Software engineering challenges in many-core computing

Paul Kelly Group Leader, Software Performance Optimisation Department of Computing Imperial College London

Thalesians Workshop on GPUs in Finance

J Presper Eckert (1919-1995)

Co-inventor of, and chief engineer on, the ENIAC, arguably the first storedprogram computer (first operational Feb 14th 1946)

27 tonnes, 150KW, 5000 cycles/sec



J.G. Brainerd & T.K. Sharpless. "The ENIAC." pp 163-172 Electrical Engineering, Feb 1948.



ENIAC was designed to be set up manually by plugging arithmetic units together (reconfigurable logic)

You could plug together quite complex configurations **Parallel** - with multiple units working at the same time



Gloria Gorden and Ester Gerston: programmers on ENIAC

THEORY AND TECHNIQUES FOR DESIGN OF ELECTRONIC DIGITAL COMPUTERS

> Lectures given at the Moore School 8 July 1946-31 August 1946

Volume IV Lectures 34-48



UNIVERSITY OF PENNSYLVANIA Moore School of Electrical Engineering PHILADELPHIA, PENNSYLVANIA

June 30, 1948

The Moore School Lectures

The first ever computer architecture conference

July 8th to August 31st 1946, at the Moore School of Electrical Engineering, University of Pennsylvania

A defining moment in the history of computing

To have been there....

A PARALIEL CHANNEL COMPUTING MACHINE

Lecture by J. P. Eckert, Jr. Electronic Control Company

... Again I wish to reiterate the point that all the arguments for parallel operation are only valid provided one applies them to the steps which the built in or wired in programming of the machine operates. Any steps which are programmed by the operator, who sets up the machine, should be set up only in a serial fashion. It has been shown over and over again that any departure from this procedure results in a system which is much too complicated to use.



But the free linch is over

Philip E Ross, Why CPU Frequency Stalled - http://www.spectrum.ieee.org/apr08/6106/CPU

- But "It has been shown over and over again..." that this results in a system too complicated to use
- How can we get the speed and efficiency without suffering the complexity?
- What have we learned since 1946?

But "It has been shown over and over again..." that this results in a system too complicated to use

HOMMONTOFATTOR

- How can we get the speed and efficiency without suffering the complexity?
- What have we learned since 1946?
 - We really need parallelism

- But "It has been shown over and over again..." that this results in a system too complicated to use
- How can we get the speed and efficiency without suffering the complexity?
- What have we learned since 1946?
 - Compilers and out-of-order processors can extract some instruction-level parallelism
 - Explicit parallel programming in MPI, OpenMP, VHDL are flourishing industries they can be made to work
 - SQL, TBB, Cilk, Ct (all functional...), many more speculative proposals
 - No attractive general-purpose solution

- But "It has been shown over and over again..." that this results in a system too complicated to use
- How can we get the speed and efficiency without suffering the complexity?
- What have we learned since 1946?
 - Some discipline for controlling complexity
 - Program generation....
 - Programs that generate programs
 - That are correct by construction
 - The generator encapsulates parallel programming expertise

Example:

for (i=0; i<N; ++i) {
 points[i]->x += 1;

Can the iterations of this loop be executed in parallel?

ASTRALGIST



No problem: each iteration is independent



ASUTAR GIST

Oh no: not all the iterations are independent!

You want to re-use piece of code in different contexts

Whether it's parallel depends on context!

Example:

for (i=0; i<N; ++i) {
 points[i]->x += 1;

Can the iterations of this loop be executed in parallel?

ASUTARIA



"Balloon types" or "ownership types" ensure that each cell is reached only by it's owner pointer

```
Variable s of
                                                          function g might
int *f(int *p) {
                                                          point to variable
                                  (1) f_* \supseteq f_p
  return p;
                                                          p of function g
int q() {
                                                          R might point to
  int x,y,*p,*q,**r,**s;
                                                         anything s might
                                  (2) g_s \supseteq \{g_p\}
  s=&p;
                                                               point to
                                  (3) g_p \supseteq \{g_x\}
  if(...) p=&x;
                                                          f's p might point
                                  (4) \ g_p \supseteq \{g_y\}
  else p=&y;
                                                            to anything r
                                                           might point to
                                  (5) g_r \supseteq g_s
  r=s;
                                  (6) f_p \supseteq *g_r
                                                          q might point to
  q=f(*r);
                                                             anything f
                                  (7) g_q \supseteq f_*
                                                               returns
```

TAS A MALIRI

Goal: for each pointer variable (p,q,r,s), find the set of objects it might point to at runtime

We have quite a large constraint graph

Eg for 126.gcc from SPEC95:

194K lines of code (132K excl comments)

51K constraint variables (22K of them heap)

7.4K "trivial" constraints

39K "simple" constraints

25K "complex" constraints (due to dereferencing)

Need to bring together several tricky techniques to get sensible solution times

Difference-sets: propagate only changes so you can track what has changed

Topological sort: visit nodes in order that maximises solution propagation

Cycle detection: zero-weighted cycles can be collapsed

Dynamically: dereferencing pointers adds new edges

0.61s for the whole program (900MHz Athlon)

Histogram of pointsto set size at dereference sites for 126.gcc:



Shared memory makes parallel programming much easier:

> for(i=0; I<N; ++i) par_for(j=0; j<M; ++j) A[i,j] = (A[i-1,j] + A[i,j])*0.5; par_for(i=0; I<N; ++i) for(j=0; j<M; ++j) A[i,j] = (A[i,j-1] + A[i,j])*0.5;

- First loop operates on rows in parallel
- Second loop operates on columns in parallel
- With distributed memory we would have to program message passing to transpose the array in between
- With shared memory... no problem!



celf-optimising linear algebra library

Olav Beckmann's PhD thesis:

- Each library function comes with metadata describing data layout constraints
- Solve for distribution of each variable that minimises redistribution cost



Finding parallelism is usually easy

Imperial College

London

- Very few algorithms are inherently sequential
 - But if you want a large speedup you need to parallelise almost all of your program

Parallelism breaks abstractions:

- Whether code should run in parallel depends on context
- How data and computation should be distributed across the machine depends on context
- *"Best-effort", opportunistic parallelisation is almost useless:*
 - Robust software must robustly, predictably, exploit large-scale parallelism

How can we build robustly-efficient multicore software

Easy parallelism - tricky engineering

While maintaining the abstractions that keep code clean, reusable and of long-term value?

- The Foundry is a London company building visual effects plug-ins for the movie/TV industry (http://www.thefoundry.co.uk/)
- Core competence: image processing algorithms
- Core value: large body of C++ code based on library of image-based primitives

Opportunity 1:

 Competitive advantage from exploitation of whatever platform the customer may have - SSE, multicore, vendor libraries, GPUs

Opportunity 2:

Redesign of the Foundry's Image Processing Primitives Library

Risk:

- Premature optimisation delays delivery
- Performance hacking reduces value of core codebase

Imperial College Londo I SURI EFFECTS IN MOVIE DOST-DIOLUCION



Nuke compositing tool (http://www.thefoundry.co.uk)

- Visual effects plugins (Foundry and others) appear as nodes in the node graph
- We aim to optimise individual effects for multicore CPUs, GPUs etc
- In the future: tunnel optimisations across node boundaries at runtime.

(c) Heribert Raab, Softmachine. All rights reserved. Images courtesy of The Foundry

Imperial Colleg Visual effects: degrain example London

Image degraining effect – a complete Foundry plug-in

Random texturing noise introduced by photographic film is removed without compromising the clarity of the picture, either through analysis or by matching against a database of known film grain patterns

Based on undecimated wavelet transform

Up to several seconds per frame

Image DeGrainRecurse(Image input, int level = 0) {
 Image HY,LY,HH,HL,LH,LL,HHP,HLP,LHP,LLP,pSum1,pSum2,out;



A'ENI

The recursive wavelet-based degraining visual effect in C++
 Visual primitives are chained together via image temporaries to form a DAG
 DAG construction is captured through delayed evaluation.

```
Imperial College
```

- Functor represents function over an image
- Kernel accesses image via indexers
- Indexers carry metadata that characterises kernel's data access pattern

```
class DWT1D : public Functor<DWT1D, <u>eParallel</u>> {
   Indexer<<u>eInput</u>, <u>eComponent</u>, <u>e1D</u>> Input;
   Indexer<<u>eOutput</u>, <u>eComponent</u>, e0D> HighOutput;
   Indexer<<u>eOutput</u>, <u>eComponent</u>, e0D> LowOutput;
   mFunctorIndexers(Input, HighOutput, LowOutput);
   DWT1D(Axis axis, Radius radius) : Input(axis, radius) {}
```

```
HighOutput() = high;
LowOutput() = centre - high;
```



ITTI BYCAT HINGLO

^{1;} One-dimensional discrete wavelet transform, as indexed functor
Compilable with standard C++ compiler
Operates in either the horizontal or vertical axis
Input indexer operates on RGB components separately
Input indexer accesses ±radius elements in one (the axis) dimension

Use of indexed functors is optimised using a source-to-source compiler (based on ROSE, www.rosecompiler.org)

SOURCE ATCHICEUTE



Goal:

- single source code, high-performance code for multiple manycore architectures
- Proof-of-concept: two targets
 - Very different, need very different optimisations



Lots of cache per threadLower DRAM bandwidth

SIMD Multicore CPU



WO GENERER REFER

- Very, very little cache per thread
- Very small scratchpad RAM shared by blocks of threads
- Higher DRAM bandwidth

SIMT Manycore GPU

Code generation for conventional PC with SSE ("SIMD") instructions:

TRASHIDEREE

Aggressive loop fusion and array contraction

Using the CLooG code generator to generate the loop fragments

Vectorisation and Scalar promotion

Correctness guaranteed by dependence metadata

If-conversion

Generate code to use masks to track conditionals

Memory access realignment:

In SIMD architectures where contiguous, aligned loads/stores are faster, placement of intermediate data is guided by metadata to make this so

Contracted load/store rescheduling:

- Filters require mis-aligned SIMD loads
- After contraction, these can straddle the end of the circular buffer we need them to wrap-around
- We use a double-buffer trick...

Imperial College London SINT – COLLE GENERATION for INITIA's BUDA

Constant/shared memory staging

Where data needed by adjacent threads overlaps, we generate code to stage image sub-blocks in scratchpad memory

Maximising parallelism

- Moving-average filters are common in VFX, and involve a loopcarried dependence
- We catch this case with a special "eMoving" index type
- We create enough threads to fill the machine, while efficiently computing a moving average within each thread

Coordinated coalesced memory access

- We shift a kernel's iteration space, if necessary, to arrange an thread-to-data mapping that satisfies the alignment requirements for high-bandwidth, coalesced access to global memory
- We introduce transposes to achieve coalescing in horizontal moving-average filters

Choosing optimal scheduling parameters

Resource management and scheduling parameters are derived from indexed functor metadata, and used to select optimal mapping of threads onto processors.



Domain-specific "active" library Visual effects Finite element encapsulates specialist performance expertise

- Each new platform requires new performance tuning effort
- So domain-specialists will be doing the performance tuning
- Our challenge is to support them



A ATTA TOTATAS

GPU Multicore FPGA Quantum?

Parallelism is everywhere Parallelism is essential

Parallelism is disruptive – it breaks abstractions



Eckert was wrong – we just need the right...

- Language
- Machine
- Discipline
- Abstractions
- Education

So what of the future?

Eckert was right –

- Avoid parallel programming!
- Isolate ordinary software from parallelism

Tools to build really clever parallel implementations Tools to deliver them And protect us from what lurks below

http://www.ralphclevenger.com