

class 2

Data Systems 101

Prof. Manos Athanassoulis

<https://midas.bu.edu/classes/CS591A1/>



some reminders

no smartphones



no laptop



class summary

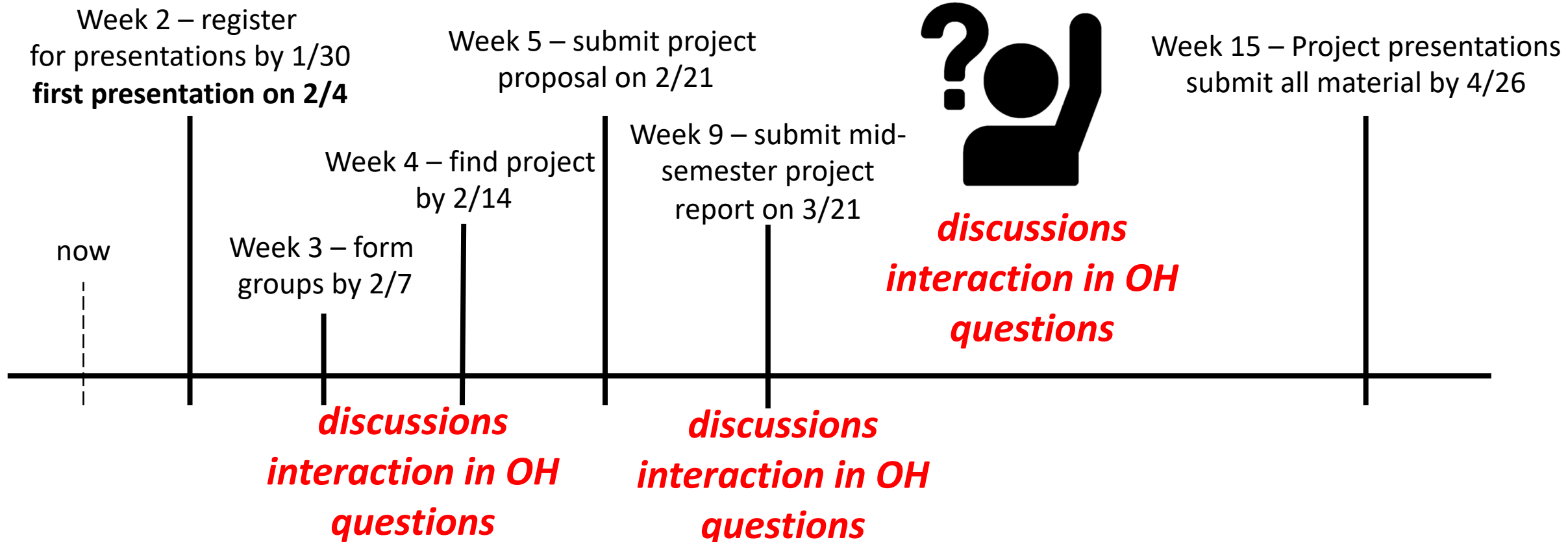
2 classes per week / OH 5 days per week

each student

1 presentation/discussion lead + 2 reviews/questions per week

systems or research project + proposal + mid-semester report

class timeline



Piazza



all discussions & announcements

<http://piazza.com/bu/spring2020/cs591a1/>

also available on class website

10 already registered!

register so we can reach you easily

size (volume)

rate (velocity)


sources (variety)

big data

(it's not only about size)

The 3 V's

+ our ability to collect ***machine-generated*** data

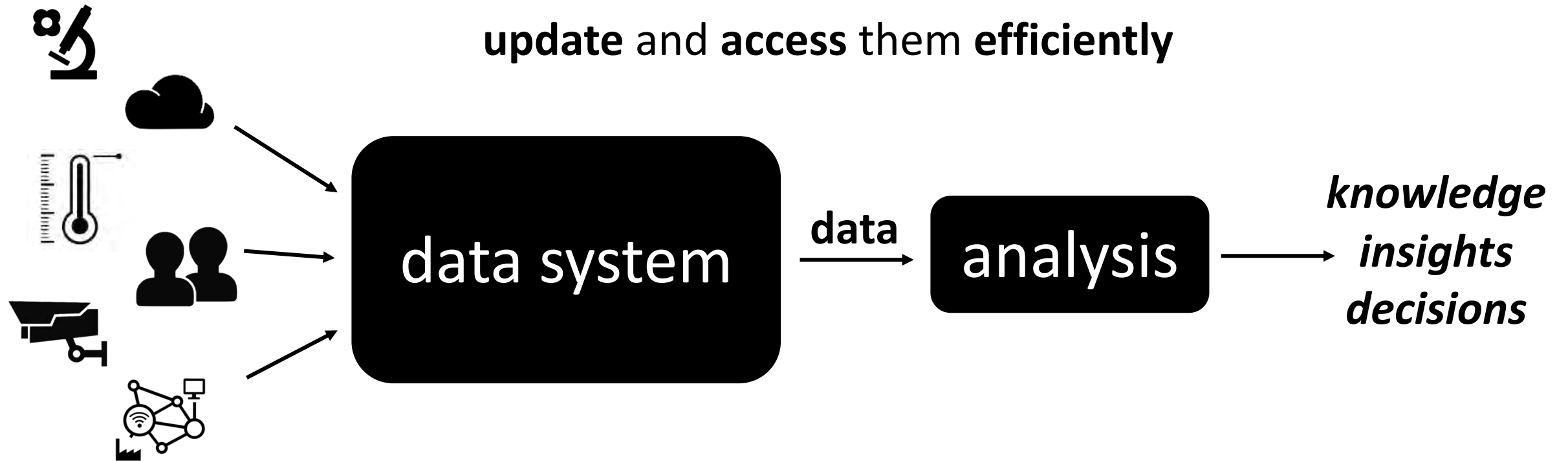
 scientific experiments

 sensors

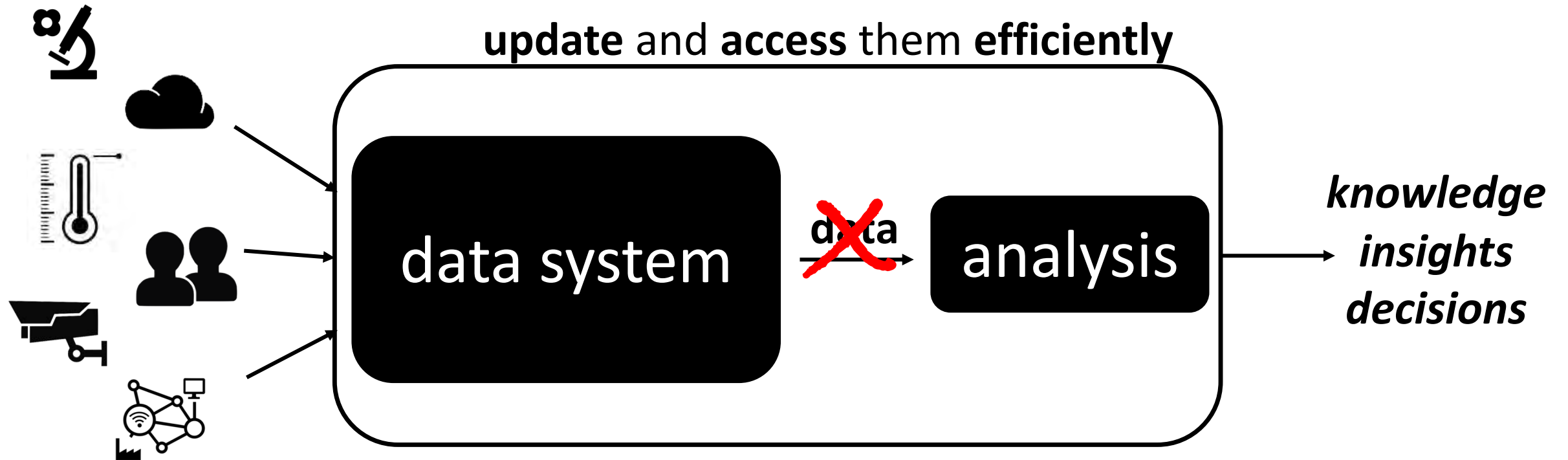
social 

Internet-of-Things 

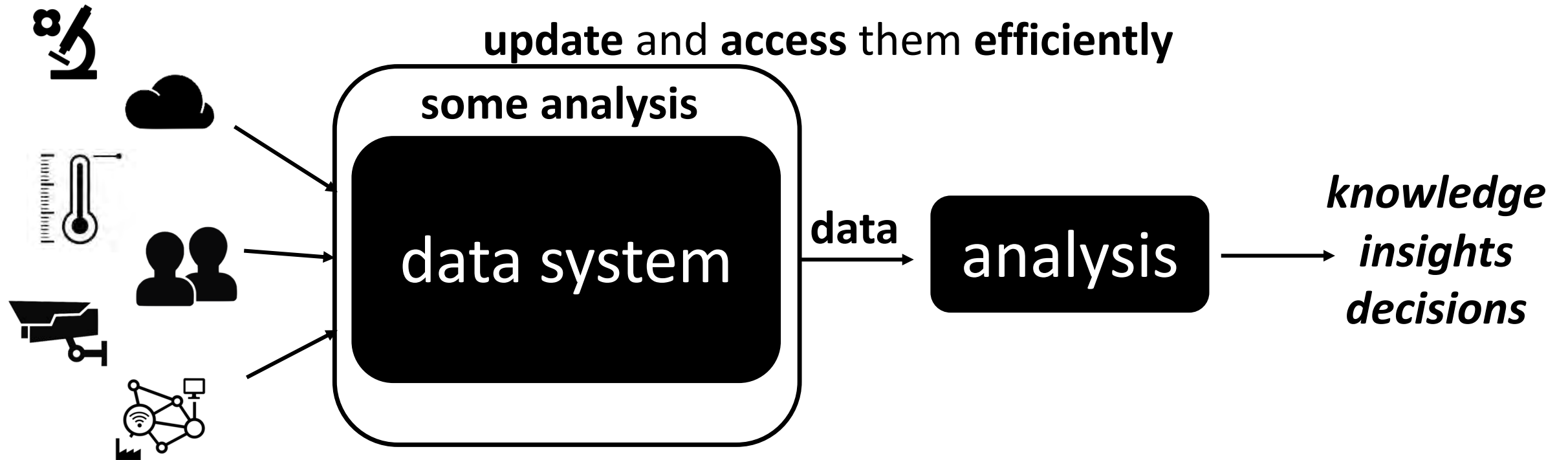
a **data system** is a large software system that **stores data**, and provides the **interface** to **update** and **access** them **efficiently**



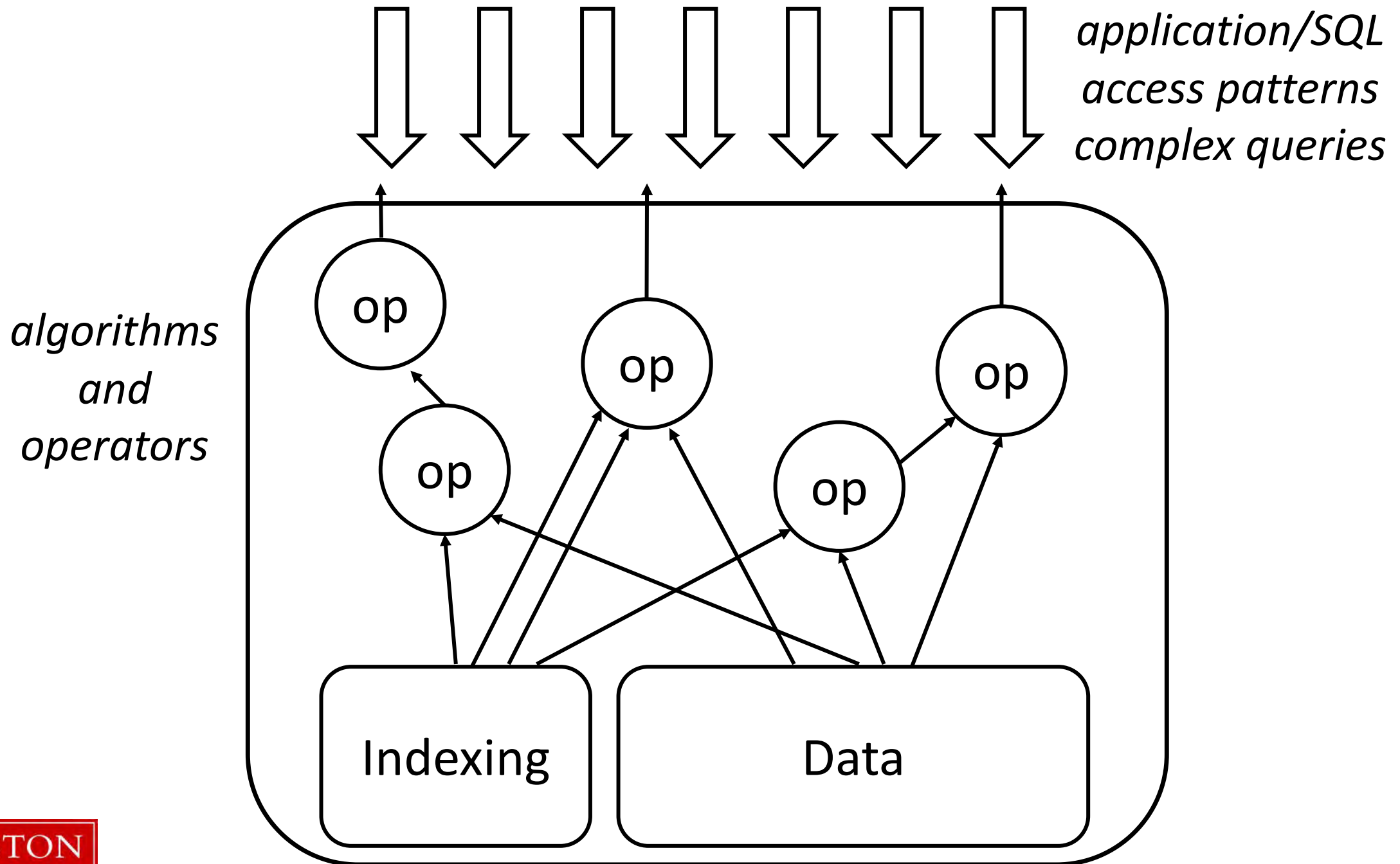
a **data system** is a large software system that **stores data**, and provides the **interface** to **update** and **access** them **efficiently**

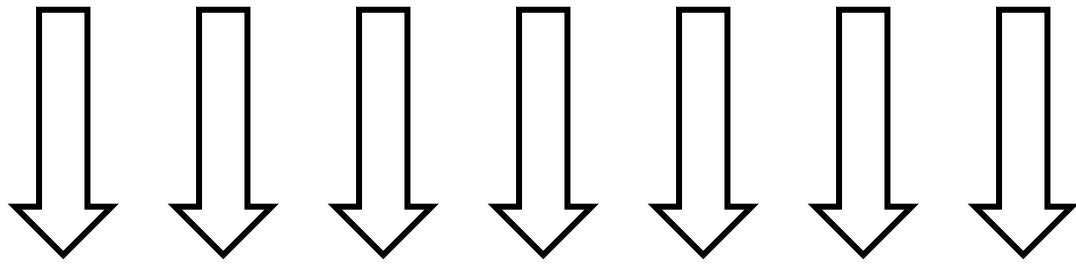


a **data system** is a large software system that **stores data**, and provides the **interface** to **update** and **access** them **efficiently**



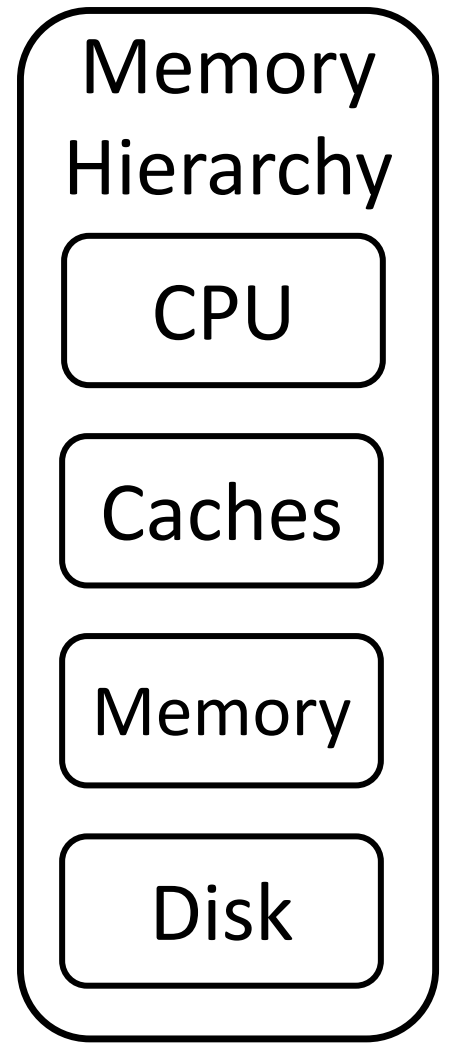
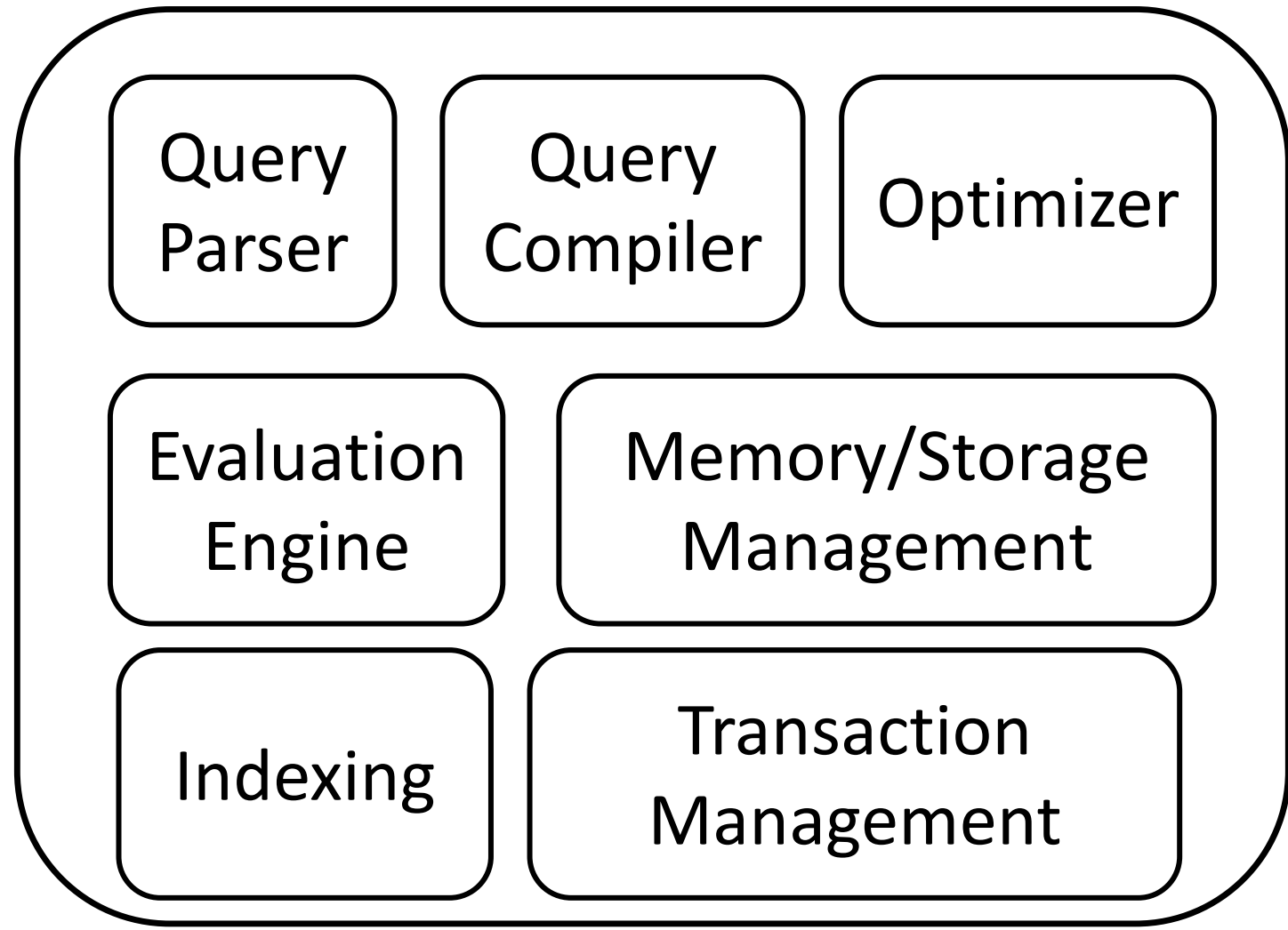
data system, what's inside?





*application/SQL
access patterns
complex queries*

modules



growing environment

ORACLE®

facebook

db
large systems
complex
lots of tuning
legacy

noSQL
simple, clean
“just enough”

Microsoft

IBM

Google

twitter

SAP

>\$200B by 2020, growing at 11.7% every year
[The Forbes, 2016]

more **complex**
applications

need for
scalability

newSQL

[noSQL]

\$3B by 2020, growing at 20% every year
[Forrester, 2016]

growing need for tailored systems



new applications



new hardware

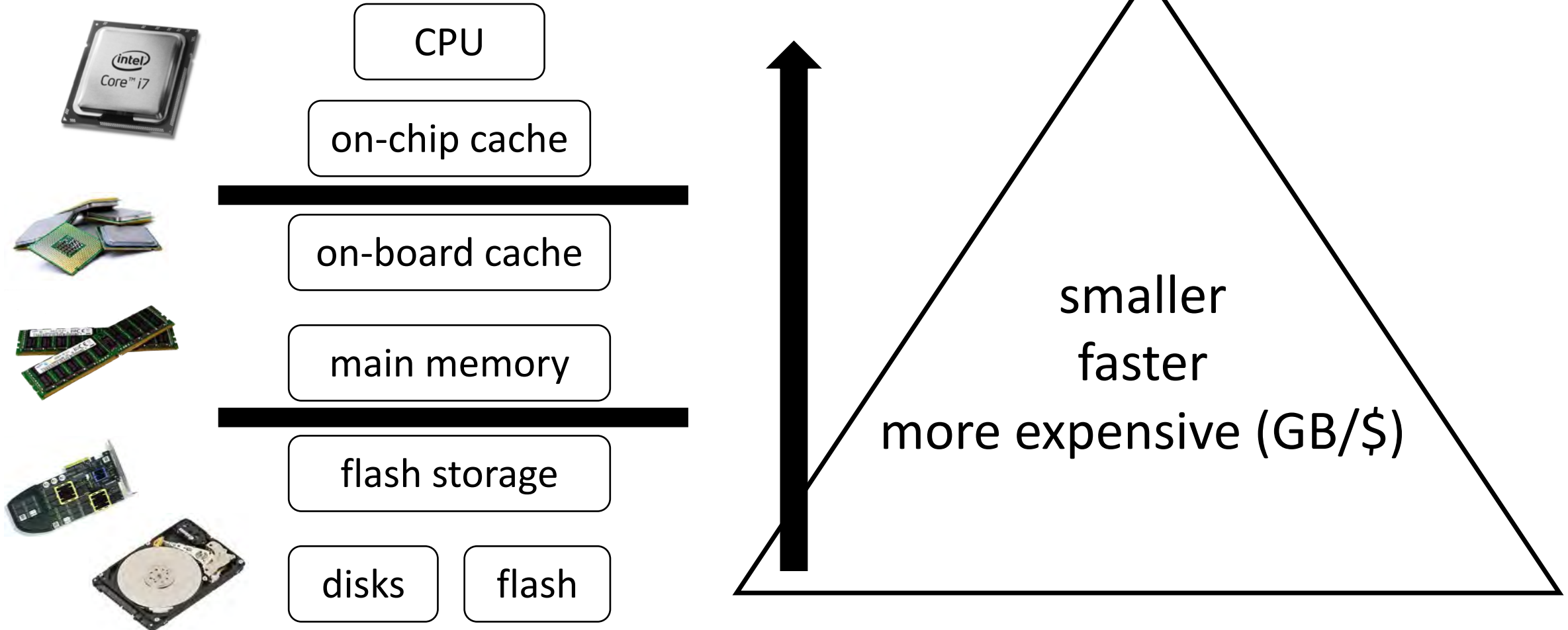


more data

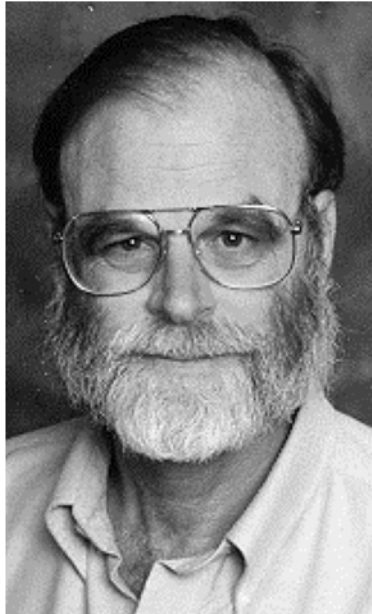


data system, what's underneath?

memory hierarchy



memory hierarchy (by Jim Gray)

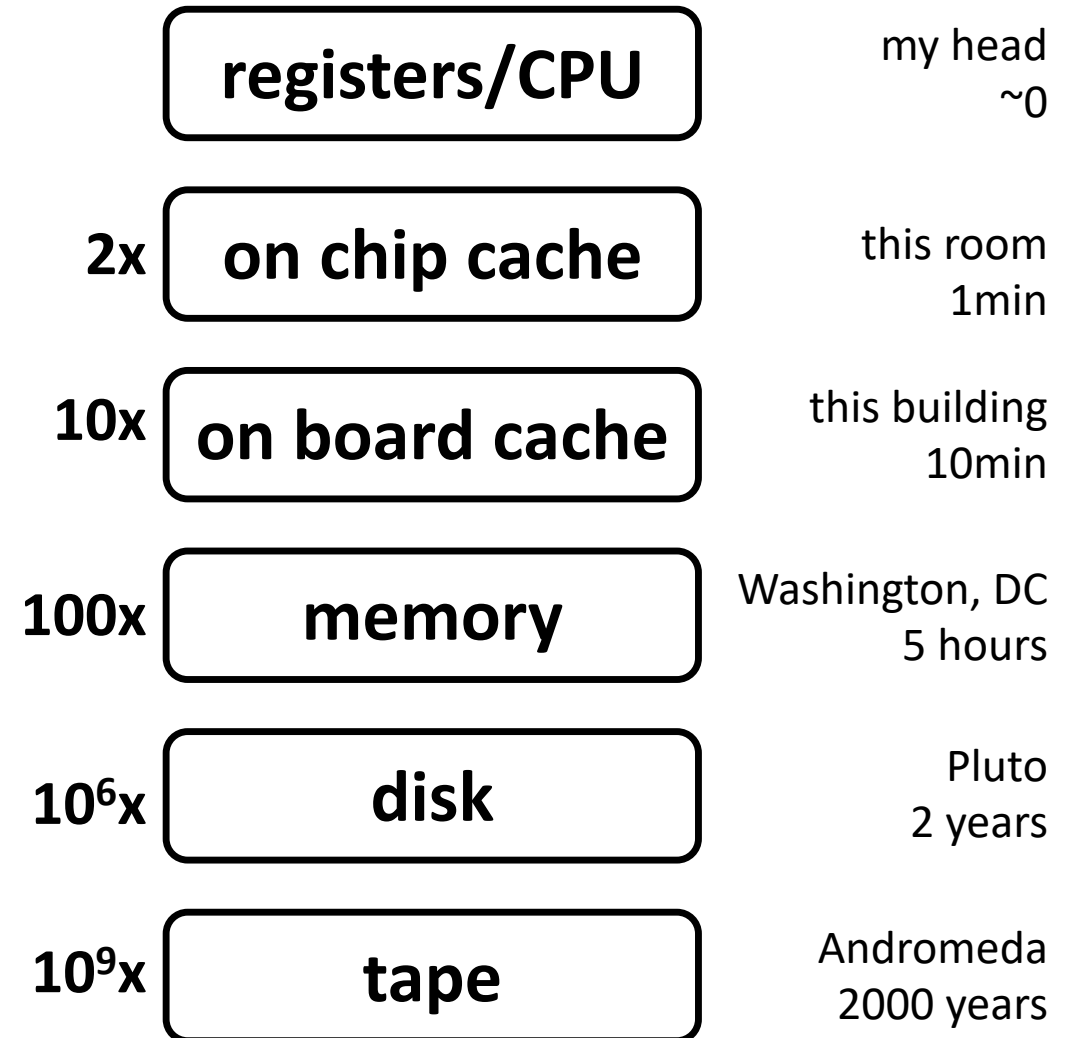


Jim Gray, IBM, Tandem, Microsoft, DEC

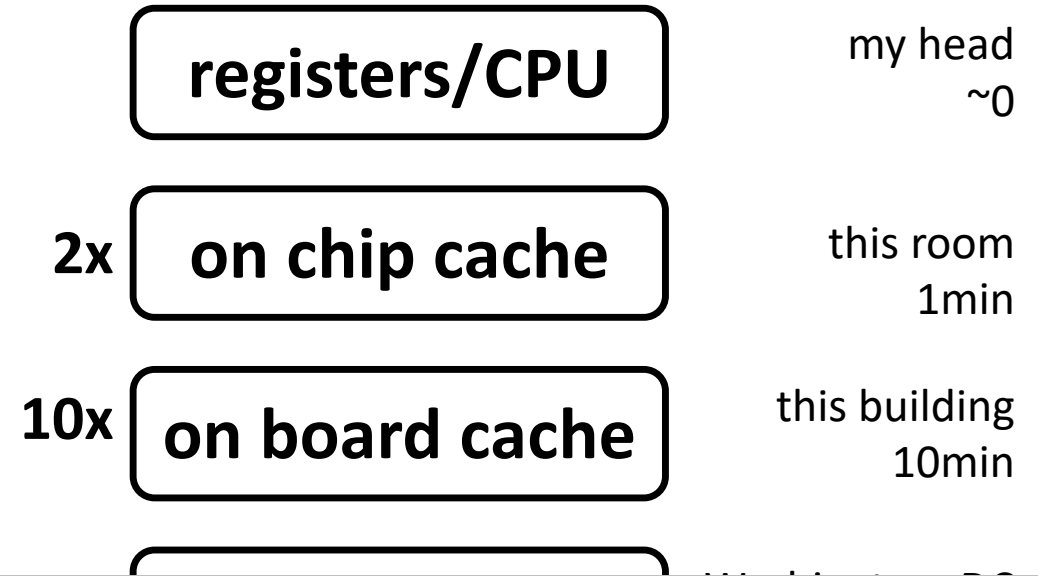
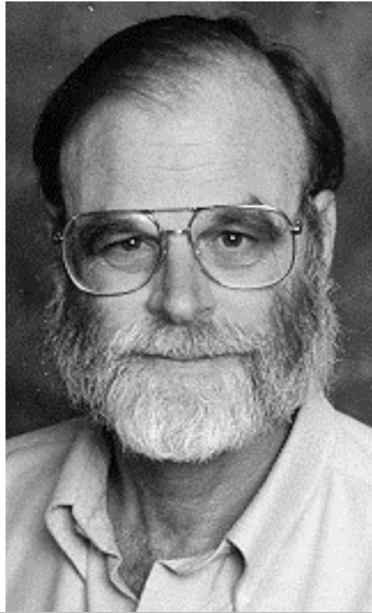
“The Fourth Paradigm” is based on his vision

ACM Turing Award 1998

ACM SIGMOD Edgar F. Codd Innovations award 1993



memory hierarchy (by Jim Gray)

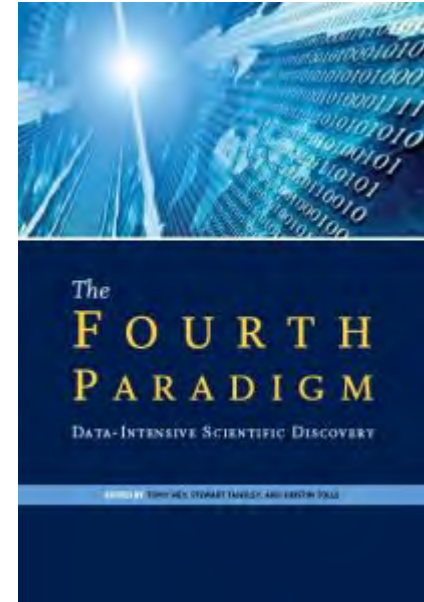
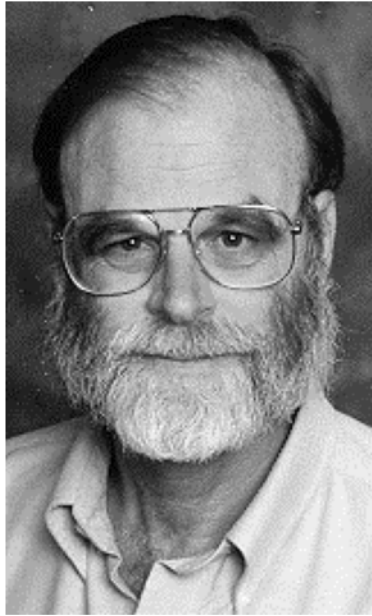


tape?

sequential-only magnetic storage
still a multi-billion industry



Jim Gray (a great scientist and engineer)



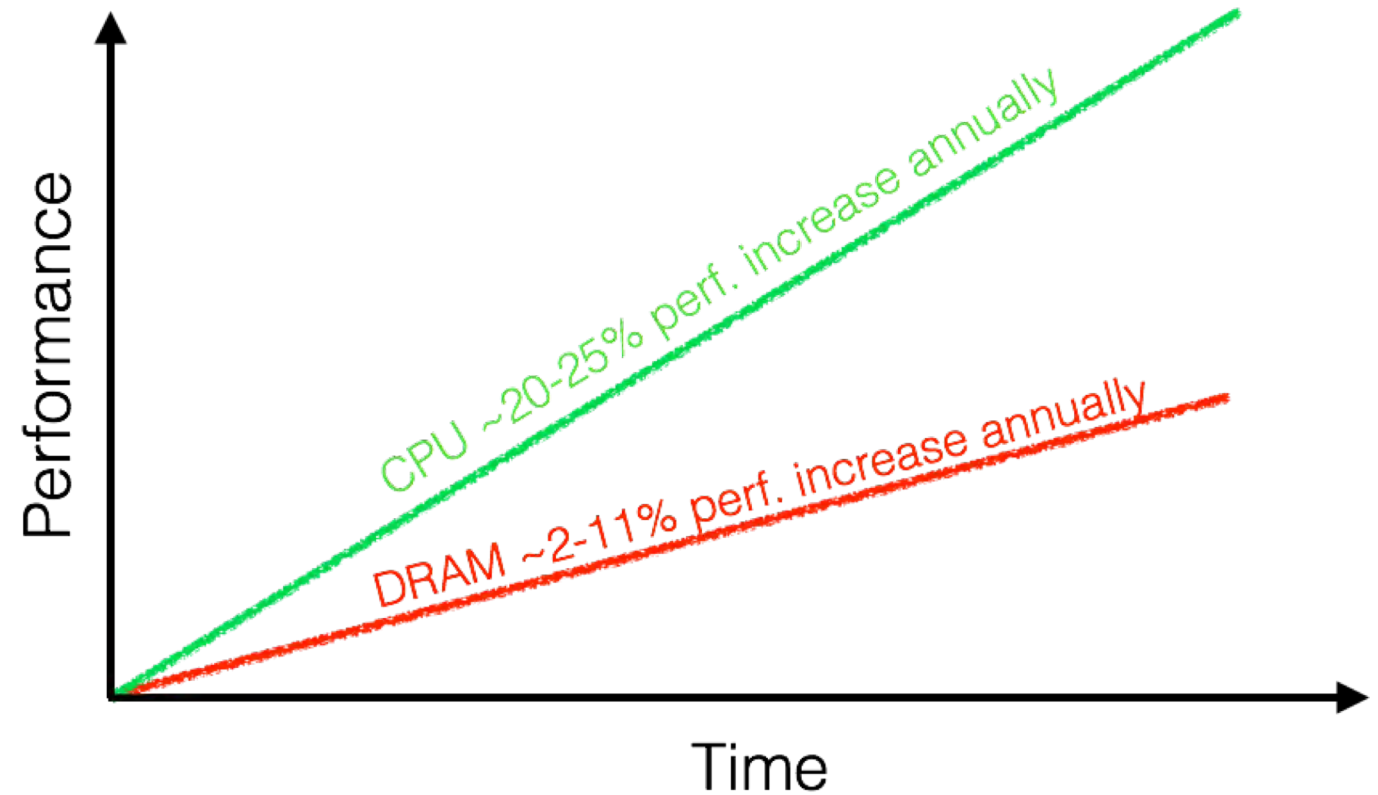
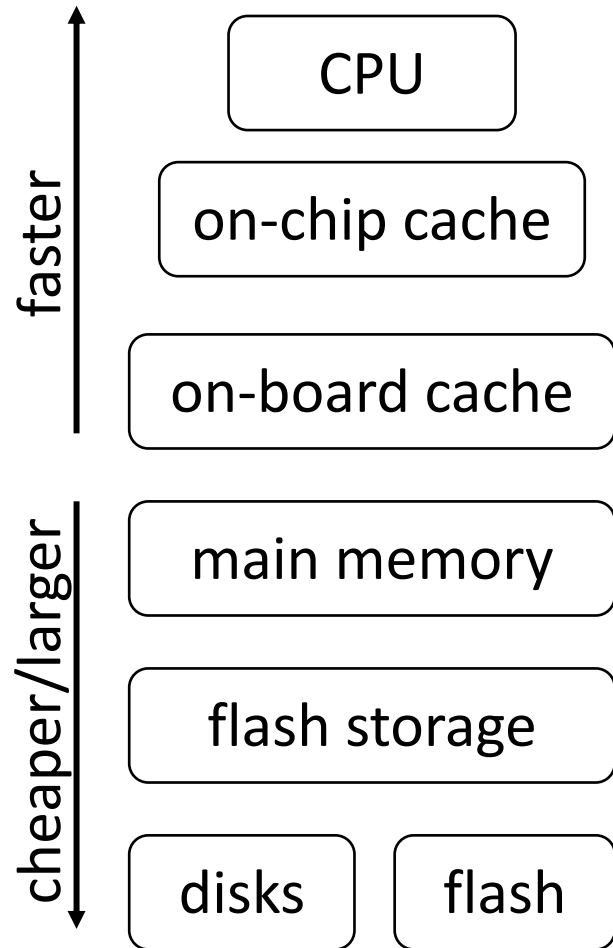
*the first collection of
technical visionary research on
a data-intensive scientific discovery*

Jim Gray, IBM, Tandem, Microsoft, DEC
“The Fourth Paradigm” is based on his vision

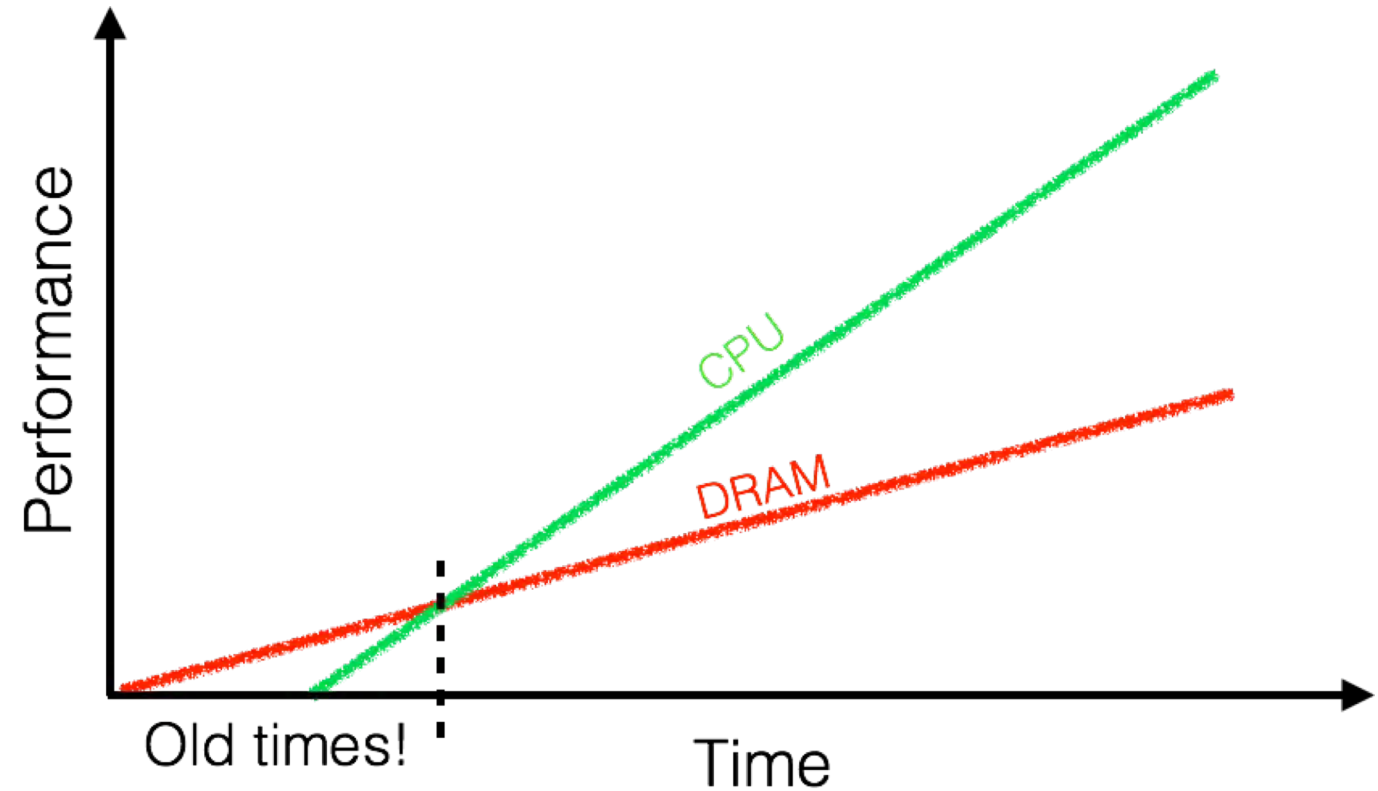
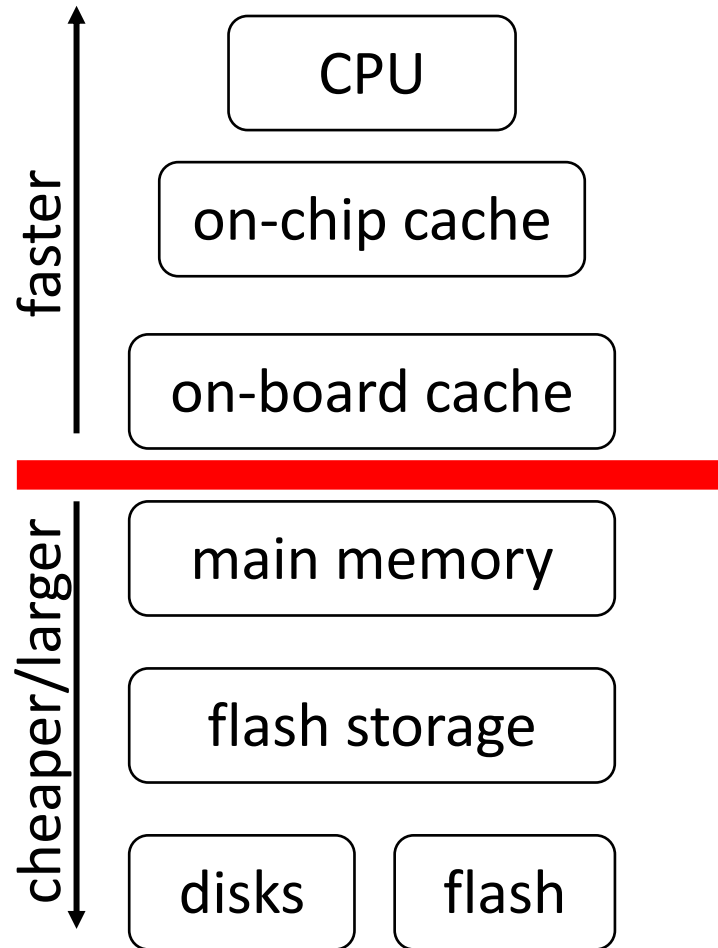
ACM Turing Award 1998

ACM SIGMOD Edgar F. Codd Innovations award 1993

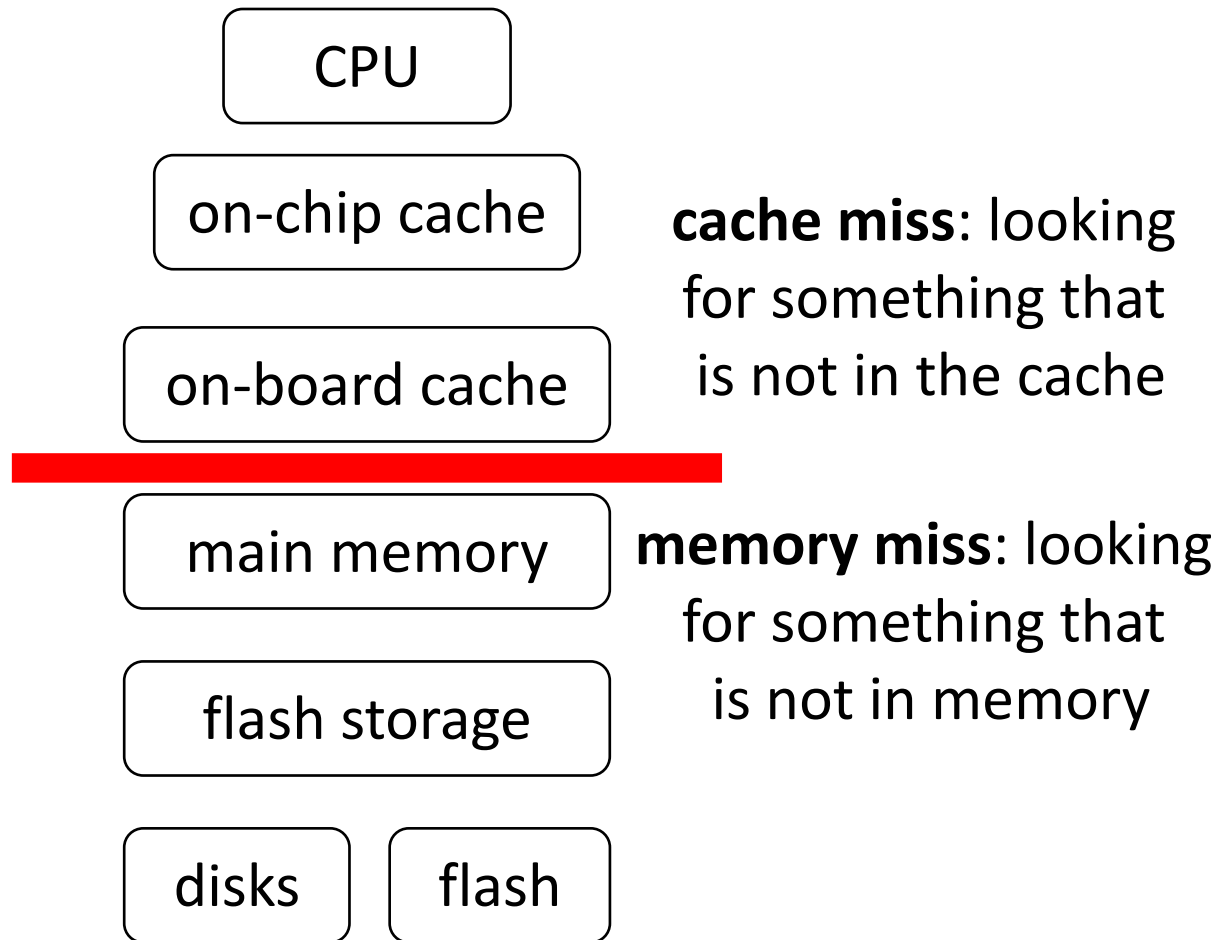
memory wall



memory wall



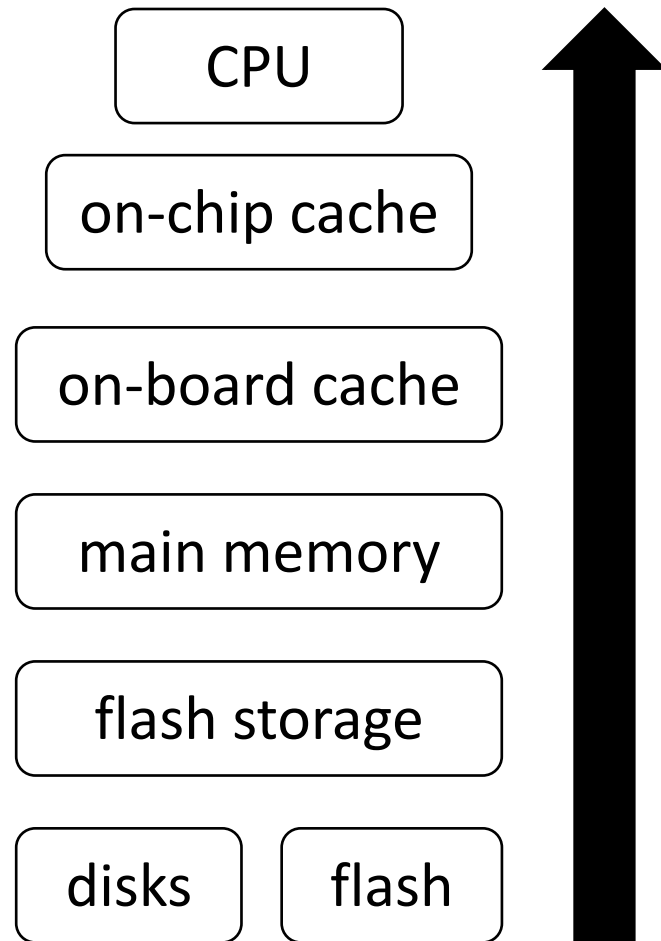
cache/memory misses



what happens if I miss?



data movement



data go through
all necessary levels

also read
unnecessary data

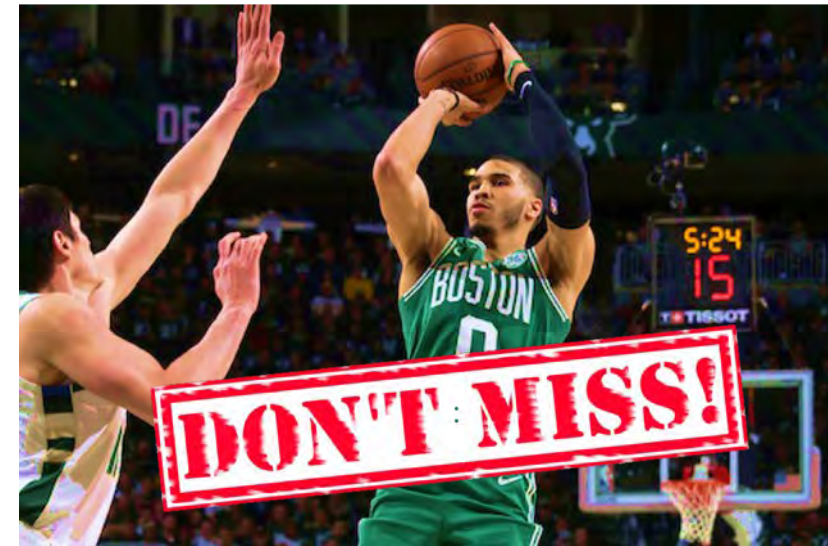
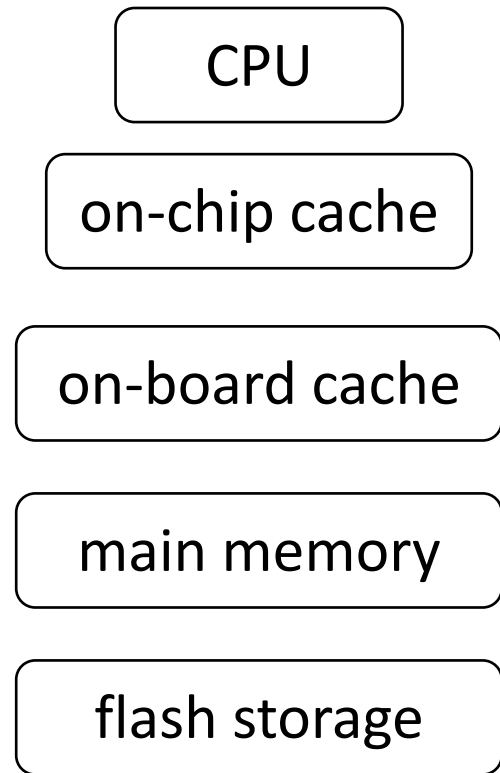


Photo by Gary Dineen/NBAE via Getty Images

need to read only X
read the whole page



data movement



data go through
all necessary levels

also read
unnecessary data



Photo by Gary Dineen/NBAE via Getty Images

need to read only X
read the whole page



remember!

disk is millions (mem, hundreds) times slower than CPU

page-based access & random access

query $x < 7$



size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

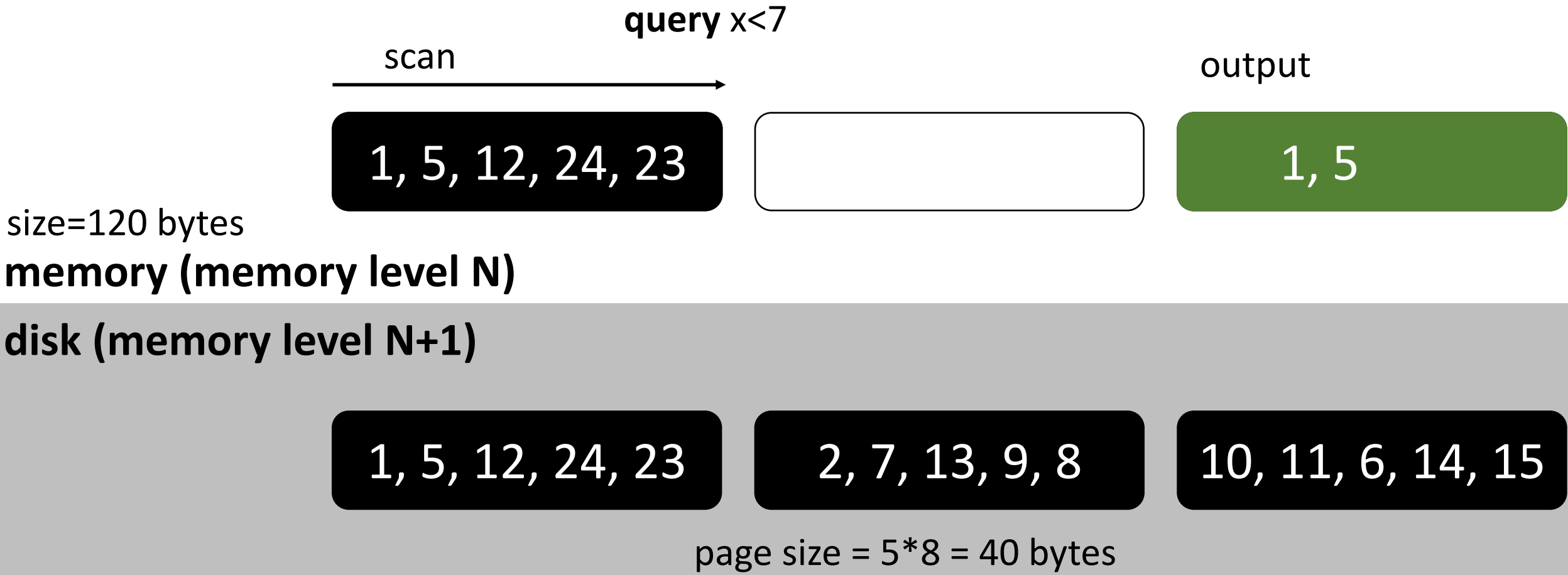
2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 40 bytes

page-based access & random access



\$ 40 bytes

page-based access & random access

query $x < 7$

scan

output

1, 5, 12, 24, 23

2, 7, 13, 9, 8

1, 5

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 40 bytes

page-based access & random access

query $x < 7$

scan

output

1, 5, 12, 24, 23

2, 7, 13, 9, 8

1, 5, 2

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 80 bytes

page-based access & random access

query $x < 7$

scan

output

1, 5, 12, 24, 23

2, 7, 13, 9, 8

1, 5, 2

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

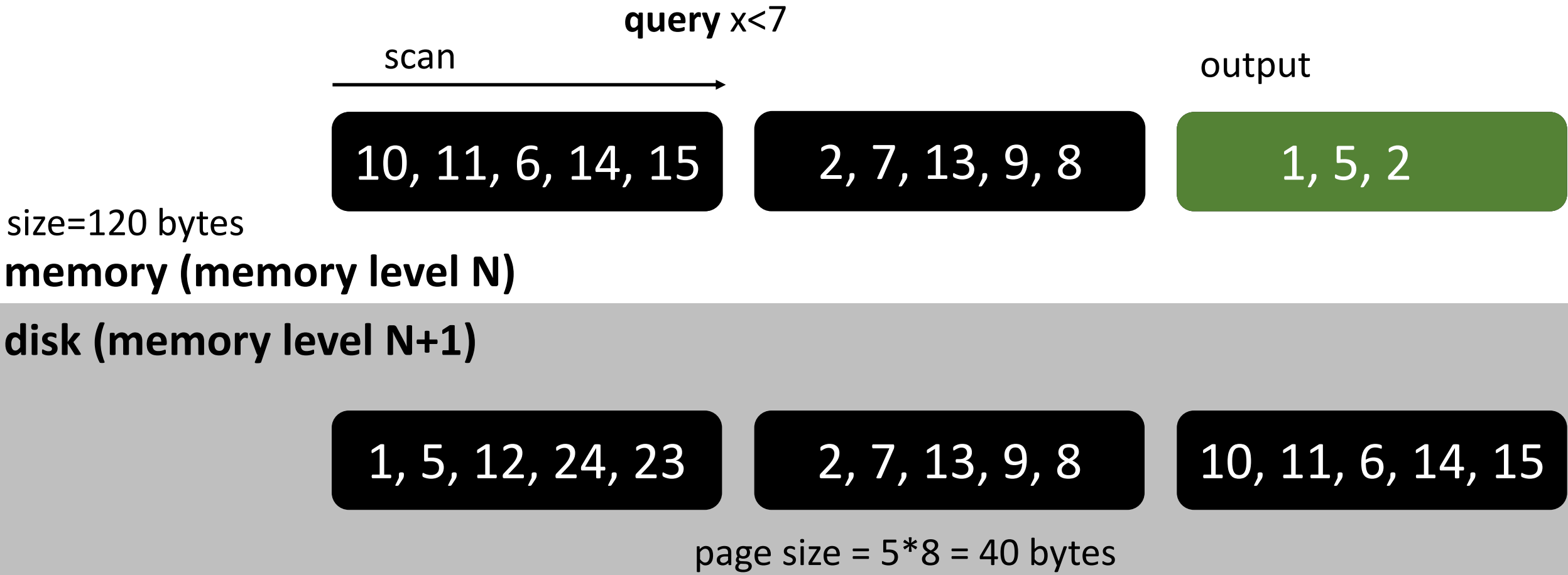
2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

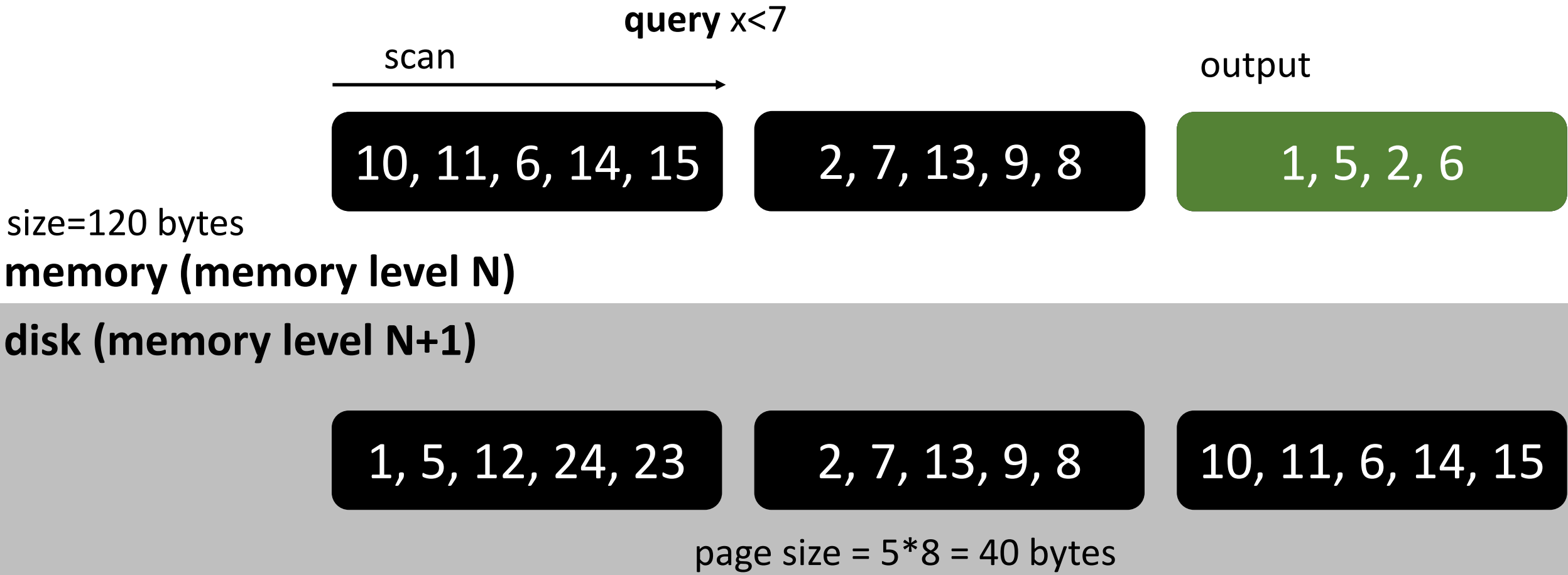
\$ 80 bytes

page-based access & random access



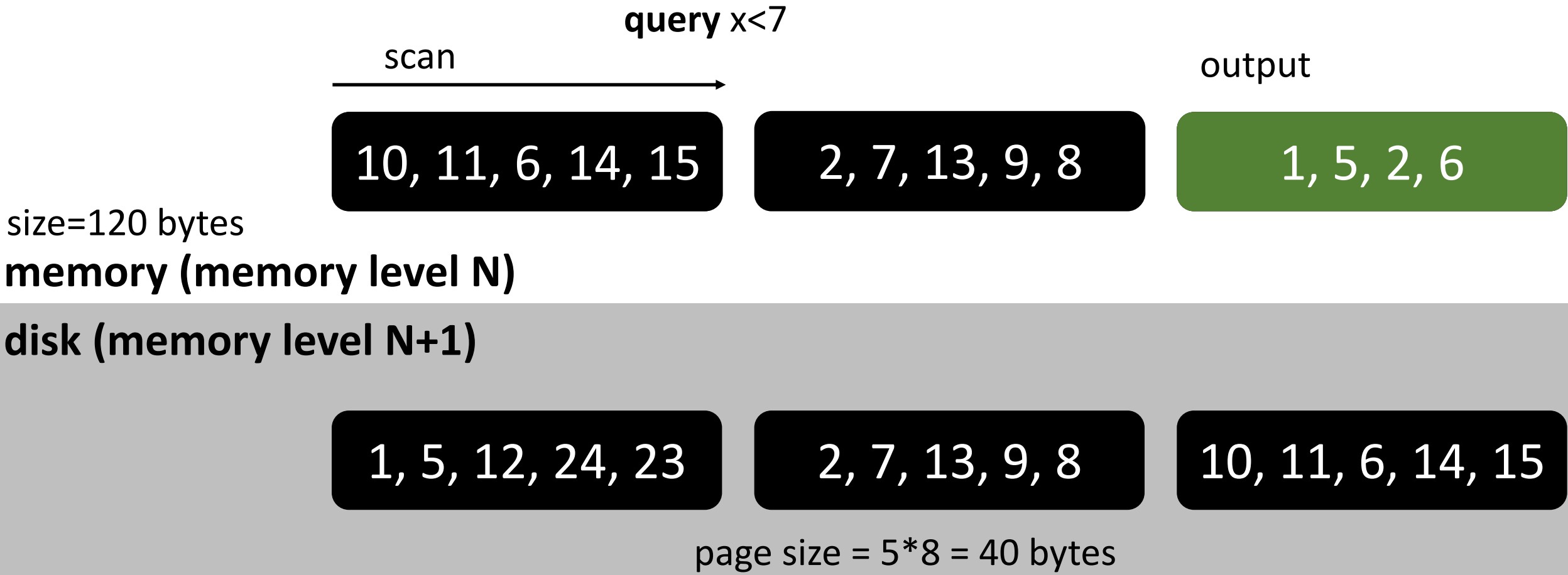
\$ 80 bytes

page-based access & random access



\$120 bytes

page-based access & random access



what if we had an oracle (perfect index)?



page-based access & random access

query $x < 7$



size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 40 bytes

page-based access & random access

query $x < 7$

oracle

1, 5, 12, 24, 23

output

1, 5

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 40 bytes

page-based access & random access

query $x < 7$

oracle

output

1, 5, 12, 24, 23

2, 7, 13, 9, 8

1, 5

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 40 bytes

page-based access & random access

query $x < 7$

oracle

output

1, 5, 12, 24, 23

2, 7, 13, 9, 8

1, 5, 2

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 80 bytes

page-based access & random access

query $x < 7$

oracle

output

1, 5, 12, 24, 23

2, 7, 13, 9, 8

1, 5, 2

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 80 bytes

page-based access & random access

query $x < 7$

oracle

output

10, 11, 6, 14, 15

2, 7, 13, 9, 8

1, 5, 2

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

\$ 80 bytes

page-based access & random access

query $x < 7$

oracle

output

10, 11, 6, 14, 15

2, 7, 13, 9, 8

1, 5, 2, 6

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

page-based access & random access

\$120 bytes



query $x < 7$

was the oracle helpful?

oracle

output

10, 11, 6, 14, 15

2, 7, 13, 9, 8

1, 5, 2, 6

size=120 bytes

memory (memory level N)

disk (memory level N+1)

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

when is the oracle helpful?



for which query would an oracle help us?

how to decide whether to use the oracle?

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

how we store data

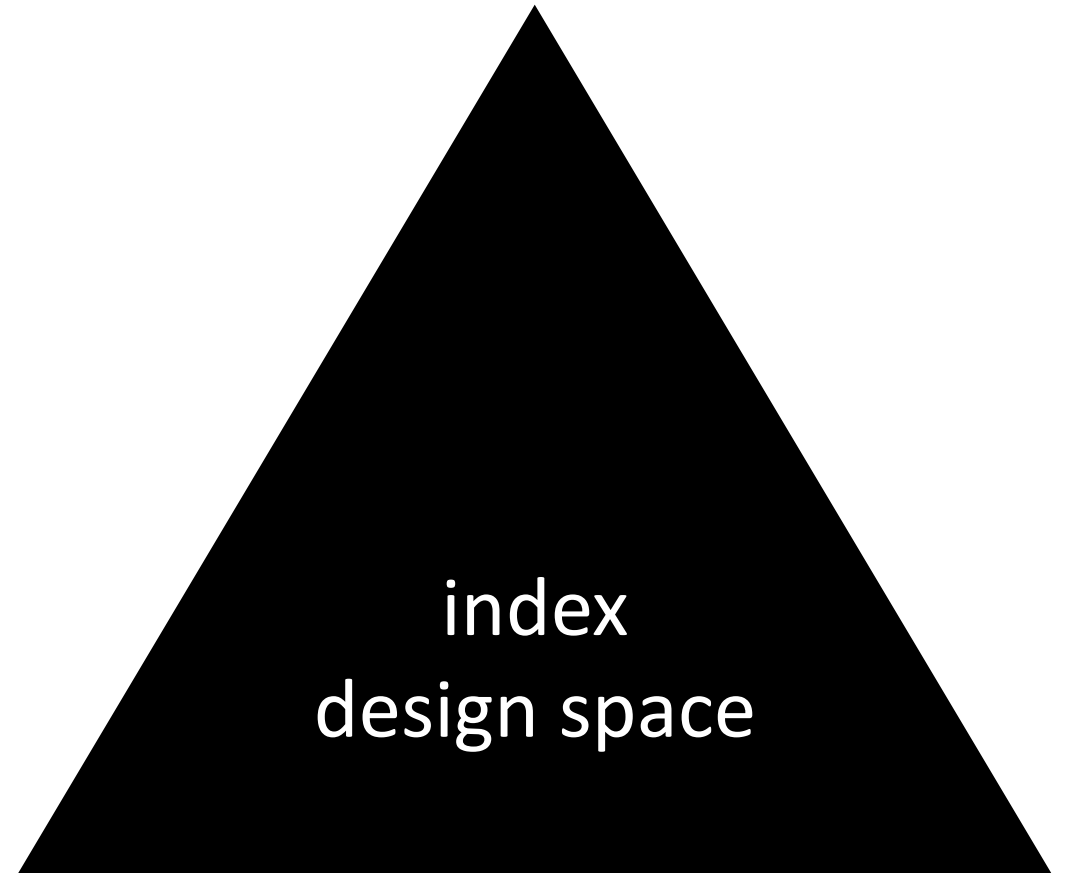
layouts, indexes

every **byte** counts

overheads and tradeoffs

know the **query**

access path selection



rules of thumb

sequential access

read one block; consume it completely; discard it; read next;

hardware can predict and start prefetching

prefetching can exploit full memory/disk bandwidth

random access

read one block; consume it partially; discard it; (may re-use);

read random next;



ideal random access?

the one that helps us **avoid a large number of accesses** (random or sequential)

the language of efficient systems: C/C++

why?

low-level control over hardware

make decisions about physical data placement and consumptions

fewer assumptions

the language of efficient systems: C/C++

why?

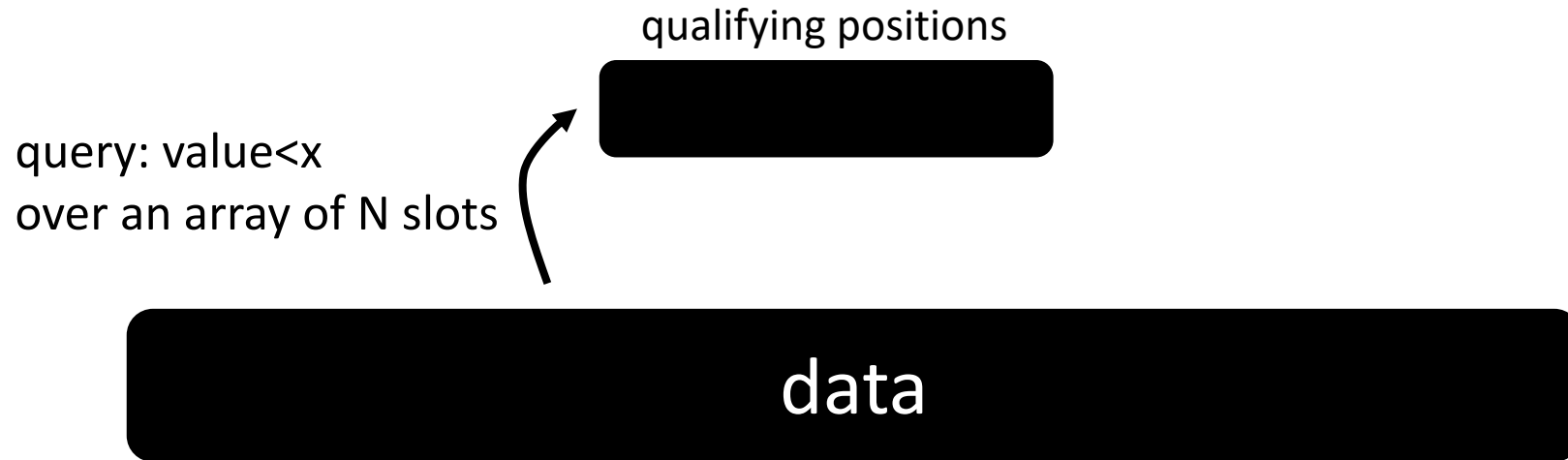
low-level control over hardware

we want you in the project to make low-level decisions

main-memory optimized-systems

a “simple” database operator

select operator (scan)





how to implement it?

```
result = new array[data.size];  
j=0;  
for (i=0; i<data.size; i++)  
  if (data[i]<x)  
    result[j++]=i;
```

qualifying positions

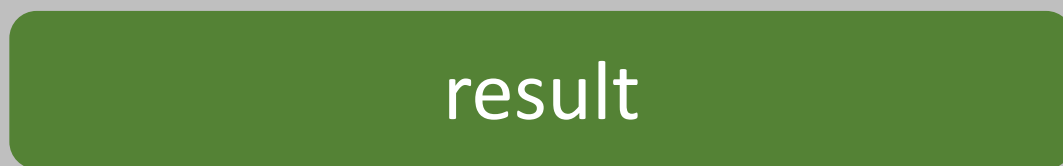


query: value<x
over an array of N slots



what if only 0.1% qualifies?

memory





how to implement it?

```
result = new array[data.size];  
j=0;  
for (i=0; i<data.size; i++)  
  if (data[i]<x)  
    result[j++]=i;
```

qualifying positions

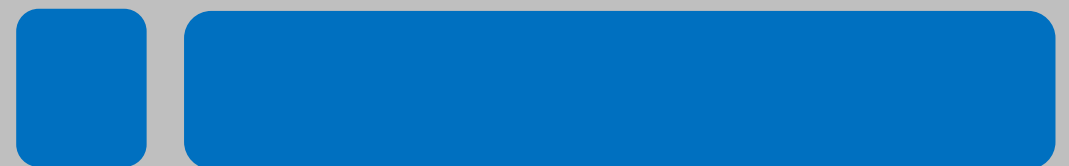
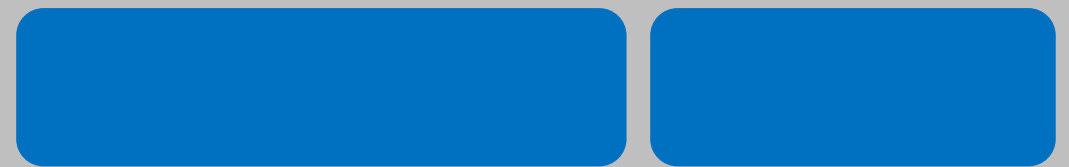


query: value<x
over an array of N slots



what if only 0.1% qualifies?

memory





how to implement it?

```
result = new array[data.size];  
j=0;  
for (i=0; i<data.size; i++)  
  if (data[i]<x)  
    result[j++]=i;
```

```
result = new array[data.size];  
j=0;  
for (i=0; i<data.size; i++)  
  result[j+=(data[i]<x)]=i;
```

qualifying positions



query: value<x
over an array of N slots

data

what if 99% qualifies?



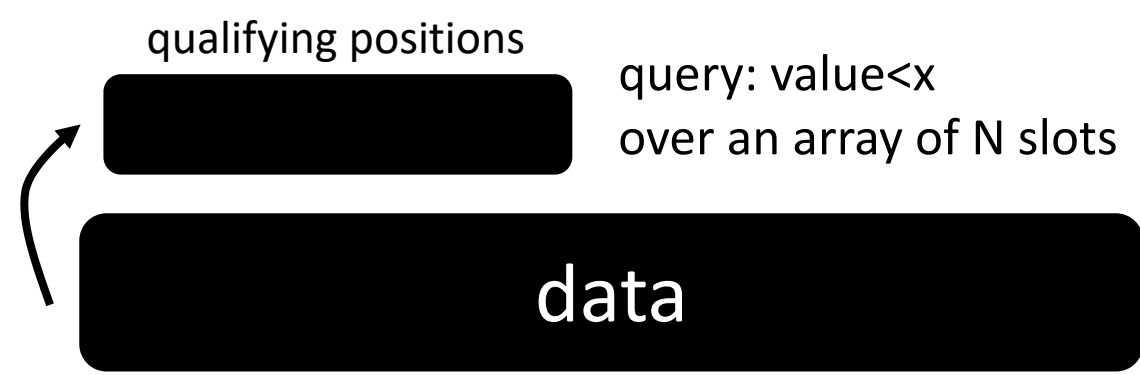
how can we know?

branches (if statements)
are bad for the processors,
can we avoid them?

how to bring the values?
(remember we have the positions)

```
result = new array[data.size];
j=0;
for (i=0; i<data.size; i++)
  if (data[i]<x)
    result[j++]=i;
```

needs coordination!
what about result writing?



what about multi-core?
NUMA? SIMD? GPU?

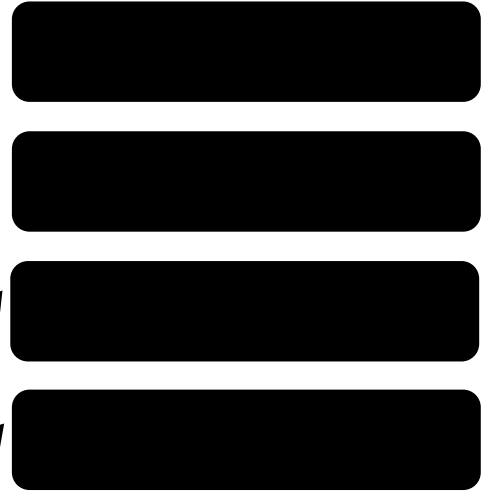


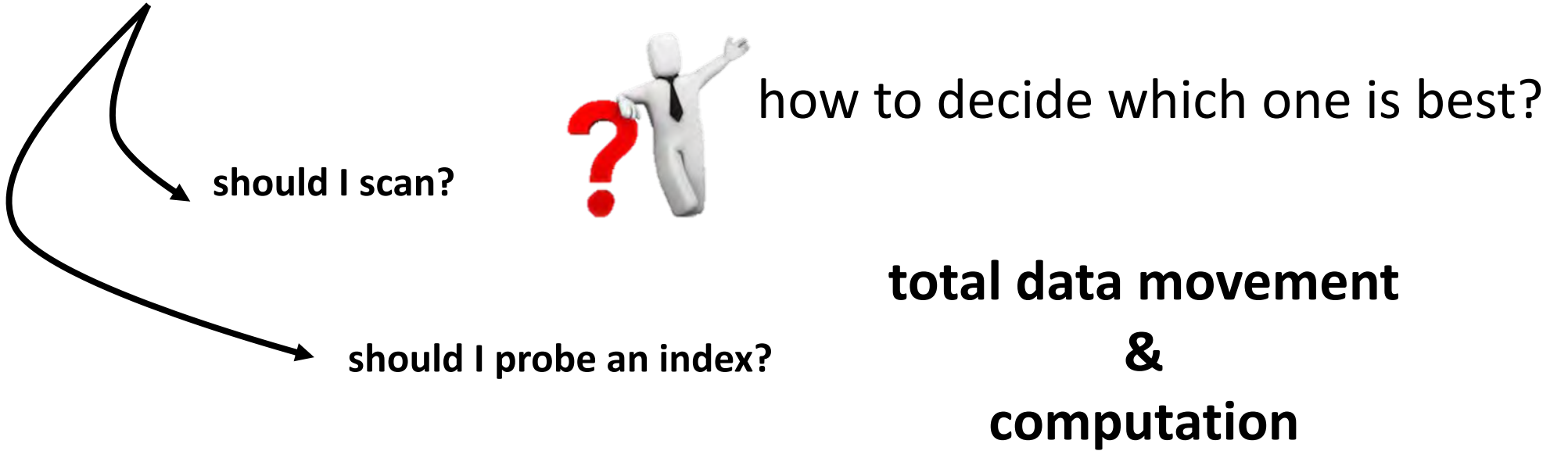


what about having multiple queries?

query1: value<x1
query2: value<x2 ...

```
result = new array[data.size];  
j=0;  
for (i=0; i<data.size; i++)  
  if (data[i]<x)  
    result[j++]=i;
```





how can I prepare?

1) Read background research material

- **Architecture of a Database System.** By J. Hellerstein, M. Stonebraker and J. Hamilton. Foundations and Trends in Databases, 2007
- **The Design and Implementation of Modern Column-store Database Systems.** By D. Abadi, P. Boncz, S. Harizopoulos, S. Idreos, S. Madden. Foundations and Trends in Databases, 2013
- **Massively Parallel Databases and MapReduce Systems.** By Shivnath Babu and Herodotos Herodotou. Foundations and Trends in Databases, 2013

2) Start going over the papers

what to do now?

- A) read the syllabus and the website
- B) register to piazza
- C) register to gradescope
- D) register for the presentation (**early next week!**)
- E) start submitting paper reviews (week 3)
- F) go over the project (next week will be available)
- G) start working on the proposal (week 3)

survival guide

class website: <https://midas.bu.edu/classes/CS591A1/>

piazza website: <http://piazza.com/bu/spring2020/cs591a1/>

presentation registration: <https://tinyurl.com/S2020-CS591-presentations>

gradescope entry-code: 9568G3

office hours: Manos (Tu/Th, before class)

Andy (M/W 3-4pm), Ju Hyoung (M 11am-noon / F 3-4pm)

material: papers available from BU network

class 2

Data Systems 101

next week: modern main-memory data systems
&
semester project