

HPC Advisory Council Switzerland Conference 2012

March 13, 2012 Lugano

#### Pattern-based Parallel Edge Preserving Algorithm for Salt-and-Pepper Image Denoising

UNIVERSITÀ DEGLI STUDI DI TORINO



PROGRAMME

Marco Aldinucci - Uni. of Torino, Italy, HPCAC University Award 2011 Winner Maurizio Drocco - Uni. of Torino, Italy, supported by HPCAC University Award 2011

Massimo Torquati and Marco Danelutto - Uni. Pisa, Italy Concetto Spampinato - Uni. of Catania, Italy Massimiliano Meneghin - IBM Research, Ireland Peter Kilpatrick - Queen's Uni. Belfast, U.K.



## Outline

- \* ParaPhrase
  - Parallel Patterns for Heterogeneous Multicore Systems
    - EC-STREP FP7
- \* Concurrency and multi-core, the theoretical background
  - a personal perspective
- \* FastFlow
  - A programming model (and a library)
  - Fast (<10 nS) core-to-core lock-free messaging</li>
  - Supporting multi-core and accelerators
- \* Testbed: a novel two-phase image/video denoiser
  - clean salt-and-pepper noise up to 90%
  - reconstruction quality at 50% close to jpeg compression

PARAPHRASE

2

♦ demo

## ParaPhrase

- \* ParaPhrase: Parallel Patterns for Heterogeneous Multicore Systems
  - ✤ 3 Year targeted research project (FP7 STREP)
    - Runs from 1/10/11 to 30/9/14
    - Funded by objective 3.4, "Computing Systems"
    - Project Number ICT-2011-288570
  - 9 partners from five countries
  - €3.5M budget, €2.6M EU contribution
  - Coordinated by the U. of St Andrews, Scotland, UK

#### PARAPHRASE

## Project Consortium

Country
UK
UK
Israel
Austria
UK
Italy
Germany
Italy
UK

### Multi-cores: programming issues

### \* We can muddle through on 2-16 cores

- modified sequential code may work
- we may be able to use multiple programs to soak up cores
- BUT larger systems are much more challenging

### \* Fundamentally, programmers must learn to "think parallel"

- this requires new high-level programming constructs
- 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

PARAPHRASE

5

• this needs to be packaged/abstracted!

### The need of a programming model

- \* Applications programmers must be systems programmers
  - insufficient assistance with abstraction, too much complexity to manage
- \* Difficult/impossible to scale, unless the problem is simple
- \* Future multi/many-core
  - Probably not just scaled versions of today's multi/many-core
    - Highly heterogeneous: message-passing, shared-memory cc-NUMA, variable cache consistency, specialized cores (XML, crypto, ...), soft-cores (FPGA), ...
- \* The approaches provide libraries
  - they need to provide abstractions, i.e. programming models

### The Implications for Programming

- \* We must program heterogeneous systems in an integrated way
- \* It will be impossible to program each kind of core differently
- \* It will be impossible to take static decisions about placement etc.

Overall, we aim to produce a new pattern-based approach to programming parallel systems.

- I. develop high-level design and implementation patterns that are capable of easily exposing useful parallelism on heterogeneous multicore/manycore systems.
- 2. develop new dynamic mechanisms to support adaptivity for heterogeneous multicore/manycore systems
- 3. verify that these patterns and adaptivity mechanisms can be used easily and effectively to develop a wide range of real-world applications

Concurrency and multi-core theoretical background: a personal perspective

PARAPHRASE

## Nowadays

#### \* E.g. Intel Sandybridge, AMD Opteron

- \* cache-coherent
- \* 10 core per socket (20 contexts)
- \* cc-NUMA (as matter of a fact)

#### \* NVidia/AMD GPGPU/Hybrid

- \* SIMD, no global synch
- \* performance only with proper and **not automatic** memory hierarchy management

#### \* Intel MIC CPU/GPGPU

- \* ring-based interconnection, variable coherency
- \* apparently even more NUMA

#### \* IBM powerEN

- \* general purpose cores
- \* specialised cores, soft cores?



|0|

### From programming/tuning viewpoint ... the simplest is already too complex ...

- \* Exploit cache coherence
- \* Memory fences are expensive
  - Increasing core count will make it worse
- \* Fine-grained parallelism is hard to achieve
  - I/O bound problems, High-throughput, Streaming, Irregular DP problems
  - Automatic and assisted parallelisation solves uniform&easy cases
    - OpenMP, Par4all, ...
- \* From programming model viewpoint, SIMD/GPGPU (maybe) worsen the scenario
  - Well-understood programming model
  - Atomic ops in memory (i.e. fences) are still needed
  - Not everything can be described with do independent (a.k.a. map)

### Micro-benchmarks: farm of tasks

Used to implement: parameter sweeping, master-worker, etc.

```
void Emitter () {
                                       int main () {
  for ( i =0; i <streamLen;++i){</pre>
                                         spawn thread( Emitter ) ;
                                         for ( i =0; i <nworkers;++i){</pre>
    task = create_task ();
    queue=SELECT_WORKER_QUEUE();
                                           spawn thread(Worker);
    queue ->PUSH(task);
                                         }
                                        wait_end () ;
                                       }
void Worker() {
 while (!end_of_stream) {
 myqueue ->POP(&task);
                                                  W_2
  do work(task) ;
                                                  Wn
```

PARAPHRASE

### Task farm with POSIX lock/unlock



### Can we avoid locks?

- \* Under relaxed memory models, using CAS/atomic ops
  - "lock-free" data structures
  - they perform better than lock-based
  - they fence the memory and pay cache coherency reconciliation overhead

PARAPHRASE

|4

- ✦ in GPUs ...
  - CAS/atomic ... you have to go to the global memory

### Lock vs CAS at fine grain (0.5 $\mu$ S)



## Re-starting from the basics

### \* Reducing the problem to the bare bones

- Producer-Consumer model (streaming)
- Directly control thread blocking using non-blocking synchronisations
- Directly design the "data channel"
  - Having clear how data move in the whole memory hierarchy

PARAPHRASE

\* Restarting from the FIFO queue



# Interaction models: theoretical background (in a nutshell)

- \* low-level synchronisation in the shared memory model
  - Mutual Exclusion (mutex)
    - typically used as basic building block of synchronisations
  - Producer Consumer

### \* they are not equally demanding

- Mutual Exclusion is inherently more complex since requires deadlock-freedom
  - require interlocked ops (CAS, ...), that induces memory fences, thus cache invalidation
  - Dekker and Bakery requires Sequential Consistency (++)
- Producer Consumer is a cooperative (non cyclic) process

### Recap: coherence and consistency

### \* Memory/Cache Coherence

Deal with multiple replicas of the same location in different caches

Thread 1



 Deal with the ordering in which writes and reads at different locations take effect in memory (issued by either the same or different processors/cores)



write(A,3)

write(A,1) read(A,?)

## FastFlow SPSC queues

```
push_nonbocking(data) {
  if (NEXT(head) == tail) {
    return EWOULDBLOCK;
  }
  buffer[head] = data;
  head = NEXT(head);
  return 0;
}
pop_nonblocking(data) {
  if (head == tail) {
    return EWOULDBLOCK;
  data = buffer[tail];
  tail = NEXT(tail);
  return 0;
}
```

#### \* Proved to be correct under SC

- doesn't work under weaker models
  - TSO, e.g. x86

PARAPHRASE

• WO, e.g. PPC, ARM

19

#### Lamport FIFO - 1983

## FastFlow SPSC queues

```
push_nonbocking(data) {
  if (NEXT(head) == tail) {
    return EWOULDBLOCK;
  }
  buffer[head] = data;
  head = NEXT(head);
  return 0;
}
pop_nonblocking(data) {
  if (head == tail) {
    return EWOULDBLOCK;
  data = buffer[tail];
  tail = NEXT(tail);
  return 0;
}
```

Lamport FIFO - 1983



FastFlow FIFO derived from PICI (Higham and Kavalsh, ISPAN 1997)

## FastFlow SPSC queues



PARAPHRASE

Lamport FIFO - 1983

derived from PICI (Higham and Kavalsh, ISPAN 1997)

### Lock vs CAS vs SPSC FastFlow (50 µS)



### Lock vs CAS vs SPSC FastFlow (5 µS)



APHPASI

### Lock vs CAS vs SPSC FastFlow (0.5 µS)



P

## Medium grain (5 µS workload)



## Fine grain (0.5 µS workload)



#### PARAPHRASE

### Producer-Cons with FF queue: latency UNIVER

4 sockets  $\times$  8 core  $\times$  2 contexts

Xeon E7-4820 @2.0GHz Sandy Bridge 18MB L3 shared cache, 256K L2

DK Panda said today MPI is ~190 ns







# Wait-free fence-free unbound dynamic queue (uSPSC)



M.Aldinucci, M. Danelutto, P. Kilpatrick, M. Meneghin, M. Torquati. An Efficient Synchronisation Mechanism for Multi-Core Systems. EuroPar 2012. Submitted.

PARAPHRASE

### Bandwidth test (with uSPSC)

http://www.1024cores.net/home/technologies/fastflow





### Related Work: Lock-free, CAS-free, wait-free

- \* Single-Producer-Single-Consumer FIFO queues
  - Lamport et al. 1983 Trans. PLS (Sequential consistency only passive)
  - Higham and Kavalsh. 1997 ISPAN (PICI TSO + proof passive)
  - Giacomoni et al. 2008 PPoPP (TSO + cache slipping passive)
- \* Multiple-Producers-Multiple-Consumers FIFO queues
  - with CAS (two of them) Michael and Scott (PODC96)
    - Also implemented in FastFlow, require deferred reclamation/hazard pointers to avoid ABA problem
  - without CAS passive I Cannot be done
  - without CAS active I FastFlow
- \* Extending the taxonomy with locking algorithms is clearly useless ....

## Why elaborating on queues?

- \* Queues are the base of Producer-Consumer
  - reference paradigm for data-centric and data-flow models
  - orthogonal and complementary to Transactional Memories
    - that address competition
  - can be used in both shared memory and message passing models
- \* Queues are the primary (maybe the only) global synchronisation mechanism in current GPGPUs
  - take a look to NVidia/CUDA programming manual
  - efficient global synch can pave the way to a broader usage of GPUs

30

PARAPHRASE

• is it possible an atomic-free global memory queue on GPUs?





## Lock-free and CAS-free?

### \* Mutex cannot be done Single-Producer-Single-Consumer (SPSC) can be done

- Producer-Consumer is inherently weaker with respect to Mutex
- It does require the cooperation of partners whereas Mutex does not
- Expressive enough to build a streaming (or dataflow) programming framework
  - MPMC = SPSC + mediator threads
- \* But what about productivity at large scale?
  - Write a program is defining a graph encoding true dependencies ... not really easy

32

### FastFlow is based on producer-consumer



- Lock-free/fence-free non-blocking synchronisations
- C++ STL-like implementation
- thread-model agnostic (pthreads, QT, windows threads, ...)
- compliant with other synchronisation mechanisms in the business code (e.g. locks and semaphores)

PARAPHRASE

### Pattern-based approach: rationale

\* Abstract parallelism exploitation pattern by parametric code

- E.g. higher order function, code factories, C++ templates, ...
- Can composed and nested as programming language constructs + offloading
- Stream and Data Parallel
- \* Platform independent
  - Implementations on different multi/many-cores
  - Support for hybrid architectures thanks to pattern compositionality
- \* Provide state-of-the-art, parametric implementation of each parallelism exploitation pattern
  - With natural way of extending patterns, i.e. OO
  - Functional (seq code) and tunable extra-functional (QoS) parameters

# Patterns, their implementation, and their purpose



## Composition

- \* Composition via C++ template meta-programming
  - CPU: Graph composition
  - GPU: CUDA streams
  - CPU+GPU: offloading
- \* farm{ pipe }

\*

- \* pipe(farm, farm)
- \* pipe(map, reduce)



## On my personal view of MPI ...

- \* I love it, I've using it for 15 years
  - as a target language for building compilers for high-level parallel languages
- \* MPI collective operations are good
- \* they are an example of patterns (a.k.a. skeletons)
  - even if not particularly high-level
- \* the problem are Send/Rcv
  - and the other over hundred (and counting) operations?
    - are they perceived as really useful from anybody?
  - they are on the language as well, and people use them
- \* one problem of "cars" are drivers (and crashes)
  - how much expertise do you need to write/tune/debug something barely decent?





#### \* Patterns at the high-level

- Currently as C++ templates
- Set of patterns can be extended, semantics of patterns can be changed, complexity gracefully increase with semantic distance
- \* Used to generate cyclic streaming networks (of threads, ...)
  - Graphs, describing true data dependencies. Can be composed and transformed as graphs
  - Cyclic graphs need unbound queue
  - Heterogeneous cores, thread affinity, memory affinity, NUMA, can be managed while mapping graph onto the metal

### Applications (see <a href="http://di.unito.it/fastflow">http://di.unito.it/fastflow</a> )

- \* Smith-Waterman (based on SWPS3/SSE2, ParCo 10)
  - Recently discovered it is used as baseline for FPGA design at it seems it is among the fastest SW solutions
- \* C4.5 (EMCL-PKDD 10, Concurrency prat. exp. submitted)
- \* StochKit-ff (HiBB I I)
- \* MonteCarlo sims & Gillespie's biological sims (PDP 11, HiBB 11)
- \* Cholesky (PDP 12)
- \* nProbe-ff (part of nTop software, ParCo II)
- \* pbzip-ff (Wiley book 12)
- \* nqueens, k-means, fibonacci, Mandelbrot QT, and many microbenchmarks
- \* Two-phase edge-preserving denoiser (IEEE IPTA submitted)

### 2012: Cholesky fact vs (very efficient) PLASMA libraries - IEEE PDP 2012

#### Targeting multi cores by structured programming and data flow

M. Aldinucci°, L. Anardu<sup>\*</sup>, M. Danelutto<sup>\*</sup>, P. Kilpatrick<sup>†</sup>, M. Torquati<sup>\*</sup>

Dept. Computer Science, Dept. Computer Science, U Dept. Computer Science, Queen

#### Abstract

Data flow techniques have been around sin were used in compilers for sequential language duction they were also considered as a possible ing, although the impact here was limited. Re has been identified as a candidate for efficient programming models on multi-core architectur the burden of determining data flow "macro" programmer, while the compiler/run time sys ficient scheduling of these instructions. We di programming approach supporting automatic



Figure 10:  $mdf^3$  vs PLASMA library. Cholesky factorization for a single 1024x1024 complex matrix (Intel Nehalem).

### C4.5 (Fine grain - irregular D&C) PKDD 2011

#### Decision Tree Growing and Pruning on General-Purpose Multicore

Marco Aldinucci, Salvatore Ruggieri, and Massimo Torquati

Abstract—The whole computer hardware industry embraced multicores. The extreme optimisation of sequential algorithms is then no longer sufficient to squeeze the real machine power, which can be only exploited via thread-level parallelism. Decision tree algorithms exhibit natural concurrency that makes them suitable to be parallelised. This paper presents an in-depth study of the parallelisation of an implementation of the C4.5 algorithm for multicore architectures. In addition to the tree growing phase, we cover also the so far unaddressed problem of parallelising the error-based pruning with grafting phase. We characterise elapsed time lower bounds for the forms of parallelisations adopted, and achieve close to optimal performances. Our implementation is based on the FastFlow para<sup>11</sup>

Index Terms—Parallel class

ideal Forest Cover 7 Census-Income KDD Cup 99 6 SvD10M9A U.S. Census 5 Speedup 2 0 2 3 4 5 6 7 8 9 10 11 1 FastFlow's worker threads

Fig. 11: NAP strategy speedup.



Fig. 12: Speedup of the NAP strategy vs no. of attributes (for 1M sample subset of SvD10M9A).

### Smith Waterman vs TBB vs OpenMP vs Cilk Parallel Computing 2010



UNIVERSITÀ DEGLI STUDI DI TORINO



## From users: Intel or AMD

### \* We always test on both platforms

### \* Assembler

- MONITOR/MWAIT instruction appear interesting for non-blocking synchronisations
  - accessible only in Ring 0, why?
  - IBM PowerEN ...
- \* Memory bandwidth is the key
  - but memory affinity tools appear to be quite primitive in both platforms

- how data structures are related with access pattern?
  - we are trying to address the problem with patterns ...

## A novel two-phase edge preserving parallel de-noising

Maurizio Drocco University of Torino (MSc. student) with the support of HPC AC university award

PARAPHRASE

### Two-phase denoising + Demo



PARAPHRASE

PNSR 43.29dB MAE 0.35

PNSR 32.75dB MAE 2.67

PNSR 23.4 MAE 11.21

## Denoising explained



## On the algorithm

#### \* Adaptive median filters

- well-known
- smooth the image if used in rebuilding phase

#### \* Variational methods

- edge-preserving, can be used for any kind of noise
- the formula is well-known, our algorithm is particularly fast
  - because we found several loop invariants in the process
- if used alone can destroy "good" pixels
- not used up to now because too slow
  - > I h on Matlab to rebuild a 256x256 image with 30% of noise

#### \* Coupling them

- detect: increase the speed by reducing false positives
- denoise: guarantee a good restoration because "touches" only noisy points







# Offloading on soft (i.e. not used cores) and HW accelerators

PARAPHRASE



- \* offloading
  - onto other cores and accelerators

## Speedup (multi-core)



4 sockets x 8 core x 2 contexts Xeon E7-4820 @2.0GHz Sandy Bridge 18MB L3 shared cache, 256K L2

PARAPHRASE

### Parallelism exploitation and patterns

### \* Application programmed high-level with FastFlow

- using Intel OpenCV
  - be careful, it is not fully thread-safe, all non thread-safe operation should be managed by a single thread (i.e. main)
- using offloading support, i.e. wrap a pattern-based code into a soft or HW "accelerator"
  - nice technique to port existing legacy code with very low effort
  - support for "semi-automatic" loop streamization

### \* Let us understand how through an example

Two-phase denoising - video version

### Edge-preserving denoiser: video

```
#include <opencv/highgui.h>
#include <opencv/cv.h>
int main(int argc, char *argv[]) {
  CvCapture *capture;
  IplImage * frame, clean_frame;
  char key;
  vector<noisy_t> noisy;
  cvNamedWindow("Video", CV_WINDOW_AUTOSIZE);
  capture = cvCreateCameraCapture(CV_CAP_ANY);
  //capture = cvCreateFileCapture("/path/to/your/video/test.avi");
  while(true) {
    frame = cvQueryFrame(capture);
                                          // get a frame from device
    noisy = myDetect(frame);
                                          // detect noisy pixels
    clean_frame = myDenoise(frame, noisy); // denoise the frame
    cvShowImage( "Video", clean_frame); // show the denoised frame
    key = cvWaitKey(100);
  }
  cvReleaseCapture(&capture);
  cvDestroyWindow("Video");
}
```

## Denoising explained



### Single design - many implementations



Parallel - Farm(Pipeline(myDetect,myDenoise))

Changing the structure does not require re-writing business code (gray)

Sequential

### Single design - many implementations



Parallel - Pipeline(myDetect,Map(myDenoise))

Changing the structure does not require re-writing business code (gray)

Sequential

### Single design - many implementations



Parallel - Farm(Pipeline(myDetect, Map(myDenoise)))

Changing the structure does not require re-writing business code (gray)

Sequential

## FF multi-core vs FF hybrid

noise	FF 32 cores Intel 4x8x2 2GHz	FF 8 cores (detect) + Tesla C2050 (denoise)	Seq Intel 4x8x2 2GHz	
Lena 512x512				
10	I.8 s	1.9 s	32 s	
50	6.5 s	2.3 s	162 s	
90	10.9 s	2.8 s	290 s	
Space 4096x4096				
10	78 s	12 s	2093 s	
50	373 s	46 s	10400 s	
90	665 s	77 s	18571 s	





# Nondeterministic variants & Convergence speed

#### \* flat (used in CUDA version)

- \* use block halo, block size = I
- \* do independent
- \* deterministic slow convergence
- \* easy CPU and GPU can be linearized

#### \* border

- \* use block halo
- \* do independent on tiles, do across within tiles
- \* nondeterministic fast convergence
- \* easy on CPU and GPU cannot be linearized

#### \* **std** (used in multicore version)

- \* don't use a block halo
- \* do independent on tiles, do across within tiles
- \* nondeterministic fast convergence
- \* easy CPU and GPU cannot be linearized
- \* cluster
  - \* do independent on tiles, do across within tiles
  - \* deterministic fast convergence
  - \* easy on CPU, difficult on GPU can be linearized



PARAPHRASE

## Conclusions

#### \* FastFlow C++ pattern-based framework

- A tiny, lightweight & open research framework
  - 2 years old 6K downloads 30K contacts x86/PPC/ARM(?) + Linux/Mac/Win
  - Adopted as one run-time technology in ParaPhrase
  - Laboratory to experiment new run-time solutions
    - GPGPU integration (working on), Infiniband RDMA integration (working on), HW blocking reads (thinking on), HW transactional mem (thinking on) ...
    - Stream-specific parallel memory allocator: faster than TBB (testing)
- ✦ A step forward in parallel programming models
  - Data-centric, focus on messaging and synchronization, thread model agnostic
  - support both message passing and shared memory
- High-level = performance & portability
  - Speedup starting from ~20 clock cycles workload on standard x86\_64 (TBB ~ 50K)
  - FF/AVX Smith-Waterman among fastest existing SW solutions
  - Tested on dozen of apps, comparable or faster than TBB/OpenMP

### Conclusions

### \* GPGPUs

- Needs high-level, CUDA/OpenCL too close to the metal
- Well integrate with functional style and higher order partterns
- \* Two-phase denoising
  - Novel, fast and efficient
    - Works up to 95% or noise, comparable to jpeg on 50% of noise
  - CPU/GPGPUs/Hybrid
  - Edge-preserving restoration works also for other kinds of noise

PARAPHRASE

6

• Working on detection and color images





## FF-allocator (written in FF)



- Faster than posix, often faster than hoard and TBB
  - unpublished, but available on sourceforge
  - needs lot of comparative testing to be published
- \* Implements deferred deallocation to avoid ABA problem (i.e. garbage collection)

The graph is now cyclic and needs unbound queues to avoid deadlocks

## On programming model

- \* Shared memory or message passing (data)
  - + message passing (synchronisations)
  - Graphs exactly describes the (true) data dependency pattern
  - Additional synchronisations can be added (e.g. locks) in the user code

### \* Queue can pass pointers or data

- Passing data means copying it. A proper usage of allocator might significantly enhance locality. Copying non strictly needed data is overhead.
- The balance depends on the application. Should be studied more.

## On programming model



The balance depends on the application. Should be studied more.

Thanks to Massimo Torquati & Marco Danelutto Computer Science Dept. - University of Pisa - Italy Massimiliano Meneghin IBM Research, Ireland Peter Kilpatrick Queen's University Belfast, U.K. Maurizio

