# Caching

*DS 5110/CS 5501: Big Data Systems*
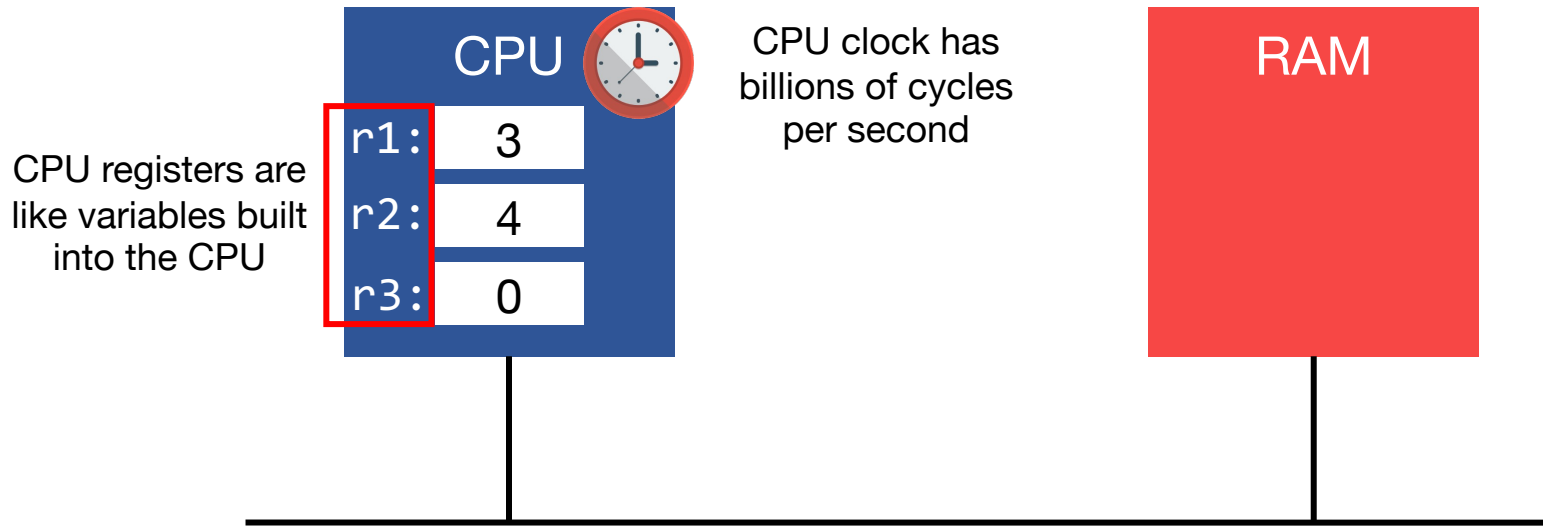
*Spring 2024*

Lecture 2d

Yue Cheng

# Learning objectives

- Describe the cache hierarchy

- Understand spatial locality and temporal locality

- Trace through access patterns with FIFO and LRU caching policies

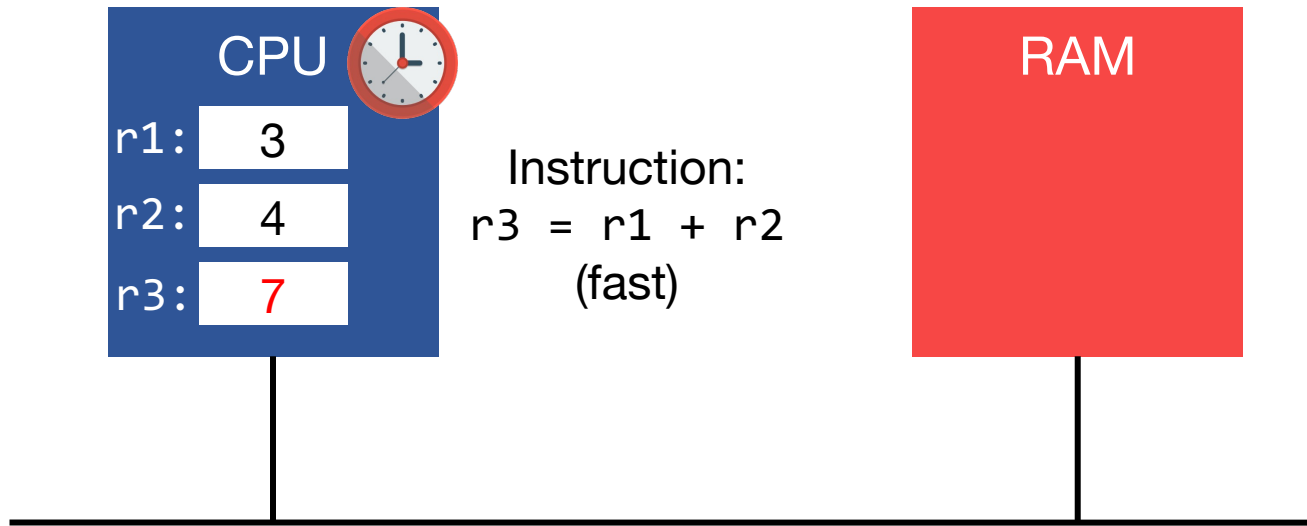  - Calculate cache performance metrics

# Outline

- Challenge: latency
- Cache hierarchy
  - CPU, RAM, SSD, Disk, Network
  - Tradeoffs
- Data access patterns, data locality, data access granularity
  - Spatial locality
  - Temporal locality
  - Cache lines and locality optimization
- What data should be cached?
  - Eviction policies: FIFO, LRU

# Interaction between CPU and RAM

**CPU**

r1: 3

r2: 4

r3: 0

CPU registers are like variables built into the CPU

CPU clock has billions of cycles per second

**RAM**

# Interaction between CPU and RAM



CPU

r1:  3

r2:  4

r3:  7

Instruction:
```
r3 = r1 + r2
```
(fast)

RAM

# Load and store

CPU
r1: 3
r2: 4
r3: 7

RAM
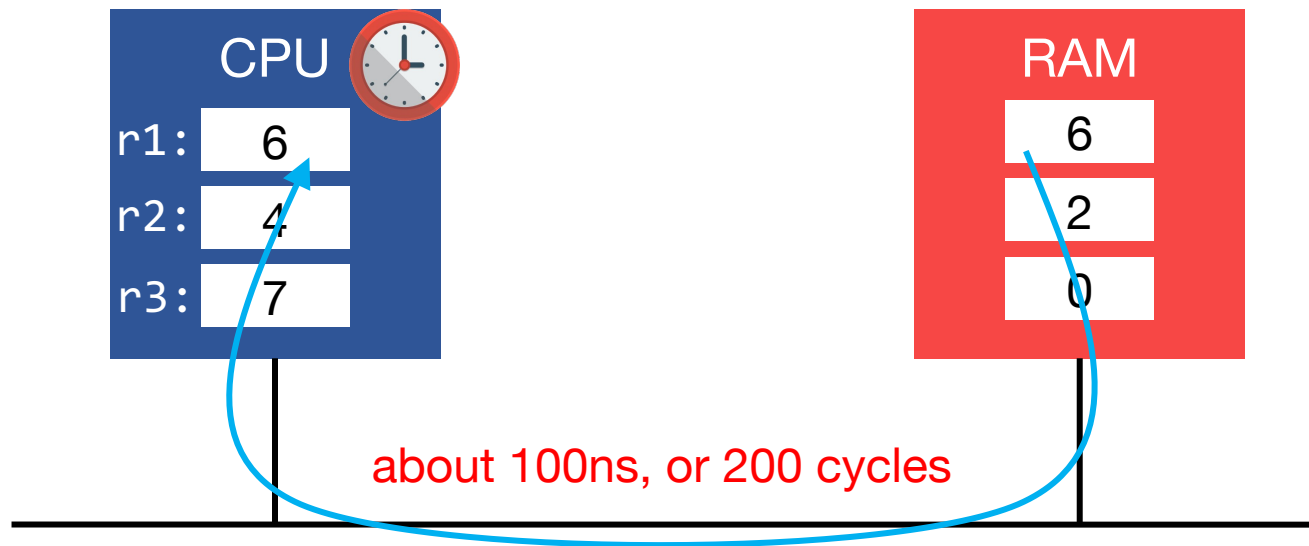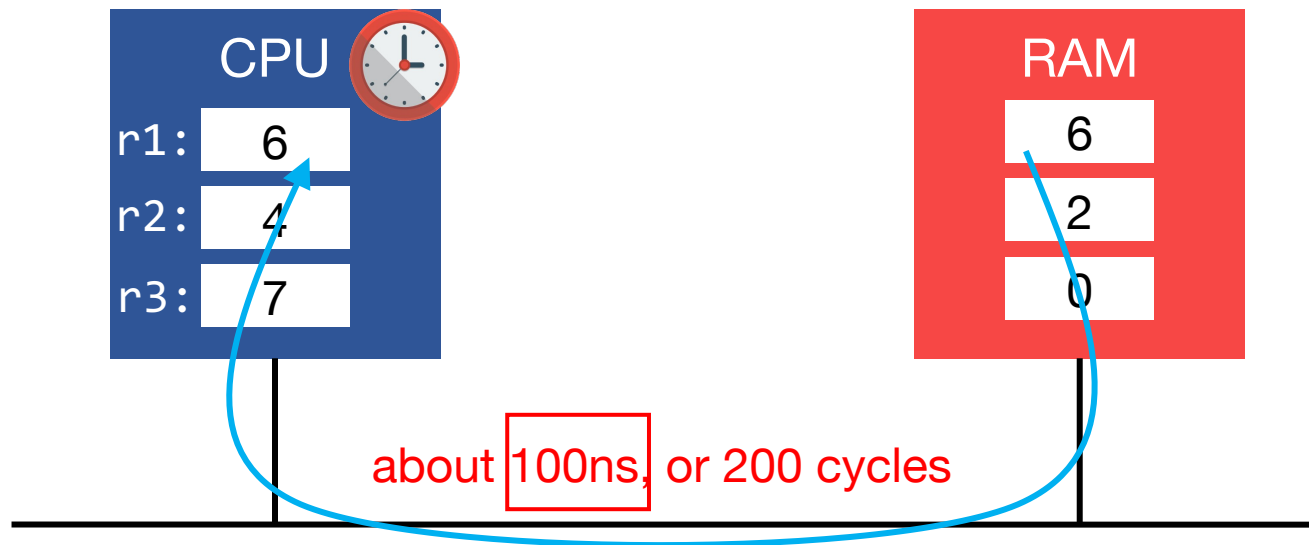6
2
0

**Challenge:** If we want to add some numbers stored in RAM, we need to **load** before adding and **store** after

# Latency to load from RAM

CPU

r1: 6
r2: 4
r3: 7

RAM

6
2
0

about 100ns, or 200 cycles

Very slow, but not long enough to switch to
a different thread…

# Latency

CPU
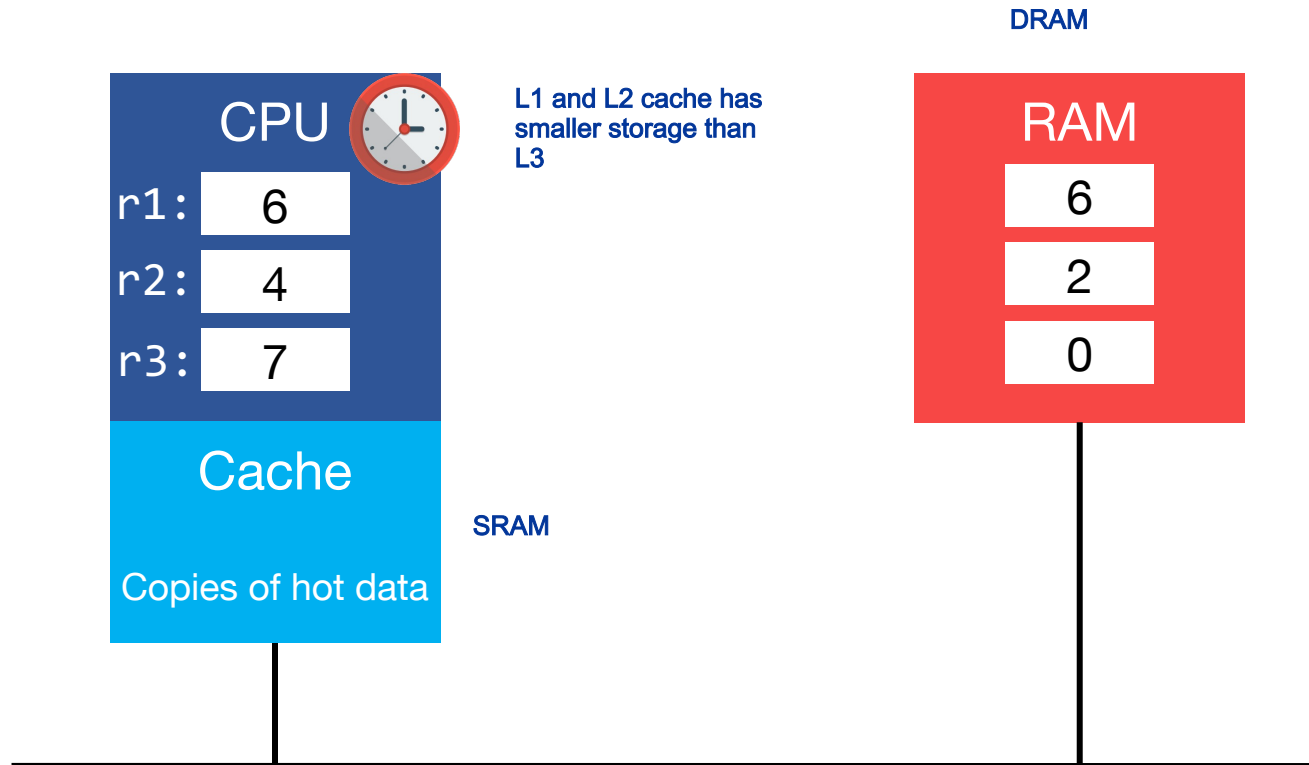| r1: | 6 |
| r2: | 4 |
| r3: | 7 |

RAM
| 6 |
| 2 |
| 0 |

about 100ns, or 200 cycles

"How much time" is a **latency** measure.

**Throughput** (bytes/second) depends on
how many loads we can do simultaneously.

# CPU Cache

**CPU**

L1 and L2 cache has smaller storage than L3

r1: 6
r2: 4
r3: 7

**Cache**

SRAM

Copies of hot data

DRAM

**RAM**

6
2
0

**Idea:** CPUs can have a small but very fast memory built in for data that is frequently accessed
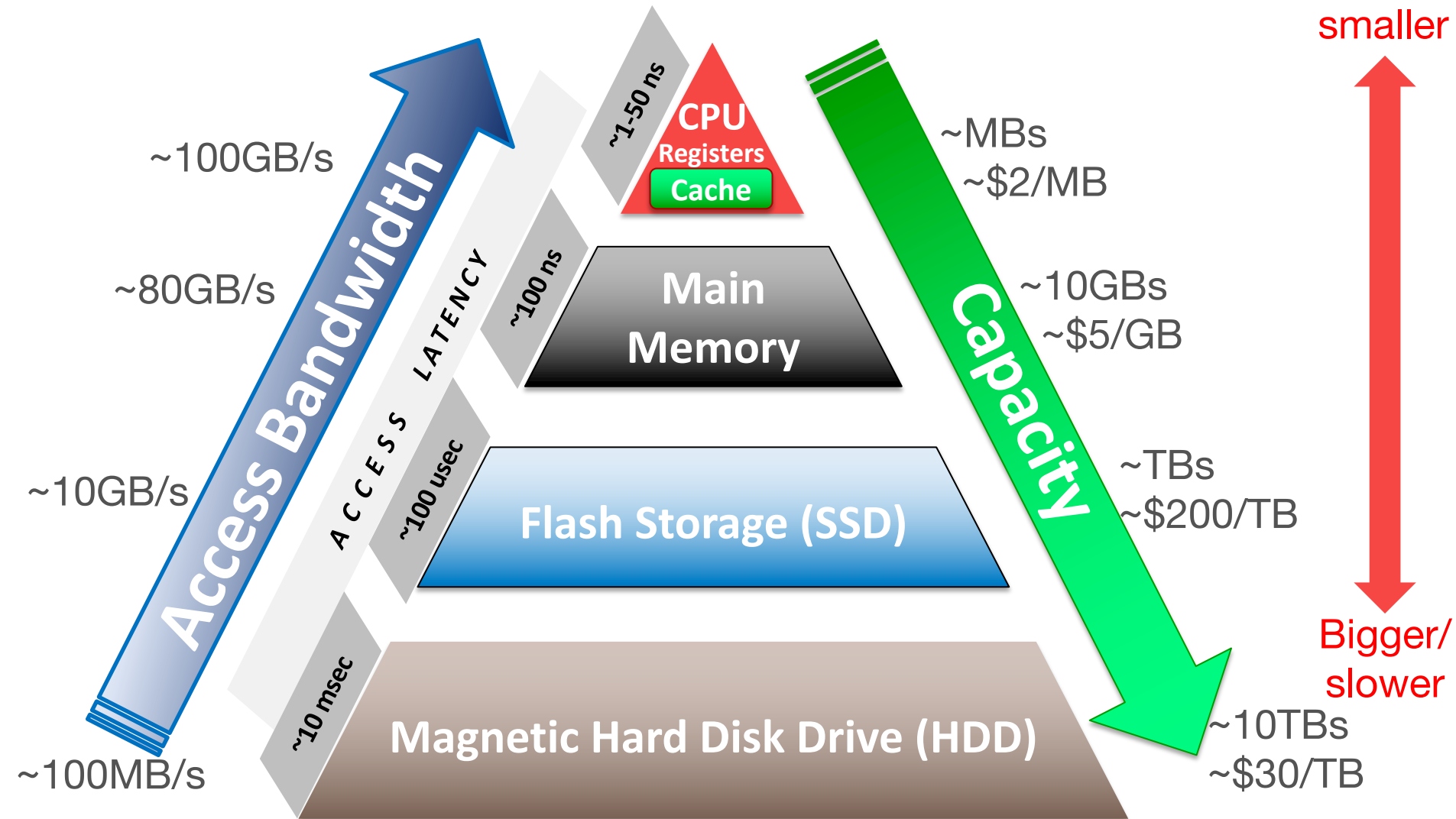
# Latency measurements

- Latency metrics
  - Average latency
  - Median latency
  - "Tail" latency ($99^{th}$ percentile, $99.9^{th}$ percentile, etc.)
    
    looking at the extreme occurances of latency

- Which metrics do we expect **caching** to help with the most? average and sometimes median! Caching is a hit or miss so it can optimize the average latency (or sometimes median)

# Cache hierarchy



Faster/ smaller

Bigger/ slower

Access Bandwidth

~100GB/s

~80GB/s

~10GB/s

~100MB/s

ACCESS LATENCY

~1-50 ns

~100 ns

~100 usec

~10 msec

CPU
Registers
Cache

Main Memory

Flash Storage (SSD)

Magnetic Hard Disk Drive (HDD)

Capacity

~MBs
~$2/MB

~10GBs
~$5/GB

~TBs
~$200/TB

~10TBs
~$30/TB

*UCSD DSC 102: Systems for scalable analysis. Arun Kumar

# Resource tradeoffs

- File system caches file data in RAM
    - Uses memory
    - Avoids storage reads

- Browser caches recently visited pages as disk files
    - Uses local storage space
    - Avoids network transfers

- Python dictionary caches return values in a `dict` (`key=args, val=return`)
    - Uses memory space
    - Avoids repeated compute

```
cache = {}
def f(x):
    if not x in cache:
        cache[x] = g(x)
    return cache[x]
```

# Workload characteristics

Application A

```
sum = 0
for i in range(0,1024):
  sum += a[i]
```

# Workload characteristics

Application A

```
sum = 0
for i in range(0,1024):
  sum += a[i]
```
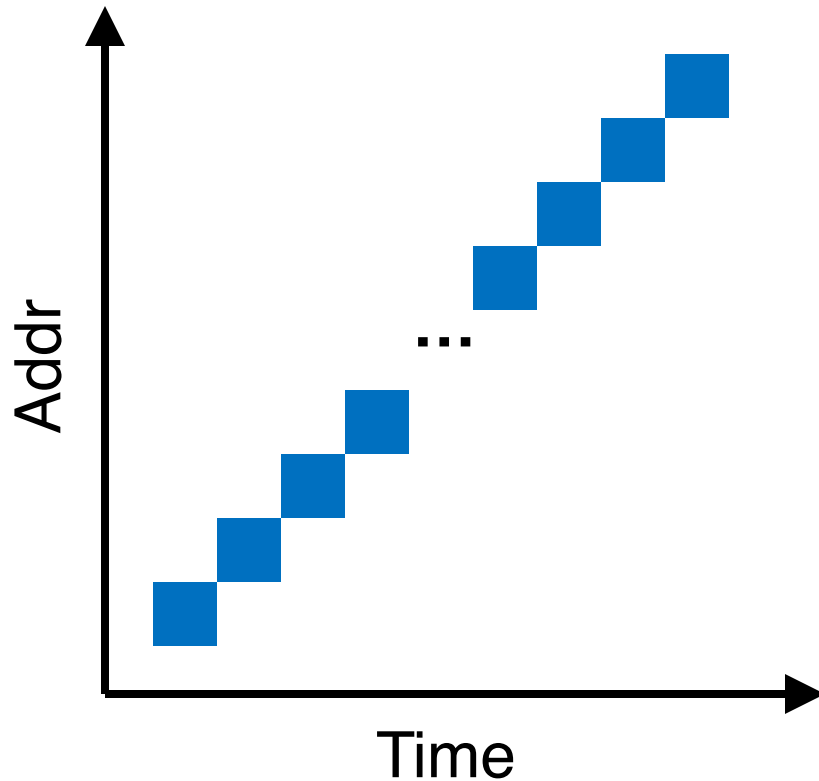
Application B

```
import random

sum = 0
random.seed(1234);
for i in range(0,512):
  sum += a[random.randint(0,1023)]

random.seed(1234) # same seed
for i in range(0,512):
  sum += a[random.randint(0,1023)]
```
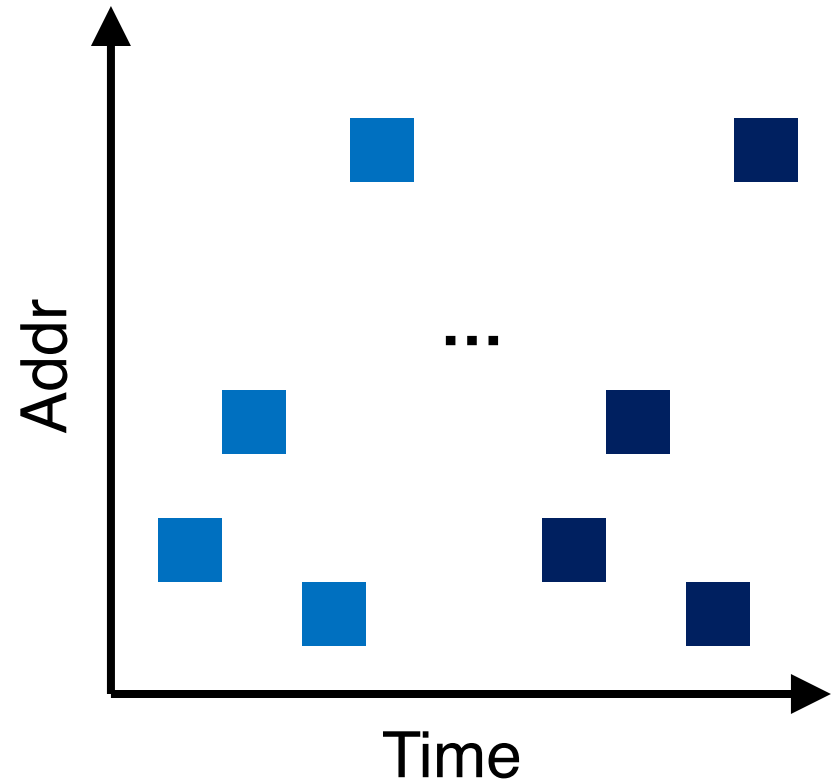
# Access patterns

Application A

Application B

Addr

Time

Addr

Time

# Access patterns

Application A

Application B



Addr → Time

**Spatial Locality**

Addr → Time

**Temporal Locality**

# Locality of data accesses

- **Spatial locality:**
  - Future access will be to nearby addresses

- **Temporal locality:**
  - Future access will be repeated to the same data

# Locality of data accesses

- **Spatial locality:**
  - Future access will be to nearby addresses


- **Temporal locality:**
  - Future access will be repeated to the same data


- Q: What is the **implication of data locality** to data systems applications?

# Locality optimization in Data Science

- Consider a matrix named data with `16*16` elements
- Each row is of size 16 floats and **prefetching+caching** means 1/2 row of accessed data item is brought to CPU cache at a time

# Locality optimization in Data Science

- Consider a matrix named data with `16*16` elements
- Each row is of size 16 floats and **prefetching+caching** means 1/2 row of accessed data item is brought to CPU cache at a time
- Program 1

```
for i in range(len(data[0]):
    for row in data:
        sum += row[i]
```

16 x 16 = **256** CPU cache misses

Not too hardware-efficient (not able to exploit prefetching+caching)

# Locality optimization in Data Science

- Consider a matrix named data with `16*16` elements
- Each row is of size 16 floats and <span style="color:red">prefetching+caching</span> means 1/2 row of accessed data item is brought to CPU cache at a time
- Program 1

```
for i in range(len(data[0]):
    for row in data:
        sum += row[i]
```

16 x 16 = <span style="color:red">256</span> CPU cache misses

Not too hardware-efficient (not able to exploit prefetching+caching)

- Program 2

```
for row in data:
    for element in row:
        sum += element
```

Only <span style="color:green">16*2</span> CPU cache misses

- Each time ½ row of data[i] is prefetched to cache so subsequent accesses are hits!

# Peeking behind the scene…

- Data access granularity
  - If a process reads one byte and misses, <span style="color:red">how much data should the CPU bring into the CPU cache?</span>
  - Tradeoff:
    - **Too little?** Will have many more misses if we read nearby bytes soon (recall spatial locality)
    - **Too much?** Wasteful to load data to cache that might never be accessed

- CPU caches data in units called **cache lines**
  - Typically, 64 bytes for modern CPUs (8 float64 numbers)

# Cache lines and misses

cache line

| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |

cache line

| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |

cache line

| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |

cache line

| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |

cache line

| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |

cache line

| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |
| fp64 |

How many misses?

How many misses?

How many misses?

# Memory layout of a matrix

Matrix of numbers
**Logically**, 2-dimensional

| |
|---|
| Row |
| Row |
| Row |
| Row |

**Physically**, those rows are arranged along **1-dimension** in the virtual address space

| Code | Row | Row | Row | Row | stack |
|---|---|---|---|---|---|

Virtual address space

# Memory layout of a matrix

Matrix of numbers
**Logically**, 2-dimensional

| Row |
| --- |
| Row |
| Row |
| Row |

Summing over row:
data consolidated into a few cache lines (CPU cache friendly)

| | Code | | Row | Row | Row | Row | | stack | |
|---|---|---|---|---|---|---|---|---|---|

# Memory layout of a matrix

Matrix of numbers
**Logically**, 2-dimensional

| Row |
|-----|
| Row |
| Row |
| Row |

Summing over row:
data consolidated into a few cache lines (CPU cache friendly)

| Code | Row | Row | Row | Row | stack |

Summing over column: each number is in its own cache line and triggers a **cache miss**

# Demo ...

UVA DS5110/CS5501 Spring '24

# Caching policies

- When to <span style="color:red">load</span> data to a cache?
    - Whenever the program reads something, add it to cache

- When to <span style="color:red">evict</span> data from a cache (eviction policy)? Several policies:
    - Random: select any data at random for eviction
    - **FIFO** (first-in, first-out): evict whichever data that has been in the cache the longest
    - **LRU** (least recently used): evict which data that has been used the least recently

# Worksheet ...