

CS 561: Data Systems Architectures

class 2

Data Systems 101

Prof. Manos Athanassoulis

<https://bu-disc.github.io/CS561/>



some reminders

no smartphones



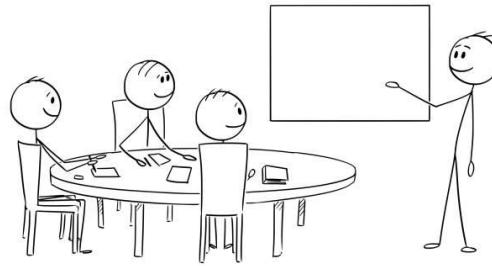
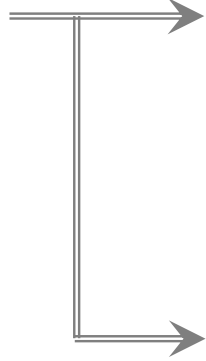
no laptop



What do we do in this class?



reading papers



presentations



quizzes + exam



projects

class summary

2 classes per week & OH/Labs 4 multiple per week

each student

1 paper discussion (as presenter or critic or proponent) + Quizzes + exam

project 0 (individual project) + project 1 (group project)

systems or research project (group project)

proposal + mid-semester report + final report/presentation

Projects

project 0

A small implementation project
to sharpen dev skills

independent project



Due on Jan 31, 2025

AND

project 1

A medium project to give you a flavor of
large-scale production system

groups of 3



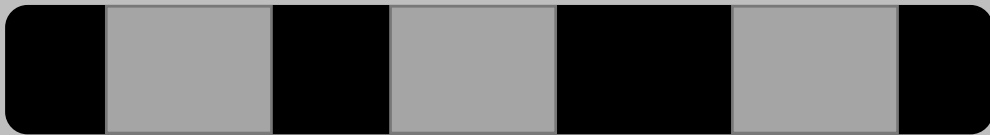
Projects

AND

project 0

A small implementation project
to sharpen dev skills

independent project



Due on Jan 31, 2025

project 1

A medium project to give you a flavor of
large-scale production system

groups of 3



Start forming groups NOW!

Due on Feb 14, 2025

Projects

systems project

groups of 3

implementation-heavy C/C++ project

[illegible]

OR

research project

groups of 3

pick a subject (list available on the website)

design & analysis

experimentation

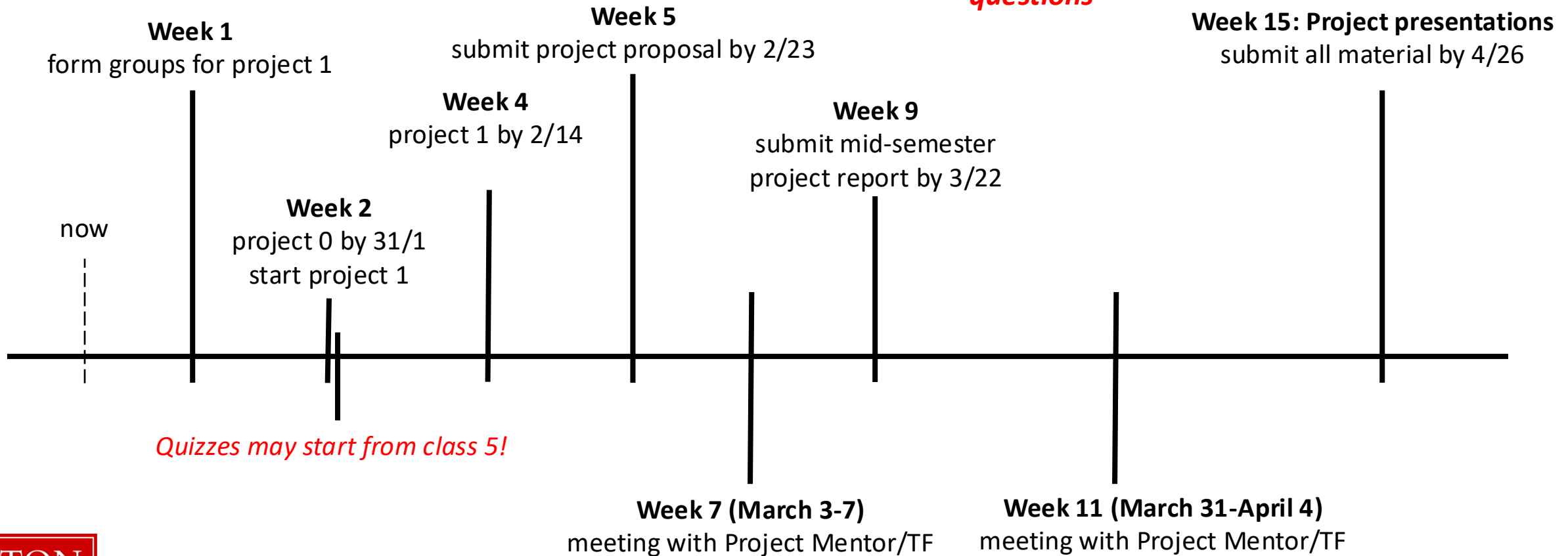




class timeline



discussions
interaction in OH & Lab
questions



Piazza



2 classes per week & OH/Labs multiple times per week

all discussions & announcements

<http://piazza.com/bu/spring2026/cs561/>

also available on class website

We have added everyone who already registered!

Please double-check!

size (volume)

rate (velocity)

sources (variety)

veracity & value

big data

(it's not only about size)

The 3 V's

size (volume)

rate (velocity)

sources (variety)


veracity & value

big data

(it's not only about size)

The 3 V's

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

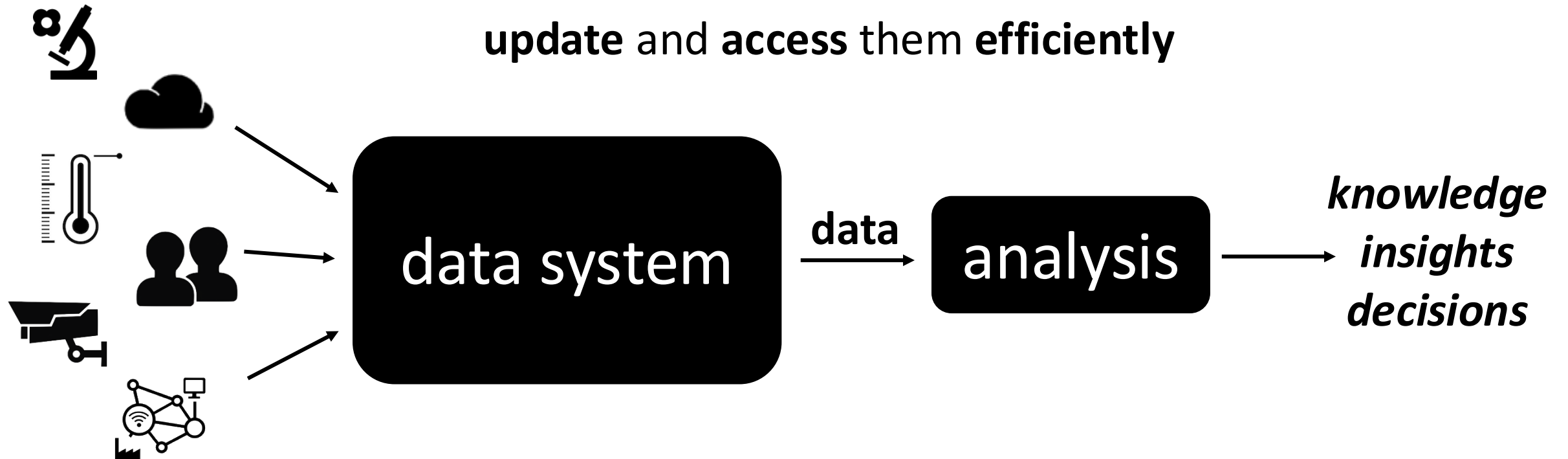
 scientific experiments

social 

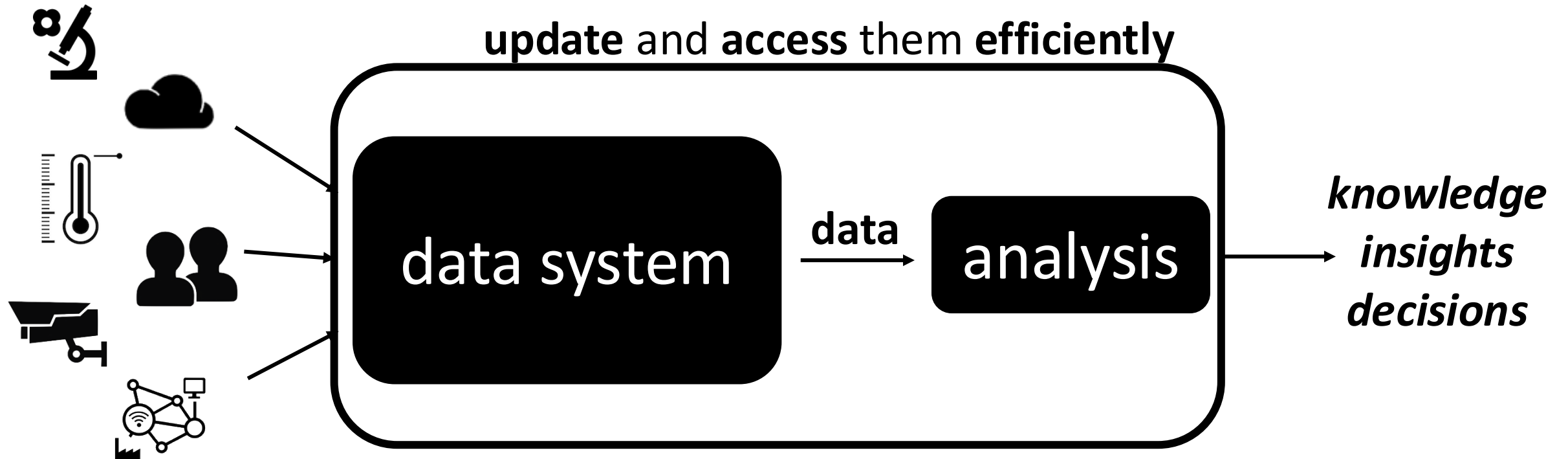
 sensors

Internet-of-things 

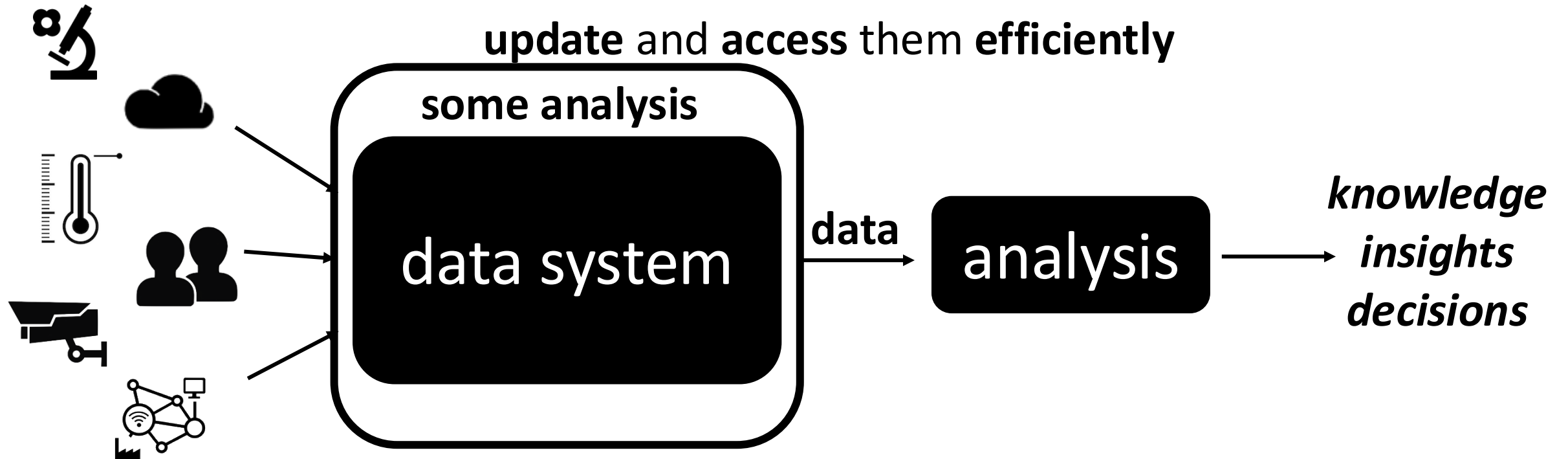
a **data system** is a large software system
that **stores data**, and provides the **interface** to
update and **access** them **efficiently**



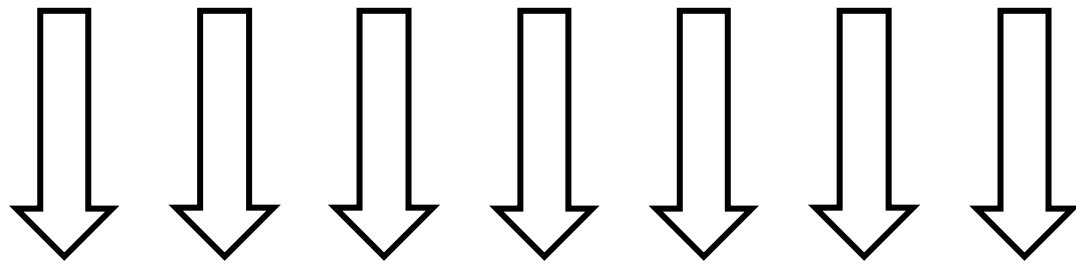
a **data system** is a large software system
that **stores data**, and provides the **interface** to
update and **access** them **efficiently**



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



data system: breaking the blackbox

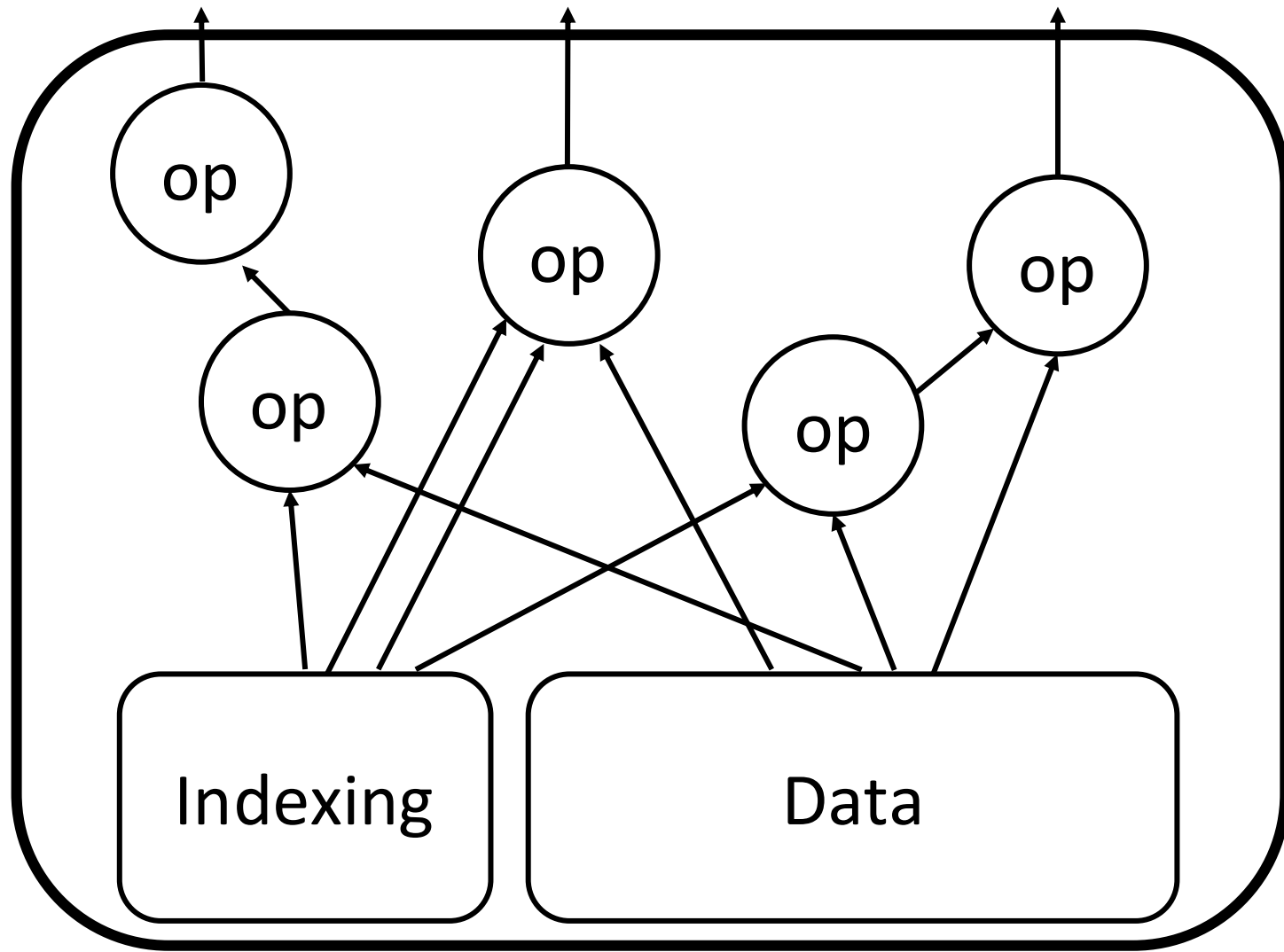


*application/SQL
access patterns
complex queries*

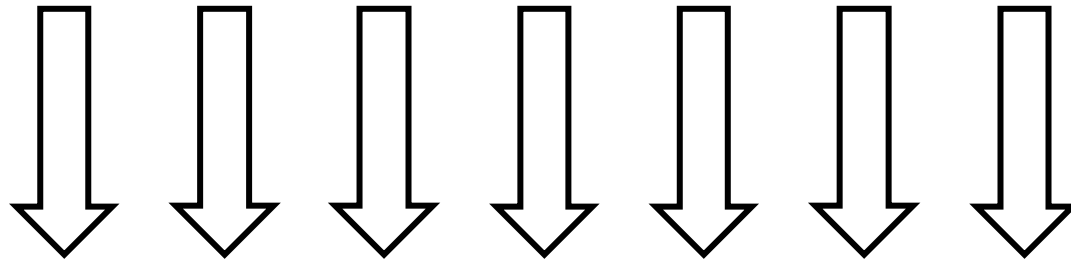


*algorithms
&
operators*

*data &
metadata*

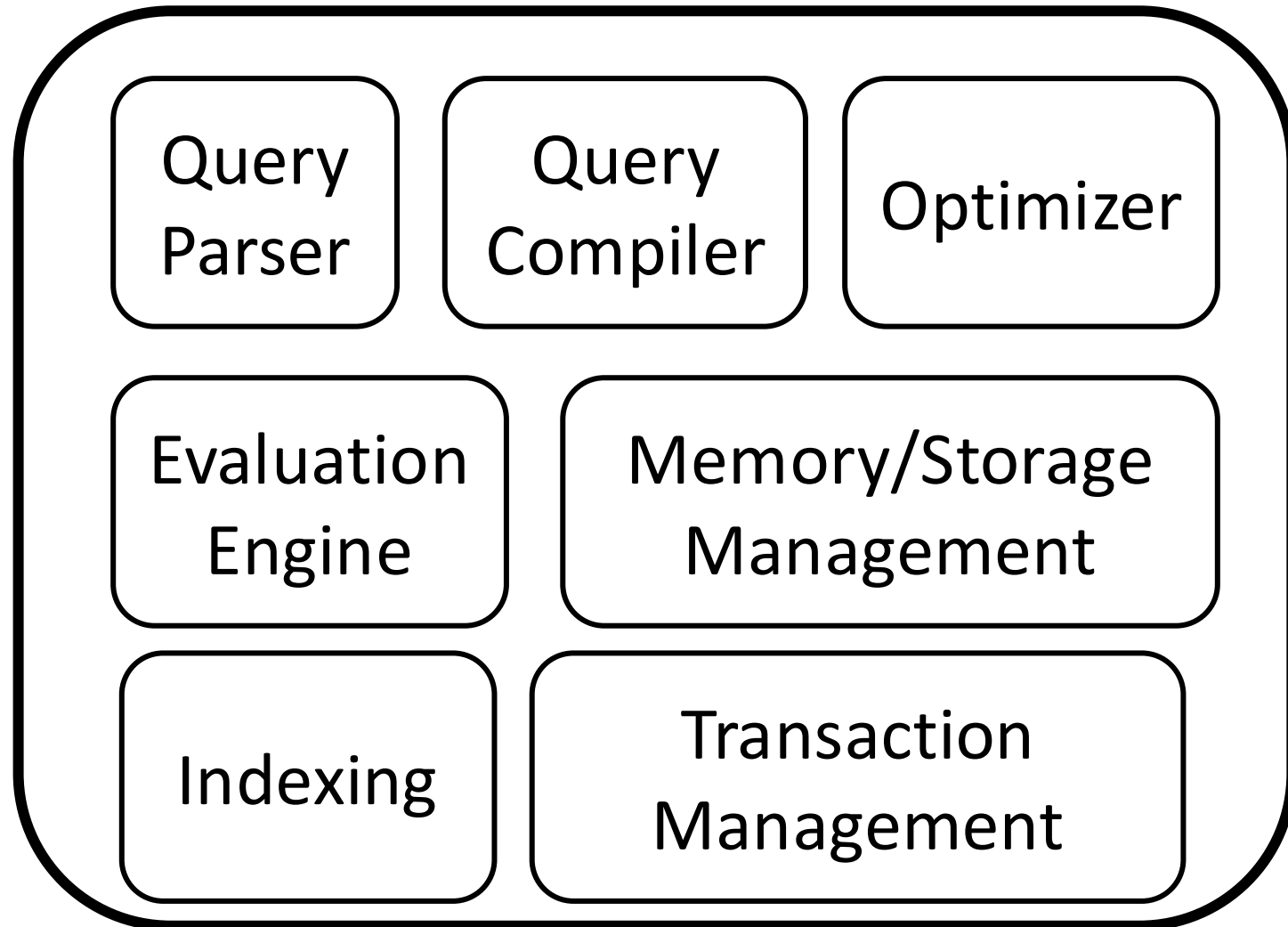


*selection
projection
join
aggregate
hashing
sorting*

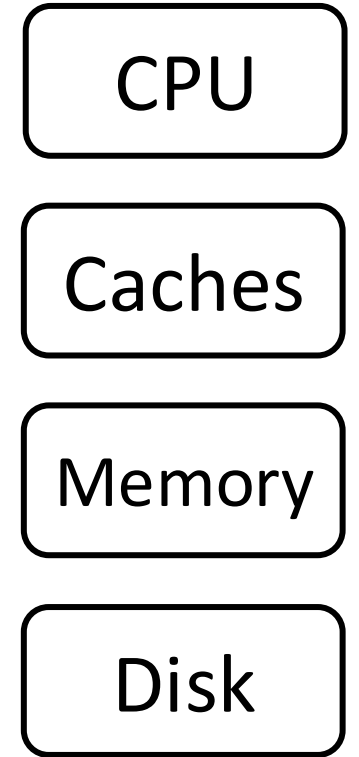


*application/SQL
access patterns
complex queries*

modules



*memory
hierarchy*



growing environment



DB

ACID

large systems

complex

lots of tuning

noSQL

BASE

simple, clean

“just enough”



®

>\$200B by 2020, growing at 11.7% every year

[The Forbes, 2016]

growing environment



DB

ACID

large systems
complex
lots of tuning

noSQL

BASE

simple, clean
“just enough”



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



APACHE



Cockroach Labs



mongoDB®



Couchbase

SCYLLA

\$3B by 2020, growing at 20% every year

[Forrester, 2016]

growing environment

ORACLE®

facebook

DB

ACID

large systems
complex
lots of tuning

noSQL

BASE

simple, clean
“just enough”

more **complex**
applications

need for
scalability

newSQL



IBM

Google



>\$200B by 2020, growing at 11.7% every year

[The Forbes, 2016]



APACHE



Cockroach Labs



mongoDB®



Couchbase

SCYLLA

\$3B by 2020, growing at 20% every year

[Forrester, 2016]

growing environment

ORACLE®

facebook

DB

ACID
large systems
complex
lots of tuning

noSQL

BASE
simple, clean
“just enough”



IBM

Google



more **complex**
applications

need for
scalability



newSQL

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



APACHE



Cockroach Labs



mongoDB®



Couchbase

SCYLLA

\$3B by 2020, growing at 20% every year

[Forrester, 2016]

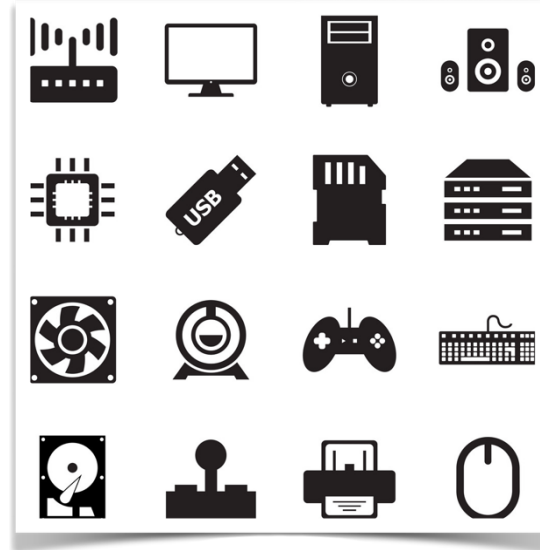
growing need for tailored systems



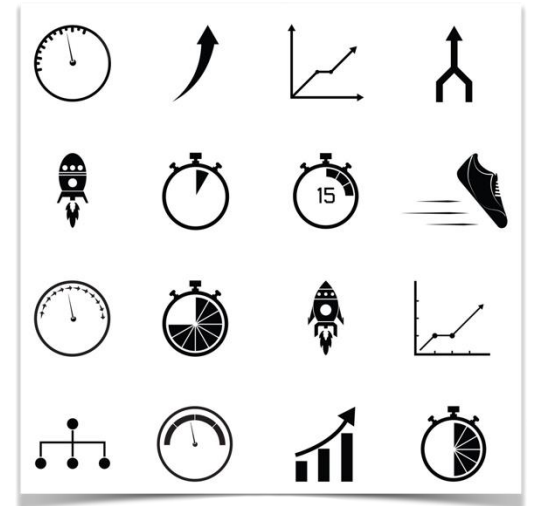
more data



new hardware



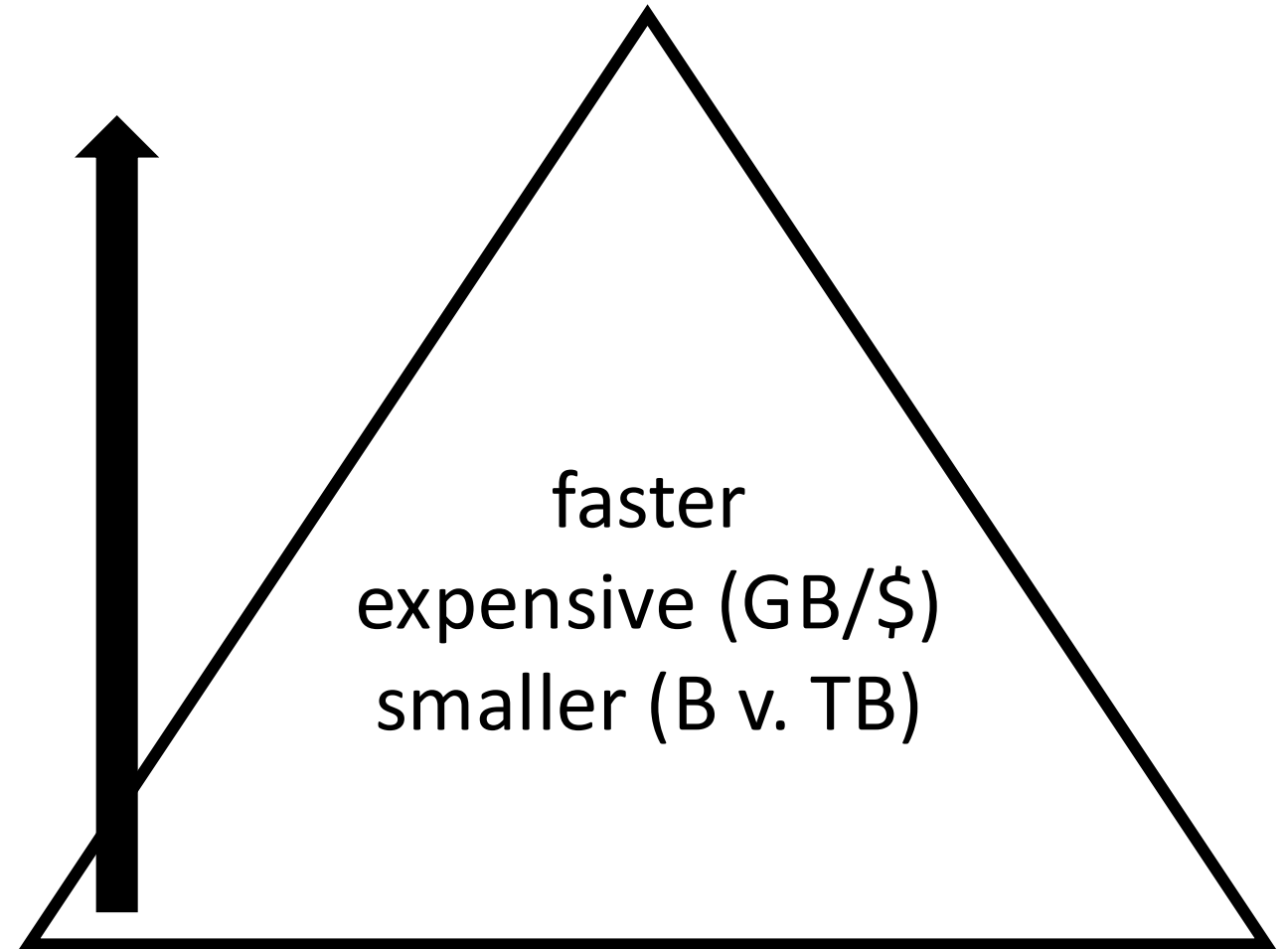
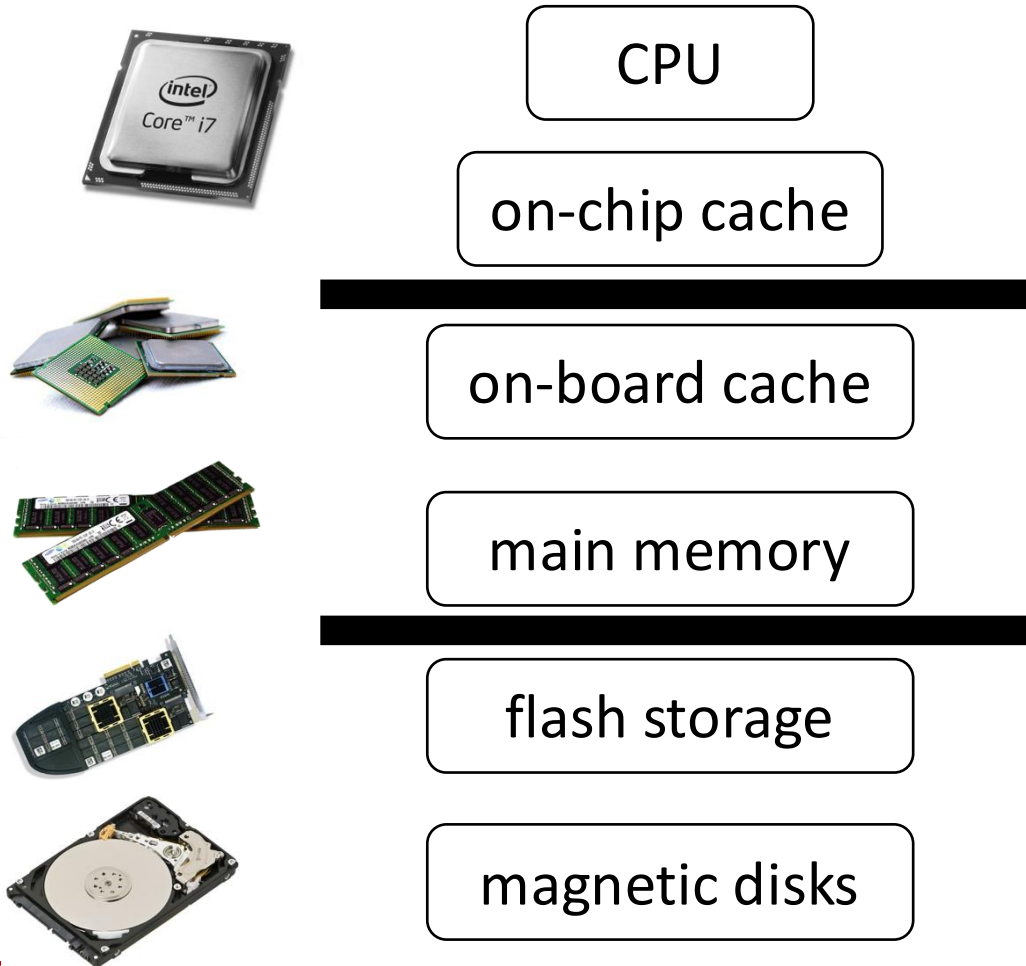
new applications



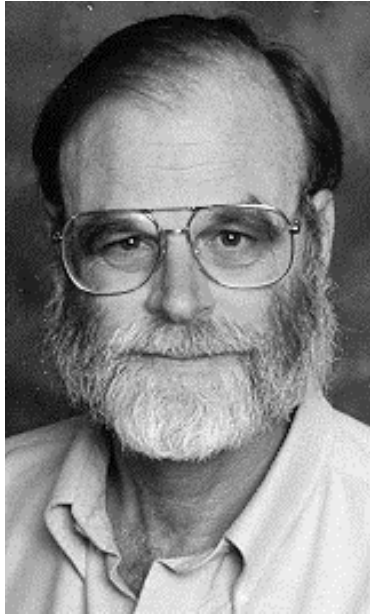
new performance
goals

data systems & the hardware

memory hierarchy



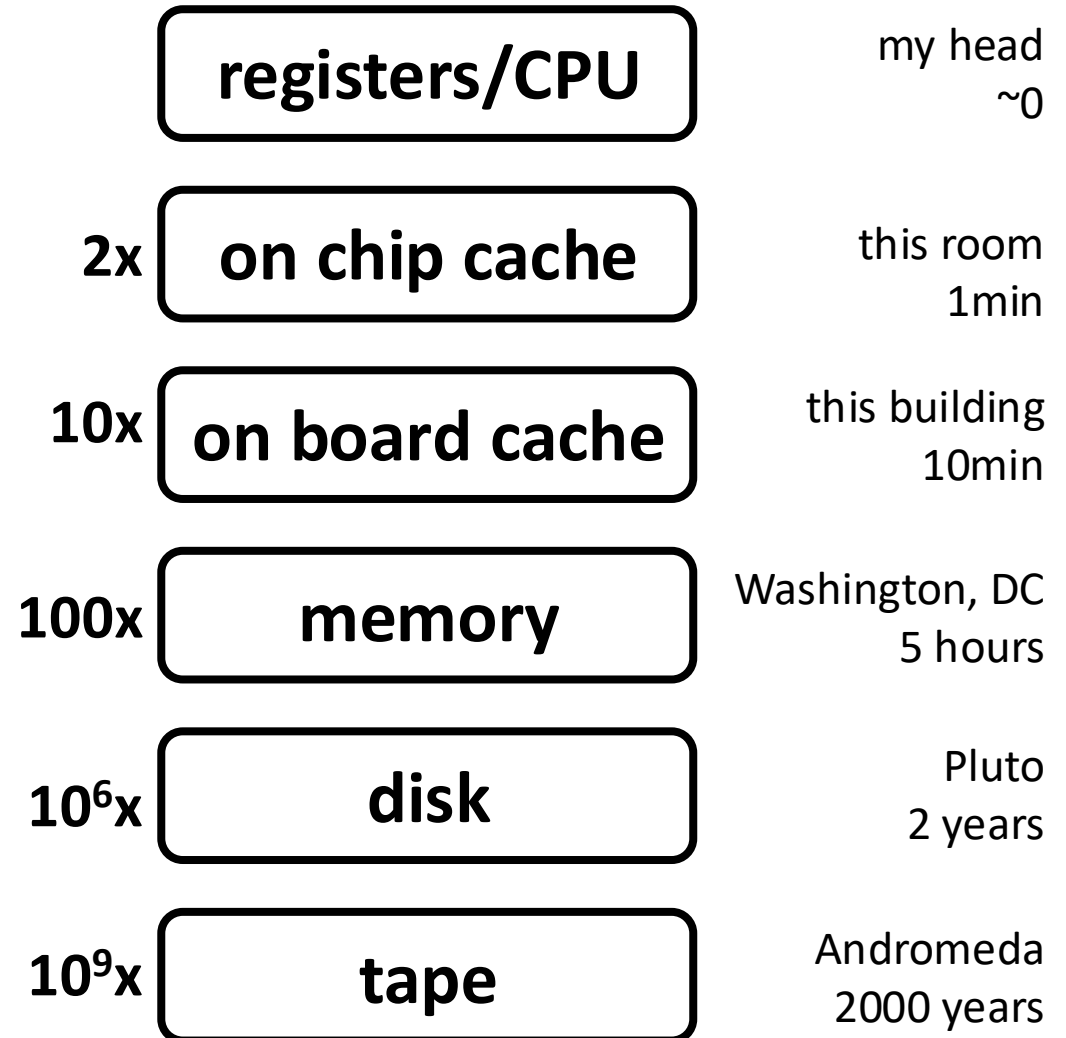
memory hierarchy (by Jim Gray)



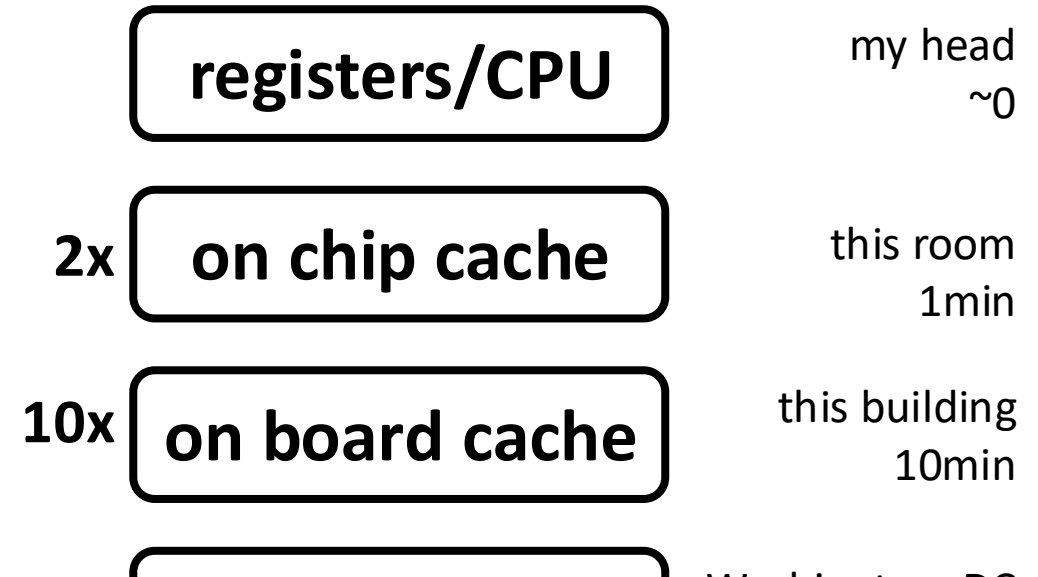
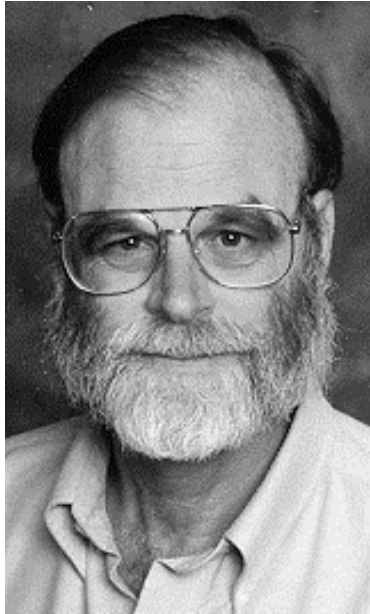
Jim Gray, IBM, Tandem, Microsoft, DEC

ACM Turing Award 1998

ACM SIGMOD Edgar F. Codd Innovations award 1993



memory hierarchy (by Jim Gray)

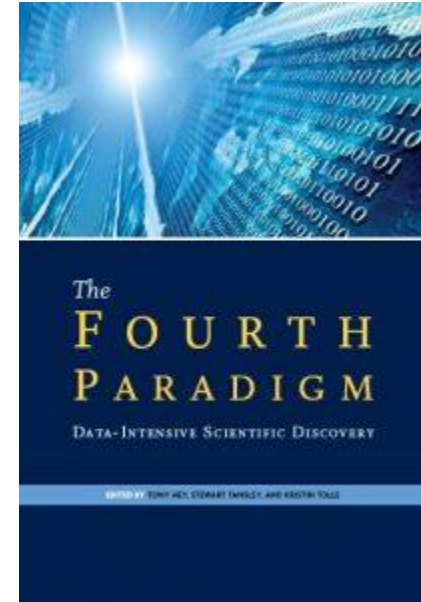
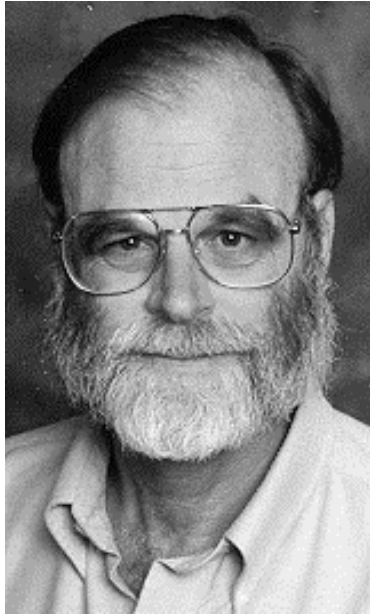


45TB @ \$150

tape?
sequential-only magnetic storage
still a multi-billion industry



Jim Gray (a great scientist and engineer)



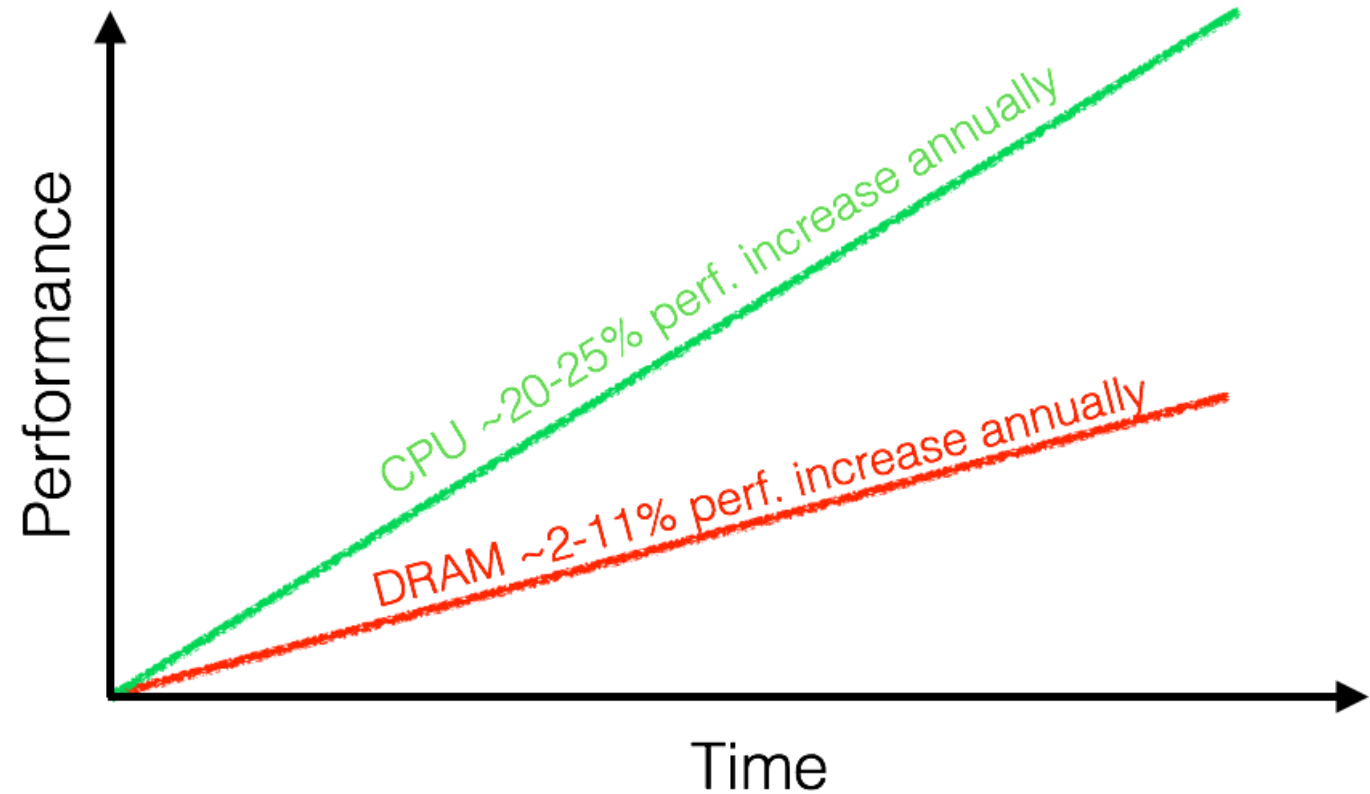
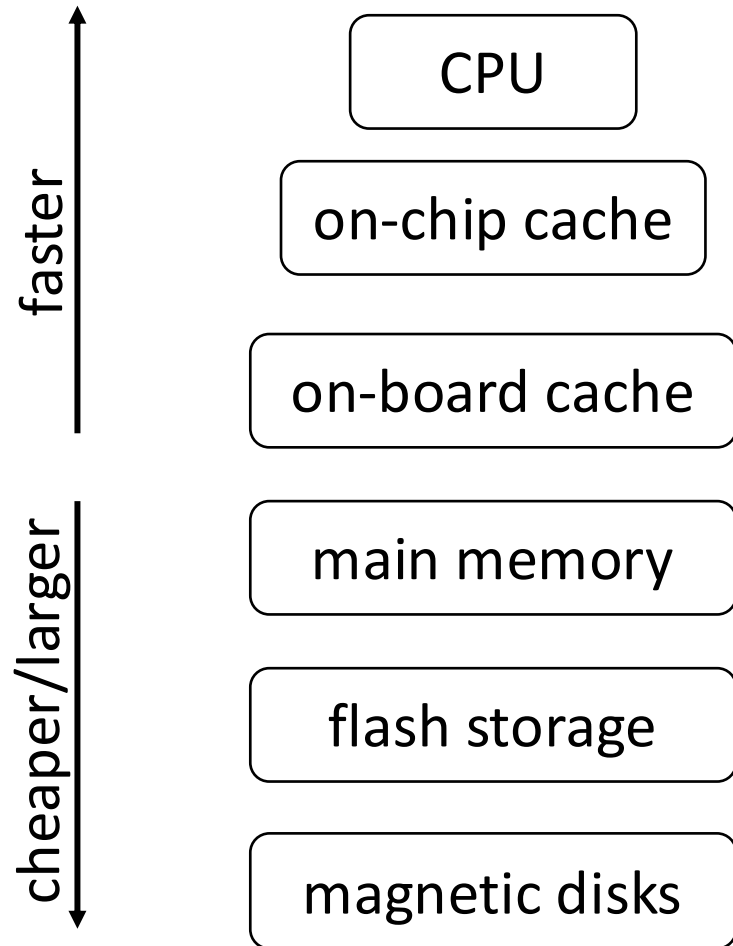
Jim Gray, IBM, Tandem, Microsoft, DEC

ACM Turing Award 1998

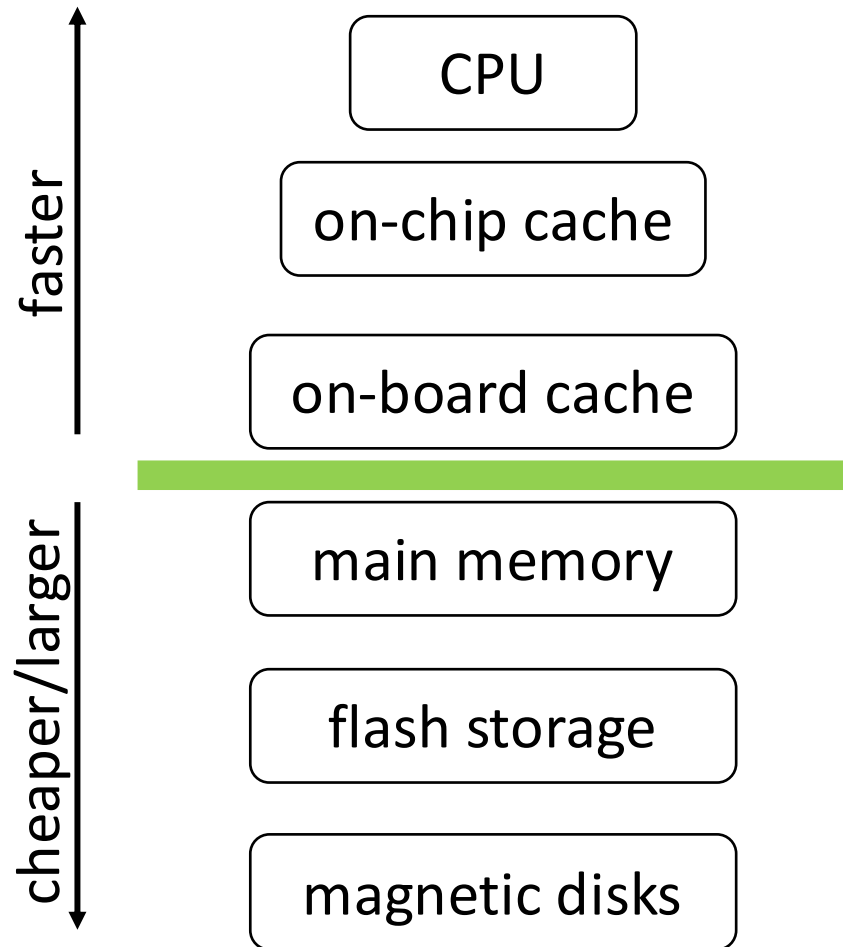
ACM SIGMOD Edgar F. Codd Innovations award 1993

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

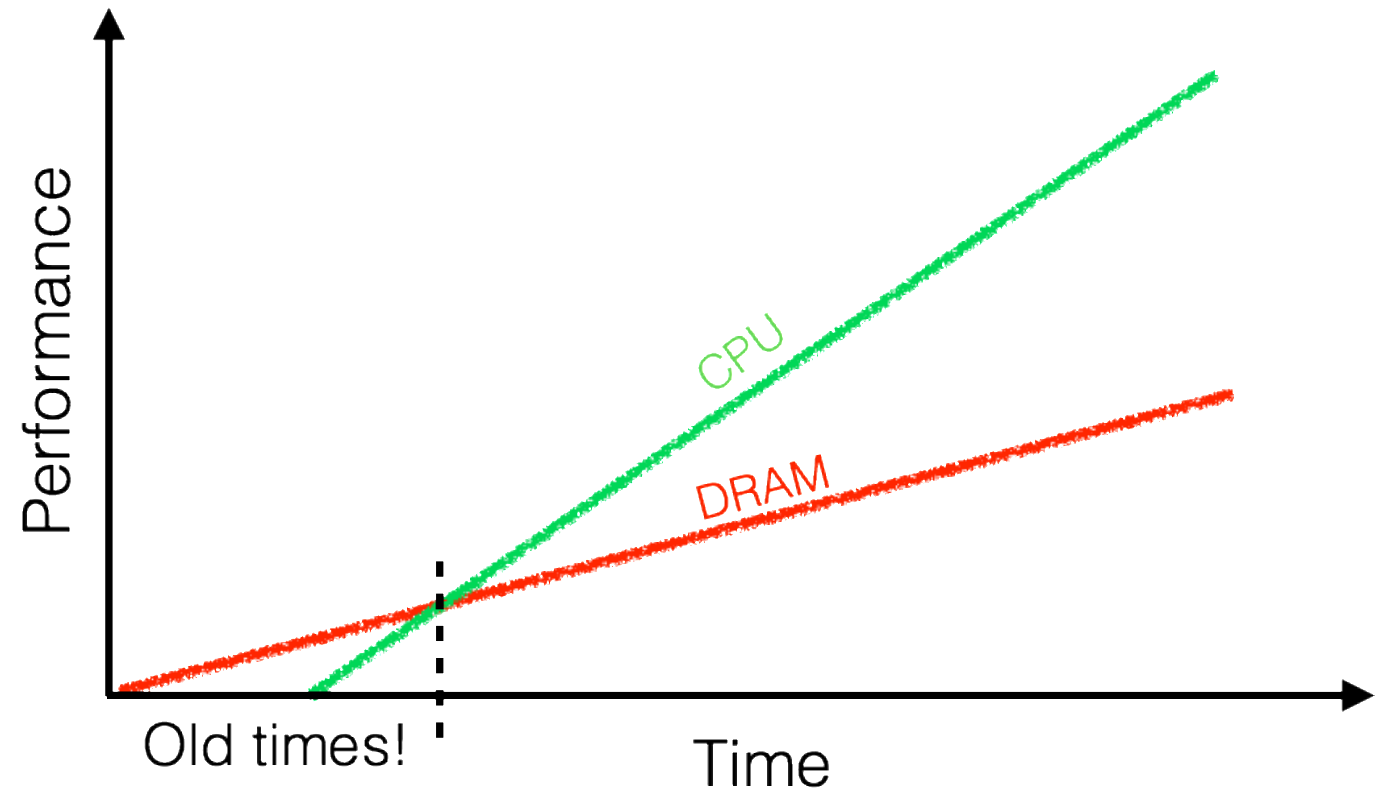
memory wall



memory wall

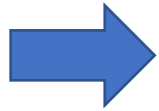


be careful when you go below the green line



cache/memory misses

*computations
happen here*



CPU

on-chip cache

on-board cache

main memory

flash storage

magnetic disks

be careful when you go below the green line

cache miss: looking
for something that
is not in the cache

what happens if I miss?

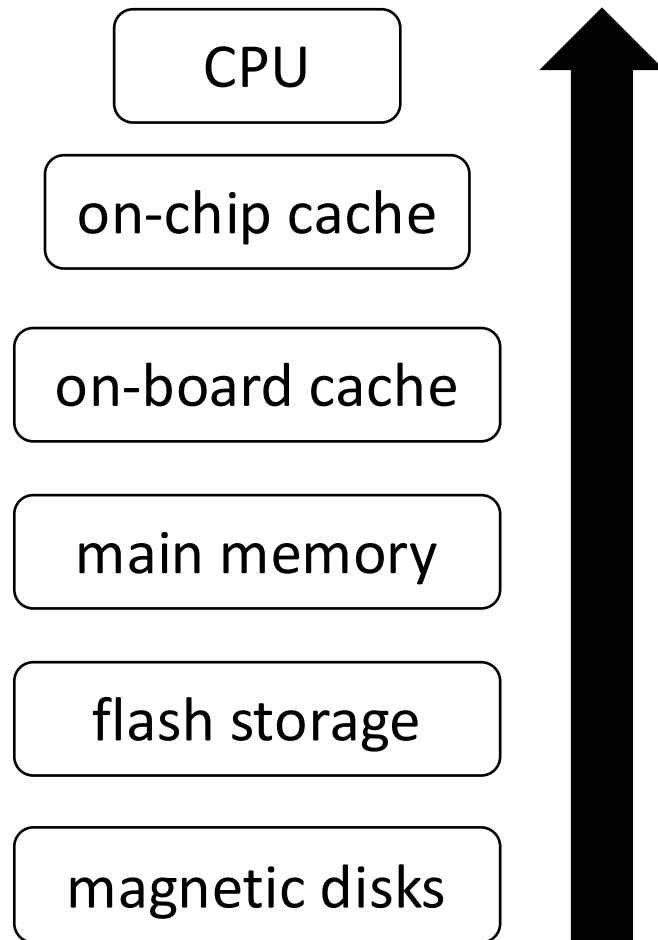
memory miss: looking
for something that
is not in memory

what happens if I miss again?

be very careful when you go below the green line



data movement



data goes through
all necessary levels

also read
unnecessary data

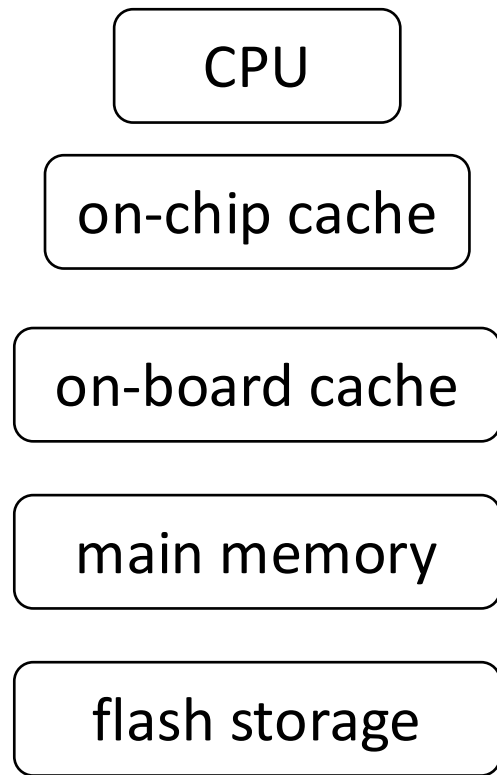


Photo by Gary Dineen/NBAE via Getty Images

need to read only X
read the whole page



data movement



data goes 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) of 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

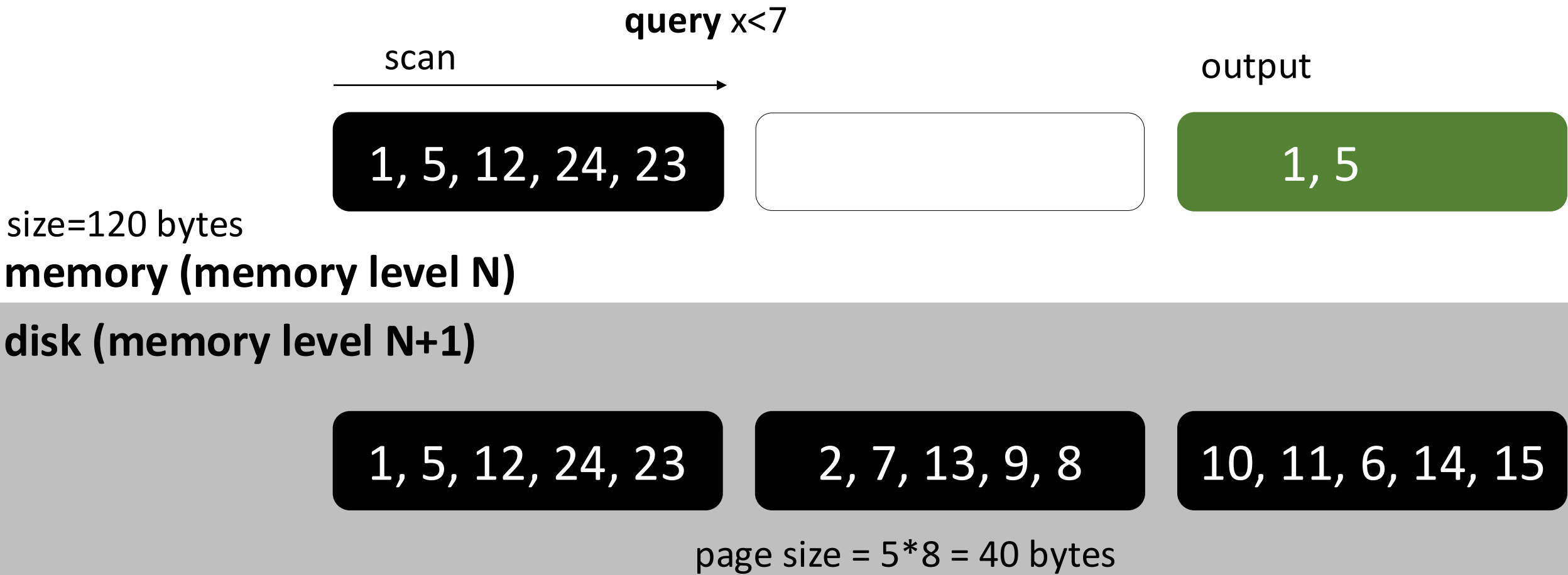
2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

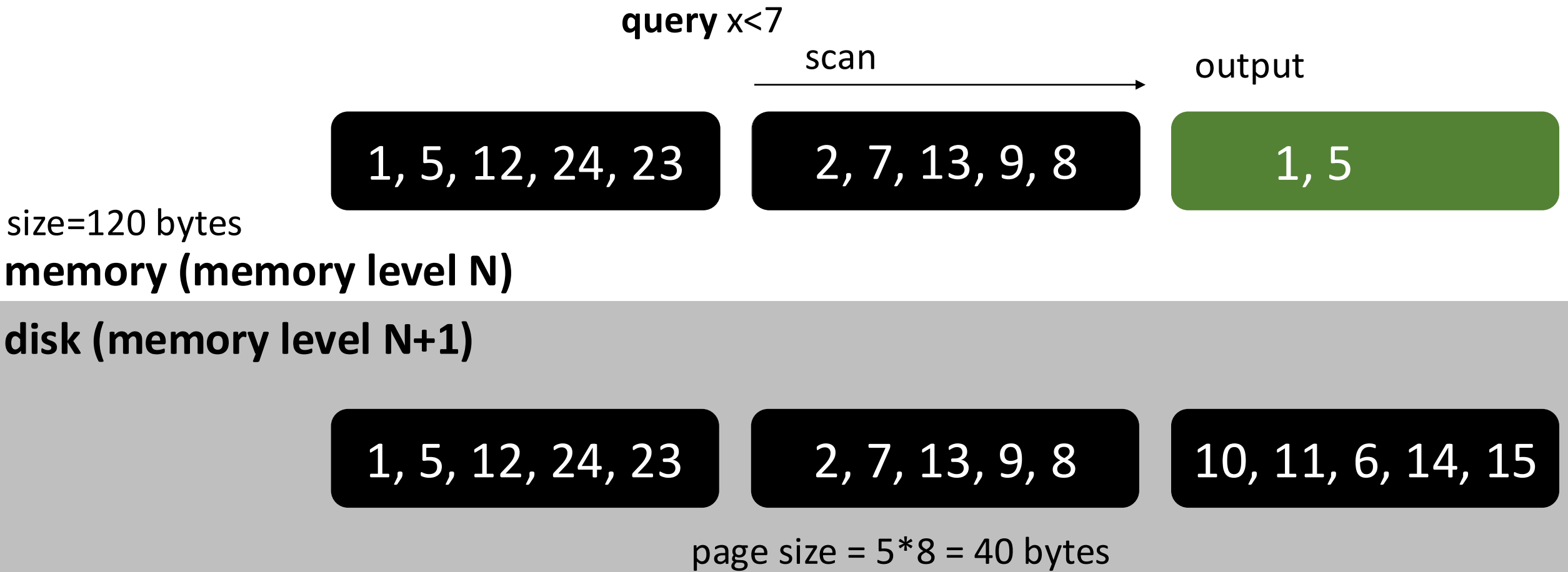
\$ 40 bytes

page-based access & random access



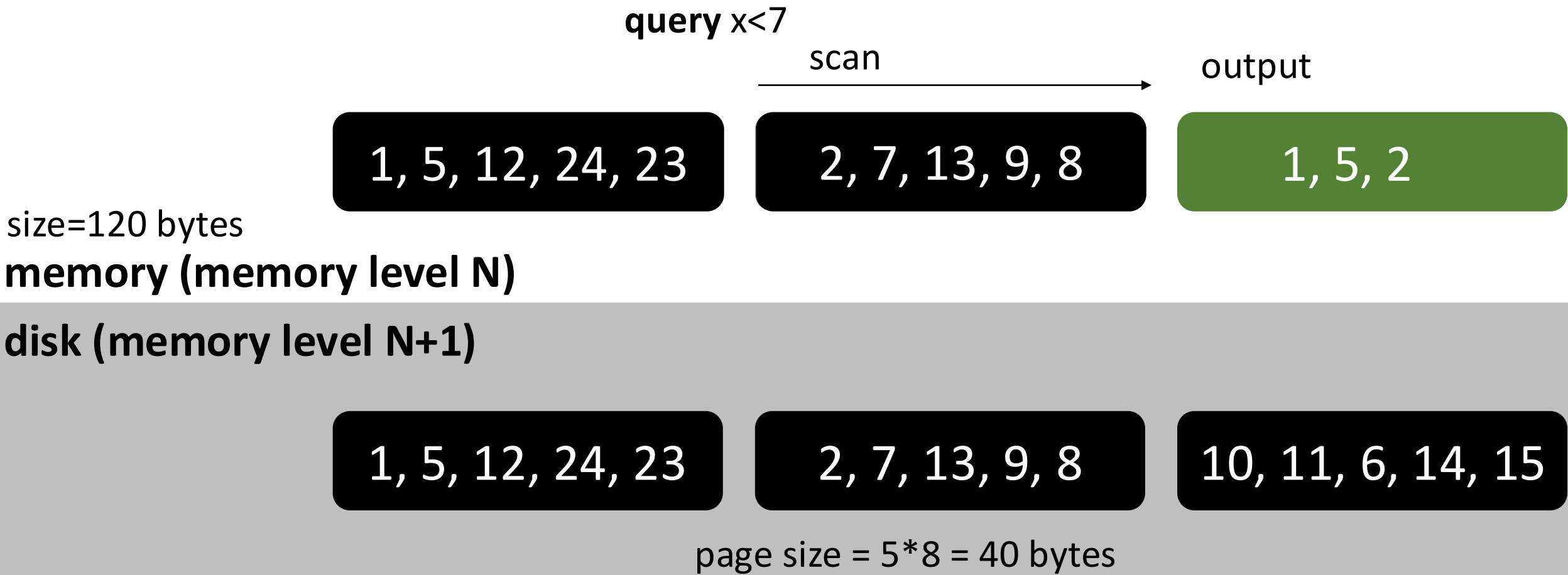
\$ 40 bytes

page-based access & random access



\$ 40 bytes

page-based access & random access



\$ 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

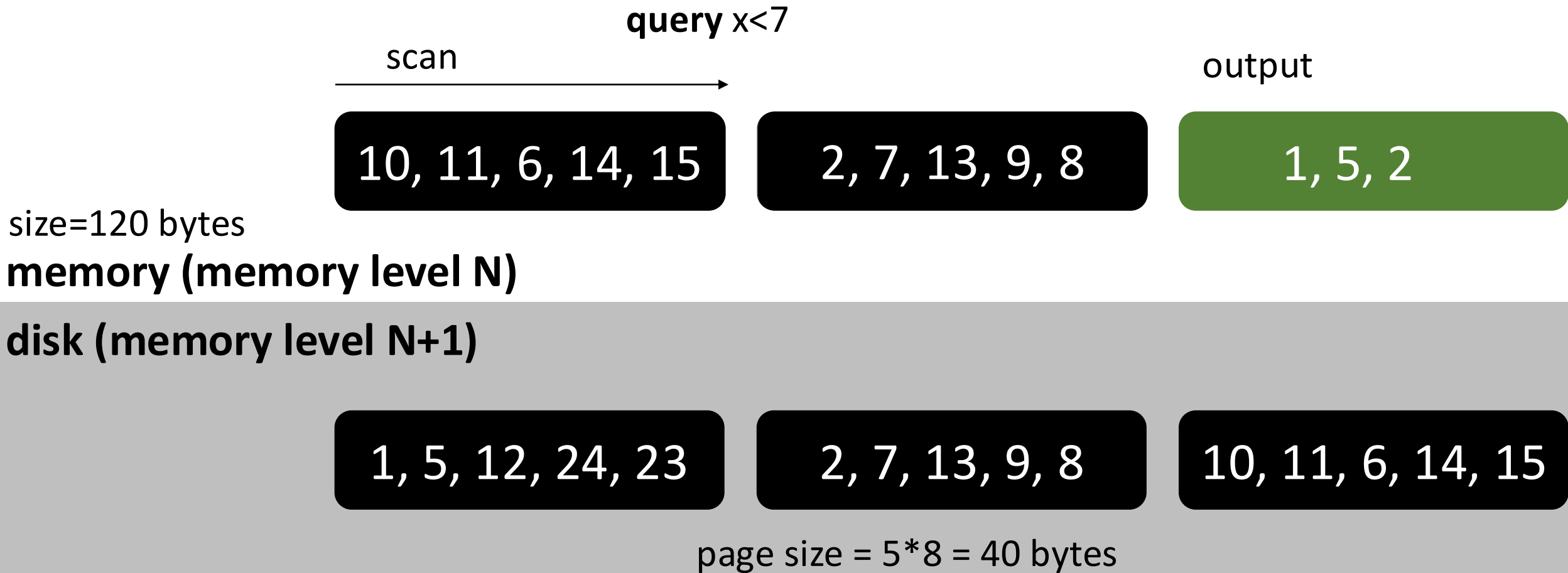
2, 7, 13, 9, 8

10, 11, 6, 14, 15

page size = $5 * 8 = 40$ bytes

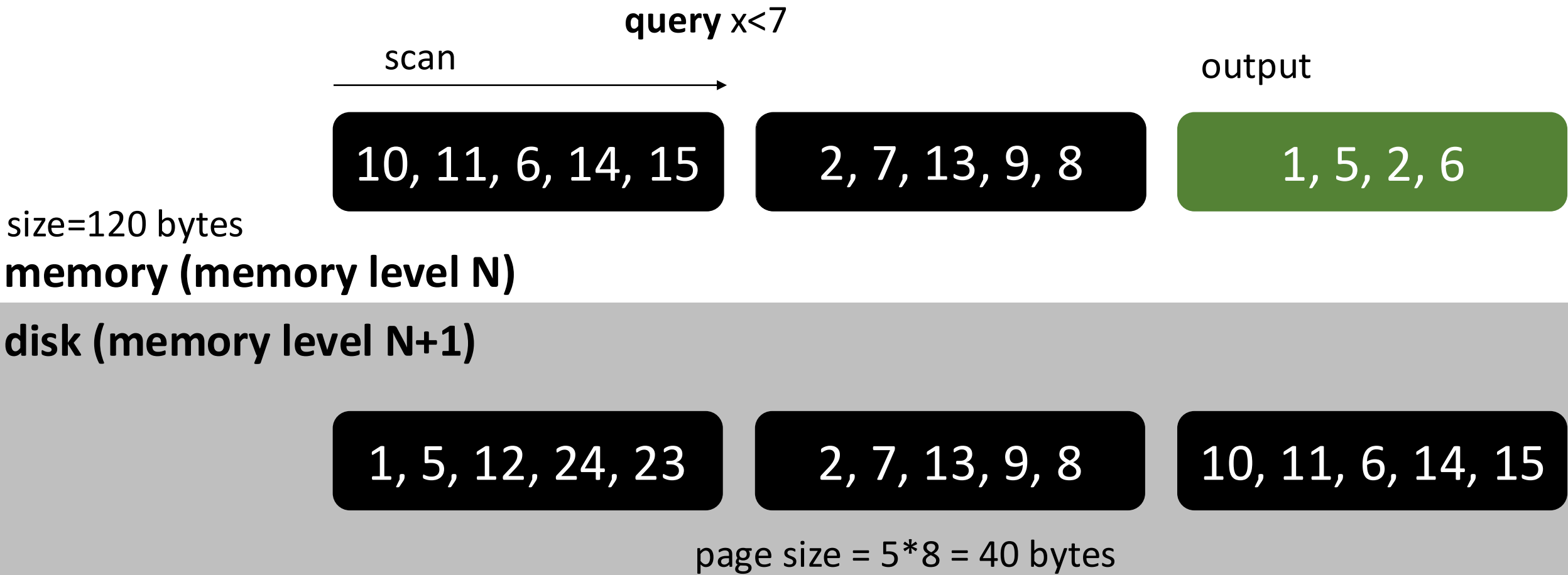
page-based access & random access

\$ 80 bytes



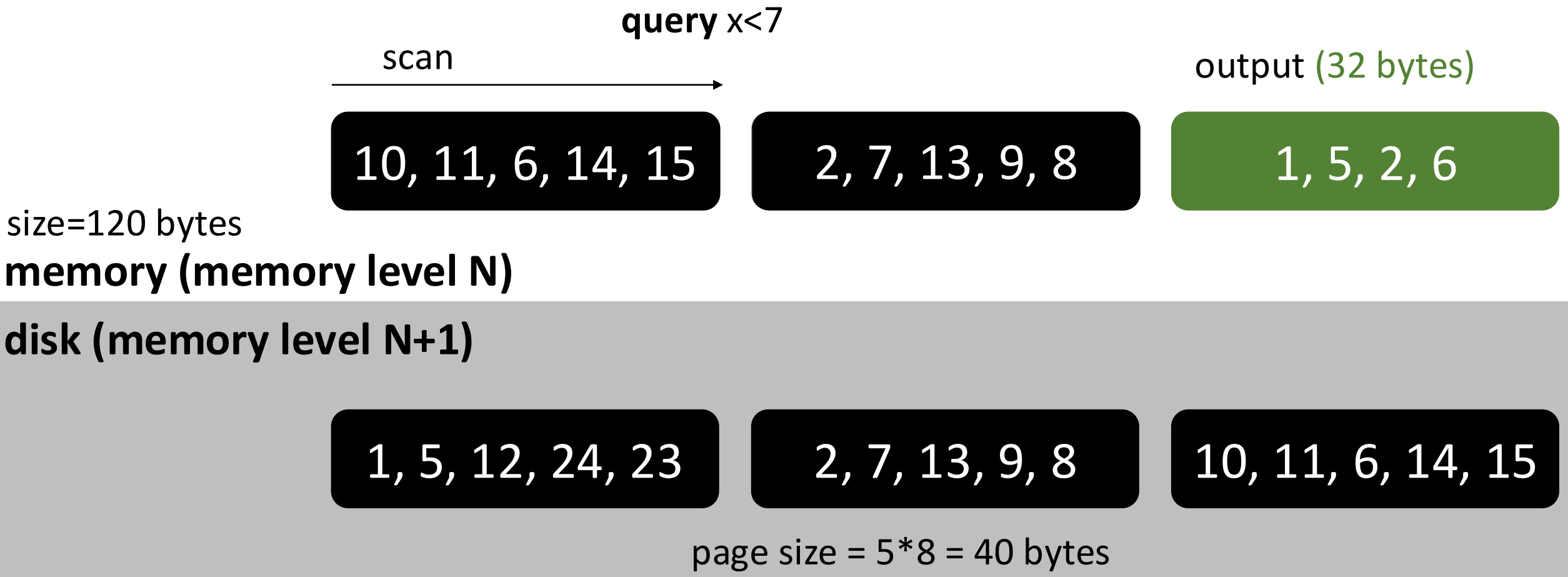
\$ 80 bytes

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

10, 11, 6, 14, 15

2, 7, 13, 9, 8

output

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



was the oracle helpful?

query $x < 7$

oracle

output (32 bytes)

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 or not?

1, 5, 12, 24, 23

2, 7, 13, 9, 8

10, 11, 6, 14, 15

every **byte** counts

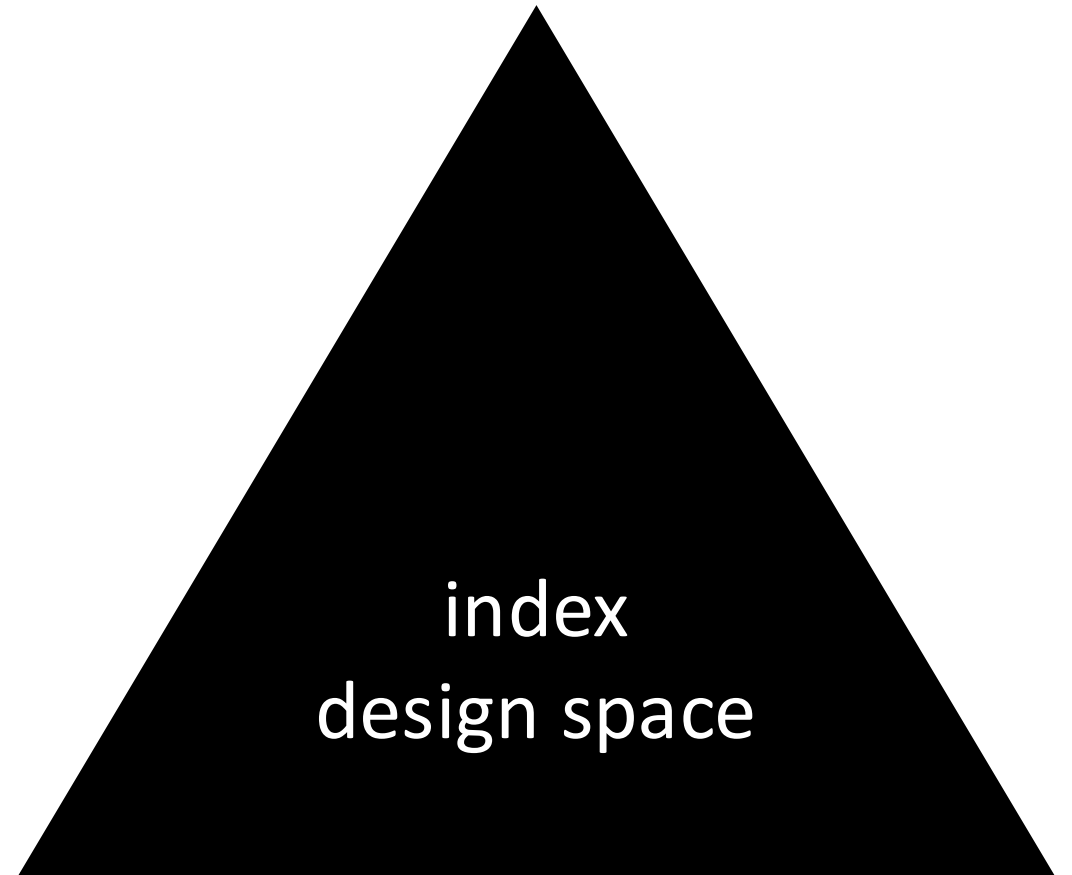
overheads and tradeoffs

how we store data

layouts, indexes

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)



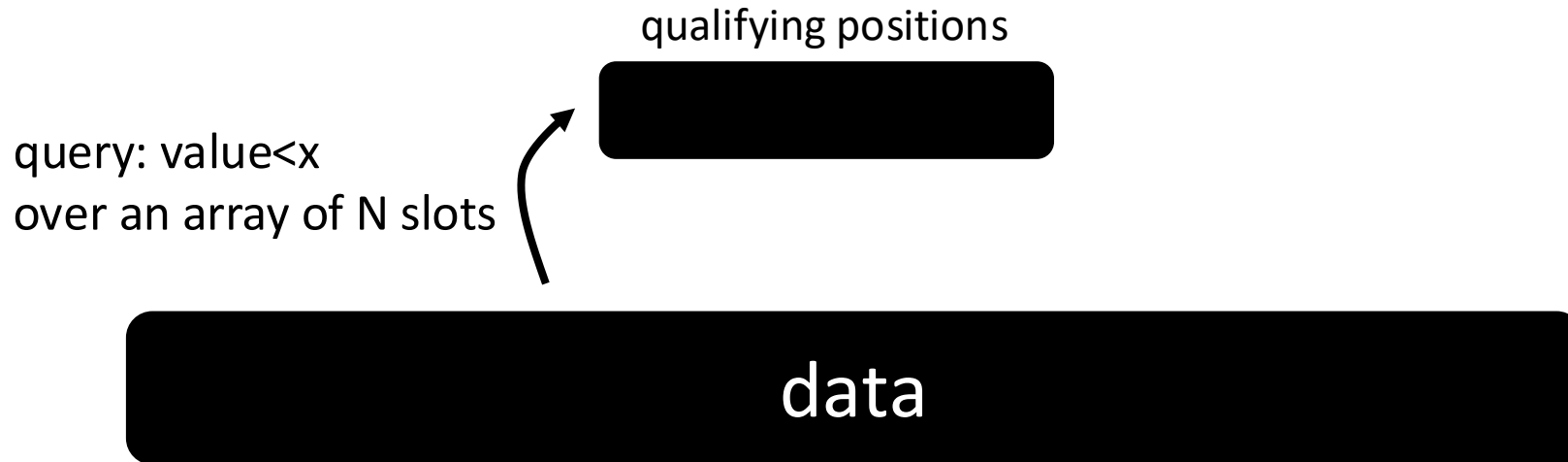
are random accesses always bad?

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

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

data

what if only 0.1% qualifies?

memory

data

result



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

data

what if only 0.1% qualifies?

memory

data



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?



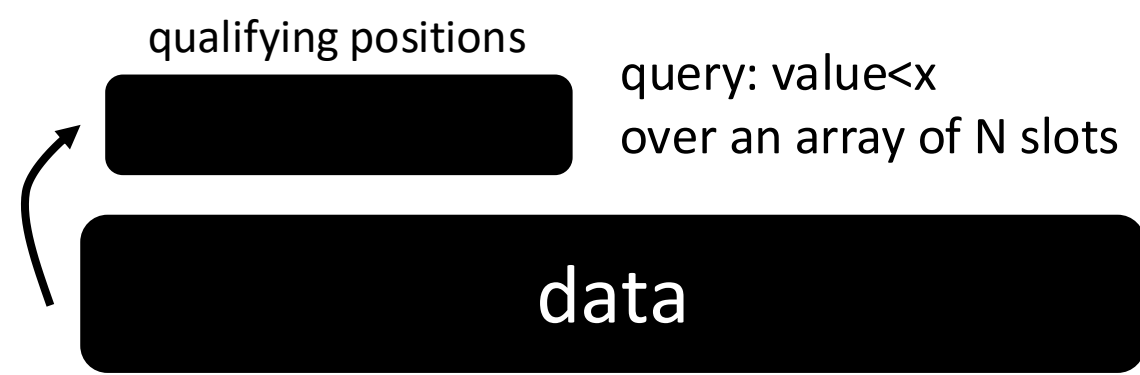
how can we know?

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?

data

core1

core2

core3

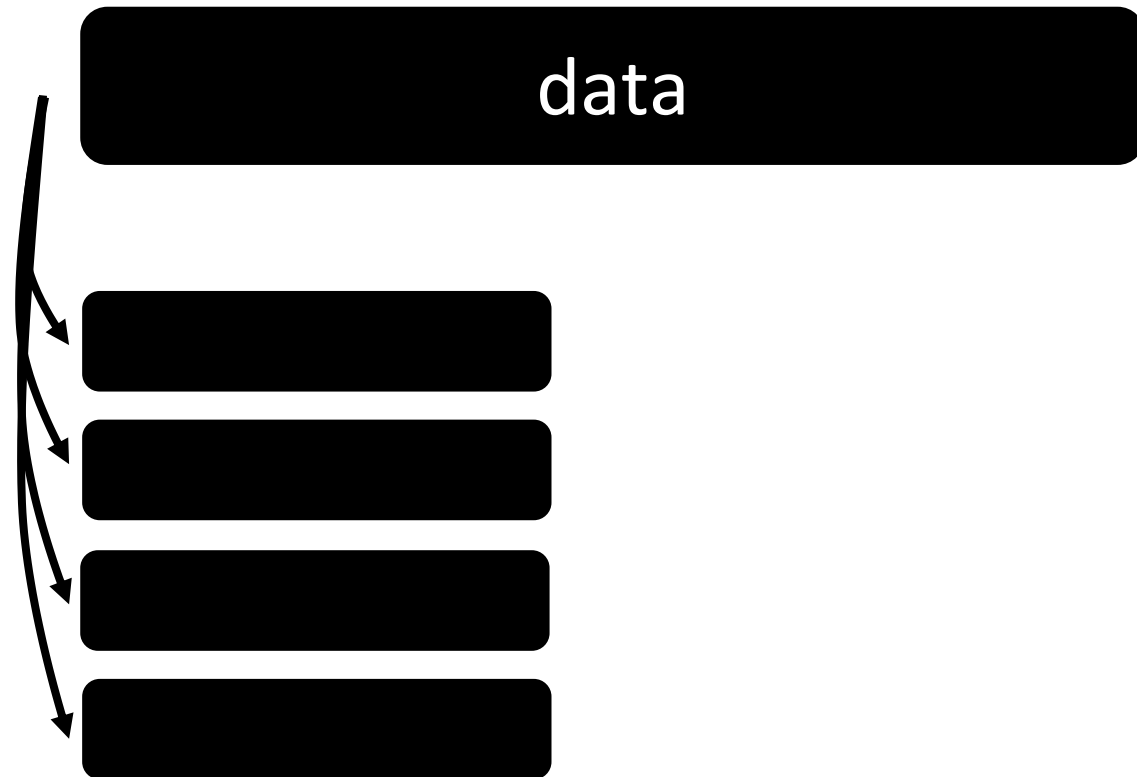
core4

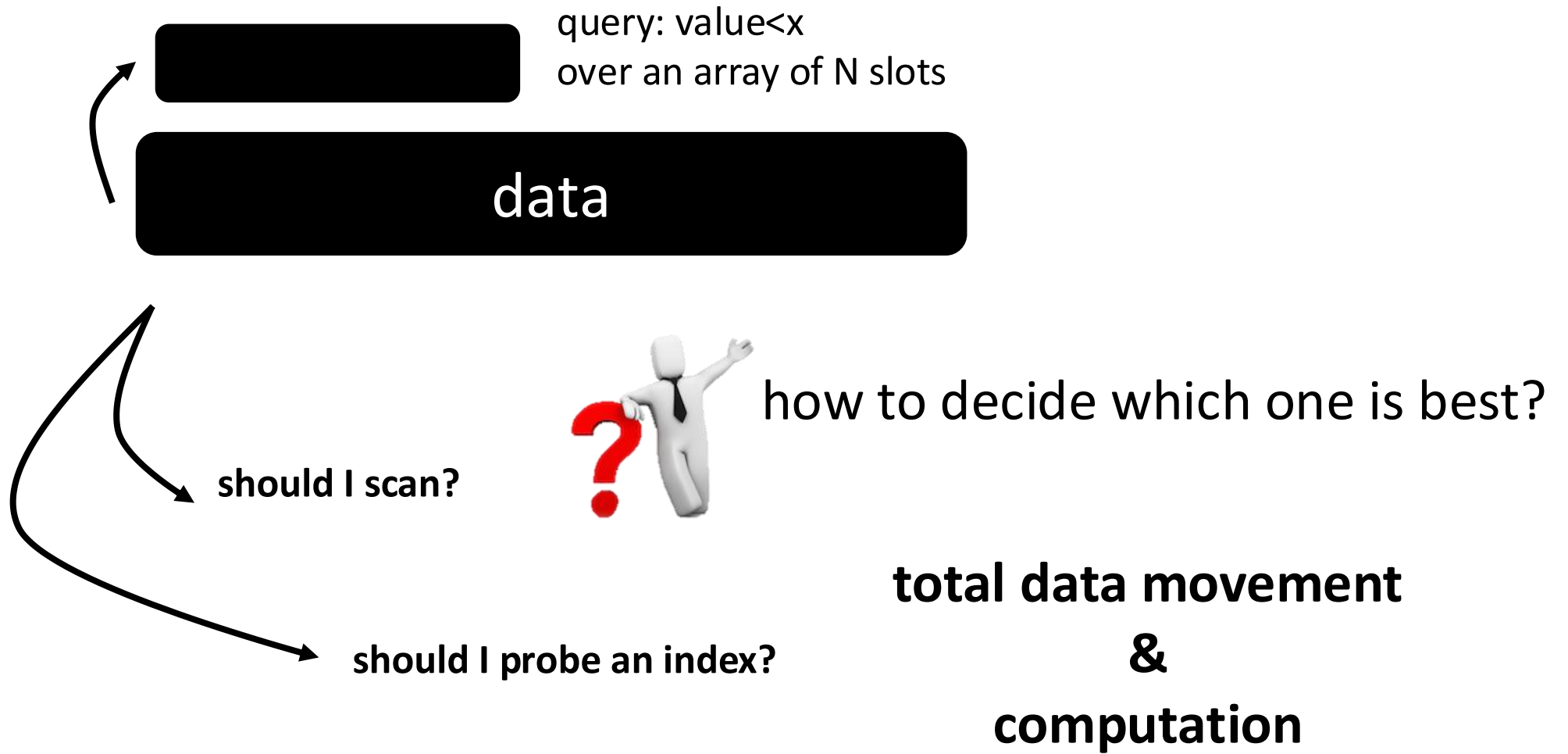


what about having multiple queries?

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;
```





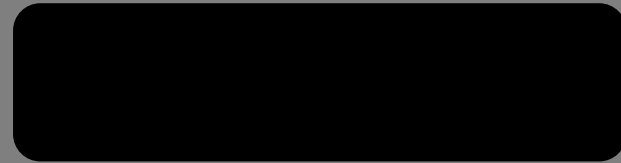
zonemaps

file = collection of pages

page 0



page 1



page 2



page 3



zonemaps

file = collection of pages

page 0

3, 16, 34, 31, 21

page 1

1, 5, 12, 24, 23

page 2

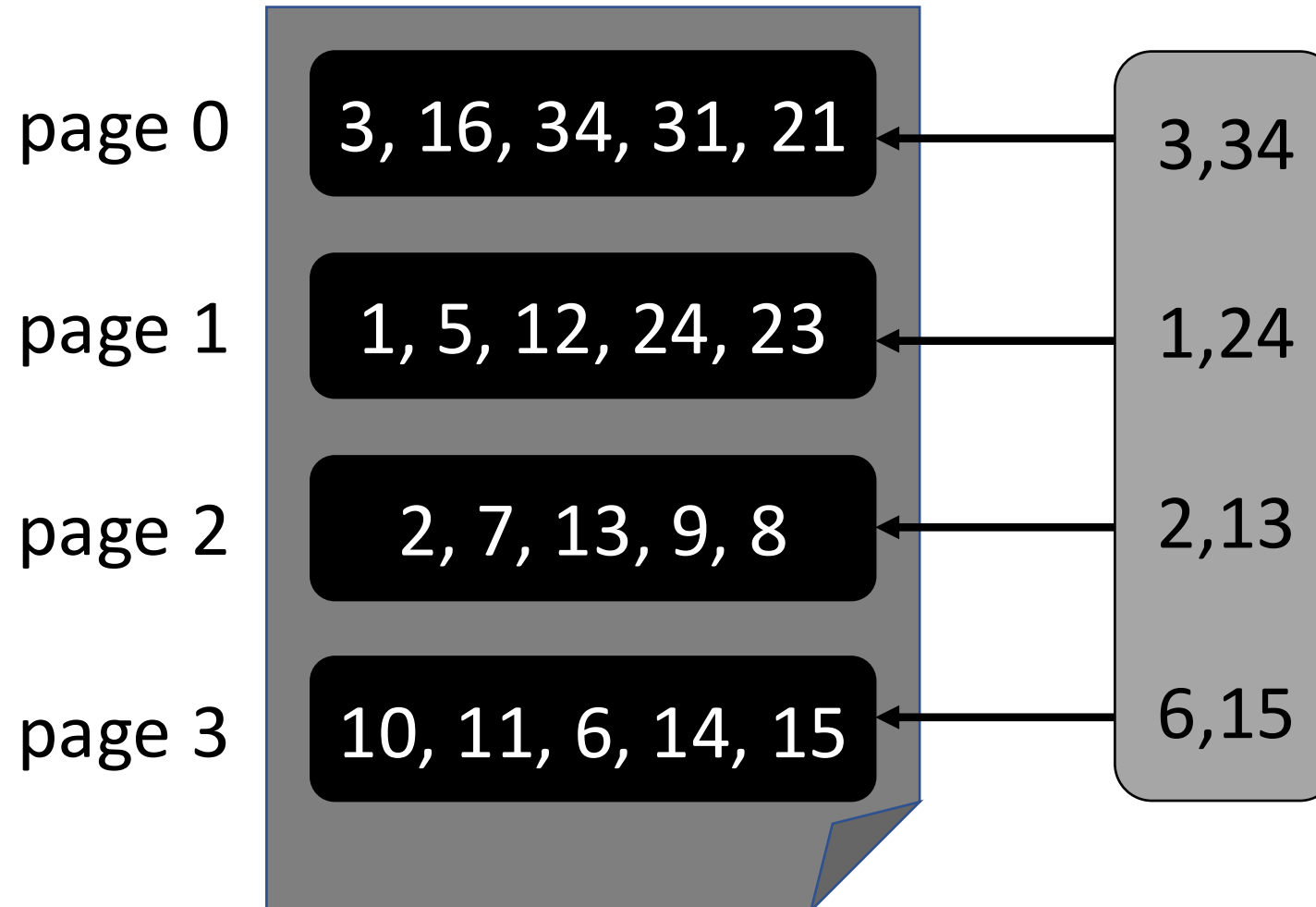
2, 7, 13, 9, 8

page 3

10, 11, 6, 14, 15

zonemaps

file = collection of pages



light-weight

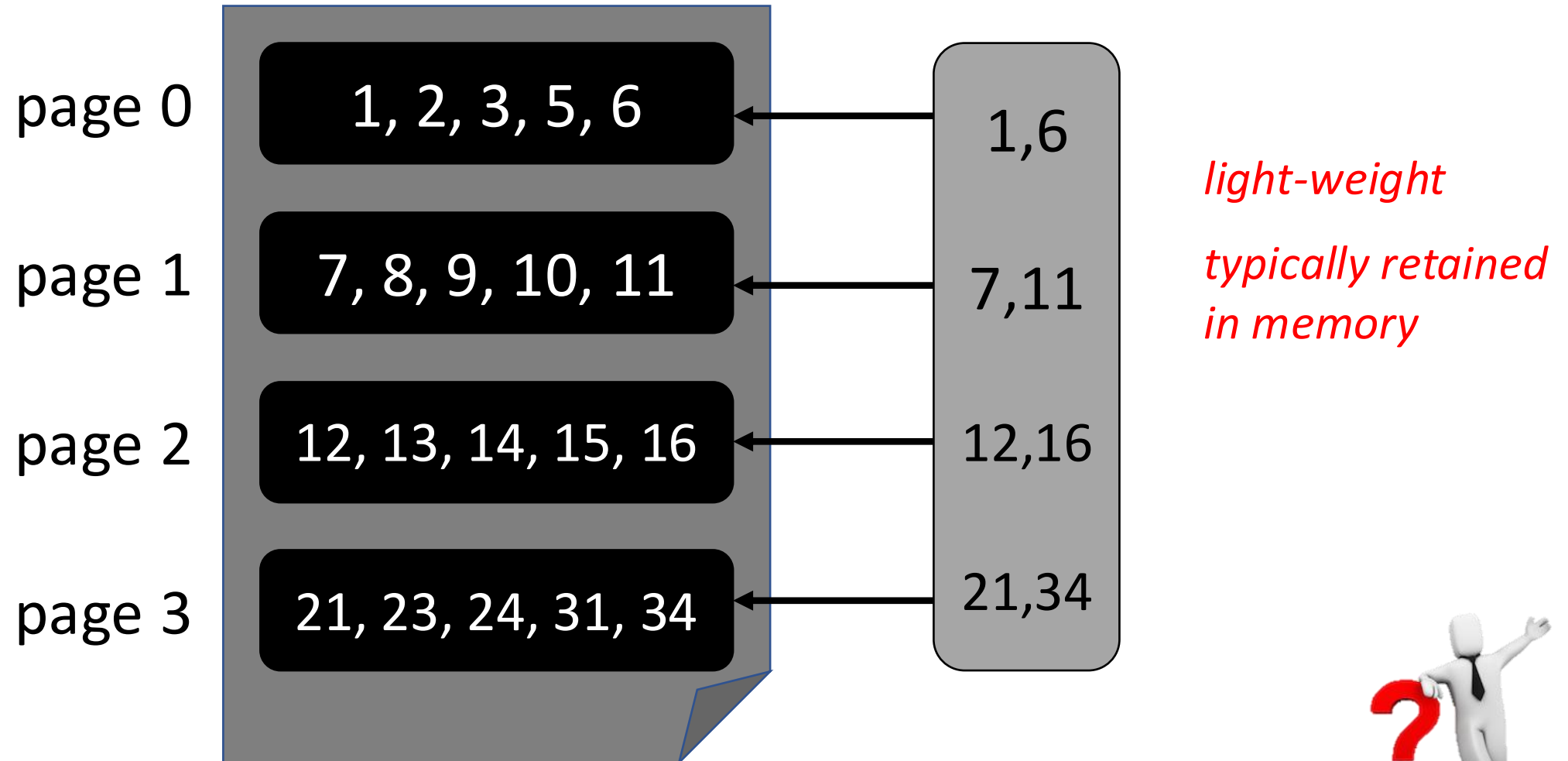
*typically retained
in memory*

But what if the data is sorted?



zonemaps

file = collection of pages



But what if the data is sorted?



the language of efficient systems: C/C++

why?

fewer assumptions

low-level control over hardware

make decisions about physical data placement and consumptions

the language of efficient systems: C/C++

why?

fewer assumptions

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

CS 561: Data Systems Architectures

class 2

Data Systems 101

next : modern main-memory data systems
&
semester project