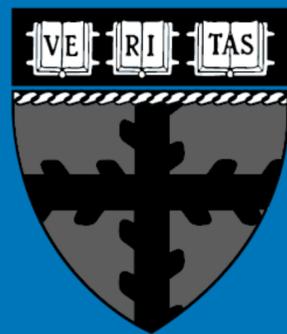
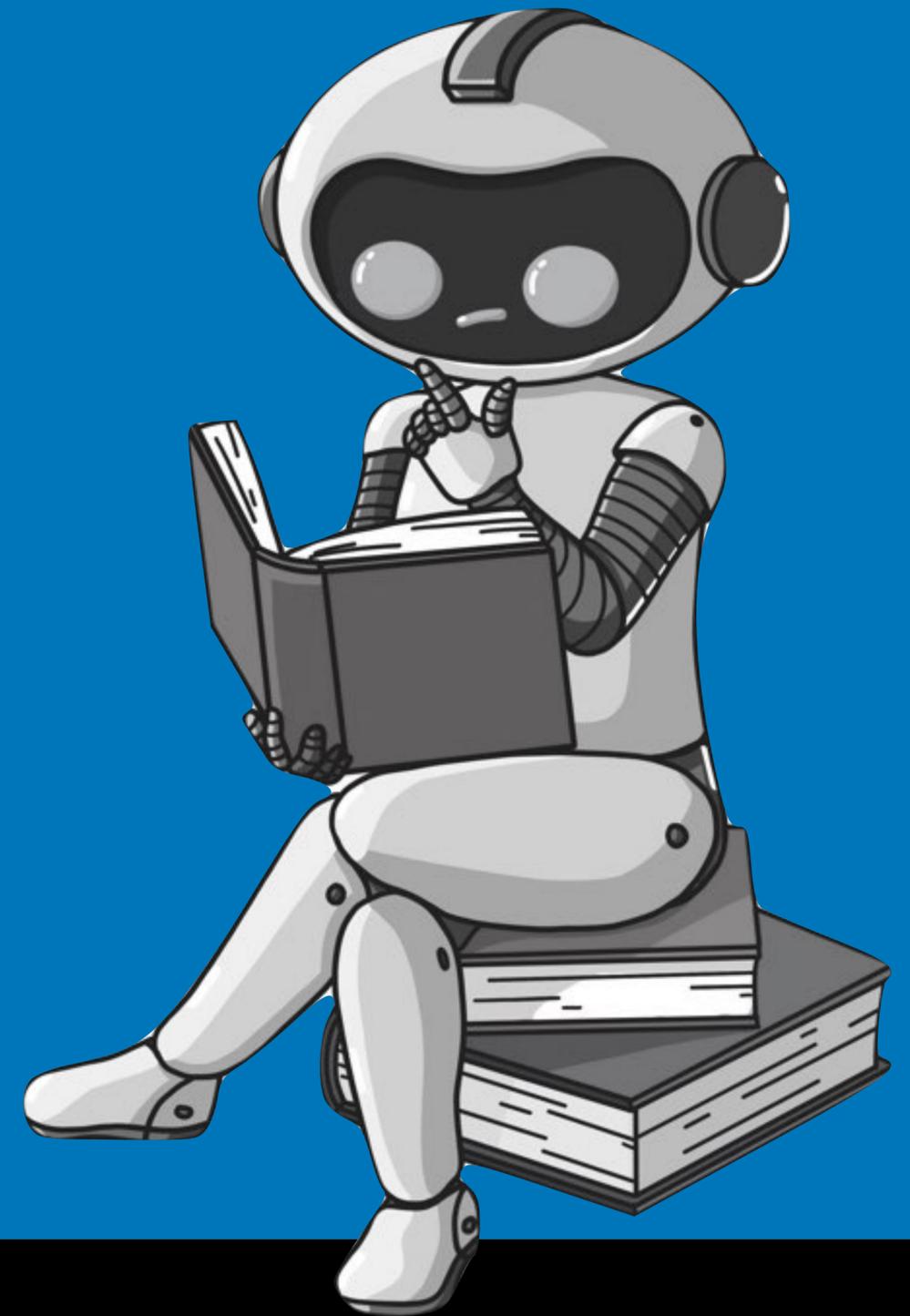


SELF-DESIGNING DATA SYSTEMS FOR THE AI ERA

Stratos
Idreos



DASlab
@ Harvard SEAS



What if we can reason about systems design?

What is a data system?

Why do we need self-designing systems?

A TYPICAL BIG DATA TASK

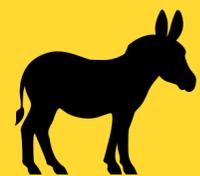
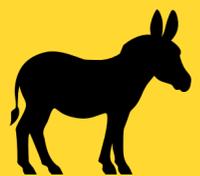
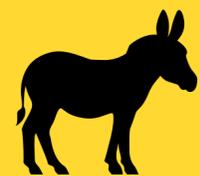
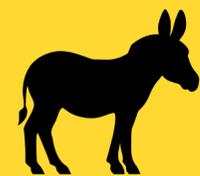
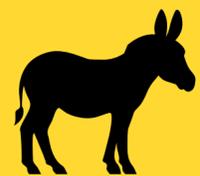
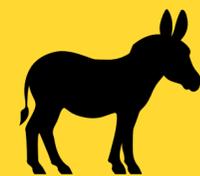
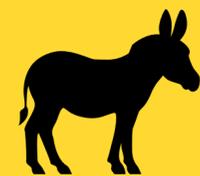
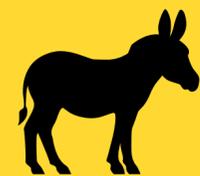
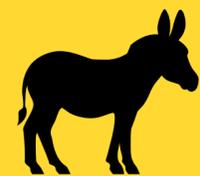
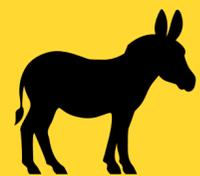
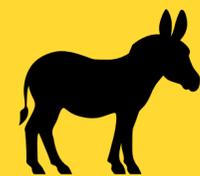
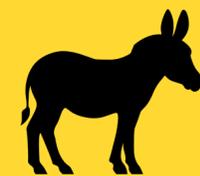
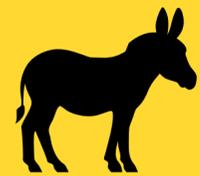
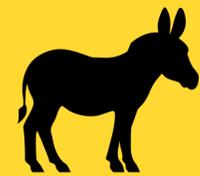
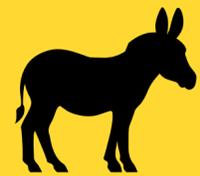
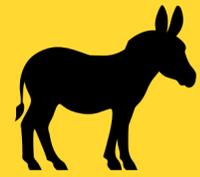
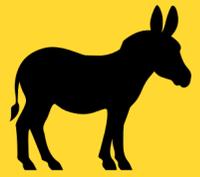
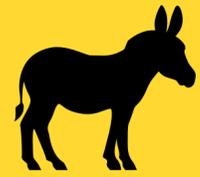
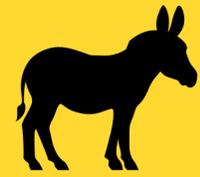
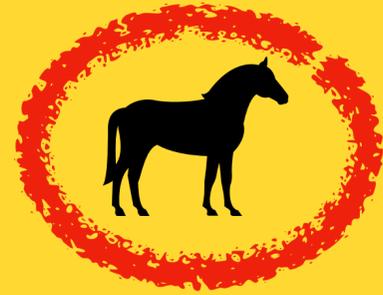
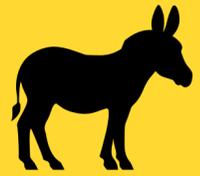
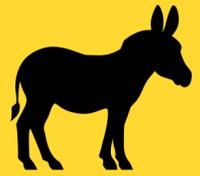
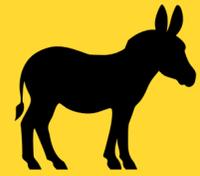
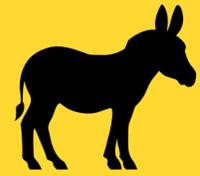
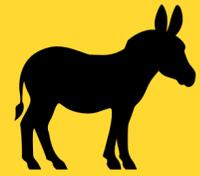
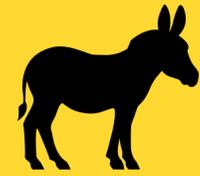
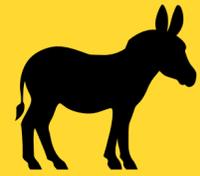
image analysis: e.g., detect the number of horses

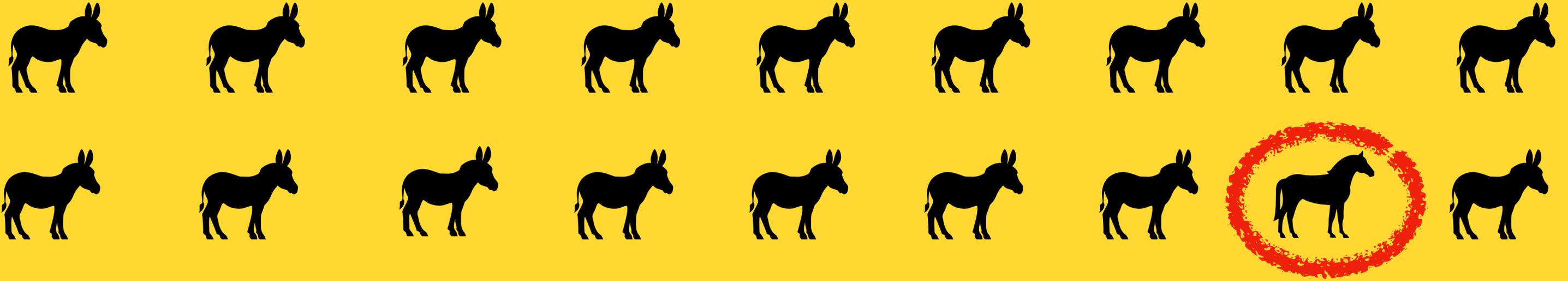


A TYPICAL BIG DATA TASK

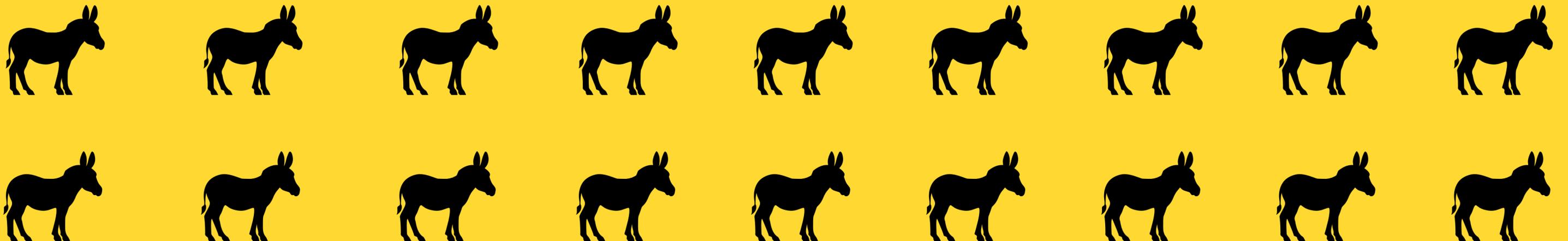
image analysis: e.g., detect the number of horses







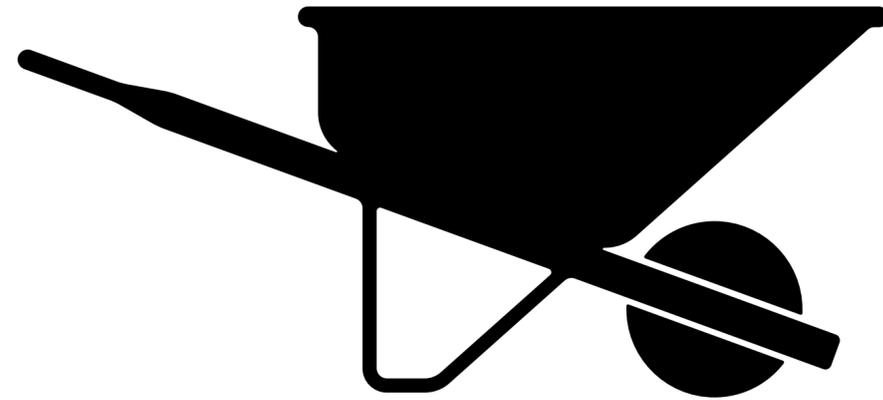
The core problem:
The size and organization of the data



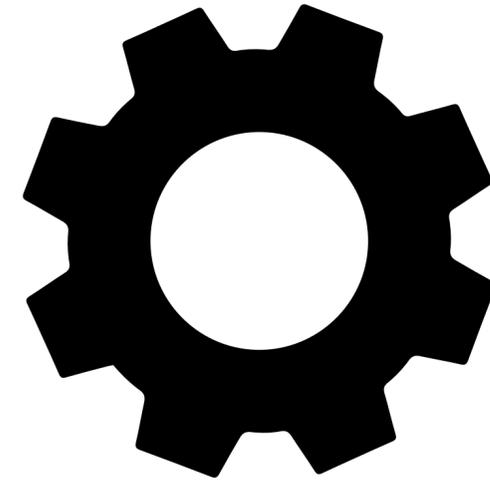
Three steps in big data/AI **regardless of application**



STORE



MOVE

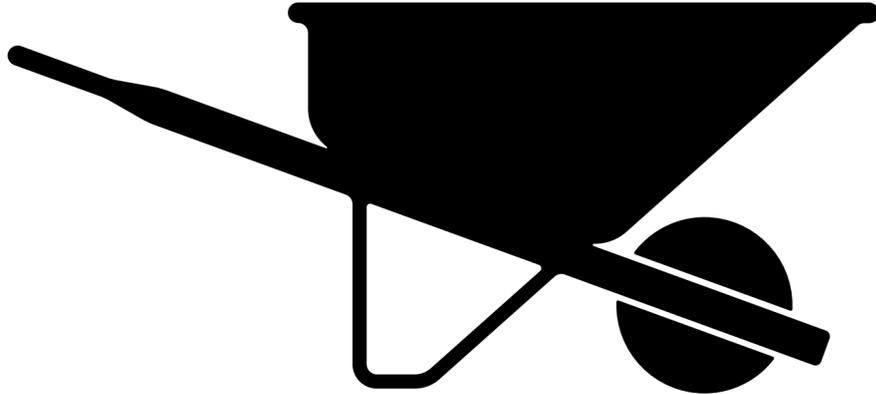


PROCESS

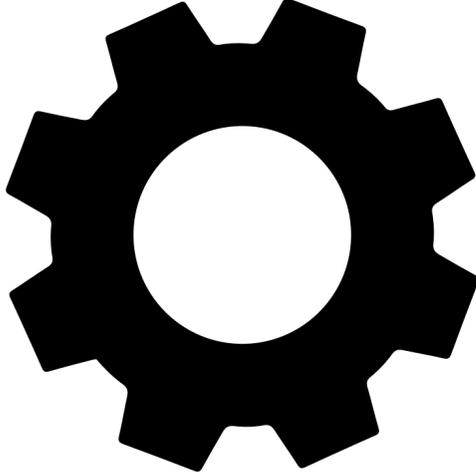
Three steps in big data/AI **regardless of application**



STORE



MOVE



PROCESS



How fast we can move and process data depends on the storage design decisions

What is a data system?

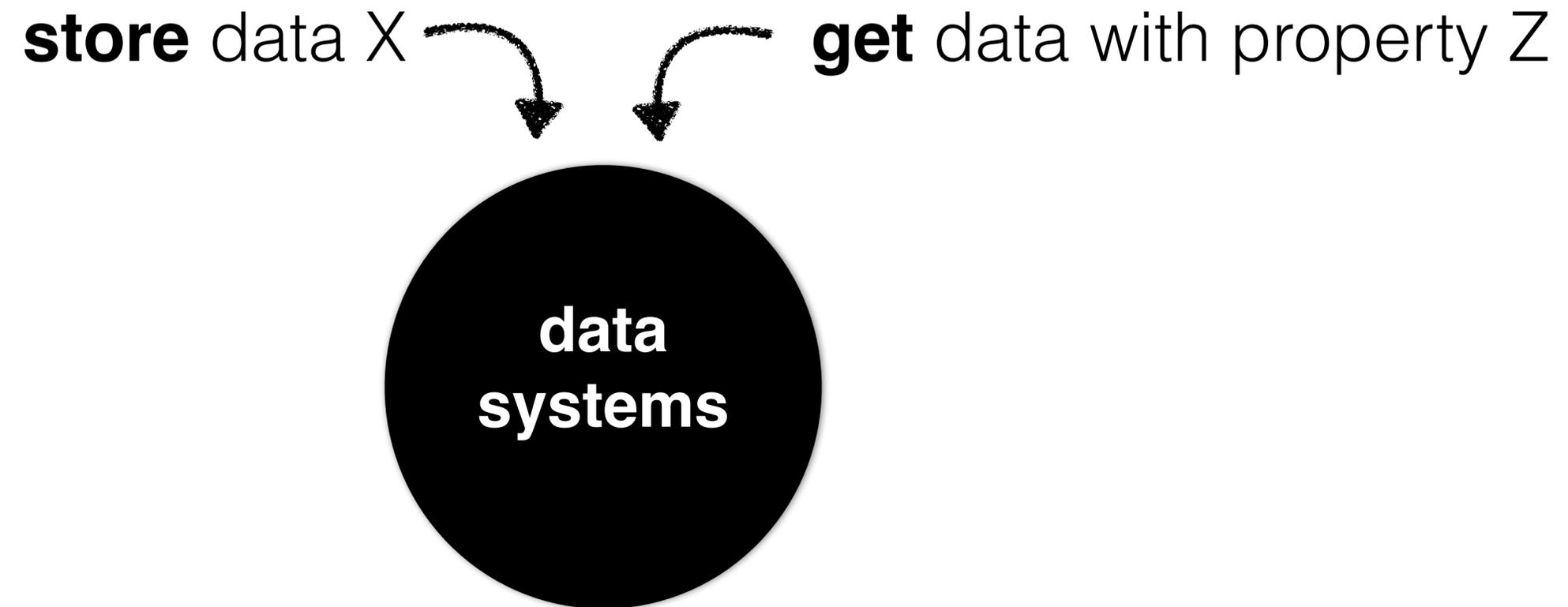
A data system is an end-to-end software system that:
manages storage, data movement, and provides access to data

What is a data system?

A data system is an **end-to-end software system** that:
manages storage, data movement, and provides access to data

1. For decades: data systems = SQL DBs
but with big data, the need for fast data systems is drastically broader than SQL

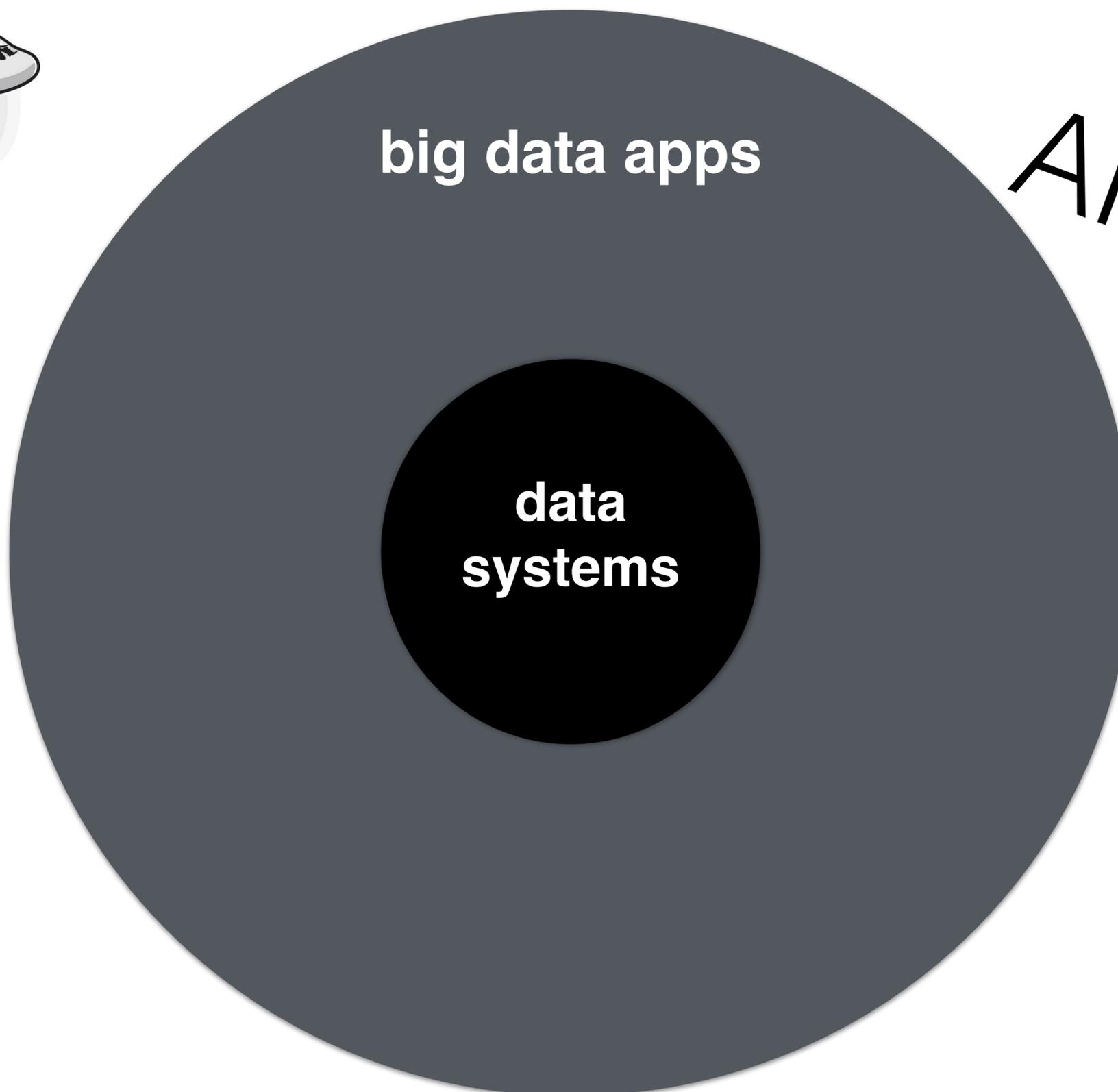
broader than SQL



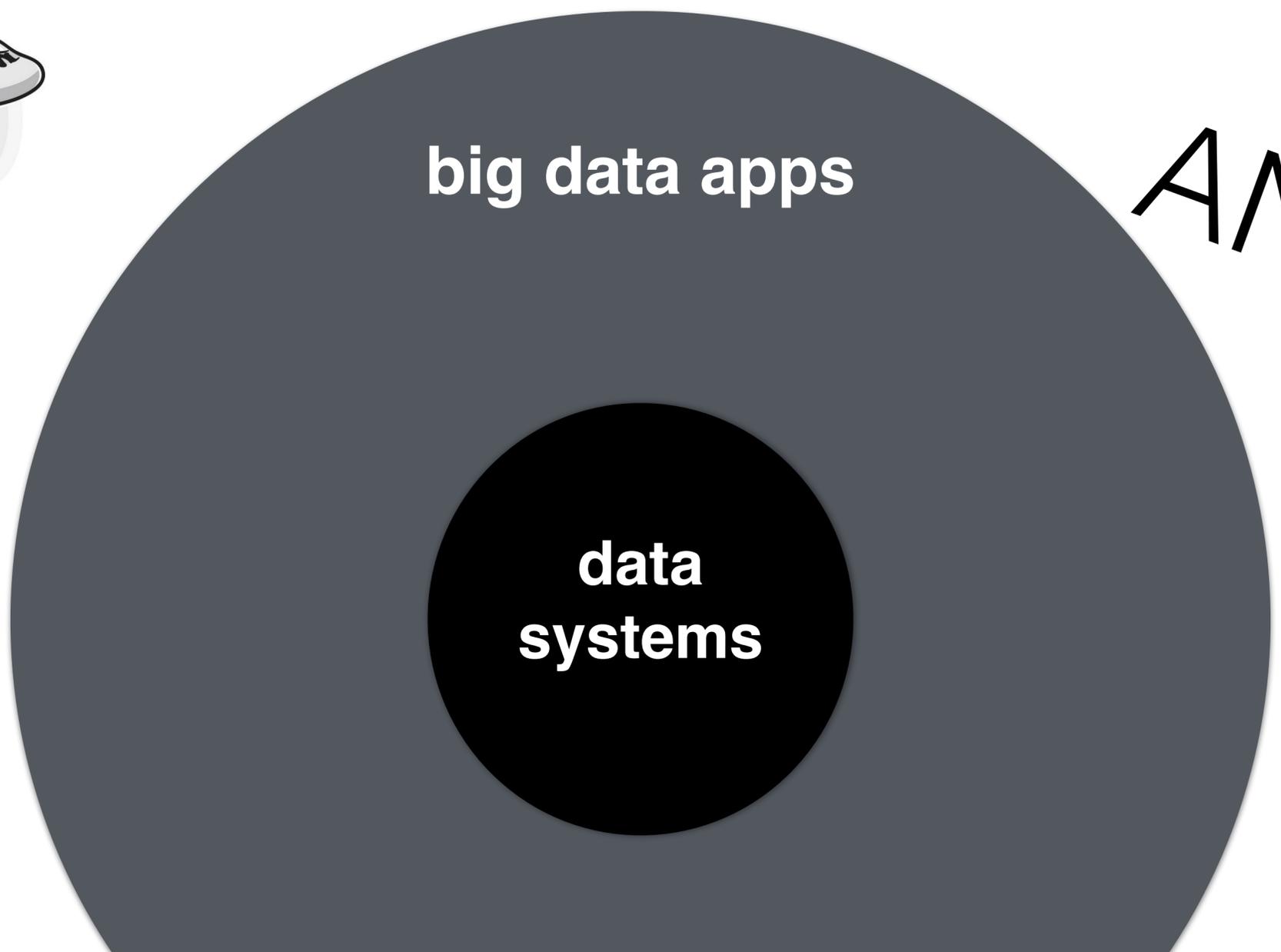
broader than SQL

ANALYTICS

AI

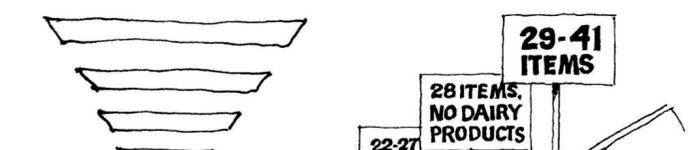


broader than SQL



ANALYTICS

AI



New data systems to handle new requirements

broader than SQL

TRANSACTIONS

Deposit money to my bank account

Transfer money from ... to...



broader than SQL

TRANSACTIONS

Deposit money to my bank account

Transfer money from ... to...



ANALYTICS

How much do customers
of X spent on average every month?

TRANSACTIONS

Deposit money to my bank account

Transfer money from ... to...



ANALYTICS

How much do customers
of X spent on average every month?

AI

Is this transaction legal?

Should we give a loan to customer X?

SOCIAL NETWORKS: REVIEWS/POSTS

How many costumers on average
leave a 4 star review or better?

SOCIAL NETWORKS: REVIEWS/POSTS

How many costumers on average
leave a 4 star review or better?

AI

Is this new review a legitimate one?

SOCIAL NETWORKS: REVIEWS/POSTS

How many costumers on average
leave a 4 star review or better?

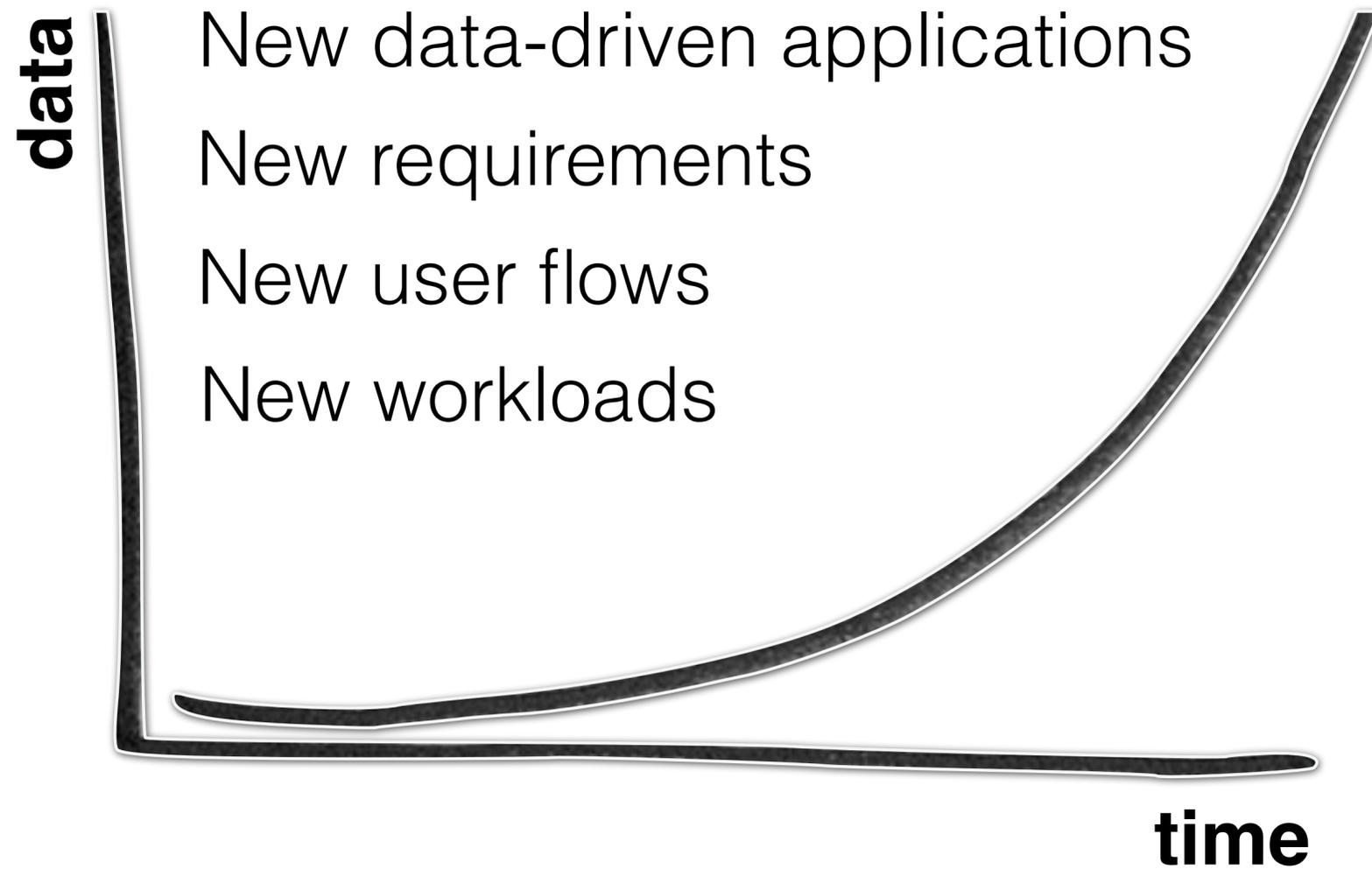
AI

Is this new review a legitimate one?

COMMUTING

Compute price for next Uber ride

broader than SQL

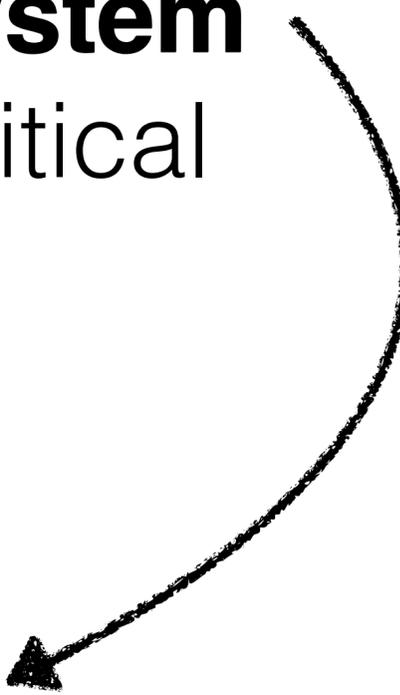


**The need for
data systems
grows with data**

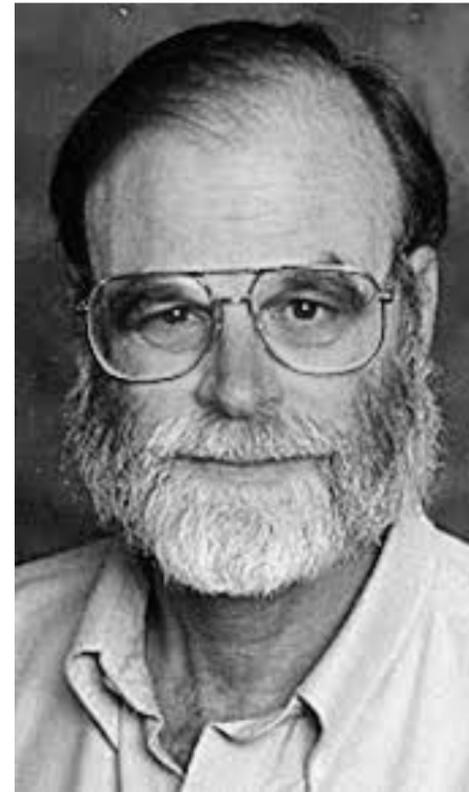
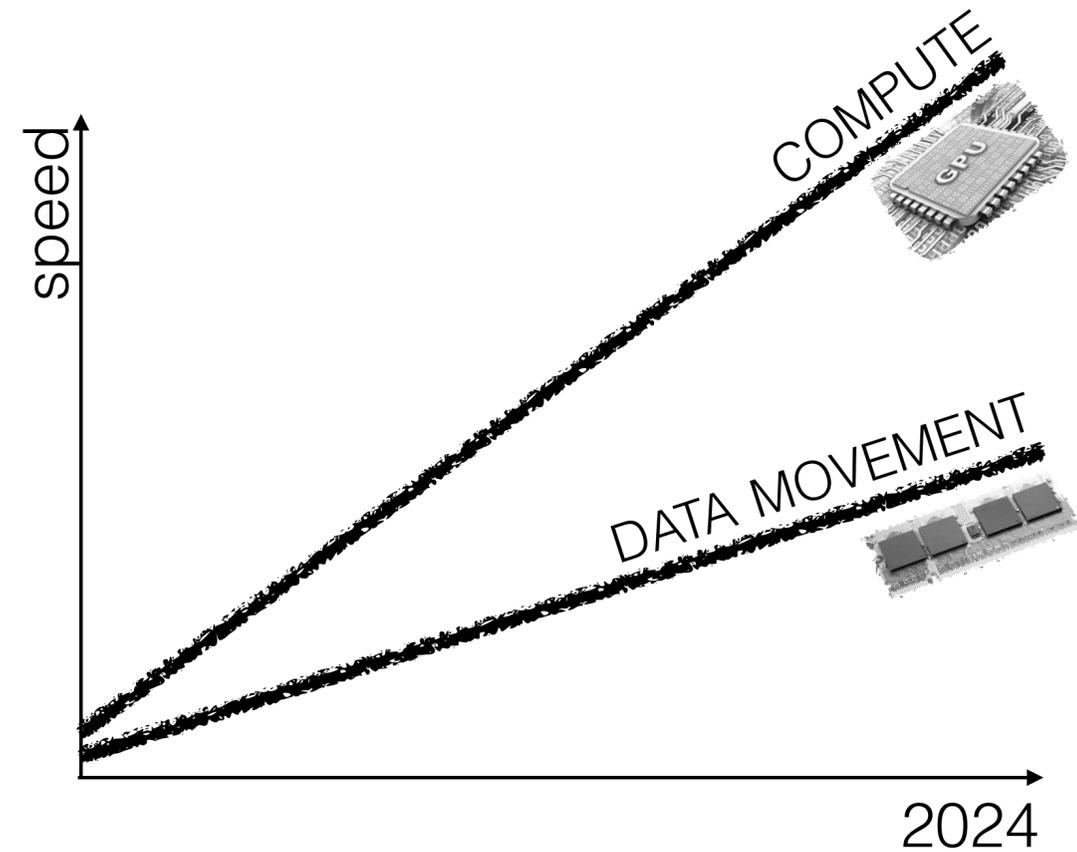
2. As data grows, having the right data system
for each application is increasingly more critical

2. As data grows, having the right data system
for each application is increasingly more critical

system architecture
it starts with storage



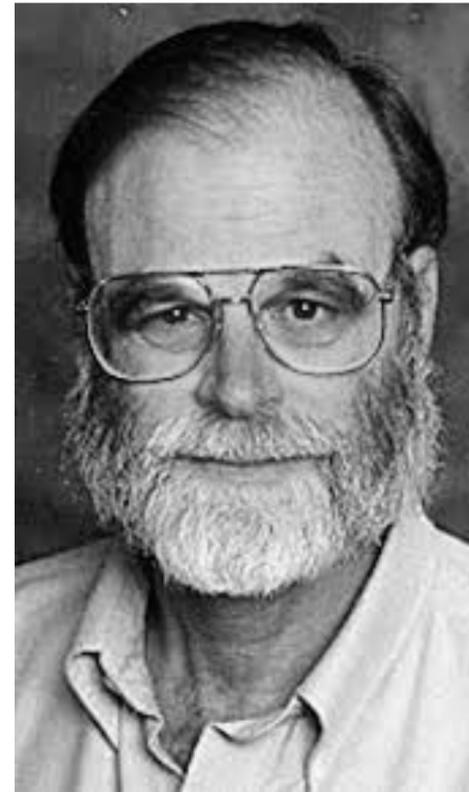
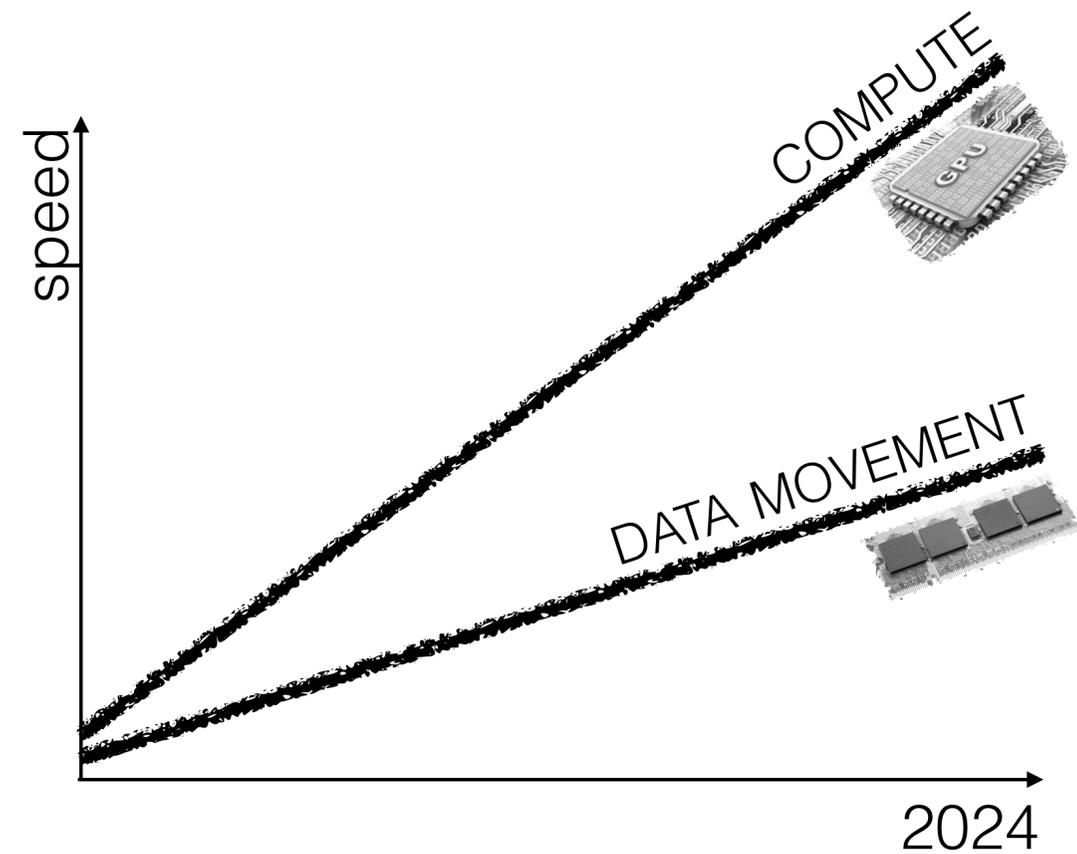
the right data system



register = this room
caches = this city
memory = nearby city
disk = Pluto

Jim Gray, Turing Award 1998

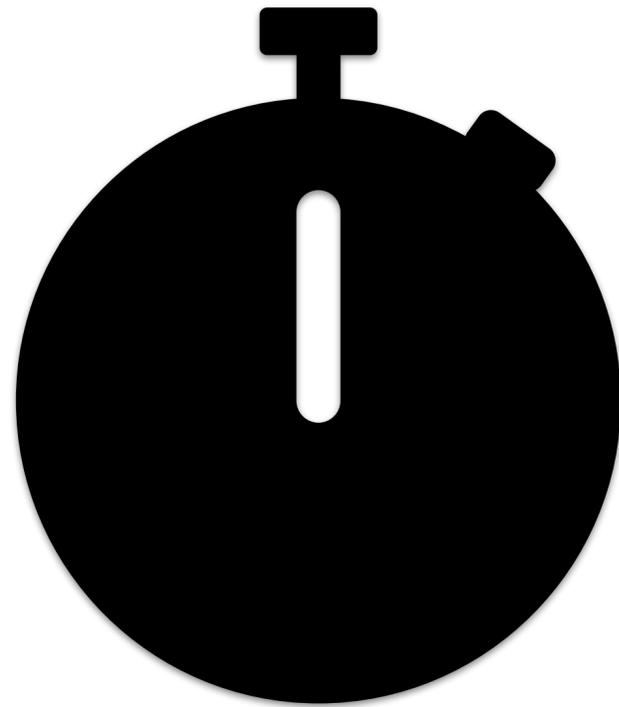
the right data system



register = this room
caches = this city
memory = nearby city
disk = Pluto

Jim Gray, Turing Award 1998

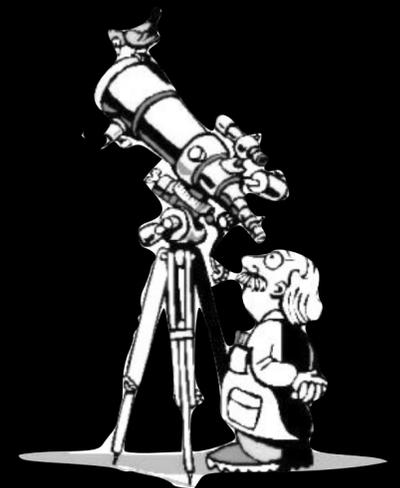
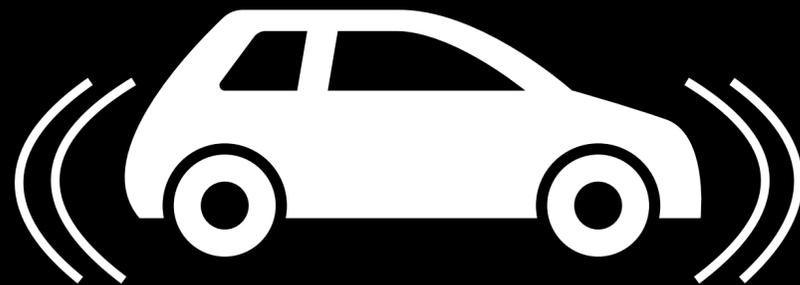
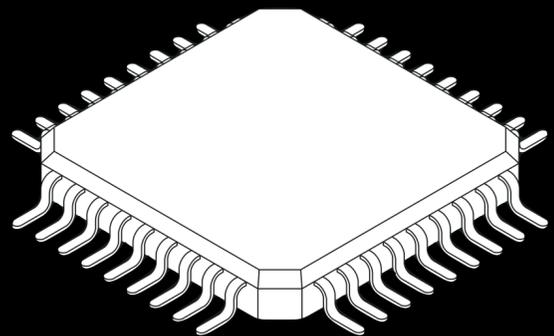
Data movement dominates everything



**70-80% of processing costs
go into data movement**

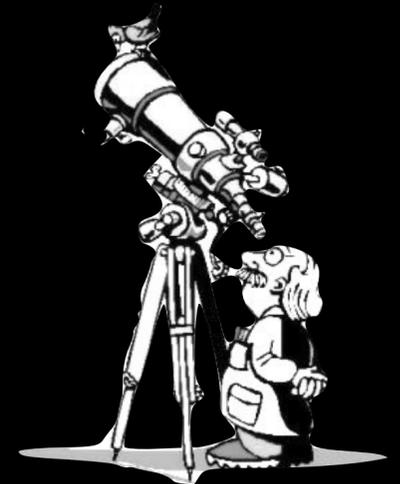
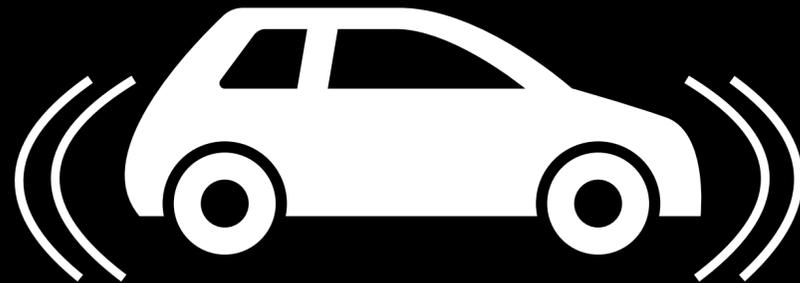
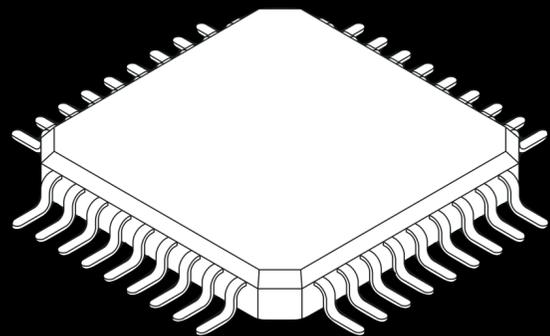
**computational hardware
utilization: only 30-50%**

The problem: as the big data/AI world keeps changing...

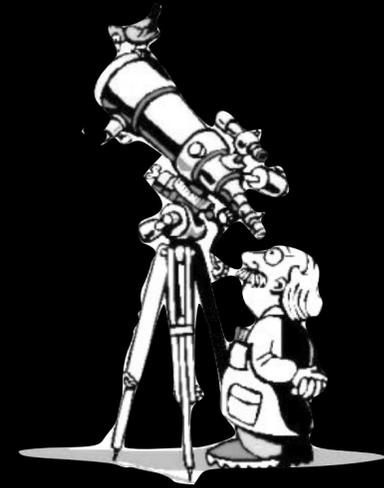
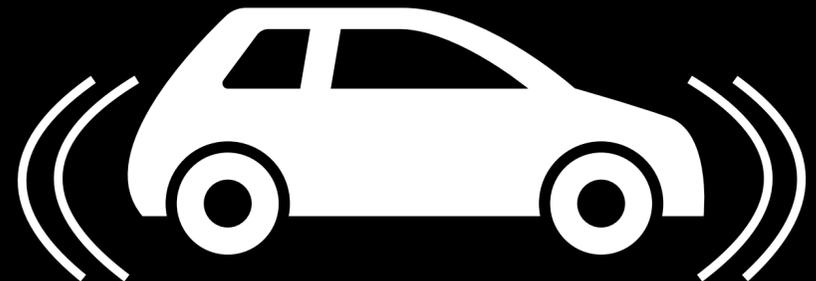
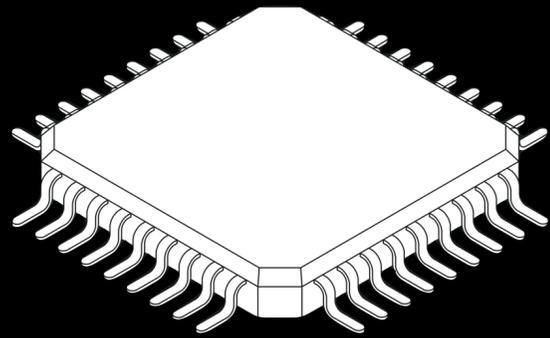


The problem: as the big data/AI world keeps changing...

there is a continuous need for new data systems
but it is **extremely hard to design & build new systems**



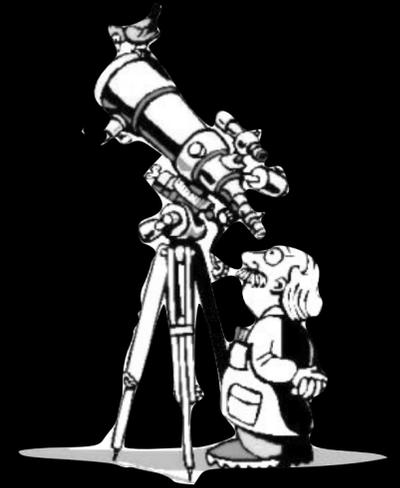
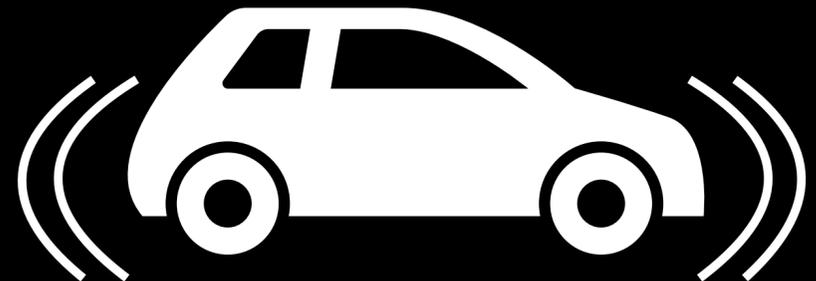
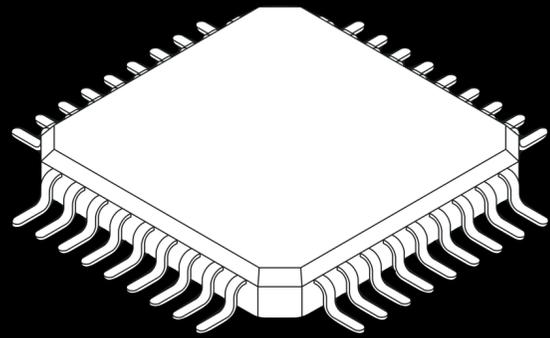
How do we design a data system that is **X times faster for a workload W?**



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?**



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

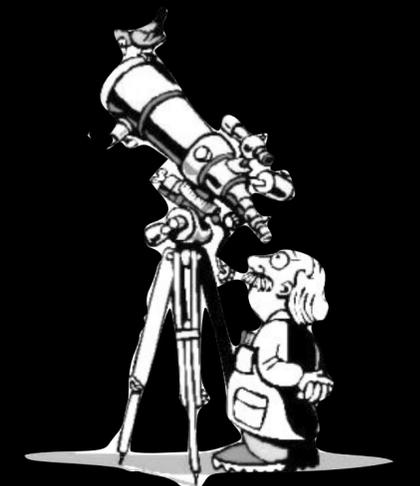
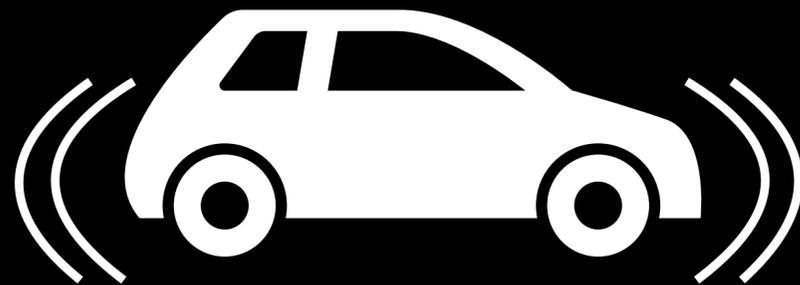
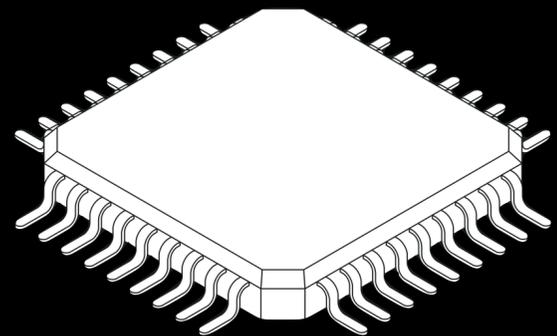
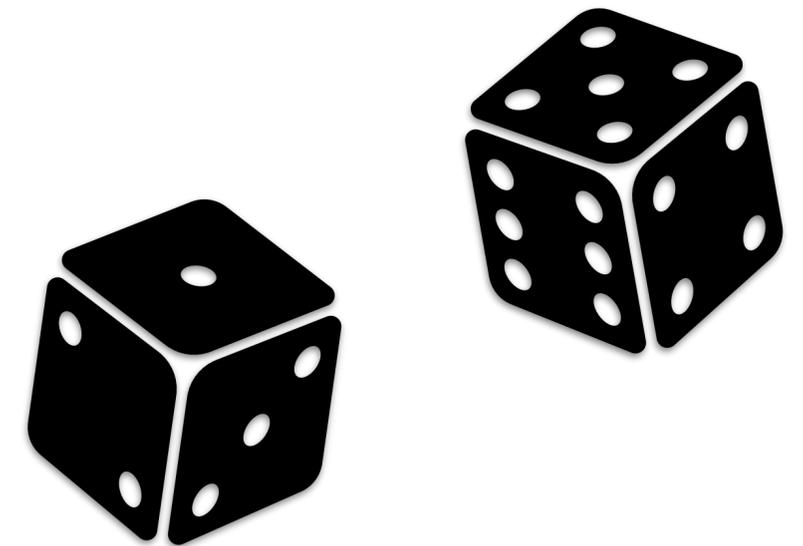


How do we design a data system that allows for control of **cloud cost?**

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

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

What will **break** our system?

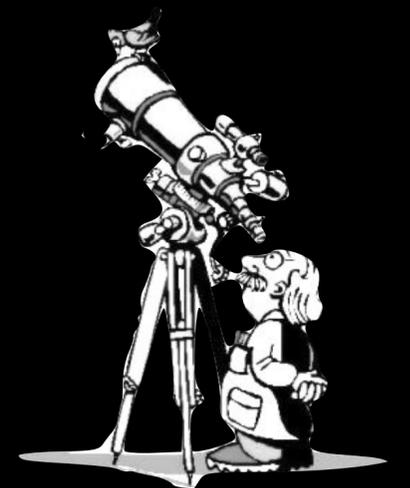
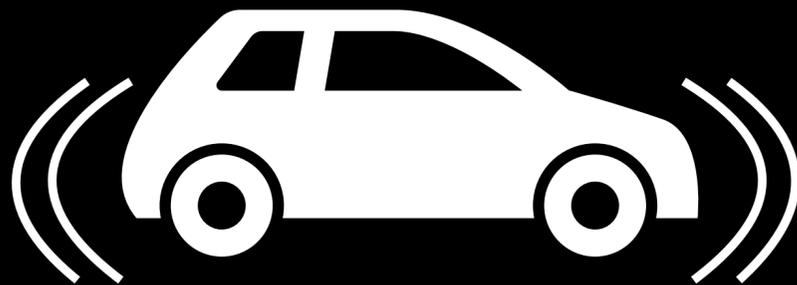
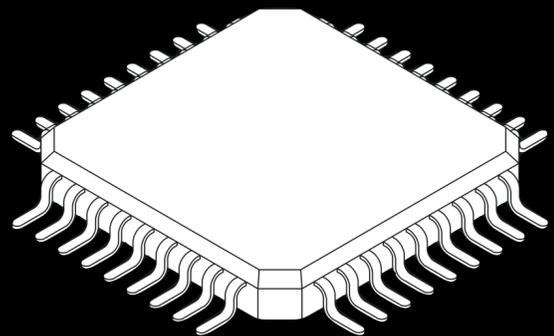
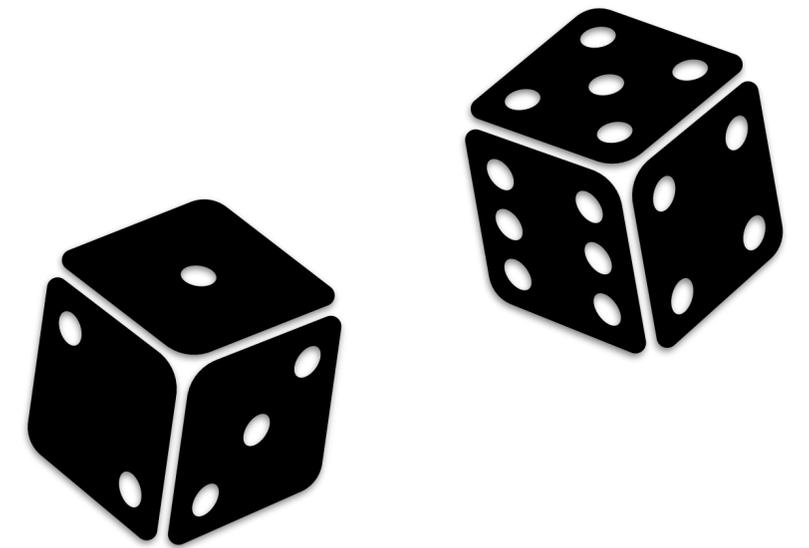


BOTTLENECK: SUB-OPTIMAL DATA SYSTEMS

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

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

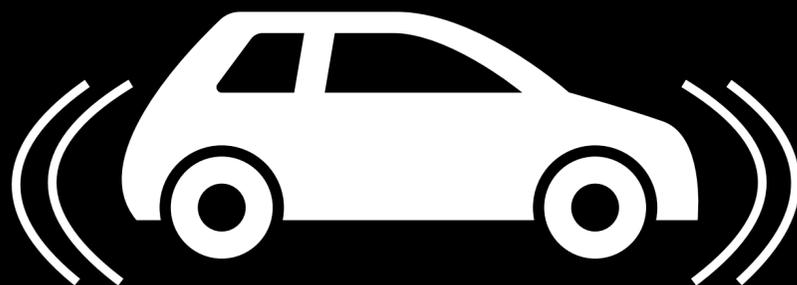
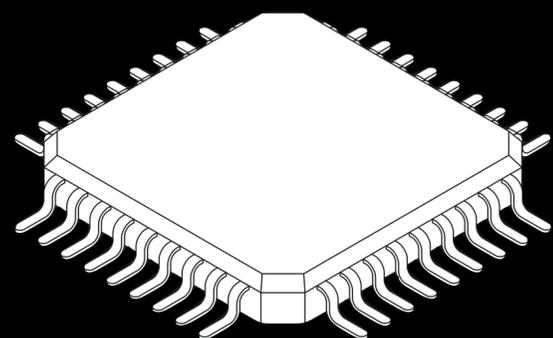
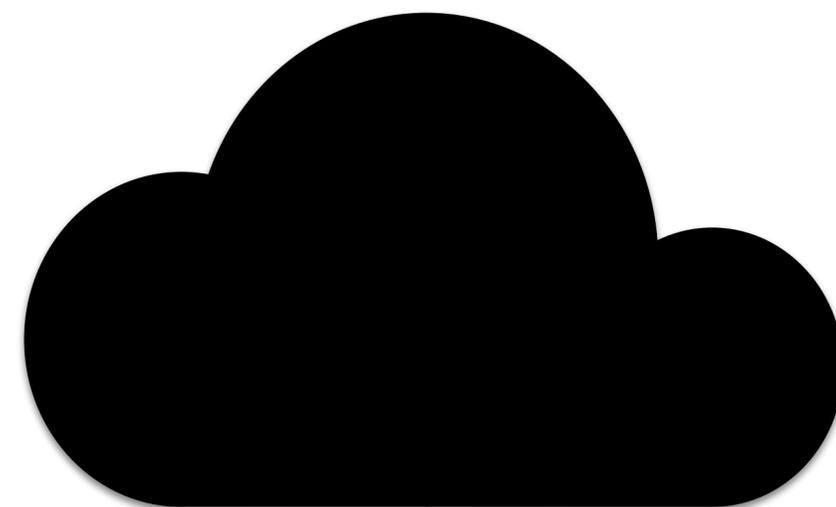
What will **break** our system?



BOTTLENECK: SUB-OPTIMAL DATA SYSTEMS

huge cloud cost

environmental impact

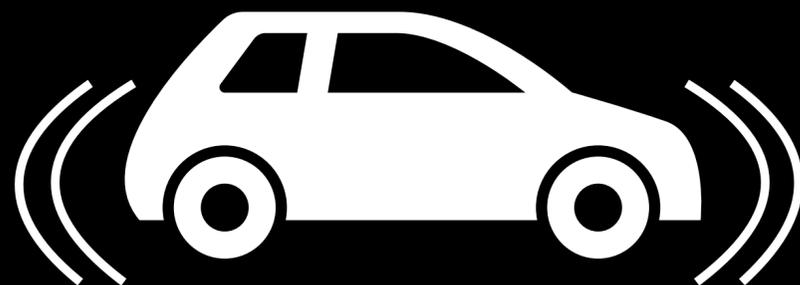
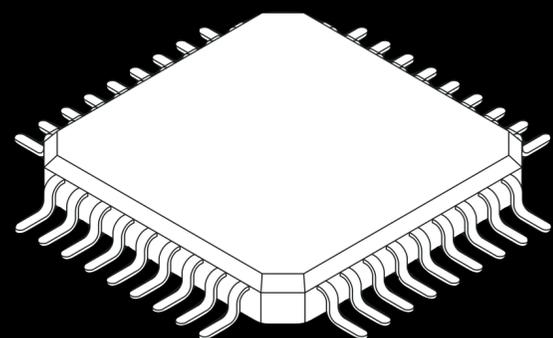


BOTTLENECK: SUB-OPTIMAL DATA SYSTEMS

huge cloud cost

expensive transitions

environmental impact



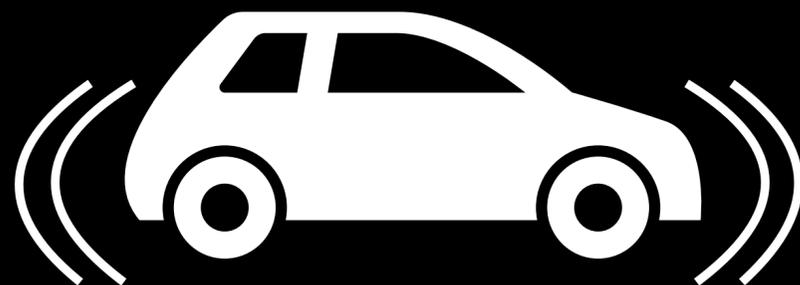
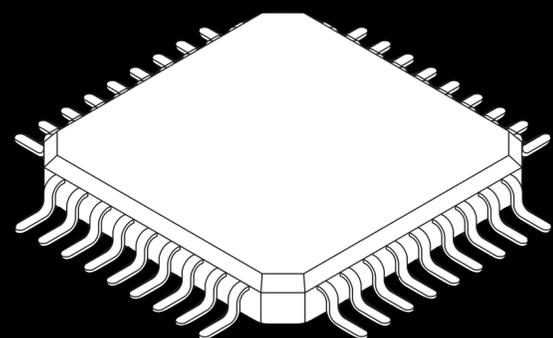
BOTTLENECK: SUB-OPTIMAL DATA SYSTEMS

huge cloud cost

expensive transitions

application feasibility

environmental impact

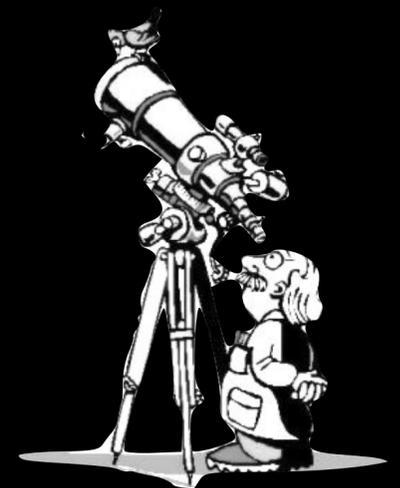
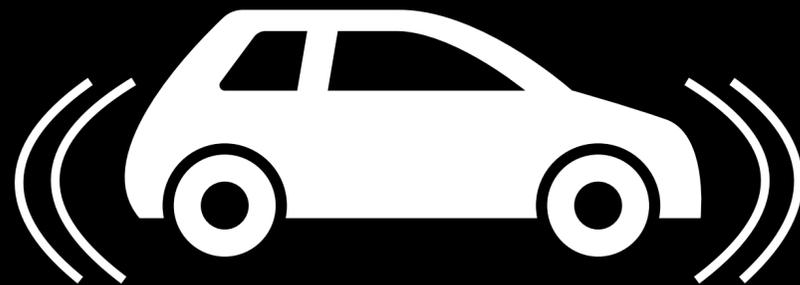
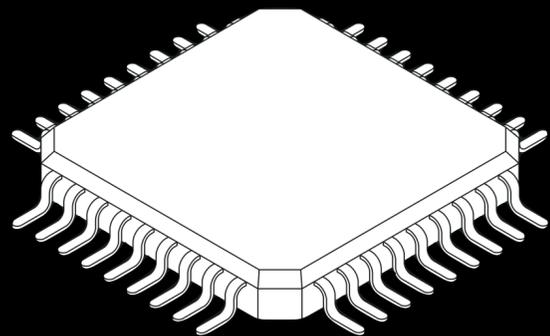


BOTTLENECK: SUB-OPTIMAL DATA SYSTEMS

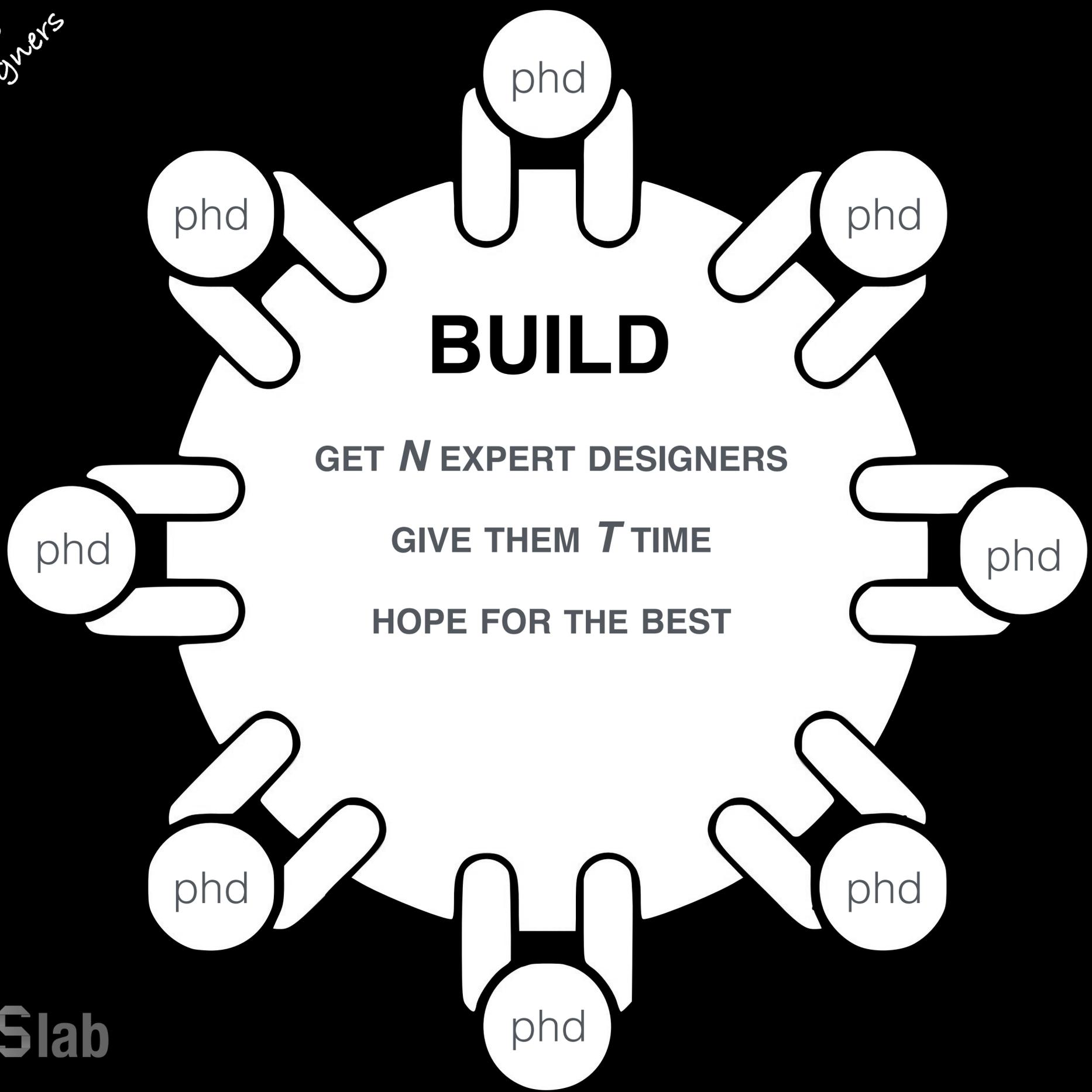
huge cloud cost expensive transitions application feasibility environmental impact

complexity

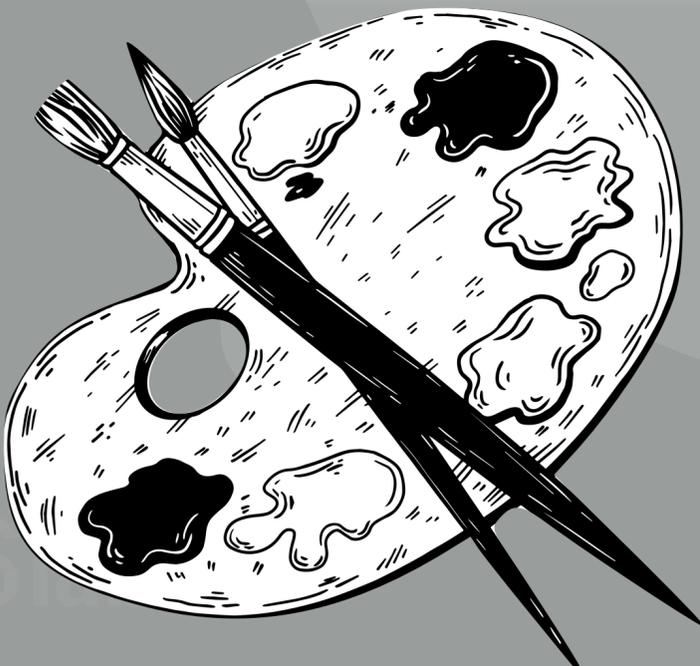
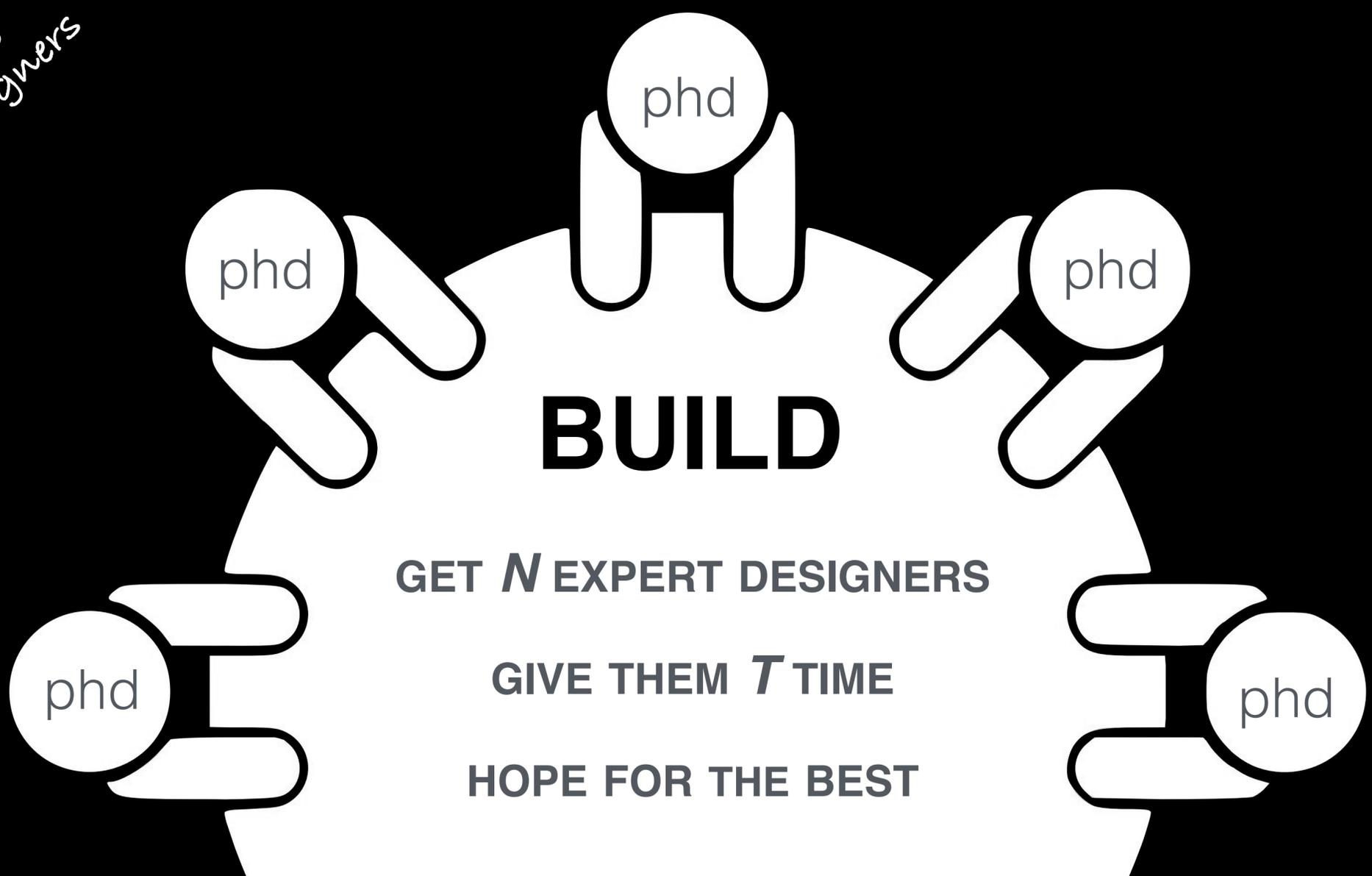
how we **BUILD** systems



the dining
systems designers

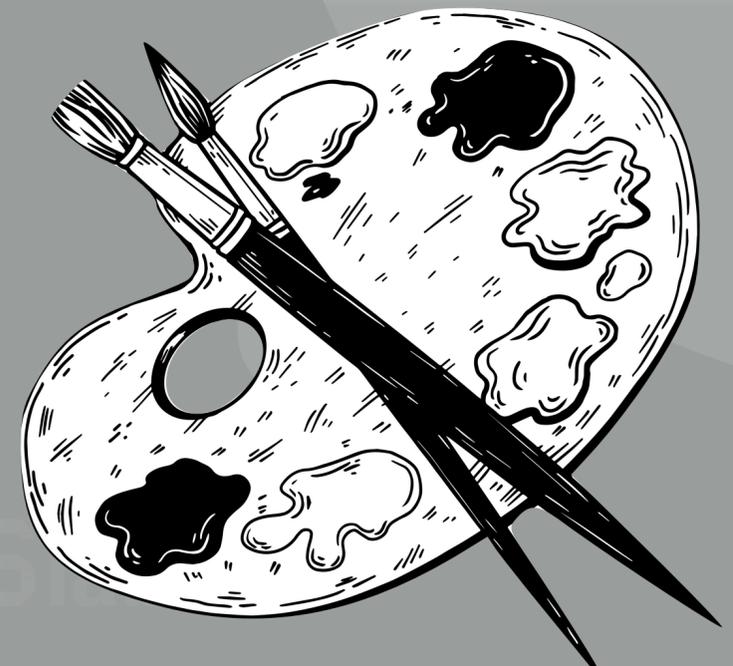
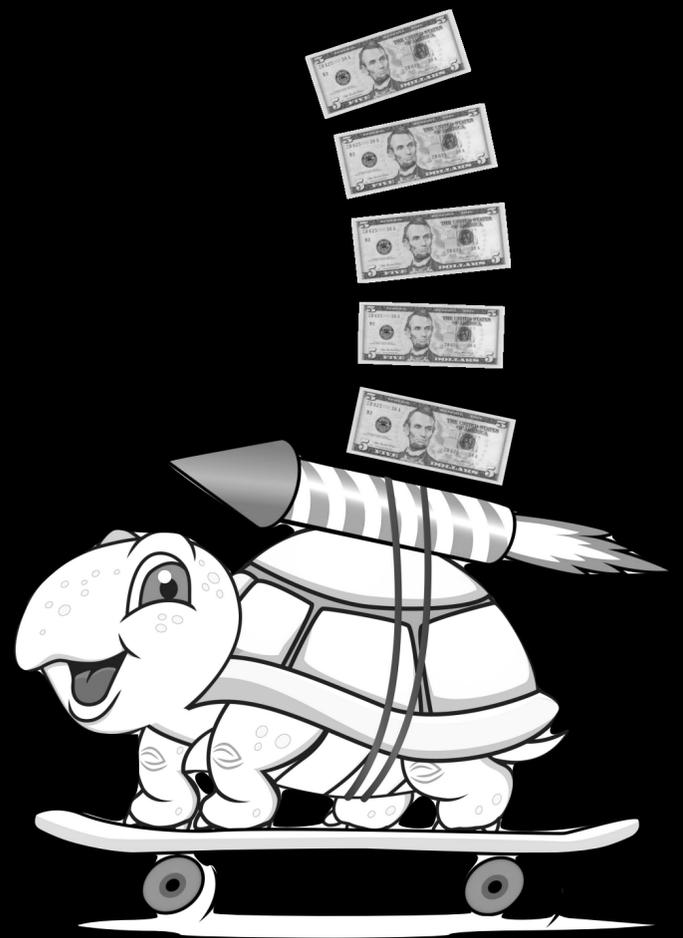
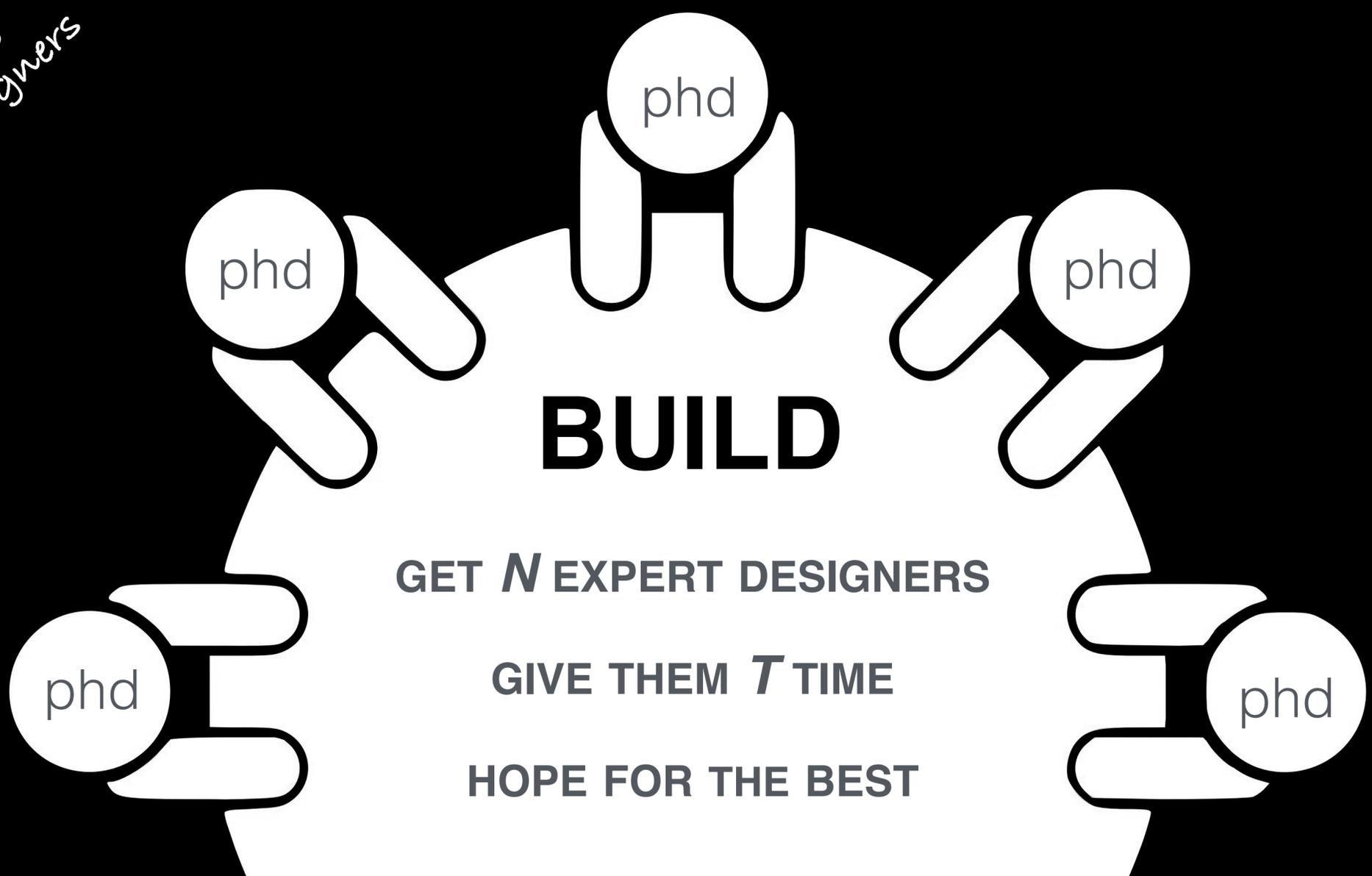


the dining
systems designers



design is an art

the dining
systems designers



design is an art

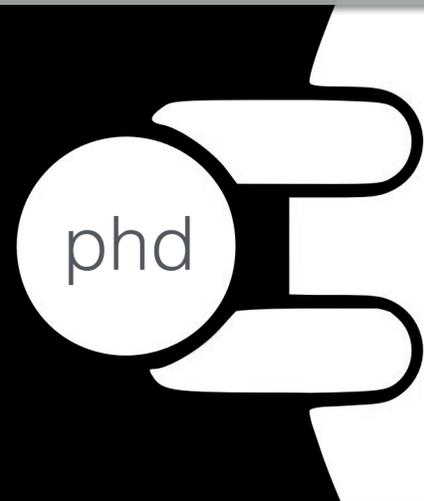
the dimming systems designers

phd

Design: 6-7 years

Reasoning: months/impossible

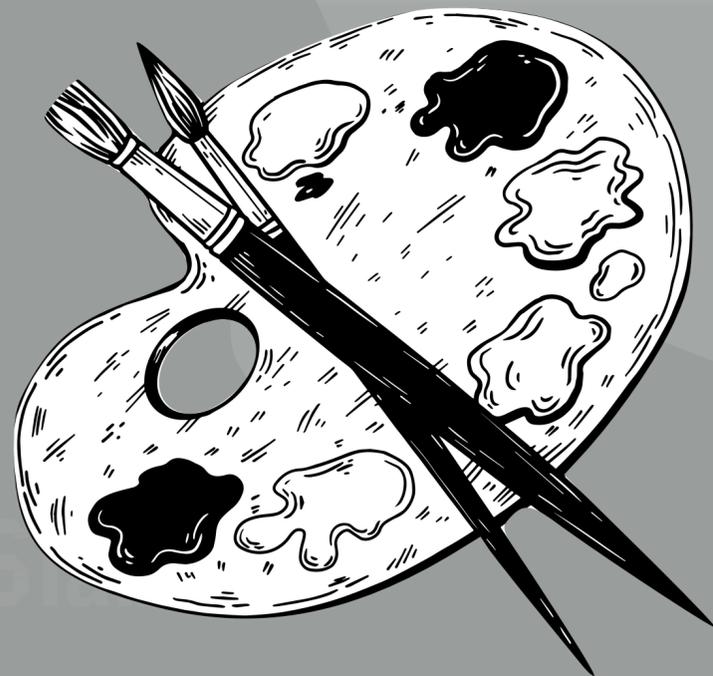
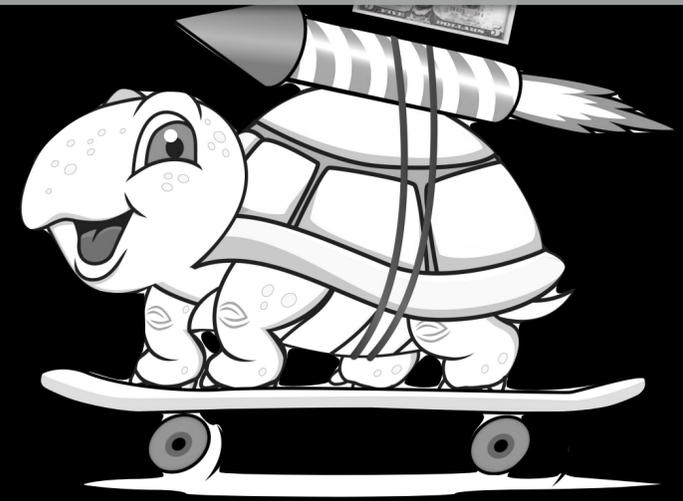
BUILD



GET N EXPERT DESIGNERS

GIVE THEM T TIME

HOPE FOR THE BEST



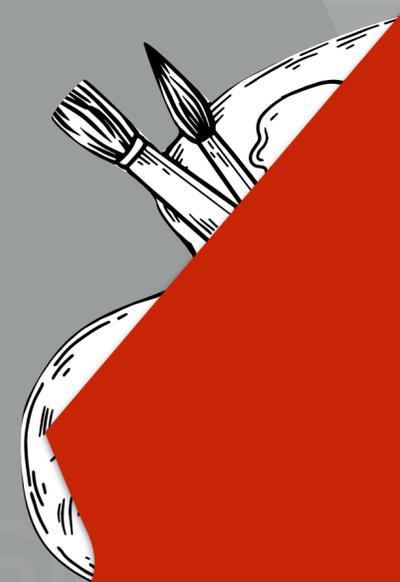
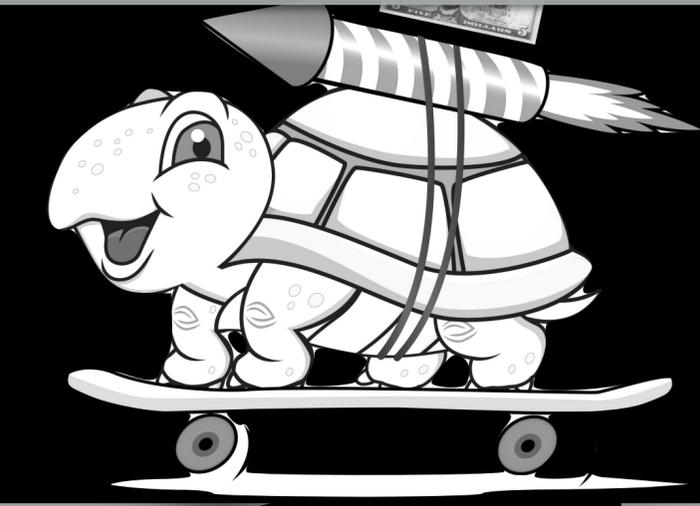
design is an art

the dimming
systems designers

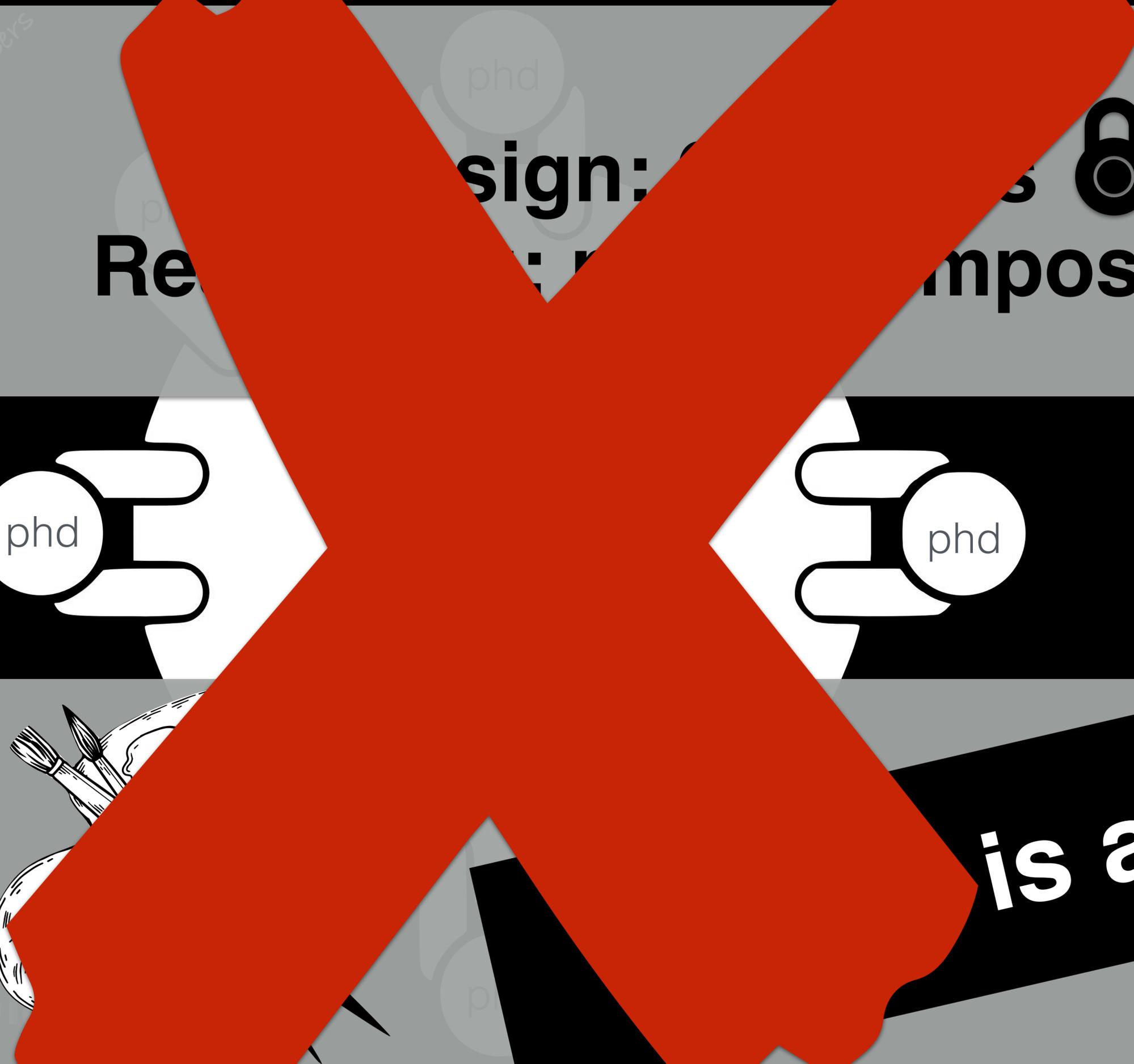
sign: s



Re: impossible



is an art



DAS

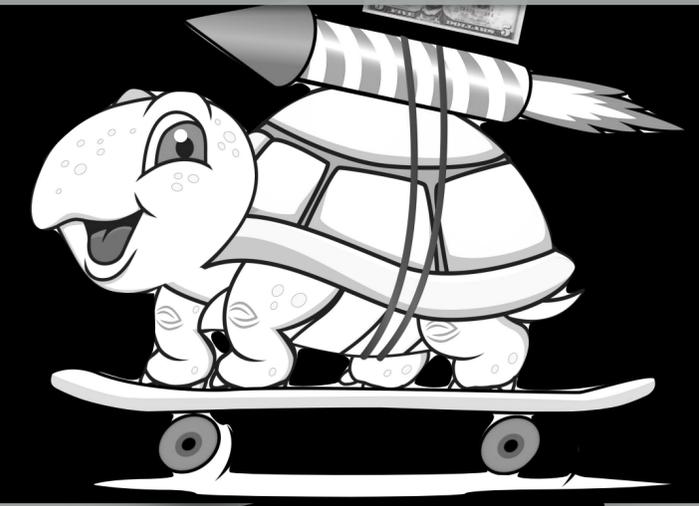
the dimming
systems designers

design: 
Re: impossible



data
hardware
applications

years

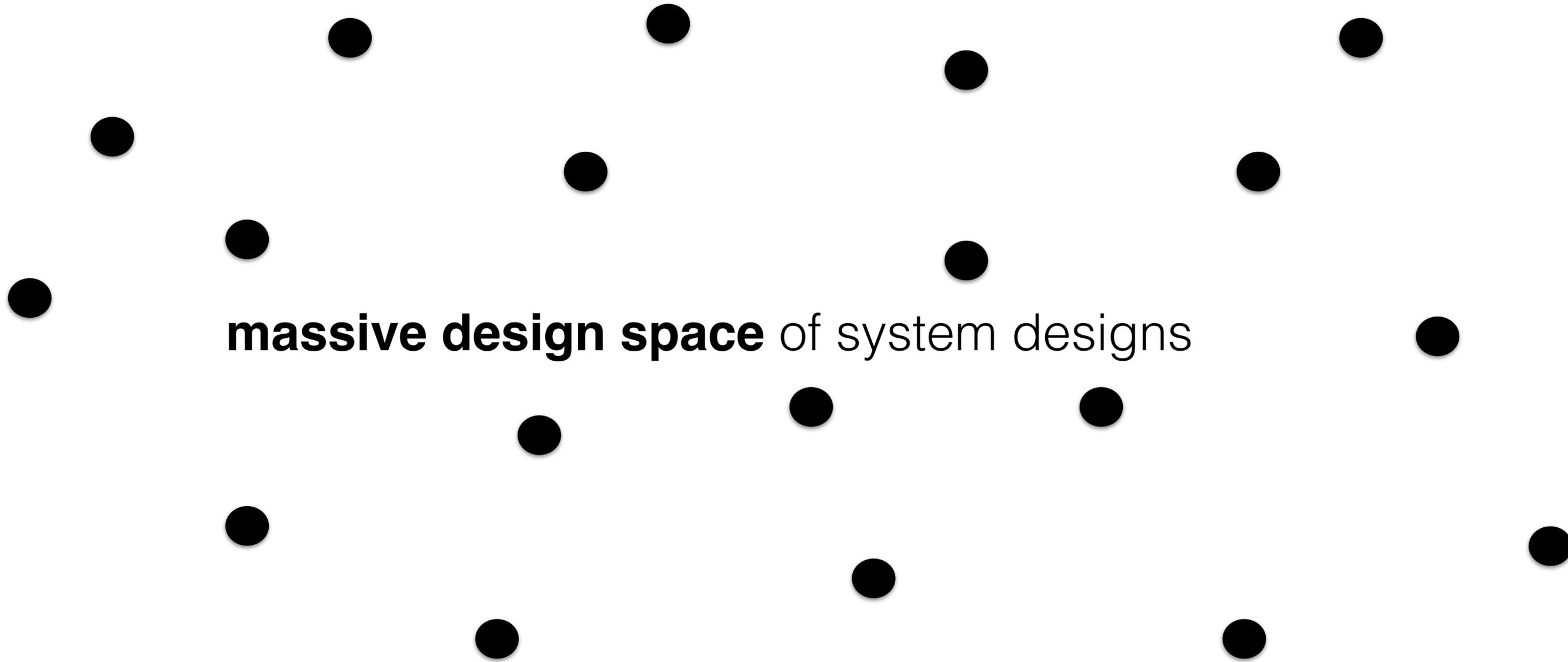


is an art

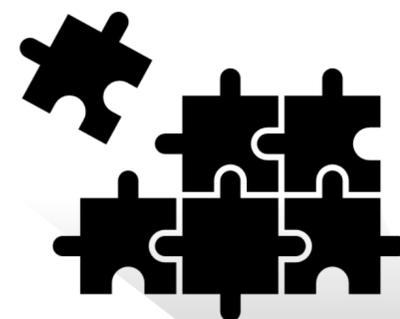
DAS

SELF-DESIGNING SYSTEMS

Automatically invent & build the perfect system for any new application



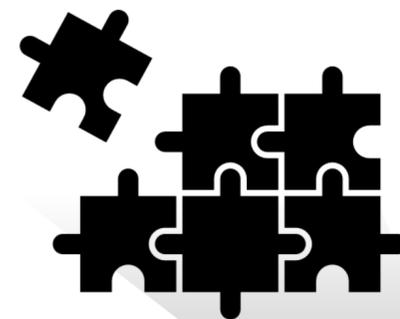
massive design space of system designs



system=
a set of low-level
design decisions

massive design space of system designs

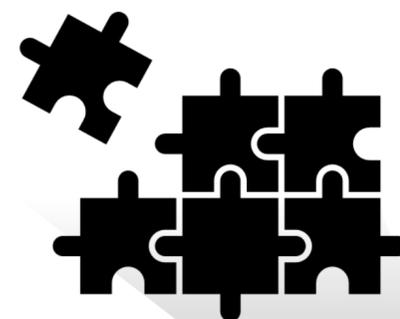
few existing designs



system=
a set of low-level
design decisions

massive design space of system designs

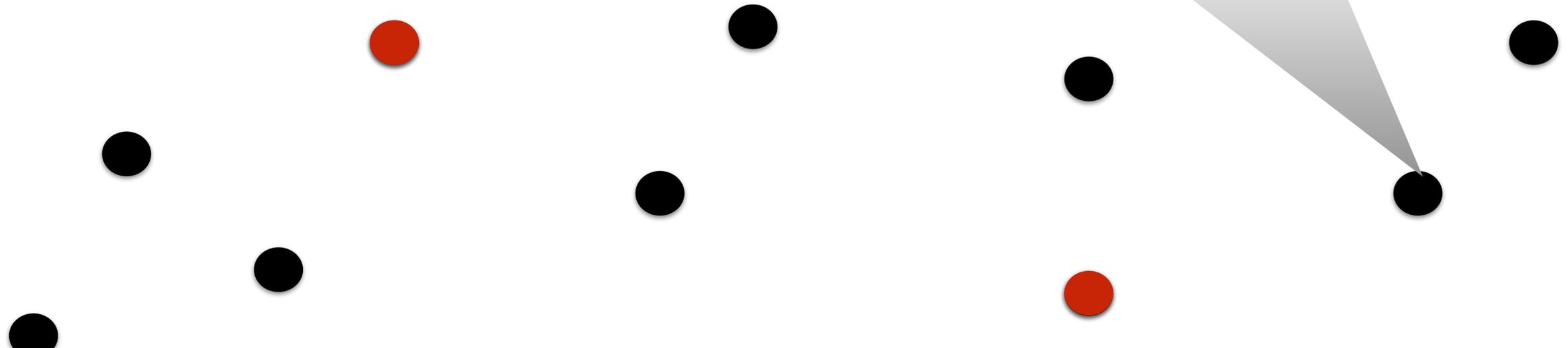
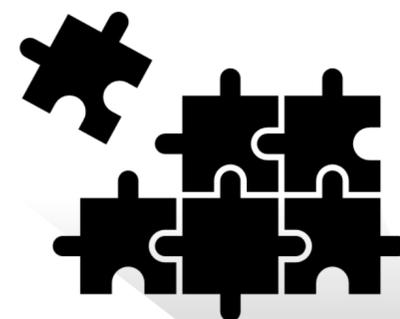
few existing designs



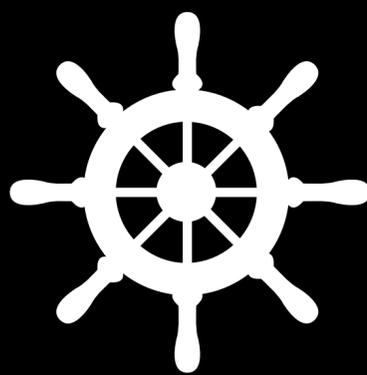
system=
a set of low-level
design decisions

massive design space of system designs



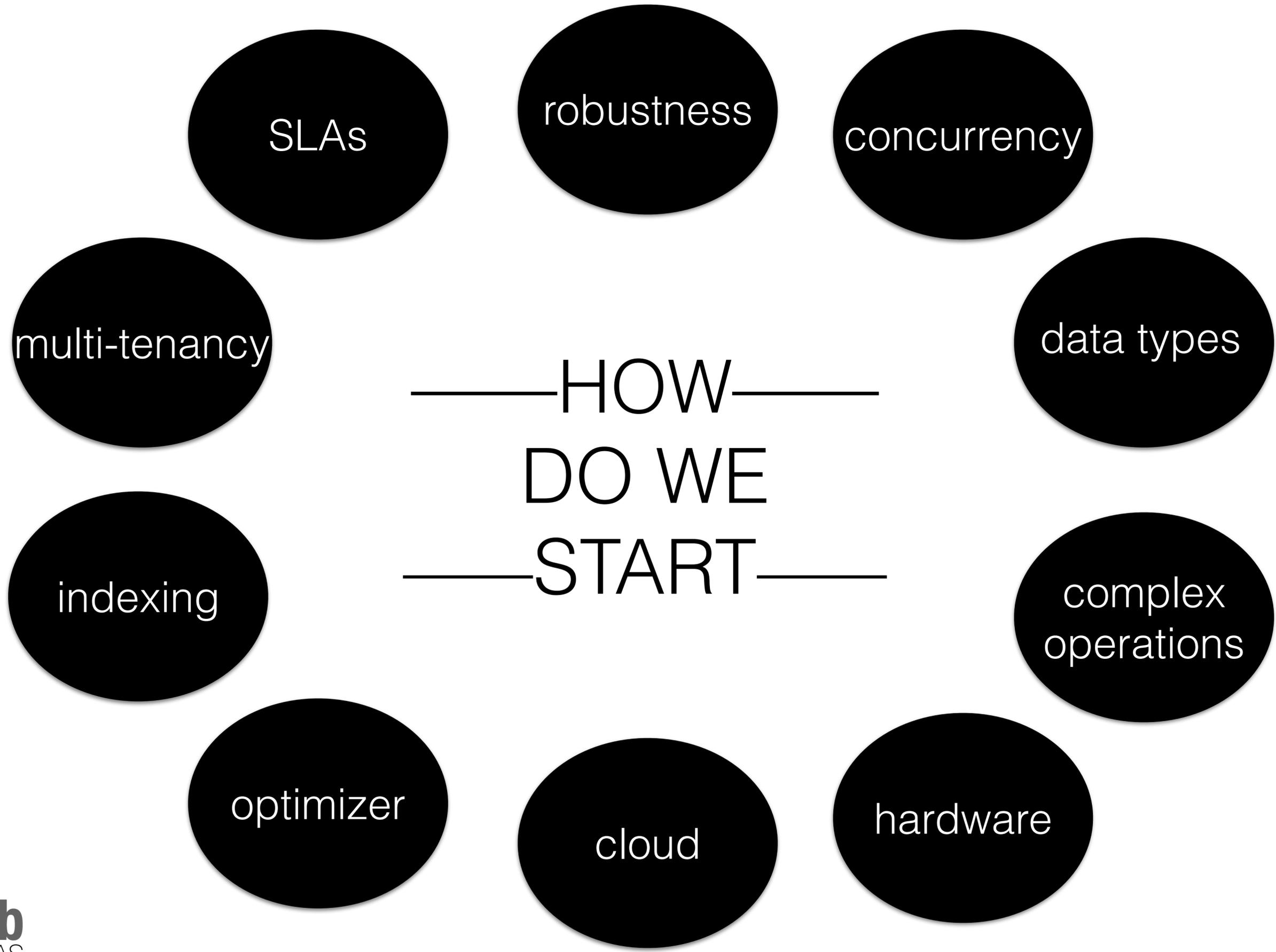


massive design space of system designs



reasoning: understand all the design decisions & their impact

—HOW—
DO WE
—START—





— HOW —
DO WE
— START —

ALGORITHMS

data structure decisions define
the algorithms that access data

INDEX

DATA

ALGORITHMS

unordered

[7,4,2,6,1,3,9,10,5,8]

INDEX

DATA

ALGORITHMS

unordered

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
[7,4,2,6,1,3,9,10,5,8]

INDEX

DATA

ALGORITHMS

unordered


[7, 4, 2, 6, 1, 3, 9, 10, 5, 8]

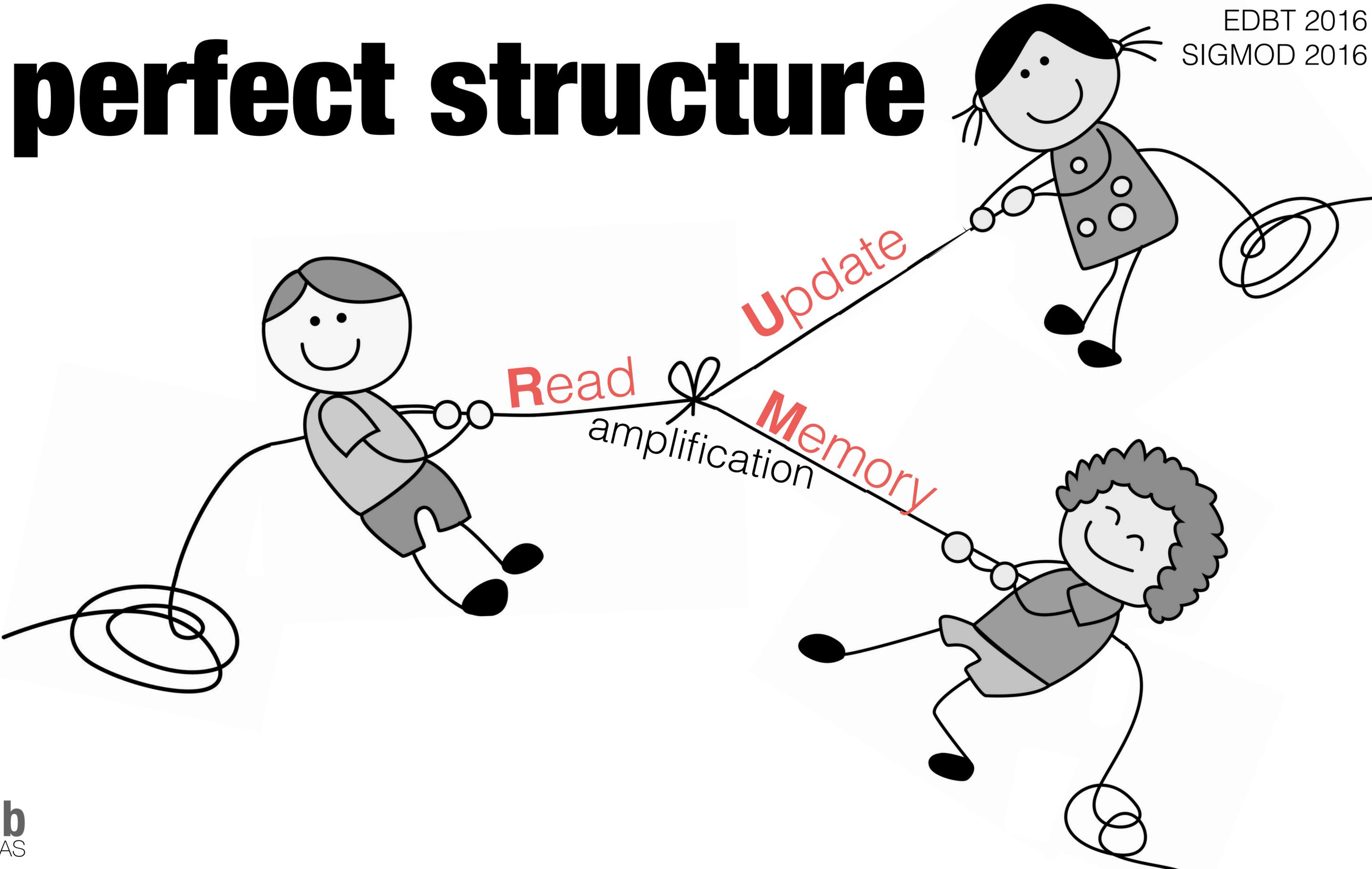
ordered


[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

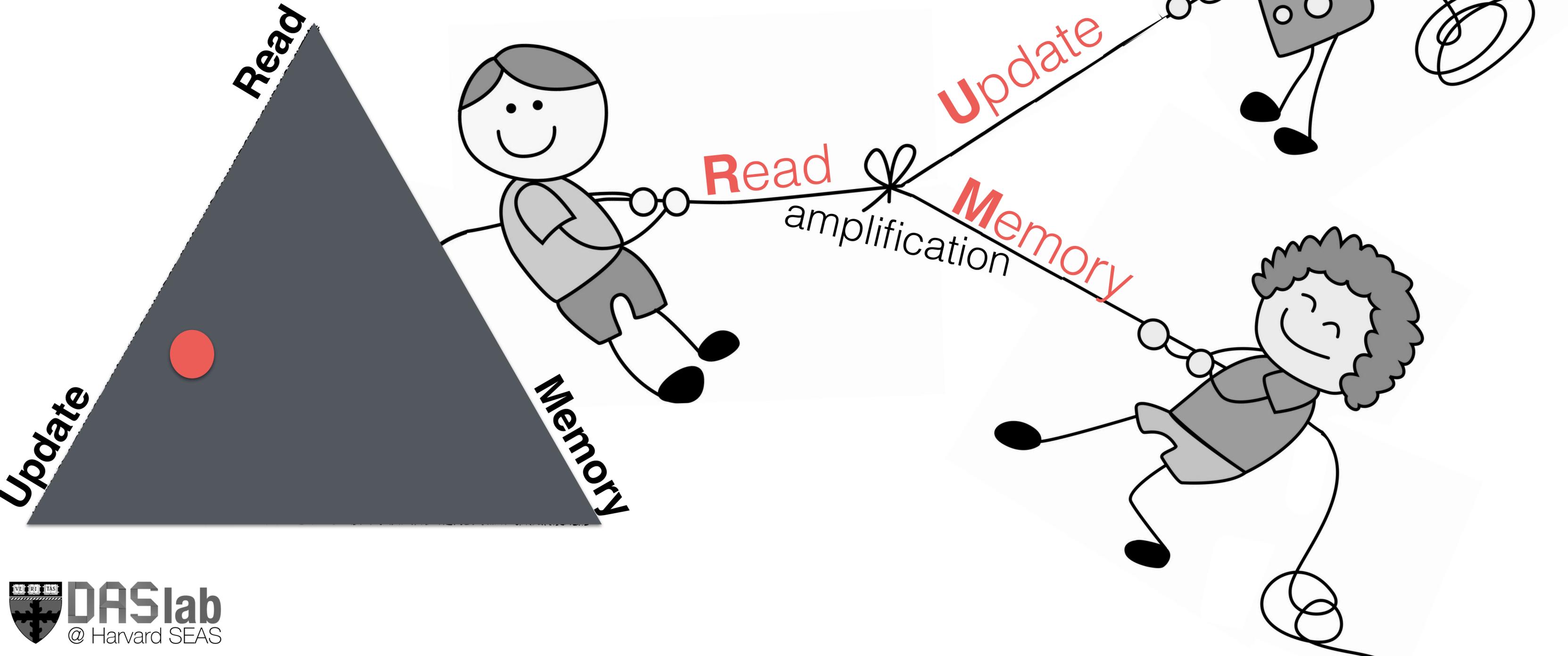
INDEX

DATA

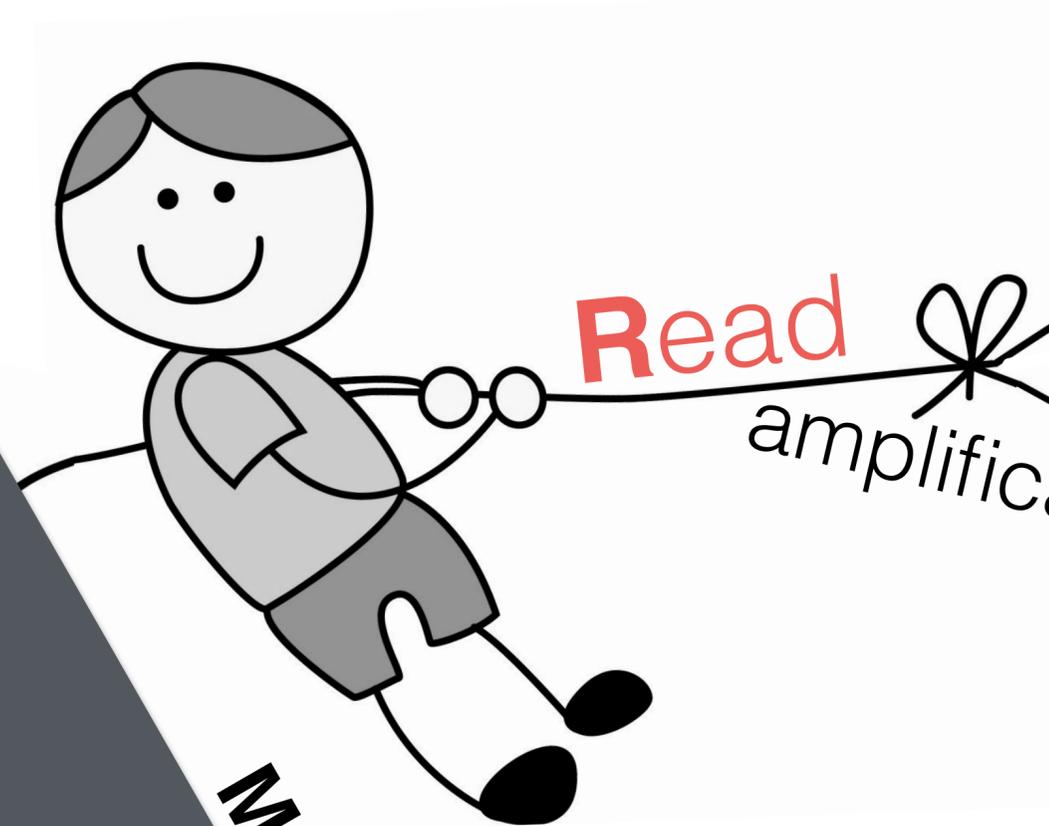
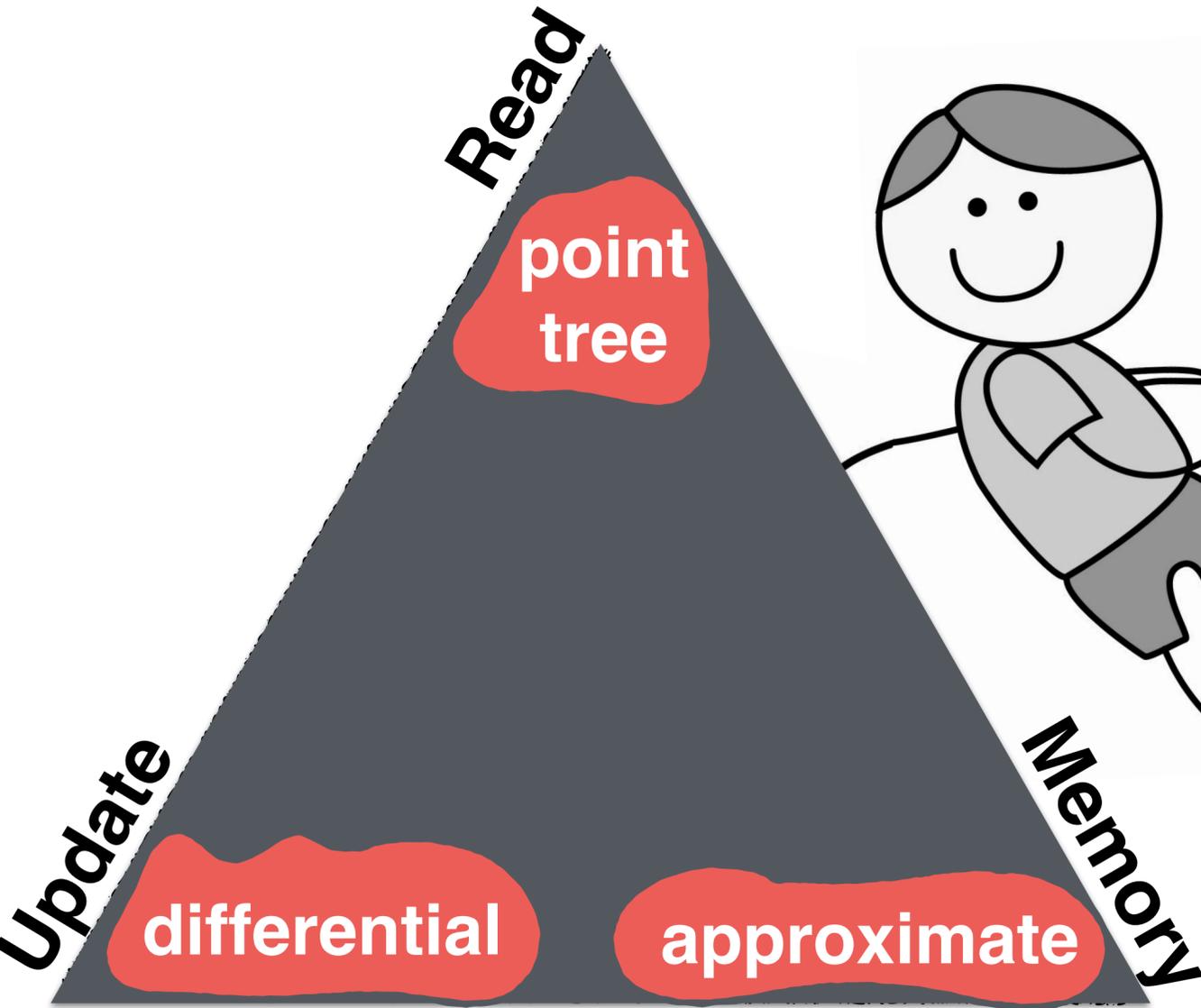
no perfect structure



no perfect structure



no perfect structure



Read

amplification

Update

Memory



read



update

memory

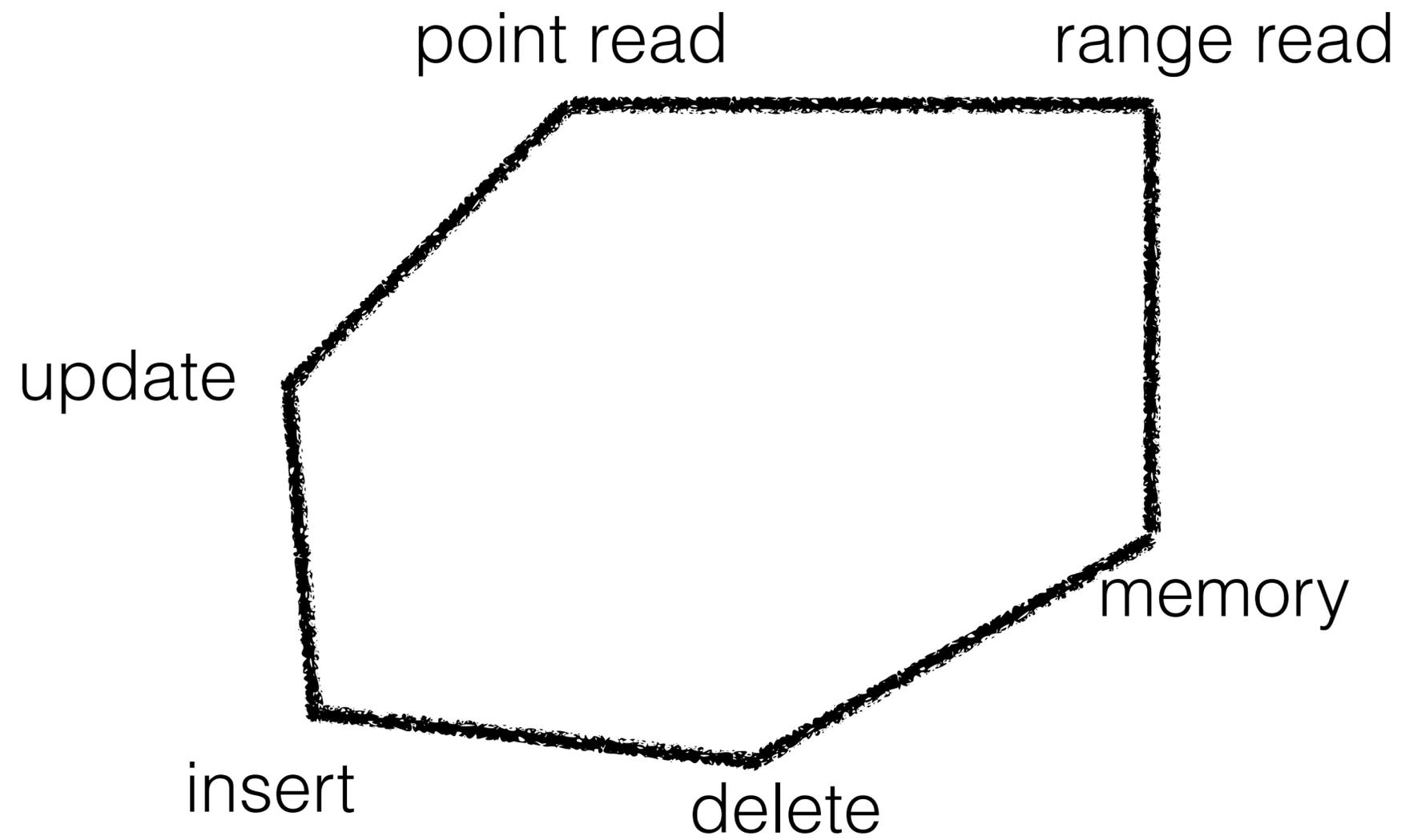
point read

range read

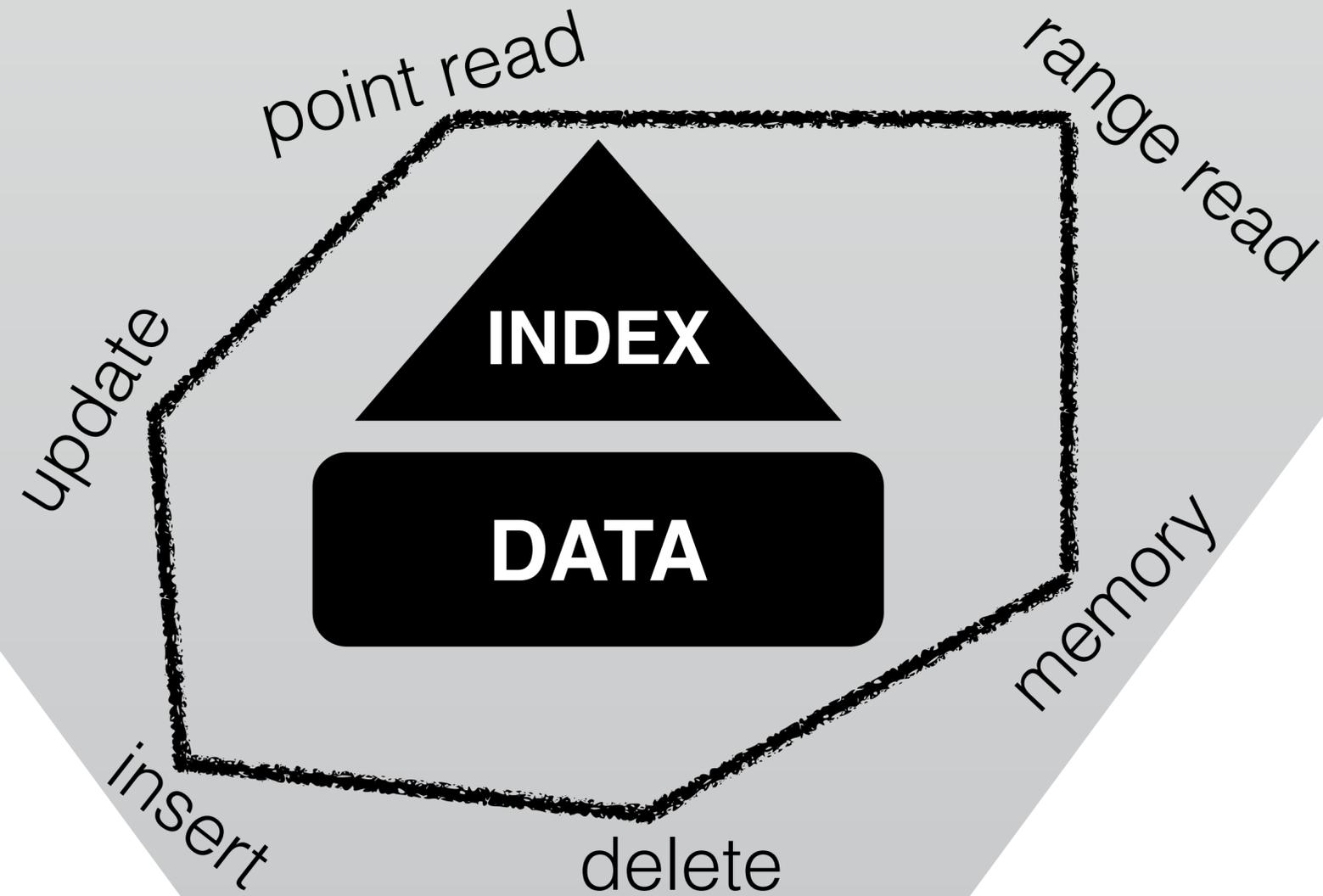


update

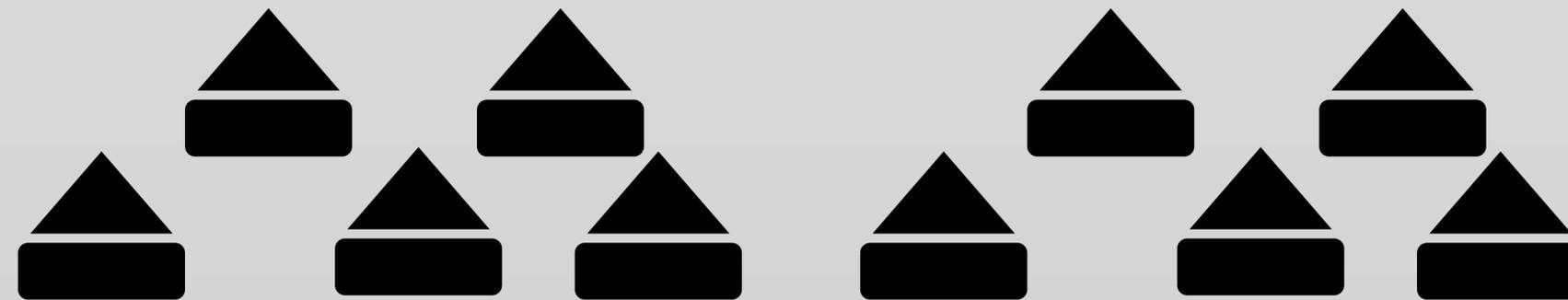
memory



ALGORITHMS



DATA SYSTEMS



ALGORITHMS

INDEX

DATA

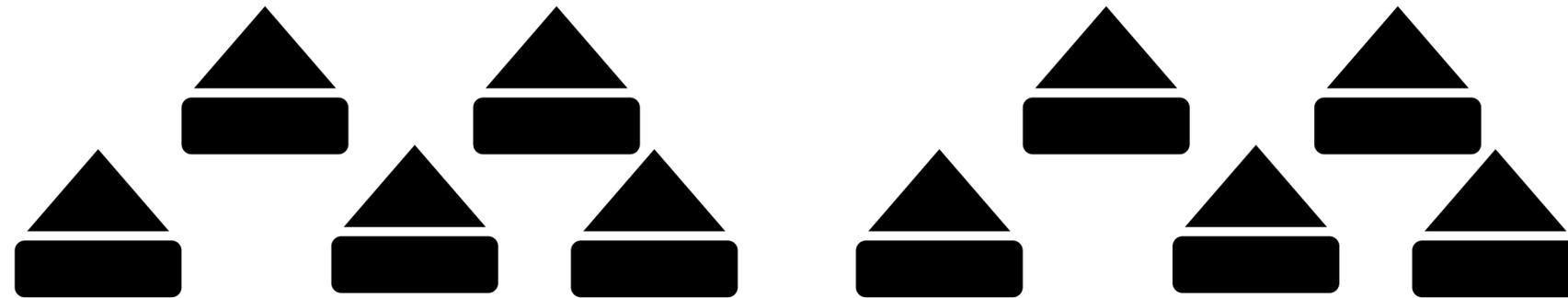
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



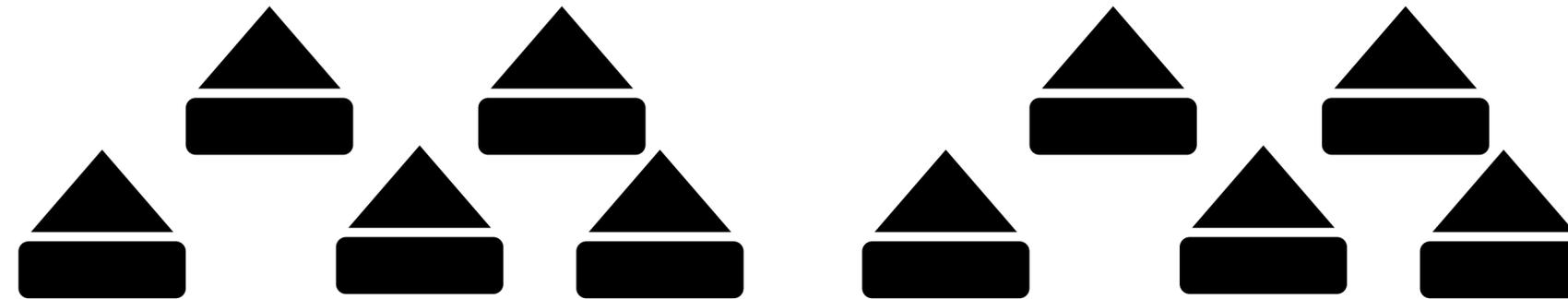
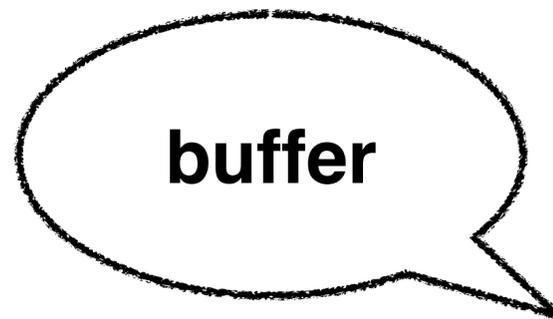
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



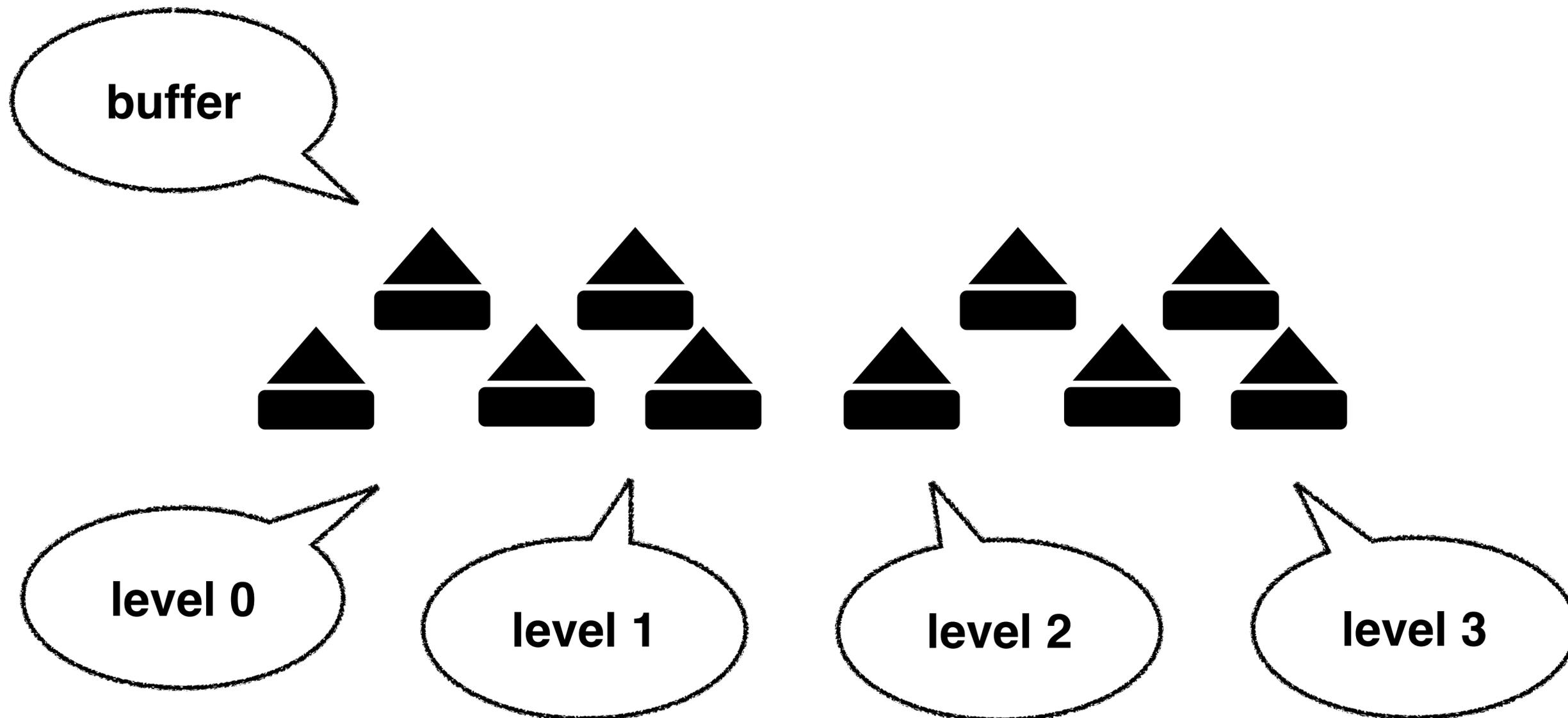
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



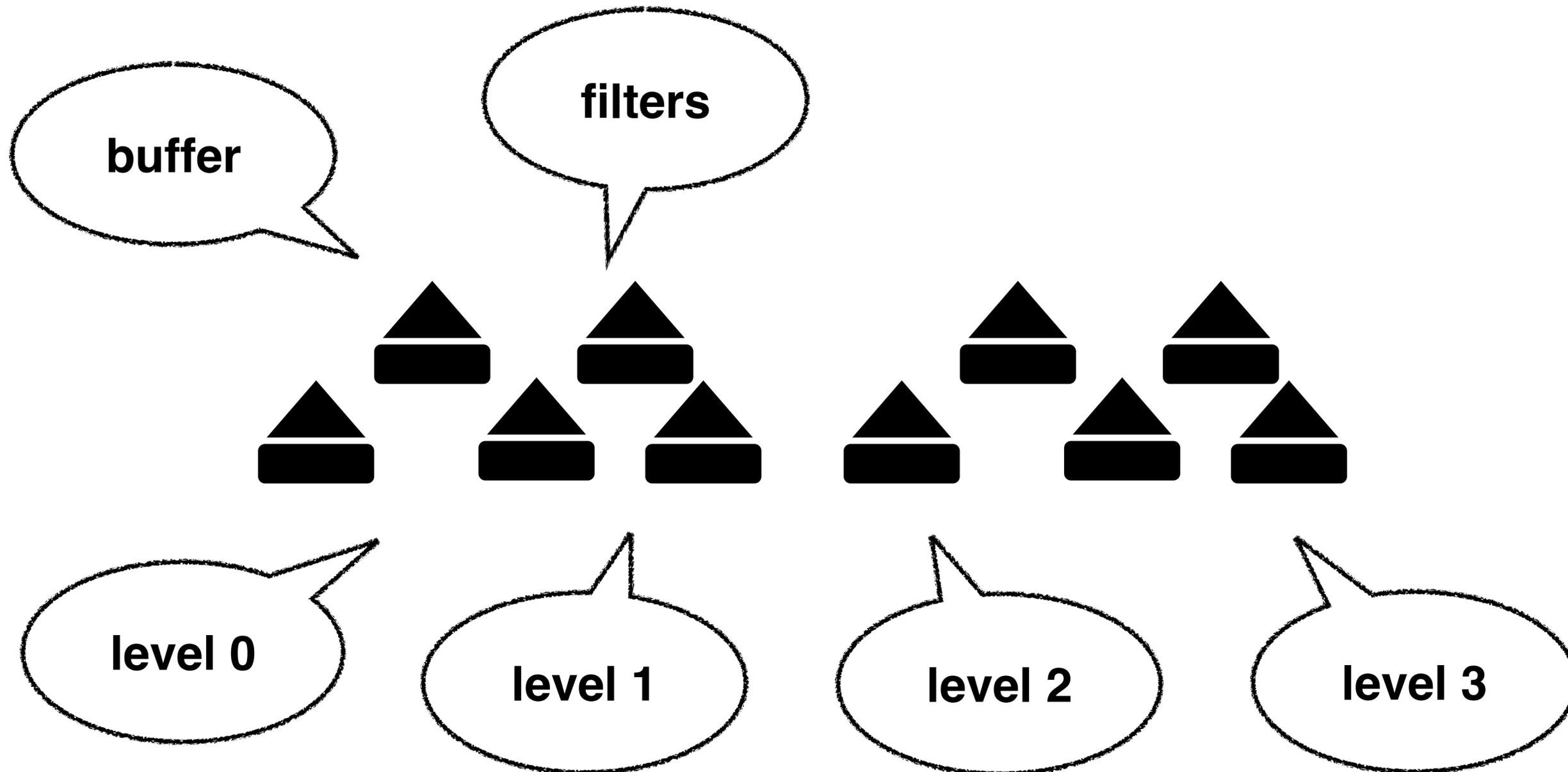
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



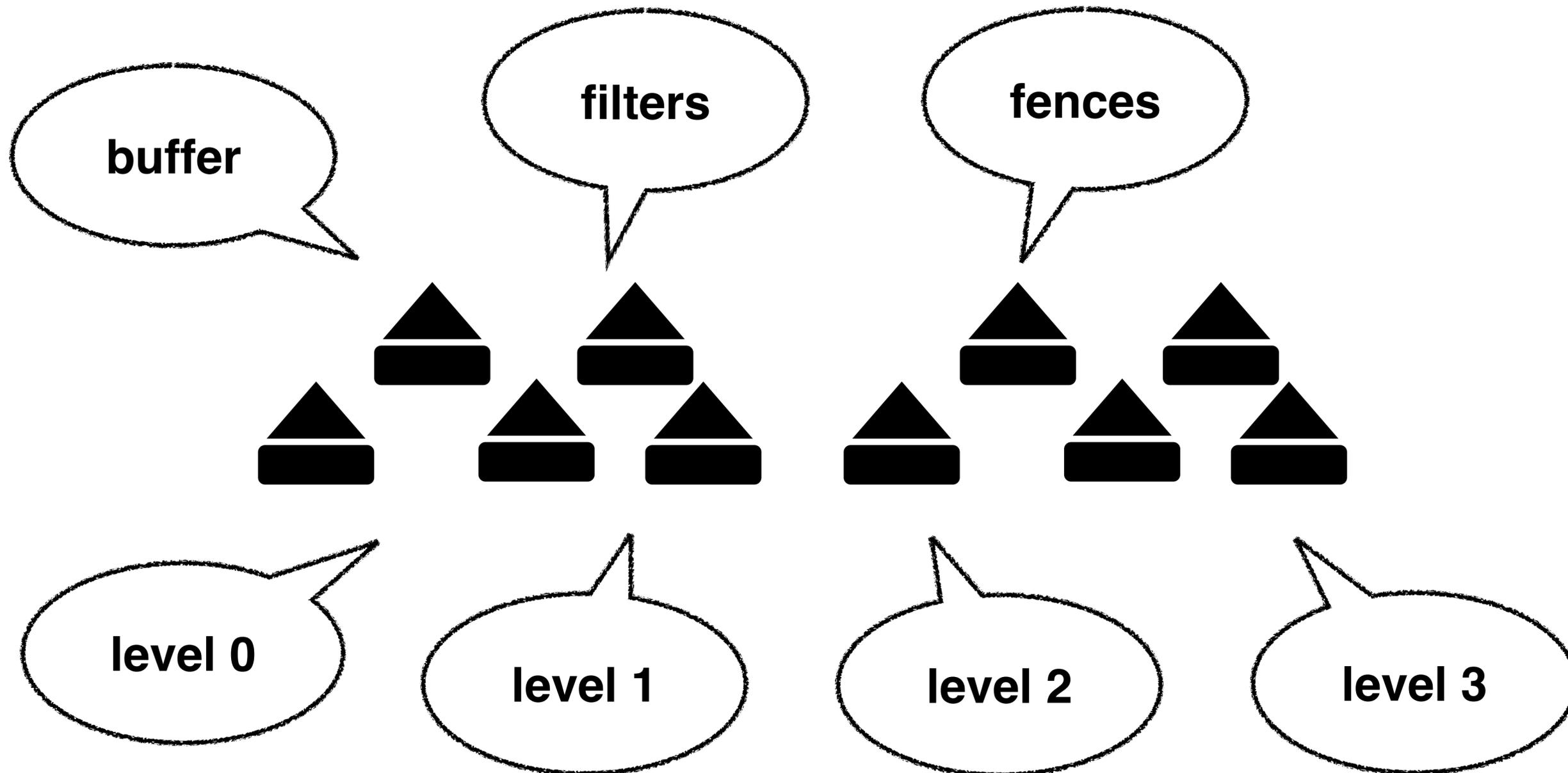
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



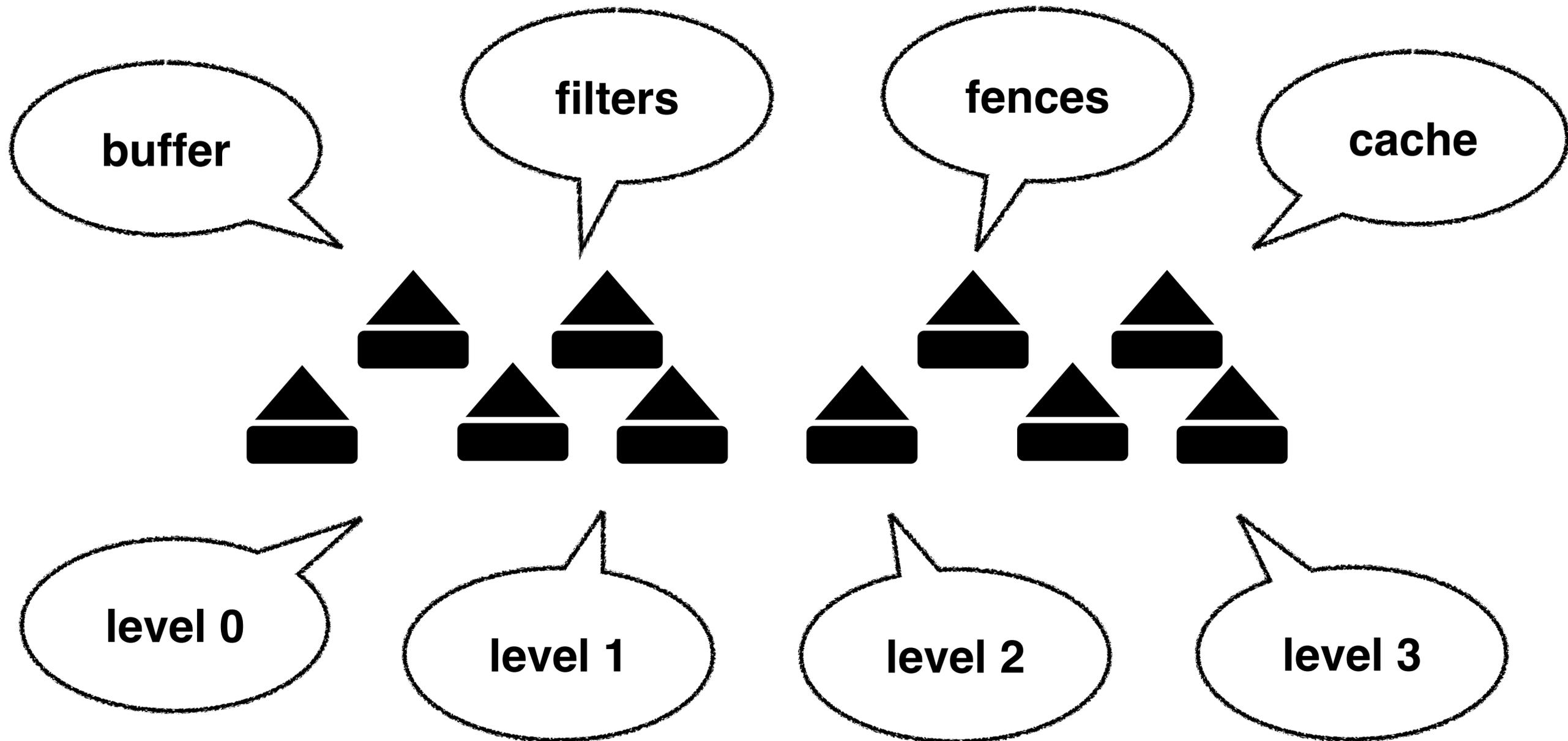
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



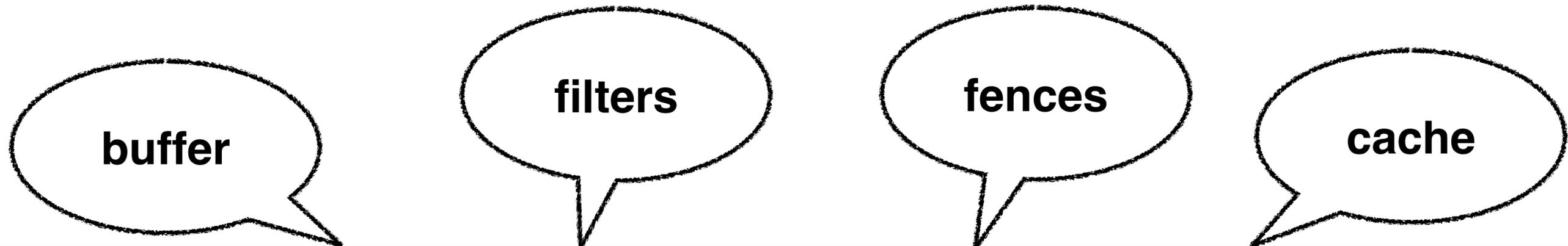
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

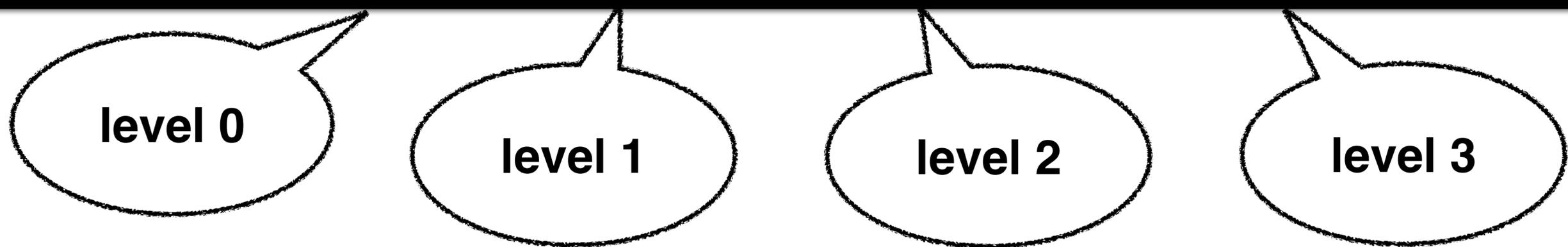
FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



diverse
data structures



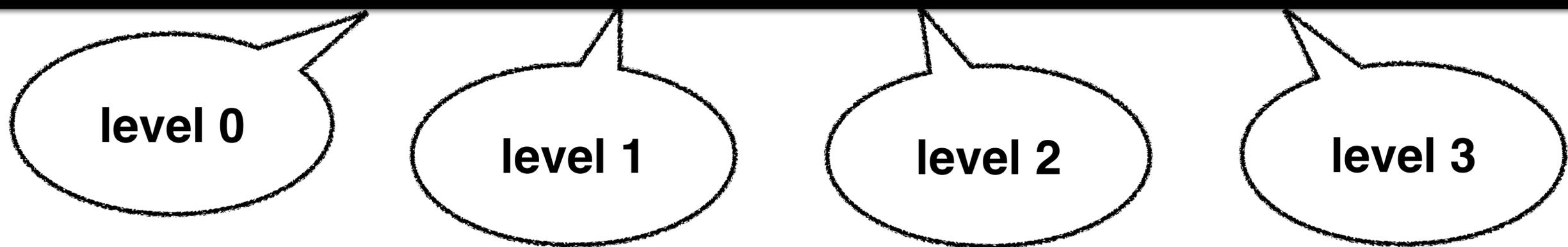
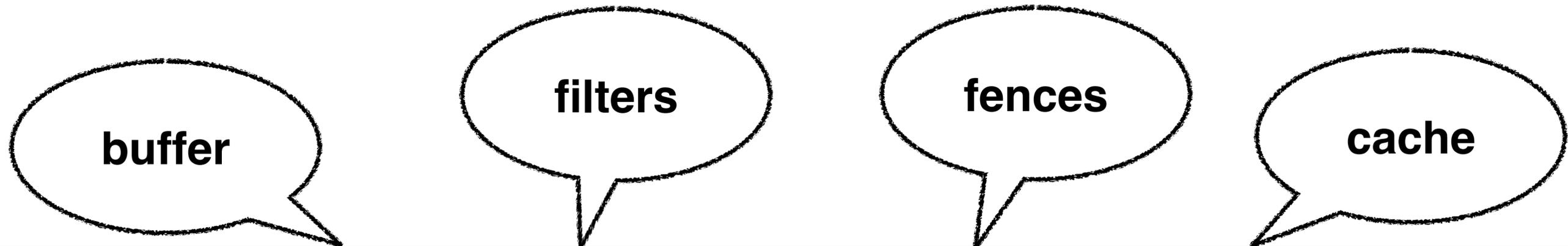
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



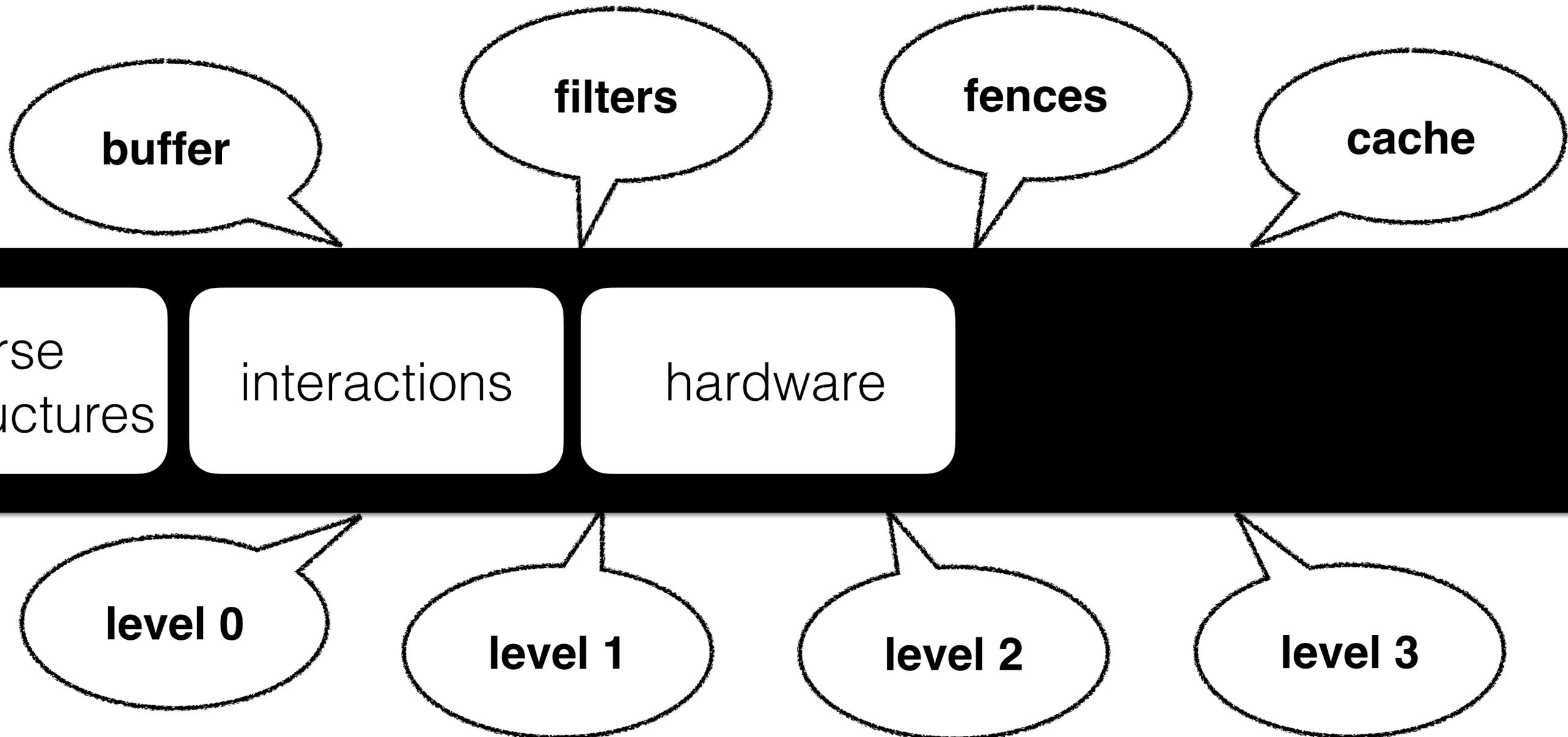
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



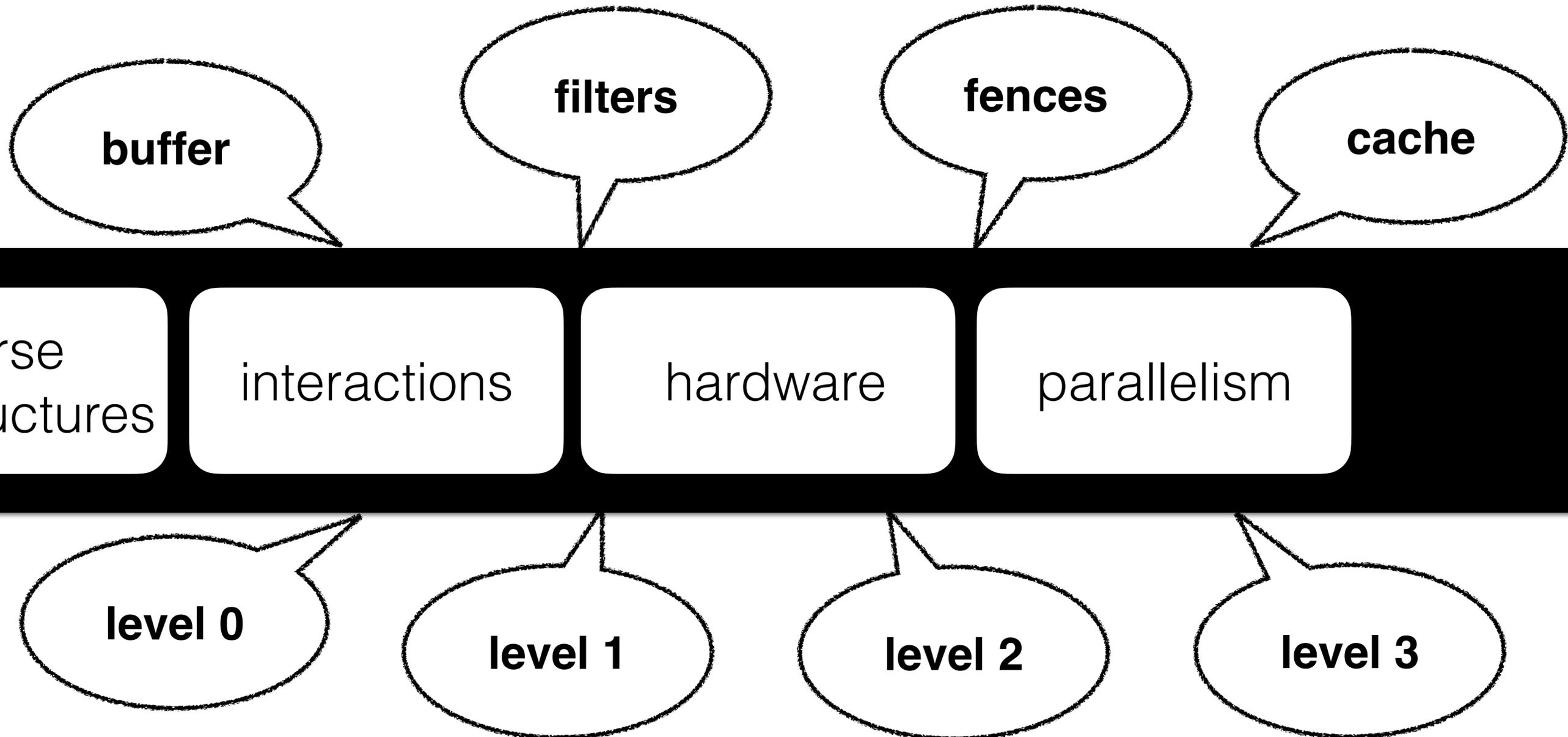
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



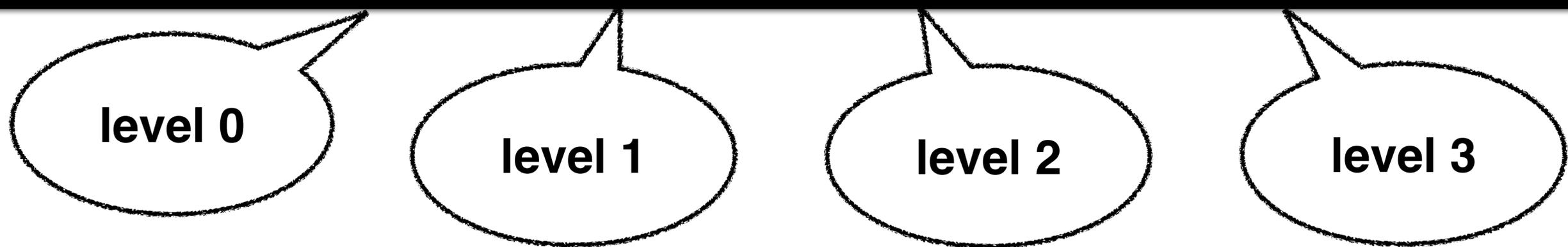
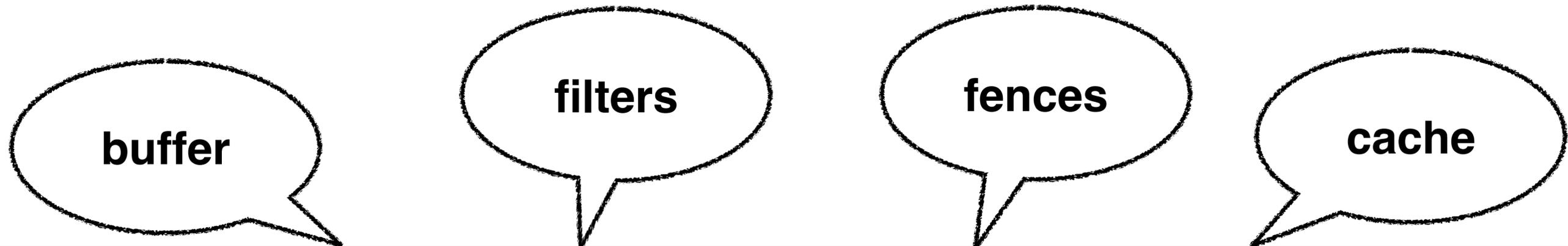
NoSQL systems are the backbone of the BigData and AI era

LSM-tree

FACEBOOK, AMAZON, GOOGLE, TWITTER, LINKEDIN

KV-stores

MACHINE LEARNING, SQL, CRYPTO, SCIENCE



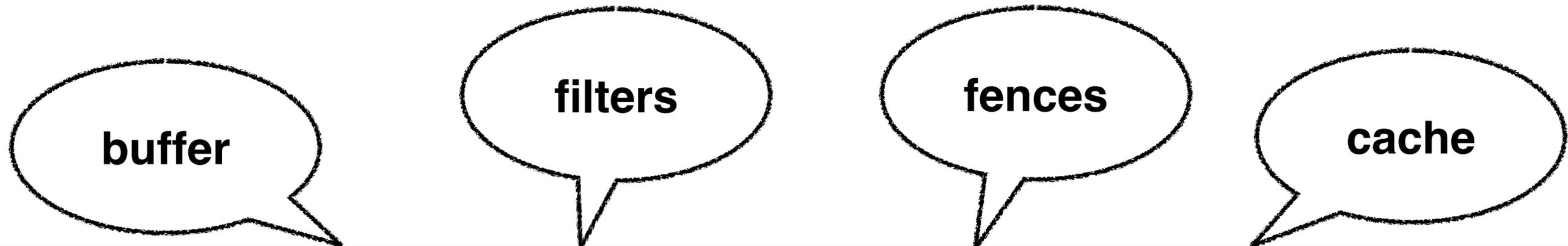
NoSQL systems are the backbone of the BigData and AI era

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 numerous variations of NoSQL KV-stores
LSM-tree variants, B-trees (MongoDB), Hash-index (Microsoft)



diverse
data structures

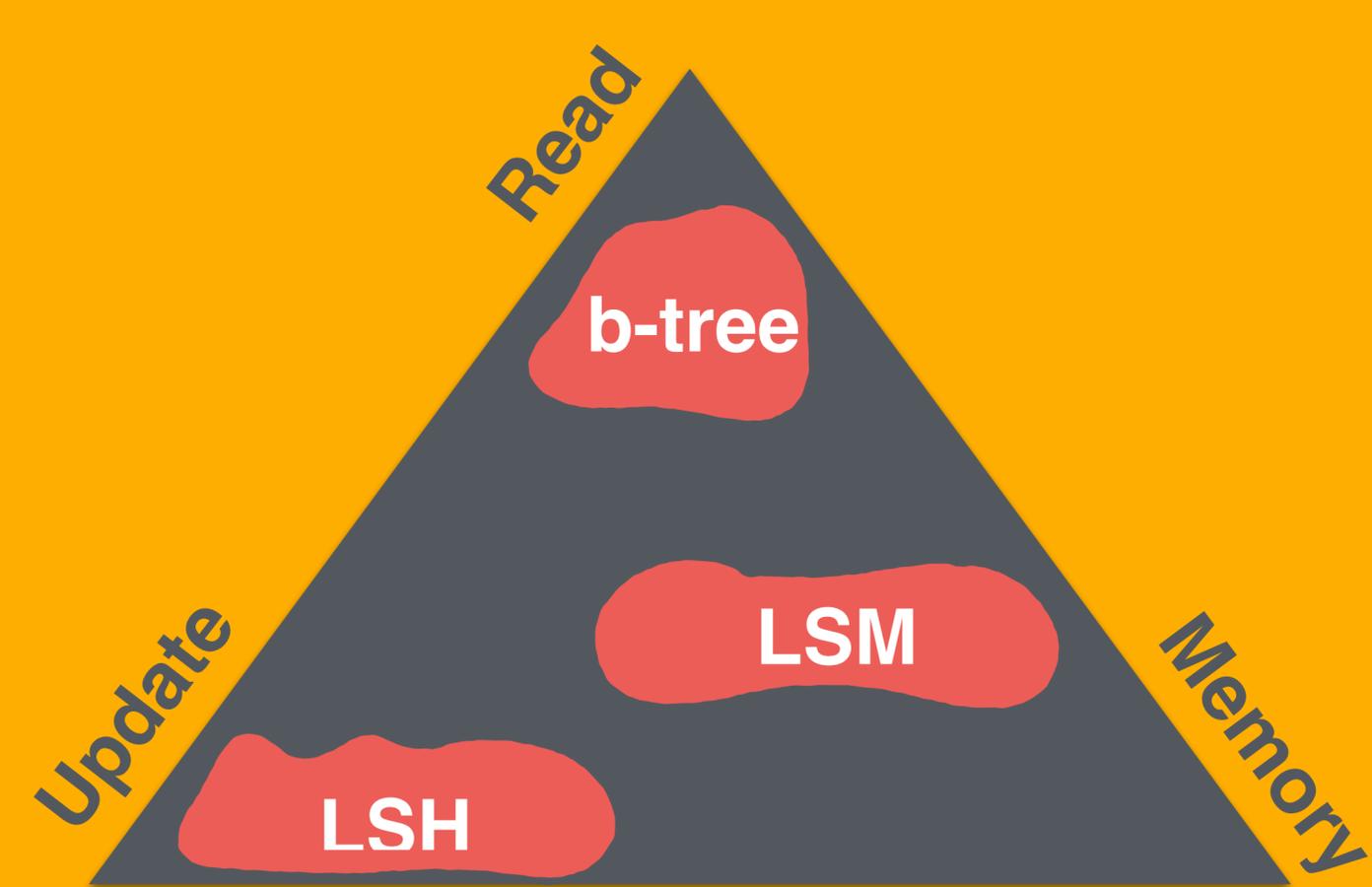
interactions

hardware

parallelism

robustness
cloud cost
SLAs

There exist numerous variations of NoSQL KV-stores
LSM-tree variants, B-trees (MongoDB), Hash-index (Microsoft)



diverse
data structures

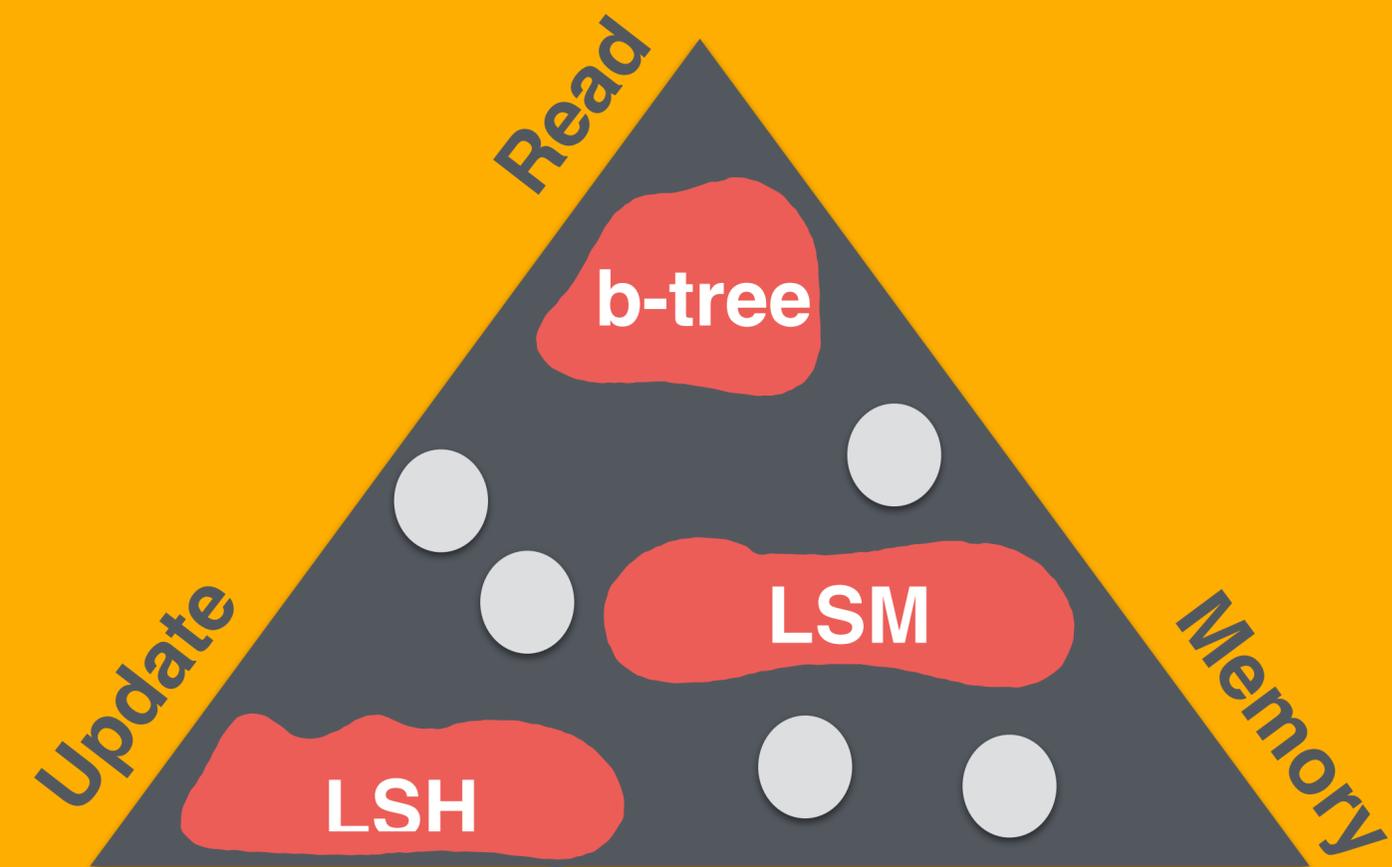
interactions

hardware

parallelism

robustness
cloud cost
SLAs

There exist numerous variations of NoSQL KV-stores
LSM-tree variants, B-trees (MongoDB), Hash-index (Microsoft)



**Constant and increasing efforts
for new system designs as
applications & hardware change**

diverse
data structures

interactions

hardware

parallelism

robustness
cloud cost
SLAs

**There exist numerous variations of NoSQL KV-stores
LSM-tree variants, B-trees (MongoDB), Hash-index (Microsoft)**

diverse
data structures

interactions

hardware

parallelism

robustness
cloud cost
SLAs

Requirements/Goals



diverse
data structures

interactions

hardware

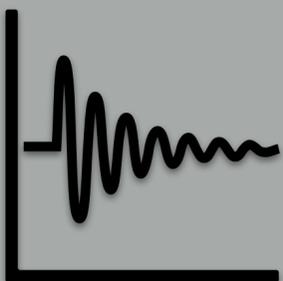
parallelism

robustness
cloud cost
SLAs

Requirements/Goals

Context

data & queries

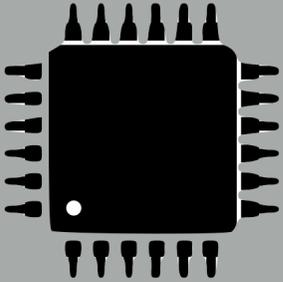


performance



budget

\$\$\$



SLA

diverse
data structures

interactions

hardware

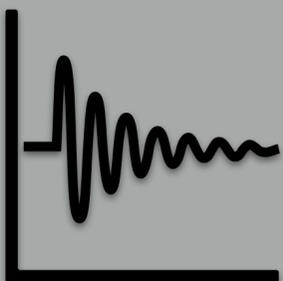
parallelism

robustness
cloud cost
SLAs

Requirements/Goals

Context

data & queries

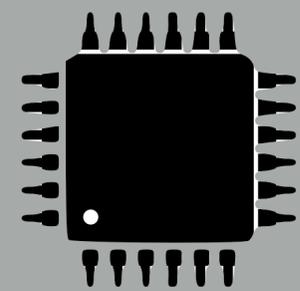


performance

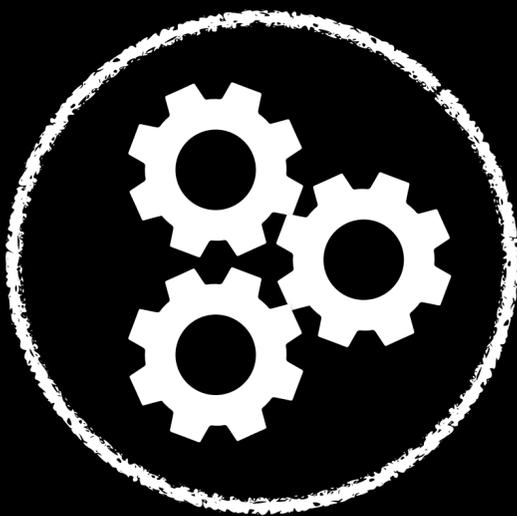
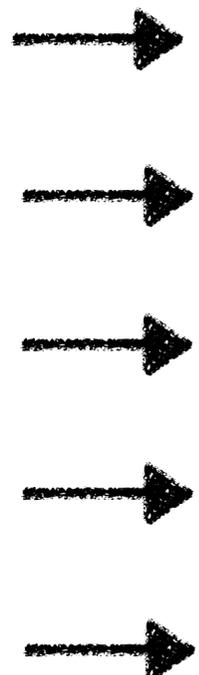


budget

\$\$\$



SLA



diverse
data structures

interactions

hardware

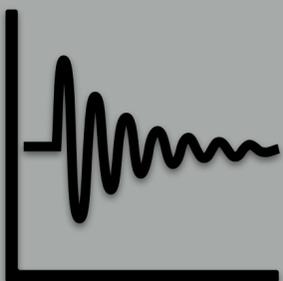
parallelism

robustness
cloud cost
SLAs

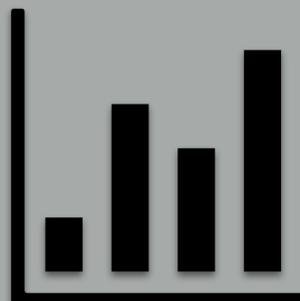
Requirements/Goals

Context

data & queries

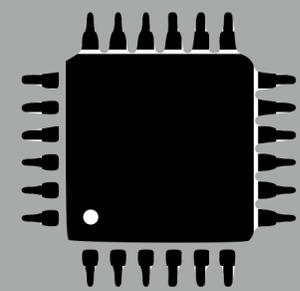


performance

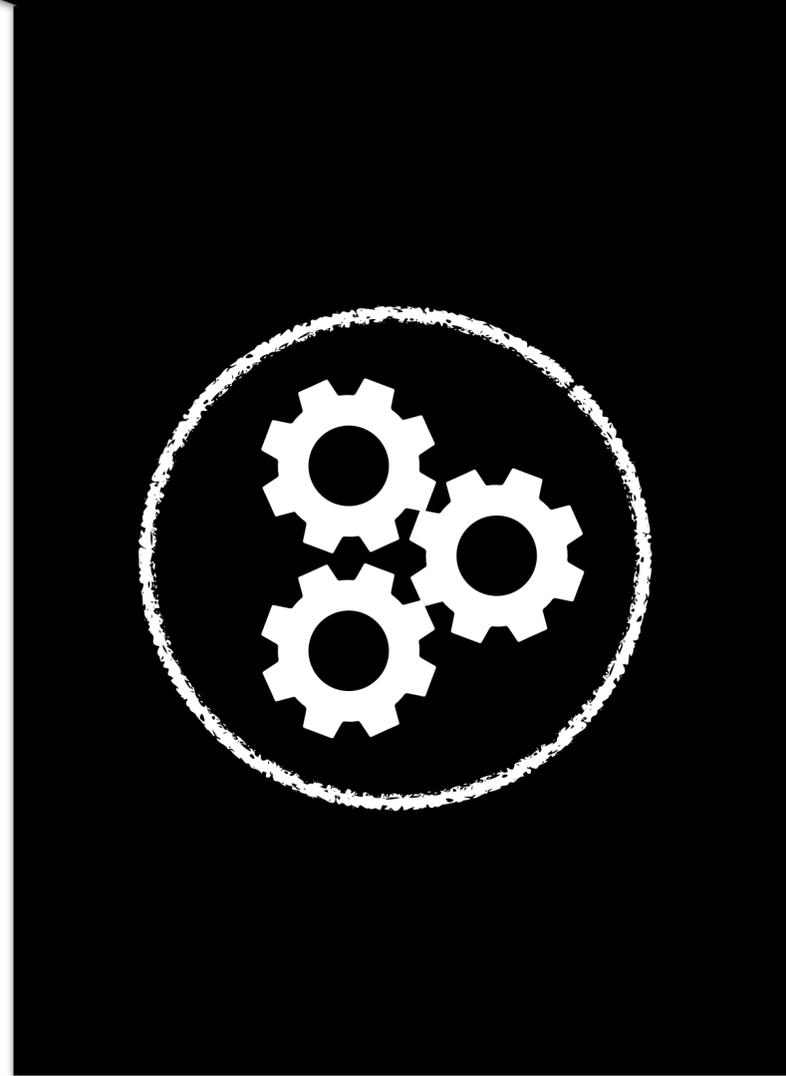
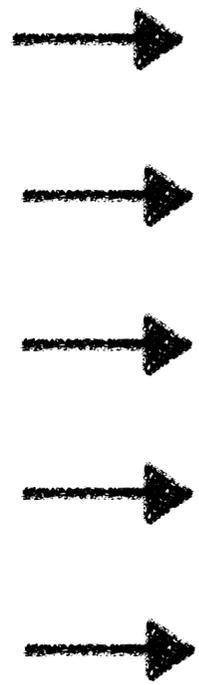


budget

\$\$\$

SLA

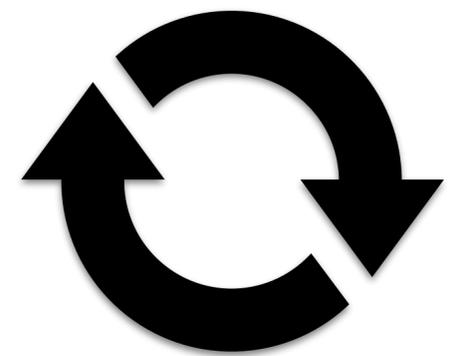
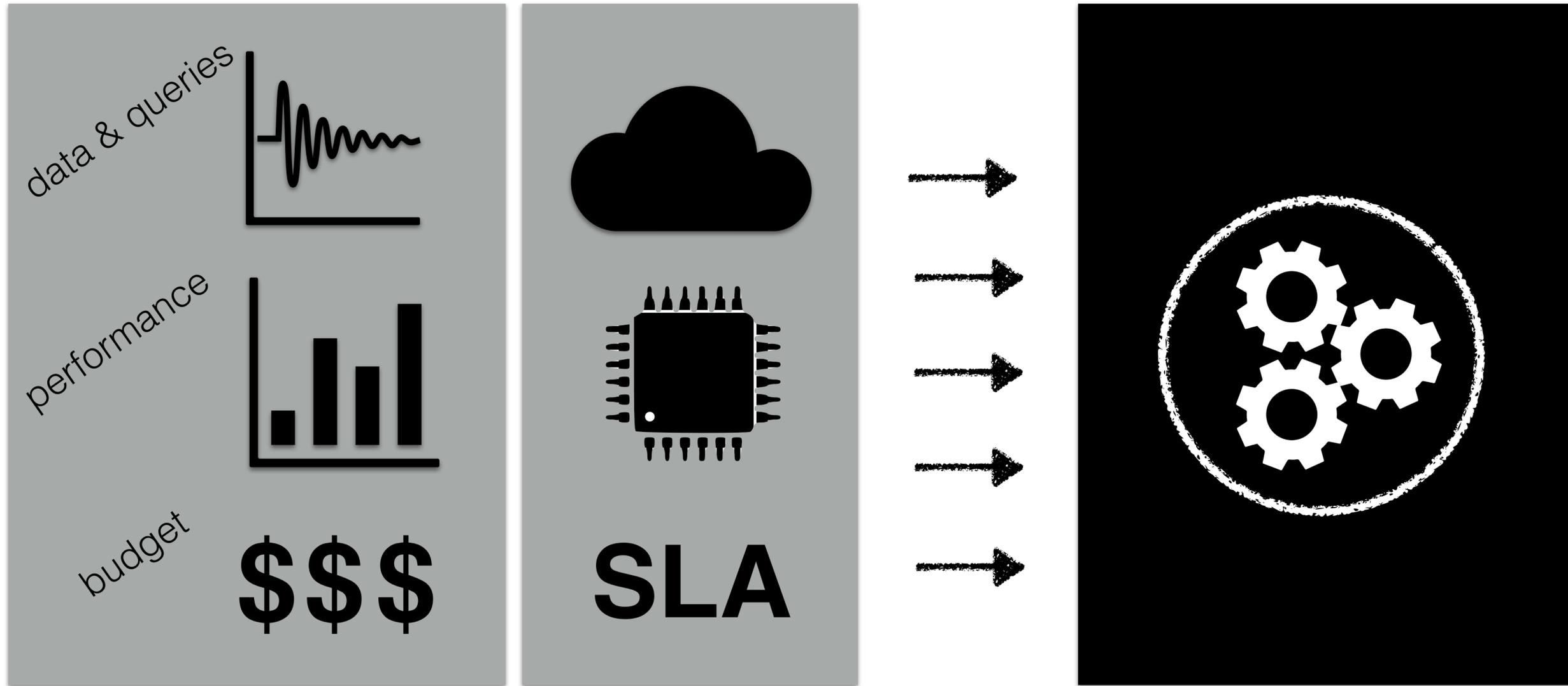


best

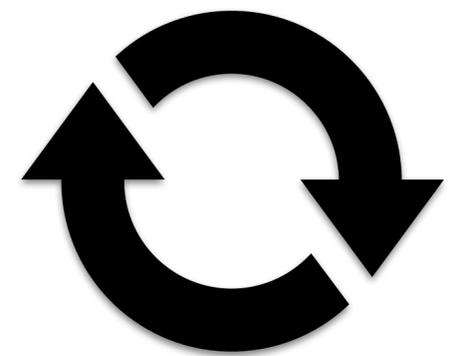
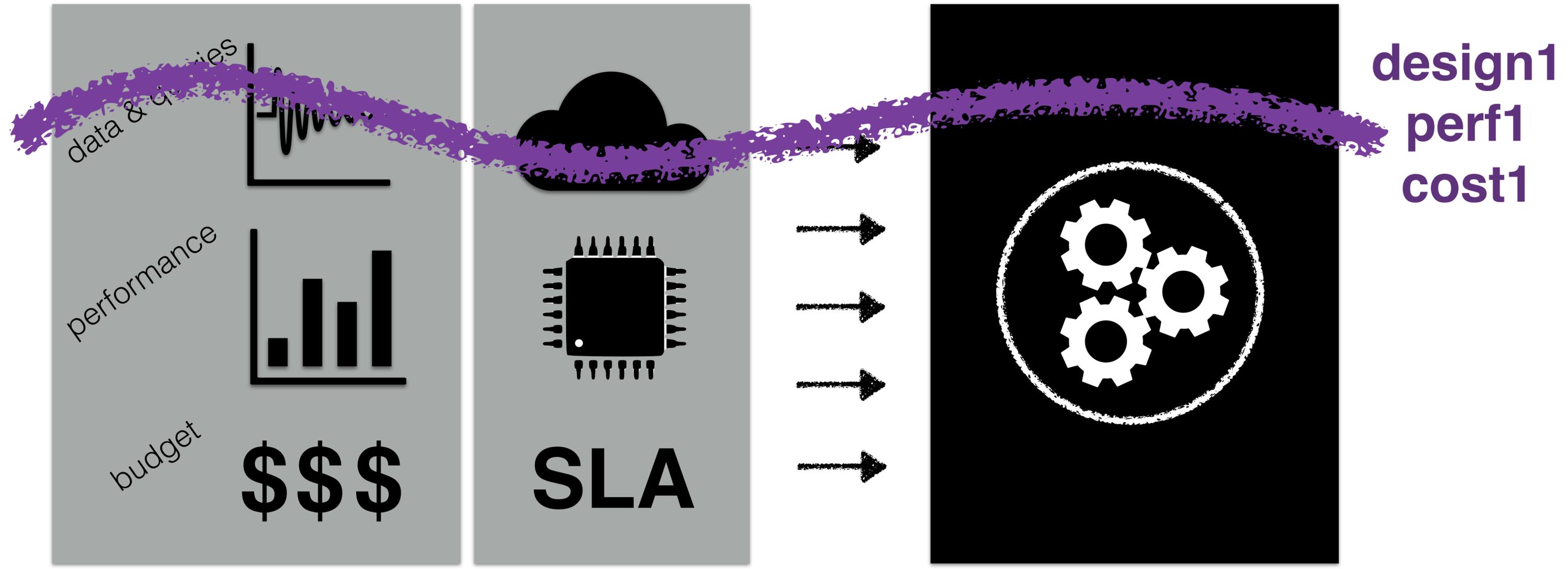


DATA
SYSTEM
design & code

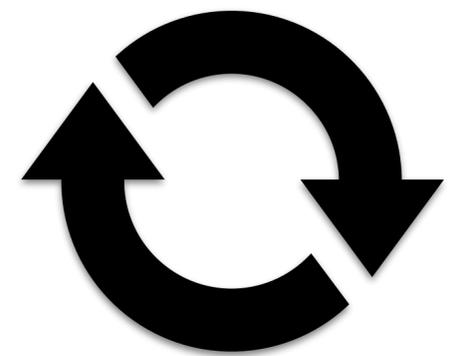
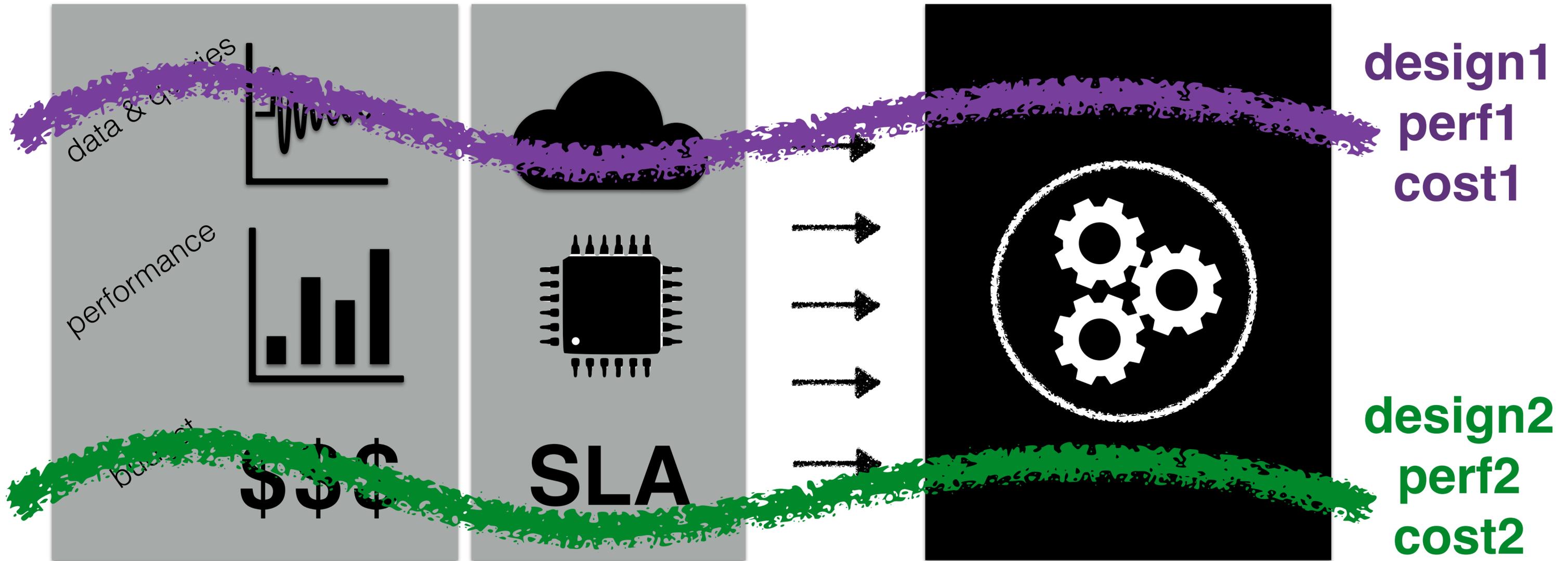




what-if reasoning



what-if reasoning



what-if reasoning

AUTO DESIGN



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
am fear free form of
ultimate I theory*

Nikos Kazantzakis, philosopher



the **grammar** of data systems design

*action is
the most holy
of form
ultimate theory*

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

Nikos Kazantzakis, philosopher



alphabet

Nikos Kazantzakis, philosopher

the **grammar** of data systems design

*action is
the most holy
ultimate form of
theory*

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



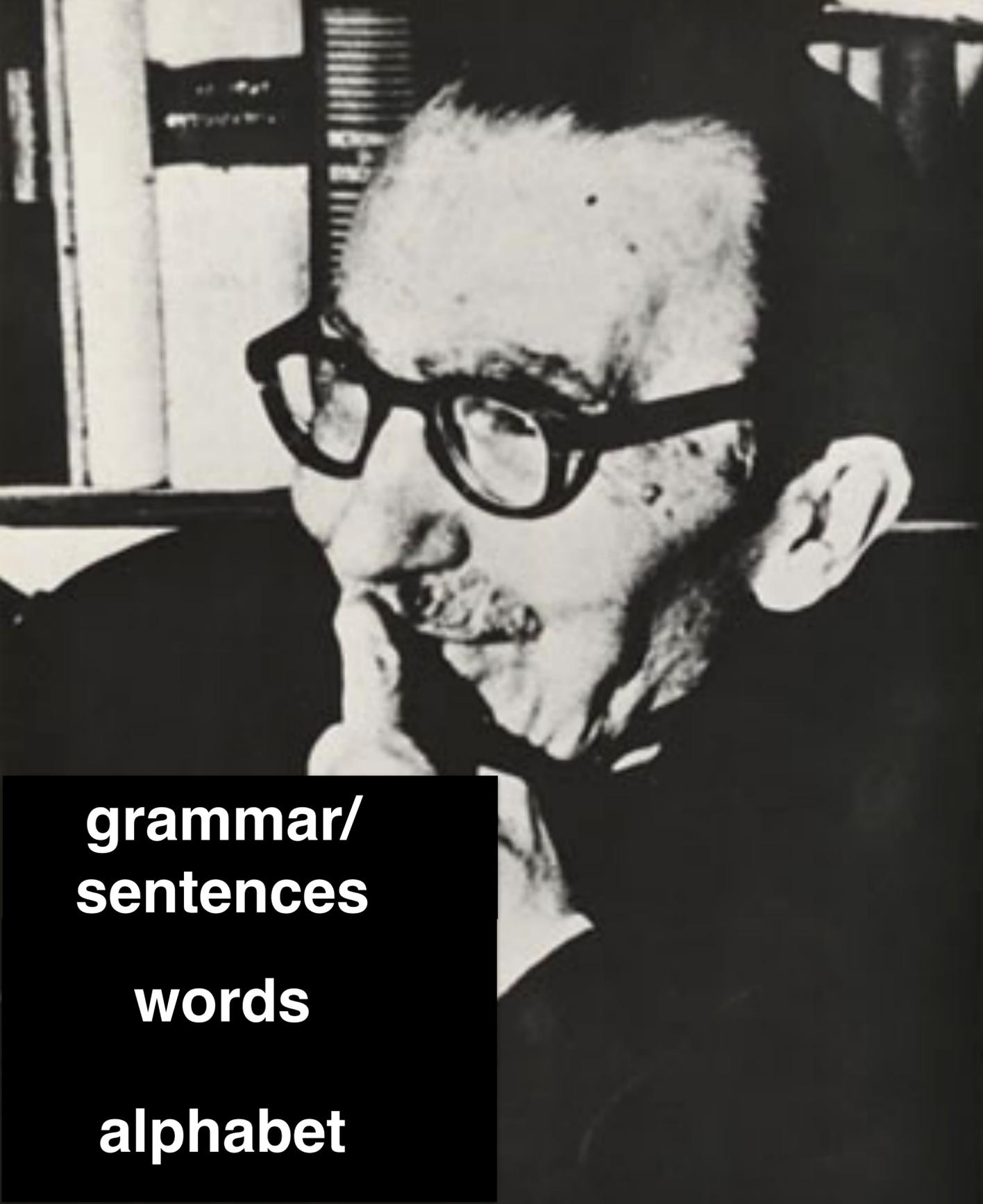
the **grammar** of data systems design

*action is
the most holy
of form
ultimate theory*

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

**words
alphabet**

Nikos Kazantzakis, philosopher



**grammar/
sentences**

words

alphabet

Nikos Kazantzakis, philosopher

the **grammar** of data systems design

action is

the

*most holy
of
form*

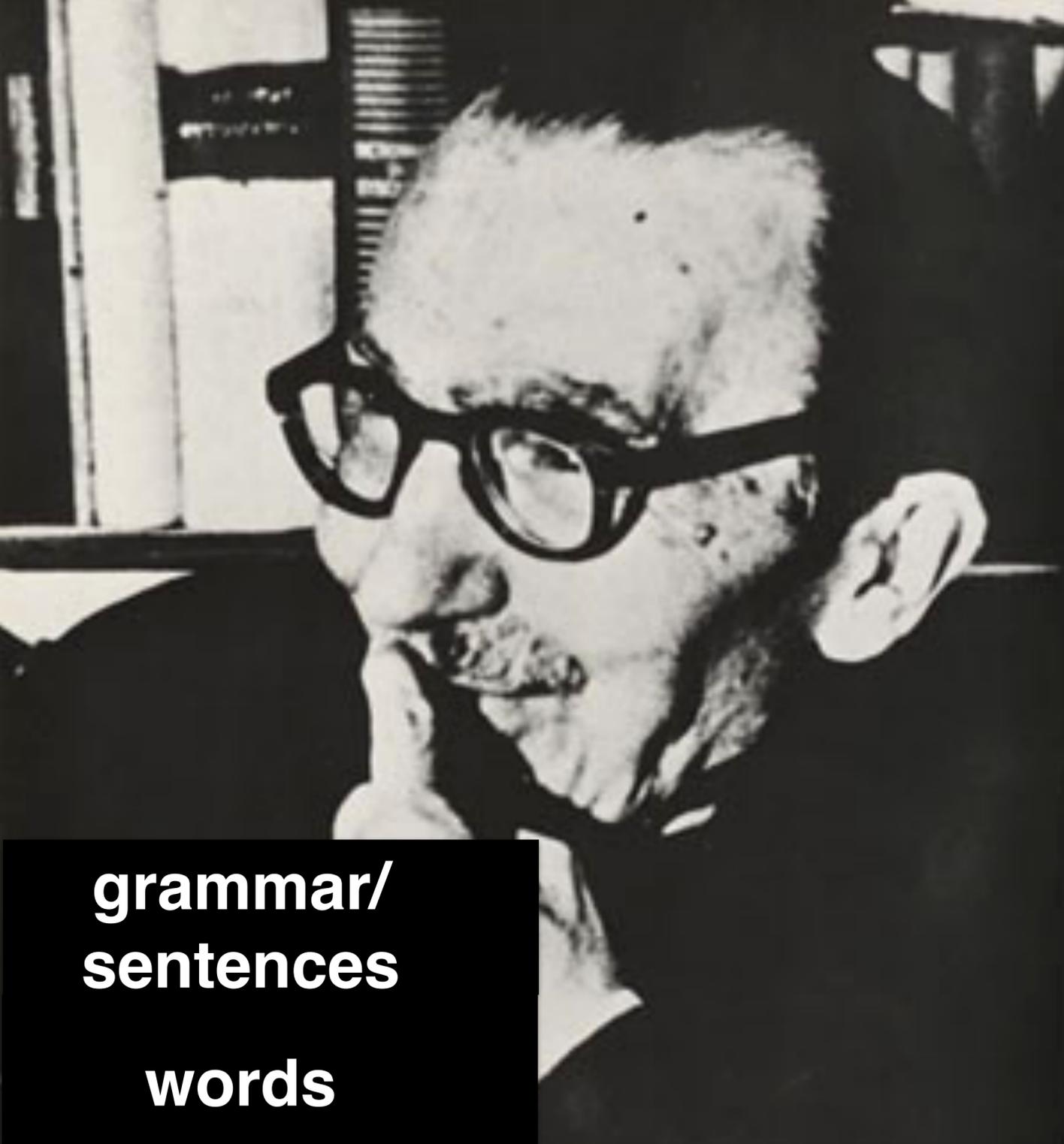
ultimate

theory

I hope for nothing

I fear nothing

I am free



**grammar/
sentences**

words

alphabet

principles

Nikos Kazantzakis, philosopher

the **grammar** of data systems design

action is

the

*most holy
of
form*

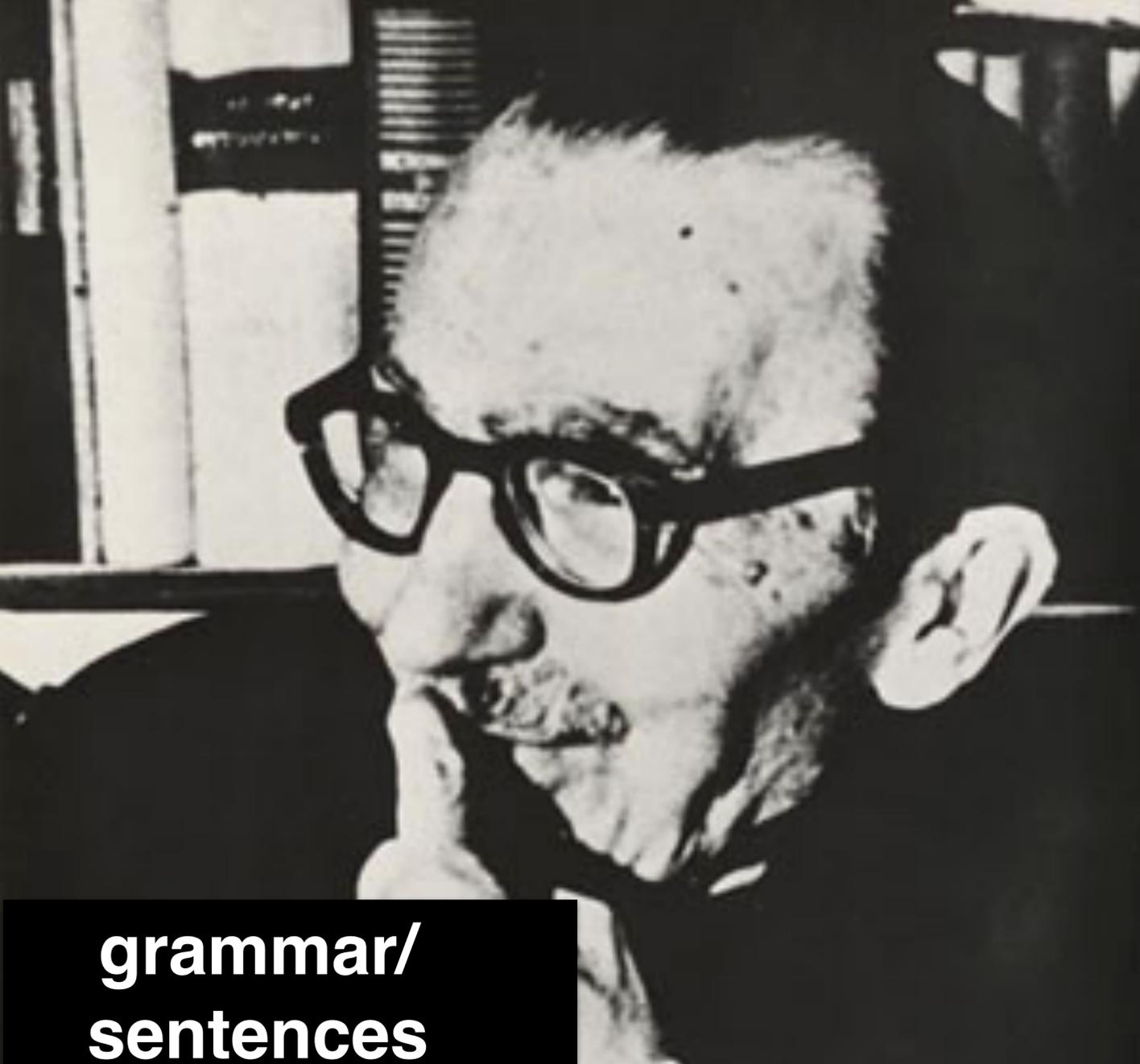
ultimate

theory

I hope for nothing

I fear nothing

I am free



the **grammar** of data systems design

*action is
the most holy
of form
ultimate theory*

**grammar/
sentences**

words

alphabet

data structures

principles

Nikos Kazantzakis, philosopher

I hope for nothing

I fear nothing

I am free



the **grammar** of data systems design

*action is
the most holy
ultimate form of
theory*

**grammar/
sentences**

words

alphabet

interactions

data structures

principles

Nikos Kazantzakis, philosopher

I hope for nothing

I fear nothing

I am free



the **grammar** of data systems design

*action is
the most holy
ultimate form of
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?

**grammar/
sentences**

words

alphabet

interactions

data structures

principles

Nikos Kazantzakis, philosopher

I hope for nothing

I fear nothing

I am free

Trillions of possible data structures

Data Calculator @SIGMOD 2018

Trillions of possible data structures

Data Calculator @SIGMOD 2018

New NoSQL systems: 1000x faster

Cosine @PVLDB 2022 and Limousine @SIGMOD 2024

Trillions of possible data structures

Data Calculator @SIGMOD 2018

New NoSQL systems: 1000x faster

Cosine @PVLDB 2022 and Limousine @SIGMOD 2024

Synthesized statistics, 10x faster ML

Data Canopy @SIGMOD 2017

Trillions of possible data structures

Data Calculator @SIGMOD 2018

New NoSQL systems: 1000x faster

Cosine @PVLDB 2022 and Limousine @SIGMOD 2024

Synthesized statistics, 10x faster ML

Data Canopy @SIGMOD 2017

10x faster Neural Networks

MotherNets @MLSys 2020, and M2 @MLSys 2023

Trillions of possible data structures

Data Calculator @SIGMOD 2018

New NoSQL systems: 1000x faster

Cosine @PVLDB 2022 and Limousine @SIGMOD 2024

Synthesized statistics, 10x faster ML

Data Canopy @SIGMOD 2017

10x faster Neural Networks

MotherNets @MLSys 2020, and M2 @MLSys 2023

10x faster Image AI

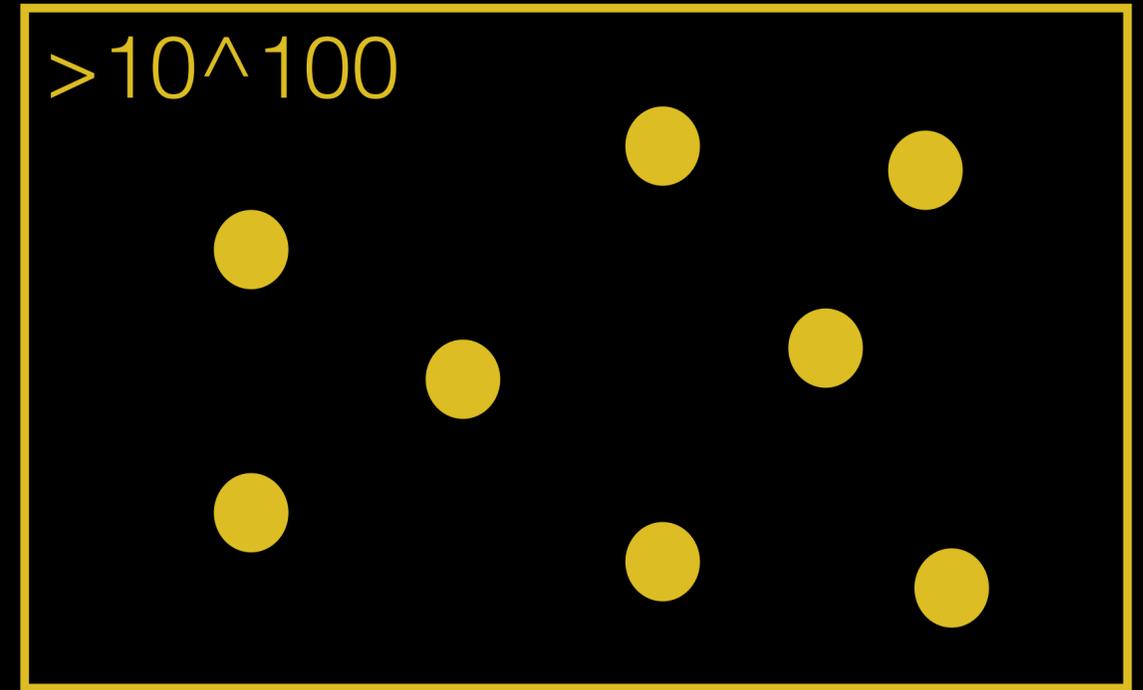
Image Calculator, SIGMOD 2024

1. DESIGN SPACE

data layout of data structures

algorithm design

systems: interactions of components

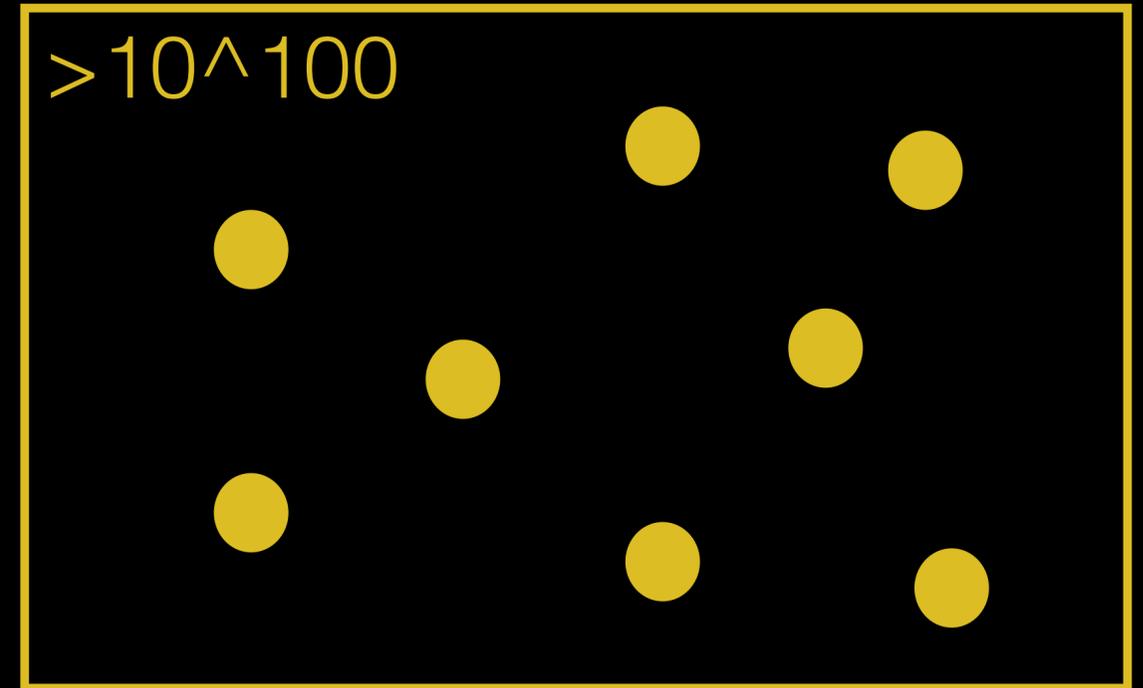


1. DESIGN SPACE

data layout of data structures

algorithm design

systems: interactions of components



2. NAVIGATE SEARCH SPACE

cost synthesis: computation and data movement

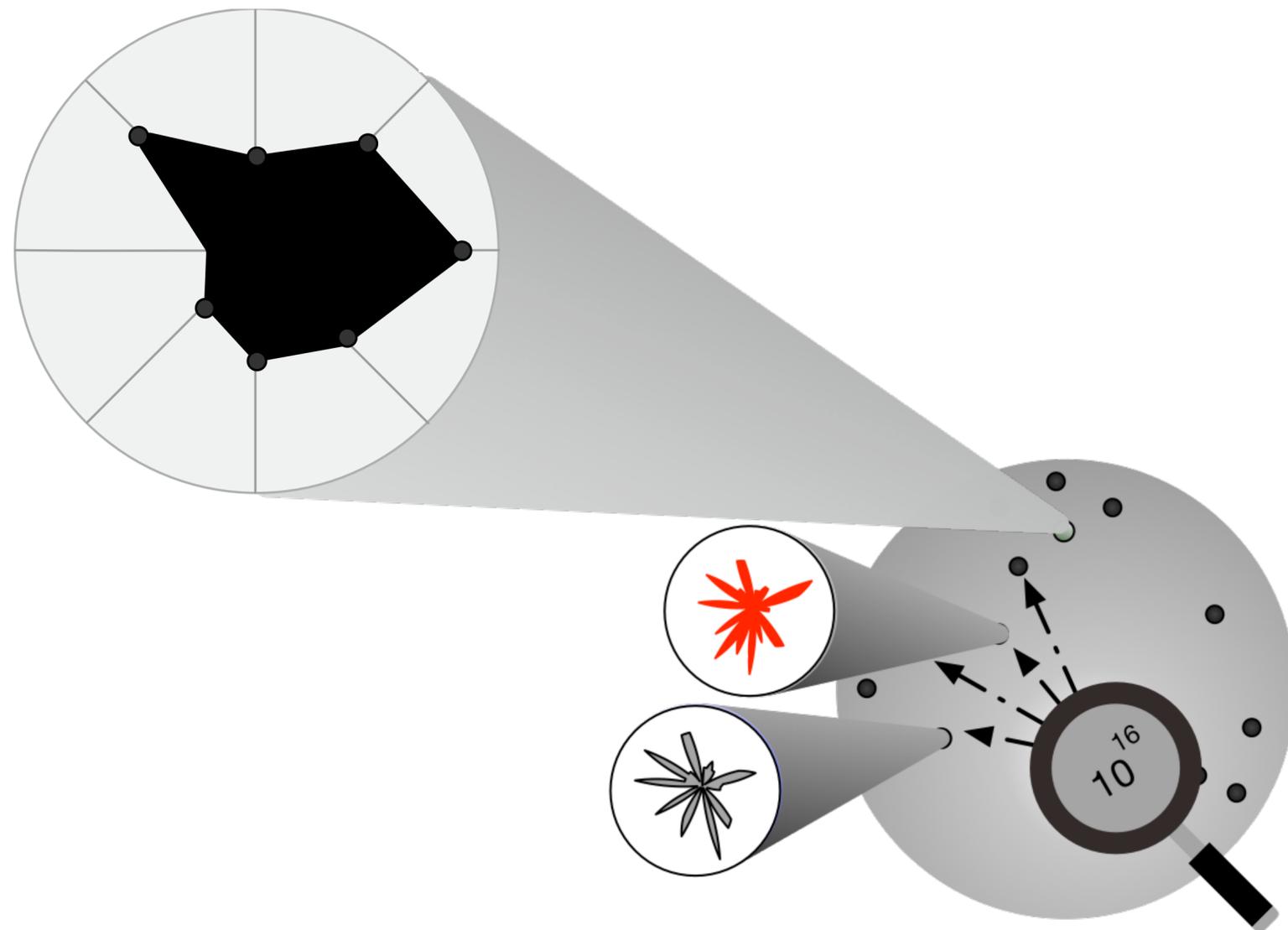
learned cost models in memory/parallelism

design continuums to shrink space

Categories

Category	ID	Description	Options	Count	Partitioning			Linking			Other
					unrestr	fixed(25)	balanced	balanced	balanced	balanced	
Partitioning		or balanced (i.e., all sub-blocks have the same size), unrestricted or functional.	unrestricted function(func) <i>(up to 10 different fixed capacity values are considered)</i>	13	unrestr	fixed(25)		balance	balance	balance	
		<u>Rules:</u> 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	
		<u>Rules:</u> 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 <i>(up to 3 different values for layer-grouping are considered)</i>	5	scatter	scatter		scatter	BFS	BFS-LL	
		<u>Rules:</u> 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	
		<u>Rules:</u> requires fanout/radix != terminal.									
	19	Sub-block consolidation. Single children are merged with their parents.	boolean	2	false	false		false	false	false	
	<u>Rules:</u> 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	
		<u>Rules:</u> 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					
		<u>Rules:</u> 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)	
		<u>Rules:</u> requires fanout/radix != terminal.			no	no					

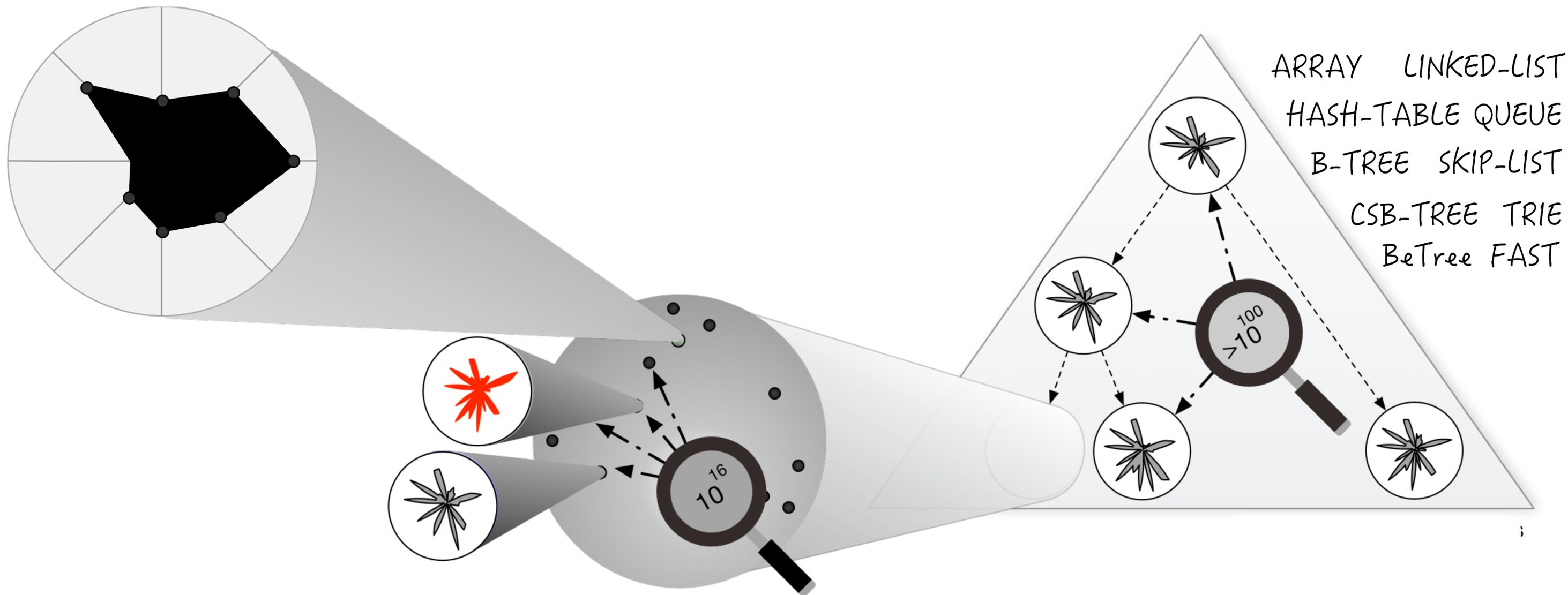
DESIGN PRINCIPLES POSSIBLE NODE DESIGNS



DESIGN PRINCIPLES

POSSIBLE NODE DESIGNS

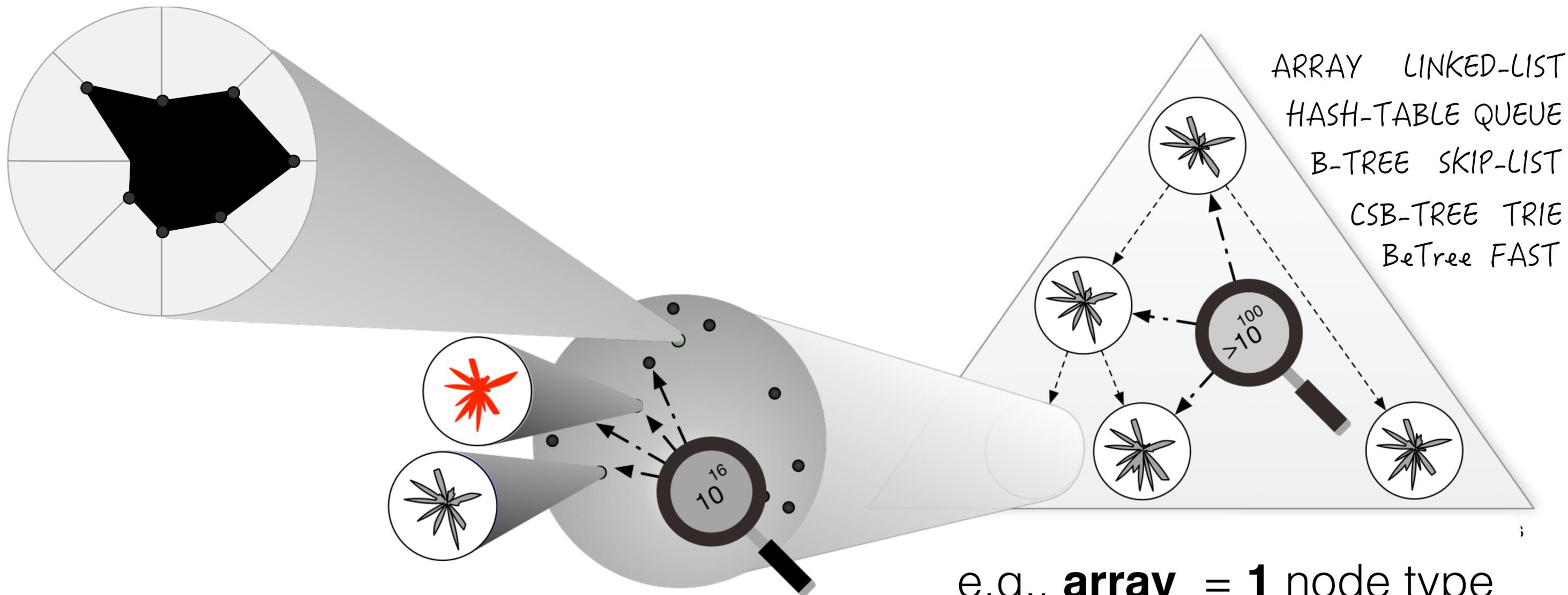
POSSIBLE STRUCTURES



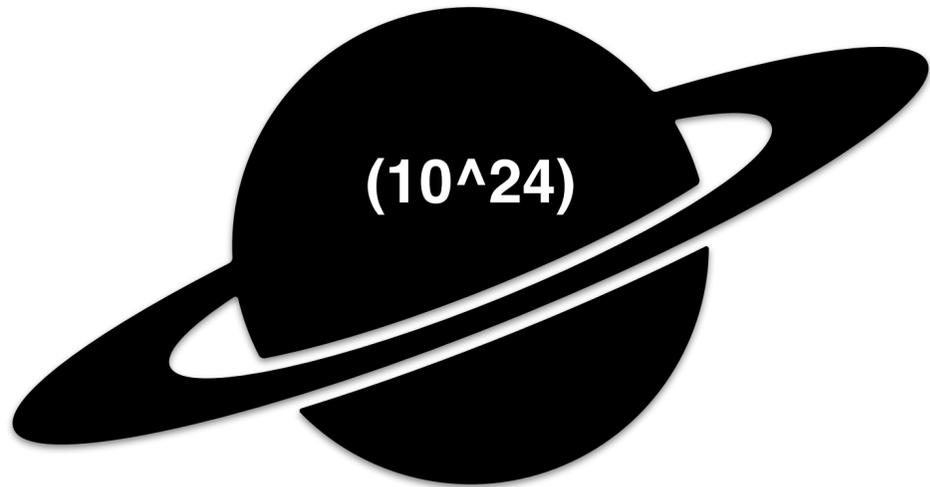
DESIGN PRINCIPLES

POSSIBLE NODE DESIGNS

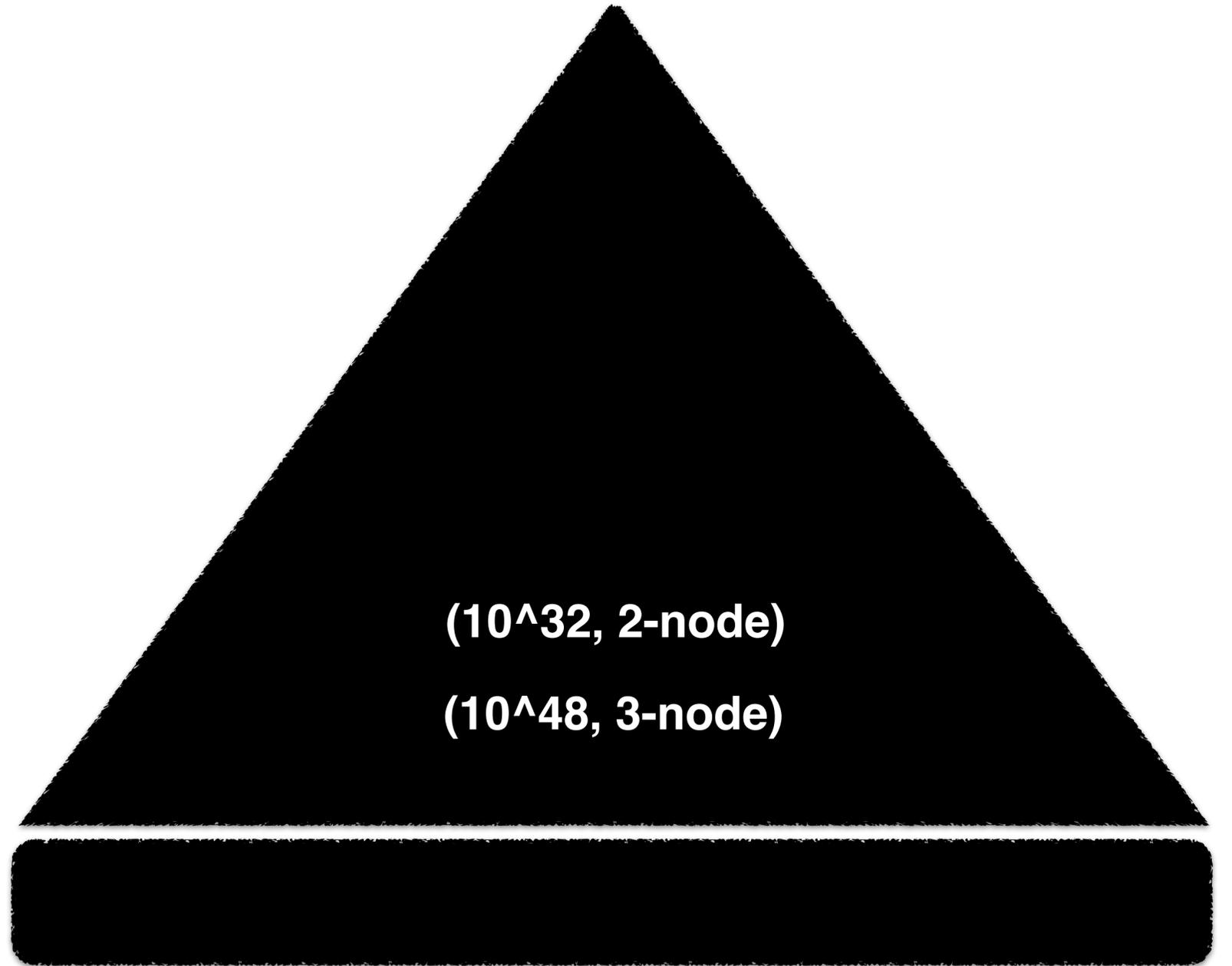
POSSIBLE STRUCTURES



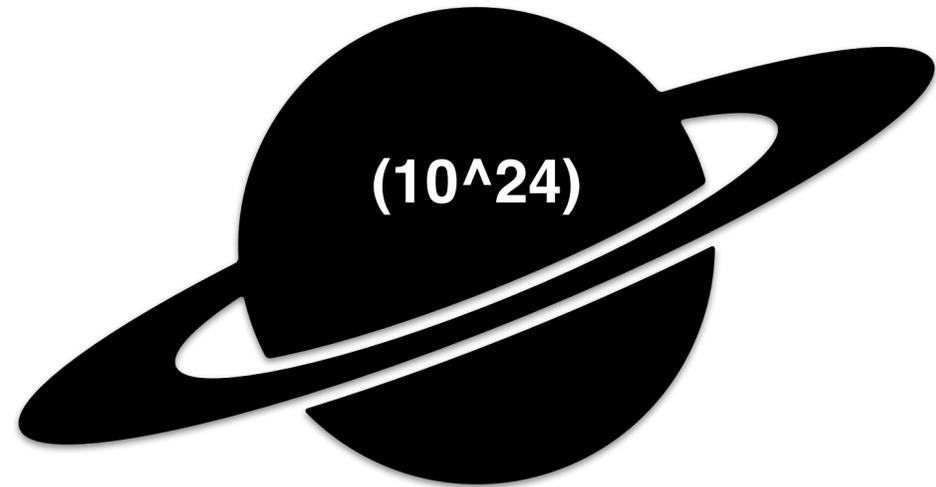
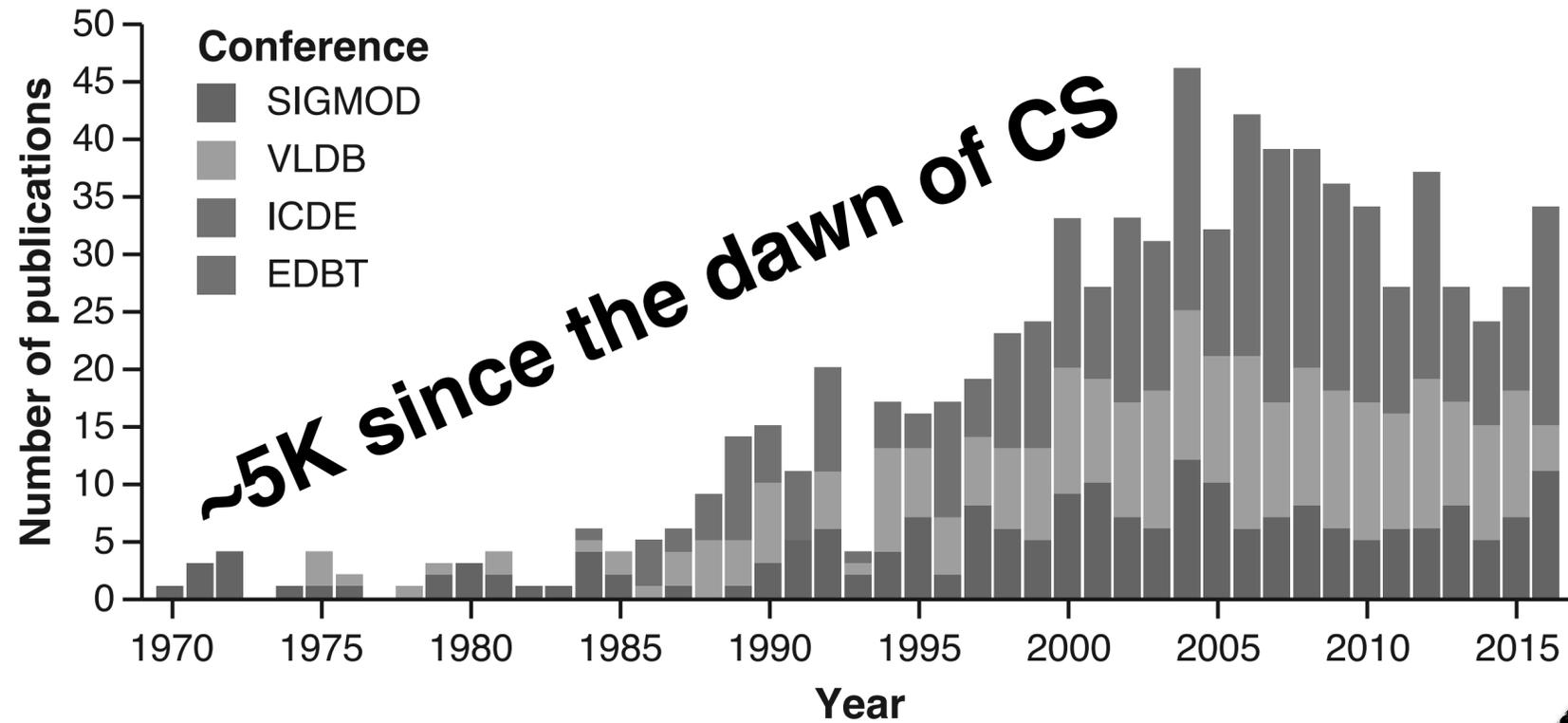
e.g., **array** = **1** node type
e.g., **b-tree** = **2** node types



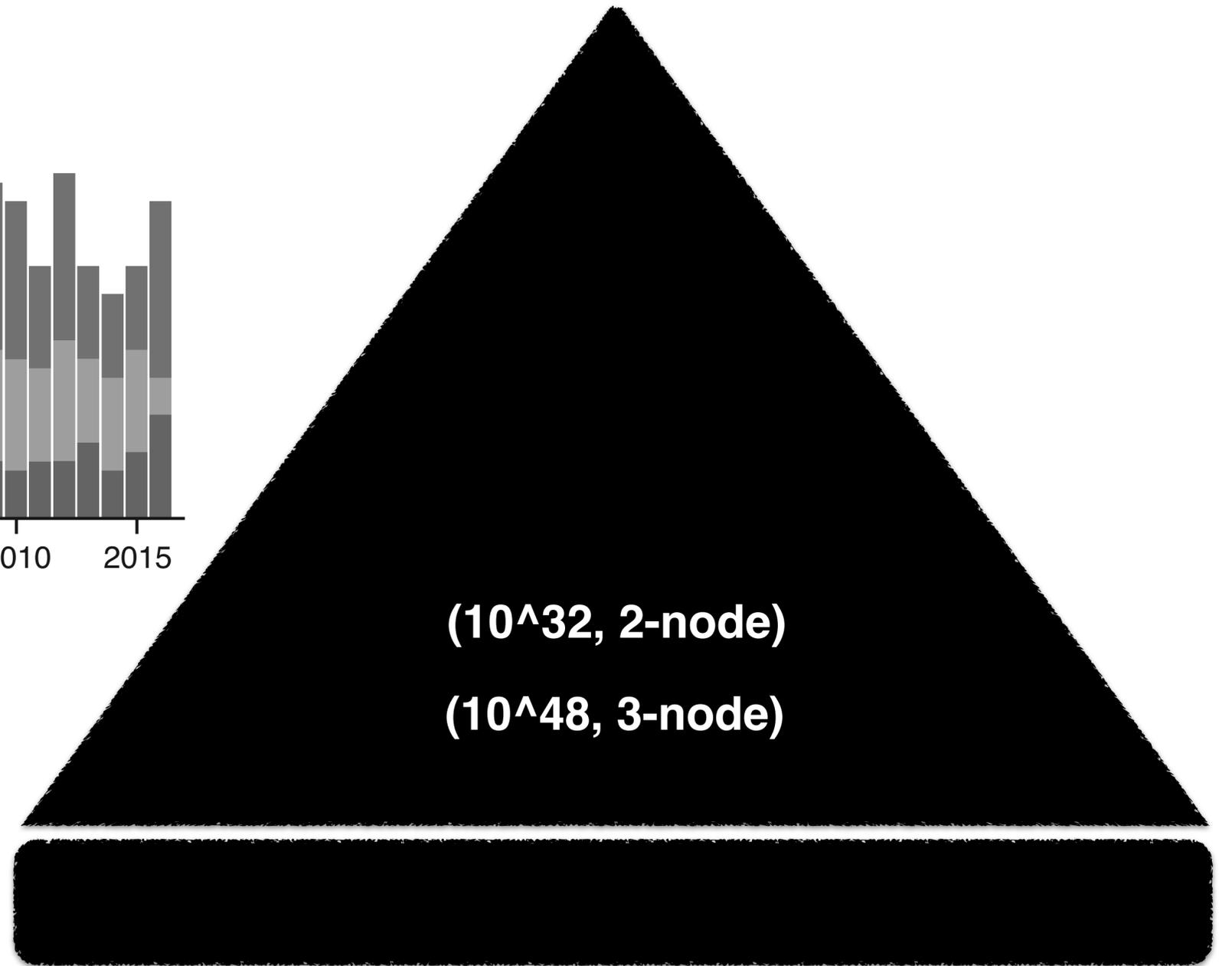
STARS ON THE SKY



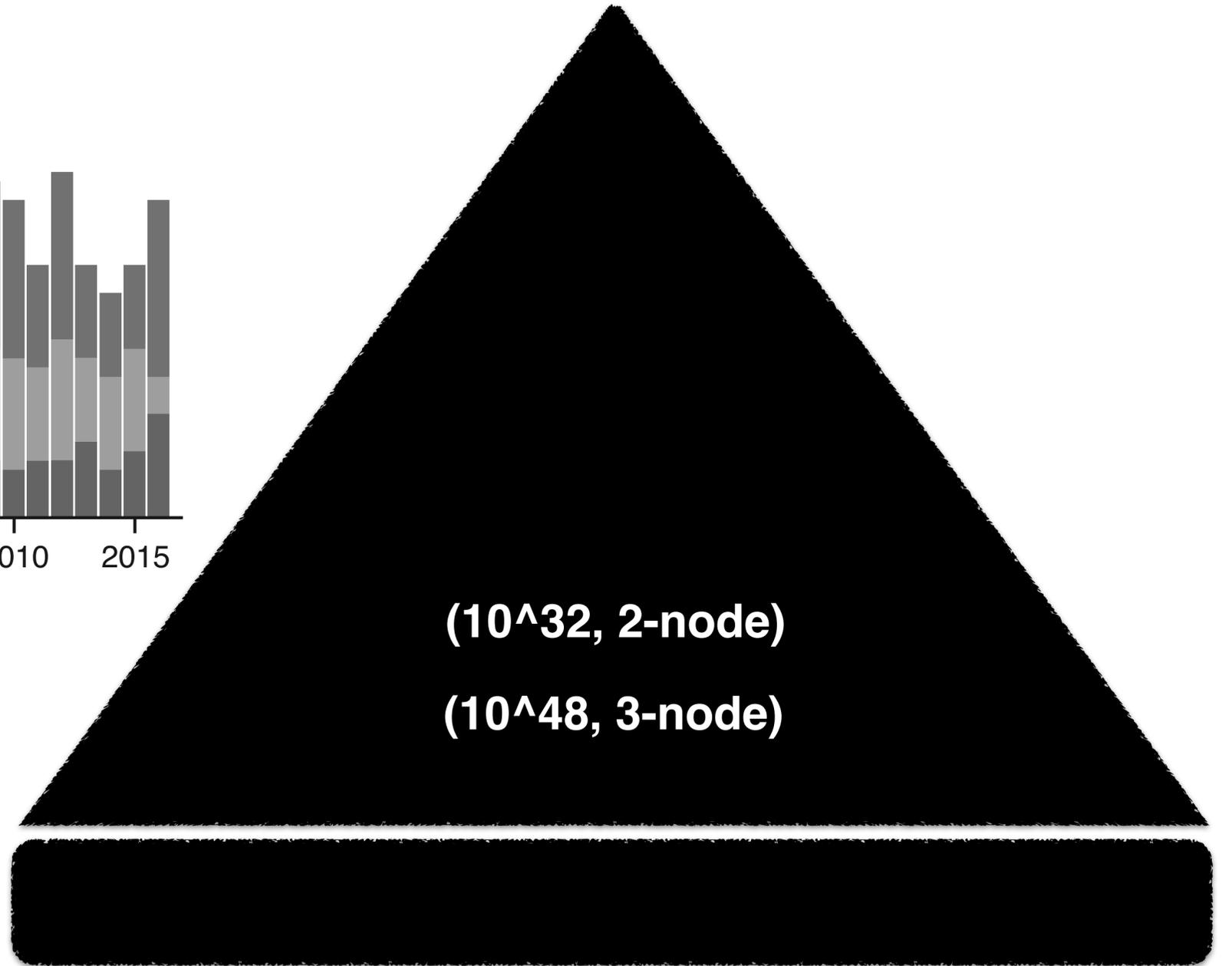
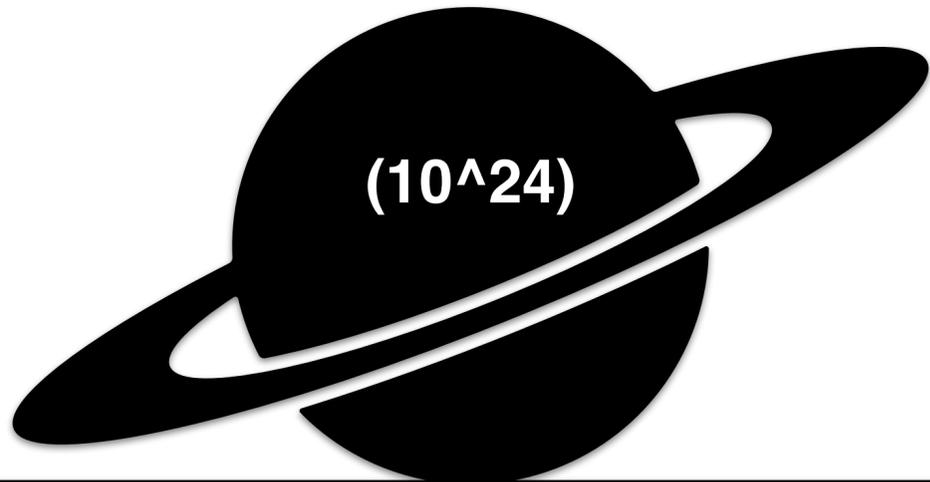
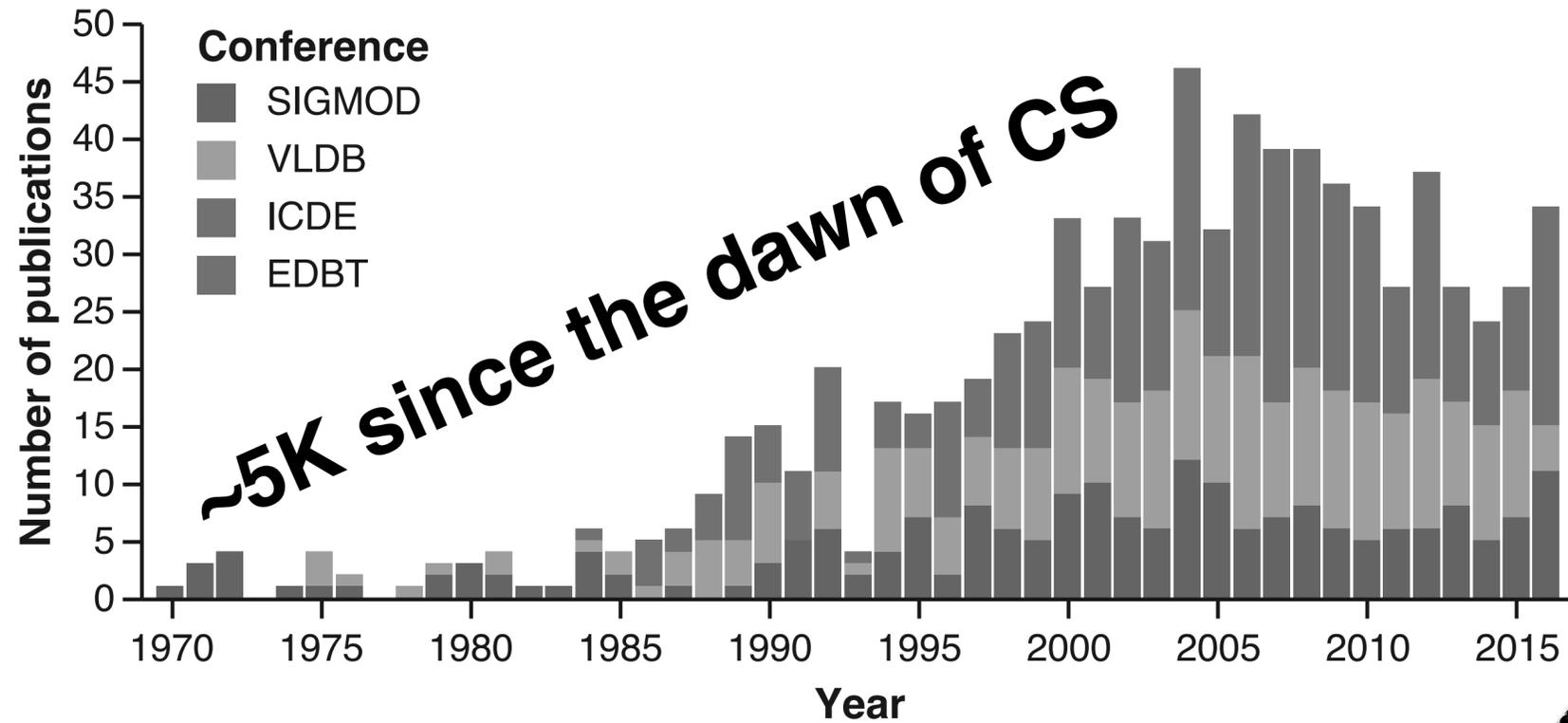
POSSIBLE DATA STRUCTURES



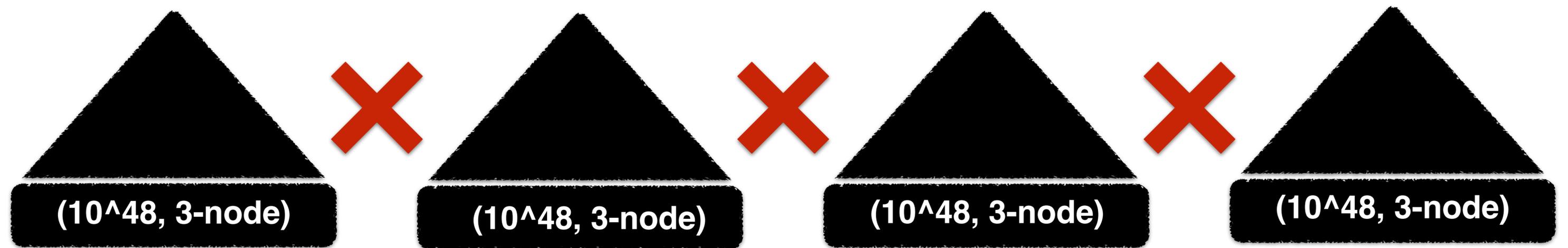
STARS ON THE SKY



POSSIBLE DATA STRUCTURES

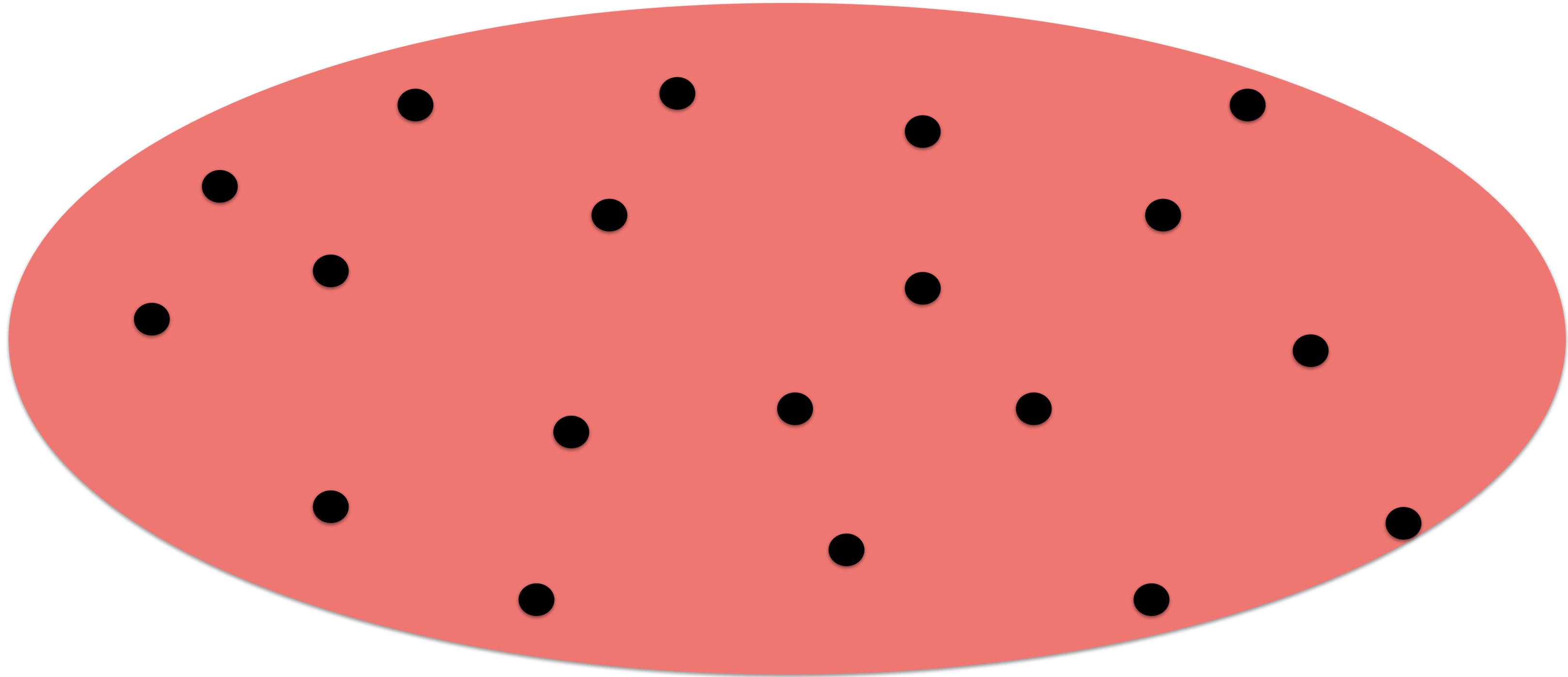


$$10^{48} - 5 \times 10^3 = 10^{48} \quad \text{zero progress}$$

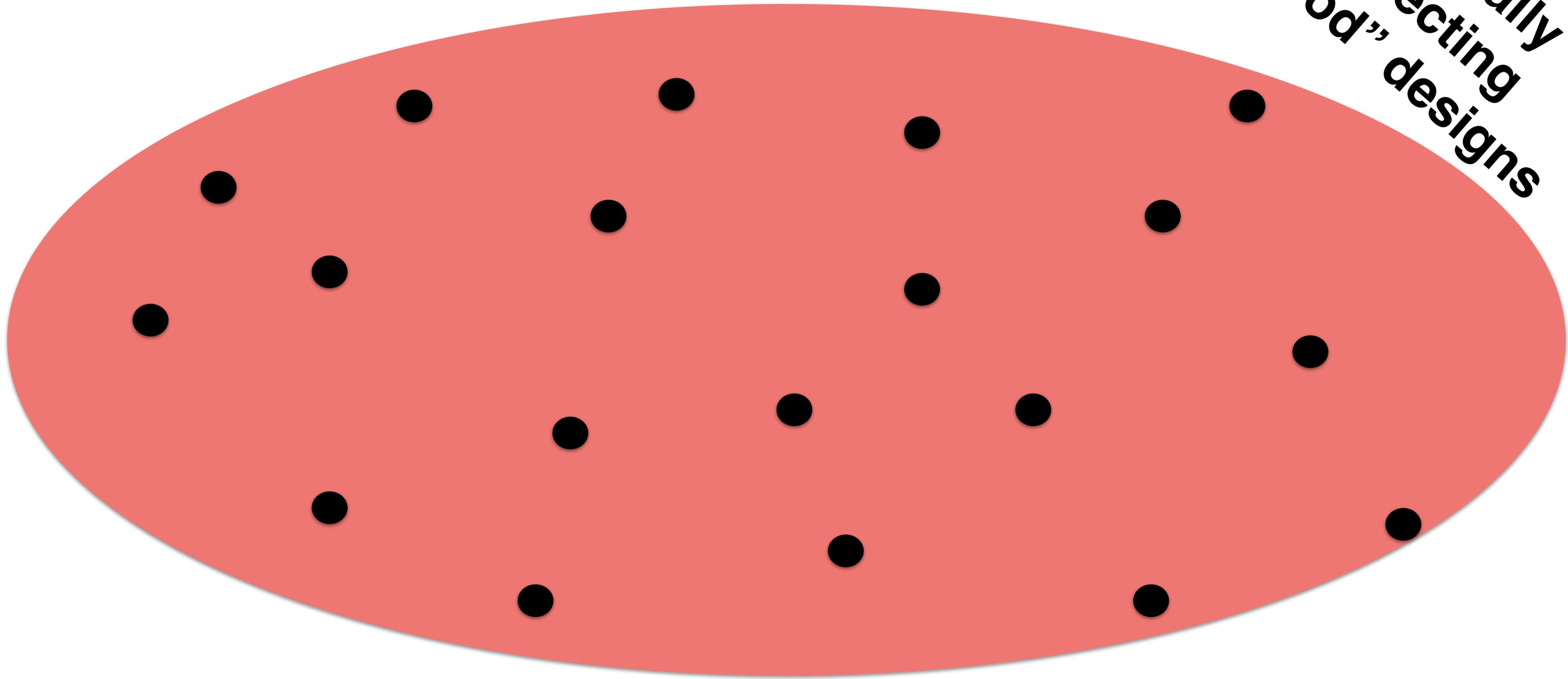


The design space of systems is even larger

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**

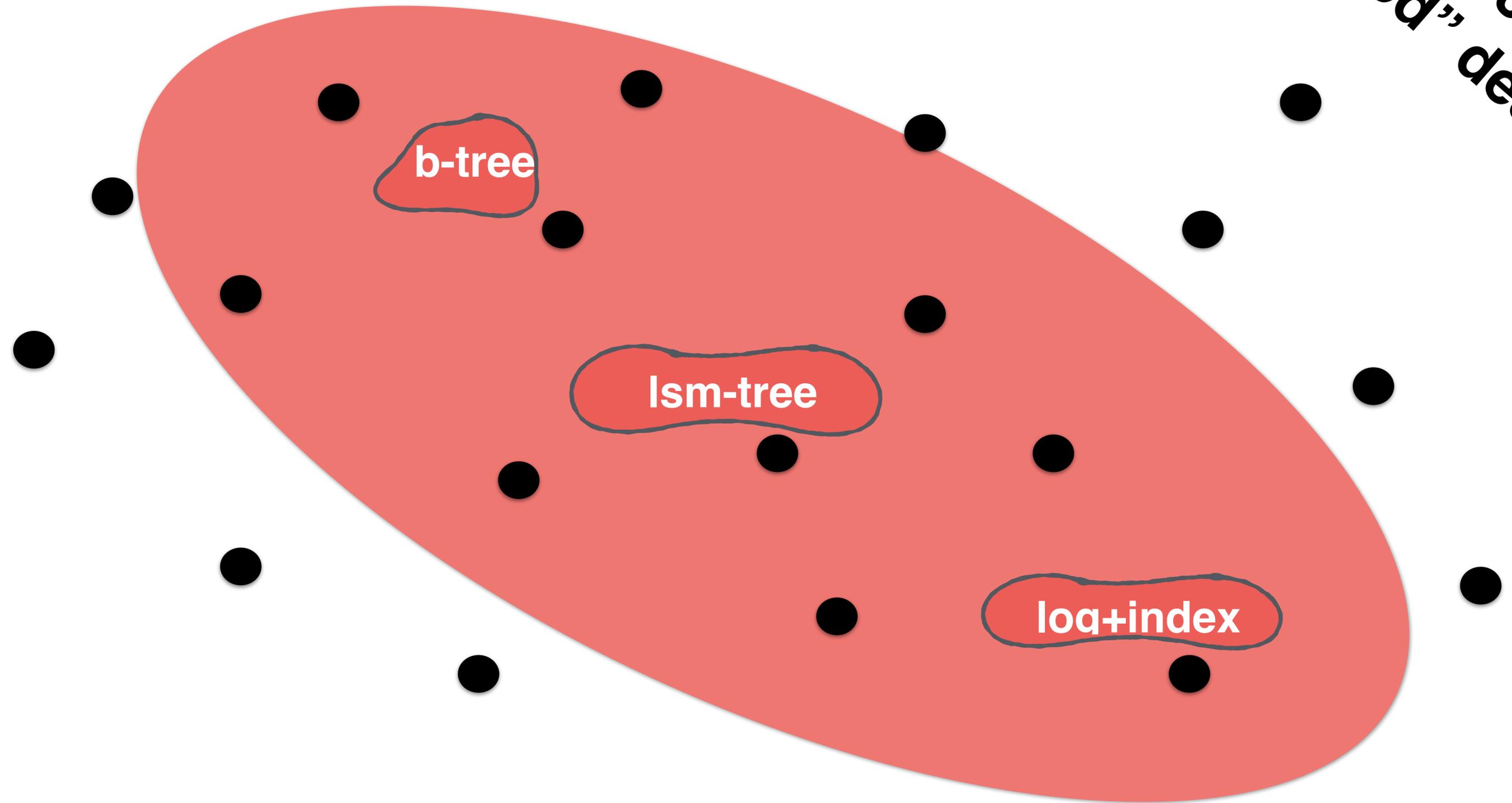
manually
selecting
“good” designs

b-tree

lsm-tree

log+index

manually
selecting
“good” designs



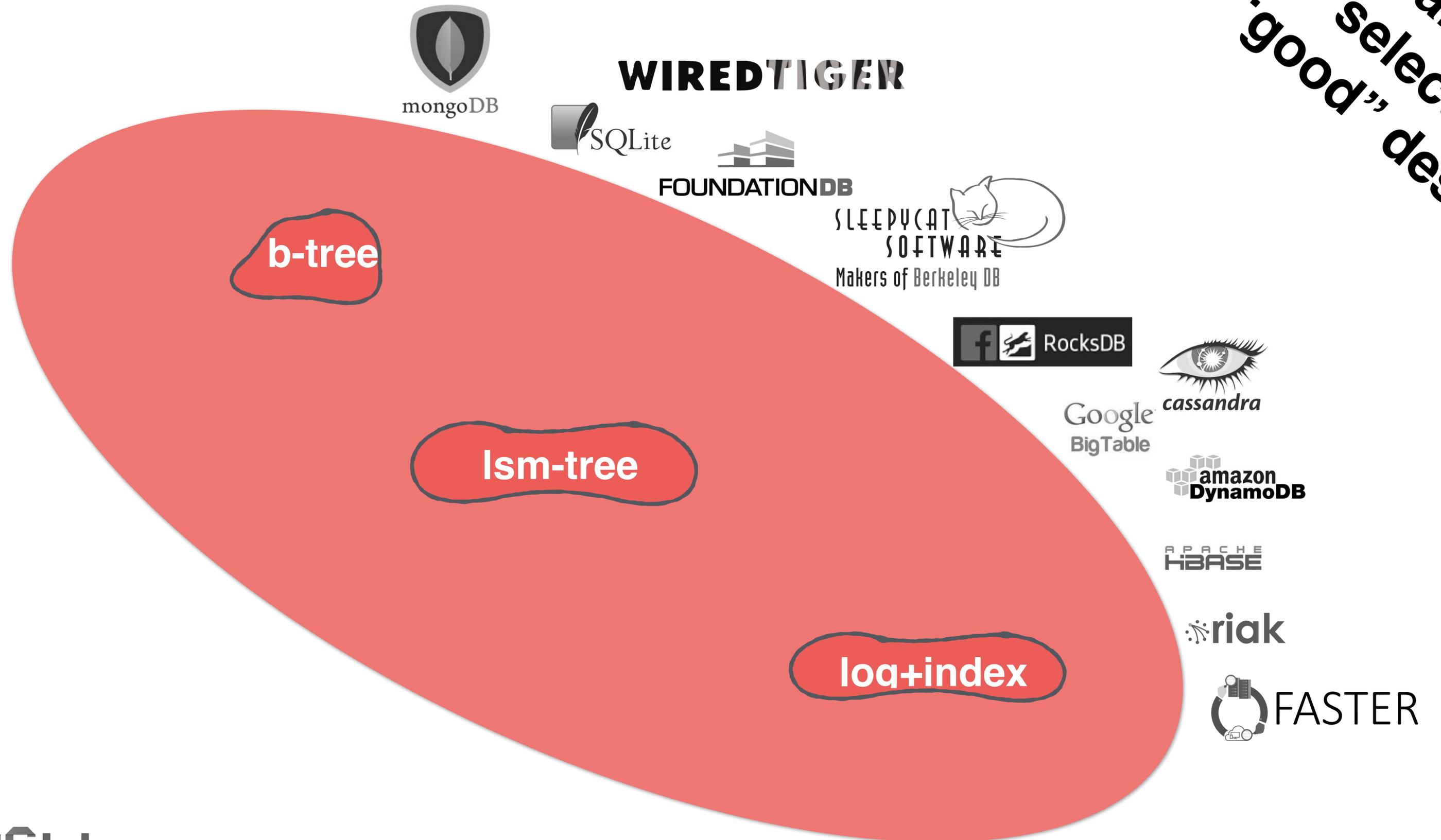
manually
selecting
“good” designs

b-tree

lsm-tree

log+index

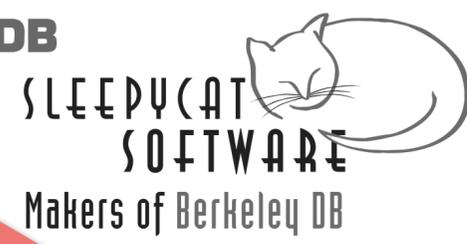
manually selecting "good" designs



WIREDTIGER



FOUNDATIONDB



Google BigTable



APACHE HBASE

riak

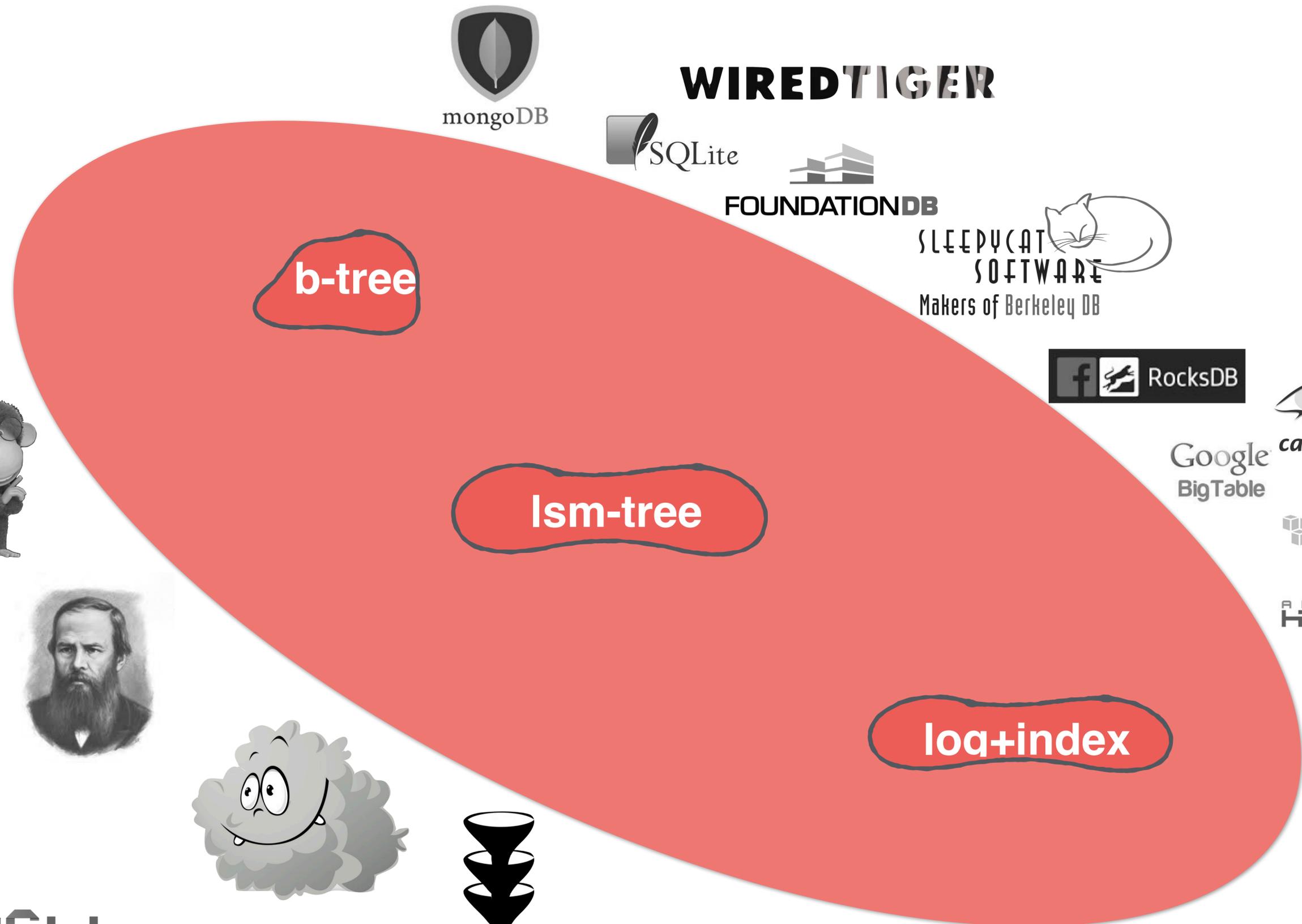


b-tree

lsm-tree

log+index

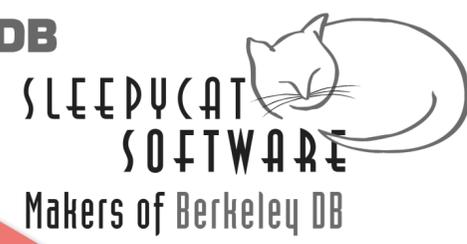
manually selecting "good" designs



WIREDTIGER



FOUNDATIONDB



Google BigTable



APACHE HBASE

riak



b-tree

lsm-tree

log+index



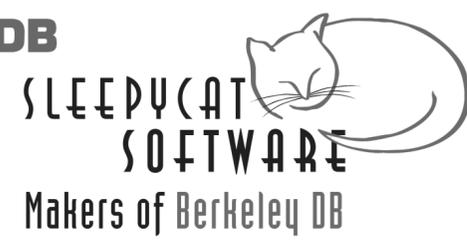
manually selecting "good" designs



WIREDTIGER



FOUNDATIONDB

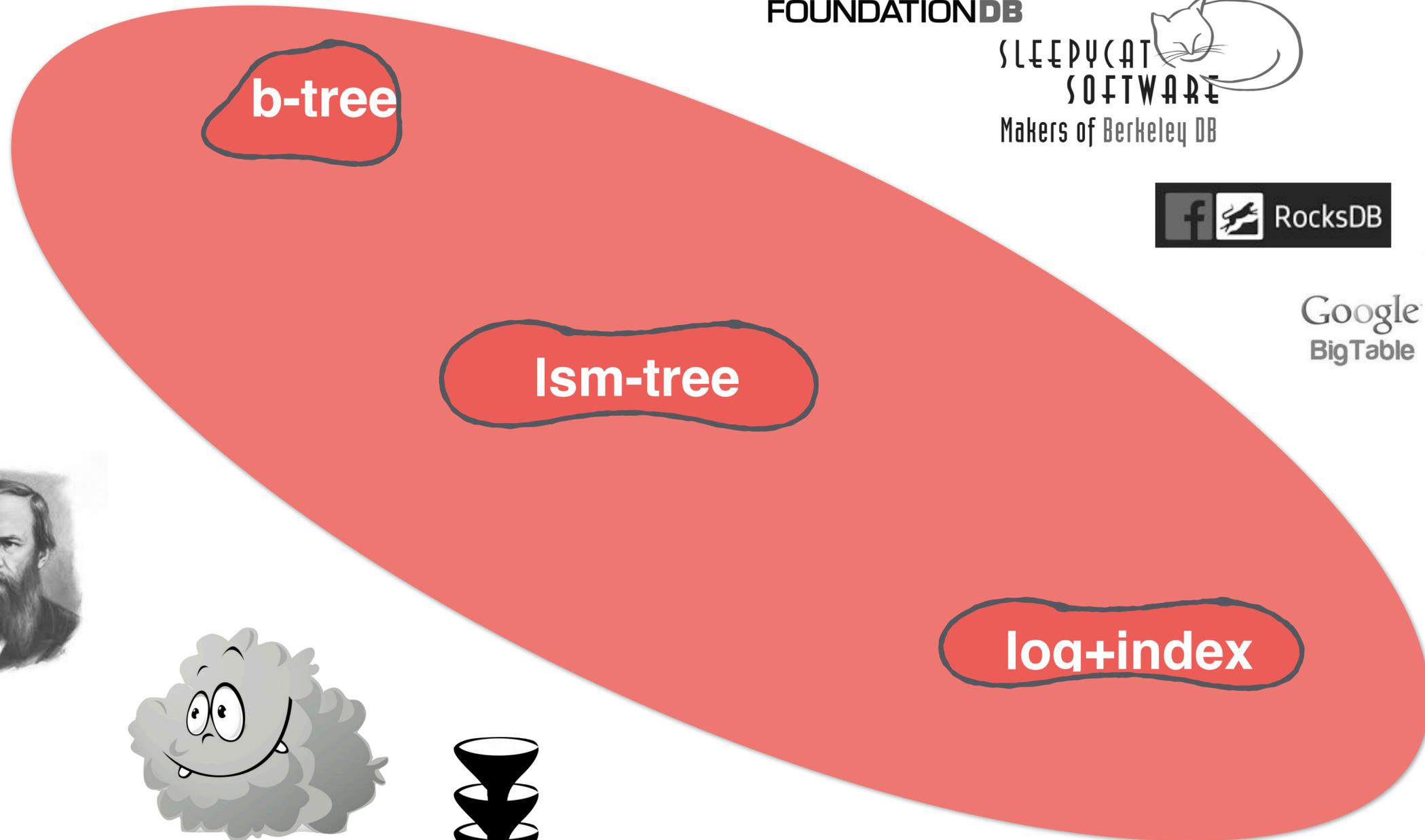


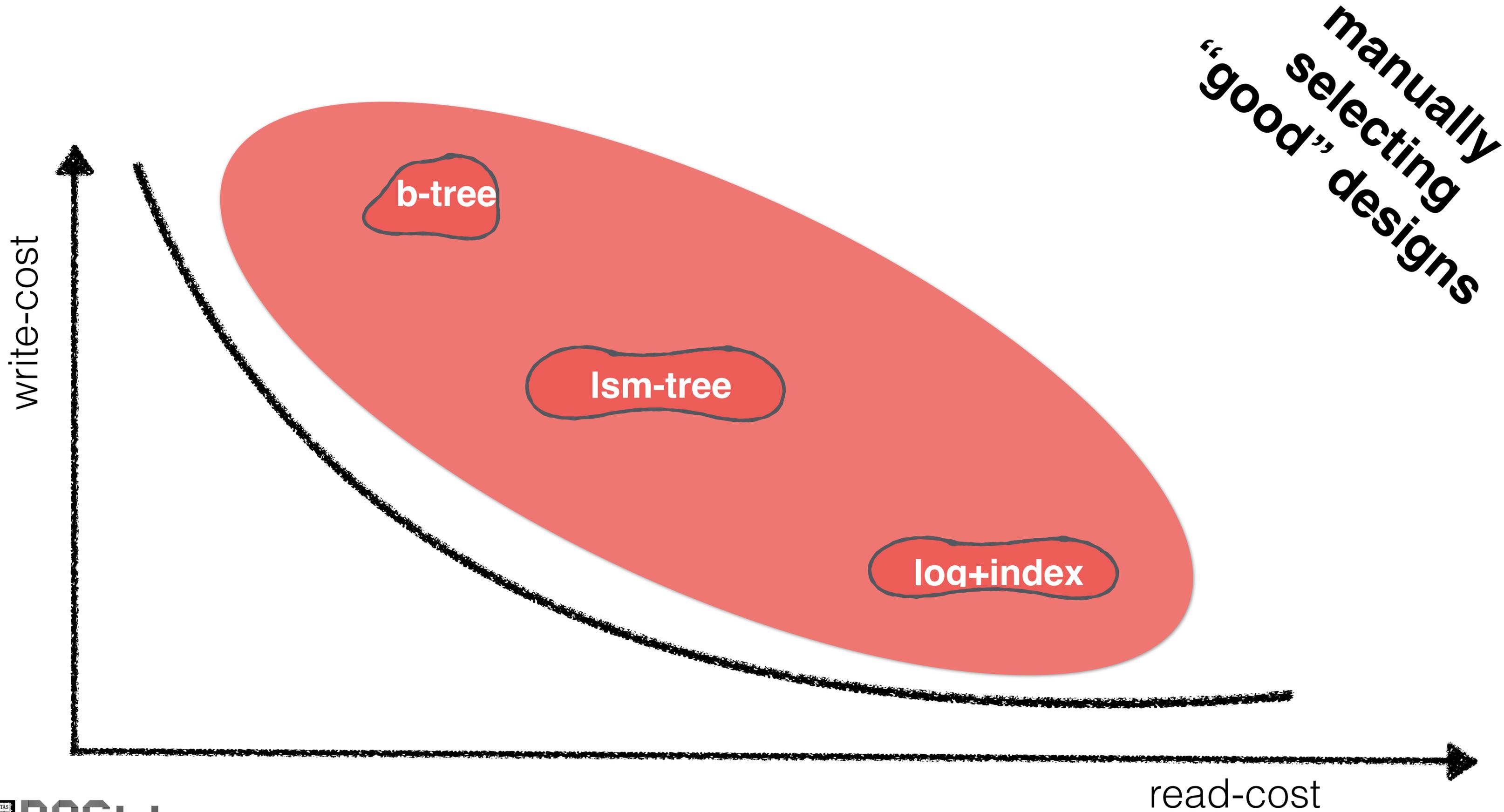
Google BigTable

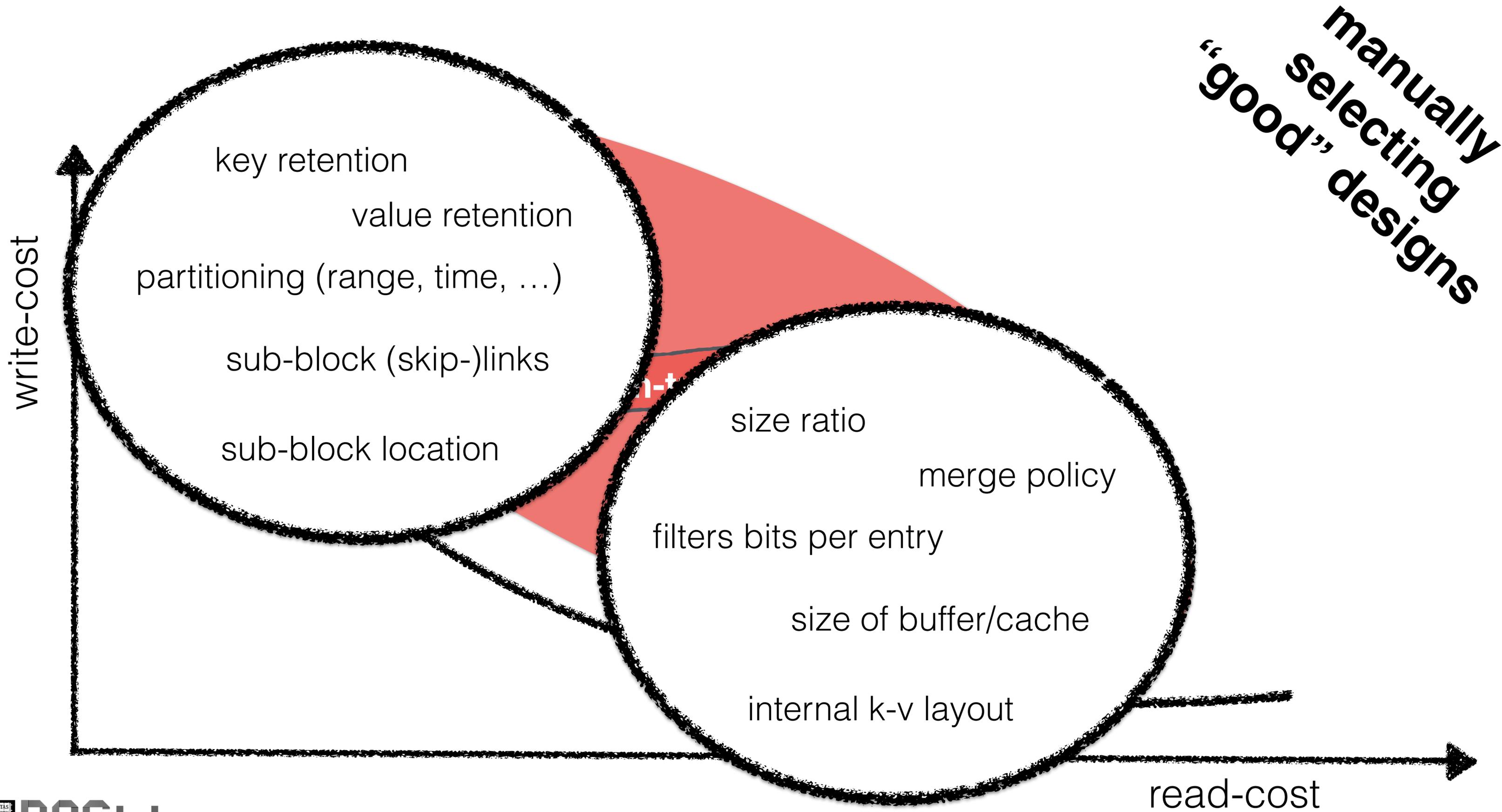


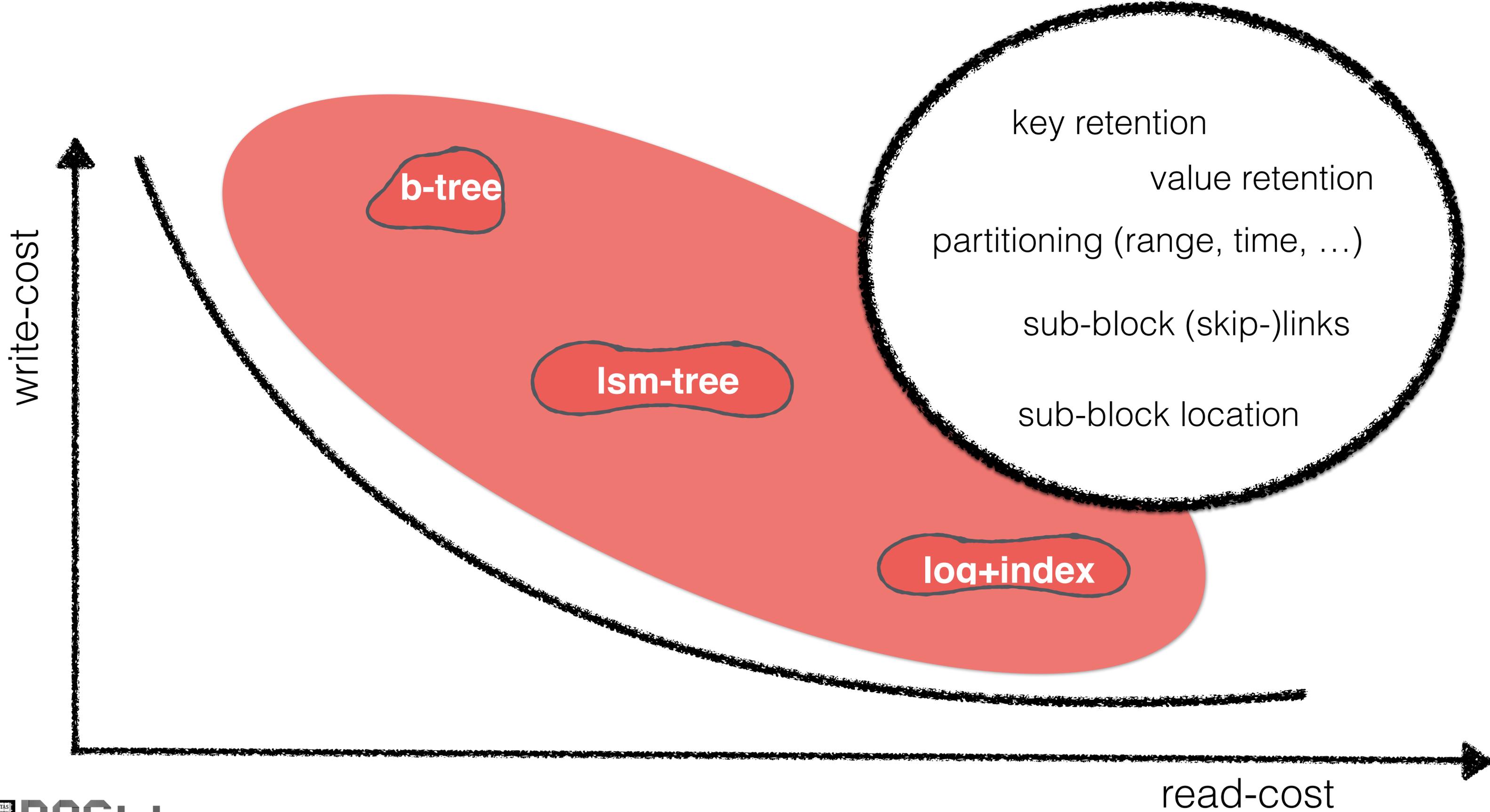
APACHE HBASE

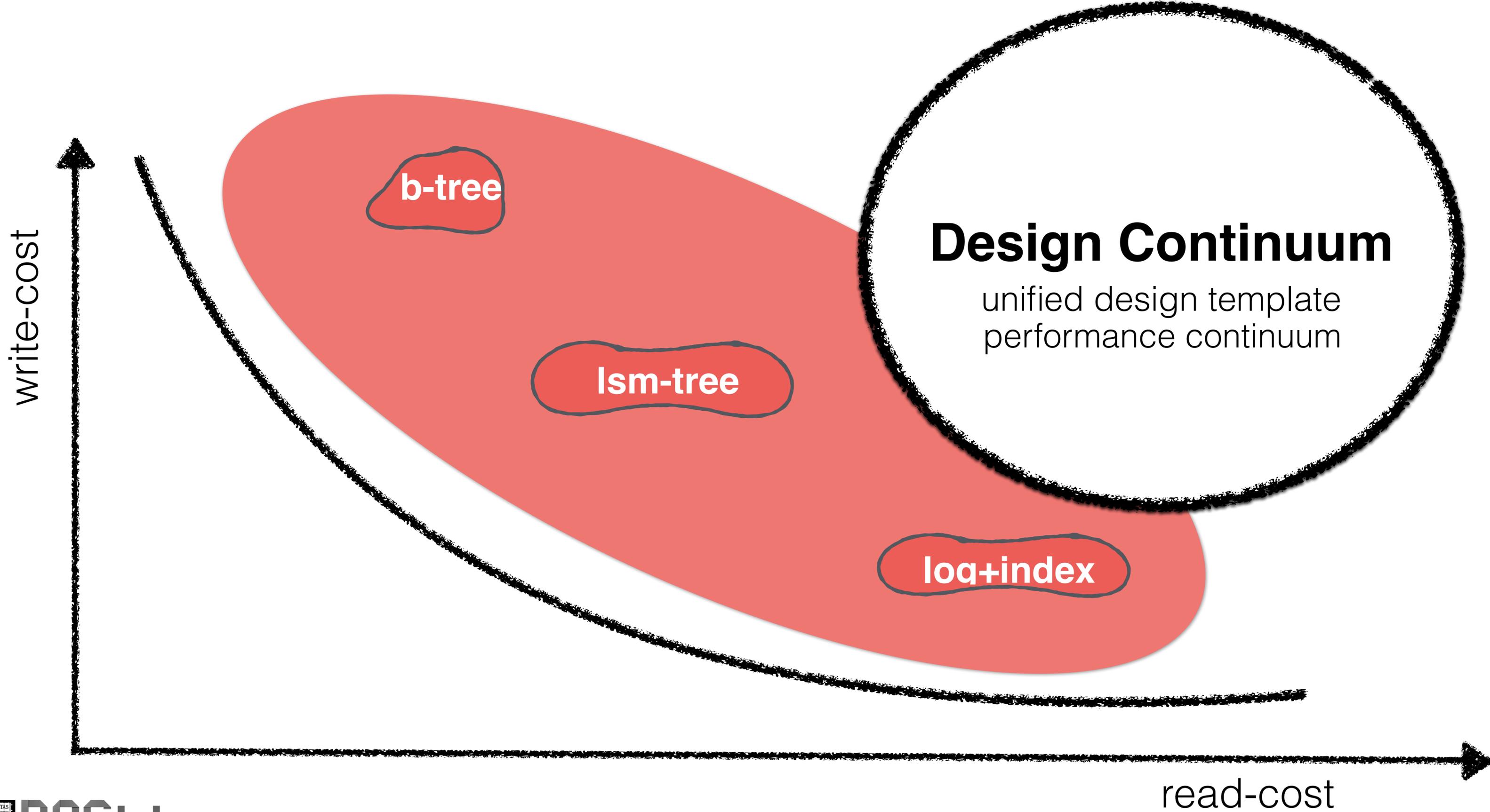
riak











Cosine
@PVLDB2022

Cloud-cost
Optimized

Self
Designing

Key-value
Store

↑ LAYOUT PRIMITIVES ↓
↑ ALGORITHMIC ABSTRACTIONS ↓

		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 <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)
Parallelism <i>derived with empirically verified rules</i>	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

↑ LAYOUT PRIMITIVES ↓
↑ ALGORITHMIC ABSTRACTIONS ↓

		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 <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)
Parallelism <i>derived with empirically verified rules</i>	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

↑ LAYOUT PRIMITIVES ↓
↑ ALGORITHMIC ABSTRACTIONS ↓

		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 <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)
Parallelism <i>derived with empirically verified rules</i>	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

↑ LAYOUT PRIMITIVES ↓
↑ ALGORITHMIC ABSTRACTIONS ↓

		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 <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)
Parallelism <i>derived with empirically verified rules</i>	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

↑ LAYOUT PRIMITIVES ↓
↑ ALGORITHMIC ABSTRACTIONS ↓

		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 <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)
Parallelism <i>derived with empirically verified rules</i>	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.

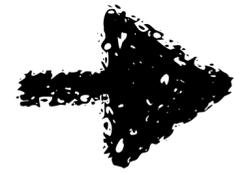
massive design space of system designs



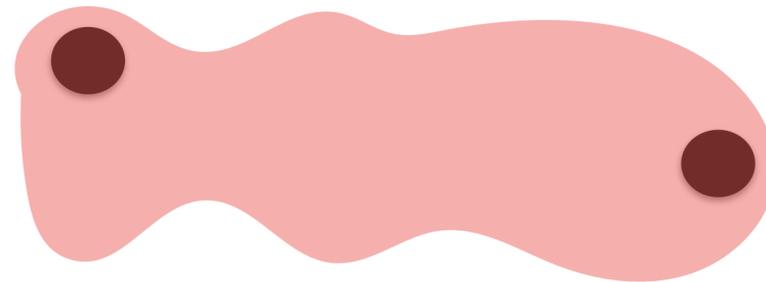
massive design space of system designs



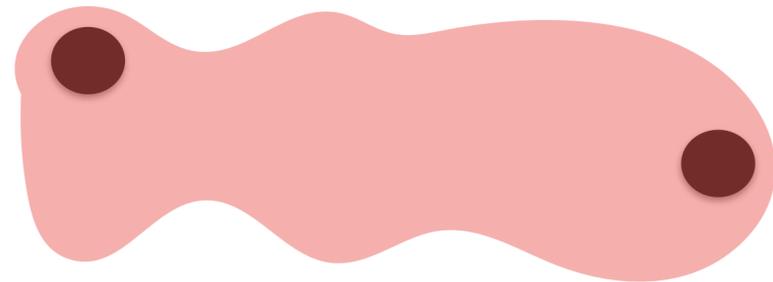
workload/hardware



performance, cloud cost, robustness

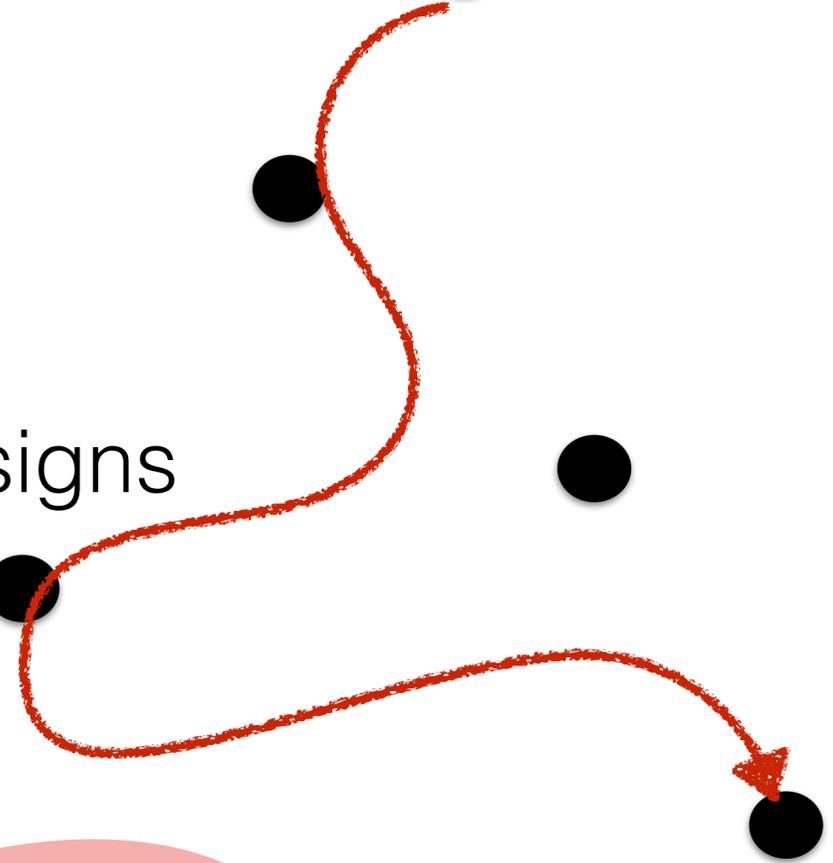
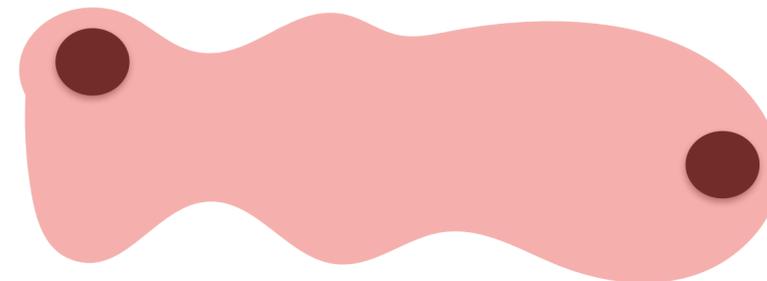


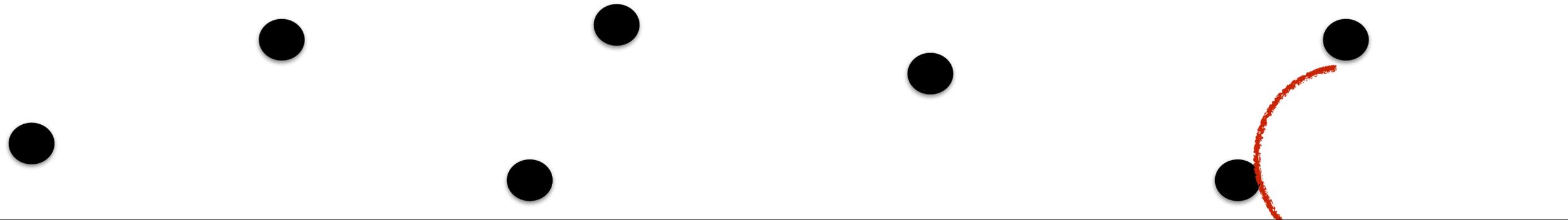
workload/hardware → performance, cloud cost, robustness



without having to code

massive design space of system designs

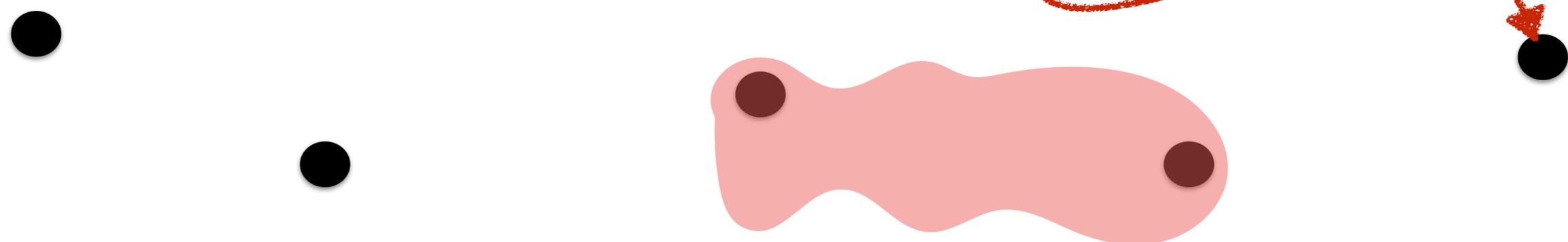




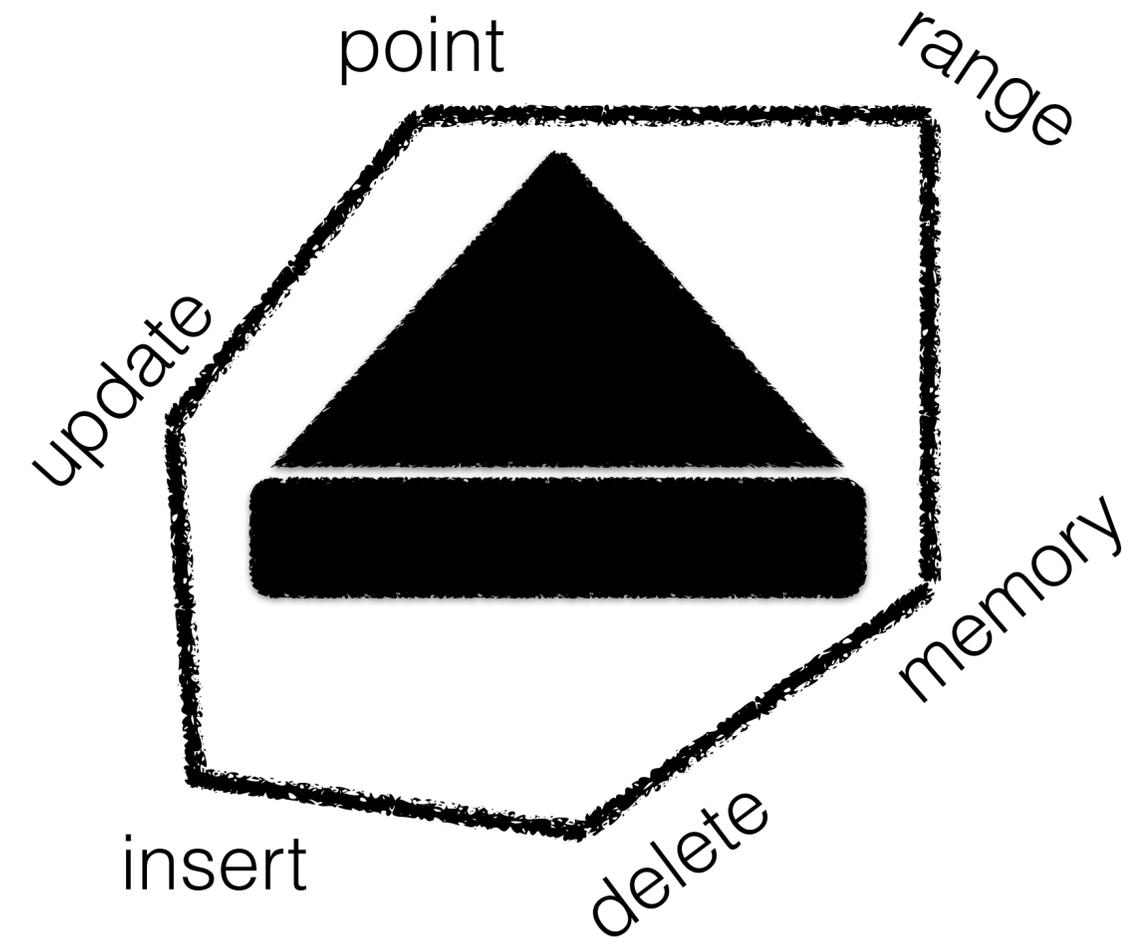
data layouts



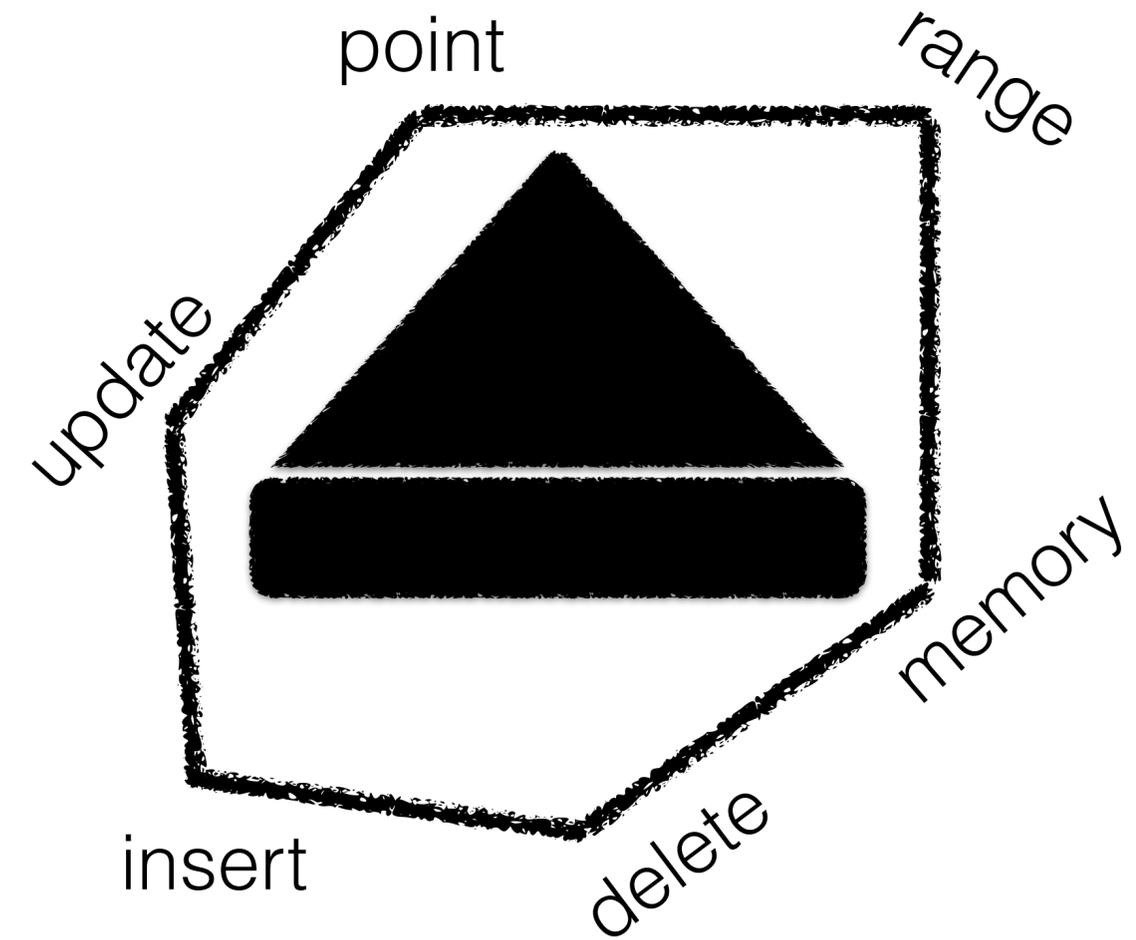
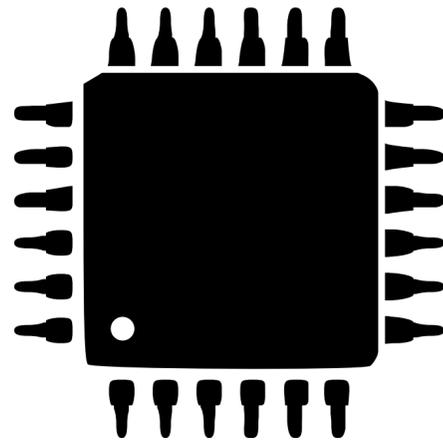
algorithm & cost synthesis



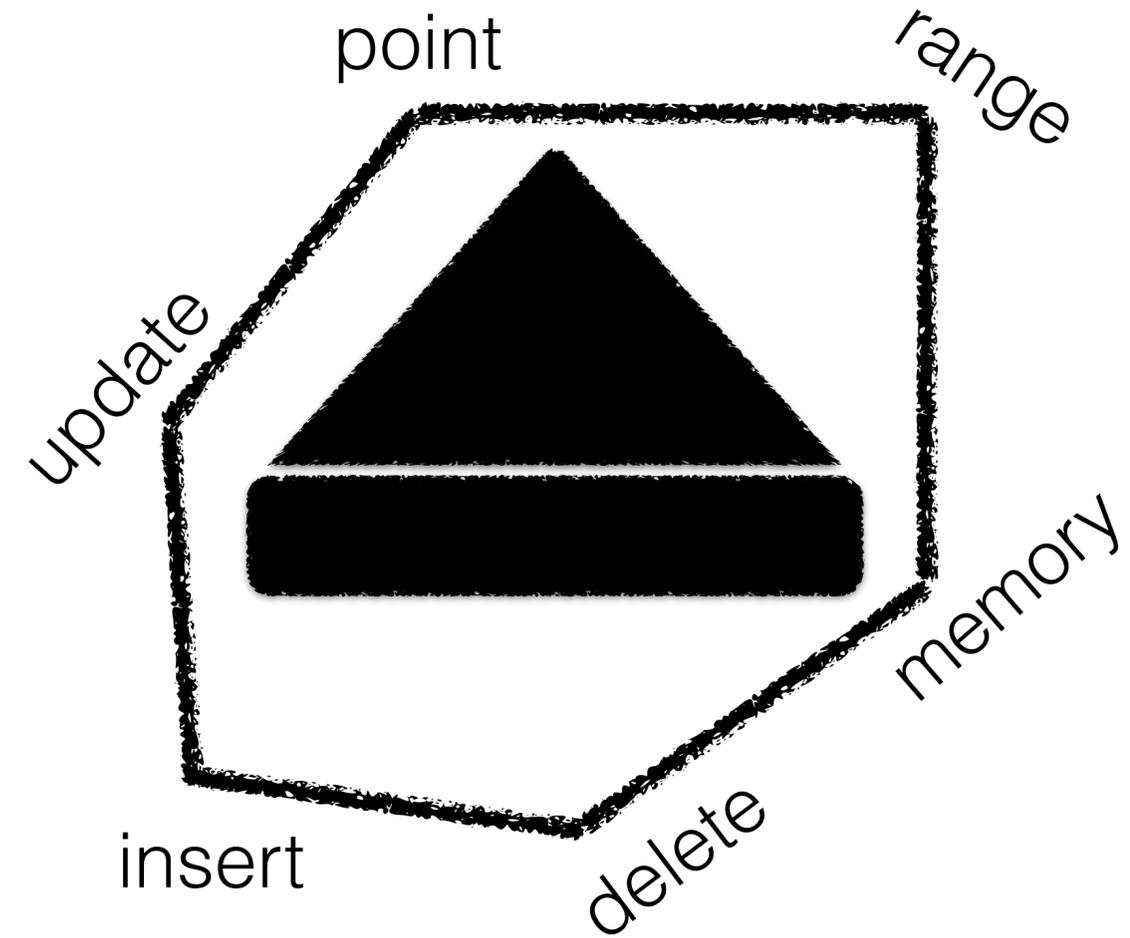
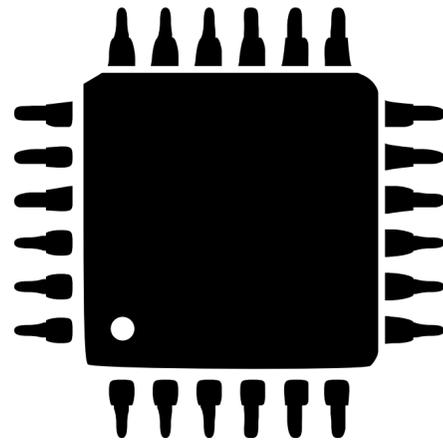
algorithm & cost synthesis



algorithm & cost synthesis



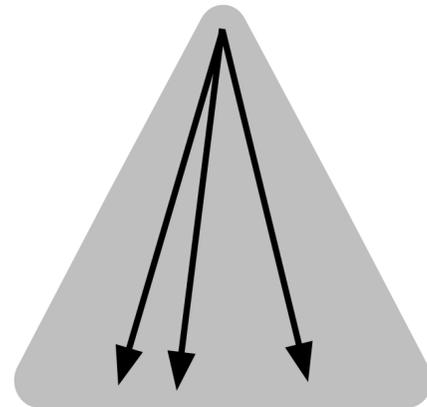
algorithm & cost synthesis



query &
hardware parallelism

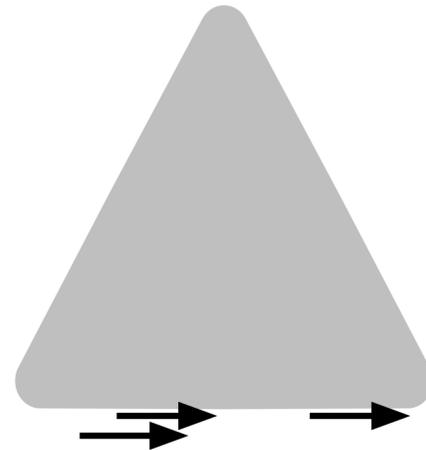
Should I scan or should I probe? @SIGMOD2016

Tree Traversal



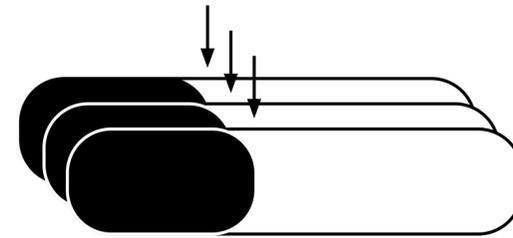
+

Leaf Traversal



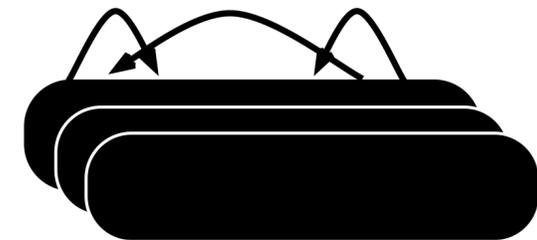
+

Result Writing



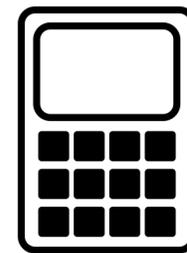
+

Sorting



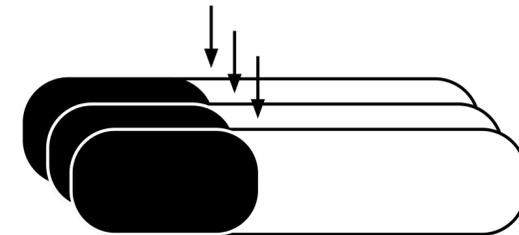
Base Scan

+



$\sum p_i$

+

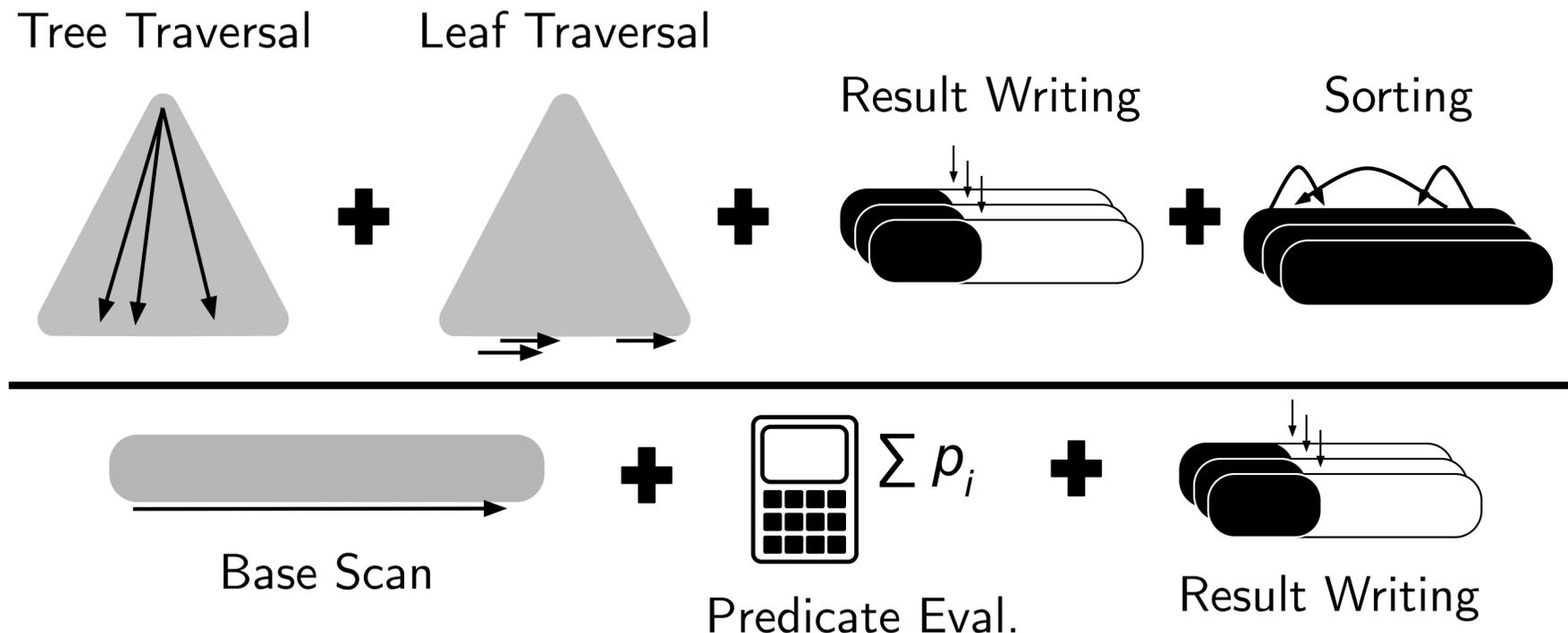


Predicate Eval.

Result Writing

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

Should I scan or should I probe? @SIGMOD2016

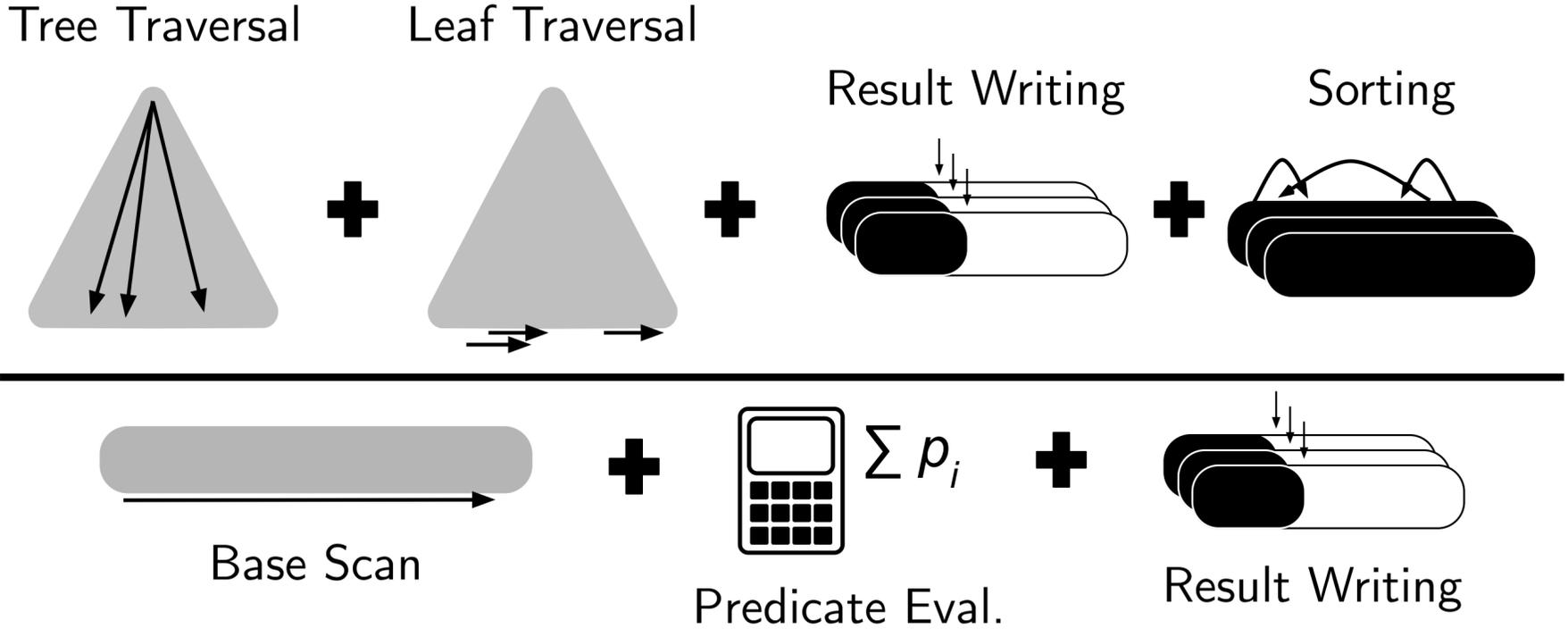


Workload	q s_i S_{tot}	number of queries selectivity of query i total selectivity of the workload
Dataset	N ts	data size (tuples per column) tuple size (bytes per tuple)
Hardware	C_A C_M BW_S BW_R BW_I p f_p	L1 cache access (sec) LLC miss: memory access (sec) scanning bandwidth (GB/s) result writing 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)

But we have a massive design space to cover.....

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

Should I scan or should I probe? @SIGMOD2016



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)

But we have
 massive des
 space to cover

$$APS(q, S_{tot}) = \frac{q \cdot \left(\frac{1}{2} \cdot (b \cdot BW_S \cdot C_A + \frac{b \cdot BW_S \cdot f_p \cdot p}{2}) \right)}{q \cdot BW_S} + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}$$

$$\frac{(w + ow) \cdot \frac{BW_S}{BW_I} + rw \cdot \frac{BW_S}{BW_R}}{p \cdot p \cdot q \cdot BW_S} + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}$$

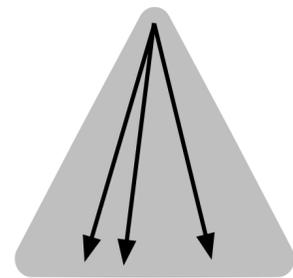
$$\frac{t \cdot \log_2(S_{tot} \cdot N) \cdot BW_S \cdot C_A}{s, 2 \cdot f_p \cdot p \cdot q \cdot BW_S} + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}$$

Should I scan or should I probe? @

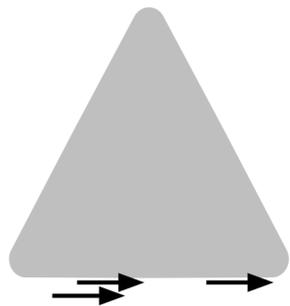
Tree Traversal

Leaf Traversal

Result



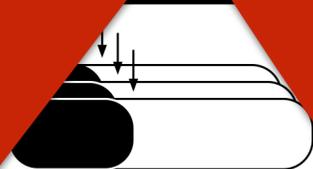
+



+



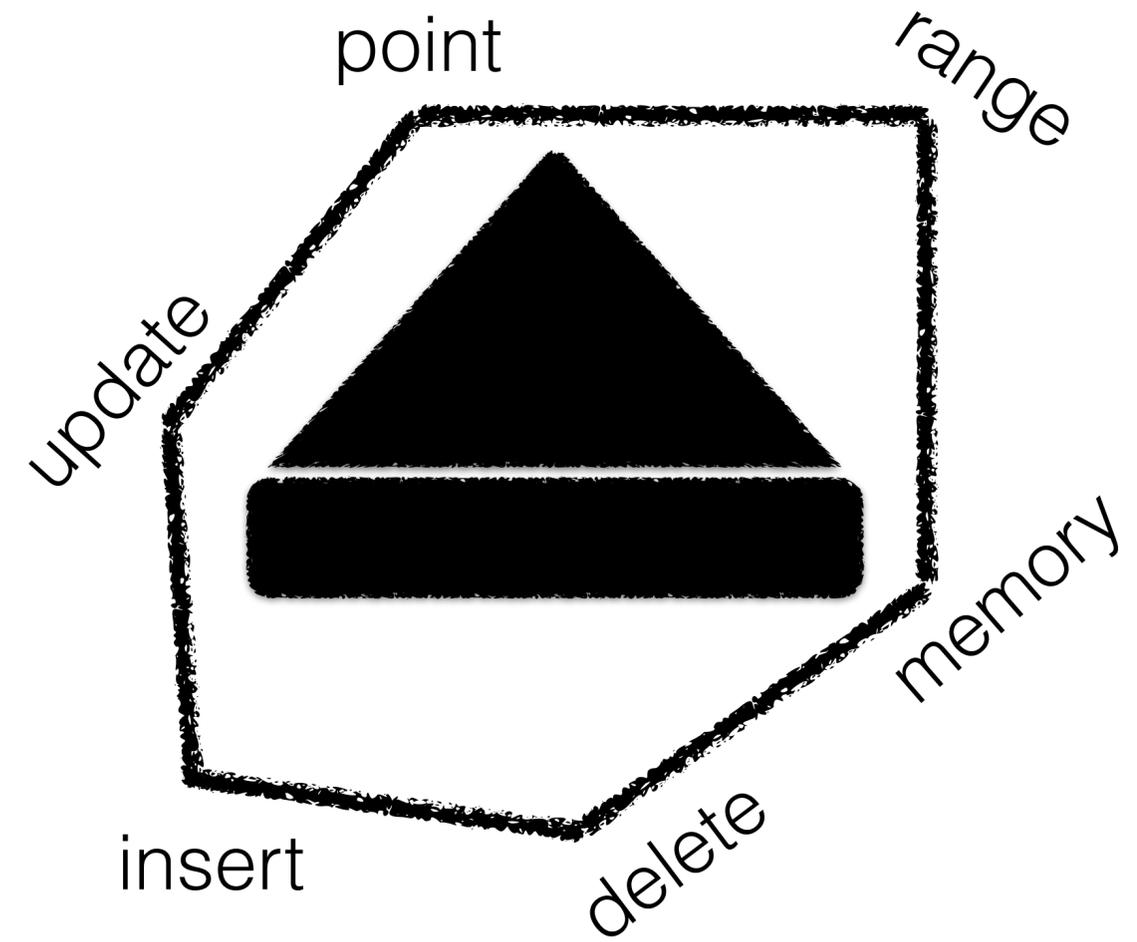
Base Scan



Result Writing

workload	q	number of queries
	s_i	selectivity of query i
	S_{tot}	total selectivity of the workload
	N	data size (tuples per column)
	ts	tuple size (bytes per tuple)
	C_A	L1 cache access (sec)
	m	LLC miss: memory access (sec)
	s	scanning bandwidth (GB/s)
	w	result writing bandwidth (GB/s)
	ow	leaf traversal bandwidth (GB/s)
	f_p	inverse of CPU frequency
	p	for accounting for pipelining
	b	result width (bytes per output tuple)
	2	tree fanout
	a	attribute width (bytes of the indexed column)
	o	offset width (bytes of the index column offset)

algorithm & cost synthesis



algorithm & cost synthesis

scan

hash
probe

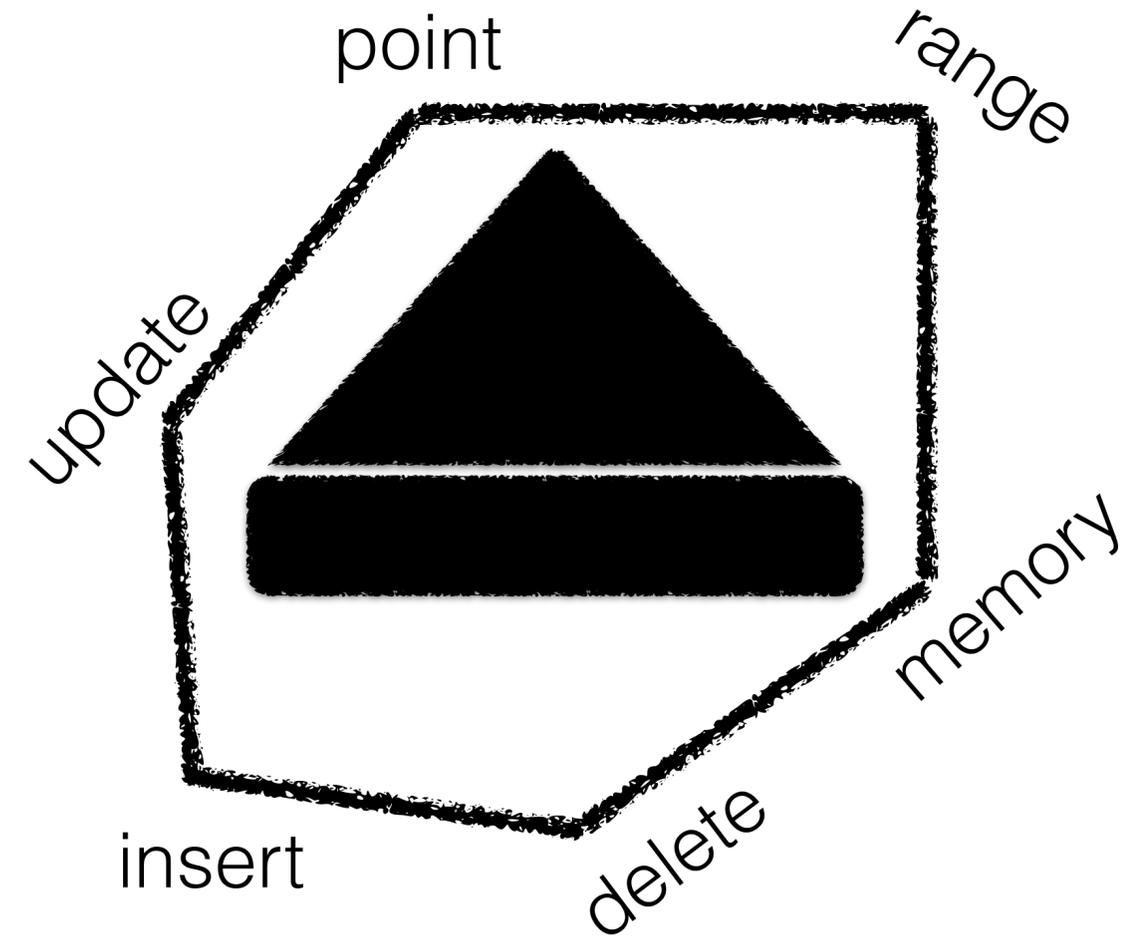
random
access

sorted
search

bloom
probe



synthesize from first principles



algorithm & cost synthesis

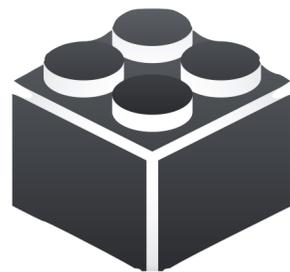
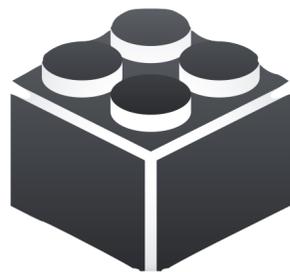
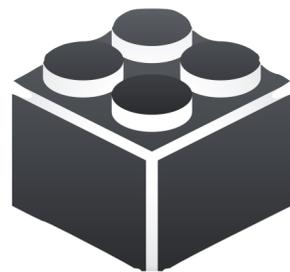
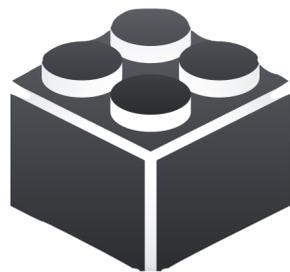
scan

hash
probe

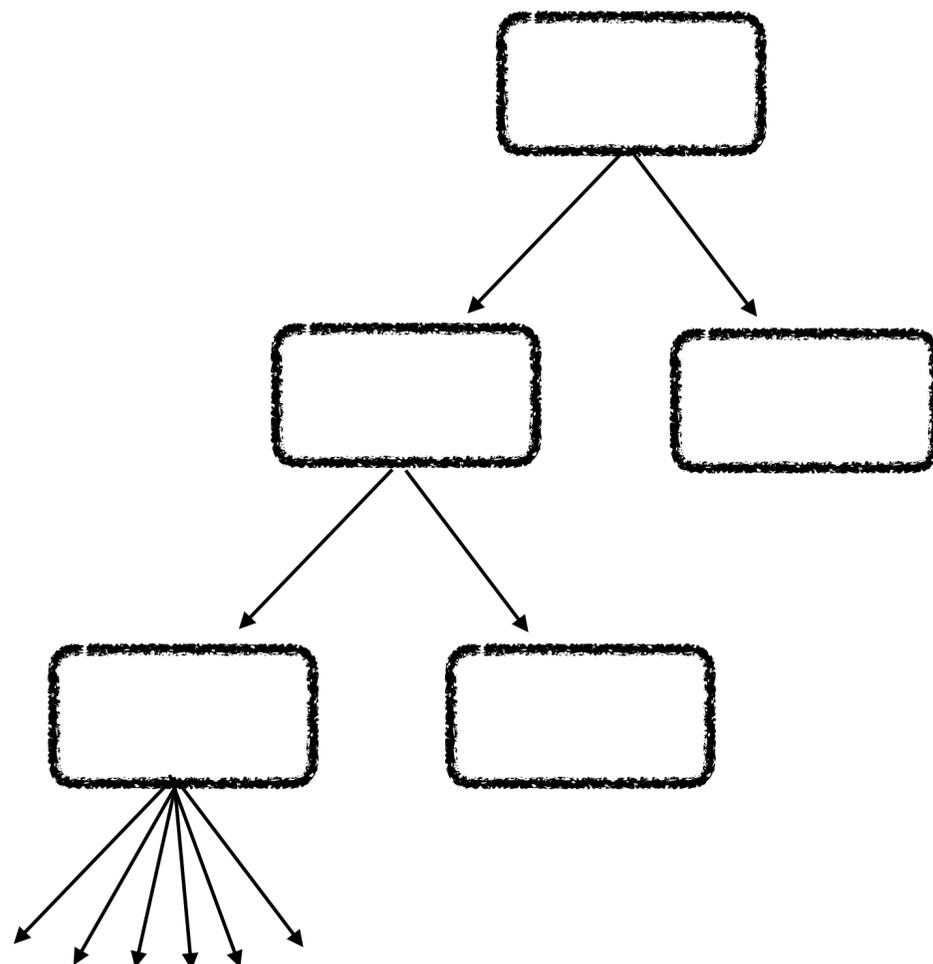
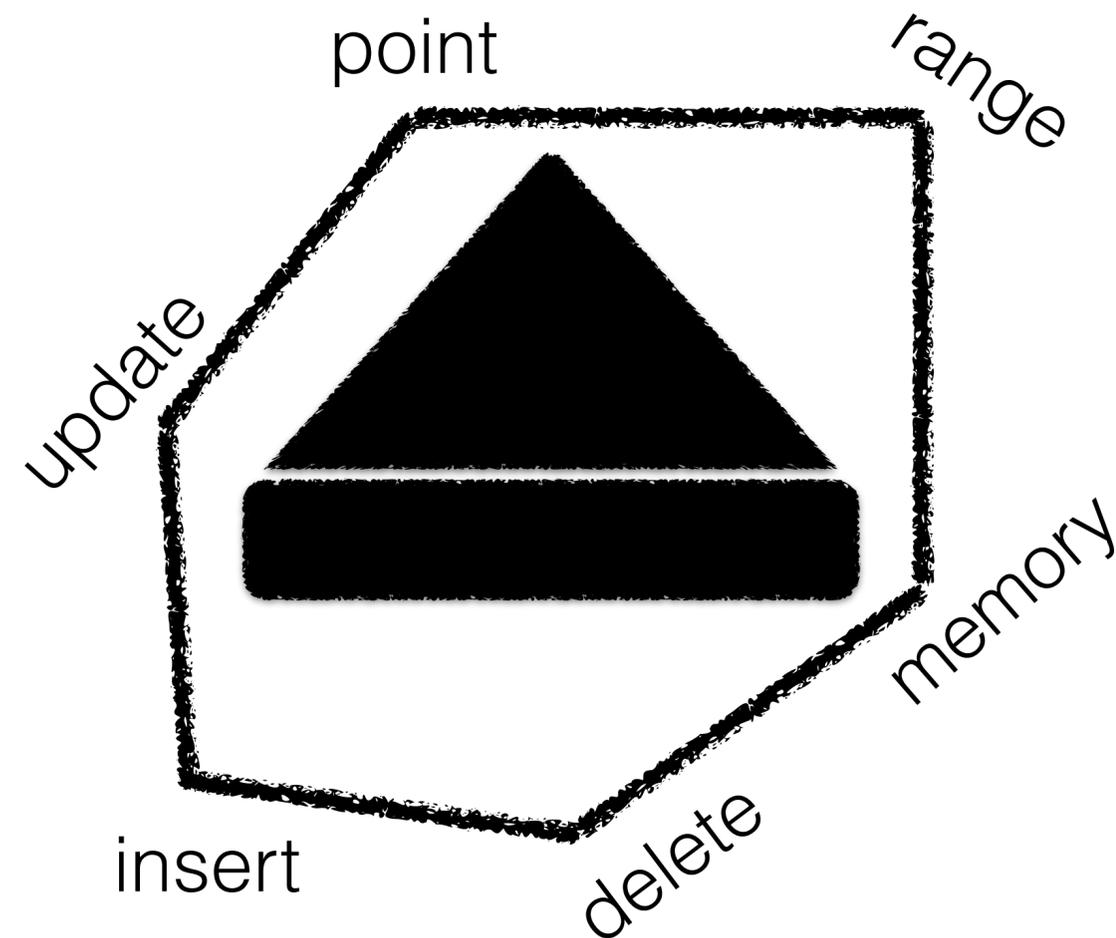
random
access

sorted
search

bloom
probe



synthesize from first principles



algorithm & cost synthesis

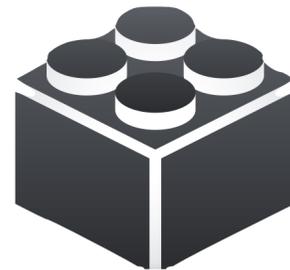
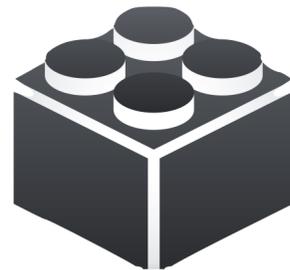
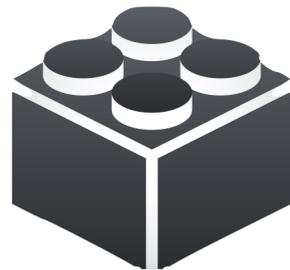
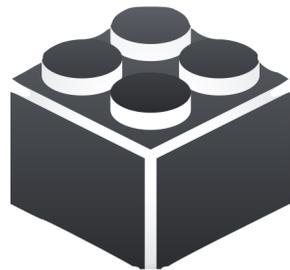
scan

hash
probe

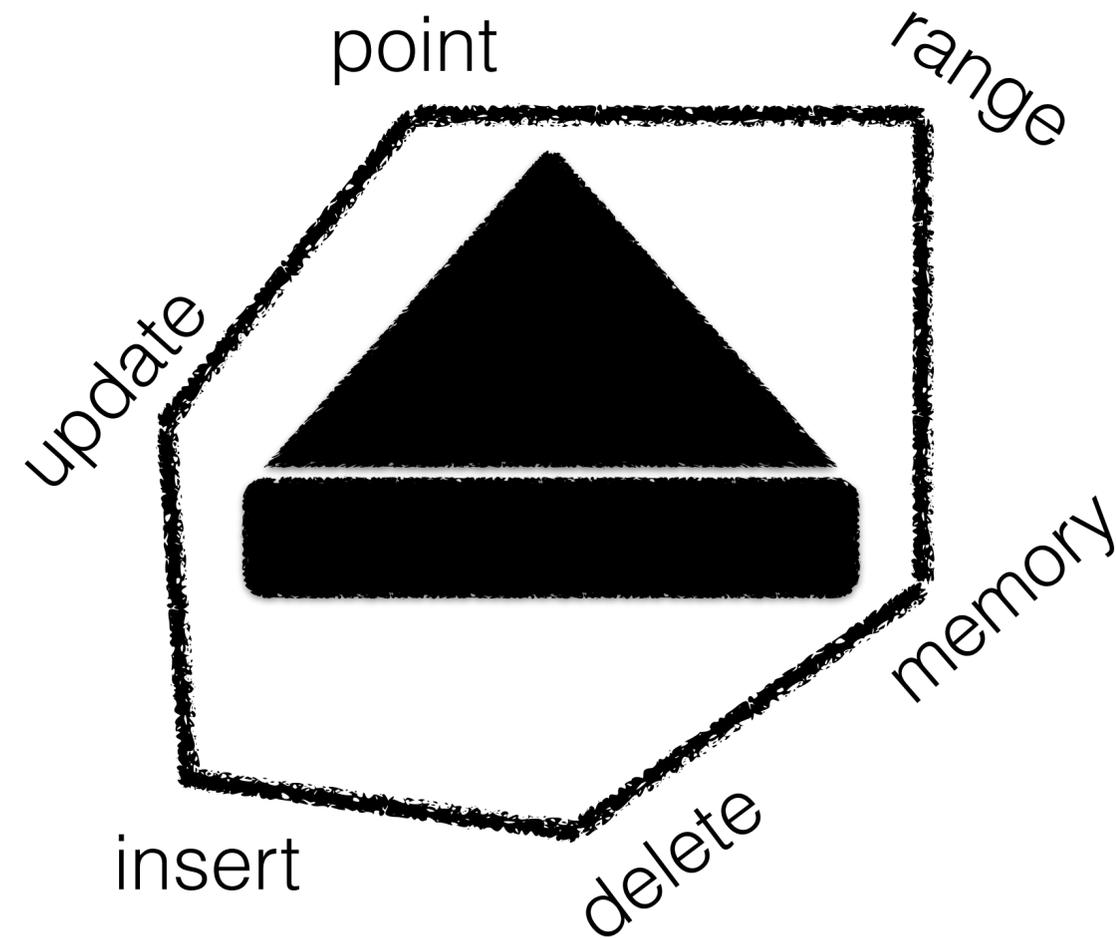
random
access

sorted
search

bloom
probe

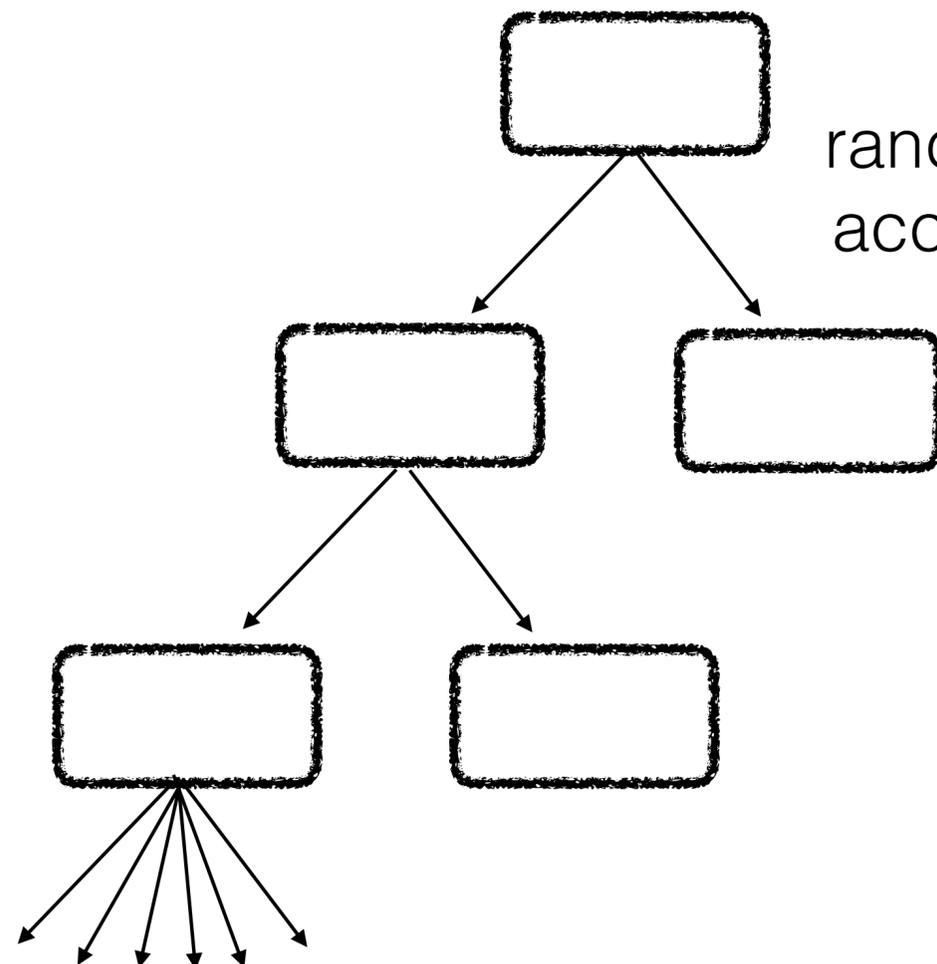


synthesize from first principles



sorted
search

random
access



algorithm & cost synthesis

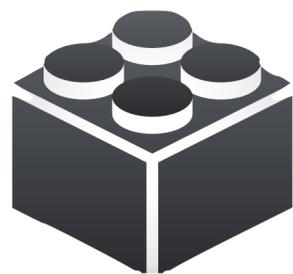
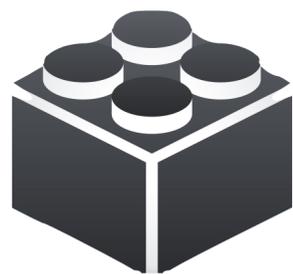
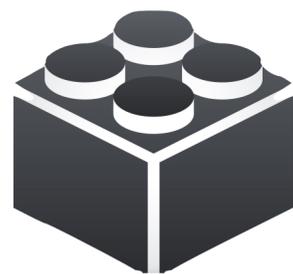
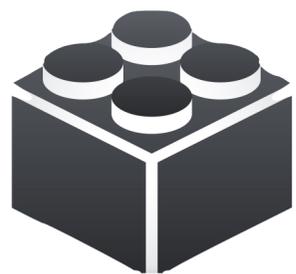
scan

hash probe

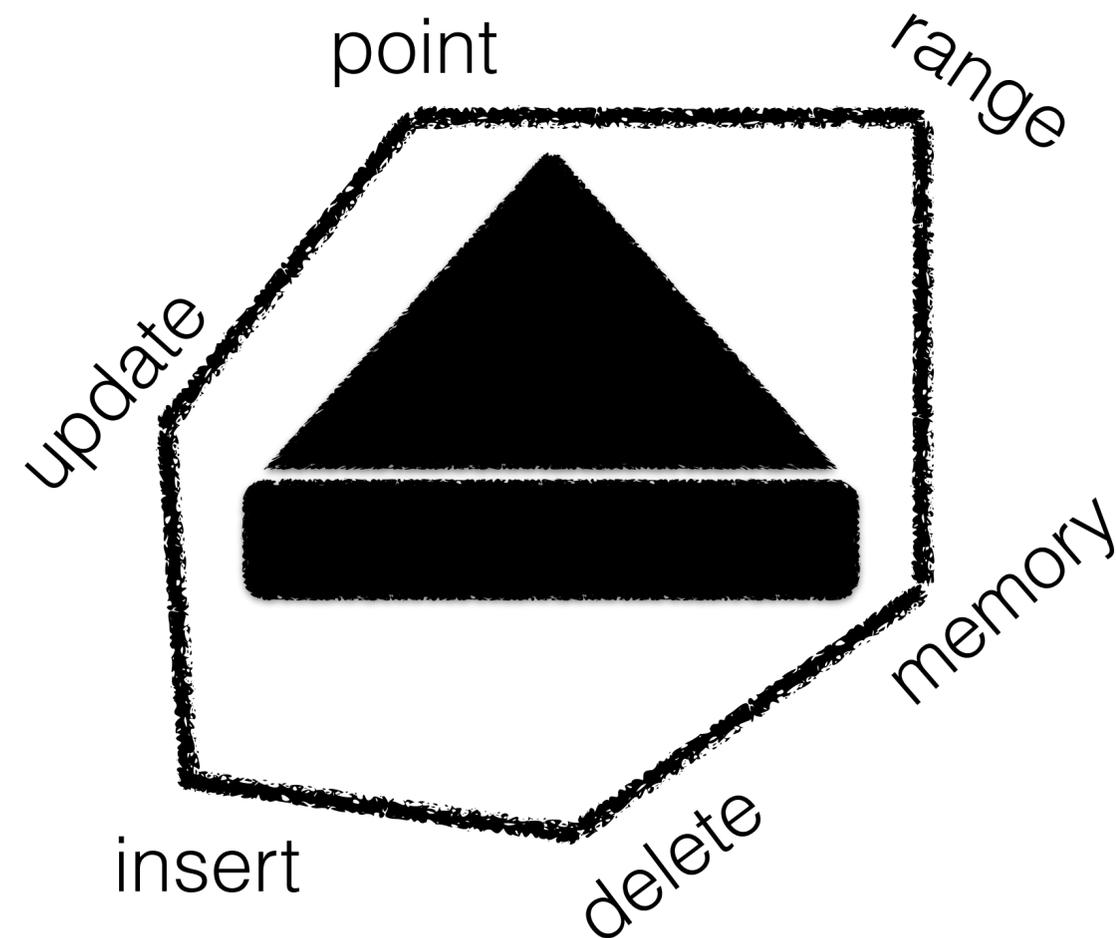
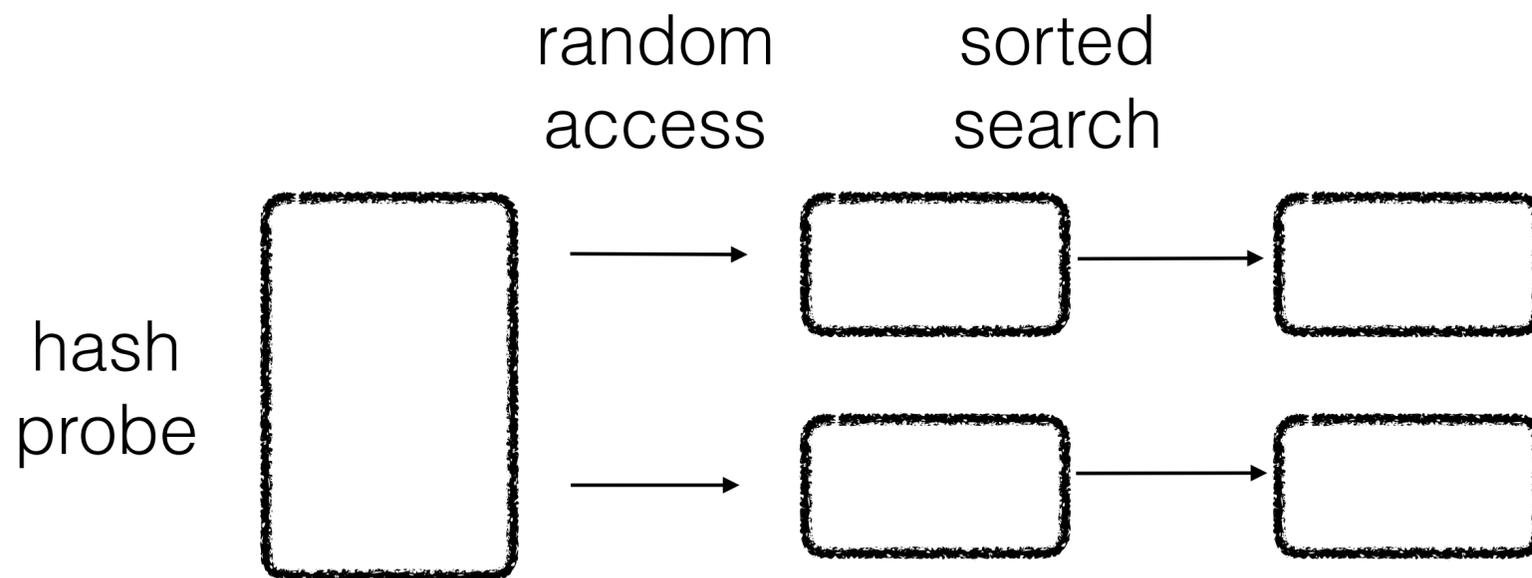
random access

sorted search

bloom probe



synthesize from first principles



algorithm & cost synthesis

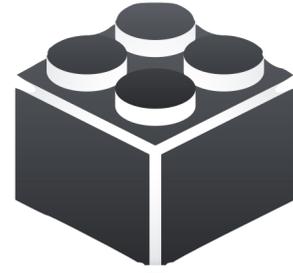
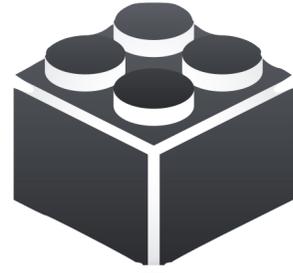
scan

hash
probe

random
access

sorted
search

bloom
probe



synthesize from first principles



1. MINIMAL CODE

e.g., binary search

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



algorithm & cost synthesis

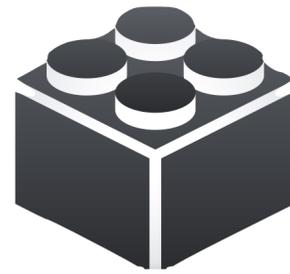
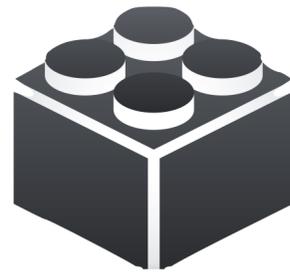
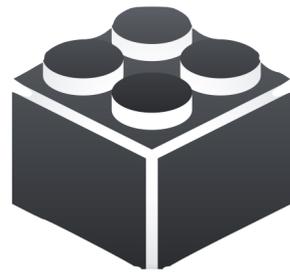
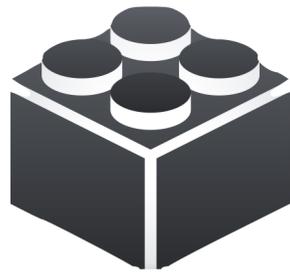
scan

hash probe

random access

sorted search

bloom probe



synthesize from first principles



1. MINIMAL CODE

e.g., binary search

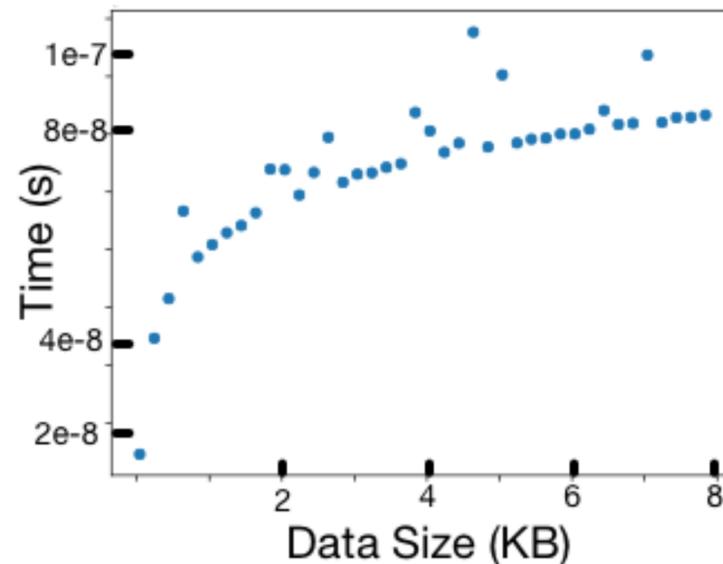
```
if (data[middle] < search_val) {  
    low = middle + 1;  
} else {  
    high = middle;  
}  
middle = (low + high)/2;
```



Run

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

2. BENCHMARK



algorithm & cost synthesis

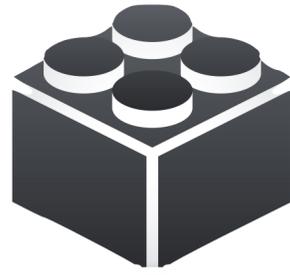
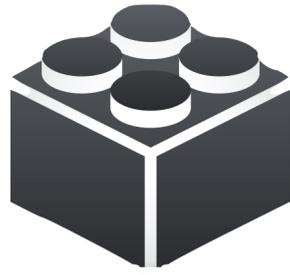
scan

hash probe

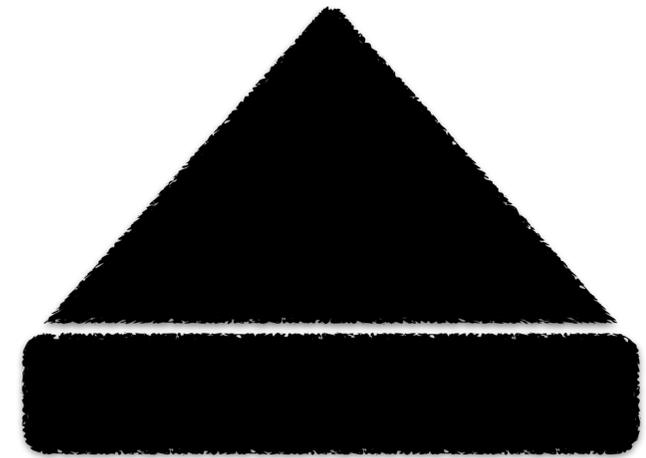
random access

sorted search

bloom probe



synthesize from first principles



1. MINIMAL CODE

e.g., binary search

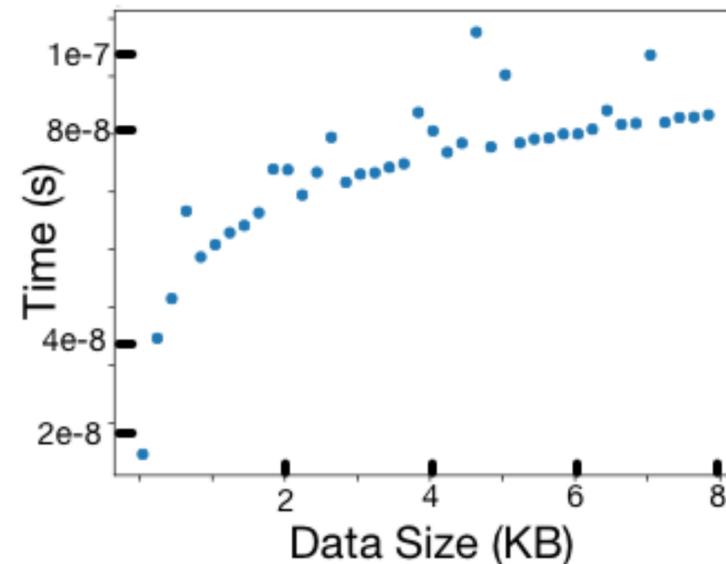
```
if (data[middle] < search_val) {  
  low = middle + 1;  
} else {  
  high = middle;  
}  
middle = (low + high)/2;
```



Run



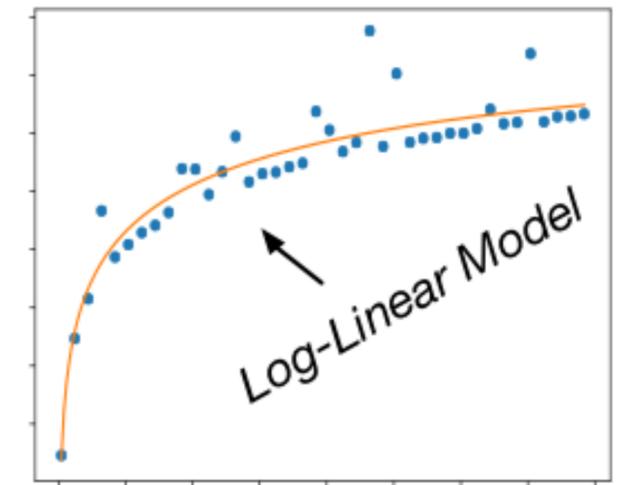
2. BENCHMARK



$f(x)$

Train

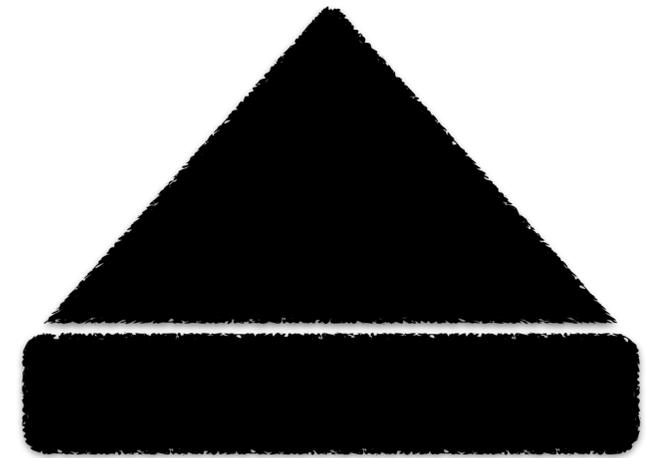
3. FIT MODEL



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

algorithm & cost synthesis

Learned Cost Models @SIGMOD2018



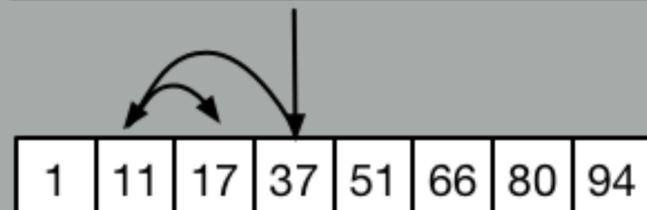
1. MINIMAL CODE

e.g., binary search

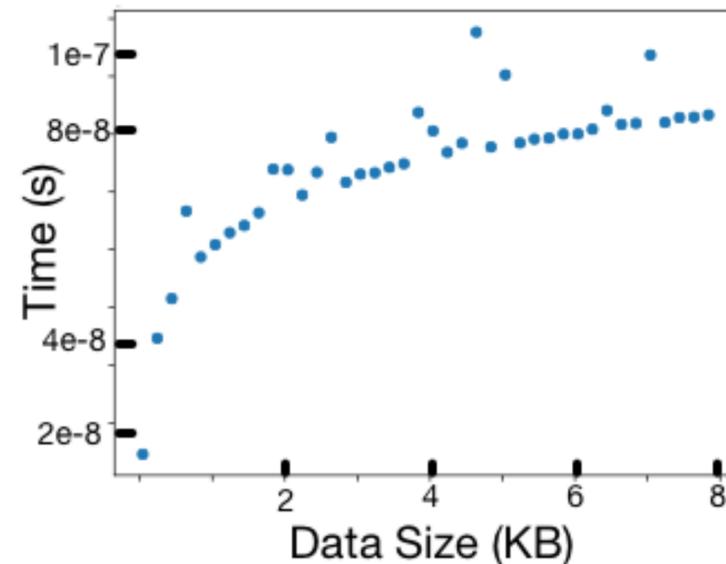
```
if (data[middle] < search_val) {  
    low = middle + 1;  
} else {  
    high = middle;  
}  
middle = (low + high)/2;
```



Run



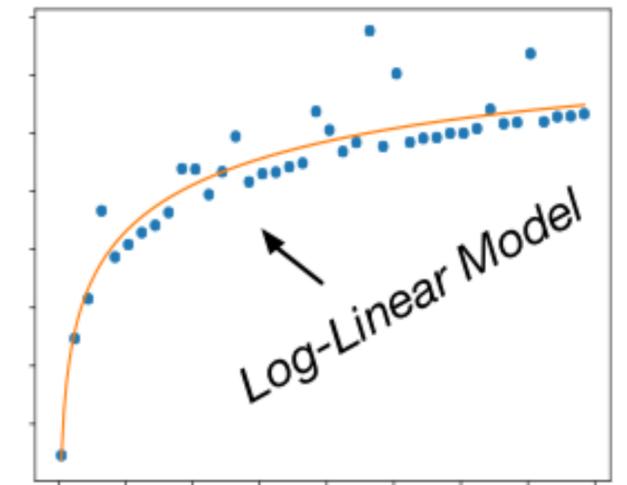
2. BENCHMARK



$f(x)$

Train

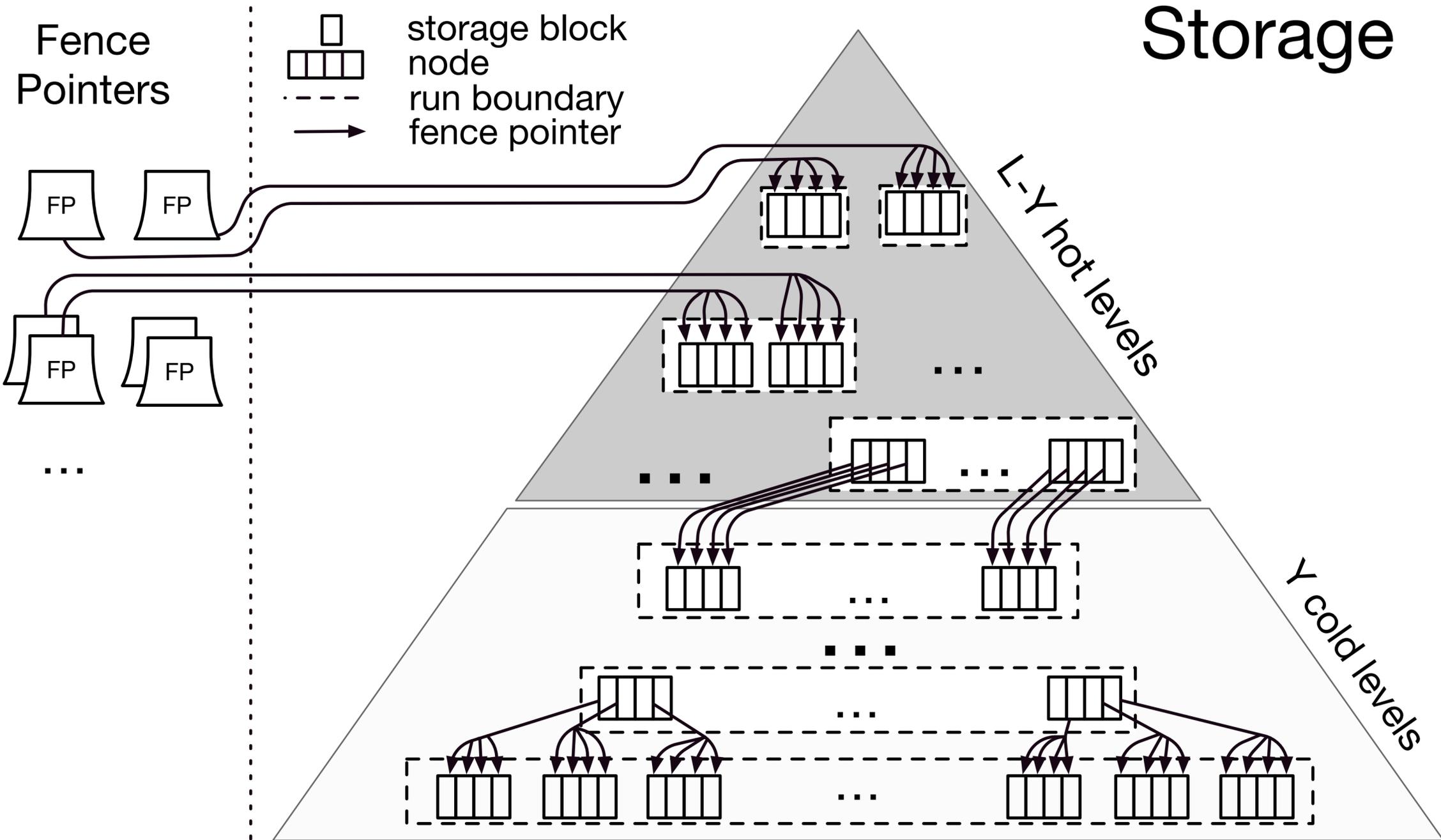
3. FIT MODEL



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

unified design storage engine template

Storage



Buffer

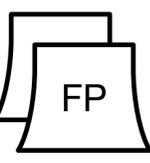
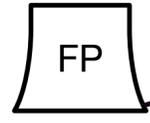


Filters



...

Fence Pointers



...

Bloom Filter

or

Hash Table

(Fewer Bits per Key)

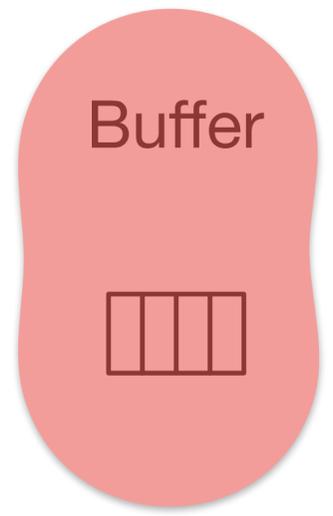
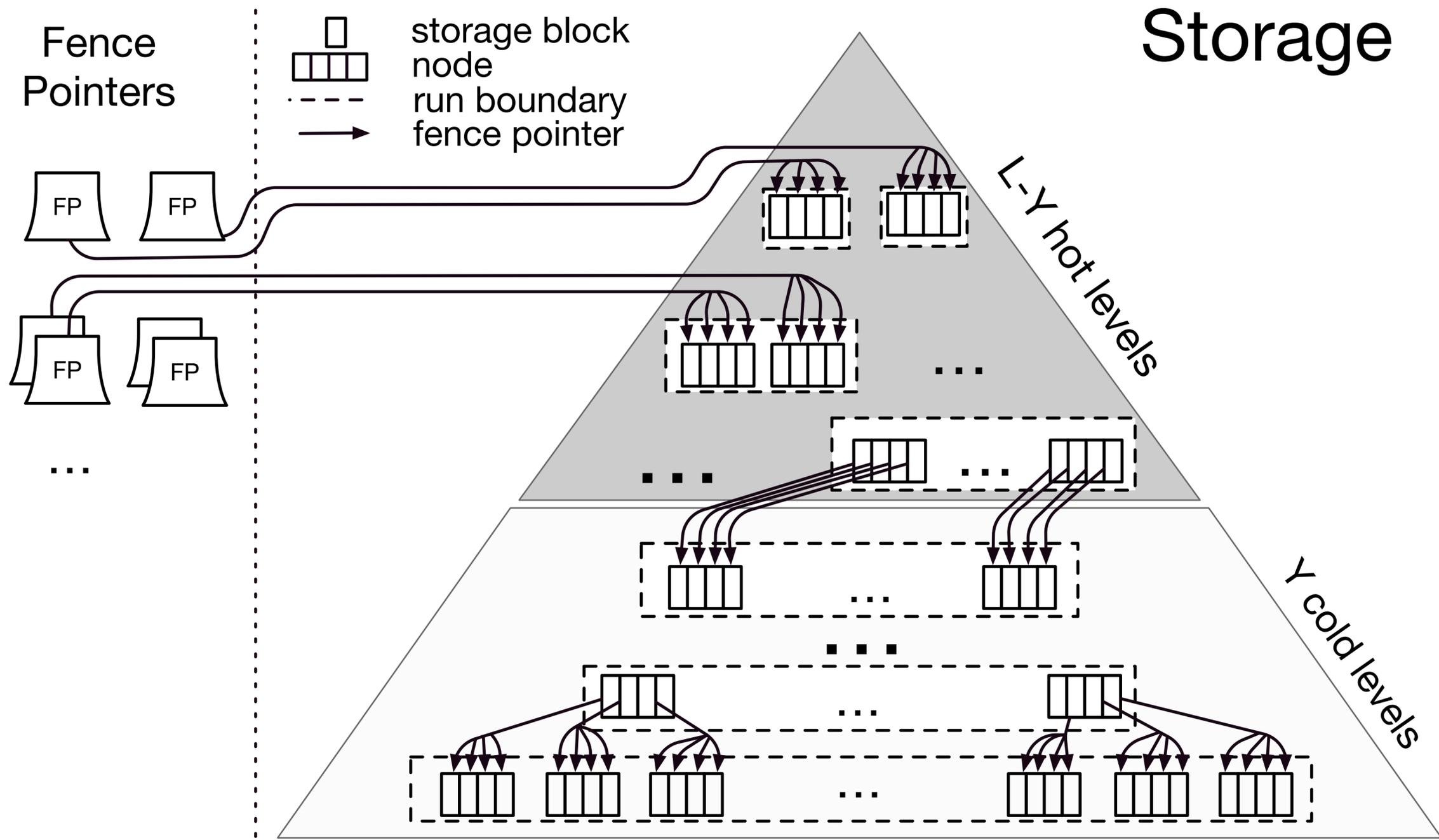
(Full Key Size)

Filter Choice by Mem. Budget

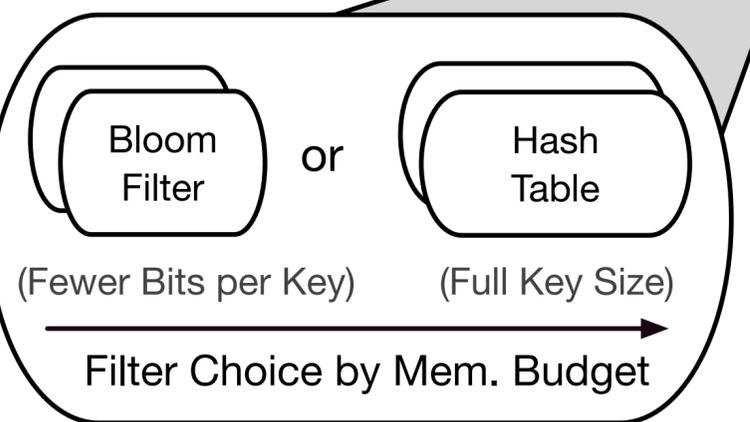
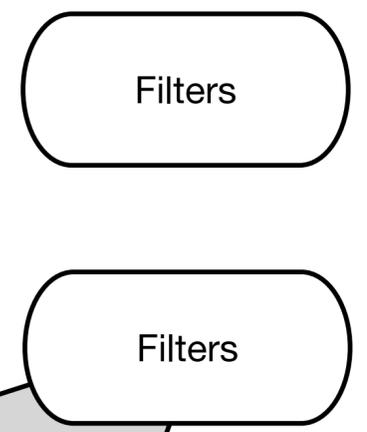
Memory

unified design storage engine template

Storage

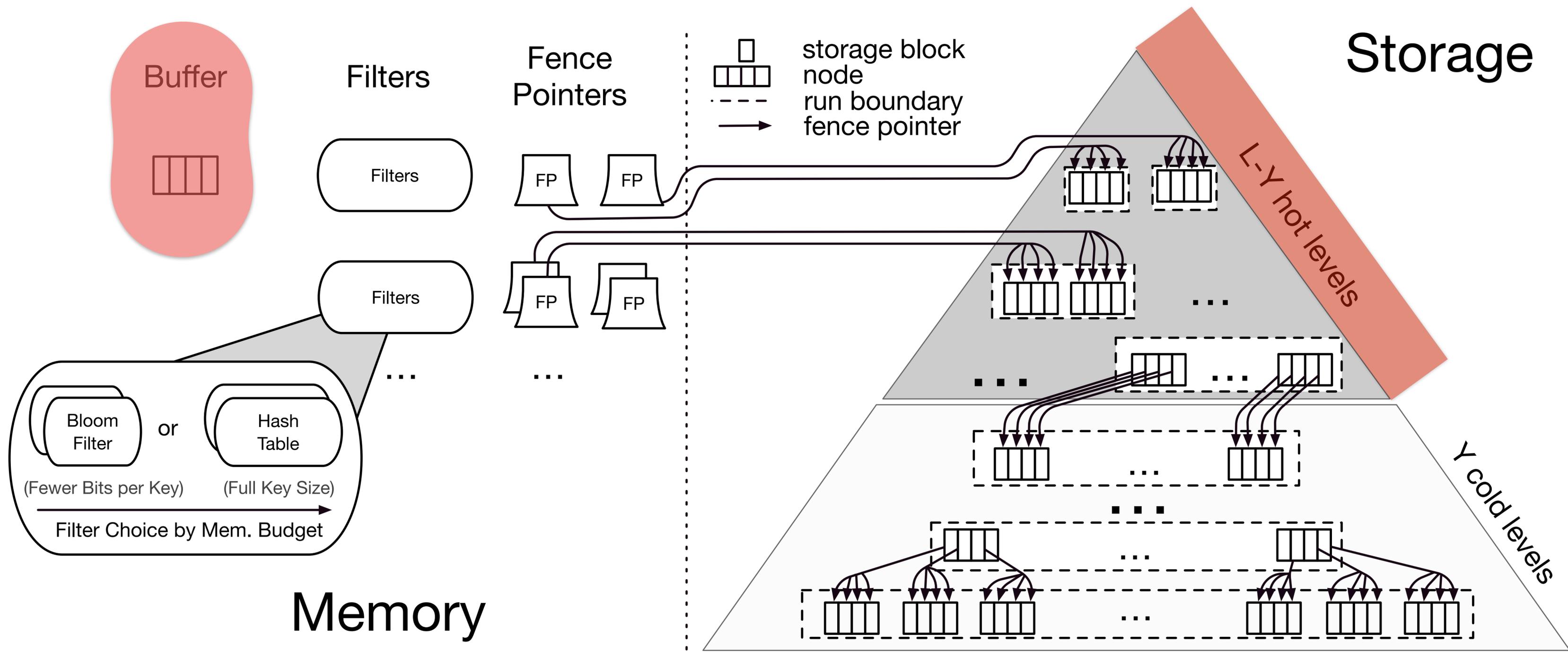


Filters

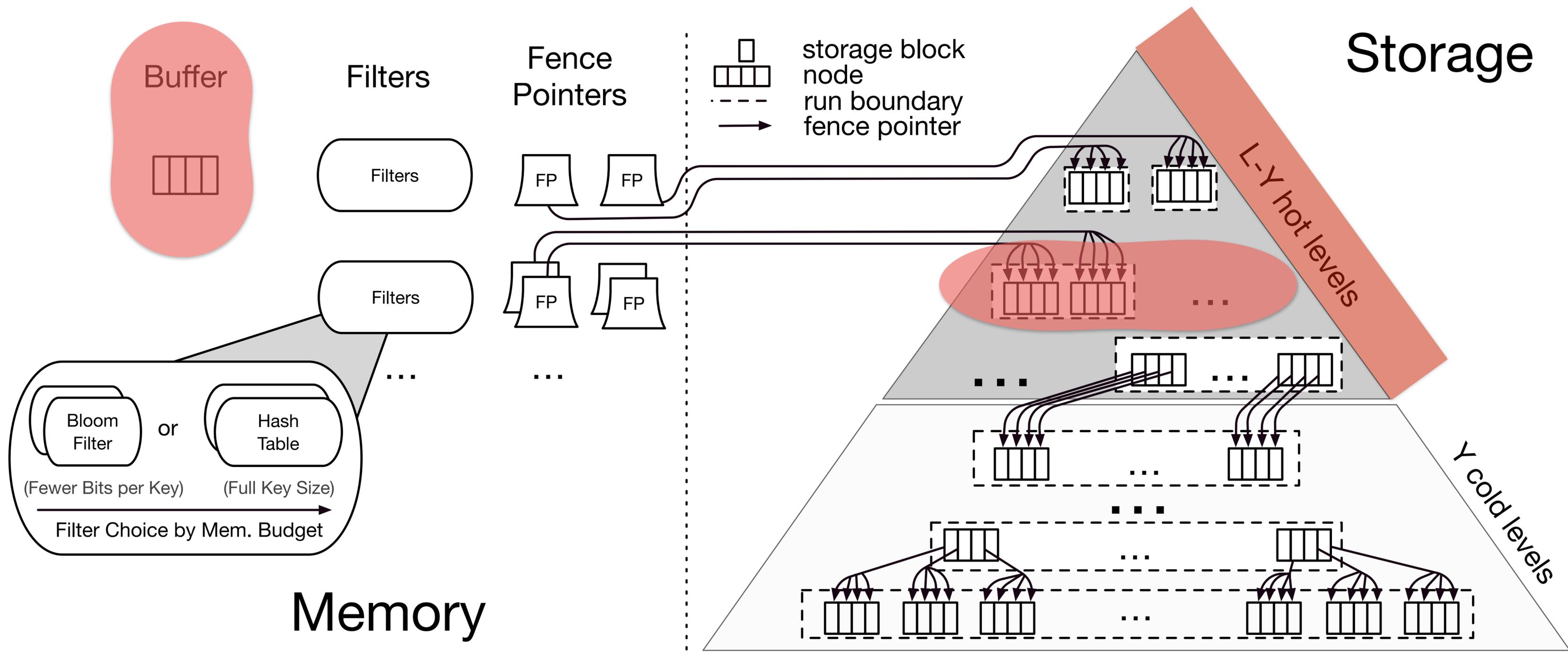


Memory

unified design storage engine template

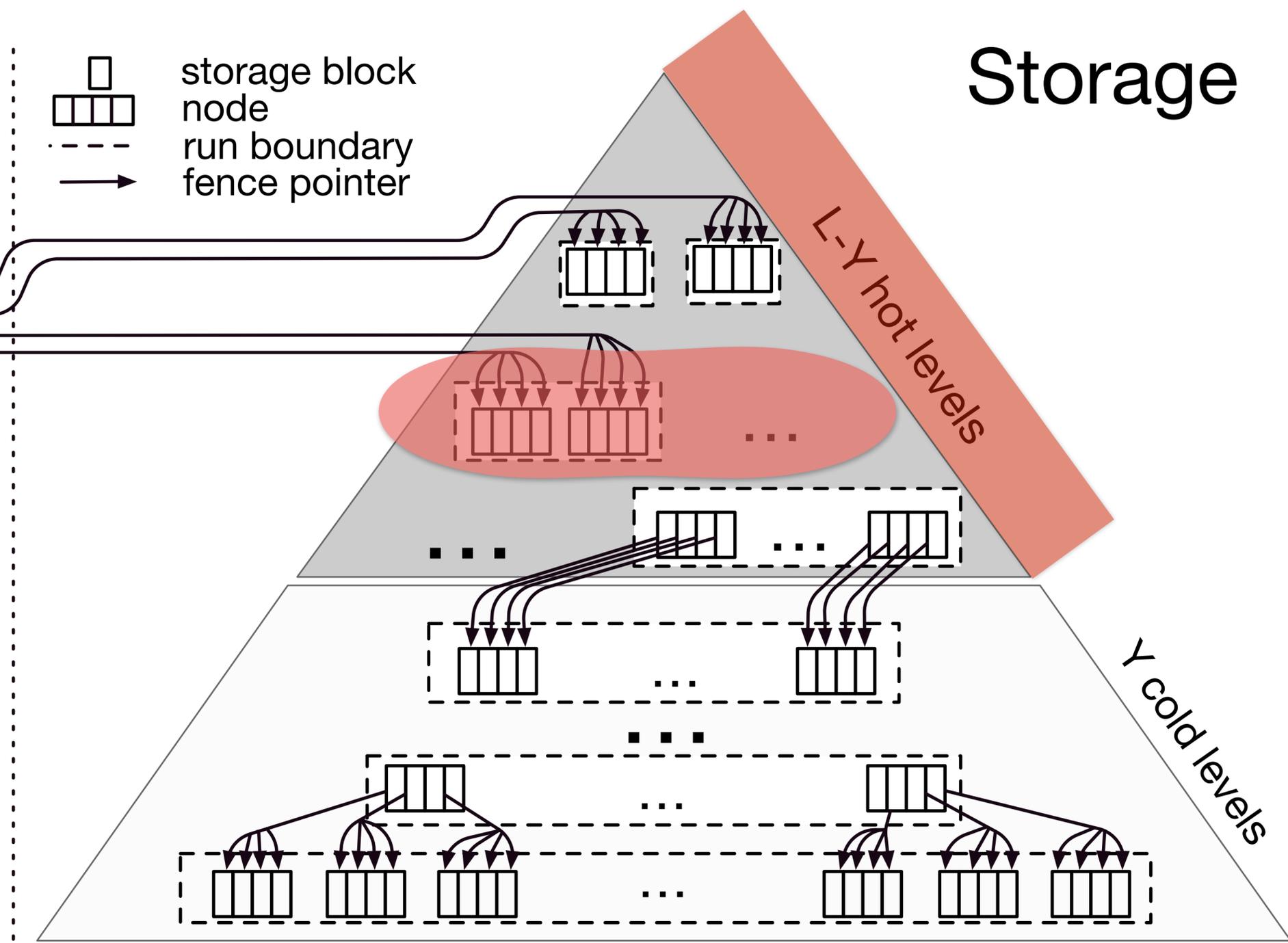


unified design storage engine template



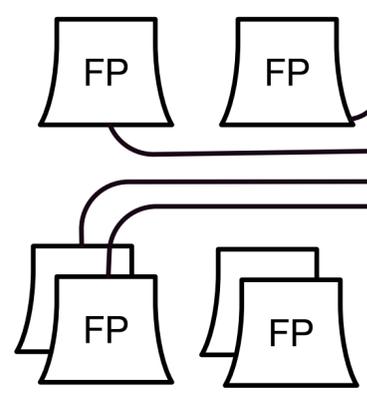
unified design storage engine template

Storage

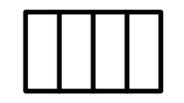


storage block
node
run boundary
fence pointer

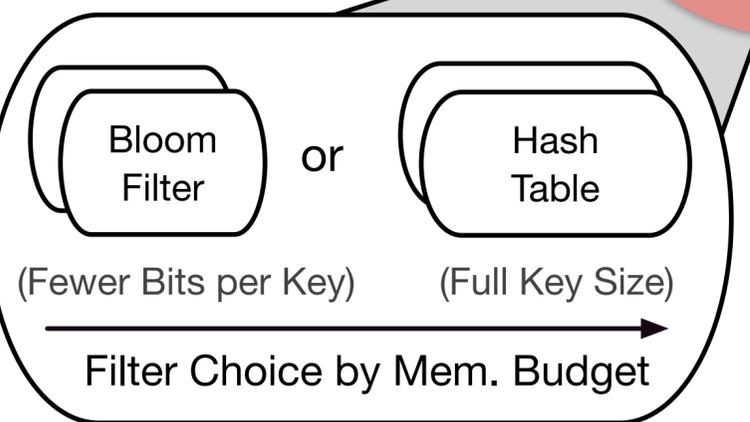
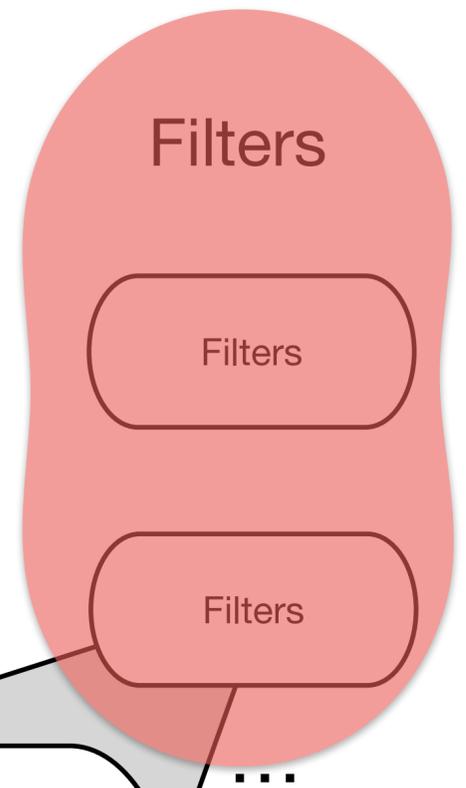
Fence
Pointers



Buffer



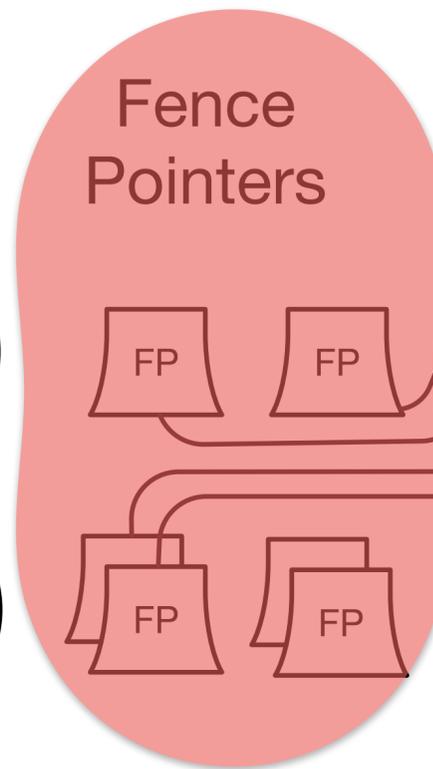
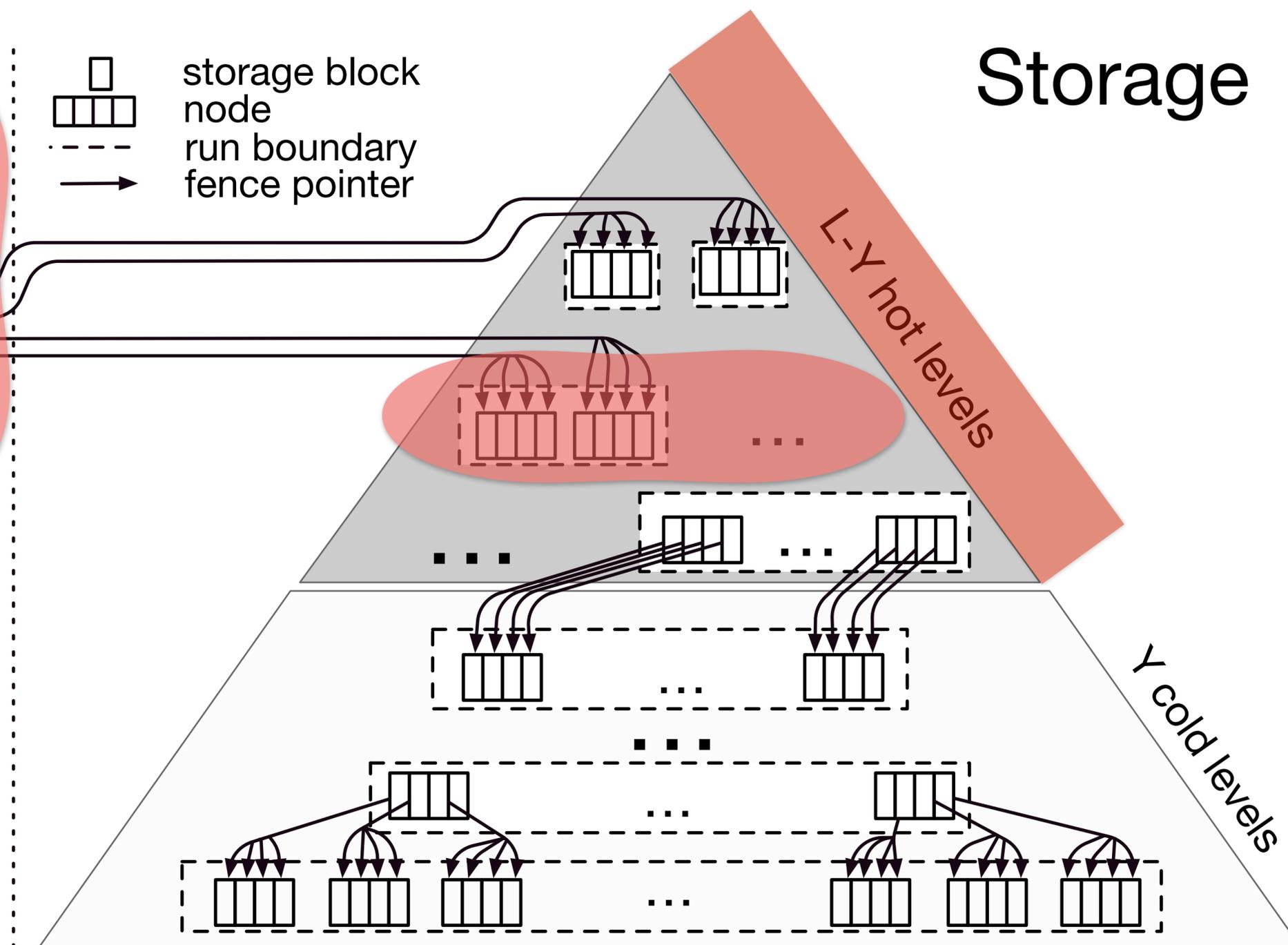
Filters



Memory

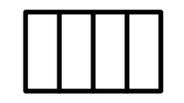
unified design storage engine template

Storage

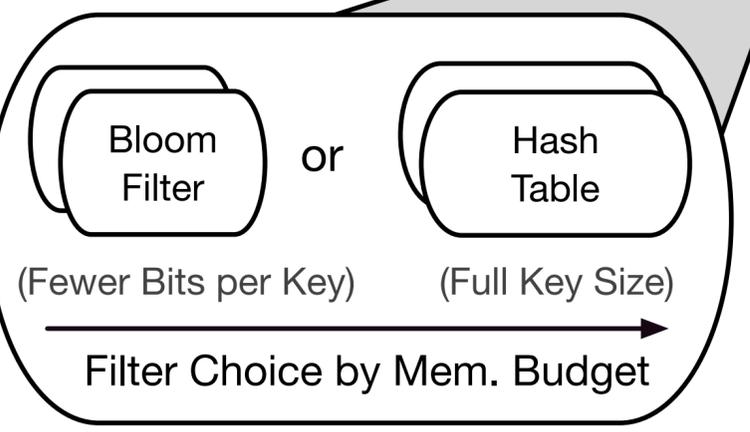


□ storage block node
- - - run boundary
→ fence pointer

Buffer



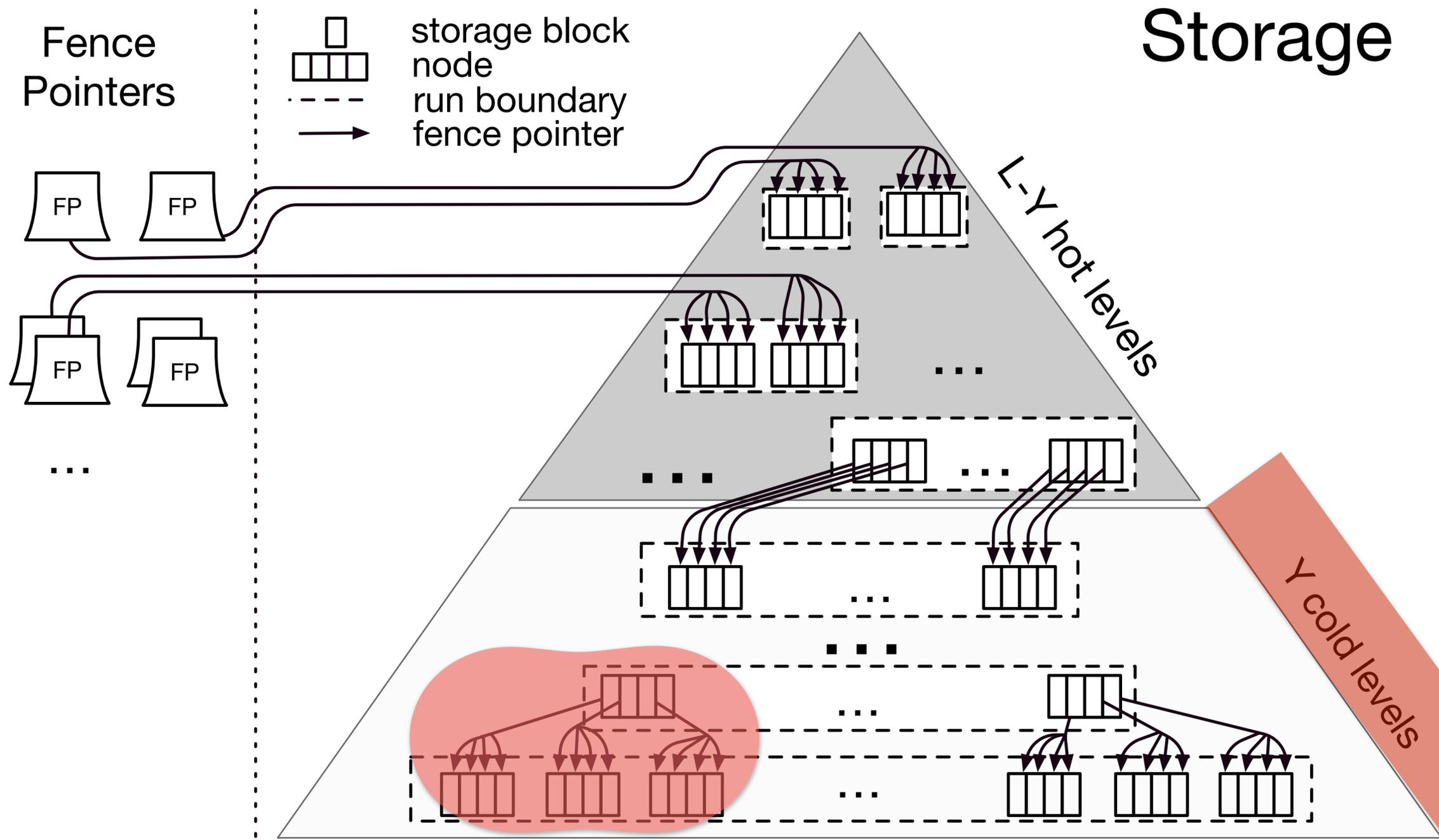
Filters



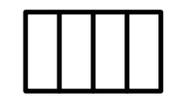
Memory

unified design storage engine template

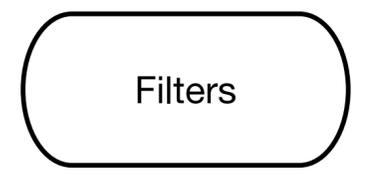
Storage



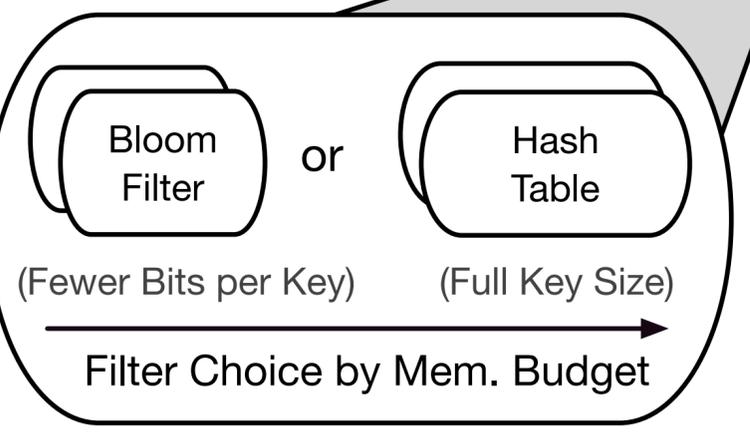
Buffer



Filters



...

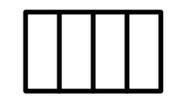


Memory

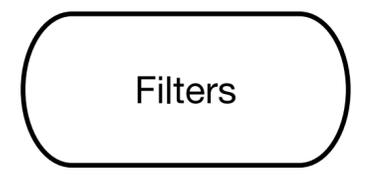
unified design storage engine template

Storage

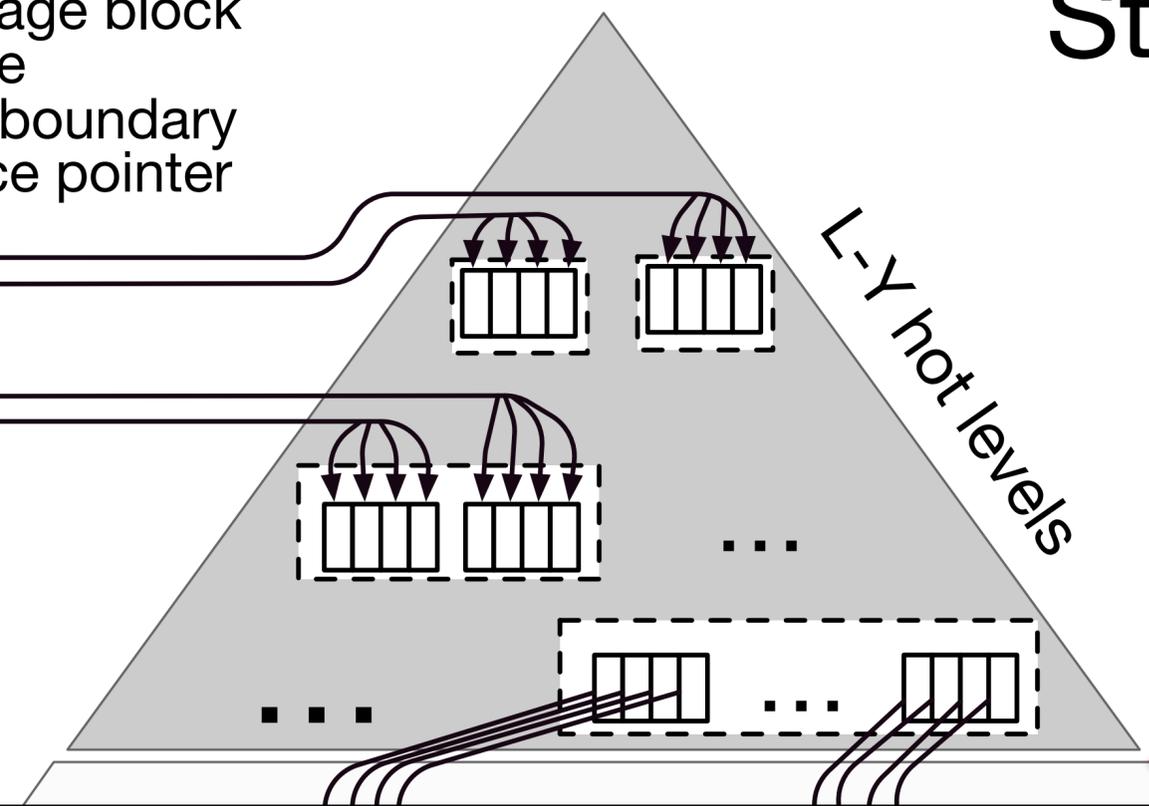
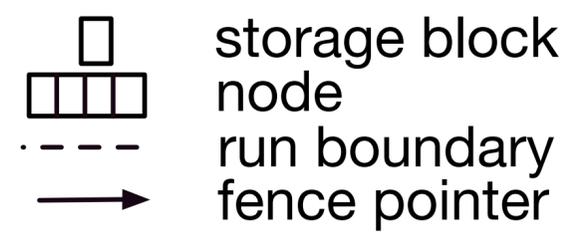
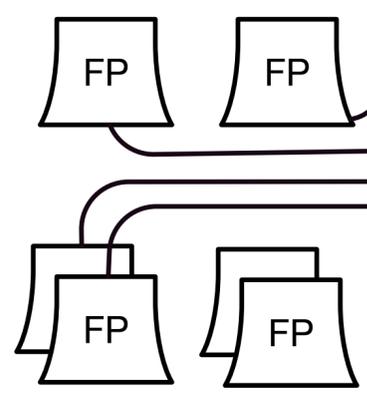
Buffer



Filters



Fence
Pointers

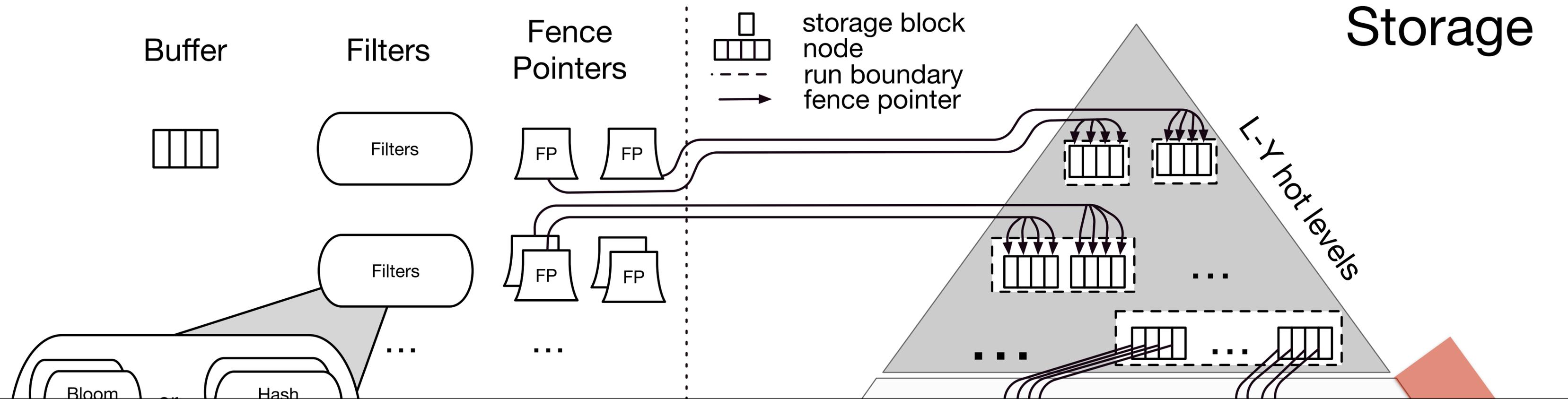


**analytical
I/O model**

unified
closed form

unified design storage engine template

Storage



**analytical
I/O model**

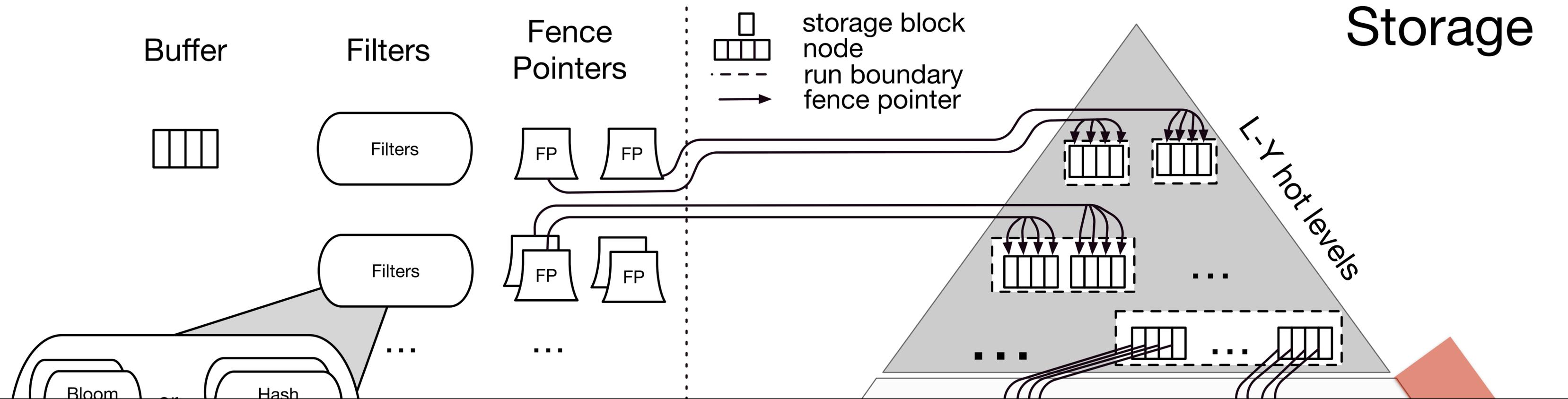
**learned
CPU model**

unified
closed form

h/w,
parallelism

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**

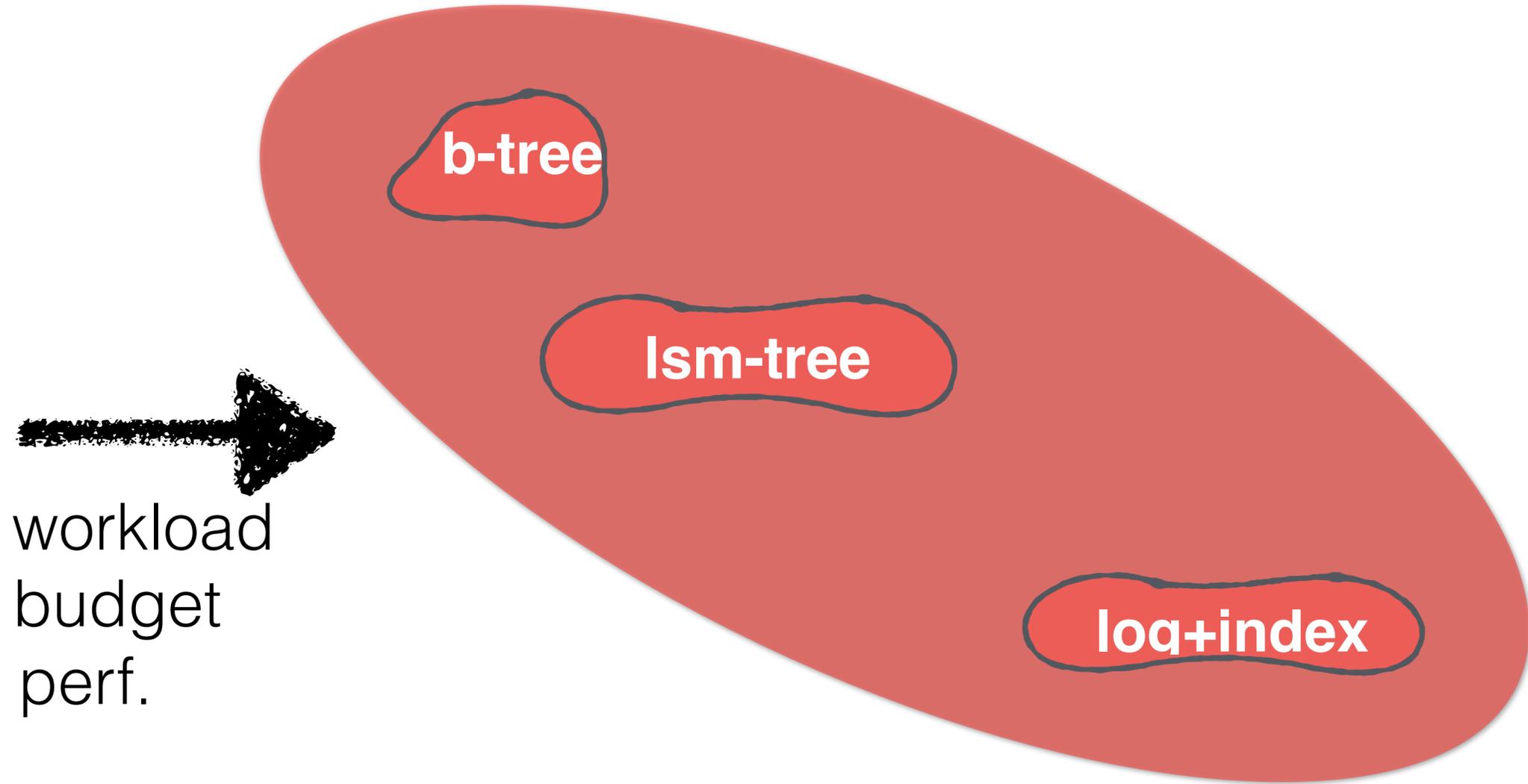
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

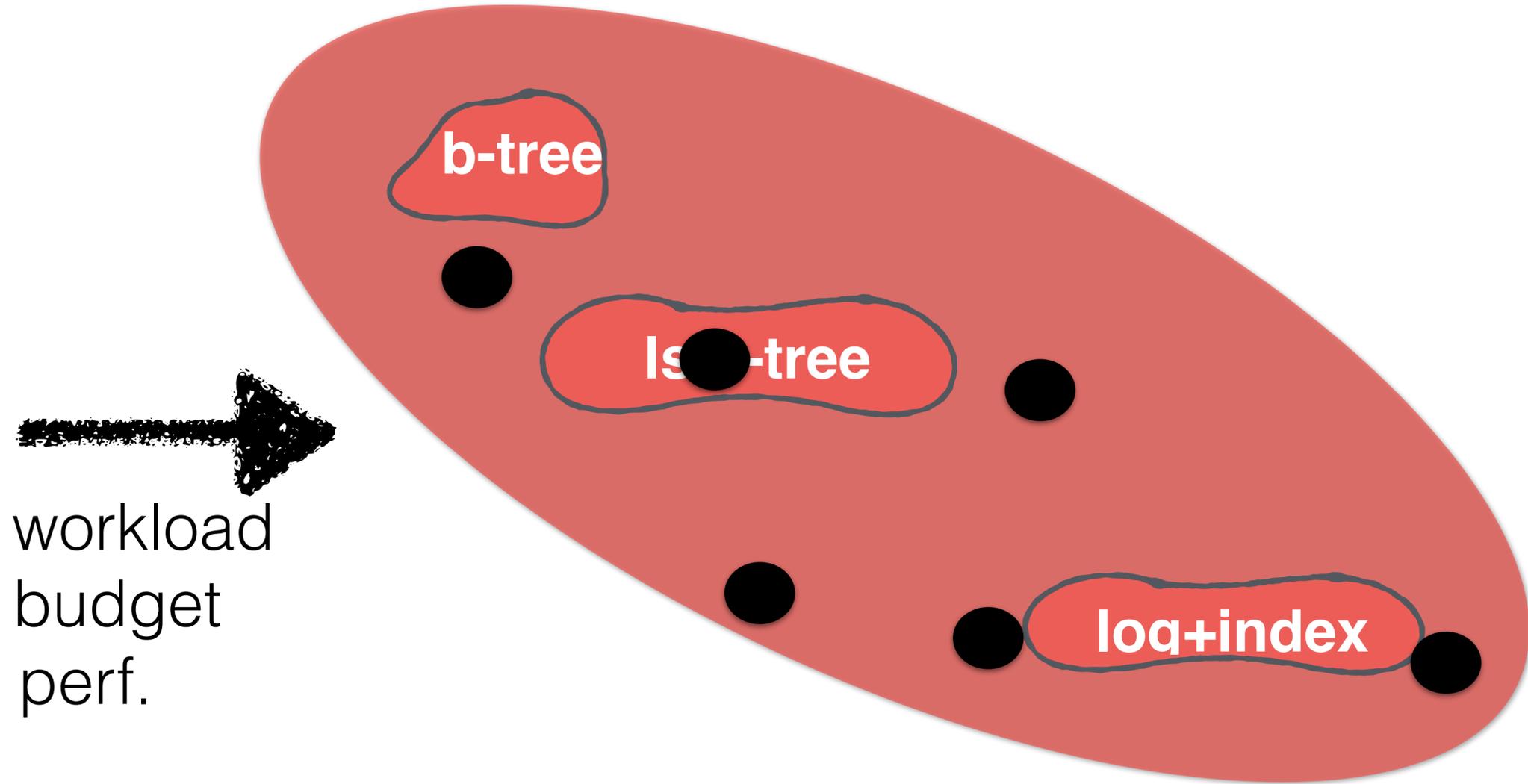
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

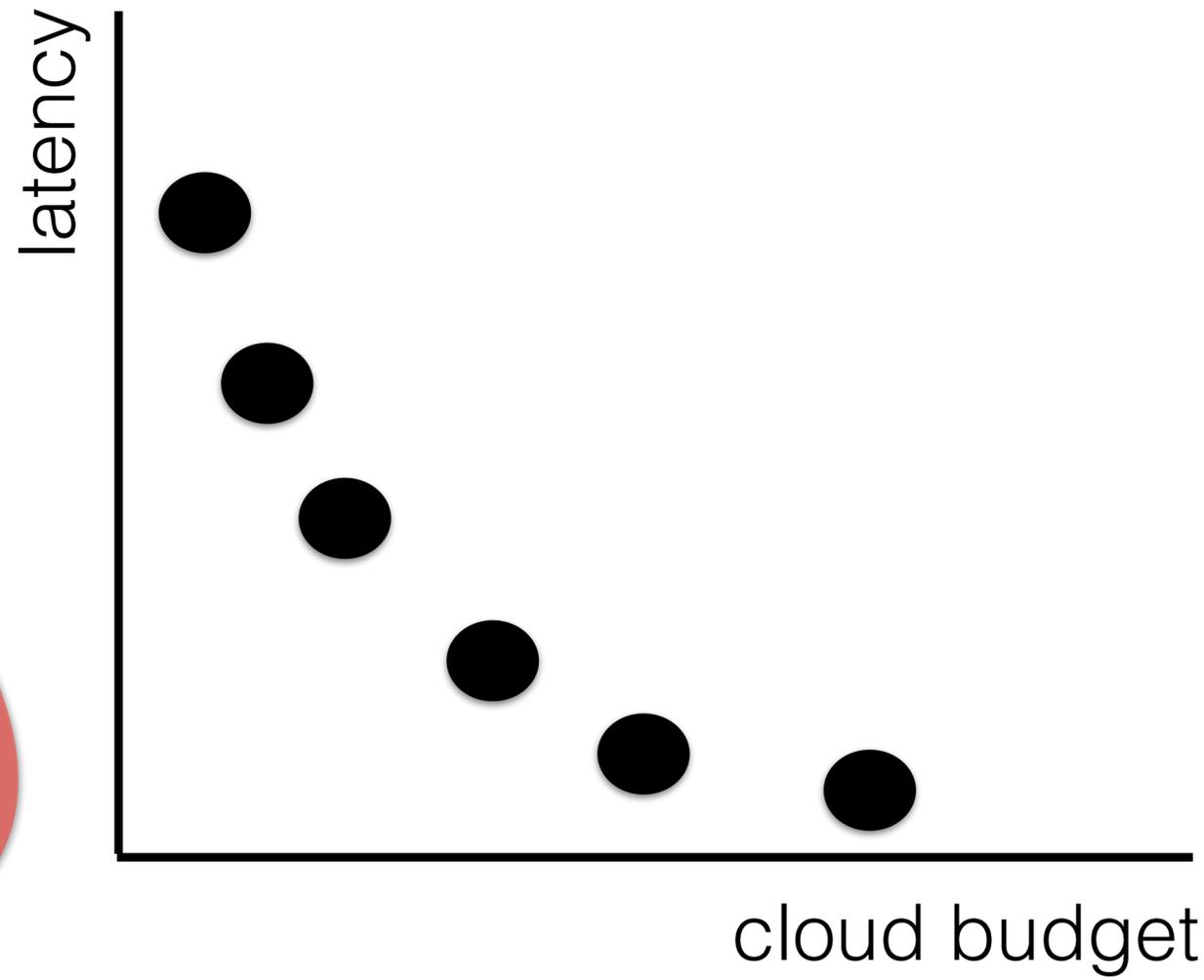
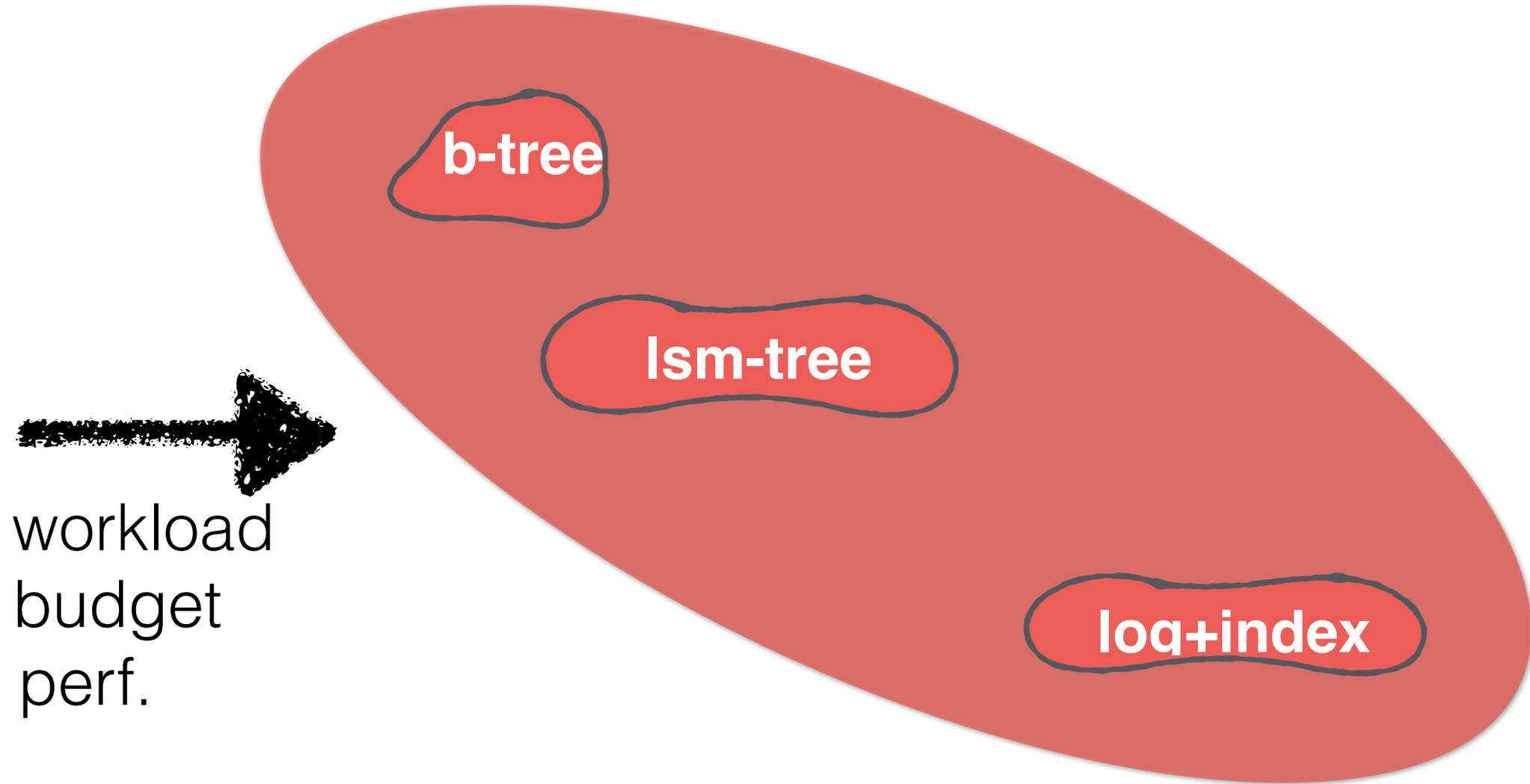
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

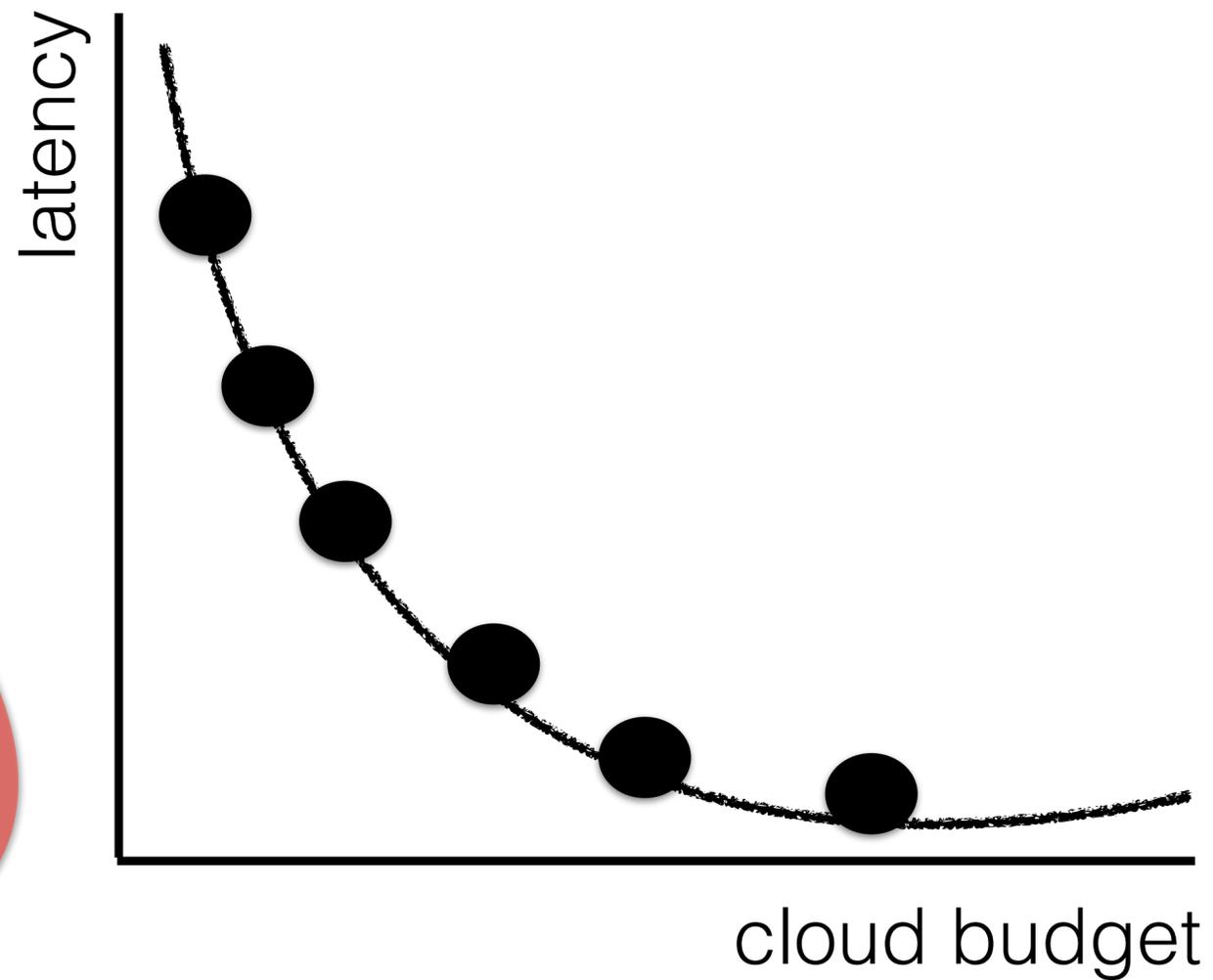
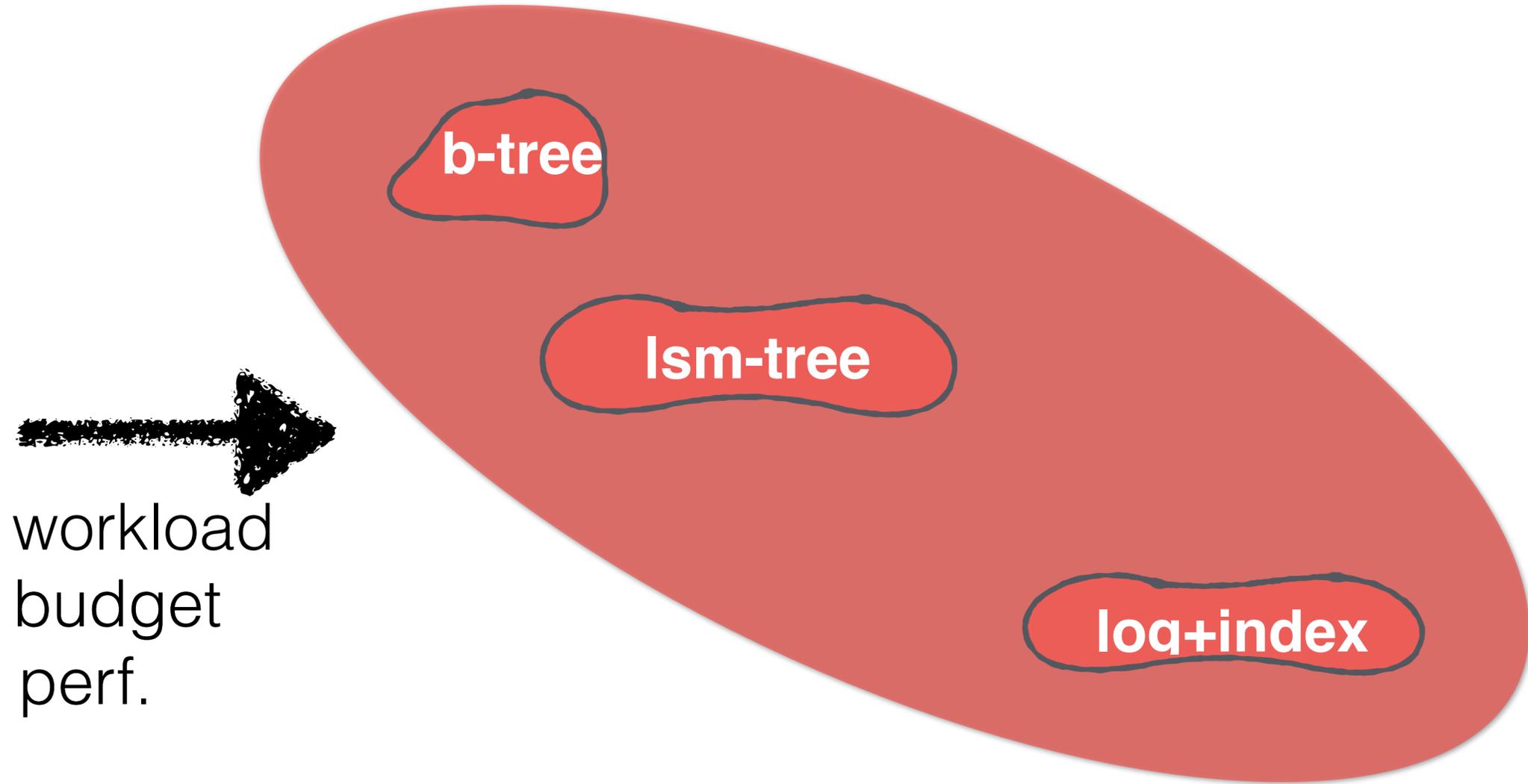
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



**analytical
I/O model**

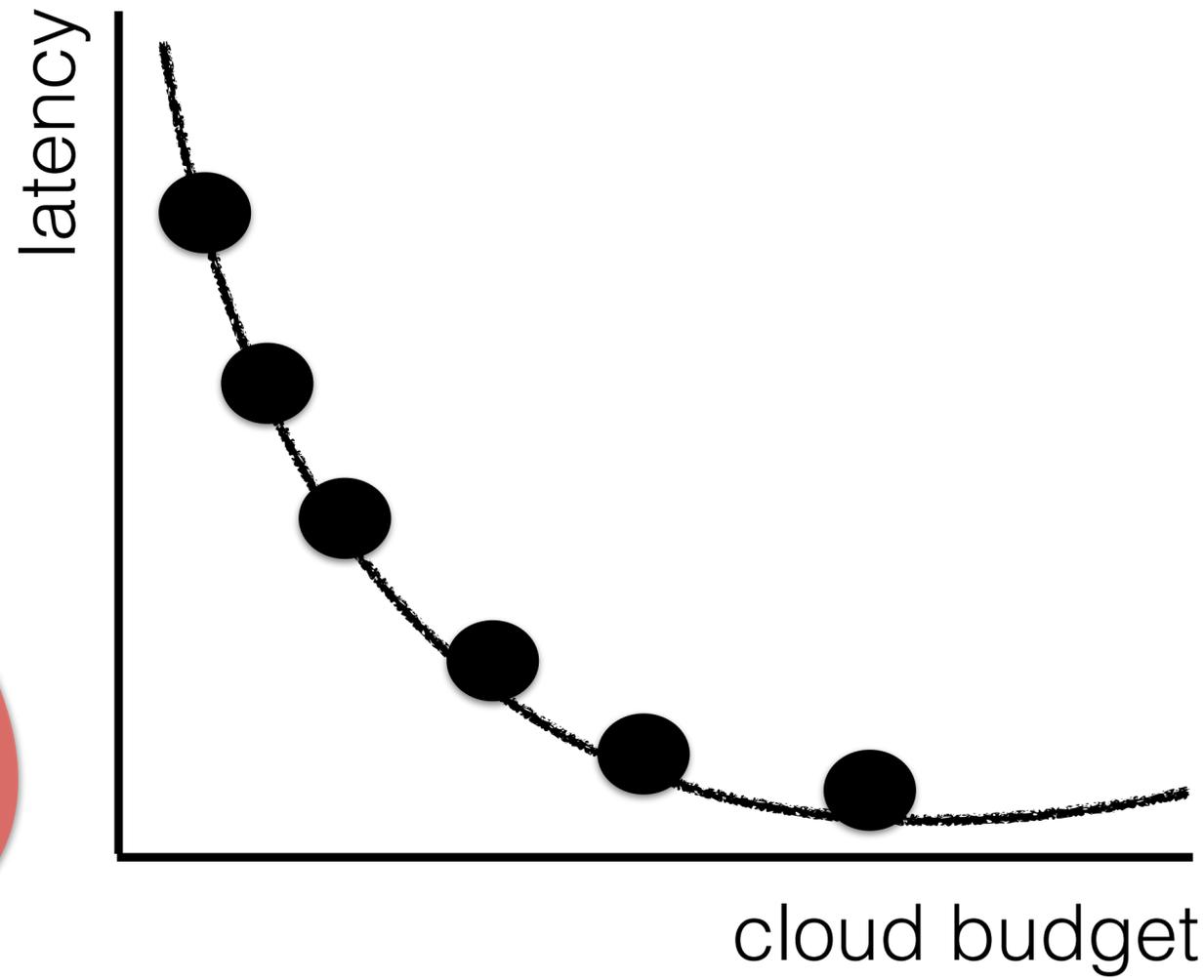
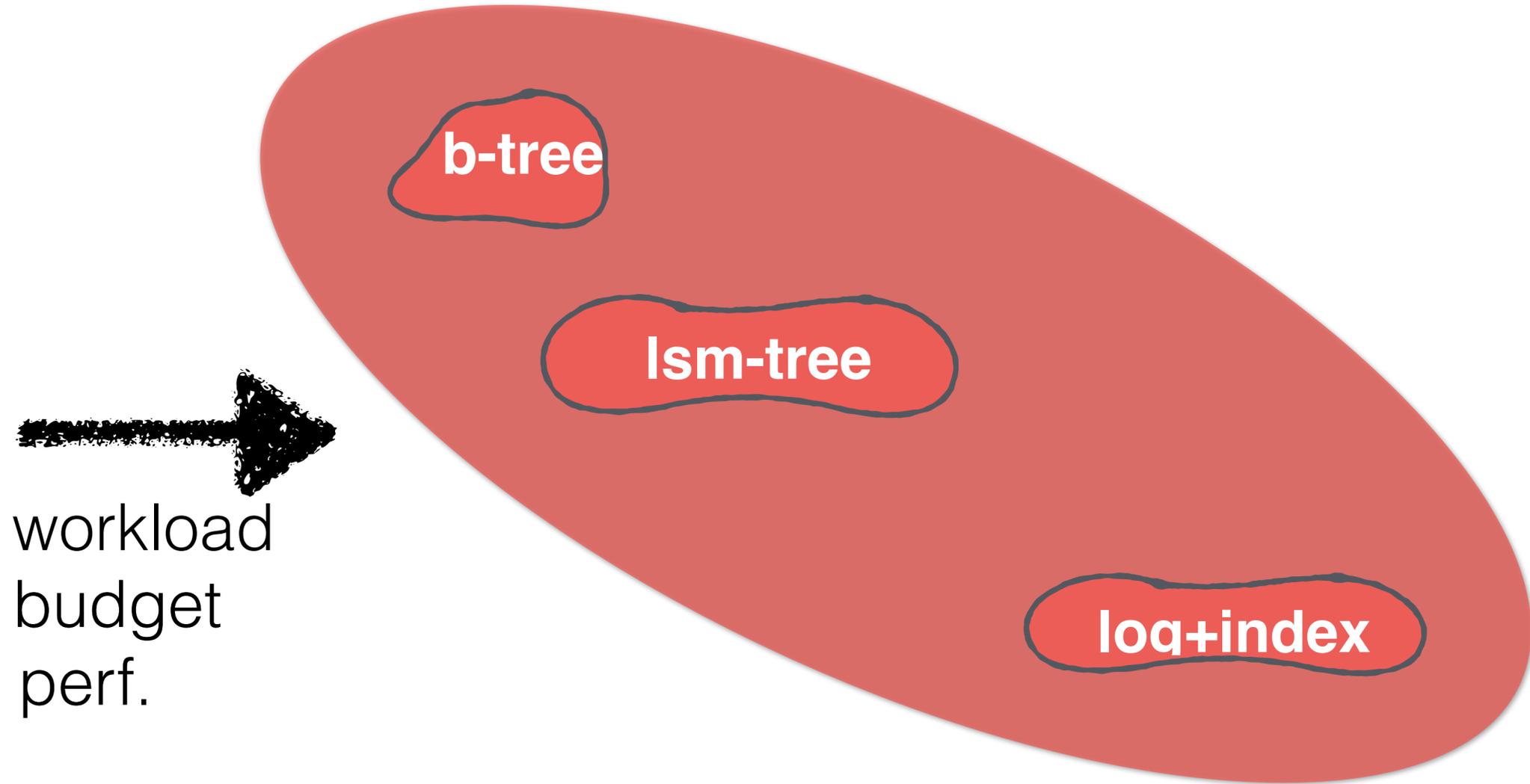
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

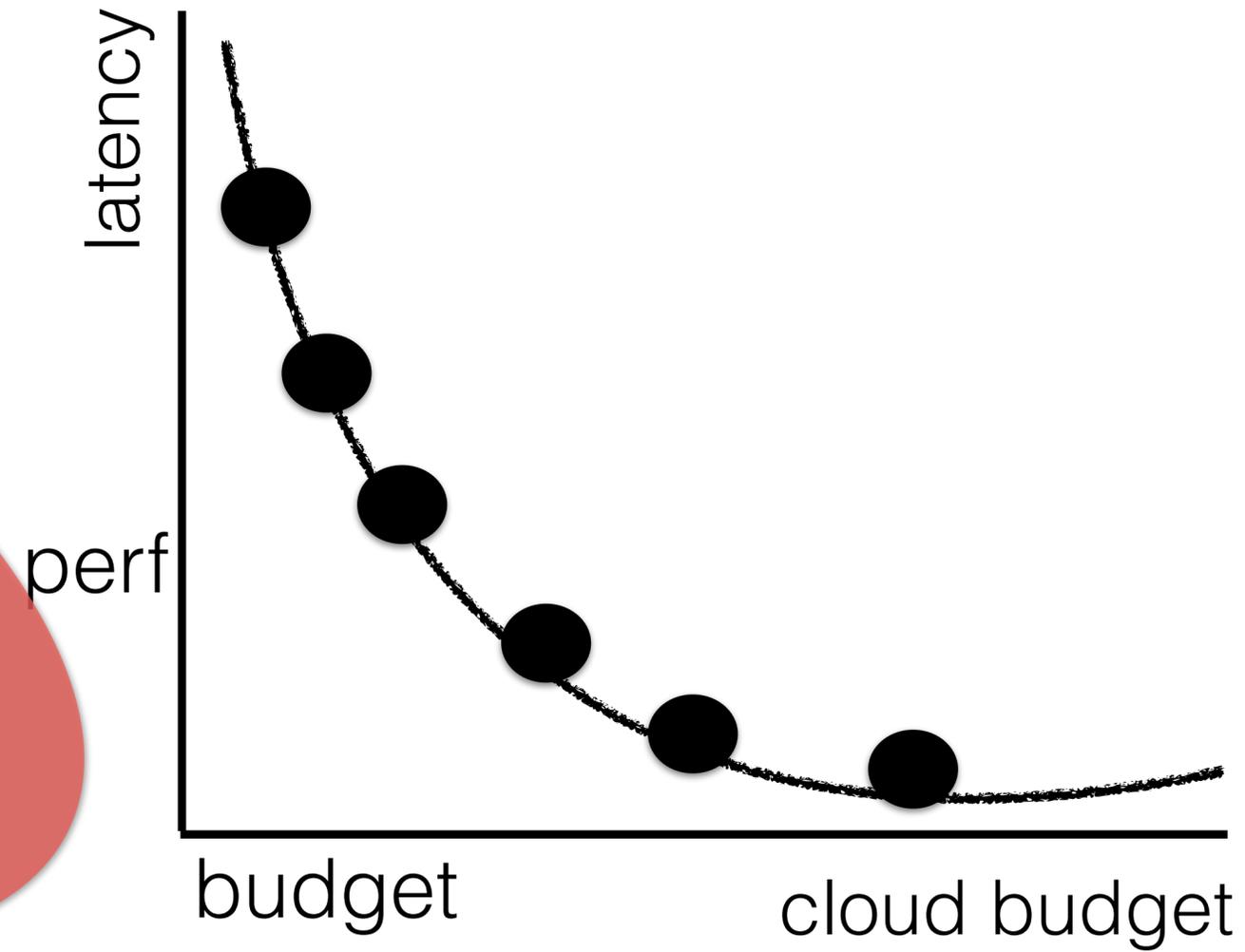
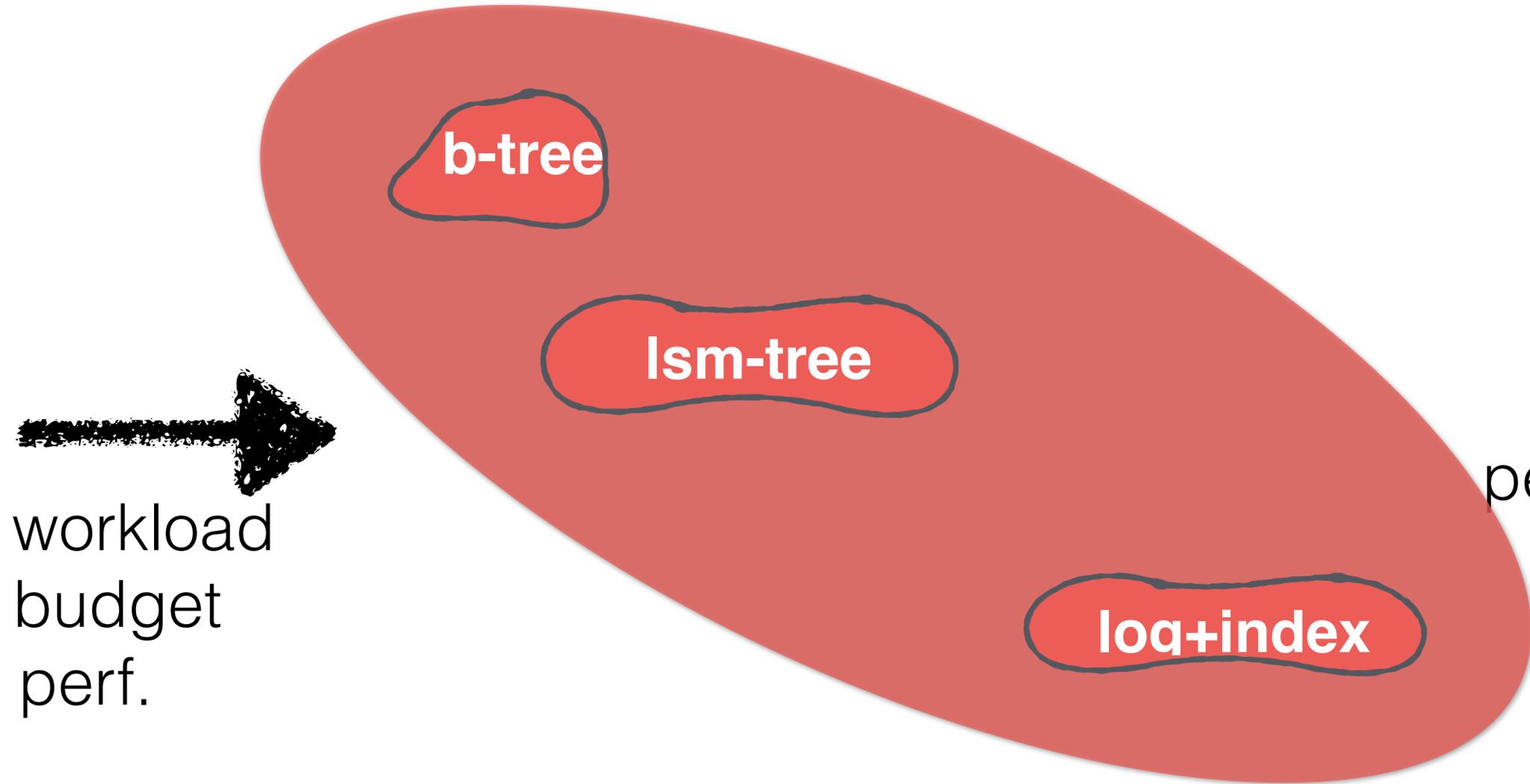
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

AWS, Azure, Google



**analytical
I/O model**

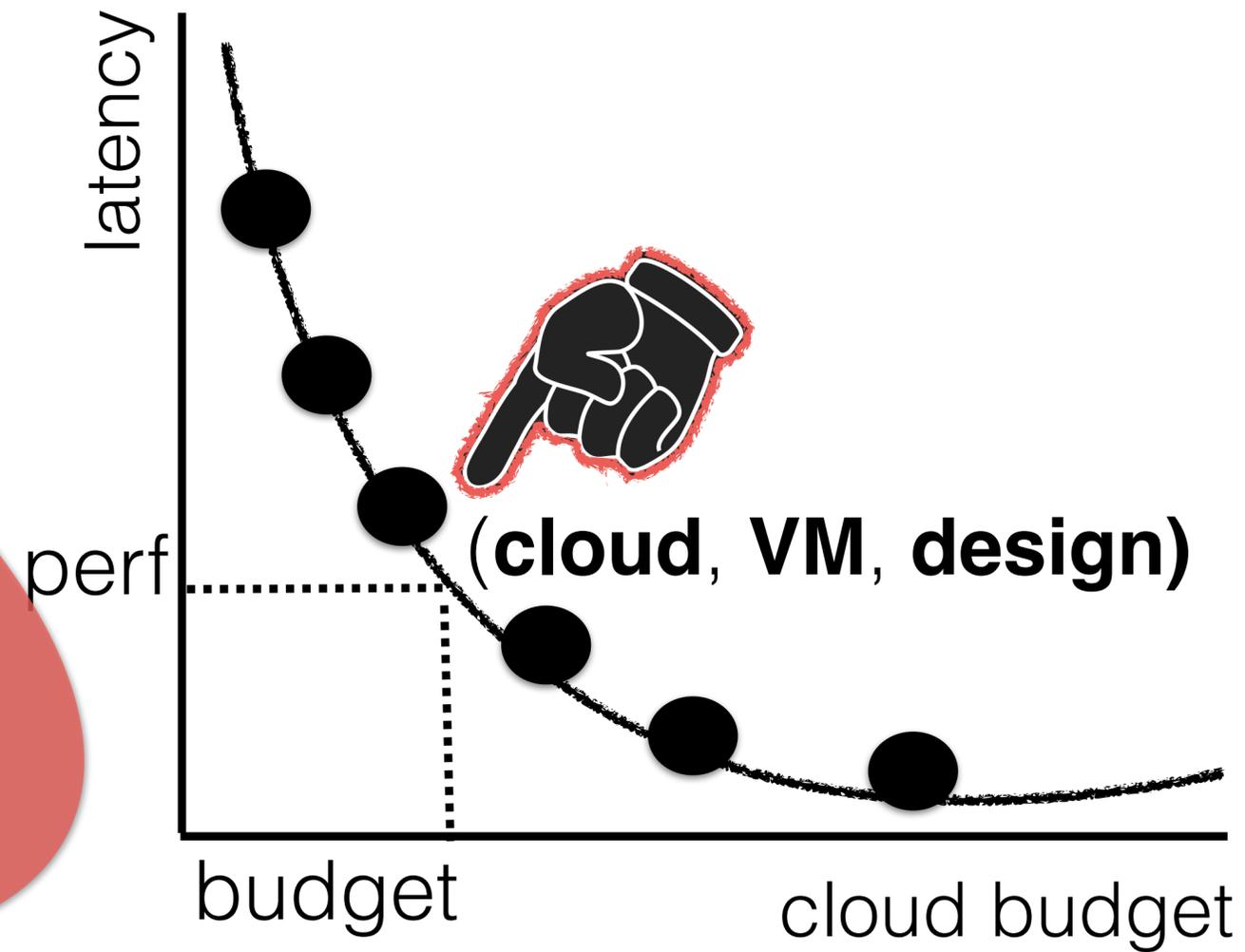
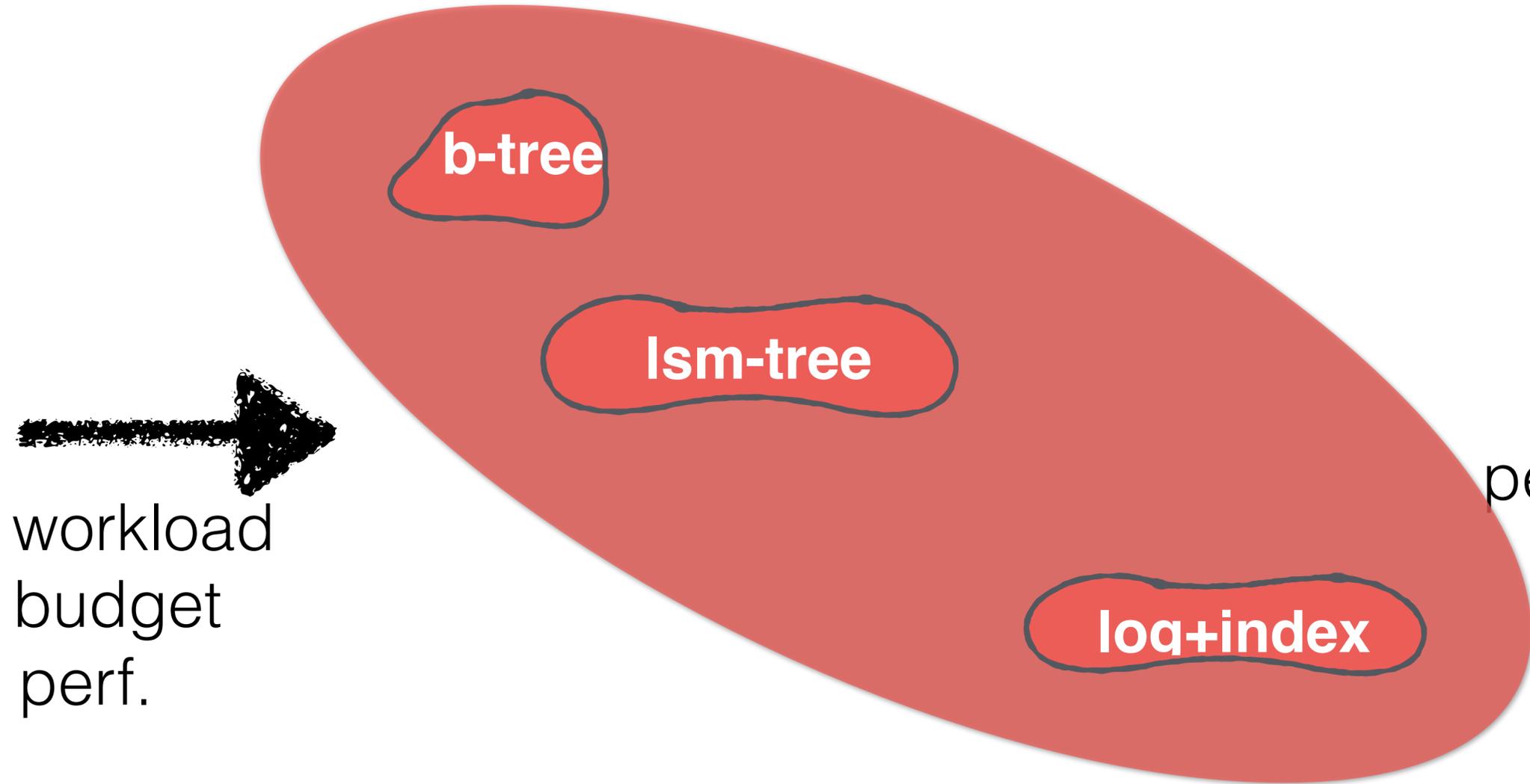
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

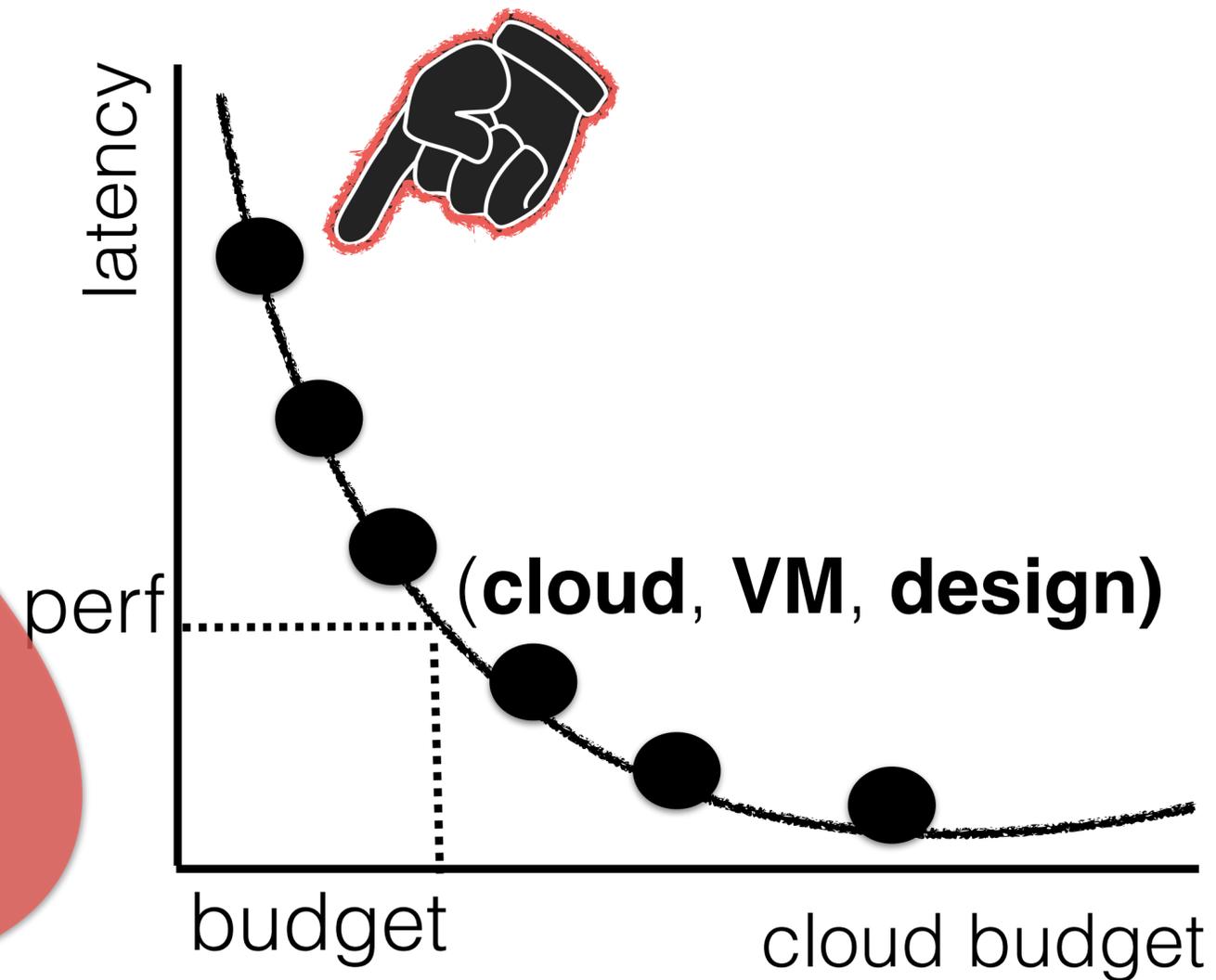
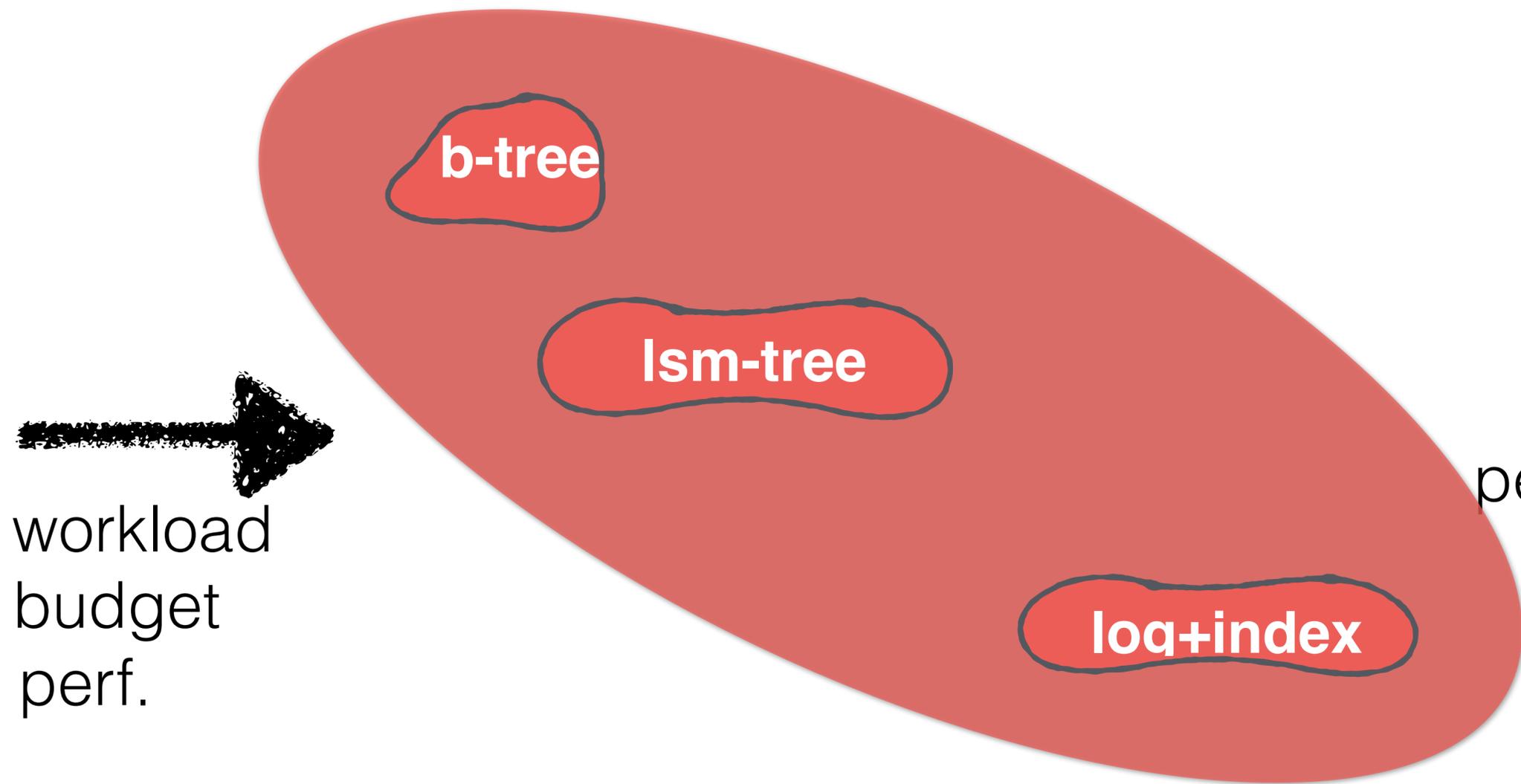
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

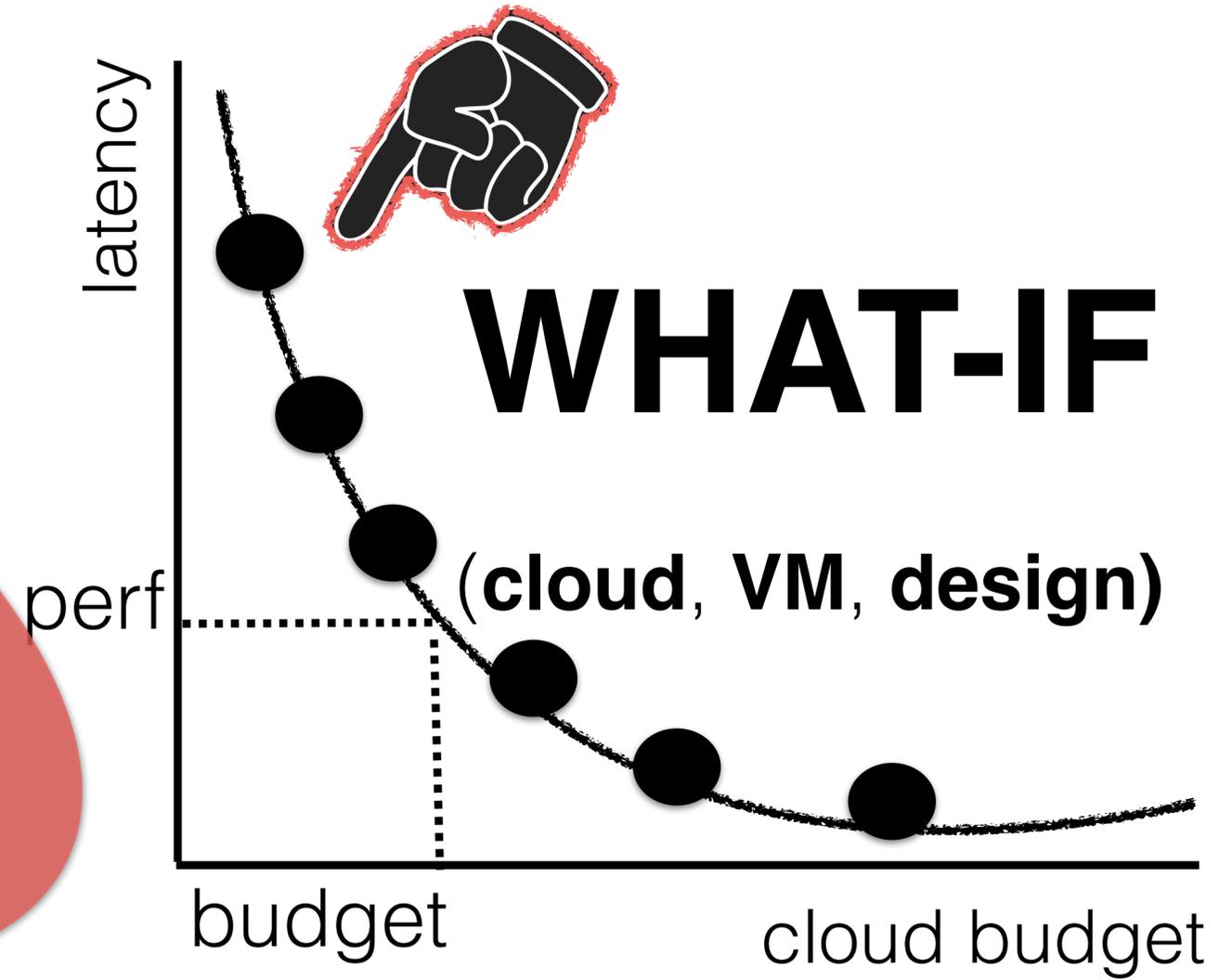
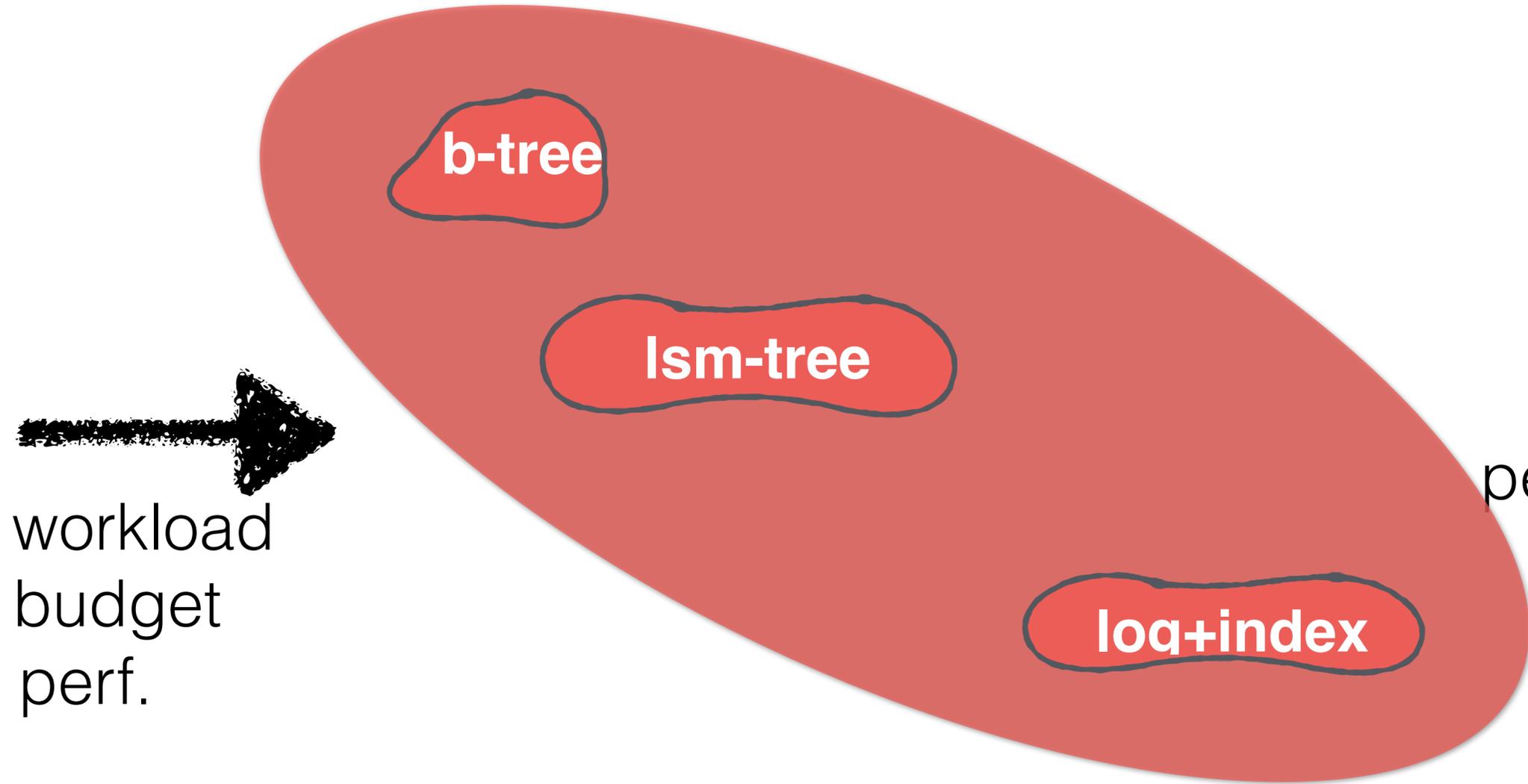
unified
closed form

**learned
CPU model**

h/w,
parallelism

**cloud cost
mapping & SLA**

AWS, Azure, Google



**analytical
I/O model**

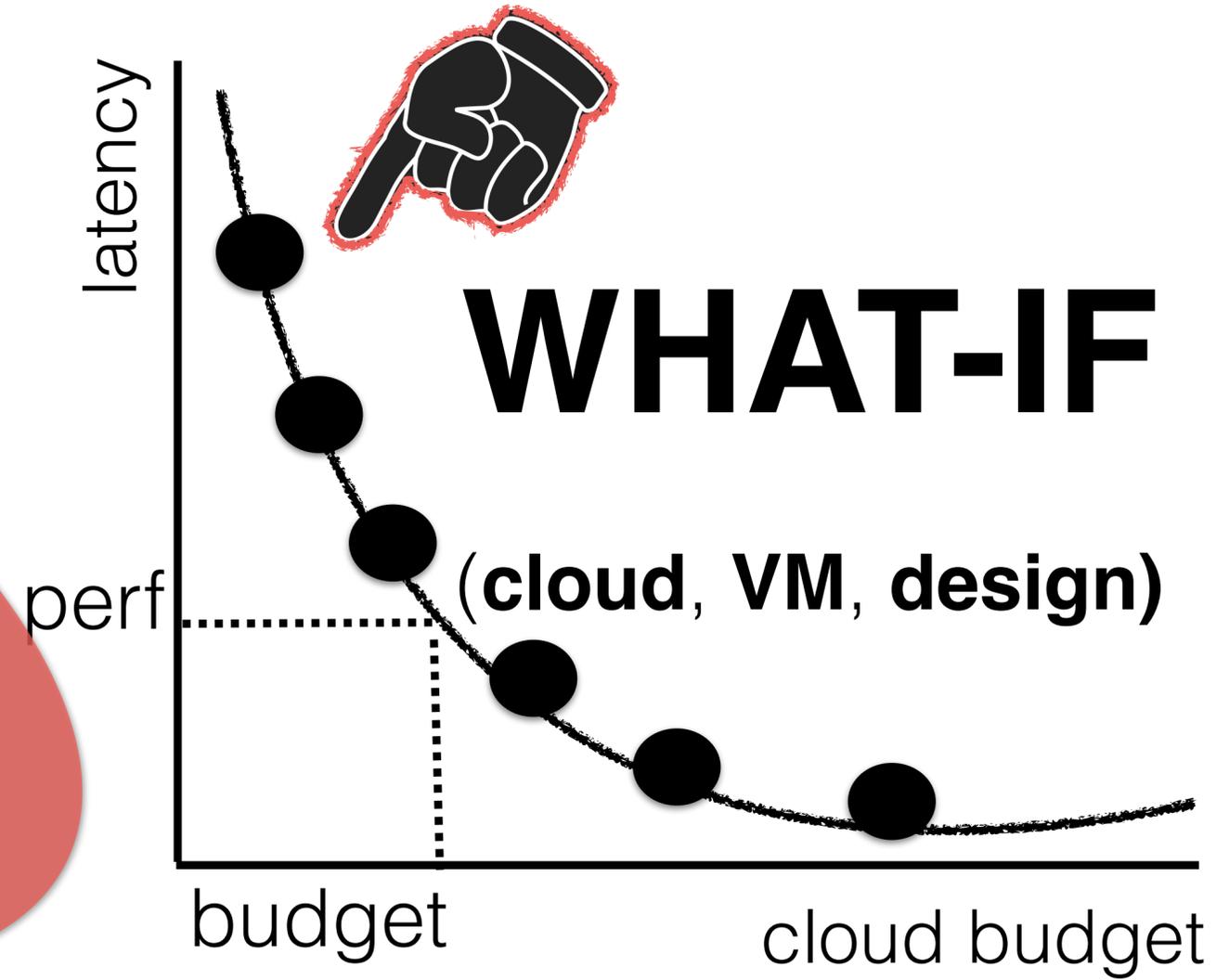
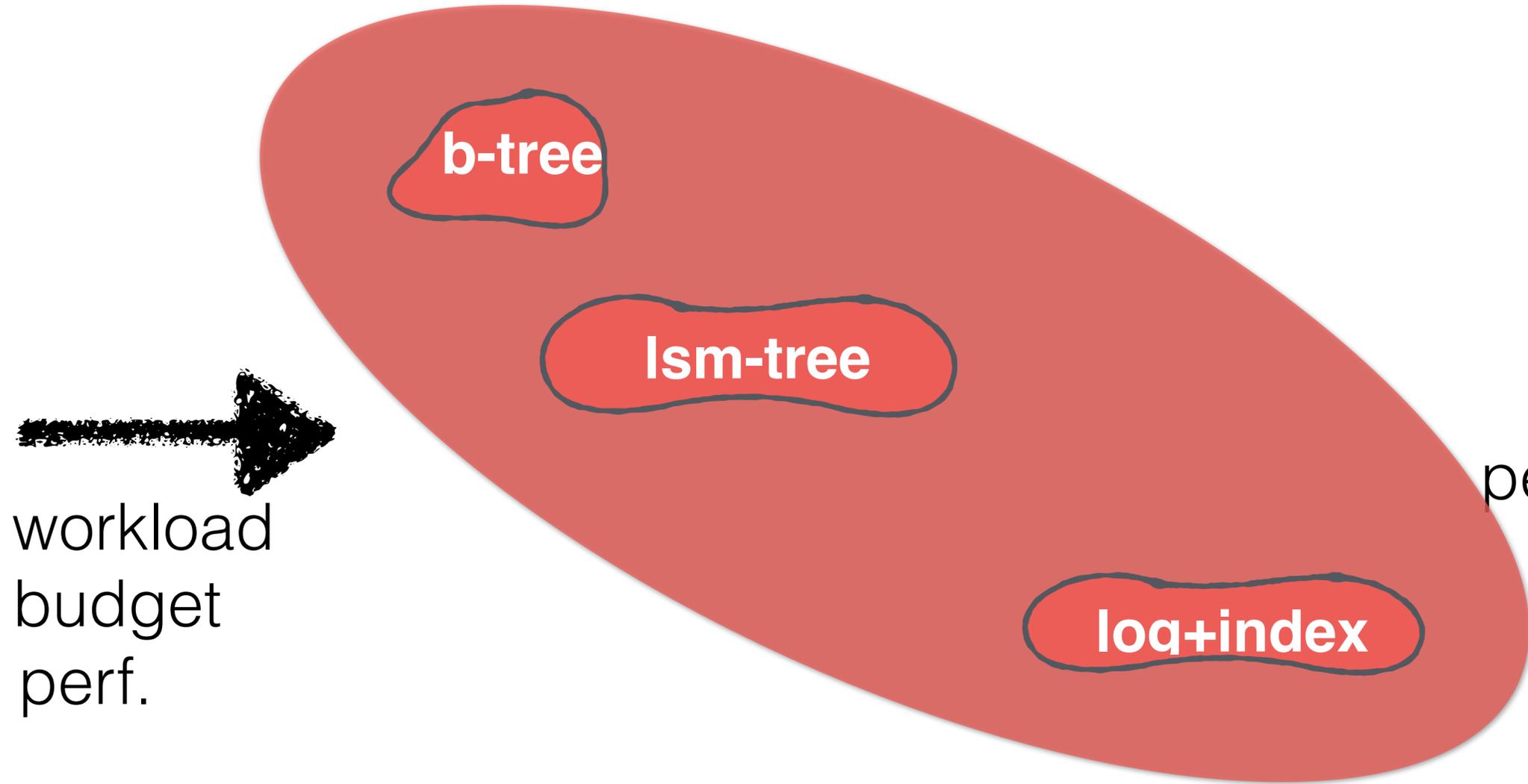
**learned
CPU model**

**cloud cost
mapping & SLA**

unified
closed form

h/w,
parallelism

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?

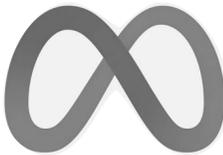
workload/budget diversity

Can we beat the best system for every workload?

How to test?

workload/budget diversity

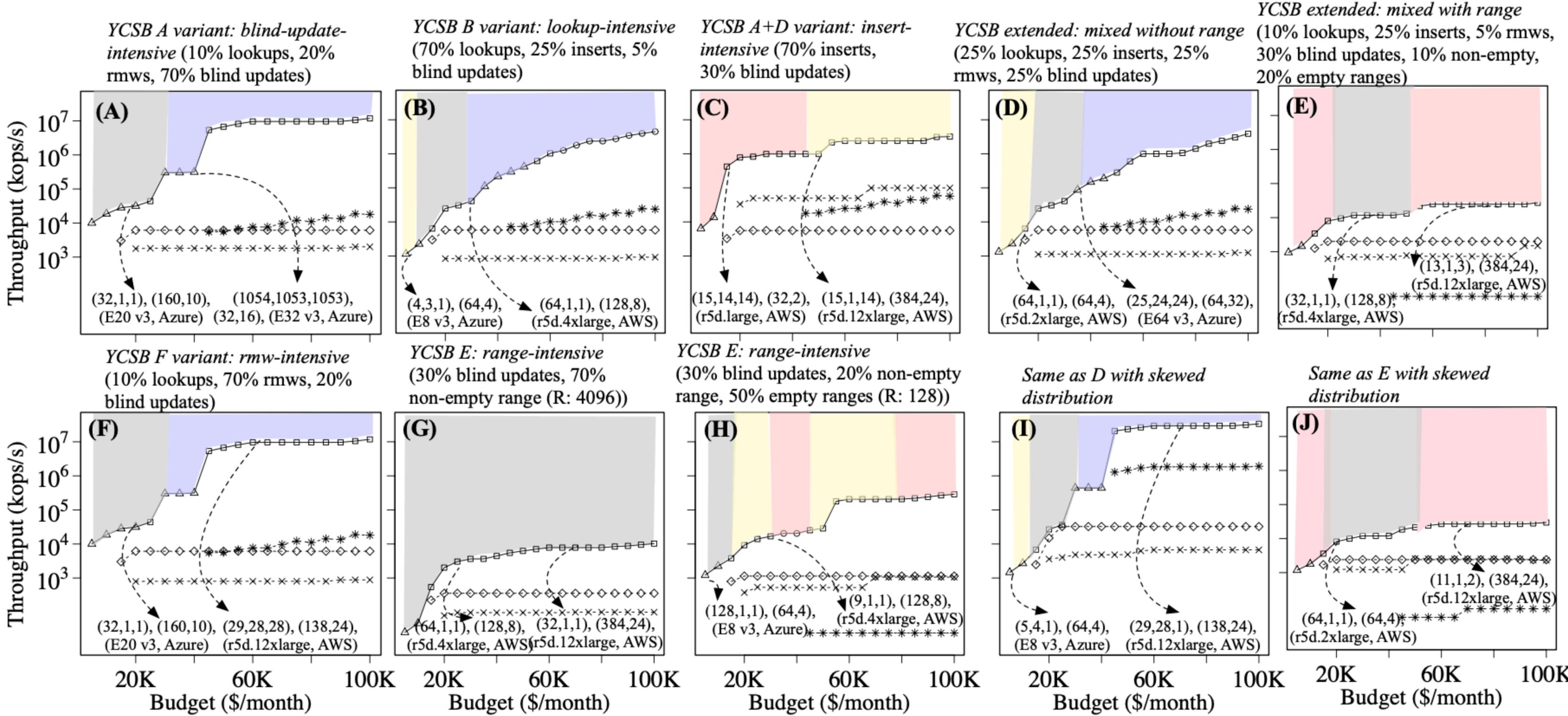
Can we beat the best system for every workload?

 **Meta**

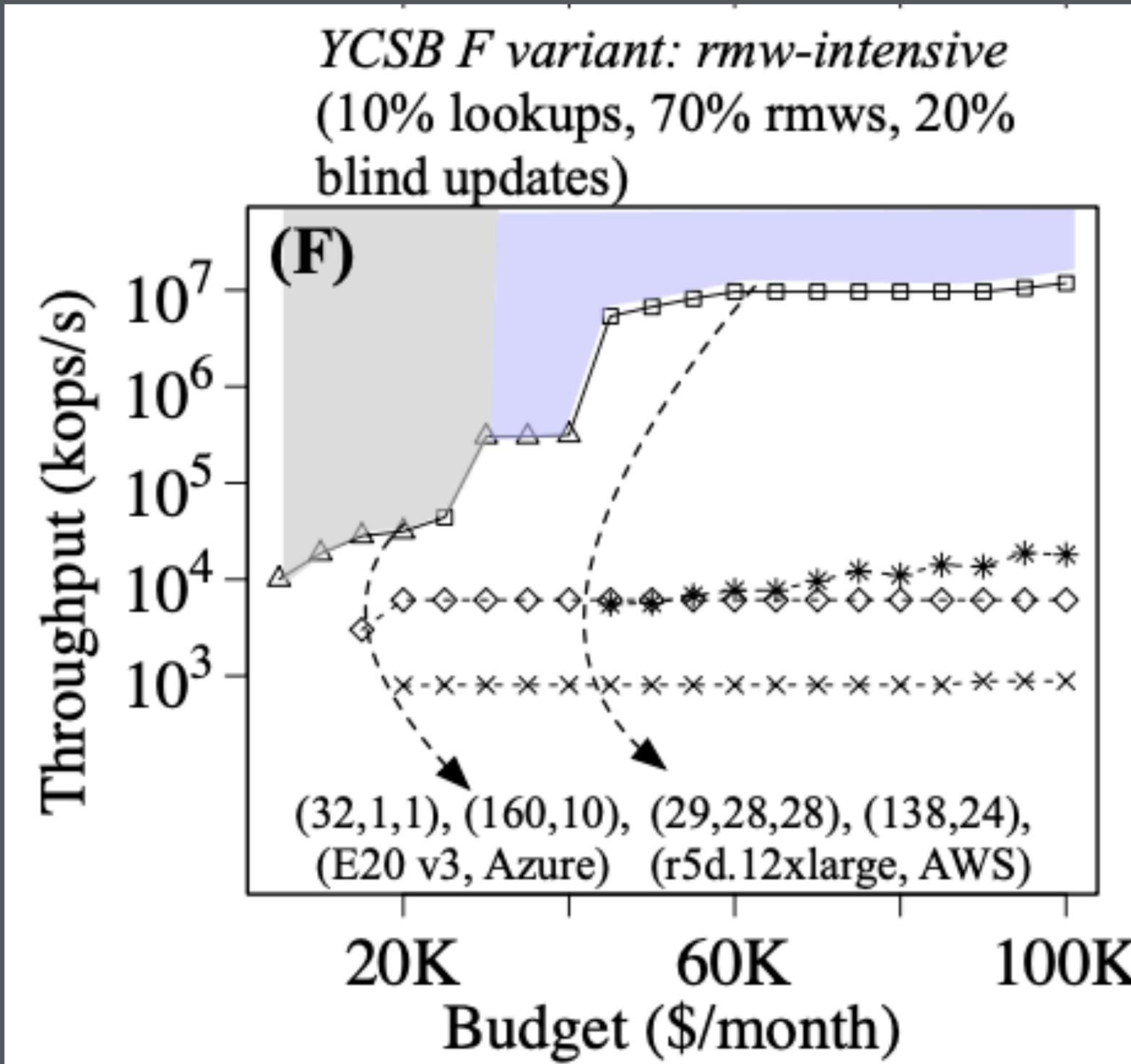
 **mongoDB[®]**

 **Microsoft**

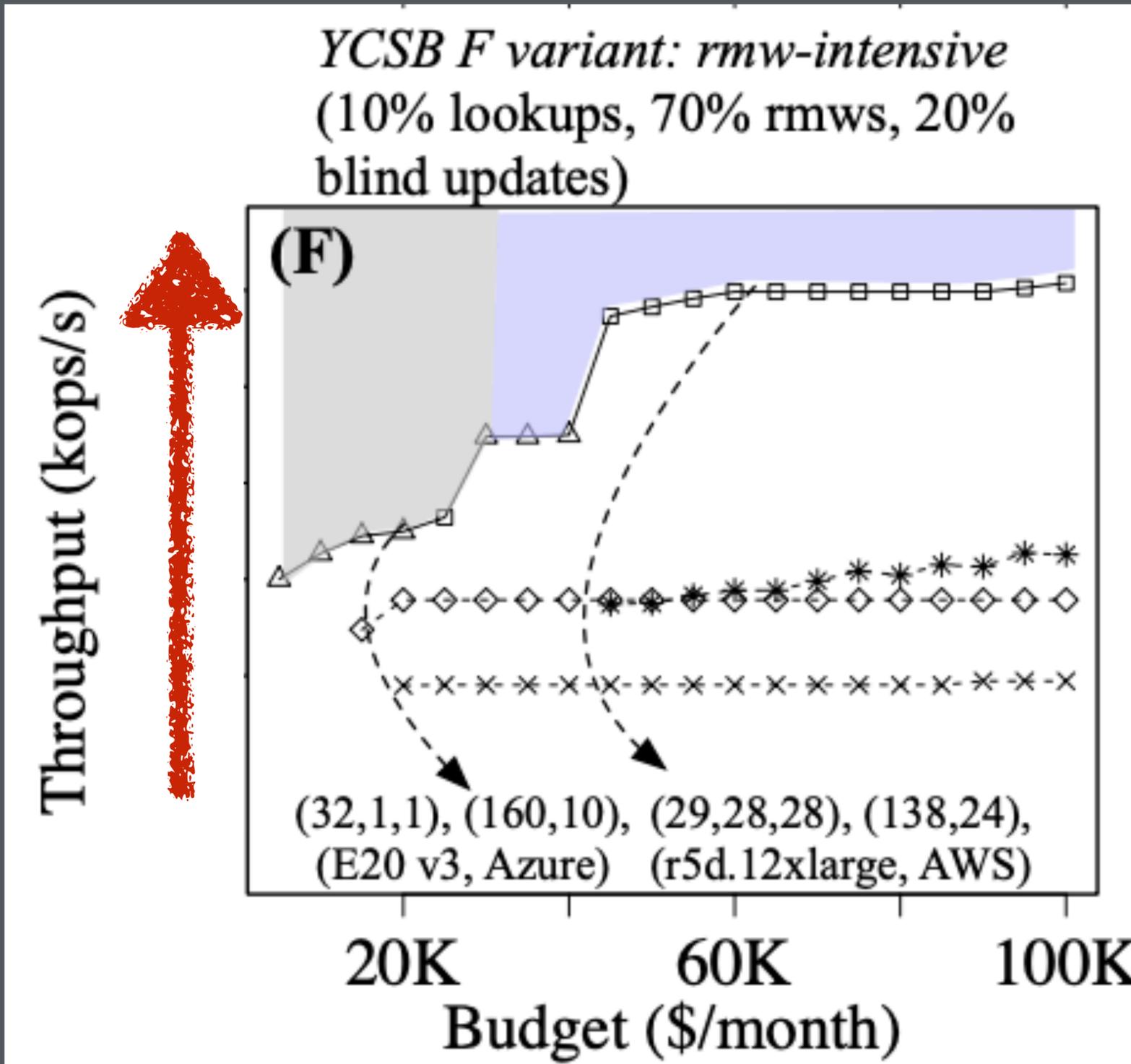
-x- RocksDB -◇- WiredTiger -*- FASTER -⊕- Cosine on AWS -⊖- Cosine on GCP -△- Cosine on Azure
 LSM class
 B-tree class
 LSH class
 Hybrid class



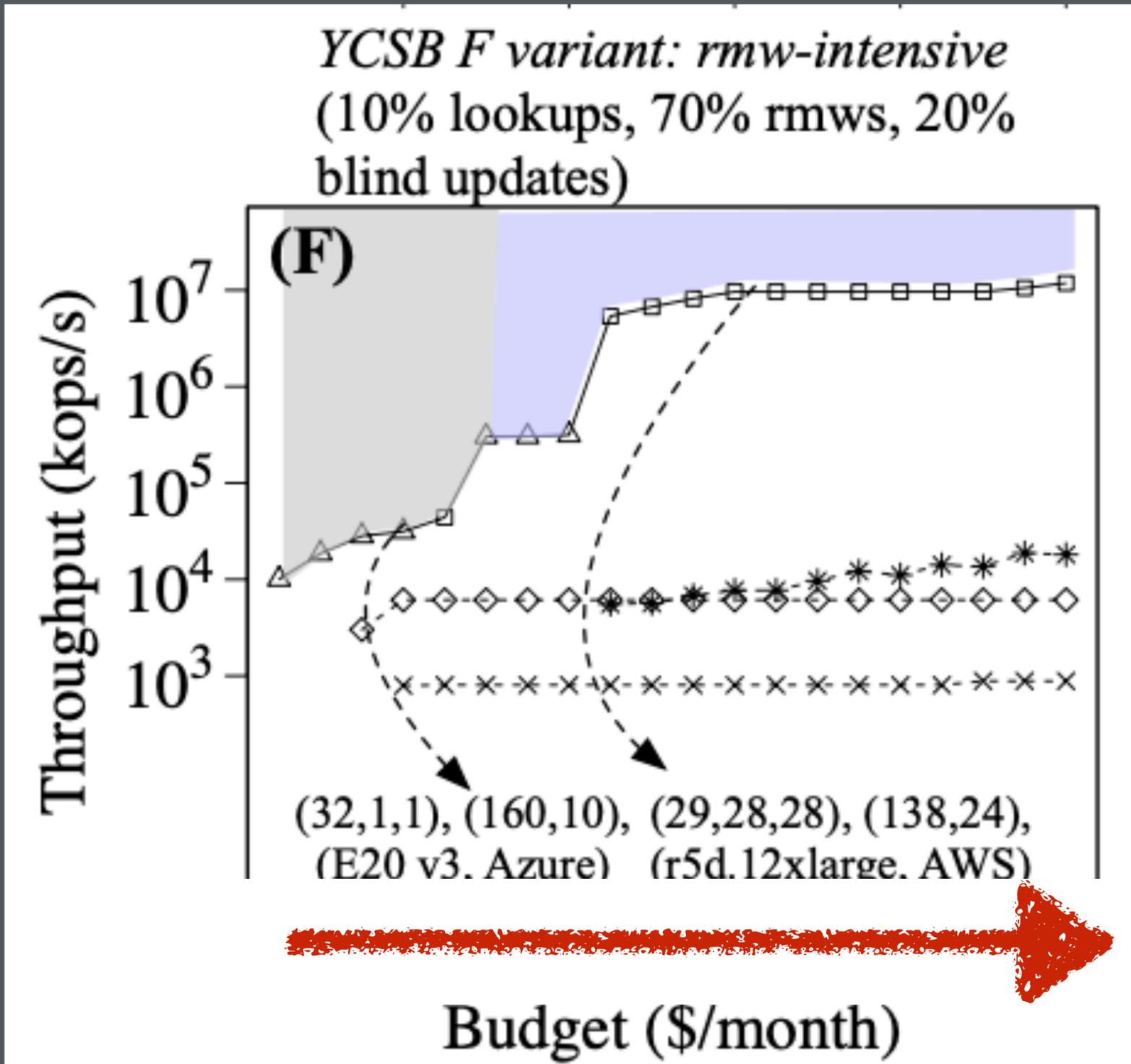
workload/budget diversity



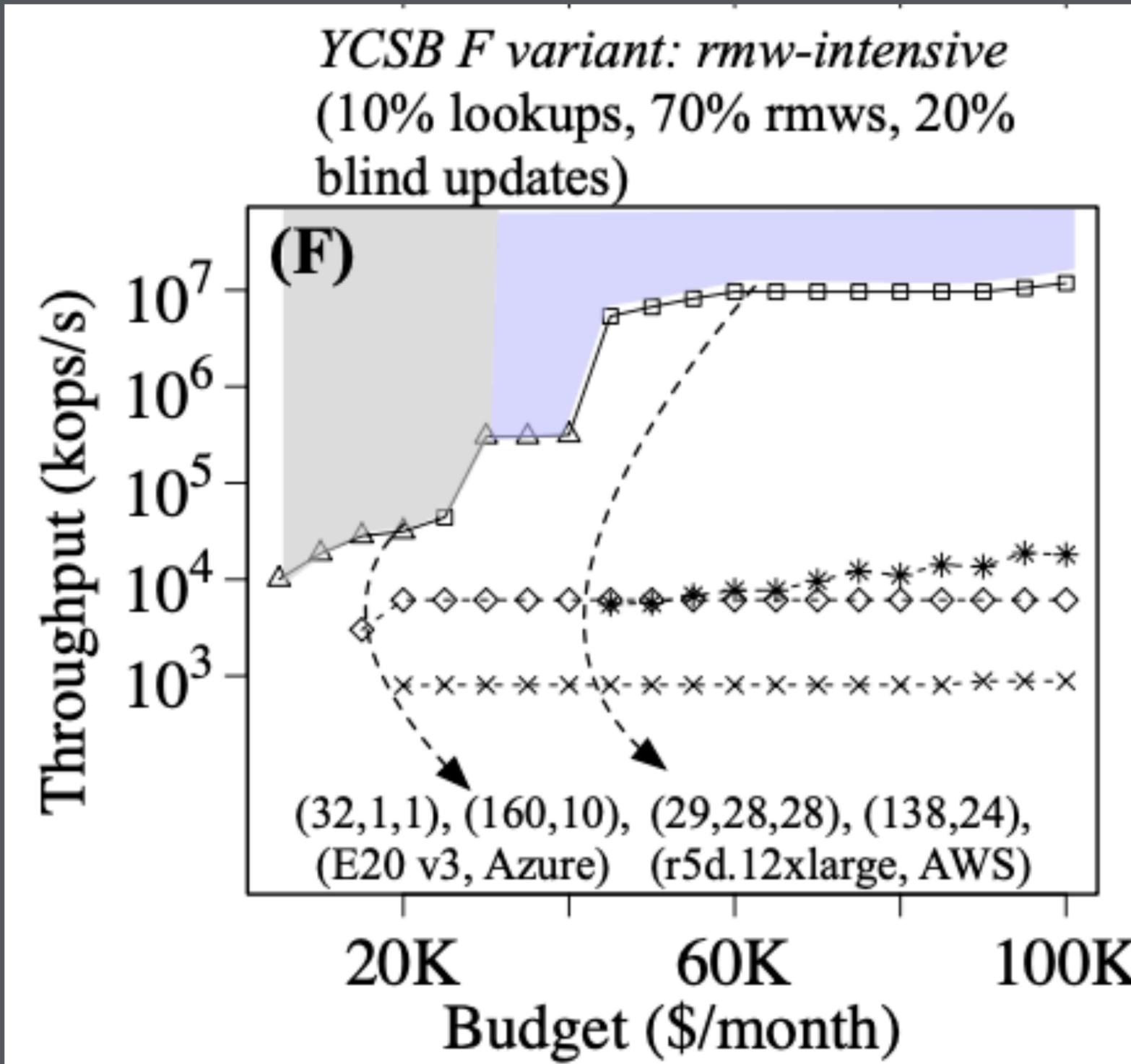
workload/budget diversity



workload/budget diversity

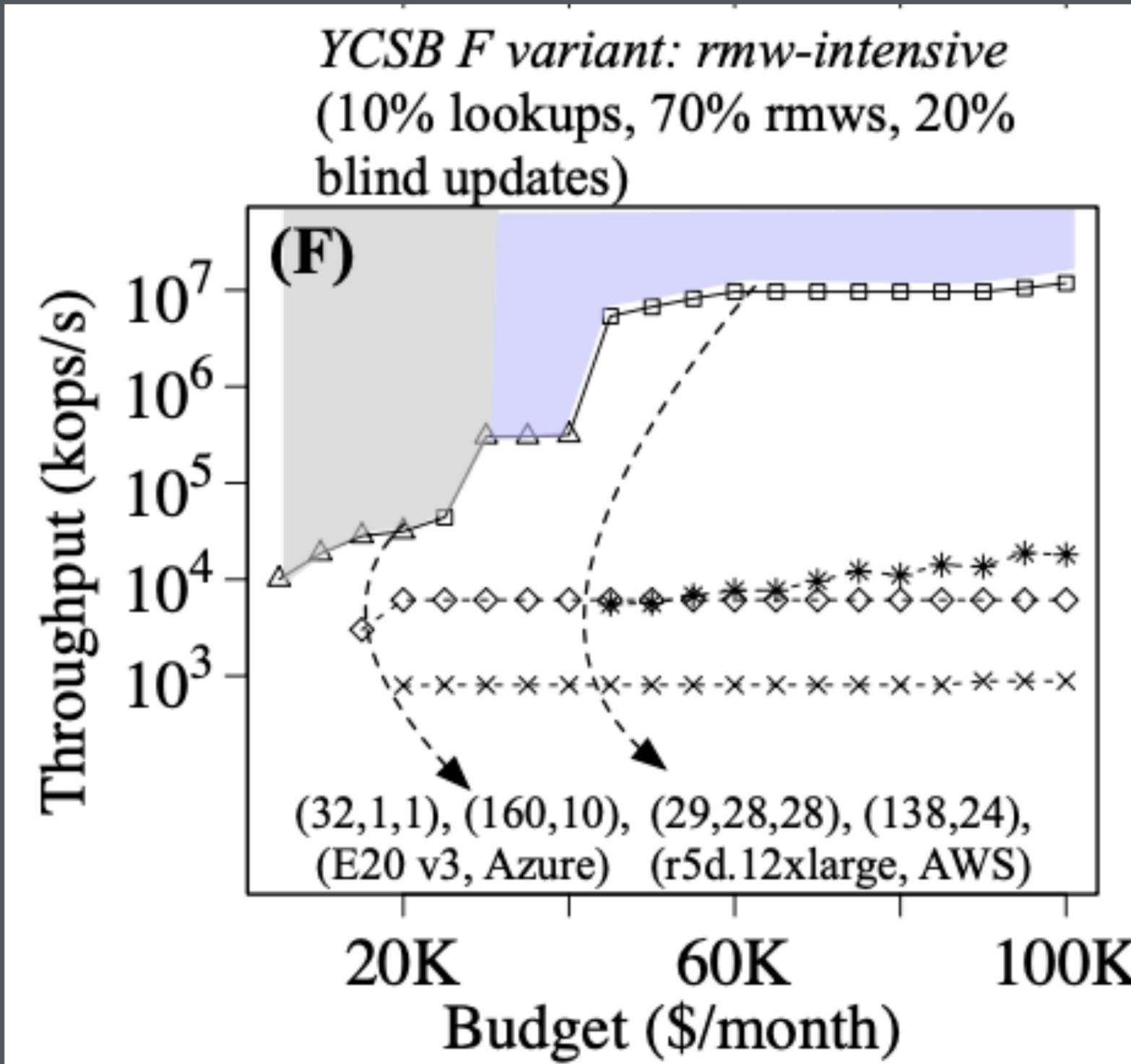


workload/budget diversity



⊕ COSINE ⊕

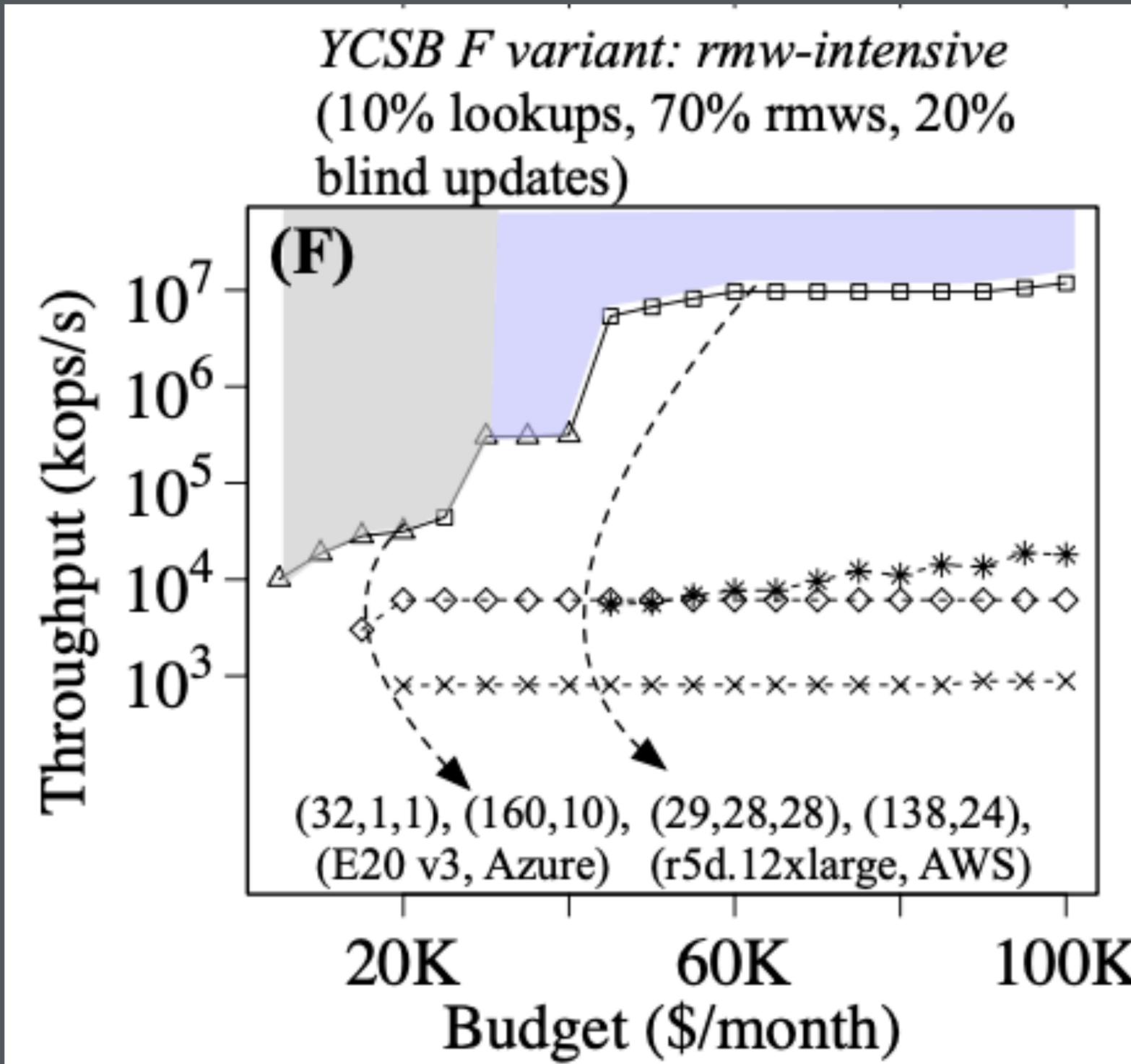
workload/budget diversity



⚙️ COSINE ⚙️

{ State of the Art
Meta, Microsoft, Mongo

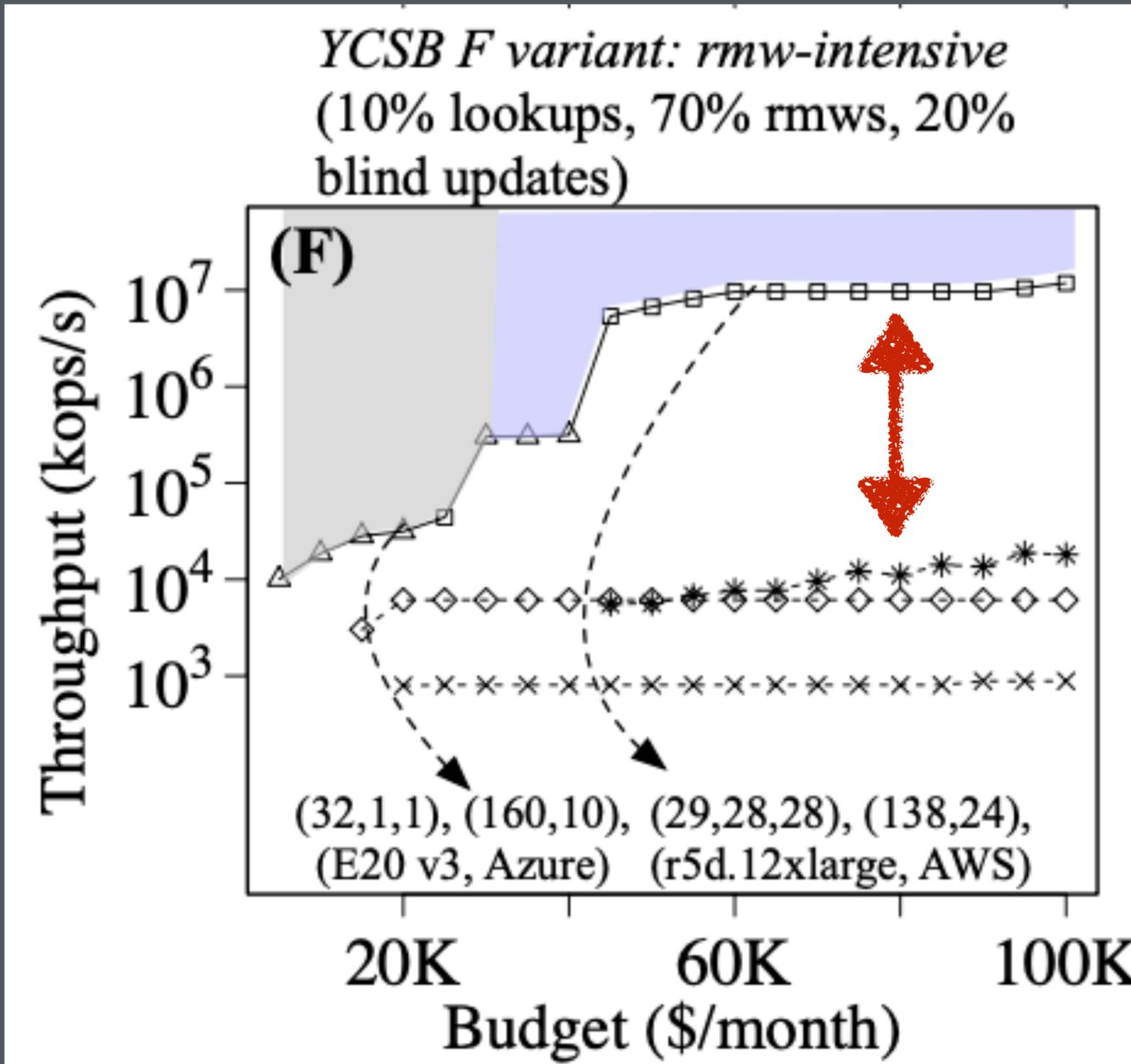
workload/budget diversity



⚙️ COSINE ⚙️

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity

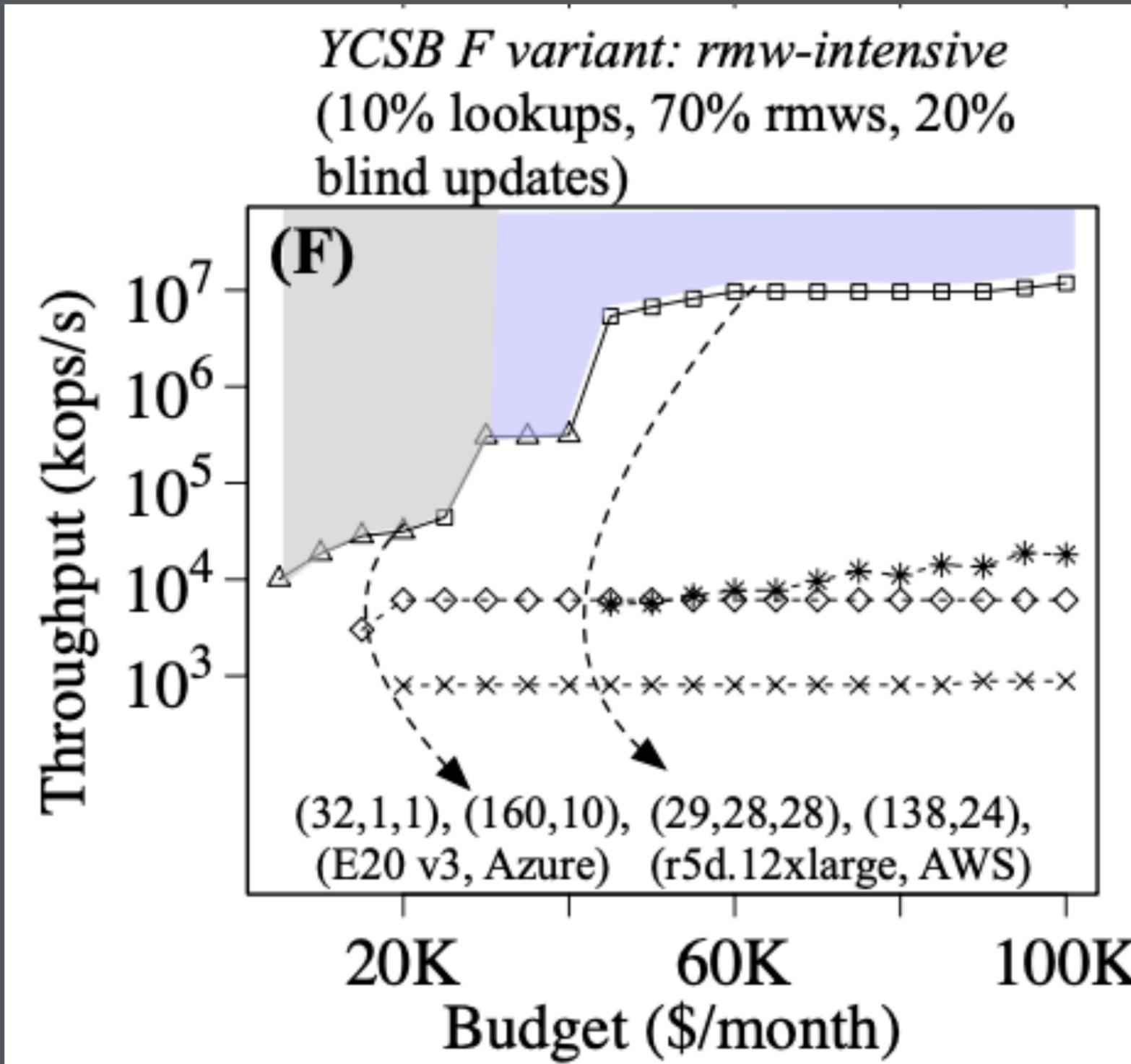


⚙️ COSINE ⚙️

Better throughput/cost

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity



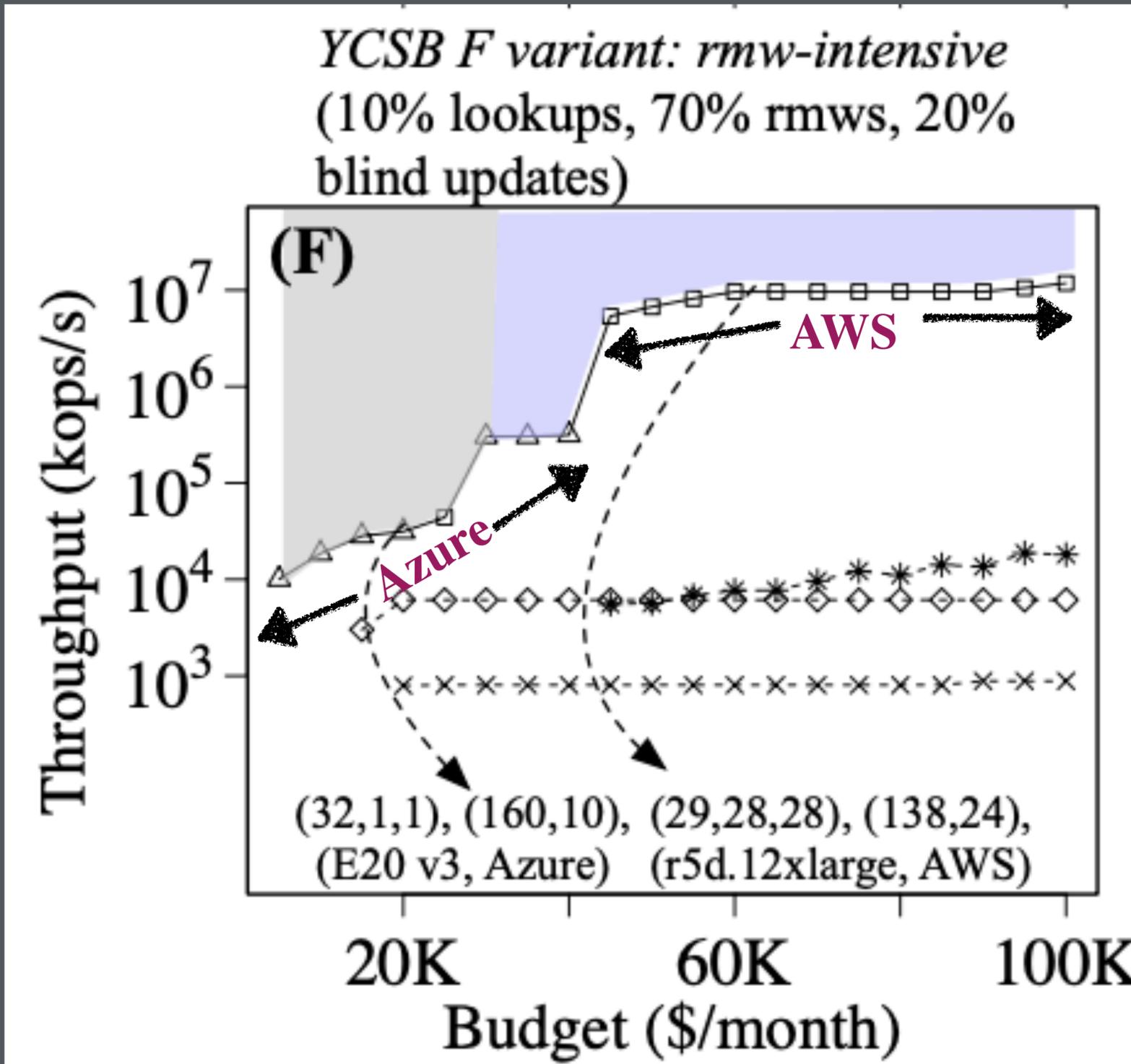
⚙️ COSINE ⚙️

Better throughput/cost

Self-designs

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity

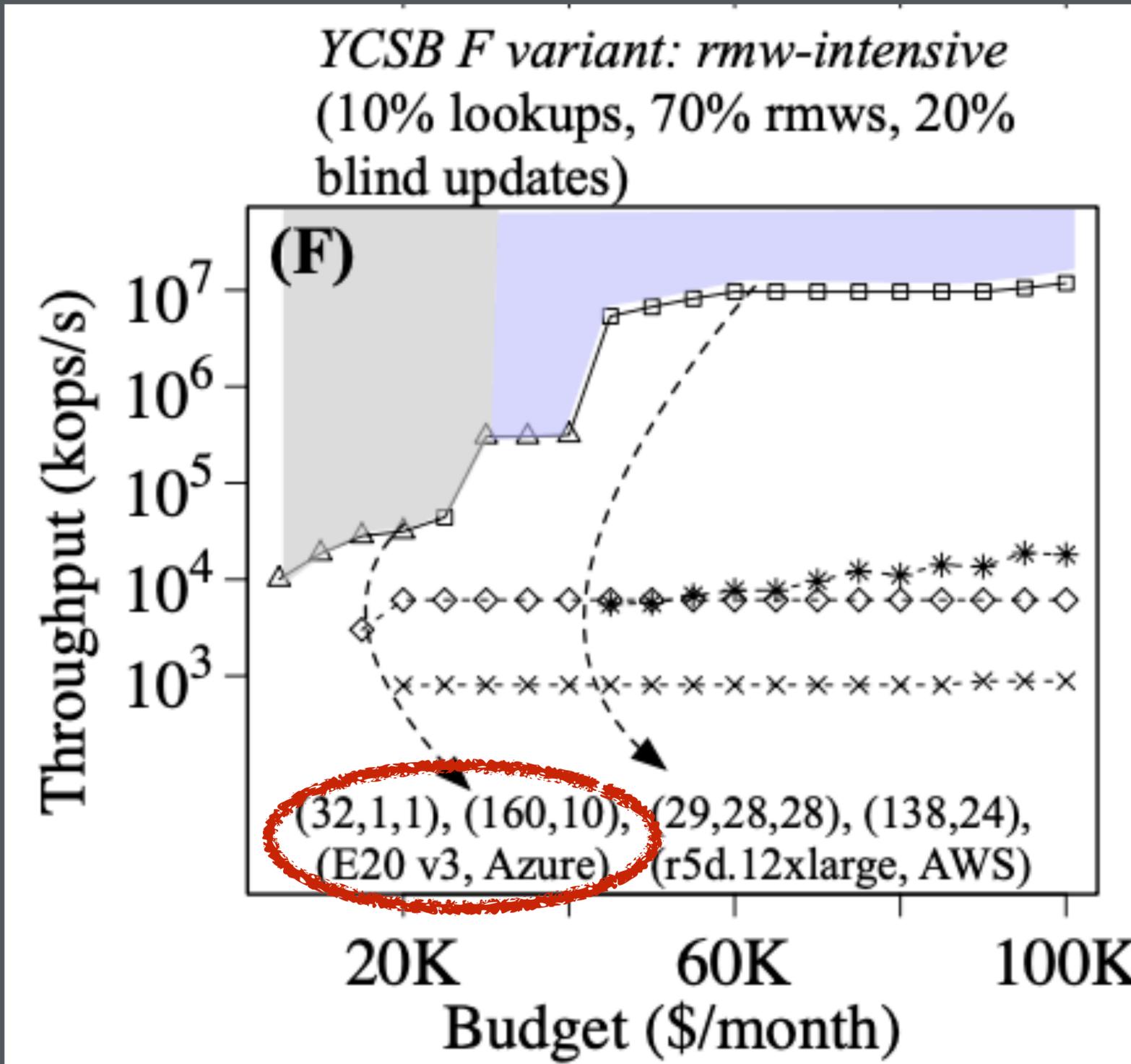


⚙️ COSINE ⚙️

Better throughput/cost
Self-designs

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity



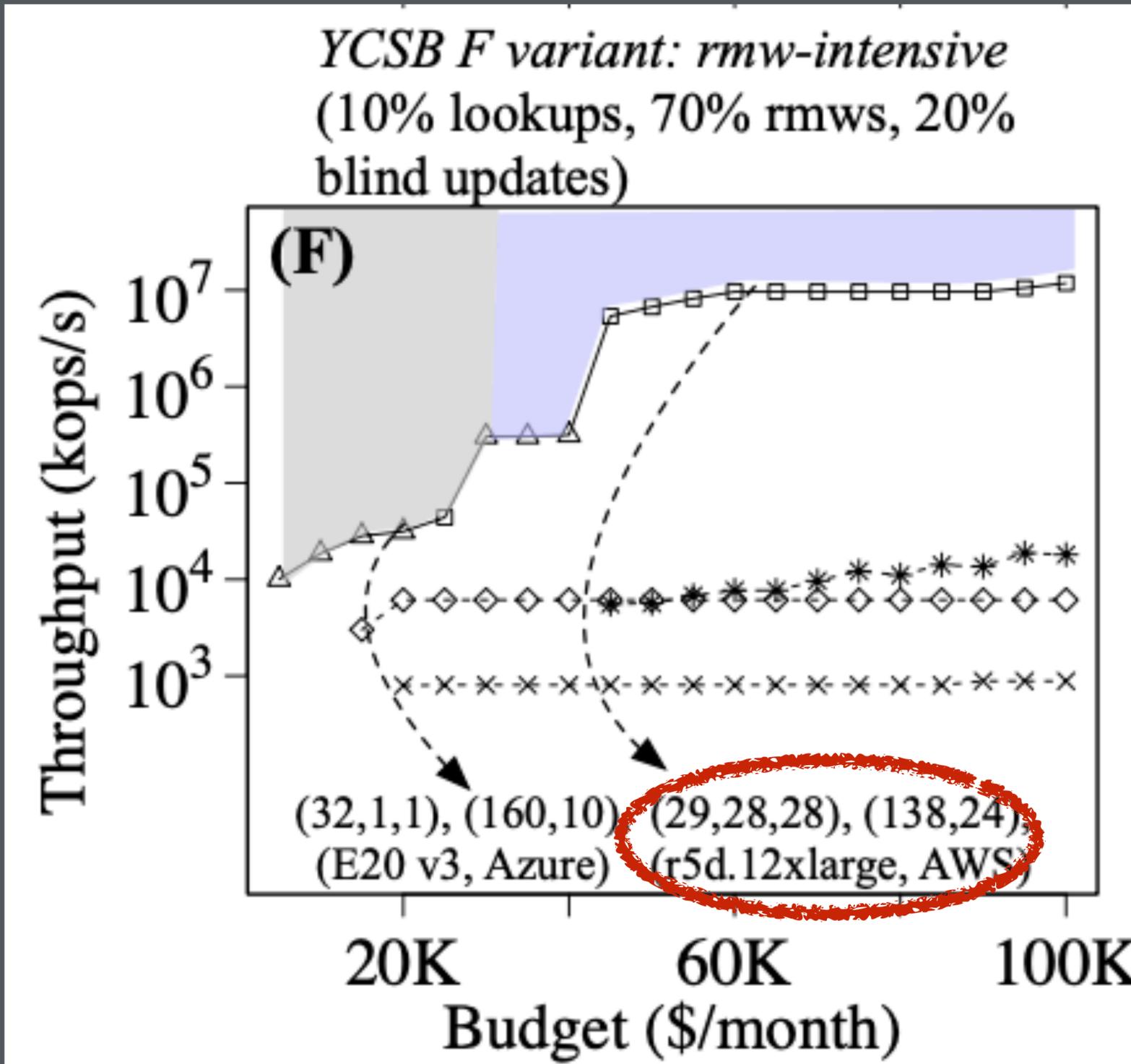
⚙️ COSINE ⚙️

Better throughput/cost

Self-designs

{ State of the Art
Meta, Microsoft, Mongo

workload/budget diversity



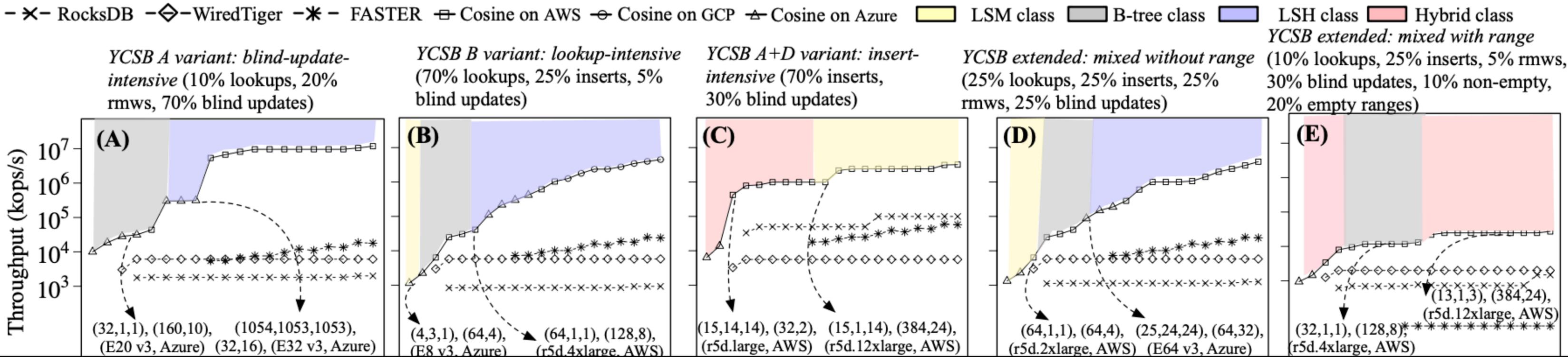
⚙️ COSINE ⚙️

Better throughput/cost

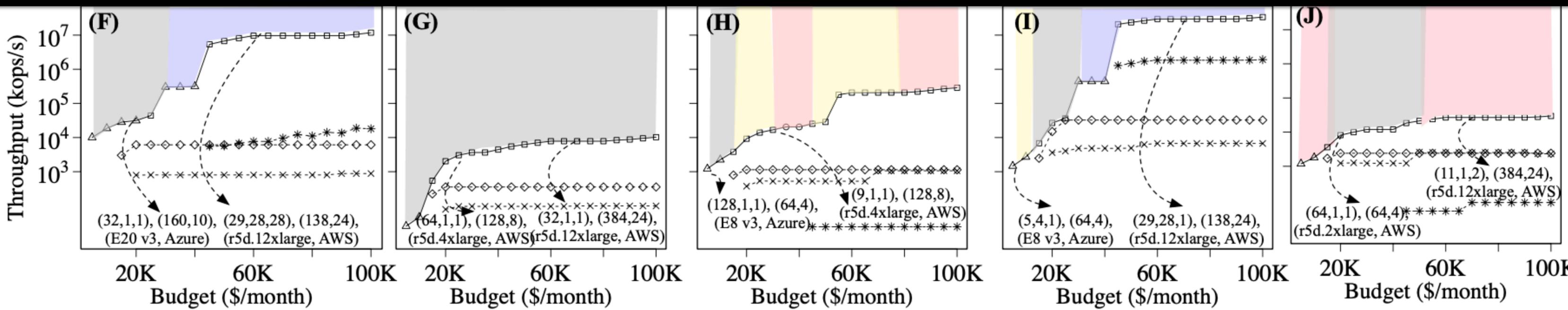
Self-designs

{ State of the Art
Meta, Microsoft, Mongo

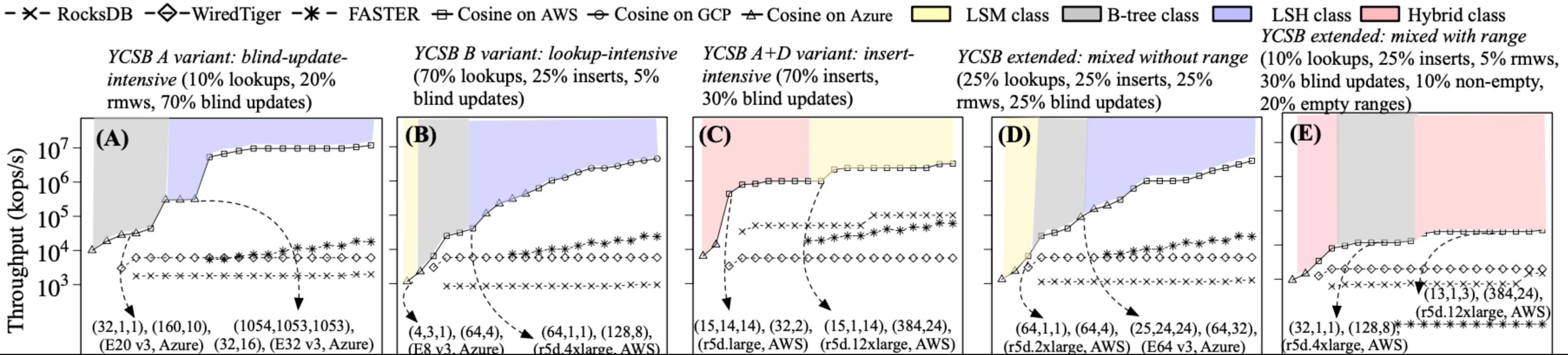
workload/budget diversity



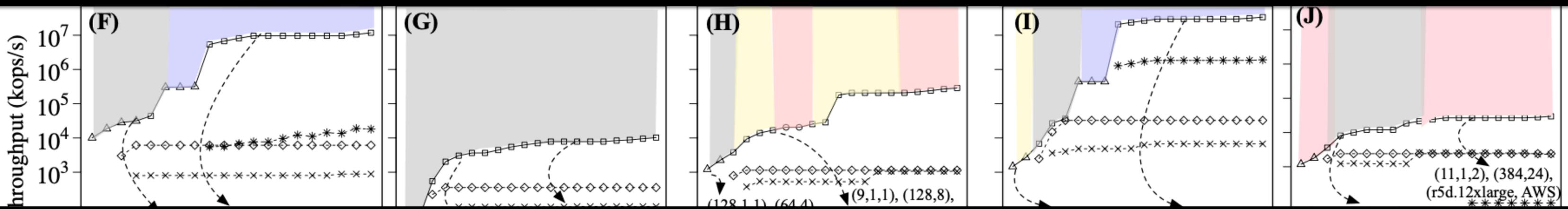
diversity beats top systems self-designs (provider, VM, new design)



workload/budget diversity



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

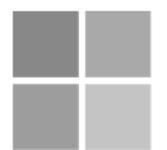


Cosine achieves the best perf. across all workloads by automatically designing a new system every time.

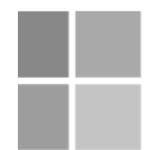
- ✕ - RocksDB - ⬠ - WiredTiger - * - FASTER - ⇆ - Cosine on AWS - ⊖ - Cosine on GCP - ⚡ - Cosine on Azure

 LSM class
 B-tree class
 LSH class
 Hybrid class

YCSB A variant: blind-update-intensive (10% lookups, 20% rmws, 70% blind updates)
 YCSB B variant: lookup-intensive (70% lookups, 25% inserts, 5% blind updates)
 YCSB A+D variant: insert-intensive (70% inserts, 30% blind updates)
 YCSB extended: mixed without range (25% lookups, 25% inserts, 25% rmws, 25% blind updates)
 YCSB extended: mixed with range (10% lookups, 25% inserts, 5% rmws, 30% blind updates, 10% non-empty, 20% empty ranges)

 Microsoft
  Microsoft
  Meta
  Microsoft
  mongo

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

 Microsoft
  mongo
  Meta
  Microsoft
  Meta

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

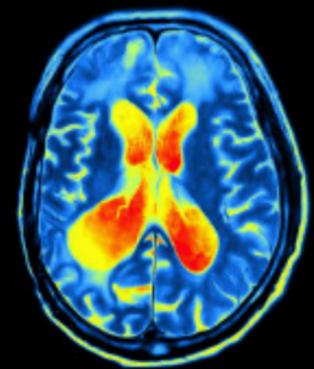
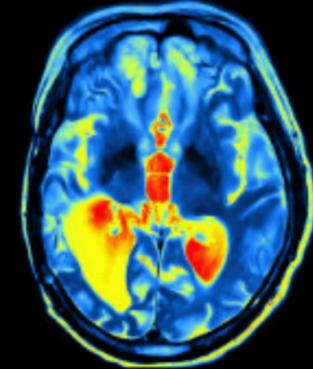
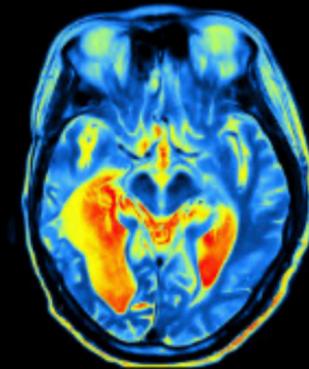
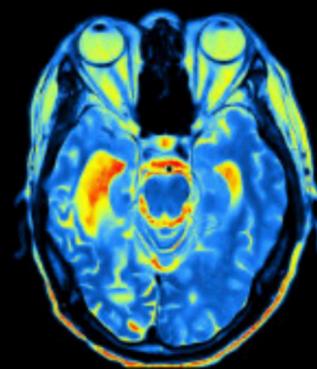
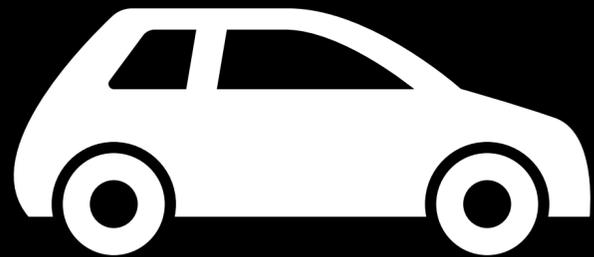
Papers: **Cosine** PVLDB 2023, and new **Limousine** at SIGMOD 2024

How do these concepts translate to the other big data areas
neural networks, image AI, Blockchain, ...?

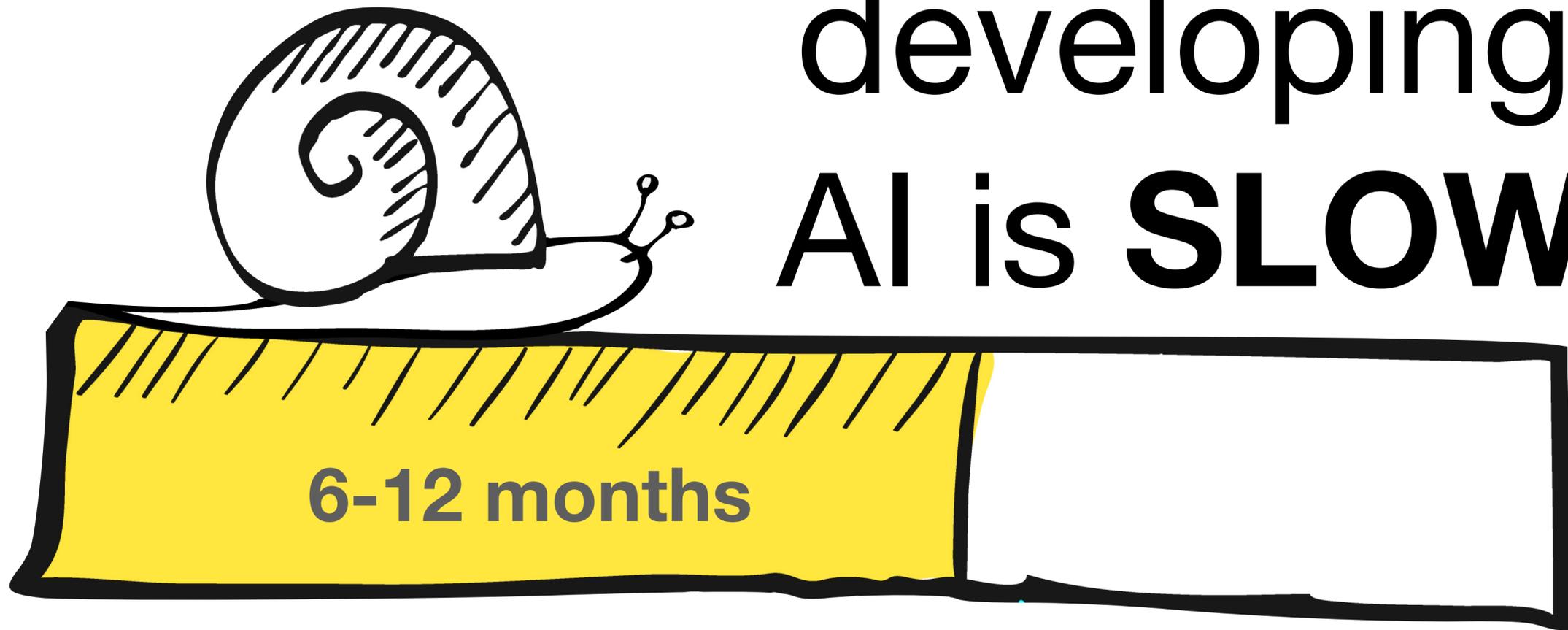
again, it all starts from the storage design space

seeing is at the very center of AI
because it is at the center of human life

image processing



developing
AI is **SLOW**

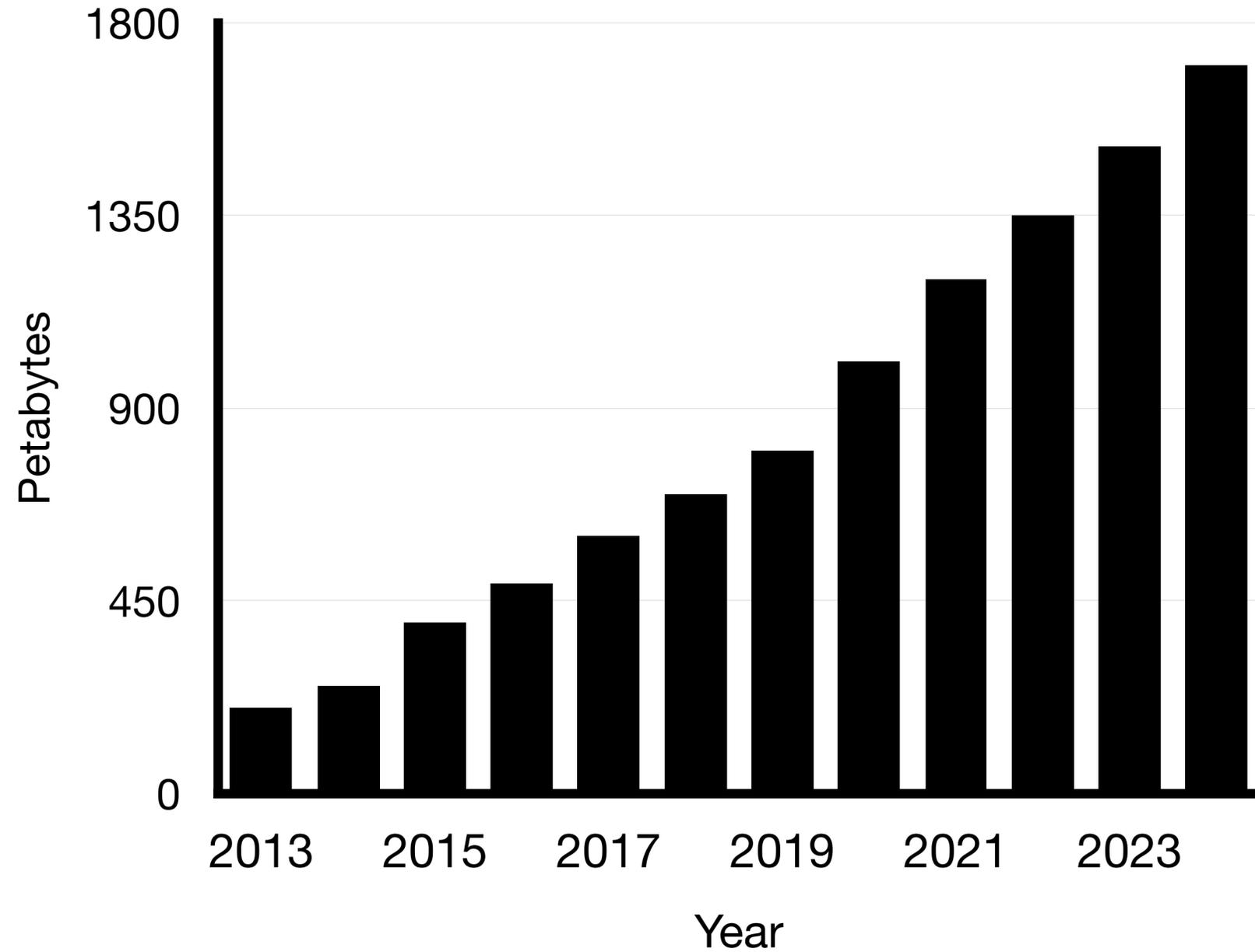


6-12 months

Data size

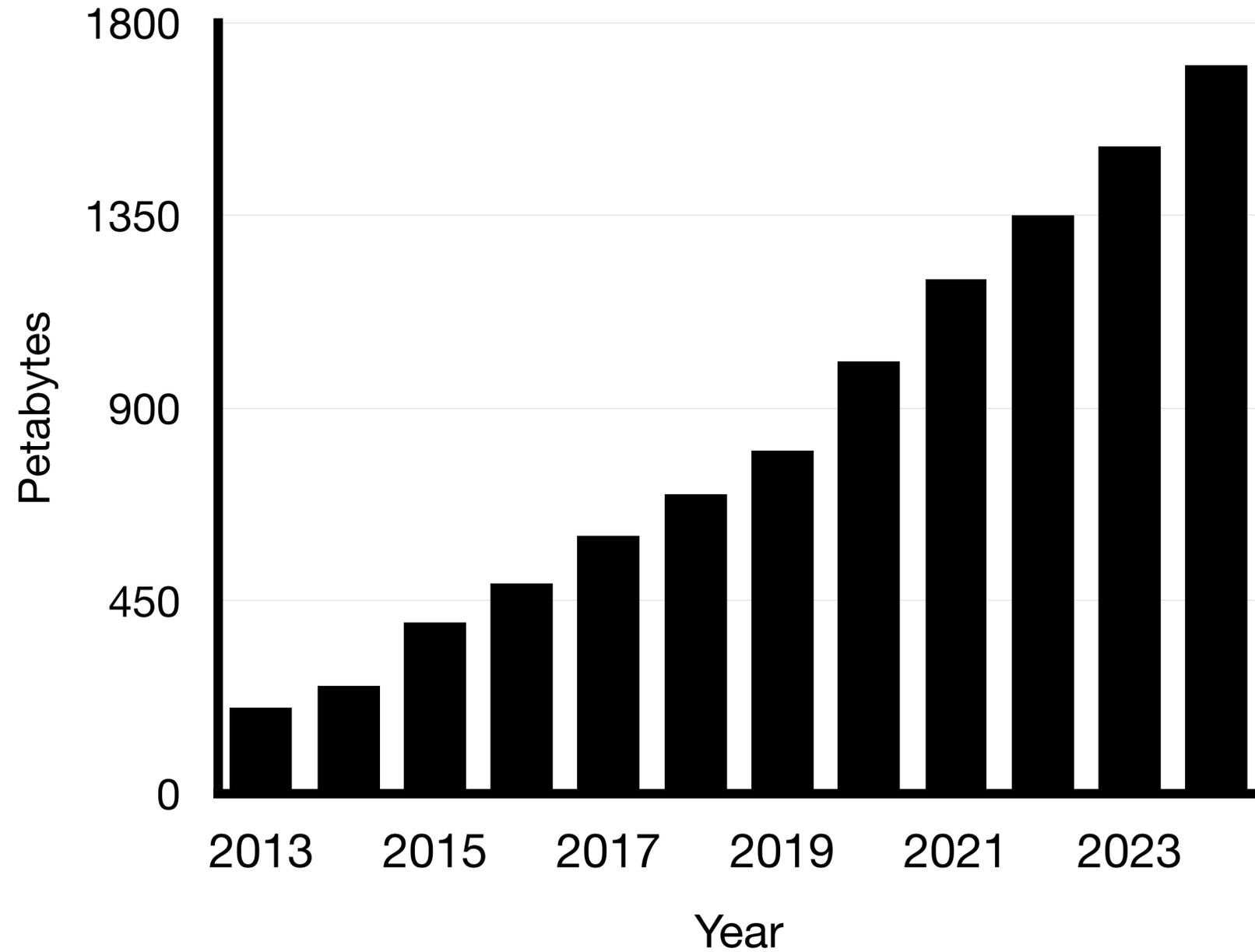
Model size

Data size

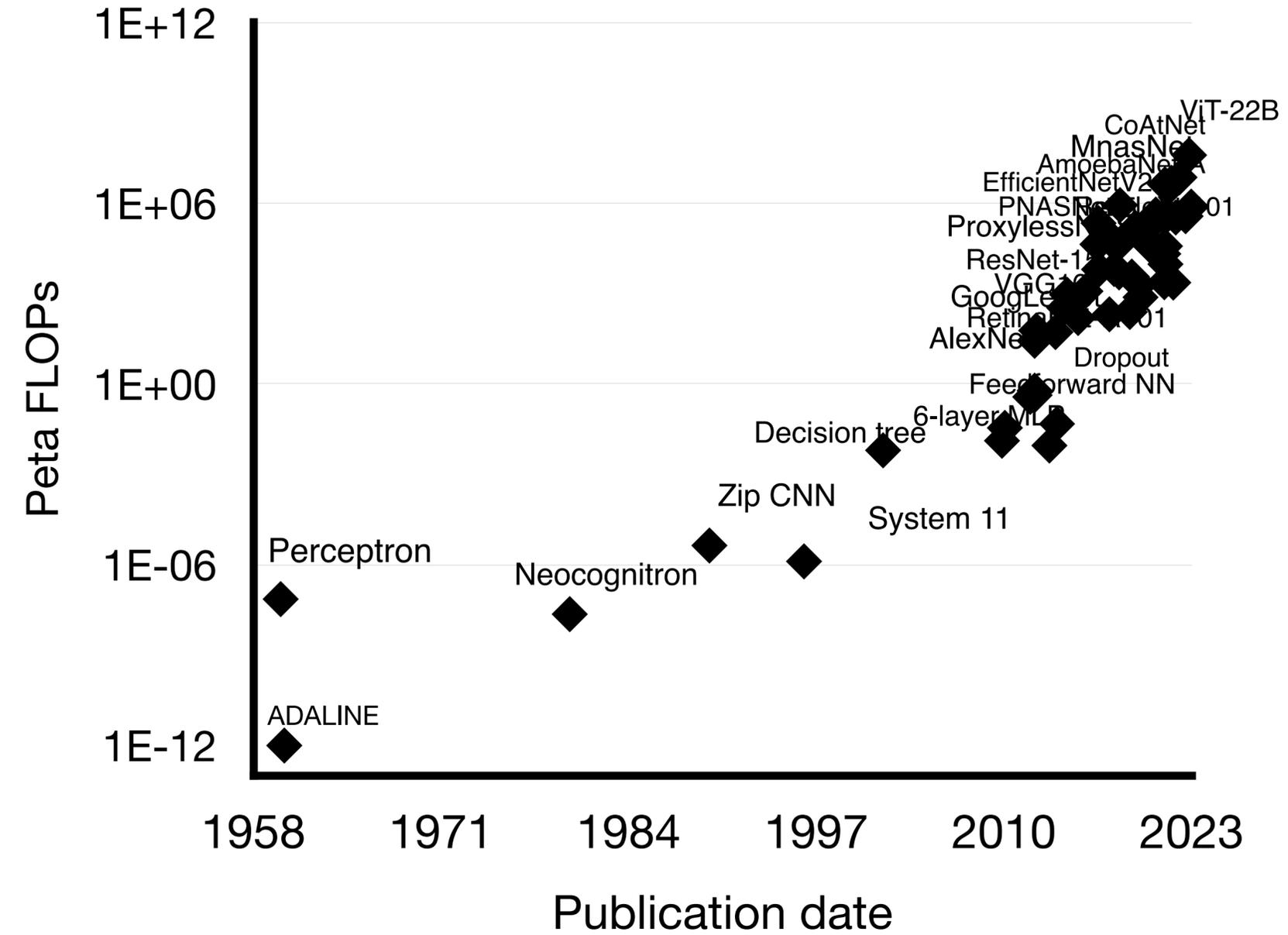


Model size

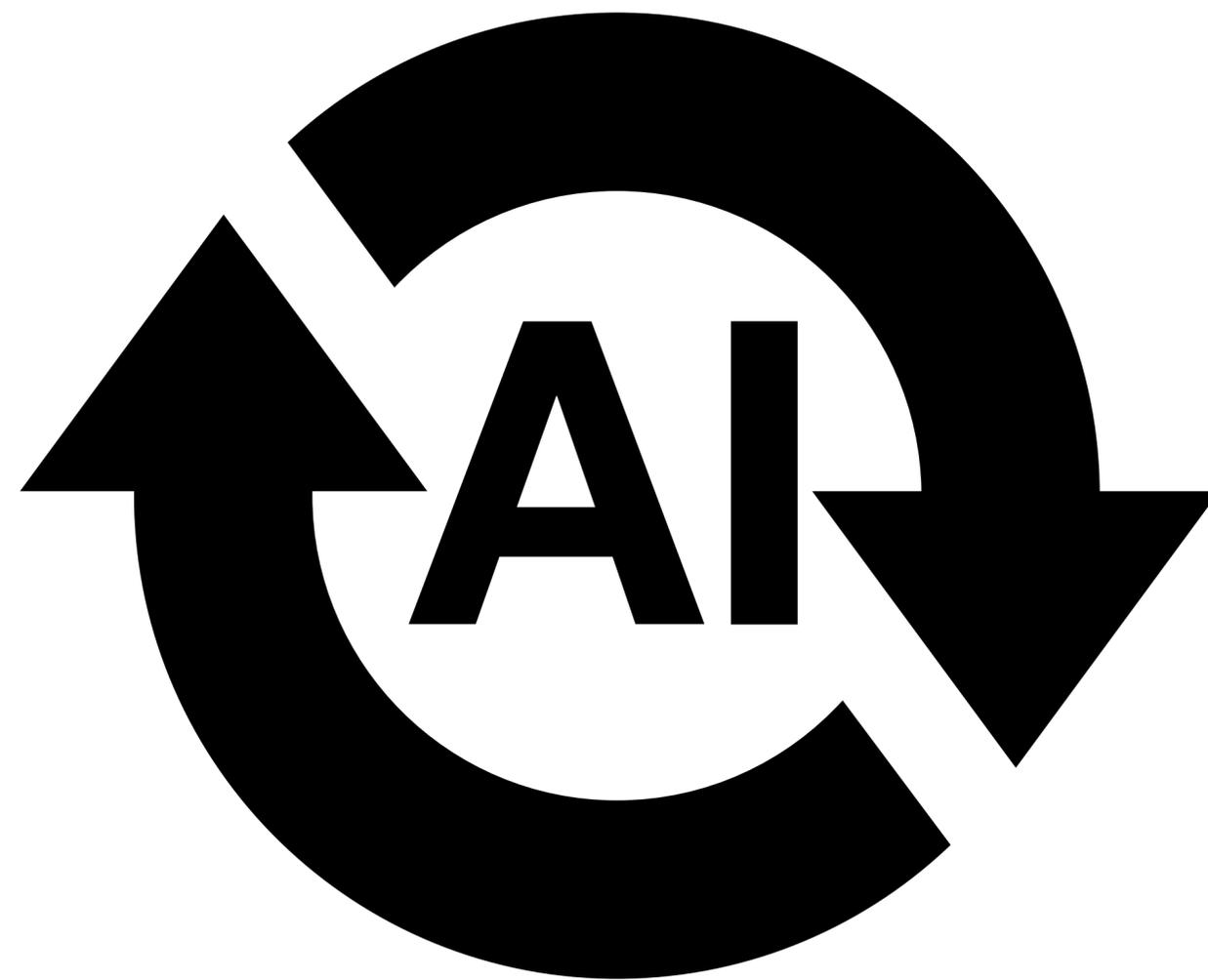
Data size



Model size

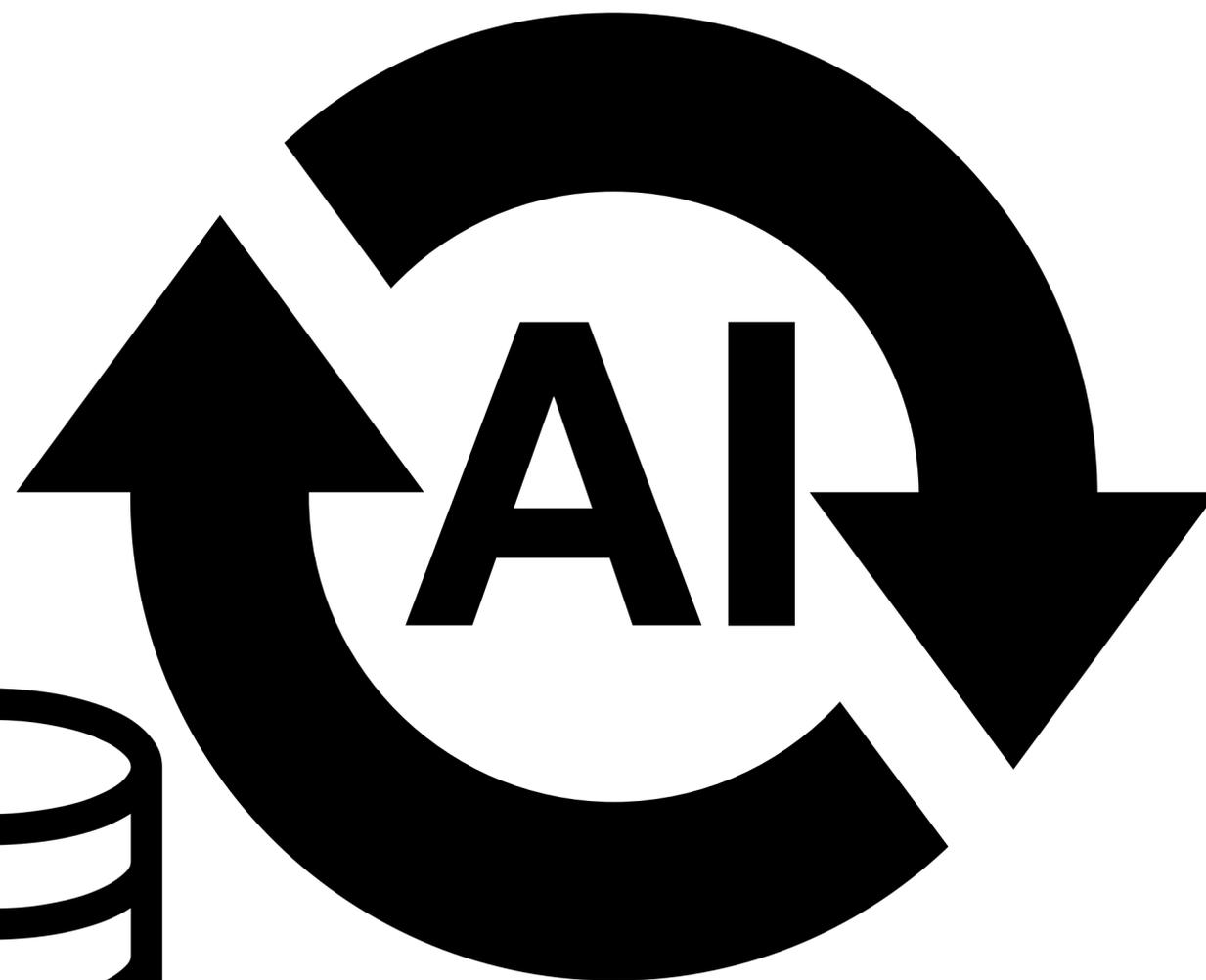
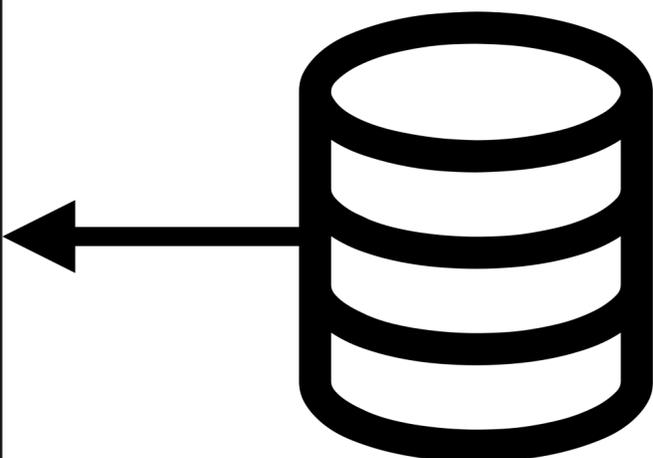
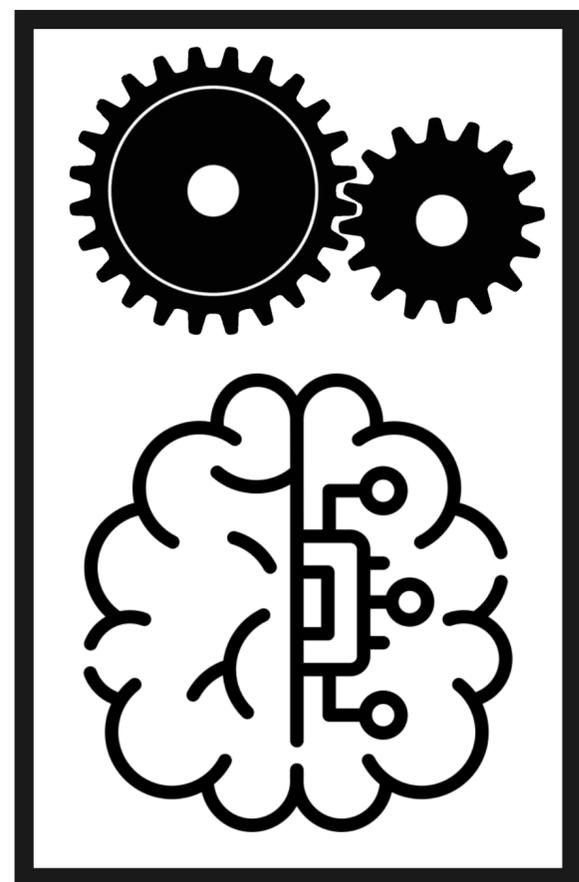


Training



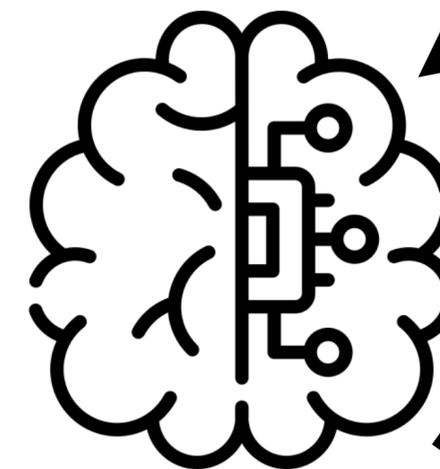
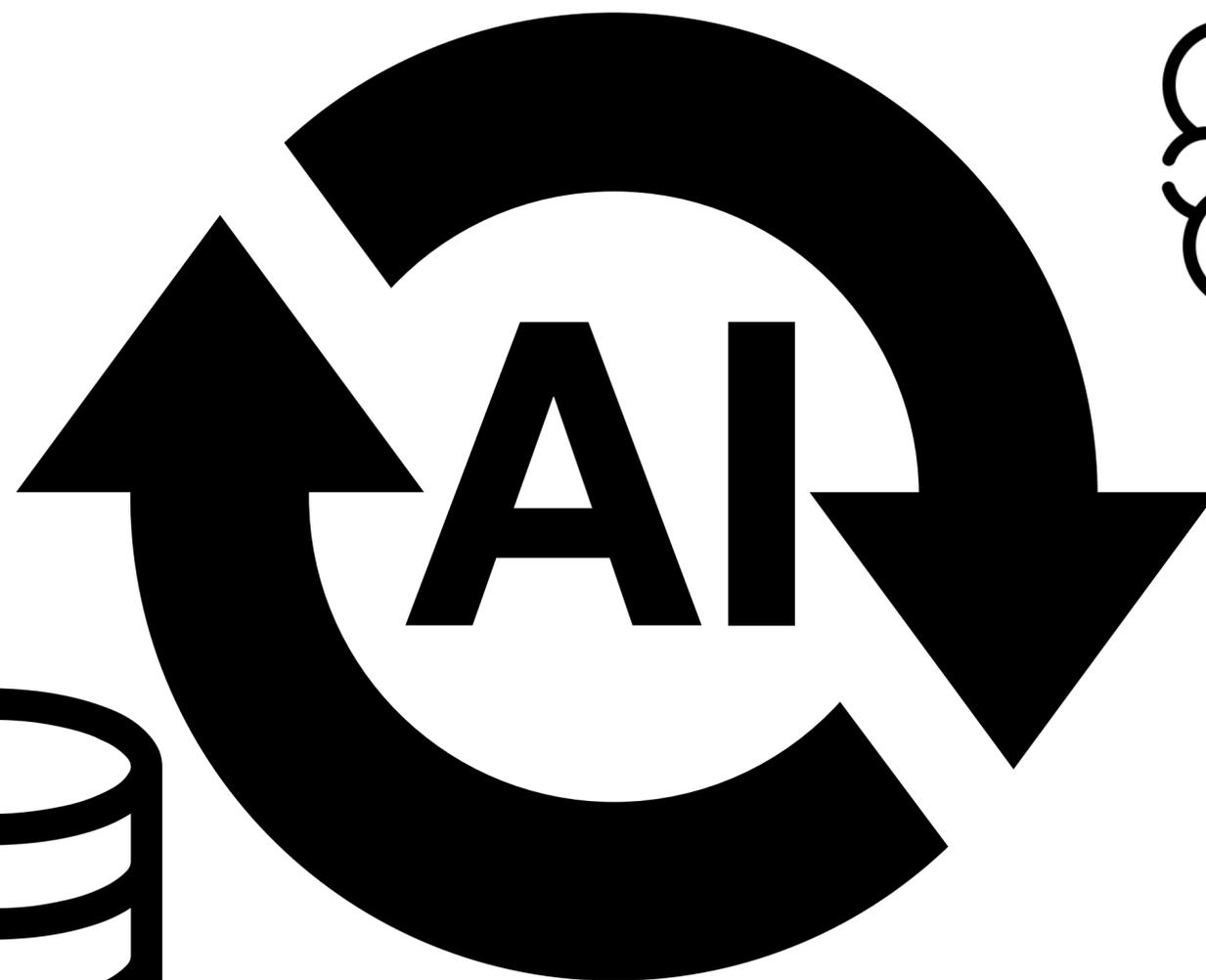
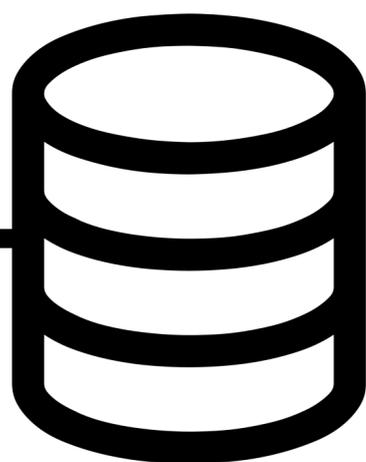
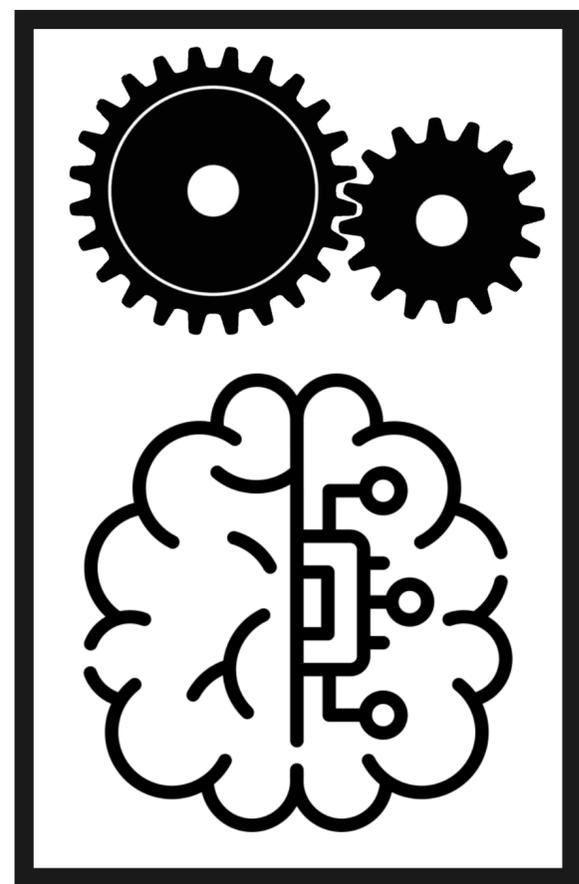
Inference

Training



Inference

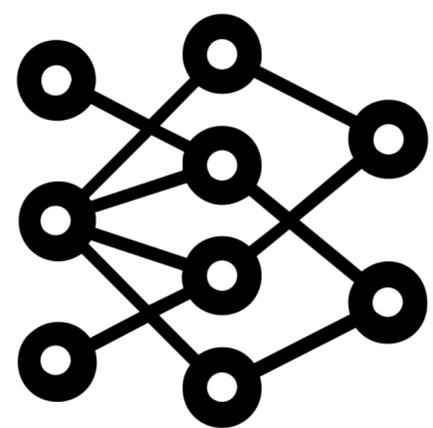
Training



What's this?

It's a dog!

Inference



Training

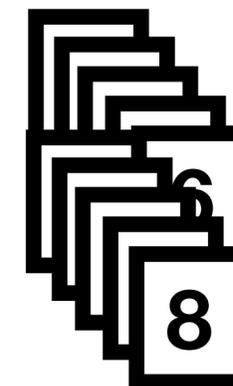


Inference



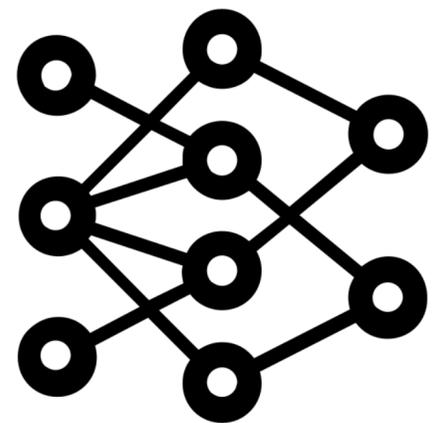
images

labels



Re-training





Training

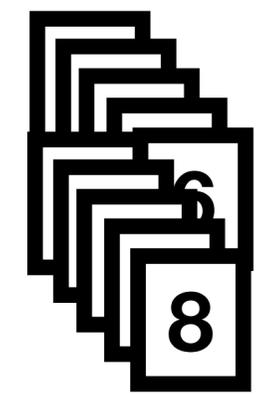


Inference



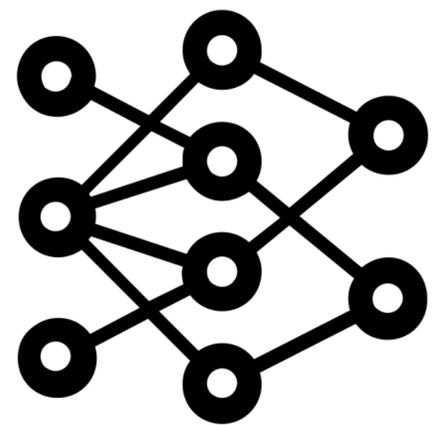
images

labels



Re-training





Training

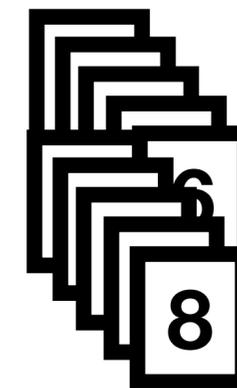


Inference



images

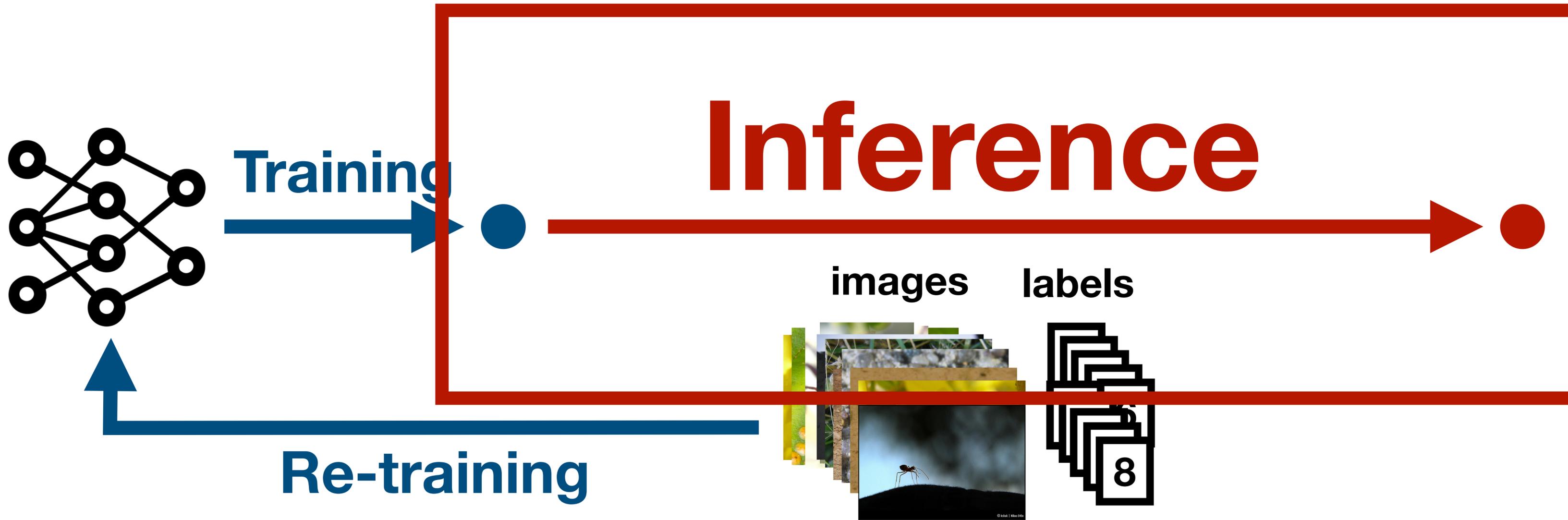
labels



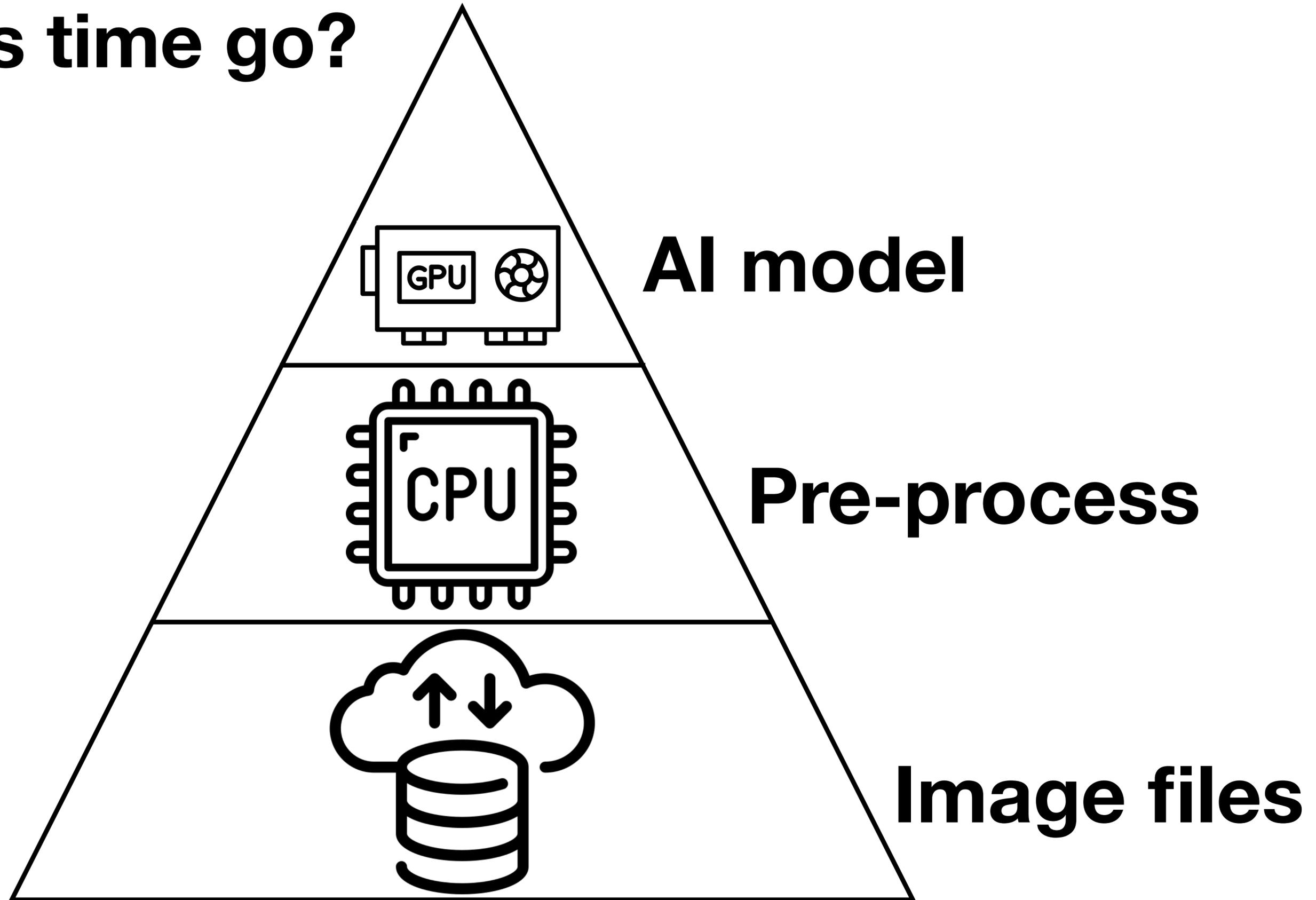
Re-training



90% of cost!

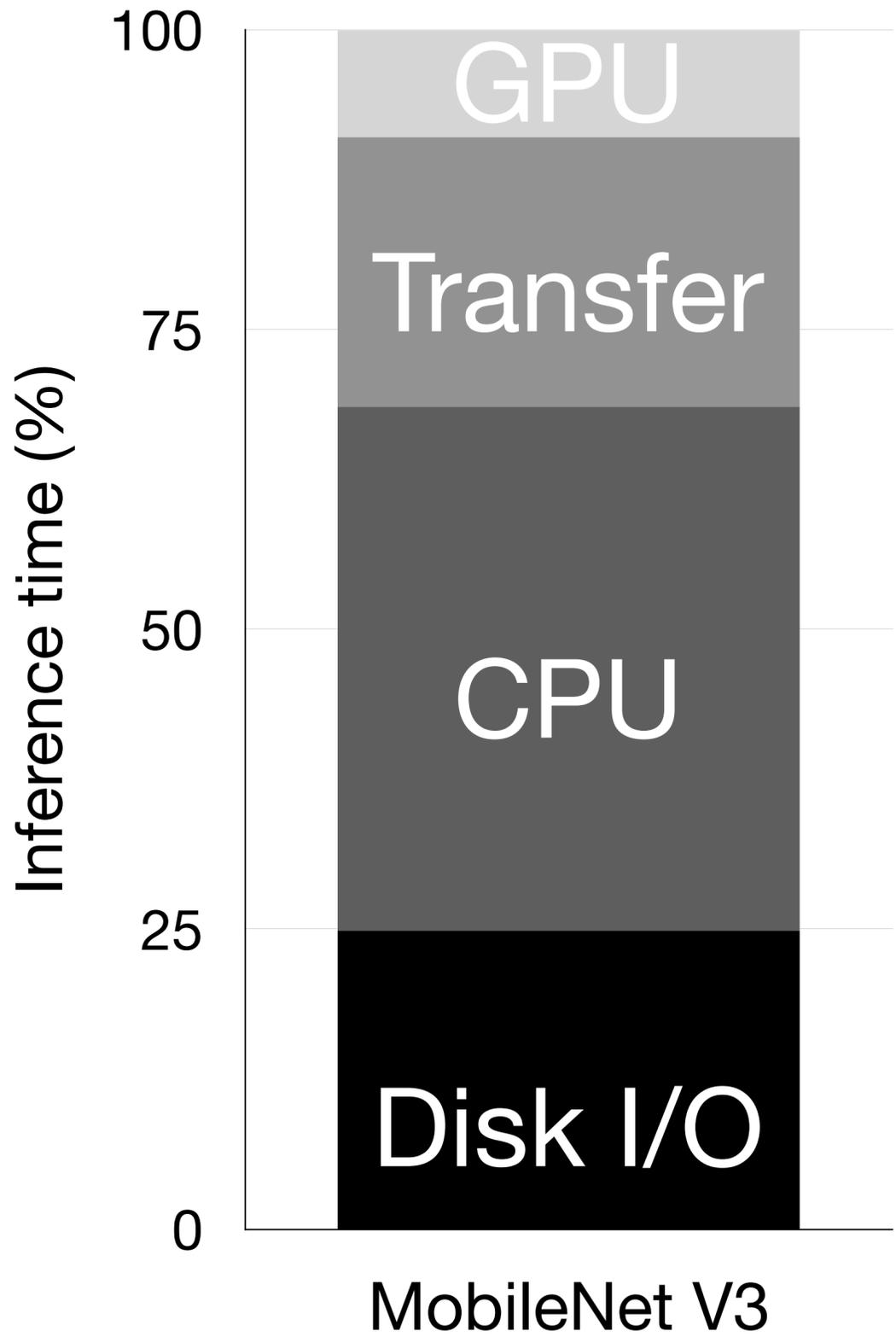


Where does time go?



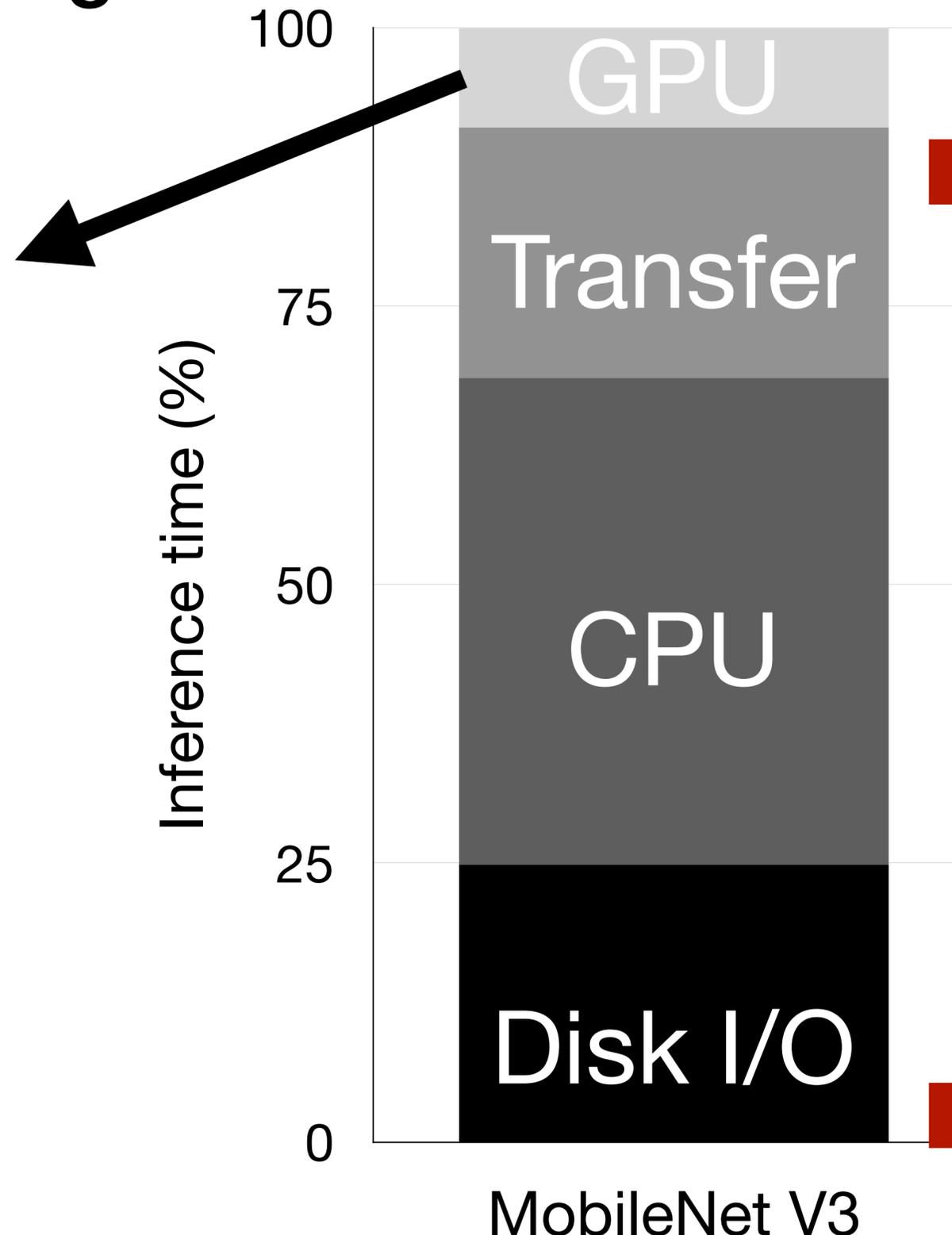
Where does time go?

Data: ImageNet
AI Model: MobileNet V3
Machine: V100, PCIe
Xeon, SSD
Framework: PyTorch v1



Where does time go?

**Only 10%
is GPU!**

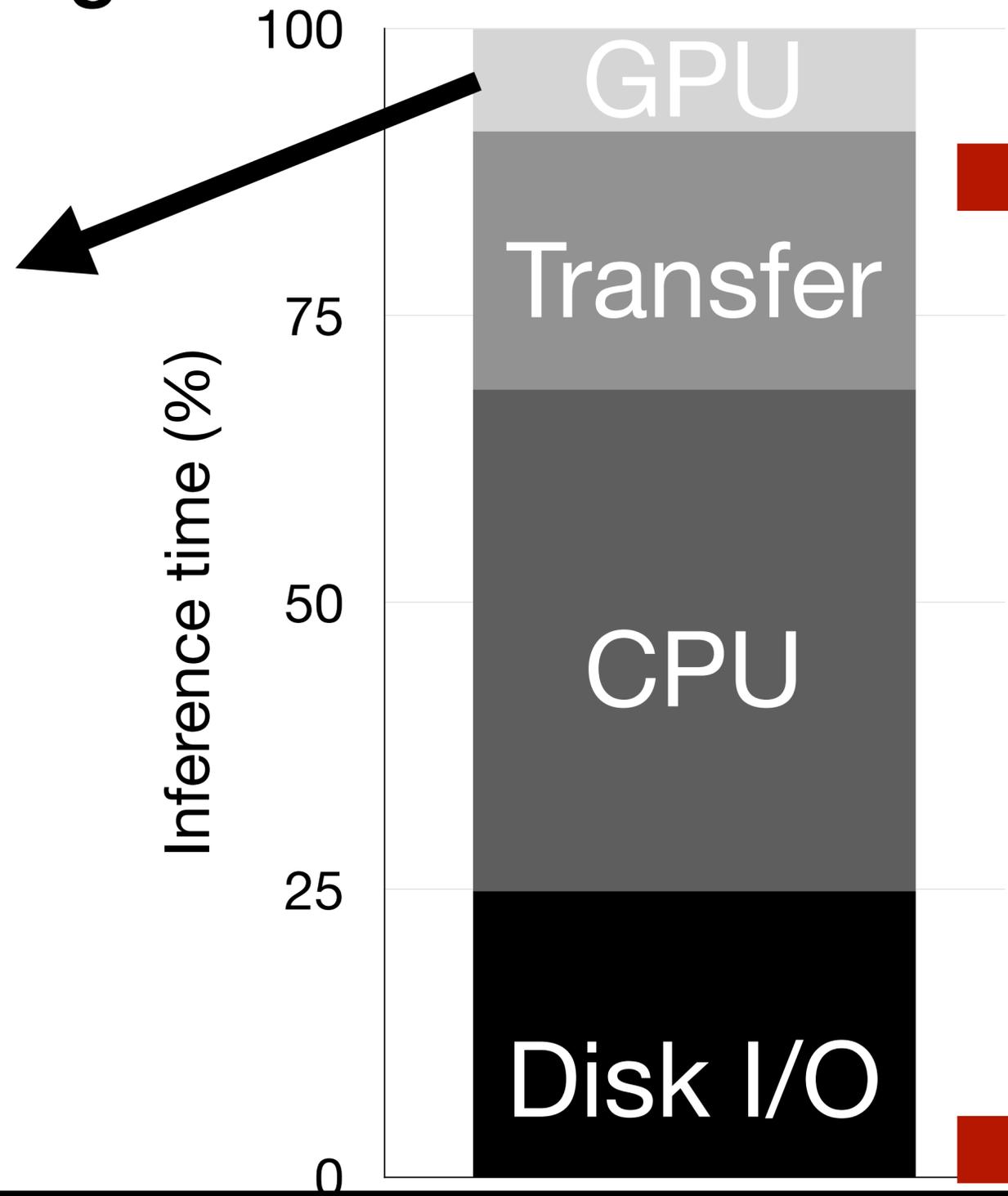


Data: ImageNet
AI Model: MobileNet V3
Machine: V100, PCIe Xeon, SSD
Framework: PyTorch v1

**Data movement/
pre-processing**

Where does time go?

**Only 10%
is GPU!**



Data: ImageNet
AI Model: MobileNet V3
Machine: V100, PCIe Xeon, SSD
Framework: PyTorch v1

**Data movement/
pre-processing**

Re-think Storage for Image AI

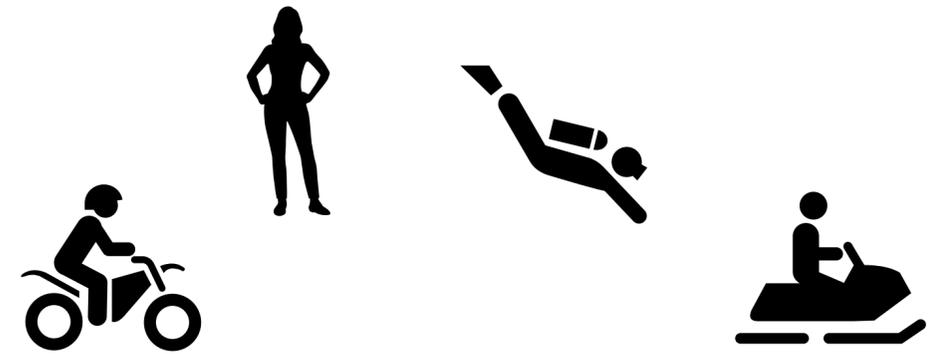
**How do machines
store images today?**

**How do machines
store images today?**

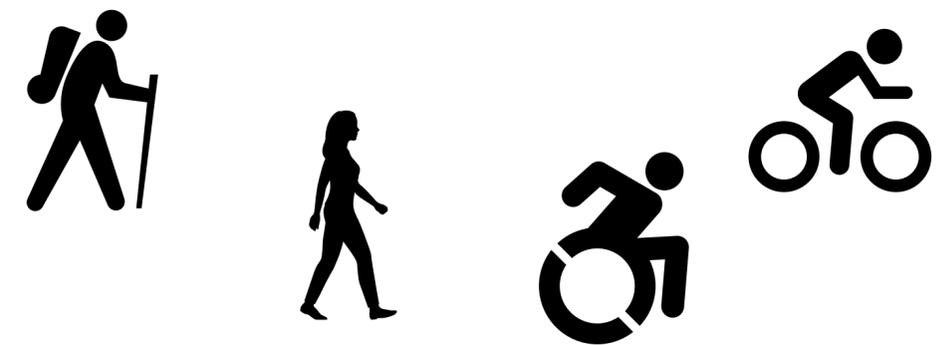
JPEG

Joint Photographic Experts Group

**How do machines
store images today?**

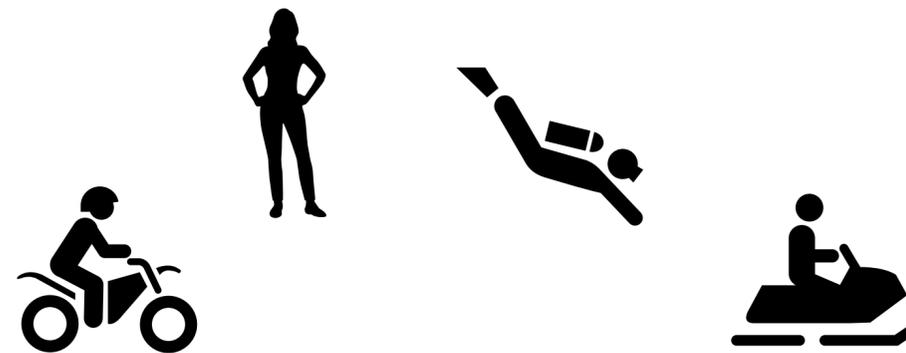


standard

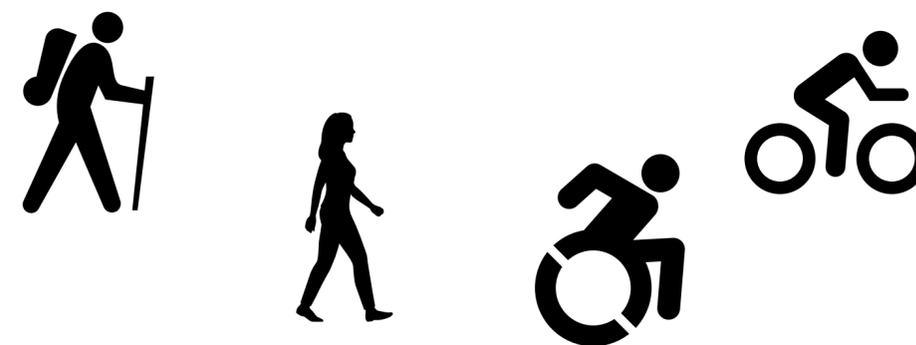




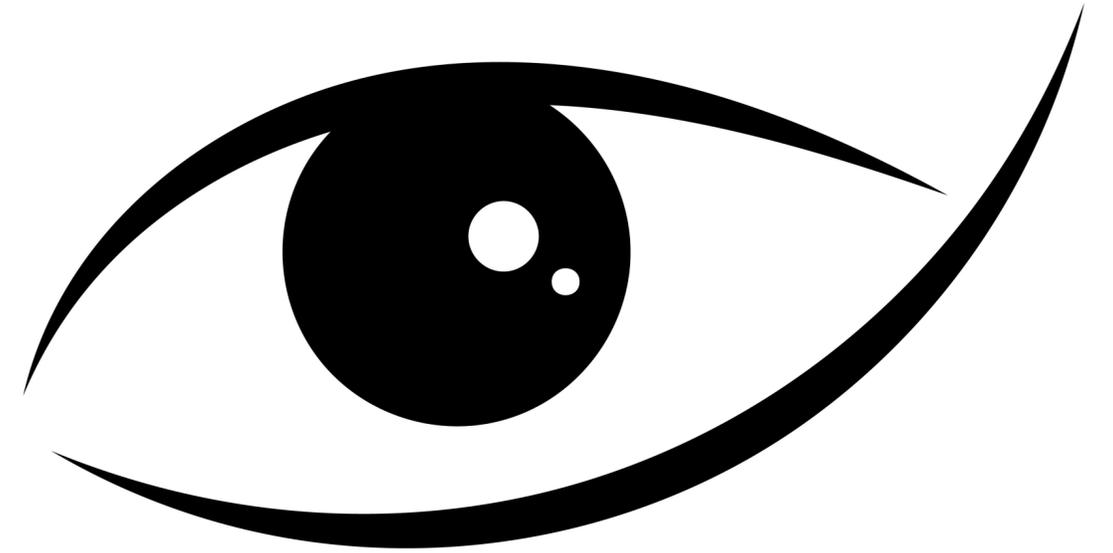
compression



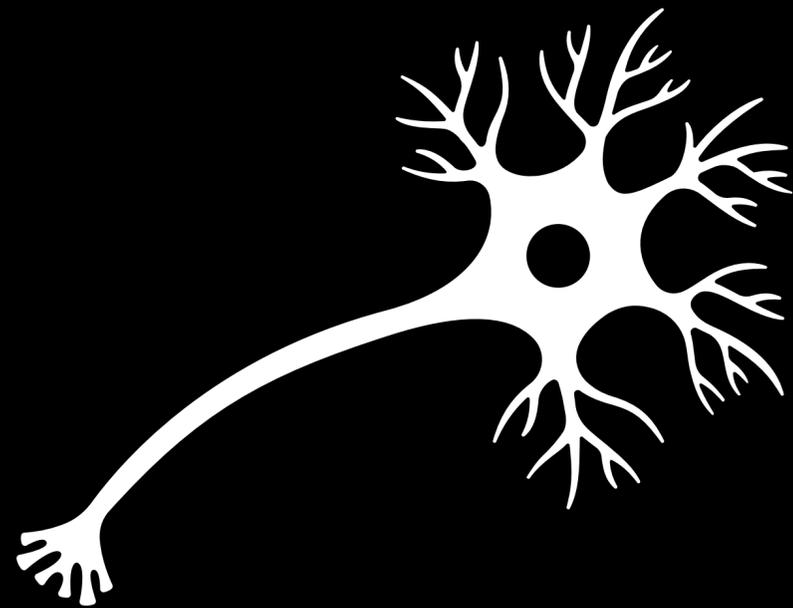
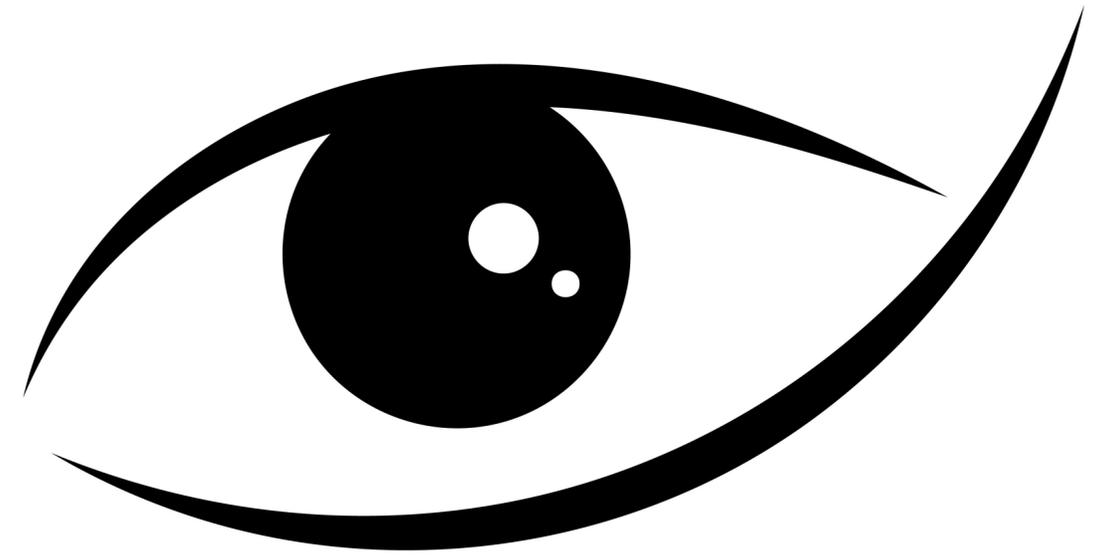
standard



JPEG is tailored for the
properties of the human eye



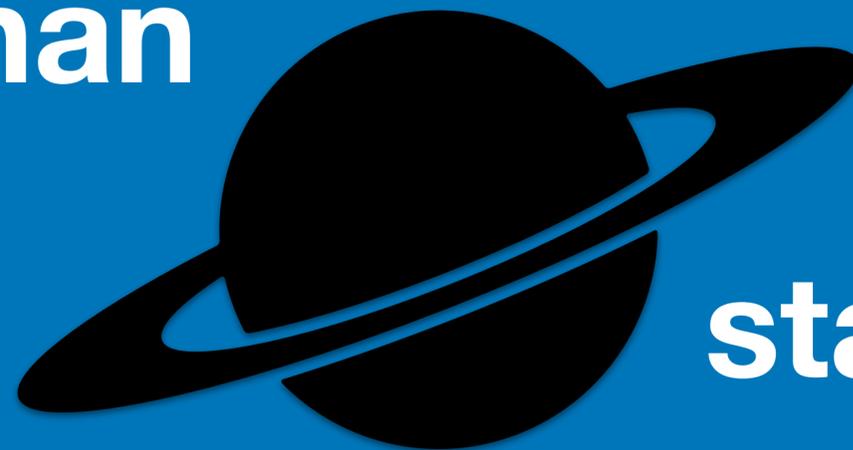
JPEG is tailored for the
properties of the human eye



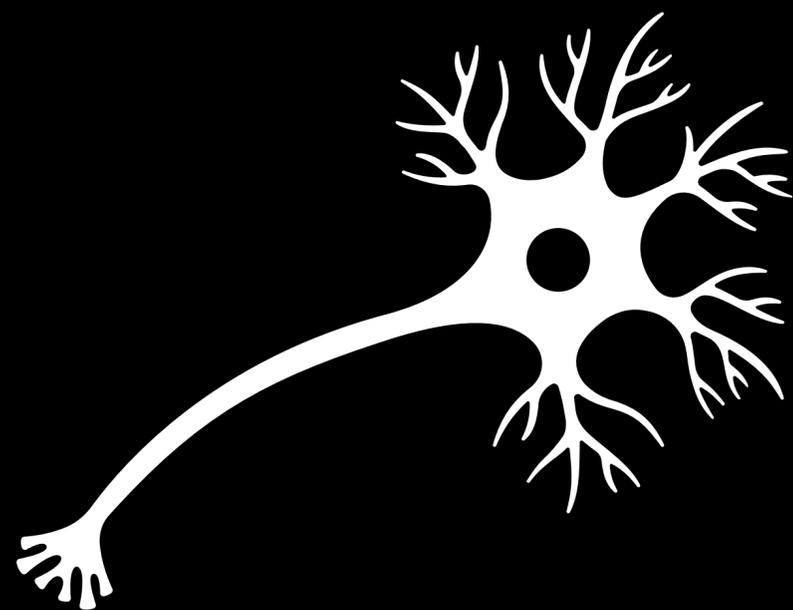
images for AI are seen by
algorithms, not humans

there are more possible ways
to store an image than

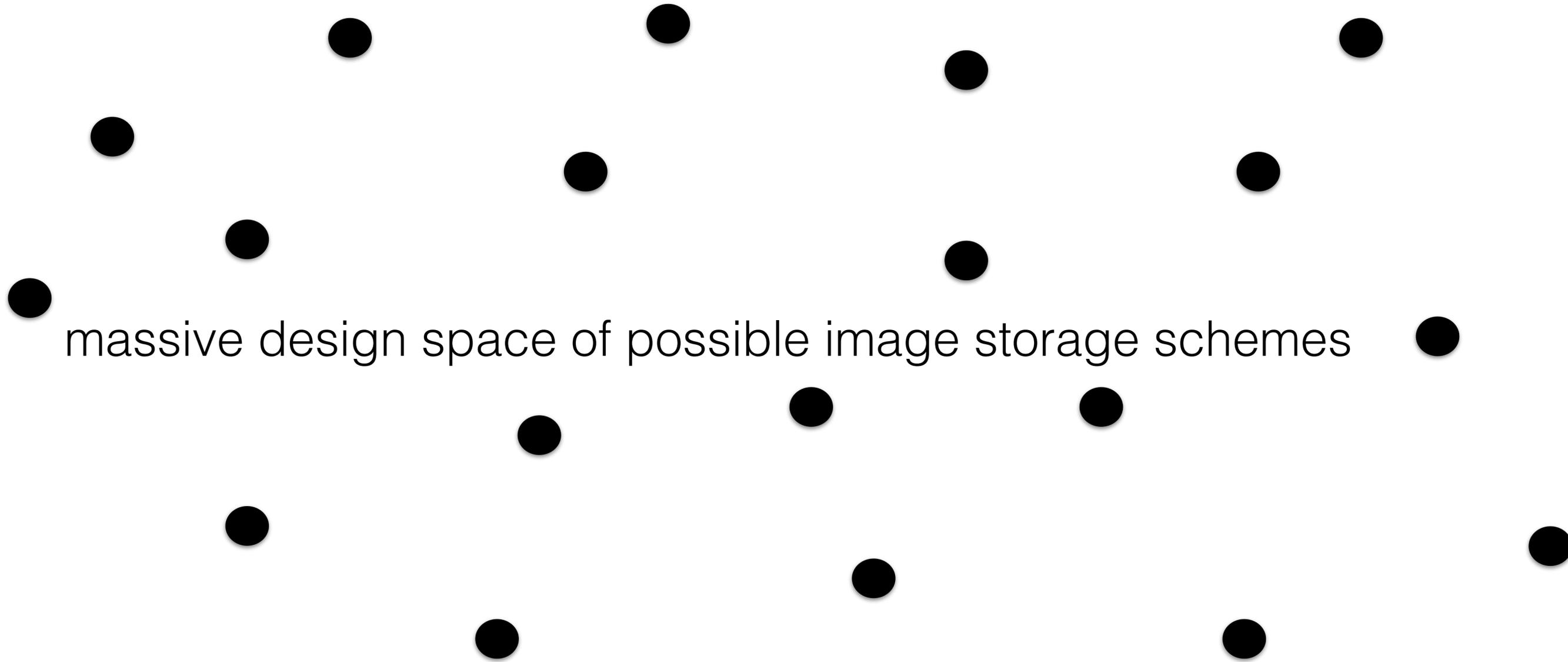
$$10^{100} > 10^{24}$$



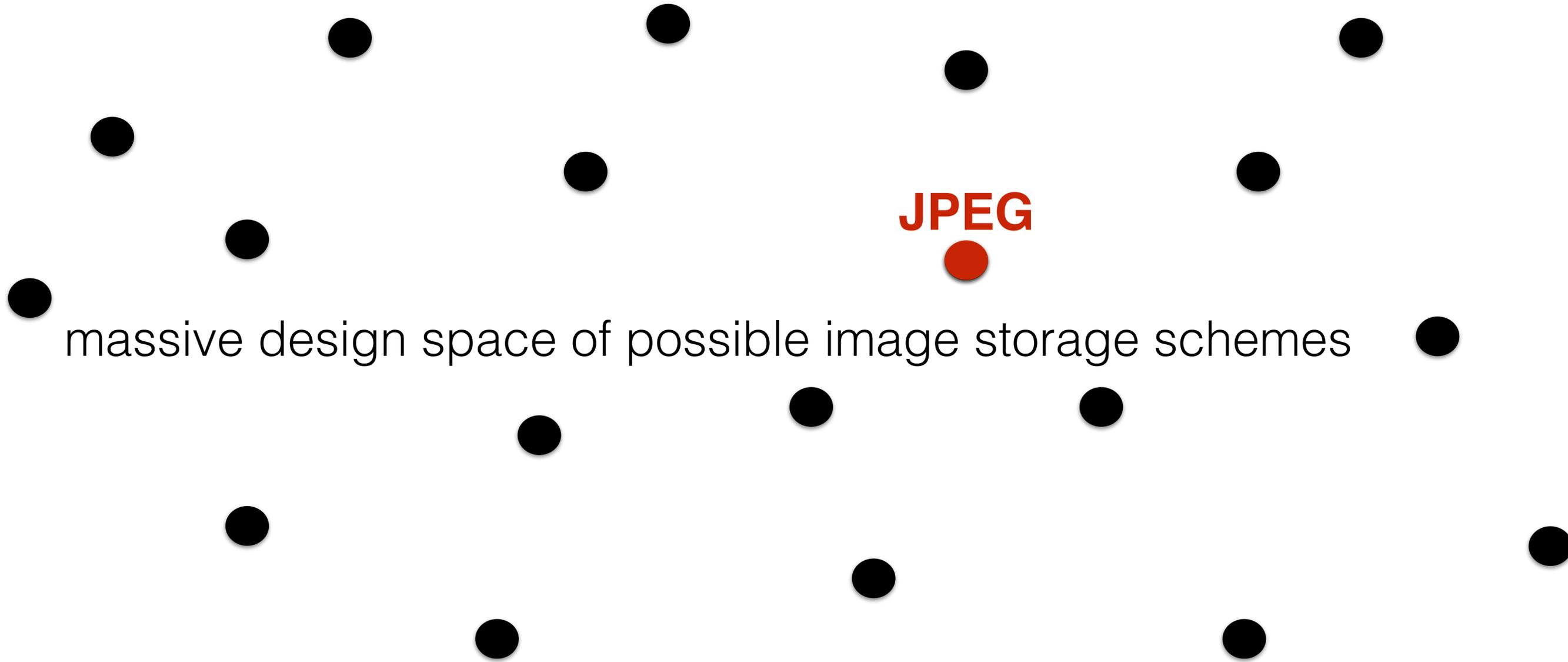
stars on the sky



images for AI are seen by
algorithms, not humans



massive design space of possible image storage schemes



JPEG

massive design space of possible image storage schemes

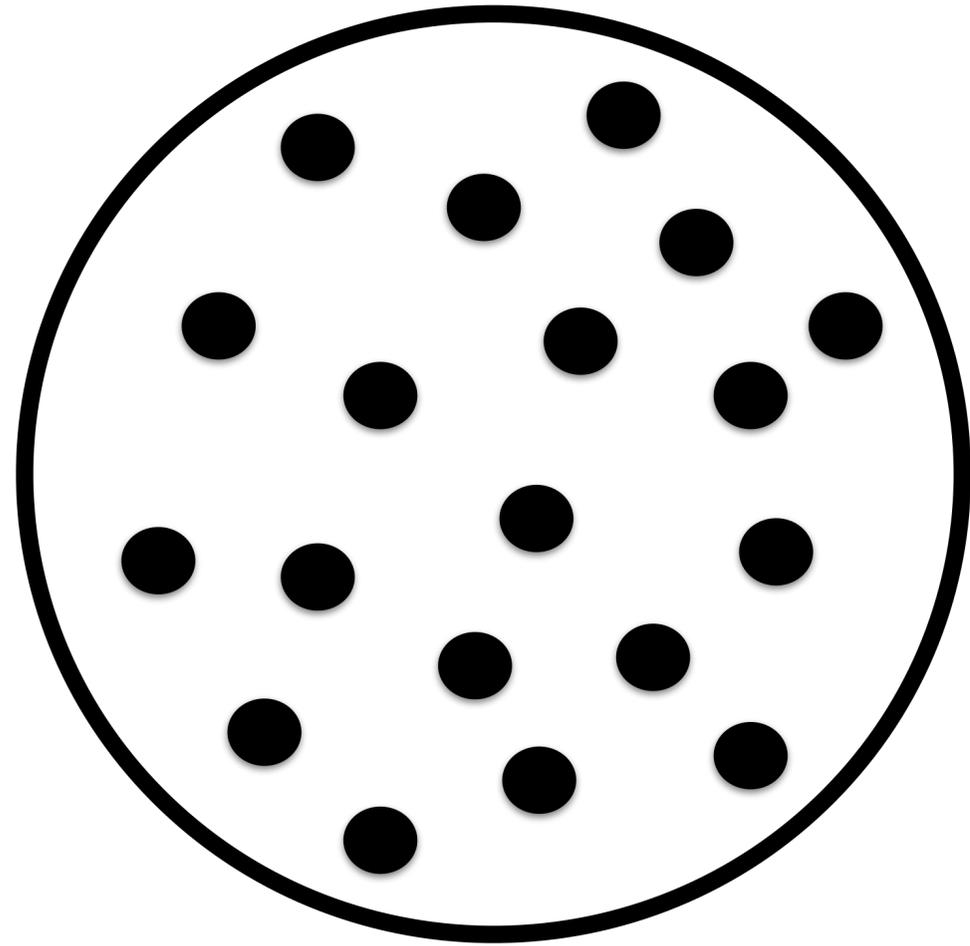
JPEG

massive design space of possible image storage schemes

accuracy cost



Design space



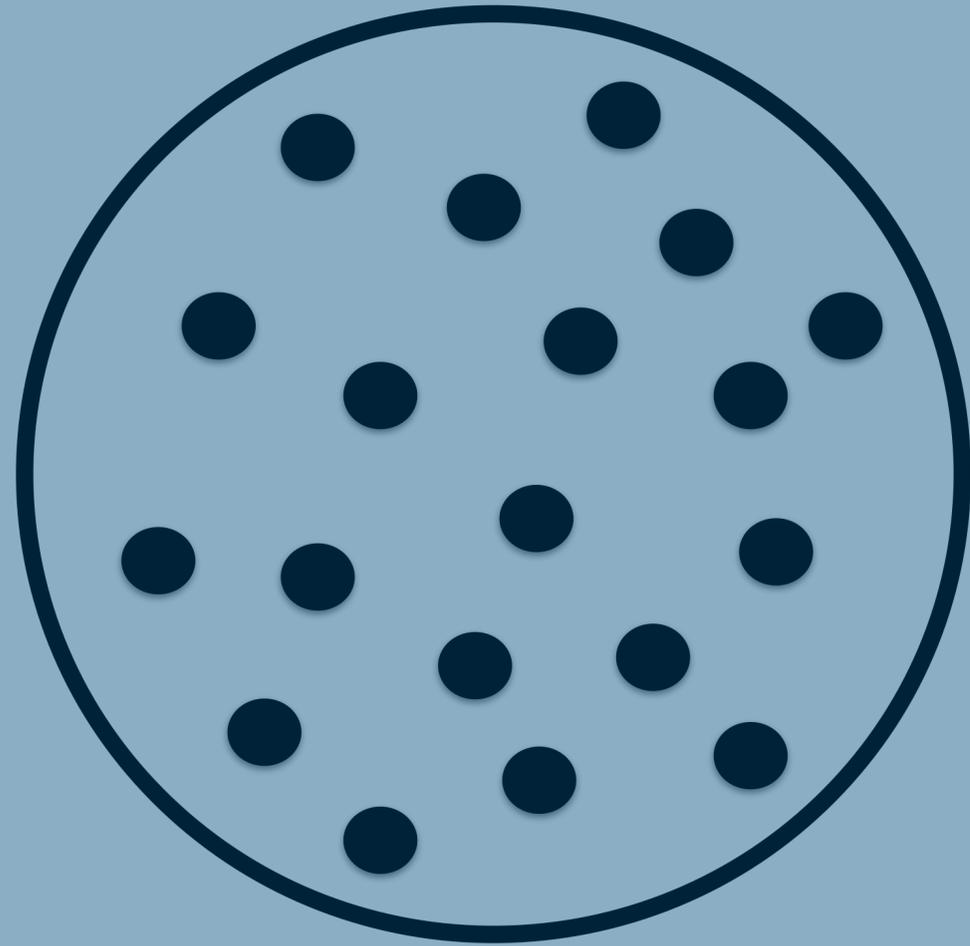
Search



14x faster



Design space



Search



14x faster



Sampling

Partitioning

Pruning

Quantization

Remove rows/
columns

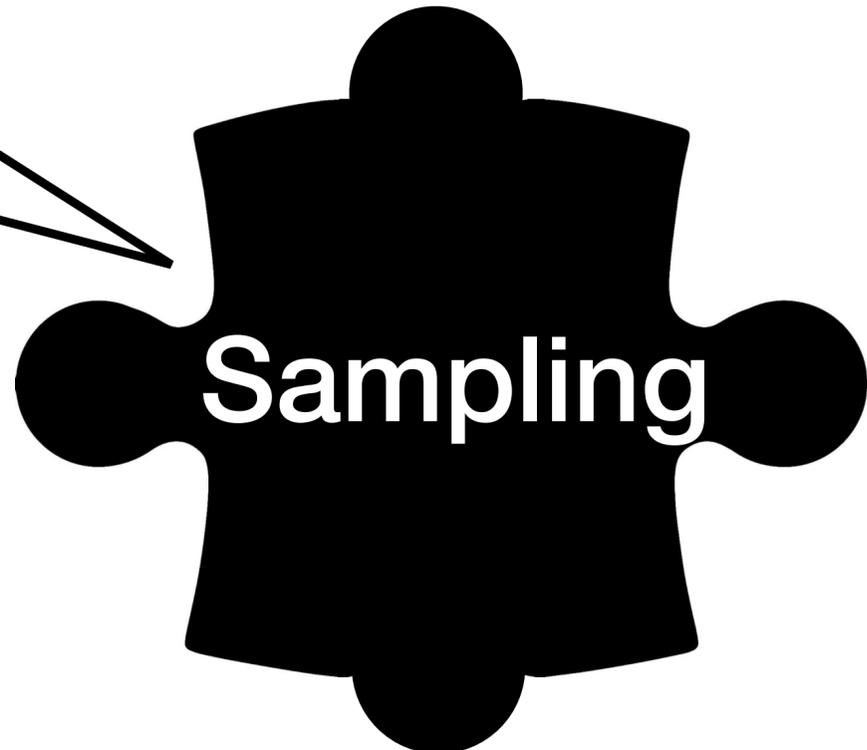
Sampling

Partitioning

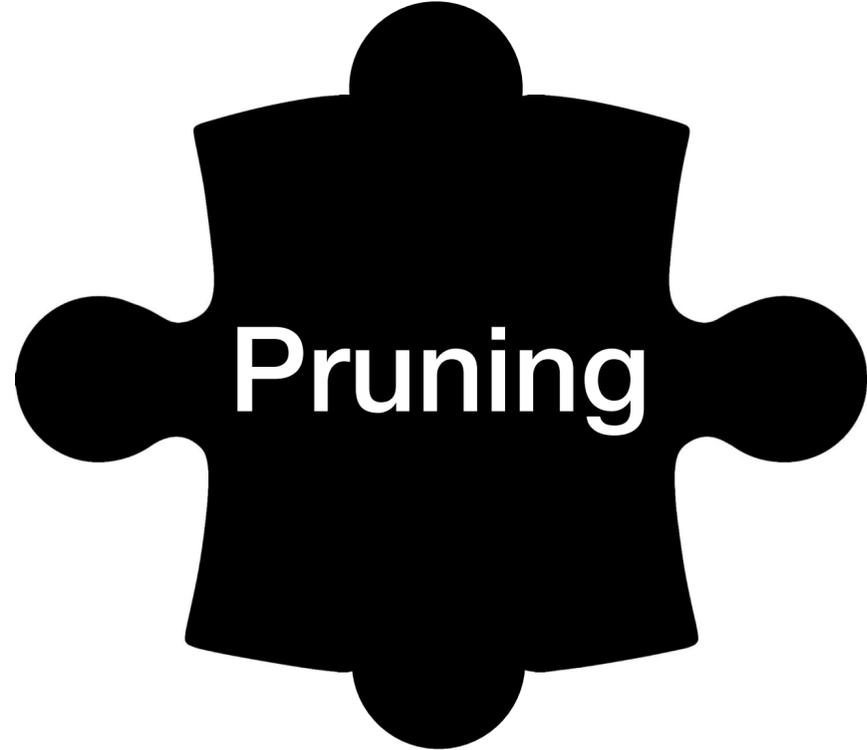
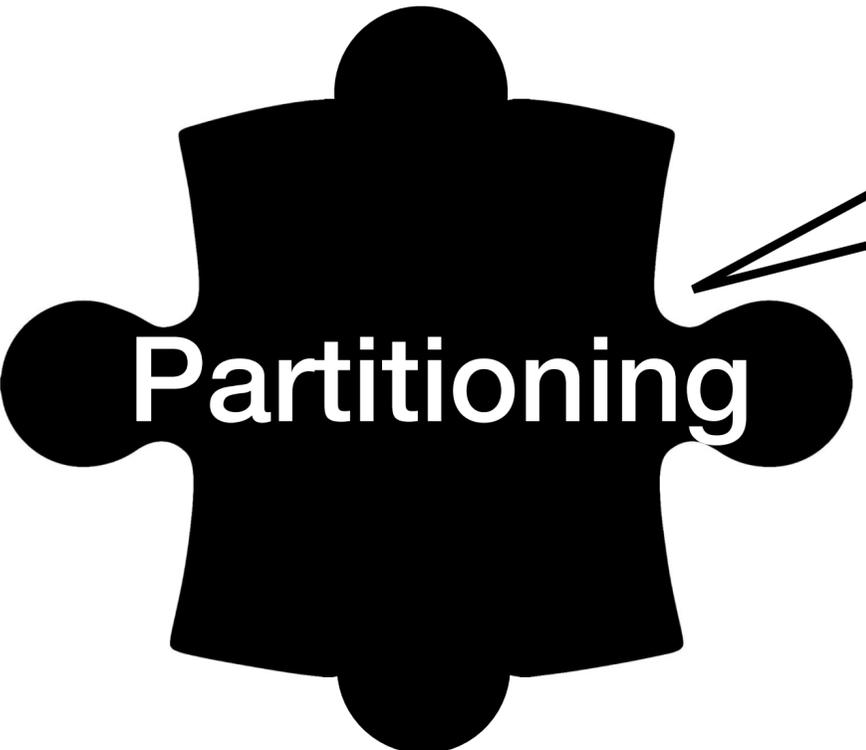
Pruning

Quantization

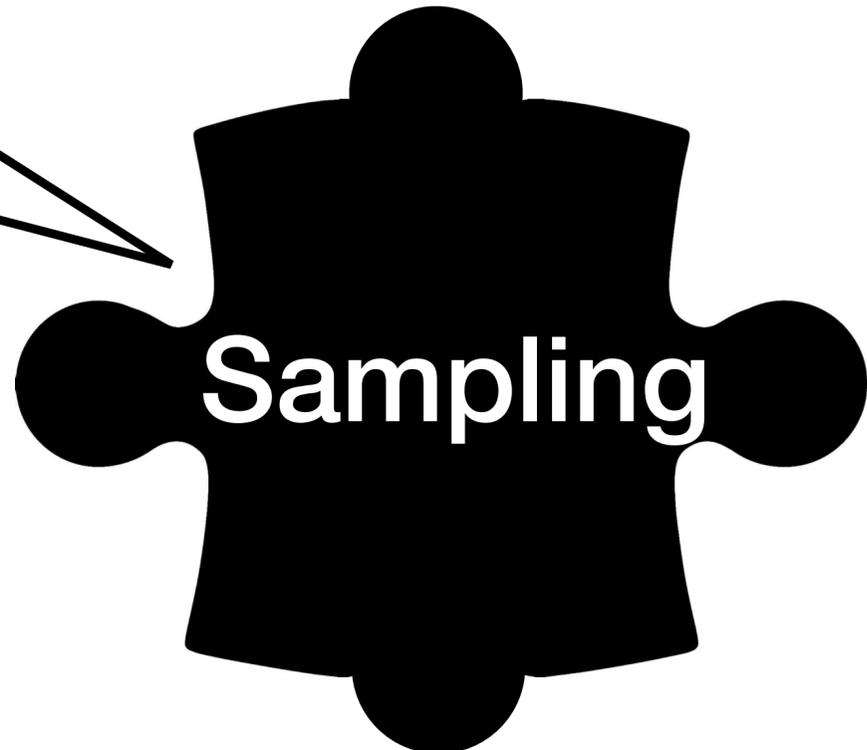
Remove rows/
columns



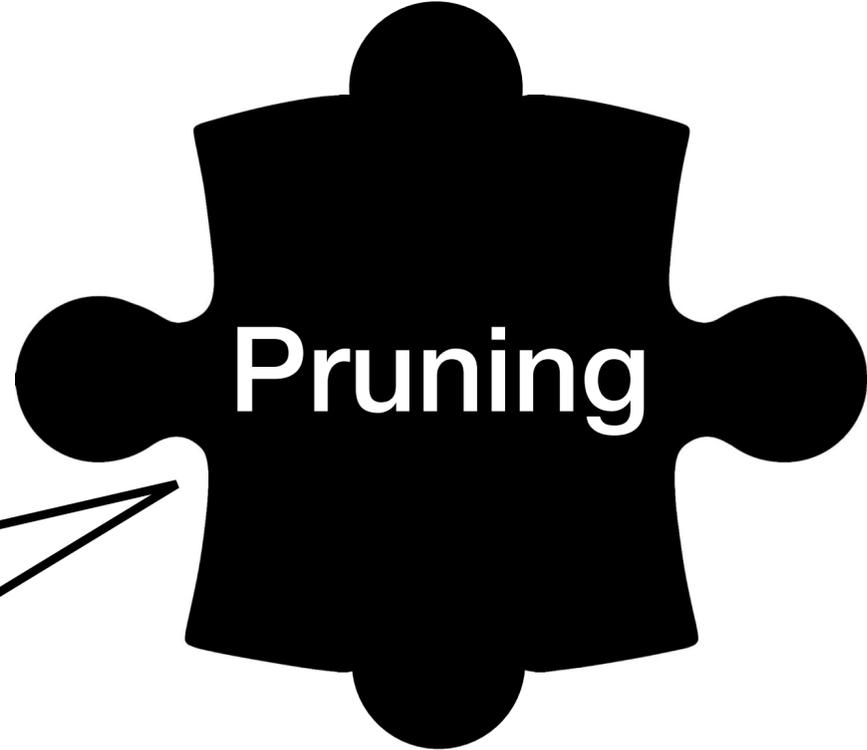
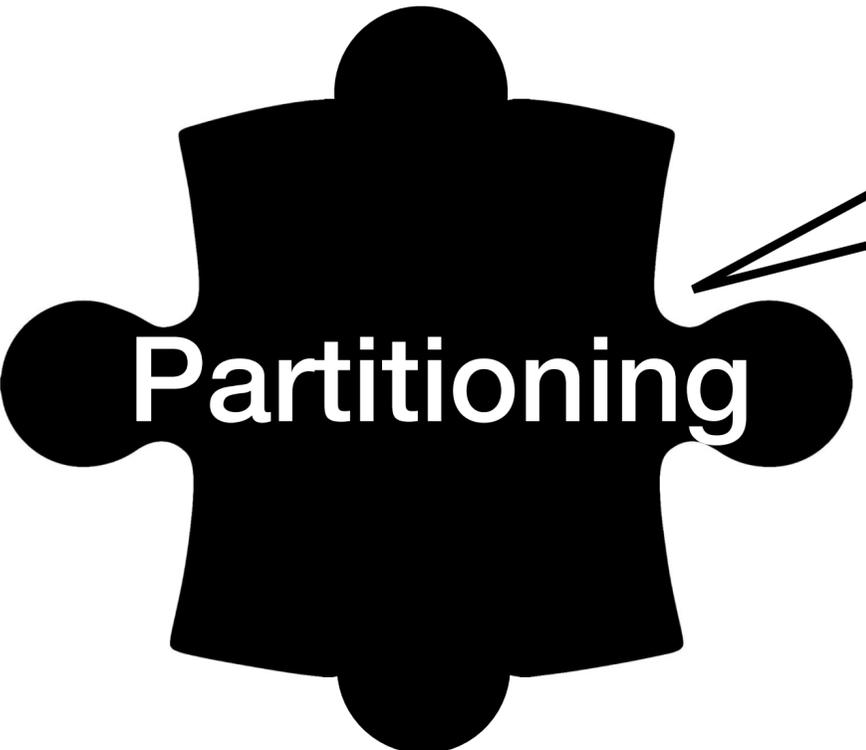
Processing
granularity



Remove rows/
columns

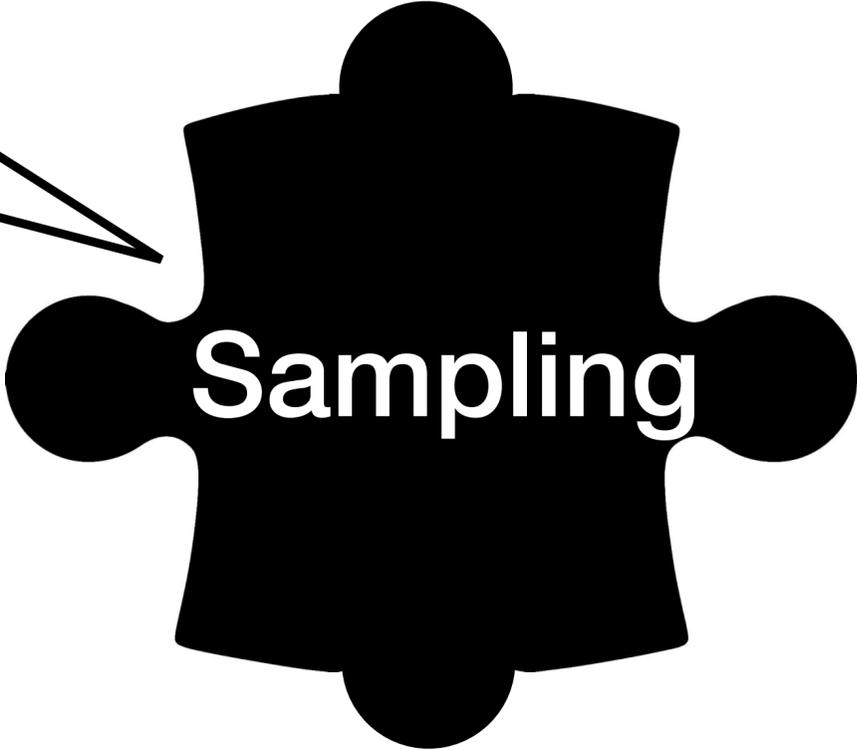


Processing
granularity

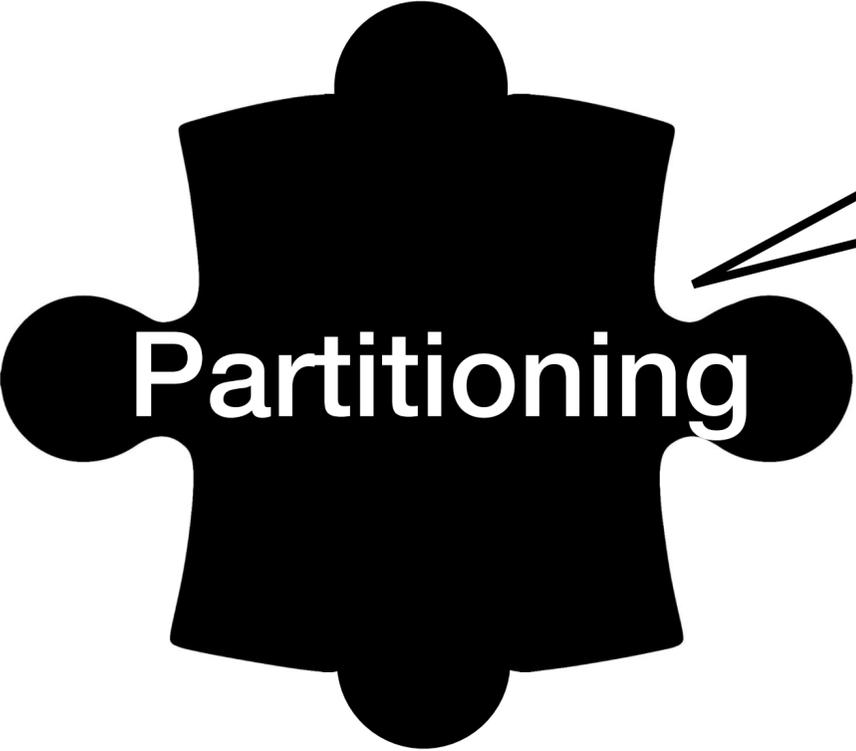


Remove
unuseful data

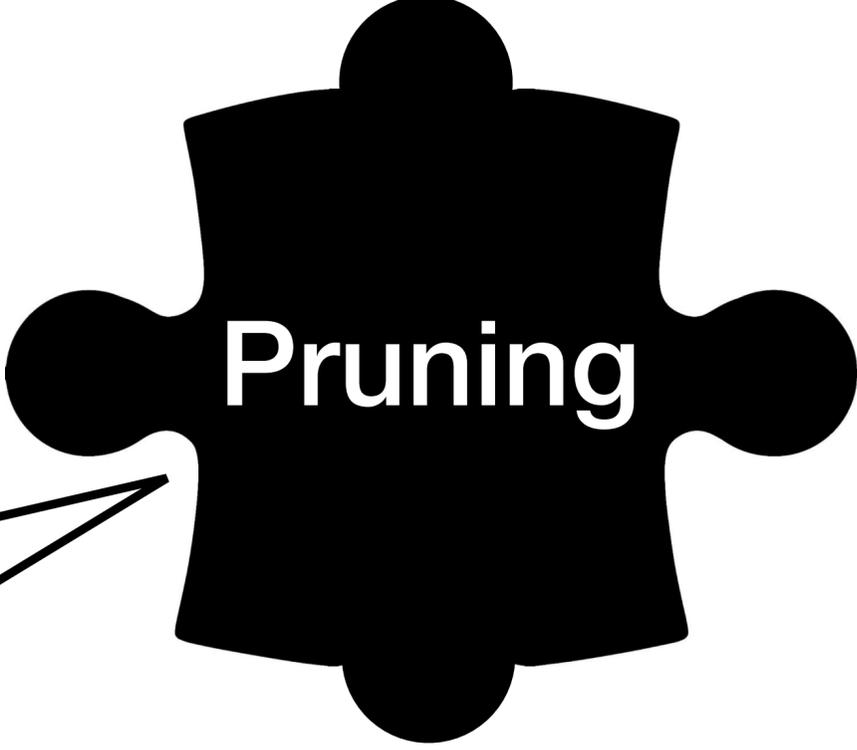
Remove rows/
columns



Processing
granularity



Remove
unuseful data



Magnitude
reduction



Remove rows/
columns

Processing
complexity

10^{150K}

possibilities

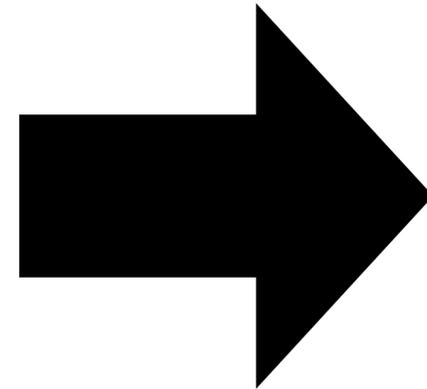
Remove
unuseful data

magnitude
reduction

**How can we prune the
design space?**

Pruning

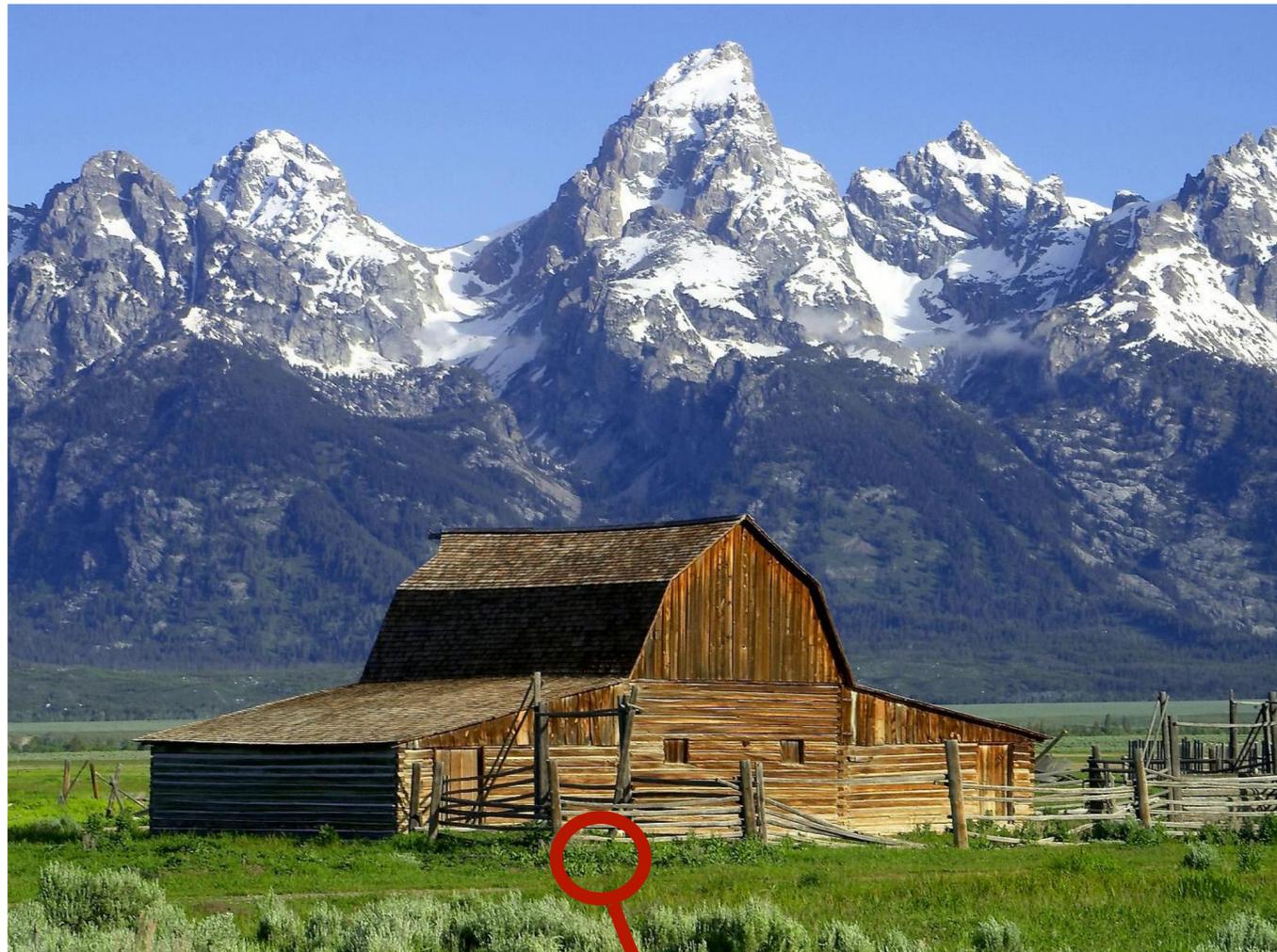
Spatial domain



Frequency domain

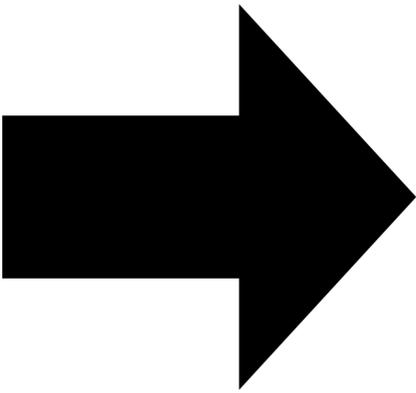
102	90	80	72	152	95	70	44	82	90	60	14	92	90	50	12
76	89	72	39	99	70	65	45	80	65	40	21	82	60	32	20
56	58	49	30	78	58	30	20	77	40	24	20	60	40	30	4
54	49	32	5	40	32	20	2	70	60	59	30	20	10	8	8
200	150	100	98	180	154	134	102	90	86	82	76	120	102	64	50
180	172	160	150	120	134	103	95	85	80	72	45	90	82	70	54
170	165	120	112	103	90	83	74	70	63	53	45	80	65	55	30
92	90	70	49	94	89	80	24	60	54	40	30	70	50	30	22
112	108	100	78	90	85	80	78	110	95	85	78	104	92	82	70
97	85	65	48	82	70	65	49	98	90	76	40	95	80	70	35
90	78	74	30	78	65	54	32	68	49	31	20	68	52	32	24
72	45	24	5	74	45	30	21	31	20	10	4	52	32	15	2

Spatial domain



A pixel

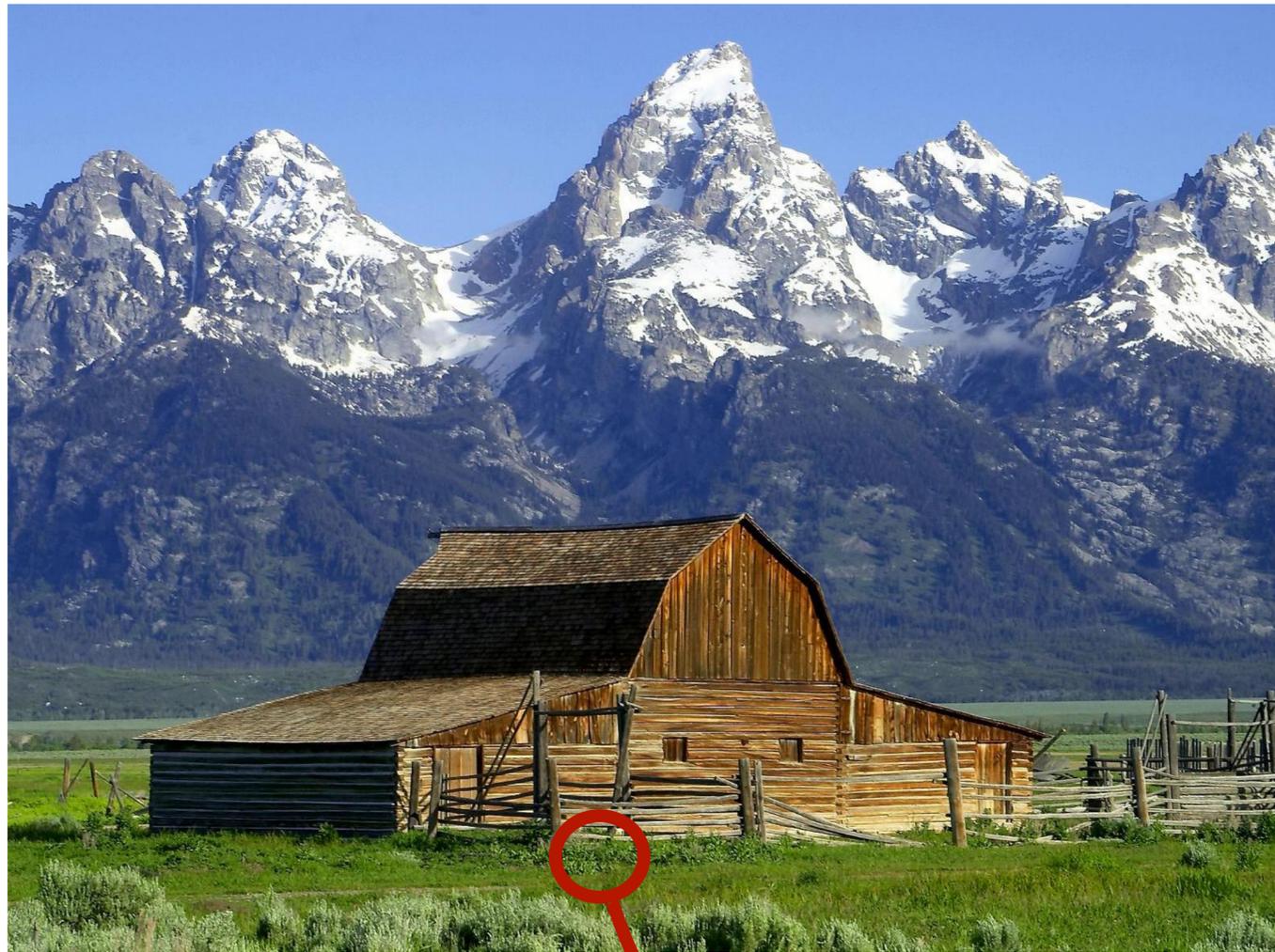
Frequency domain



102	90	80	72	152	95	70	44	82	90	60	14	92	90	50	12
76	89	72	39	99	70	65	45	80	65	40	21	82	60	32	20
56	58	49	30	78	58	30	20	77	40	24	20	60	40	30	4
54	49	32	5	40	32	20	2	70	60	59	30	20	10	8	8
200	150	100	98	180	154	134	102	90	86	82	76	120	102	64	50
180	172	160	150	120	134	103	95	85	80	72	45	90	82	70	54
170	165	120	112	103	90	83	74	70	63	53	45	80	65	55	30
92	90	70	49	94	89	80	24	60	54	40	30	70	50	30	22
112	108	100	78	90	85	80	78	110	95	85	78	104	92	82	70
97	85	65	48	82	70	65	49	98	90	76	40	95	80	70	35
90	78	74	30	78	65	54	32	68	49	31	20	68	52	32	24
72	45	24	5	74	45	30	21	31	20	10	4	52	32	15	2

A frequency coefficient

Spatial domain



A pixel

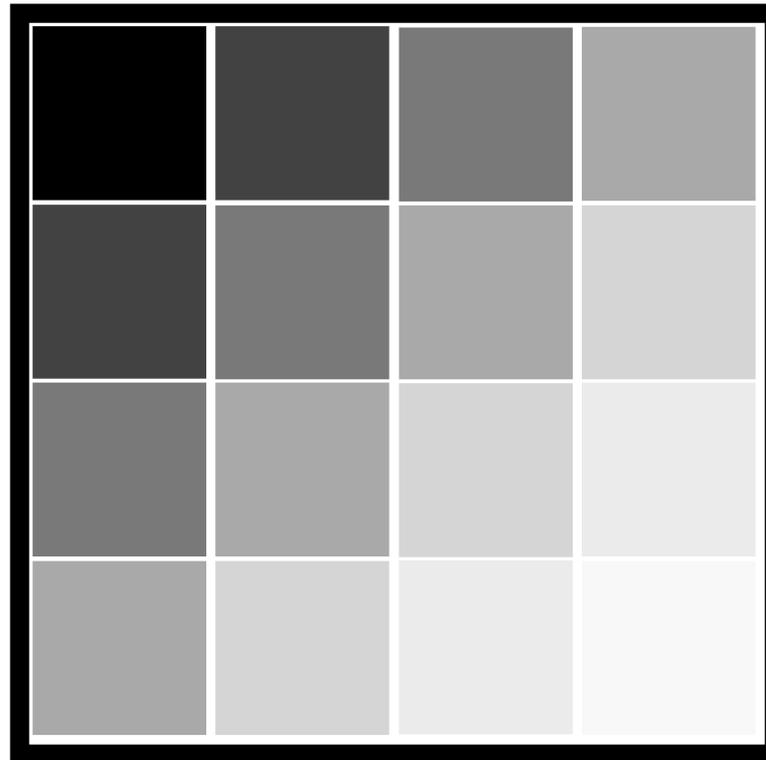
A block

Frequency domain

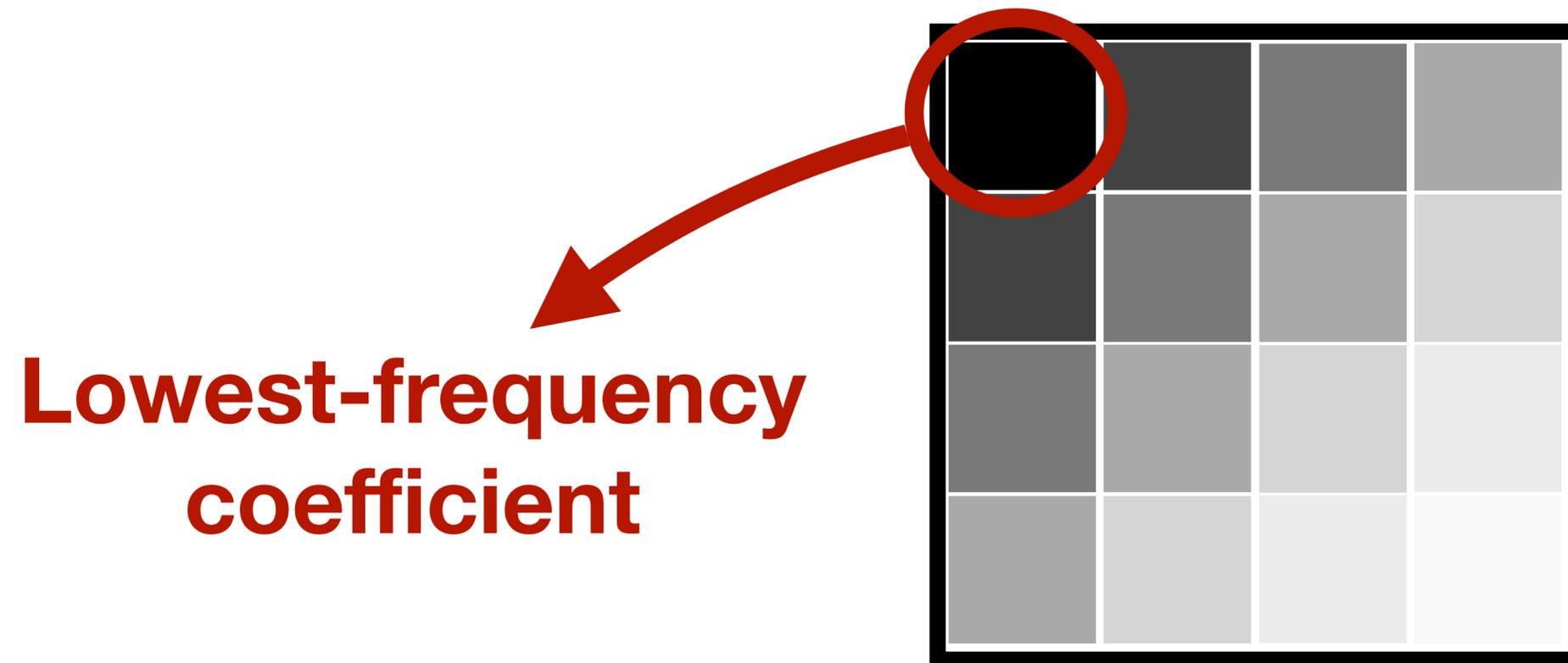
102	90	80	72	52	95	70	44	82	90	60	14	92	90	50	12
76	89	72	39	99	70	65	45	80	65	40	21	82	60	32	20
56	58	49	30	78	58	30	20	77	40	24	20	60	40	30	4
54	49	32	5	40	32	20	2	70	60	59	30	20	10	8	8
200	150	100	98	180	154	134	102	90	86	82	76	120	102	64	50
180	172	160	150	120	134	103	95	85	80	72	45	90	82	70	54
170	165	120	112	103	90	83	74	70	63	53	45	80	65	55	30
92	90	70	49	94	89	80	24	60	54	40	30	70	50	30	22
112	108	100	78	90	85	80	78	110	95	85	78	104	92	82	70
97	85	65	48	82	70	65	49	98	90	76	40	95	80	70	35
90	78	74	30	78	65	54	32	68	49	31	20	68	52	32	24
72	45	24	5	74	45	30	21	31	20	10	4	52	32	15	2

A frequency coefficient

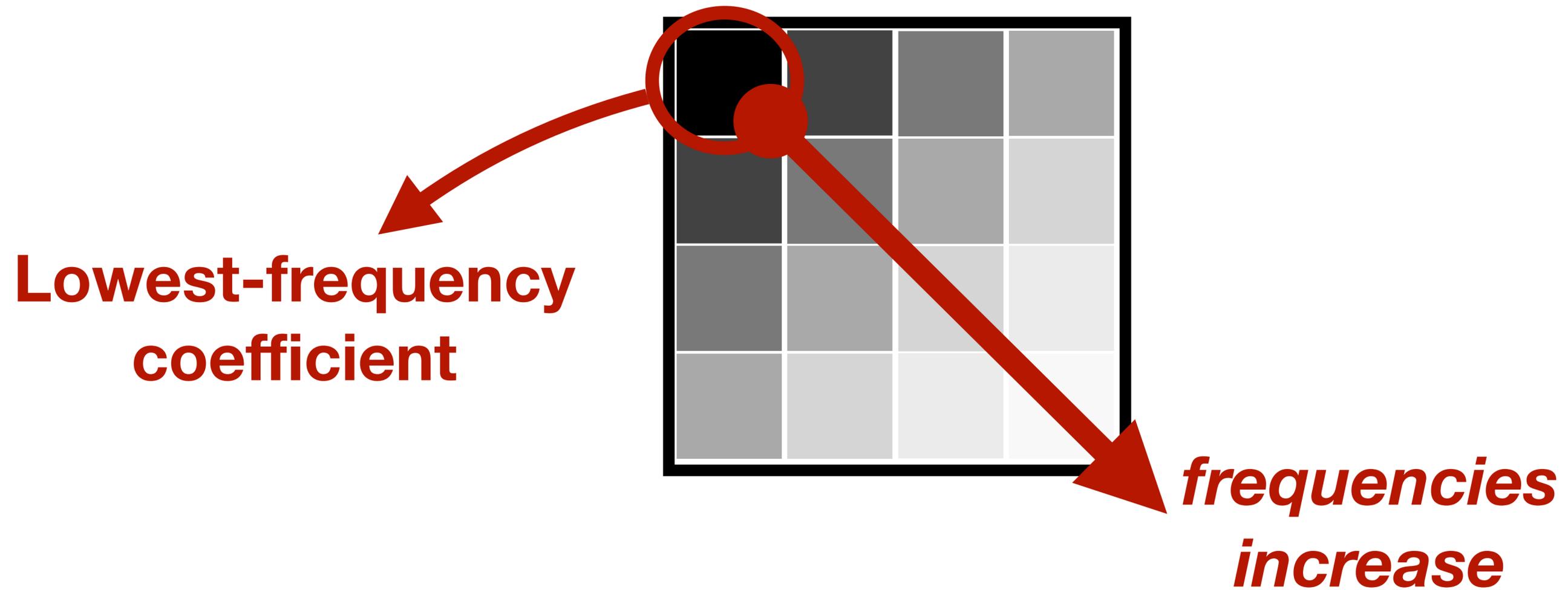
A single block



A single block

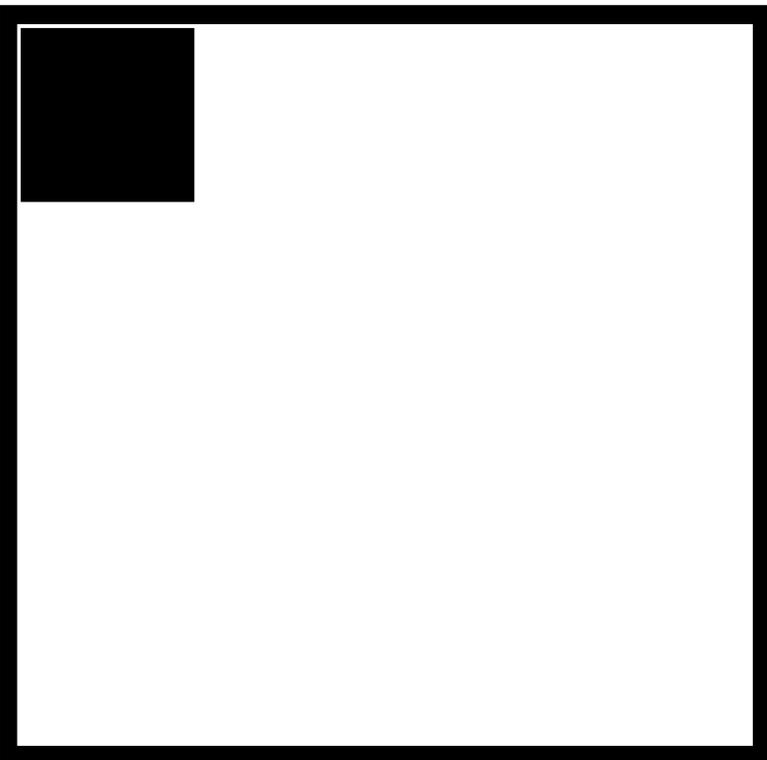


A single block

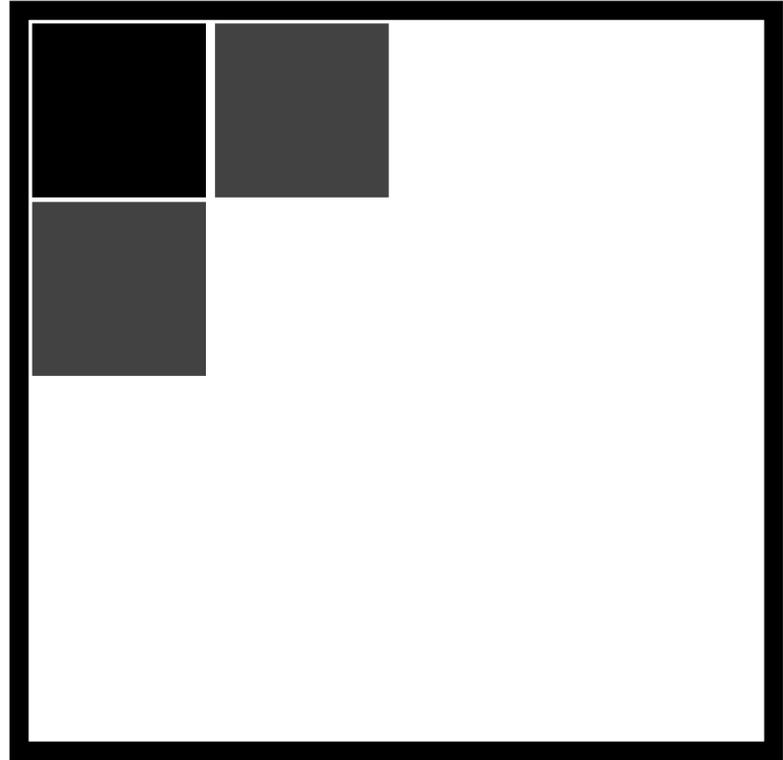


Pruning strategy #1

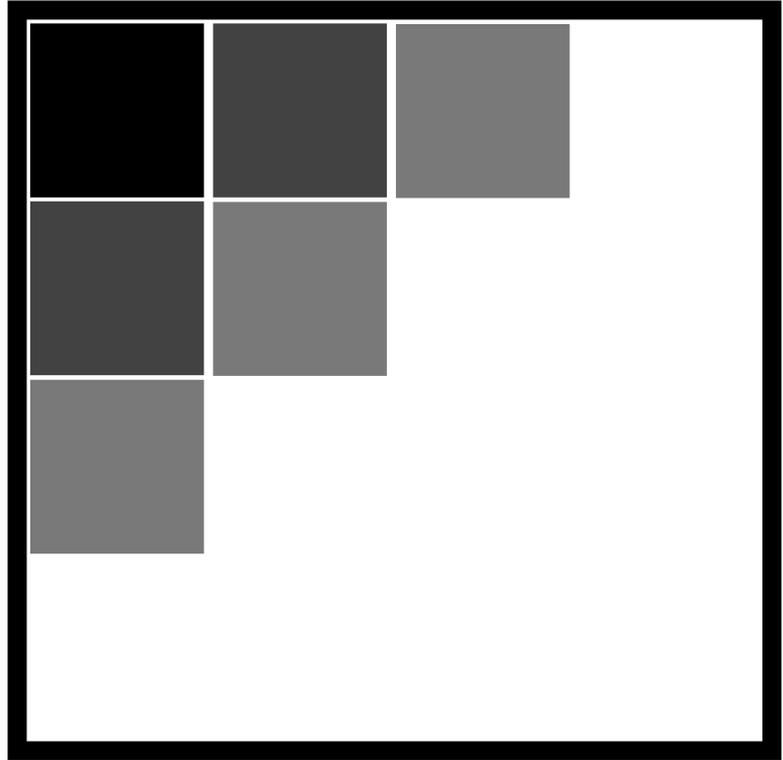
Value 1



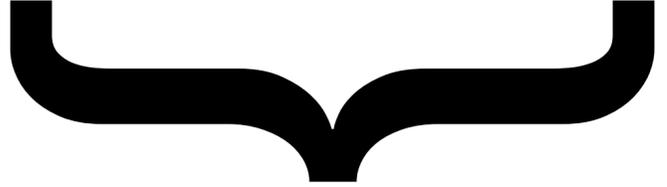
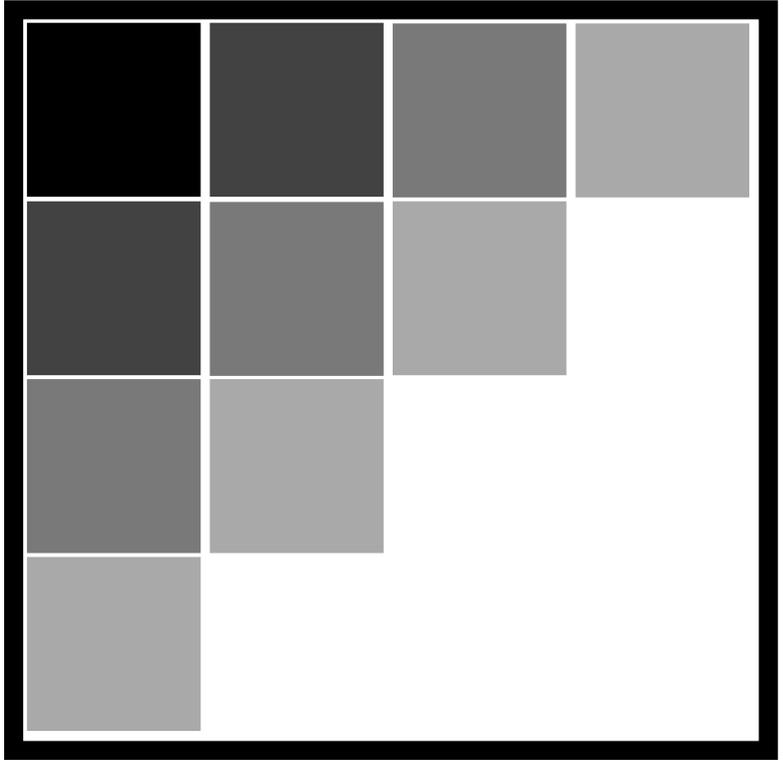
Value 2



Value 3



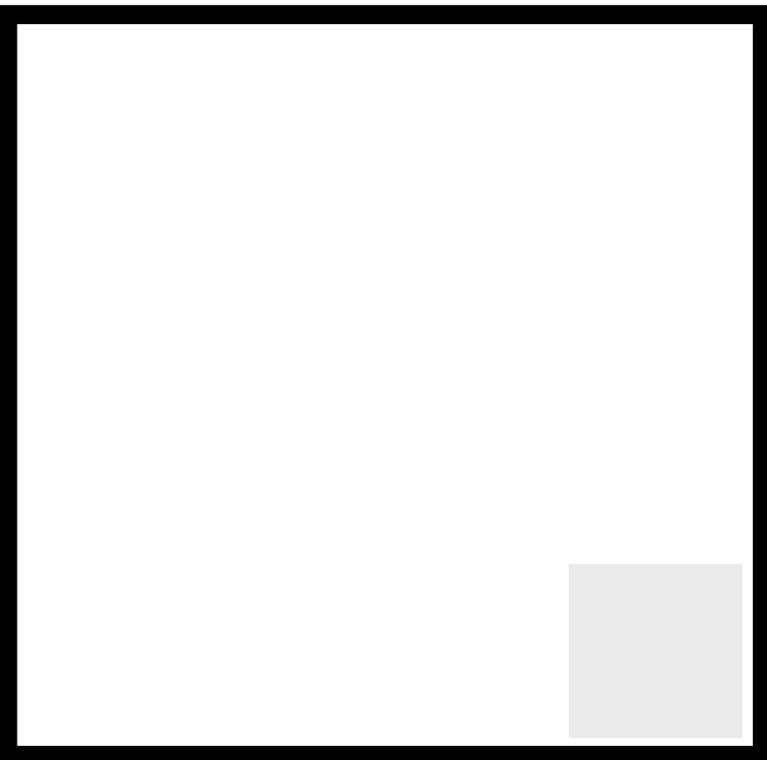
Value 4



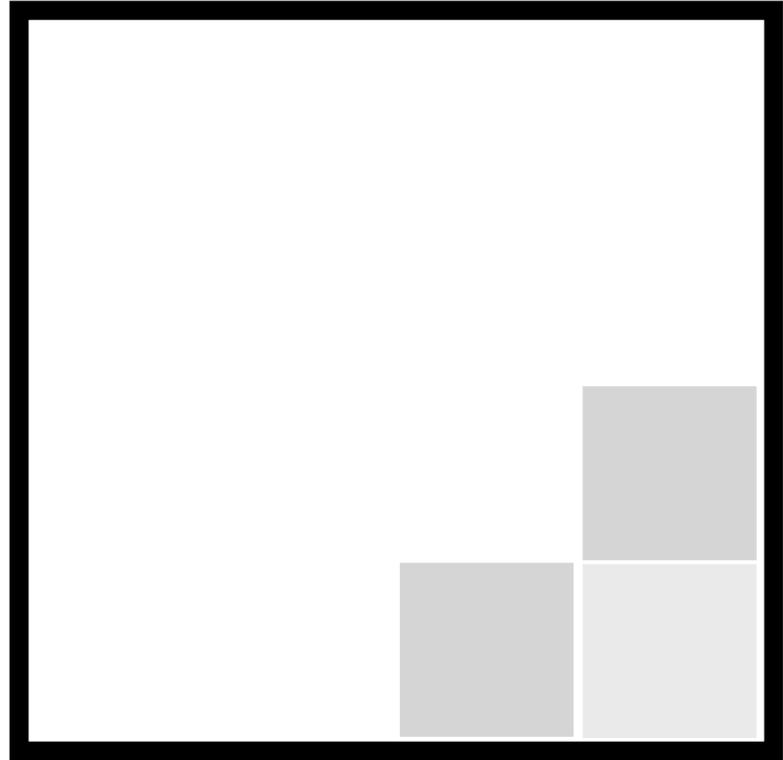
A single block

Pruning strategy #2

Value 1



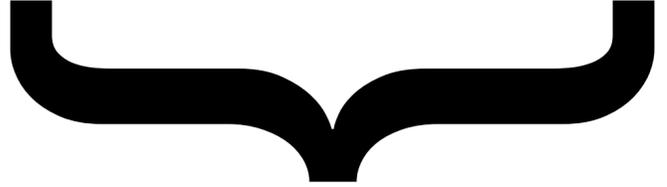
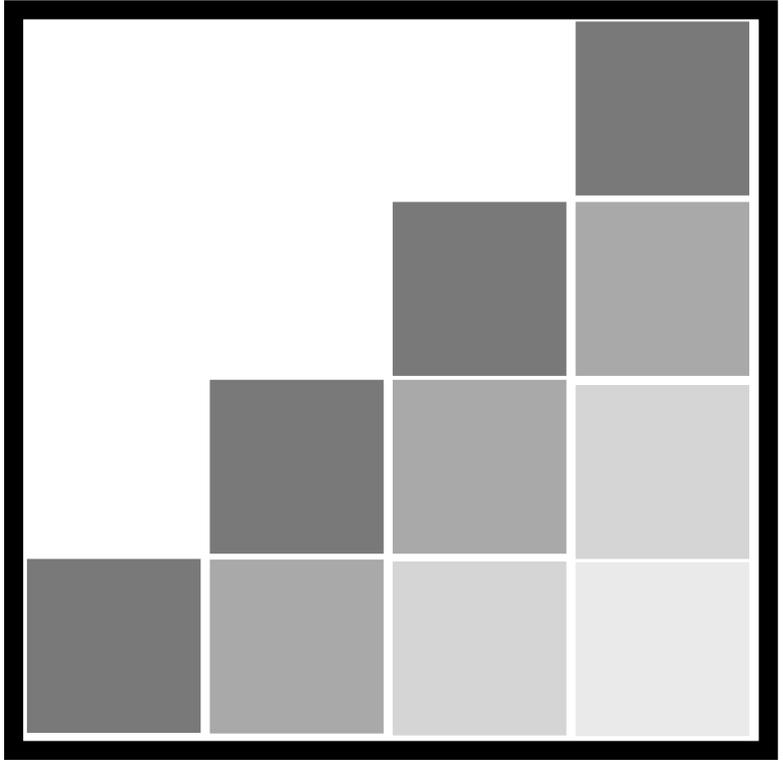
Value 2



Value 3

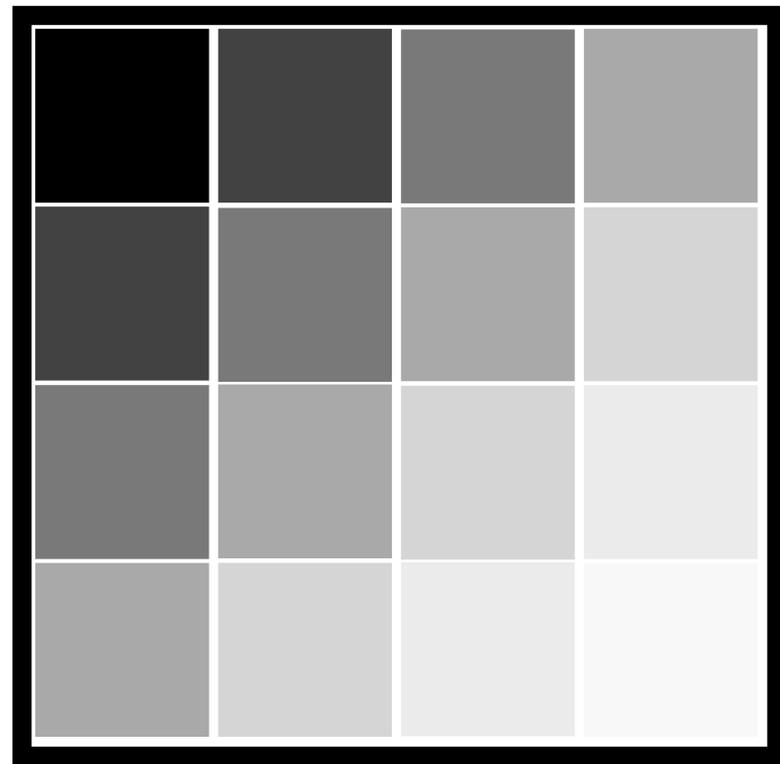


Value 4

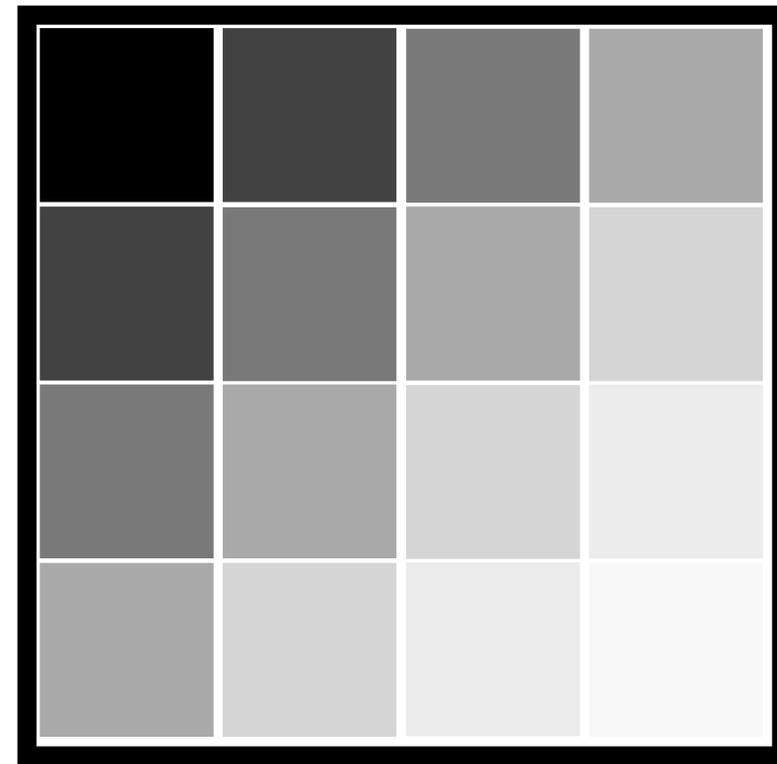


A single block

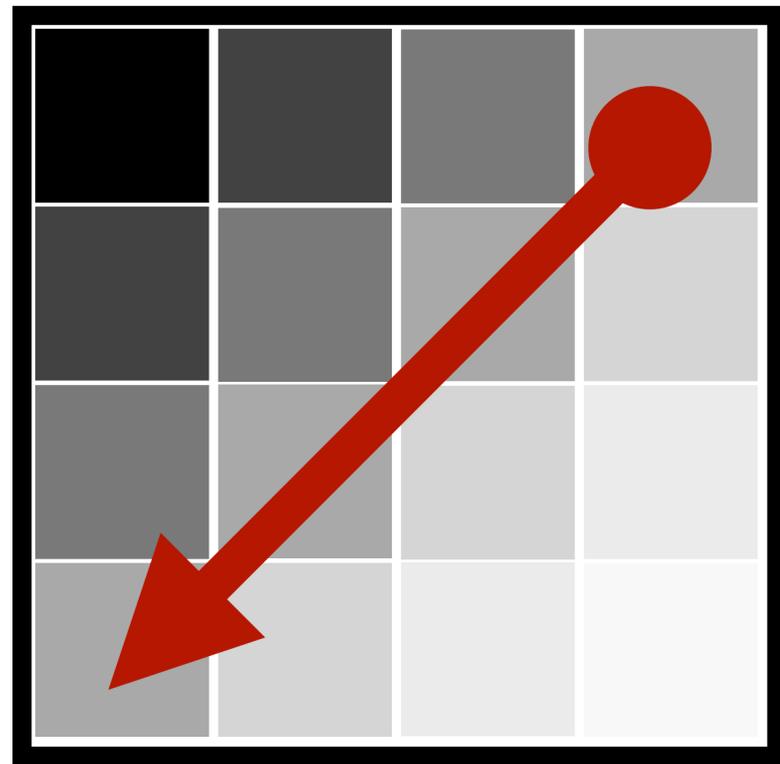
Strategy 3



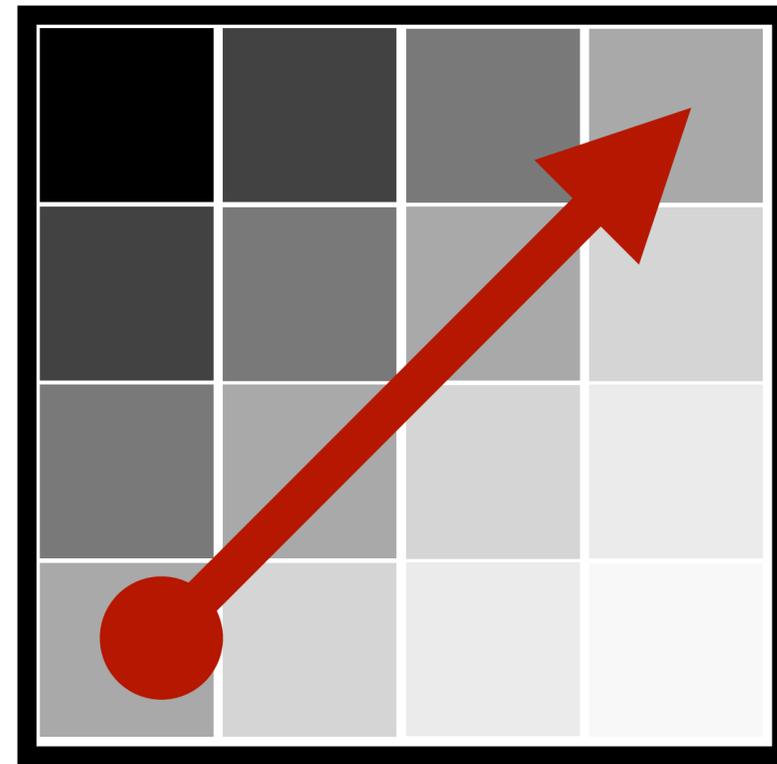
Strategy 4



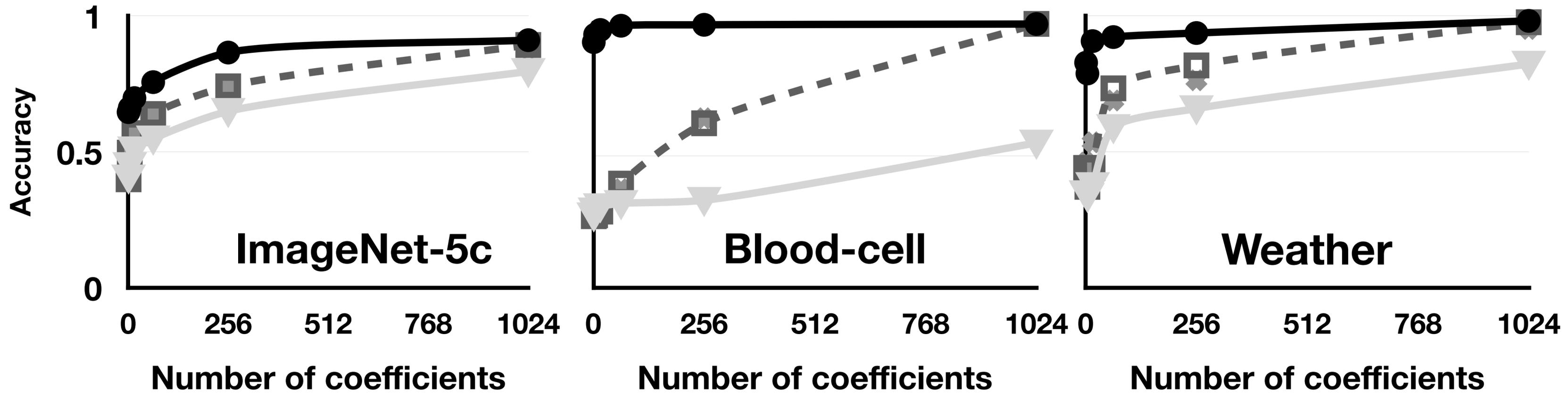
Strategy 3

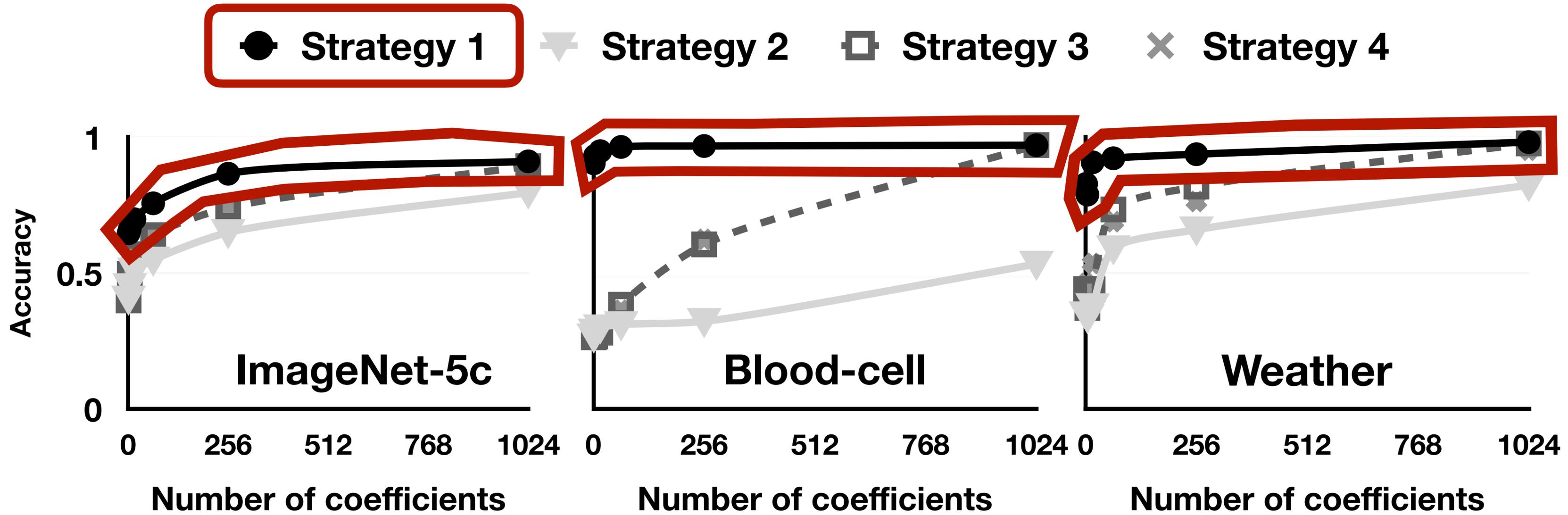


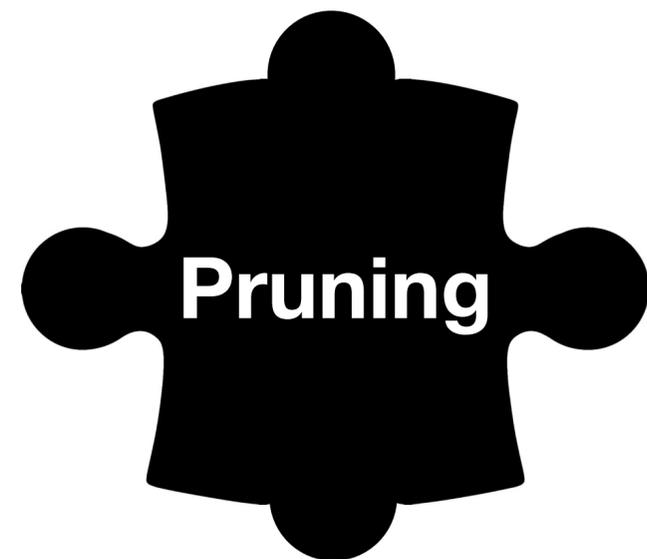
Strategy 4



● Strategy 1 ▼ Strategy 2 □ Strategy 3 × Strategy 4

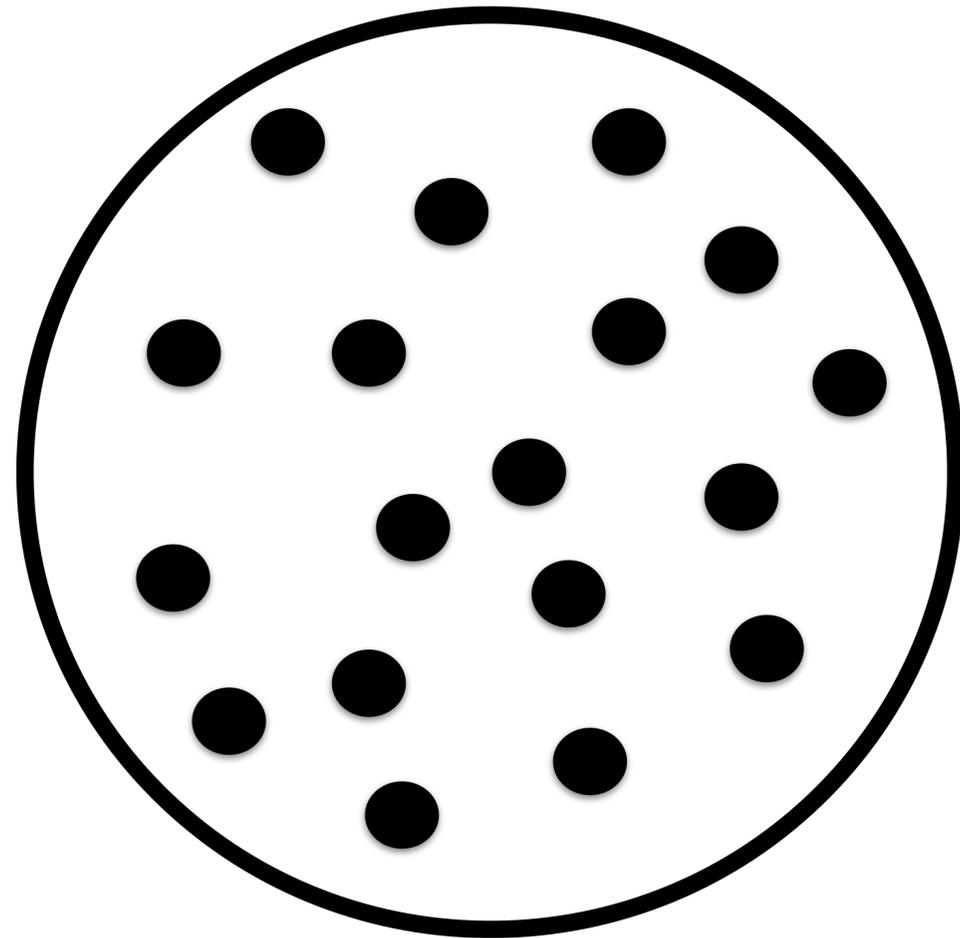




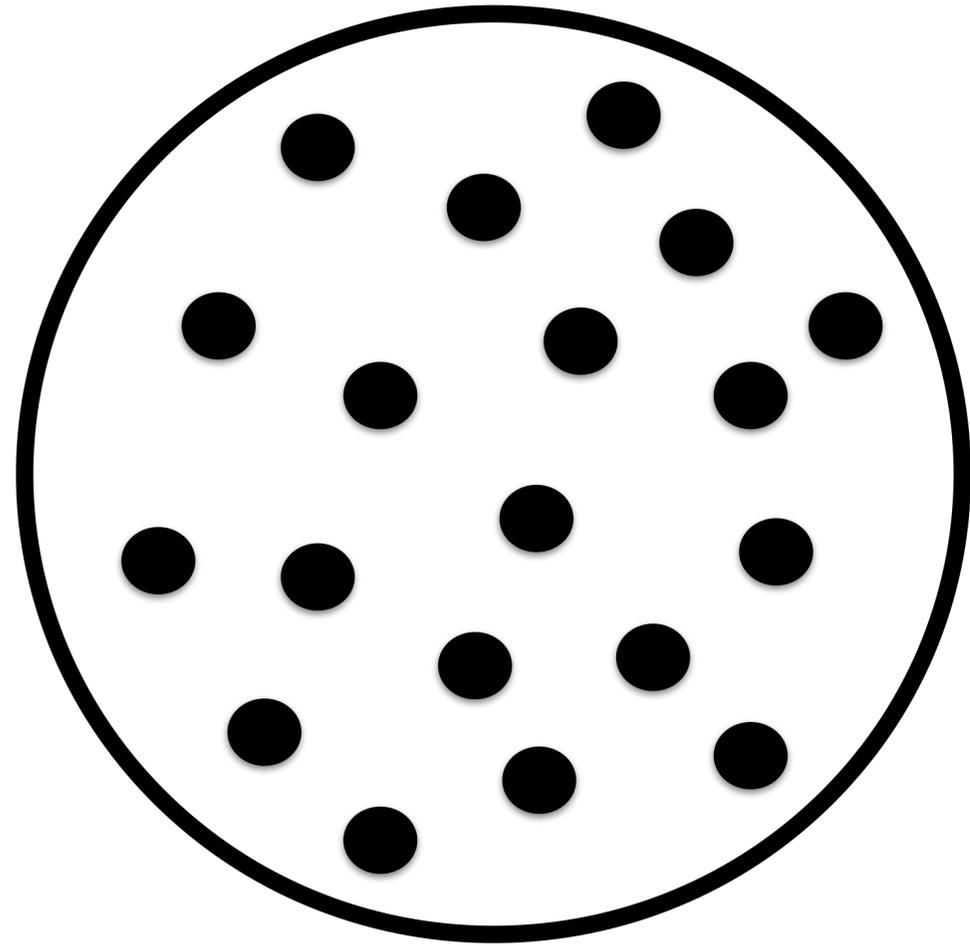


10150K → 6048

6048 new designs



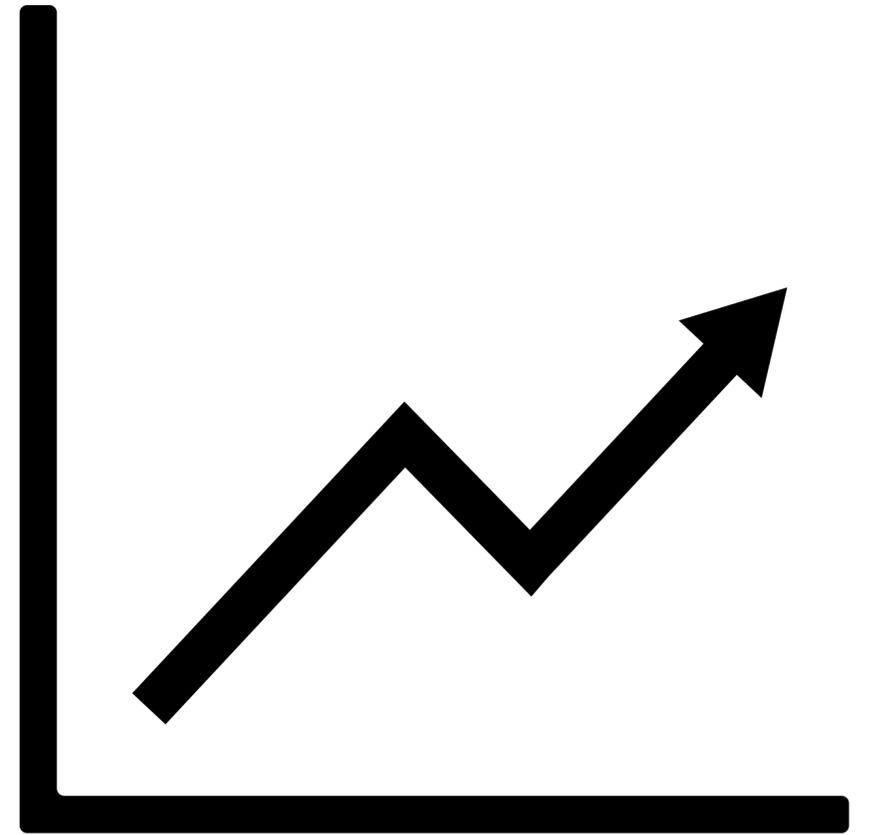
Design space



Search

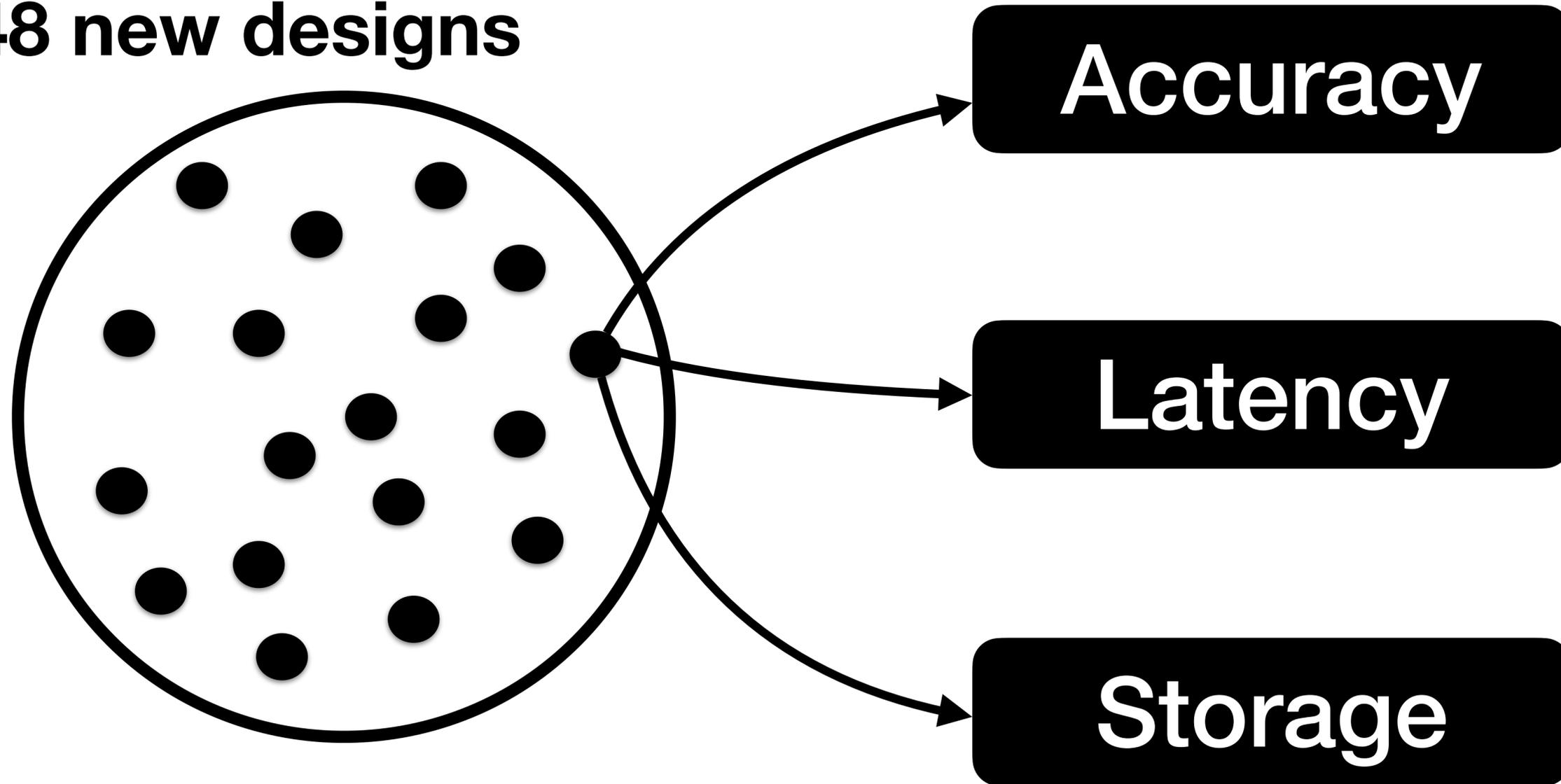


14x faster



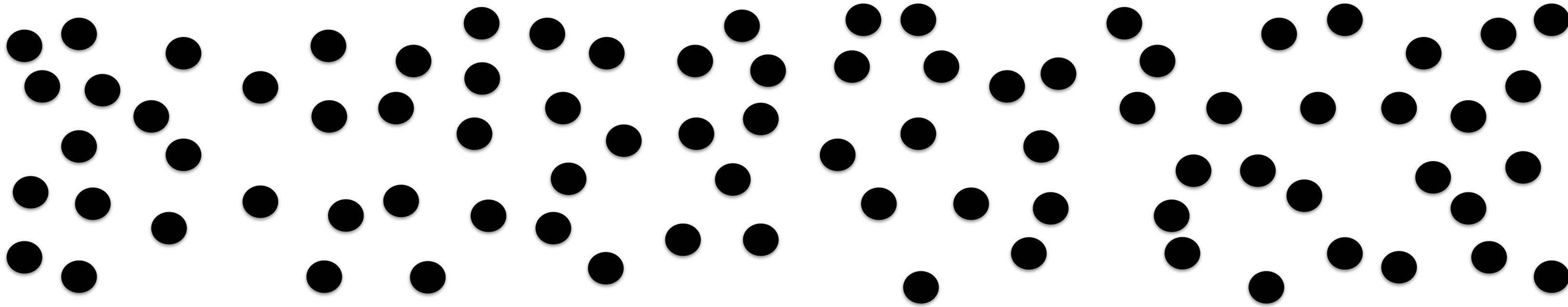
Performance models

6048 new designs



Accuracy model

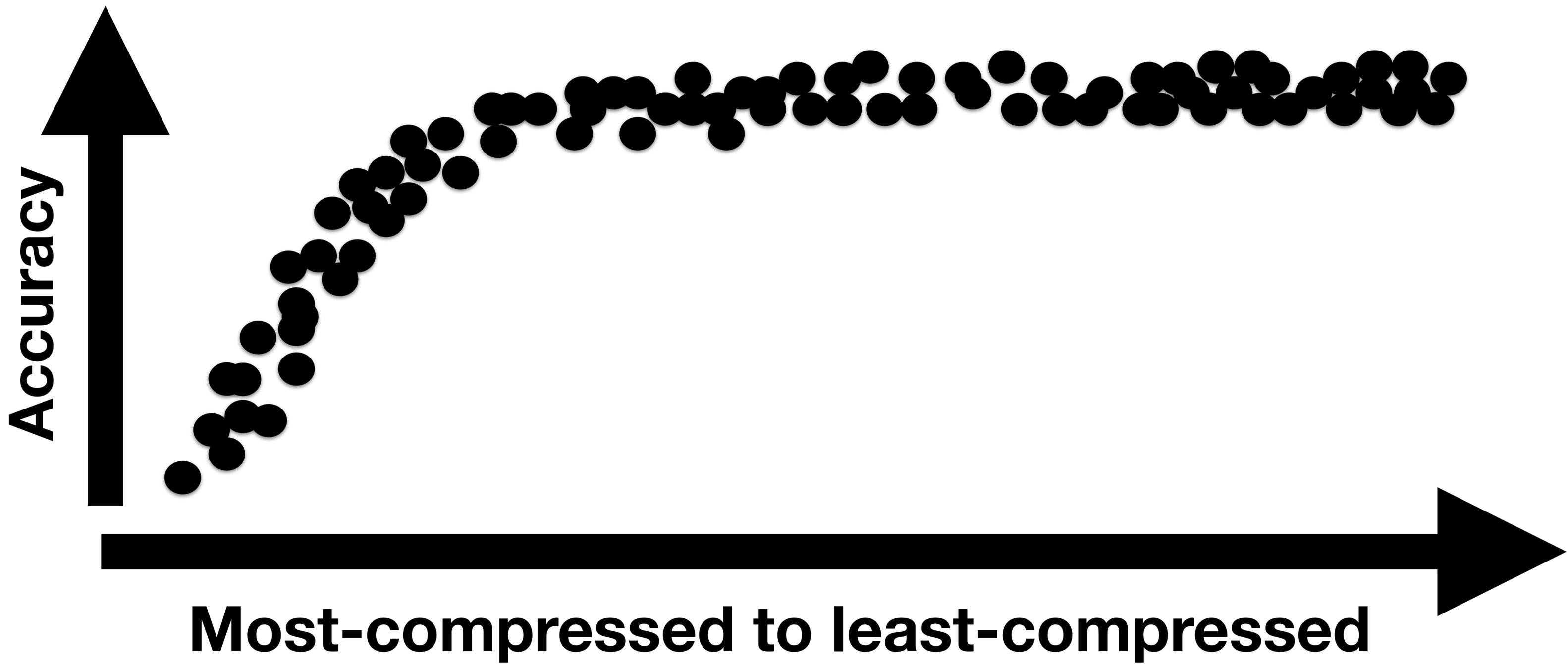
sort



Most-compressed to least-compressed

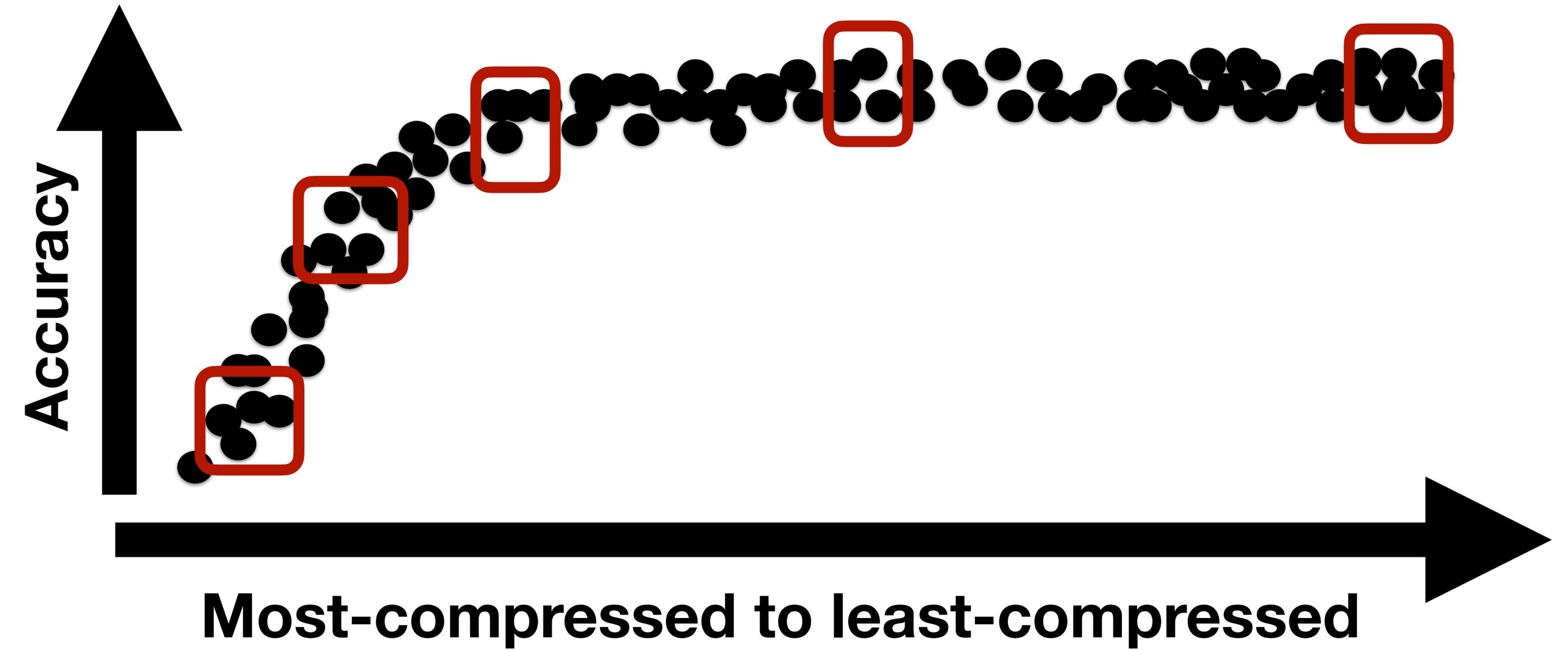
sort & analyze

*ImageNet-5c, ResNet50
A100, PyTorch v1*



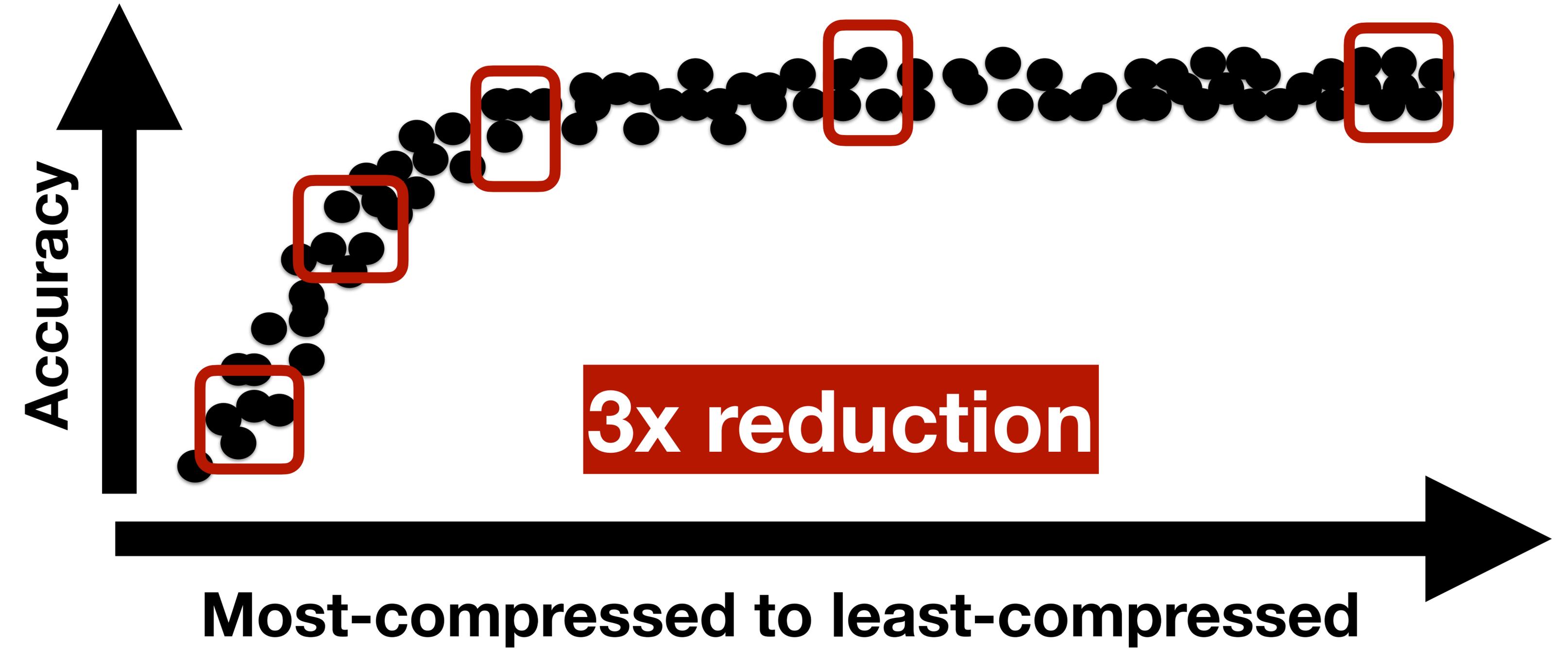
sort & analyze — sample & interpolate

*ImageNet, ResNet50
A100, PyTorch v1*

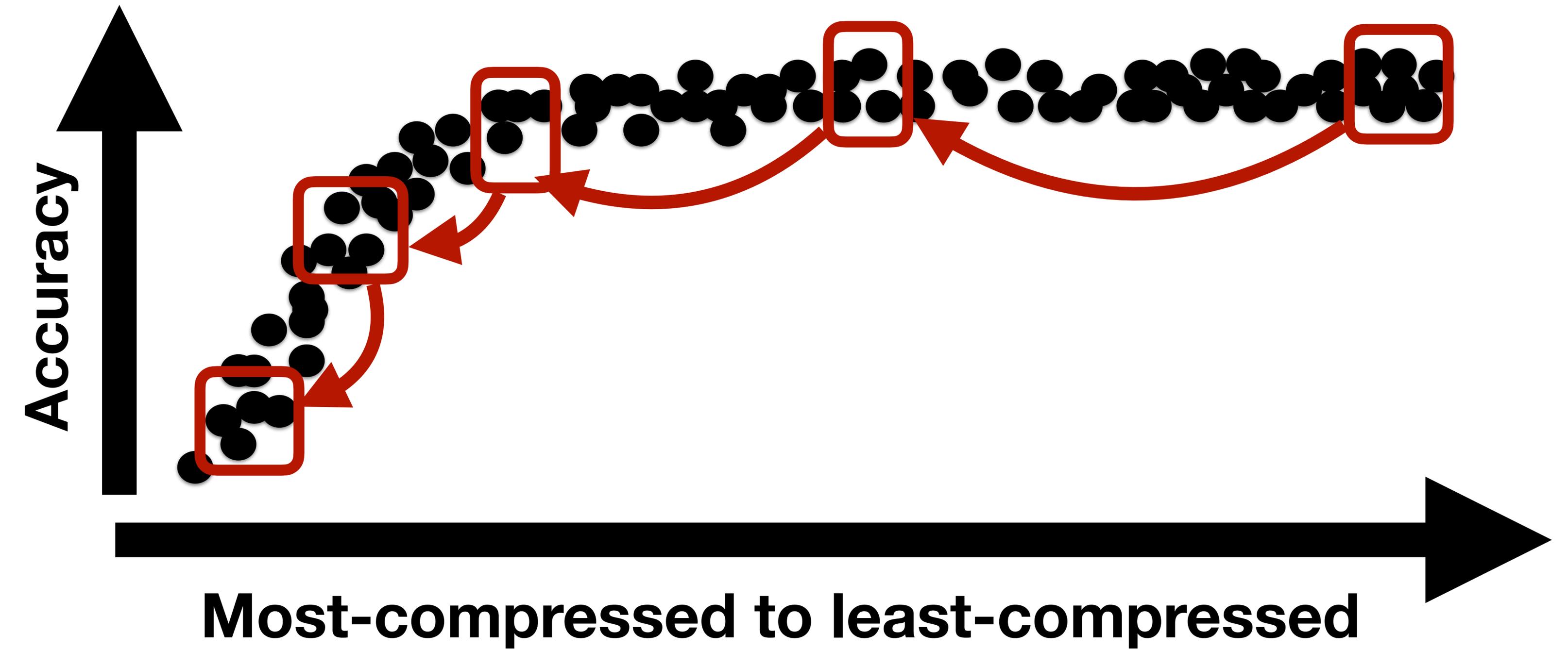


sort & analyze — sample & interpolate

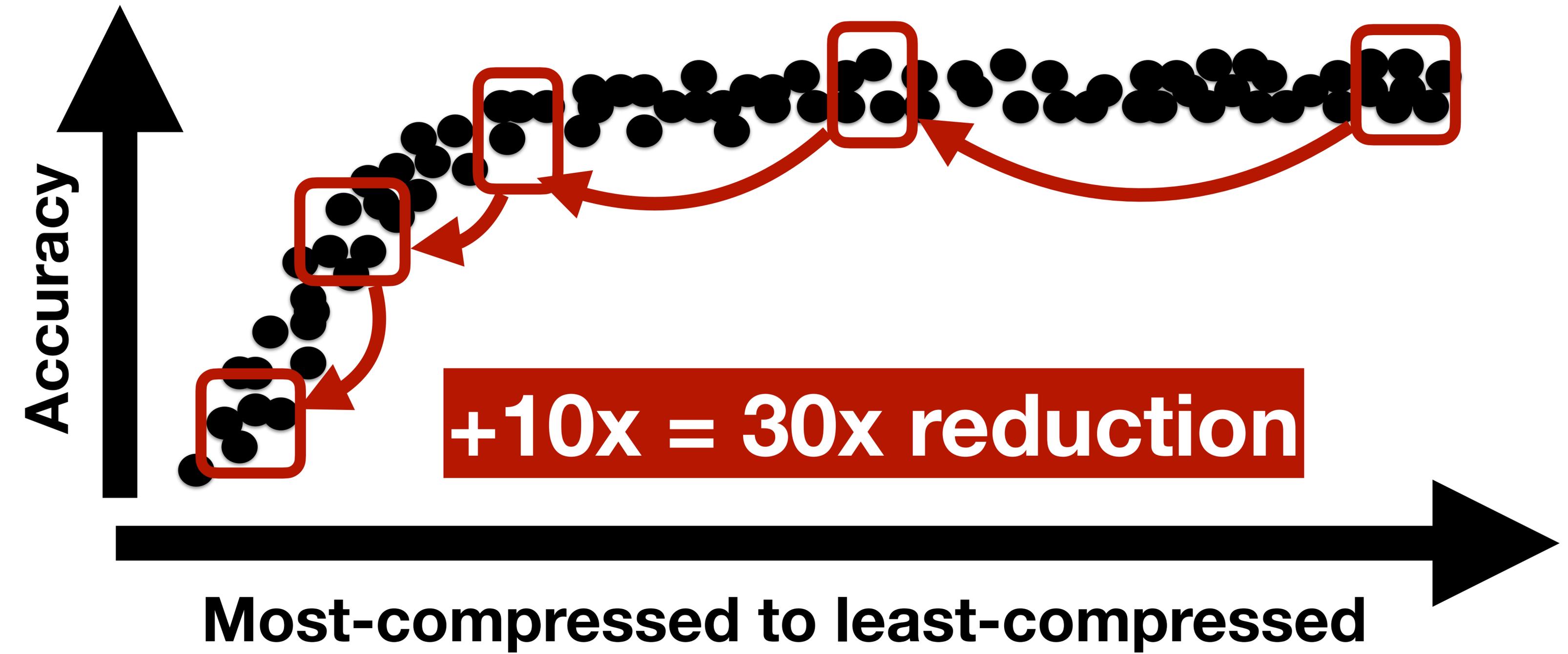
*ImageNet, ResNet50
A100, PyTorch v1*



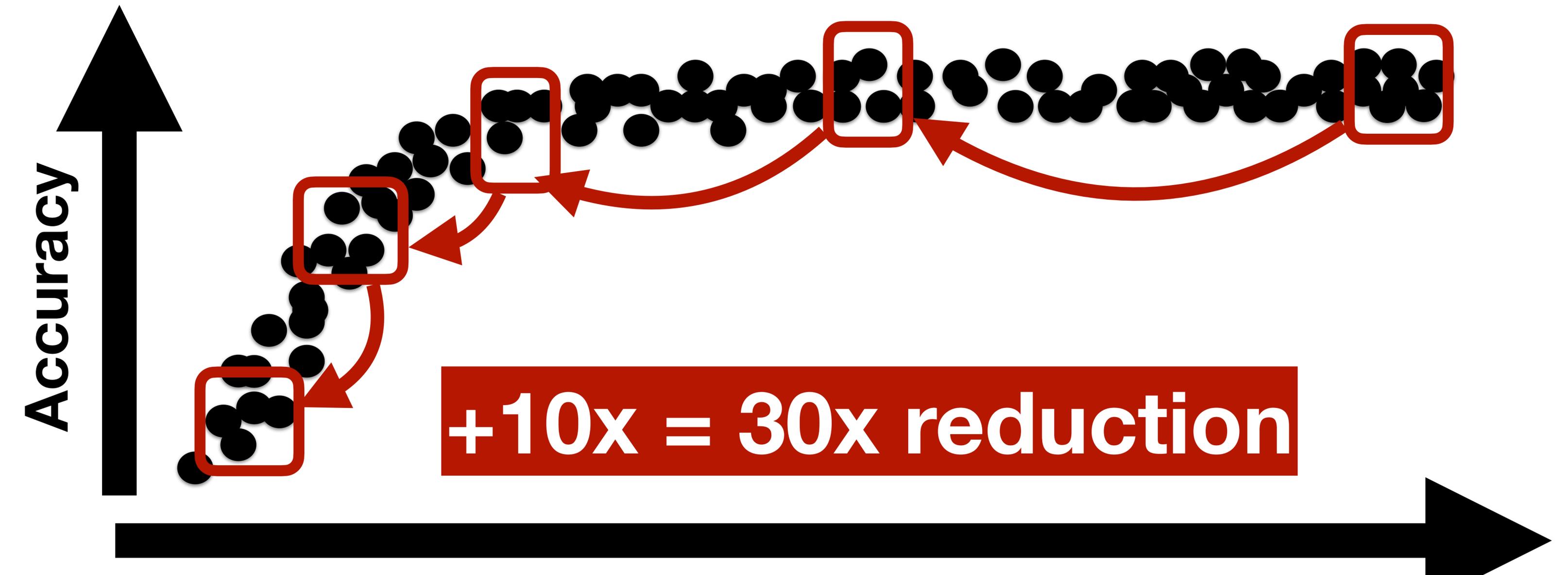
sort & analyze – sample & interpolate + transfer learn



sort & analyze – sample & interpolate + transfer learn

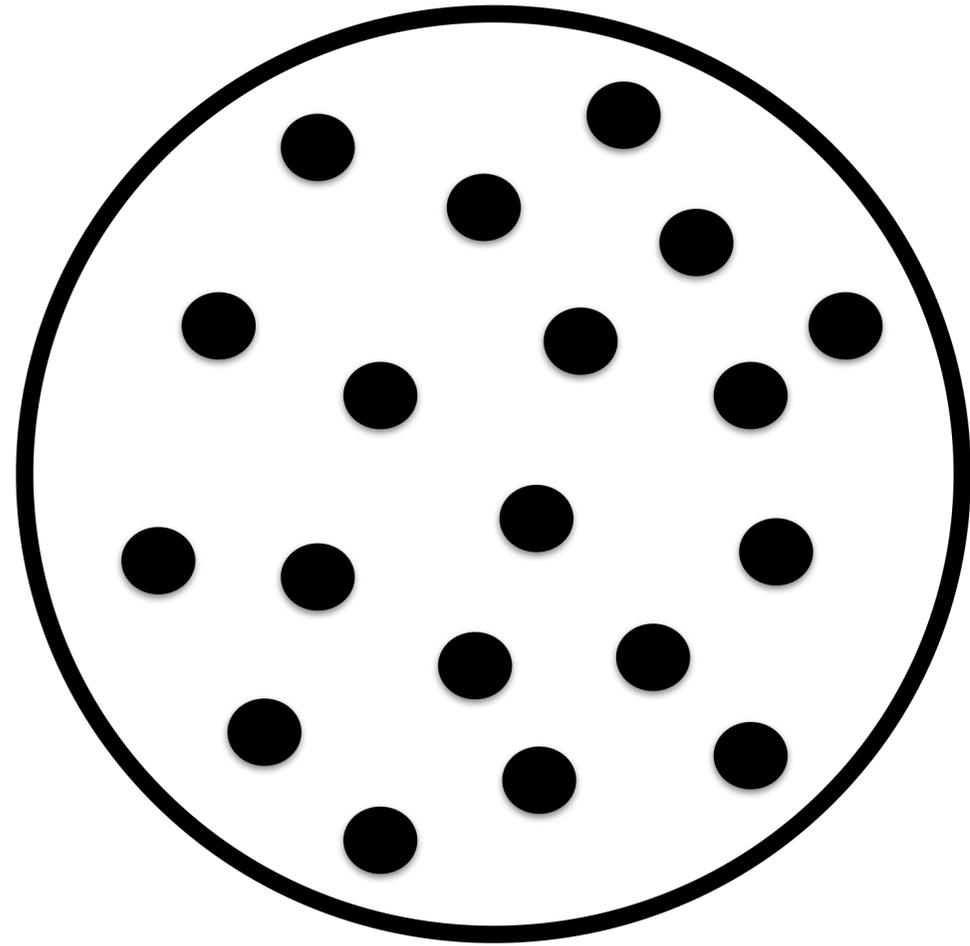


sort & analyze – sample & interpolate + transfer learn



6048 → 200 training

Design space



Search



14x faster

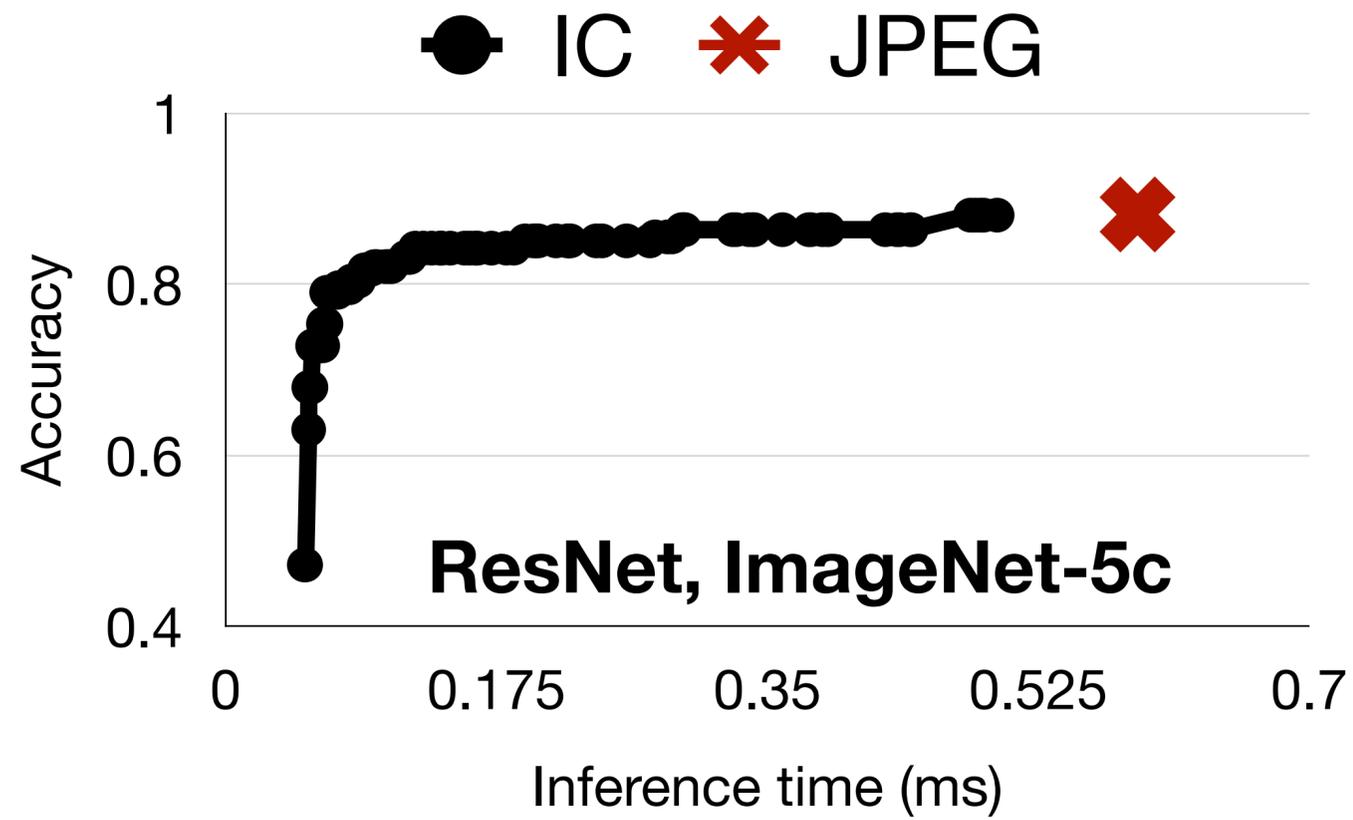


IC brings benefits on diverse datasets

ResNet50

A100

PyTorch v1

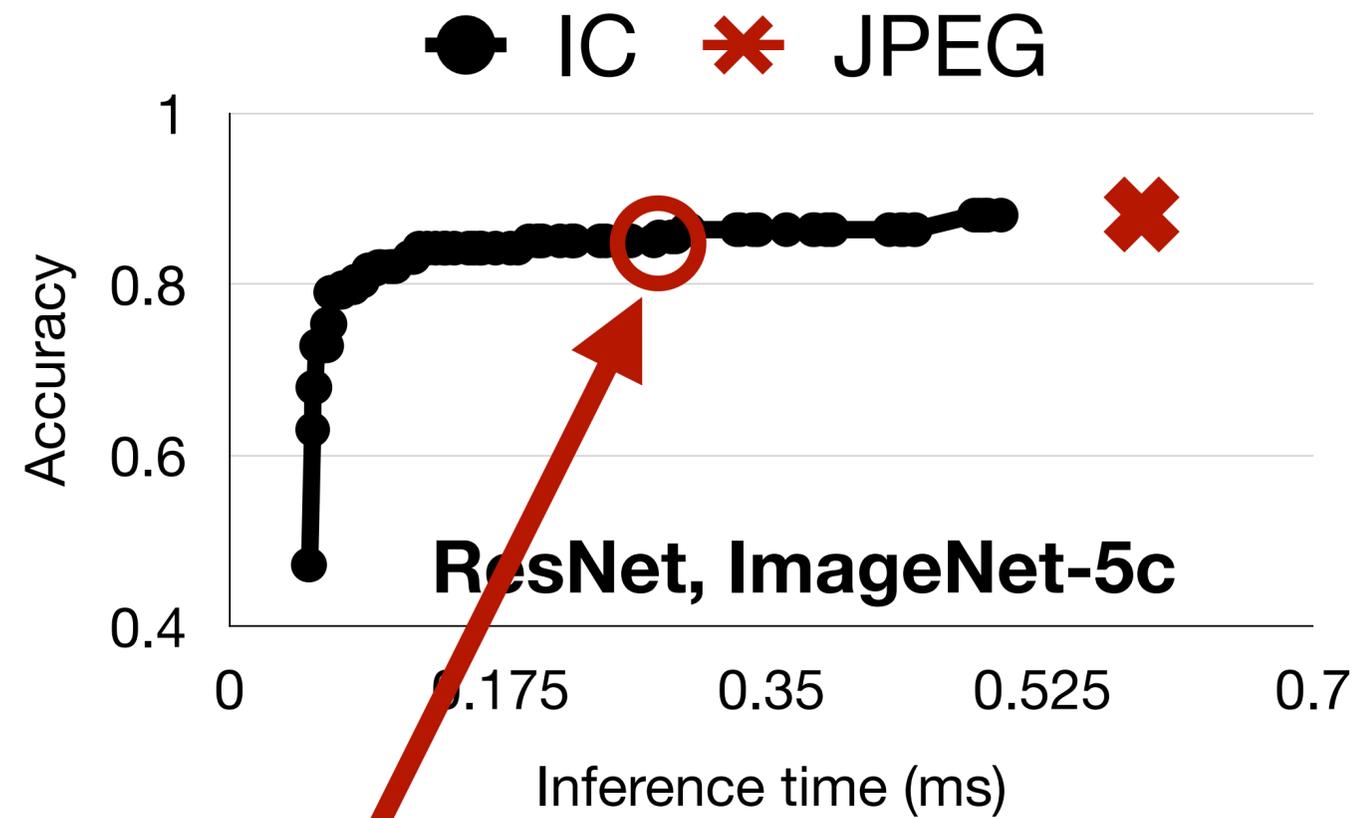


IC brings benefits on diverse datasets

ResNet50

A100

PyTorch v1



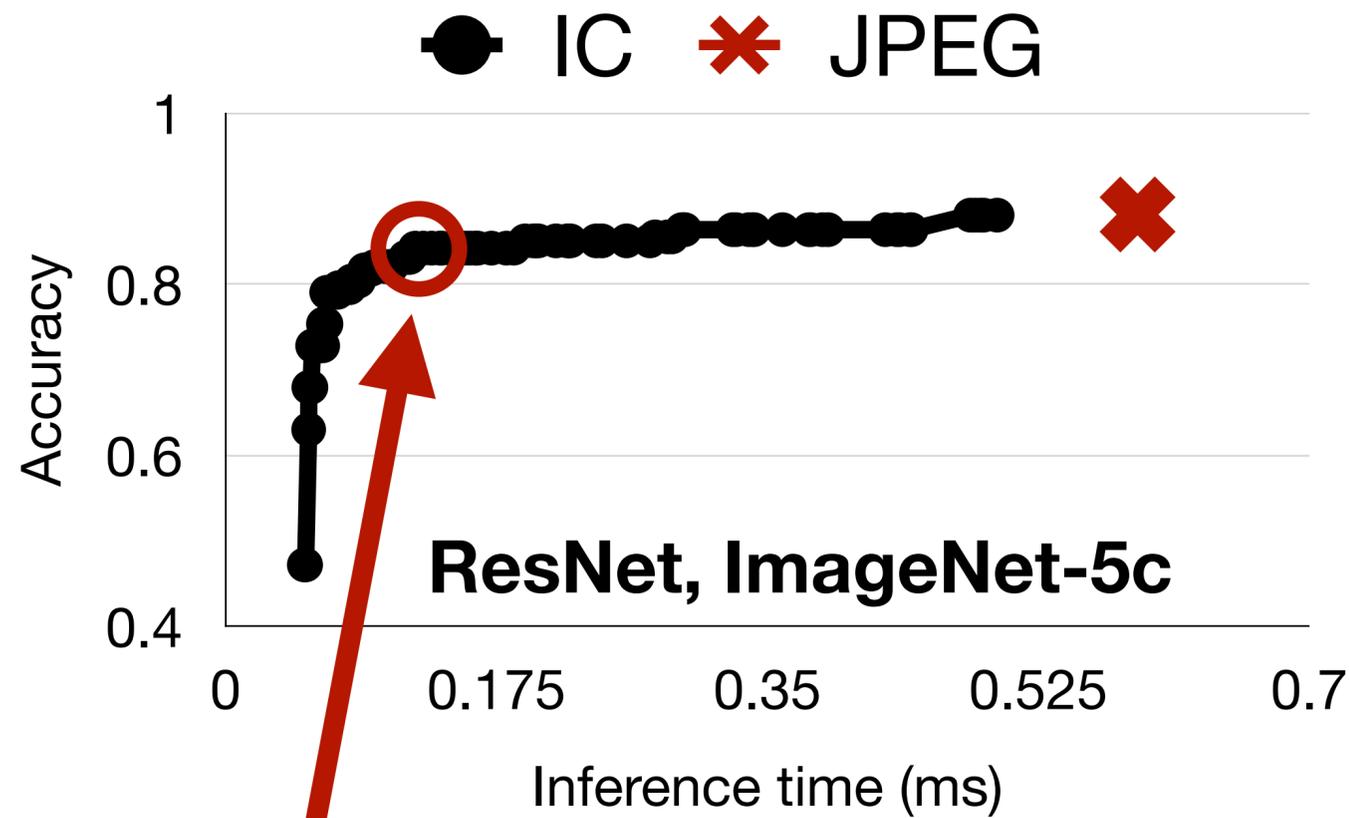
**2% loss vs.
2x gain**

IC brings benefits on diverse datasets

ResNet50

A100

PyTorch v1



**4% loss vs.
4x gain**

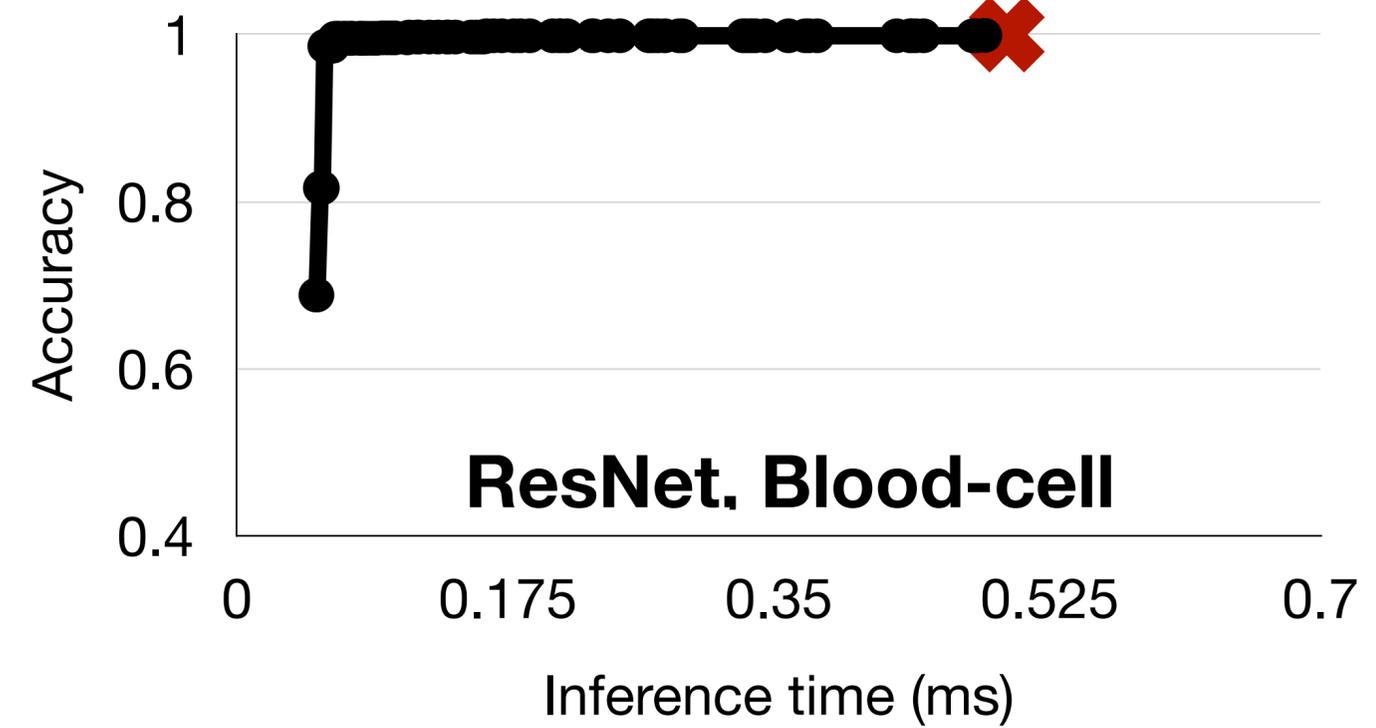
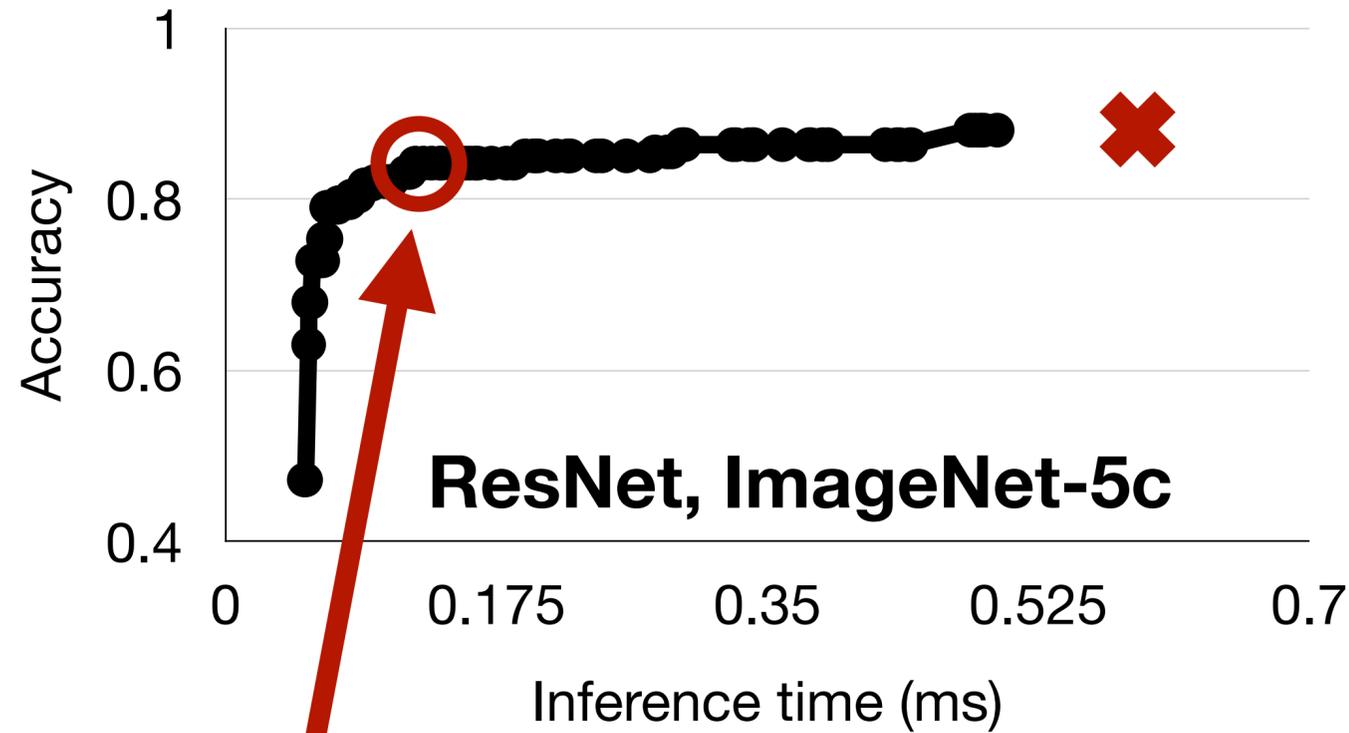
IC brings benefits on diverse datasets

ResNet50

A100

PyTorch v1

● IC ✖ JPEG



**4% loss vs.
4x gain**

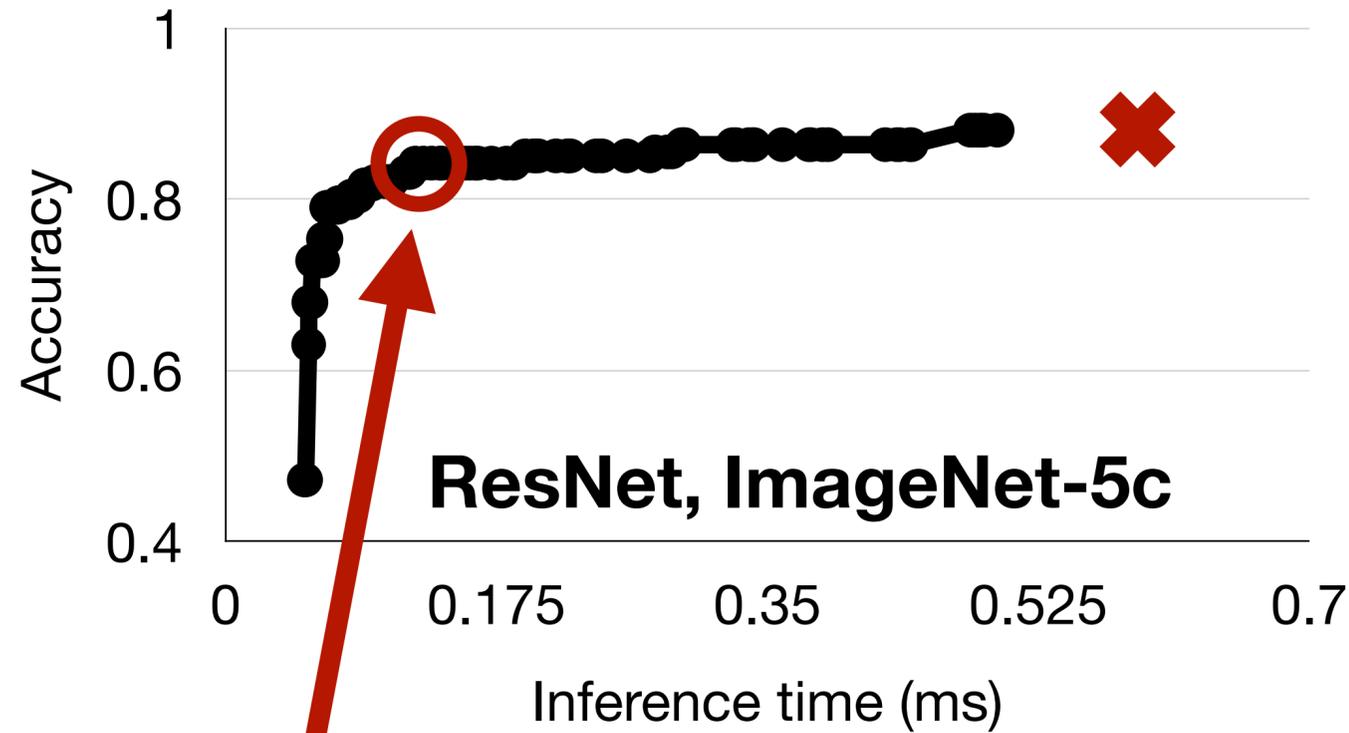
IC brings benefits on diverse datasets

ResNet50

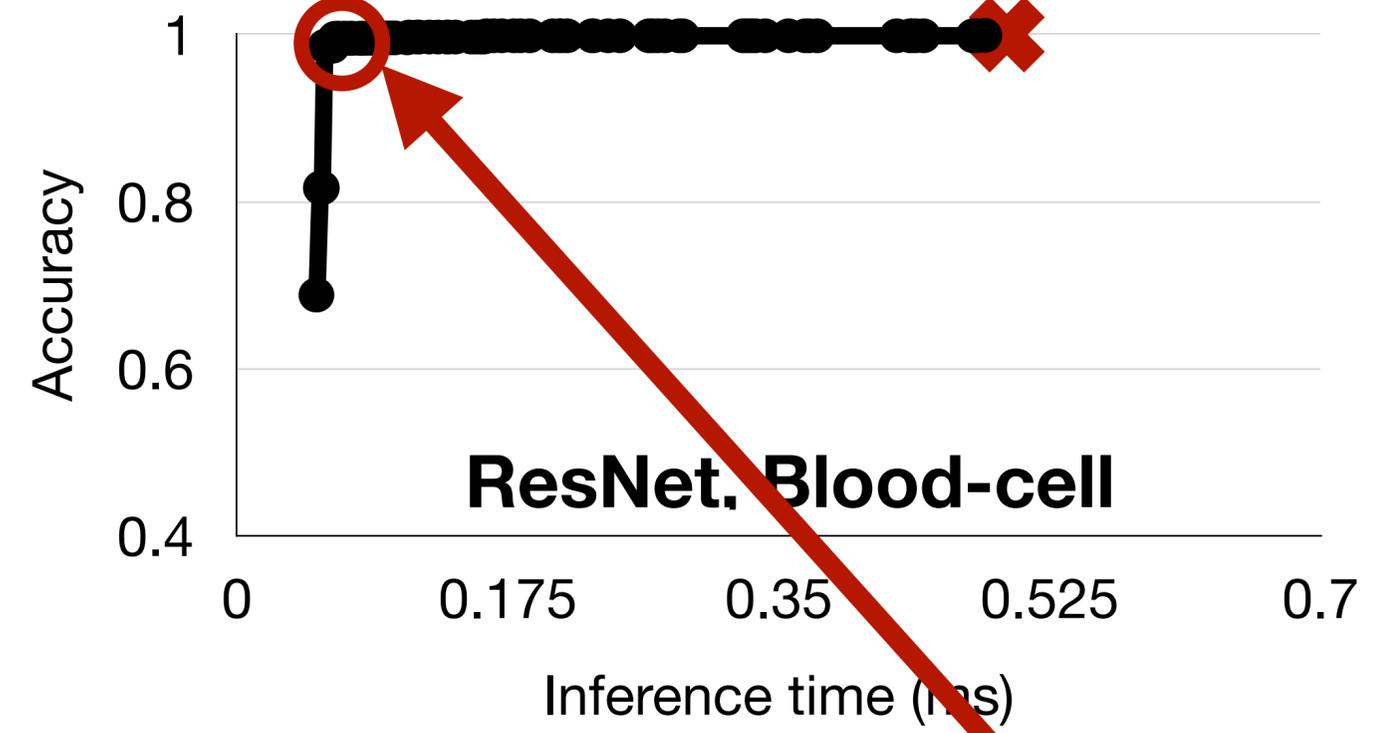
A100

PyTorch v1

● IC ✖ JPEG



**4% loss vs.
4x gain**



**No loss w/
9x gain**

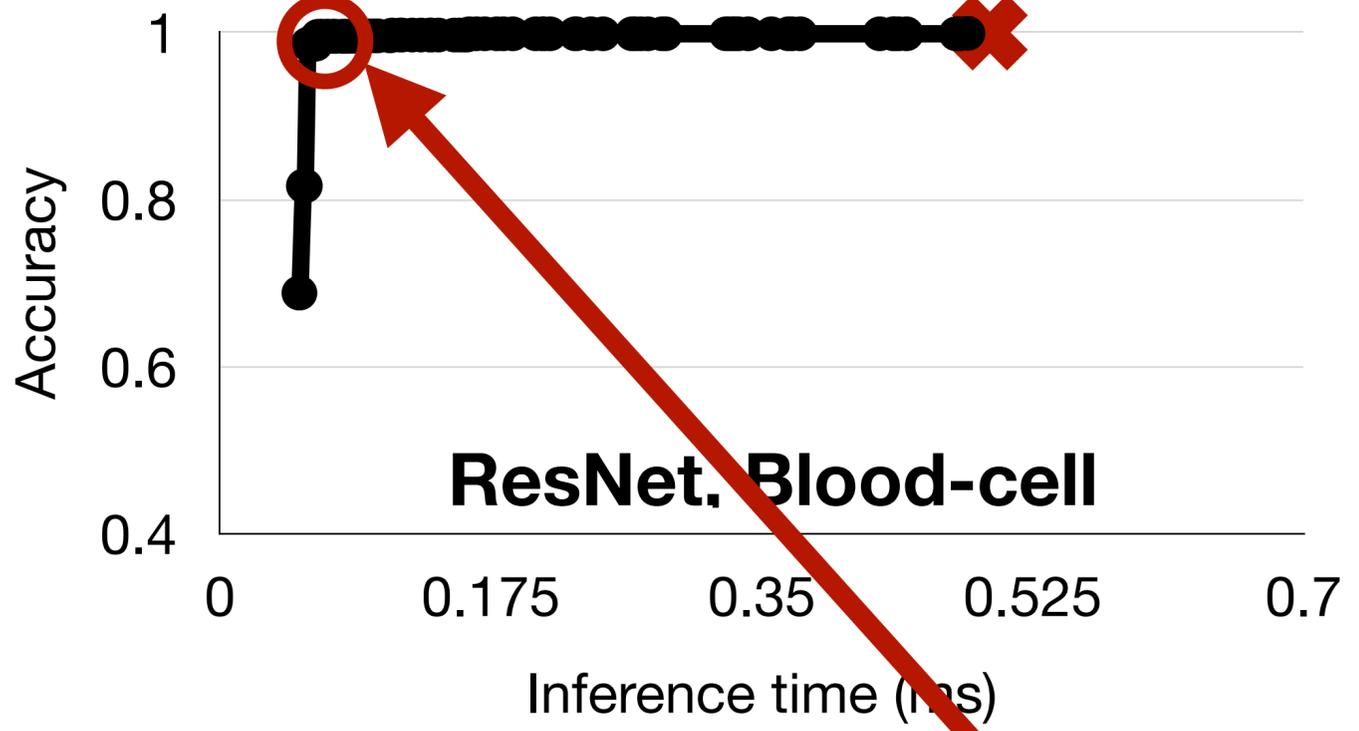
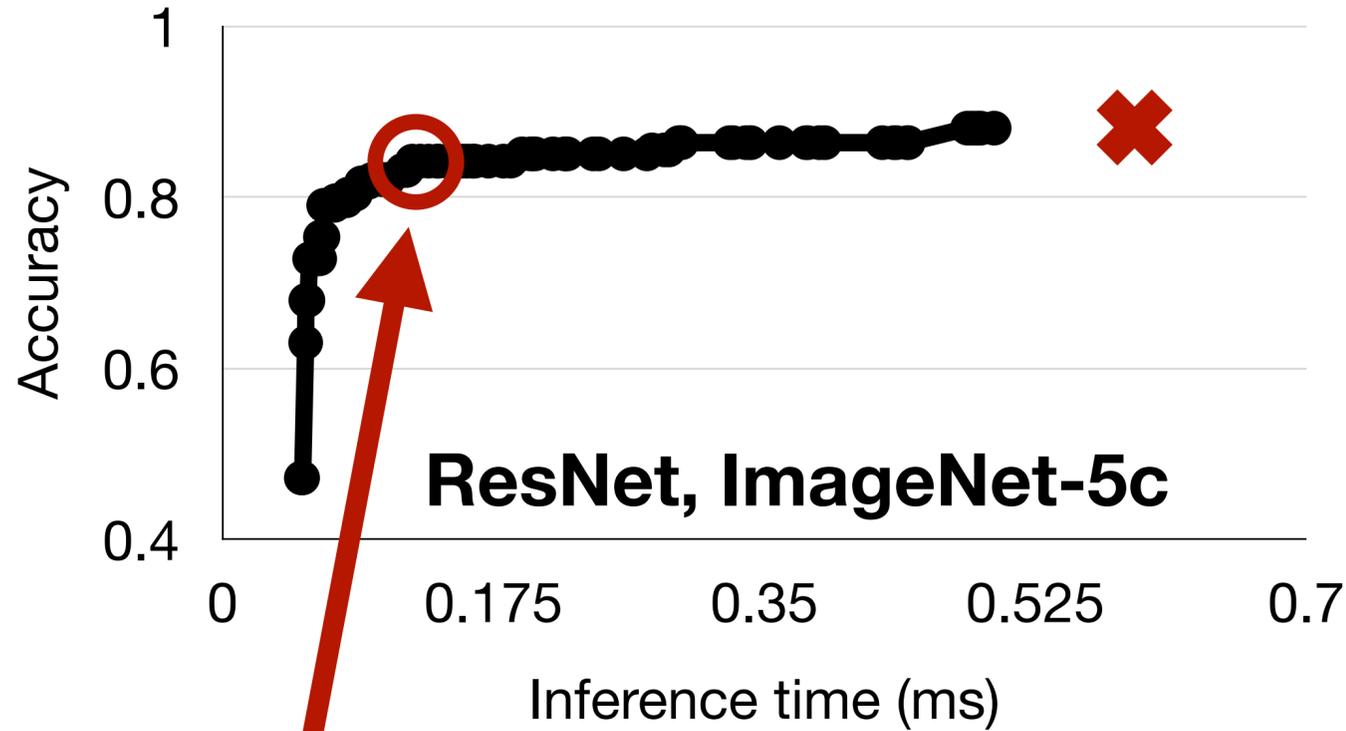
IC brings benefits on diverse datasets

ResNet50

A100

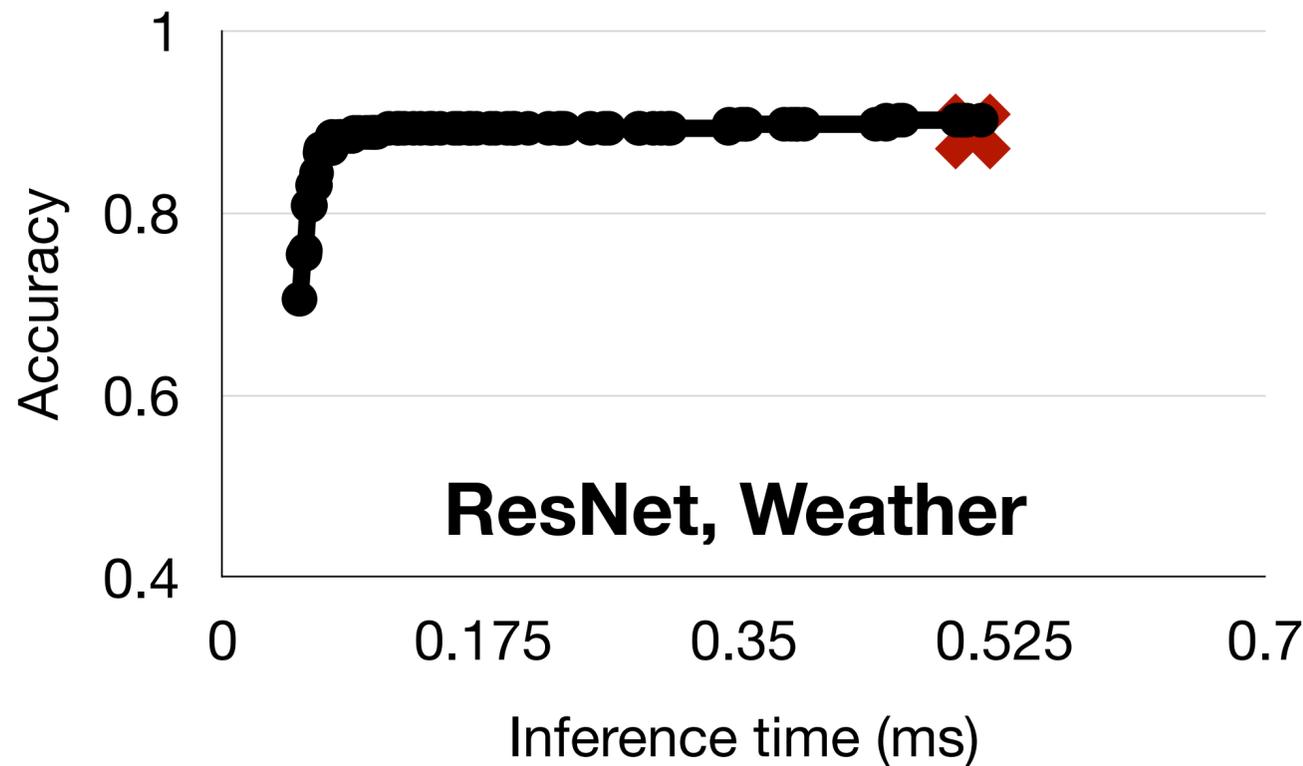
PyTorch v1

● IC ✖ JPEG



**4% loss vs.
4x gain**

**No loss w/
9x gain**



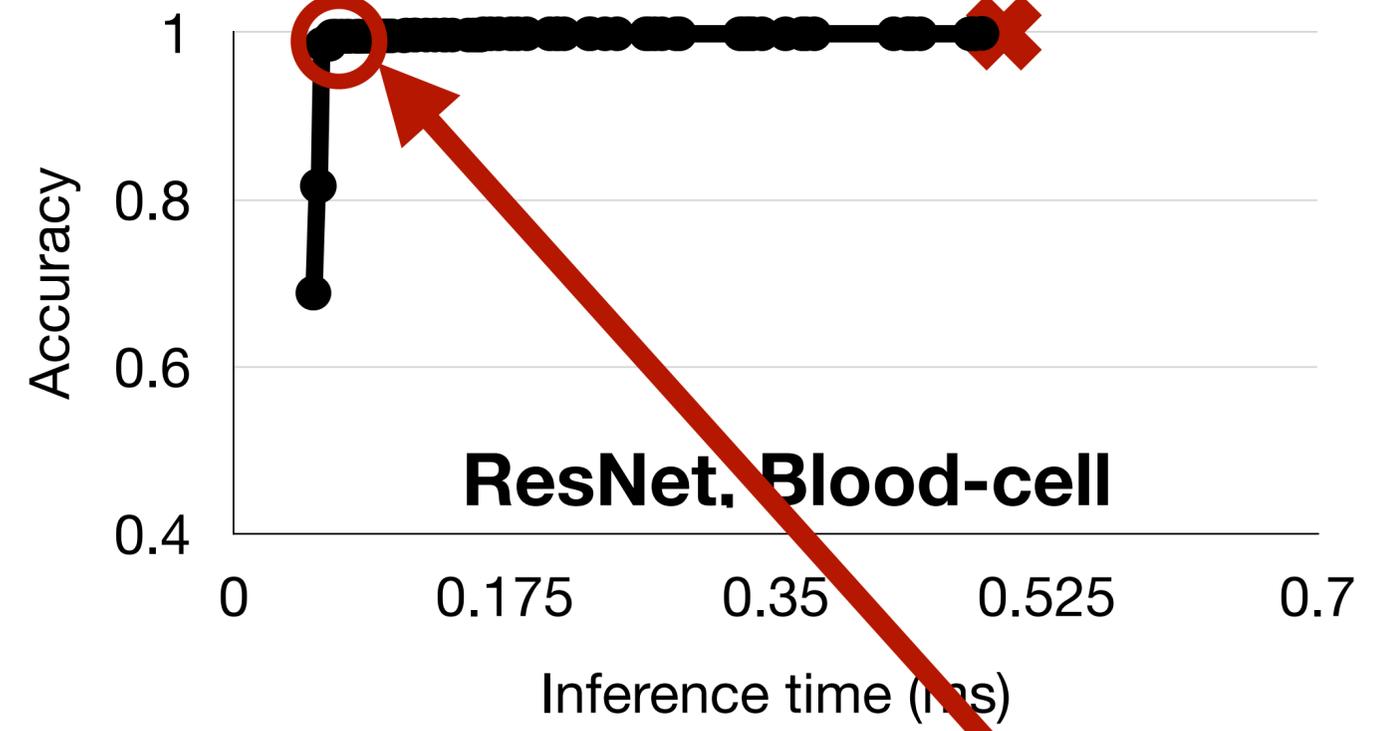
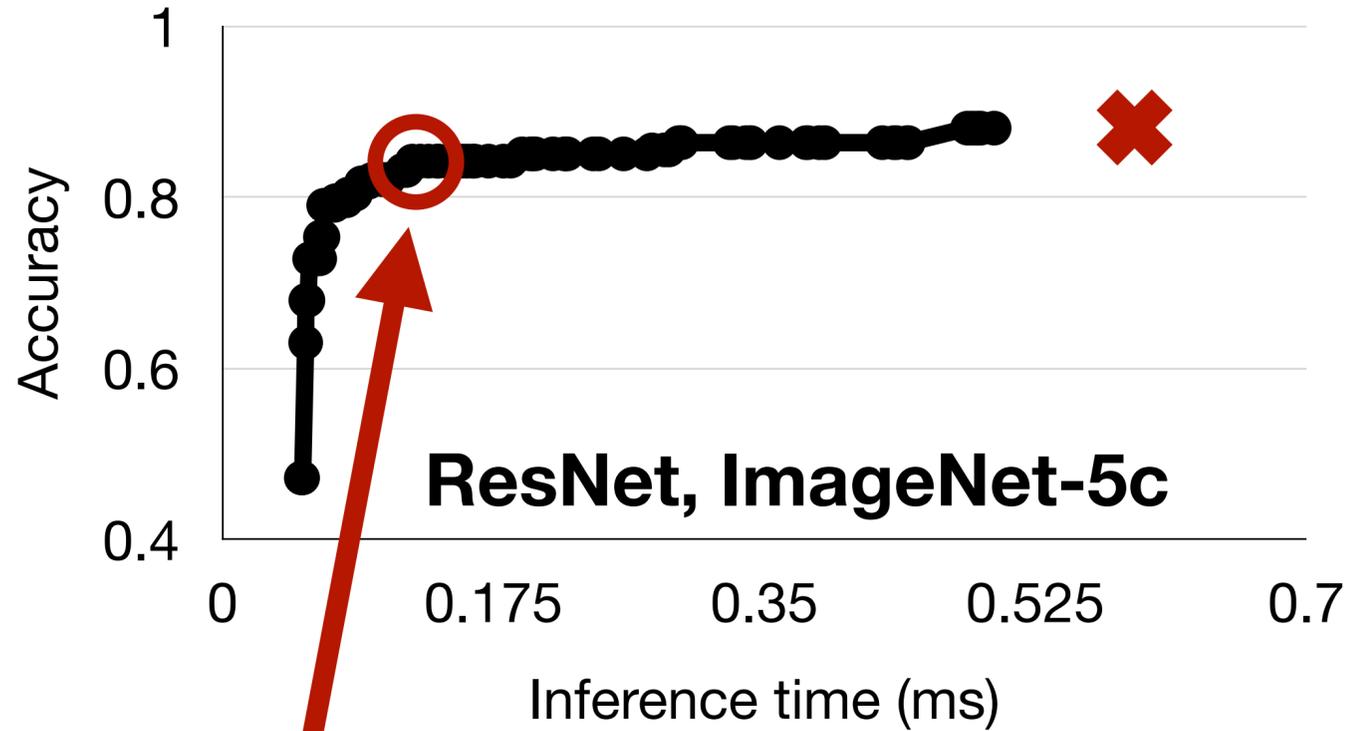
IC brings benefits on diverse datasets

ResNet50

A100

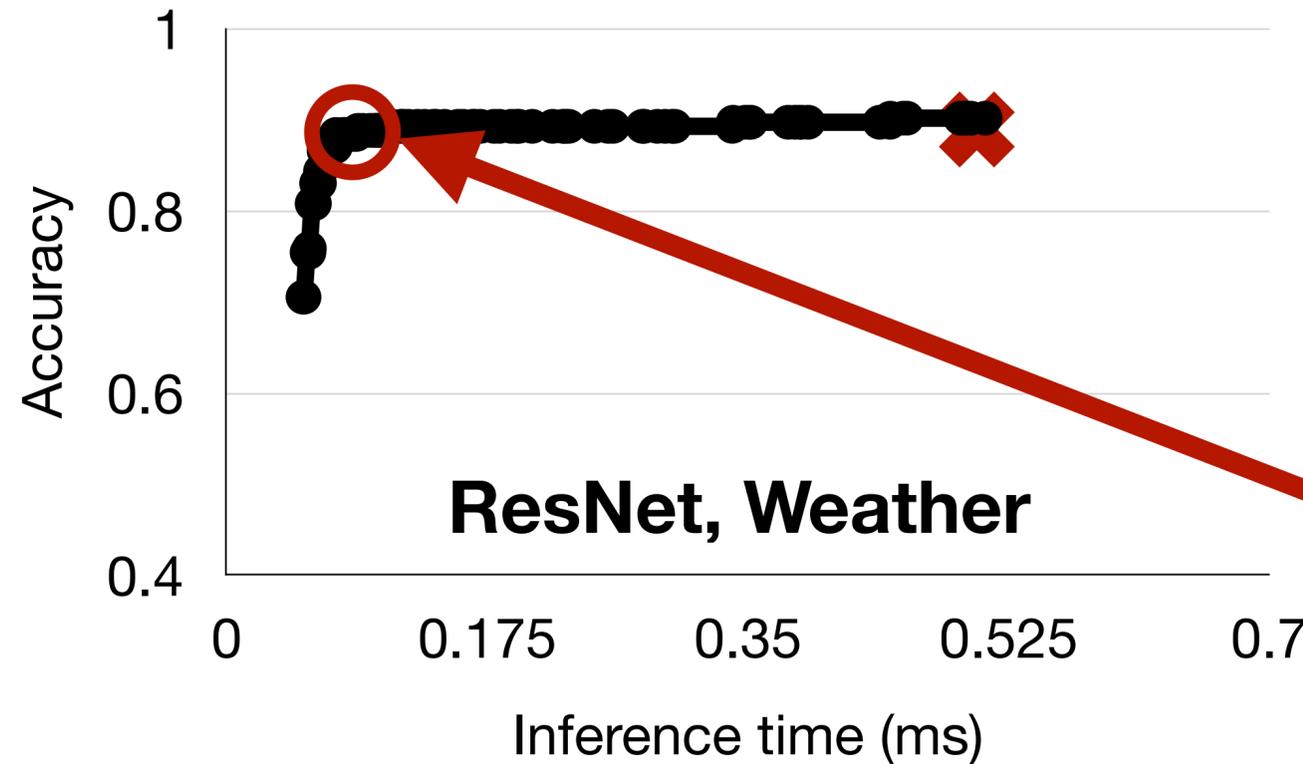
PyTorch v1

● IC ✖ JPEG



**4% loss vs.
4x gain**

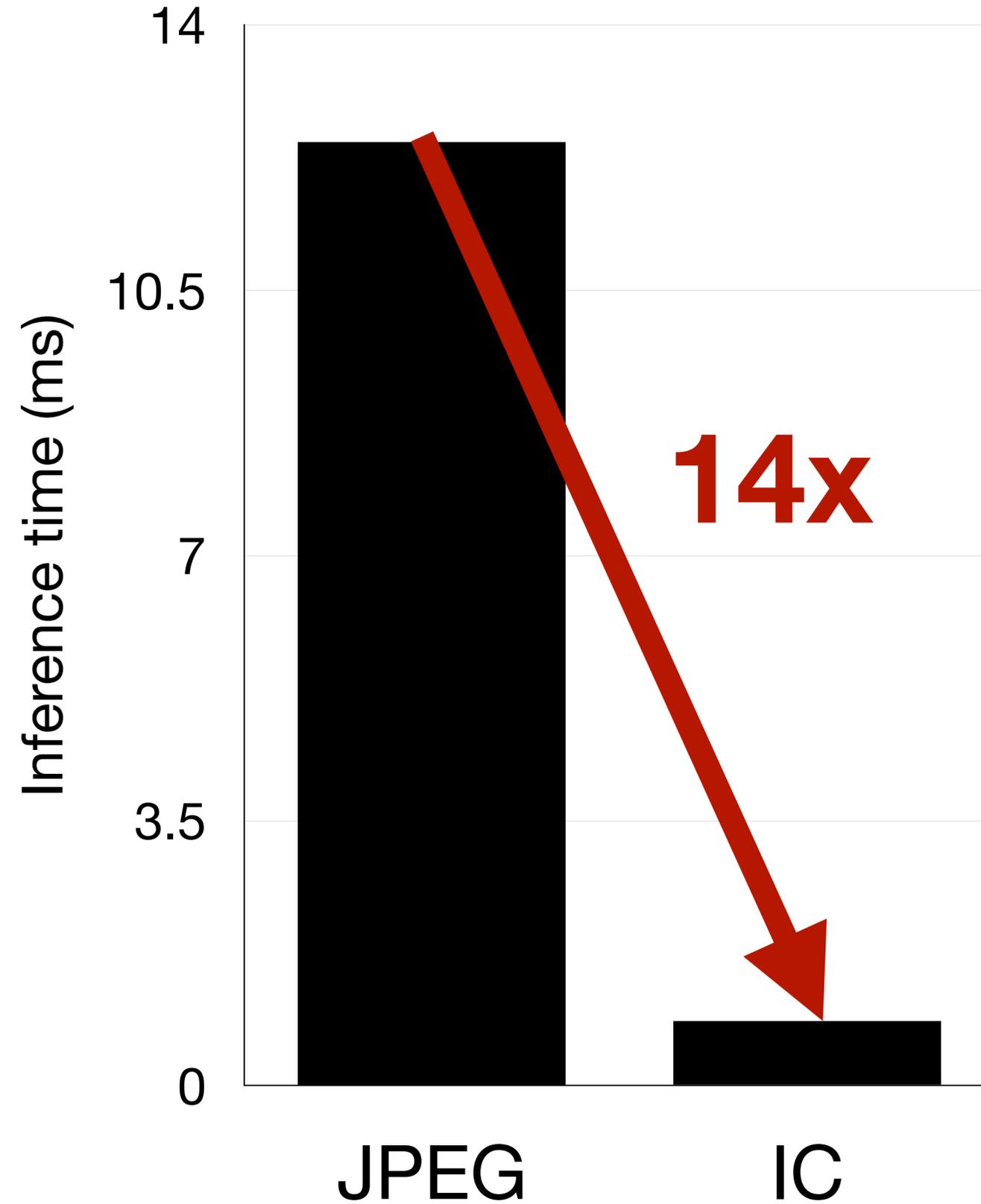
**No loss w/
9x gain**



**No loss w/
7x gain**

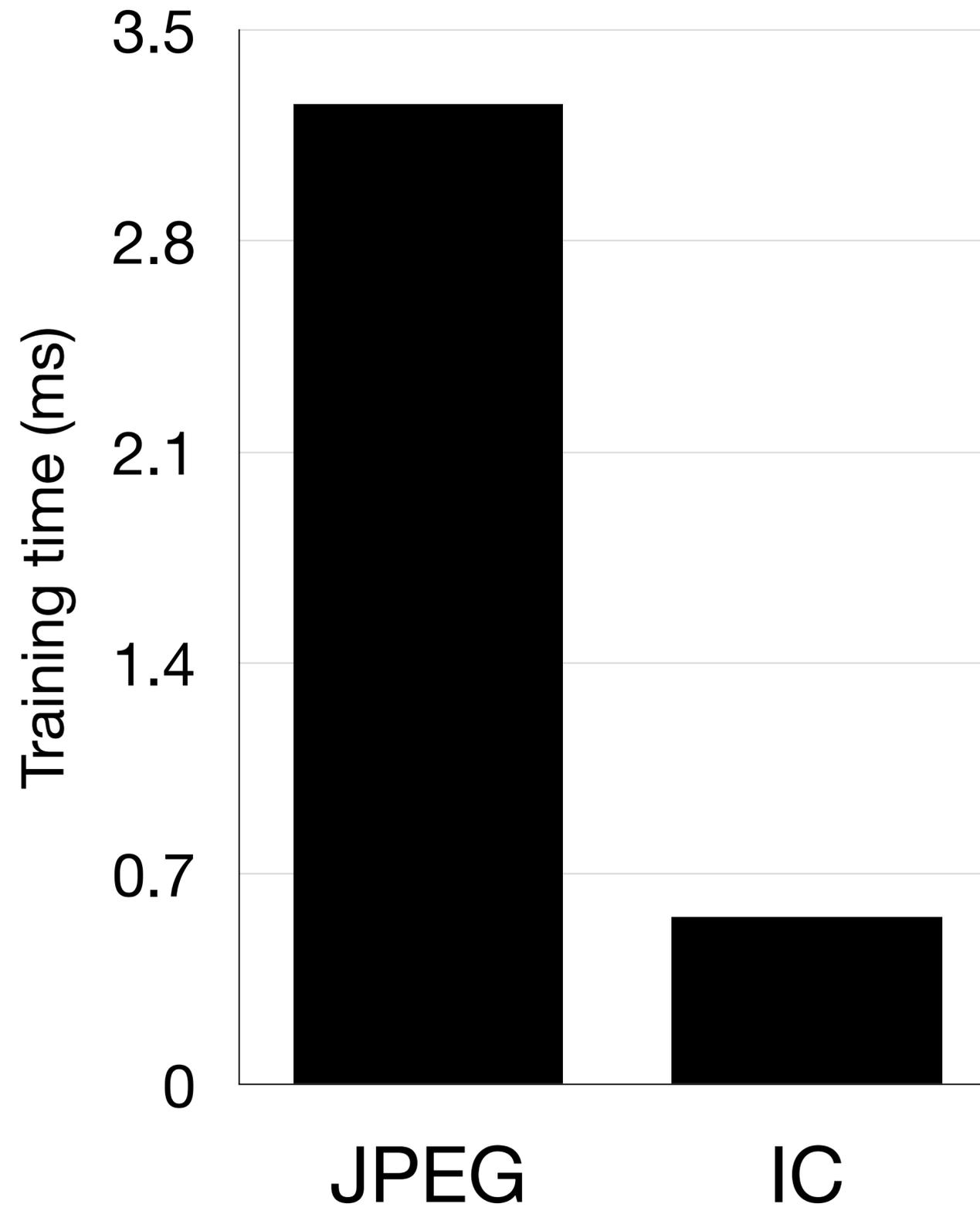
IC improves on cheap CPUs

*Blood-cell
MobileNetV3
Intel CPU
PyTorch v1*



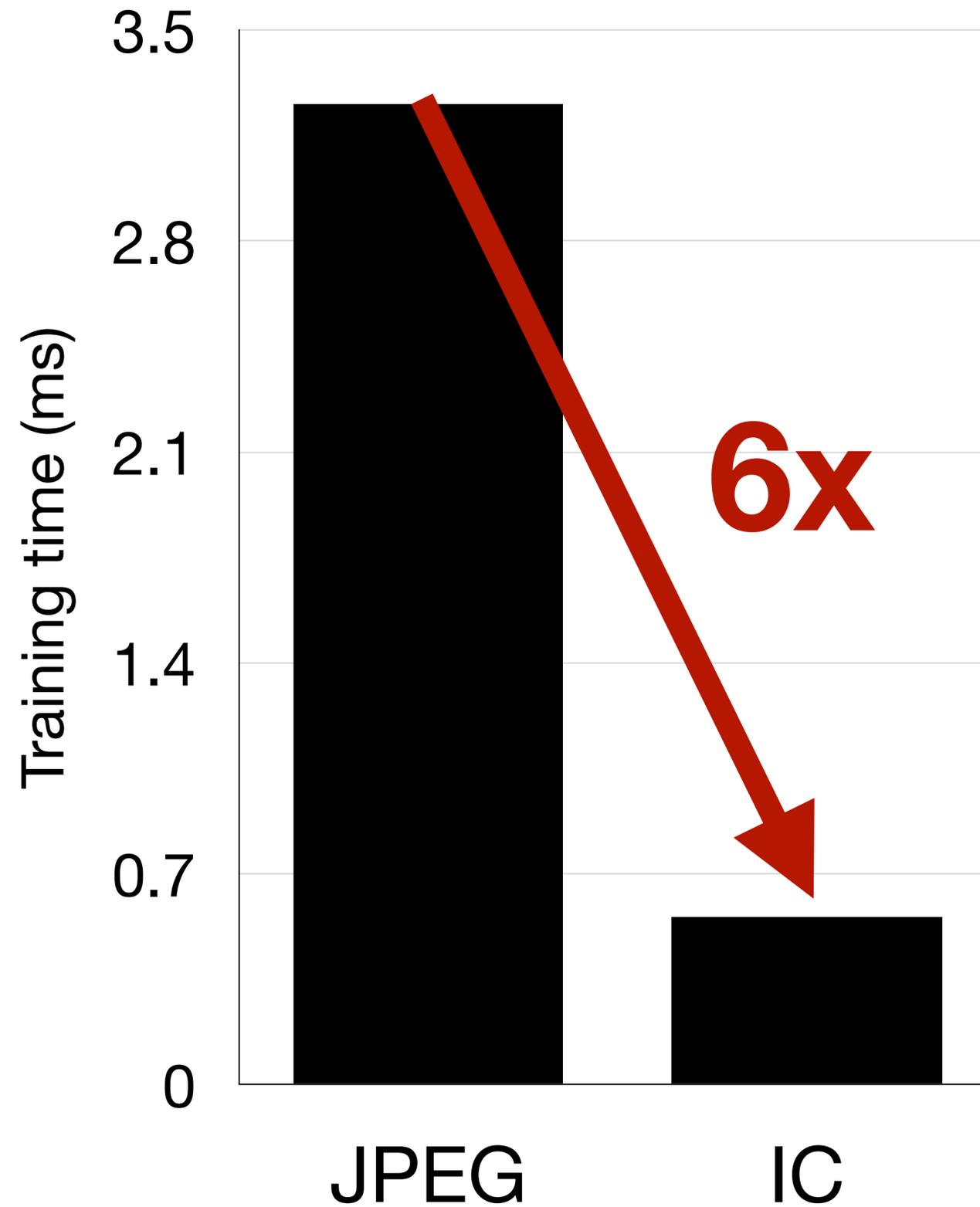
IC reduces training time

*Blood-cell, ResNet50
A100, PyTorch v1*



IC reduces training time

*Blood-cell, ResNet50
A100, PyTorch v1*



Generating design spaces for whole systems.

Reasoning: rules, math and ML to create entirely new designs

Generating design spaces for whole systems.

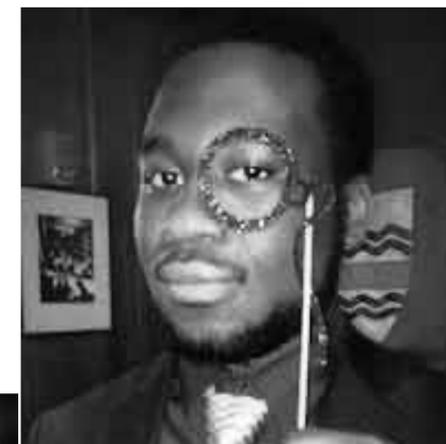
Reasoning: rules, math and ML to create entirely new designs

Primer: **The Periodic Table of Data Structures**

IEEE Data Eng. Bull. 41(3), 2018




DASlab
 @ Harvard SEAS
daslab.seas.harvard.edu





DASlab
@ Harvard SEAS

daslab.seas.harvard.edu

THANKS!

