

This works because of locality

2 types of locality

Spatial locality

↳ if I'm accessed it's likely my neighbor will be

- instruction stream
- stack accesses
- reading an array

Temporal locality

↳ I'm likely to be reused

- access data repeatedly
- loops of instructions
- stack accesses

sum = 0

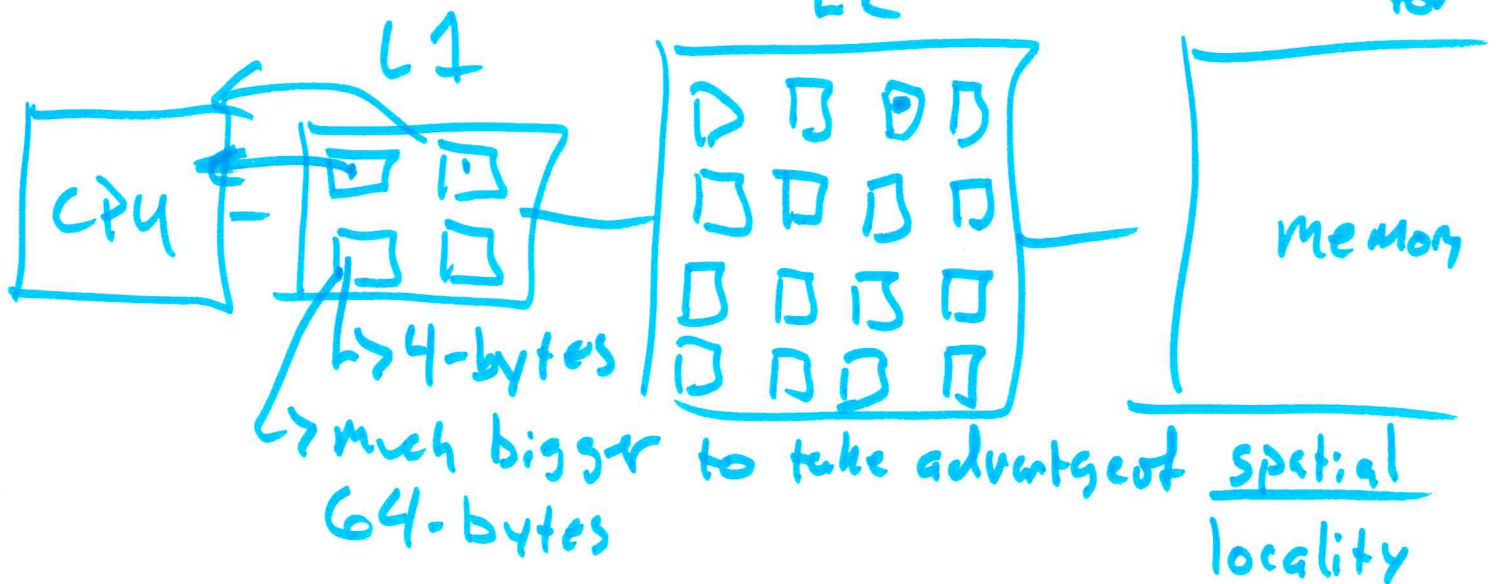
```
for (int i = 0; i < n; i++) {  
    sum += a[i];  
}  
return sum;
```

Annotations: An arrow points from `a[i]` to the word "spatial". Another arrow points from the curly braces to the word "temporal".

Spatial locality in instruction stream
+ temporal locality in the loop

Caches take advantage of locality

- Temporal locality → keep recently used things close
- Spatial locality → load more than was asked for



CPU checks cache \rightarrow if it's there: HIT
if not there: MISS

3 different kinds of cache misses (3 C's)

- Cold/compulsary misses

\hookrightarrow ~~first~~ miss because cache is empty

\hookrightarrow "warm up" caches

- Conflict misses

- cache could hold the data, but
because of limited "slots" it's a miss

- Capacity misses

- cache is full

- Working set is too big for cache

\hookrightarrow current data program is using

Caches are managed by hardware

\hookrightarrow no software / OS / etc involved

\hookrightarrow logically transparent to the programmer

in contrast to registers which are managed by the program

Average memory access time or miss ratio
 $(T_{hit}) + (T_{miss}) \cdot (\text{miss rate}) \rightarrow$

Cache hit time 1ns

miss time 100ns

Hit rate \rightarrow 99%

\rightarrow Miss rate = 1%

$\hookrightarrow \frac{\# \text{ hits}}{\text{total } \# \text{ accesses}}$

$\frac{\# \text{ misses}}{\text{total } \# \text{ of accesses}}$

$$\begin{aligned} \text{AMAT} &= 1\text{ns} + 100\text{ns} \cdot (.01) \\ &= 2\text{ns} \end{aligned}$$

$$T_{\text{miss}} = \text{AMAT}_{L2} = T_{\text{hit}L2} + T_{\text{miss}L2} \cdot \left(\frac{\text{L2 miss rate}}{\text{rate}} \right)$$

L2 miss rate = 25%

L2 hit time = 10ns

L2 miss time = 100ns

$$\begin{aligned} \text{AMAT} &= 1\text{ns} + (.01) \cdot [10\text{ns} + .25 \cdot 100] \\ &= 1 + .01 \cdot (35) \\ &= 1.35\text{ns} \end{aligned}$$

$$AMAT = T_{hit} + T_{miss} \cdot \left(\frac{\text{miss}}{\text{rate}} \right)$$

$$T_{miss} = AMAT_{L2} = T_{hitL2} + T_{missL2} \cdot \left(\frac{L2 \text{ miss}}{\text{rate}} \right)$$

$$AMAT = T_{hitL1} + \left(\frac{L1 \text{ miss}}{\text{rate}} \right) \left[T_{hitL2} + T_{missL2} \cdot \left(\frac{L2 \text{ miss}}{\text{rate}} \right) \right]$$