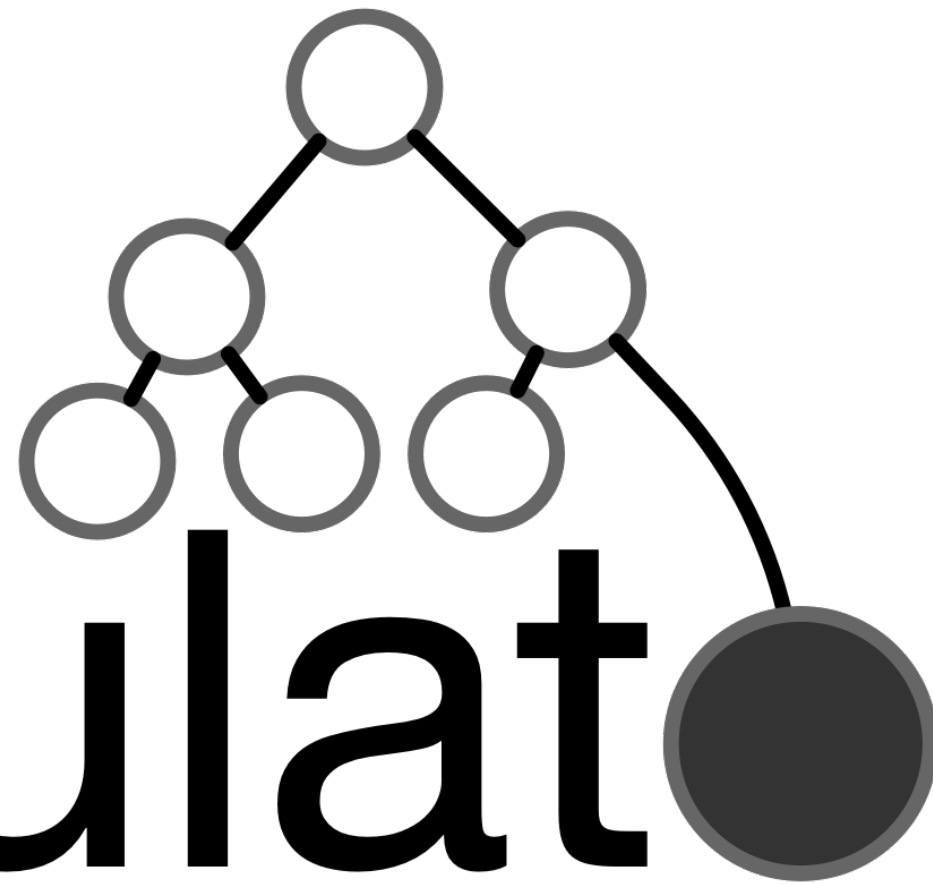Data Calculator
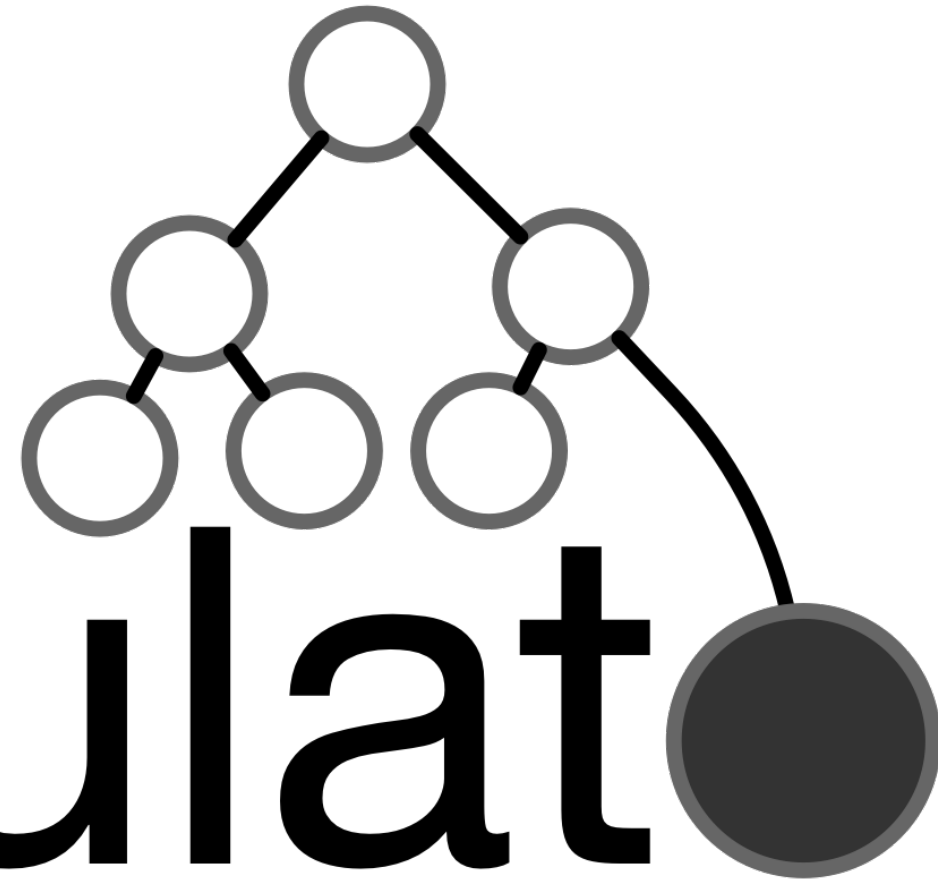
Stratos Idreos & Data Systems Lab

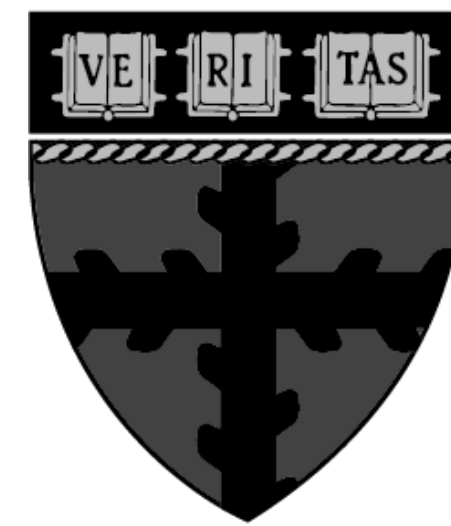@ Harvard SEAS

# Data Calculator

Stratos Idreos & Data Systems Lab
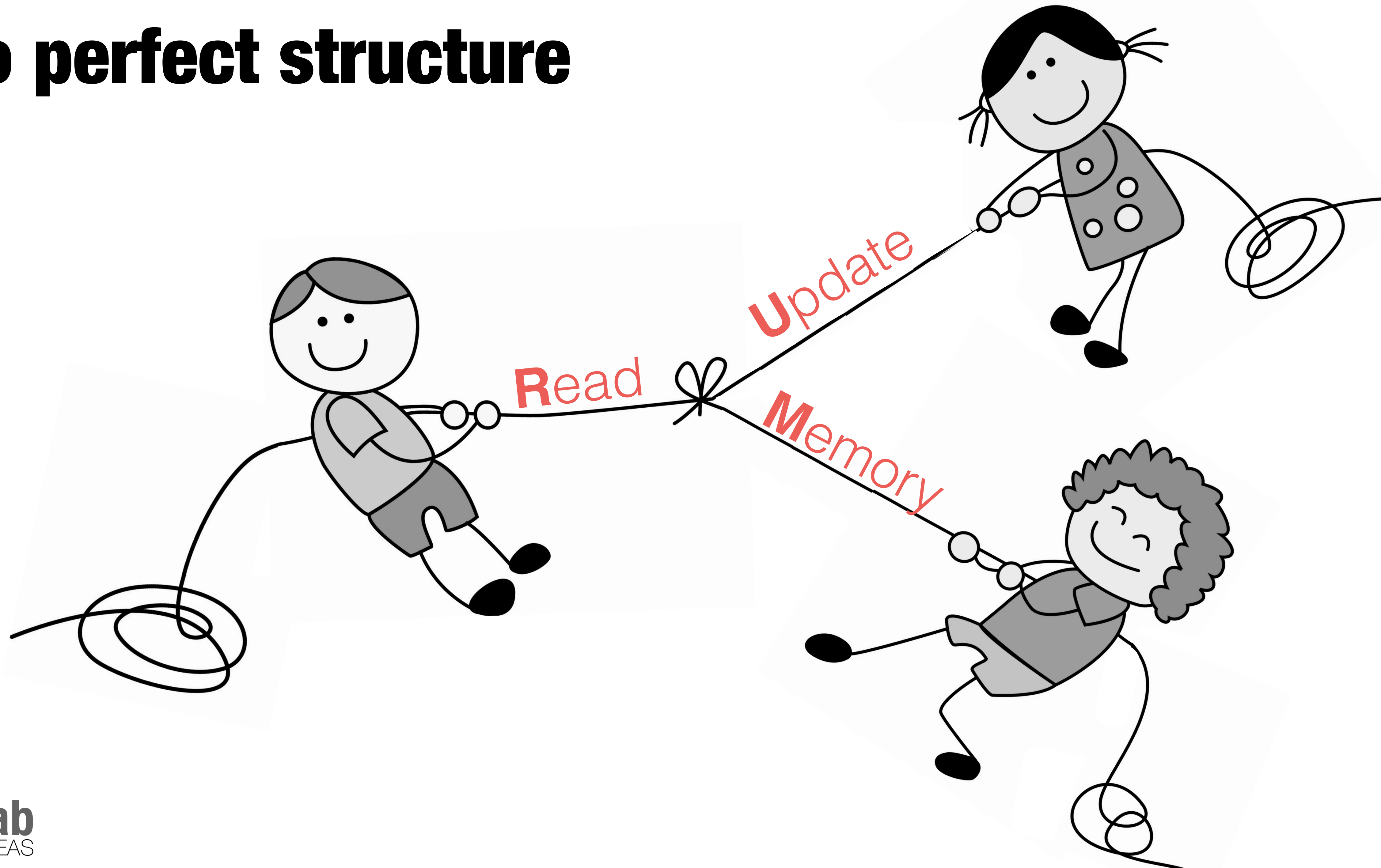
DASlab @ Harvard SEAS

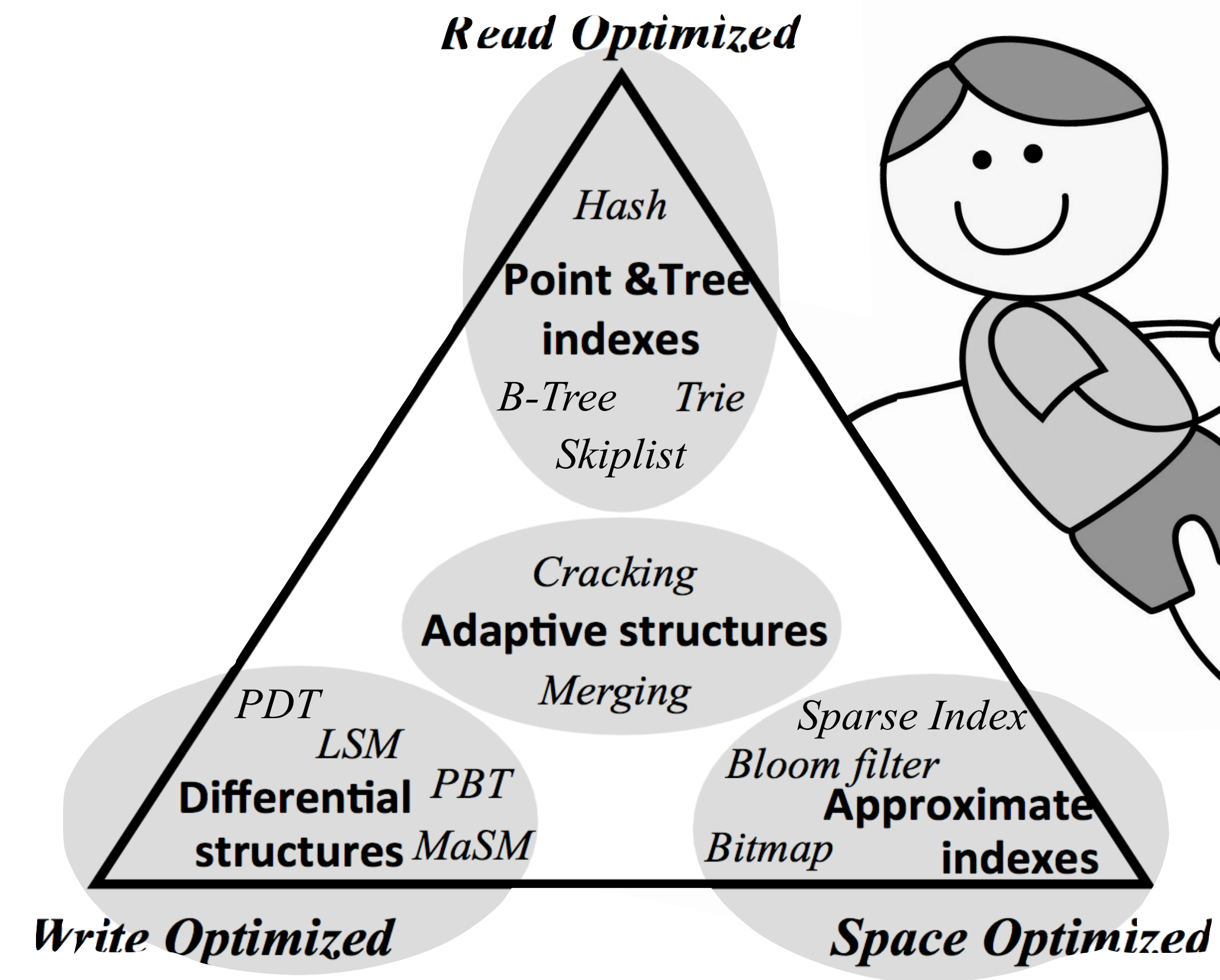WHAT IF WE COULD **REASON** ABOUT THE **DESIGN SPACE** OF DATA STRUCTURES?

# 1. no perfect structure



Read · Update · Memory
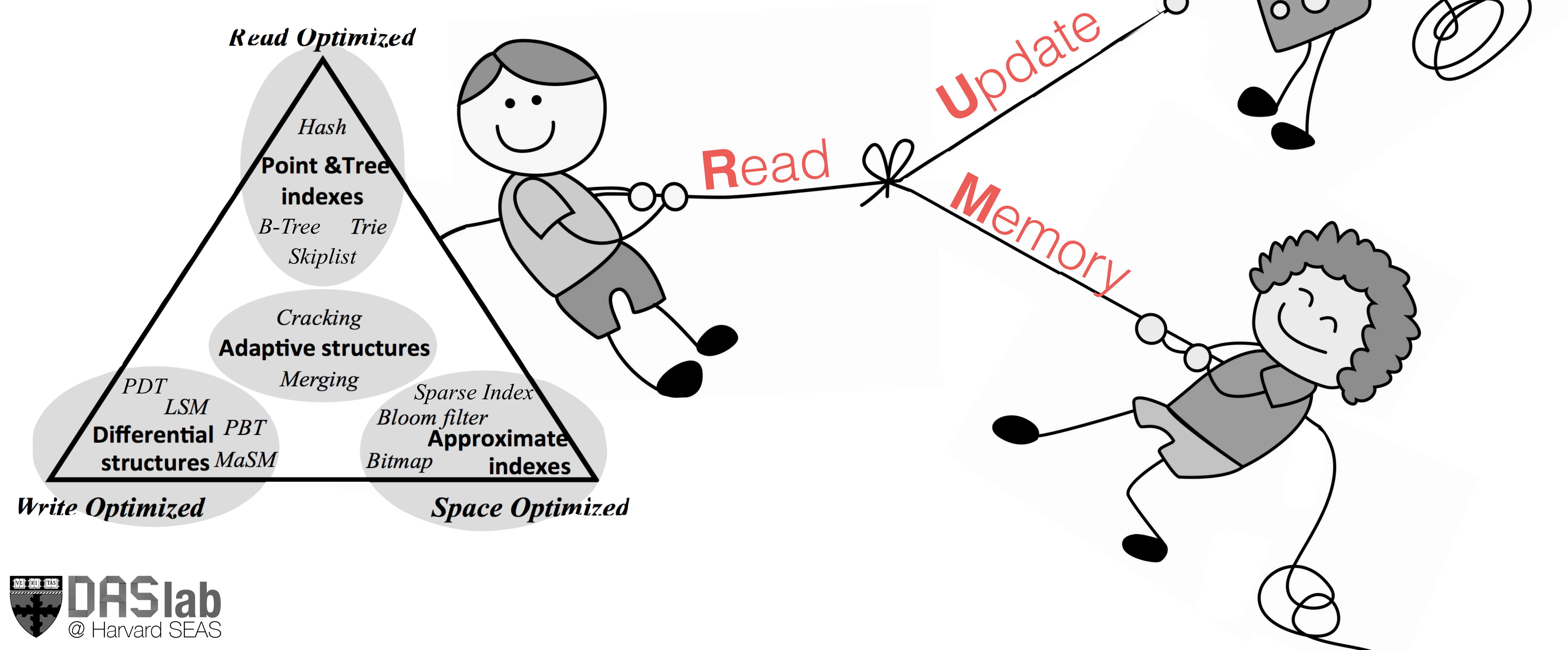
# 1. no perfect structure



Read Optimized

Hash

**Point &Tree indexes**

B-Tree    Trie
Skiplist

Cracking
**Adaptive structures**
Merging

PDT
LSM
**Differential** PBT
**structures** MaSM

Sparse Index
Bloom filter
**Approximate indexes**
Bitmap

Write Optimized

Space Optimized

**R**ead

**U**pdate

**M**emory

# 1. no perfect structure
# 2. apps & h/w evolve



Read Optimized

Hash
**Point &Tree indexes**
B-Tree    Trie
Skiplist

Cracking
**Adaptive structures**
Merging

PDT
LSM
**Differential structures**    PBT
MaSM

Sparse Index
Bloom filter
**Approximate indexes**
Bitmap

Write Optimized

Space Optimized

**R**ead

**U**pdate

**M**emory

DASlab
@ Harvard SEAS

the dinning systems designers

phd
phd
phd
phd
phd
phd
phd
phd

GET **N** SMART PEOPLE

GIVE THEM **T** TIME

HOPE FOR THE BEST

DASlab

phd

phd

phd

phd

phd

phd

phd

phd

GET *N* SMART PEOPLE

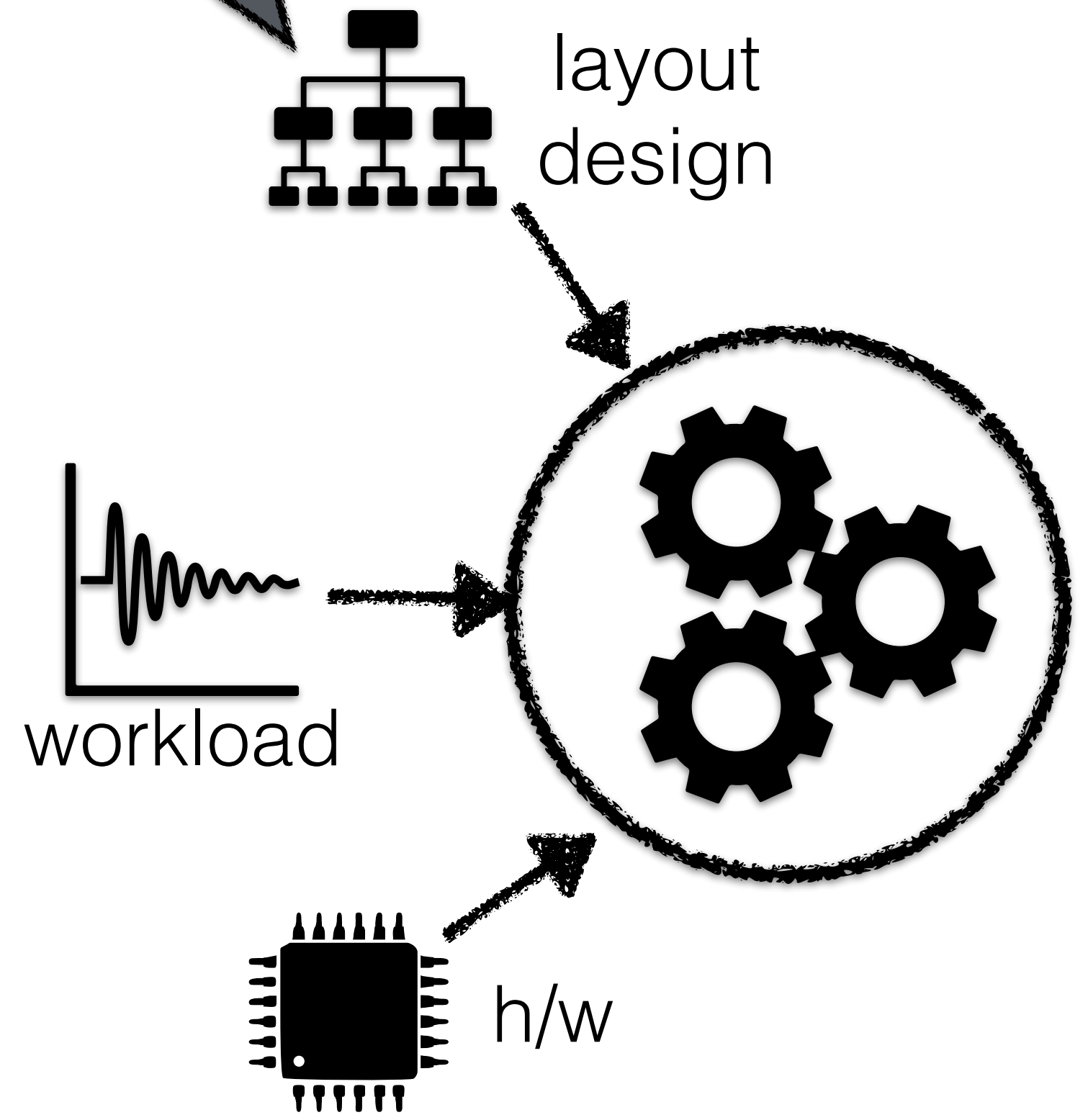GIVE THEM *T* TIME

HOPE FOR THE BEST

DASlab

# AUTO DESIGN

"**IS THERE A CALCULUS OF DATA STRUCTURES** by which one can choose the appropriate representation and techniques for a given problem?" (SIAM,1978)
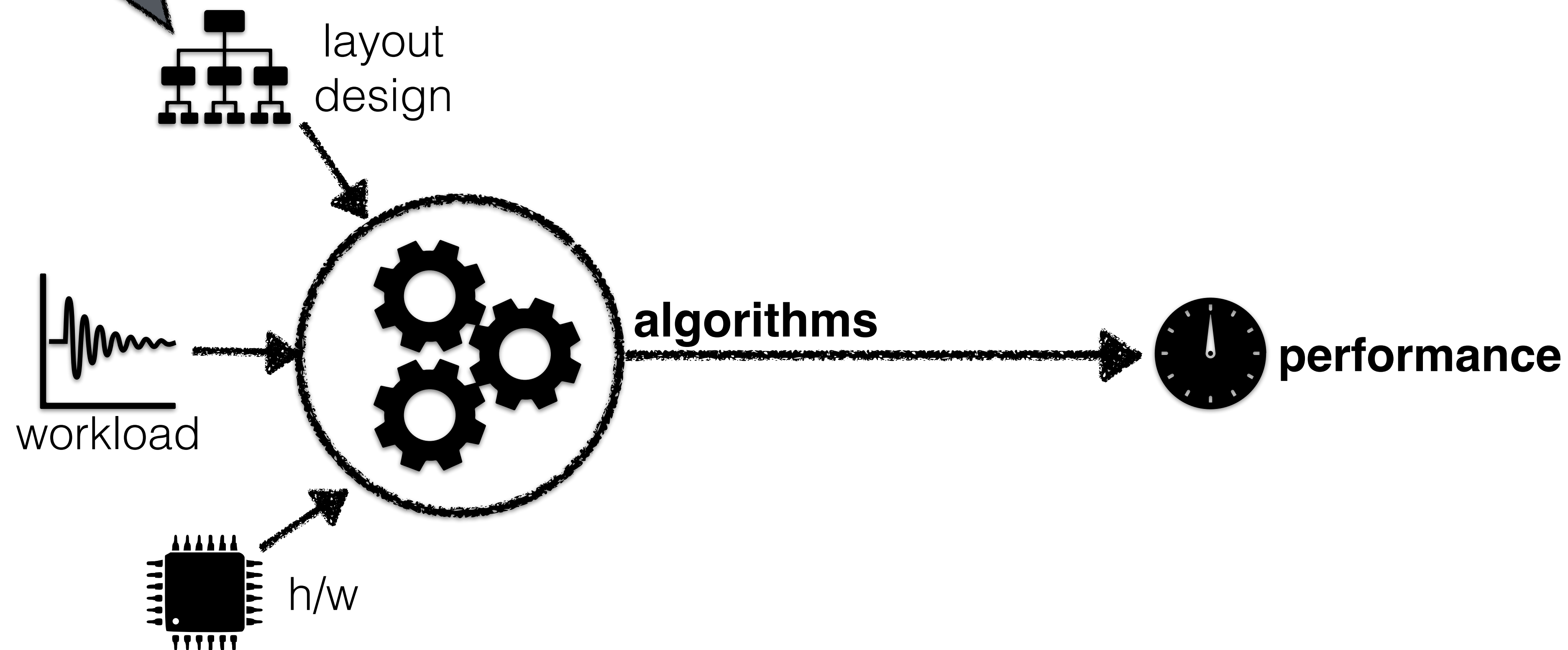
ROBERT TARJAN

DASlab
@ Harvard SEAS

# HOW MANY AND WHICH?

# COMPUTE PERFORMANCE?

DASlab
@ Harvard SEAS

DESIGN SPACE
OF POSSIBLE
STORAGE LAYOUTS

Data Calculator

layout design

workload

h/w

DASlab
@ Harvard SEAS

Data Calculator

What-if Design    Auto-Design    Self-designing Systems

DASlab
@ Harvard SEAS

DESIGN SPACE | COST SYNTHESIS | HOW TO USE

DASlab
@ Harvard SEAS

# 1

## EACH DESIGN: A **SET** OF CONCEPTS

{partitioning, links, fence pointers,…}

# 2

DESIGN {
**COMBINATION**
of existing concepts

**TUNING**
of existing concepts

**NEW**
concept

# 3

(ALMOST) ALL DESIGNS ARE A COMBINATION/TUNING OF **EXISTING CONCEPTS**

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

Nikos Kazantzakis

if we know the **fundamental** building blocks,

if we know the **fundamental** building blocks,
how they combine and their properties,

if we know the **fundamental** building blocks,
how they combine and their properties,

then we can **automate** the discovery of
novel combinations and tunings

DESIGN SPACE

if we know the **fundamental** building blocks, how they combine and their properties,

then we can **automate** the discovery of novel combinations and tunings
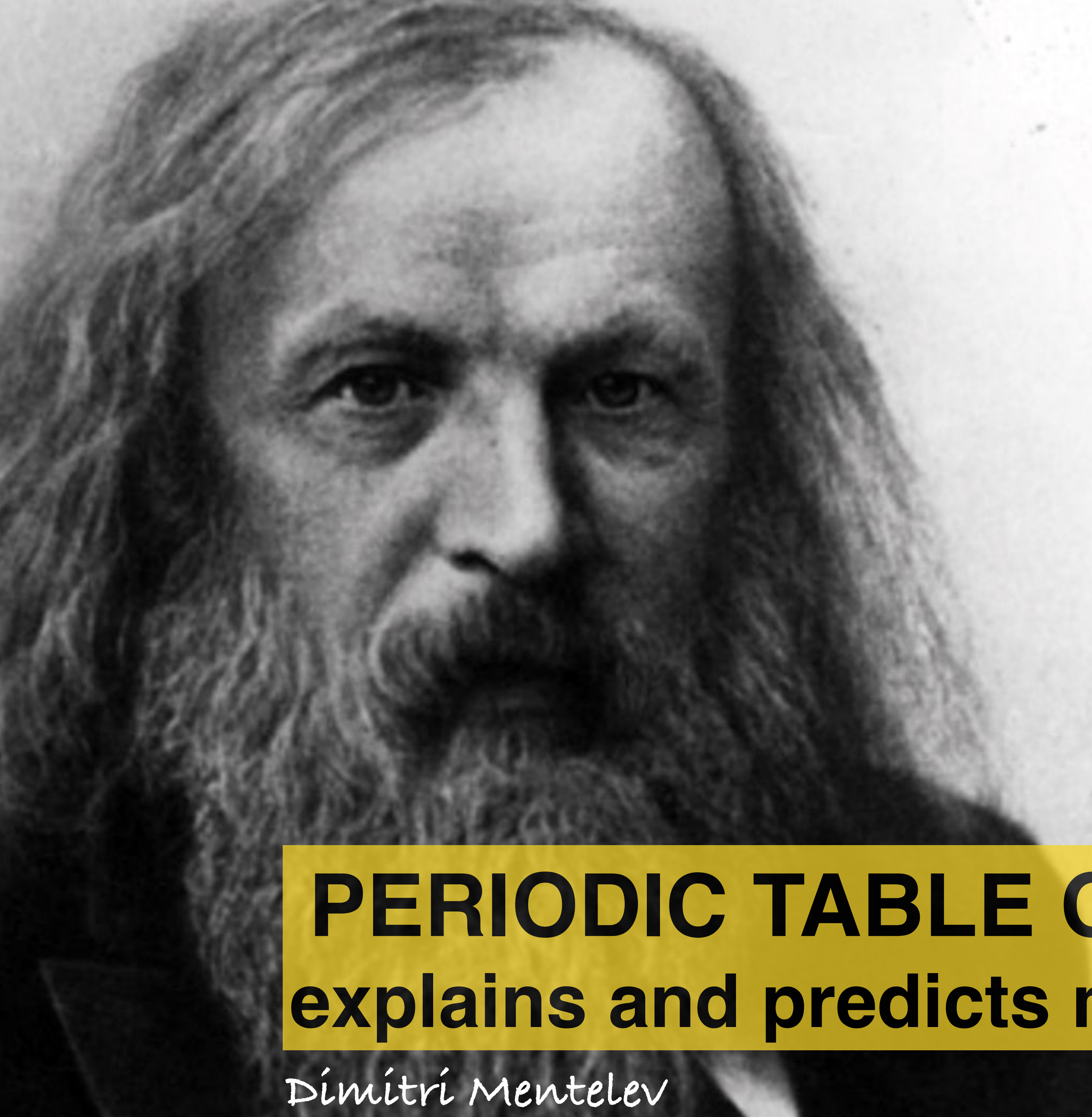
NAVIGATION

**structures elements** based on atomic number, electron configuration, and recurring chemical properties

**PERIODIC TABLE OF ELEMENTS**
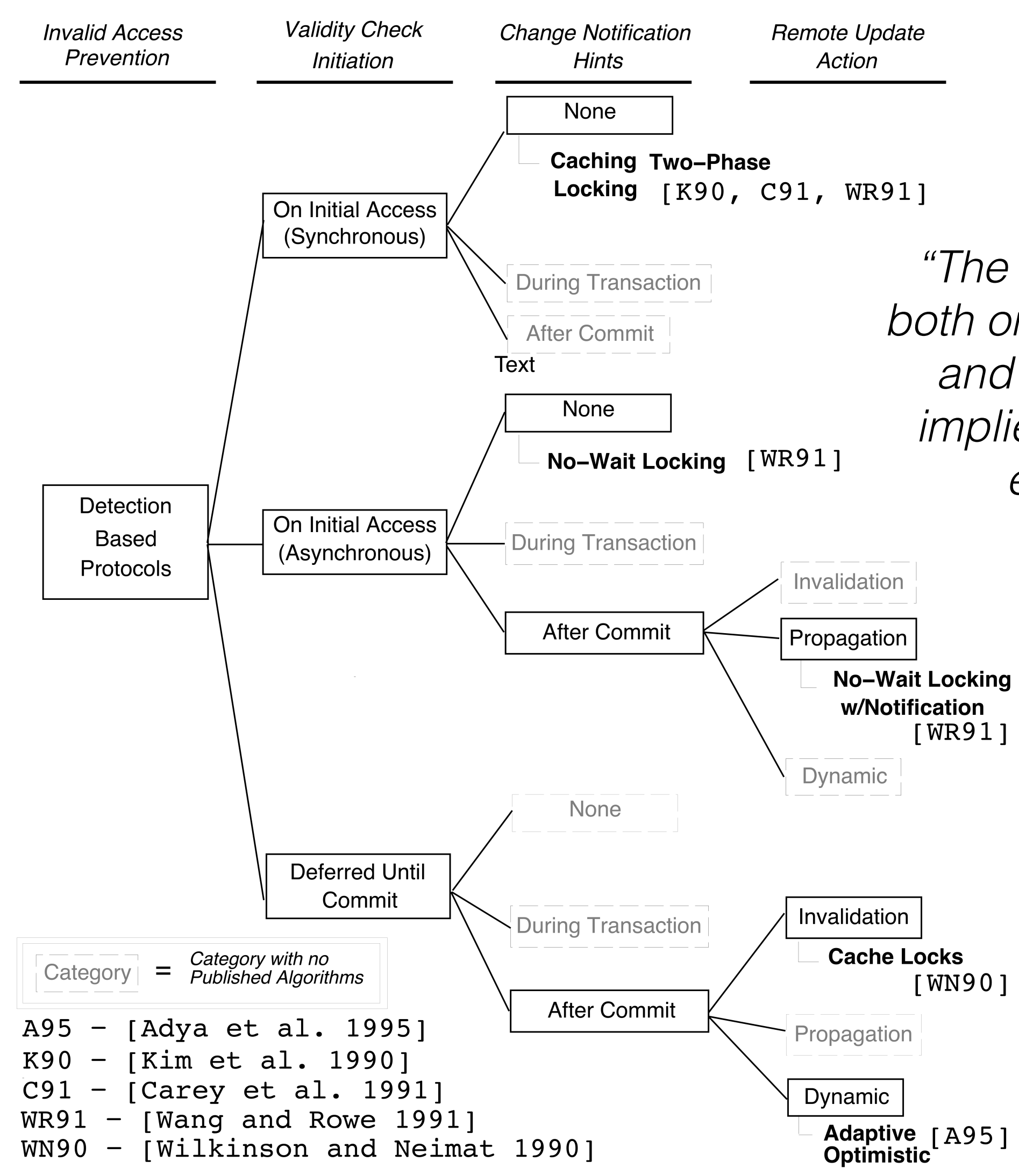**explains and predicts missing elements**

Dimitri Mentelev

Kosuke Morita

TAXONOMY OF COMPLEX ALGORITHMS
transactional cache consistency maintenance

Mike Franklin

trial and error

DASlab
@ Harvard SEAS

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

trial and error

DASlab
@ Harvard SEAS

# MAP LAYOUT FIRST

**(1)** KNOWN DESIGNS

**(2)** OPEN QUESTIONS

*trial and error*

# **{arrays, logs, lsm-trees, b-trees}, filters, bitmaps, compression, stats**

e.g., 1000x NoSQL k-v: bloom filter bits, merging policy
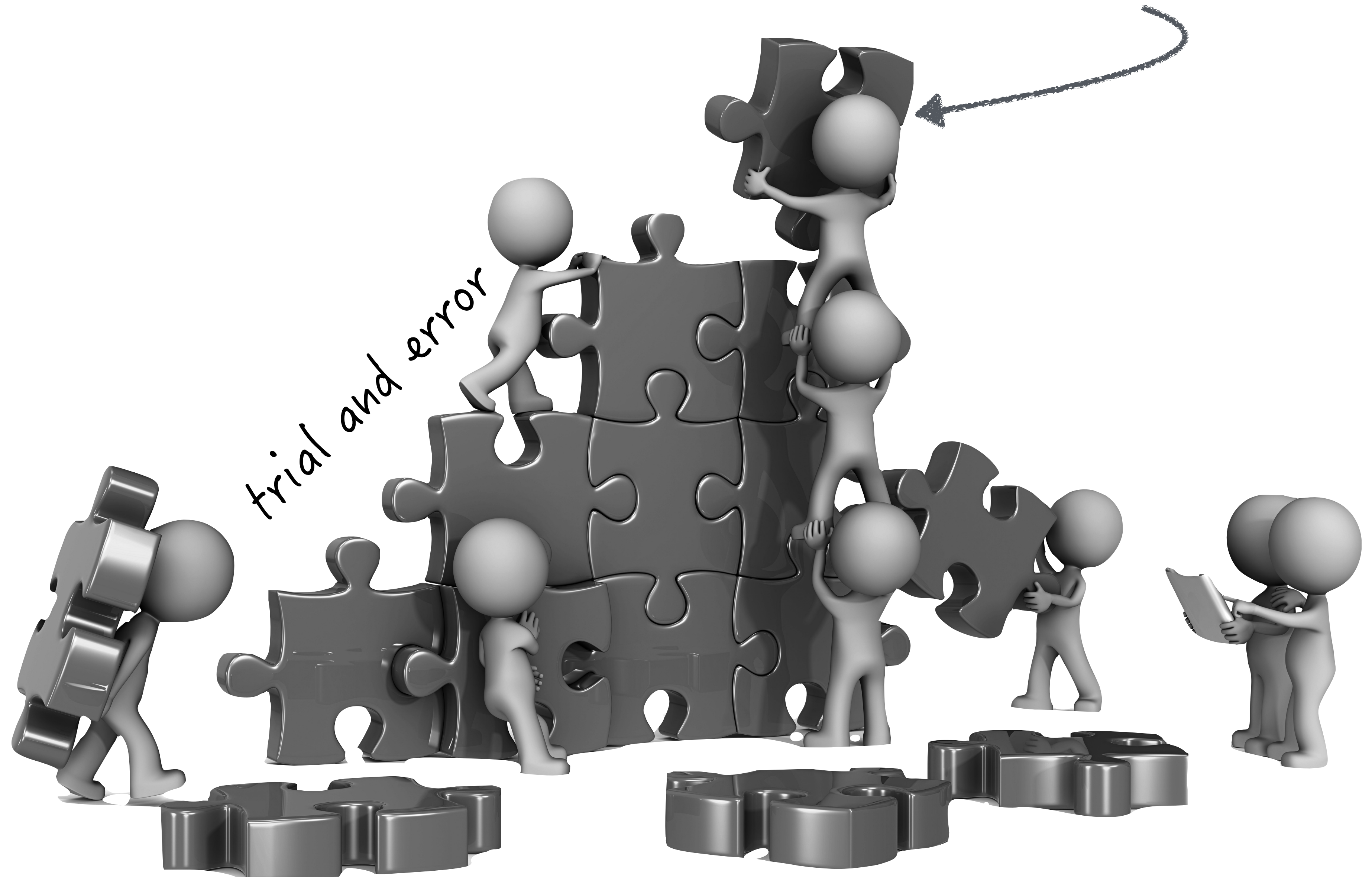e.g., access path selection: scans vs b-tree depends on concurrency
e.g., robust scans with value by value lossy compression
e.g., updatable bitmap indexes
e.g., fast statistics/ML
…

# EXAMPLE: The design space of **NoSQL Key-value Stores**

MEMORY

DISK

**insert (key-value)**

buffer

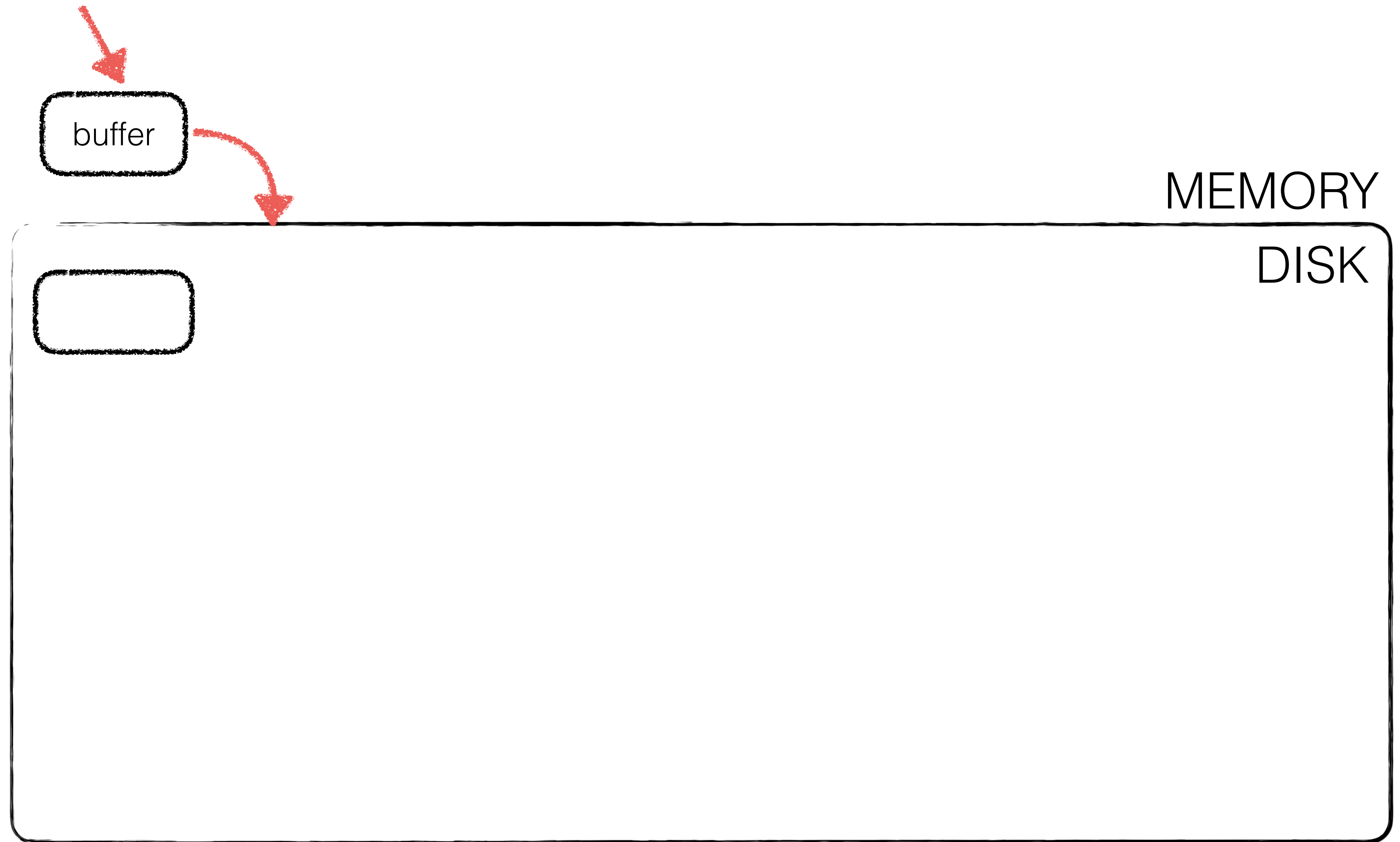MEMORY

DISK

**insert (key-value)**

buffer

MEMORY

DISK

**insert (key-value)**

buffer

MEMORY

DISK

Level 1

**insert (key-value)**

buffer

MEMORY

DISK

Level 1

Level 2

**insert (key-value)**

buffer

MEMORY

DISK

Level 1

Level 2

Level 3

**insert (key-value)**

buffer

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

**insert (key-value)**

buffer

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

tiered

leveled

sorted

**insert (key-value)**

buffer

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

DASlab
@ Harvard SEAS

buffer

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

[1,0,0,1,1,1]
hash fun.

[min-max]

**get (key)**

buffer

bloom
filters

fence
pointers

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

...

...

DASlab
@ Harvard SEAS

# Instant writes, but overtime have to be merged >1

[1,0,0,1,1,1]
hash fun.

[min-max]

buffer

bloom
filters

fence
pointers

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

Instant writes, but overtime have to be merged >1
Temporal partitioning = recent data to the top

[1,0,0,1,1,1]
hash fun.

[min-max]

buffer

bloom
filters

fence
pointers

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

DASlab
@ Harvard SEAS

Instant writes, but overtime have to be merged >1
Temporal partitioning = recent data to the top
But duplicates across levels (memory amplification)

[1,0,0,1,1,1]
hash fun.

[min-max]

bloom
filters

fence
pointers

...

...

buffer

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

Instant writes, but overtime have to be merged >1
Temporal partitioning = recent data to the top
But duplicates across levels (memory amplification)
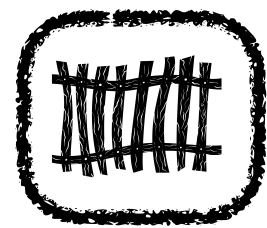Reads can get away with 0-1 disk page read per level
Run layout affect read/write costs

[1,0,0,1,1,1]
hash fun.

[min-max]

buffer

bloom
filters

fence
pointers

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

DASlab
@ Harvard SEAS

Instant writes, but overtime have to be merged >1
Temporal partitioning = recent data to the top
But duplicates across levels (memory amplification)
Reads can get away with 0-1 disk page read per level
Run layout affect read/write costs
Size ratio affects everything
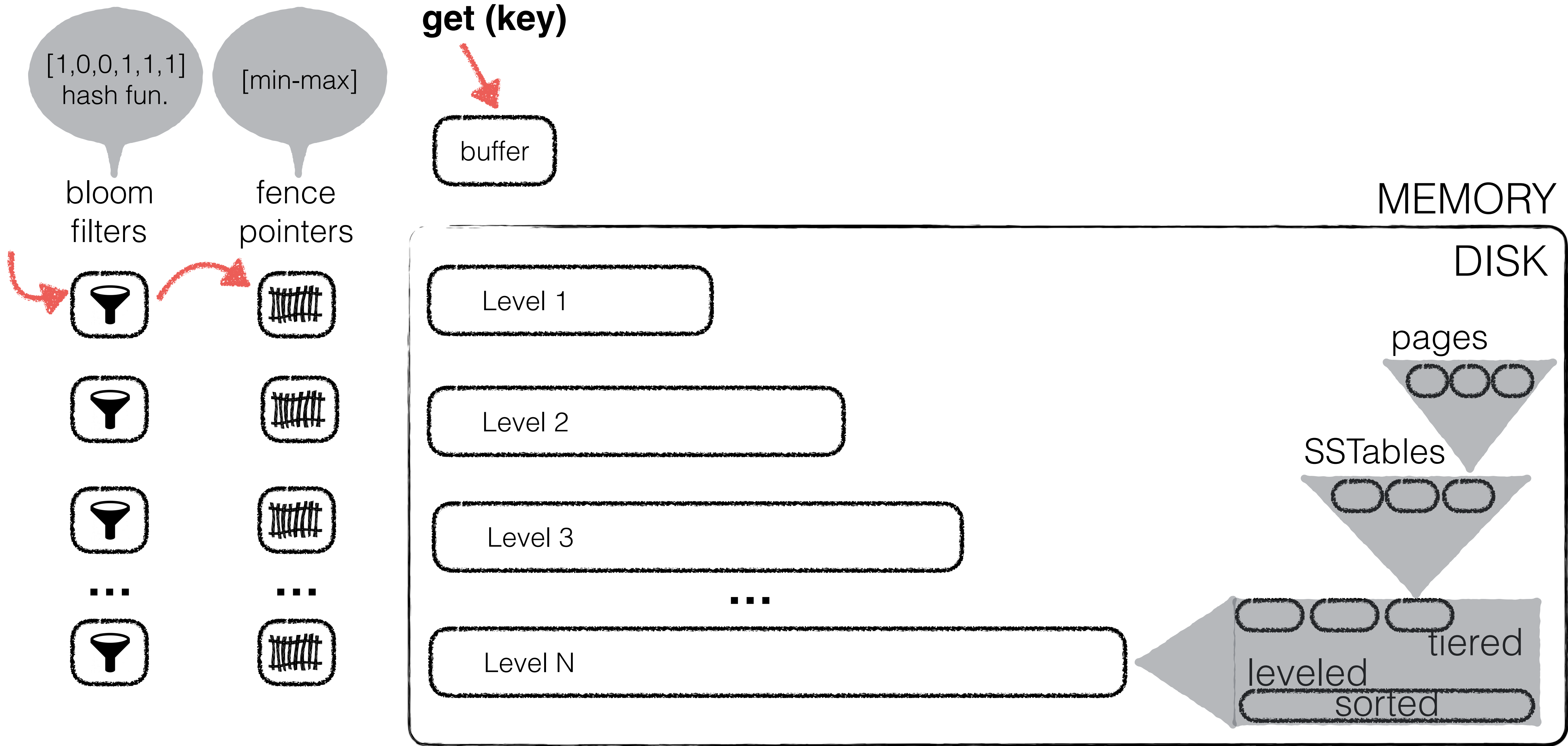
bloom filters

[1,0,0,1,1,1] hash fun.

fence pointers

[min-max]

buffer

MEMORY

DISK

Level 1

Level 2

Level 3

...

Level N

pages

SSTables

tiered

leveled

sorted

DASlab
@ Harvard SEAS

size ratio

merge policy

filters bits per entry

size of (page, buffer, ..)

internal k-v layout

**For every design principle X:**

1. Which are X's meaningful values?

2. How does X affect read, update and memory amplification?

3. Should X be a design principle or can it be optimized out?

**bits per entry in filters**

monkey@SIGMOD2017

**M**onkey: **O**ptimal **N**avigable

**Key**-Value Store

**bits per entry in filters**
**monkey**@SIGMOD2017
**M**onkey: **O**ptimal **N**avigable
**Key**-Value Store

standard design: fixed per run

at most one I/O per level

worst case I/O: sum of false positive rates

buffer

Level 1

Level 2

. . .

Level N

**Bits per entry in bloom filter**

**bits per entry in filters**
**monkey**@SIGMOD2017
**M**onkey: **O**ptimal **N**avigable
**Key**-Value Store

minimize sum of FPRs
***"move" filters memory budget
from big to small levels***

buffer

Level 1

Level 2

. . .

Level N

**Bits per entry in bloom filter**

bits per entry in filters

monkey@SIGMOD2017
Monkey: Optimal Navigable
Key-Value Store

minimize sum of FPRs
*"move" filters memory budget
from big to small levels*

buffer

Level 1

Level 2

. . .

Level N

**Bits per entry in bloom filter**

bits per entry in filters

monkey@SIGMOD2017

Monkey: Optimal Navigable

Key-Value Store

minimize sum of FPRs
*"move" filters memory budget
from big to small levels*

buffer

Level 1

Level 2

. . .

Level N

**Bits per entry in bloom filter**

lookup cost

WiredTiger

Cassandra, HBase

RocksDB, LevelDB

*monkey*

update cost

minimize sum of FPRs
***"move" filters memory budget
from big to small levels***

buffer

Level 1

Level 2

. . .

Level N

**Bits per entry in bloom filter**

**bits per entry in filters**
**monkey**@SIGMOD2017
**M**onkey: **O**ptimal **N**avigable
**Key**-Value Store

Lookup latency (ms)

□ LevelDB
○ Monkey

$\approx 0.2$ I/Os per lookup

data entries

uniform, zero result, point queries, entry size=1KB

**merge policy**

**dostoevsky**@SIGMOD2018

**Do**stoevsky: **S**pace-**T**ime **O**ptimized

**Ev**olvable **S**calable **K**ey-Value Store

standard design: fixed across levels

Leveled helps reads; tiered helps writes

buffer

Level 1

Level 2

. . .

Level N

Bits per entry in bloom filter

**merge policy**
**dostoevsky**@SIGMOD2018
**Do**stoevsky: **S**pace-**T**ime **O**ptimized
**Ev**olvable **S**calable **K**ey-Value Store

merge policy

**dostoevsky**@SIGMOD2018

**Do**stoevsky: **S**pace-**T**ime **O**ptimized
**Ev**olvable **S**calable **K**ey-Value Store

*hybrid merge policy*
*tiered for small levels*

buffer

Level 1

Level 2

. . .

Level N

Bits per entry in bloom filter

*hybrid merge policy*
*tiered for small levels*

buffer

Level 1

Level 2

. . .

Level N

Bits per entry in bloom filter

**hybrid merge policy**
*tiered for small levels*

buffer

Level 1

Level 2

. . .

Level N

Bits per entry in bloom filter

**merge policy**
**dostoevsky**@SIGMOD2018
**Do**stoevsky: **S**pace-**T**ime **O**ptimized
**Ev**olvable **S**calable **K**ey-Value Store

lookup cost

WiredTiger

Cassandra, HBase

RocksDB, LevelDB

**monkey**

**dosto**

update cost

arrays

b-trees

size ratio

**merge policy**

**filters bits per entry**

size of (page, buffer, ..)

internal k-v layout

...

@SIGMOD18

key-value
(read, bulk load)

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

DASlab
@ Harvard SEAS

Labels around the wheel:
valueRetention.function
valueRetention.type
zeroElementNullable
blockAccess.direct
blockAccess.headLink
blockAccess.tailLink
bloom.active
bloom.hashFunctionsNum
bloom.numberOfBits
capacity.function
capacity.type
capacity.value
external.links.next
external.links.prev
fanout.fixedValue
fanout.function
fanout.type
filters.filtersMemLayout
filters.zoneMaps.exact
filters.zoneMaps.max
filters.zoneMaps.min
indirectedPointers
keyRetention.function
keyRetention.type
links.linksMemLayout
links.next
links.prev
links.skipLinks.probability
links.skipLinks.type
log.filtersPerLevel
log.filtersPerRun
log.initialRunSize
log.maxRunsPerLevel
log.mergeFactor
orderMaintenance.type
partitioning.function
partitioning.type
recursionAllowed
retainedDataLayout
sub-block.phys.homog.
sub-block.phys.layout
sub-block.phys.location
utilization.constraint
utilization.function

**Are keys retained?** (yes, no, function)

**Are values retained?**

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



Columnar | Row-wise | Row-Groups | None

Key and value layout

**Are keys retained?** (yes, no, function)

**Are values retained?**

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



Key and value layout

**Fanout** (fixed/functional | unlimited | terminal |)

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

**Are keys retained?** (yes, no, function)

**Are values retained?**

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



Columnar | Row-wise | Row-Groups | None

Key and value layout

**Fanout** (fixed/functional | unlimited | terminal |)

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

**Intra node access** (direct | head_link | tail_link | link_function(func))

**Are keys retained?** (yes, no, function)

**Are values retained?**

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



Key and value layout

**Fanout** (fixed/functional | unlimited | terminal |)

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

**Intra node access** (direct | head_link | tail_link | link_function(func))

**Sub block links** (next | previous | both | none)

**Sub block skip links** (perfect | randomized(prob: double) | function(func) | none)



Immediate: Both

Logical Block

**Are keys retained?** (yes, no, function)

**Are values retained?**

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



Columnar | Row-wise | Row-Groups | None

Key and value layout

**Fanout** (fixed/functional | unlimited | terminal |)

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

**Intra node access** (direct | head_link | tail_link | link_function(func))

**Sub block links** (next | previous | both | none)

**Sub block skip links** (perfect | randomized(prob: double) | function(func) | none)

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

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

**Filters layout** (consolidate | scatter)

**Links layout** (consolidate | scatter)



Consolidate | Scatter

**Are keys retained?** (yes, no, function)

**Are values retained?**

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



Columnar | Row-wise | Row-Groups | None

**Key and value layout**

**Fanout** (fixed/functional | unlimited | terminal |)

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

**Intra node access** (direct | head_link | tail_link | link_function(func))

**Sub block links** (next | previous | both | none)

**Sub block skip links** (perfect | randomized(prob: double) | function(func) | none)

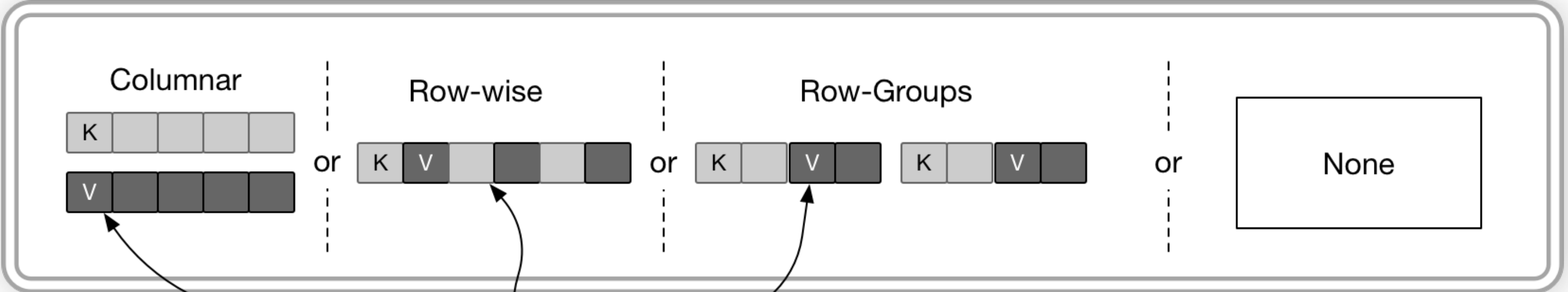**Zone Maps** (min | max | both | exact | off)

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

**Filters layout** (consolidate | scatter)

**Links layout** (consolidate | scatter)

**Physical location** (inline | pointed | double- pointed)
**Physical layout** (BFS | scatter)



Inline — Logical Block — Logical Block

Pointed — Logical Block — Logical Block

# SETS OF CONCEPTS

# SETS OF CONCEPTS    POSSIBLE NODE DESIGNS

SETS OF CONCEPTS     POSSIBLE NODE DESIGNS     POSSIBLE STRUCTURES

sorted
bloom
filter bits
zone map
link
children
layout

ARRAY    LINKED-LIST
HASH-TABLE QUEUE
B-TREE  SKIP-LIST
BINARY-TREE
MASSTREE
CSB-TREE
FAST
TRIE

all

DASlab
@ Harvard SEAS

# the periodic table of data structures

*classes of designs*

*classes of primitives*

| | B-trees & Variants | Tries & Variants | LSM-Trees & Variants | Differential Files | Membership Tests | Zone maps & Variants | Bitmaps & Variants | Hashing | Base Data & Columns | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Partitioning** | DONE | DONE | DONE | | | | | DONE | DONE | R U M |
| **Logarithmic Design** | DONE | DONE | DONE | | | | | | | R U M |
| **Fractional Cascading** | DONE | | DONE | DONE | | | | | | R U M |
| **Log-Structured** | DONE | | DONE | DONE | | | | | | R U M |
| **Buffering** | DONE | | | DONE | | | DONE | | | R U M |
| **Differential Updates** | DONE | | | DONE | | | | | | R U M |
| **Sparse Indexing** | DONE | | | | DONE | DONE | | | | R U M |
| **Adaptivity** | DONE | | | | | | | | DONE | |

DASlab
@ Harvard SEAS

# the periodic table of data structures

*classes of designs*

*classes of primitives*

| | B-trees & Variants | Tries & Variants | LSM-Trees & Variants | Differential Files | Membership Tests | Zone maps & Variants | Bitmaps & Variants | Hashing | Base Data & Columns | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Partitioning** | DONE | DONE | DONE | | | | | DONE | DONE | ↓↑↑↑ RUM |
| **Logarithmic** | DONE | DONE | DONE | | | | | | | ↓↓↑ |

**PAPER MACHINE**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Differential Updates** | DONE | | | DONE | | | | | | ↑↓↓ RUM |
| **Sparse Indexing** | DONE | | | | DONE | DONE | | | | ↓↕↑ RUM |
| **Adaptivity** | DONE | | | | | | | | DONE | |

DASlab
@ Harvard SEAS

(10^9)

(10^24)

(10^32, 2-node)
(10^48, 3-node)

PEOPLE ON EARTH

STARS IN THE SKY

POSSIBLE DATA STRUCTURES

DASlab
@ Harvard SEAS

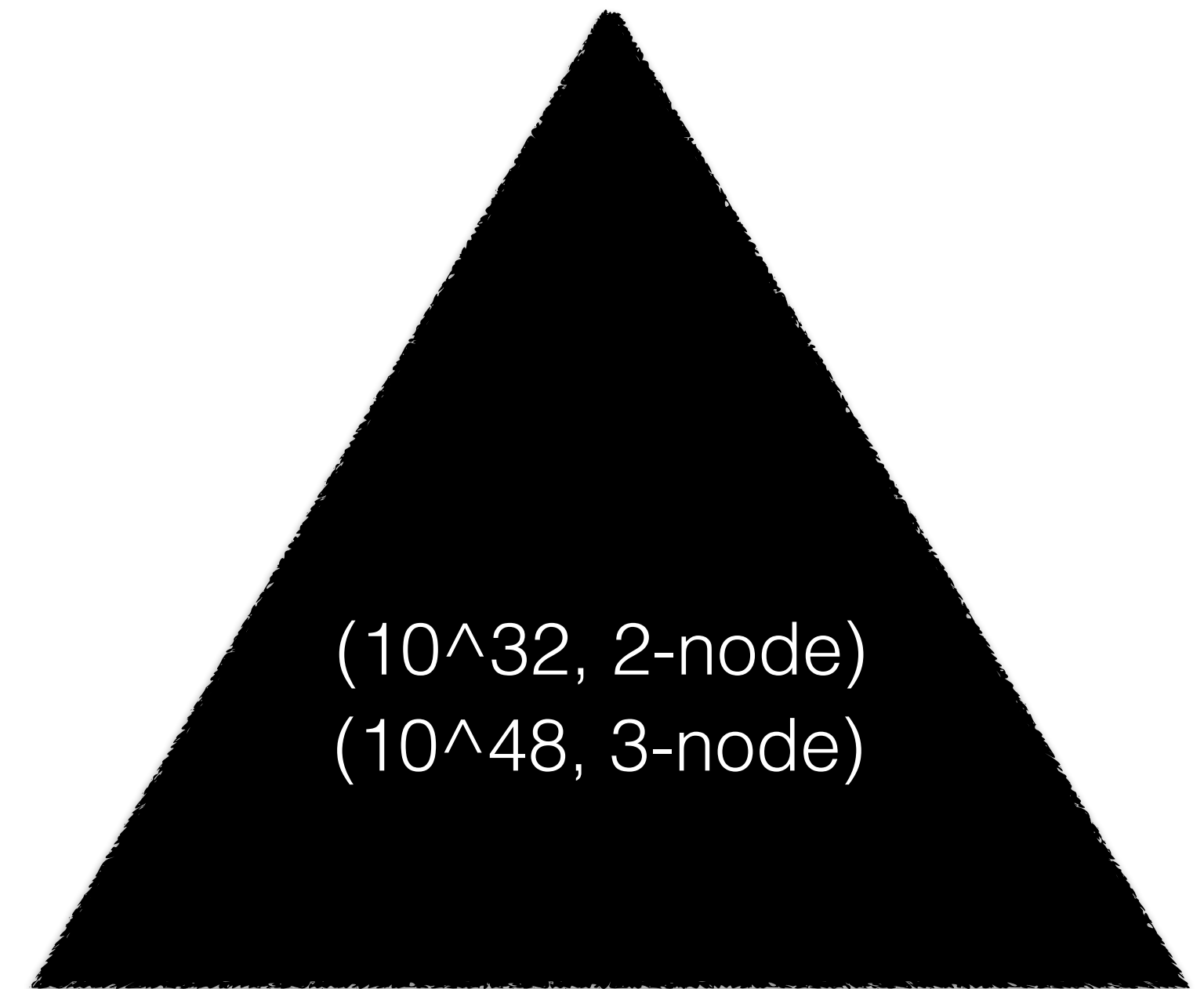~5K since the dawn of CS

(10^32, 2-node)
(10^48, 3-node)

(10^9)

(10^24)

PEOPLE ON EARTH

STARS IN THE SKY

POSSIBLE DATA STRUCTURES

DESIGN SPACE

if we know the fundamental building blocks,
how they combine and their **properties**,

then we can automate the discovery of
novel combinations and tunings

DESIGN SPACE

DASlab
@ Harvard SEAS

DESIGN SPACE

COST SYNTHESIS

HOW TO USE

# HOW TO JUDGE A DESIGN?

**1**

**COMPLEXITY ANALYSIS**

**2**

**IMPLEMENTATION & TESTING**

**3**

**GENERALIZED MODELS**

# HARD & SLOW

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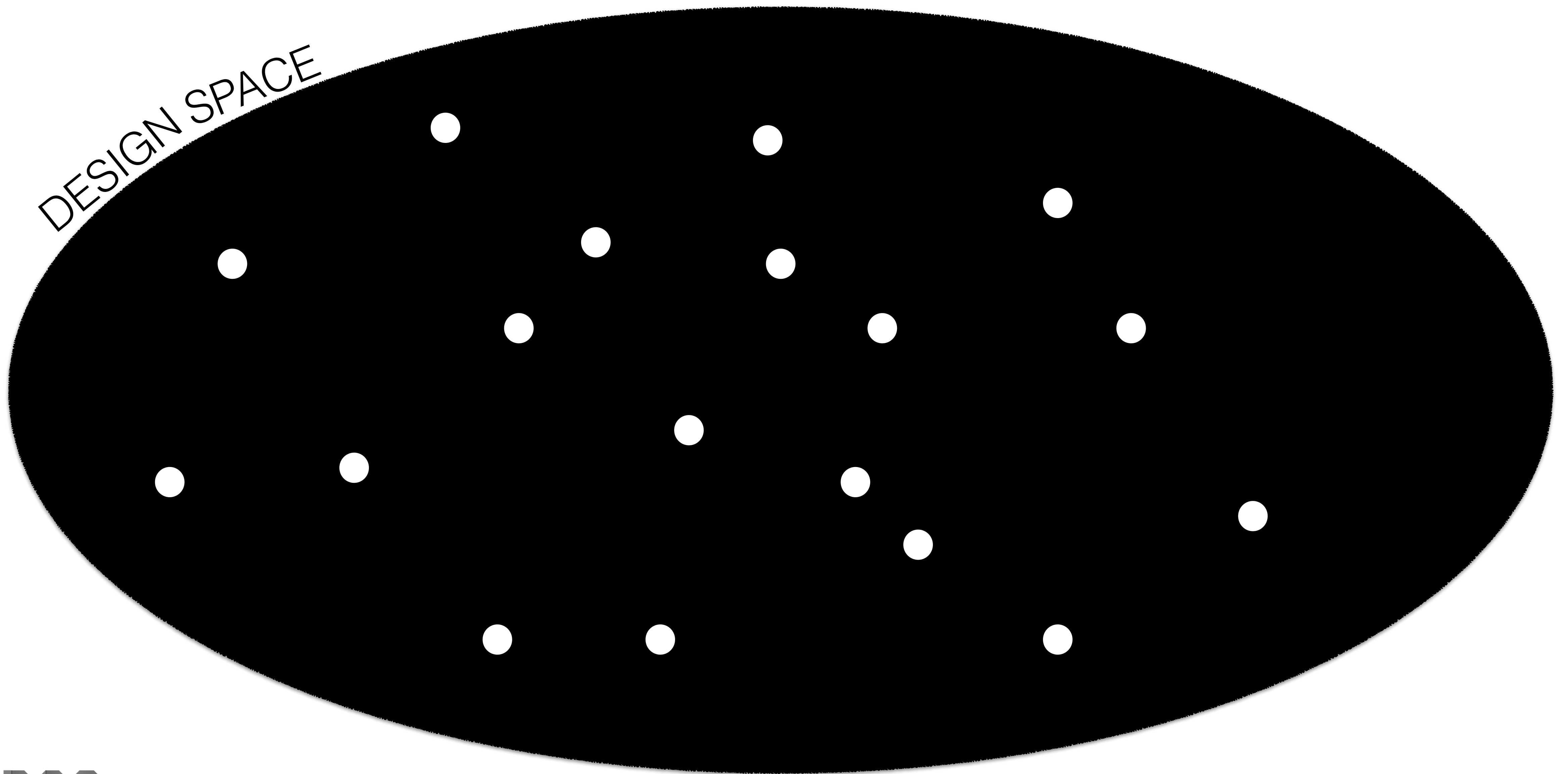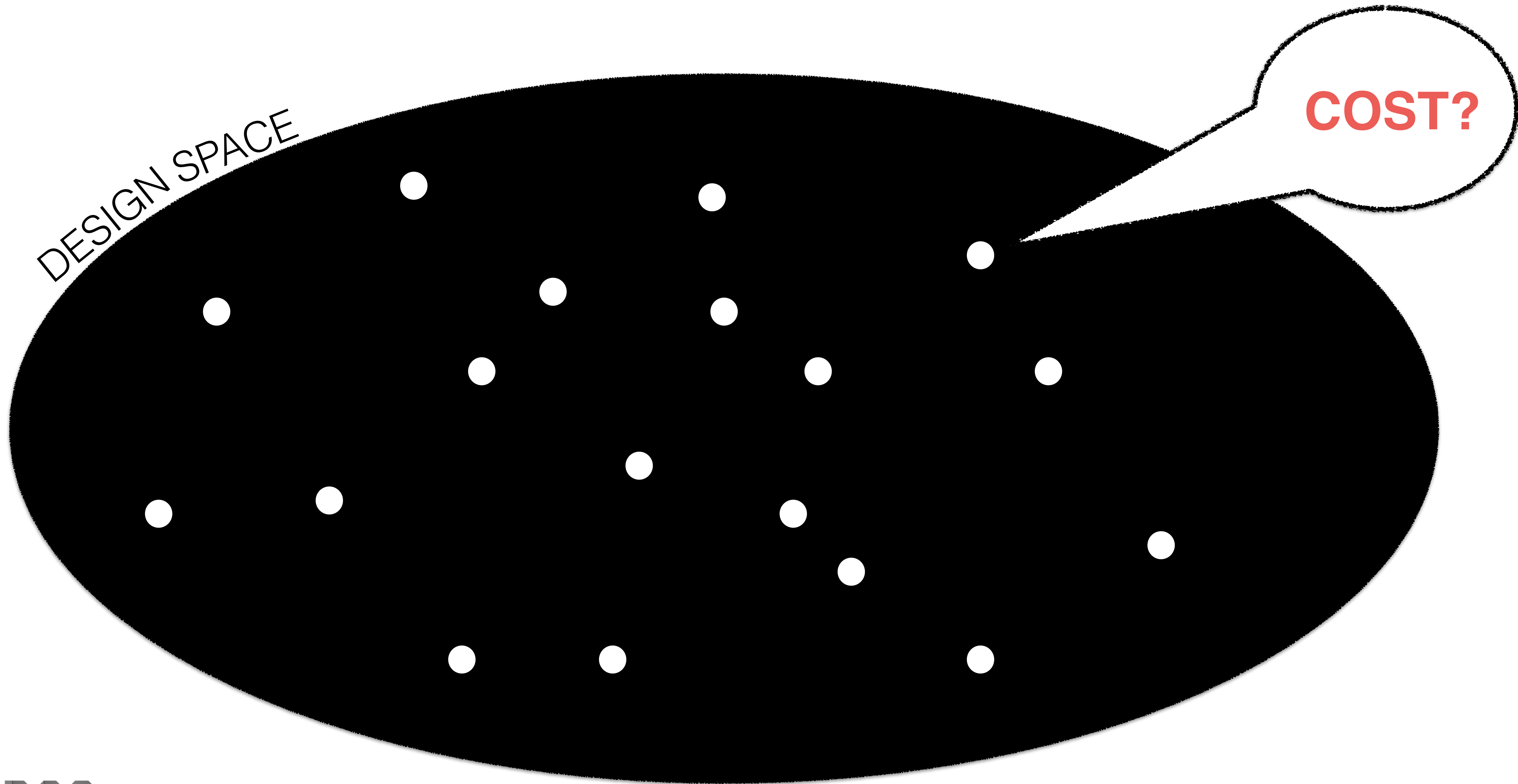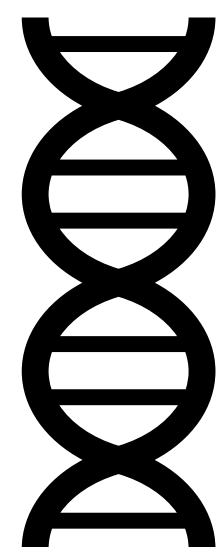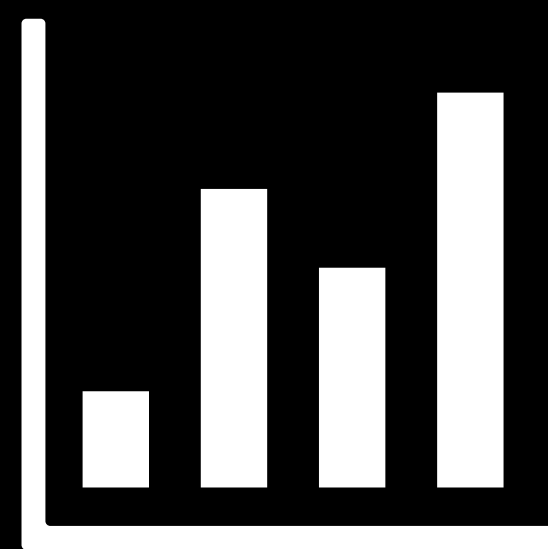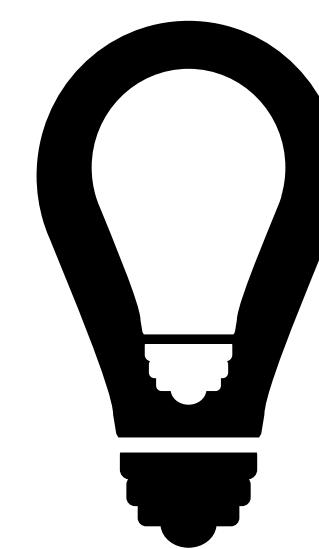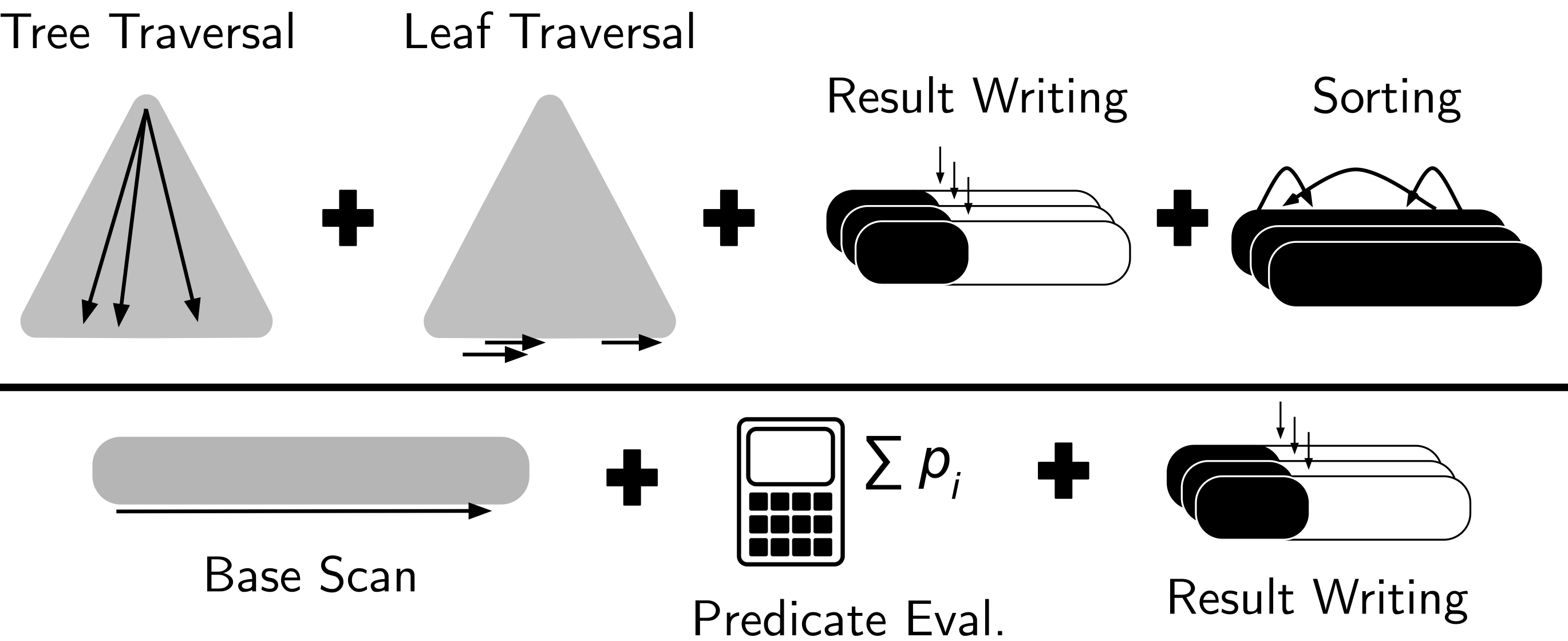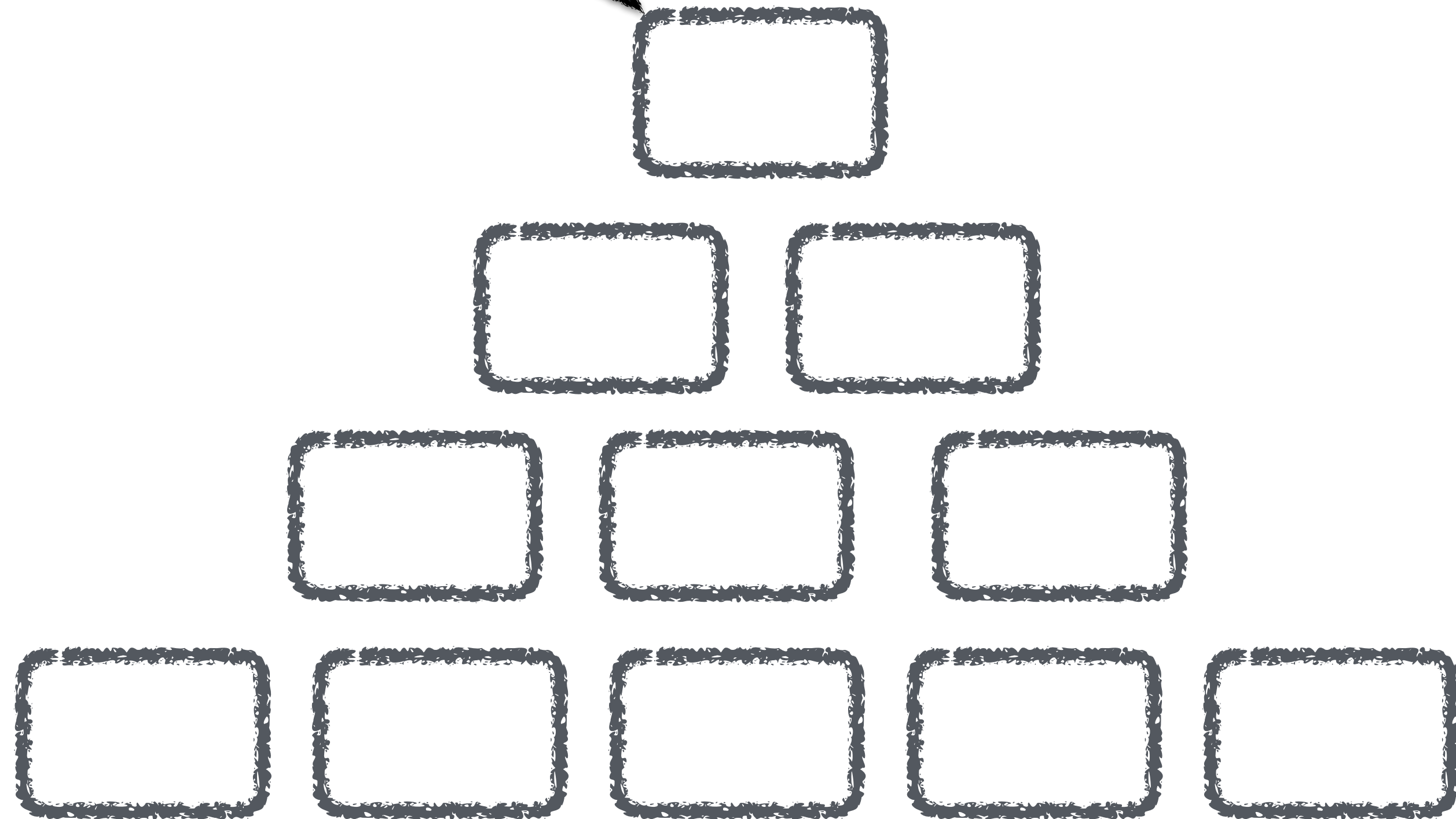\cdot p}{2}\right)}{max\left(ts, 2 \cdot f_p \cdot p \cdot q \cdot BW_S\right) + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}}$$

$$+ \frac{S_{tot}\left(\frac{BW_S \cdot C_M}{b} + (aw + ow) \cdot \frac{BW_S}{BW_I} + rw \cdot \frac{BW_S}{BW_R}\right)}{max\left(ts, 2 \cdot f_p \cdot p \cdot q \cdot BW_S\right) + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}}$$

$$+ \frac{S_{tot} \cdot log_2\left(S_{tot} \cdot N\right) \cdot BW_S \cdot C_A}{max\left(ts, 2 \cdot f_p \cdot p \cdot q \cdot BW_S\right) + S_{tot} \cdot rw \cdot \frac{BW_S}{BW_R}}$$

## Access path selection @SIGMOD2017



Tree Traversal + Leaf Traversal + Result Writing + Sorting

Base Scan + Predicate Eval. $\sum p_i$ + Result Writing

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

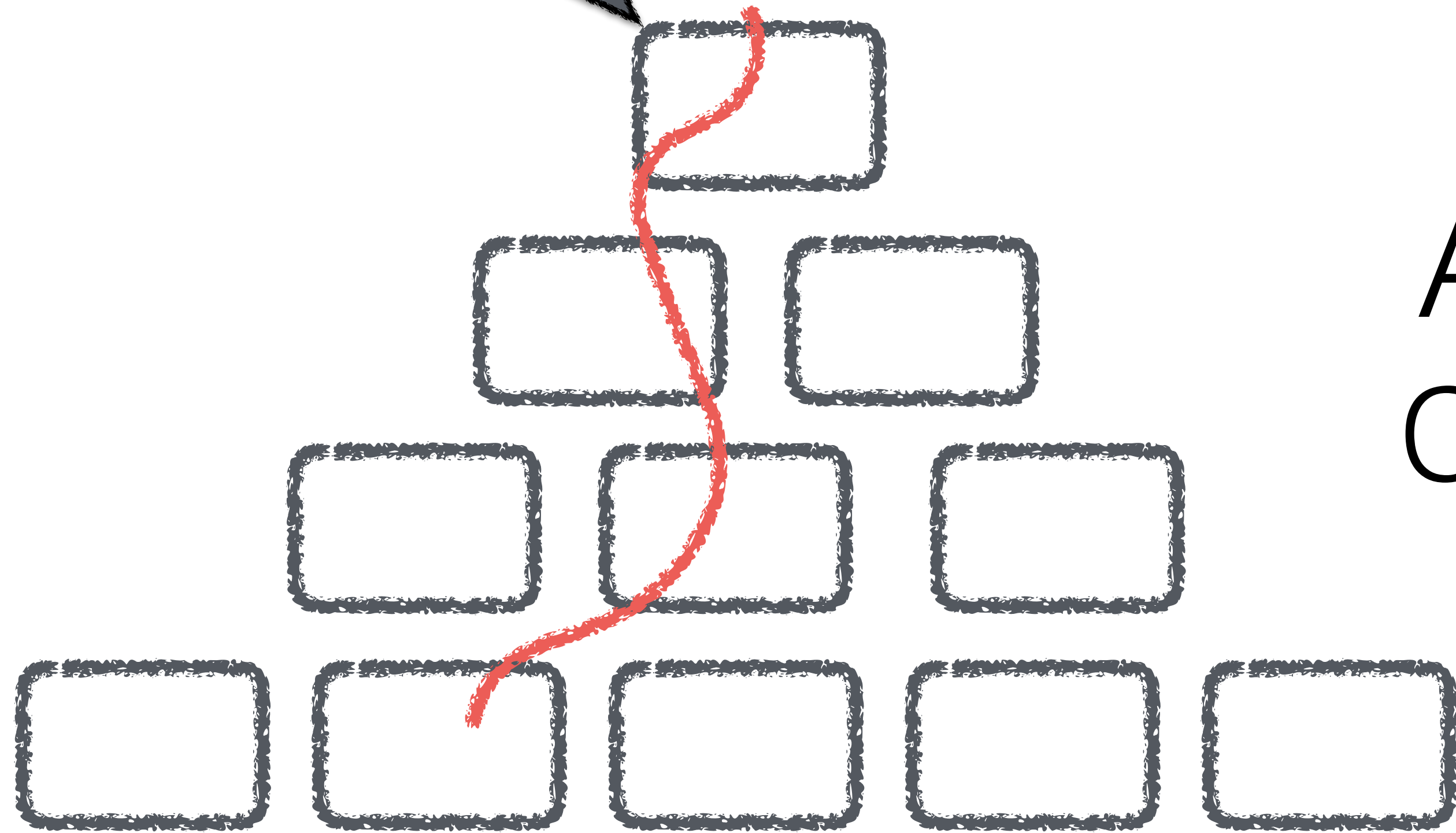DESIGN SPACE OF POSSIBLE STORAGE LAYOUTS
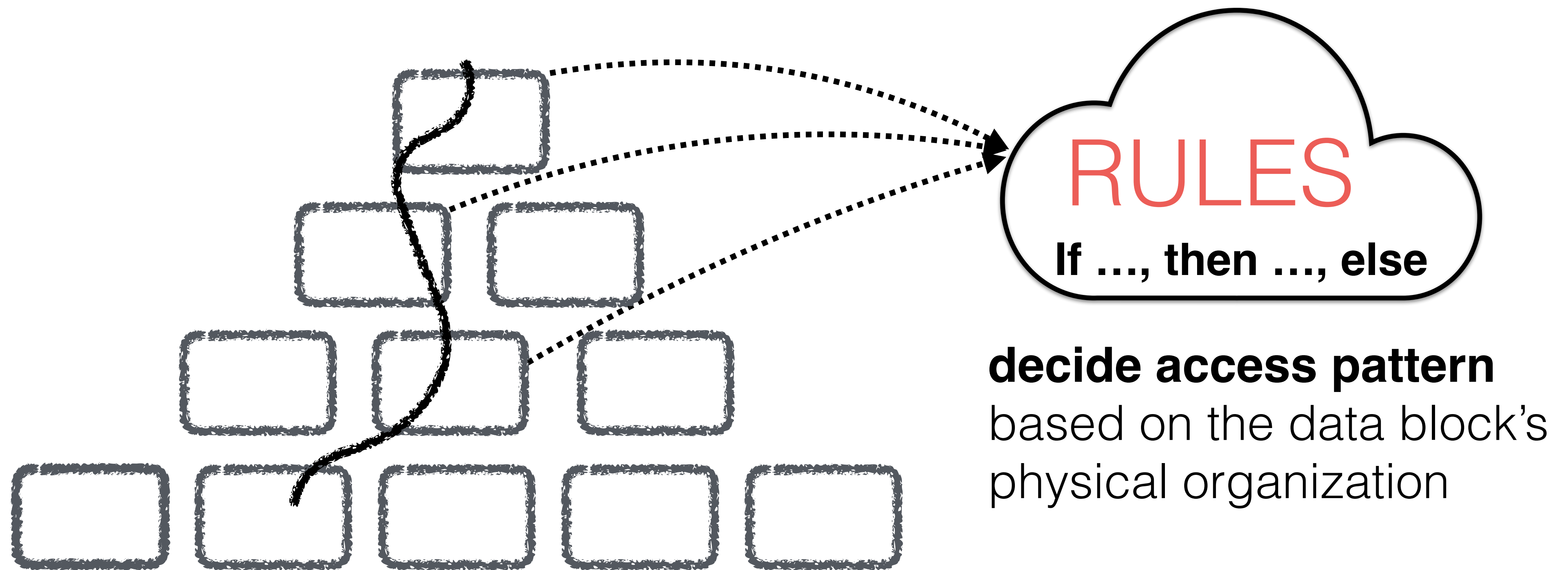
DESIGN SPACE
OF POSSIBLE
STORAGE LAYOUTS

operation

ALGORITHM &
COST SYNTHESIS

# RULES

**If …, then …, else**

**decide access pattern**
based on the data block's
physical organization

DASlab
@ Harvard SEAS

sorted keys
columnar layout

sorted keys
columnar layout

RULES →

sorted
search

# SYNTHESIS FROM LEARNED MODELS
## coding, modeling, generalized models, and a touch of ML

**1.** **MINIMAL CODE**

e.g., binary search



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

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

DASlab
@ Harvard SEAS

# SYNTHESIS FROM LEARNED MODELS
## coding, modeling, generalized models, and a touch of ML

**1. MINIMAL CODE**

**2. BENCHMARK**

e.g., binary search

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

C++

Run

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



DASlab
@ Harvard SEAS

# SYNTHESIS FROM LEARNED MODELS
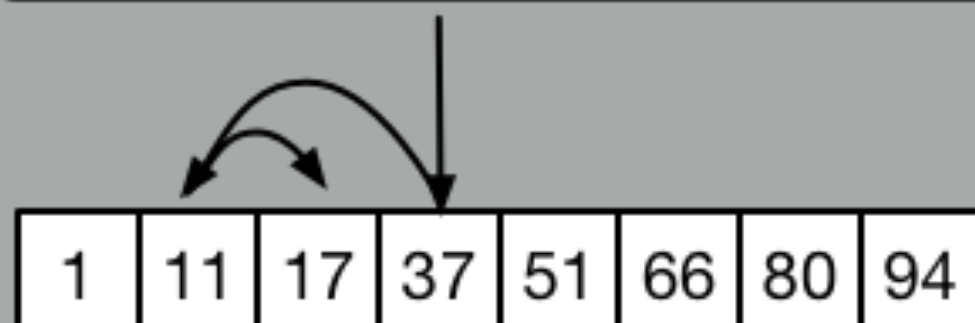## coding, modeling, generalized models, and a touch of ML

**1.** MINIMAL CODE

e.g., binary search
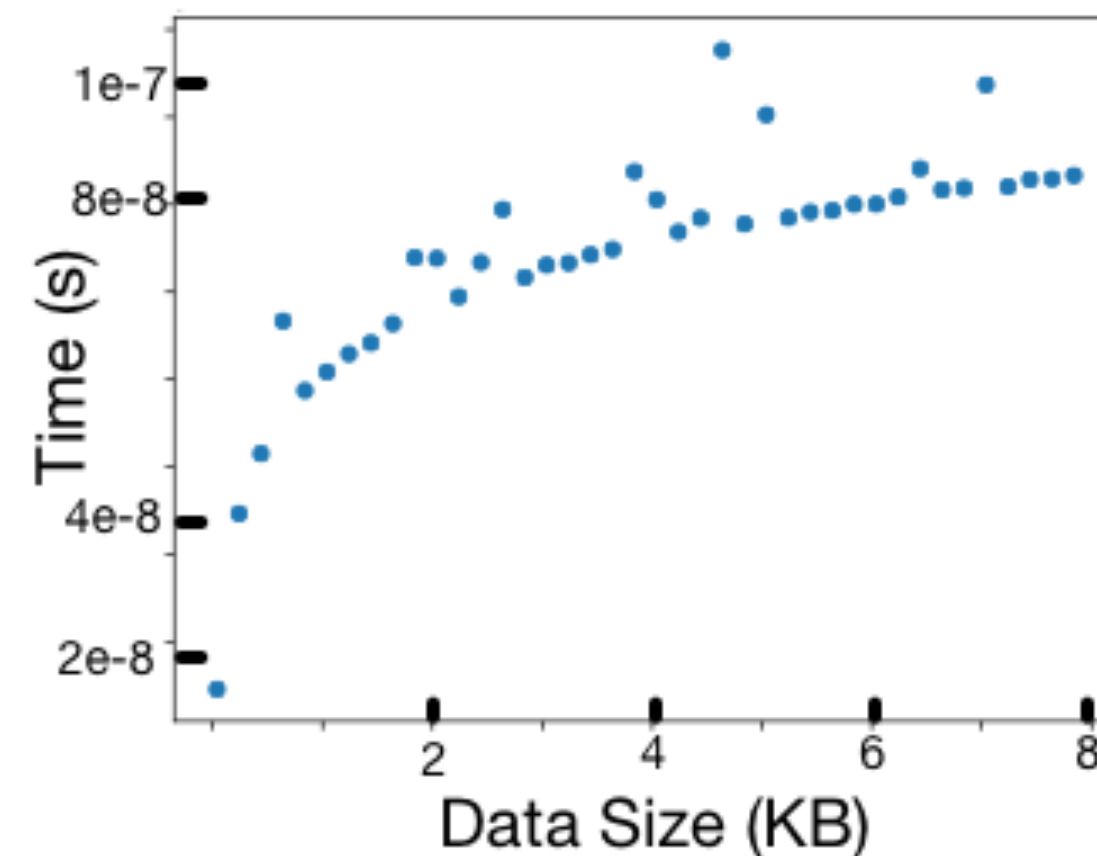
```
if (data[middle] < search_val) {
    low = middle + 1;
} else {
    high = middle;
}
middle = (low + high)/2;
```
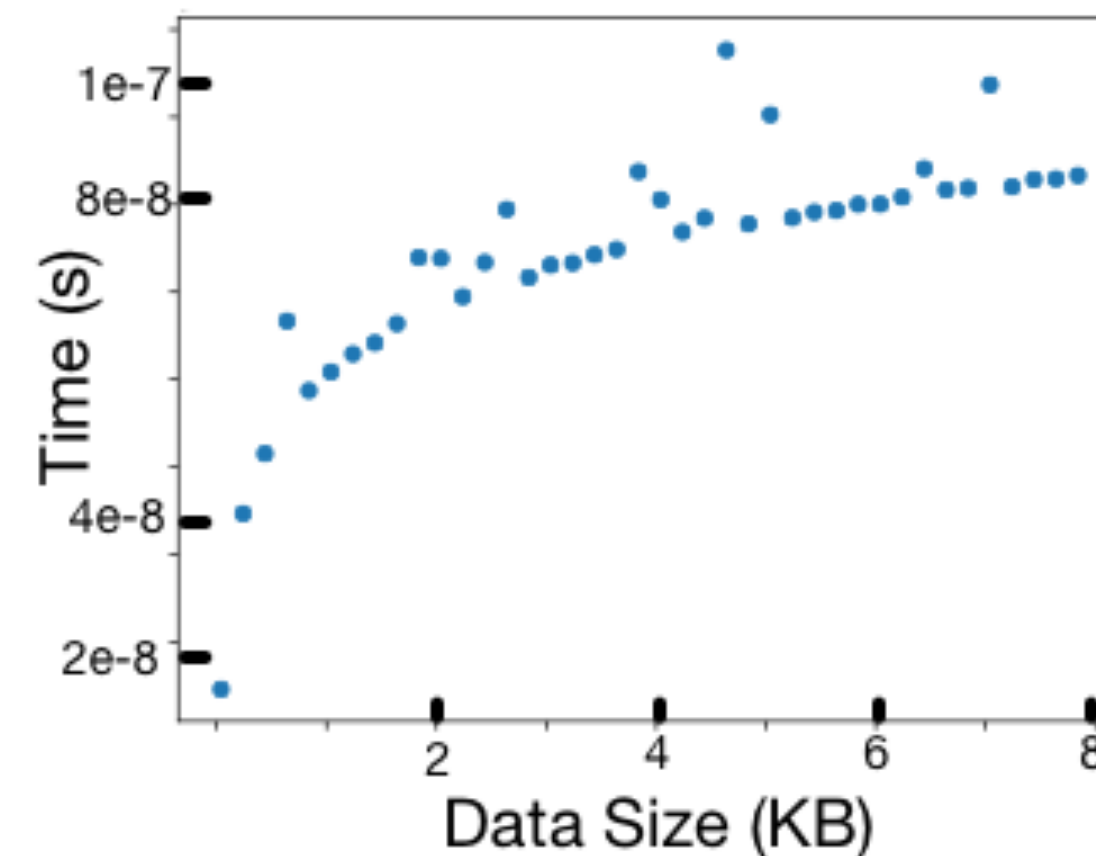
C++

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

Run

**2.** BENCHMARK



Time (s) vs Data Size (KB)

$f(x)$

Train

**3.** FIT MODEL



Log-Linear Model

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

DASlab
@ Harvard SEAS

Data Calculator

cost

Binary
Search

**CPP**

Scan

**CPP**

Probe

**CPP**

Hardware

First Principles (access)

Data
Calculator

. . .

cost

# TRAINING

Binary
Search

Scan

Probe

**CPP**

**CPP**

**CPP**

Hardware

First Principles (access)

Data
Calculator

Learned Models

cost

DASlab
@ Harvard SEAS

QUEYRYING

Binary Search
Scan
Probe
CPP CPP CPP

Hardware

First Principles (access)

Learned Models

Data Calculator

cost

key retention
val retention
key order
...

First Principles (layout)

DASlab
@ Harvard SEAS

QUEYRYING

Binary
Search

Scan

Probe

CPP

CPP

CPP

Hardware   Workload Specification

First Principles (access)

...

Data
Calculator

Learned Models

cost

key
retention   val
retention   key
order   ...

First Principles (layout)

DASlab
@ Harvard SEAS

# QUEYRYING

Binary Search

Scan

Probe

CPP

CPP

CPP

...

Hardware

Workload Specification

**First Principles (access)**

Data Calculator

cost

**Learned Models**

key retention

val retention

key order

...

**First Principles (layout)**

FOR EACH OPERATION

FOR EACH NODE

1. Decide access **strategy** (L1) based on node design

2. Decide exact access strategy **implementation** (L2) based on available models
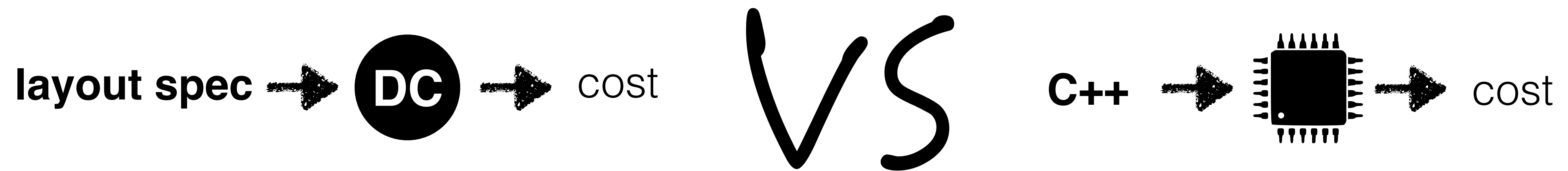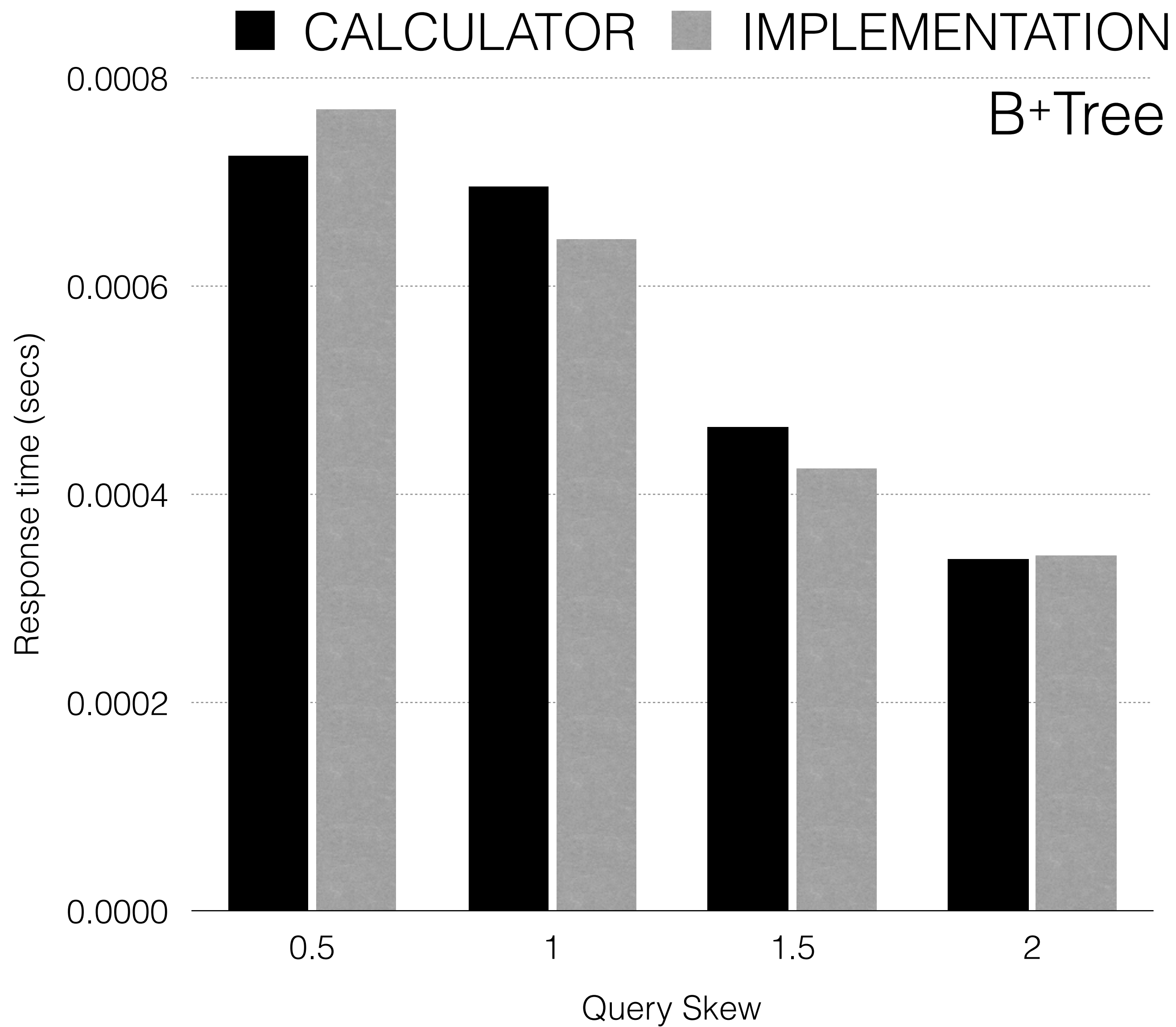
3. **Get cost** for chosen model

DASlab
@ Harvard SEAS

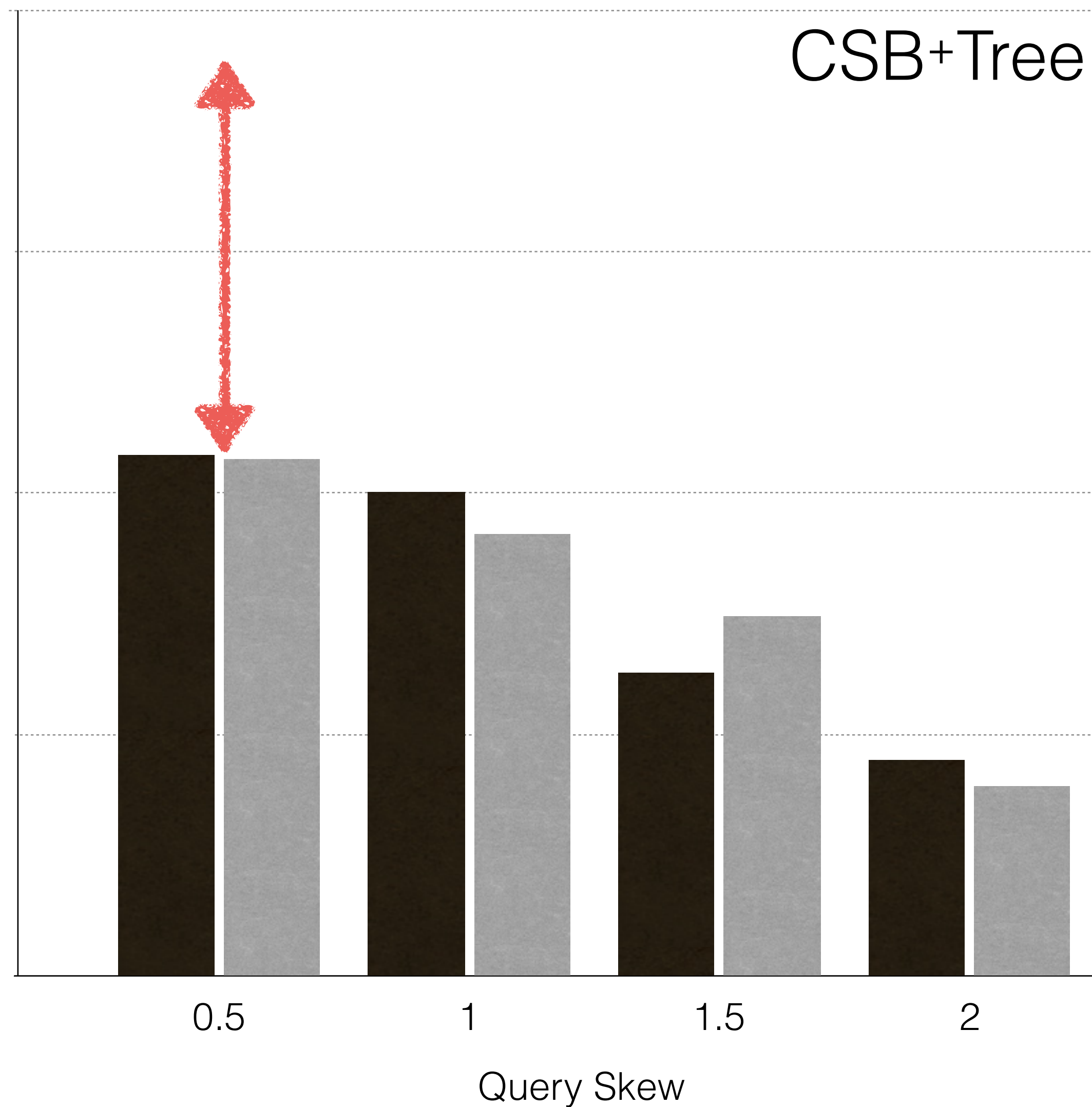# CAN WE COMPUTE PERFORMANCE ACCURATELY?
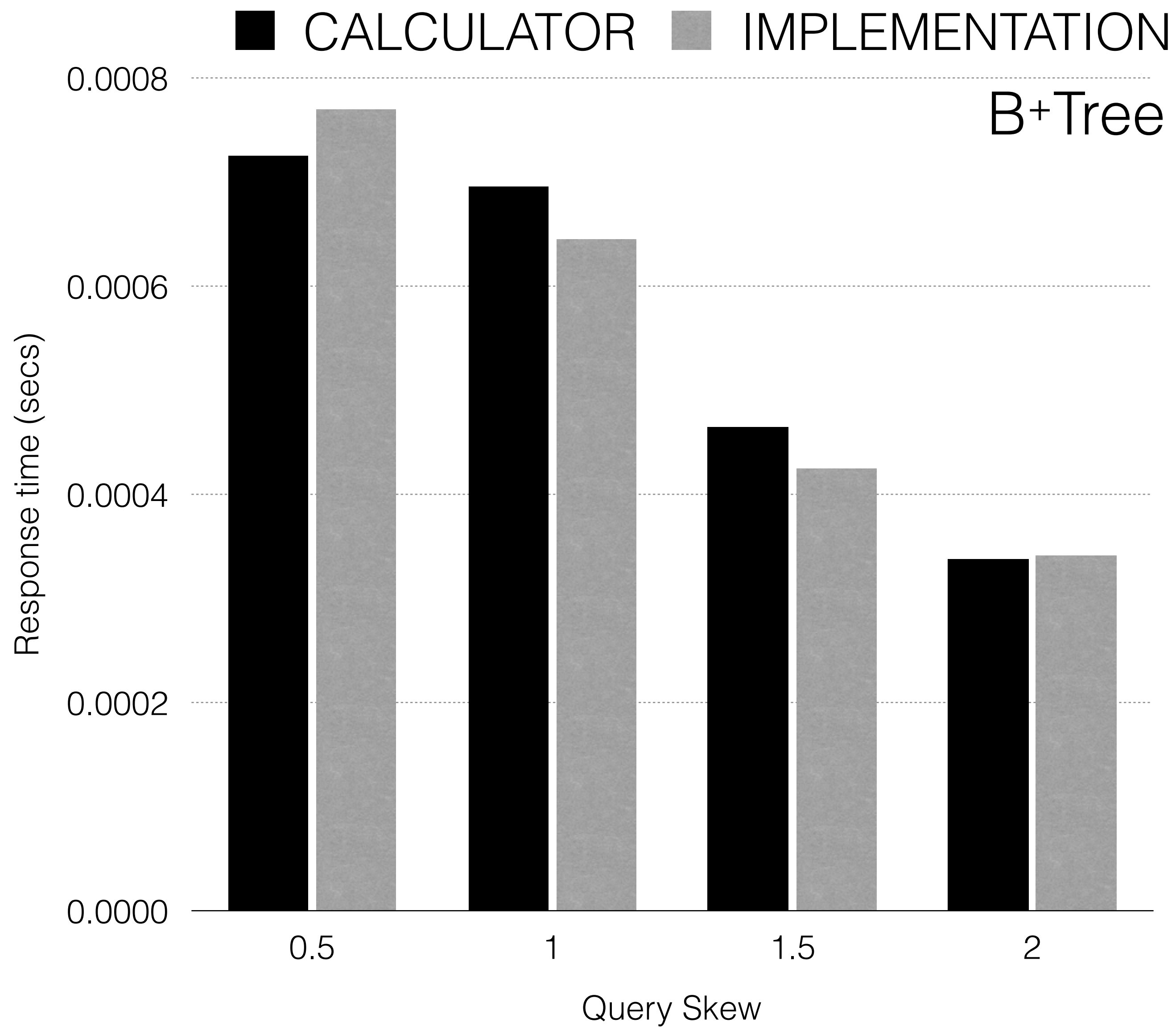
**layout spec** ➜ **DC** ➜ cost
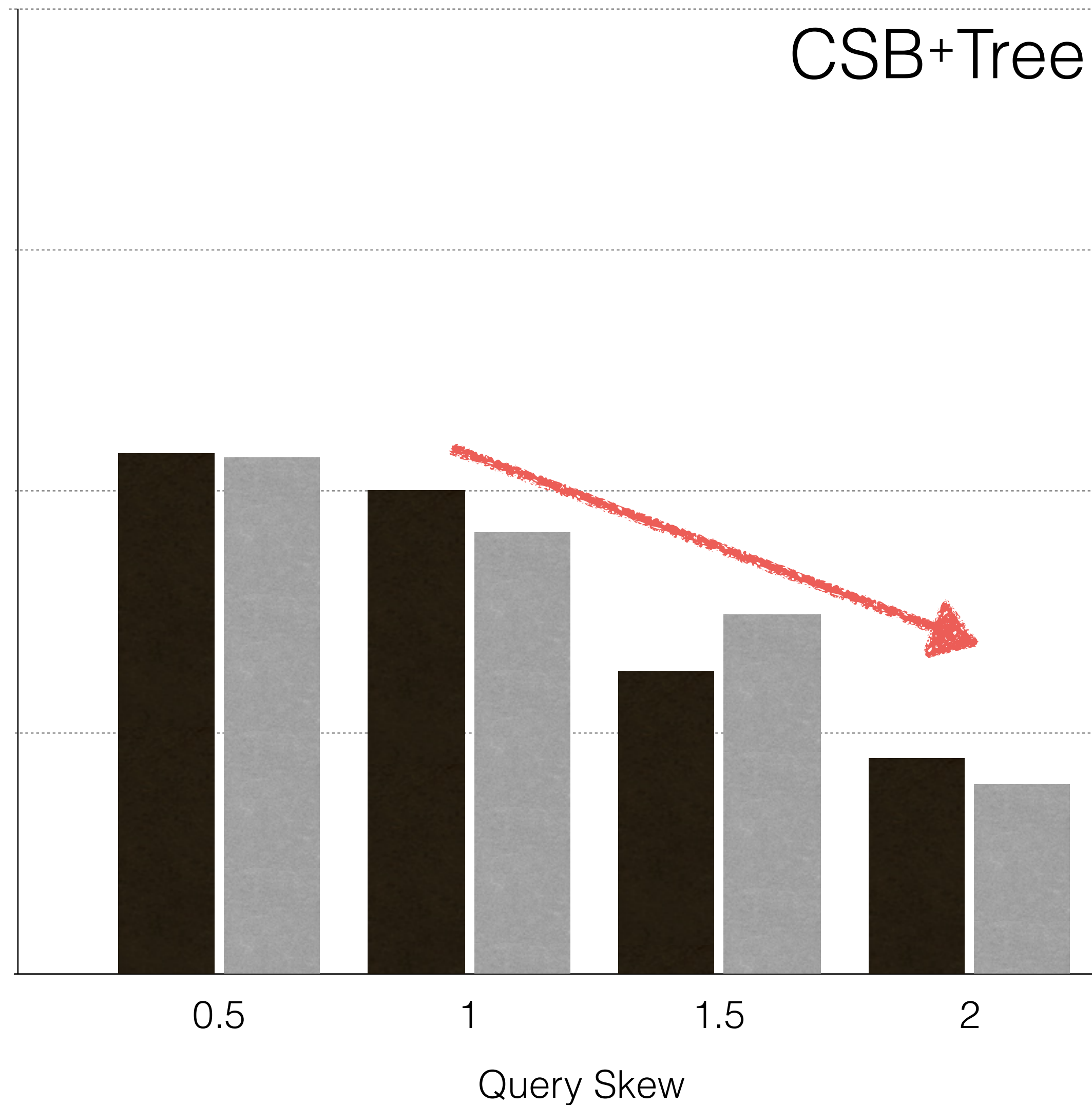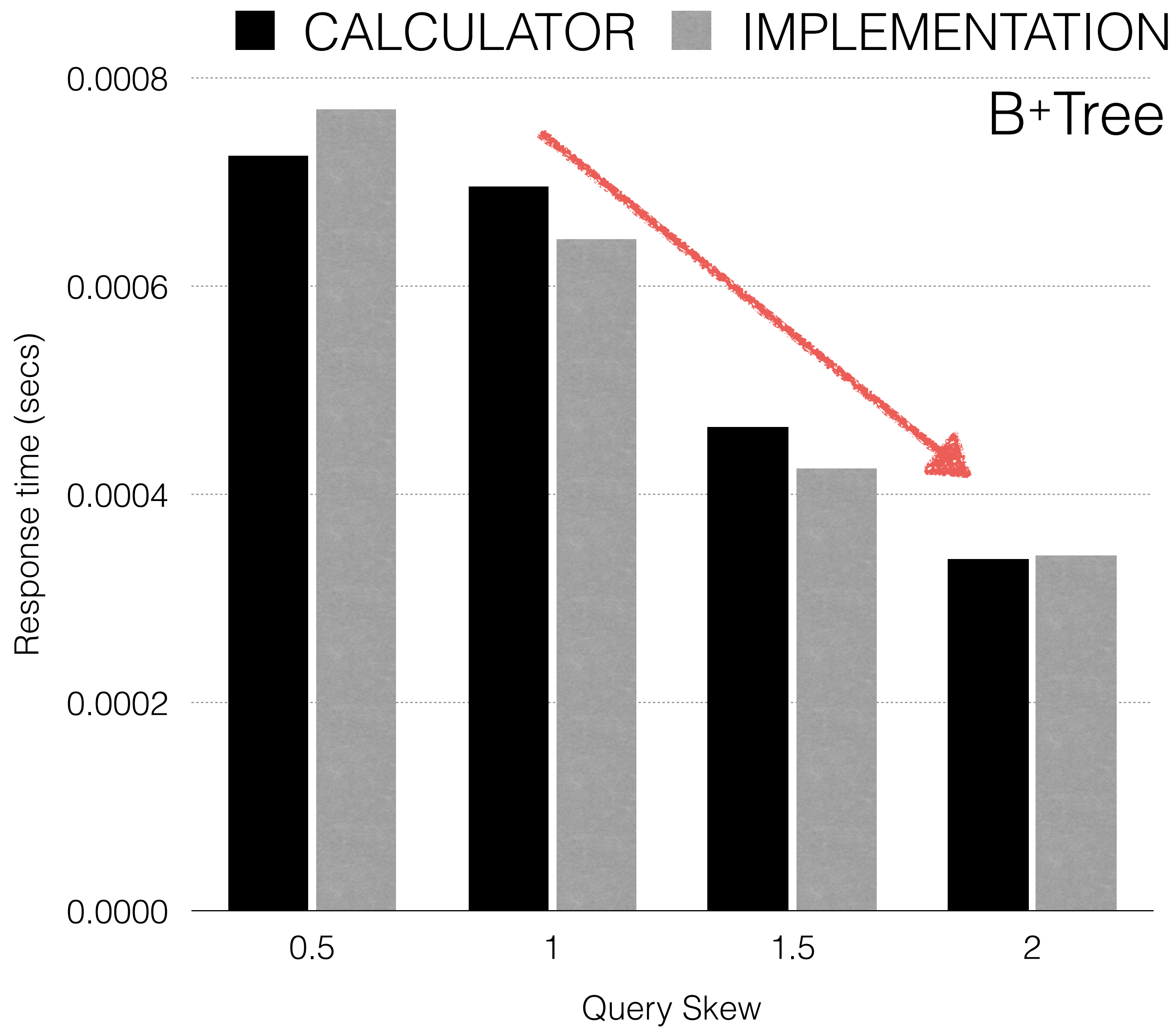
# VS

**C++** ➜ [chip] ➜ cost

(same workload, hardware, data)

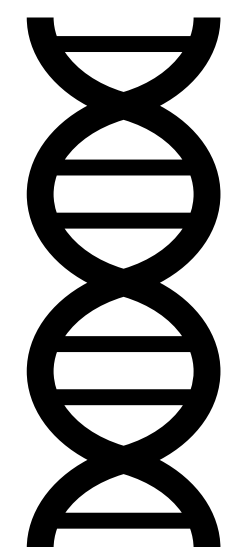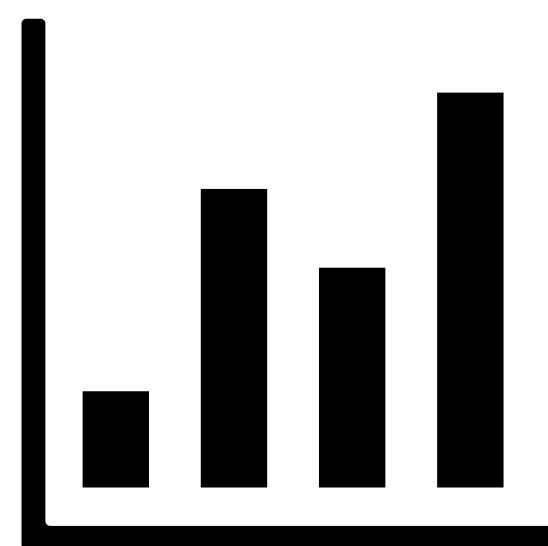{10M (uniform) k-v pairs, 100 point queries (skewed)}

{10M (uniform) k-v pairs, 100 point queries (skewed)}

DESIGN SPACE | COST SYNTHESIS

WHO AND HOW TO USE

WHO AND
HOW TO USE

What-if we **add bloom filters** in the hash-table buckets?

**~20 SECONDS**

(workload:10 Million entries, 100 queries)

*What-if the workload changes to **90% writes**?*

**~20 SECONDS**

(workload:10 Million entries, 100 queries)

What-if we **buy faster CPU X**?

# ~20 SECONDS
(workload:10 Million entries, 100 queries)

auto.design

dynamic programming
genetic algorithms
reinforcement learning

DASlab
@ Harvard SEAS

B+TREE
ELEMENT

Synthesis cost (min.)

30

20

10

0

$10^5$          $10^6$          $10^7$

# of inserts

Hybrid B+Tree / Hash Table / Array

Point get
intensive

Range
intensive

Only writes

HASH
PARTITIONING

B+TREE
ELEMENT

DATAPAGE
(system page size)

DATAPAGE
(large chunks)

DATAPAGE
(system page size)

DASlab
@ Harvard SEAS

lookup cost

update cost

*log*

○○○○○

lookup cost

update cost

DASlab
@ Harvard SEAS

*log*

lookup cost

*LSM-tree\**

update cost

DASlab
@ Harvard SEAS

**log**

lookup cost

**LSM-tree\***

**B-tree\***

update cost

DASlab
@ Harvard SEAS

**log**

lookup cost

**LSM-tree***

**B-tree***

**sorted array**

update cost

log

one model

LSM-tree*

B-tree*

sorted array

lookup cost

update cost

DASlab
@ Harvard SEAS

log

one model

LSM-tree*

B-tree*

sorted array

lookup cost

update cost

@ Harvard SEAS

**few continuous**
parameters
&
**closed form**
formulas for metrics

**few continuous**
parameters
&
**closed form**
formulas for metrics

near **instant** design
space navigation

*log*

one model

lookup cost

*LSM-tree\**

tiering

size
ratio

leveling

*B-tree\**

*sorted array*

update cost

DASlab
@ Harvard SEAS

log

lookup cost

one model

LSM-tree*

tiering

size
ratio

leveling

B-tree*

sorted array

update cost

DASlab
@ Harvard SEAS

**few continuous**
parameters
&
**closed form**
formulas for metrics

near **instant** design
space navigation

Crimson DB

a self-designing key-value store

S. BING YAO
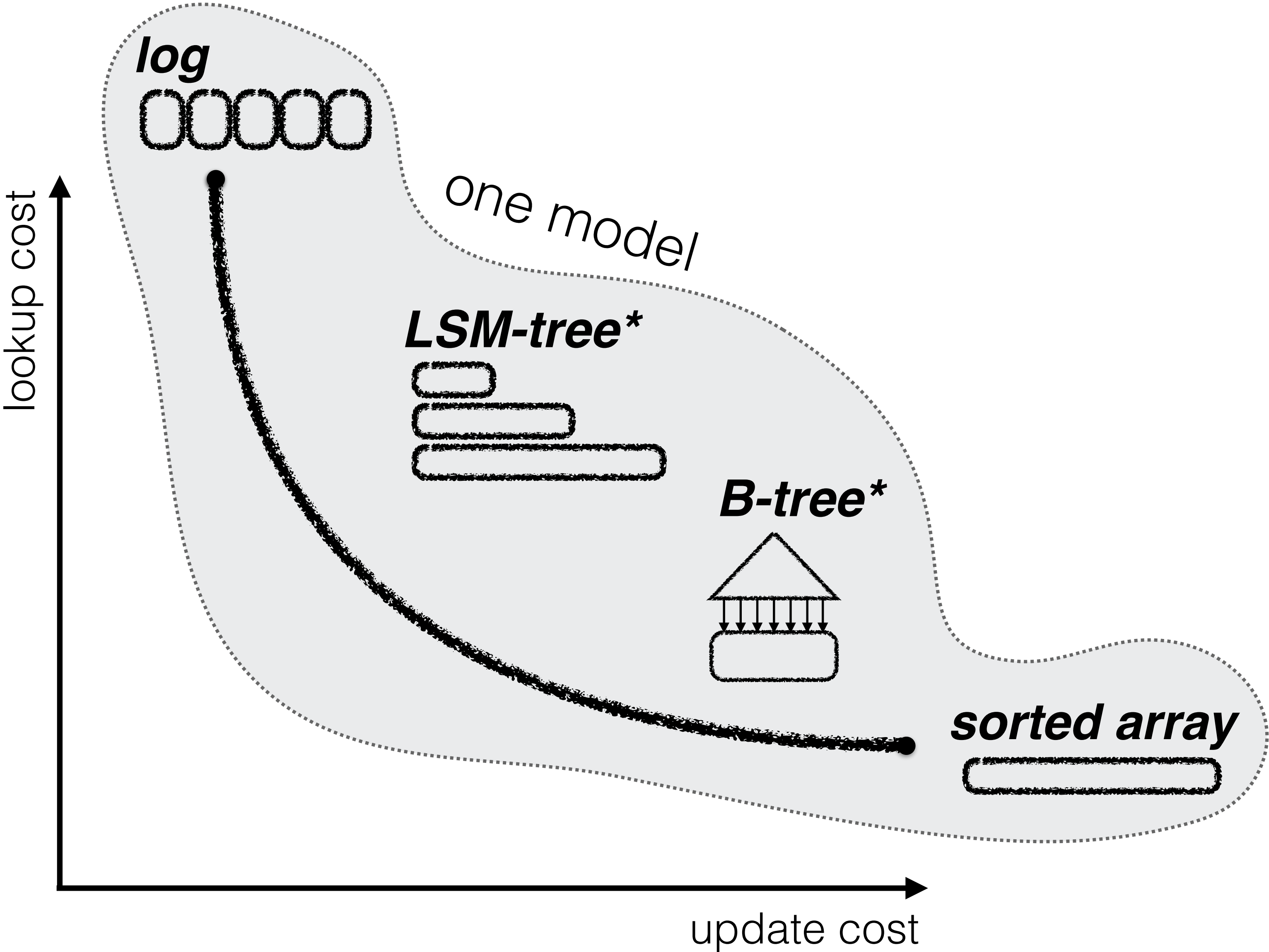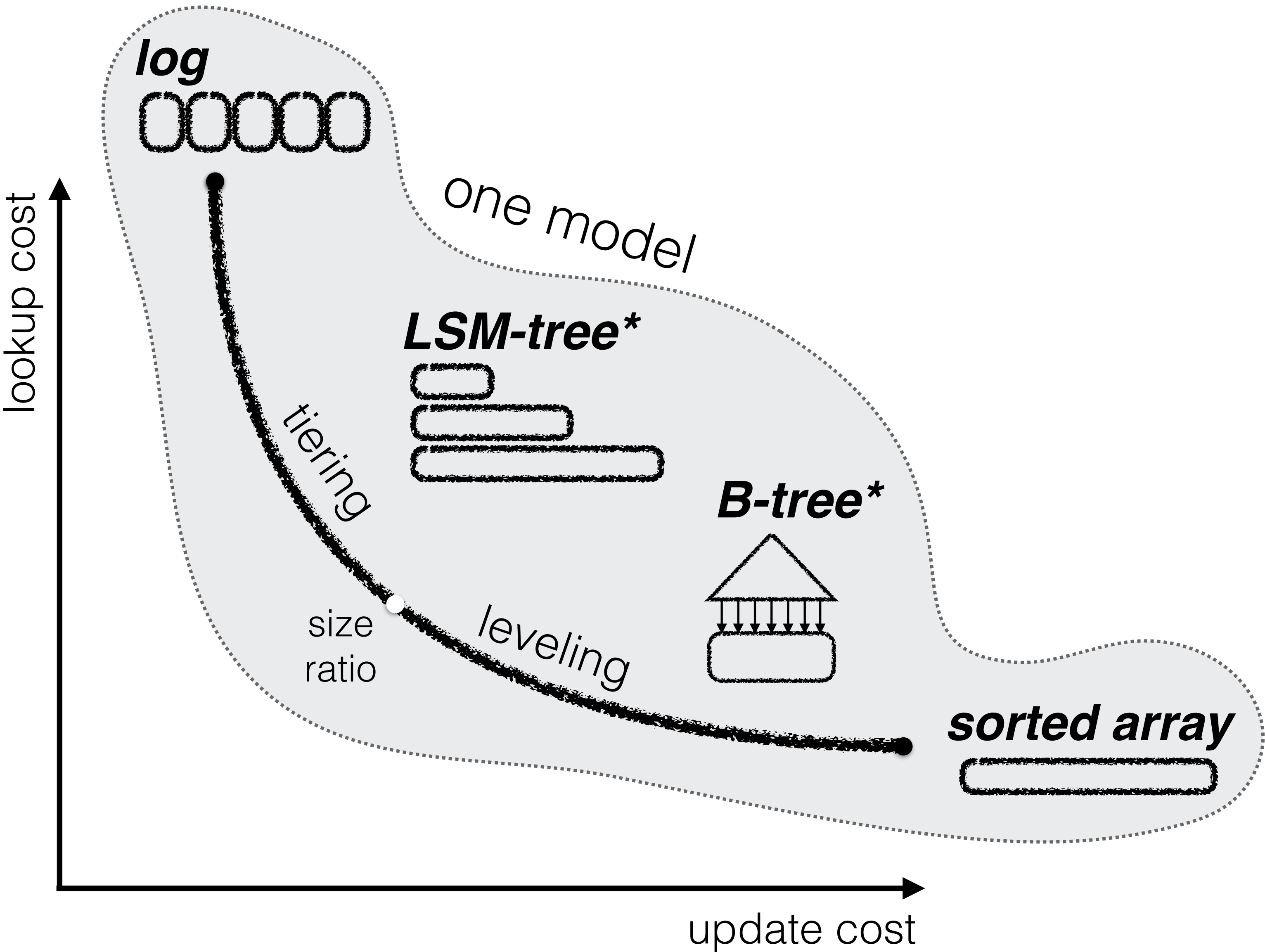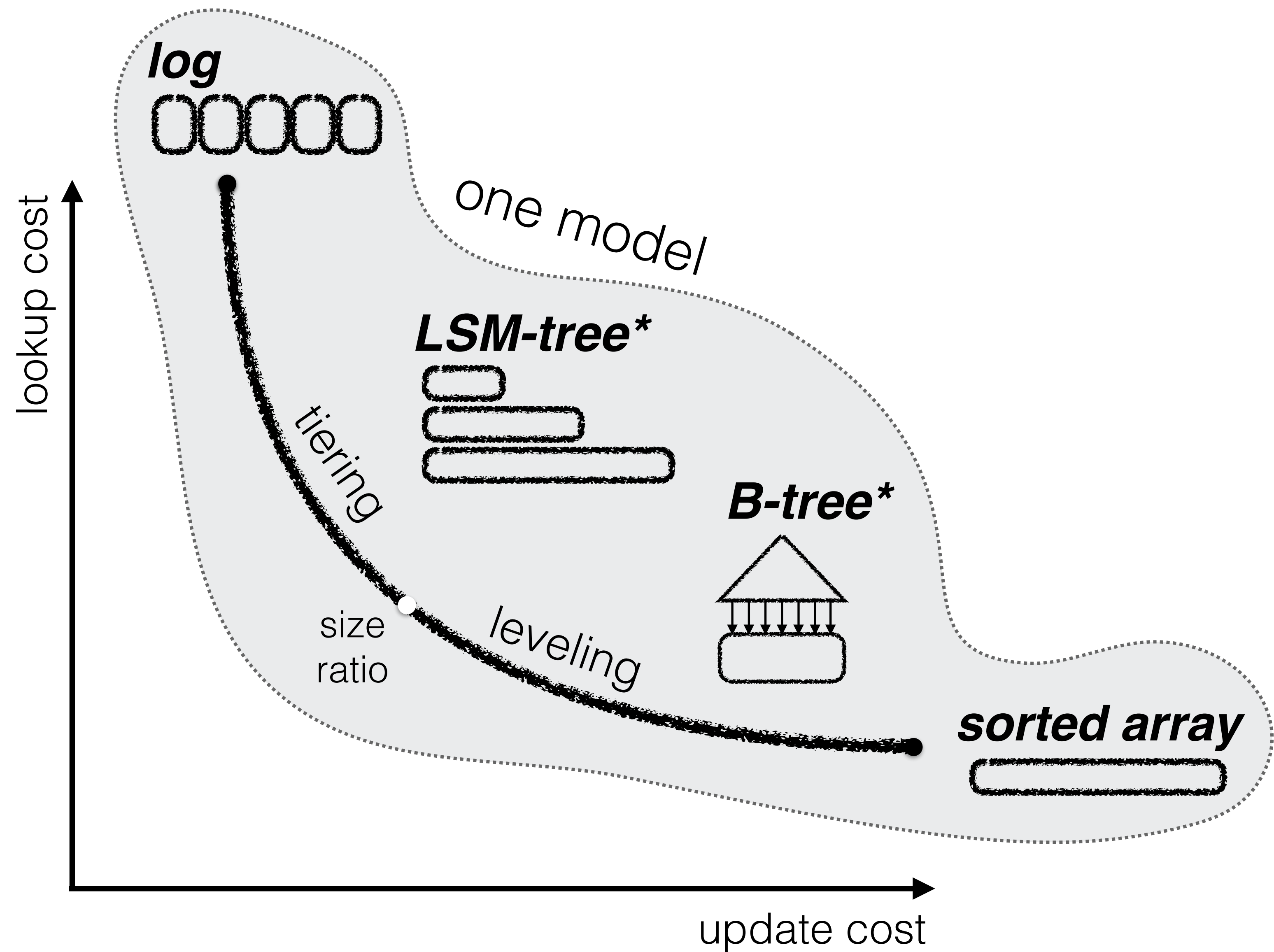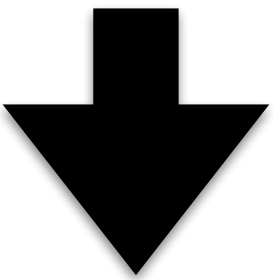models/advisors

DON BATORY
modular synthesis

JOE HELLERSTEIN
extensible indexing

STEFAN MANEGOLD
model synthesis

## calculator infrastructure

DESIGN SPACE (updates, concurrency)

COST SYNTHESIS (accuracy, scalability)

EASY EXTENSIBILITY (plug & play rules)

## study structure & gaps

DESIGN CONTINUUMS (optimizations)

GO AFTER GAPS (new design opport.)

DESIGN GUIDE (static design rule sys)

## building more on top of

AUTO-SEARCH (ML/algo hybrids, hints)

SELF-DESIGNING (log to sorted arrays)

DSL & COMPILERS (productivity & perf)

# DASlab
## @ Harvard SEAS

daslab.seas.harvard.edu

# THANKS!

**Manos Athanassoulis** — *Postdoctoral Researcher*

**Niv Dayan** — *Postdoctoral Researcher*

**Kostas Zoumpatianos** — *Postdoctoral Researcher*

**Michael Kester** — *Ph.D. Researcher*

**Lukas Maas** — *Ph.D. Researcher*

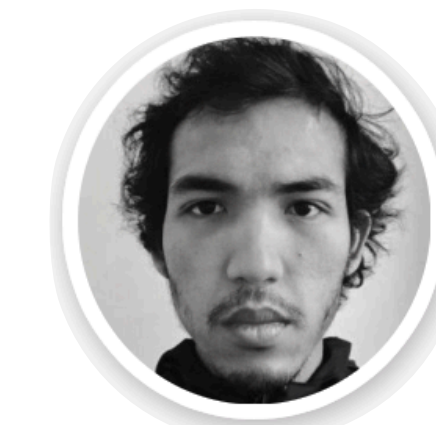**Abdul Wasay** — *Ph.D. Researcher*

**Brian Hentschel** — *Ph.D. Researcher*

**Wilson Qin** — *Ph.D. Researcher*

**Angelo Kastroulis** — *Graduate Researcher*

**Franco Solleza** — *Graduate Researcher*

**Mali Akmanalp** — *Graduate Researcher*

**Rachna Sha** — *Graduate Researcher*

**Dhruv Gupta** — *Undergraduate Researcher*

**Demi Guo** — *Undergraduate Researcher*

**Yiyou Sun** — *Undergraduate Research Intern*

**Mo Sun** — *Undergraduate Research Intern*

**Yuze Liao** — *Undergraduate Research Intern*

## *calculator infrastructure*

DESIGN SPACE (updates, concurrency)

COST SYNTHESIS (accuracy, scalability)

EASY EXTENSIBILITY (plug & play rules)

## *study structure & gaps*

DESIGN CONTINUUMS (optimizations)

GO AFTER GAPS (new design opport.)

DESIGN GUIDE (static design rule sys)

## *building more on top of*

AUTO-SEARCH (ML/algo hybrids, hints)

SELF-DESIGNING (log to sorted arrays)

DSL & COMPILERS (productivity & perf)