Chapter 3 Design for Parallel

2 Design for Parallel

Single processor design

fastest; smallest; algorithm

You have a single processor so it corresponds to how you might solve the problem.

Multi-processor design: need to design in parallelism.

Communications are the biggest problem.

Bandwidth is much more limited than cpu cycles.

The CPUs are running independently and it is difficult to provide a single consistent time - and yet communication is normally relevant only when both processors have reached a particular point, which requires synchronisation.

May wish to write to a common data structure, but then we need to limit access to the data structure to a single process.

Both problems have solutions, but neither are simple and both tend to slow down the calculation. Simulations are embarrassingly parallel

3 Design for Parallel

Minimal communication – ideally zero between processors.

Example

For a film you need 25 frames a second. Each frame takes 4 seconds to generate. Assuming perfect scaling with 100 processors you can generate a frame every 1/25 of a second on average.

Trivial for creating a film. One frame per free processor: first frame takes 4 seconds and after that they arrive at 25 per second.

Put them in order and display



Hard for real time generation in say a game – the image is not ready until the user has made a decision – a frame cannot be generated until user action, thus we need to split the image into 100 bits and generate them all separately.

Light source in one part of the image, reflects off object in a second part modifying colour and illuminates something in a third part.

Communication – potentially non-local

Simulations are embarrassingly parallel

Data Parallelism

different parts of the data are processed on different units. Same code, different data, but not SIMD where the processors are executing instructions in lockstep.

Functional Parallelism

different subtasks are executed on different processors *MPMD* (multiple programme multiple data).

A particular processing step executed for all the data,

Results passed on to further processors for further processing

(Must be possible to pass on part of the data set)

Historical note

Prior to invention of digital electronic computers – "computers" were people who worked together in large rooms doing calculations. On the Los Alamos Project Richard Feynman speeded up calculations by working out a way to parallelise them

Continuous system (differential equation) – replaced by a discrete system – Geometrical Grid – and Finite difference equations.

Electric fields, magnetic fields, temperature distribution, stress calculations, diffusion, weather, fluid flow.

 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

A cpu calculates the values for a range of grid points.

These require input from a cell and some of its neighbours.

Not a grid in the module title

6 Data Parallel

Continuous system (differential equation) – replaced by a discrete system, difference equation.

Every step in the calculation requires input from a cell and some of its neighbours. Input changes at every step.

Domain Decomposition

Normally each core will deal with more than one point.

How to decide on the mapping of the points to the cores?

Communication is needed between cores to provide updated values.



Parallel Design

 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

7 Topology

Each core needs memory for points not calculated but needed for calculation *halo* or *ghost* values.

Moving through all points by a core is a *sweep*.

Updating occurs when all cores have finished a sweep.

Communication is much slower than calculation

Loads must be balanced across processors – waiting is inefficient (and wastes power).

Communication overheads depends on the size of the **stencil.** Value, first derivative and second derivative depend on nearest neighbour. Higher derivatives require longer range communications.



More complex update schedules may be possible

 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

 0
 0
 0
 0
 0
 0
 0
 0

8 Longer Range

Overhead on longer range interactions depend on topology.

Fully connected



9 Functional

Master Slave

Break the problem into distinct independent problems.

CGI movie – each frame

Weather predictions – each starting point.

Master knows about the work units.

Distributes to slaves.

Collects and collates results.

Problems – bandwidth/congestion/single point of failure.

Lost jobs

This works especially well with clusters – collections of independent processors, rather than the tightly connected special purpose systems.

So cloud – although there are still significant problems to an efficient solution.

Trivially parallelisable

10 Design for Parallel	Going to talk about a way of approaching the task of solving a problem on a distributed system.	
	Might think that defining how a disperforms in relation to a single proeasy.	stributed system cessor system is
	There are subtleties	
	In order to quantify improvements understand metrics – both to "sell" you come up with and to understan people might describe their solutio	you need to any solution that nd the way other ns.
Design always needs	Parallel design: hard	
intelligent application	An approach – ideas to bear in mind.	
	Partition Communication	Concurrency and scalability
	Agglomeration Mapping	Locality

1	1	Partition

May need to agglomerate later

Looking for possibilities for parallelism. Maximum number – smallest unit of work.

Fine grained decomposition.

Greatest flexibility in terms of building solution. Easier to consolidate than divide.

Divide computation and data

OO design naturally produced good targets for parallelism.

Object includes data and code.

Objects should be as small as possible (they do one task) and with minimal coupling.

Domain decomposition: break up the data and associate the code with it.

Functional decomposition: break up the operations and associate the data.

Clearly increases communication

One object per processor is always possible.

Parallel Design

Old technique

12	Partition	Check
----	-----------	-------

Will not scale

If tasks get larger then solution will not scale

Following checklist

Does your decomposition have at least an order of magnitude more tasks than processors?

Does it avoid repeated computation and multiple storage?

Are the tasks equal size?

Does number of tasks should scale with problem size?

Do you have alternatives?

A large task which cannot be decomposed behaves like a serial part of the execution.

Need flexibility

Allocation hard

13 Communication check

Communication is via a **channel**

One end sends, other end receives

Define

channels – links between producers and consumers.

Messages – data which travels down the channel.

Communication should be spread widely A central communication hub is a bottleneck Many messages should be possible *concurrently*

Communication modes

Local: only with "near" neighbours

Structured: communicating nodes for a regular structure.

Static: does not change with time. Might be known at compile or submit time. *Dynamic* changes.

(A)Synchronous: both parties have to be on the line at the same time.

Send is asynchronous Receive is synchronous

For grid "near" means high bandwidth, low latency.

14 Example

Look at an example of a problem.

Definitely suitable for distributed solution

Consider communication patterns and their implication

15 Local communication

Consider a calculation done on a rectangular grid. (stick to 2d)



Another meaning of grid computing

16 Finite Element



Another meaning of grid computing

Many applications

Many problem domains

17 Difference from differential



Another meaning of grid computing

Many problem domains

18 Stencil



Shown are the channels. Each arrow is two channels

Boundary conditions

The **stencil** is the pattern of communication channels to the neighbours.

Rectangular grid and nearest neighbour leads to the **5 point stencil**.

 $\psi^{t}(i,j) = 1/8 \left[\psi^{t-1}(i,j-1) + \psi^{t-1}(i-1,j) + \psi^{t-1}(i+1,j) + \psi^{t-1}(i,j+1) + 4*\psi^{t-1}(i,j+1) \right]$

The important thing to look at here is that the values of all the points depend on the values of the four nearest neighbour at the previous step.

19 Traversal pattern

The initial update strategy seems to be a sensible one.

We update all nodes in turn. We can of course think of various orders in which to do the update.C

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	Obvious is x then y (or y then x)	Trivial extrapolation to higher dimensions
	A wavefront travelling from bottom leftt to top right	
	Red-black or <i>checkerboard</i>	Parallel Design

20 Simple replace	Which ever order we calculate the nodes for the simple replace <i>Jacobi</i> algorithm.	
	We conceptually just replace all elements at step t, with all elements at t+1 <i>simultaneously</i> .	
<u> </u>	Some coordination needed.	
	Increment a flag when you complete and look for a value. <i>Global.</i>	
	Send to the output channels when complete and listen on input channels for neighbours to finish. <i>Local</i> .	Are we guaranteed not to get out of step?
	Convergence is slow.	
		Parallel Design







No red node depends in value on any other red node No black node depends in value on any other black node.

Update all red, then update all black. Like a leapfrog update.

Need to look at – efficiency of speed-up.

Communication overhead.

Efficiency of algorithm.

If A speeds up less than B, but B needs more step to produce the answer. The correct algorithm is not trivial.

Trivial extrapolation to higher dimensions

Parallel Design

Looking for answer not efficiency

What if communication is asynchronous?	
Producer doesn't know when result is required.	
FEA result is always required as produced.	
Large data structure – all of which can be written by many tasks.	
Data structure I/O can be structured as a single task.	Bottleneck
Structure distributed among computational tasks.	
Requirement to check for I/O requests complicates code	
Separate set of communication tasks. No data locality	
	Parallel Design

25 Communication check

Checklist for communication.

All tasks perform the same number of communication operations. If unbalanced worry – can we duplicate data to reduce number of operations.

Are communication partners local or global. Local is best.

Possibility of concurrent communications?

Can the computation in different tasks proceed concurrently. Consider order of communication and computation tasks.

26 A	Agglor	neration
------	--------	----------

Agglomeration

Stages 1 and 2 produce an abstract solution. Stages 3 and 4 look at the particular architecture.

How do we put the bits together to create an efficient and scaleable solution?

Do we need to duplicate data?

Do we aim at one task per processor?

Consider granularity, flexibility, engineering costs Balance sometimes conflicting requirements.

<u>Granularity</u>

Communications can form a big part of parallel computation.

Reducing granularity may reduce message volumes beneficially.

It will reduce total task creation time.

May be constant message volume but reduced message number is desirable.

Alternately leave to mapping

When message start up is a significant overhead





Only true for local communications

28 Agglomeration

Agglomeration on the Grid

Traditional parallel asks about multiple tasks on one machine versus spreading through a cluster.

Parallel v. communication.

Cloud is normally in one data centre - may be communication problems – is bandwidth in a rack as good as communication between racks – can you specify location?

What if we want to run on several sites.

Communication overheads are large between sites.

Problem of ensuring co-ordination between sites.

Booking bandwidth between sites.

Remember looking for optimum solution. Not the most scalable algorithm. (Normally)

29 Mapping	Allocating tasks to processors	
	Common sense.	
Can you specify with Cloud provider?	Place tasks which need to communicate close. Separate tasks which can run concurrently	
	Mapping problem is NP complete.	
		In effect it is unsoluble
	Of then the load per task is stochastic and some sort of balancing during execution is required.	
	Not enough to send 100 tasks to 100 processors, one may end up taking 10 time as long as the others.	
Load balancing should not use more time than it saves!	Solution then arrives much later.	
	In most complex cases	
	number of tasks changes dynamically	
	size of tasks appears dynamically.	
	Multiple operations – local algorithms cause less overhead.	
		Parallel Design

30 Mapping algorithms

Send part of problem out and sub-divide later.

Recursive bisection

Divide the problem into subdomains

Aims

Equal computational cost

Minimum communication overhead.

Repeat – the algorithm can itself be executed in parallel.

Bisection may not give best results in complex geometries. *Unbalanced recursive bisection.* Break into P parts. Same aim.

Balance computational cost, minimise communications.

Other more complex "recursive spectral bisection"

Pothen et al SIAM J. Mat. Anal. Appl. 11(3) 1990

Reviews

H. Simon. Computing Systems in Engineering 2(2/3) 135-148 1991

R. Williams. Concurrency: Practice and Experience 3(5)457-481 1991

Local algorithms

For instance compare load with other local processors Transfer load to more lightly loaded processors.

Multi-level – distribute loads at high level to local "controllers". Controllers communicate with "local" nodes to balance load.

(For some problems "Controllers" can "request" more load after start.)

Reduces communications overheads.

Random

Allocate tasks at random.

Many tasks and expect the load to be equal (probably) computationally very light.

32 Scheduling	
---------------	--

Scheduling

If we have only weak locality constraints.

Chop the task up into a pool of tasks. Free processors are allocated tasks.

Works on the grid for particle physics. Works on cloud for similar problems

Manager - worker. Very simple to implement.

LHC use "pilot" jobs – investigate the performance of nodes and adjust submission appropriately.

Hierarchical manager-worker. Looks like distributed controllers. Controllers can re-allocate after start. Managers allocate to free.

Split the task into tasks with guaranteed "longest time" – manager works.

Managers simpler than controllers.

Parallel Design

Easy for suitable problems

A single manager is susceptible to failure.

"Elections"

33 Termination

How does a processor know when it is finished? Central co-ordination is easy – ask the boss.

Completely independent tasks simple – finished the allocated task and that is your job done.

Decentralised – much harder ...

Particle Physics problems – are completely independent and so finish and exit works.

Some of the LHC experiments have found the overheads in sitting in the queue unacceptable.

Run a task – which contacts central server and asks for events for this node.

Blocks a slot in a processor farm – causes problems.

Local management lose control of their resources.

Tasks creating other tasks have security implications

34 Mapping checklist

Have you considered both a static solution when all tasks are created and then run

and

a solution which includes dynamic task creation and deletion

If you have centralised load balancing schemes have you ensured the central point is not a bottleneck.

For dynamic load balancing have you considered several different strategies.

If using probabilistic or cyclic methods do you have enough tasks to distribute. At least an order of magnitude more tasks than processors

35 Failure	Designing for Failure	
	Distributed Computing	
	"Where your can fail due to the action(s) of a computer of which you have no knowledge."	
	Failures are inherent in solving large problems on Cloud	For instance paper on RGMA testing
And is	Infrastructure needs to be resilient – and is designed as such.	
	Need to confirm cloud infrastructure of the provider will is resilient. Connections to the cloud are as resilient as the connections to the internet.	
	Need to ensure workflow system can cope with failures	
	The infrastructure does not protect individual computing tasks.	
	Failures can be classified according to their effect.	
	Understand the morphology and protect your jobs.	
	Protection may involve manual intervention.	Parallel Design

36 Failure morphology





Crash

processor halts. Does nothing else *in particular no illegal messages.* Never restarts.

Crashes are not detectable in asynchronous systems crashes are indistinguishable from slow connections.

Link failure

A communication link fails. It stays failed. Link failure may reduce communication bandwidth **OR** it may partition the network, some pairs of nodes can never communicate

Omission

Only a proper subset of the required messages is *sent* or *received*.

Byzantine fault

System fails with arbitrary behaviour.

A system which can tolerated Byzantine behaviour can tolerate any faults.

NETWORK PARTITONED

THAT is no country for old men. The young In one another's arms, birds in the trees - Those dying generations - at their song, The salmon-falls, the mackerel-crowded seas, Fish, flesh, or fowl, commend all summer long Whatever is begotten, born, and dies. Caught in that sensual music all neglect Monuments of unageing intellect.

37 Failure mo	rphology
(i)	

Link failures

Some analyses assume that there are no link failures on the basis that they can always be modelled by other forms of failures.

Useful in GC especially if we are looking at failure probabilities.

Heterogeneous processor/network – failure probabilities can be shuffled between sources.

GC with multiple CPU and link types the allocation can be made but would be rather artificial.

Two identical CPUs – different sites so different failure probabilities.

Byzantine fault

Hardest to protect against. Various theorems.

- N processors with t failures need to run correctly.
- t < N/2 for benign failures
- t < N/3 for malign failures

NETWORK PARTITONED



Synchronous v Asynchronous

What is the difference?

Asynchronous means that a process can ask for some communication and then go away and do something else.

Synchronous means that having asked the process waits.

Implication – synchronous system there is a maximum wait time at which point the end-point is assumed to have failed in some way.

Asynchronous – there is no maximum time.

A slow asynchronous system is indistinguishable from a system where one part has died.



Computer theory tells you reasoning about asynchronous systems is much harder.

Making a reliable system is likely to be easier if you insist on synchronicity.