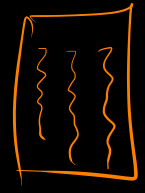


Sai

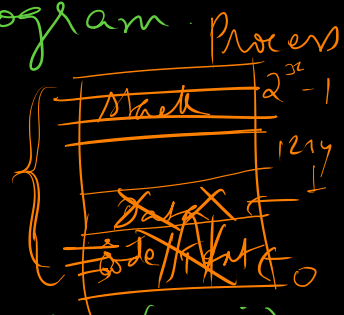
Alternate Approaches for hiding Memory Latency

I. Multithreading

Thread : Single Stream of control in the flow of a program



```
for (i = 0; i < NN; i++)  
    for (j = 0; j < N; j++)  
        c[i][j] = dot-product(get-row(a, i),  
                               get-col(b, j));
```



n^2 dot-product (...) \rightarrow Each can be executed independently.

Main thread
Can spawn n^2 or n^2-1 threads for each to take one iteration

Can be assigned to different compute abstractions (Threads)

↓
One process, multiple threads
Main thread to wait until all threads have finished

→ Program Complete

↓
What if there were n^2 cores?

n^2 speedup

Speedup?

What if there was only one core?

```
for (i = 0; i < M; i++)
```

```
for (j = 0; j < N; j++)
```

```
c[i][j] = create_thread(dot_product(
    get_row(a, i), get_col(b, j)));
```

(Assume only one core)

⇒ Hides Memory latency provided,

→ There is enough Concurrency to keep the processor busy

→ Program explicitly specifies it

→ Fast Context Switching of threads

II. Prefetching

```
for (i = 0; i < N; i++)
```

```
    c[i] = a[i] + b[i];
```

```
load R1, a[0]
```

```
add R1, b[0]
```

```
store R1, c[0]
```

```
⋮
```

(Load just before use)

```
load R1, a[0]
```

```
load R2, b[0]
```

```
load R3, a[1]
```

```
load R4, b[1]
```

```
⋮
```

```
add R1, R2
```

```
store R1, c[0]
```

```
⋮
```

(Load much ahead of use)

→ Even if there is a cache miss, data is likely to arrive by the time it is used.

→ Fresh load needed if a data item is overwritten between load and use (but no worse than the original)

- Effectively, same as multithreading, in identifying no resource or data dependencies
- Done by compiler optimizations

Issues in Multithreading & Prefetching

I. 1 GHz CPU, 4-word cache line, 1 cycle access to cache, 100 ns latency to DRAM

Assume a computation:

- CPU requests 1 word every cycle
- For 1 KB cache, Cache Hit Ratio = 25%.
- For 32 KB cache, Cache Hit Ratio = 90%.

Two cases: a) Single thread, 32 KB cache

b) 32 threads, 1 KB cache to each

Case a) 1 word every cycle and 90% Cache Hit Ratio

⇒ 1 in 10 requests hits memory

4 Bytes ^{10 words} ⇒ On an average, 1 word in 10 ns from DRAM

⇒ Memory bandwidth requirement = 400 MB/s

Case b) 1 word every cycle and 25% Cache Hit Ratio

⇒ 3 in 4 requests hits memory

⇒ On an average, 3 words in 4 ns from DRAM

⇒ Memory bandwidth requirement = 3 GB/s

12 bytes 4
? 10

// Bandwidth requirements of multithreaded program may increase significantly due to smaller cache residency of each thread.

// Tend to become bandwidth bound instead of latency bound.

// Prefetching does not worsen than the original, but if needed to fetch again, it means, fetching data item twice, resulting in doubling memory bandwidth requirement.

II. Additional hardware needed to effectively utilize multithreading or prefetching.
(Need large Register files and Caches)