

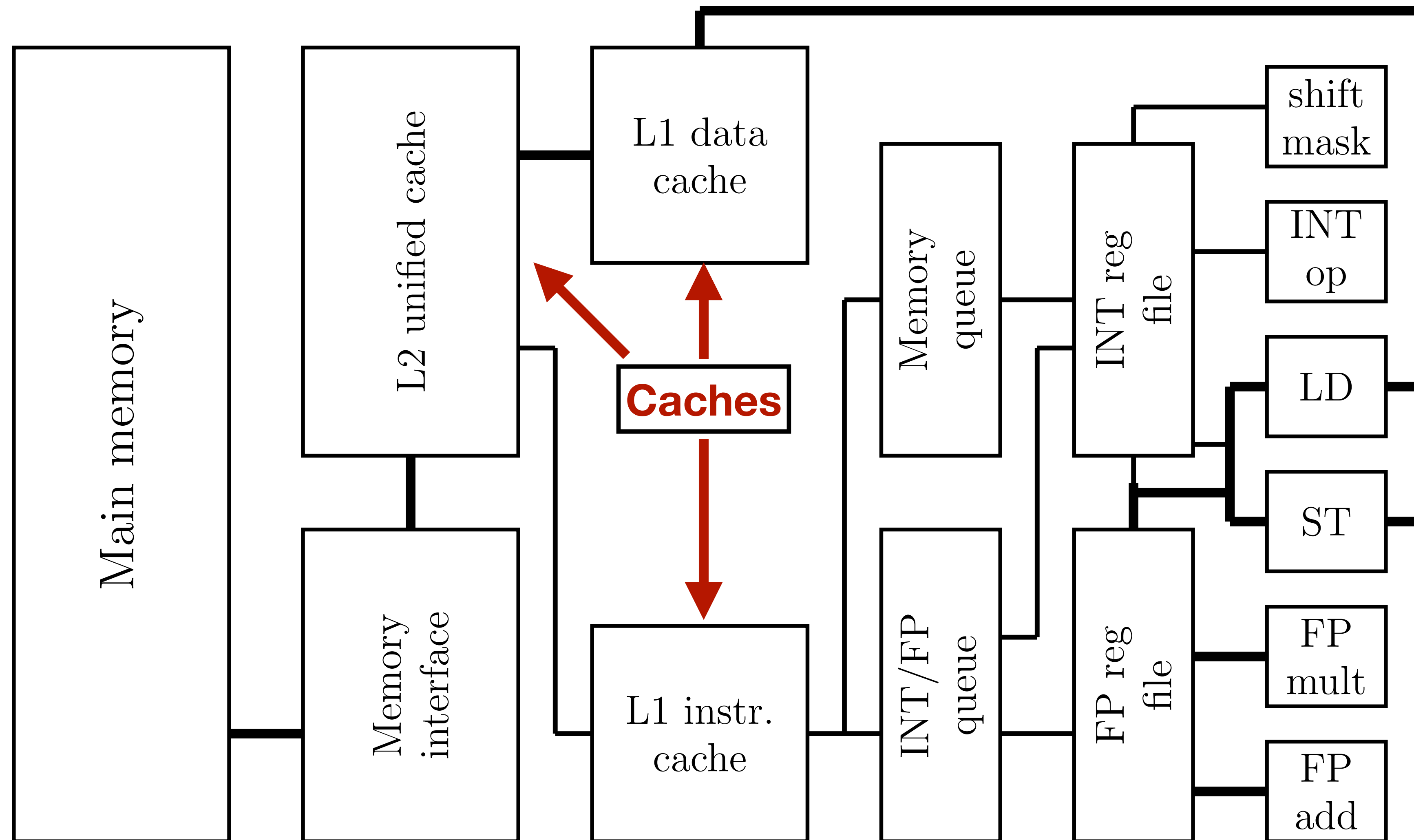
# Introduction to Parallel Processing

Lecture 3 : Serial Optimizations

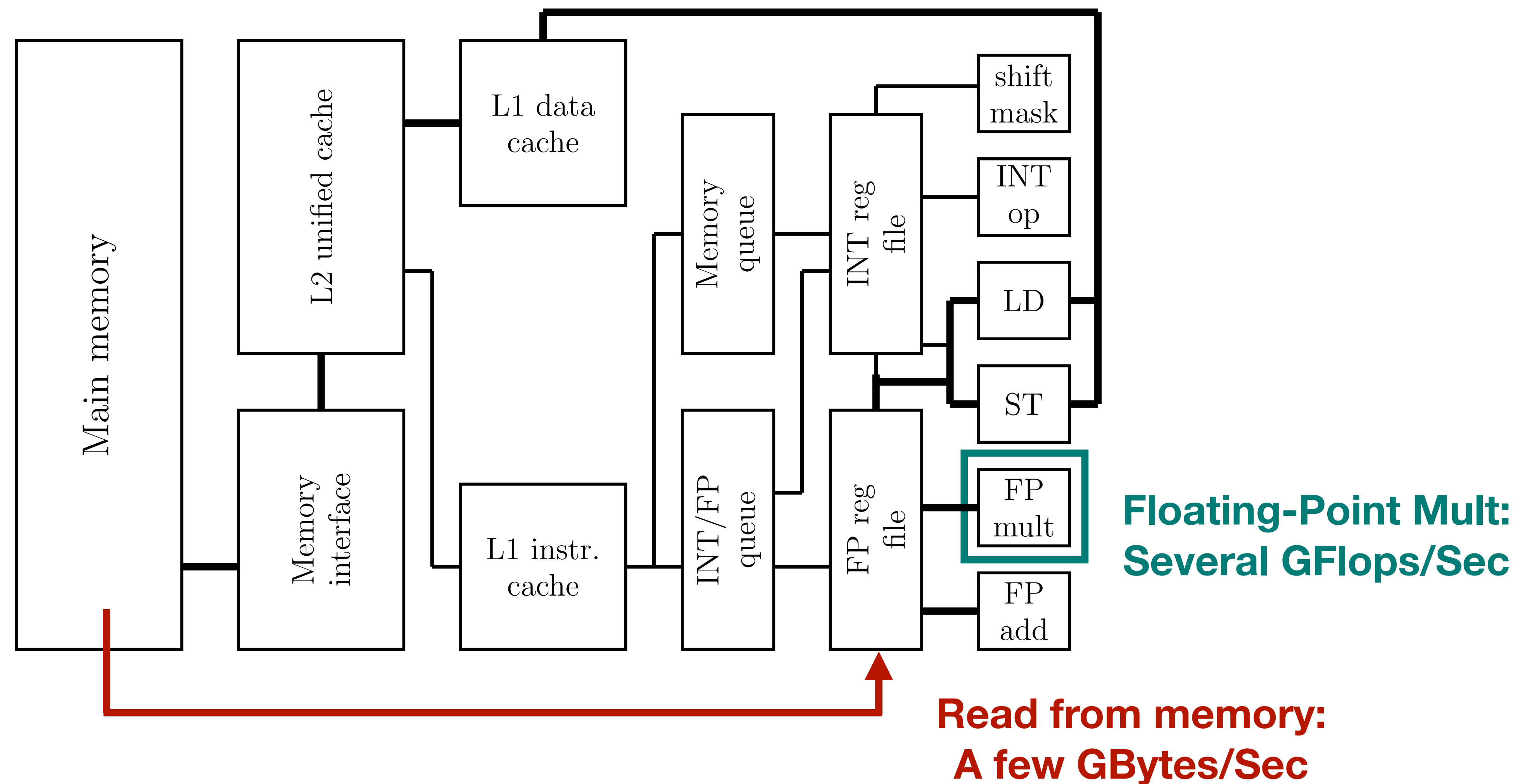
08/31/2022

Professor Amanda Bienz

# Cache-Based Microprocessor

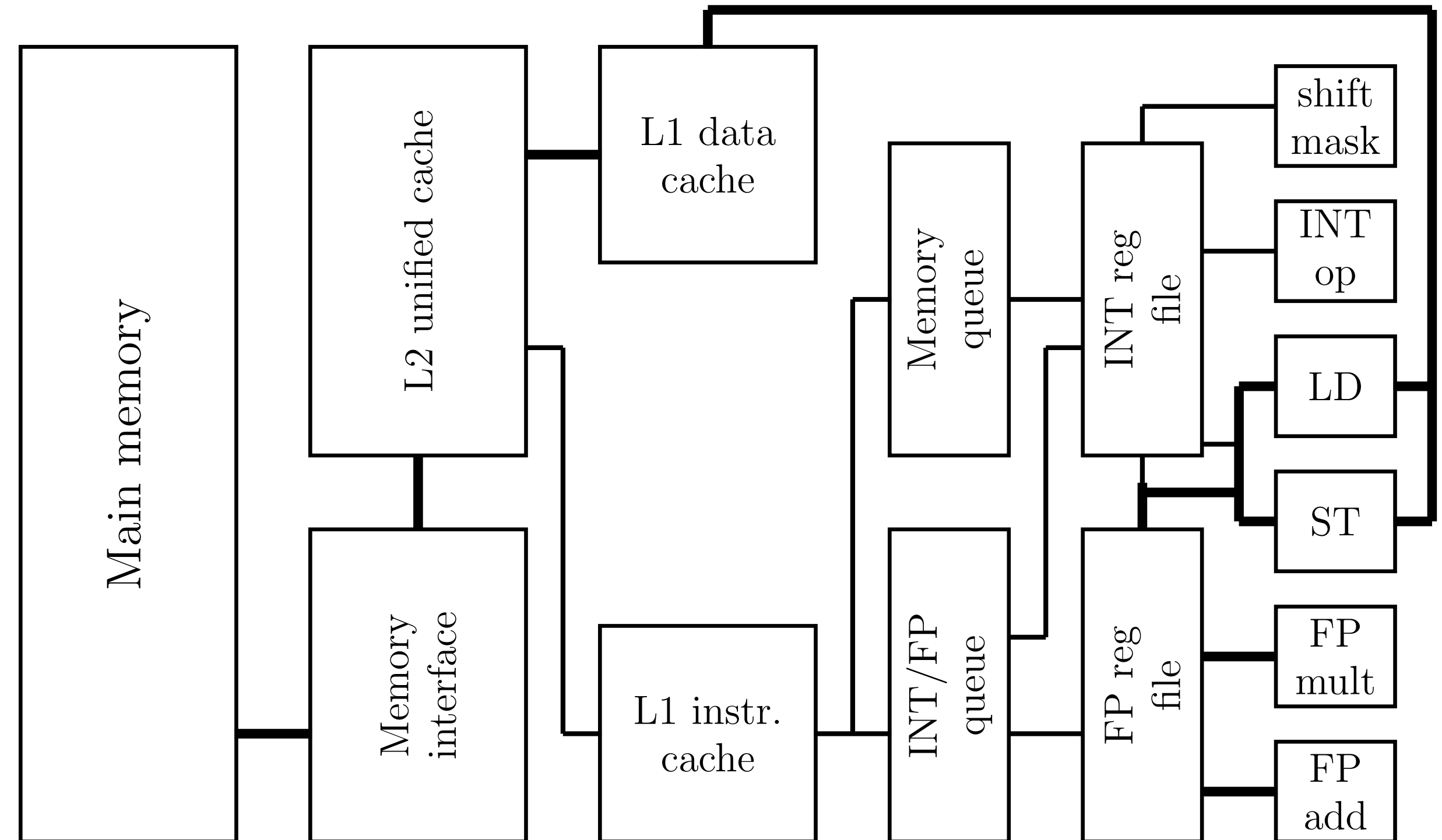


# Why Do We Have Caches



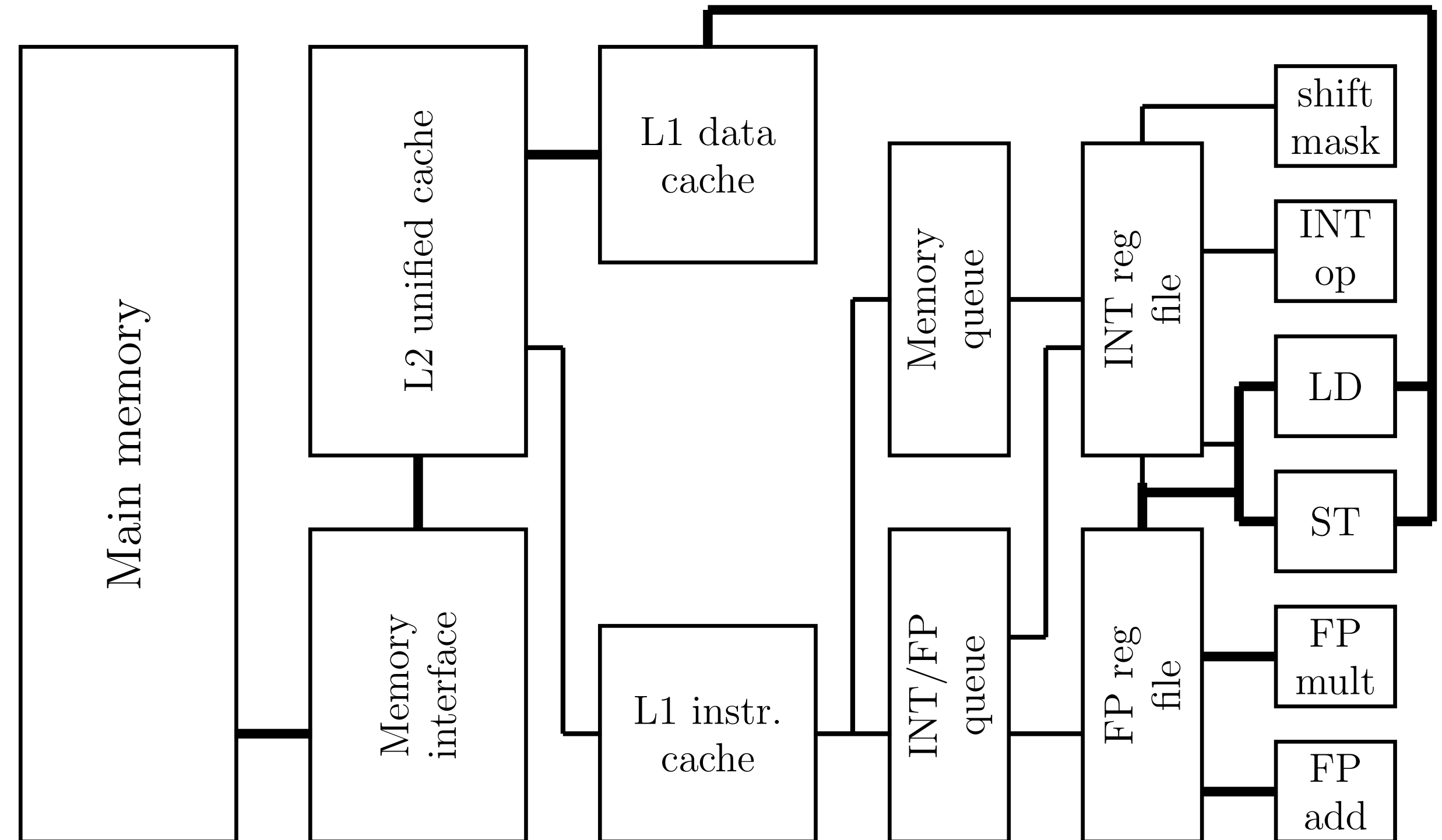
# Cache Overview

- Typically at least 2 caches
- L1 is split (data, instr.)
- Others are unified
- Cost of reading data:
  - Latency
  - Bandwidth
- **Closer the cache is to register, lower latency, higher bandwidth**



# Cache Overview

- When reading data:
  - **Cache hit** : data in cache
  - **Cache miss** : not in cache
    - Outer cache or main memory
- If cache is full, data is evicted



# Memory / Cache Latency

- Overhead for reading a single byte
- Not dependent on the number of bytes being read
- **Cache Lines** : a group of data is read from memory (say 64 bytes)
- During a cache miss, fetch entire line from memory
  - Neighboring items will then be read from cache

# Writing Data is More Complicated

- **Write Hit** : data to be written already in cache
  - **Write-back** : modify cache line in cache, write to memory whole cache line when evicted
- **Write-Miss** : Cache line first transferred from memory to cache before being modified
  - Increased data transfer cost

# How Much Does Cache Help?

Component	Required Cycles	Cost, relative to previous
Register	1	
L1 Cache	A few	1-3x
L2 Cache	10	3-10x
Memory	250	25x

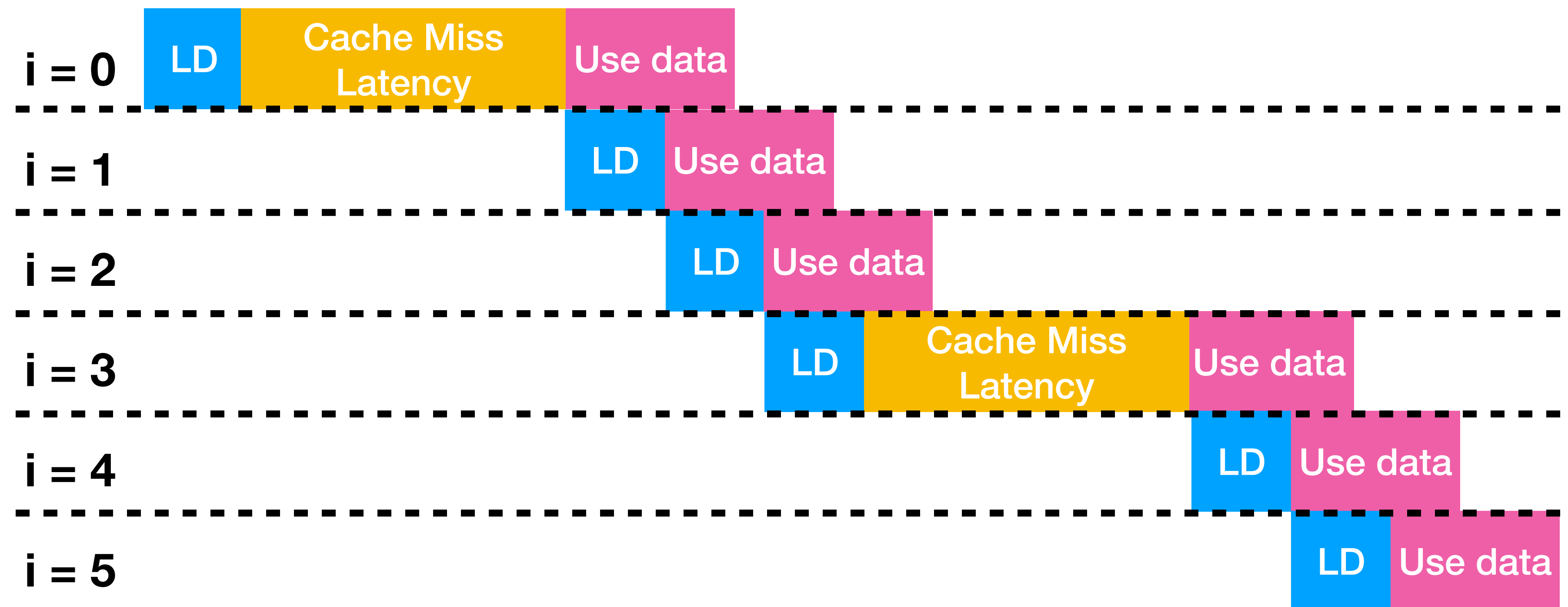


# Prefetching

Cache miss has high latency, even with use of cache lines  
Let's assume cache line size is 4 doubles below

Vector Norm :

```
for i = 0...n:  
  sum += v[i]*v[i]
```



# Prefetching

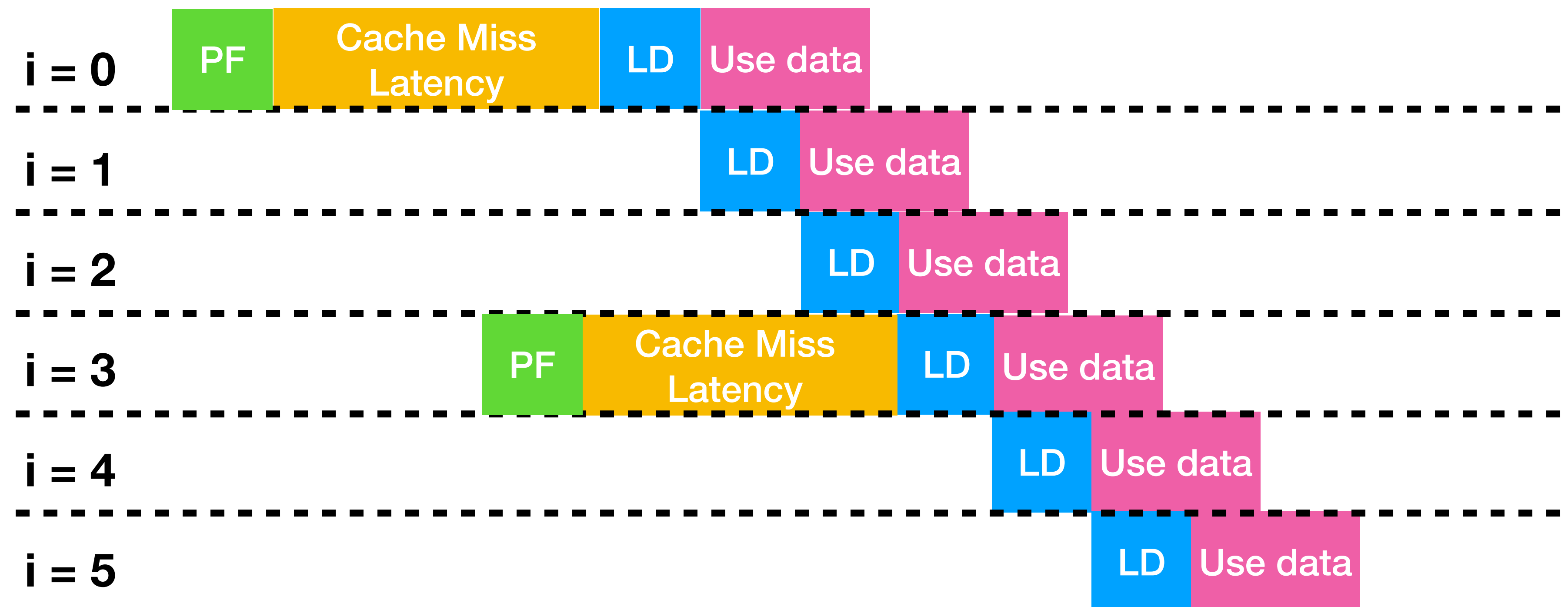
Load data into cache before it is needed by the application

Compiler adds instructions to touch cache lines early

*Hardware prefetcher* detects regular access patterns

**Vector Norm :**

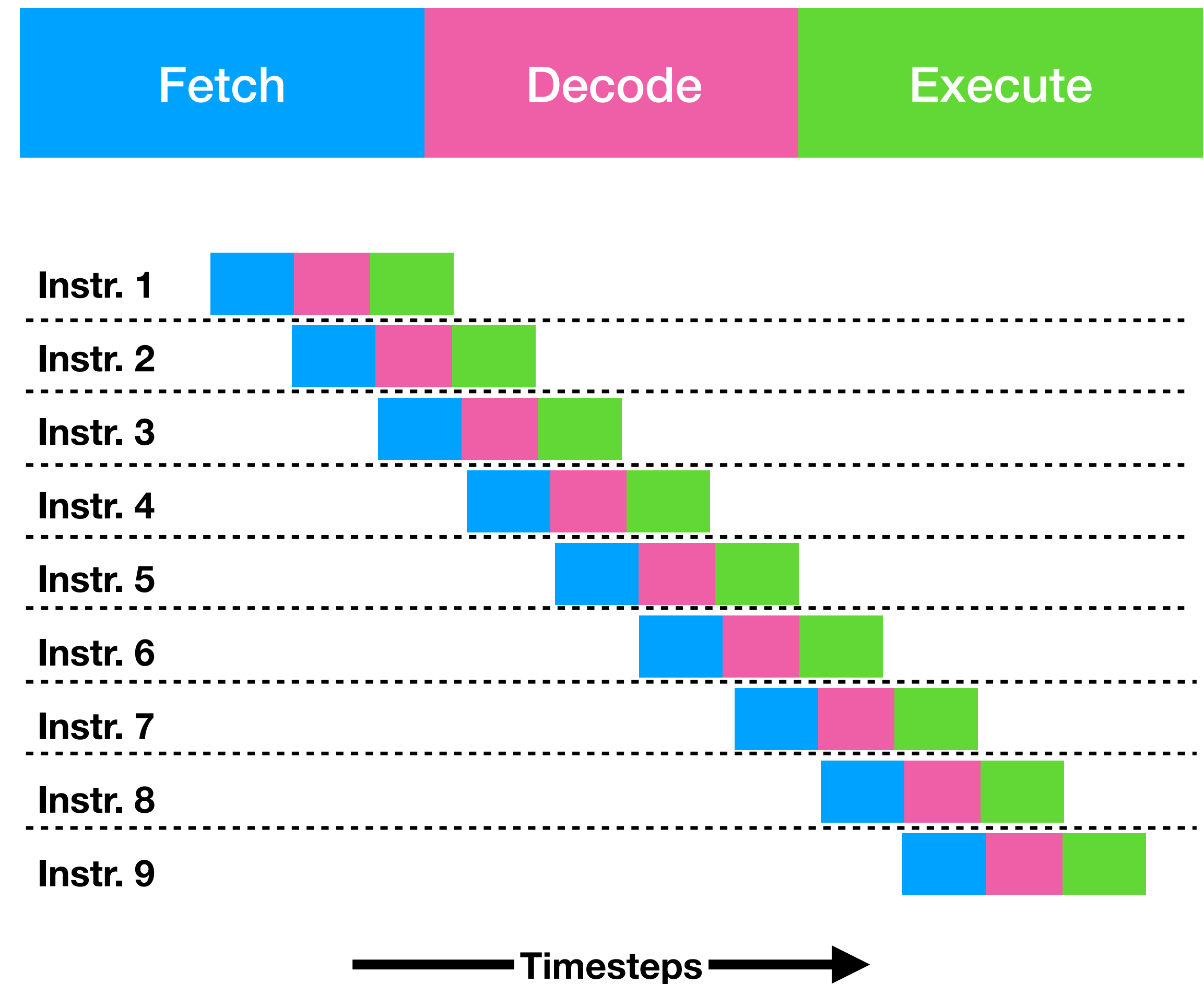
```
for i = 0...n:
  sum += v[i]*v[i]
```



**Vastly improves performance, but depends on regular access patterns**

# Branch Prediction

- Pipelining instructions can also be described as 'instruction prefetching'
- Can run into bottlenecks
- If/Else statements : which path to take?



# If you want to learn more..

- When prefetching works, when it doesn't, and why
- Some great slides on caches by Bill Gropp
- A Study of Branch Prediction Strategies
- <https://ieeexplore.ieee.org/abstract/document/628926/>