

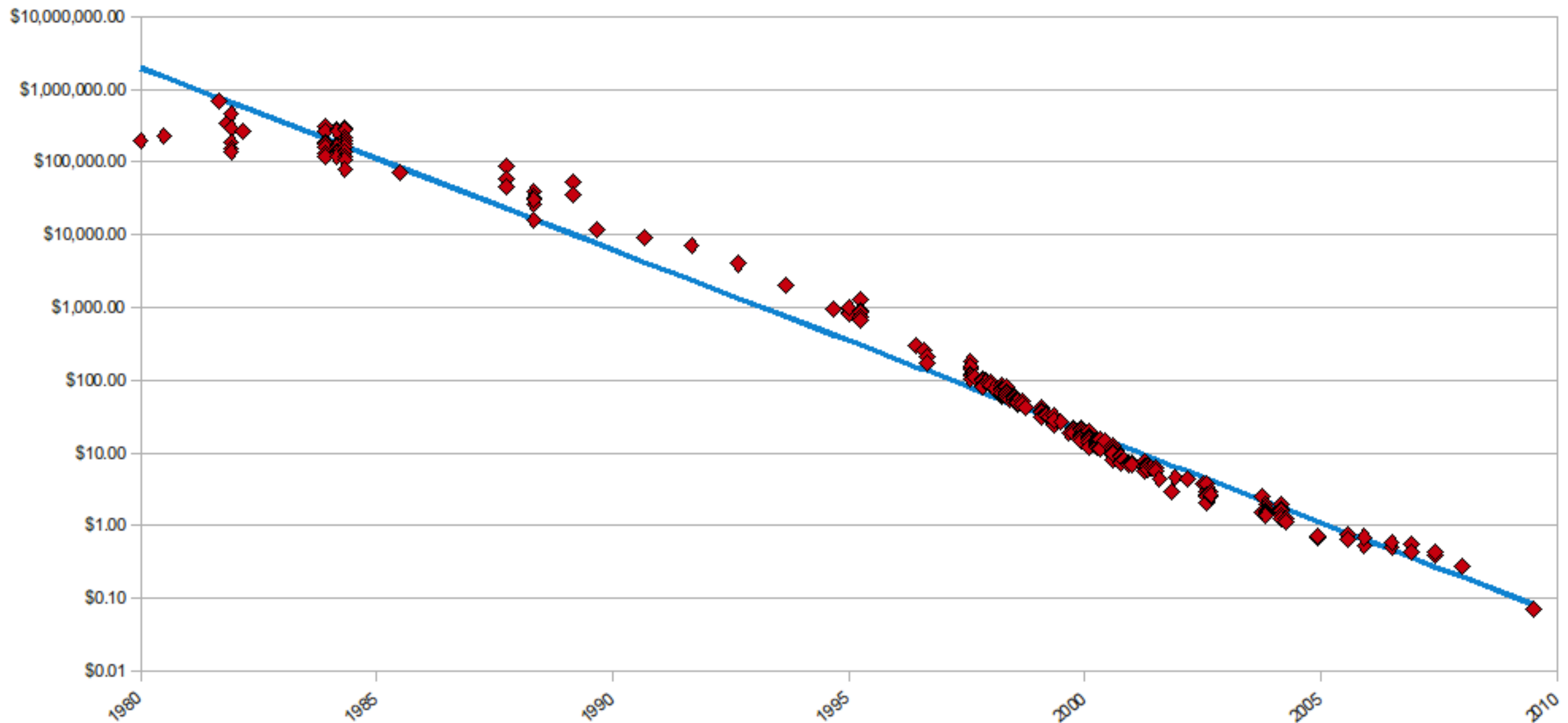
Computer Architecture Trends That Affect Data Analytics



CHIDATA

Last Lecture

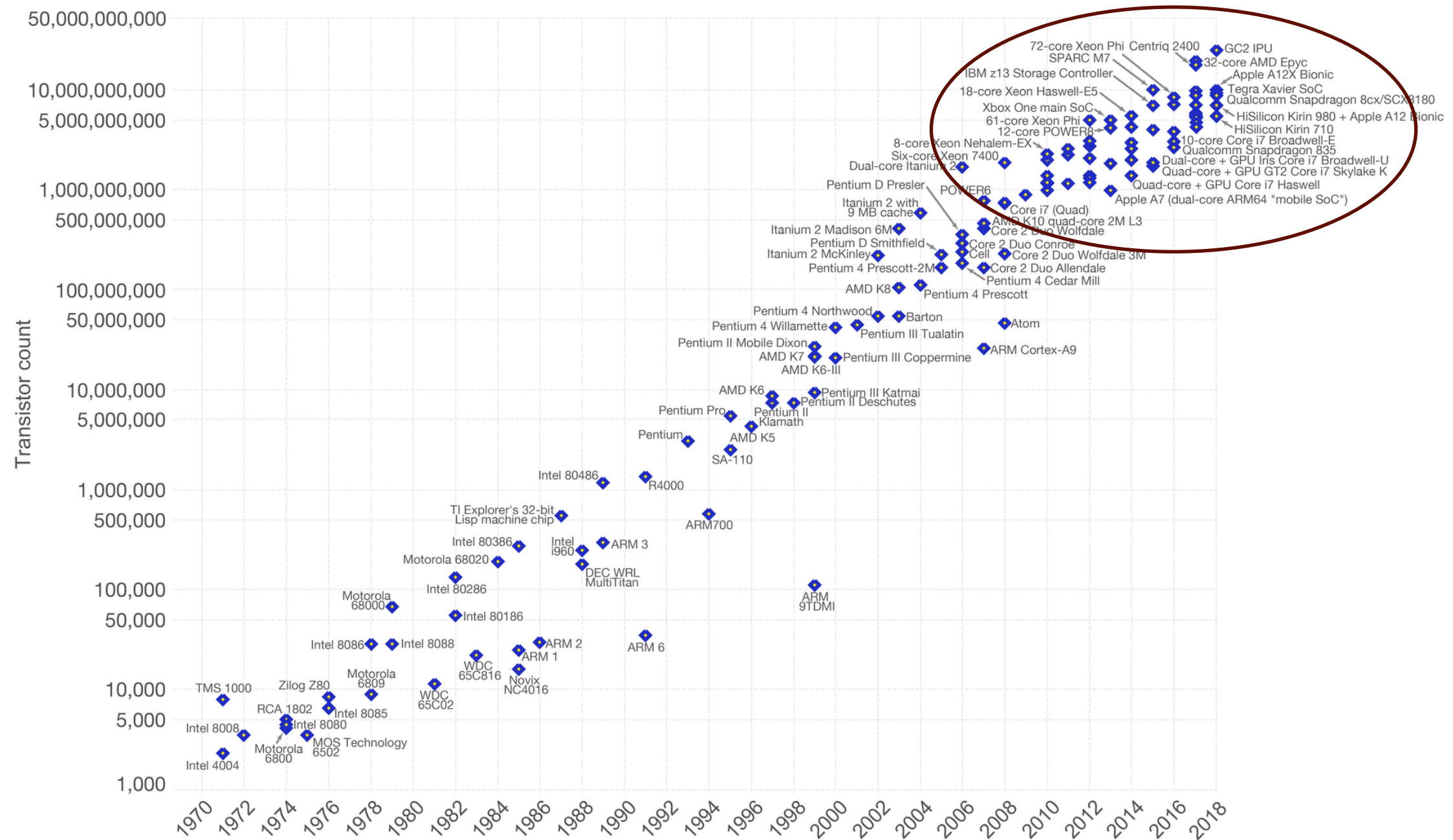
Hard Drive Cost per Gigabyte
1980 - 2009



Last Lecture

Moore's Law – The number of transistors on integrated circuit chips (1971-2018)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are linked to Moore's law.

Our World
in Data

This Lecture

CPU = Fast

Moving Data = Slow

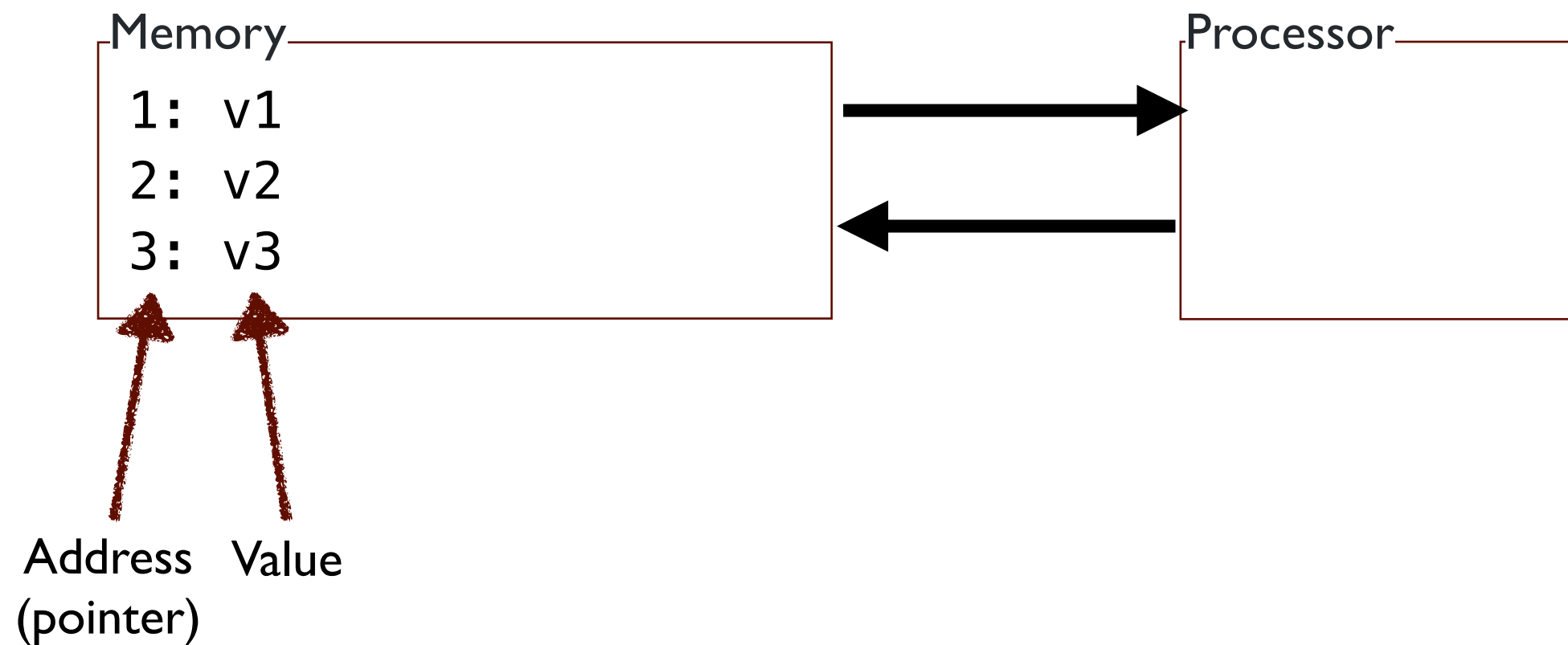
This Lecture

The “storage” story is actually much more complicated.

Data are typically stored in a hierarchy of slow but big to fast but small storage systems.

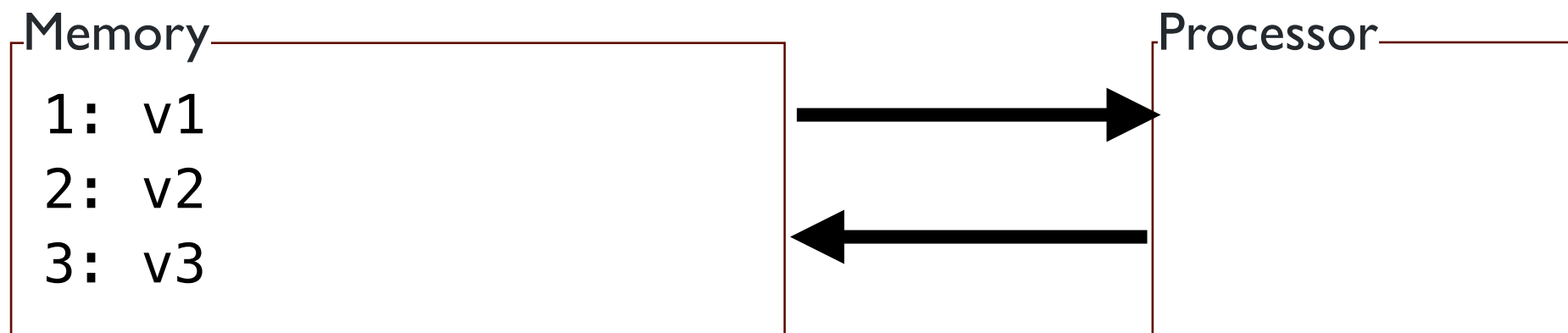
Overview of the changes on the horizon.

A Basic Computing Model



A Basic Computing Model

read(address)

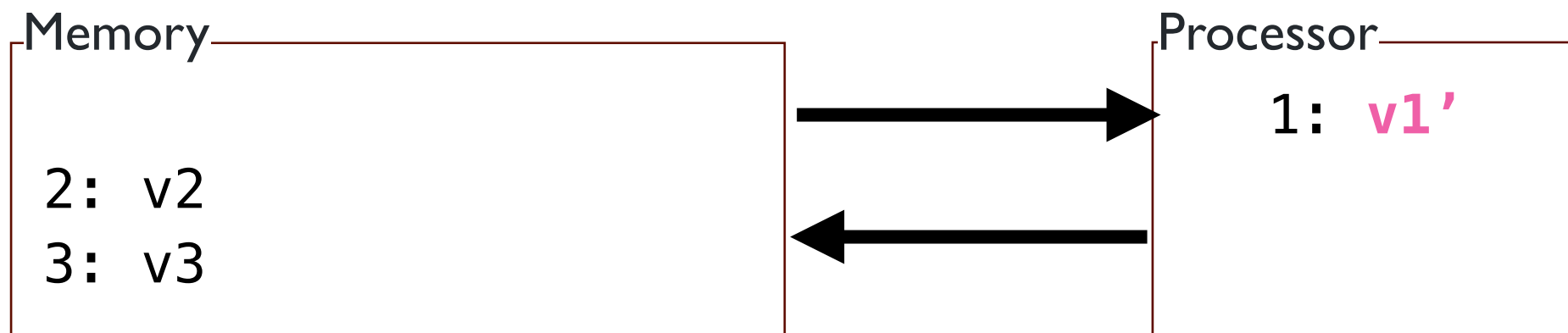


Move data from memory to processor

Processor has a limited amount of native memory

A Basic Computing Model

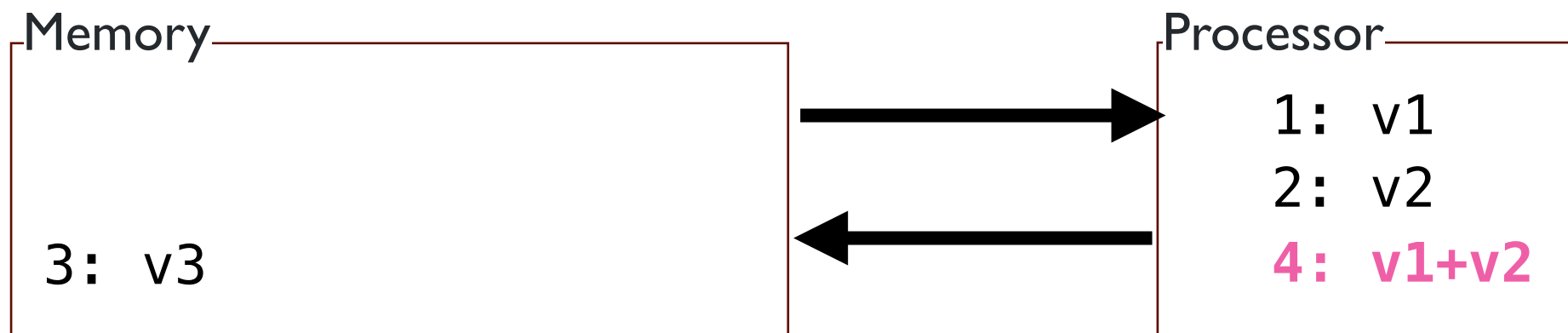
`inst(value)`



Apply a basic instruction to the data

A Basic Computing Model

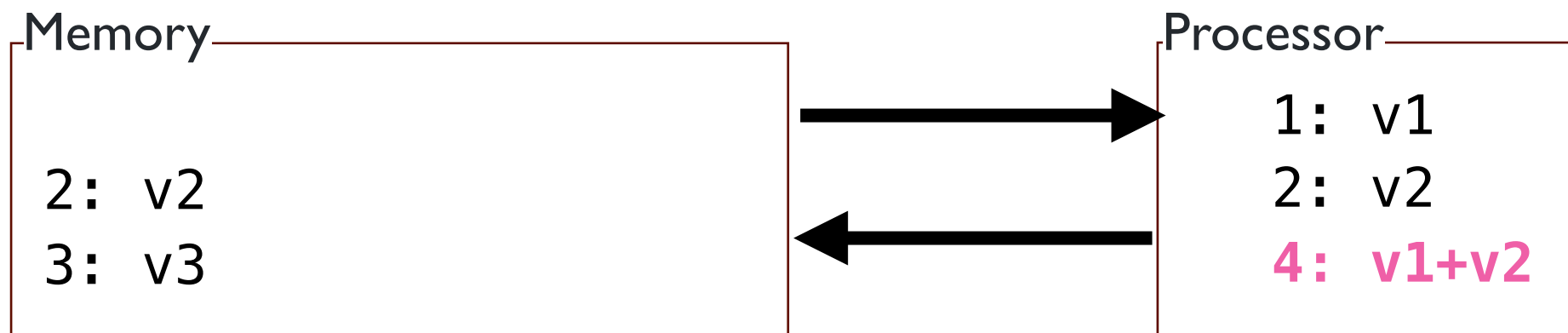
`inst(value)`



Apply a basic instruction to the data

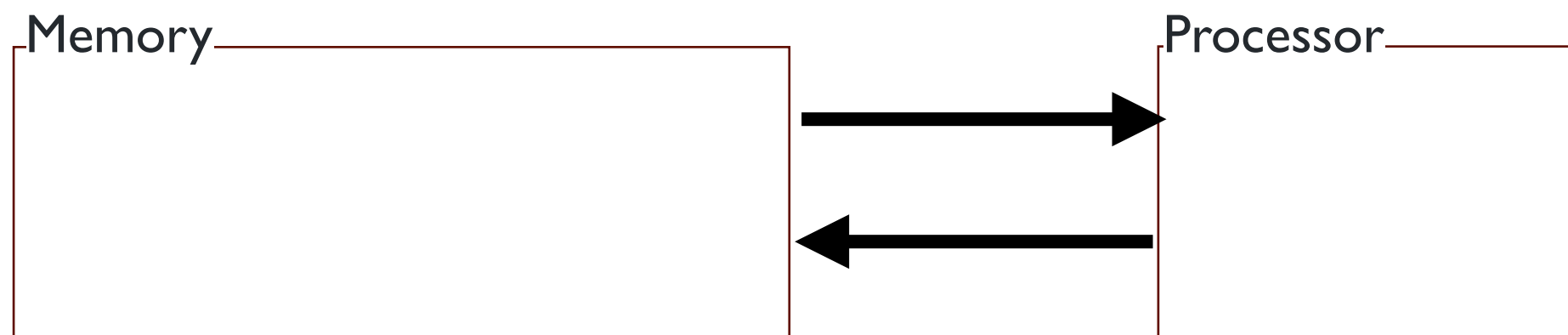
A Basic Computing Model

`write(address)`



Write data back to memory

Run Time



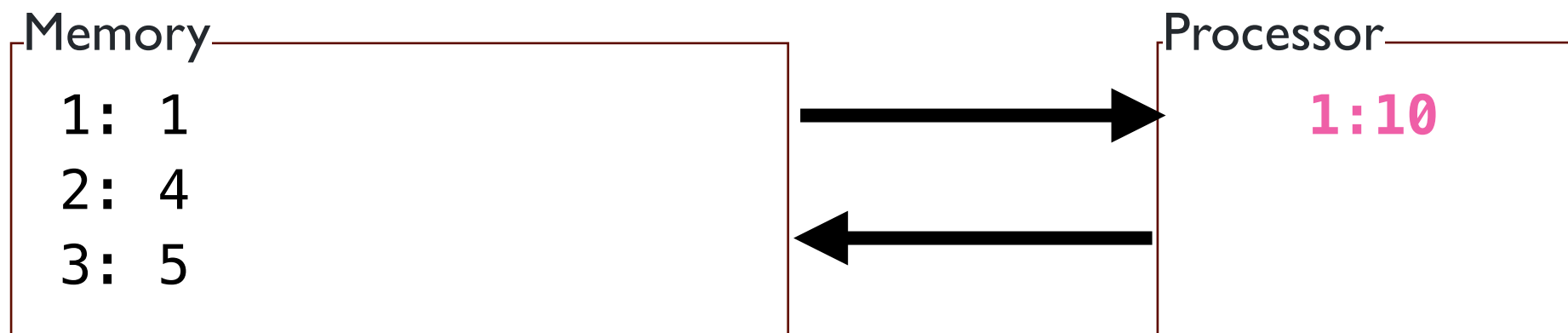
$$\text{Runtime} = \text{I/O Time} + \text{Compute Time}$$

↑
Total
read()/write()

↑
Total
inst()

Calculate I/O and CPU Costs

Given a list of numbers [1,4,5,...,6] calculate the sum:



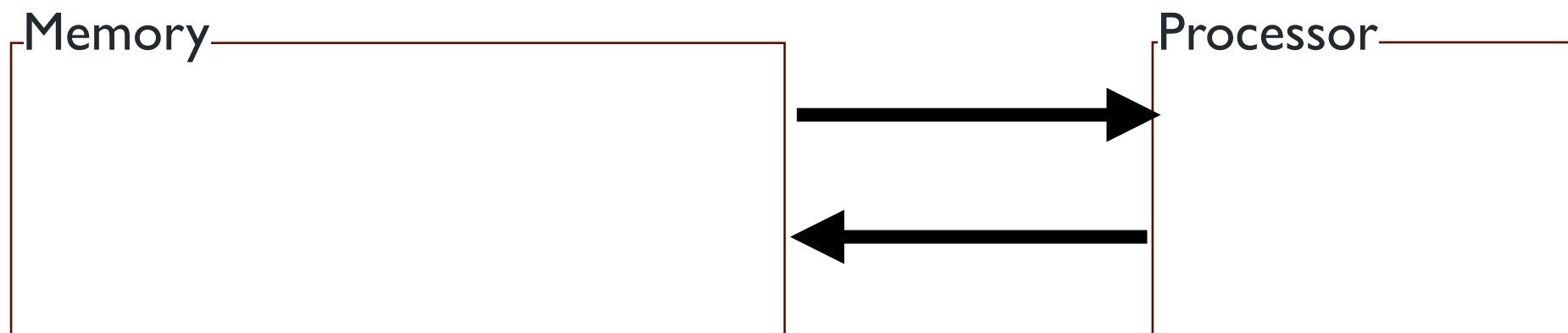
I read() + I sum() per data item

I write() for final answer

I/O Cost: $N+1$, CPU Cost: $N-1$

Calculate I/O and CPU Costs

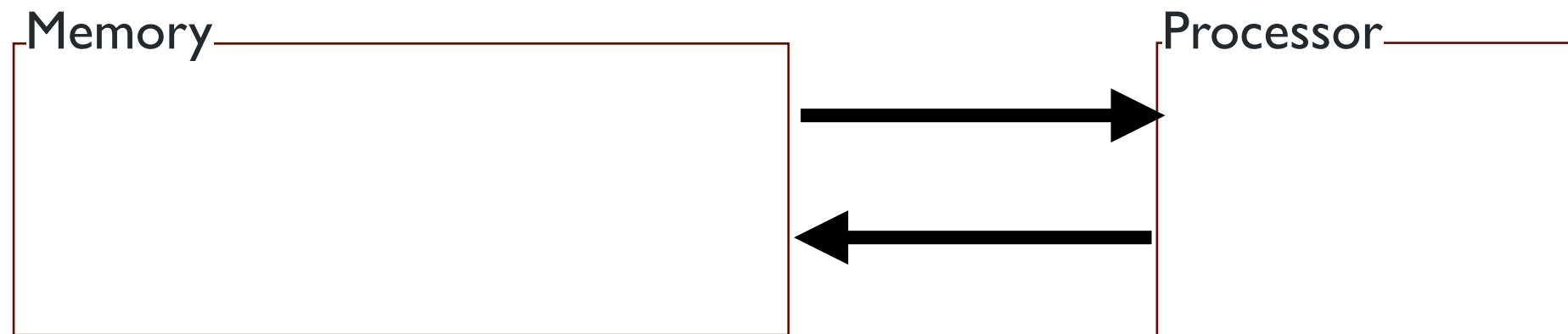
Given a list of numbers $[1, 4, 5, \dots, 6]$ calculate the sum:



I/O Cost: $_N$, CPU Cost: $_N$

Generally care about orders and ignore small constants

Rate-Limiting Operations



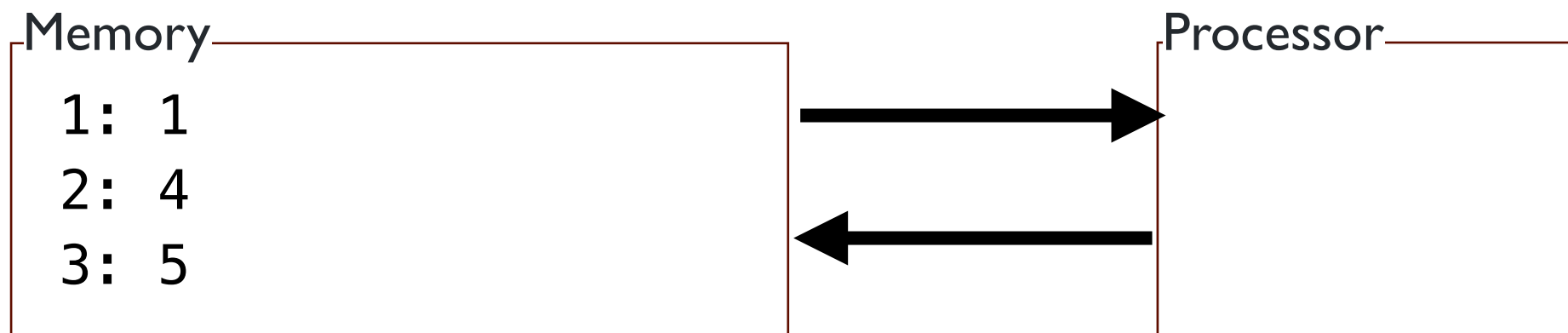
Runtime = I/O Time + Compute Time

Runtime = $a \cdot \text{IO Steps} + b \cdot \text{CPU Steps}$

a is usually MUCH bigger than b ! (thousands of times!)

Calculate I/O and CPU Costs

Given a list of numbers [1,4,5,...,6] calculate the sum of the even numbers:

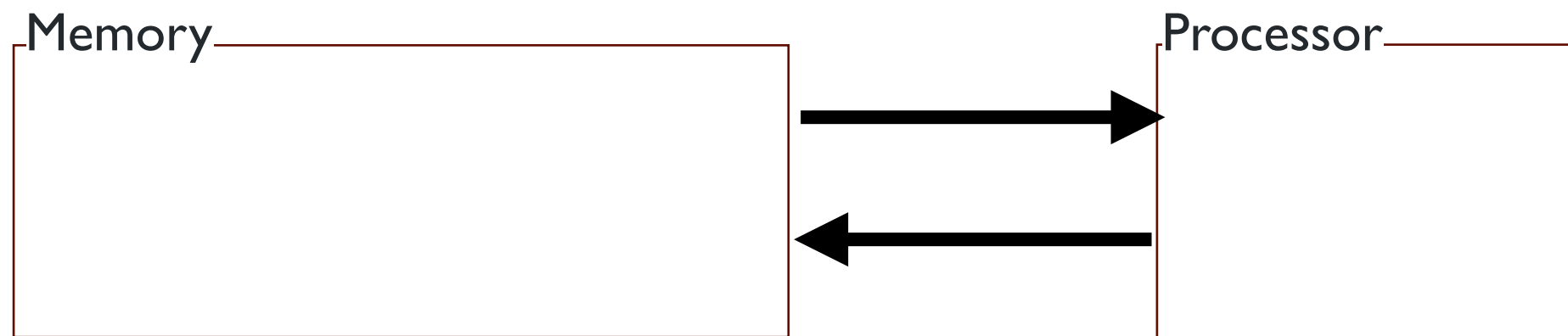


$I \text{ read}() + I \text{ mod}() + I \text{ branch}() + \text{at most } I \text{ sum}() \text{ per data item}$

$I \text{ write}() \text{ for final answer}$

I/O Cost: N , CPU Cost: $3N$

Rate-Limiting Operations



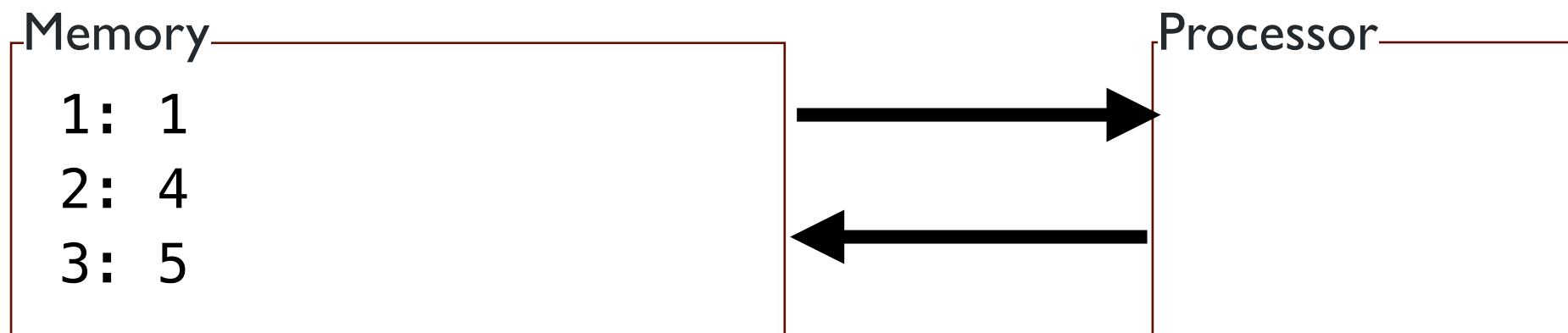
Runtime = I/O Time + Compute Time

“I/O Bound” : I/O Time \gg Compute Time

“CPU Bound” : Compute Time \gg I/O Time

Calculate I/O and CPU Costs

Given the median of $[1, 4, 5, \dots, 6]$



What if the processor can store $>N/2$ data points?

I/O V.S CPU Bound

“I/O Bound” : Simple operations over lots of data

- Sorting a list of numbers. (Operation: Comparison)
- Summing a list of numbers (Operation: Sum)
- Finding a number in the list that is less than 5 (Operation: Comparison)

“CPU Bound” : Complex repetitive over small working sets

- Matrix multiplication
- Cryptographic hashing
- Image processing

Data Analytics is (usually) I/O Bound!

“I/O Bound” : Simple operations over lots of data

- Sorting a list of numbers. (Operation: Comparison)
- Summing a list of numbers (Operation: Sum)
- Finding a number in the list that is less than 5 (Operation: Comparison)

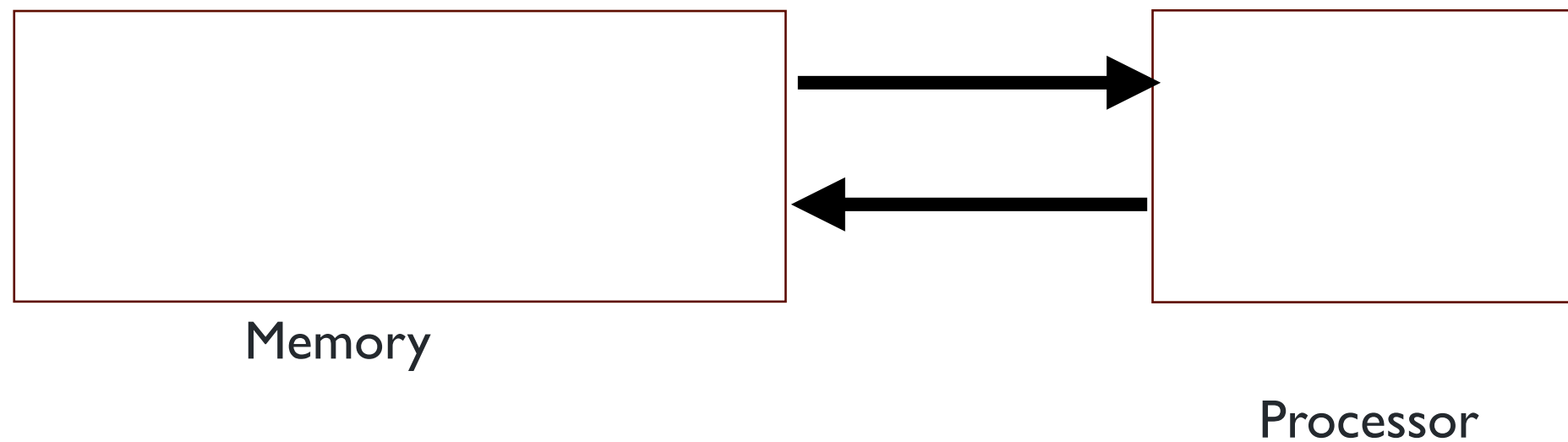
“CPU Bound” : Complex repetitive over small working sets

- Matrix multiplication
- Cryptographic hashing
- Image processing

Optimizing I/O Bound Processes

“I/O Bound” : Simple operations over lots of data

Finding a number in the list that is less than 5 (Operation: Comparison)

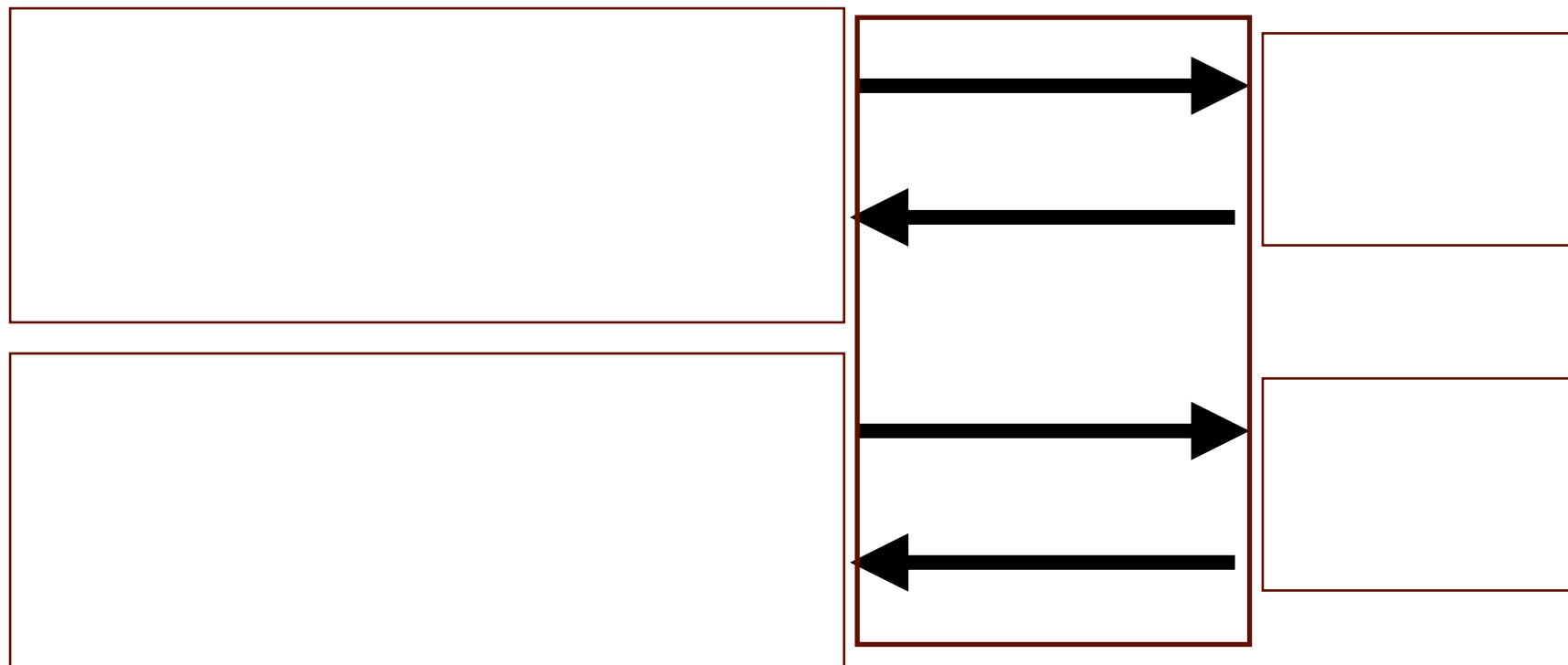


Not as useful! I/O Cost Still N

Optimizing I/O Bound Processes

“I/O Bound” : Simple operations over lots of data

Finding a number in the list that is less than 5 (Operation: Comparison)



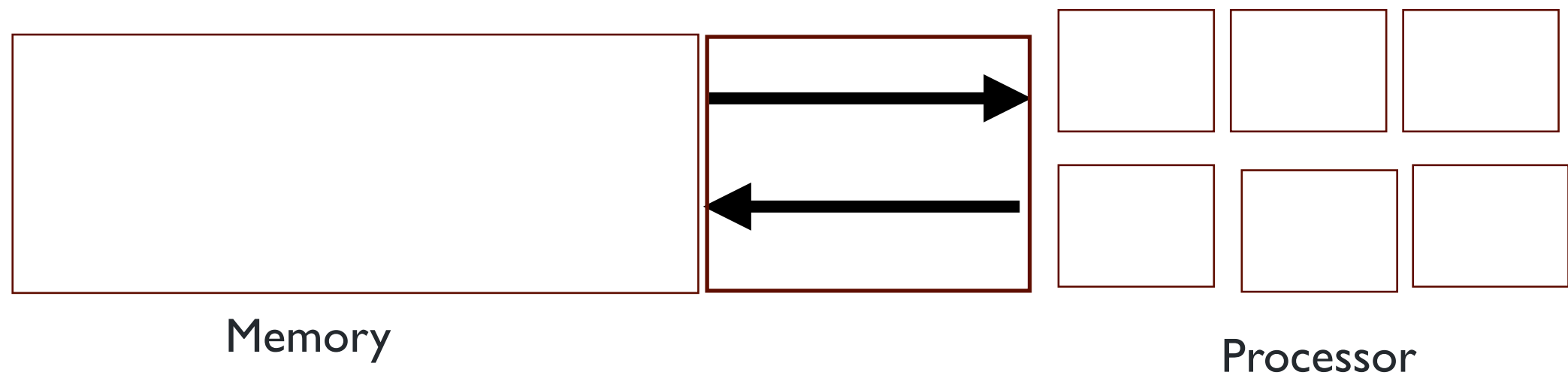
Better Strategy, I/O constant goes down!

Equivalent Serial I/O Cost of $N/2^$*

Optimizing I/O Bound Processes

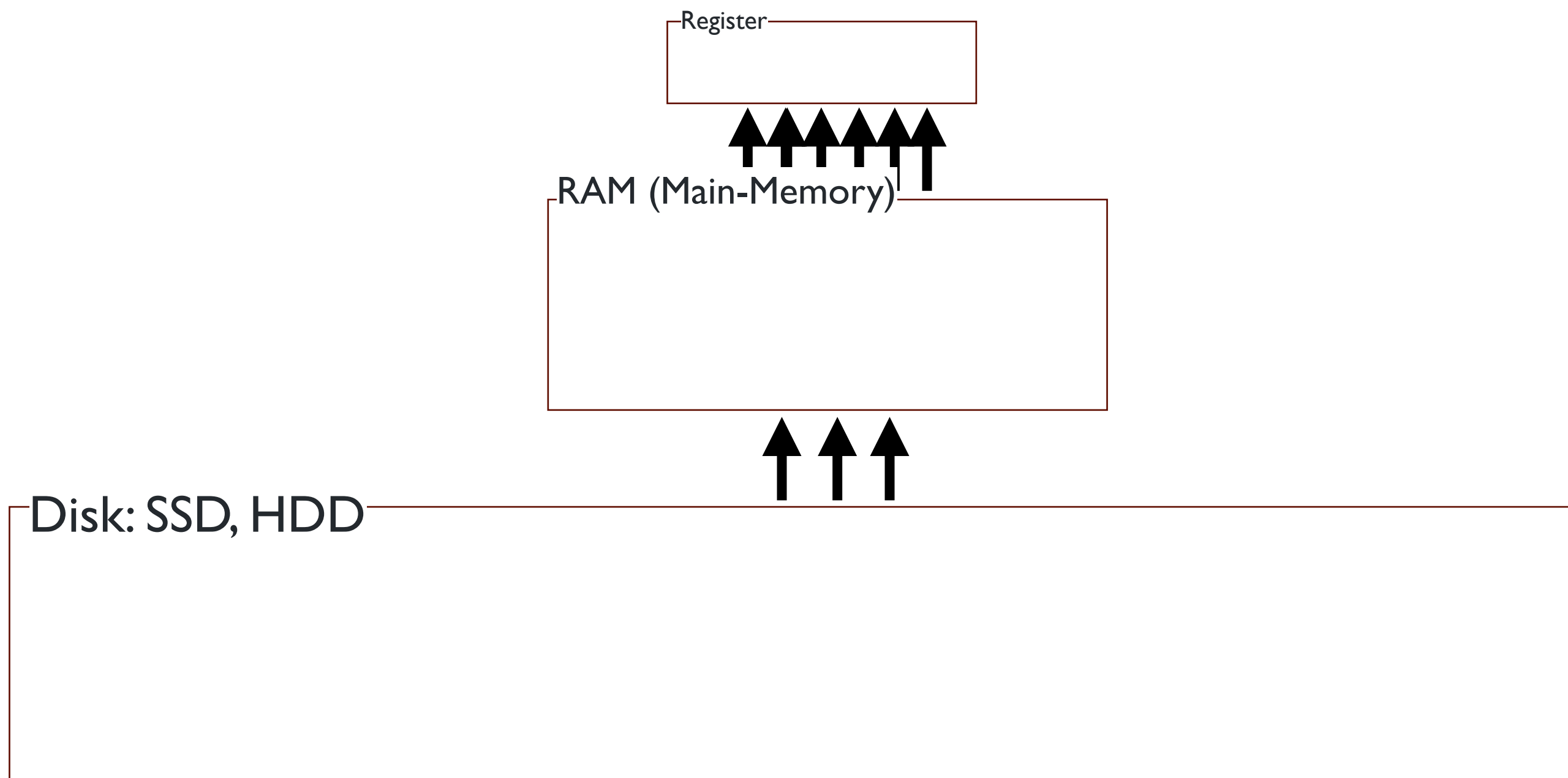
“I/O Bound” : Simple operations over lots of data

Finding a number in the list that is less than 5 (Operation: Comparison)

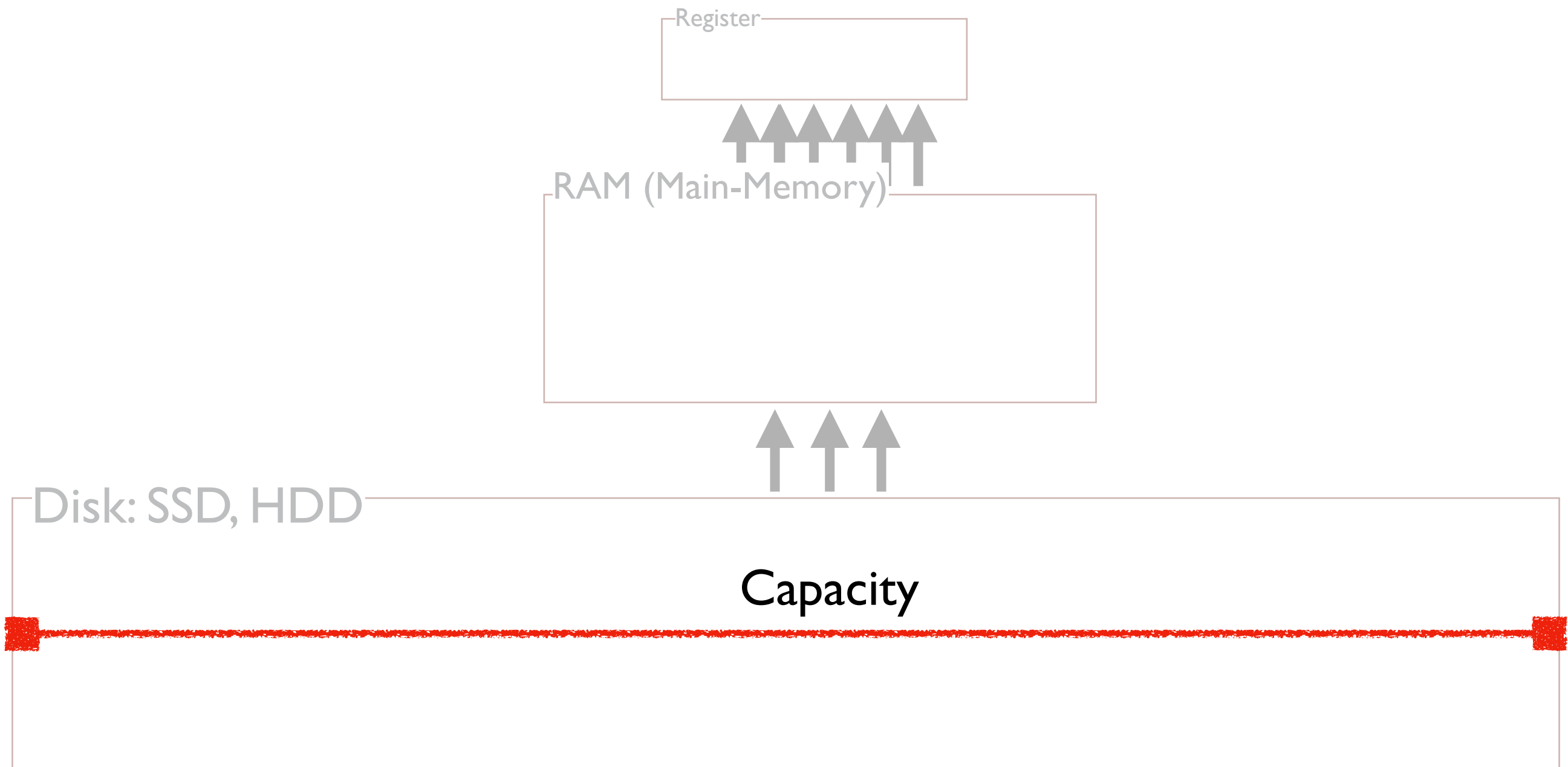


Depends on how much I/O can be parallelized

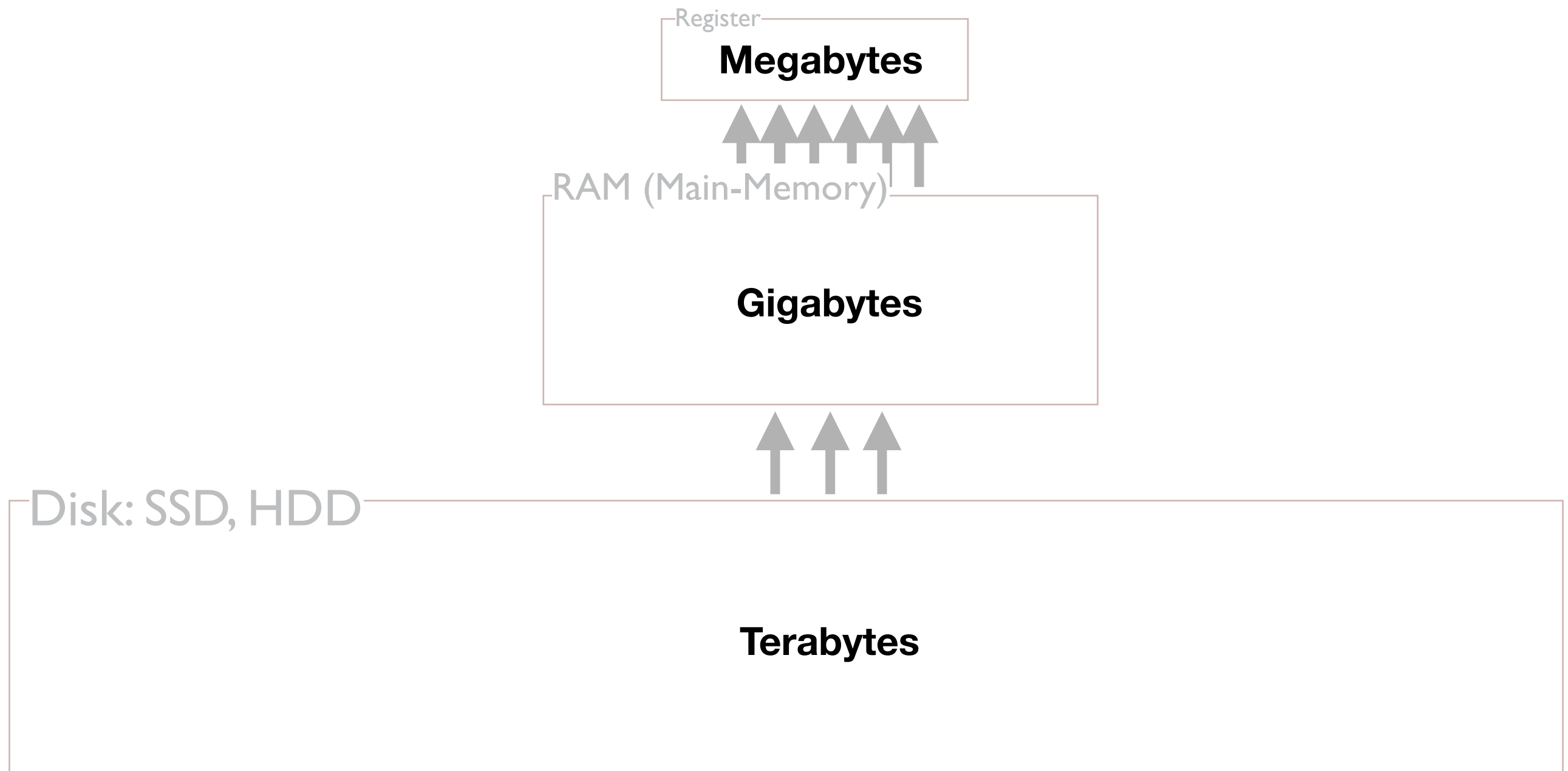
There are more levels!



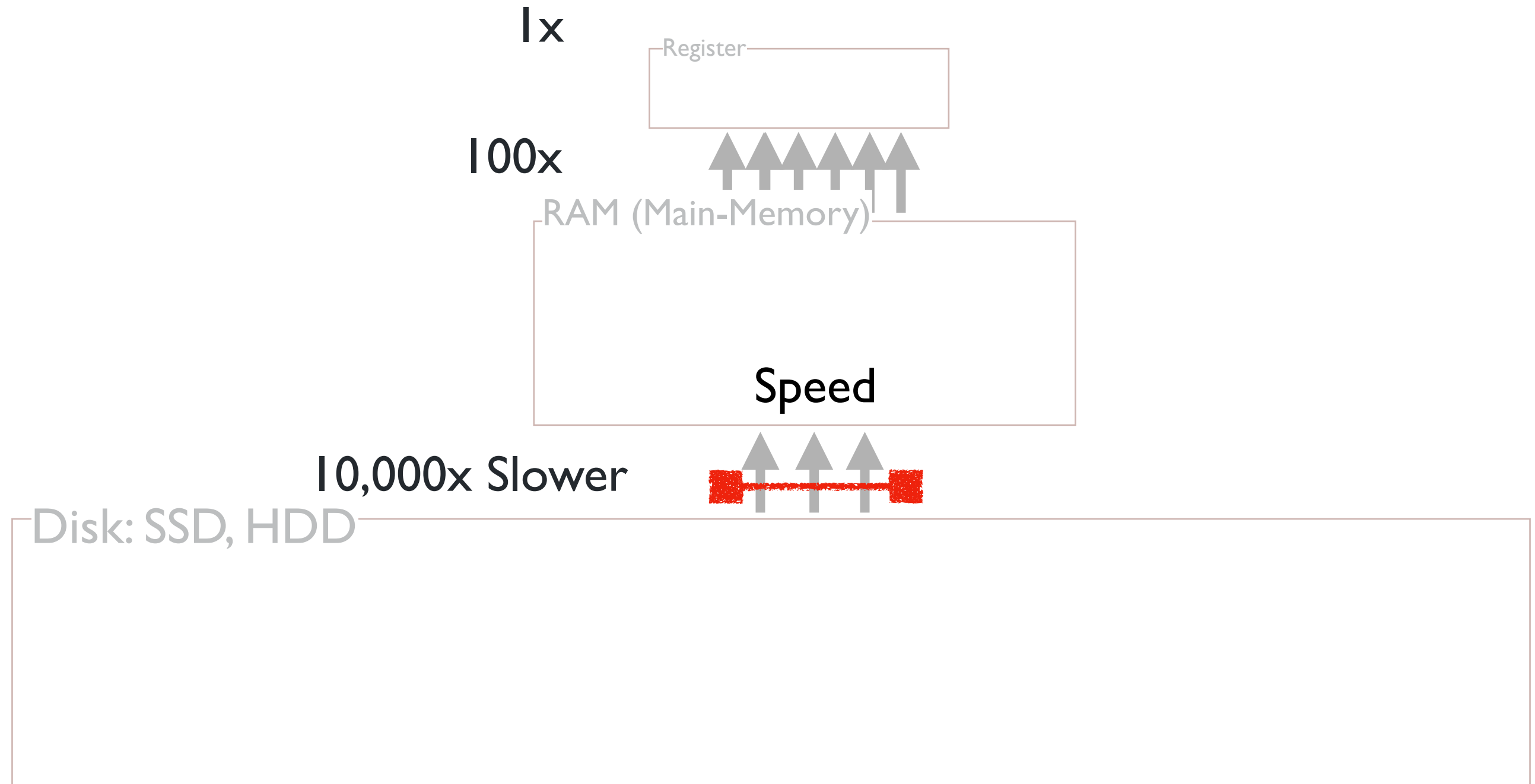
Schematic of Computer Storage



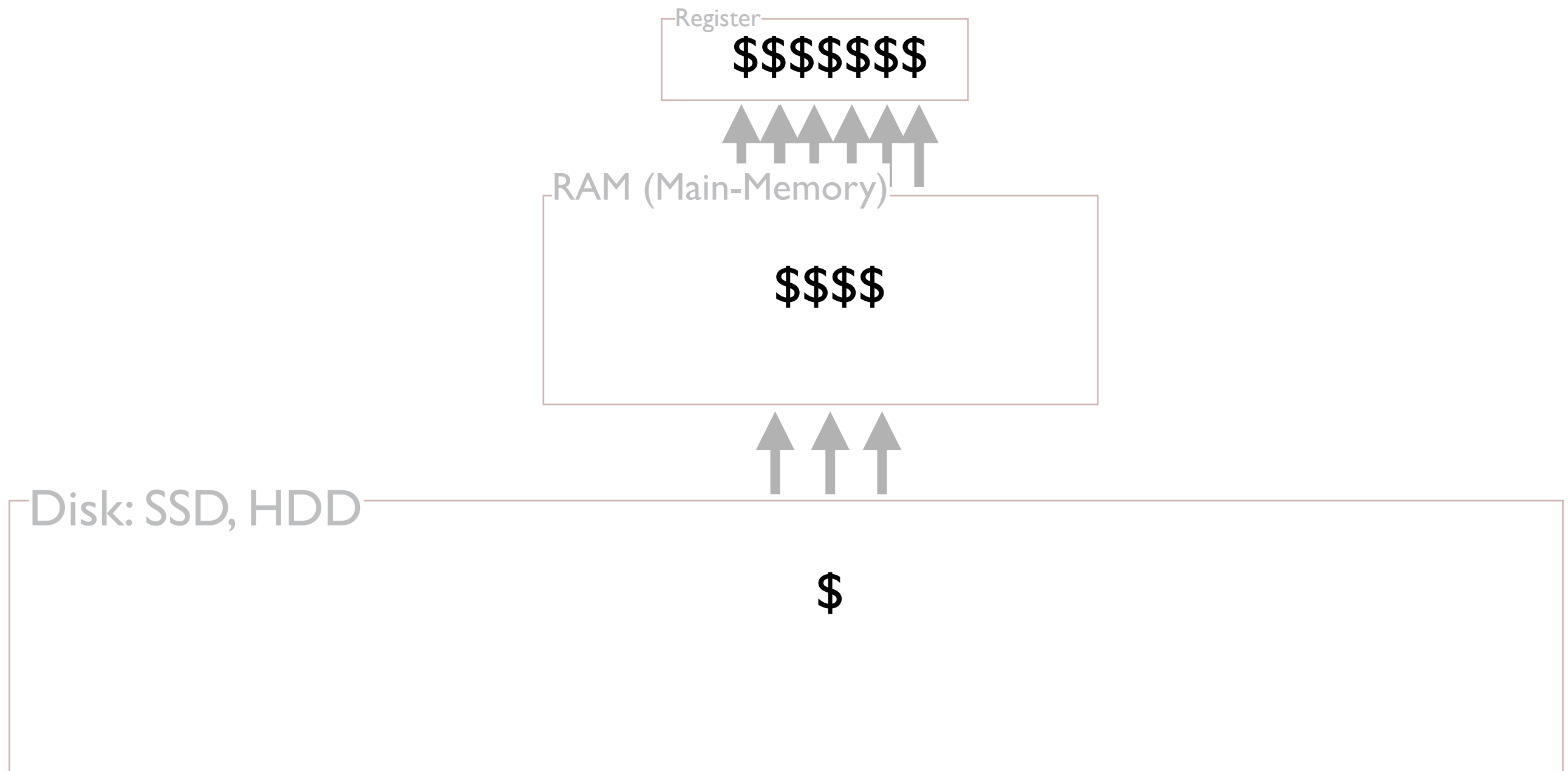
Schematic of Computer Storage



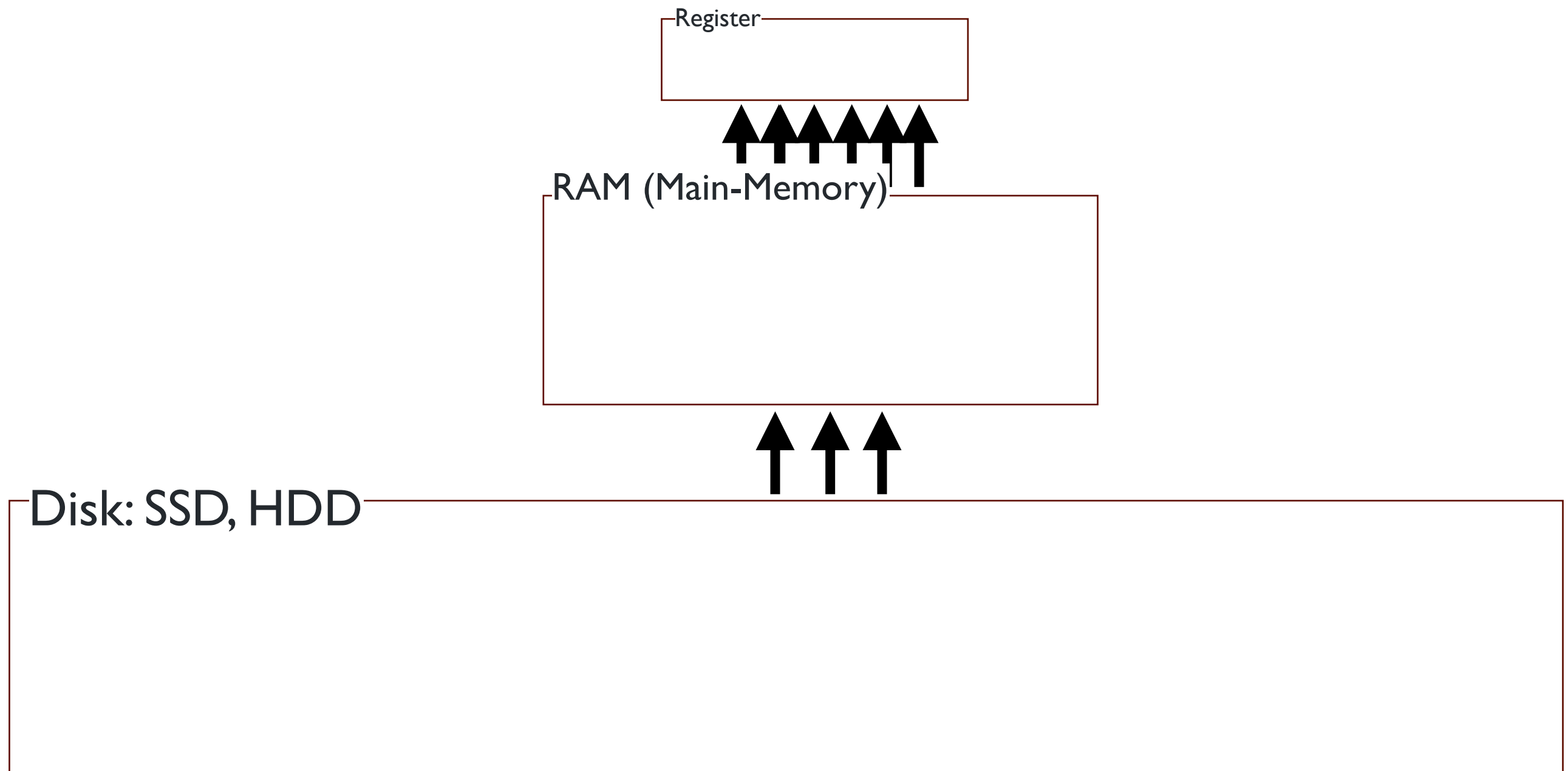
Schematic of Computer Storage



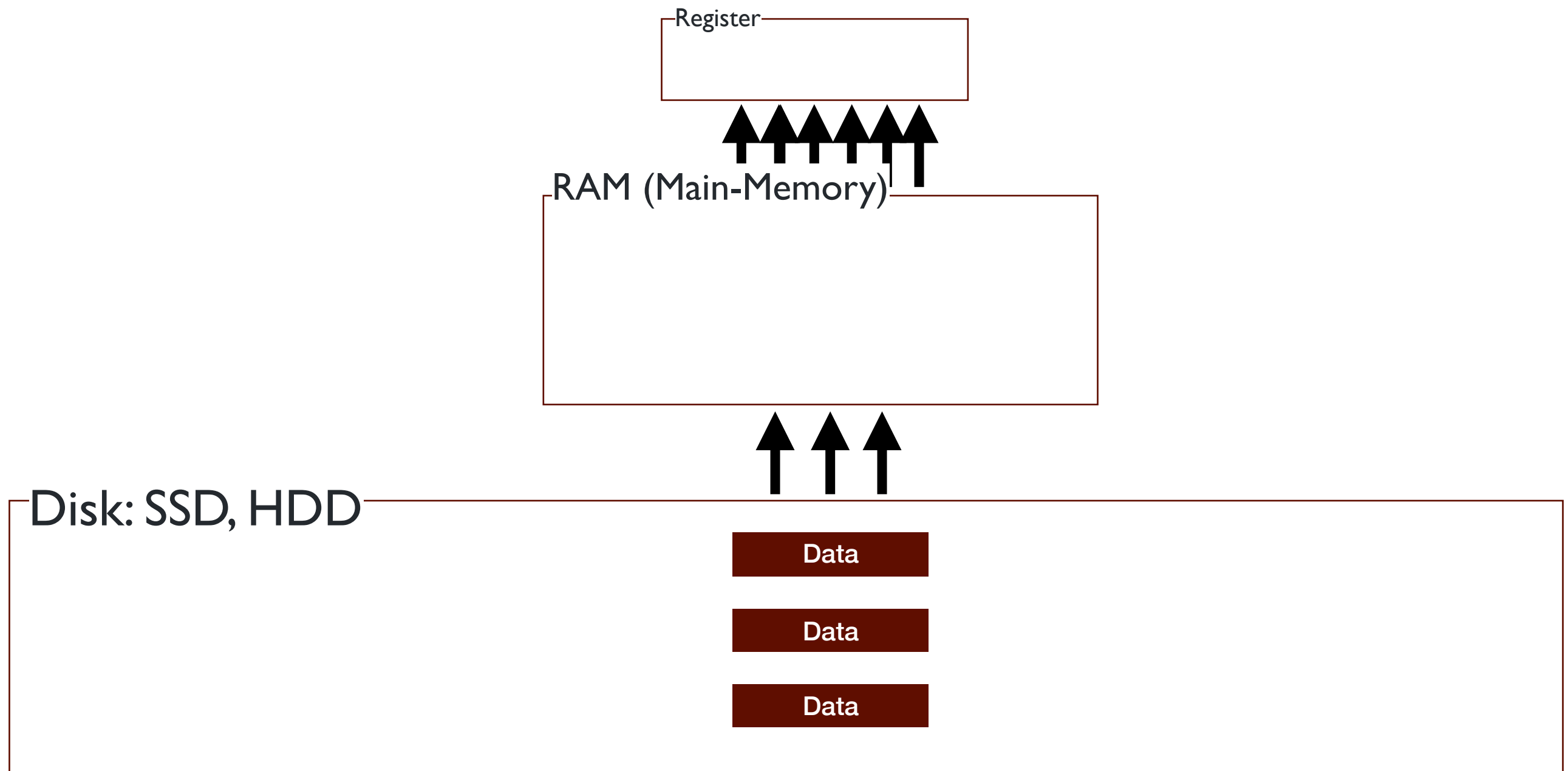
Schematic of Computer Storage



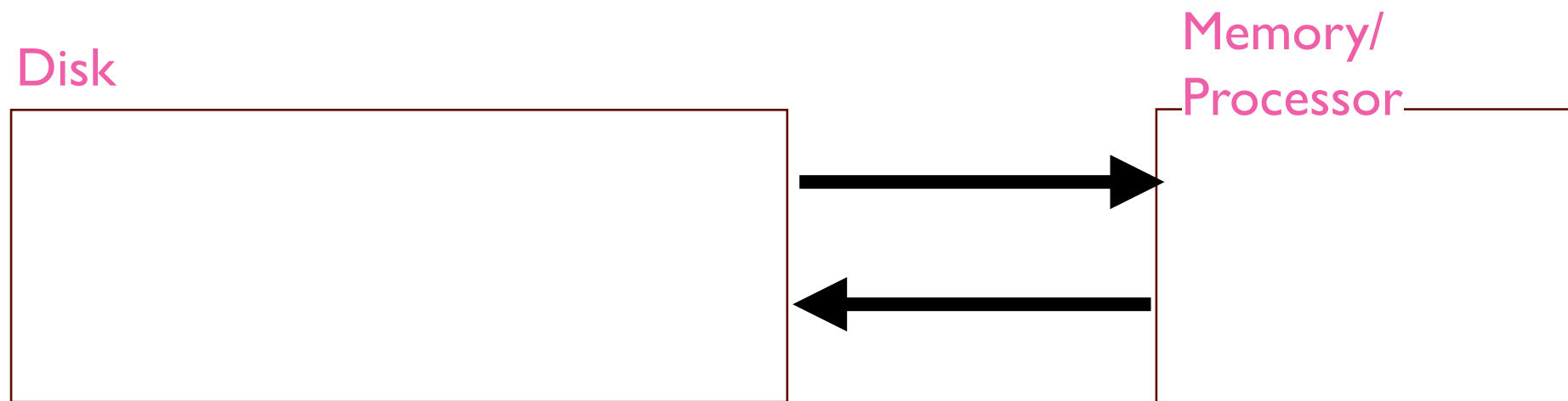
Schematic of Computer Storage



Data Must Move For Analysis



Rate-Limiting Operations

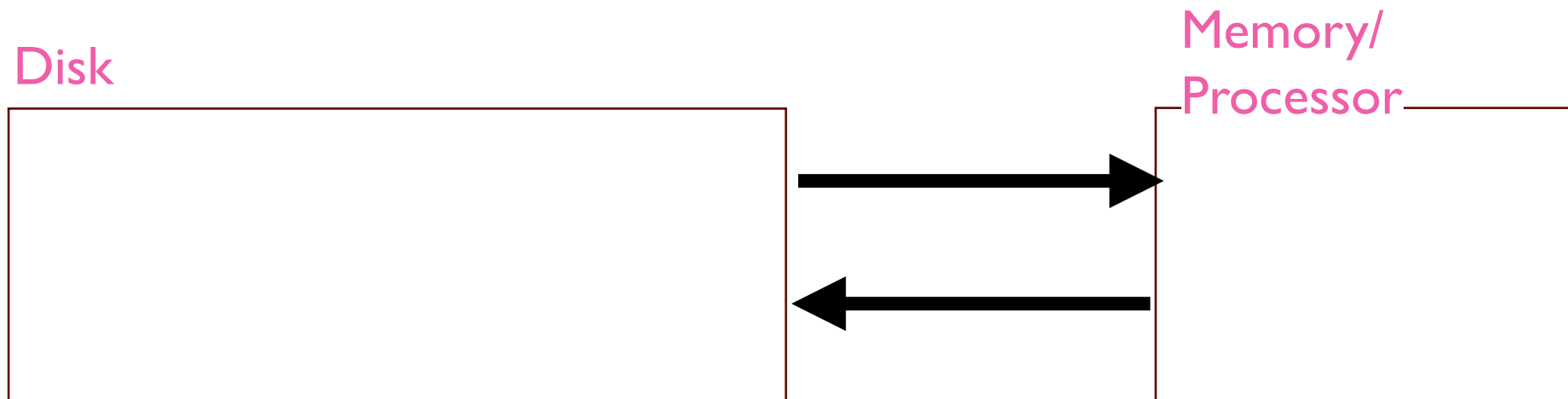


Runtime = I/O Time + Compute Time

“I/O Bound” : I/O Time \gg Compute Time

“CPU Bound” : Compute Time \gg I/O Time

Rate-Limiting Operations



Given the median of $[1, 4, 5, \dots, 6]$

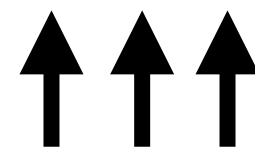
Memory $> N/2 \Rightarrow$ I/O Cost $\sim N$

Memory $< N/2 \Rightarrow$ I/O Cost $\sim N \log N$

Effect is called “spilling”

Physical Design

Finding a number in the list that is less than 5 (Operation: Comparison)
Median of a list of numbers



10,000x Slower

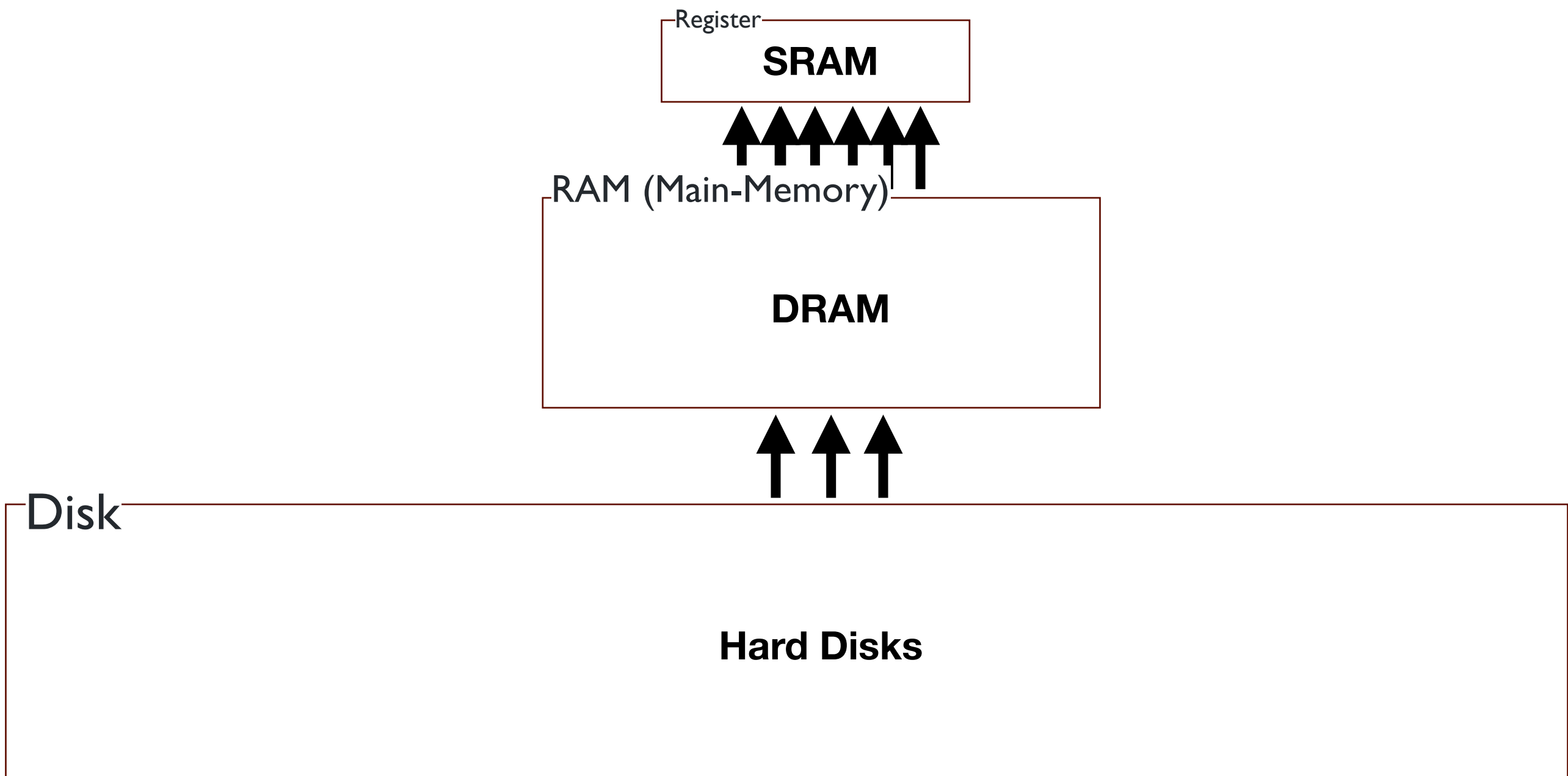
Disk: SSD, HDD

Stored data sorted!

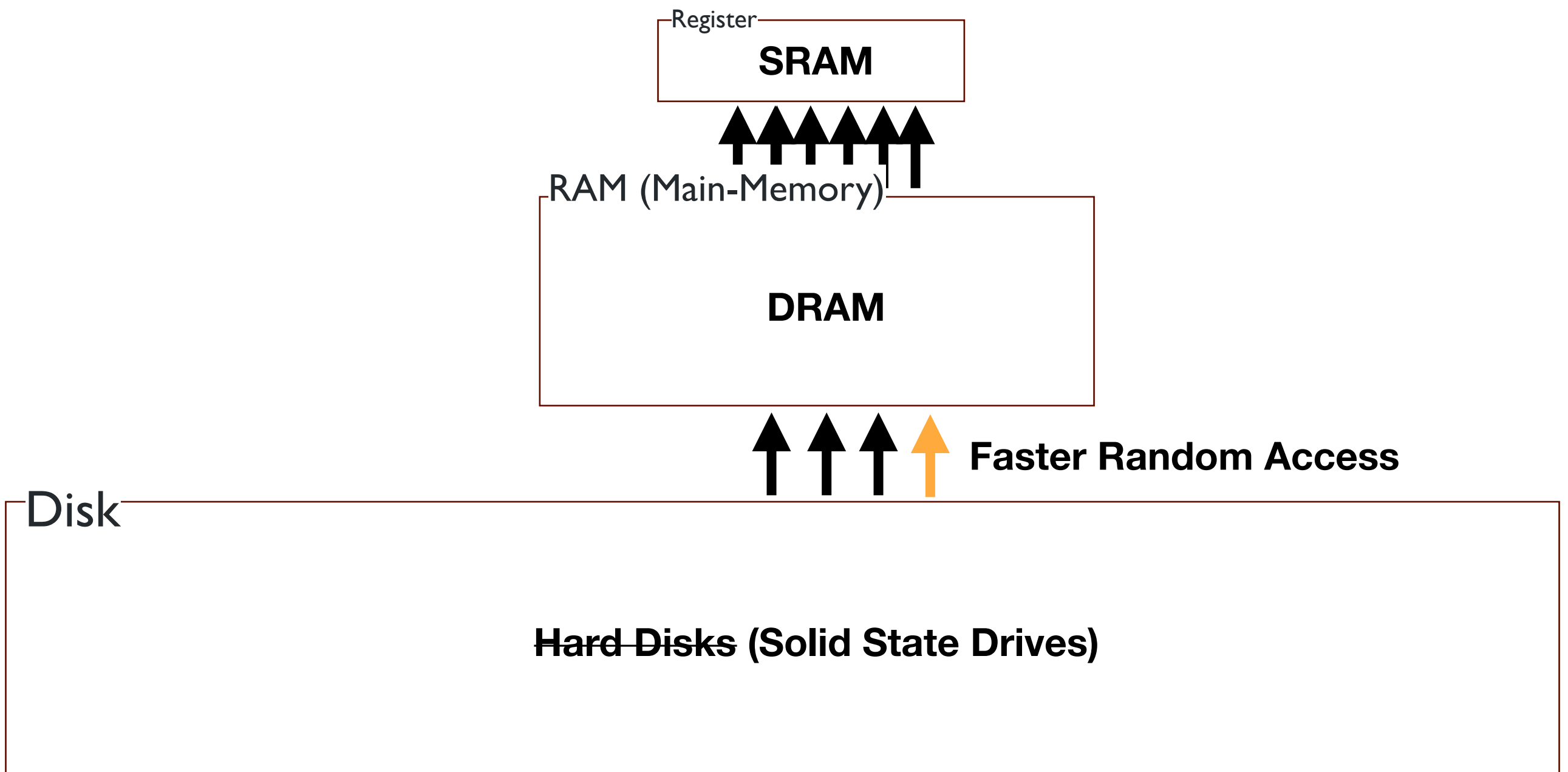
Indexing, Partitioning, and Sorting

Avoid transferring data you don't need

Technologies



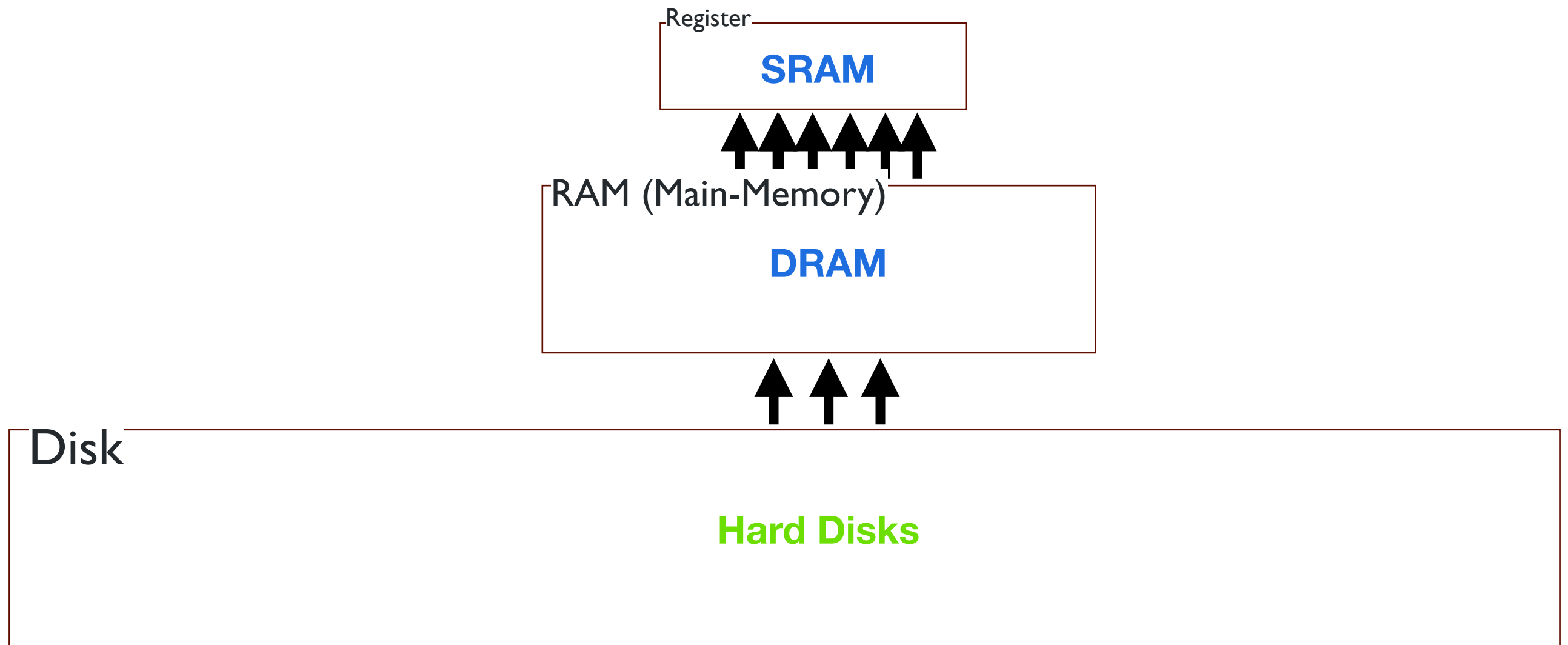
Technologies



Volatile v.s. Non-Volatile

Volatile: Power needs to be applied to store data

Non-Volatile: Data are persistent



Technologies

