

Parallel and Distributed Computing

Scalability

- ↳ Property of system to handle growing amount of work by adding resources to the system
- ↳ adding resources to a system to handle much more work e.g. 1000 customers 2 staff
10,000 customers, so need to increase staff

SCALABILITY DIMENSIONS

Physical/Load Scalability

- ↳ a system scalable w.r.t its size
- ↳ easily add/remove users and resources

if work load decreases
remove resources

Administrative Scalability

- ↳ ability for an increasing no of organisations/users to a system

↳ if made scalable then has the administration/organisation increased along with it or not

Functional Scalability

- ↳ ability to enhance the system by adding new functions w/o disrupting existing activities

Generational Scalability

- ↳ ability of a system to scale by adopting new generations of components

↳ if system of 8 gen
can we add a component of a newer gen?

Horizontal Scalability

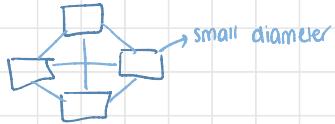
→ scale out

Vertical Scalability

→ scale up

NOT completed

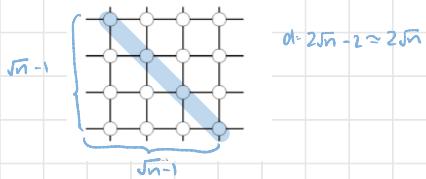
→ COMPLETE NETWORK



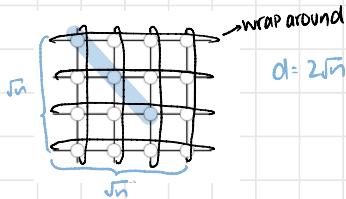
→ STRAIGHT PATH



→ 2D MESH



→ 2D TORUS



Parallel computing systems

- ↳ doing things simultaneously

- ↳ task is split into parts

- ↳ each part is further broken into stream of instruction

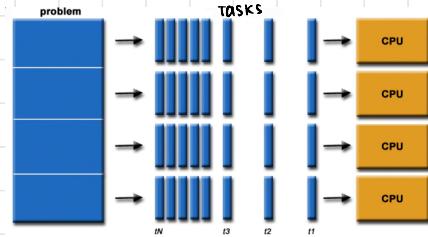
- ↳ instructions of different parts are then executed concurrently

PROS

- ↳ save time

- ↳ solve large problems

- ↳ provide concurrency



USES

- ↳ processors aren't getting faster

- ↳ solves problems of huge memory requirements

APPS

- ↳ data mining

- ↳ web search engines

- ↳ advanced graphics

- ↳ AI, ML, DL

- ↳ real time systems

Implicit parallelism refers to the ability of the compiler to perform certain operations in parallel without "any" hint from the programmer. This is achievable if the computation you want to perform has certain properties (e.g. no data dependencies among the parallel jobs)

Serial computing

- ↳ task split into stream of instructions

done after the other

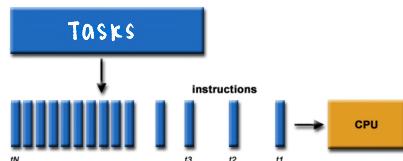
- ↳ they are executed sequentially on a single processor

CONS

- ↳ one instruction executes at a time

- ↳ slow

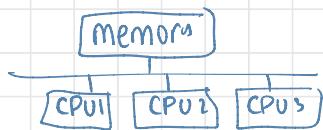
- ↳ expensive to make a single processor faster



PARALLEL ARCHITECTURE

Shared Memory

- ↳ large memory unit
 - ↳ multiple processors accessing it
 - ↳ tasks share a common address space
- CON**
- ↳ not scalable → limit to no. of CPU's that can be added



PROGRAMMING MODELS

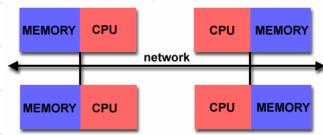
Shared memory → OPEN MP

- ↳ Public



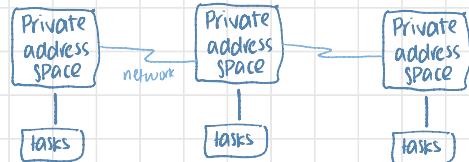
Distributed System

- ↳ multiple processors with their own memory
- ↳ can be scalable



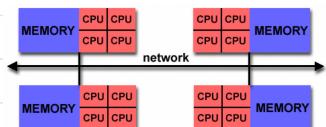
MESSAGE PASSING → MPI

- ↳ Private
- ↳ exchange messages via send/receive



Hybrid

- ↳ shared + distributed memory



→ typically used
↳ Cache coherent

→ if 1 processor updates location in shared memory, all other processors know

- ↳ largest and fastest computers

Shared Memory

- ↳ all processors have direct access to common physical memory

PROS

- ↳ user friendly
- ↳ fast and uniform data sharing

CONS

- ↳ non scalable memory
- ↳ expensive

Distributed Memory

- ↳ network based memory access to non common physical memory
- ↳ each processor has its own local memory
- ↳ no cache coherency

PROS

- ↳ scalable memory
- ↳ fast memory access
- ↳ cost effective

CONS

- ↳ difficult to map some ds
- ↳ NUMA

Uniform
Memory
Access

UMA

- ↳ identical processors
- ↳ equal access time to memory
- ↳ cache coherent
 - if 1 processor updates location in shared memory, all other processors know

NON
UNIFORM
MEMORY
ACCESS

Symmetric
multi processor

NVMA

- ↳ 1 SMP can access another SMP
- ↳ not all processors have equal access time
- ↳ slow memory access

Parallel Overhead

- ↳ amount of time to coordinate parallel tasks as opposed to doing useful work
- ↳ include factors
 - ↳ Task start up time
 - ↳ synchronizations
 - ↳ Data communications
 - ↳ Task termination time

Massively Parallel

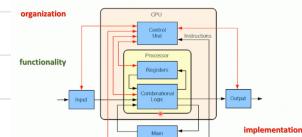
- ↳ hardware that has many processors

Scalability

- ↳ property of system to handle growing amount of work by adding resources to the system
- ↳ adding resources to a system to handle much more work
 - e.g. 1000 customers → staff
 - 10,000 customers, so need to increase staff
- ↳ factors contribute to scalability
 - ↳ hardware
 - ↳ Application algorithm
 - ↳ Parallel overhead

Architecture

- ↳ set of rules that describe functionality, organisation
- ↳ balancing performance ↳ efficiency ↳ cost ↳ reliability



↳ categorisation of things

FLYNN'S TAXONOMY

- ↳ 4 classifications based on no of concurrent instructions and data streams available

SISD → Single instruction Single data

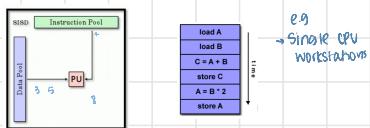
- ↳ no parallelism
- ↳ single control unit
- ↳ one operation at a time

PROS

- ↳ no issue of complex comm b/w no of cores
- ↳ requires less power

CONS

- ↳ speed is limited
- ↳ not suitable for larger applications

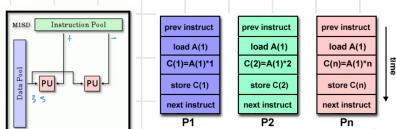


MISD multiple instruction single data

- ↳ performs different instructions on same data set

USES

- ↳ multiple cryptography algos to crack code



↳ pipelined computers

SIMD → Single instruction multiple data

- ↳ instructions can be executed sequentially or parallelly

PROS

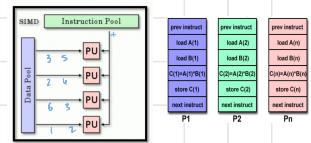
- ↳ increasing no of cores increases throughput
- ↳ processing speed > SISD

CONS

- ↳ complex comm b/w
- ↳ cost > SISD

USES

- ↳ image processing



↓
vector
processors

common
fastest!

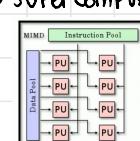
MIMD → multiple instruction multiple data

- ↳ performs different instruction on different data

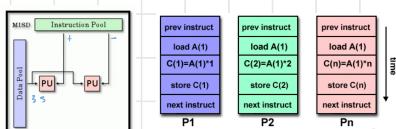
USES

- ↳ most modern computers

USES

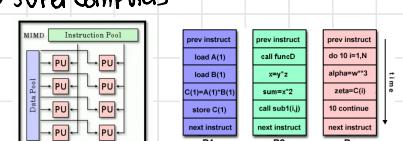


synchronous/asynchronous
deterministic/non-deterministic



↳ pipelined
computers

a + b



Parallel execution

- ↳ executing multiple statements at the same time

Serial Execution

- ↳ executing one statement at a time sequentially

Task

- ↳ set of instructions executed by a processor

Data Parallel

- ↳ same operation on diff data
- ↳ Scalable

Task Parallel

- ↳ diff operations on same/diff data
- ↳ not scalable

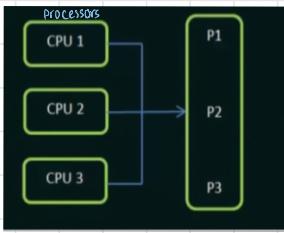
Communication

- ↳ data exchange for parallel tasks through
 - ↳ shared memory bus
 - ↳ network

Synchronization

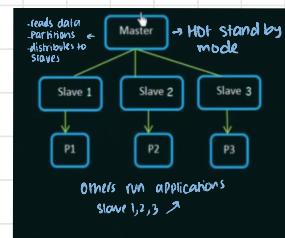
- ↳ coordination of parallel tasks in real time

Symmetric Multiprocessing



- ↳ shared memory
- ↳ all CPUs similar

Asymmetric Multiprocessing



- ↳ NO shared memory
- ↳ independent memory
- ↳ master and slave

ishma hafeez
notes
reprt

Granularity

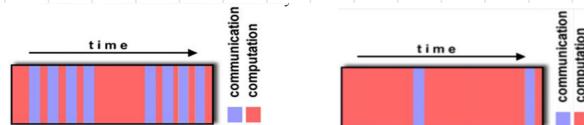
↳ ratio of computation to communication

1. Fine grain Parallelism

- ↳ relatively small amount of computational work
- ↳ facilitates load balancing
- ↳ high communication overhead
- ↳ less performance enhancement

2. Coarse grain Parallelism

- ↳ relatively large amount of computational work
- ↳ hard to facilitate load balancing



Amdahl's law

↳ used in parallel computing

↳ to predict theoretical speed up when using multiple processors

$$S = \frac{1}{(1-P) + P \cdot \frac{1}{S}}$$

Annotations:
S latency
= S
P → Parallelised process
no. of processes

- For example, if a program needs 20 hours using a single processor core, and a particular part of the program which takes one hour to execute cannot be parallelized.

can be Parallelized

19 HOURS

$$P = \frac{19}{20} = 0.95$$

can't be Parallelized

1 hour → critical

$$\frac{1}{20} = 0.05$$

min execution > that critical 1 hour

For example, suppose that we use our strategy to search for primes using 4 processors, and that 90% of the running time is spent checking 2k-digit random numbers for primality (after an initial 10% of the running time computing a list of k-digit primes). Then $P = .90$ and $S = 4$ (for 4-fold speedup). According to Amdahl's Law,

$$\text{overall speedup} = \frac{1}{(1 - 0.90) + \frac{0.90}{4}} = \frac{0.10}{0.225} = 3.077$$

This estimates that we will obtain about 3-fold speedup by using 4-fold parallelism.

Pipelining

↳ each instruction is broken down into stages

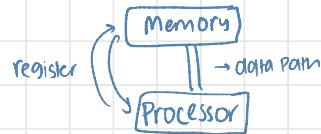
↳ process executes 1 stage in each cycle

CONS

↳ if stage slow, speed slow

↳ if too many conditional jumps, accurate branch prediction needed

↳ misprediction causes large no of instructions to be flushed



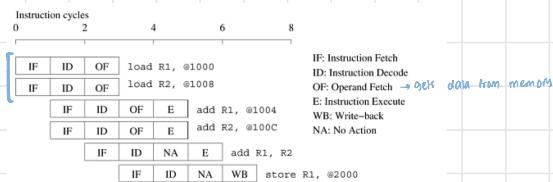
SUPER SCALAR EXECUTION (SUPER PIPELINING)

↳ processor checks which instructions can be run together → done by processor at runtime

↳ if instructions have nothing to do with each other

it can run them parallelly

↳ needs multiple logical units → ALU IF/ID/OF... multiple instructions at a time



CONS

↳ data dependency

↳ branching ??



↳ memory latency

↳ expensive

↳ complex

VERY LONG INSTRUCTION WORDS (VLIW)

↳ instructions that can run concurrently

are packed into groups

↳ then dispatched together

→ solves data dependency issue
done by compiler at compile time
during compile time
group these instructions
→ check next slide

SUPER PIPELINE

CONS [↳ more complex h/w
↳ real time → can't take time as small window]

PRO [↳ view of dynamic state]

VLIW

PROS [↳ more simple h/w
↳ offline → can take time]

CON [↳ no view of dynamic state → reduces accuracy of branch and memory prediction
↳ more difficult branch and memory prediction]

How 2-way superscalar remain idle

Horizontal waste

- not all issue slots can be filled in a cycle

Vertical waste

- Processor issues no instructions in a cycle

(b) Execution schedule for code fragment (i) above.



(c) Hardware utilization trace for schedule in (b).

DATA DEPENDENCY

result of one operation is output to next

in order issue

issue instruction on the order they came in



if 2nd instruction is dependent on 1st only one instruction is issued in cycle

out of order

issue instruction in any order



1st and 3rd can be coscheduled and not 2nd as it is dependent on 1

BRANCH DEPENDENCY

- scheduling instructions across conditional branch statements can't be done



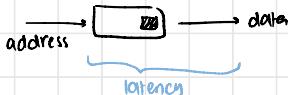
RESOURCE DEPENDENCY

- two instructions require same resource

Limitations of memory system

MEMORY LATENCY

- gap b/w rate at which processor works and time taken to fetch data from memory



SOLUTION IS
CACHE

memory Bandwidth

- rate at which data is moved b/w memory to processor
- increasing blocksize increase bandwidth

CACHE

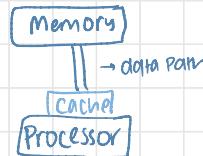
- ↳ smaller and faster memory
- ↳ sits b/w processor and main memory

WORKS

- ↳ if data is in cache
 - get from cache
- ↳ if data is in main memory
 - store in cache

PROS

- ↳ cache has lower latency
 - even lower if repeated data



TEMPORAL LOCALITY

- ↳ notion of repeated reference to a data item in a small time window

CACHE HIT RATIO

- ↳ % of memory access found in cache
- ↳ it defines performance

- Given:
- Fetching the two matrices into the cache corresponds to fetching 2K words, which takes approximately 200 μ s.
- Multiplying two $n \times n$ matrices takes $2n^3$ operations. For our problem, this corresponds to 64K operations, which can be performed in 16K cycles (or 16 μ s) at four instructions per cycle.
- Solution :
- The total time for the computation is therefore approximately the sum of time for load/store operations and the time for the computation itself, i.e., $200 + 16 \mu$ s.
- This corresponds to a peak computation rate of 64K/216 or 303 MFLOPS.

Storing 2 matrices into cache
= 200 μ s → load

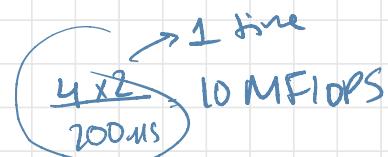
Multiplying 2 matrices → muladd
= $2(32)^3 = 64K$ operations
= $64K / 4$ operations in 1 cycle
= 16K μ s → muladd

Total time

$$\text{FLOPS} = \frac{64K}{200\mu\text{s} + 16\mu\text{s}} = 303 \text{ MFLOPS}$$

$$1 \text{ word} = 4 \text{ byte}$$

$$4 \text{ instructions per cycle} = \text{cache line} = 4$$



load $R_1, [R_2]^{a_{1k}}$ → 100 ns
load $R_3, [R_4]^{b_{1k}}$ → 100 ns
mult, add R_5, R_1, R_3
↓
2 operations



no of operations: $2n^3$
no of iterations: n^3

LATENCY

MATRIX MULTIPLY

$C = 0$

for($i=0; i < n; i++$)

 for($j=0; j < n; j++$)

 for($k=0; k < n; k++$)

$C_{ij} += a_{ik} * b_{kj}$

$$C = A \times B$$

64x64 matrices

4 bytes

latency

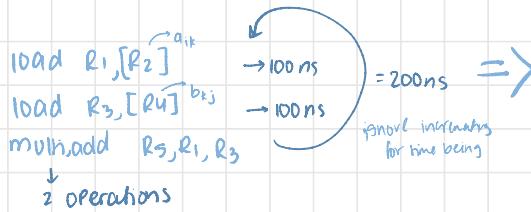
processor: 1 GHz

memory access time: 100 ns

word size: 4 bytes

Whenever we
data from some
location

→ if size not
given of matrix
THEN USEFUL



$$FLOPS = \frac{2}{200 \mu s} = 10 \times 10^6 \text{ FLOPS}$$

↓
10 MFLOPS

execution time

CACHE

$$C = A \times B$$

64x64 matrices

4 bytes

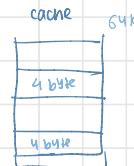
latency

$$MM = 100 \text{ ns}$$

$$\text{Cache} = 1 \text{ ns}$$

cache size = 64K

word size = 4 byte



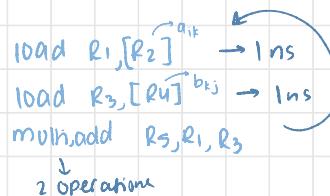
$$\text{matrix A size} = 64 \times 64 \times 4 = 16 \text{ K ns}$$

↑ elem per row
in bytes size

round off

$$\text{B size} = " 16 \text{ K}$$

$$\text{C size} = " 16 \text{ K}$$



$$\text{Total data size} = 16 \text{ K} + 16 \text{ K} + 16 \text{ K}$$

$$= 48 \text{ K}$$

↑ can all fit
in cache

$$\text{total no of iterations} = n^3 = 64^3$$

$$\text{total time} = 64^3 \times 20 \text{ ns} + 800 \mu \text{s}$$

$$FLOPS = \frac{2 \times 64^3}{64^3 \times 20 \text{ ns} + 800 \mu \text{s}} = \approx 90 \text{ MFLOPS}$$

Load all data into cache

$$= 100 \text{ ns} \times 64 \times 64 \times 2 \rightarrow \text{no need to load matrix C}$$

= 800 μs

Storing 2 matrices into cache

$$= 200\text{ }\mu\text{s} \rightarrow \text{load}$$

Multiplying 2 matrices $\rightarrow \text{multadd}$

$$= 2(64)^3 = 524288 \approx 524\text{K}$$

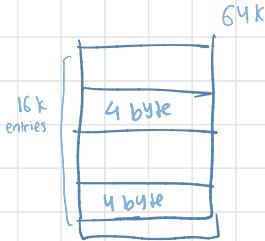
$$= 524\text{K}/4 \rightarrow 4 \text{ operations in 1 cycle}$$

$$= 131\text{K }\mu\text{s} \rightarrow \text{multadd}$$

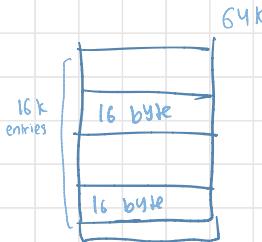
Total time

$$\text{FLOPS} = \frac{2(64)^3}{200\mu\text{s} + 131\text{K }\mu\text{s}} = 4\text{ GFLOPS}$$

BANDWIDTH



increasing blocksize



Dot product \rightarrow doesn't use cache

\rightarrow latency

$$\text{MM} = 100\text{ns}$$

word size: 4 byte

NO CACHE

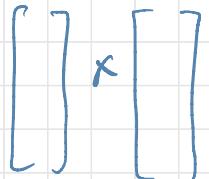
load $R_1, [R_2]^{a_{ik}} \rightarrow 100\text{ns}$
 load $R_3, [R_4]^{b_{kj}} \rightarrow 100\text{ns}$
 mult/add R_5, R_1, R_3
 \downarrow
 2 operations

$$\text{FLOPS: } \frac{2}{200\text{ns}} = 10 \text{ MFLOPS}$$

\rightarrow latency

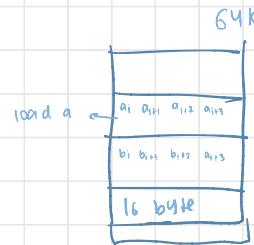
$$\text{MM} = 100\text{ns}$$

word size: 16 byte



CACHE

load $R_1, [R_2]^{a_{ik}} \rightarrow 100\text{ns}$
 load $R_3, [R_4]^{b_{kj}} \rightarrow 100\text{ns}$
 mult/add R_5, R_1, R_3
 \downarrow
 2 operations
 ...
 load $R_1 \rightarrow 2\text{ns}$
 load $R_3 \rightarrow 1\text{s}$ $\rightarrow 2\text{ns} \rightarrow$ ignore for now
 ...
 2 more times



$$\text{FLOPS: } \frac{2 \times 4}{200\text{ns}} = 40 \text{ MFLOPS}$$

SOLVE USING CACHE HIT RATIO

1 \rightarrow MM 3 \rightarrow cache

cache hit ratio = 75%.

$$\text{avg memory access time} = \underbrace{75 \times 1\text{ns}}_{\text{hit ratio}} + \underbrace{25 \times 100\text{ns}}_{\text{miss ratio}} = 25.75\text{ns}$$

$$\text{FLOPS: } \frac{2}{50\text{ns}} = 40 \text{ MFLOPS}$$

50ns

\downarrow

$$25.75 \times 2$$

\downarrow
 2 access to
 the memory

2.2 Consider a memory system with a level 1 cache of 32 kB and DRAM of 512 MB with the processor operating at 1 GHz. The latency to L1 cache is one cycle and the latency to DRAM is 100 cycles. In each memory cycle, the processor fetches four words (cache line size is four words). What is the peak achievable performance of a dot product of two vectors? Note: Where necessary, assume an optimal cache placement policy.

```
1  /* dot product loop */  
2  for (i = 0; i < dim; i++)  
3      dot_prod += a[i] * b[i];
```

DOT PRODUCT OF 2 VECTORS

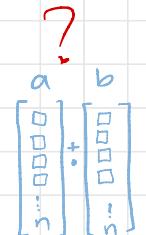
cache size = 32 kB → latency = 1 ns

MM size = 512 MB → latency = 100 ns

Processor = 1 GHz

10^9

cache line size = 4 words



↳ load from Memory

= 100 ns + 100 ns

no of miss = 2
 ↳ load a
 ↳ load b

load a

load b

muladd a, b, c

cache 32kB

time for miss = $2 \times 100 \text{ ns} = 200 \text{ ns}$

no of operations = 2

no of iteration = 4 ← cache size

total no of operation = $4 \times 2 = 8$

$$= \frac{8}{200 \text{ ns}} = 40 \text{ MFLOPS}$$

$\downarrow 10^9$

2.3 Now consider the problem of multiplying a dense matrix with a vector using a two-loop dot-product formulation. The dimension is 4 × 4K. (Each row of the matrix takes 16 kB of storage.) What is the peak achievable performance of this technique using a two-loop dot-product based matrix-vector product?

```
1  /* matrix-vector product loop */  
2  for (i = 0; i < 4; i++)  
3      for (j = 0; j < 4k; j++)  
4          c[i] += a[i][j] * b[j];
```

MULTIPLE MATRIX X VECTOR

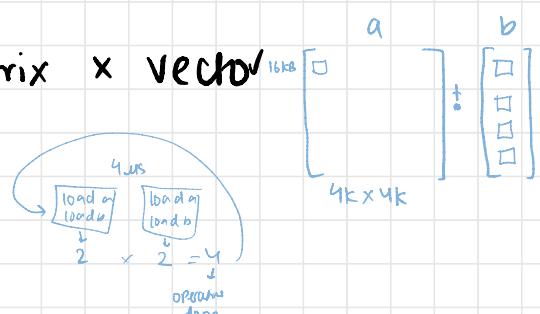
matrix = 4K + 4K

each row takes = 16 kB

Total storage matrix takes

= 16 kB × 4K

= 64 000 kB → can't fit inside cache



no of miss in a = 1

no of operation = 2

time for miss = 100 ns

no of iterations = 4

total operations = 8

storing B in cache

= 100 ns × 4K/U

= 4,000,000 ns

= 1×10^{-4}

$$\text{FLOPS} = \frac{8}{100 \text{ ns}} = 80 \text{ MFLOPS}$$

2.4 Extending this further, consider the problem of multiplying two dense matrices of dimension $4K \times 4K$. What is the peak achievable performance using a three-loop dot-product based formulation? (Assume that matrices are laid out in a row-major fashion.)

```

1      /* matrix-matrix product loop */
2      for (i = 0; i < dim; i++)
3          for (j = 0; j < dim; j++)
4              for (k = 0; k < dim; k++)
5                  c[i][j] += a[i][k] * b[k][j];

```

MULTIPLY MATRIX INTO MATRIX

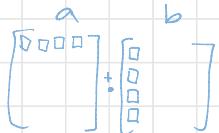
cache size = 32 kB \rightarrow latency = 1 ns

cache line size = 4 words

MM size = 512 MB \rightarrow latency = 100 ns

Processor = 1 GHz

10^9



Total storage matrix takes

$$= 16 \text{ kB} \times 4K$$

$= 64000 \text{ kB} \rightarrow$ can't fit inside cache

no of miss = 5

time for miss = 500 ns

How

no of operations = 2

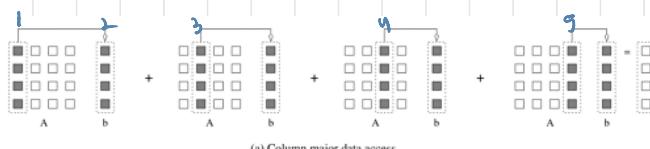
no of iterations = 4

total operations = 8

$$\text{FLOPs} = 8$$

$$= 16 \text{ MFLOPs}$$

$$500 \text{ ns}$$



load a

load b

muladd a, b, c

Example 2.4 Effect of block size: dot-product of two vectors

Consider again a memory system with a single cycle cache and 100 cycle latency DRAM with the processor operating at 1 GHz. If the block size is one word, the processor takes 100 cycles to fetch each word. For each pair of words, the dot-product performs one multiply-add, i.e., two FLOPs. Therefore, the algorithm performs one FLOP every 100 cycles for a peak speed of 10 MFLOPS as illustrated in Example 2.2.

$$\frac{1}{100 \text{ ns}} = 10 \text{ MFLOPs}$$

Now let us consider what happens if the block size is increased to four words, i.e., the processor can fetch a four-word cache line every 100 cycles. Assuming that the vectors are laid out linearly in memory, eight FLOPs (four multiply-adds) can be performed in 200 cycles. This is because a single memory access fetches four consecutive words in the vector. Therefore, two accesses can fetch four elements of each of the vectors. This corresponds to a FLOP every 25 ns, for a peak speed of 40 MFLOPS. Note that increasing the block size from one to four words did not change the latency of the memory system. However, it increased the bandwidth four-fold. In this case, the increased bandwidth of the memory system enabled us to accelerate the dot-product algorithm which has no data reuse at all.

$$\frac{4 \times 2}{200 \text{ ns}} = 40 \text{ MFLOPs}$$

Another way of quickly estimating performance bounds is to estimate the cache hit

Example 2.5 Impact of strided access

Consider the following code fragment:

```
1  for (i = 0; i < 1000; i++)  
2      column_sum[i] = 0.0;  
3      for (j = 0; j < 1000; j++)  
4          column_sum[i] += b[j][i];
```

The code fragment sums columns of the matrix `b` into a vector `column_sum`. There are two observations that can be made: (i) the vector `column_sum` is small and easily fits into the cache; and (ii) the matrix `b` is accessed in a column order as illustrated in [Figure 2.2\(a\)](#). For a matrix of size 1000 x 1000, stored in a row-major order, this corresponds to accessing every 1000th entry. Therefore, it is likely that only one word in each cache line fetched from memory will be used. Consequently, the code fragment as written above is likely to yield poor performance. ■

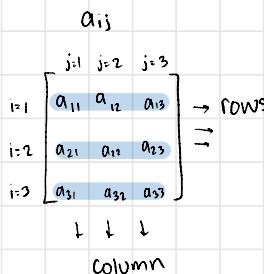
A strided access pattern accesses a sequence of addresses with a uniform skip between each referenced address. For example, sequence 1, 1001, 2001, 3001, 4001, 5001, ... is a strided access pattern with a stride of +1000.

Strided access patterns refer to memory access patterns where data elements are accessed with a fixed step or stride between consecutive accesses. Instead of accessing consecutive memory locations, the address of each accessed element is computed by adding a fixed offset or stride value to the previous address.

The above example illustrates problems with strided access (with strides greater than one). The lack of spatial locality in computation causes poor memory system performance. Often it is possible to restructure the computation to remove strided access. In the case of our example, a simple rewrite of the loops is possible as follows:

Example 2.6 Eliminating strided access

- Strided access patterns can have a significant impact on performance, especially in memory-bound applications.
1. Cache misses: When accessing data with a large stride, the likelihood of cache misses increases because the cache may not be able to prefetch or store the required data in advance. This can lead to additional latency and inefficient memory access.
 2. Reduces Spatial Locality: When data elements are not accessed consecutively but instead with a stride, it introduces gaps in memory access patterns, reducing spatial locality. This can result in less effective utilization of cache lines and decreased cache hit rates, negatively impacting performance.
 3. Reduce Memory Bandwidth utilization
 4. Increase memory latency



Alternate Approaches for Hiding Memory Latency

- Consider the problem of browsing the web on a very slow network connection. We deal with the problem in one of three possible ways:
 - we anticipate which pages we are going to browse ahead of time and issue requests for them in advance;
 - we open multiple browsers and access different pages in each browser, thus while we are waiting for one page to load, we could be reading others; or
 - we access a whole bunch of pages in one go - amortizing the latency across various accesses.

- **Prefetching**
- **Multi threading**
- **Spatial locality**

Prefetching

- ↳ Used to speed up fetch operations
- ↳ It involves
 - ↳ Loading of a resource before it's required to decrease waiting time of that resource

CON

- ↳ might need more space to store advanced loads

Multi Threading

- ↳ run multiple instruction threads at the same time

PROS

- ↳ concurrent execution of multiple threads
 - ↳ allows processor to work on other tasks while waiting for memory operation to complete

Spatial Locality

- ↳ accessing nearby memory locations which can improve
 - ↳ cache utilization
 - ↳ memory performance

May become
bandwidth bound

Tradeoffs of Multithreading and Prefetching

- Bandwidth requirements of a multithreaded system may increase very significantly because of the smaller cache residency of each thread.
- Multithreaded systems become bandwidth bound instead of latency bound.
- Multithreading and prefetching only address the latency problem and may often exacerbate the bandwidth problem.
- Multithreading and prefetching also require significantly more hardware resources in the form of storage.

Interconnection Networks

↳ carry data b/w processors and memory

Static

↳ links → wires

↳ links connect nodes together

Dynamic

↳ links + switches

↳ a device which redirects data from any input port to output ports

↳ establish paths among nodes and memory banks

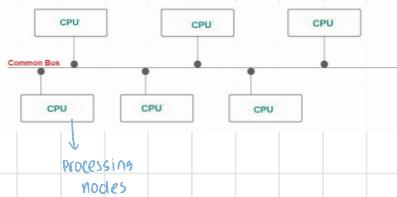
→ BUS Based Network

↳ processors access common bus to exchange data

↳ distance b/w 2 nodes is $O(1)$ in bus

↳ provides broadcast media

↳ scalable cost



→ CROSSBAR NETWORKS

→ dynamic $O(P^2)$

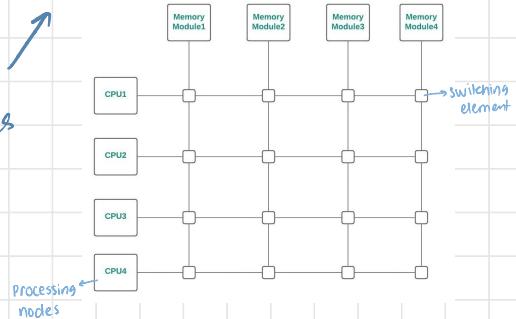
↳ connects processors to memory banks

↳ a grid of switches

↳ non blocking network → doesn't block other

↳ scalable performance

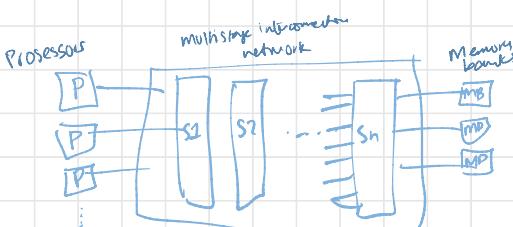
↗
NETWORK
TOPOLOGIES
↘



→ Multistage Network

↳ lies b/w bus and crossbar

↳ more scalable ↳ cost performance



Interconnection Networks: Network Interfaces

- Processors talk to the network via a network interface.
- The network interface may hang off the I/O bus or the memory bus.
- In a physical sense, this distinguishes a cluster from a tightly coupled multicomputer.
- The relative speeds of the I/O and memory buses impact the performance of the network.

HOW GOOD A NETWORK IS

↳ diameter → latency

↳ no of links → cost

↳ degree of nodes → practicality/cost

↳ bisection bandwidth

↳ amount of data that can pass
through 2 halves of system

* degree of switch = no of ports