



University
of
St Andrews



Predictable Timing Analysis of x86 Multicores using High-Level Parallel Patterns



Kevin Hammond, Susmit Sarkar and Chris Brown

University of St Andrews, UK

T: @paraphrase_fp7

E: kh@cs.st-andrews.ac.uk

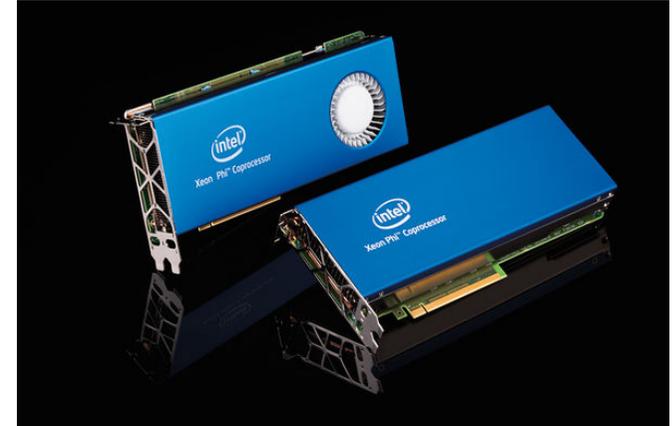
W: <http://www.paraphrase-ict.eu>





Motivation

- **No future system will be single-core**
 - parallel programming will be essential
- **It's not just about performance**
 - **it's also about energy usage**
- **If we don't solve the multicore challenge, then no other advances will matter!**
- **We need to produce predictable timing models for widely used multicores (e.g. x86, ARM)**





University
of
St Andrews

Even Mobile Phones are Multicore!



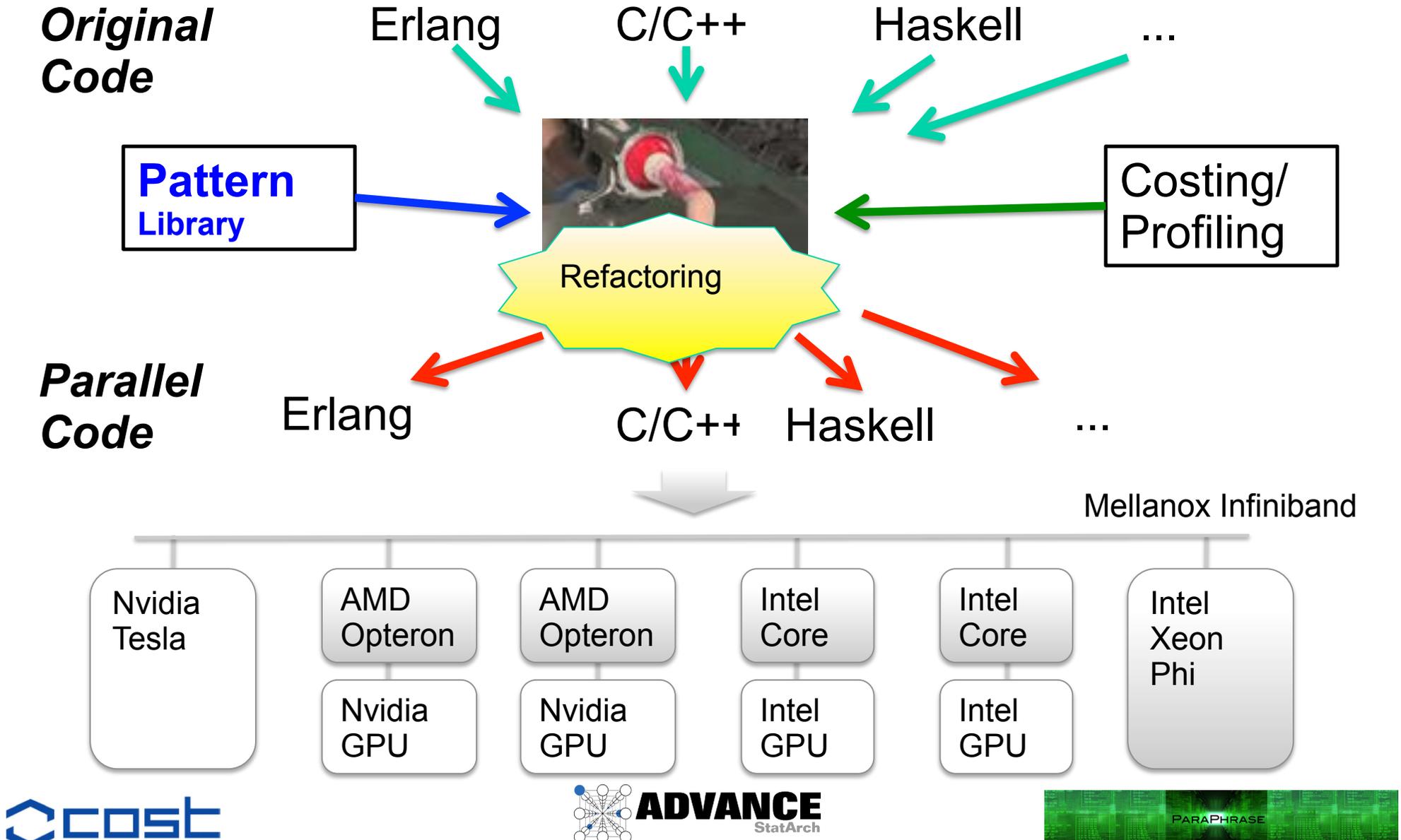
Current Parallel Methodologies

- **Applications programmers *must be systems programmers***
 - insufficient assistance with abstraction
 - too much complexity to manage
- **Difficult/impossible to scale, unless the problem is simple**
- **Difficult/impossible to change fundamentals**
 - scheduling
 - task structure
 - migration
- **Many approaches provide *libraries***
 - they need to provide abstractions

Thinking Parallel

- **Fundamentally, programmers must learn to “think parallel”**
 - this requires new *high-level* programming constructs
 - perhaps dealing with large numbers of threads
- **You cannot program effectively while worrying about deadlocks etc.**
 - *they must be eliminated from the design!*
- **You cannot program effectively while fiddling with communication etc.**
 - *this needs to be packaged/abstracted!*
- **You cannot program effectively without performance information**
 - *this needs to be included as part of the design!*

The ParaPhrase Approach



Components and Abstraction

- **Components give some of the advantages of functional programming**
 - clean abstraction
 - pure computations, easily scheduled
 - dependencies can be exposed
- **Hygiene/discipline is necessary**
 - no unwanted state leakage
(e.g. in terms of implicit shared memory state)

The ParaPhrase Approach

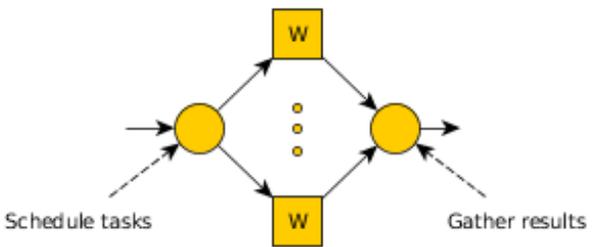
- **Start bottom-up**
 - identify (strongly hygienic) **COMPONENTS**
 - *using semi-automated refactoring*
- **Think about the PATTERN of parallelism**
 - e.g. map(reduce), task farm, parallel search, parallel completion, ...
- **STRUCTURE the components into a parallel program**
 - *turn the patterns into concrete (skeleton) code*
 - Take performance, **energy** etc. into account (multi-objective optimisation)
 - also using refactoring
- **RESTRUCTURE/TUNE if necessary! (also using refactoring)**

*both legacy and
new programs*

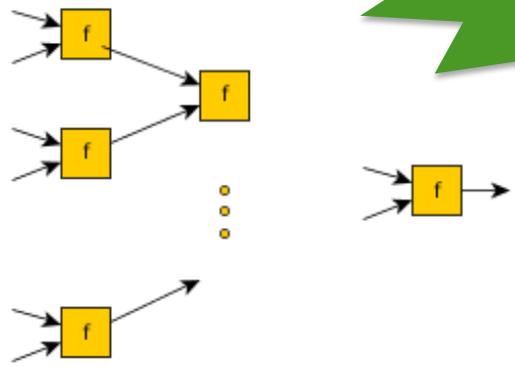


Some Common Parallel Patterns

Farm



Reduce

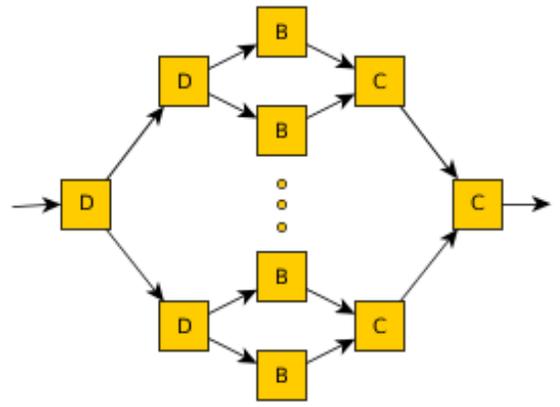


Generally, we need to nest/combine patterns in arbitrary ways

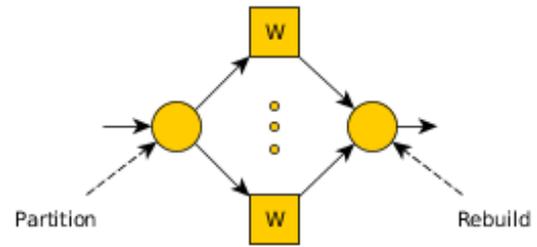
Pipeline



Divide&Conquer



Map



Skeletons

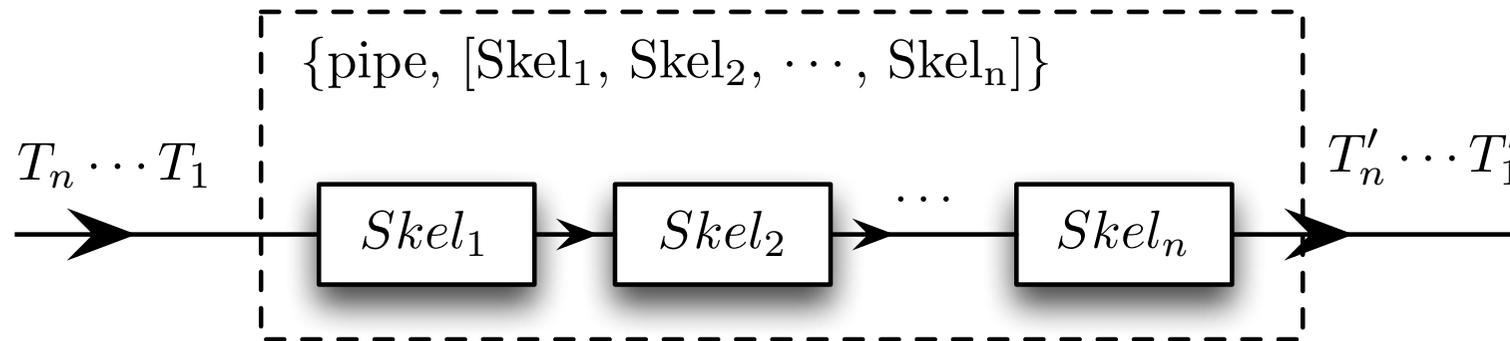
- Skeletons are *implementations* of parallel patterns
- A skeleton is a template
 - pluggable *higher-order* functions
 - can be instantiated with concrete worker functions
- Skeletons *avoid* deadlock, race conditions
 - communication is implicit and structured

Murray Cole, "Algorithmic Skeletons: structured management of parallel computation" MIT Press, 1989

Horacio González-Vélez and Mario Leyton:
A survey of algorithmic skeleton frameworks: high-level structured parallel programming enablers.
Softw., Pract. Exper. 40(12): 1135-1160 (2010)

Parallel Pipeline Skeleton

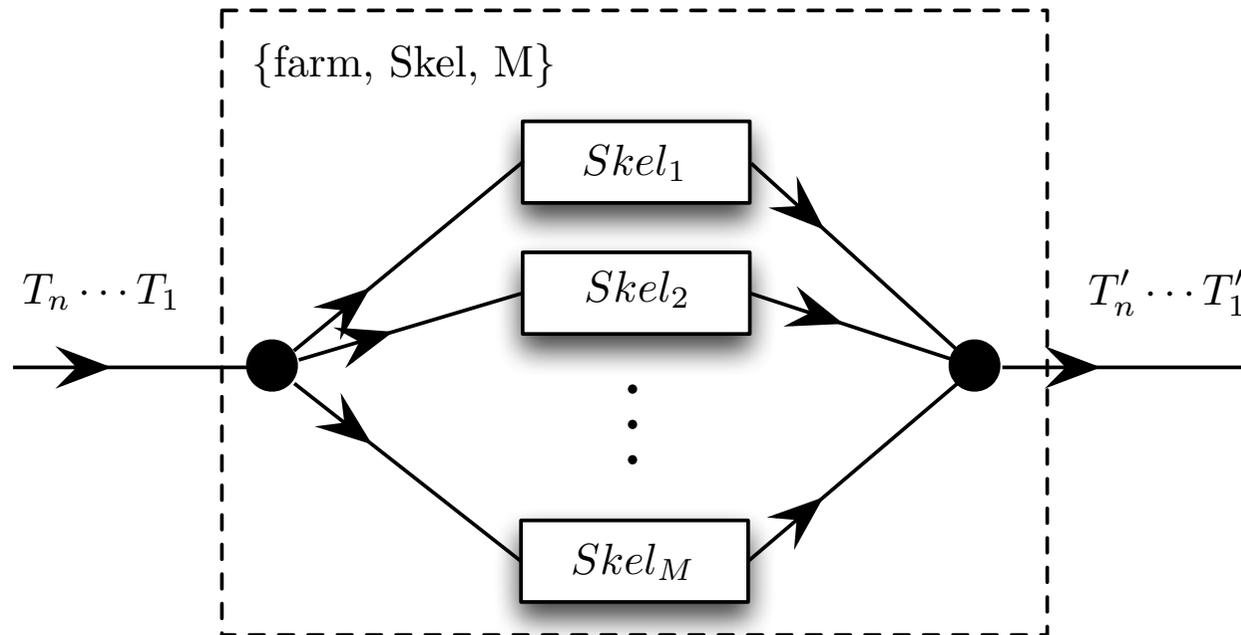
- Each stage of the pipeline can be executed in parallel
- The input and output are streams
- Each stage is itself an instance of a pattern (Skel)



```
skel:do([{pipe, [Skel1, Skel2, ..., SkelN]}], Inputs).
```

Parallel Task Farm Skeleton

- Each worker is executed in parallel
- A bit like a 1-stage pipeline



```
skel:do([ {farm, Skel, M} ], Inputs).
```

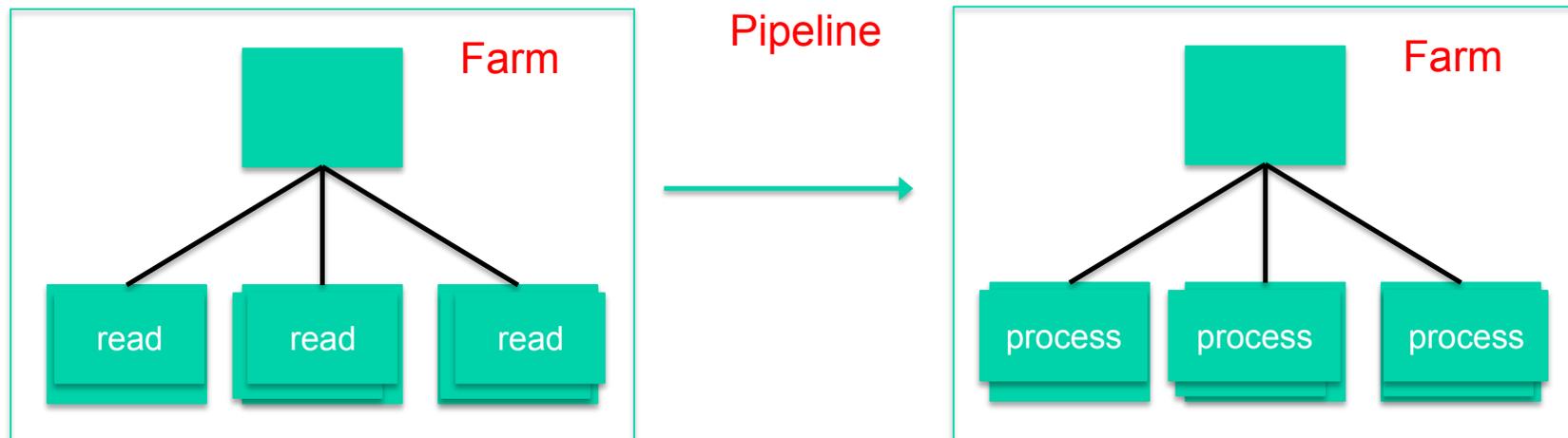
Example Parallel Structure

Sequential

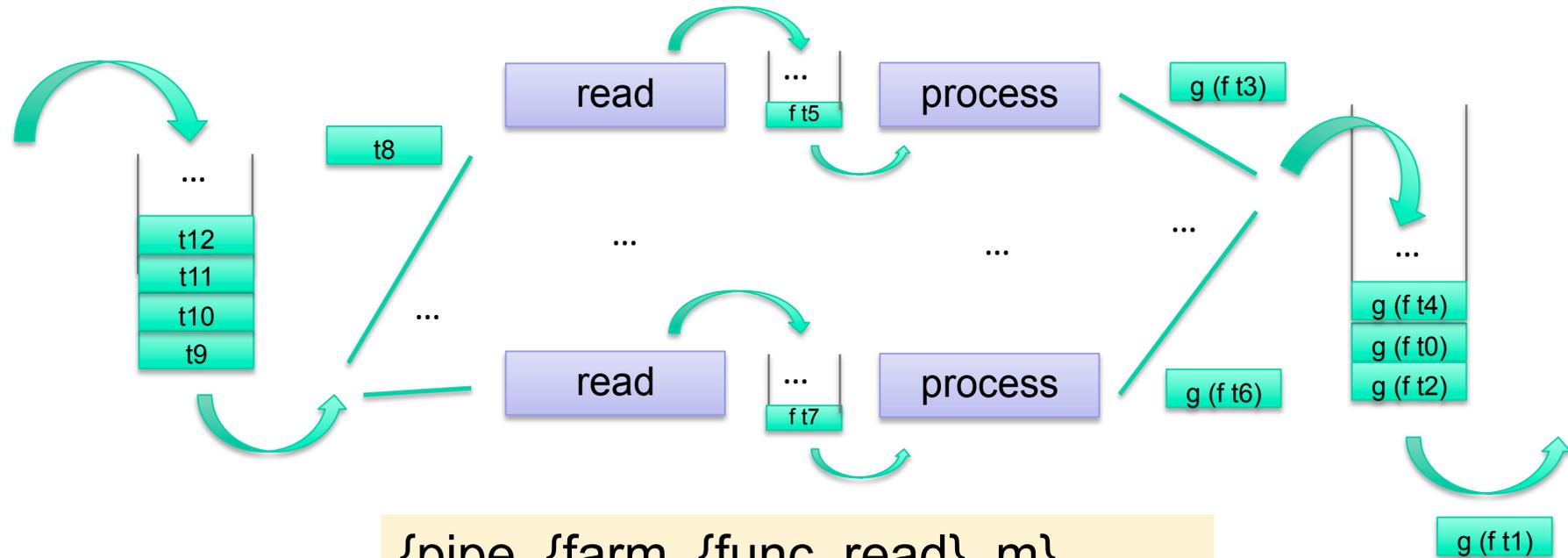
for each image, i .
process(read i)

Parallel

{pipe, {farm, {func, read}, m },
{farm, {func, process}, n }}



Composing Skeletons

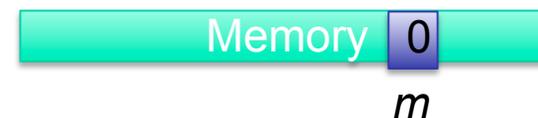
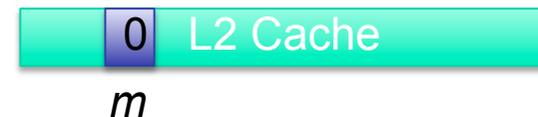
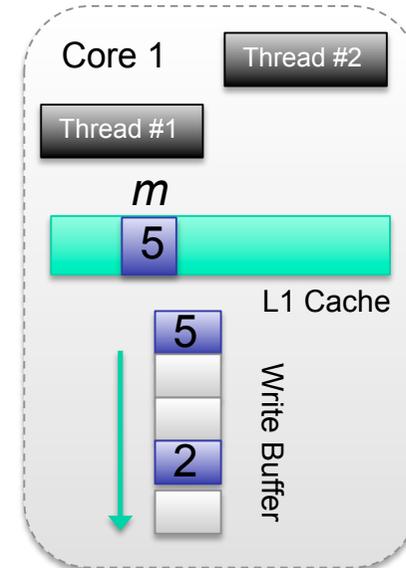


$\{\text{pipe}, \{\text{farm}, \{\text{func}, \text{read}\}, m\},$
 $\{\text{farm}, \{\text{func}, \text{process}\}, n\}\}$

- Queues link skeletons

x86 Multicore Cache Design

- **Each core has**
 - a local write-back cache
 - a FIFO-ordered *write buffer*
- **A core may run many threads**
- **Cores share**
 - level 2 (and 3) cache
 - global memory



Sequential Consistency (SC)

```
1 int x = 0, y = 0;
2
3 ...
4
5 { // thread 1           { // thread 2
6   x = 1;                y = 2;
7   return(x+y);          return(x+y);
8 }                       }
```

- **Memory accesses are effectively interleaved**
 - as if run by a single processor
- **Either**
 - both threads return 3
 - thread 1 returns 1, thread 2 returns 3
 - thread 1 returns 3, thread 2 returns 2
- **Not**
 - thread 1 returns 0, thread 2 returns 0

x86 Total Store Order (TSO)

- **On a multicore, SC can be inefficient**
- **Intel uses a weaker (relaxed memory) consistency model**
 - Total store order (TSO) guarantees that the order that **Writes** are seen by a location is the same as the order they were issued
- **ARM uses an even weaker consistency model**

Basic TSO Rules

- **The basic rules are:**
 - (1) **Reads** are not reordered with other **Reads**.
 - (2) **Writes** are not reordered with older **Reads**.
 - (3) **Writes** are not reordered with other **Writes**.
 - (4) **Reads** may be reordered with older **Writes** to different memory locations but not with older **Writes** to the same memory location
- An **Exchange** is treated as an indivisible Read/Write pair to a specific memory location
- A **Fence** is treated as both a Read and Write to all memory locations, except that no actual memory transfer occurs



Simple Spin Lock Implementation

```
1 void lock( volatile char *lockcell ) {  
2     char old_value ;  
3  
4     do {  
5         old_value = exchange(lockcell,1);  
6     } while ( 1 == old_value ) ;  
7 }  
8  
9  
10  
11 void unlock( volatile char *lockcell ) {  
12     *lockcell = 0 ;  
13 }
```



x86 Assembly code for spin lock

```
1 lockr:  
2   push ebp           ; Start new stack frame  
3   mov  ebp, esp  
4   mov  ebx, [ebp+8]  ; Get address of lock cell  
5  
6   trylock:  
7   mov  eax, 1        ; Set EAX register to 1 (locked)  
8   xchg eax, [ebx]    ; Exchange EAX and lock cell  
9   test eax, eax      ; Test whether the cell is already locked  
10  jnz  trylock       ; Retry the lock if so  
11  
12  pop  ebp           ; revert stack frame  
13  ret                ; The lock has been acquired  
14  
15  unlockr:  
16  push ebp           ; Start new stack frame  
17  mov  ebp, esp  
18  mov  ebx, [ebp+8]  ; Get address of lock cell  
19  
20  mov  eax, 0        ; Set EAX register to 0 (unlocked)  
21  mov  [ebx], eax    ; Release the lock  
22  
23  pop  ebp           ; revert stack frame  
24  ret                ; The lock has been released.
```



Simple Queue using spin lock

```
1 Value qget(Queue q) {  
2   Value v;  
3  
4   do {  
5     lock(&q.lock_cell);  
6  
7     if (qempty(q))  
8       break;  
9  
10    unlock(&q.lock_cell);  
11  } while (1);  
12  
13  /* lock is held */  
14  v = front(q);  
15  
16  unlock(&q.lock_cell);  
17  
18  return(v);  
19 }
```

```
21 void qput(Queue q, Value v) {  
22   lock(&q.lock_cell);  
23  
24   addtoq(q, v);  
25  
26   unlock(&q.lock_cell);  
27 }
```

We have used HOL to
prove that this is sound wrt
the TSO relaxed-memory
model



Simple Timing Model

- The worst-case costs if n threads contend a lock are

$$T_{\text{qput}} = n \cdot T_{\text{Exchange}} + T_{\text{Write}} + T_{\text{Write}}$$

$$T_{\text{qget}} = n \cdot T_{\text{Exchange}} + T_{\text{Read}} + 2T_{\text{Write}}$$

Timing Model for a Farm

- The amortised average cost for each farm operation is

$$T_{qget} + T_f + T_{qput}$$

which simplifies to:

$$2 \cdot (n + 1) \cdot T_{Exchange} + 5 \cdot T_{Write} + T_{Read} + T_f$$

Timing Model for a Pipeline

- If the first stage dominates (function f), its cost is

$$T_{qget} + T_f + T_{qput}$$

which simplifies to:

$$2 \cdot (|f| + 1) \cdot T_{Exchange} + 5 \cdot T_{Write} + T_f$$

- The total cost for both stages is therefore:

$$2 \cdot (|f| + 1) \cdot T_{Exchange} + 5 \cdot T_{Write} + T_f + T_g$$

or, if the second stage dominates (function g)

$$T_f + (2 \cdot |g| + |f| + 1) \cdot T_{Exchange} + 5 \cdot T_{Write} + T_g$$

Including Store-Buffer Flushing

- The cost of an **exchange** depends on items to be flushed, **b**

$$T_{Exchange} = b \cdot T_{Fl} + T_{JustX}$$

- The cost of a spin-lock on **t** contending threads is

$$b \cdot T_{Fl} + t \cdot T_{JustX}$$

- The costs of queue operations change slightly

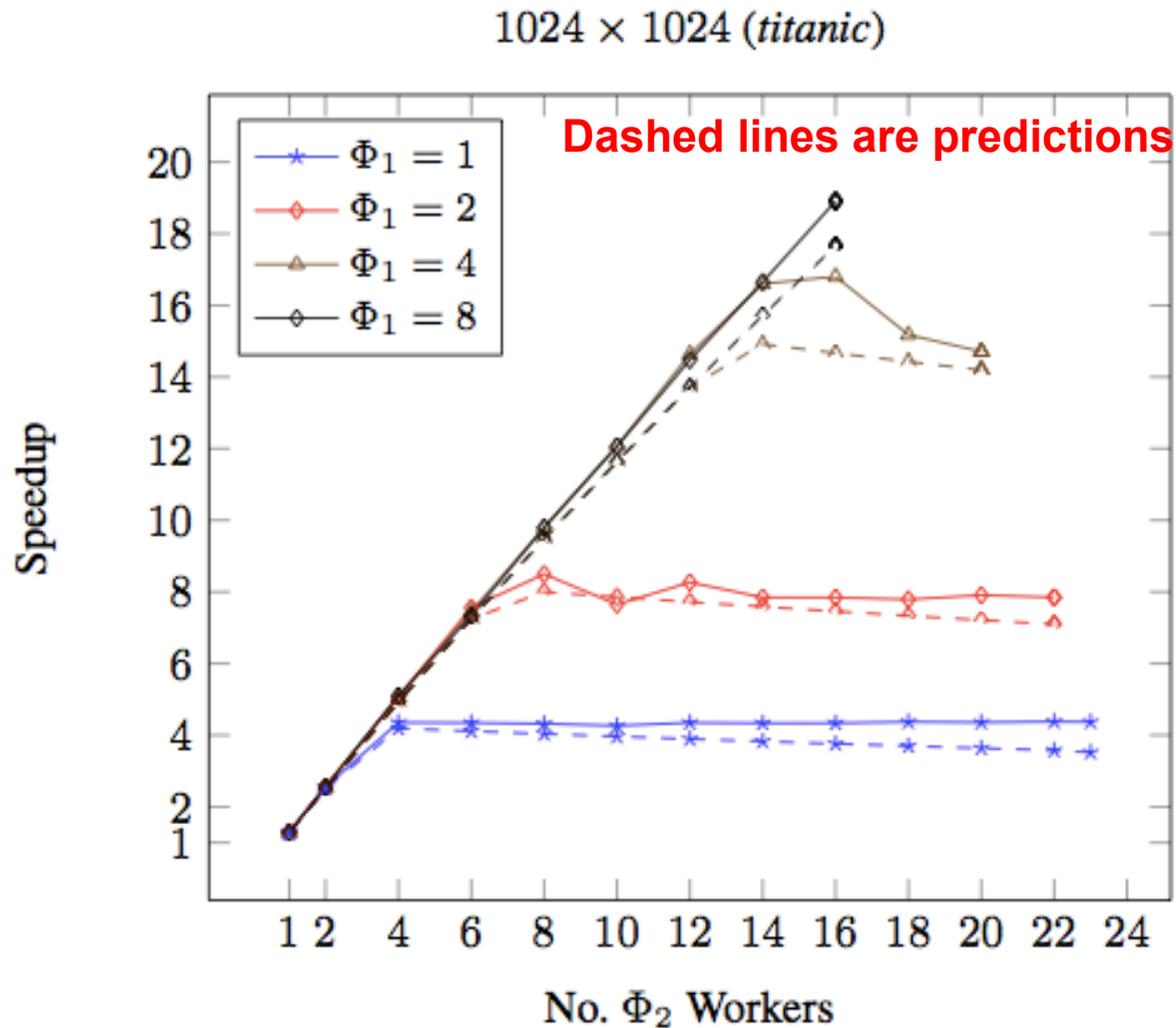
$$T_{qput} = b \cdot n \cdot T_{Fl} + n \cdot T_{JustX} + T_{Write} + T_{Write};$$

$$T_{qget} = b \cdot n \cdot T_{Fl} + n \cdot T_{JustX} + T_{Read} + 2 \cdot T_{Write};$$

- The cost of a farm is:

$$2 \cdot (n + 1) \cdot b \cdot T_{Fl} + (n + 1) \cdot T_{JustX} + 4T_{Write} + T_{Read} + T_f$$

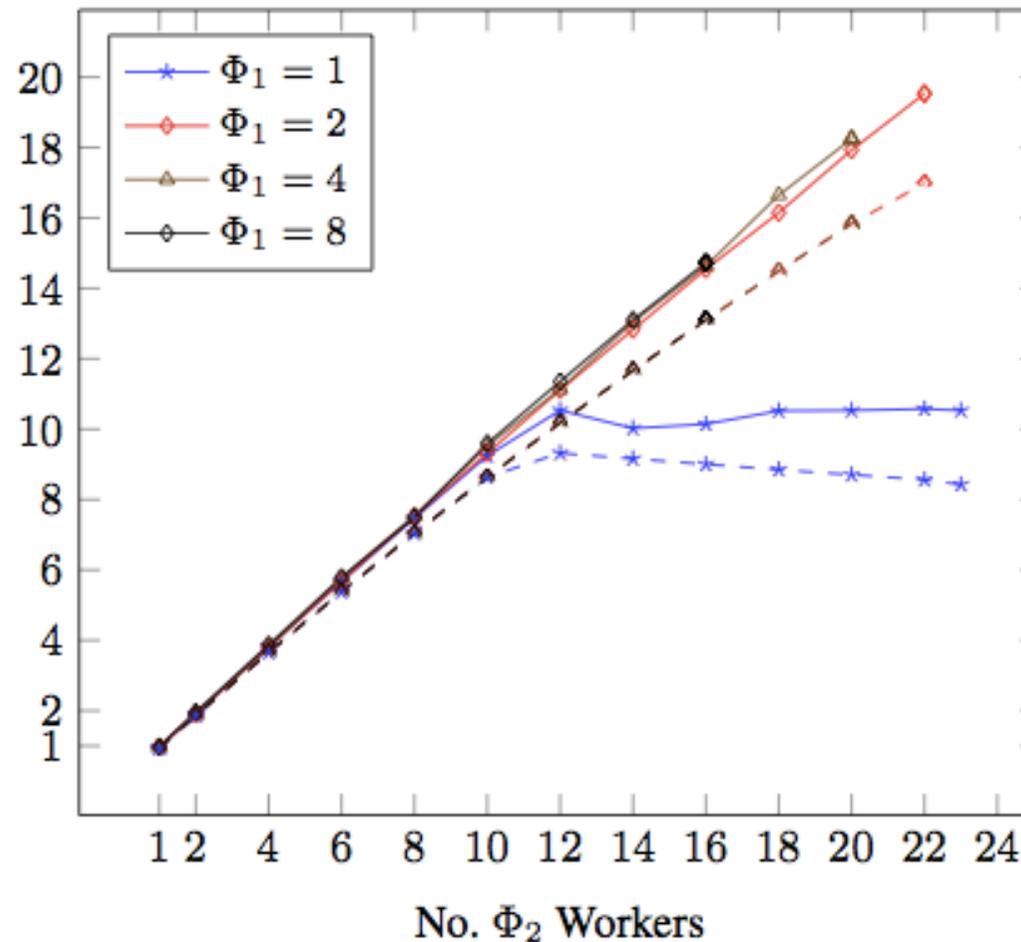
Performance Predictions (Image Convolution, 1024x1024)



24 core machine at Uni. Pisa
2xAMD Opteron 6176. 800 Mhz
32GB RAM
1 x NVidia Tesla C2050 GPU

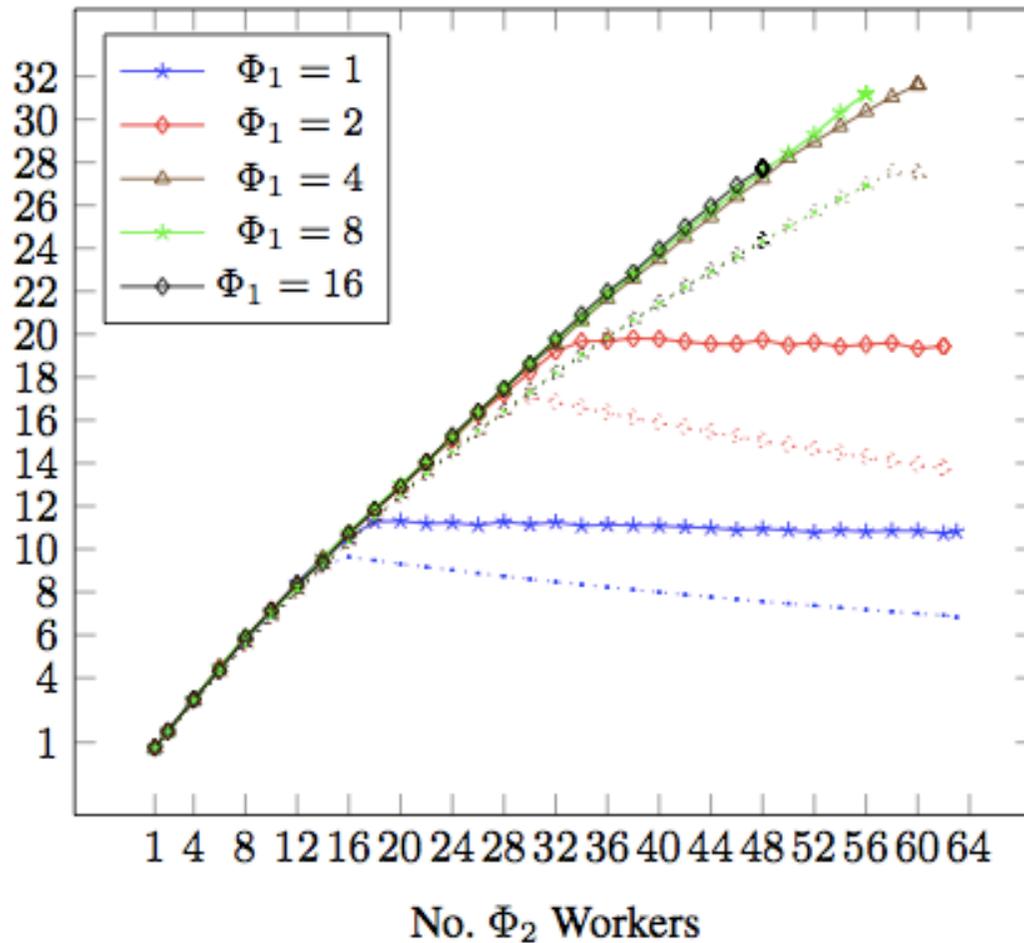
Performance Predictions (Image Convolution, 2048x2048)

2048 × 2048 (*titanic*)



Performance Predictions (Image Convolution, 2048x2048)

2048 × 2048 (*lovelace*)



64-core machine at Uni. St Andrews

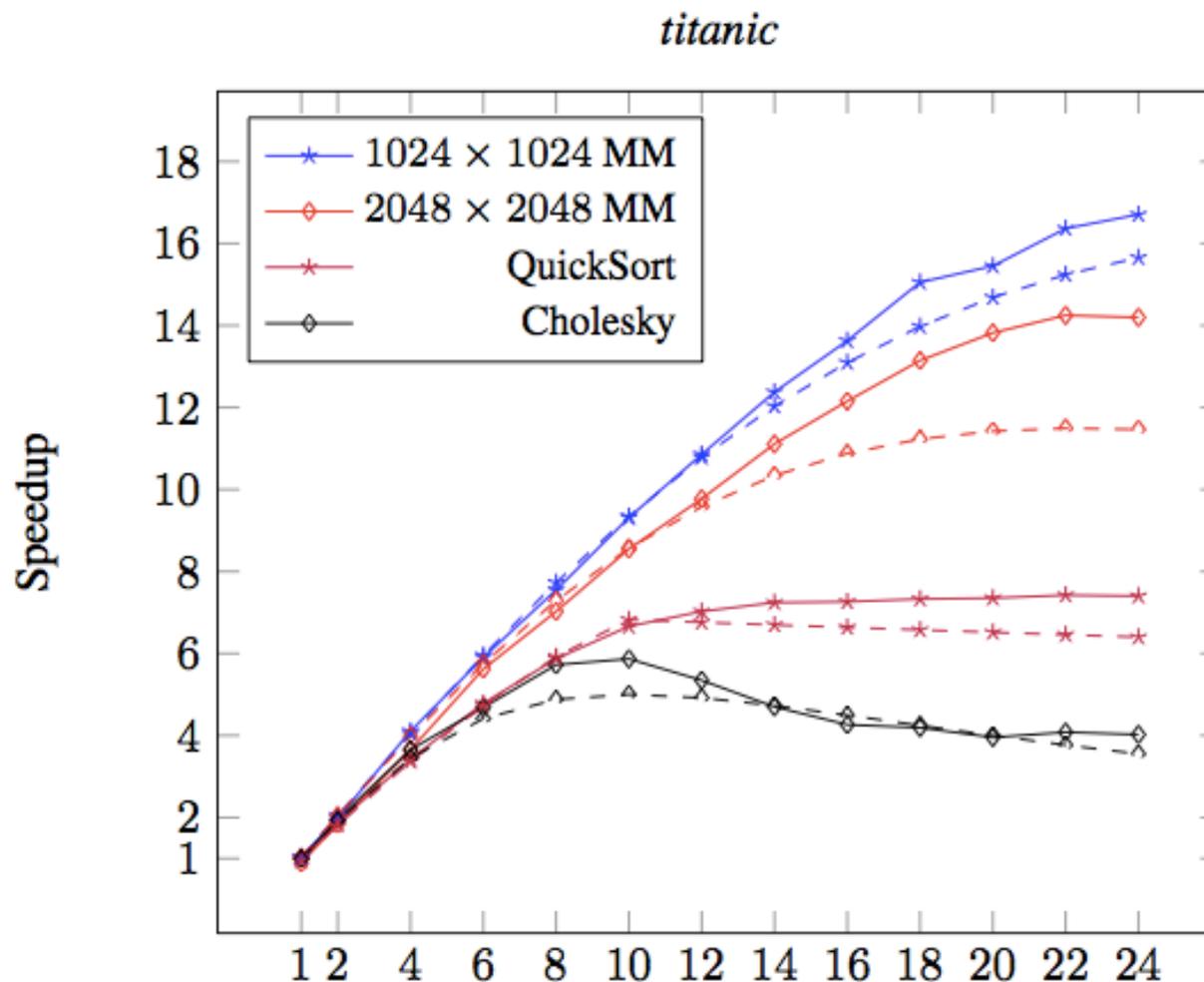
8xAMD Opteron 6376. 2.3Ghz

32GB RAM

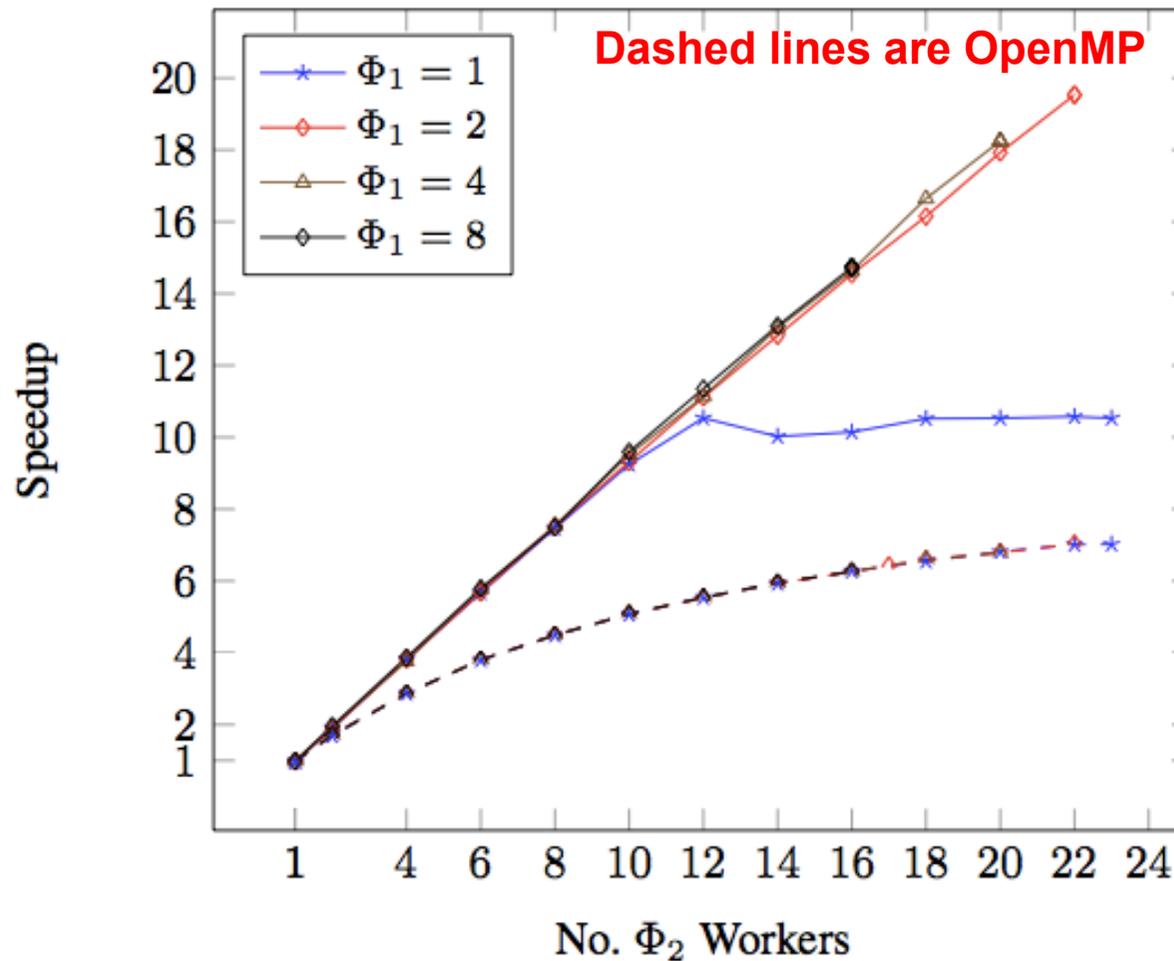
Performance Predictions (Matrix Multiplication etc)



University
of
St Andrews



Comparison with OpenMP



Combining CPUs and GPUs

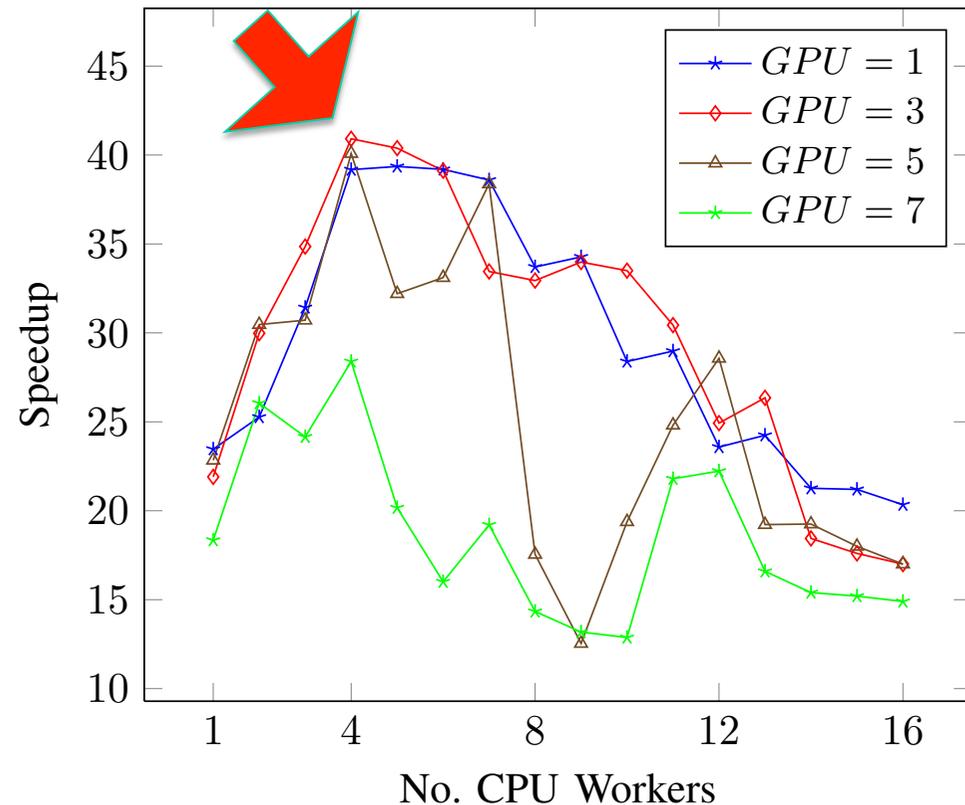
- **Machine Learning chooses**

- best combination of patterns
- CPU/GPU allocations

- **Excellent Results**

- within 5% of optimal
- > 40x speedup
over sequential CPU

Speedups for $\Delta(r) \parallel \Delta(p)$



Lock-Free Queue Implementation

```
--
16 Value qgetlf(Queue q) {
17     NodePtr first;
18
19     do {
20         if (qempty(q))
21             continue;
22
23         first = q->first;
24
25         if (first==NULL)
26             continue;
27
28     } while (!dcas(q->first, first, first->next));
29
30     return (first->value);
31 }
```

- **This uses a double compare-and-swap variant of Exchange**
 - atomically swaps two values
 - allows us to avoid ABA errors by including a count field



Lock-Free Queue Implementation

```
34 void qputlf(Queue q, Value v) {
35     Node n = new Node(v);           // new queue node
36     NodePtr np = new NodePtr(n);    // queue pointer
37     Queue last;
38
39     do {
40         last = q->last;
41
42         np.count = last.count+1;
43     } while(!dcas(q->last, last, np));
45 }
```

- **At the pattern level, this is plug-replaceable with a lock**
 - The cost model needs to change but most details are the same
 - All proof is the same above the lock-free level



Comparison of Development Times

	Man.Time	Refac. Time
Convolution	3 days	3 hours
Ant Colony	1 day	1 hour
BasicN2	5 days	5 hours
Graphical Lasso	15 hours	2 hours

Conclusions

- **High-level Patterns help structure parallel computations**
 - avoid deadlock, race conditions etc (formal proof in paper!)
 - reduce development time by an order of magnitude
 - allow us to construct predictable cost models
- **Cost model for x86 constructed from first principles**
 - Predictable timings for x86 (provably correct from TSO semantics)
 - Highly Accurate
 - All previous formal models have been for much simpler memory models (e.g. PPC)
- **Proved to be deadlock-free**
- **Applicable to energy as well as time**

Funded by



University
of
St Andrews

- **ParaPhrase (EU FP7), Patterns for heterogeneous multicore,**

€4.2M, 2011-2014



- **SCIENCE (EU FP6), Grid/Cloud/Multicore coordination**

•€3.2M, 2005-2012



- **Advance (EU FP7), Multicore streaming**

•€2.7M, 2010-2013



- **HPC-GAP (EPSRC), Legacy system on thousands of cores**

•£1.6M, 2010-2014

- **Islay (EPSRC), Real-time FPGA streaming implementation**

•£1.4M, 2008-2011



- **TACLE: European Cost Action on**

•€300K, 2012-2015





University of St Andrews

Some of our Industrial Connections

Mellanox Inc.



Erlang Solutions Ltd



SAP GmbH, Karlsruhe



BAe Systems



Selex Galileo



Biold GmbH, Stuttgart

Philips Healthcare

Software Competence Centre, Hagenberg



Microsoft Research



Well-Typed LLC

Microsoft Research





- ▶ Home
- ▶ Sitemap
- ▶ Restricted Area
- ▶ FAQ
- ▶ Glossary
- ▶ Links
- ▶ Contact
- ▶ Jobs

- About COST
- Domains and Actions**
- Participate
- Events
- Media
- e-COST**



Home | Domains and Actions | Information and Communication Technologies (ICT) | Actions | IC1202

- ▶ All Actions
- ▶ Biomedicine and Molecular Biosciences (BMBS)
- ▶ Chemistry and Molecular Sciences and Technologies (CMST)
- ▶ Earth System Science and Environmental Management (ESSEM)
- ▶ Food and Agriculture (FA)
- ▶ Forests, their Products and Services (FPS)
- ▶ Individuals, Societies, Cultures and

ICT COST Action IC1202

Timing Analysis on Code-Level (TACLe)

Descriptions are provided by the Actions directly via e-COST.

Embedded systems increasingly permeate our daily lives. Many of those systems are business- or safety-critical, with strict timing requirements. Code-level timing analysis (used to analyse software running on some given hardware w.r.t. its timing properties) is an indispensable technique for ascertaining whether or not these requirements are met. However, recent developments in hardware, especially multi-core processors, and in software organisation render

Information and Communication Technologies COST Action IC1202

- ▶ Description
- ▶ Parties
- ▶ Management Committee



General Information*





University
of
St Andrews

THANK YOU!

<http://www.paraphrase-ict.eu>

<http://www.project-advance.eu>

@paraphrase_fp7