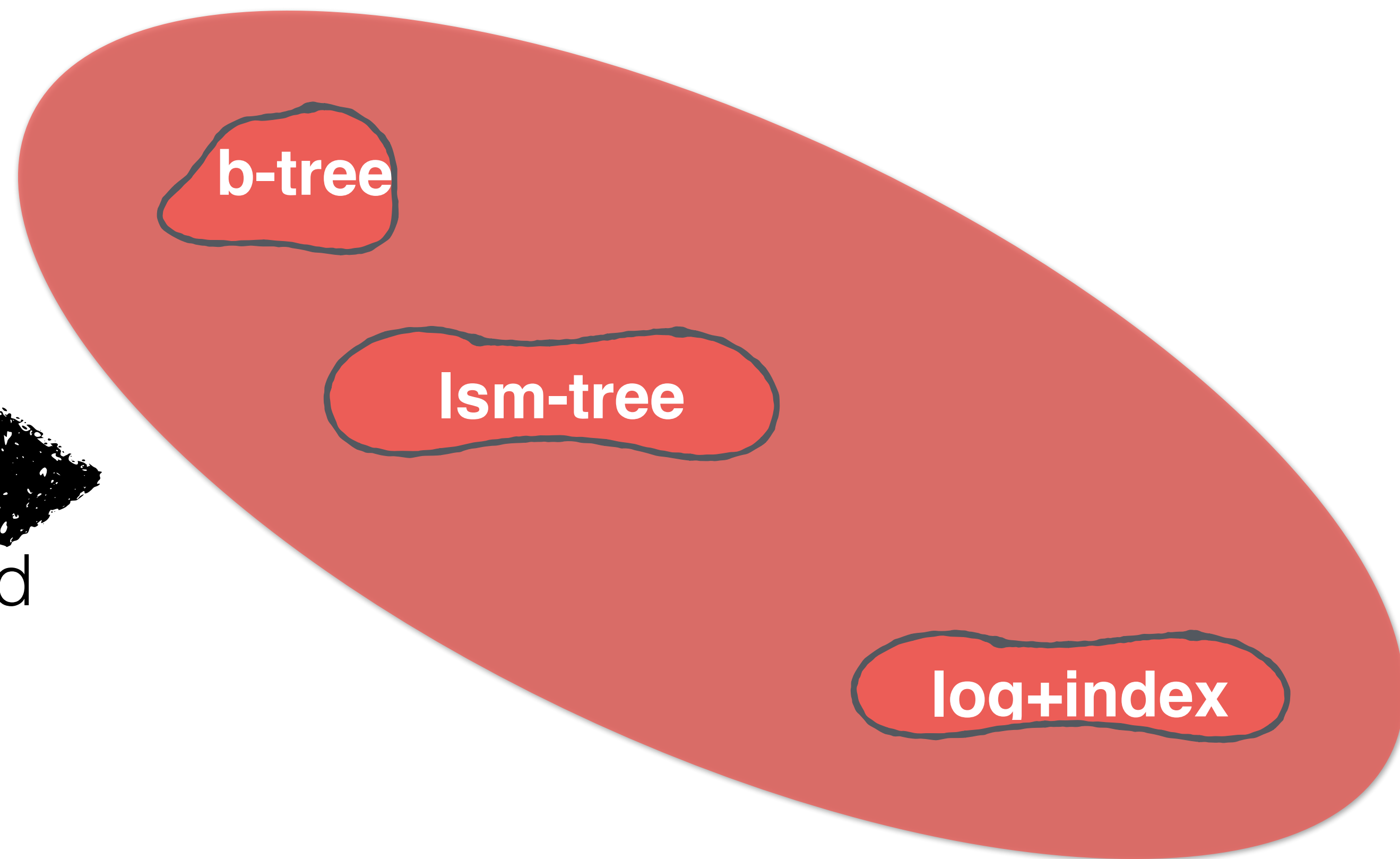
**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google

workload
budget
perf.



**analytical
I/O model**

**learned
CPU model**

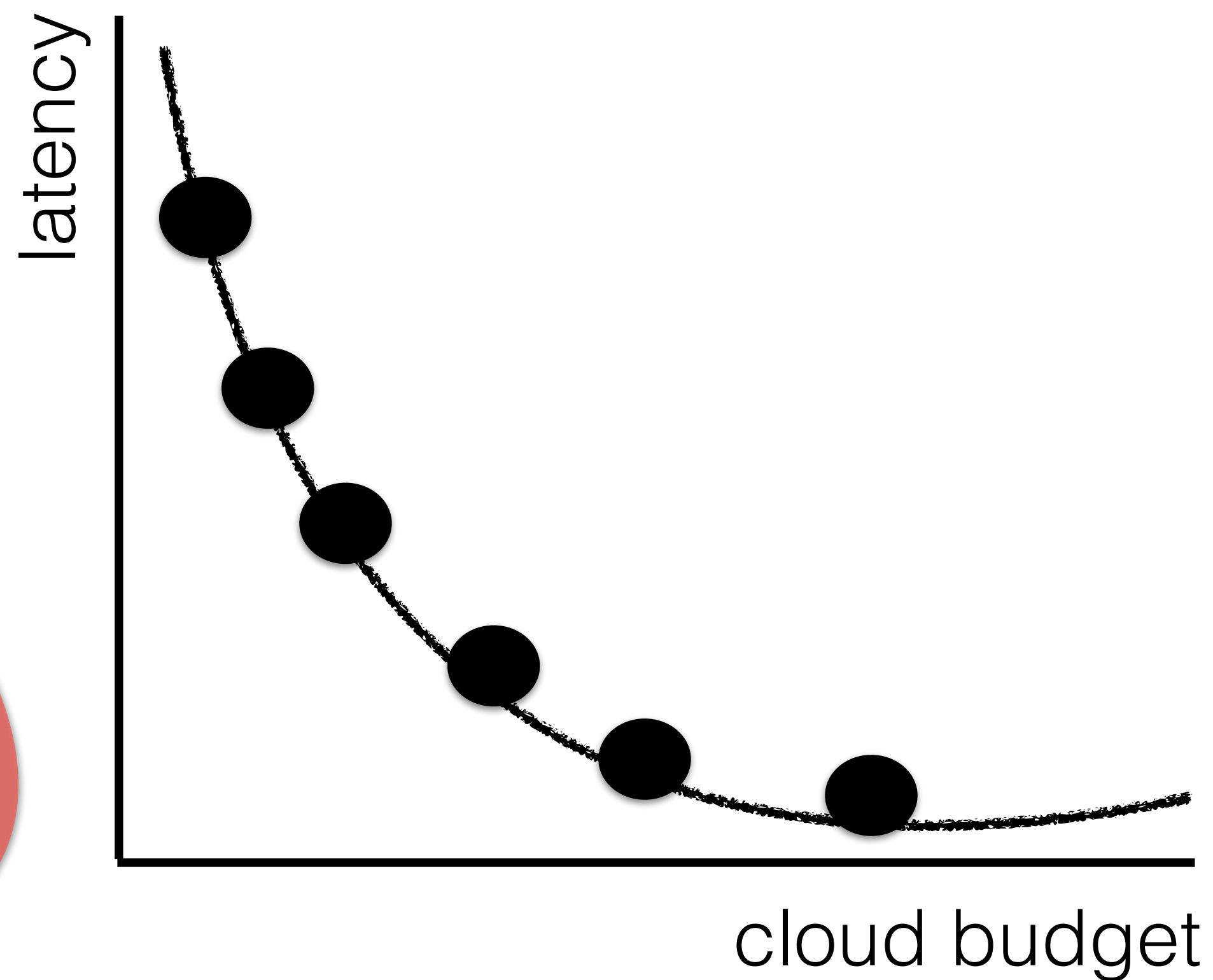
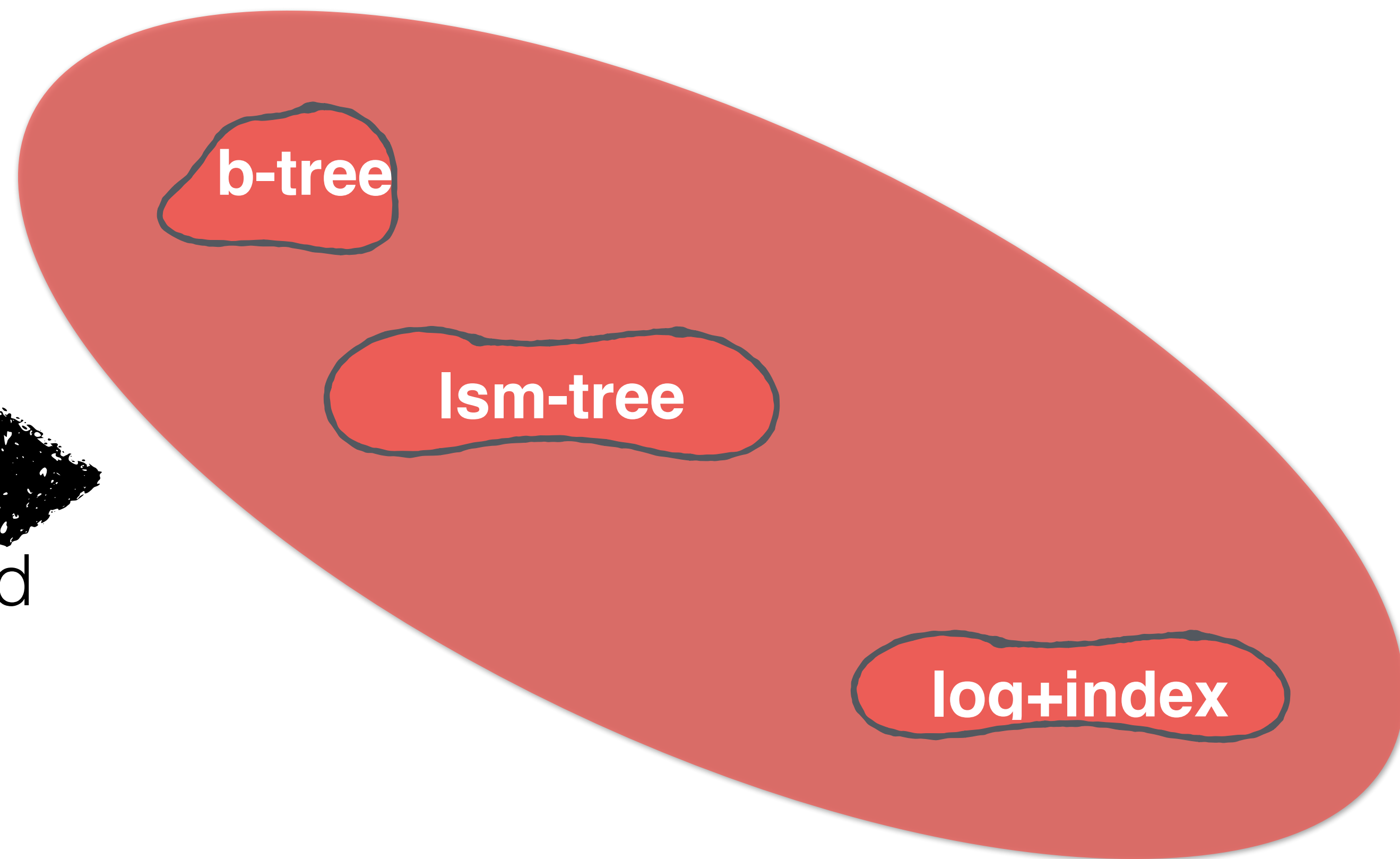
**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google

workload
budget
perf.



**analytical
I/O model**

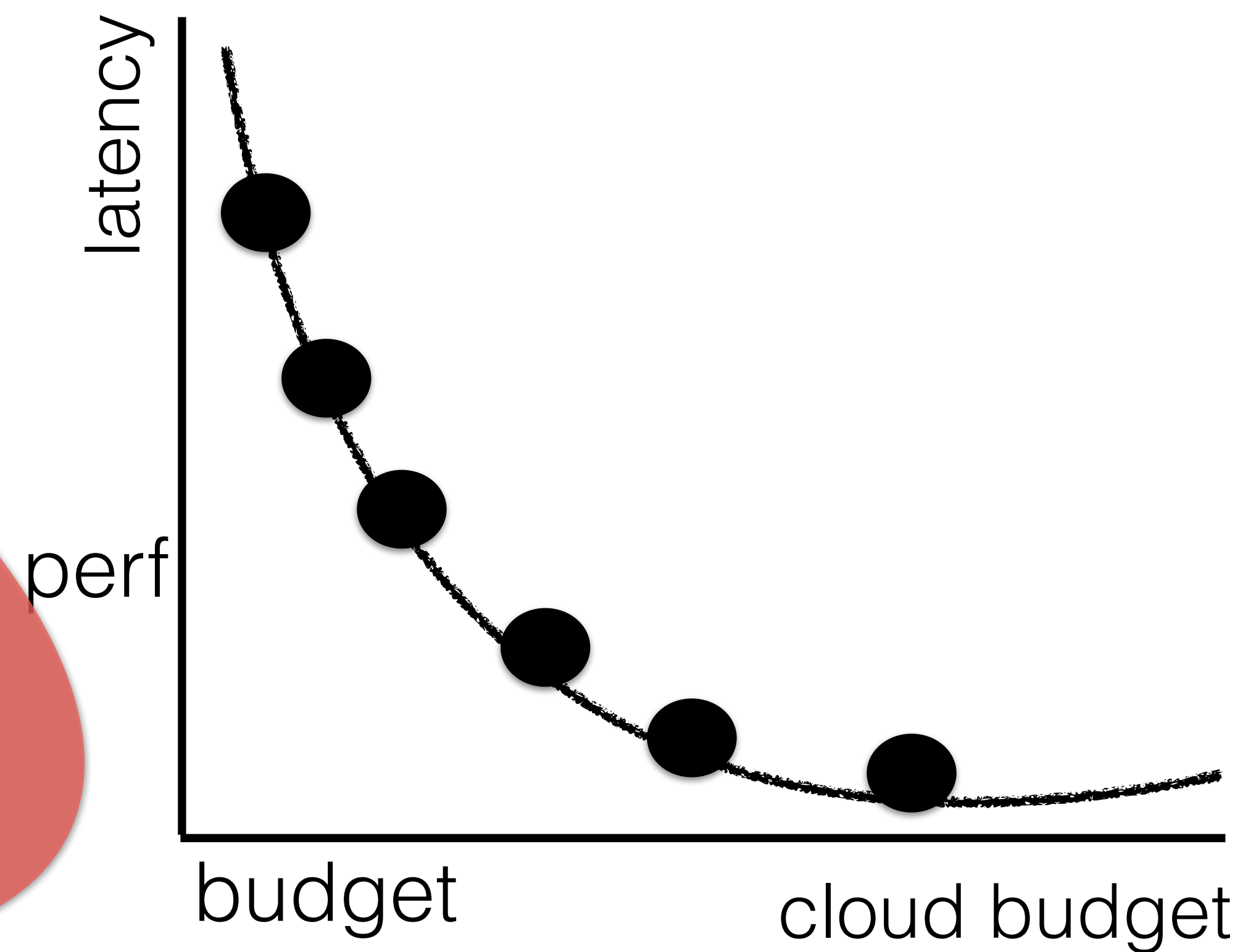
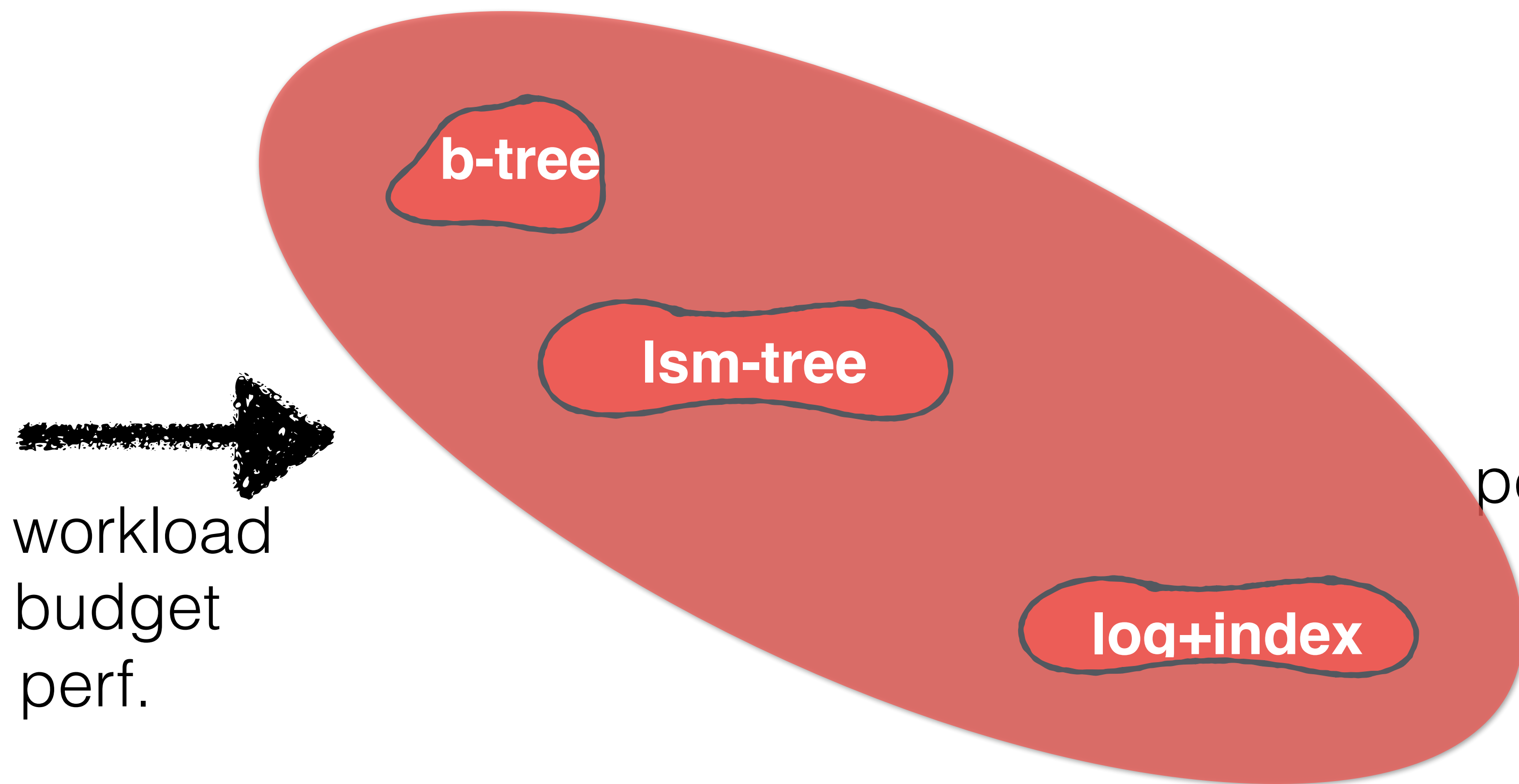
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



**analytical
I/O model**

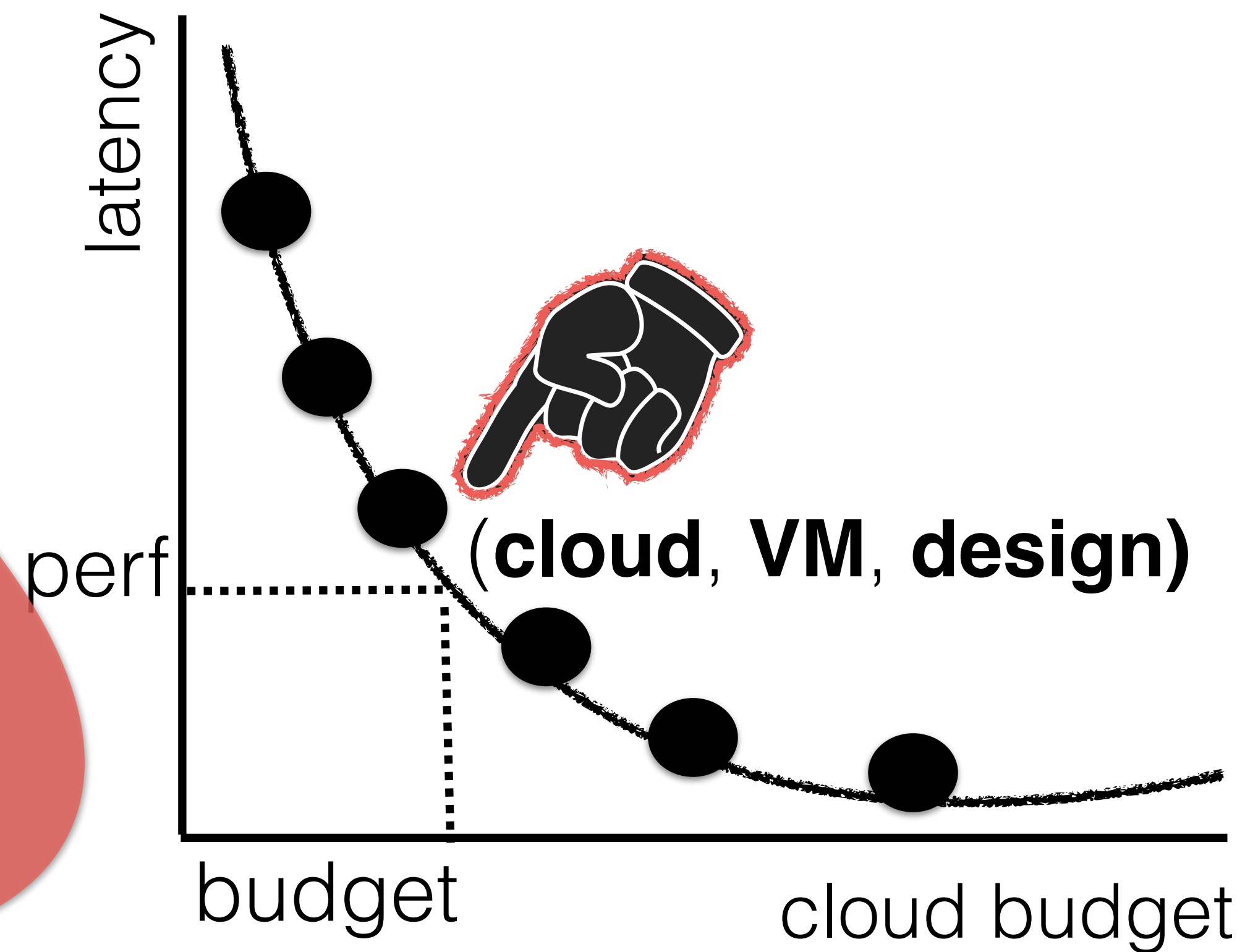
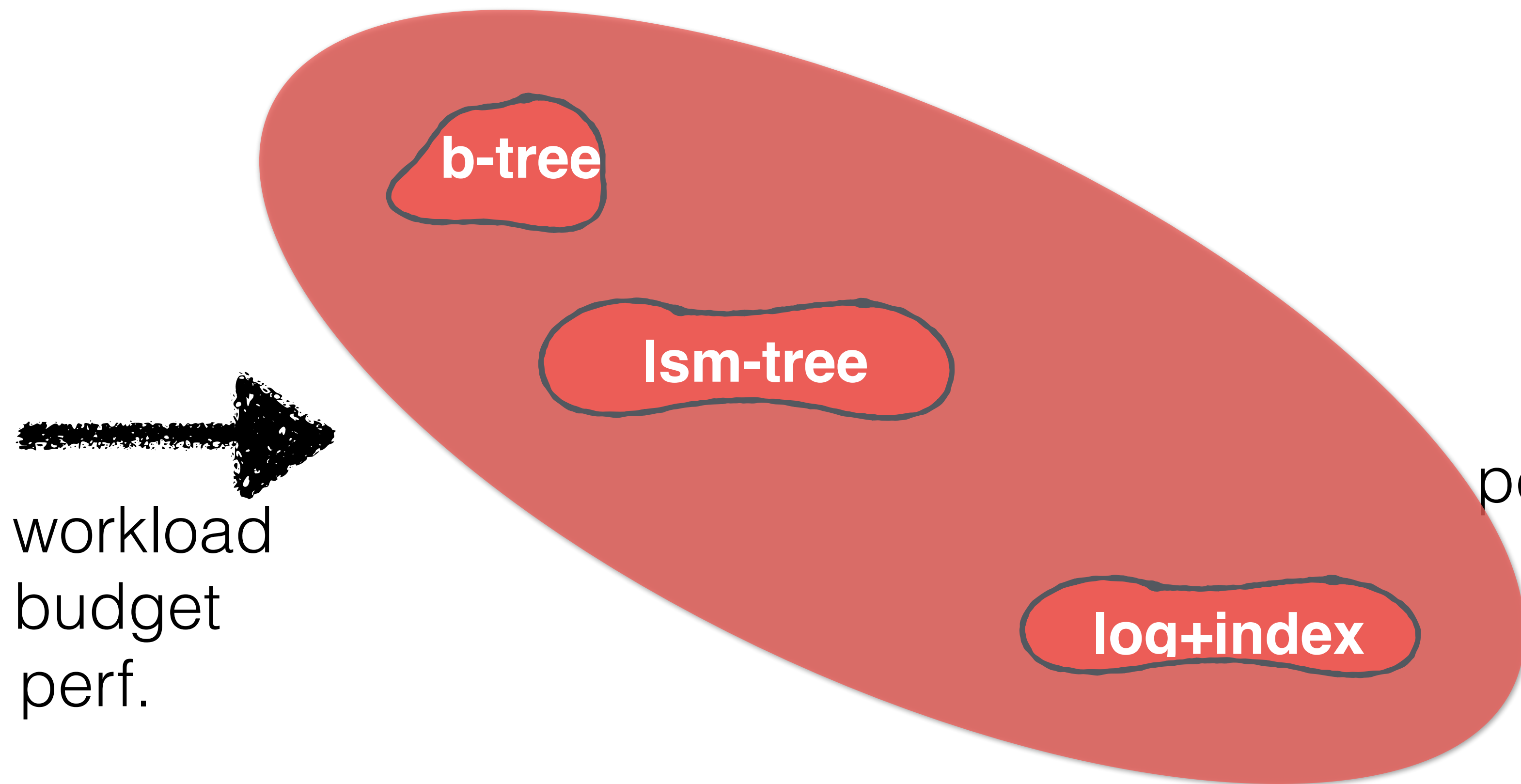
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



**analytical
I/O model**

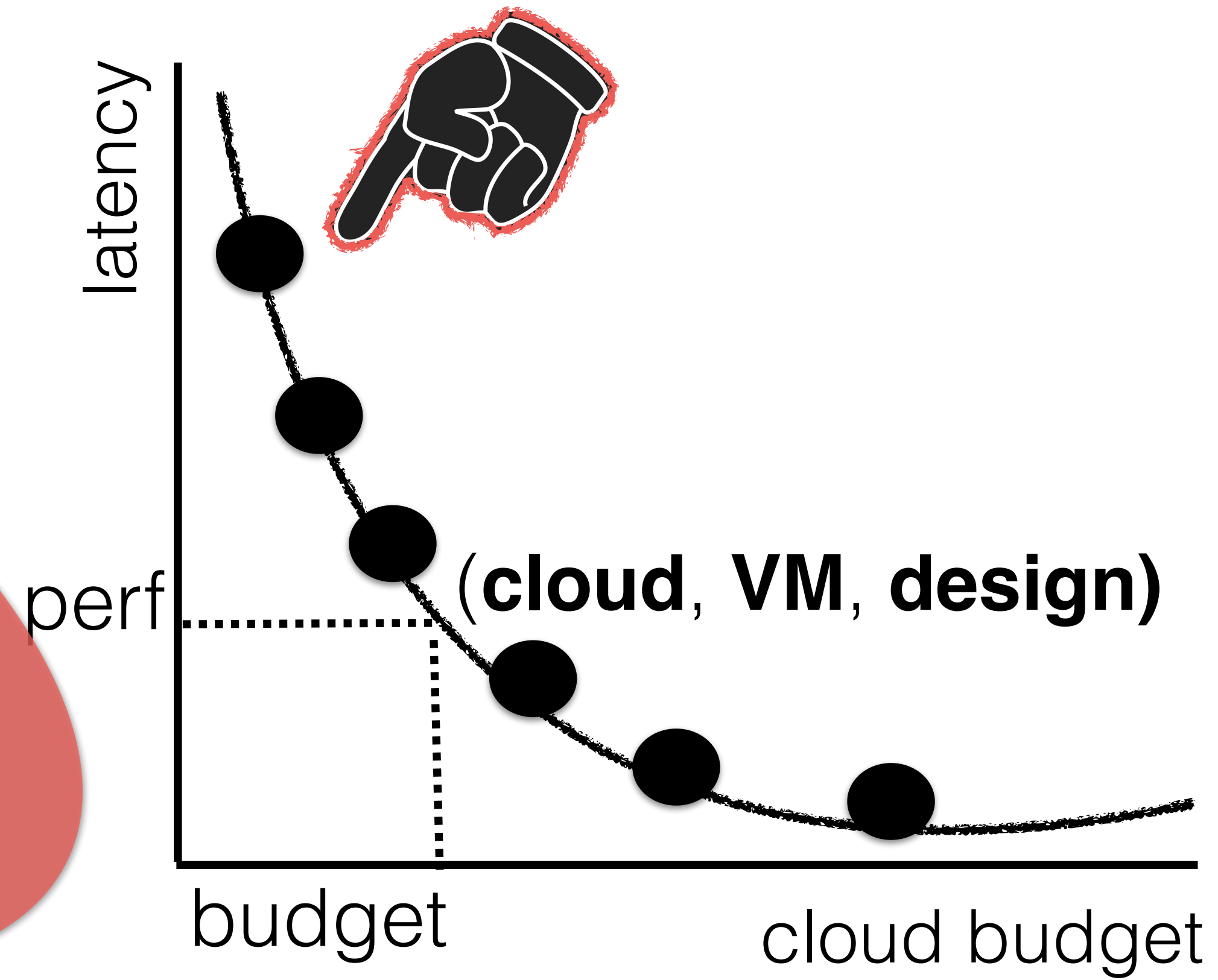
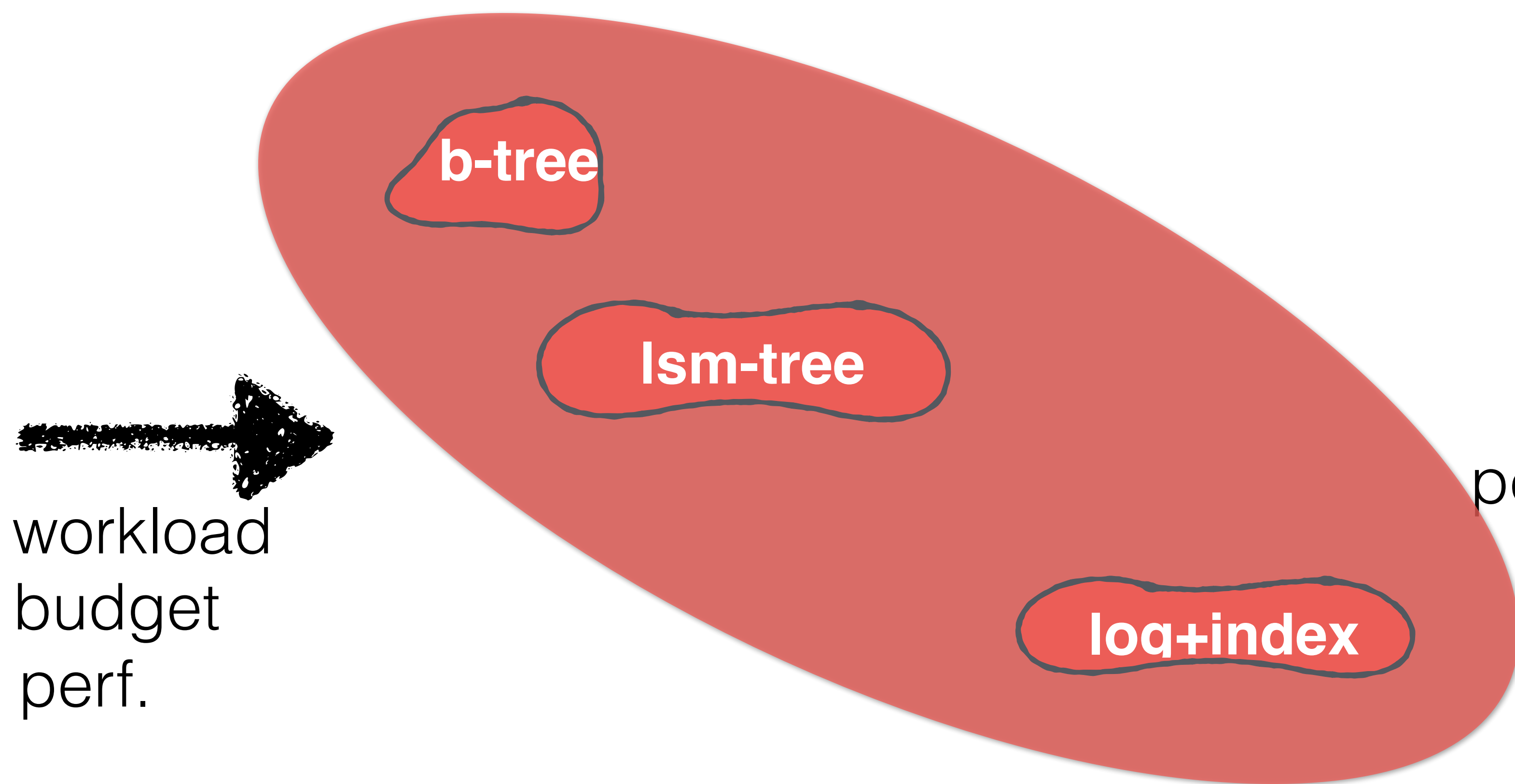
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



**analytical
I/O model**

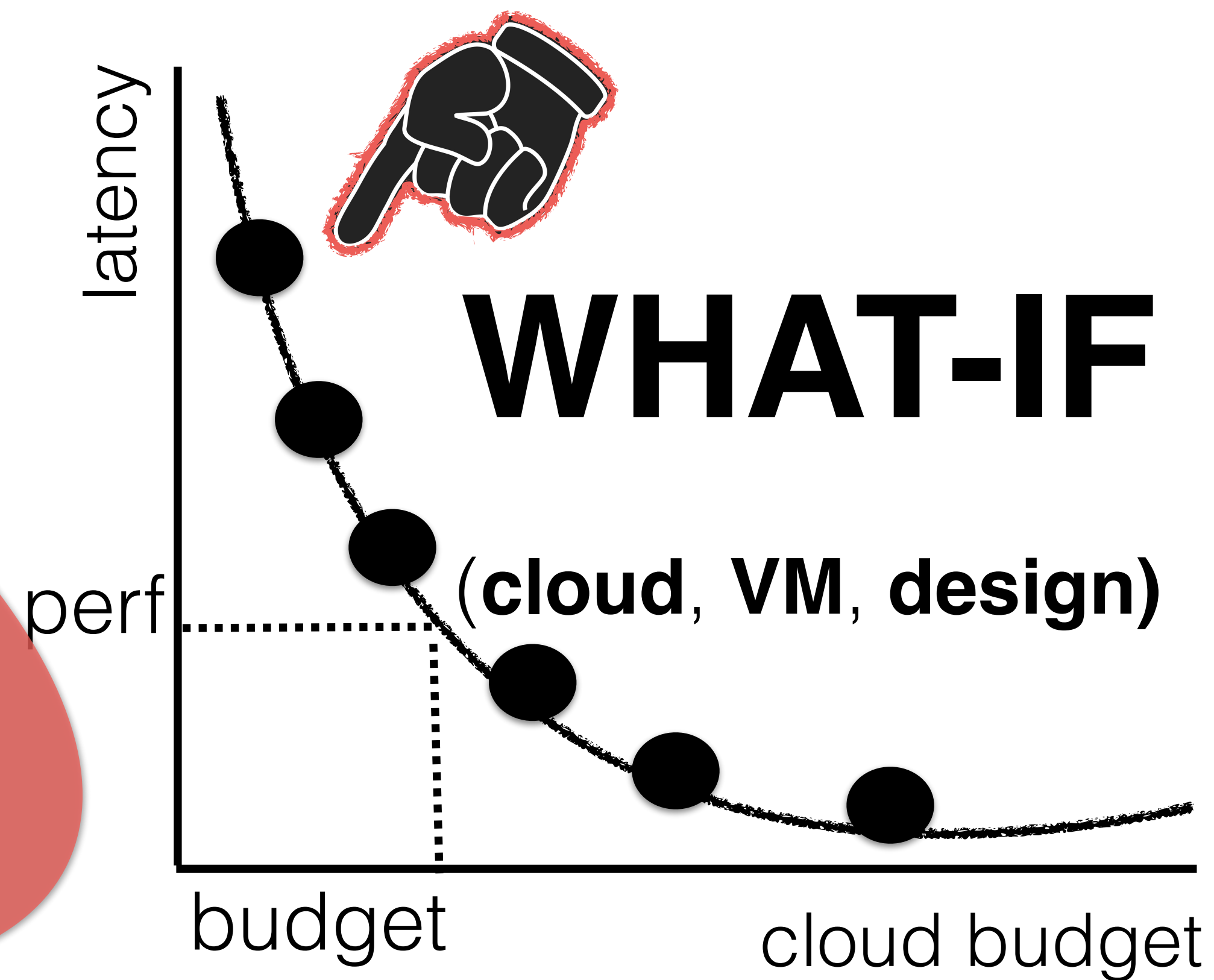
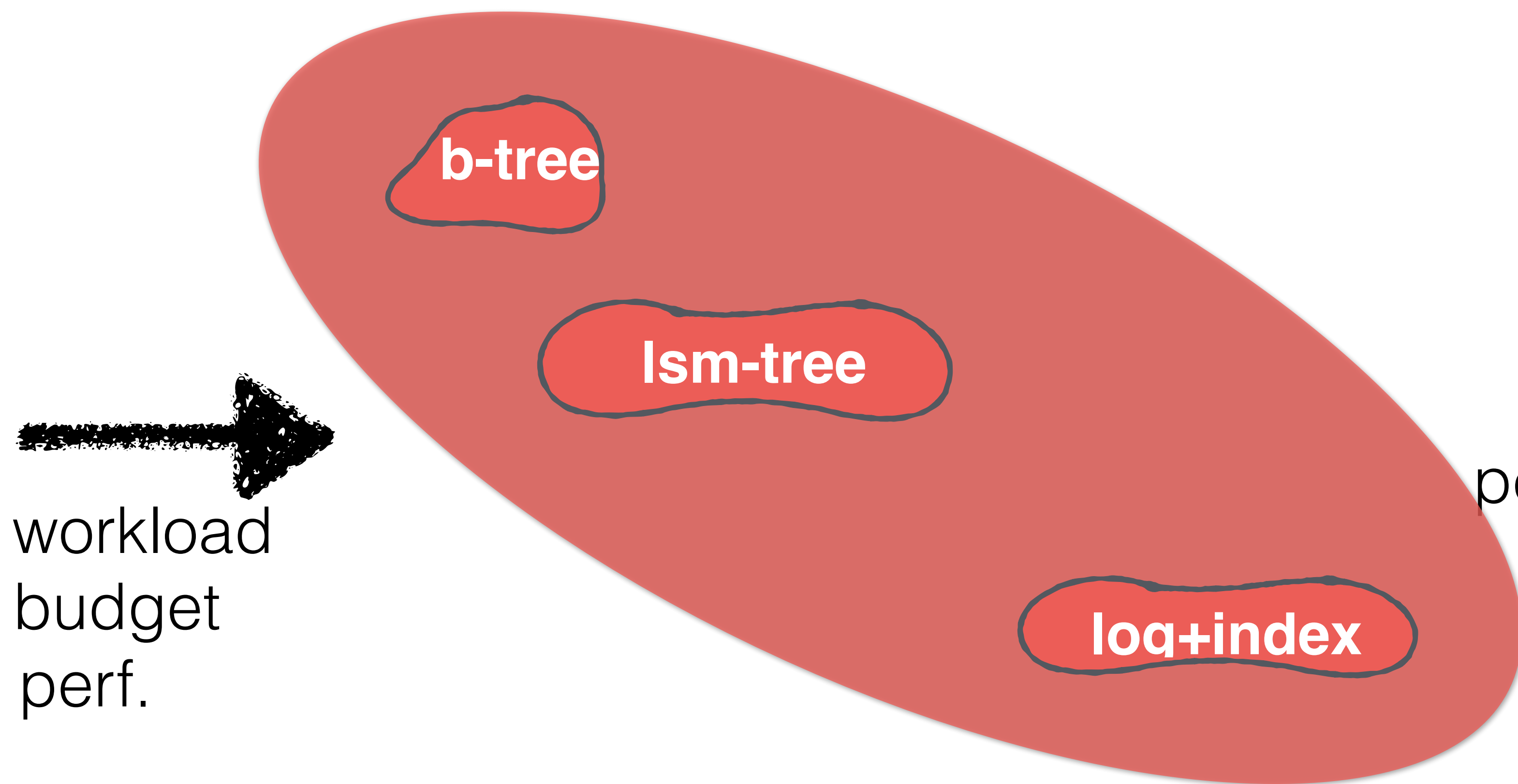
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



**analytical
I/O model**

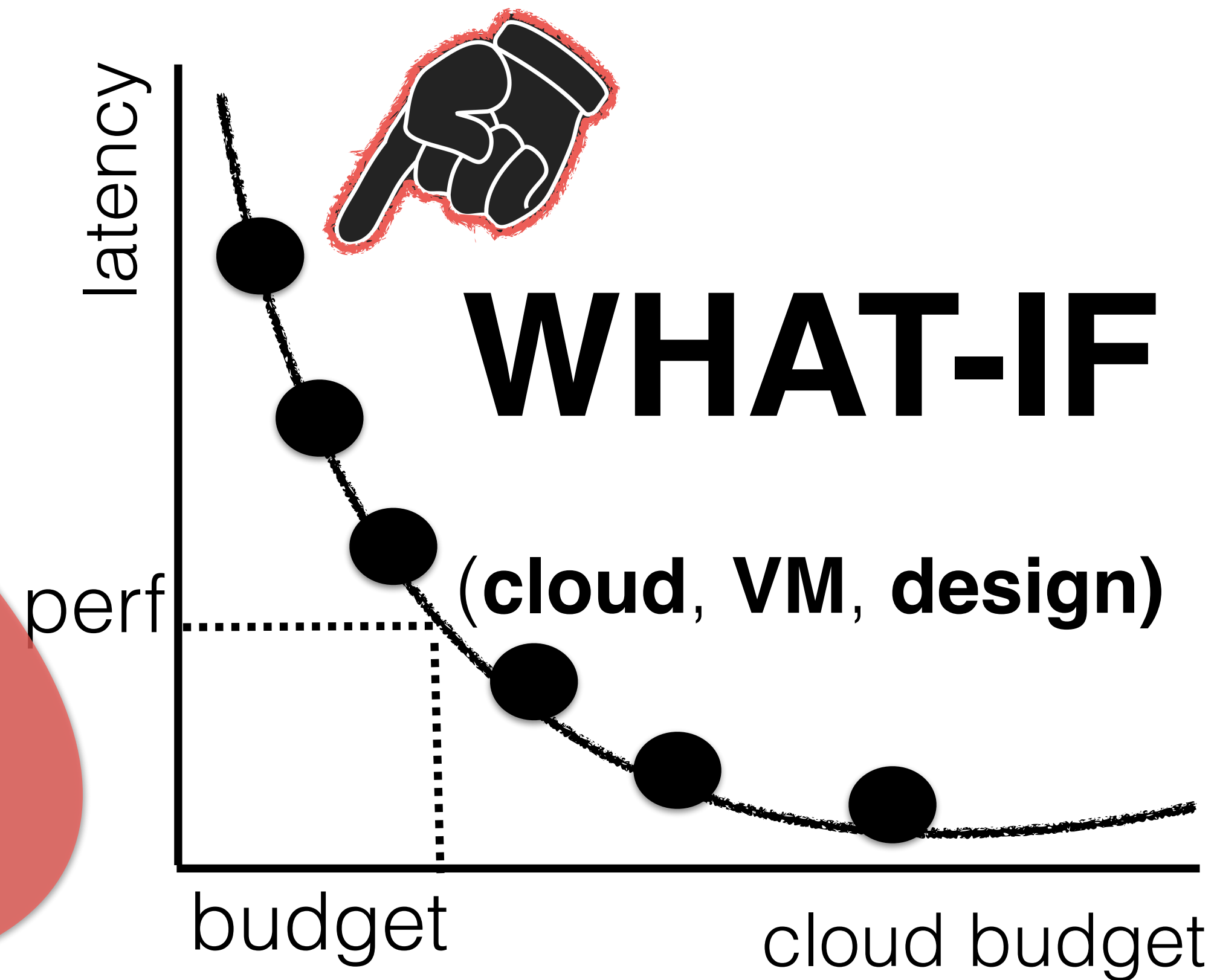
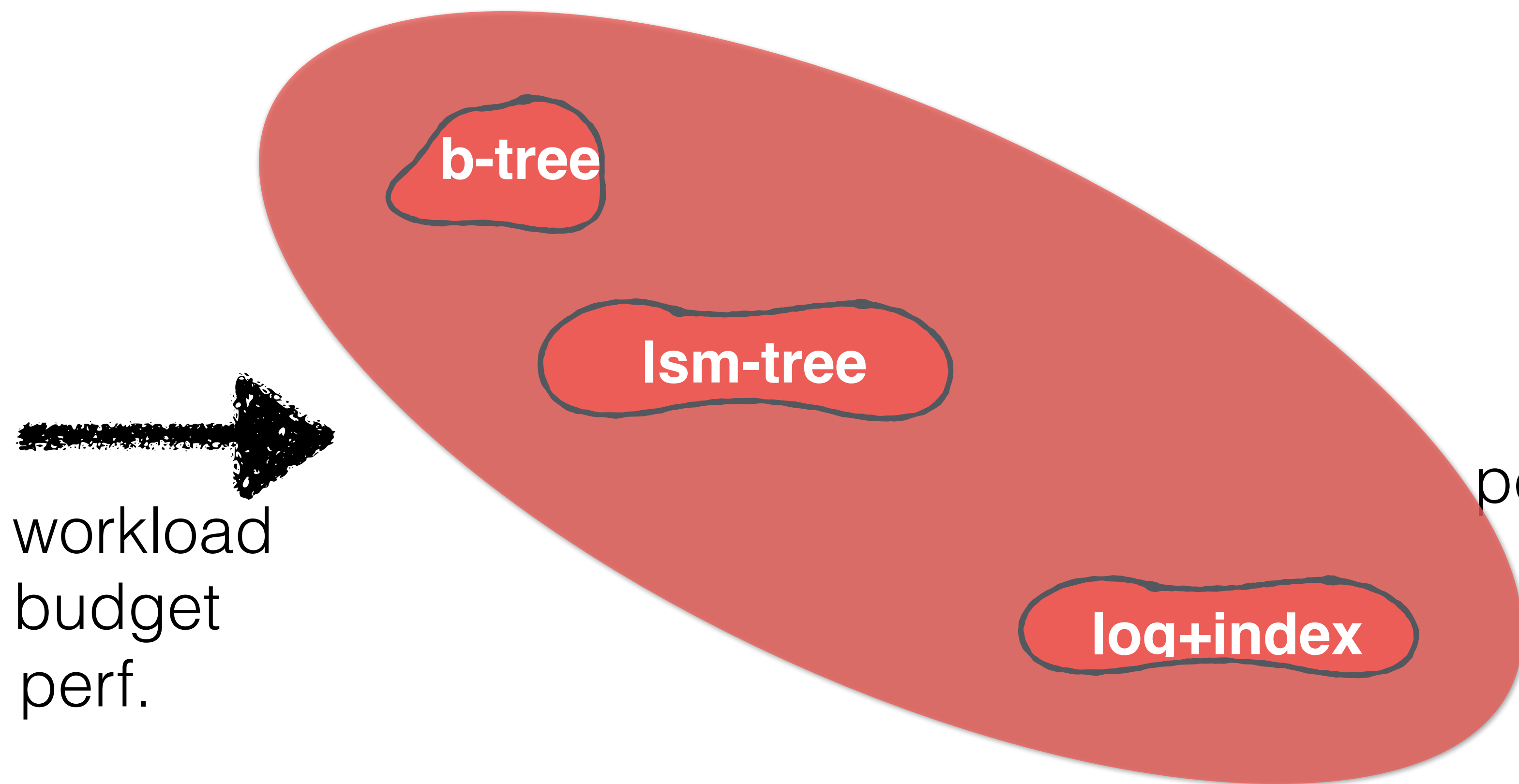
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



**analytical
I/O model**

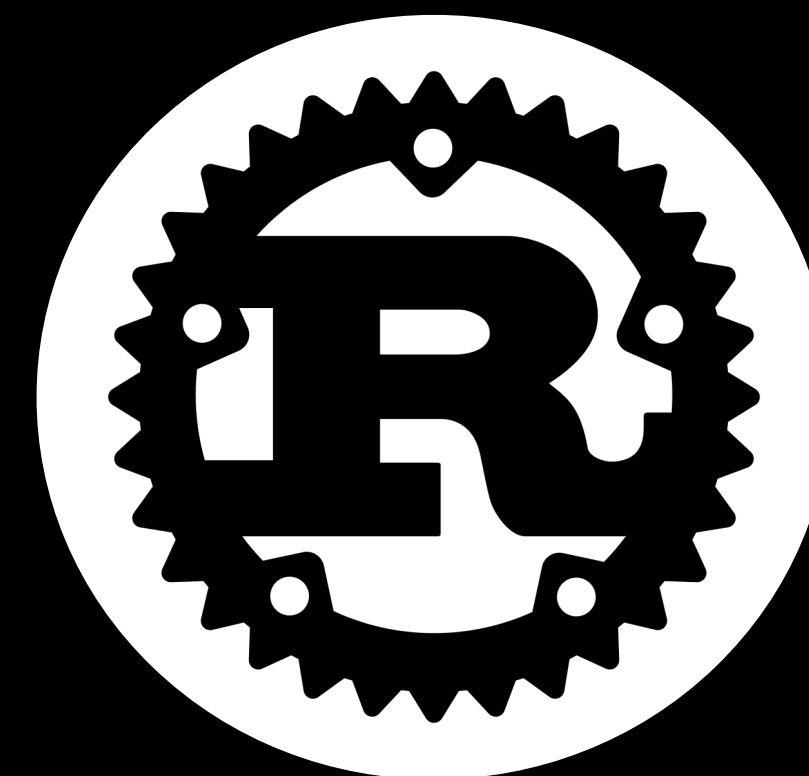
**learned
CPU model**

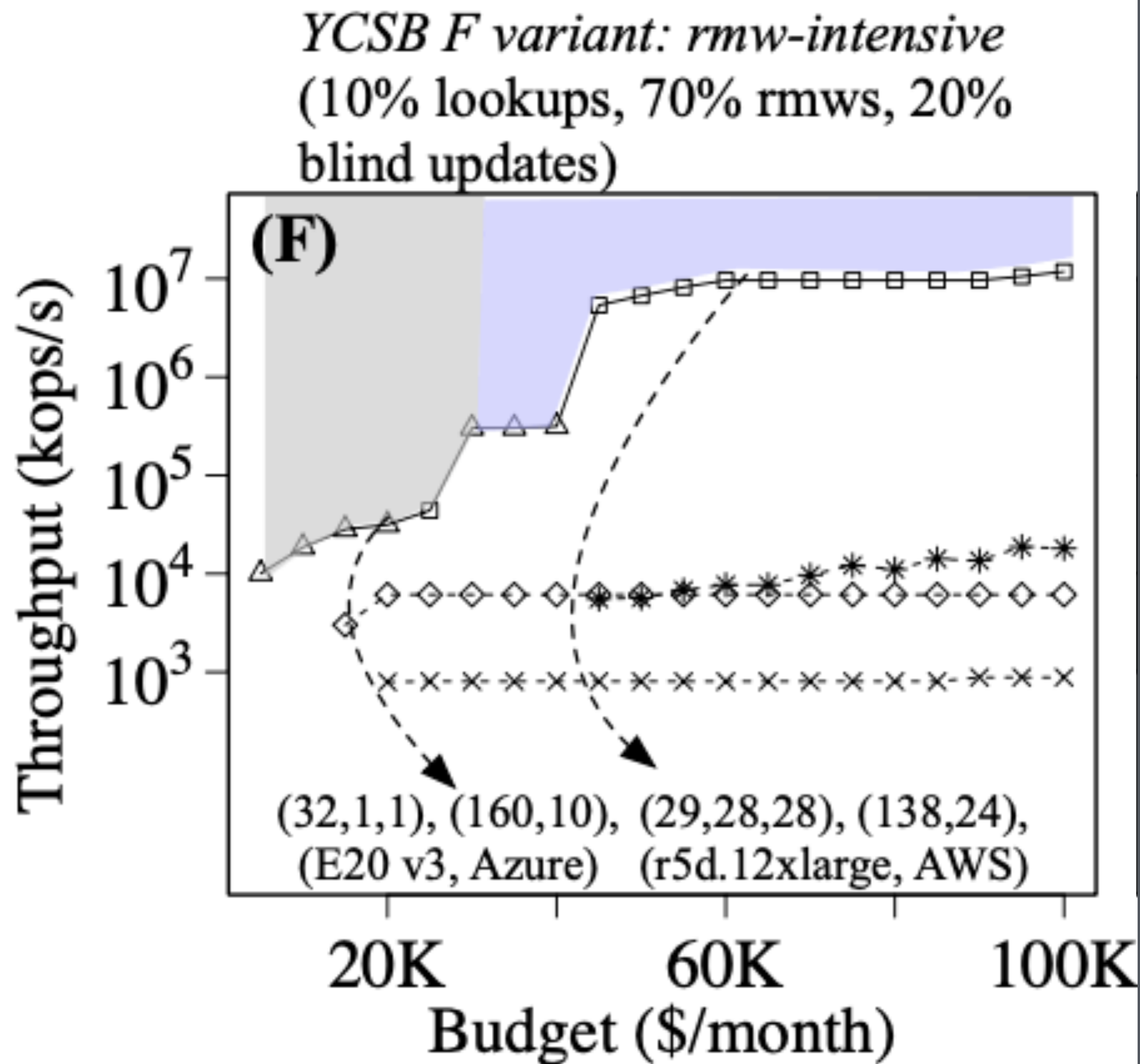
**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

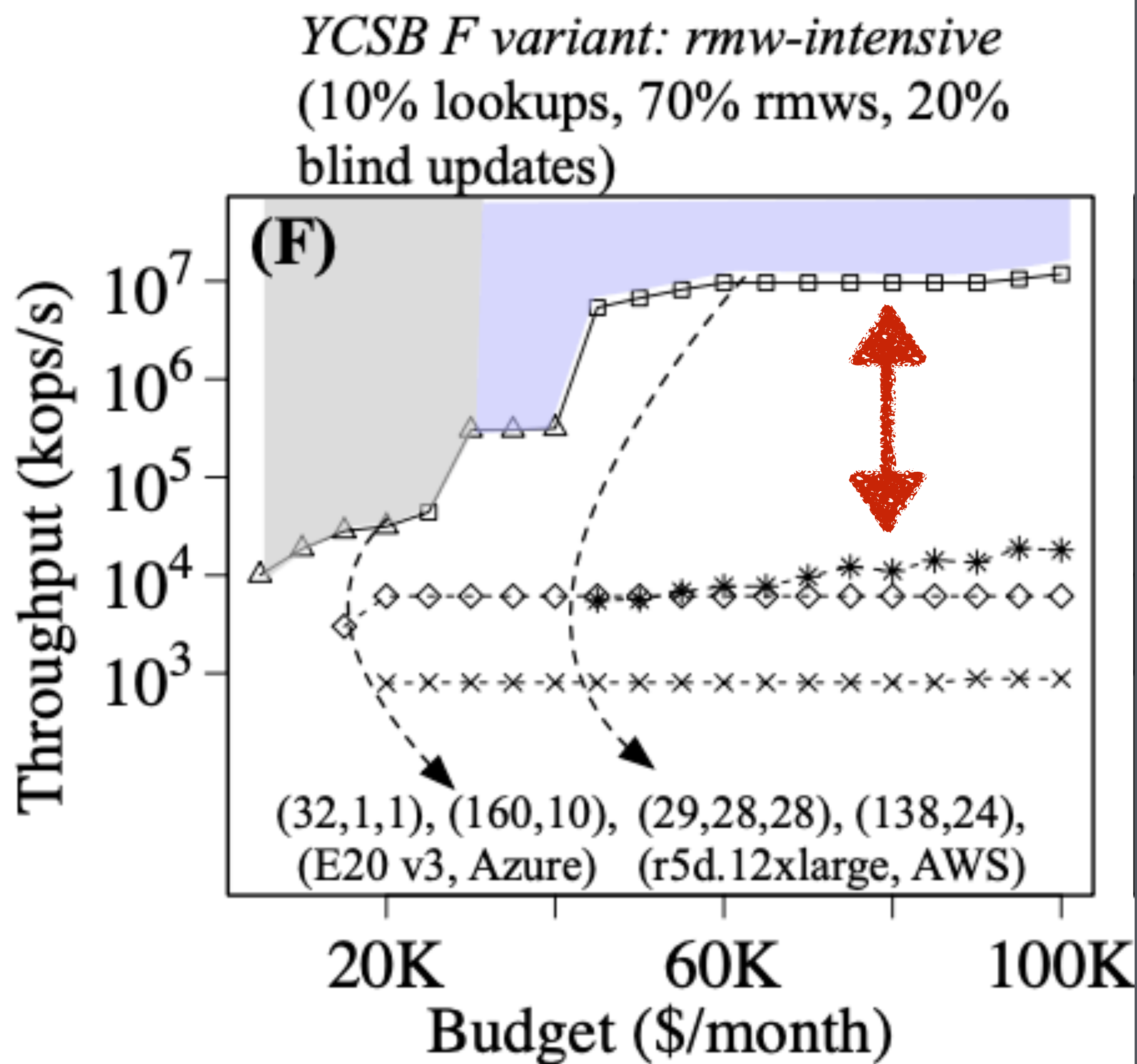
AWS, Azure, Google





✚ COSINE ✚

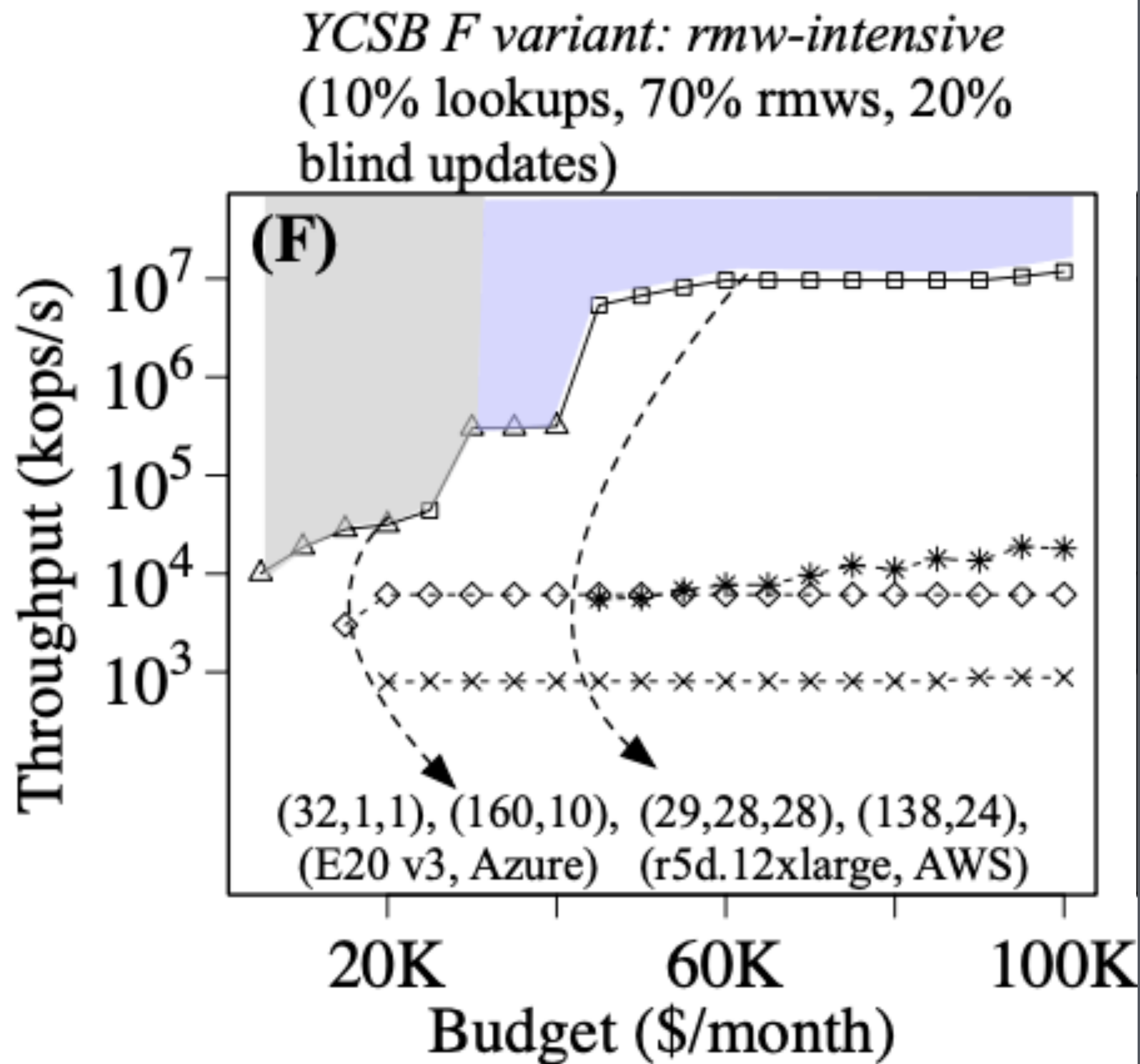
{ State of the Art
Meta, Microsoft, Mongo



✚ COSINE ✚

Better throughput/cost

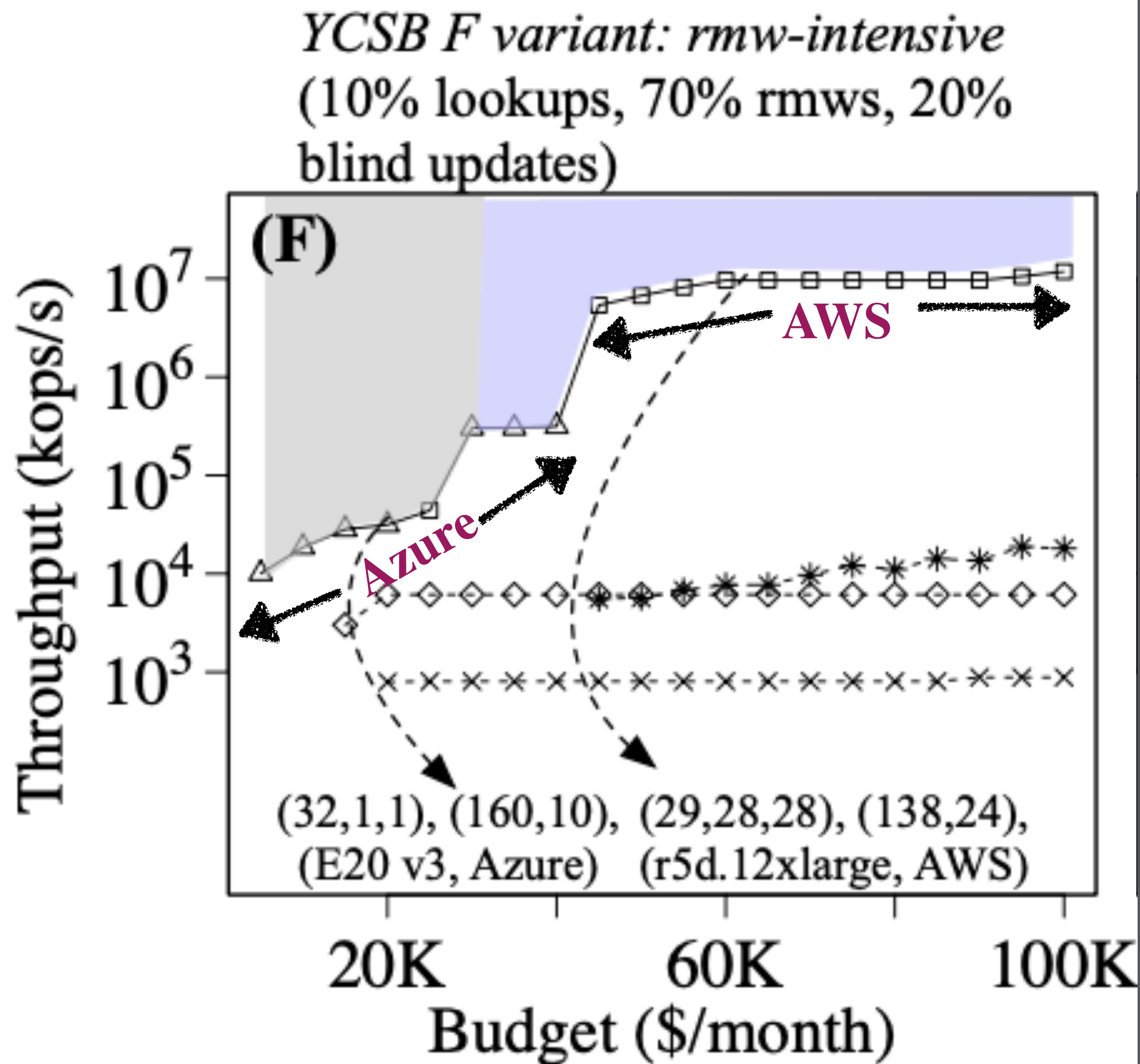
{ State of the Art
Meta, Microsoft, Mongo



✚ COSINE ✚

Better throughput/cost
Self-designs

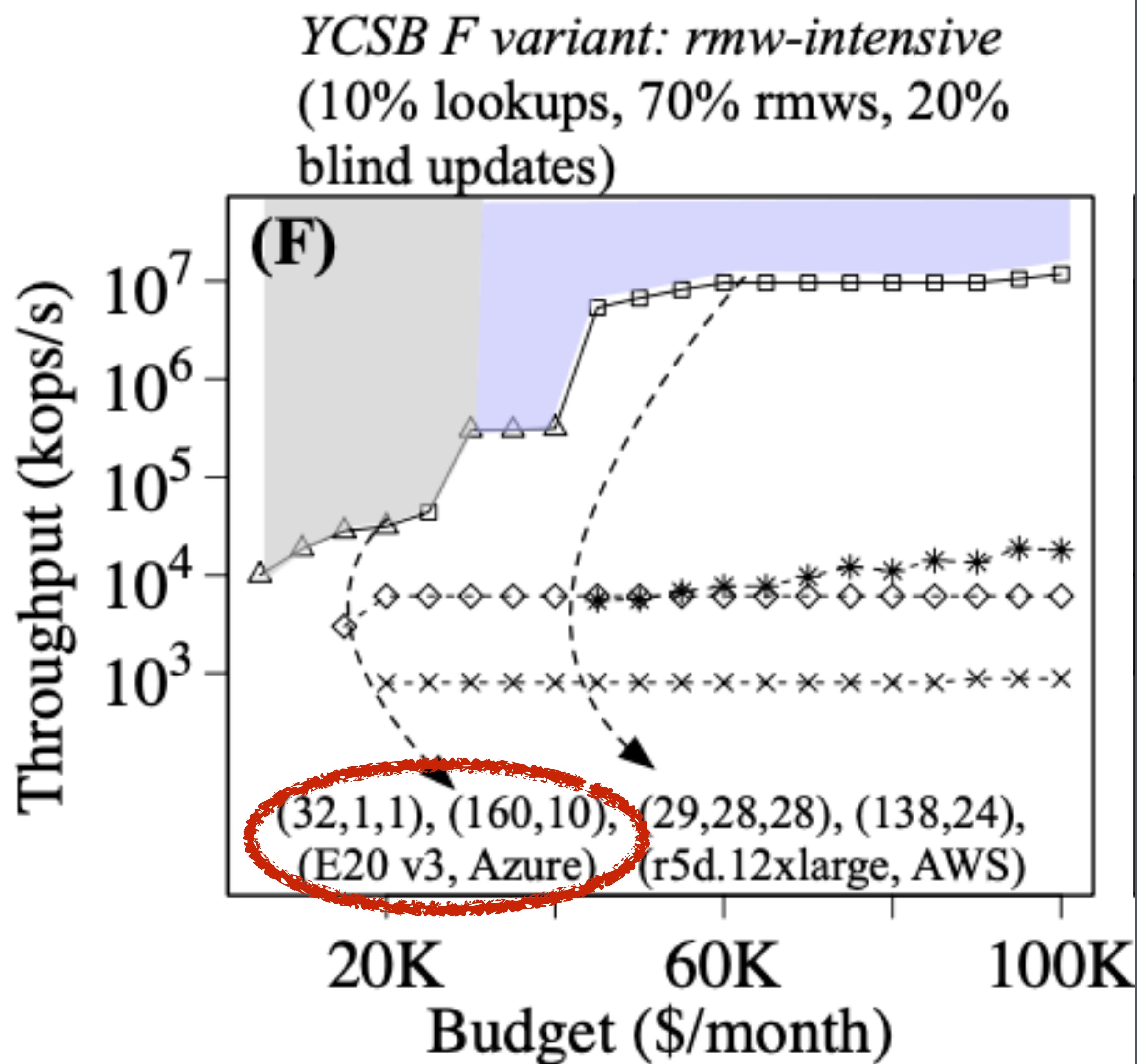
{ State of the Art
Meta, Microsoft, Mongo



⚙️ COSINE ⚙️

Better throughput/cost
Self-designs

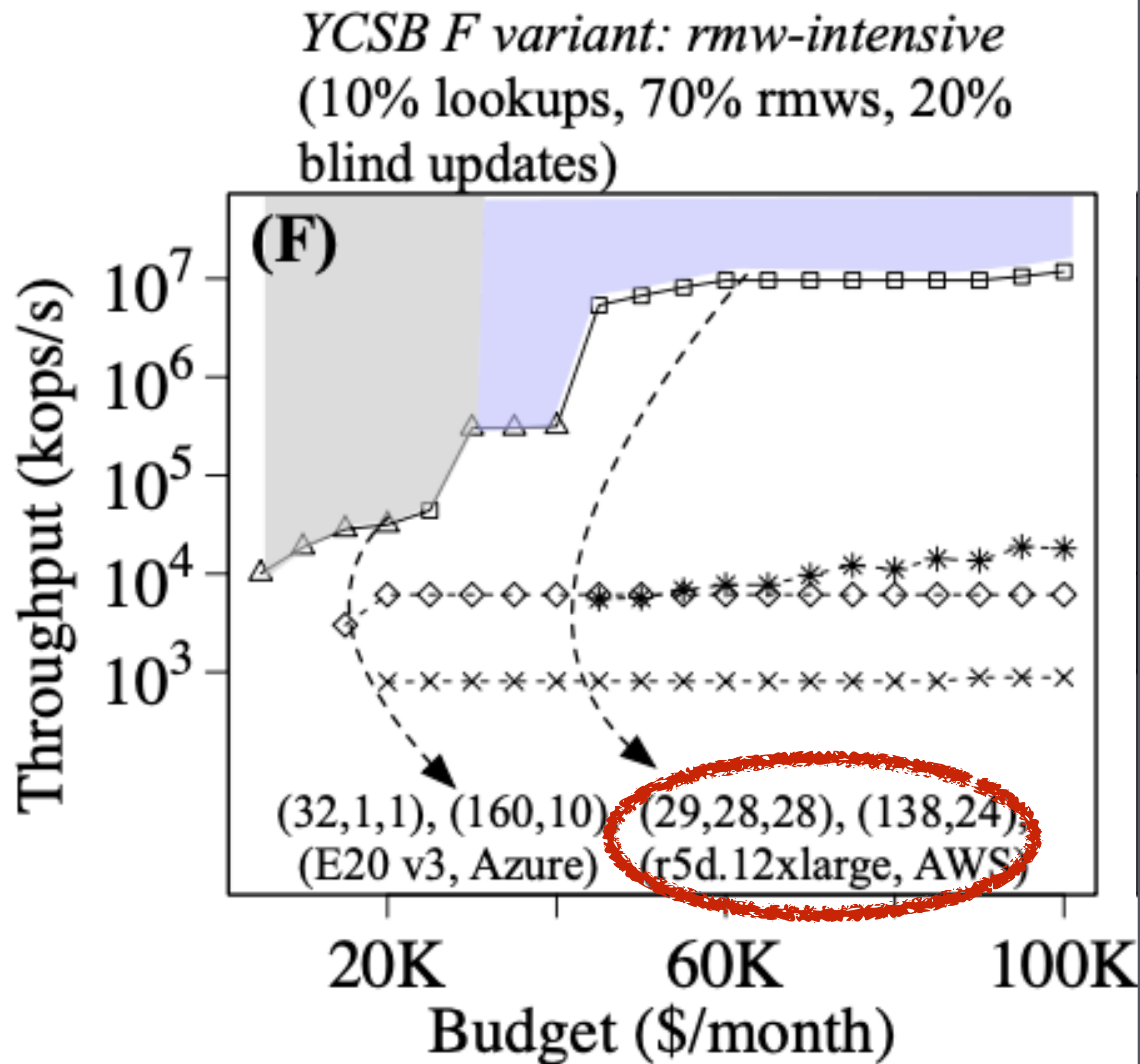
{ State of the Art
Meta, Microsoft, Mongo



✚ COSINE ✚

Better throughput/cost
Self-designs

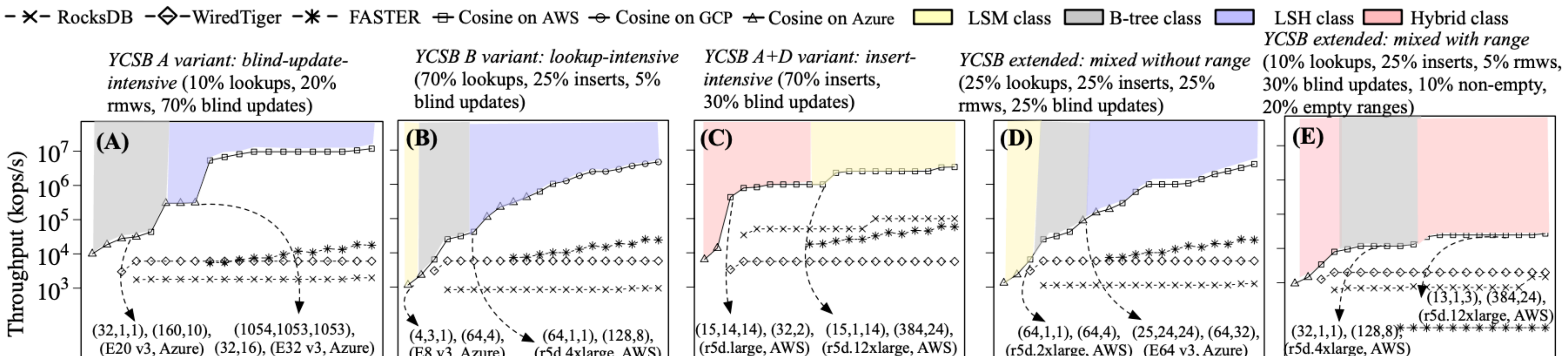
{ State of the Art
Meta, Microsoft, Mongo



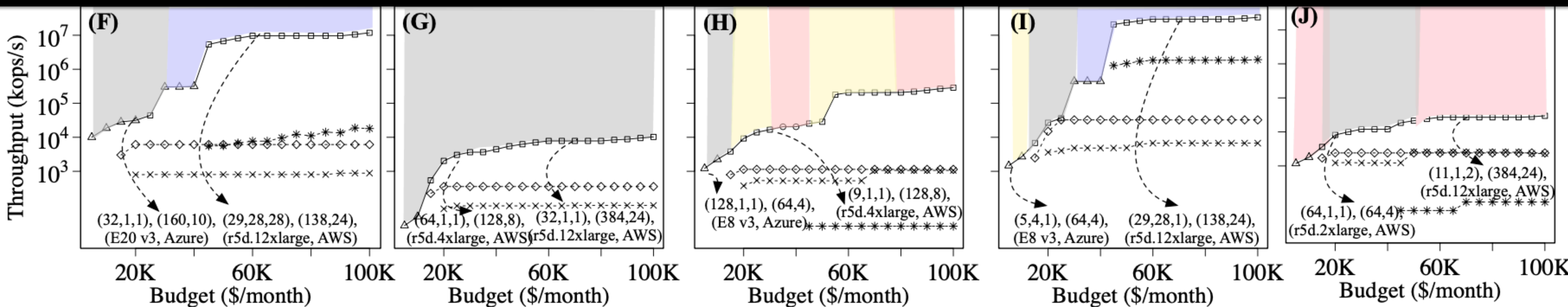
✚ COSINE ✚

Better throughput/cost
Self-designs

{ State of the Art
Meta, Microsoft, Mongo



beats top systems across all workloads by self-designing new systems

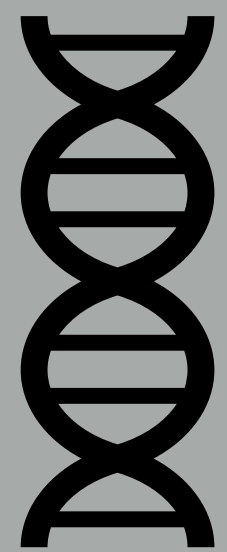


workload/budget diversity

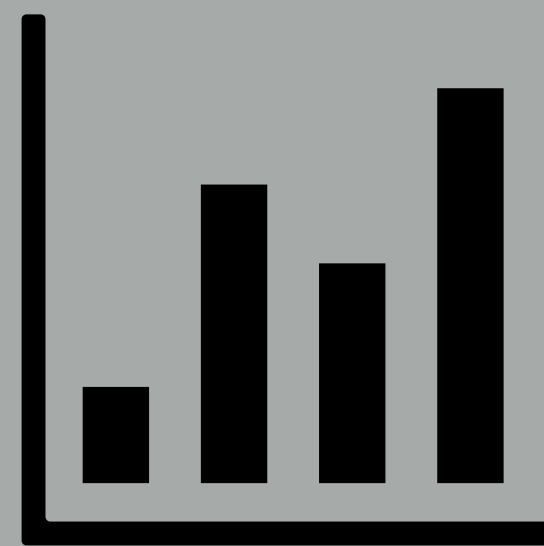
IMAGE AI STORAGE

new primitives

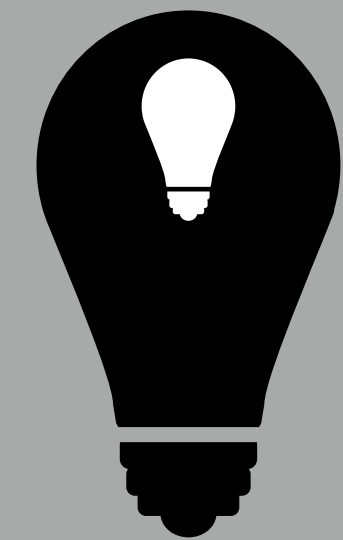
accuracy



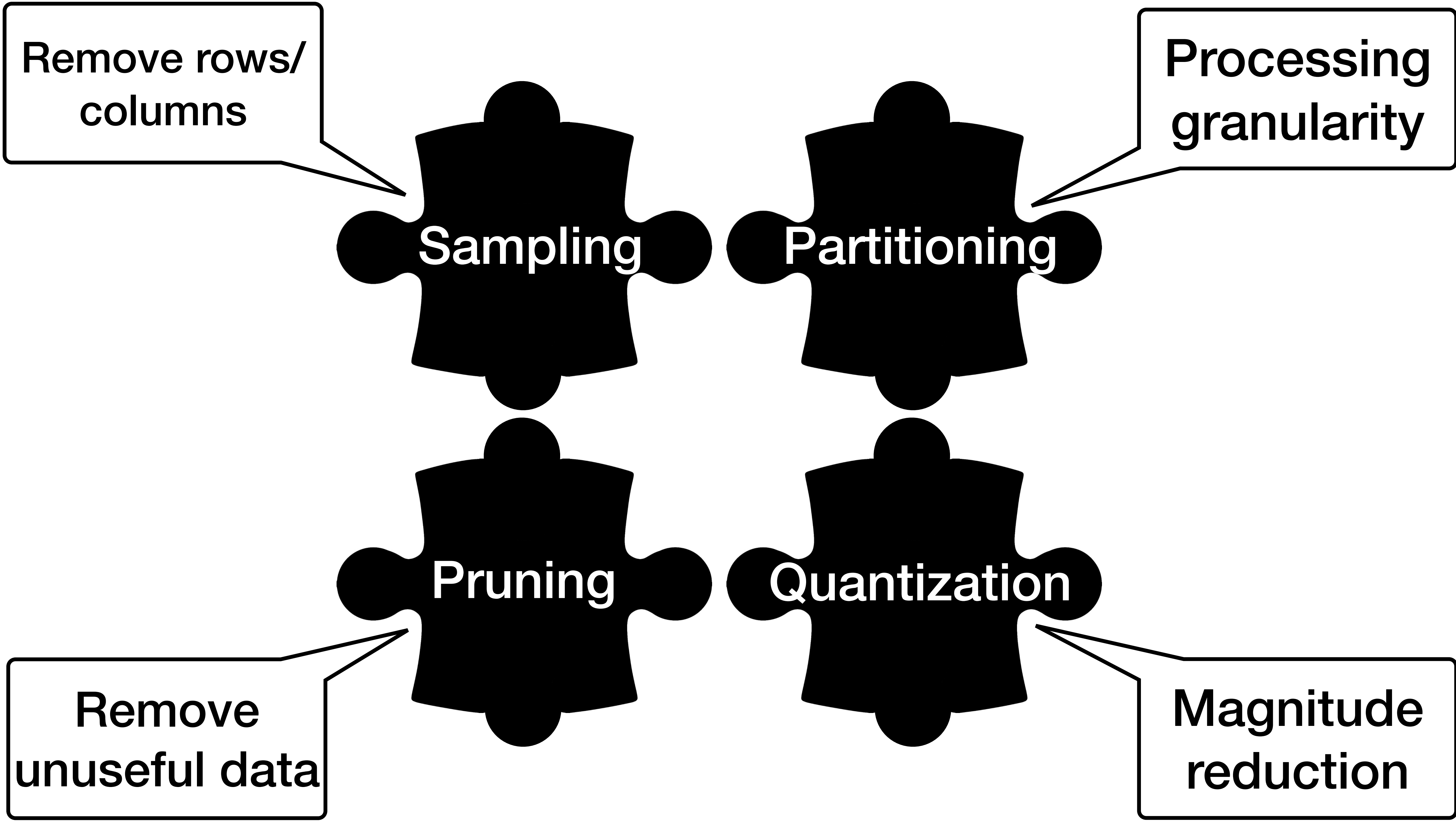
DESIGN SPACE



COST ESTIMATION



SEARCH



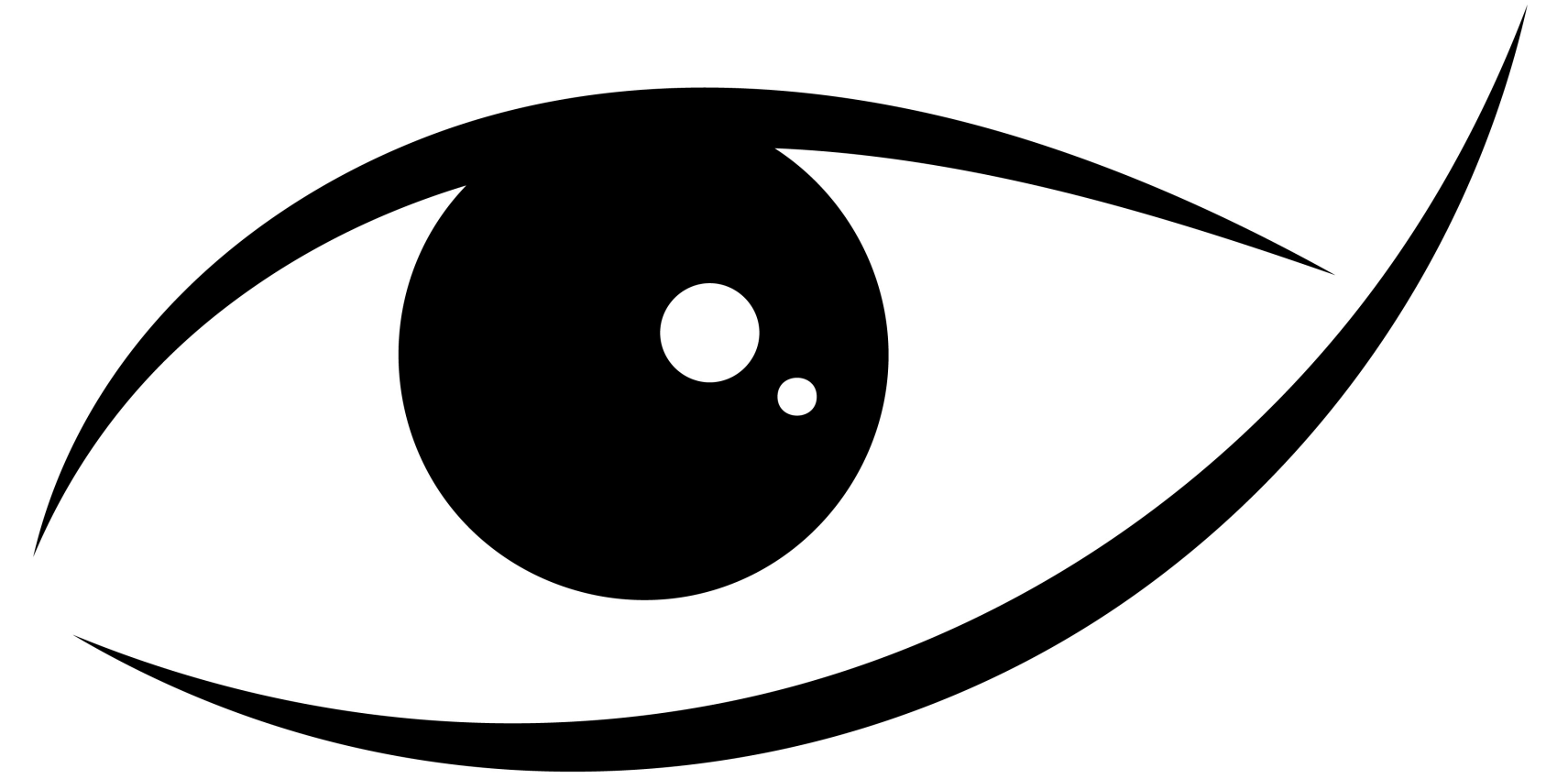
**How do machines
store images today?**

**How do machines
store images today?**

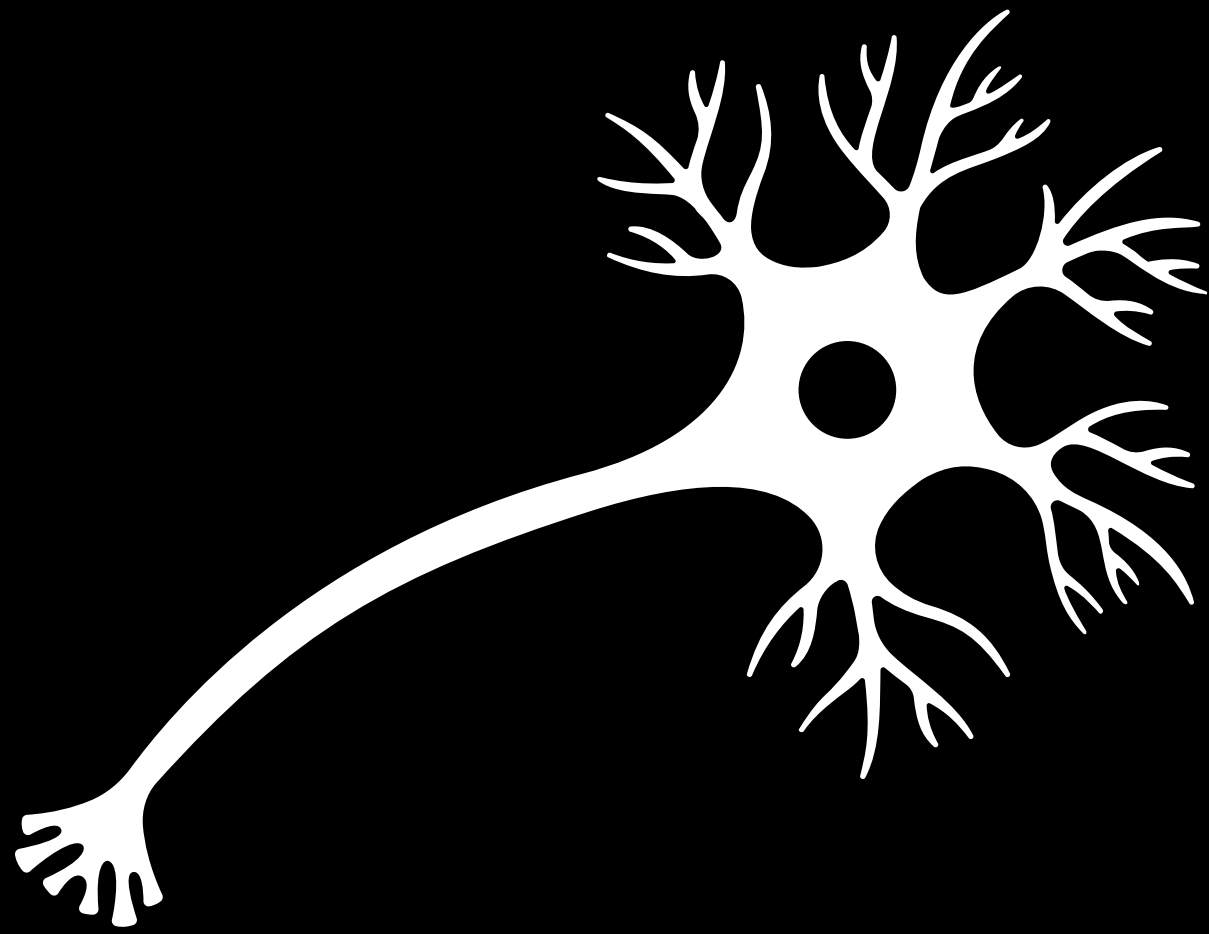
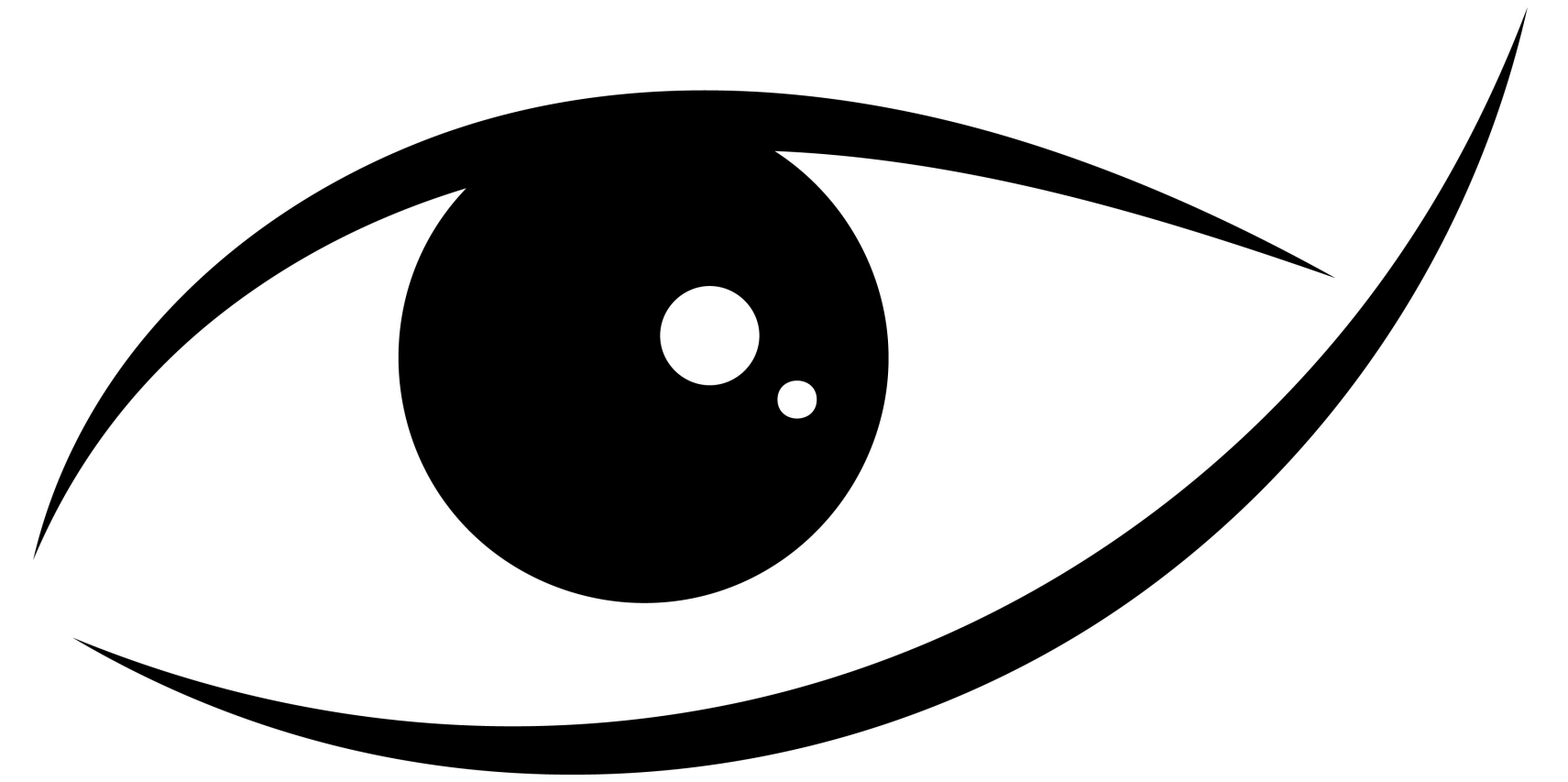
JPEG

Joint Photographic Experts Group

JPEG is tailored for the
properties of the human eye



JPEG is tailored for the
properties of the human eye



images for AI are seen by
algorithms, not humans

10^{150} designs

massive design space of possible image storage schemes

10^{150} designs

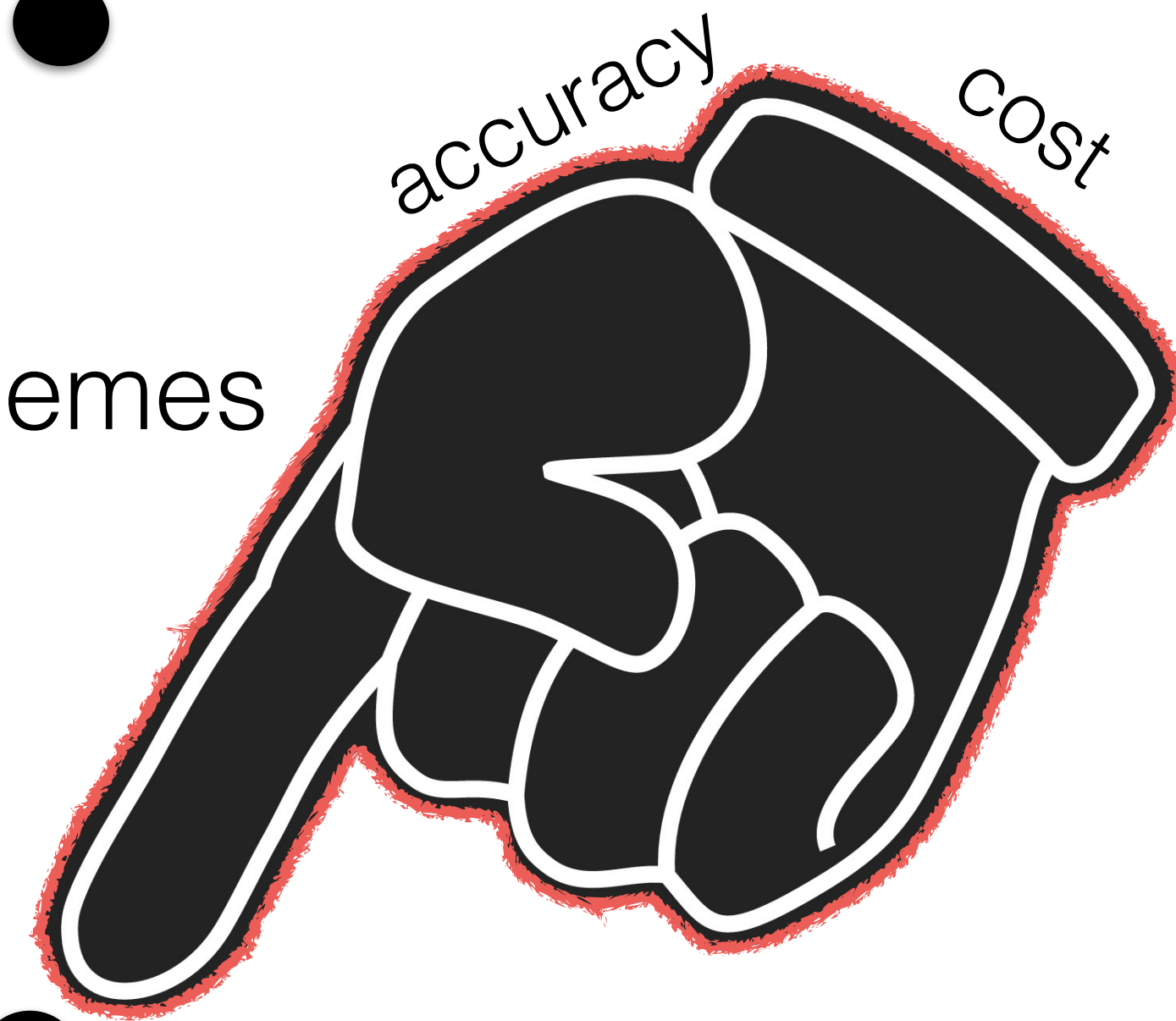
JPEG

massive design space of possible image storage schemes

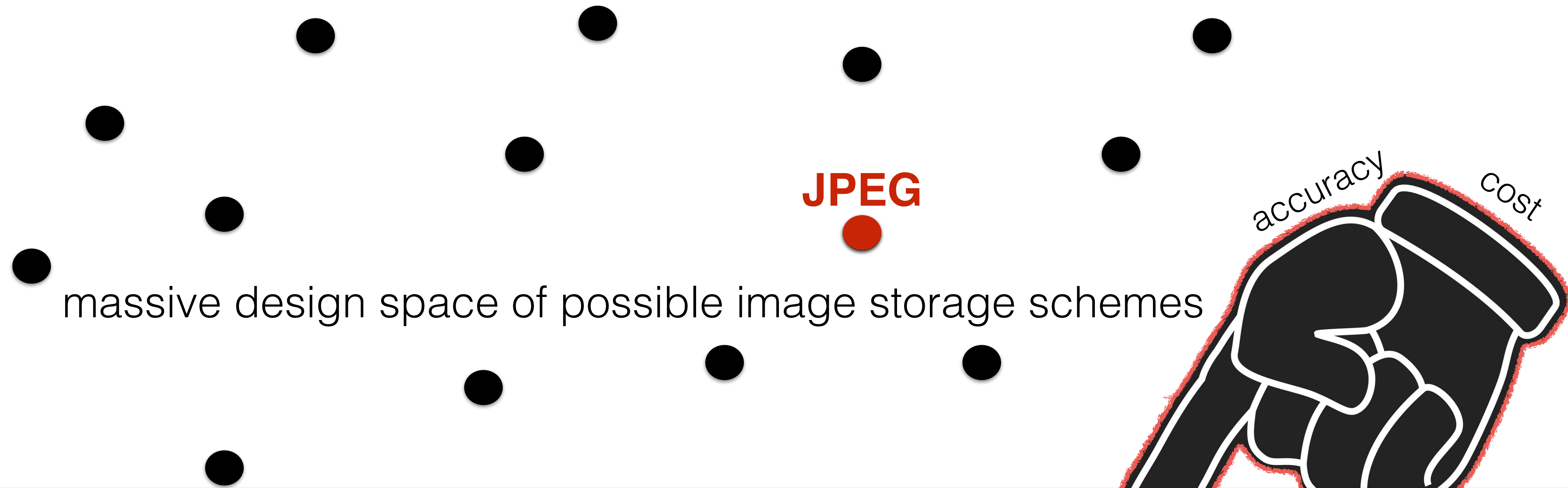
10^{150} designs

JPEG

massive design space of possible image storage schemes

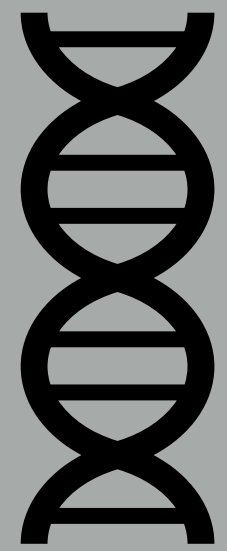


10^{150} designs

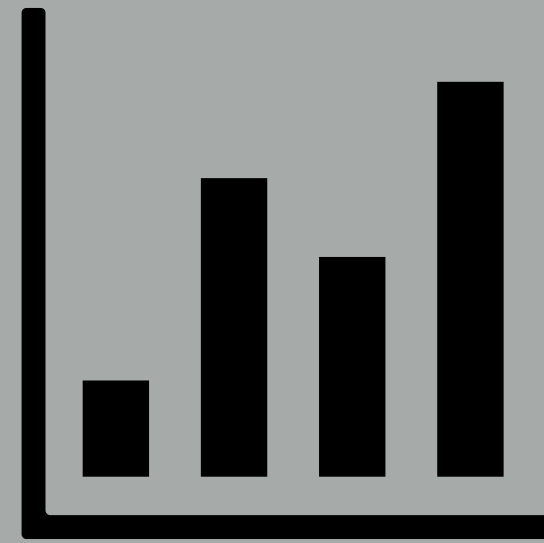


keep only the required bits
10x faster inference

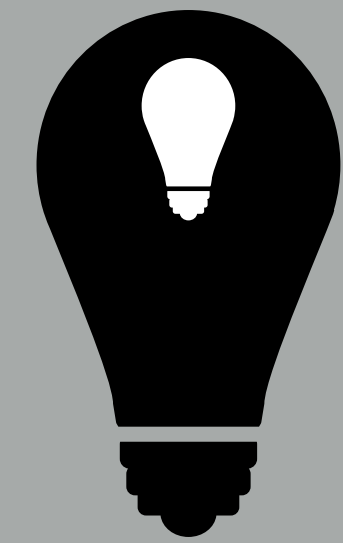
LARGE MODEL TRAINING



DESIGN SPACE

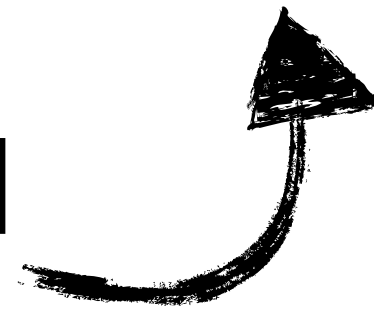


COST ESTIMATION



SEARCH

**GPU-based
Distributed**



TorchTitan

- No failed training
- New algos at scale



S. BING YAO
models/advisors



DON BATORY
modular synthesis

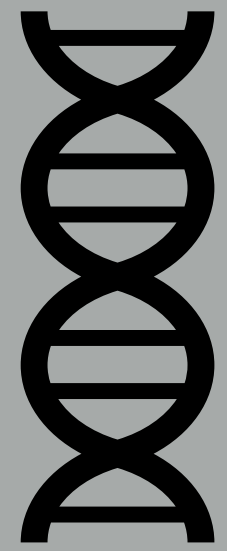


JOE HELLERSTEIN
extensible indexing

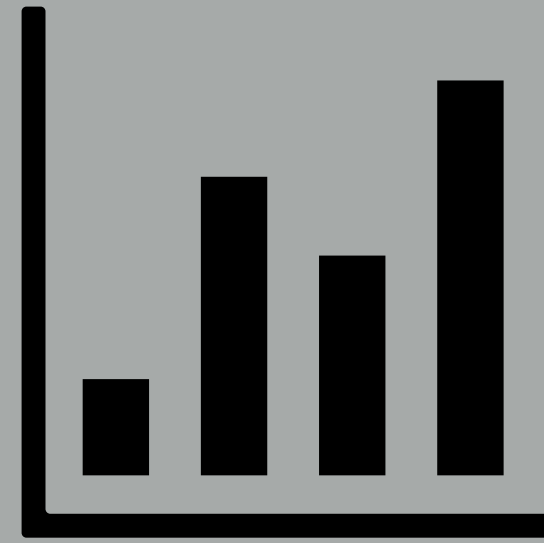


STEFAN MANEGOLD
model synthesis

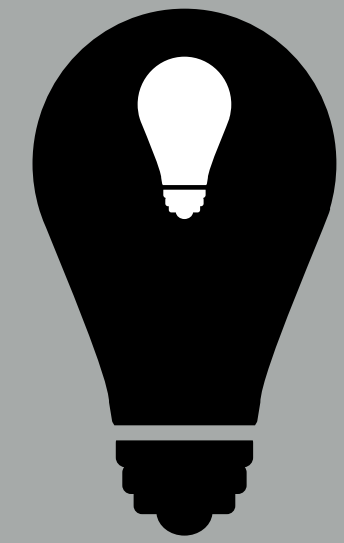
from one design at a time to design spaces



DESIGN SPACE



COST ESTIMATION

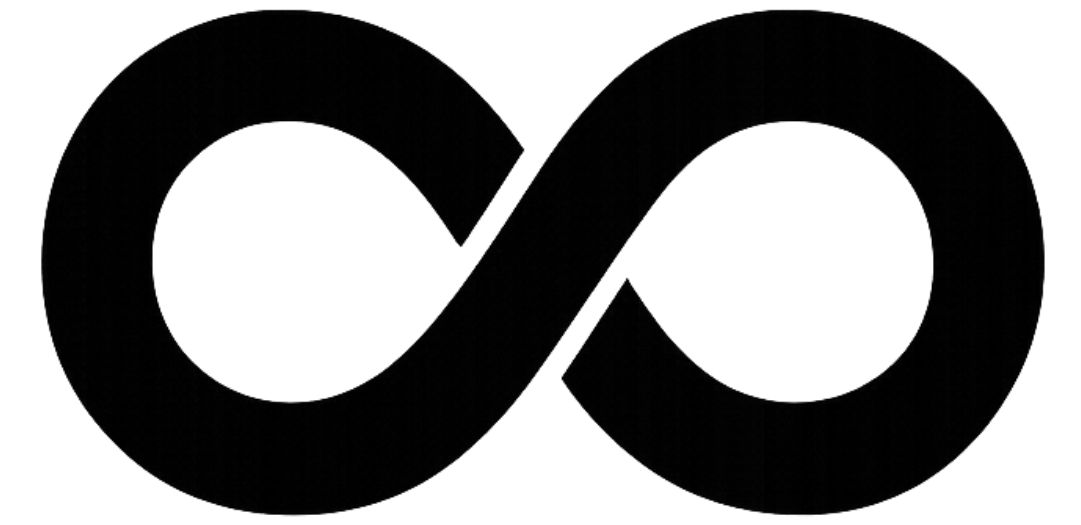


SEARCH

Tailored Systems, Accelerate Research

Next
Steps

Research Topics for all data fields



Long-term: A Generalized Design Space Engine

Short-term: LLMs Design Space

Startup: E2E Data-centric AI Platform



THANKS!

