

Processor-oblivious parallel algorithms with provable performances

- Applications

Jean-Louis Roch

MOAIS

Lab. Informatique Grenoble, INRIA, France

Overview

- Introduction : interactive computation, parallelism and processor oblivious
 - Overhead of parallelism : parallel prefix
- Machine model and work-stealing
- Scheme 1: Adaptive parallel algorithms
- Scheme 2: Amortizing the overhead of synchronization (Nano-loop)
- Scheme 3: Amortizing the overhead of parallelism (Macro-loop)
- Putting things together: processor-oblivious prefix computation

Interactive parallel computation?

Any application is "parallel":

- composition of several programs / library procedures (possibly concurrent) ;
- each procedure written independently and also possibly parallel itself.

Interactive Distributed Simulation
3D-reconstruction + simulation + rendering
[B Raffin & E. Boyer]
- 1 monitor
- 5 cameras,
- 6 PCs

Video

[B Raffin, MOAIS & E Boyer, PERCEPTION]

http://www-id.imag.fr/~jlroch/perso_html/talks/2007-01-MSRL-Berkeley/iceevr06-mov1.mpg

New parallel supports from small too large

- **Parallel chips & multi-core architectures:**
 - MPPSoCs (Multi-Processor Systems-on-Chips)
 - GPU : graphics processors (and programmable: Shaders, Cuda SDK)
 - Dual Core processors (Opteron, Itanium, etc.)
 - Heterogeneous multi-cores : CPUs + GPUs + DSPs+ FPGAs (Cell)
- **Commodity SMPs:**
 - 8 way PCs equipped with multi-core processors (AMD Hypertransport) + 2 GPUs
- **Clusters:**
 - 72% of top 500 machines
 - Trends: more processing units, faster networks (PCI-Express)
 - Heterogeneous (CPUs, GPUs, FPGAs)
- **Grids:**
 - Heterogeneous networks
 - Heterogeneous administration policies
 - Resource Volatility
- **Dedicated platforms:** eg Virtual Reality/Visualization Clusters:
 - Scientific Visualization and Computational Steering
 - PC clusters + graphics cards + multiple I/O devices (cameras, 3D trackers, multi-projector displays)

Parallelism induces overhead : e.g. Parallel prefix on fixed architecture

- **Prefix problem :**
 - input : a_0, a_1, \dots, a_n
 - output : π_1, \dots, π_n with $\pi_i = \prod_{k=0}^i a_k$
- **Sequential algorithm :**
 - for ($\pi[0] = a[0]$, $i = 1$; $i \leq n$; $i++$) $\pi[i] = \pi[i-1] \cdot a[i]$; performs only n operations
- **Fine grain optimal parallel algorithm :**

Critical time = $2 \cdot \log n$
but performs $2 \cdot n$ ops

Parallel requires twice more operations than sequential !!
- **Tight lower bound on p identical processors:**

Optimal time $T_p = 2n / (p+1)$
but performs $2 \cdot n \cdot p / (p+1)$ ops

Lower bound(s) for the prefix

Prefix circuit of depth d
 \Downarrow [Fitch80]
 #operations $> 2n - d$

$$\text{parallel time} \geq \frac{2n}{(p+1) \cdot \Pi_{ave}}$$

The problem

To design a single algorithm that computes efficiently prefix(a) on an arbitrary dynamic architecture

Sequential algorithm parallel P=2 ... parallel P=100 ... parallel P=max

Which algorithm to choose ?

Heterogeneous network Multi-user SMP server Grid

Dynamic architecture : non-fixed number of resources, **variable speeds**
eg: grid, ... but not only: SMP server in multi-users mode

Processor-oblivious algorithms

Dynamic architecture : non-fixed number of resources, variable speeds
eg: grid, SMP server in multi-users mode,

=> motivates the design of «processor-oblivious» parallel algorithm that:

- + is **independent** from the underlying architecture:
no reference to p nor $\Pi_i(t)$ = speed of processor i at time t nor ...
- + on a given architecture, has **performance guarantees** :
behaves as well as an optimal (off-line, non-oblivious) one

2. Machine model and work stealing

- Heterogeneous machine model and work-depth framework
- Distributed work stealing
- Work-stealing implementation : work first principle
- Examples of implementation and programs:
Cilk , Kaapi/Athapascan
- Application: Nqueens on an heterogeneous grid

Heterogeneous processors, work and depth

Processor speeds are assumed to change arbitrarily and adversarially:
model [Bender,Rabin 02] $\Pi_i(t)$ = **instantaneous speed** of processor i at time t
(in μ nit operations per second)

Assumption : $\Pi_{max}(t) < \text{constant} \cdot \Pi_{min}(t)$

Def. for a computation with duration T

- **total speed**: $\Pi_{tot} = \sum_{i=0..p} \sum_{t=0..T} \Pi_i(t)$
- **average speed** per processor: $\Pi_{ave} = \Pi_{tot} / P$

“Work” W = #total number operations performed

“Depth” D = #operations on a critical path
(~parallel “time” on ∞ resources)

For any greedy **maximum utilization** schedule:
[Graham80, Jaffe80, Bender-Rabin02]

$$\text{makespan} \leq \frac{W}{p \Pi_{ave}} + \left(1 - \frac{1}{p}\right) \frac{D}{\Pi_{ave}}$$

The work stealing algorithm

- A distributed and randomized algorithm that computes a greedy schedule :
 - > Each processor manages a local task (depth-first execution)

The work stealing algorithm

- A distributed and randomized algorithm that computes a greedy schedule :
 - > Each processor manages a local stack (depth-first execution)

> When idle, a processor steals the topmost task on a remote -non idle- victim processor (randomly chosen)

> **Theorem**: With good probability, [Acar,Bliech, Blumofe02, BenderRabin02]

- > #steals $< p \cdot D$
- > execution time $\leq \frac{W}{p \Pi_{ave}} + O\left(\frac{D}{\Pi_{ave}}\right)$

> **Interest**:
if W independent of p and D is small, work stealing achieves **near-optimal** schedule

Work stealing implementation

13

efficient policy (close to optimal) ← Scheduling → control of the policy (realisation)

Difficult in general (coarse grain)
But easy if D is small (work-stealing)
 Execution time $\leq \frac{W}{P1_{loc}} + O\left(\frac{D}{P1_{loc}}\right)$ (fine grain)

Expensive in general (fine grain)
But small overhead if a small number of tasks (coarse grain)

If D is small, a work stealing algorithm performs a **small number of steals**

=> **Work-first principle**: "scheduling overheads should be borne by the critical path of the computation" [Frigo 98]

Implementation: since all tasks but a few are executed in the local stack, overhead of task creation should be as close as possible as sequential function call

At any time on any non-idle processor, efficient local **degeneration** of the parallel program in a sequential execution

Work-stealing implementations following the work-first principle: Cilk

14

▪ **Cilk-5** <http://supertech.csail.mit.edu/cilk/>: C extension

- **Spawn f(a)**: sync (serie-parallel programs)
- Requires a shared-memory machine
- Depth-first execution with synchronization (on sync) with the end of a task:
 - Spawned tasks are pushed in double-ended queue
- "Two-clone" compilation strategy [Frigo-Leiserson-Randall98]:
 - on a successful steal, a thief executes the continuation on the topmost ready task;
 - When the continuation hasn't been stolen, "sync" = nop; else synchronization with its thief

```

01 cilk int fib (int n)
02 {
03   if (n < 2) return n;
04   else
05     {
06       int x, y;
07
08       x = spawn fib (n-1);
09       y = spawn fib (n-2);
10
11       sync;
12
13       return (x+y);
14     }
15 }
  
```

```

1  int fib (int n)
2  {
3    fib_frame *f;
4    f = (fib_frame *) malloc(sizeof(fib_frame));
5    f->n = n;
6    if (n < 2)
7      return n;
8    return n;
9  }
10
11 int x, y;
12
13 f->spawn = 1;
14 f->pc = 0;
15
16
17 if (sync() == FAILED)
18   ...
19   ...
20   ...
21   ...
22   ...
23   ...
24 }
  
```

▪ won the 2006 award "Best Combination of Elegance and Performance" at HPC Challenge Class 2, SC'06, Tampa, Nov 14 2006 [Kuszmaul] on SGI ALTIX 3700 with 128 bi-lithium

Work-stealing implementations following the work-first principle: KAAPI

15

▪ **Kaapi / Athapascan** <http://kaapi.gforge.inria.fr/>: C++ library

- Fork-<>(a, ...) with **access mode** to parameters (value,read,write;r/w,cw) **specified in f prototype** (macro dataflow programs)
- Supports distributed and shared memory machines; heterogeneous processors
- Depth-first (reference order) execution with synchronization on data access:
 - Double-end queue (mutual exclusion with compare-and-swap)
 - on a successful steal, one-way data communication (write&signal)

```

1 struct sum {
2   void operator()(Shared_L< int > a,
3                 Shared_L< int > b,
4                 Shared_W< int > r) {
5     { r->write(a->read() + b->read()); }
6   };
7 };
8
9 struct fib {
10  void operator()(int n, Shared_W<int> r)
11  { if (n < 2) r->write(n);
12    else
13      { int r1, r2;
14        Fork< fib >( n-1, r1 );
15        Fork< fib >( n-2, r2 );
16        Fork< sum >( r1, r2, r );
17      }
18  };
  
```

▪ Static scheduling won the 2006 award "Prix special du Jury" for the best performance at NQueens contest, Plugtests-Grid&Work'06, Nice, Dec. 1, 2006 [Luscau-Chevalier] on Grid 5000 1458 processors with different speeds.

N-queens: Takaken C sequential code parallelized in C++ / Kaapi

16

▪ T. Gautier & S. Guellon won the 2006 award "Prix special du Jury" for the best performance at NQueens contest, Plugtests-Grid&Work'06, Nice, Dec. 1, 2006

- Some facts [on on Grid'5000, a grid of processors of heterogeneous speeds]
 - NQueens (21) in 78 s on about 1000 processors
 - NQueens (22) in 502.9s on 1458 processors
 - NQueens(23) in 4435s on 1422 processors [-24.10¹³ solutions]
 - 0.625% idle time per processor
 - < 20s to deploy up to 1000 processes on 1000 machines [Taktak, Huard]
 - 15% of improvement of the sequential due to C++ (template)

Grid'5000 utilization during contest

6 instances NQueens(22)

Competitor X, Competitor Y, Grid'5000 free, N-Queens(23)

Experimental results on SOFA [CIMIT-ETZH-INRIA]

17

[Allard 06]

Bar-fem-implicit-32

Speedup GPU Bar-spring-euler

Kaapi (C++, ~500 lines) Cilk (C, ~240 lines)

Preliminary results on GPU NVIDIA 8800 GTX

- speed-up ~9 on Bar 10x10x46 to Athlon64 2.4GHz
- 128 "cores" in 16 groups
- CUDA SDK: "BSP"-like, 16 X [16 .. 512] threads
- Supports most operations available on CPU
- ~2000 lines CPU-side + 1000 GPU-side

3. Work-first principle and adaptability

18

- **Work-first principle**: -implicit- dynamic choice between two executions:
 - a sequential "depth-first" execution of the parallel algorithm (local, default);
 - a parallel "breadth-first" one.
- Choice is performed at runtime, depending on resource idleness: rare event if Depth is small to Work
- **WS adapts parallelism to processors with practical provable performances**
 - Processors with changing speeds / load (data, user processes, system, users,
 - Addition of resources (fault-tolerance [Cilk/Perch, Kaapi, ...])
- **The choice is justified only when the sequential execution of the parallel algorithm is an efficient sequential algorithm:**
 - Parallel Divide&Conquer computations
 - ...

-> **But**, this may not be general in practice

How to get both optimal work W_1 and W_∞ small?

19

- General approach: to mix both
 - a sequential algorithm with optimal work W_1
 - and a fine grain parallel algorithm with minimal critical time W_∞
- Folk technique : parallel, then sequential
 - Parallel algorithm until a certain « grain »; then use the sequential one
 - Drawback : W_∞ increases « » ... and, also, the number of steals
- Work-preserving speed-up technique [Blau-Pan04] sequential, then parallel Cascading [Laf09] : Careful interplay of both algorithms to build one with both W_∞ small and $W_1 = O(W_{seq})$
 - Use the work-optimal sequential algorithm to reduce the size
 - Then use the time-optimal parallel algorithm to decrease the time
 - Drawback : sequential at coarse grain and parallel at fine grain «

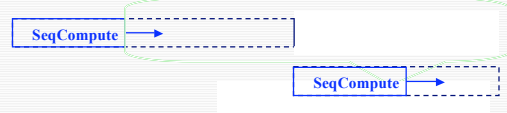
Extended work-stealing: concurrently sequential and parallel

20

Based on the work-stealing and the Work-first principle :
Instead of optimizing the sequential execution of the best parallel algorithm, let optimize the parallel execution of the best sequential algorithm

Execute always a sequential algorithm to reduce parallelism overhead
⇒ parallel algorithm is used only if a processor becomes idle (ie workstealing) [Roch&al2005,...] to extract parallelism from the remaining work a sequential computation

Assumption : two concurrent algorithms that are complementary:
- one sequential : SeqCompute (always performed, the priority)
- the other parallel, fine grain : LastPartComputation (often not performed)



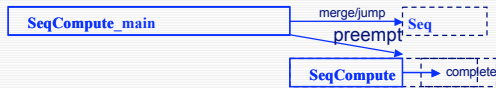
Extended work-stealing : concurrently sequential and parallel

21

Based on the work-stealing and the Work-first principle :
Instead of optimizing the sequential execution of the best parallel algorithm, let optimize the parallel execution of the best sequential algorithm

Execute always a sequential algorithm to reduce parallelism overhead
⇒ parallel algorithm is used only if a processor becomes idle (ie workstealing) [Roch&al2005,...] to extract parallelism from the remaining work a sequential computation

Assumption : two concurrent algorithms that are complementary:
- one sequential : SeqCompute (always performed, the priority)
- the other parallel, fine grain : LastPartComputation (often not performed)



Note:

- merge and jump operations to ensure non-idleness of the victim
- Once SeqCompute_main completes, it becomes a work-stealer

Extended work-stealing and granularity

22

• Scheme of the sequential process : nanoloop

```
while (not completed(Wrem) ) and (next_operation hasn't been stolen)
{
    atomic { extract_next k operations ; Wrem -= k ; }
    process the k operations extracted ;
}
```

• Processor-oblivious algorithm

- Whatever p is, it performs $O(p \cdot D)$ preemption operations (« continuation faults »)
-> D should be as small as possible to maximize both speed-up and locality
- If no steal occurs during a (sequential) computation, then its arithmetic work is optimal to the one W_{opt} of the sequential algorithm (no spawn/lock/copy)
-> W should be as close as possible to W_{opt}

- Choosing $k = \text{Depth}(W_{rem})$ does not increase the depth of the parallel algorithm while ensuring $O(W/D)$ atomic operations ;
since $D > \log_2 W_{em}$, then if $p = t$: $W \sim W_{opt}$

- Implementation : atomicity in nano-loop based on efficient local lock

- Self-adaptive granularity

Interactive application with time constraint

23

Anytime Algorithm:

- Can be stopped at any time (with a result)
- Result quality improves as more time is allocated

In Computer graphics, anytime algorithms are common:
Level of Detail algorithms (time budget, triangle budget, etc...)
Example: Progressive texture loading, triangle decimation (Google Earth)

Anytime processor-oblivious algorithm:

On p processors with average speed Π_{ave} , it outputs in a fixed time T a result with the same quality than a sequential processor with speed Π_{ave} in time $p \cdot \Pi_{ave}$

Example: Parallel Octree computation for 3D Modeling

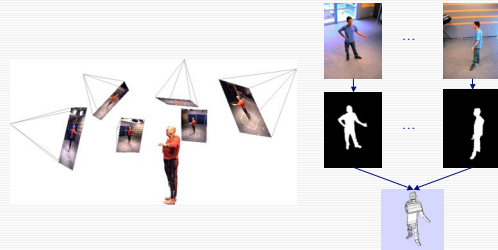
Parallel 3D Modeling

24

3D Modeling :

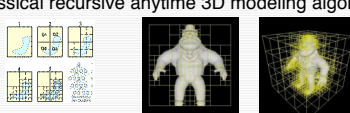
build a 3D model of a scene from a set of calibrated images

On-line 3D modeling for interactions: 3D modeling from multiple video streams (30 fps)

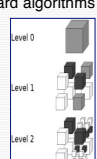


Octree Carving [L. Soares 06]

A classical recursive anytime 3D modeling algorithm.



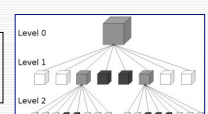
Standard algorithms with time control:



Depth first
+ iterative deepening

State of a cube:

- Grey: mixed => split
- Black: full : stop
- White: empty : stop



Width first

At termination: quick test to decide all grey cubes time control

Width first parallel octree carving ²⁶


Well suited to work-stealing

- Small critical path, while huge amount of work (eg. $D = 8, W = 164\,000$)
- non-predictable work, non predictable grain :

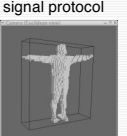
For cache locality, each level is processed by a self-adaptive grain :
"sequential iterative" / "parallel recursive split-half"

Octree needs to be "balanced" when stopping:

- Serially computes each level (*with small overlap*)
- Time deadline (30 ms) managed by signal protocol



Unbalanced



Balanced

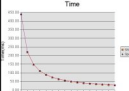
Theorem: W.r.t the adaptive in time T on p procs., the sequential algorithm:

- goes at most one level deeper : $|d_s - d_p| \leq 1$;
- computes at most : $n_s \leq n_p + O(\log n_s)$.

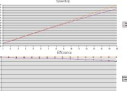
Results [L. Soares 06]

16 core Opteron machine, 64 images

- Sequential: 269 ms, 16 Cores: 24 ms
- 8 cores: about 100 steals (167 000 grey cells)



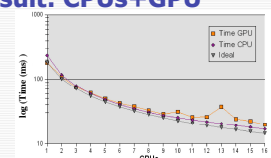
8 cameras, levels 2 to 10



64 cameras, levels 2 to 7

Preliminary result: CPUs+GPU

- 1 GPU + 16 CPUs
- GPU programmed in OpenGL
- efficient coupling till 8 but does not scale



4. Amortizing the arithmetic overhead of parallelism ²⁸

Adaptive scheme : `extract_seq/nanoloop // extract_par`

- ensures an optimal number of operation on 1 processor
- but no guarantee on the work performed on p processors

Eg (C++ STL): `find_if (first, last, predicate)`
locates the first element in [First, Last) verifying the predicate

This may be a drawback :

- unneeded processor usage ;
- undesirable for a library code that may be used in a complex application, with many components
- (or not fair with other users)
- increases the time of the application :
 - any parallelism that increases the execution time should be avoided

Motivates the building of **work-optimal** parallel adaptive algorithm (processor oblivious)

4. Amortizing the arithmetic overhead of parallelism (cont'd) ²⁹

Similar to nano-loop for the sequential process :

- that balances the -atomic- local work by the depth of the remaindering one

Here, by **amortizing** the work induced by the `extract_par` operation, ensuring this **work to be small** enough :

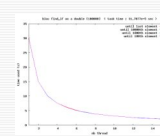
- Either w.r.t the -useful- **work already performed**
- Or with respect to the - useful - **work yet to performed** (if known)
- or both.

Eg : `find_if (first, last, predicate)` :

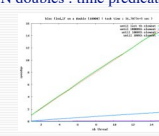
- only the work already performed is known (on-line)
- then prevent to assign more than $\alpha(W_{done})$ operations to work-stealers
- Choices for $\alpha(n)$:
 - $n/2$: similar to **Floyd's iteration** (approximation ratio = 2)
 - $n/\log n$: to ensure optimal usage of the work-stealers

Results on find_if [S. Guelton]

N doubles : time predicate ~ 0.31 ms



With no amortization macroloop



With amortization macroloop

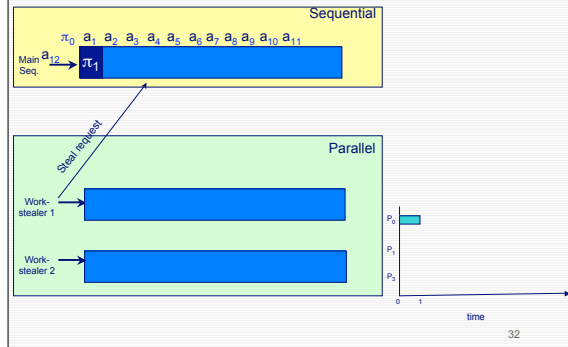
5. Putting things together *processor-oblivious prefix computation*

31

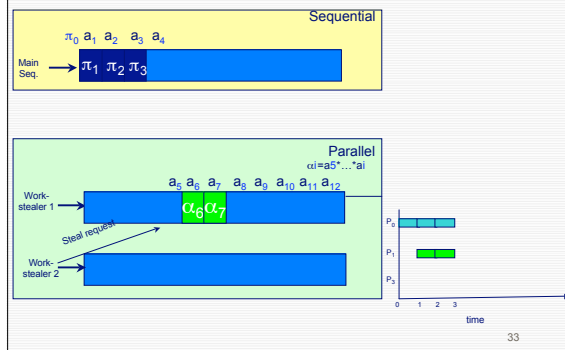
Parallel algorithm based on :

- compute-seq / extract-par scheme
- nano-loop for compute-seq
- macro-loop for extract-par

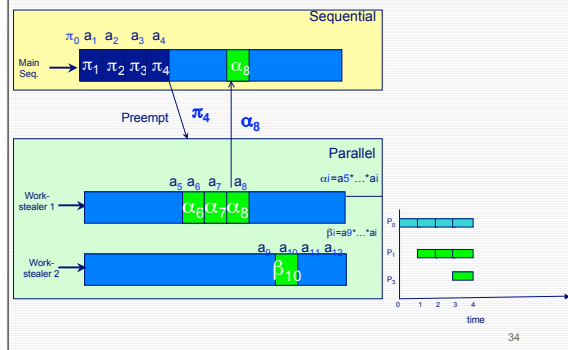
P-Oblivious Prefix on 3 proc.



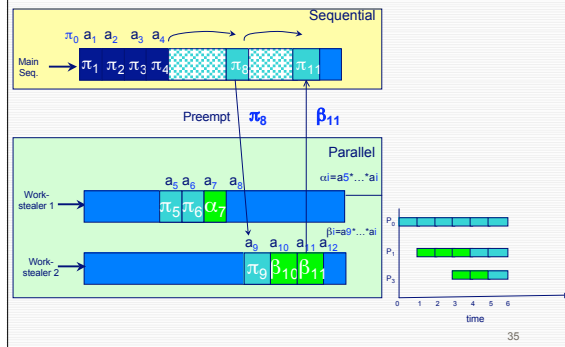
P-Oblivious Prefix on 3 proc.



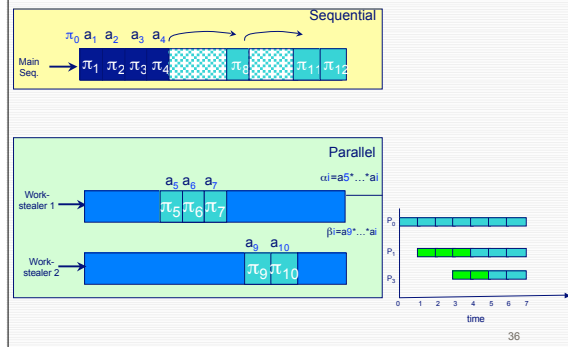
P-Oblivious Prefix on 3 proc.

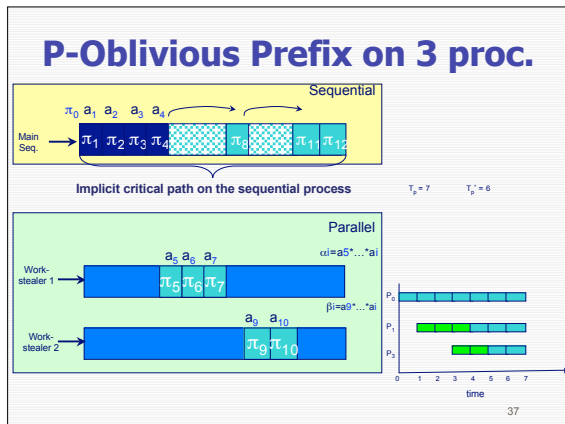


P-Oblivious Prefix on 3 proc.



P-Oblivious Prefix on 3 proc.





Analysis of the algorithm

Execution time $\leq \frac{2n}{(p+1) \cdot \Pi_{ave}} + O\left(\frac{\log n}{\Pi_{ave}}\right)$

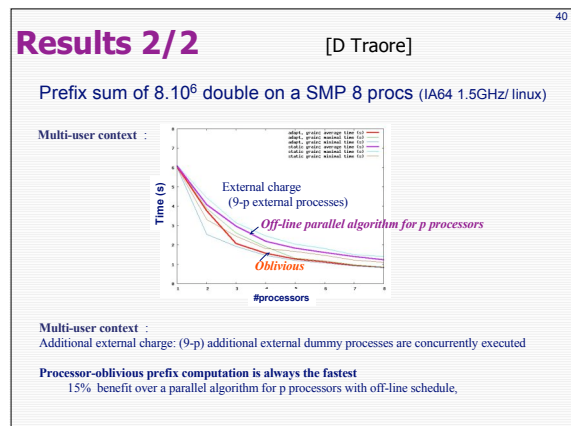
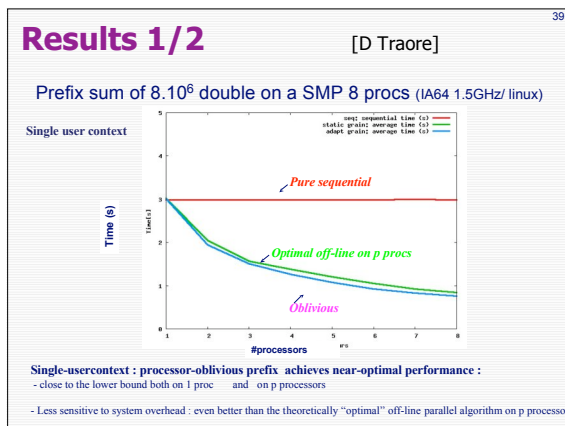
Sketch of the proof:

- Dynamic coupling of two algorithms that complete simultaneously:
 - Sequential: (optimal) number of operations S on one processor
 - Extract_par: work stealer perform X operations on other processors
 - dynamic splitting always possible till finest grain BUT local sequential
 - Critical path small (eg: $\log X$ with a $W_i = n / \log^2 n$ macroloop)
 - Each non constant time task can potentially be splitted (variable speeds)

$T_s = \frac{S}{\Pi_{ave}}$ and $T_p = \frac{X}{(p-1) \cdot \Pi_{ave}} + O\left(\frac{\log X}{\Pi_{ave}}\right)$

Algorithmic scheme ensures $T_s = T_p + O(\log X)$

\Rightarrow enables to bound the whole number X of operations performed and the overhead of parallelism $= (p+X) \cdot \beta_{ops_optimal}$.



Conclusion

- Fine grain parallelism enables efficient execution on a small number of processors**
 - Interest: portability; mutualization of code;
 - Drawback: needs work-first principle \Rightarrow algorithm design
- Efficiency of classical work stealing relies on work-first principle:**
 - Implicitly delinearates a parallel algorithm into a sequential efficient ones;
 - Assumes that parallel and sequential algorithms perform about the same amount of operations
- Processor Oblivious algorithms based on work-first principle**
 - Based on anytime extraction of parallelism from any sequential algorithm (may execute different amount of operations);
 - Oblivious: near-optimal whatever the execution context is.
- Generic scheme for stream computations:**
 - parallelism introduce a copy overhead from local buffers to the output
 - gzip / compression, MPEG-4 / H264

Kaapi (kaapi.gforce.inria.fr)

- Work stealing / work-first principle
- Dynamics Macro-dataflow: partitioning (Metis, ...)
- Fault Tolerance (add/del resources)

FlowVR (flowvr.sf.net)

- Dedicated to interactive applications
- Static Macro-dataflow
- Parallel Code coupling

Video: [B Raffin, MOAIS & E Boyer, PERCEPTION]

http://www.id.imag.fr/~jlrosch/perso_html/talks/2007-01-MSRL-Berkeley/jecevr06.mov

[E Boyer, B Raffin 2006]

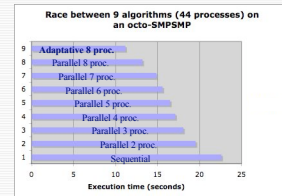
Thank you !

Back slides

43

The Prefix race: sequential/parallel fixed/ adaptive

44



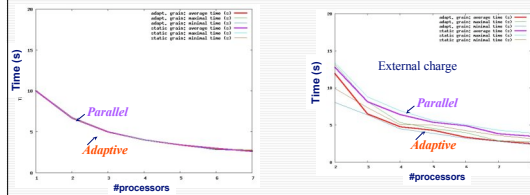
	Sequential	pa=2	pa=4	pa=6	pa=7	pa=8	Adaptatif
Minimum	21.83	18.16	15.89	14.99	13.02	12.51	8.76
Maximum	23.34	20.73	17.66	16.51	15.73	14.43	12.70
Moyenne	22.57	19.50	17.10	15.58	14.84	13.17	11.14
Mediane	22.58	19.64	17.38	15.57	14.63	13.11	11.01

On each of the 10 executions, adaptive completes first

Adaptive prefix : some experiments

45

Prefix of 10000 elements on a SMP 8 procs (IA64 / linux)



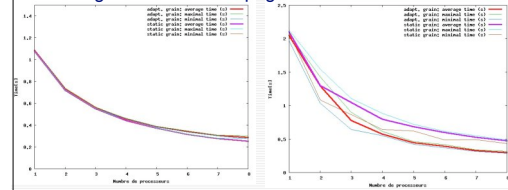
Single user context
Adaptive is equivalent to:
- sequential on 1 proc
- optimal parallel-2 proc. on 2 processors
- ...
- optimal parallel-8 proc. on 8 processors

Multi-user context
Adaptive is the fastest
15% benefit over a static grain algorithm

With * = double sum (r[i]=r[i-1] + x[i])

46

Finest "grain" limited to 1 page = 16384 octets = 2048 double

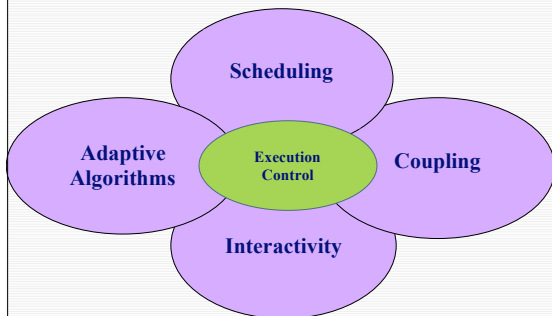


Single user

Processors with variable speeds

Remark for n=4,096,000 doubles :
- "pure" sequential : 0,20 s
- minimal "grain" = 100 doubles : 0.26s on 1 proc
and 0.175 on 2 procs (close to lower bound)

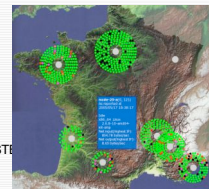
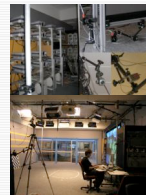
The Moais Group



Moais Platforms

46

- Cluster 2 :
 - 110 dual Itanium bi-processors with Myrinet network
- Grillage ("Grappe" and Image):
 - Camera Network
 - 54 processors (dual processor cluster)
 - Dual gigabits network
 - 16 projectors display wall
- Grids:
 - Regional: Ciment
 - National: Grid5000
 - Dedicated to CS experiments
- SMPs:
 - 8-way Itanium (Bull novascale)
 - 8-way dual-core Opteron + 2 GPUs
- MPSoCs
 - Collaborations with ST Microelectronics on ST...



Parallel Interactive App.



- Human in the loop
- Parallel machines (cluster) to enable large interactive applications
- Two main performance criteria:
 - Frequency (refresh rate)
 - Visualization: 30-60 Hz
 - Haptic: 1000 Hz
 - Latency (makespan for one iteration)
 - Object handling: 75 ms
- A classical programming approach: data-flow model
 - Application = static graph
 - Edges: FIFO connections for data transfert
 - Vertices: tasks consuming and producing data
 - Source vertices: sample input signal (cameras)
 - Sink vertices: output signal (projector)
- One challenge:
 - Good mapping and scheduling of tasks on processors

