

CS 265

Big Data & AI Systems

NoSQL | Neural Networks | Image AI | LLMs | Data Science

Scope: End-to-end AI systems

Topics: LLMs, Context, Agents, RAG

Inspiration: Research + Industry

Technical: Storage/Computation/Self-designing Projects:

Systems (LLM core, or design)

Research (LLM compiler, RAG, Image,

Fine-tuning, Context Management)

Research is open to 165 & systems students but eventually open to all

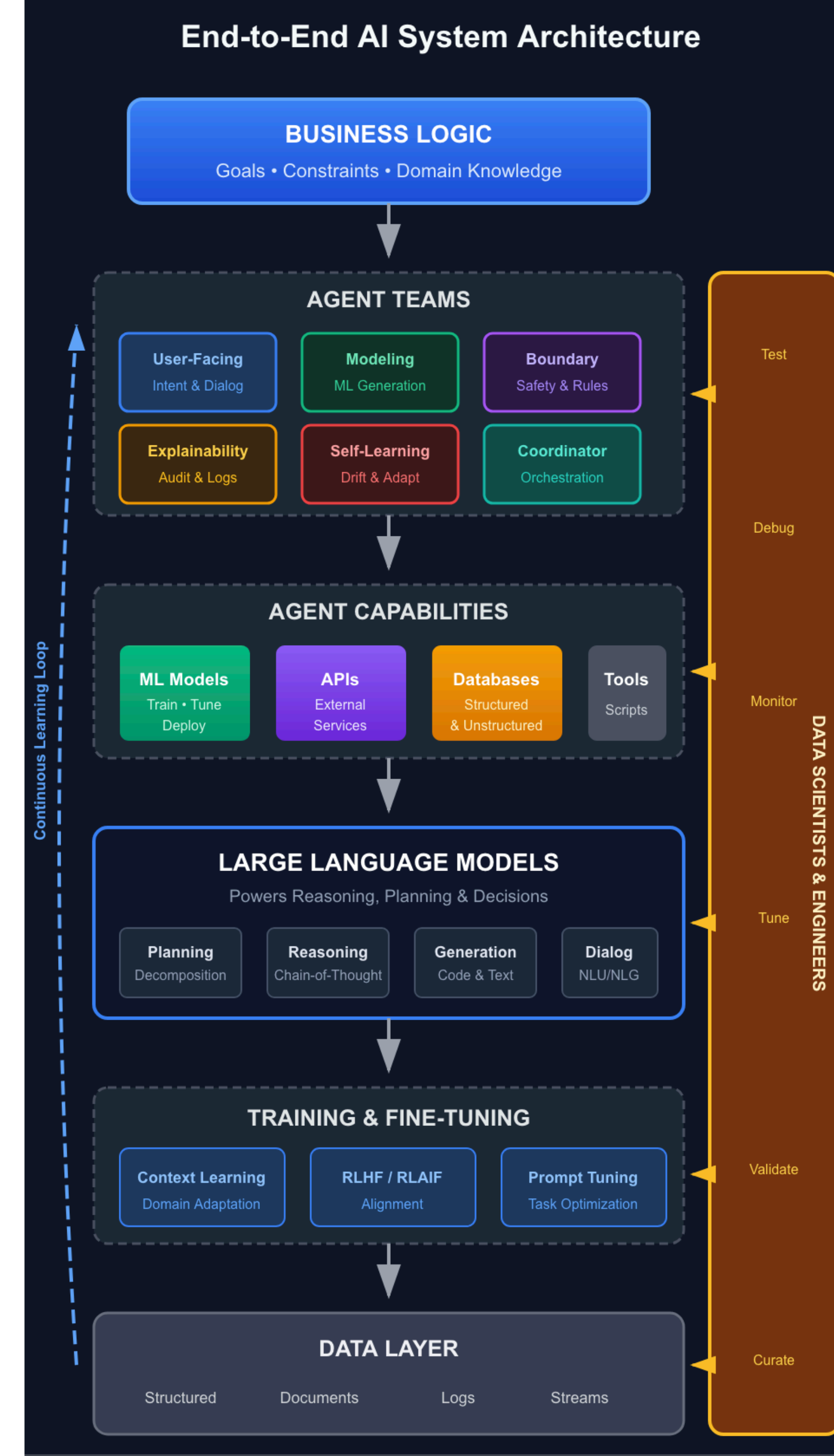
Timeline:

5 weeks of introduction

then reading research papers

Goals: Develop to an “AI systems person”

Info: <http://daslab.seas.harvard.edu/classes/cs265/>



A TYPICAL BIG DATA TASK

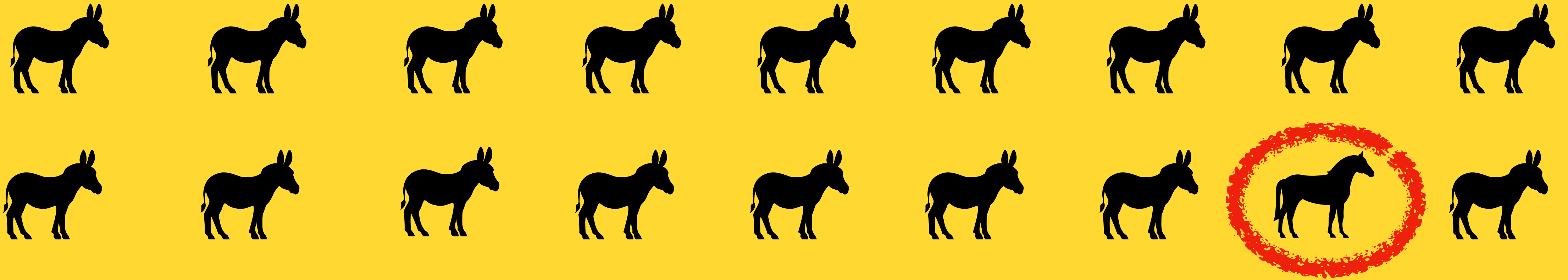
image analysis: e.g., detect the number of horses



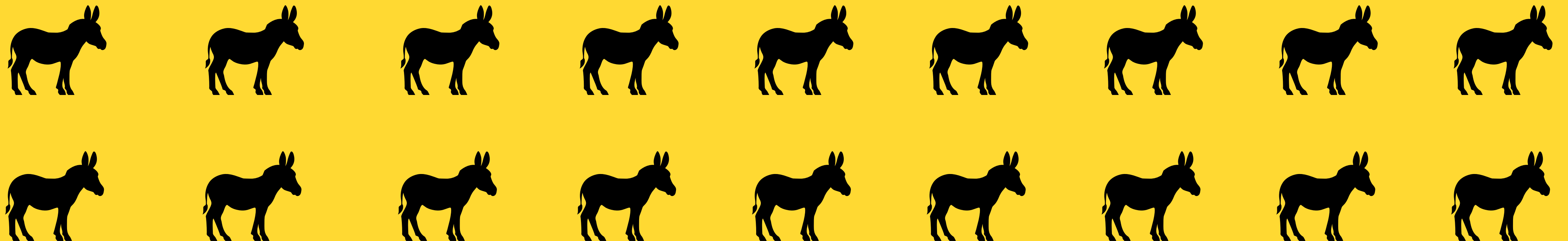
A TYPICAL BIG DATA TASK

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

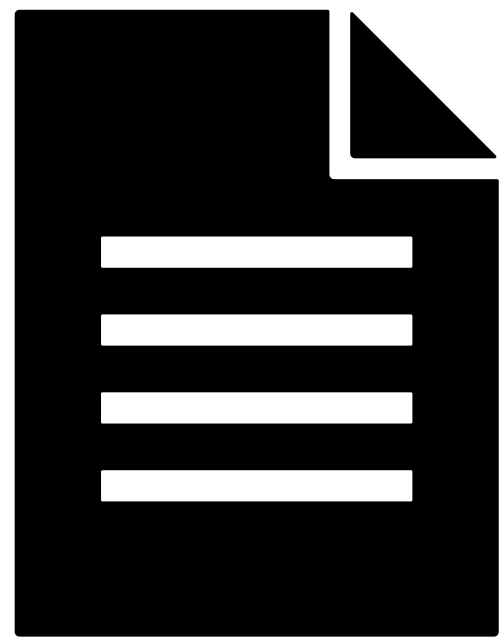




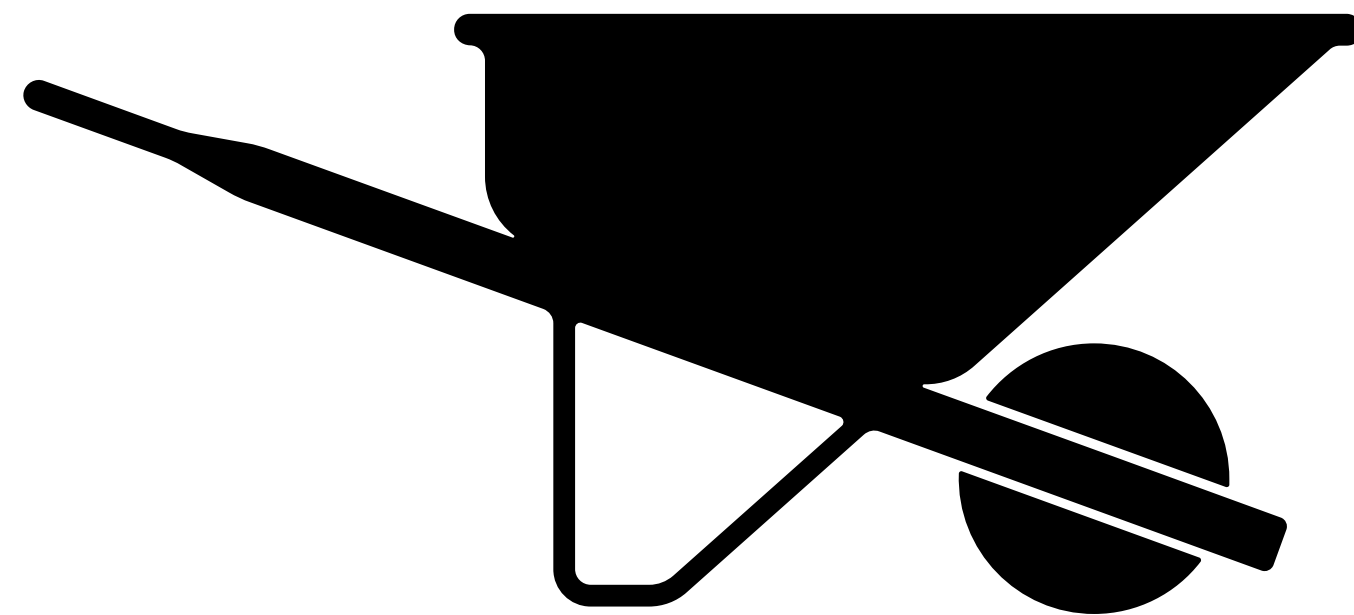
The core problem:
The size and organization of the data



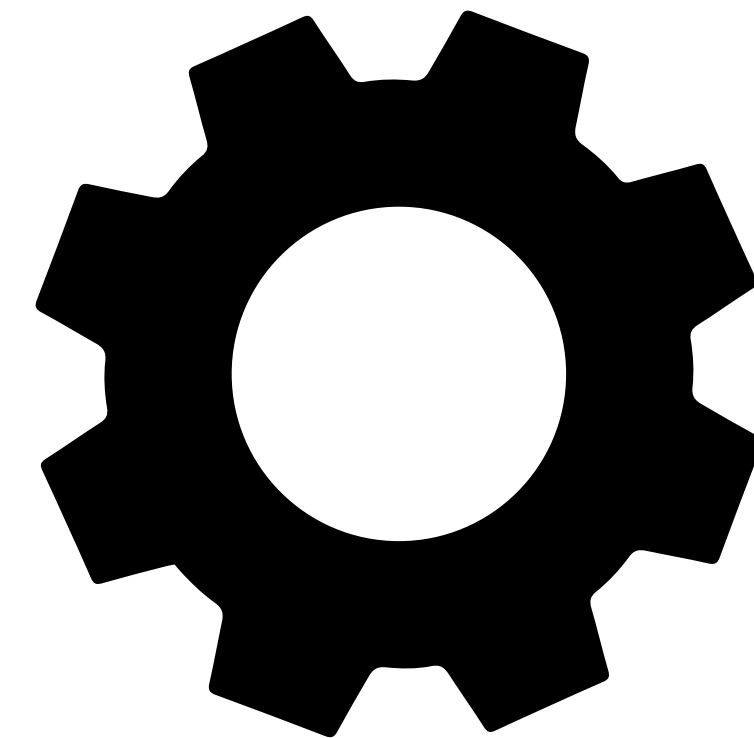
Three steps in big data **regardless of application**



STORE

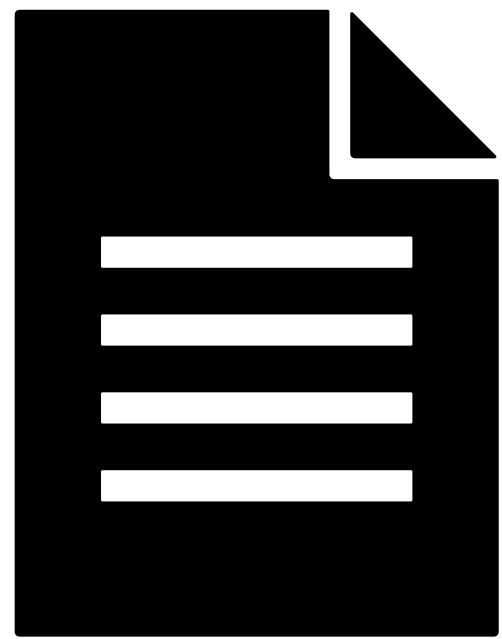


MOVE

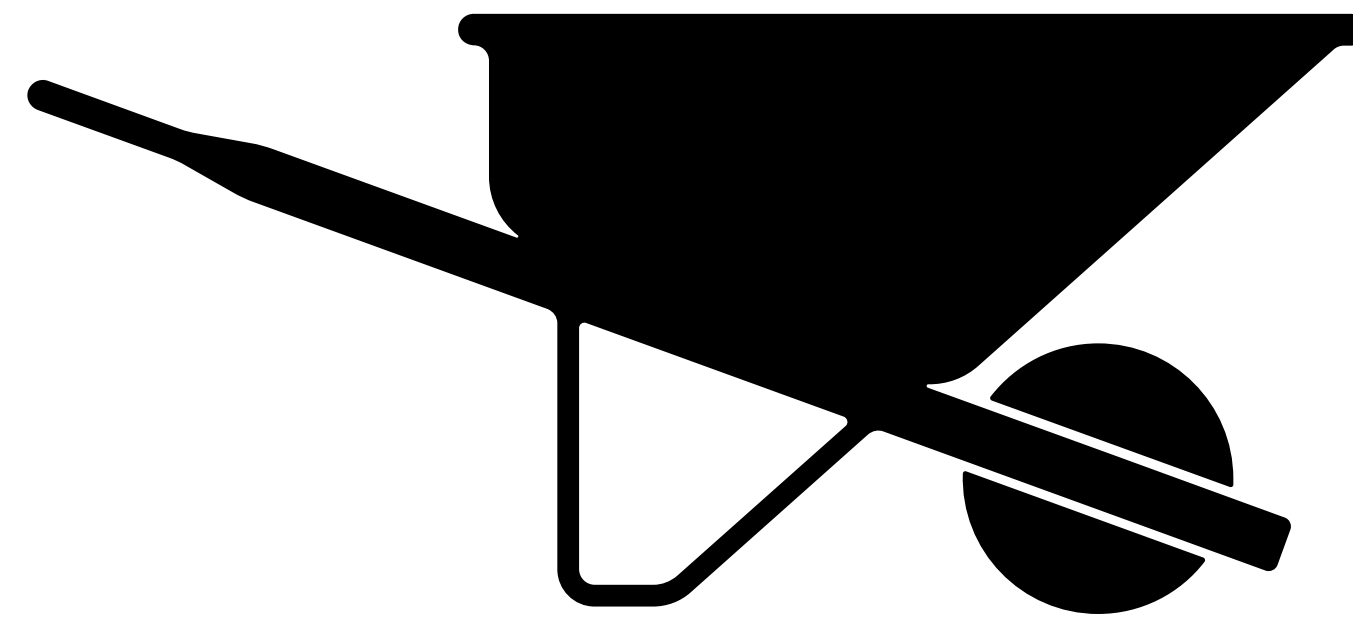


PROCESS

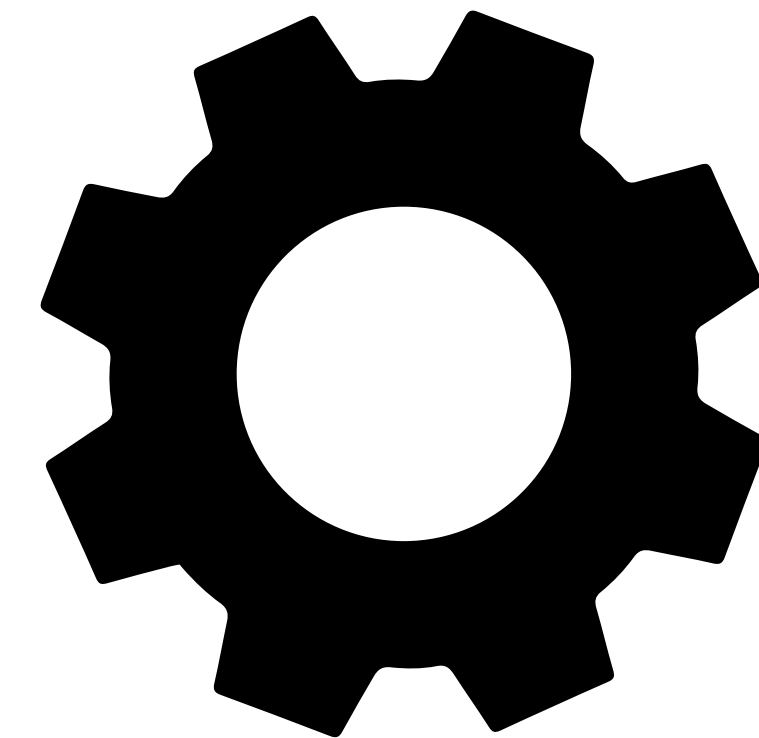
Three steps in big data **regardless of application**



STORE



MOVE



PROCESS



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

50-80% of end-to-end time is due to storage-related decisions

50-80% of end-to-end time is due to storage-related decisions



cloud cost

learning outcome

Fundamentals of storage

data structures, SQL, NoSQL, Agents, LLMs, RAG, Data Science, Image AI

learning outcome

Fundamentals of storage

data structures, SQL, NoSQL, Agents, LLMs, RAG, Data Science, Image AI
same set of principles across all fields (performance: design & implementation)

learning outcome

Fundamentals of storage

data structures, SQL, NoSQL, Agents, LLMs, RAG, Data Science, Image AI
same set of principles across all fields (performance: design & implementation)

from algorithms to systems

This class helps with:

software (systems) engineering jobs
joining data-driven startups
starting with research

This class helps with:

software (systems) engineering jobs
joining data-driven startups
starting with research

This class does not help with:

using systems

it helps with designing and building systems

First ~5 weeks:

Background on basic systems concepts (storage)
Intro into the concept of self-designing systems

Systems for: RAG, Image AI, LLMs

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

A system is a complex set of components

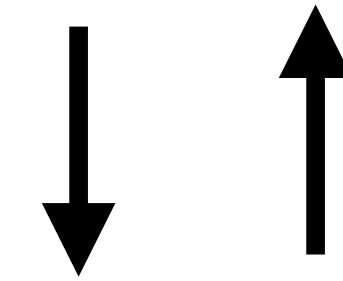
interacting in harmony depending on the context

exposing as little as possible complexity to users





declarative interface
ask “what” you want

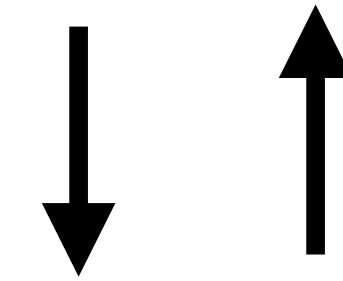


data* system

the system decides
“how” to best store
and access data



declarative interface
ask “what” you want



data* system

the system decides
“how” to best store
and access data



why is this good

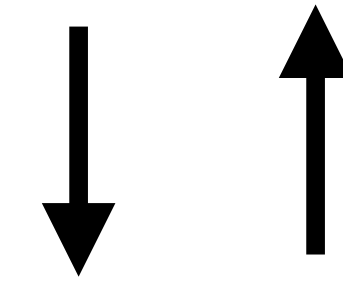


~6 decades of research

started with IBM, Microsoft, Oracle, Teradata, etc.
and a gazillion start-ups today

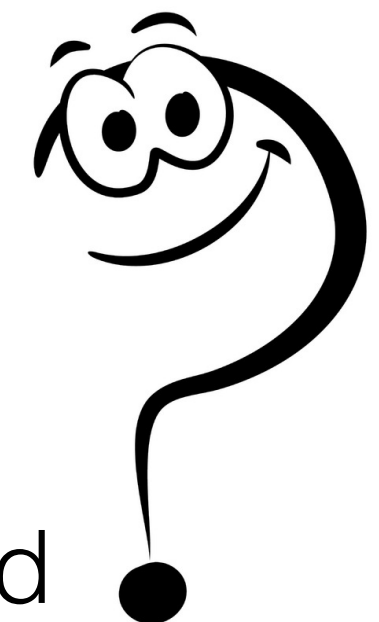
declarative interface

ask “what” you want



data* system

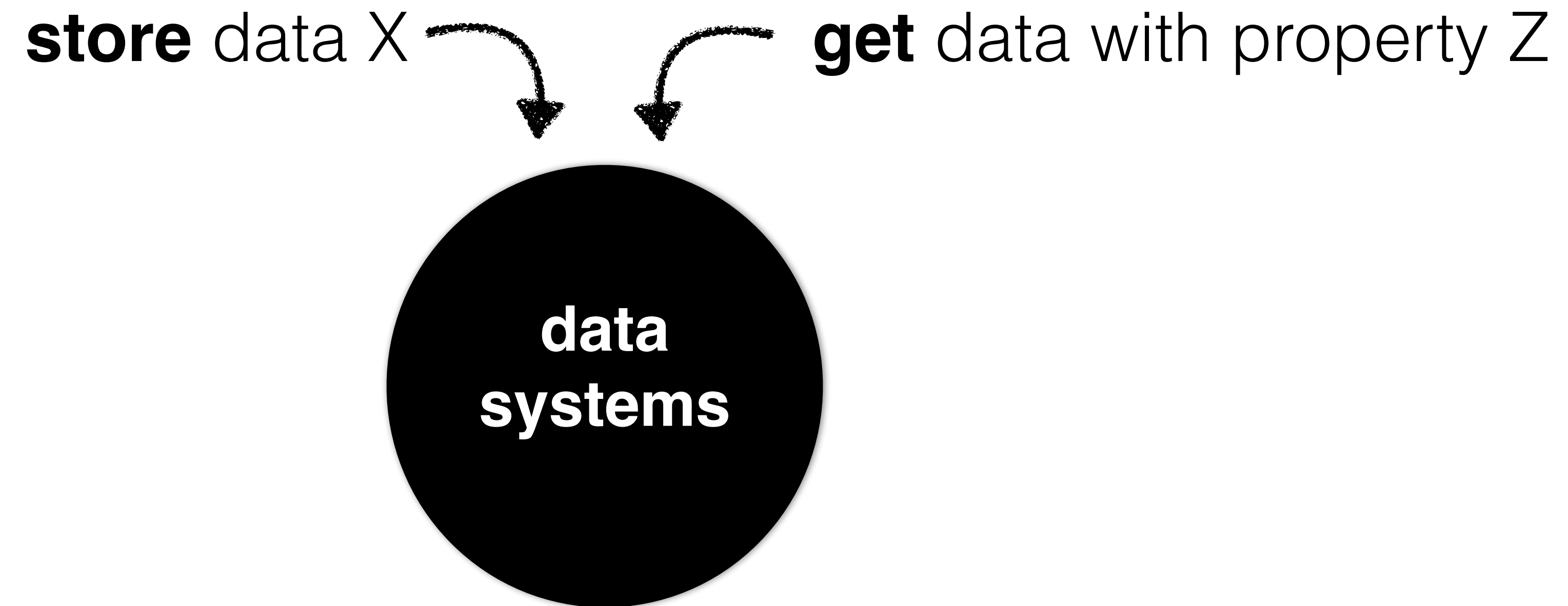
the system decides
“how” to best store
and access data



why is this good

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

big data apps

data
systems



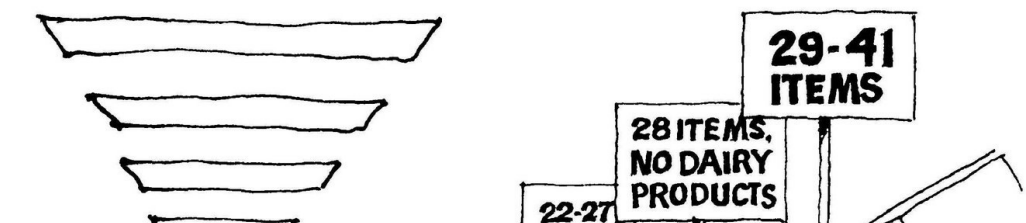
broader than SQL

big data apps

ANALYTICS

**data
systems**

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?

broader than SQL

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?

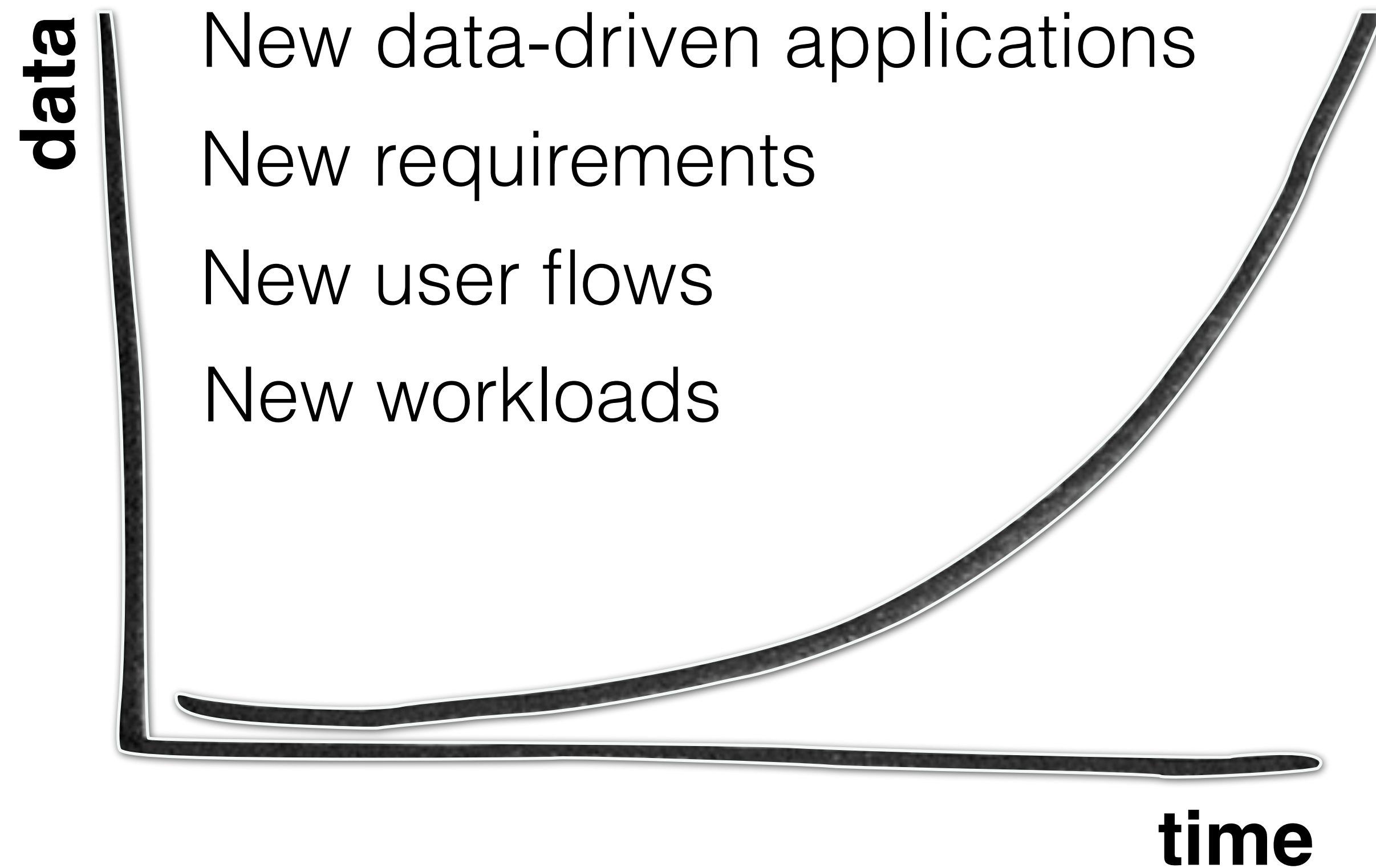
AI

Is this new review a legitimate one?

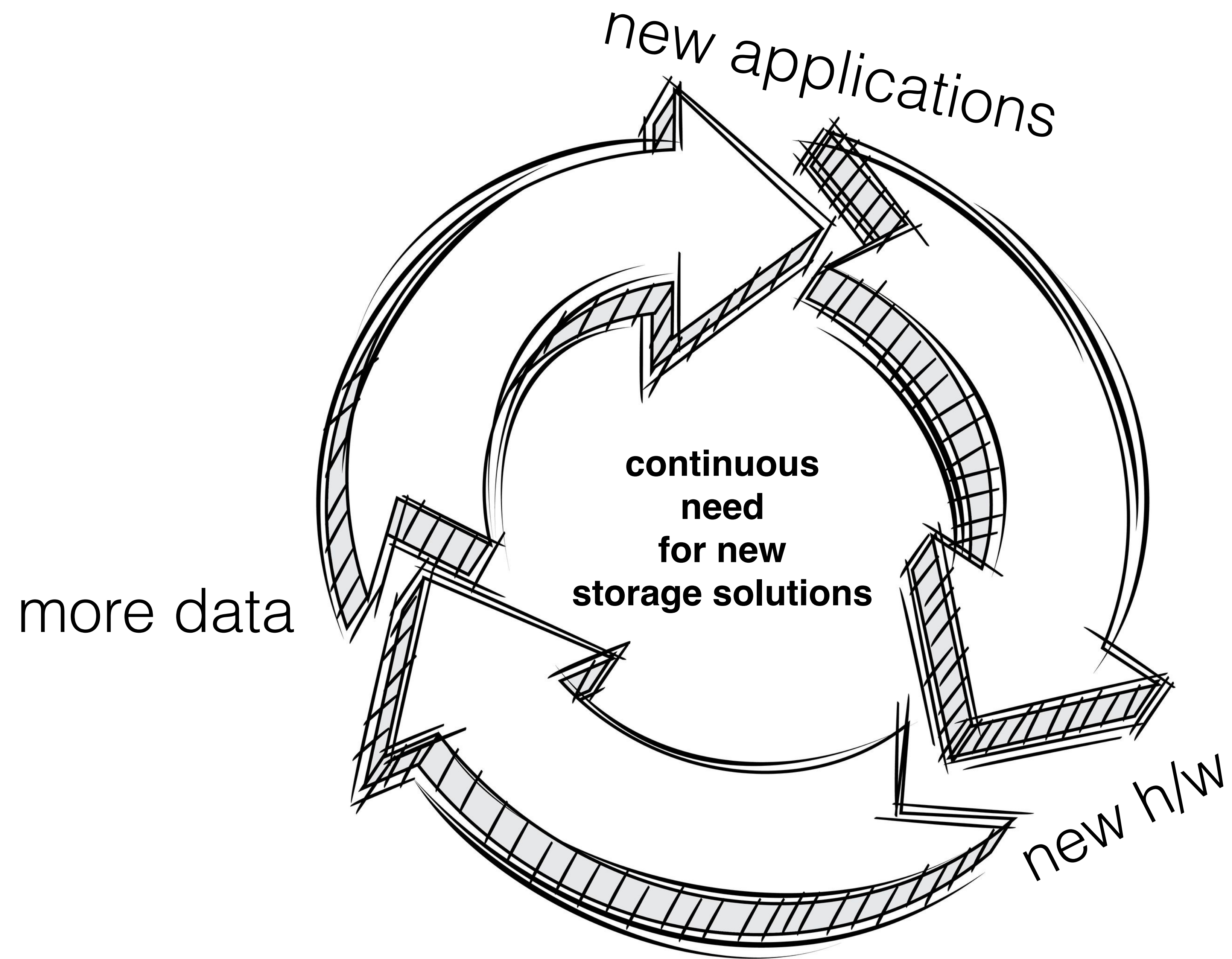
COMMUTING

Compute price for next Uber ride

broader than SQL



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



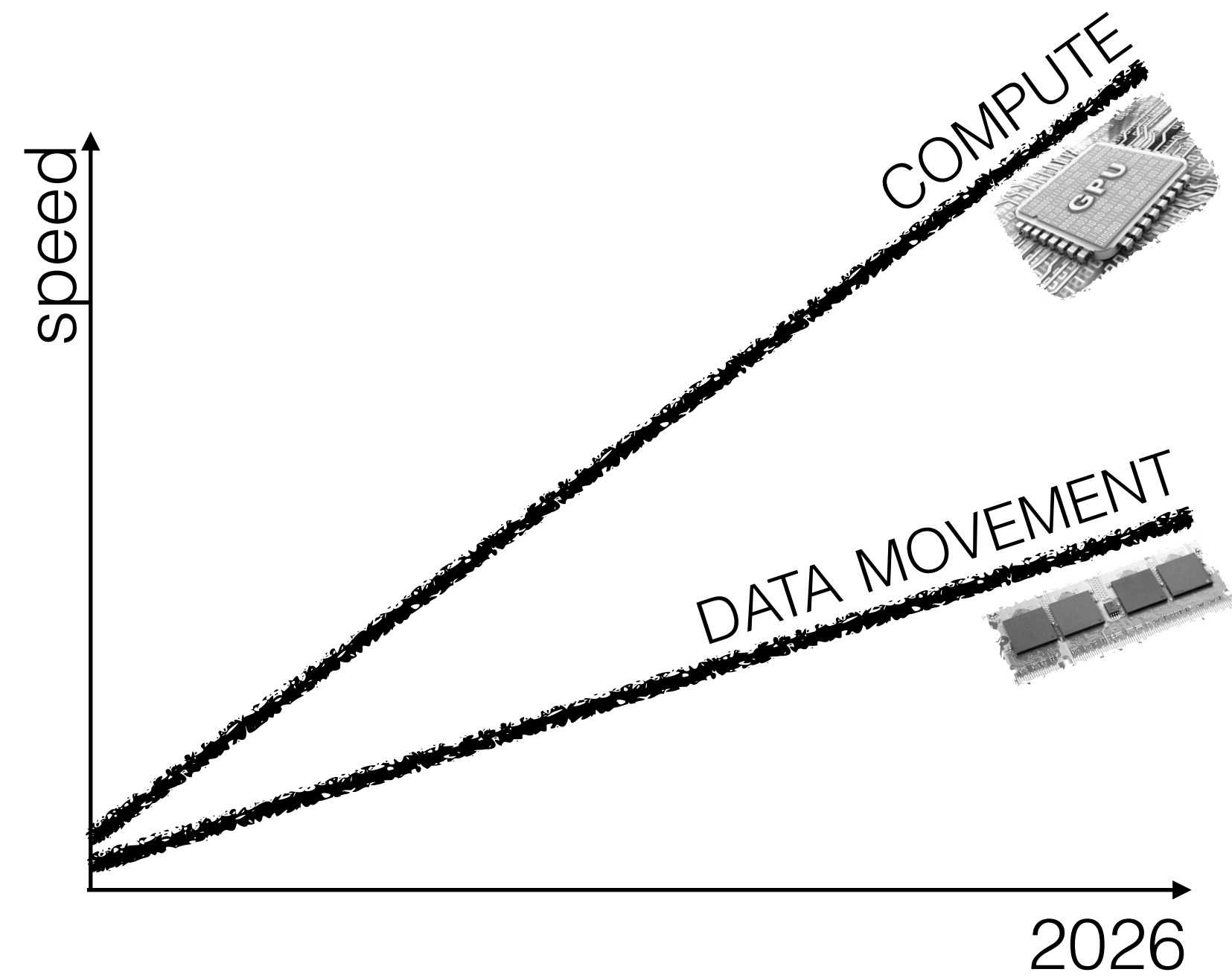
2. As data grows, having the right data system
for each application is increasingly more critical

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

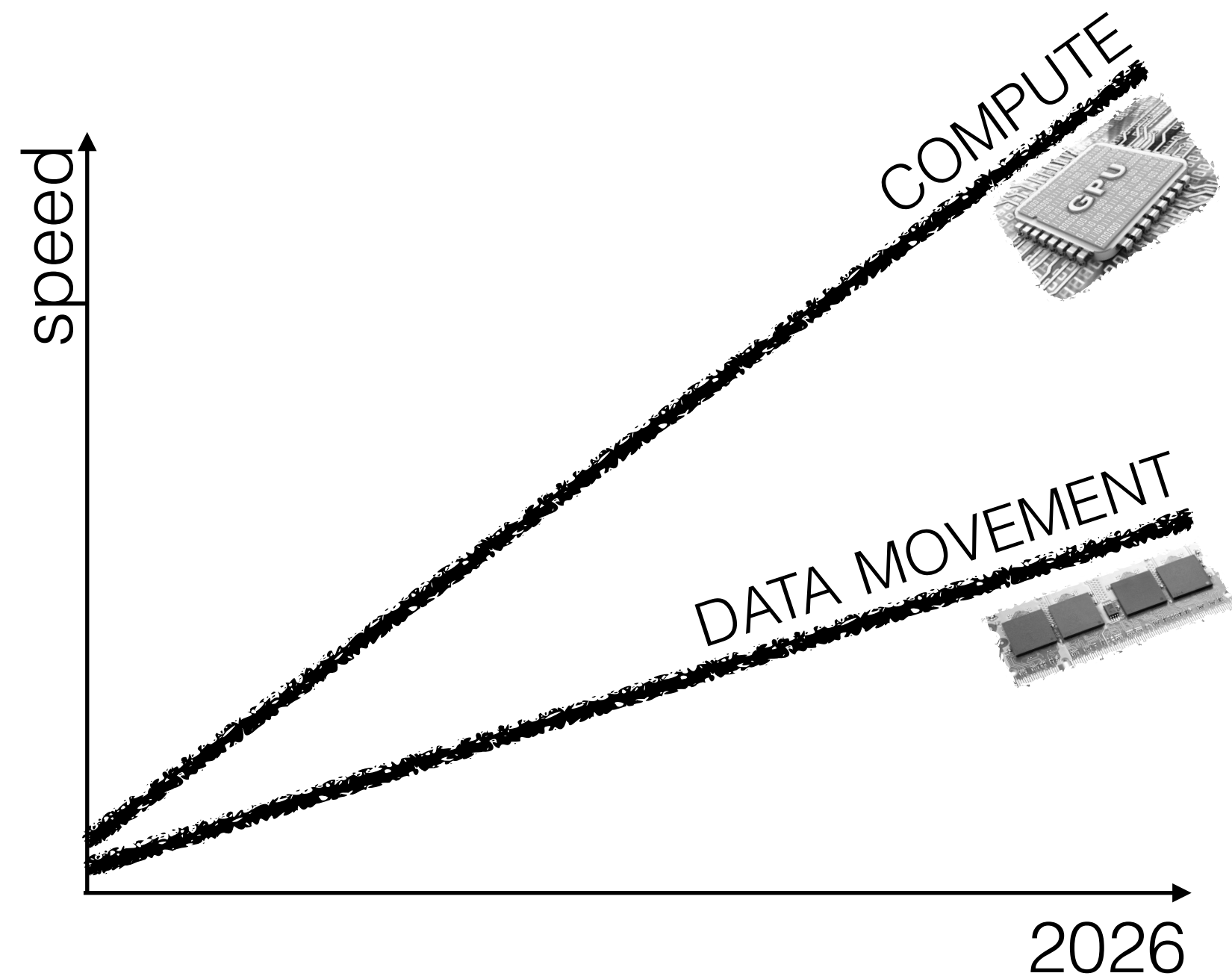
system architecture
it starts with storage



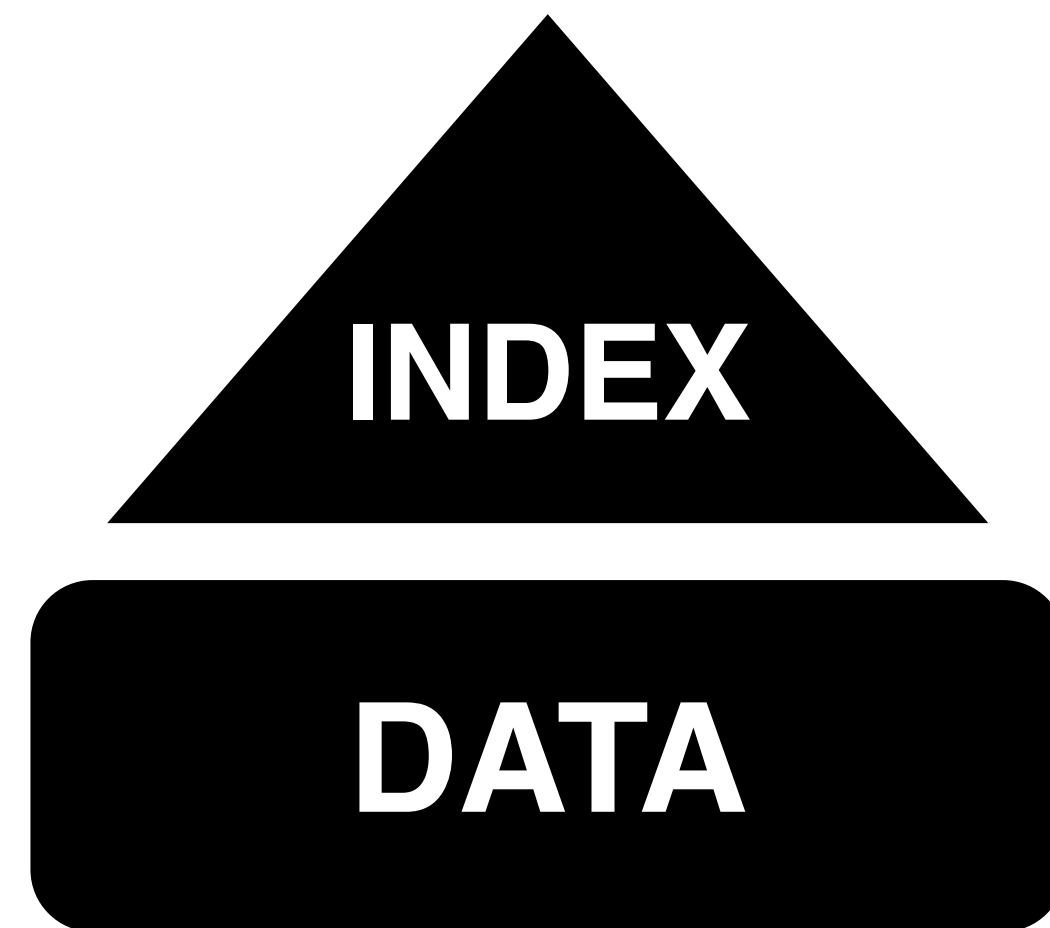
the right data system



the right data system



**System architecture design gets more complex
with bigger data and new diverse hardware**



—HOW—
TO STORE
—DATA—

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

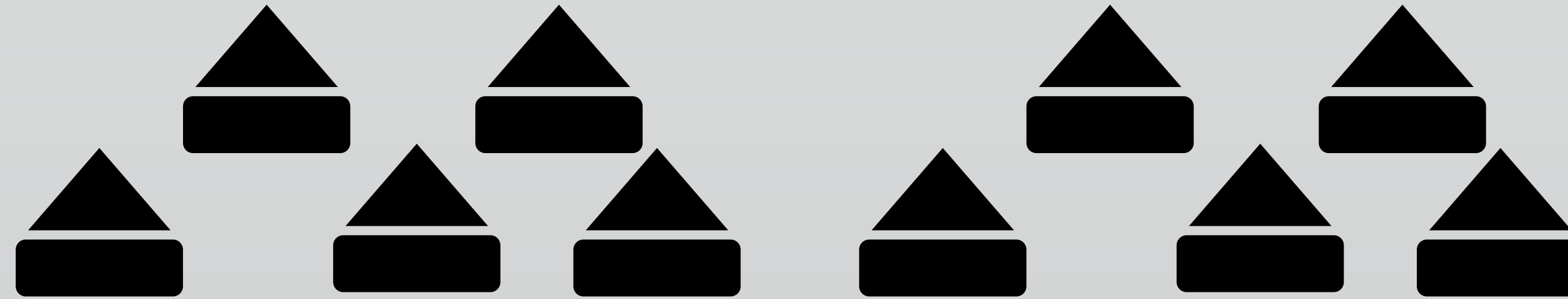


ALGORITHMS

INDEX

DATA

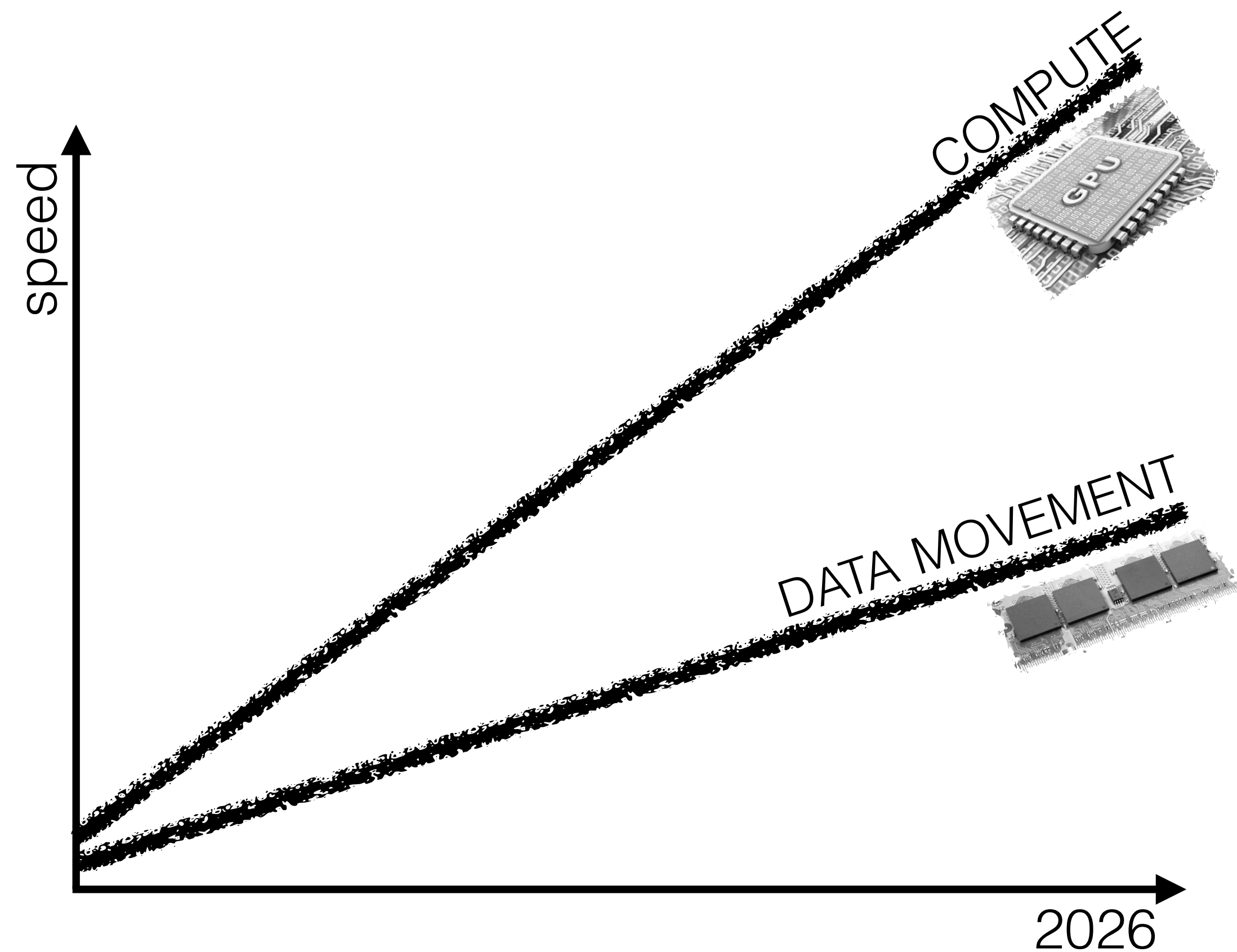
SYSTEMS



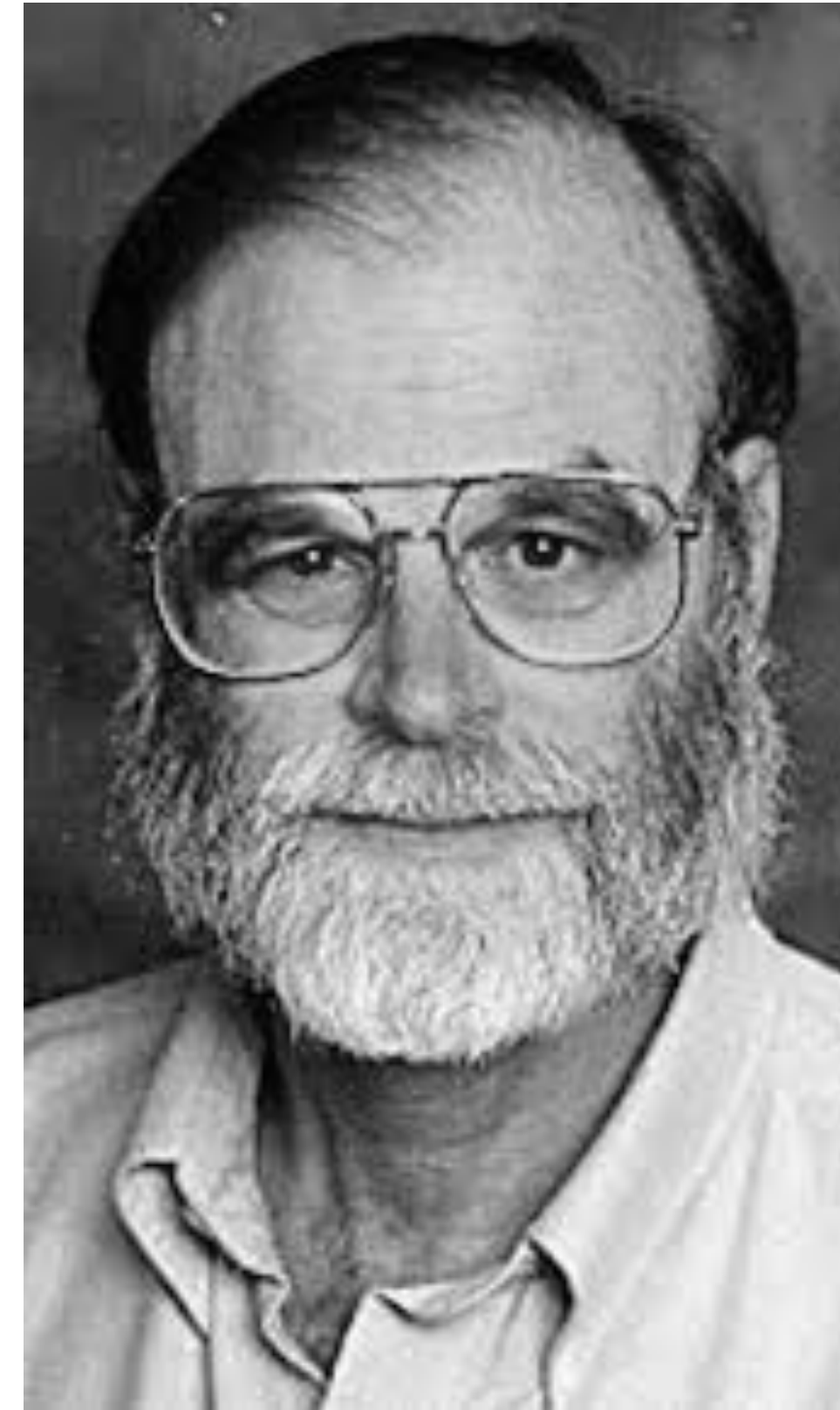
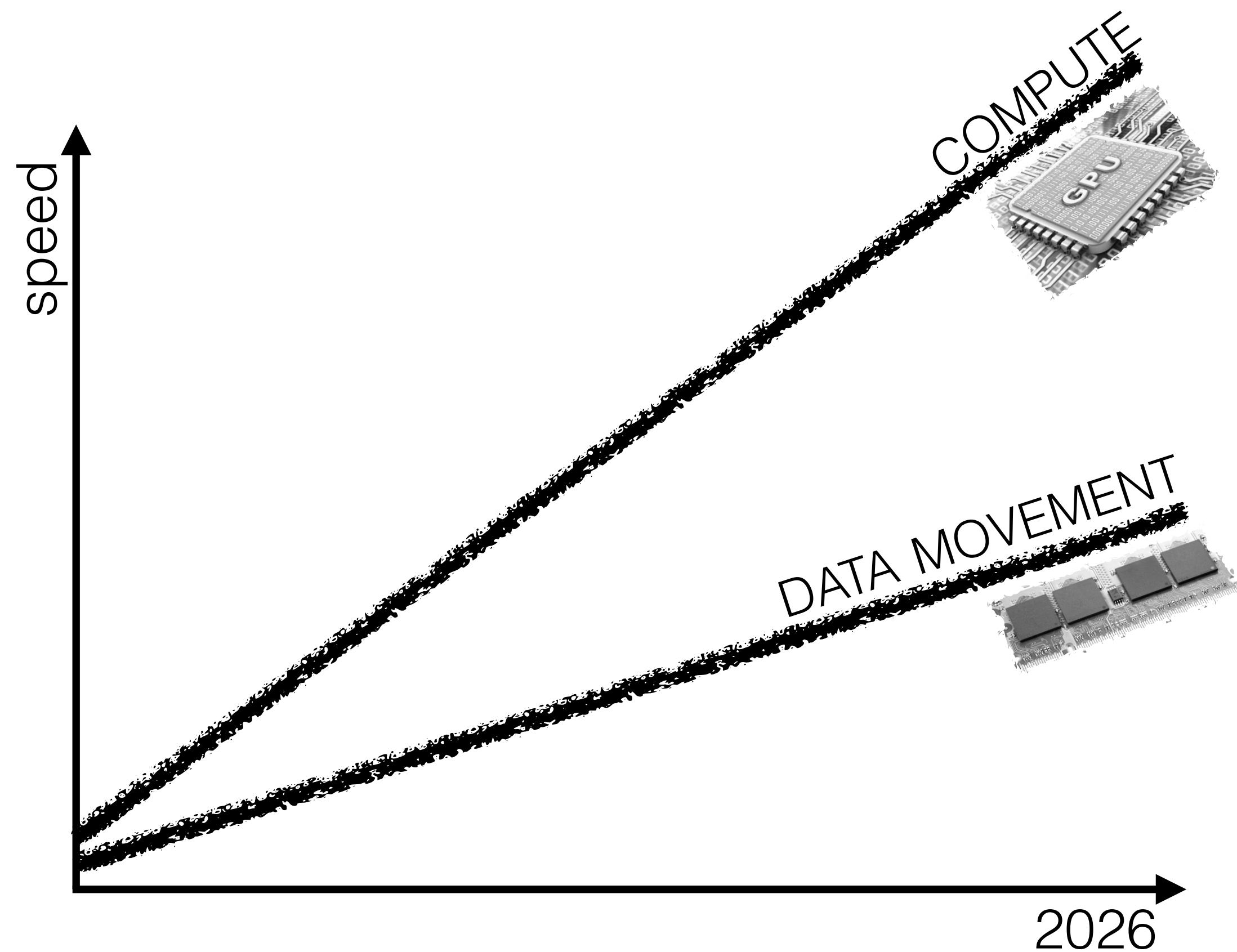
ALGORITHMS

INDEX

DATA

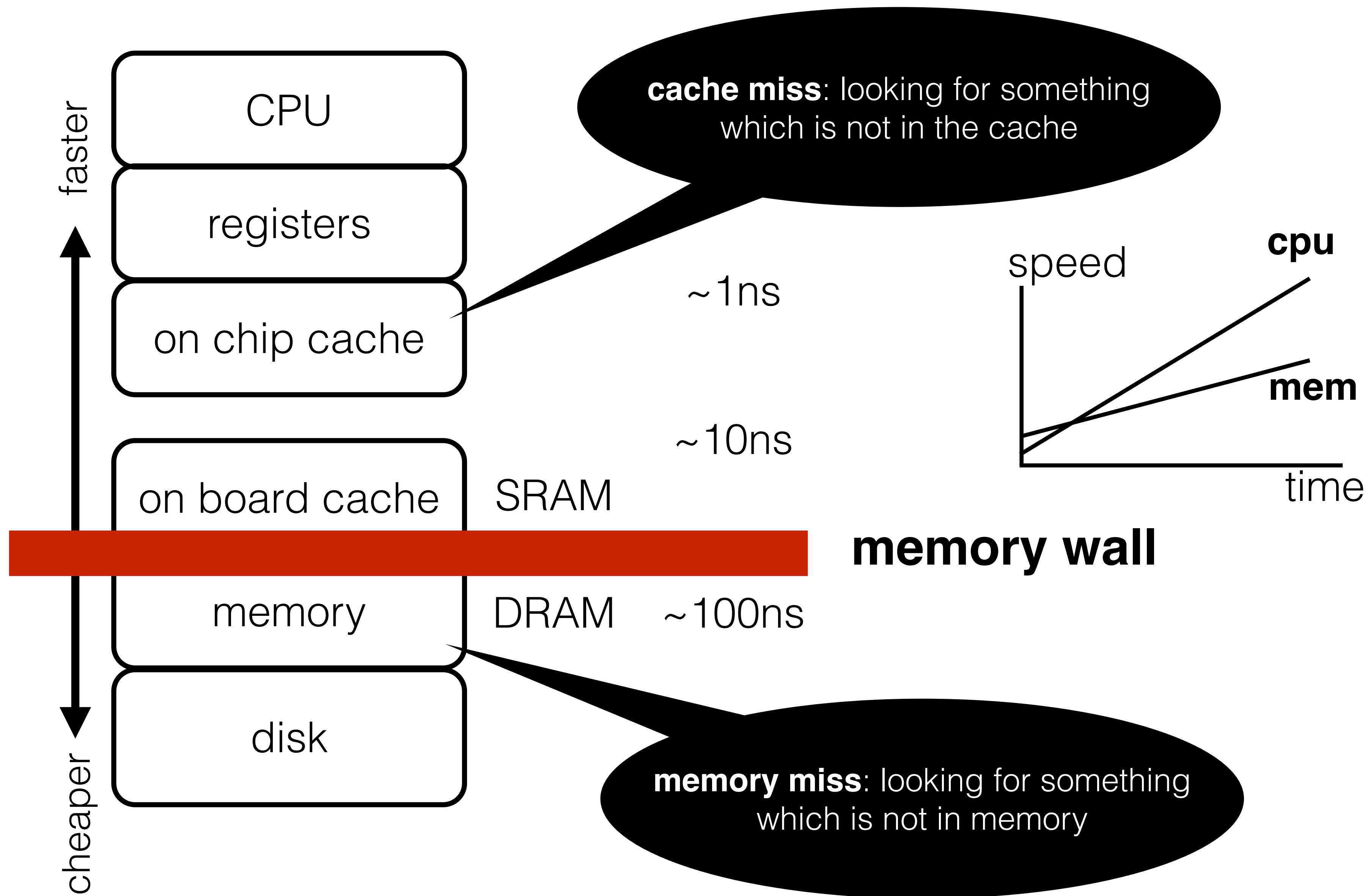


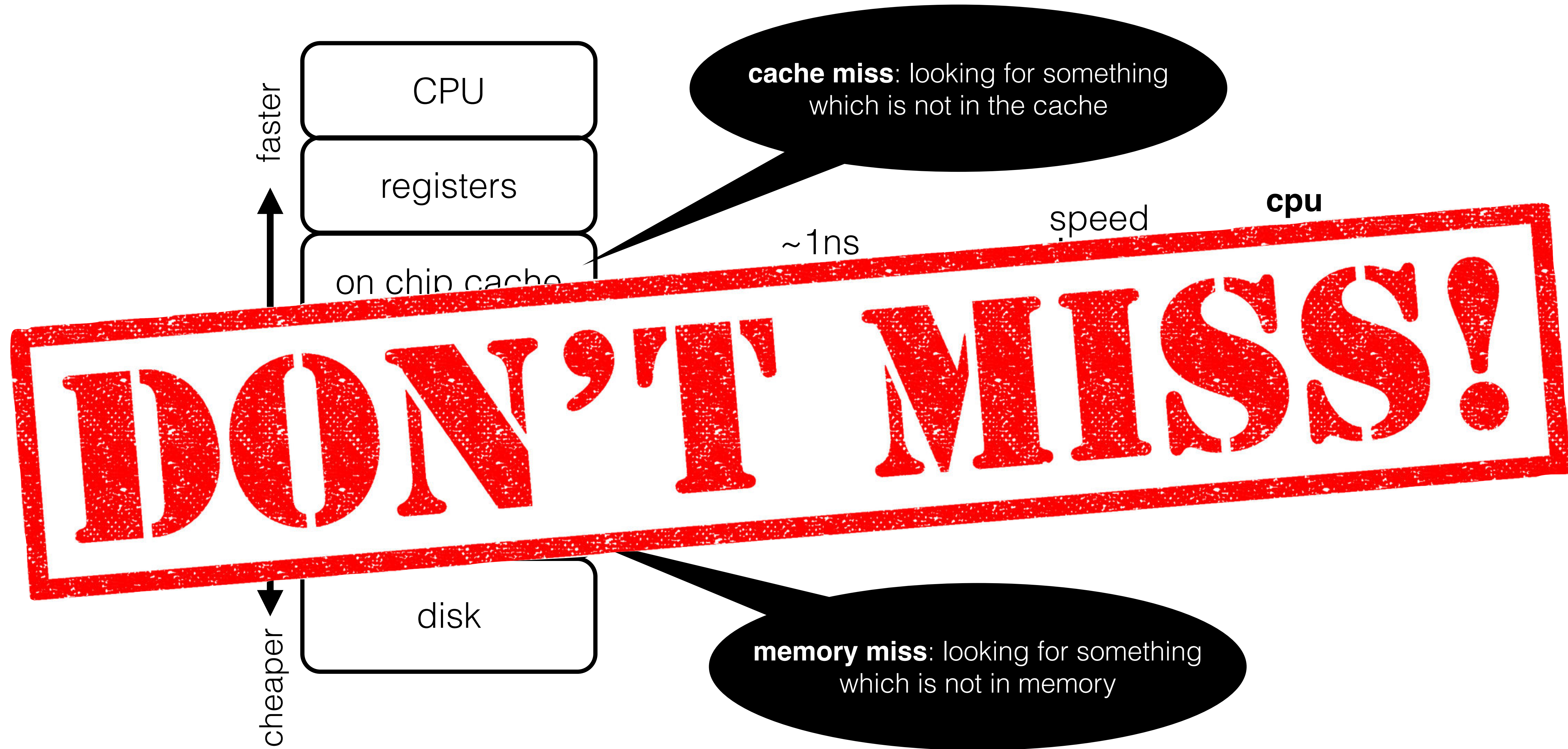
**DATA
STRUCTURES
DEFINE
PERFORMANCE**

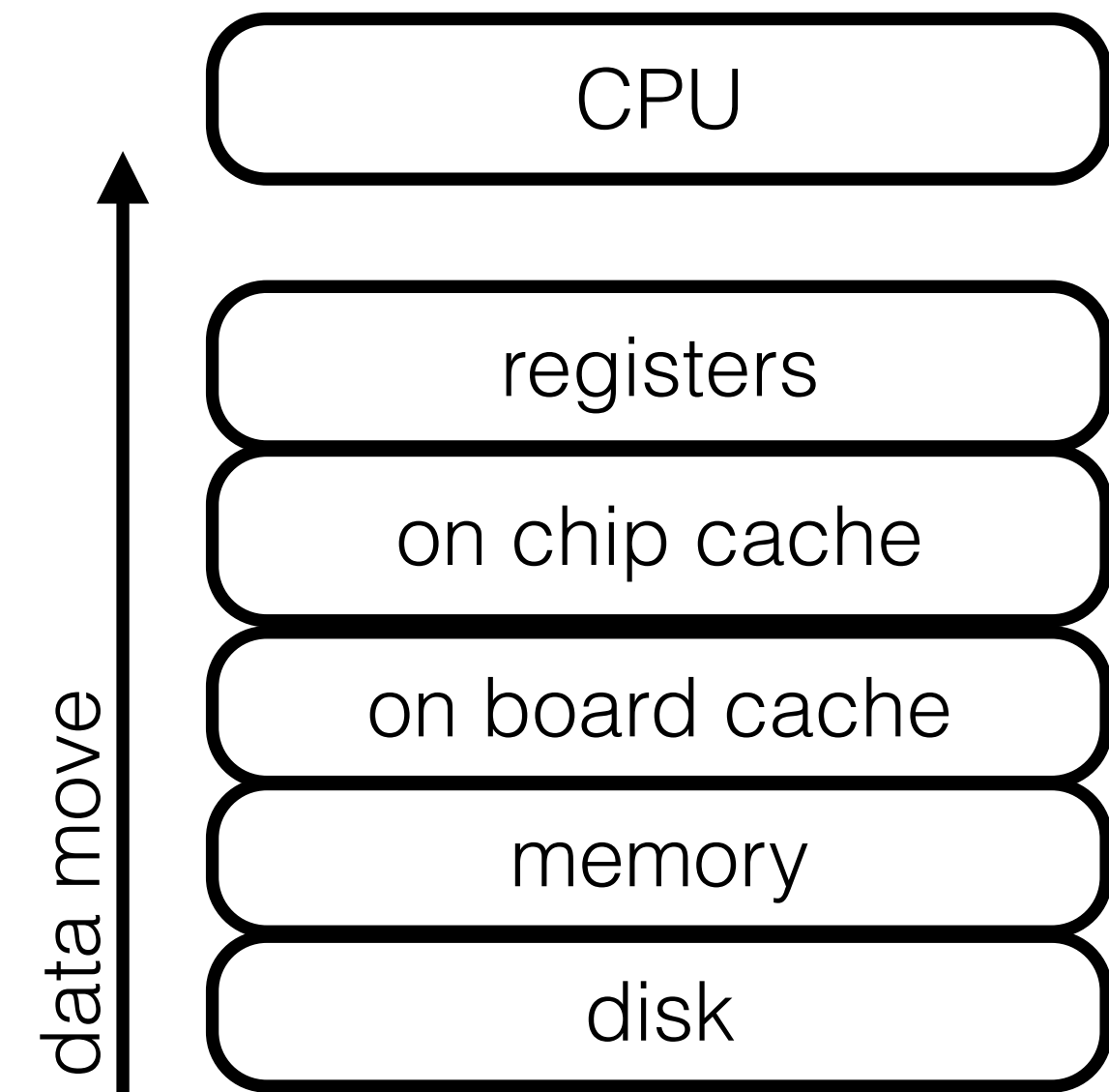
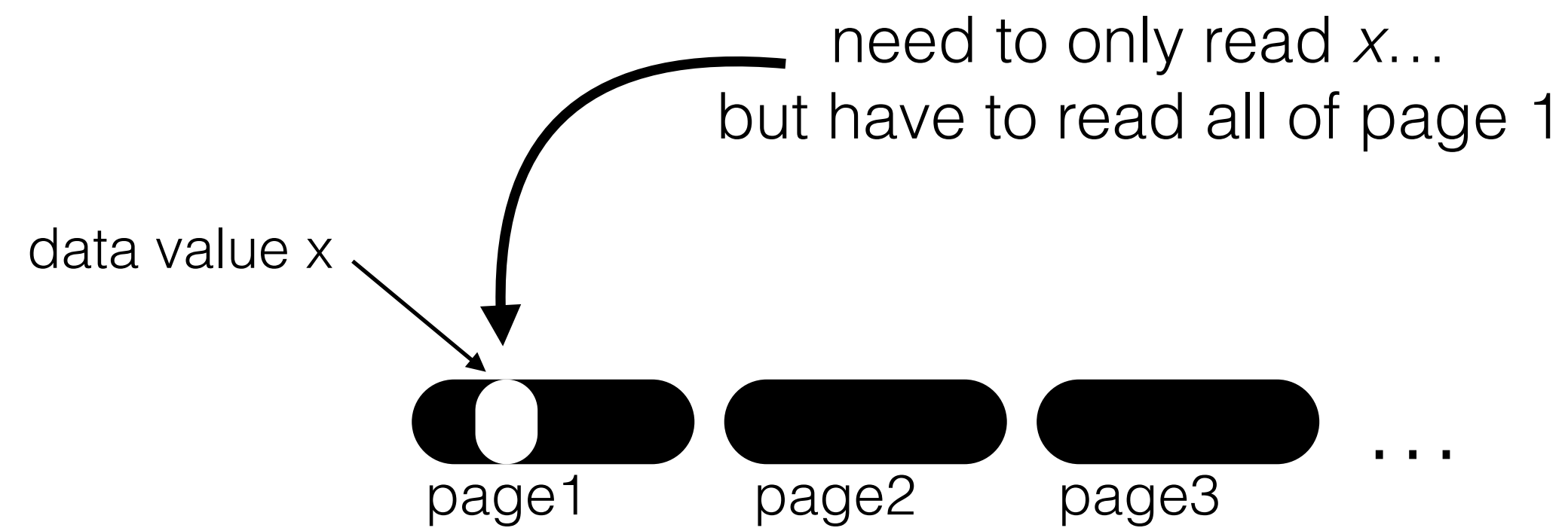


register = this room
caches = this city
memory = nearby city
disk = Pluto

Jim Gray, Turing Award 1998







query $x < 5$

(size=120 bytes)
memory level N

memory level N-1

5 10 6 4 12

2 8 9 7 6

7 11 3 9 6

...

page size: 5x8 bytes

query $x < 5$

scan

5 10 6 4 12

(size=120 bytes)
memory level N

memory level N-1

5 10 6 4 12

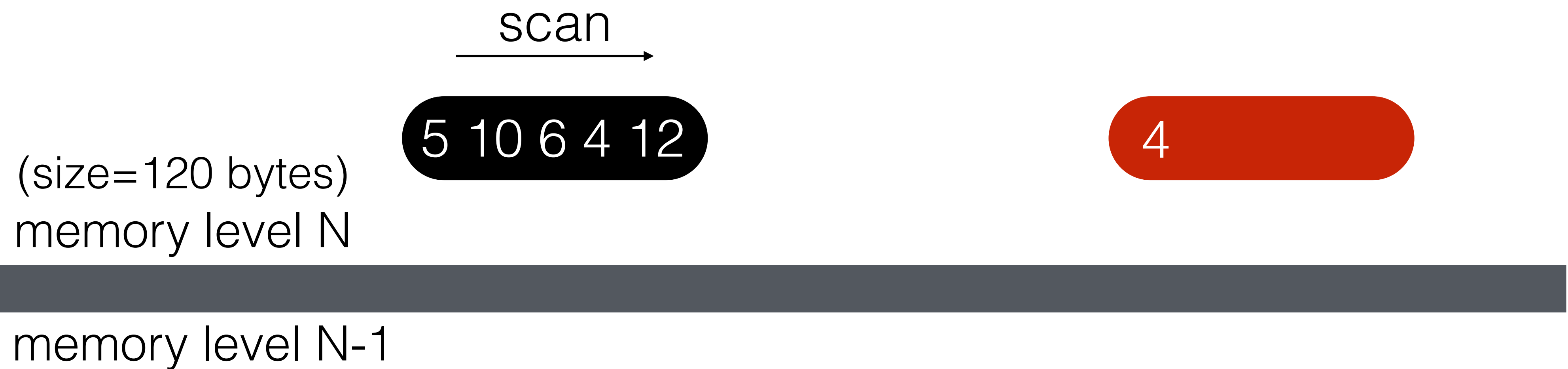
2 8 9 7 6

7 11 3 9 6

...

page size: 5x8 bytes

query $x < 5$



5 10 6 4 12 2 8 9 7 6 7 11 3 9 6 ...

page size: 5x8 bytes



40 bytes

query $x < 5$

scan →

(size=120 bytes)
memory level N

5 10 6 4 12

4

memory level N-1

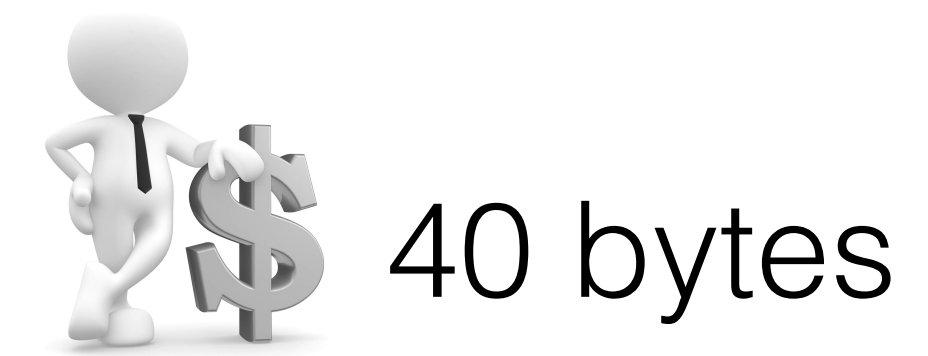
5 10 6 4 12

2 8 9 7 6

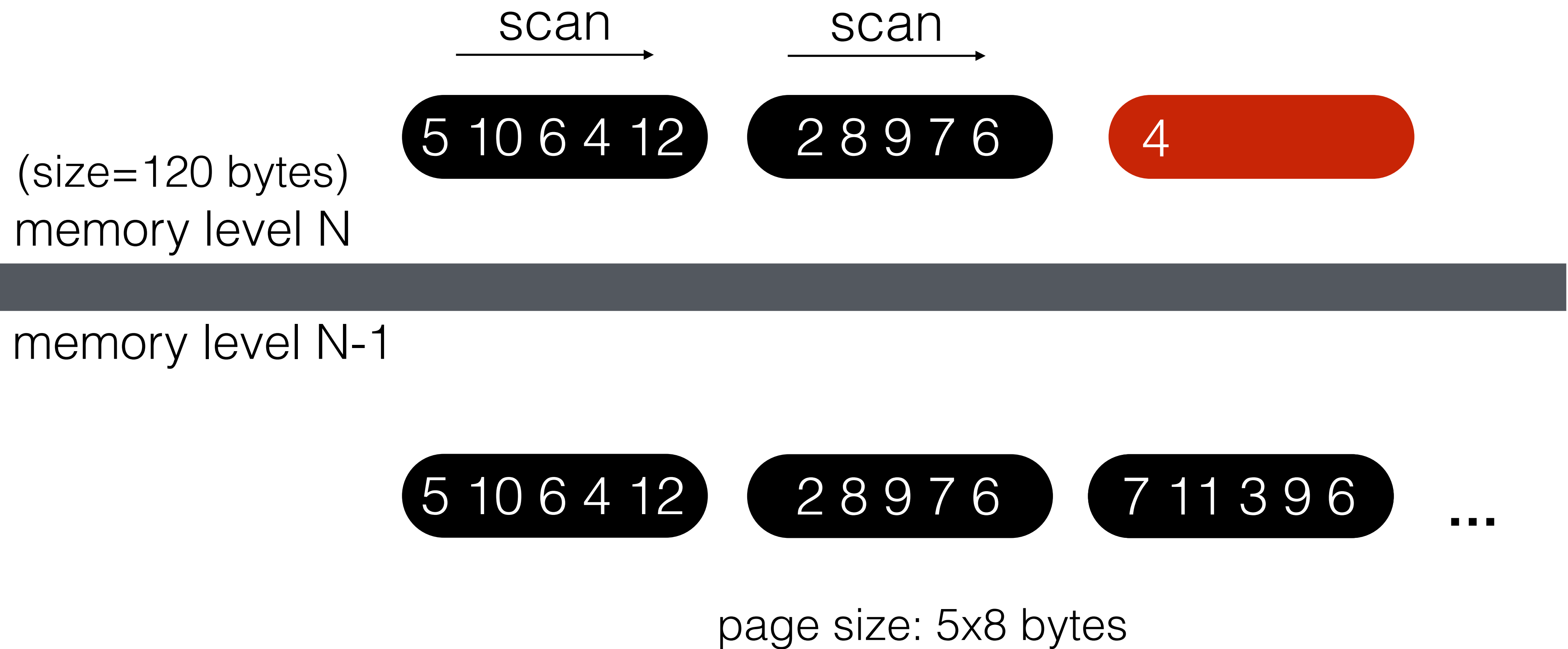
7 11 3 9 6

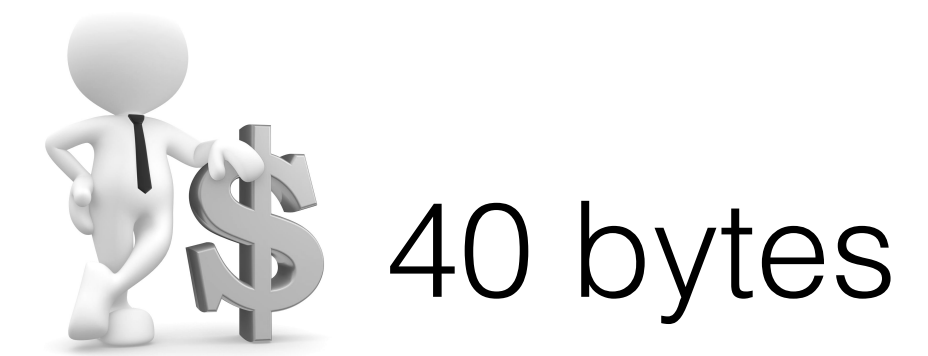
...

page size: 5x8 bytes

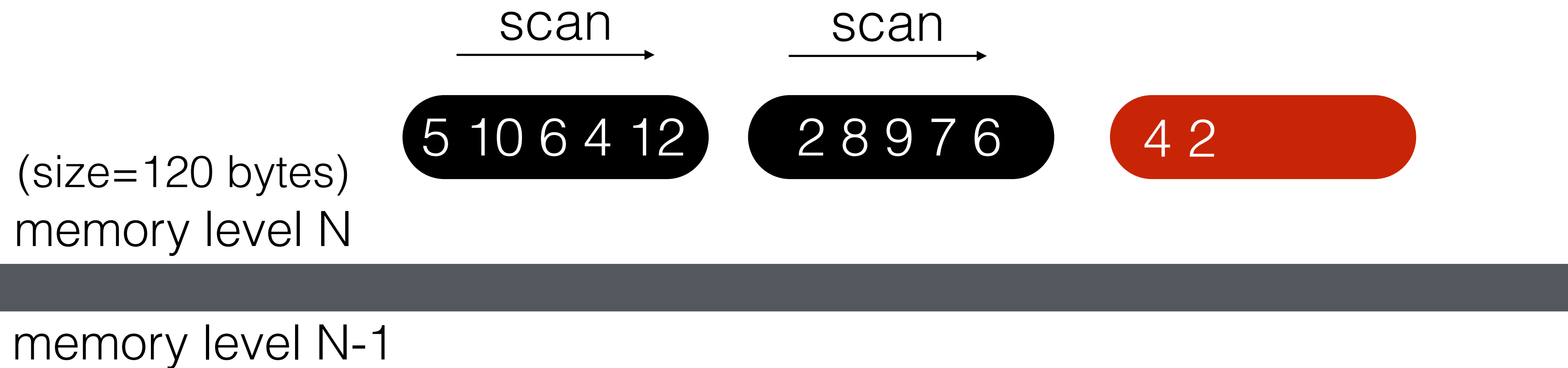


query $x < 5$

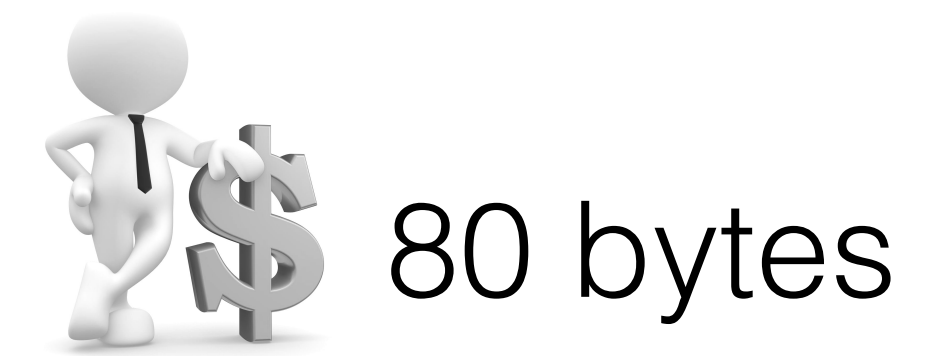




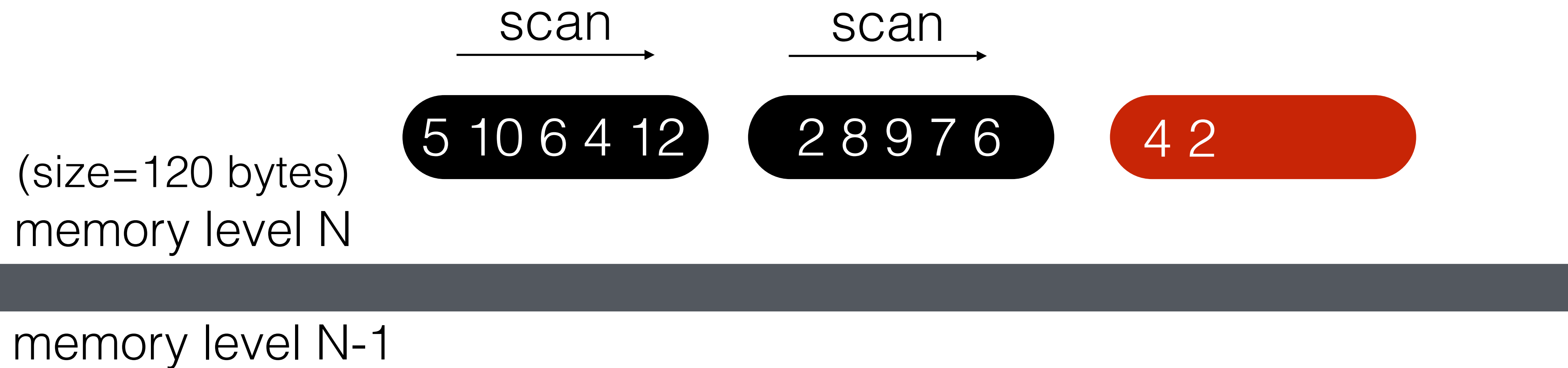
query $x < 5$



page size: 5x8 bytes



query $x < 5$



page size: 5x8 bytes



80 bytes

query $x < 5$

(size=120 bytes)
memory level N

2 8 9 7 6

4 2

memory level N-1

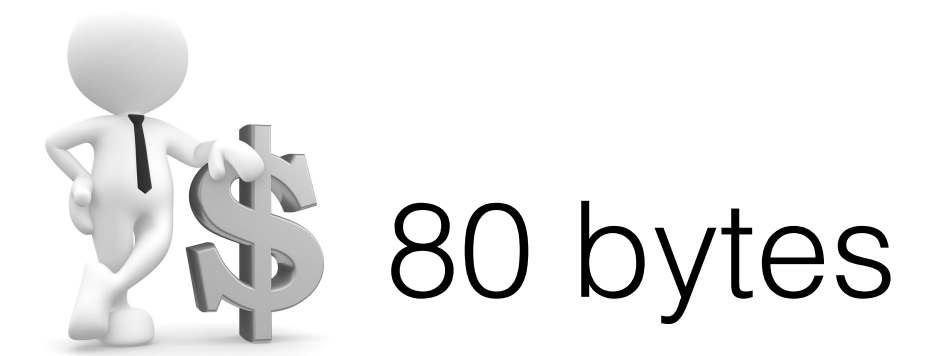
5 10 6 4 12

2 8 9 7 6

7 11 3 9 6

...

page size: 5x8 bytes



query $x < 5$

scan →

(size=120 bytes)
memory level N

7 11 3 9 6

2 8 9 7 6

4 2

memory level N-1

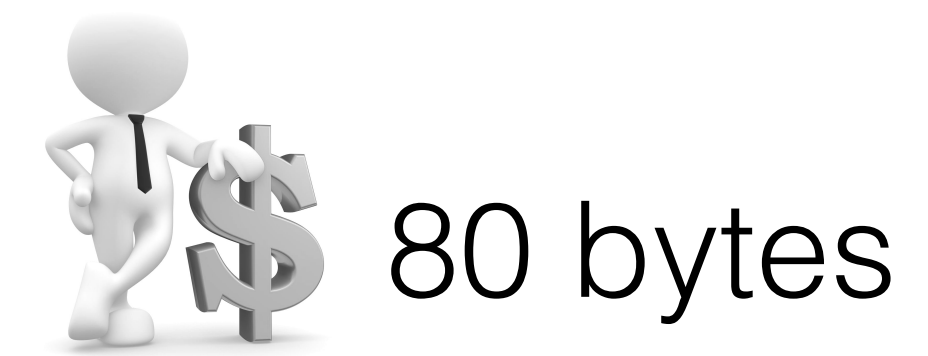
5 10 6 4 12

2 8 9 7 6

7 11 3 9 6

...

page size: 5x8 bytes



query $x < 5$

scan →

(size=120 bytes)
memory level N

7 11 3 9 6

2 8 9 7 6

4 2 3

memory level N-1

5 10 6 4 12

2 8 9 7 6

7 11 3 9 6

...

page size: 5x8 bytes



120 bytes

query $x < 5$

scan →

(size=120 bytes)
memory level N

7 11 3 9 6

2 8 9 7 6

4 2 3

memory level N-1

5 10 6 4 12

2 8 9 7 6

7 11 3 9 6

...

page size: 5x8 bytes

an oracle gives us the positions

query $x < 5$

(size=120 bytes)
memory level N



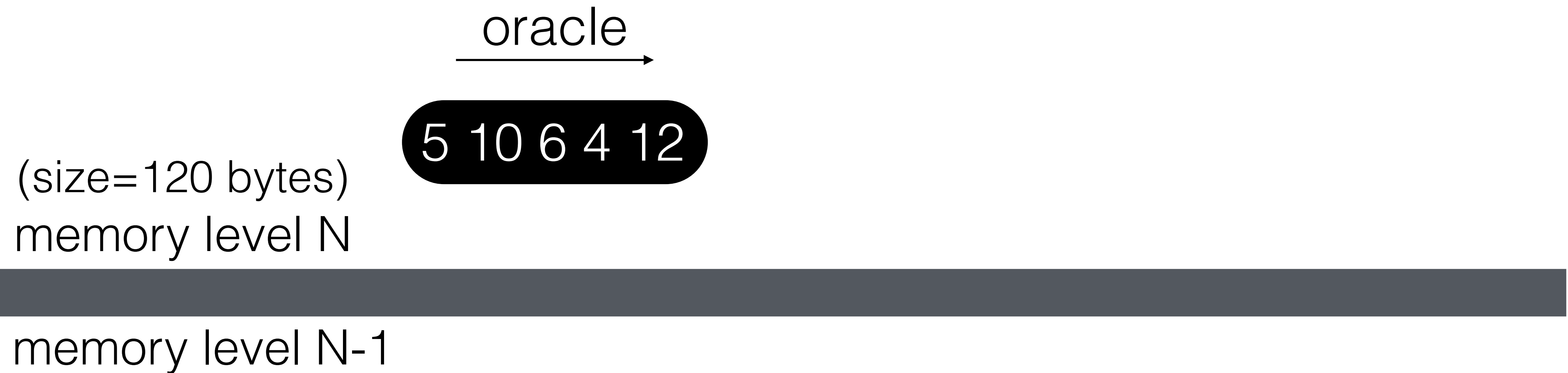
memory level N-1



page size: 5x8 bytes

an oracle gives us the positions

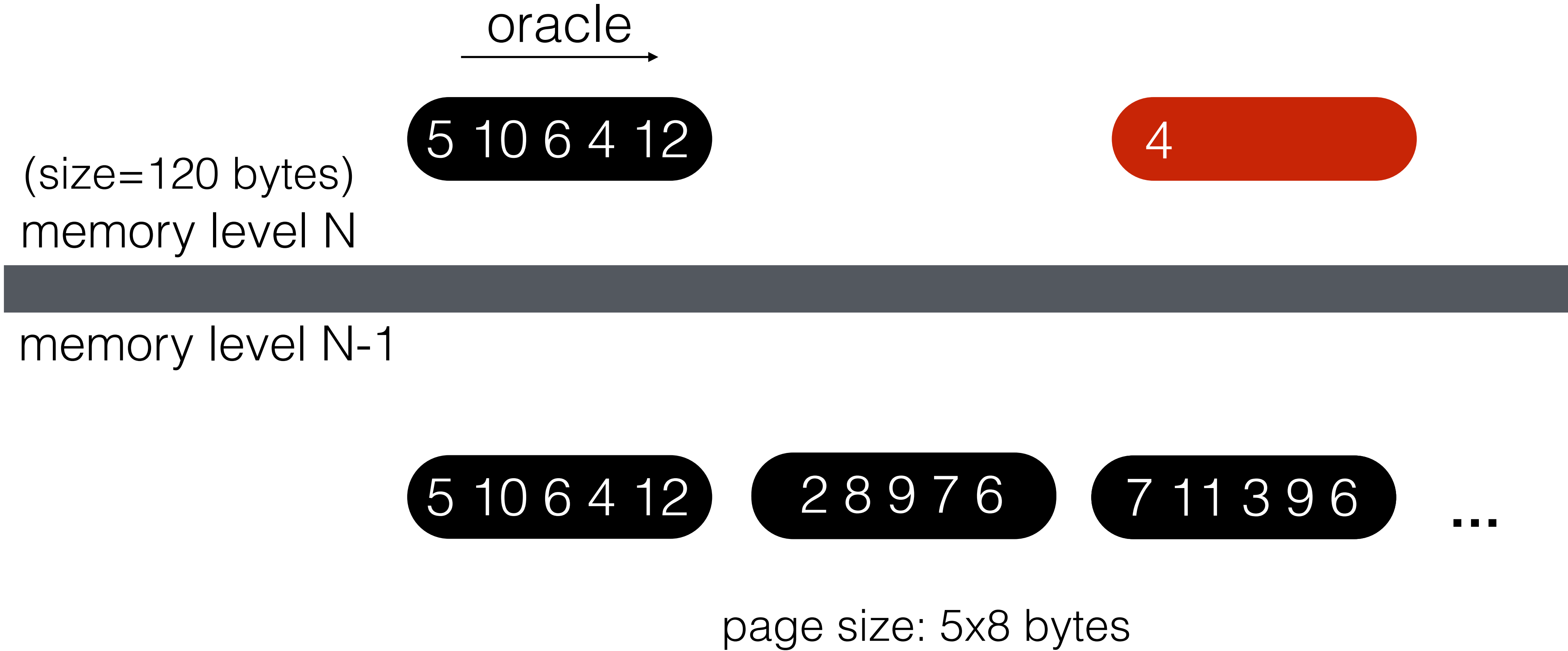
query $x < 5$



page size: 5x8 bytes

an oracle gives us the positions

query $x < 5$

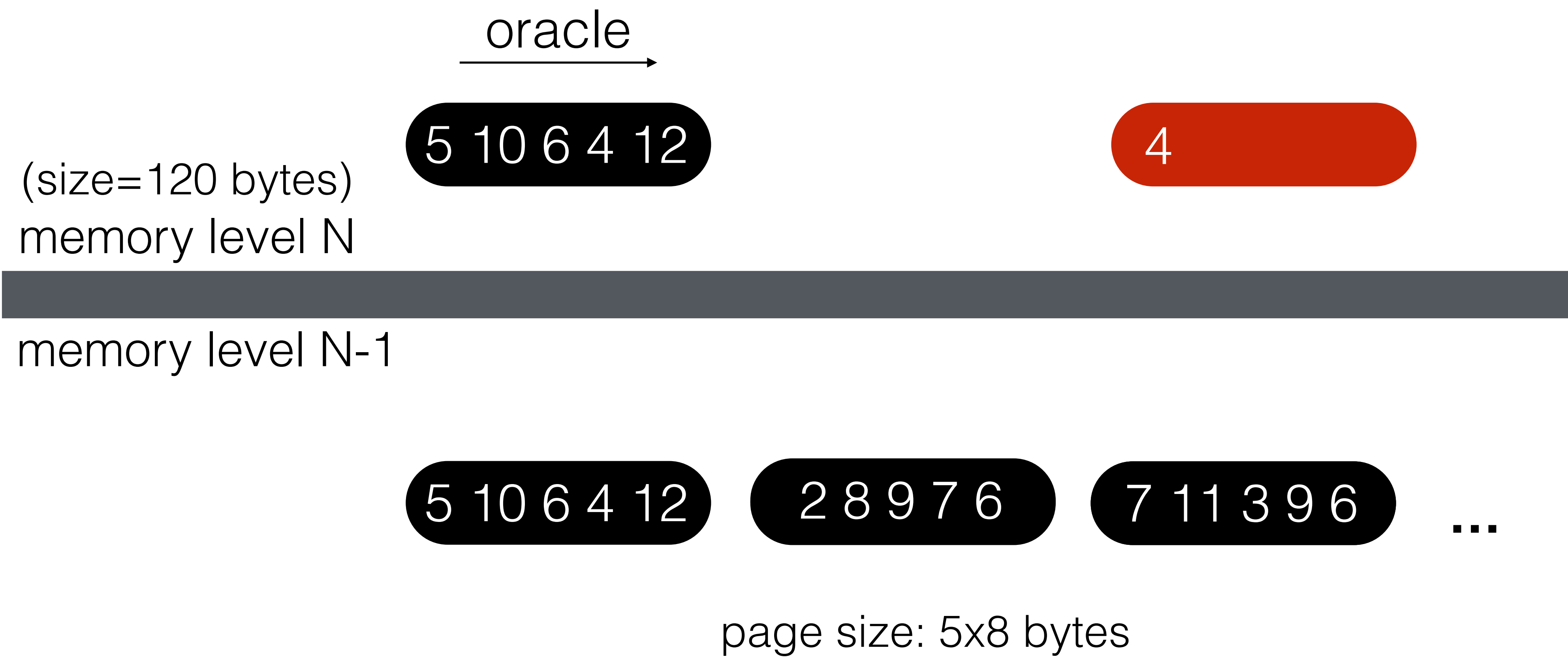


an oracle gives us the positions



40 bytes

query $x < 5$

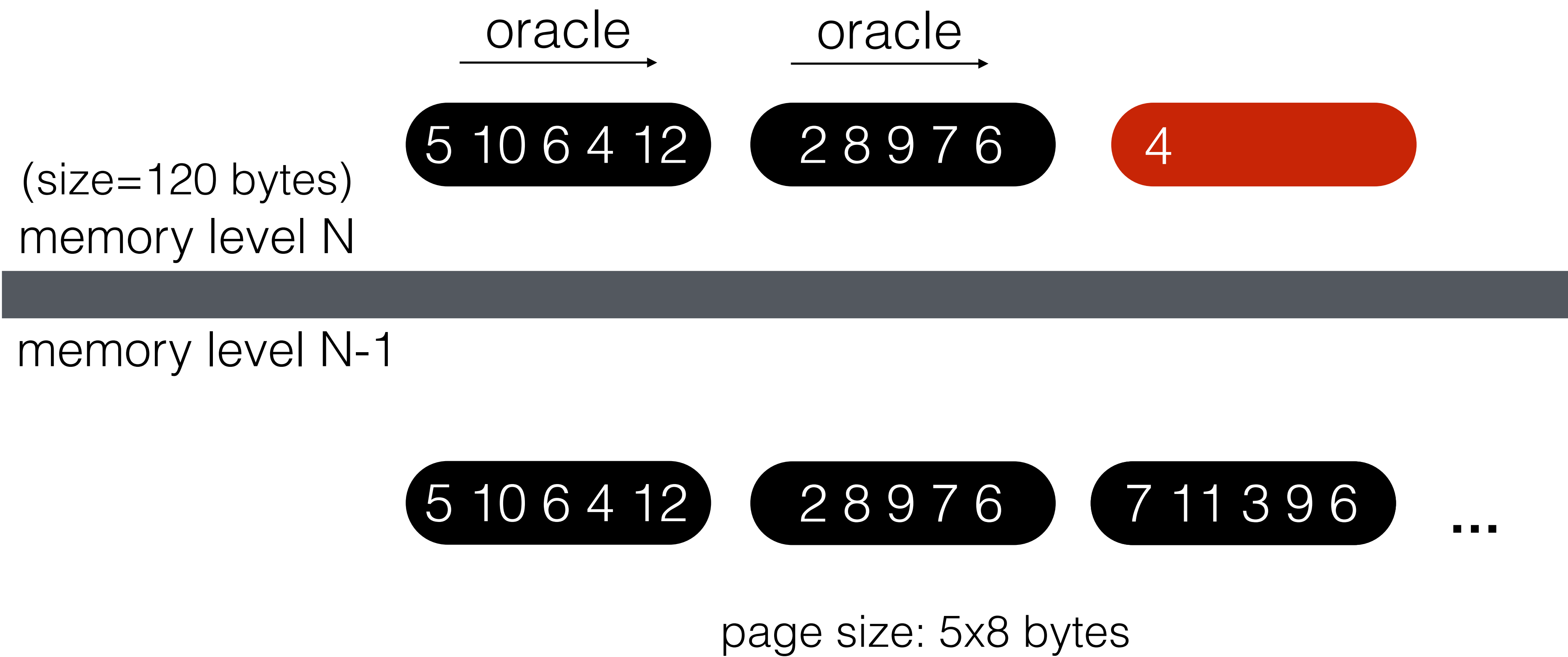


an oracle gives us the positions



40 bytes

query $x < 5$

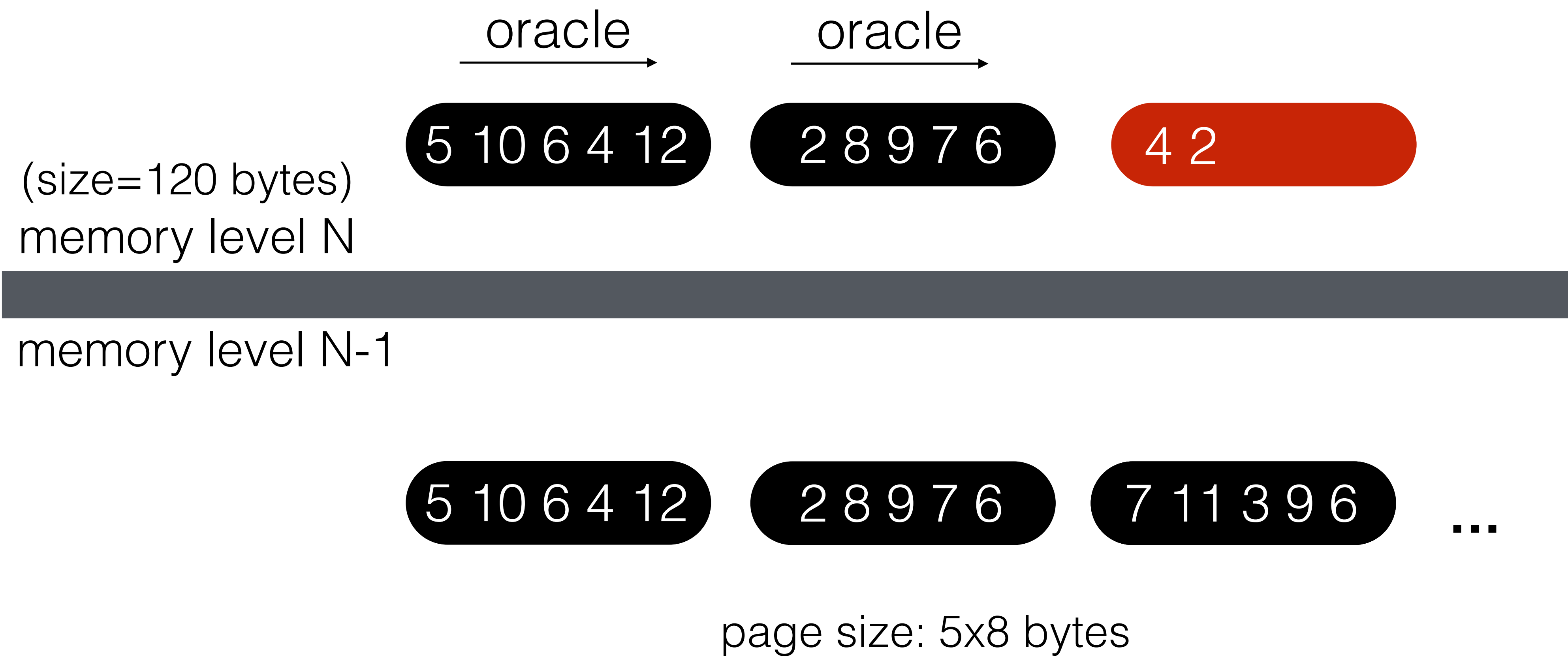


an oracle gives us the positions



40 bytes

query $x < 5$

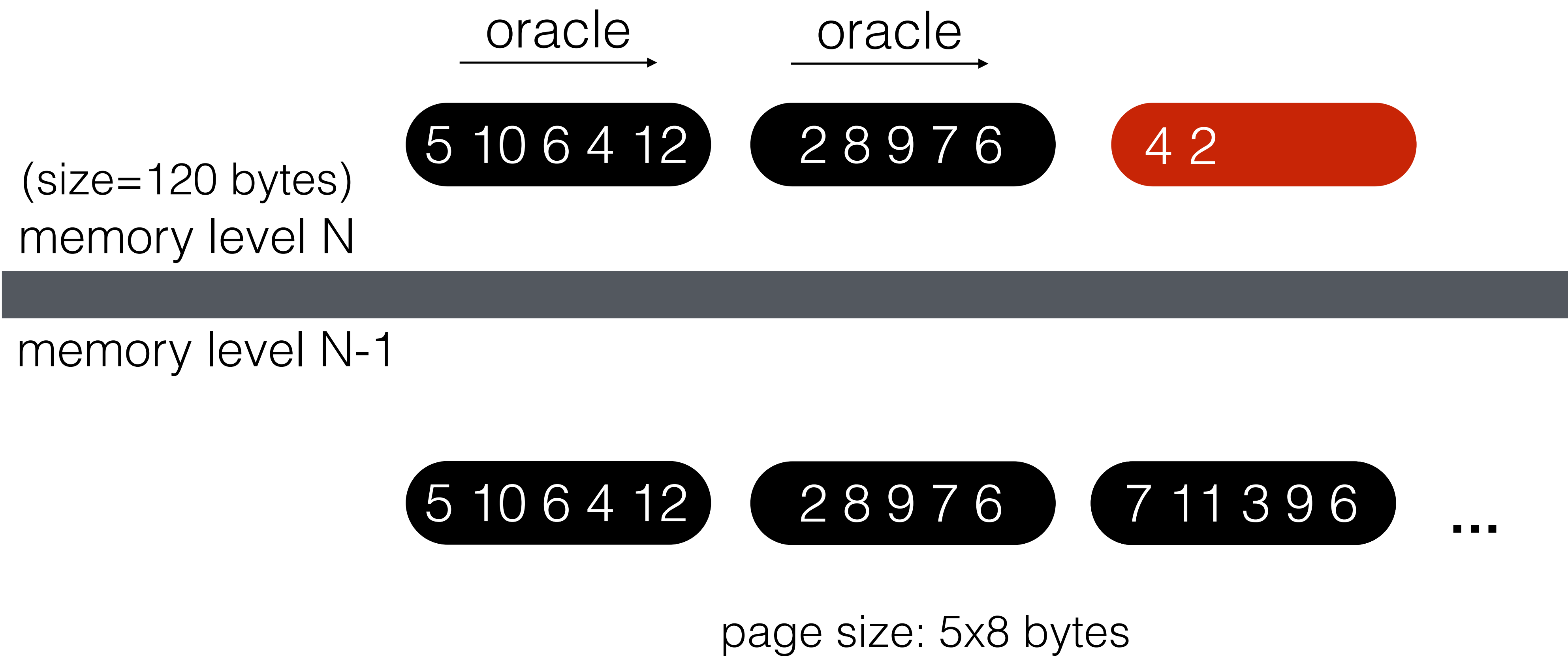


an oracle gives us the positions



80 bytes

query $x < 5$



an oracle gives us the positions



80 bytes

query $x < 5$

(size=120 bytes)
memory level N

2 8 9 7 6

4 2

memory level N-1

5 10 6 4 12

2 8 9 7 6

7 11 3 9 6

...

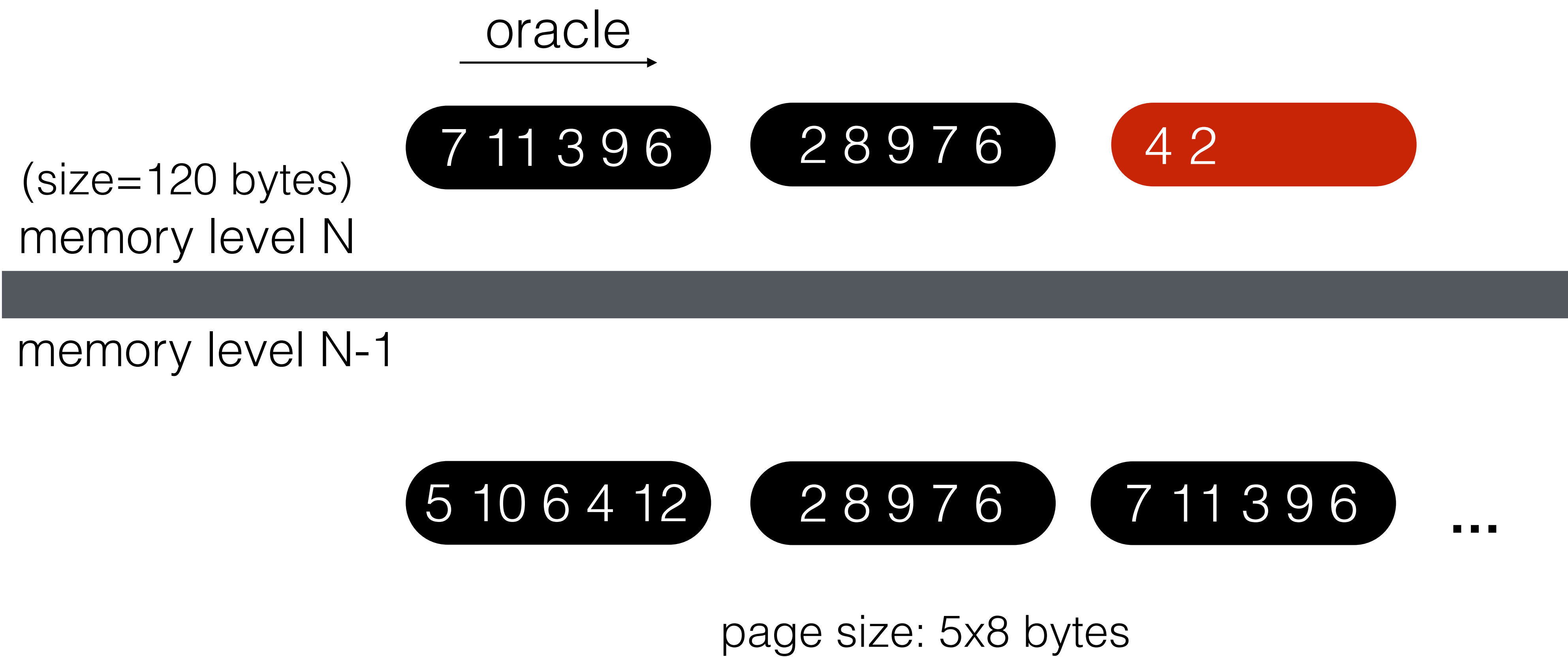
page size: 5x8 bytes

an oracle gives us the positions



80 bytes

query $x < 5$

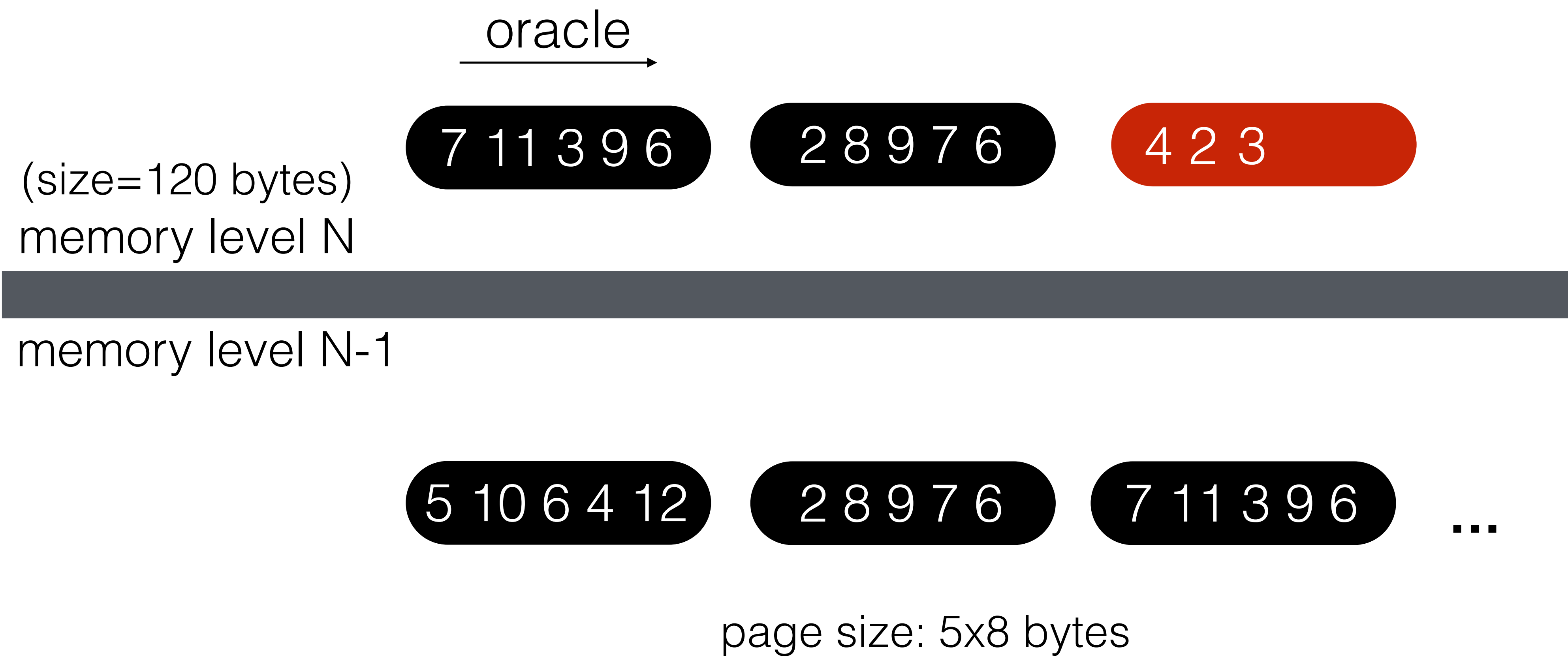


an oracle gives us the positions



80 bytes

query $x < 5$

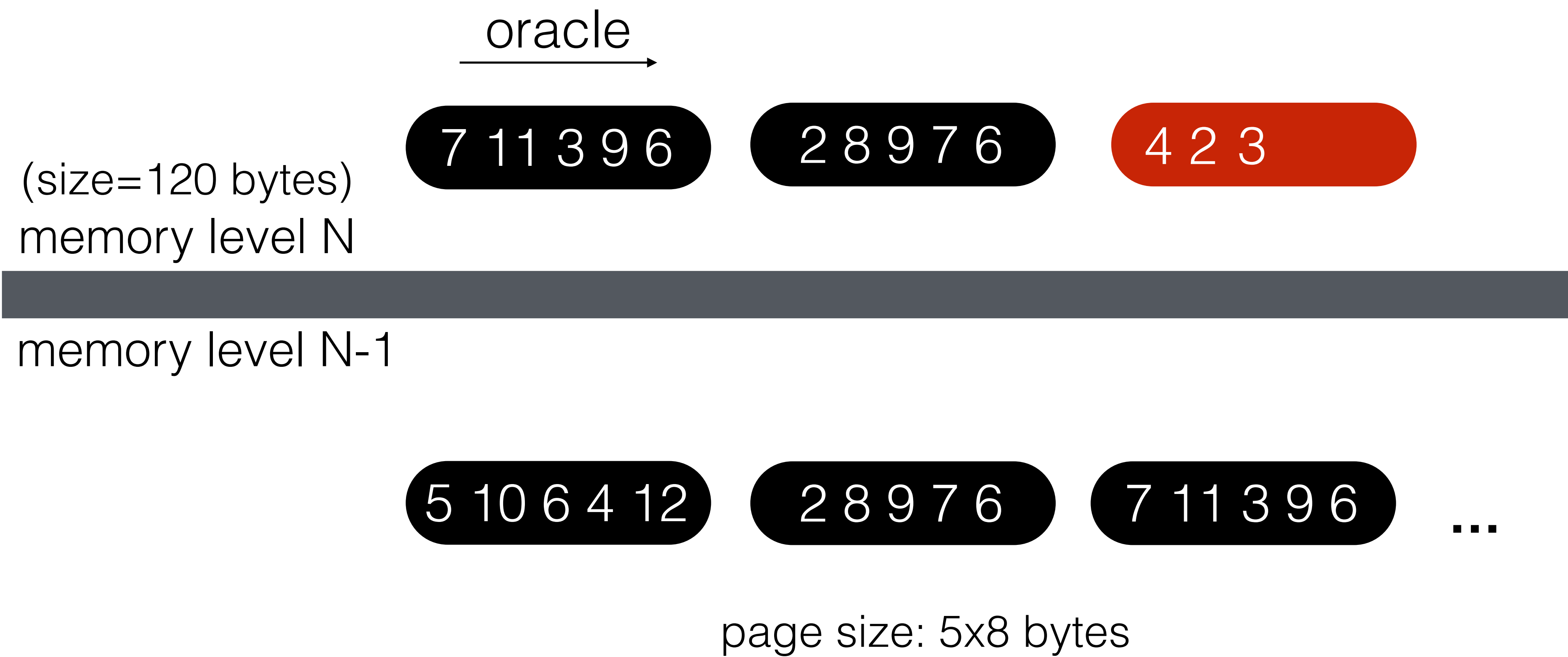


an oracle gives us the positions



120 bytes

query $x < 5$



when does it make sense to have an oracle
how can we minimize the cost



e.g., **query** $x < 5$

5 10 6 4 12

2 8 9 7 6

7 11 3 9 6

...

algorithm/system design = not just computation

algorithm/system design = not just computation

Is there maybe a perfect system? Nope...

basic CS265 logistics

learning outcome

Fundamentals of storage

data structures, SQL, NoSQL, Agents, LLMs, RAG, Data Science, Image AI

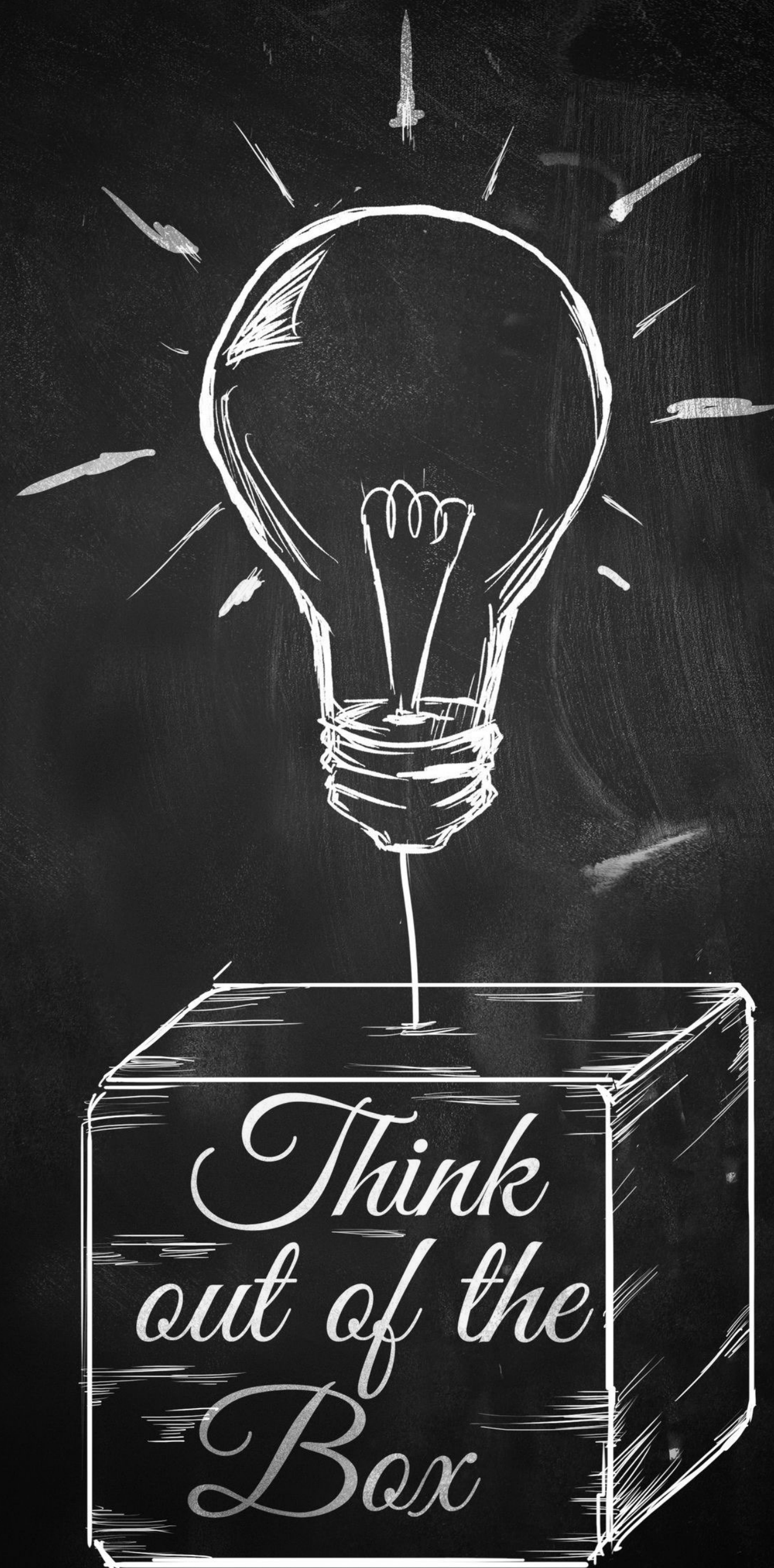
learning outcome

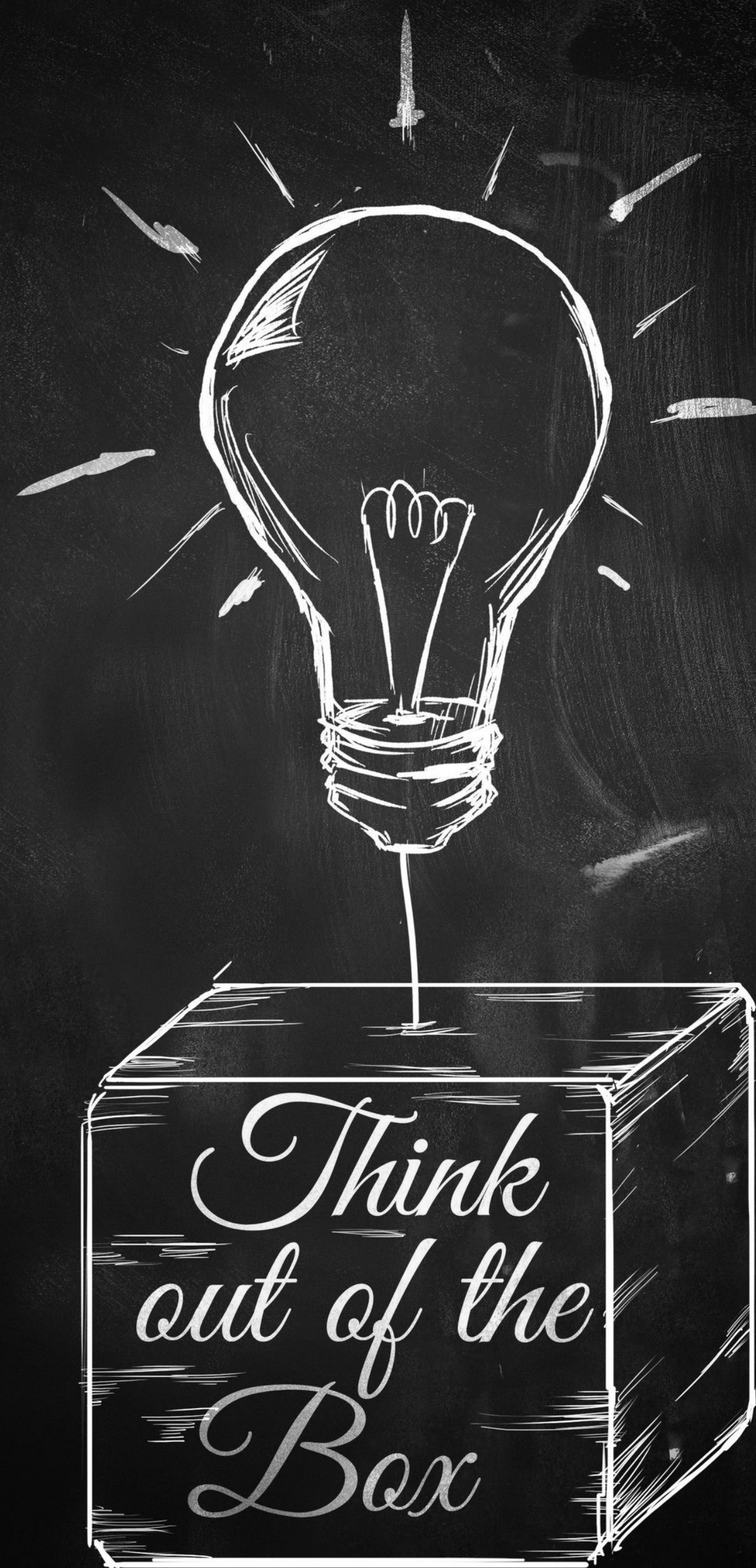
Fundamentals of storage

data structures, SQL, NoSQL, Agents, LLMs, RAG, Data Science, Image AI

Self-designing systems

Automated system design: cloud cost, hardware, data & app requirements





first ~5 weeks: Stratos & TFs

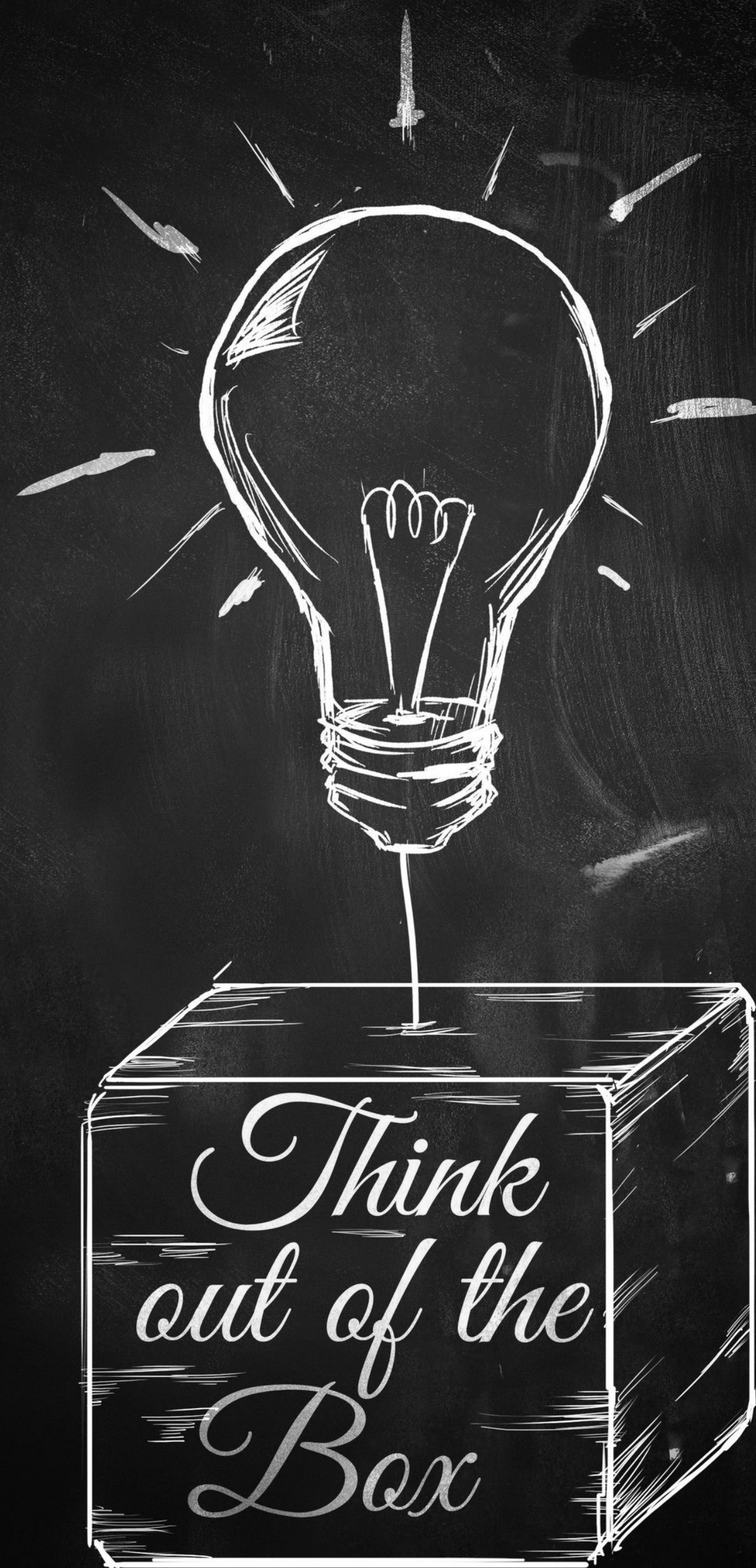
Basic background

Self-designing systems

LLM systems

Image AI systems

Research thinking



first ~5 weeks: Stratos & TFs

Basic background

Self-designing systems

LLM systems

Image AI systems

Research thinking

afterwards:

Students present research papers

One paper per class (ML systems)

In-class research/systems discussion

Research reviews

Research/systems projects



Recent Research Papers

Each student:
In-class discussions/1 presentation

review and slides should focus on

- what is the problem
- why is it important
- why is it hard
- why existing solutions do not work
- what is the core intuition for the solution
- solution step by step
- does the paper prove its claims
- exact setup of analysis/experiments
- are there any gaps in the logic/proof
- possible next steps

* follow a few citations to gain more background

learn to judge constructively

learn to present

learn to prepare slides

Each student:

In-class discussions/1 presentation

review and slides should focus on

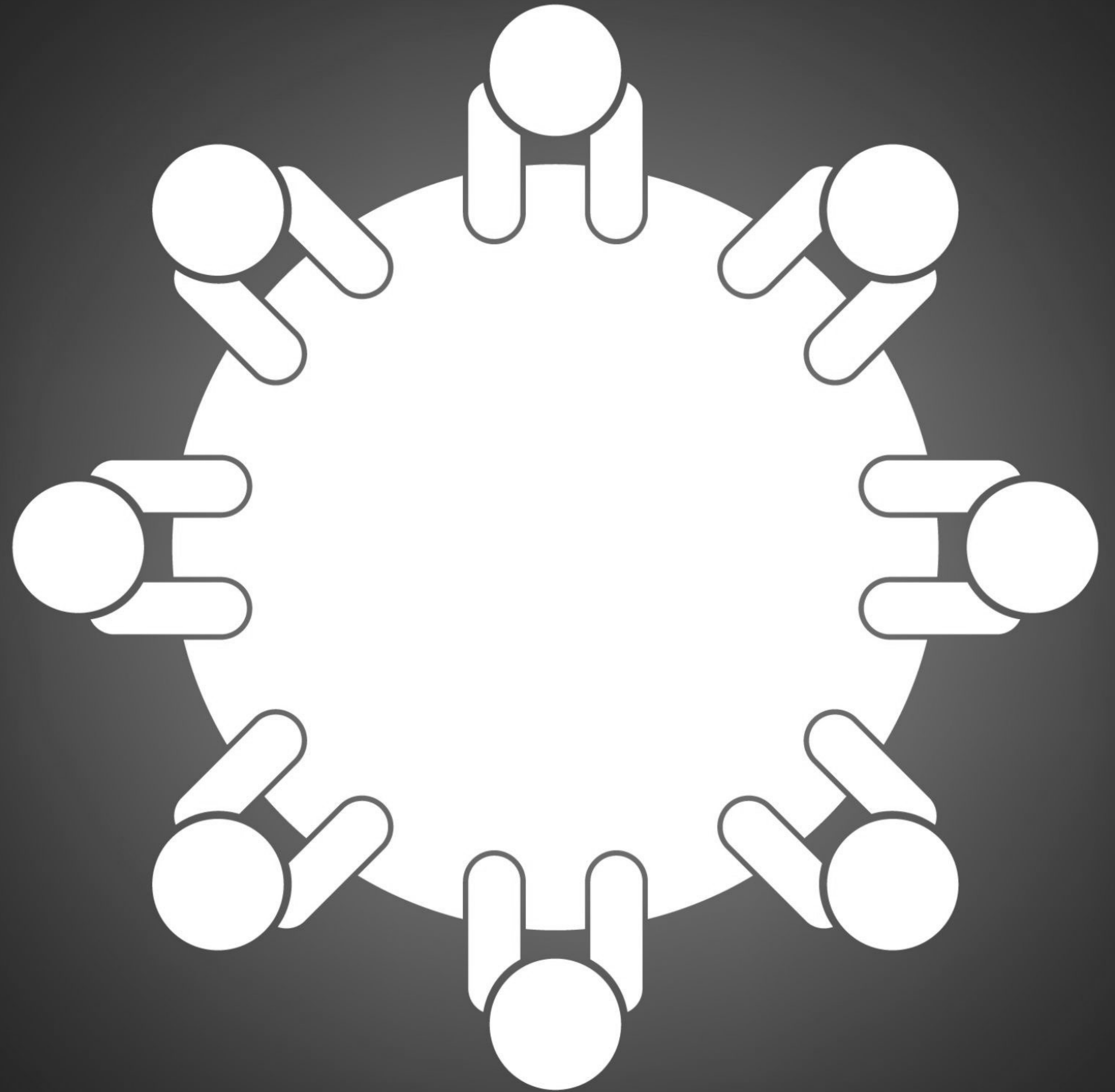
- what is the problem
- why is it important
- why is it hard
- why existing solutions do not work
- what is the core intuition for the solution
- solution step by step
- does the paper prove its claims
- exact setup of analysis/experiments
- are there any gaps in the logic/proof
- possible next steps

* follow a few citations to gain more background

In class discussions
is a critical component
and learning outcome

Think creatively
Fail quickly
Incrementally solve

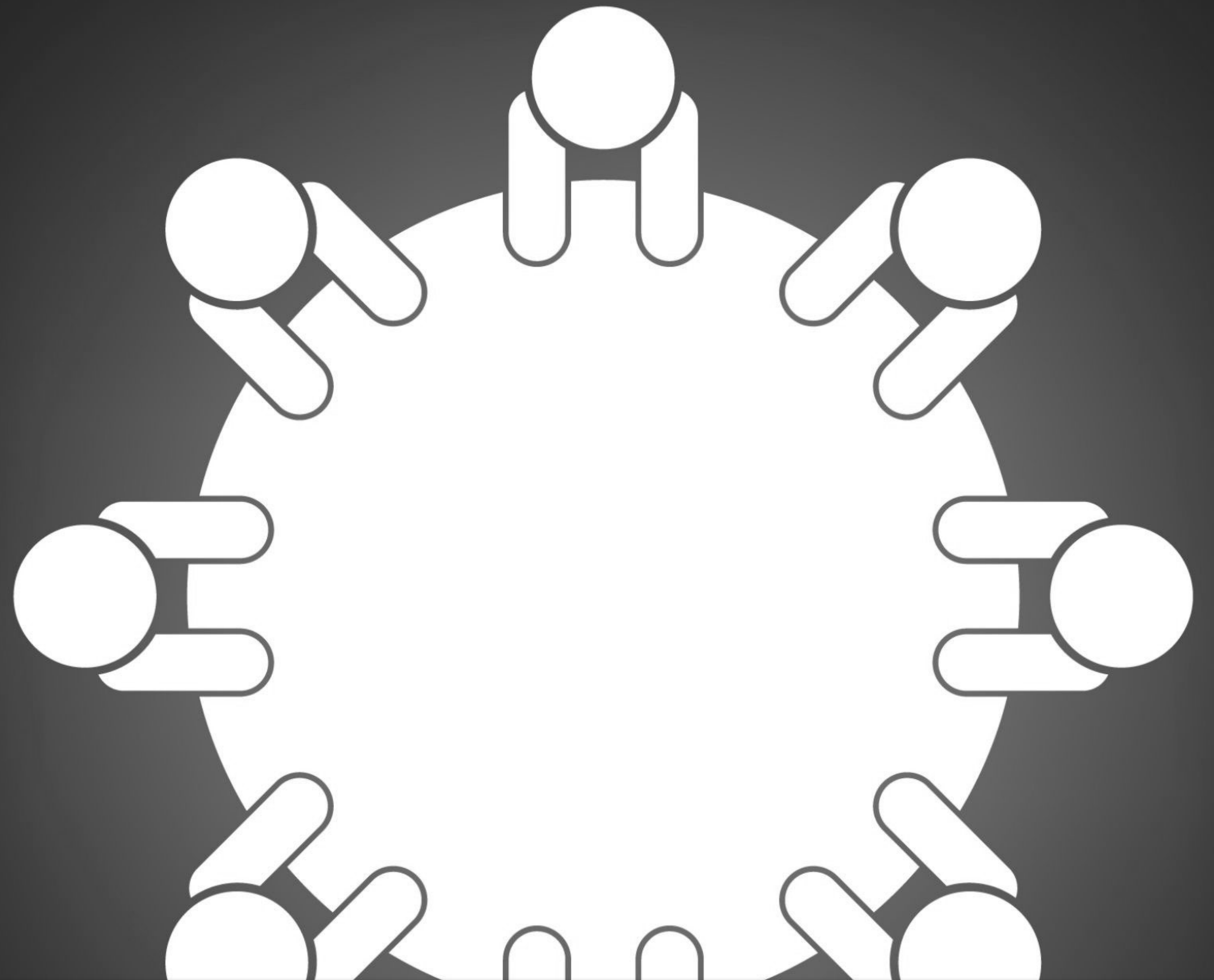
DailyOH/labs in person,
Sat/Sun remote Labs
Friday remote OH



In class discussions
is a critical component
and learning outcome

Think creatively
Fail quickly
Incrementally solve

DailyOH/labs in person,
Sat/Sun remote Labs
Friday remote OH



There is no such thing as a wrong question/answer!!!!

semester project: due in the end of semester + a midway check in (end of March, 10%)

systems project

research project (publish)

semester project: due in the end of semester + a midway check in (end of March, 10%)

systems project

individual project

LLMs, in c/c++

MLsys, in pytorch



research project (publish)

semester project: due in the end of semester + a midway check in (end of March, 10%)

systems project

individual project

LLMs, in c/c++

MLsys, in pytorch



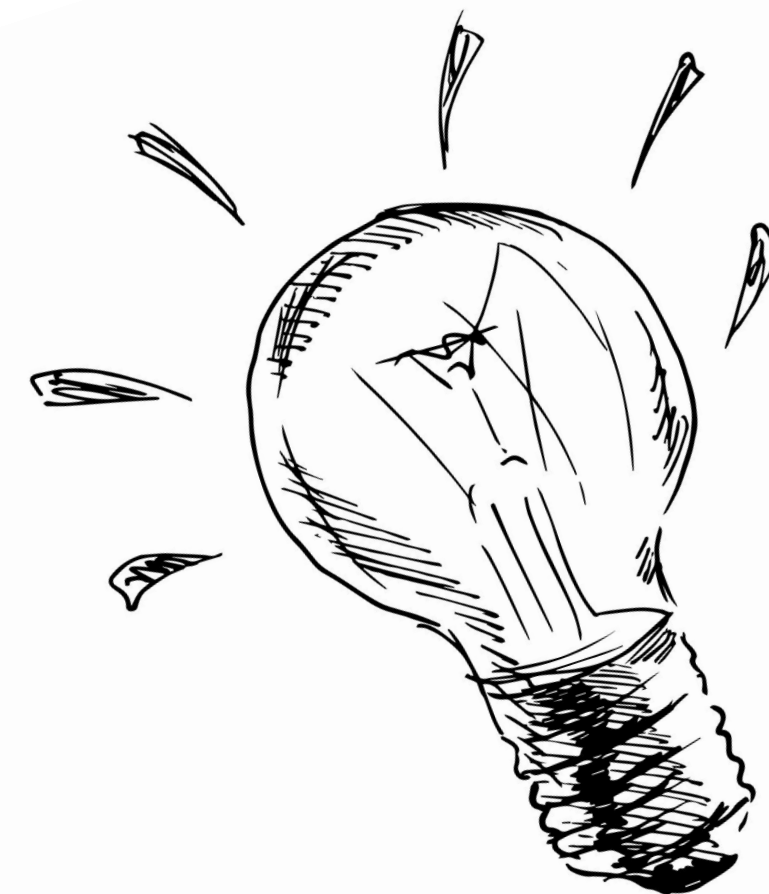
research project (publish)

groups of max three

Adaptivity/Performance

Focus this year:

LLM inference & Fine-tuning, RAG, Image AI



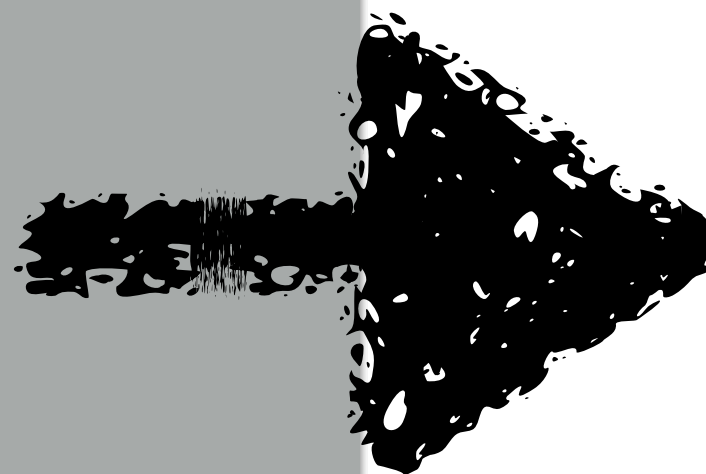
semester project: due in the end of semester + a midway check in (end of March, 10%)

systems project

individual project

LLMs, in c/c++

MLsys, in pytorch



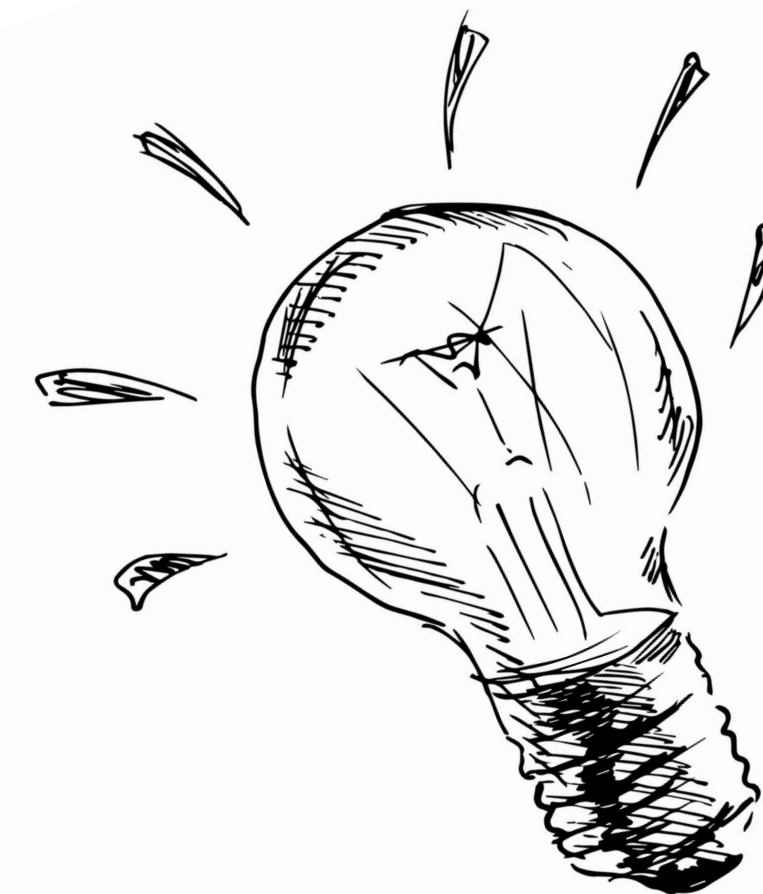
research project (publish)

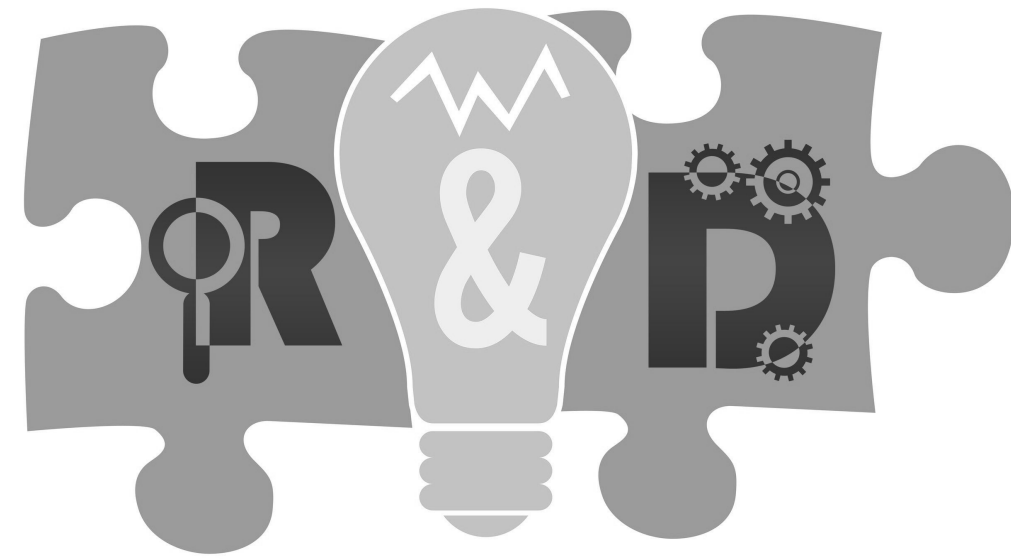
groups of max three

Adaptivity/Performance

Focus this year:

LLM inference & Fine-tuning, RAG, Image AI





ACM Special Interest Group In Data Management (SIGMOD)

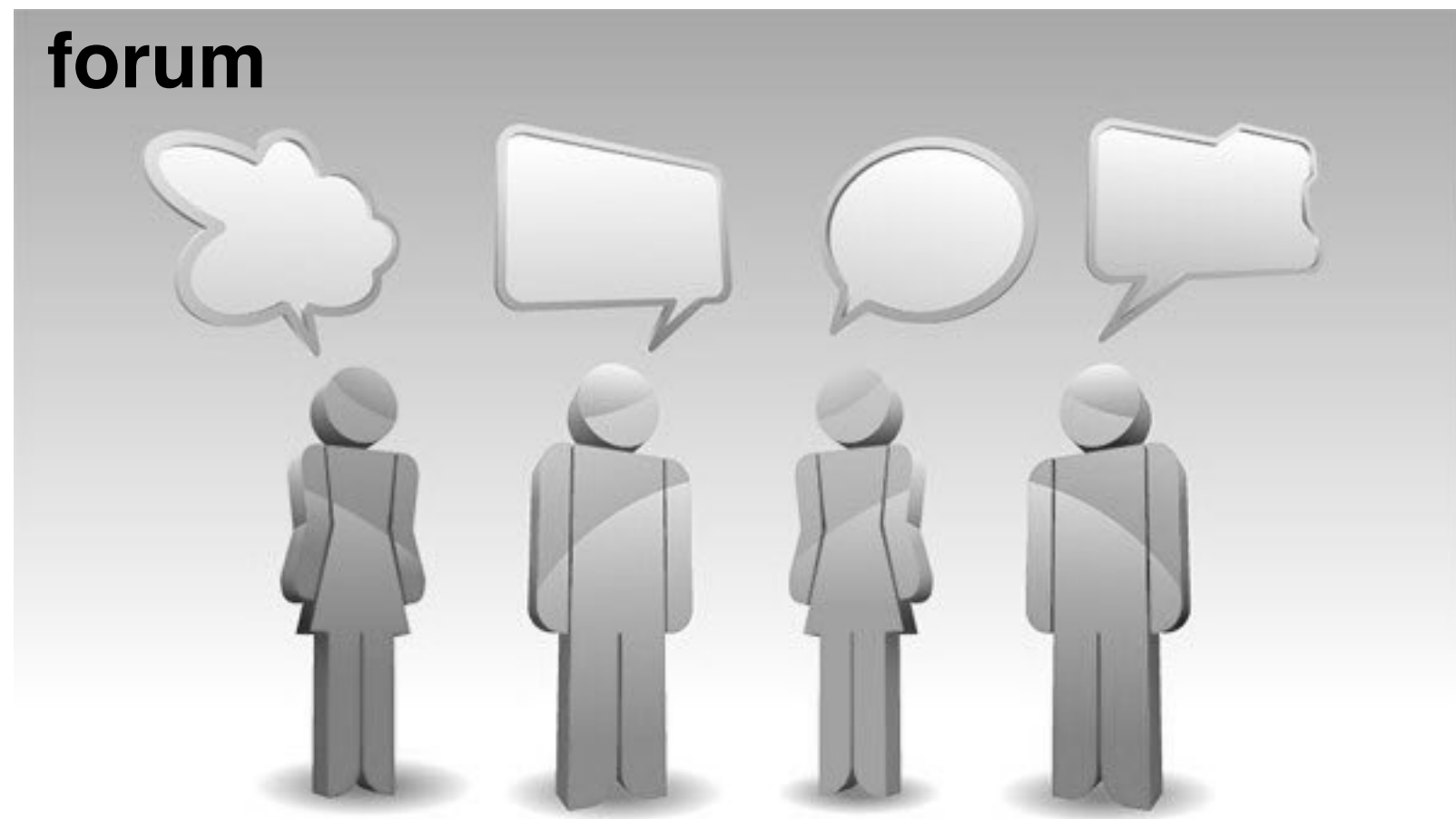
Undergrad Research Competition

first prize in 2016, 2017,
2018, 2019, 2020, 2022

Adaptive Denormalization
Evolving Trees
Splaying LSM-Trees
Adaptive NoSQL
Adaptive Filters
Distributed Deep Learning



forum

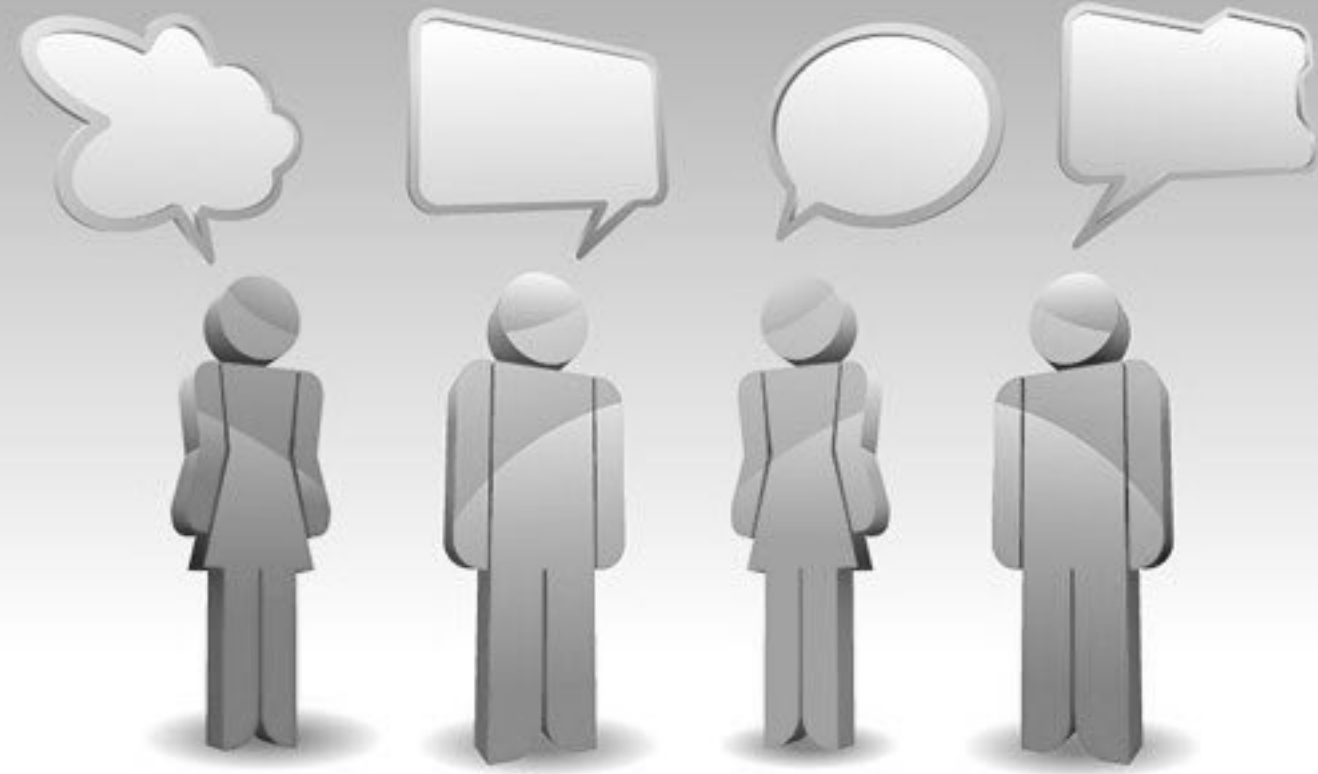


all announcements & discussions

as of week 2

link on class website - check out usage guidelines

forum



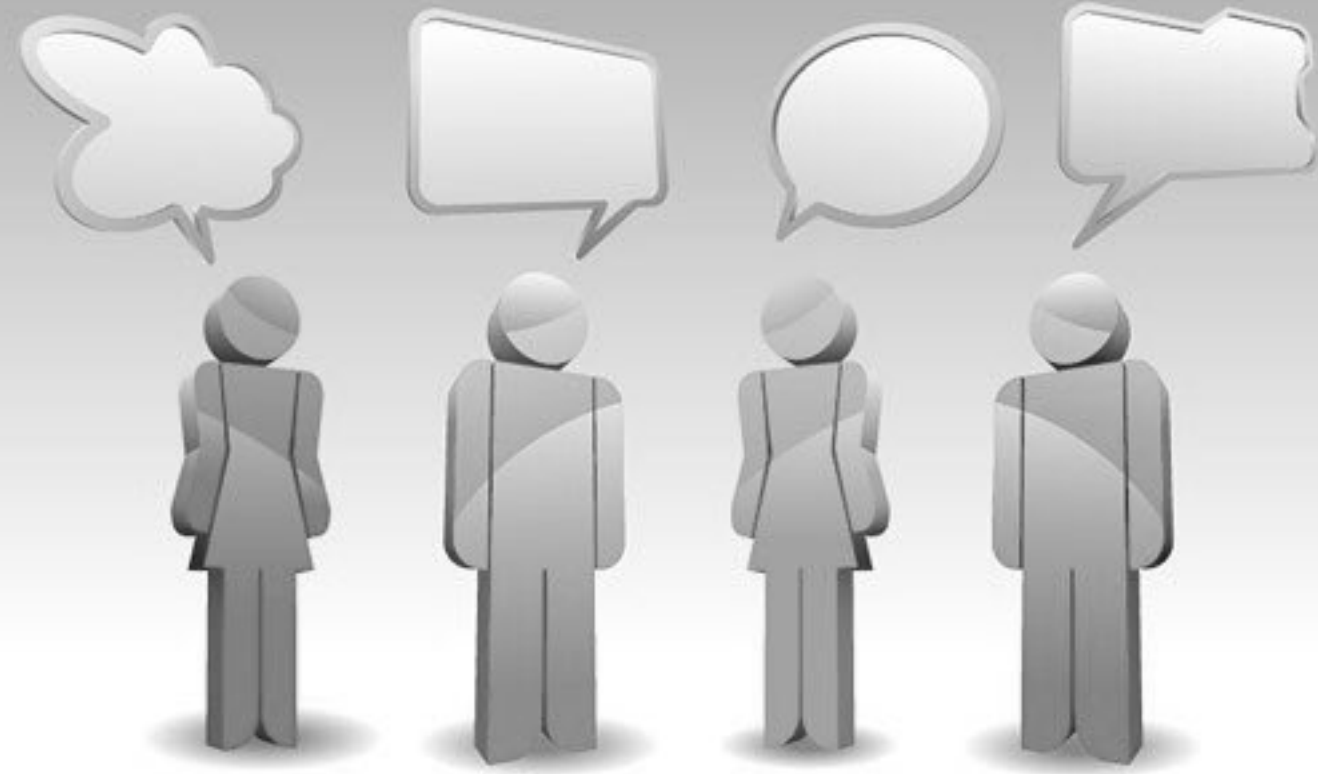
all announcements & discussions
as of week 2

link on class website - check out usage guidelines



classes are recorded
(links on canvas)

forum



all announcements & discussions
as of week 2

link on class website - check out usage guidelines

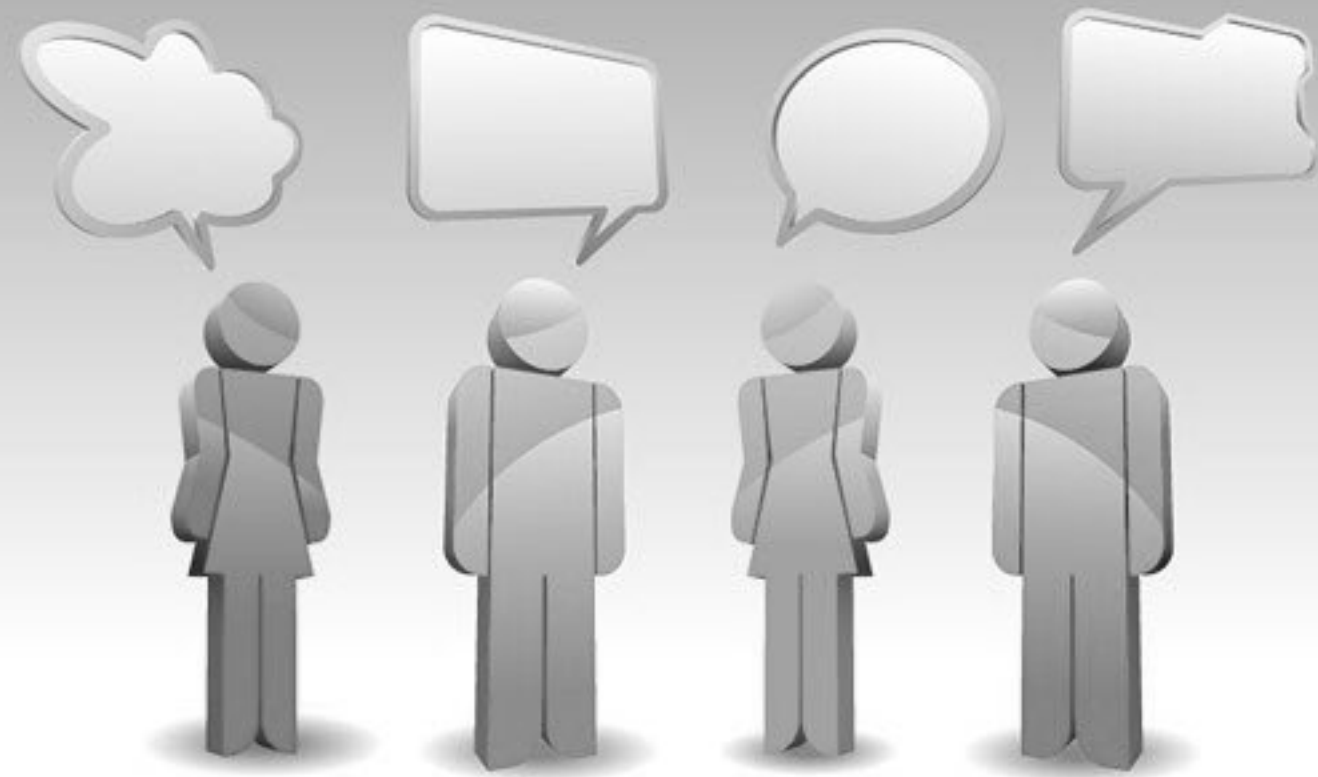


classes are recorded
(links on canvas)

Project: 45%
Midway Check-in: 10%
Discussion: 30%
Presentation: 15%

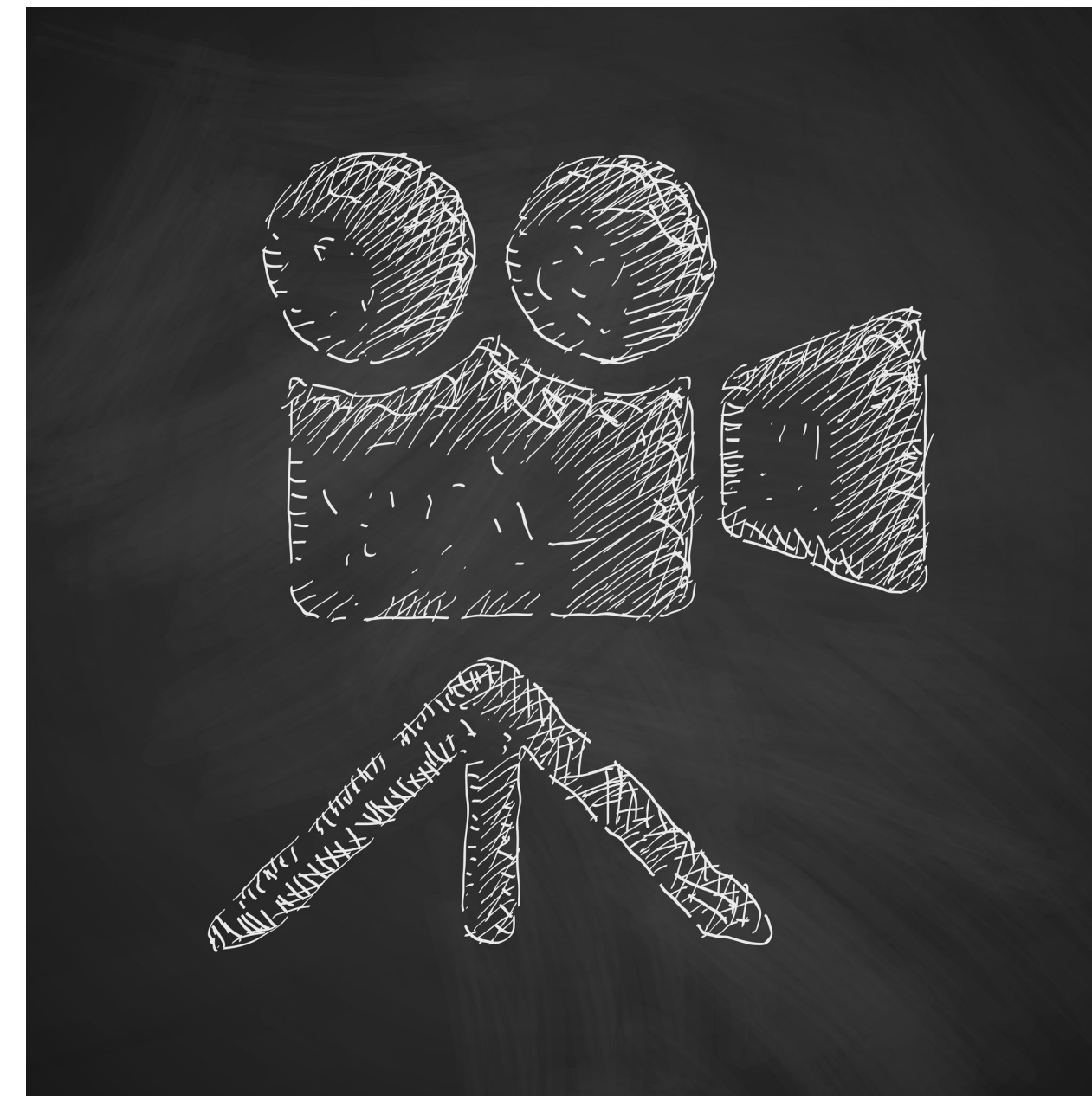


forum



all announcements & discussions
as of week 2

link on class website - check out usage guidelines



classes are recorded
(links on canvas)

Project: 45%
Midway Check-in: 10%
Discussion: 30%
Presentation: 15%



NO LAPTOP/PHONE POLICY
class is based on participation!

Teaching Fellows



Utku Sirin
Teaching Fellow
(Room: SEAS 4.435)



Konstantinos Kopsinis
Teaching Fellow
(Room: SEAS 4.435)



Qitong Wang
Teaching Fellow
(Room: SEAS 4.435)



Milad Rezaei Hajidehi
Teaching Fellow
(Room: SEAS 4.435)

Prerequisites

knowledge of algorithms, data structures, hardware, systems

Research track:

open to CS165 students

after discussion also CS161 and systems PhDs

Systems track allows taking the class without all prerequisites (but at least CS61)



Check out: syllabus,
project 0, systems projects, online sections

<http://daslab.seas.harvard.edu/classes/cs265/>

Timeline:

Research papers: 3rd week

New Systems Project: 3rd week

Research projects: 4th week

Plan to start systems/research project end of Feb

Stratos' OH start today - Labs to start on Week 3

CS 265

Big Data & AI Systems

NoSQL | Neural Networks | Image AI | LLMs | Data Science