















### A Parallel Edge Preserving Algorithm for Salt and Pepper Image Denoising

Marco Aldinucci - Uni. of Torino, Italy Parallel and distributed computing research cluster

Massimo Torquati - Uni. Pisa, Italy Maurizio Drocco - Uni. of Torino, Italy Concetto Spampinato - Uni. of Catania, Italy Simone Palazzo - Uni. of Catania, Italy





- \* A two-phase image/video denoiser
- \* FastFlow
  - A programming model (and a library)
  - Fast (<10 nS) core-to-core lock-free messaging</li>
  - Fast zero-copy distributed messaging
  - Supporting multi-core and accelerators
- \* two-phase denoiser: parallel implementation
  - clean salt-and-pepper noise up to 90%
  - ✦ demo



B

UNIVERSITÀ DEGLI STUDI DI TORINO

### A two-phase image denoiser for salt & pepper



### Salt&Pepper noise

#### \* Salt&Pepper noise

- Electronic and signal noise
- Uniform distribution of "saturated" white and black pixels
  - Often measured as percentage of affected vs overall pixels
- \* Typically restored using statistic filters
  - e.g. median, median-adaptive
  - Not satisfactory for high levels of noise, also good pixels are filtered (image is smoothed)



### Improved S&P filtering



DI TORINO

6



- median

- variational (Nikolova)

Impulse detector

- progressive switching median
- adaptive median
- neural/bayesian networks
- fuzzy/neuro-fuzzy logic

#### \* Performance

- Improved, if detection faster than restoration
- Rationale: efficient detection + high-quality denoising
  - they can be pipelined and independently parallelized

### Variational method



- \* Iterative optimization problem
  - Theoretically and practically more expensive than statistical filtering
    - Matlab w 256x256 image with 50% of noise needs hours of computation time
  - Many local minimum, fixed point criteria for termination not trivial
    - Quasi-Newtonian methods with PSNR and MAE variation can be used

$$\min_{u \in N} F(u) = \alpha \int R(u) + \beta \int D(u, d)$$

regularization term data fidelity term

$$F_d|_N(u) = \sum_{(i,j)\in N} \left[ |u_{i,j} - d_{i,j}| + \frac{\beta}{2} (S_1 + S_2) \right]$$

$$S_1 = \sum_{\substack{(m,n) \in V_{i,j} \cap N}} 2 \cdot \varphi(u_{i,j} - d_{m,n})$$
$$S_2 = \sum_{\substack{(m,n) \in V_{i,j} \cap N^c}} \varphi(u_{i,j} - u_{m,n})$$

$$\varphi(t) = |t|^{\alpha}$$
 with  $1 < \alpha \le 2$ 





PARAPHRASE

Original Baboon standard test image 1024x1024

#### Restored



PNSR 43.29dB MAE 0.35

PNSR 32.75dB MAE 2.67

PNSR 23.4 MAE 11.21

8



















		Lena	- Noi	se leve	I		Peppe	rs — N	oise lev	el		Bridg	e — No	ise leve	el		Baboon	n — Noi	se leve	1
	10%	30%	50%	70%	90%	10%	30%	50%	70%	90%	10%	30%	50%	70%	90%	10%	30%	50%	70%	90%
PSNR	41.96	35.76	32.56	29.47	24.43	42.27	35.49	31.57	28.23	23.12	36.61	31.50	28.11	25.23	21.22	35.39	29.73	26.94	24.61	22.06
MAE	0.36	1.21	2.27	3.85	8.82	0.30	1.10	2.18	3.94	9.90	0.82	2.61	4.85	8.06	15.01	0.94	3.14	5.63	8.87	13.96

#### Comparison with single phase Chan's variational method

	Method	L	ena — I	Noise le	evel	Br	idge —	Noise I	evel
		30%	50%	70%	90%	30%	50%	70%	90%
PSNR	Chan's	35.20	32.21	29.03	22.46	30.29	27.91	25.00	19.02
PSNR	Our	35.76	32.56	29.47	24.43	31.50	28.11	25.23	21.22
MAE	Chan's	1.5	2.97	4.2	12.45	3.75	5.30	8.10	21.25
	Our	1.21	2.27	3.85	8.82	2.61	4.85	8.06	15.01

PARAPHRASE

10



- \* MPI, thread synchronizations, and CUDA can be used to parallelize almost everything
- \* They give you lot freedom, you can write almost everything, move you data in any fashion
- \* Eventually "they are as a "car", you can drive where you like, when you want, ..."
  - (cit. D.K. Panda, leader of the MPI-MVAPICH group)

### On the freedom of low-level methods UNIVERSIT



EGLI STUDI

12

#### Design your algorithm



car

![](_page_13_Picture_0.jpeg)

12

### On the freedom of low-level methods UNIVERSITÀ DEGLI STUDI

#### ... and running it!

![](_page_13_Picture_3.jpeg)

![](_page_14_Picture_0.jpeg)

![](_page_14_Picture_1.jpeg)

EU FP7 Strep 3.5 M€

### http://mc-fastflow.sourceforge.net/

# FastFlow

![](_page_14_Figure_5.jpeg)

- C++ pattern-based framework, open-source LGPL
- A tiny, lightweight & open research framework for HPC
  - 5K lines of code
- 3 years old over 8K downloads over 40K contacts
- x86/PPC/ARM + Linux/Mac/Win/iOS
- Multicore, GPGPU, distributed (TCP & Infiniband)

![](_page_15_Figure_0.jpeg)

### FastFlow (multicore)

#### Applications on multicore, many-core

Efficient and portable - designed with high-level patterns

**FastFlow** 

#### Streaming network patterns

Skeletons: pipeline, map farm, reduce, D&C, ...

#### Arbitrary streaming networks

Lock-free SPSC/MPMC queues + FF nodes

#### Simple streaming networks Lock-free SPSC queues + threading model

#### Multicore and manycore SMP: cc-UMA & cc-NUMA

### Layer 1: Simple streaming networks

![](_page_17_Figure_1.jpeg)

M.Aldinucci, S. Campa, M. Danelutto, M. Torquati. An Efficient Synchronisation Mechanism for Multi-Core Systems. EuroPar 2012.

![](_page_18_Picture_0.jpeg)

PARAPHRASE

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

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

> MPI is ~190 ns (D.K. Panda)

![](_page_18_Figure_4.jpeg)

![](_page_18_Figure_5.jpeg)

![](_page_19_Picture_0.jpeg)

#### Medium grain (5 µS workload) UNIVERSITÀ DEGLI STUDI DI TORINO

![](_page_20_Figure_1.jpeg)

### FastFlow (multicore)

Applications on multicore, many-core

Efficient and portable - designed with high-level patterns

**FastFlow** 

Streaming network patterns

Skeletons: pipeline, map farm, reduce, D&C, ...

Arbitrary streaming networks Lock-free SPSC/MPMC queues + FF nodes

Simple streaming networks Lock-free SPSC queues + threading model

> Multicore and manycore SMP: cc-UMA & cc-NUMA

![](_page_22_Picture_0.jpeg)

![](_page_22_Picture_1.jpeg)

2

#### \* Data Parallel

 is a method for parallelizing a single task by processing independent data elements of this task in parallel. The flexibility of the technique relies upon stateless processing routines implying that the data elements must be fully independent. Data Parallelism also supports Loop-level Parallelism where successive iterations of a loop working on independent or read-only data are parallelized in different flows-of-control and concurrently executed.

#### \* Task Parallel

 is explicit in the algorithm and consists of running the same or different code on different executors (cores, processors, machines, etc.). Different flows-of-control (threads, processes, etc.) may communicate with one another as they work.
 Communication usually takes place to pass data from one thread to the next as part of the same data-flow graph.

#### \* Stream Parallel

 can be used when there exists a partial or total order in a computation. By processing data elements in order, local state may be maintained in each filter.

### Layer 3: streaming networks patterns

- Composition via C++ template meta-programming
  - CPU: Graph composition
  - GPU: CUDA streams
  - CPU+GPU: offloading
- farm{ pipe }
- pipe(farm, farm)
- pipe(map, reduce)

![](_page_23_Figure_8.jpeg)

![](_page_23_Figure_9.jpeg)

GPGPU

### Layer 3: streaming networks patterns

![](_page_24_Figure_1.jpeg)

# At the bottom line, also on GPUs, it is matter of abstracting & engineering well-known concepts

![](_page_25_Picture_1.jpeg)

# + distributed

Applications on multicore, many core & distributed platforms of multicores Efficient and portable - designed with high-level patterns

![](_page_26_Figure_2.jpeg)

- network channels
  - P2P or collective
  - used as frontier node of streaming graph
  - can be used to merge graphs across distributed platforms
- No changes to programming model
  - when passing pointers data is serialised

M.Aldinucci, S. Campa, M. Danelutto, M. Torquati, P. Kilpatrick. Targeting distributed systems in FastFlow. CGW/EuroPar 2012

### 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
- Rationale
  - Decrease to bare minimum synchronization overhead (speedup)
  - Provide ready-to-use patterns (productivity, time-to-market)

![](_page_28_Picture_0.jpeg)

UNIVERSITÀ DEGLI STUDI DI TORINO

27

# Two-phase edge preserving parallel de-noising

### Denoising explained (video)

![](_page_29_Picture_1.jpeg)

![](_page_29_Figure_3.jpeg)

PARAPHRASE

In the video case the two stages can be pipelined on *m*-core+*m*-core m-core+GPGPU **GPGPU+GPGPU** 

and you haven't to decide it at design time o port the code

30

### Speedup (multi-core)

![](_page_30_Picture_1.jpeg)

UNIVERSITÀ DEGLI STUDI DI TORINO

3

![](_page_30_Figure_2.jpeg)

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

![](_page_31_Picture_0.jpeg)

### 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");
```

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

![](_page_32_Picture_1.jpeg)

UNIVERSITÀ DEGLI STUDI DI TORINO

![](_page_32_Picture_3.jpeg)

- \* offloading
  - onto other cores and accelerators

![](_page_33_Picture_0.jpeg)

![](_page_33_Figure_1.jpeg)

Parallel - Farm(Pipeline(myDetect,myDenoise))

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

Sequential

![](_page_34_Picture_0.jpeg)

![](_page_34_Figure_1.jpeg)

Parallel - Pipeline(myDetect, Map(myDenoise))

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

Sequential

![](_page_35_Picture_0.jpeg)

![](_page_35_Figure_1.jpeg)

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

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

Sequential

![](_page_36_Picture_0.jpeg)

![](_page_36_Picture_1.jpeg)

UNIVERSITÀ DEGLI STUDI DI TORINO

noise	FF 32 cores Intel 4x8x2 2GHz	FF 8 cores (detect) + Tesla C2050 (denoise)	Seq Intel 4x8x2 2GHz	
Lena 512x5	12			
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 4096	×4096			
10	78 s	12 s	2093 s	17
50	373 s	46 s	10400 s	22
90	665 s	77 s	18571 s	27

![](_page_37_Picture_0.jpeg)

UNIVERSITÀ DEGLI STUDI DI TORINO

36

![](_page_37_Picture_2.jpeg)

![](_page_37_Picture_3.jpeg)

### Conclusions

![](_page_38_Picture_1.jpeg)

#### \* Two-phase denoising

- Variational methods can be made 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
- \* GPGPUs
  - Needs high-level, CUDA/OpenCL too close to the metal
  - Well integrate with functional style and higher order patterns
- \* Image analysis is just an example of usage of FastFlow
  - classification, mining, compressing, string-alignment, network inspection (also used in nTop), and many more... See <a href="http://di.unito.it/fastflow">http://di.unito.it/fastflow</a>
  - MonteCarlo simulations
    - precessed CMS@CERN channel  $H \rightarrow ZZ \rightarrow 4I$  Higg's boson for July 2012 claim

![](_page_39_Picture_0.jpeg)

UNIVERSITÀ DEGLI STUDI DI TORINO

# Thank you

![](_page_39_Picture_3.jpeg)

# Nondeterministic variants & Convergence speed

![](_page_40_Picture_1.jpeg)

DI TORINO

39

#### \* 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

![](_page_40_Figure_21.jpeg)

n. cycle