

# *Experiences in Parallel System Development*

---

POETS



Alex Rast  
University of Southampton  
UKDF 2019  
Manchester, UK

Andrew Brown	<i>adb@ecs.soton.ac.uk</i>
Graeme Bragg	<i>gmb@ecs.soton.ac.uk</i>
Alex Rast	<i>A.D.Rast@soton.ac.uk</i>
Mark Vousden	<i>m.vousden@soton.ac.uk</i>

# *What is a Parallel System?*

## Multiprocess

Multiple independent processes  
(units of execution)

## Multi-data

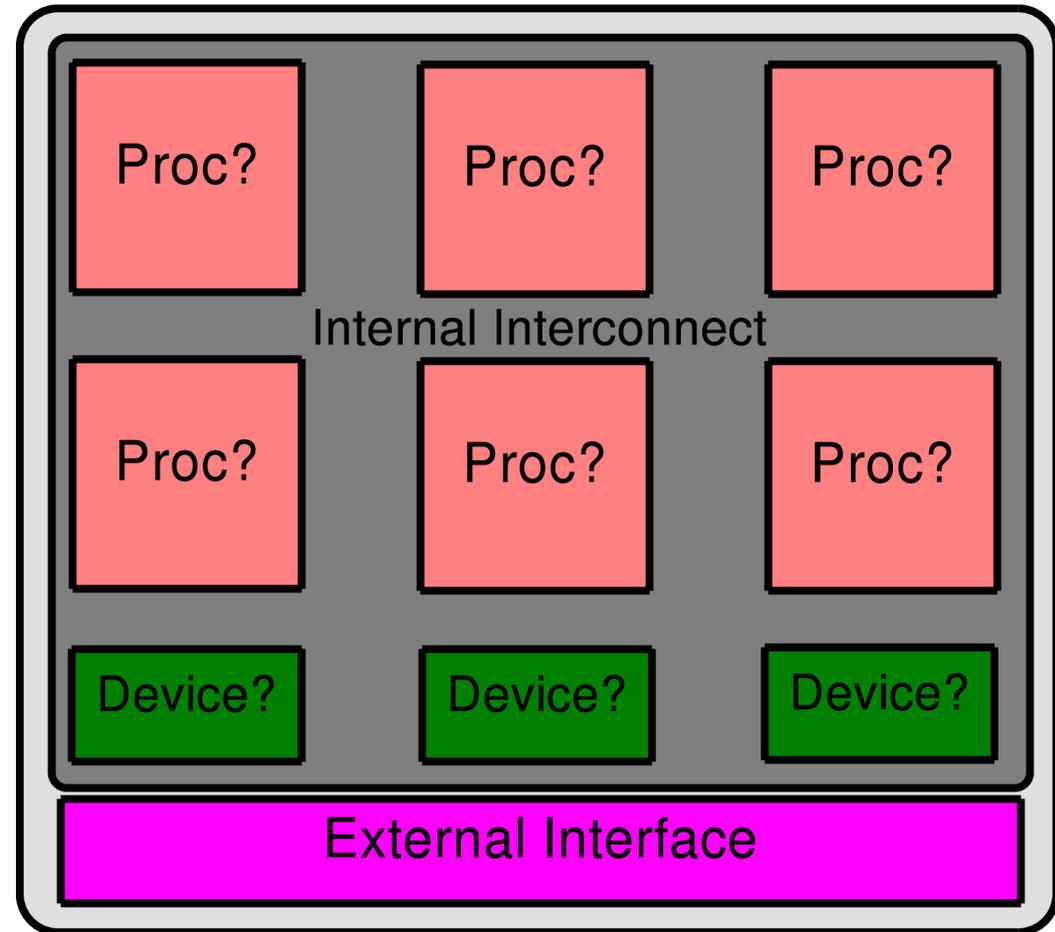
Multiple data streams  
concurrently in flight

## Multicore

Multiple execution engines  
(hardware units)

## Multipath

Multiple independent interconnect  
links capable of carrying the same  
data



# Parallel Hardware I: Synchronous Multicore

## Standard CPUs

Virtually all modern desktop PCs are multicore.

## GPUs

Massive but specialised and heterogeneous parallelism.

## Specialist Multicore

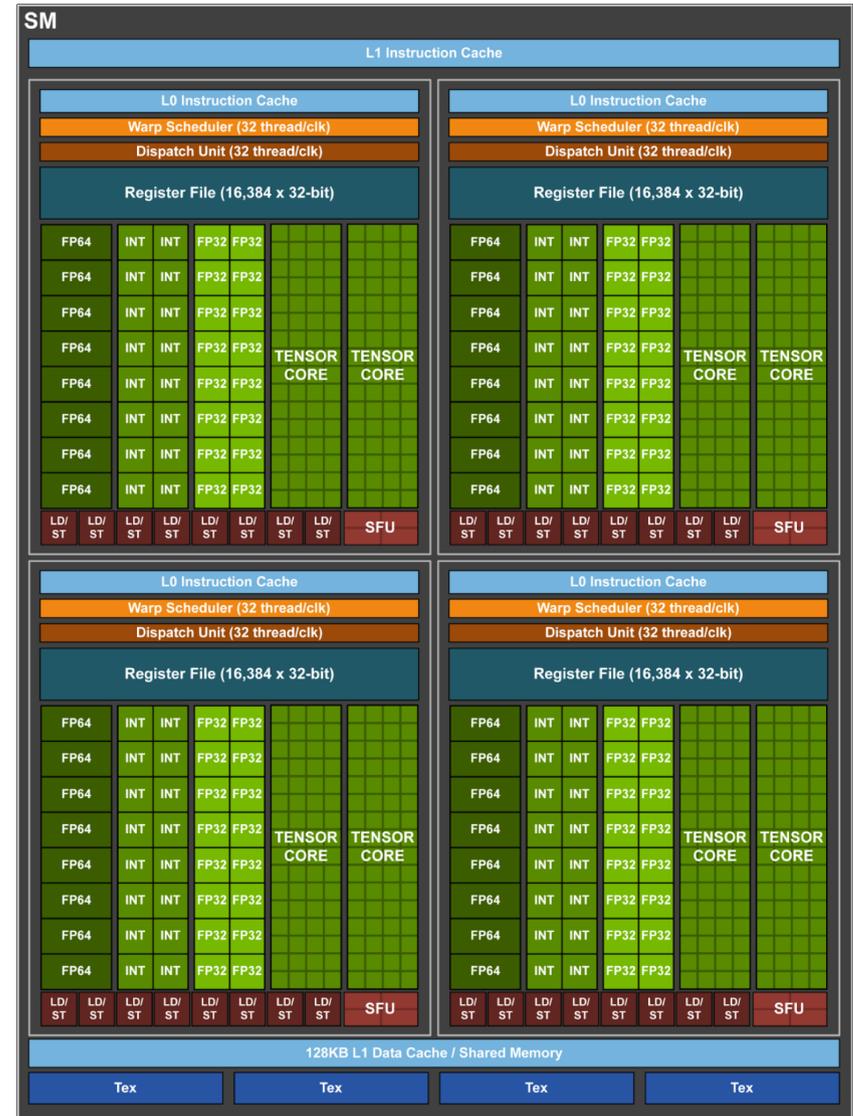
e.g. Intel Xeon PHI, Adapteva Epiphany. Specs vary.

## Synchronous Interconnect

Internal links, whether busses or NOCs, use clocked communication.

## Usually Memory Coherent

Typically, chips use shared memory and cache hierarchies.



NVIDIA Volta GV100 Tile – From 'NVIDIA Tesla V100 Architecture'

# Parallel Hardware II: (Partially) Asynchronous Multicore

## Small Cores

Fairly basic (RISC-V, ARM) RISC cores. Internally synchronous.

## Asynchronous NoC

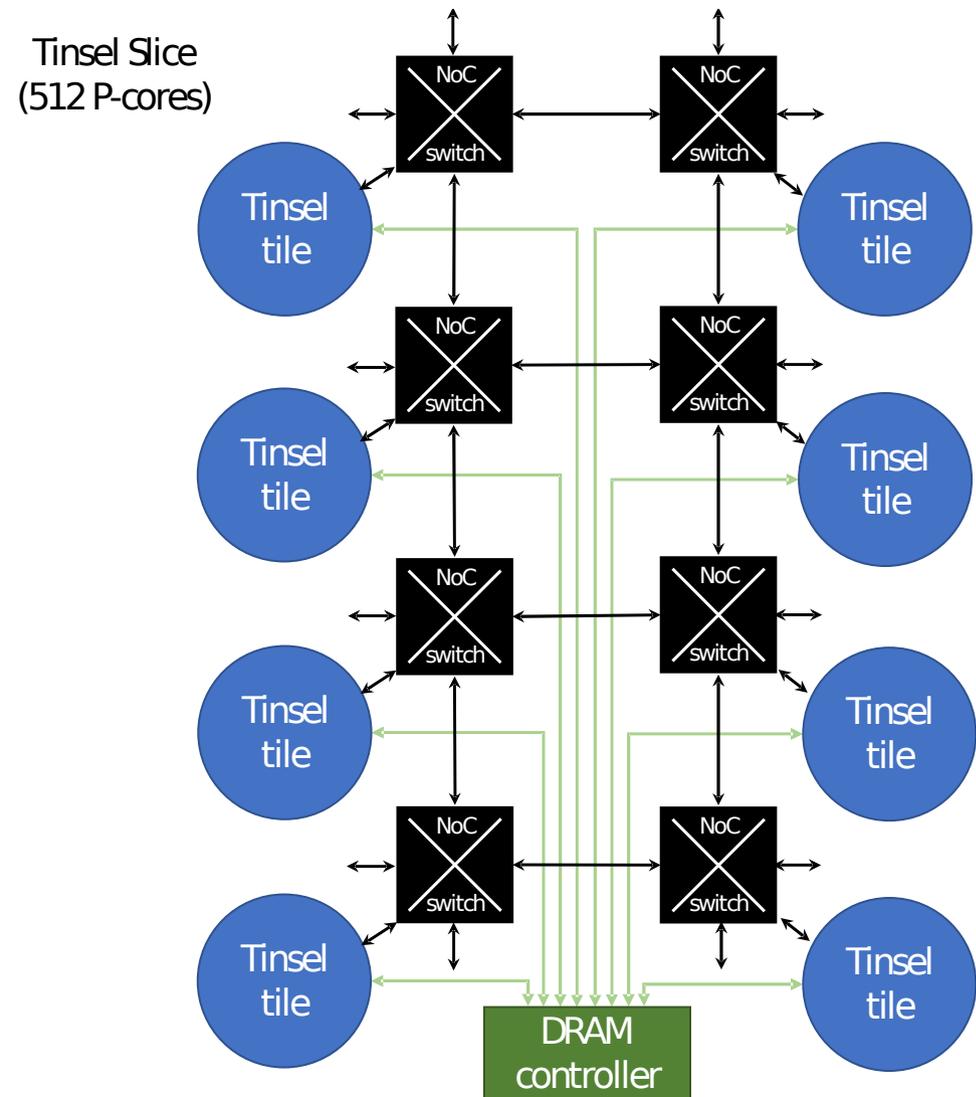
Messages are routed from device to device over an arbitrary path

## Event-Driven Control

Processors respond to events, not program sequences

## No Memory Coherency

Memory is effectively private. No attempt is made to guard against simultaneous accesses to anything.



# Horses for Courses

## Conventional Multicore

- +: Easy to program, comprehensive tool support, zero learning curve
- : Low parallelism, hard to optimise

## GPU

- +: Mature tools, high performance in the right application, NVIDIA is taking over the world
- : The problem needs to look like matrix-bashing, memory accesses must be carefully synchronised, power-hungry

## Specialist Multicore

- +: Tunable for optimal performance, relatively problem-agnostic
- : Steep learning curve, partial support, hardware obsolescence

## Asynchronous Multicore

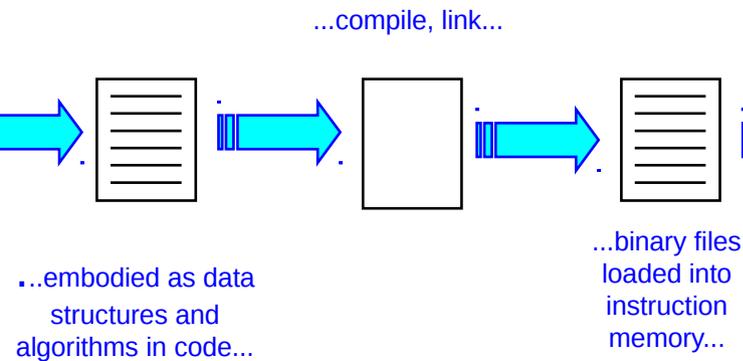
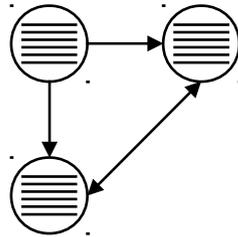
- +: Low power, intuitive application model, no need to worry about synchronisation or coherency
- : Immature tools, debugging is hard, research technology, implementation pitfalls

... *shape* of the application matters

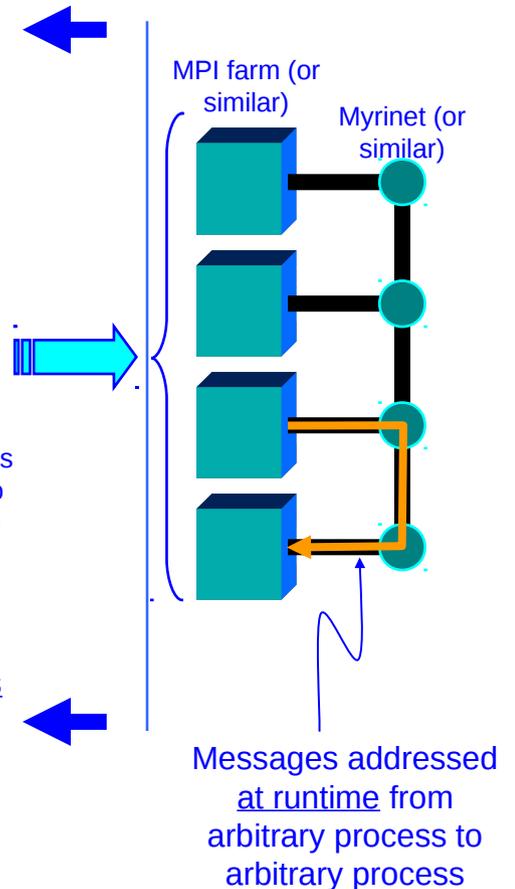
# But How to Program These Things?

A **conventional** multi-processor program:

Problem: represented as a network of programs with a certain behaviour...



Interface presented to the application is a homogenous set of processes of arbitrary size; process can talk to process by messages under application software control



- Most Existing Codebases are Sequential
- Parallel Programming is Hard
- Is Language the Solution?
- You Still Have to Map it to the Hardware

# Parallel Software I: Of Processes and Threads

## MPI

A heavy-duty standard for message-passing

## Processes and Comms

Physical ↔ Logical resource mapping.

Thread support available.

Multicomm support is spotty.

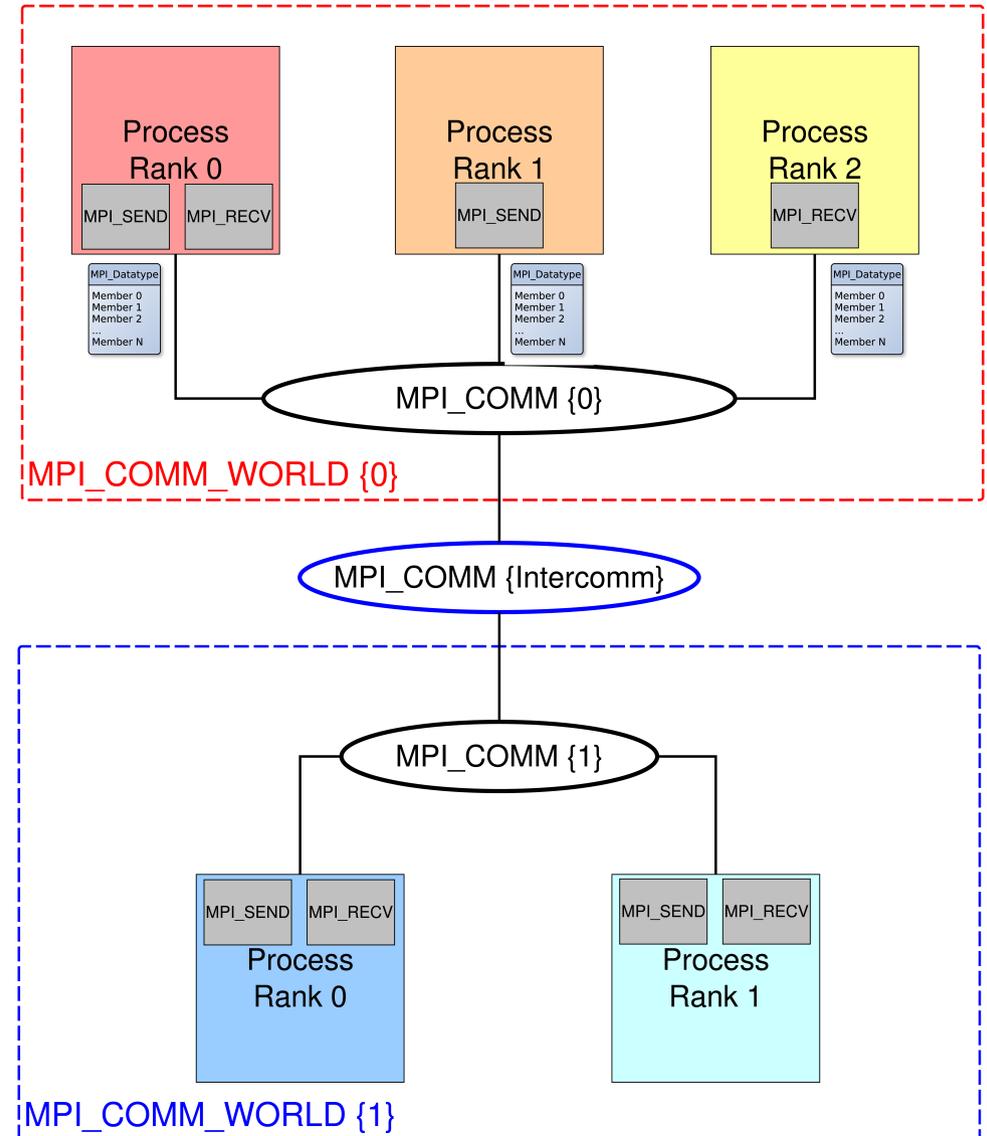
## Message Send/Receive

Relies on highly choreographed communications.

May be blocking or nonblocking.

## Diverging Implementations

MPICH, OpenMPI, MS-MPI, and others, each with different support and little cross-platform interoperability



# Parallel Software II: Of Graphs and Messages

Processing Element = Node

Performs some *simple* function

Communications Link = Edge

Represents an abstract message path between nodes

Event = Message

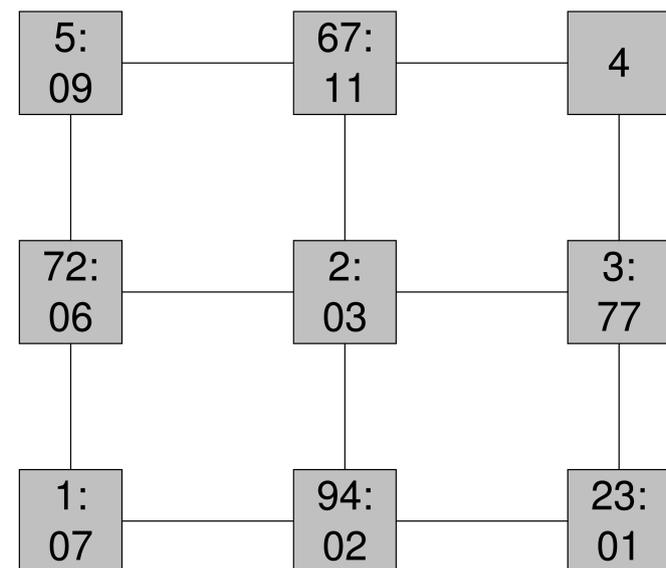
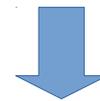
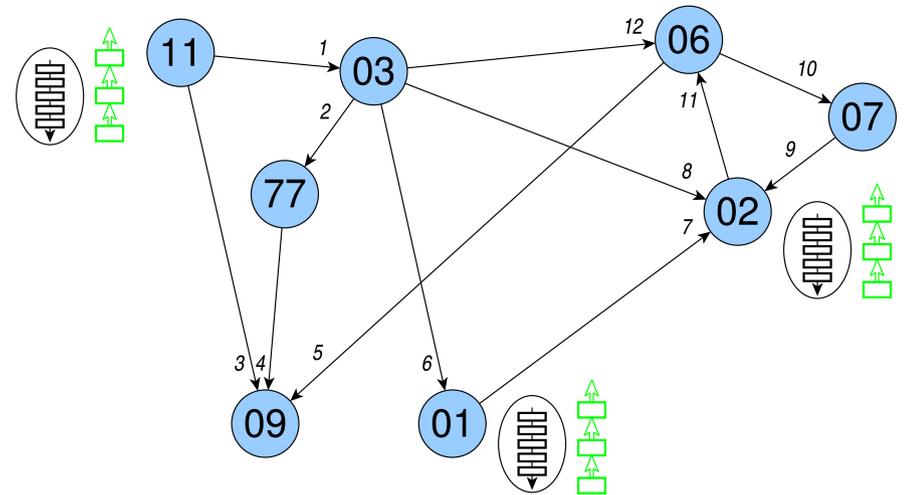
Sends a small data item and a notification

User Specifies the  
Application Graph

Describes the *program* topology

System Software Maps to a  
Hardware Graph

Application  $\leftarrow \rightarrow$  Hardware ID mappings



# An Old Problem: Clock Trees

## Parallel

Models fanout of clock buffers; in large trees there can be  $\sim 10^7$  devices

## Extremely Simple

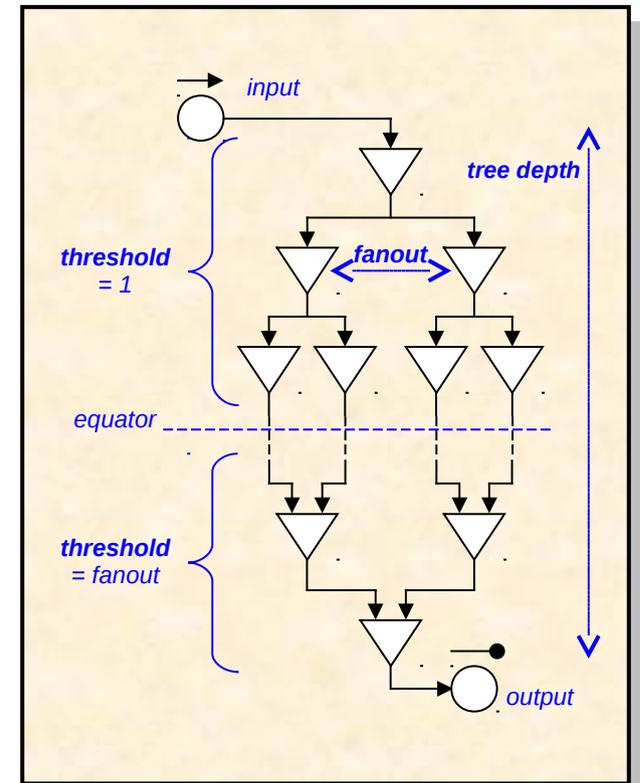
Each device simply records and forwards messages; each message is nothing more than a point event

## Synchronous or Asynchronous

Data flow does not depend on explicit synchronisation and is free from side effects

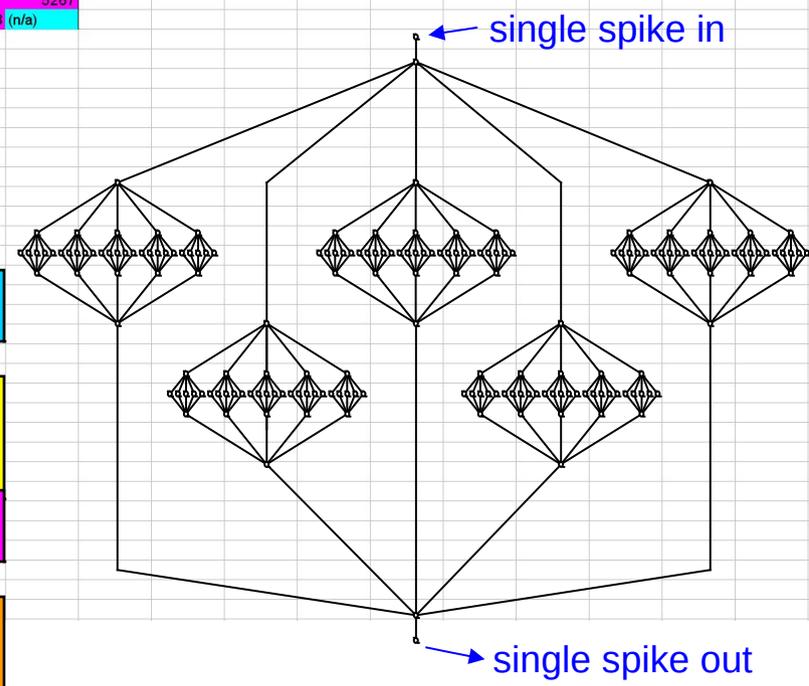
## Worst-case Modelling

The layer at the equator generates intense packet storms



# Clock Trees on SpiNNaker

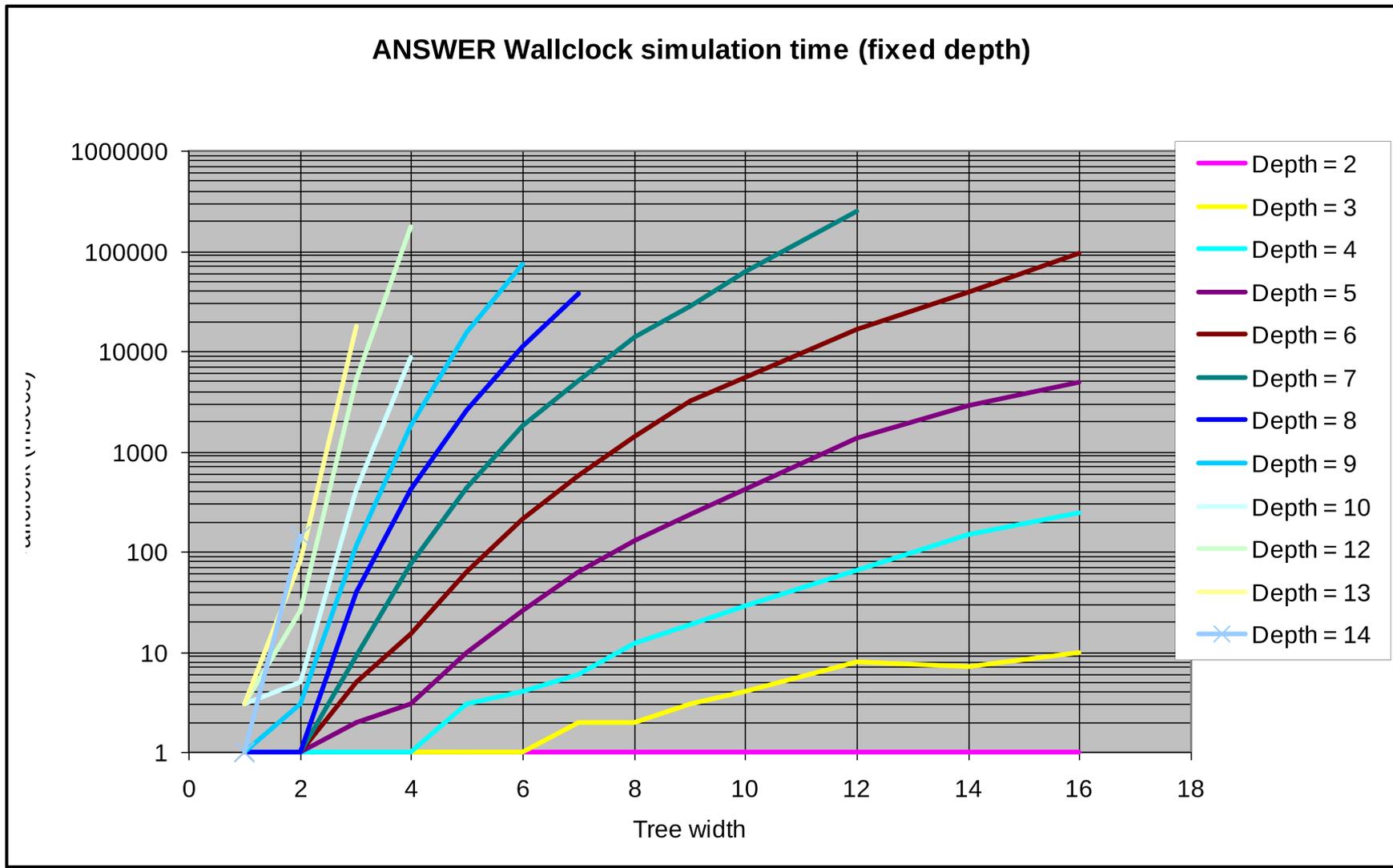
(1, 1)	1	(2, 1)	1	(3, 1)	1	(4, 1)	1	(5, 1)	1	(6, 1)	1	(7, 1)	1	(8, 1)	1	(9, 1)	1	(10, 1)	1	(12, 1)	1	(14, 1)	1	(16, 1)	1			
3																												
(1, 2)	5	0.1	10	0.1	17	0.1	26	0.1	37	0.1	50	0.1	65	0.1	82	0.1	101	0.1	122	0.1	170	0.1	226	0.1	290	0.1		
(1, 3)	7	0.1	22	0.1	53	0.1	106	0.1	187	(n/a)	302	0.1	457	0.1	658	0.1	911	0.1	1222	0.2	2042	3	(12, 3)	4	(14, 3)	7	(16, 3)	10
(1, 4)	9	0.1	46	0.1	161	0.1	426	0.1	937	0.1	1814	0.5	3201	1	5266	2.5	8201	10	12222	(n/a)	24506	(n/a)	44326	(n/a)	74274	(n/a)	146	
(1, 5)	11	0.1	94	0.1	485	0.1	1706	0.2	4687	1.5	10886	10	22409	(n/a)	42130	(n/a)	73811	(n/a)	122222	(n/a)	239	(12, 5)	575	(14, 5)	1213	(16, 5)	2322	
(1, 6)	13	0.1	190	0.1	1457	0.1	6826	2	23437	(n/a)	65318	(n/a)	156865	(n/a)	337042	(n/a)	664301	(n/a)	1298	(10, 6)	2388	(12, 6)	6893					
(1, 7)	15	0.1	382	0.1	4373	9	27306	(n/a)	117187	(n/a)	229	(6, 7)	766	(7, 7)	2145	(8, 7)	5267											
(1, 8)	17	0.1	766	0.1	13121	26	109226	(n/a)	214	(5, 8)	1145	(6, 8)	4593															
(1, 9)	19	0.1	1534	0.1	39365	(n/a)	436906	(n/a)																				
(1, 10)	21	0.1	3070	0.2	118097	(n/a)	1747626	(n/a)	3414																			
(1, 12)	25	0.1	12286	(n/a)	1062881	(n/a)																						
(1, 13)	27	0.1	24574	48																								
(1, 14)	29	0.1	49150	15																								
(1, 16)	33	0.1	196606	384																								
(1, 18)	37	0.1	786430	1536																								
(1, 20)	41	0.1																										
(1, 22)	45	0.2																										
(1, 24)	49	0.1																										
(1, 26)	53	0.1																										
(1, 28)	57	0.1																										
(1, 30)	61	0.1																										
(1, 32)	65	0.1																										



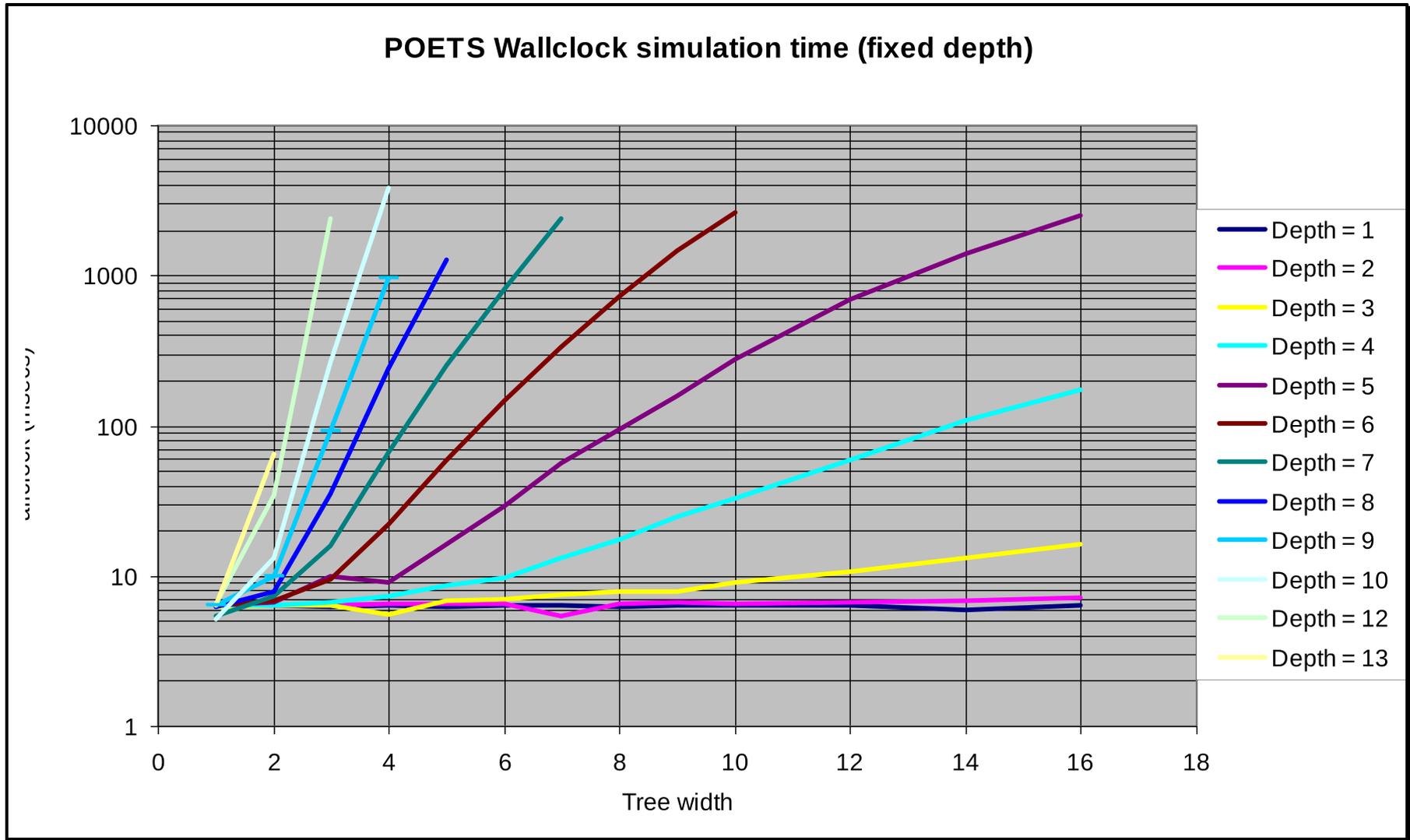
(width,depth)	dev/core
$\Sigma$ devices	min timer

- data lost
- MC table overflow
- >1000 dev/core
- no biological fidelity

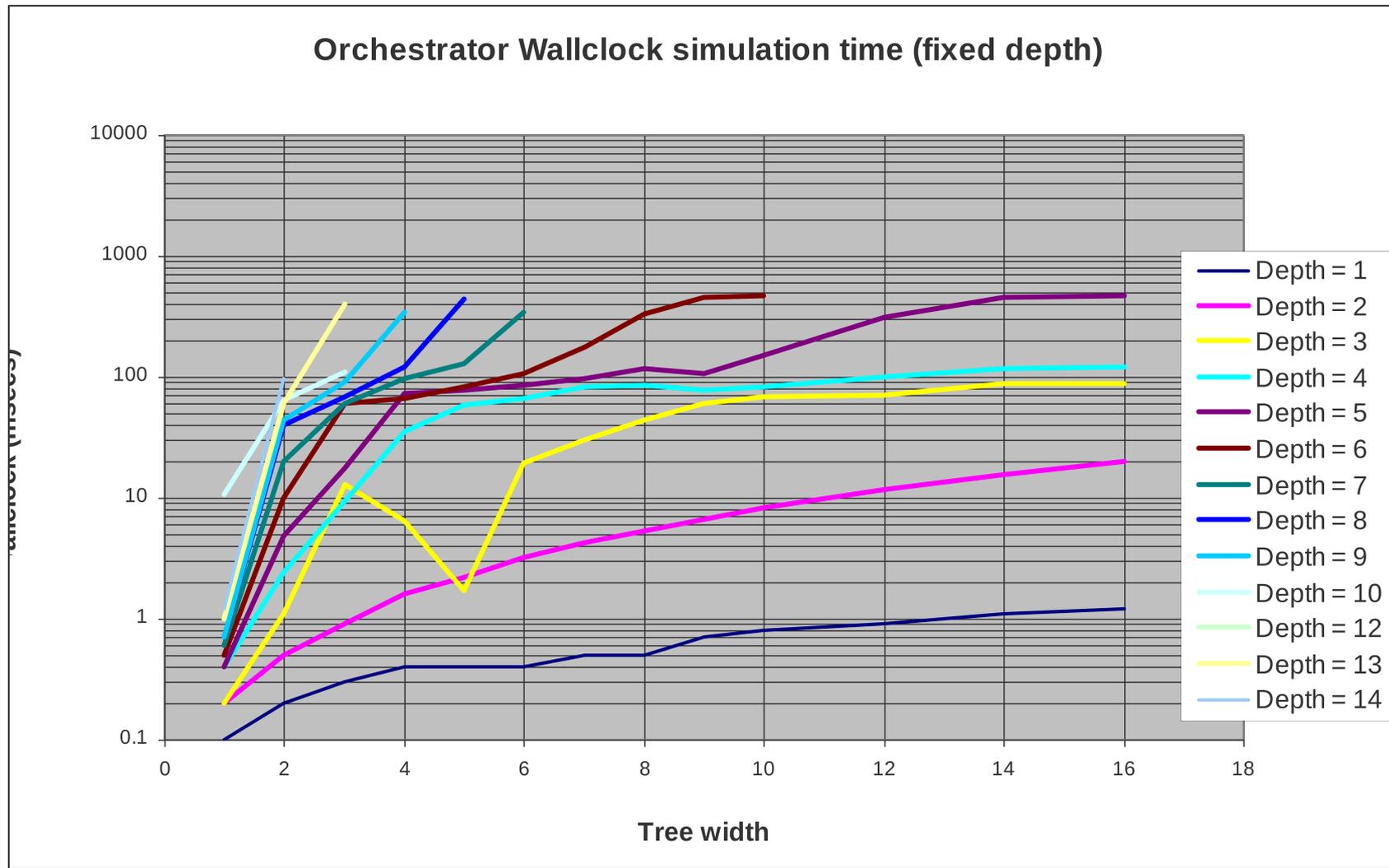
# Clock Trees on a Conventional PC



# Clock Trees on POETS

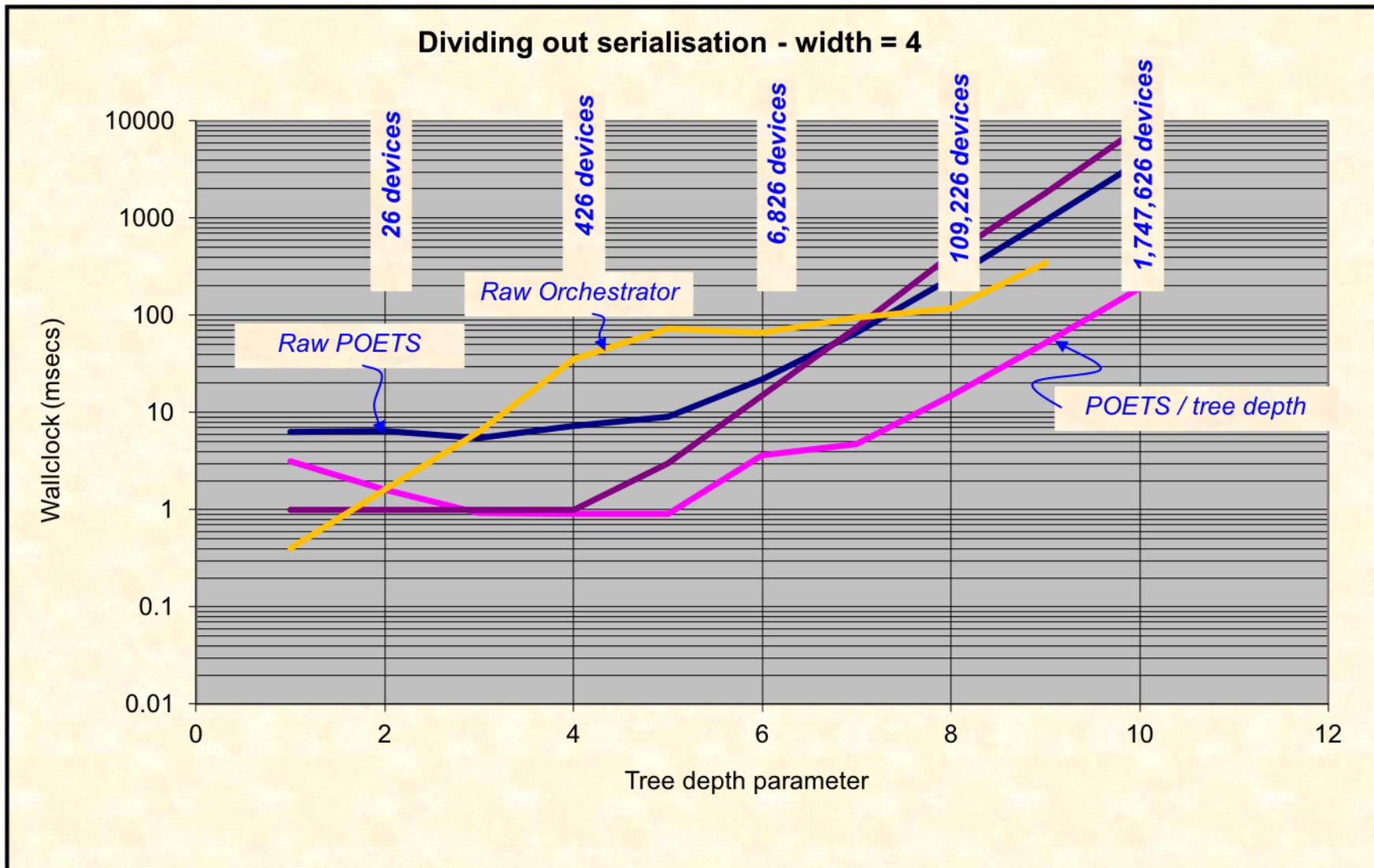


# Clock Trees on POETS (Again!)



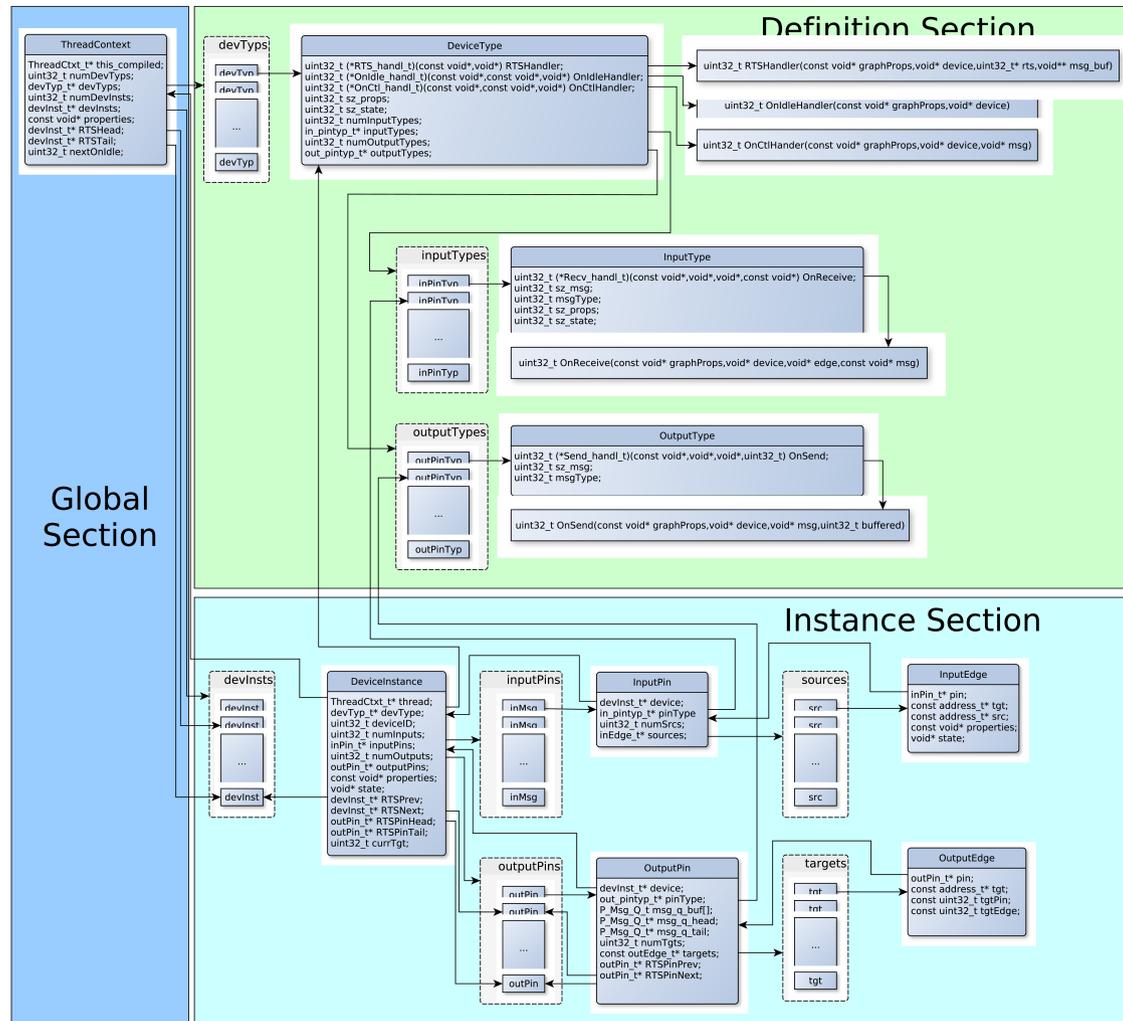
# So What's Easy (and Works)?

- Application Specification
- Scaling (up to a point!)
- Code Generation
- Basic Mapping



# And What's Hard?

- Debugging
- Recasting the Problem
- Time to Generate and Load
- Data Type Support
- Nontrivial Mapping
- Sensible I/O



# What's Next?

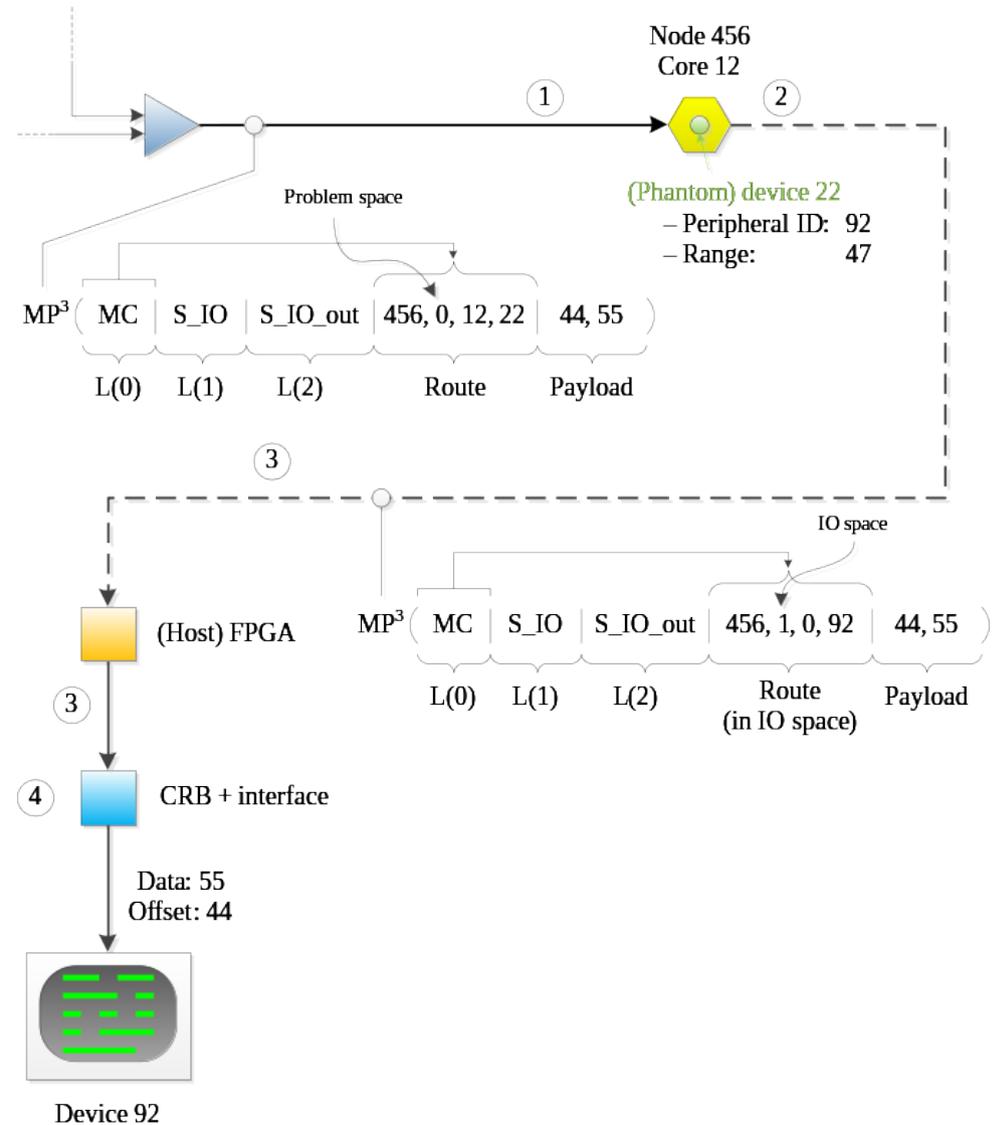
Mapping Algorithms  
(Simulated Annealing?)

Standards for Asynchronous  
Message I/O

Work on Formal Verification?

Diversifying the Applications

Better Hardware



# And What's The Future?

NVIDIA Isn't Going Away

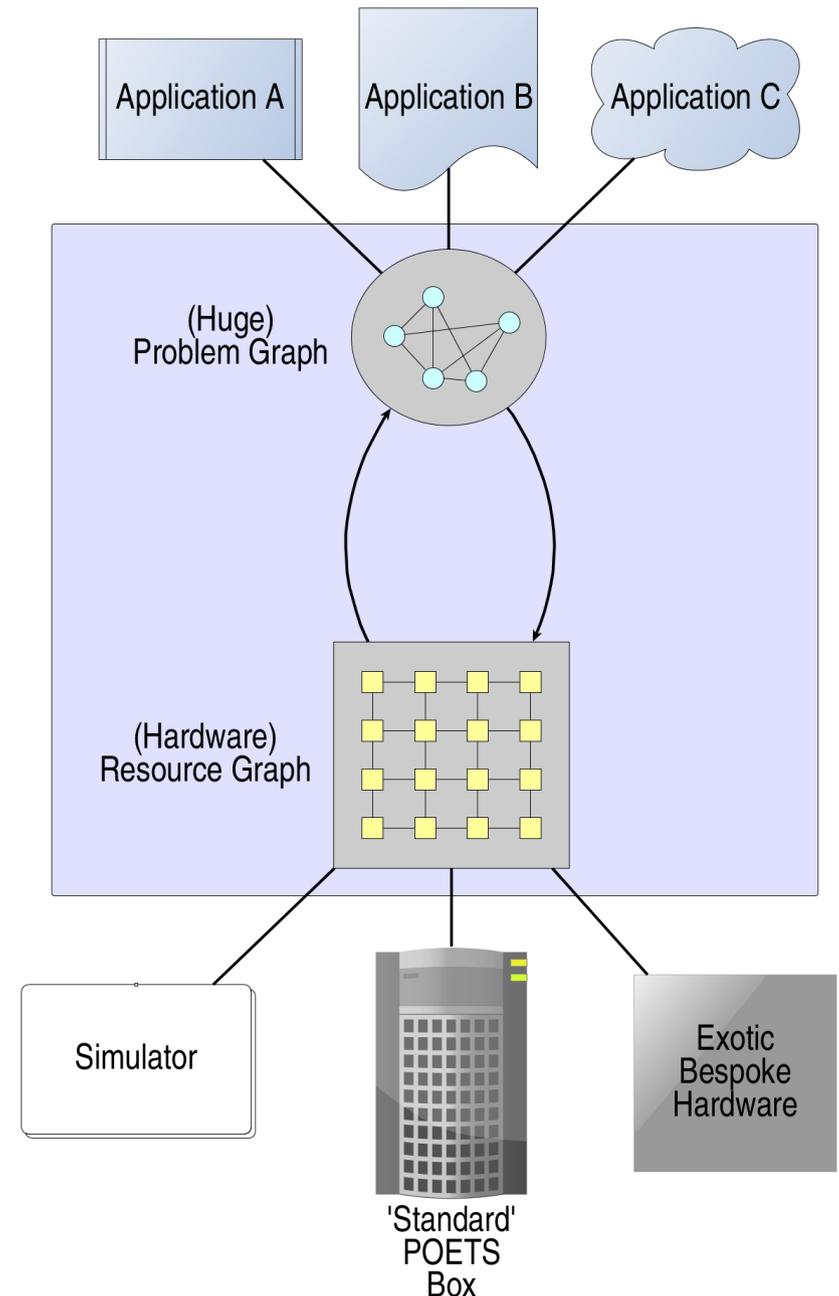
But Some Other Models  
Will Find Niches

Graph Descriptions and  
Graphical Programming Software

Automation of the  
Mapping Process

Asynchronous = True Parallel?

Embedded Parallel Systems



# Many Thanks To...

