

Parallel Algorithms

Design and Implementation

Jean-Louis.Roch at imag.fr

MOAIS / Lab. Informatique Grenoble, INRIA, France

1

Overview

- Machine model and work-stealing
- Work and depth
 - Fundamental theorem
 - Parallel divide & conquer
 - Examples
 - Accumulate
 - Monte Carlo simulations
 - Prefix/partial sum
- Work-stealing theorem
- Course 2: Work-first principle - Amortizing the overhead of parallelism
 - Sorting and merging
- Course 3: Amortizing the overhead of synchronization and communications
 - Numerical computations : FFT, matrix computations; Domain decompositions



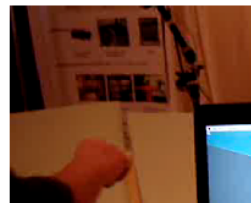
Interactive parallel computation?

Any application is "parallel":

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



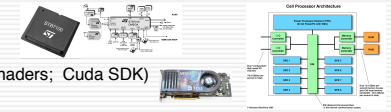
Interactive Distributed Simulation
3D-reconstruction
+ simulation
+ rendering
[B Raffin & E Boyer]
- 1 monitor
- 5 cameras,
- 6 PCs



New parallel supports from small too large

4

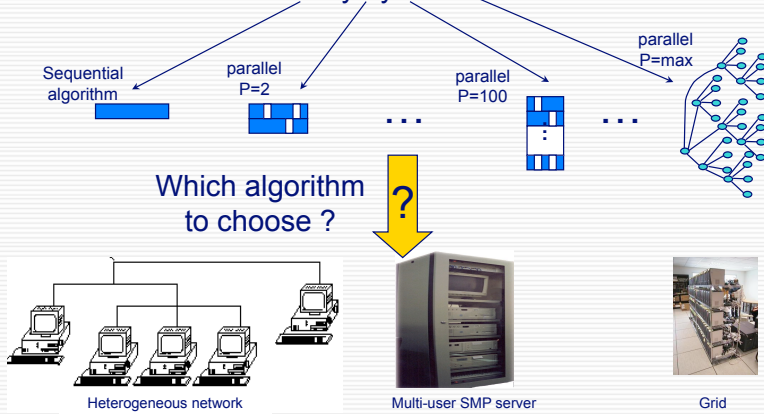
- **Parallel chips & multi-core architectures:**
 - MPSoCs (Multi-Processor Systems-on-Chips)
 - GPU : graphics processors (and programmable: Shaders; Cuda SDK)
 - Dual Core processors (Opterons, Itanium, etc.)
 - Heterogeneous multi-cores : CPUs + GPUs + DSPs+ FPGAs (Cell)
- **Commodity SMPs:**
 - 8 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)



Grimage platform

The problem

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

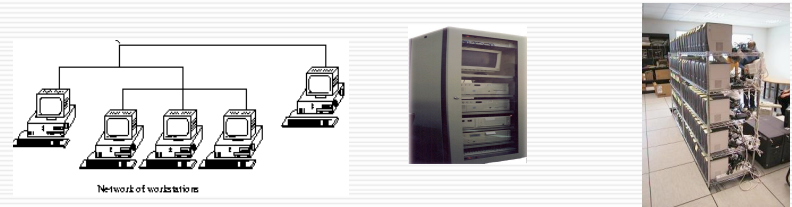


Which algorithm to choose ?

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

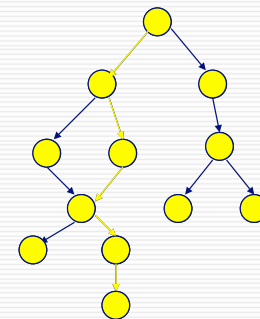
Heterogeneous processors, work and depth

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 #unit operations per second)
Assumption : $\text{Max}_{i,t} \{ \Pi_i(t) \} < \text{constant} \cdot \text{Min}_{i,t} \{ \Pi_i(t) \}$

Def: for a computation with duration T

- **total speed**: $\Pi_{tot} = \sum_{i=0, \dots, P} \sum_{t=0, \dots, T} \Pi_i(t)$
- **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:
[Graham69, Jaffe80, Bender-Rabin02]

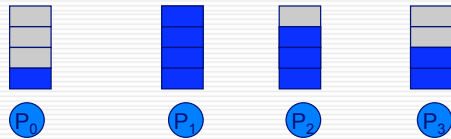
$$\text{makespan} \leq \frac{W}{p \cdot \Pi_{ave}} + \left(1 - \frac{1}{p}\right) \frac{D}{\Pi_{ave}}$$

The work stealing algorithm

9

- A distributed and randomized algorithm that computes a greedy schedule :

- Each processor manages a local task (depth-first execution)

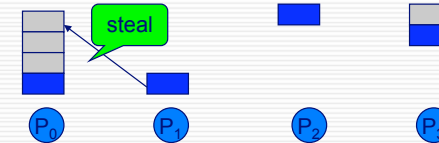


The work stealing algorithm

10

- 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,Blelloch, Blumofe02, BenderRabin02]

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

- Interest:**

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

Work stealing implementation

11

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

Difficult in general (coarse grain)

But easy if D is small [Work-stealing]

$$\text{Execution time} \leq \frac{W}{p \cdot \Pi_{ave}} + O\left(\frac{D}{\Pi_{ave}}\right) \quad (\text{fine grain})$$

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

12

- 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 = alloc(sizeof(*f));
5     f->sig = fib_sig;
6     if (sync) {
7         free(f, sizeof(*f));
8     }
9 }
10 else {
11     int x, y;
12     f->entry = 1;
13     f->n = n;
14     *f = f;
15     push(f);
16     x = fib (n-1);
17     if (sig(f) == FAILURE)
18         return 0;
19     ...
20     i;
21     free(f, sizeof(*f));
22     return (x+y);
23 }
24 }
```

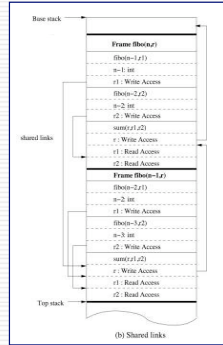
- won the 2006 award "Best Combination of Elegance and Performance" at HPC Challenge Class 2, SC'06, Tampa, Nov 14 2006 [Kuszmaul] on SGI ALTIX 3700 with 128 bi-Ithanium

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

- **Kaapi / Athapascan** <http://kaapi.gforge.inria.fr> : C++ library
 - `Fork<f>()` 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)

```

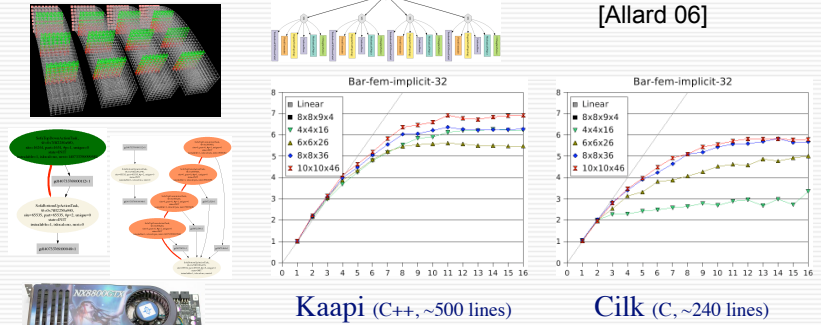
1 struct sum {
2   void operator()(Shared_r < int > a,
3                 Shared_r < int > b,
4                 Shared_w < int > r )
5   { r.write(a.read() + b.read()); }
6 };
7
8 struct fib {
9   void operator()(int n, Shared_w<int> r)
10  { if (n <2) r.write( n );
11    else
12      { int r1, r2;
13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



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

Experimental results on SOFA [CIMIT-ETZH-INRIA]

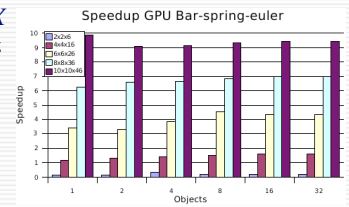
[Allard 06]



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



Algorithm design

Cascading divide & Conquer

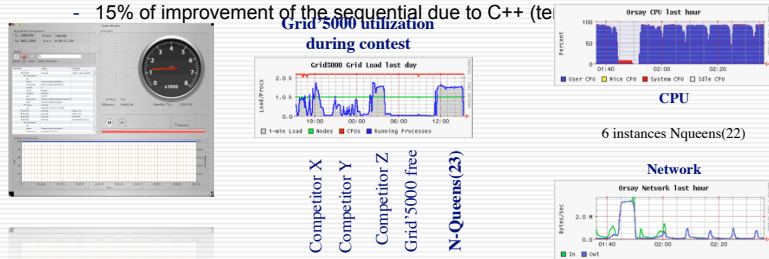
- $W(n) = a.W(n/K) + f(n)$
- $D(n) = D(n/K) + f(n)$
- $D(n) = D(\text{sqrt}(n)) + \log n$

Examples

- **Accumulate:**
 - Sequential
 - Parallel
- **Matrix-vector product – Matrix multiplication**
- **Triangular matrix inversion**
- **Maximum on CRCW**
- **Partial sum**

Example: Recursive and Monte Carlo computations

- X Besseron, T. Gautier, E Gobet, & G Huard won the nov. 2008 Plugtest-Grid&Work'08 contest – Financial mathematics application (Options pricing)
- In 2007, the team won the Nqueens contest; 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 [~24.10³³ solutions]
 - 0.625% idle time per processor
 - < 20s to deploy up to 1000 processes on 1000 machines [Taktuk, Huard]
 - 15% of improvement of the sequential due to C++ (te



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 $(\pi[0] = a[0], i = 1; i \leq n; i++) \pi[i] = \pi[i-1] * a[i];$ performs only n operations

- Fine grain optimal parallel algorithm :

