





CentraleSupélec 

SG6: High Performance Computing

## Distributed algorithms on multi-dimensional topologies

Stéphane Vialle

 universit  Paris-Saclay   Sciences et technologies de l'information et de la communication (STIC)

Stephane.Vialle@centralesupelec.fr  
<http://www.metz.supelec.fr/~vialle>

---

---

---

---

---

---

---

---

---

---

CentraleSupélec

## Distributed algorithms on multi-dimensional topologies

**Dense matrix product on a 1D ring**

- Parallelisation scheme
- Performance model

Dense matrix product on a 2D torus  
 Hyper-quicksort on a kD hypercube

---

---

---

---

---

---

---

---

---

---

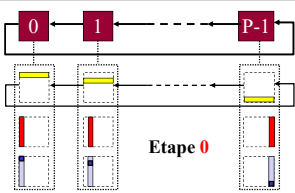
CentraleSupélec

Dense matrix product on a 1D ring

## Parallelisation scheme ( $C = A \times B$ )

**Data partitioning and circulation on a 1D ring:**

- A: partitioned into line blocks
- B, C: partitioned into column blocks
- Circulation of A slices
- Static B and C slices



Etape 0

Rig of  $P$  processes

Partitioning and circulation of A

Static partitioning of B

Static partitioning of C

**Notations:**

- Matrixes with:  $N = n \times n$  elements
- $P$  processes
- Matrixes slices of:  $\sqrt{N} \cdot \sqrt{N}/P = N/P$  elements

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 1D ring

## Parallelisation scheme ( $C = A \times B$ )

**At each of the P stages:**

- All processes are working (at any moment)
- In each process (and at each stage):
  - 1 - calculation of a small square block of C
  - 2 - circulation of a « horizontal » slice of A

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 1D ring

## Performance model

**Sequential computation of C matrix:**

A, B and C:  $\sqrt{N} \times \sqrt{N}$  elements

**$C_{ij}$  computation:**  
 $\sqrt{N}$  multiplications  
 and  $(\sqrt{N}-1)$  additions

$$T_{C_{ij} \text{ comput}}^{seq} = \sqrt{N} \cdot t_{mult} + (\sqrt{N} - 1) \cdot t_{add}$$

$$T_{C \text{ comput}}^{seq} = (\sqrt{N} \cdot \sqrt{N}) \cdot (\sqrt{N} \cdot t_{mult} + (\sqrt{N} - 1) \cdot t_{add})$$

➔  $T_{C \text{ comput}}^{seq} \approx N \cdot (2 \cdot \sqrt{N} \cdot t_{flop})$

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 1D ring

## Performance model

**Distributed computation of C matrix:**

**During one step:**

**Computation of a C block of  $\frac{\sqrt{N}}{p} \cdot \frac{\sqrt{N}}{p}$  elements**

- Time on one process:  $t_{1-comput}^{1-process} \approx \frac{\sqrt{N}}{p} \cdot \frac{\sqrt{N}}{p} \cdot (2 \cdot \sqrt{N} \cdot t_{flop})$
- Time on the P processes:

*hypothesis*  $\left\{ \begin{array}{l} \text{perfect load balancing} \\ \text{identical computing resources} \\ \text{perfect synchronization} \end{array} \right.$

➔  $t_{1-comput}^{P-process} = t_{1-comput}^{1-process}$

➔  $t_{1-comput}^{P-process} \approx \frac{N}{p^2} \cdot (2 \cdot \sqrt{N} \cdot t_{flop})$

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 1D ring

## Performance model

**Distributed computation of C matrix:**

**During one step:**

**Circulation of A slices of  $\sqrt{N} \cdot \frac{\sqrt{N}}{P}$  elements**

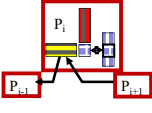
- Time of one communication:  $P_{i-1} \leftarrow P_i$

$$t_{1-comm} = t_s + \left(\sqrt{N} \cdot \frac{\sqrt{N}}{P}\right) \cdot t_w \quad \left\{ \begin{array}{l} t_s : \text{setup time or latency time} \\ t_w = \text{sizeof}(1 \text{ element})/B_w \end{array} \right.$$

- Time of all the communications across the P processes (during 1 step)

*hypothesis*  $\left\{ \begin{array}{l} \text{All communications are executed in parallel,} \\ \rightarrow \text{each process can send and receive in parallel...} \\ \dots \text{depends on the hardware \& on the MPI comm routine!} \\ \rightarrow \text{all communications cross the same « network »...} \\ \dots \text{depends on the deployment!} \end{array} \right.$

$\Rightarrow t_{1-circulation}^{P-processes} = t_{1-comm} \approx \frac{N}{P} \cdot t_w$




---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 1D ring

## Performance model

**Distributed computation of C matrix:**

**During one step:**

$$t_{1-step}^{P-process} = t_{1-comput}^{P-process} + t_{1-circulation}^{P-processes}$$

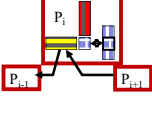
$$t_{1-step}^{P-process} \approx \frac{N}{P^2} \cdot (2 \cdot \sqrt{N} \cdot t_{flop}) + \frac{N}{P} \cdot t_w$$

**During the P steps of C distributed computation:**

$$T_{C\ comput}^{distributed} = P \cdot t_{1-step}^{P-process}$$

$\Rightarrow T_{C\ comput}^{distributed} \approx \frac{N}{P} \cdot (2 \cdot \sqrt{N} \cdot t_{flop}) + N \cdot t_w$  with:  $P > 1$

*hypothesis*  $\left\{ \begin{array}{l} \text{All computations identical \& executed in parallel} \\ \text{All communications identical \& executed in parallel} \end{array} \right.$




---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 1D ring

## Performance model

**Distributed computation of C matrix:**

**Execution times:**

$$T_{C\ comput}^{seq} \approx N \cdot (2 \cdot \sqrt{N} \cdot t_{flop})$$

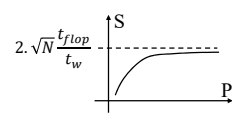
$$T_{C\ comput}^{distributed} \approx \frac{N}{P} \cdot (2 \cdot \sqrt{N} \cdot t_{flop}) + N \cdot t_w \quad \text{with: } P > 1$$

**Speedup:**

$$S(N,P) = \frac{T_{C\ comput}^{seq}}{T_{C\ comput}^{distributed}} = \frac{N \cdot (2 \cdot \sqrt{N} \cdot t_{flop})}{\frac{N}{P} \cdot (2 \cdot \sqrt{N} \cdot t_{flop}) + N \cdot t_w}$$

$\Rightarrow S(N,P) = P \cdot \frac{1}{1 + \frac{P}{2 \cdot \sqrt{N}} \cdot \frac{t_w}{t_{flop}}}$

With  $N = N_0$




---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 1D ring

**Performance model**

**Distributed computation of C matrix:**

$$S(N, P) = P \cdot \frac{1}{1 + \frac{P}{2\sqrt{N}} \cdot \frac{t_w}{t_{flop}}}$$

- When data size is fixed:  $N = N_0$

- When  $P = P_0$  is fixed, and data size increases:  
 $\lim_{N \rightarrow \infty} (S(N, P_0)) = P_0$       Speedup increases and tends to  $P_0$  (ideal  $S$ )

In fact:  $T_C^{distributed} \approx \frac{N}{P_0} \cdot \underbrace{(2 \cdot \sqrt{N} \cdot t_{flop})}_{O(N^{3/2})} + \underbrace{N \cdot t_w}_{O(N)}$

**When increasing the pb size, on a fixed number of nodes:**  
 Computations increase faster than communications  
 → Speedup will increase! .... **Nice problem!**

---

---

---

---

---

---

---

---

---

---

**Distributed algorithms on multi-dimensional topologies**

Dense matrix product on a 1D ring  
**Dense matrix product on a 2D torus**

- Parallelisation scheme
- Performance model

Hyper-quicksort on a kD hypercube

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

**Parallelisation scheme ( $C = A \times B$ )**

A, B, C :  $n \times n = N$  elements       $C = A \times B$        $c_{ij} = \sum_{k=1}^n (a_{ik} \cdot b_{kj})$

$\sqrt{N} \times \sqrt{N}$  elements       $\sqrt{P} \times \sqrt{P}$  processes

→ Which data partitioning & circulation ?

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallelisation scheme ( $C = A \times B$ )

**Initial partitioning:**

$\sqrt{N} \times \sqrt{N}$  elements  
 $\sqrt{P} \times \sqrt{P}$  blocks of  $\frac{\sqrt{N}}{\sqrt{P}} \times \frac{\sqrt{N}}{\sqrt{P}}$  elements  
 $\sqrt{P} \times \sqrt{P}$  processes

**Only one partitioning for A, B and C**  
 → Result matrix C can be reused as a left or right operand

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallelisation scheme ( $C = A \times B$ )

**Data circulation principle:**

A blocks: require a **circulation on the row** of processes ←

B blocks: require a **circulation on the column** of processes ↑

$$C_{i,j} = \sum_{k=1}^{\sqrt{P}} (A_{i,k} \times B_{k,j})$$
 → Each process will receive all data it needs to compute its C block

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallelisation scheme ( $C = A \times B$ )

**But process  $P_{i,j}$  requires **simultaneously** the blocks  $A_{i,k}$  and  $B_{k,j}$ :**

With the initial partitioning only process  $P_{i,i}$  can start their computations

$$C_{i,j} = \sum_{k=1}^{\sqrt{P}} (A_{i,k} \times B_{k,j})$$
 → We need to improve the initial partitioning

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallellisation scheme ( $C = A \times B$ )

**Approach:**

- We keep the 2D partitioning principle
- We keep the circulation scheme  
(A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
  - we install the first blocks required by  $P(0;0)$
  - and we propagate the constraints...

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallellisation scheme ( $C = A \times B$ )

**Approach:**

- We keep the 2D partitioning principle
- We keep the circulation scheme  
(A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
  - we install the first blocks required by  $P(0;0)$
  - the circulation scheme impose:
    - $A_{0,j}$  on the first line
    - $B_{i,0}$  on the first column

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallellisation scheme ( $C = A \times B$ )

**Approach:**

- We keep the 2D partitioning principle
- We keep the circulation scheme  
(A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
  - we install the first blocks required by  $P(0;0)$
  - the circulation scheme impose
    - $A_{0,j}$  on the first line
    - $B_{i,0}$  on the first column
  - we install A & B blocks required on the 1st col and line

Compatible  $A_{iK}$  and  $B_{KJ}$  blocks

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallellisation scheme ( $C = A \times B$ )

**Approach:**

- We keep the 2D partitioning principle
- We keep the circulation scheme  
(A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
  - we install the first blocks required by  $P(0;0)$
  - the circulation scheme impose  $A_{0,j}$  on the first line  
 $B_{1,0}$  on the first column
  - we install A & B blocks required on the 1st col and line
  - we install others A & B blocks, guided by circulation order**

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallellisation scheme ( $C = A \times B$ )

**Approach:**

- We keep the 2D partitioning principle
- We keep the circulation scheme  
(A blocks circulate on process rows, B blocks circulate on process columns)
- We get an initial partitioning where all  $A_{iK}$  and  $B_{Kj}$  are compatible**
- Couples of blocks remain compatible after each circulation step**

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Parallellisation scheme ( $C = A \times B$ )

**Approach:**

- We keep the 2D partitioning principle
- We keep the circulation scheme  
(A blocks circulate on process rows, B blocks circulate on process columns)
- We get an initial partitioning where all  $A_{iK}$  and  $B_{Kj}$  are compatible**
- Couples of blocks remain compatible after each circulation step**

**Finally:**

- Row  $I$  has been pre-shifted of  $I$  steps to the left**
- Column  $J$  has been pre-shifted of  $J$  steps to the top**

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

## Parallelisation scheme ( $C = A \times B$ )

**Canon algorithm:  $C = A \times B$**   
 Product of dense matrixes of  $\sqrt{N} \times \sqrt{N}$  elements  
 on a 2D torus of  $\sqrt{P} \times \sqrt{P}$  processes

---

1 : Matrix partitioning: block  $I, J$   
 on process  $I, J$

2 : Pre-shifting of the initial blocks:

- A row  $I (\in [0; \sqrt{P}-1])$  shifted of  $I$  step to the left  $\leftarrow$
- B column  $J (\in [0; \sqrt{P}-1])$  shifted of  $J$  step to the top  $\uparrow$

3 :  $\sqrt{P}$  steps: { local product  $A_{IK} \times B_{KJ}$  (partial sum of  $C_{IJ}$ );  
 1 notch shifting of A & B blocks:  $A_{IK} \leftarrow, B_{KJ} \uparrow$  }

4 : Post-shifting of A and B blocks to restore input matrix partitioning

---

---

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

## Parallelisation scheme ( $C = A \times B$ )

**Canon algorithm:  $C = A \times B$**   
 Product of dense matrixes of  $\sqrt{N} \times \sqrt{N}$  elements  
 on a 2D torus of  $\sqrt{P} \times \sqrt{P}$  processes

---

1 : Matrix partitioning: block  $I, J$   
 on process  $I, J$

2 : Pre-shifting of the initial blocks:

- A row  $I (\in [0; \sqrt{P}-1])$  shifted of  $I$  notches to the left  $\leftarrow$
- B column  $J (\in [0; \sqrt{P}-1])$  shifted of  $J$  notches to the top  $\uparrow$

3 :  $(\sqrt{P} - 1)$  steps: { local product  $A_{IK} \times B_{KJ}$  (partial sum of  $C_{IJ}$ );  
 one notch shifting of A & B blocks:  $A_{IK} \leftarrow, B_{KJ} \uparrow$  }

4 : 1 step: { local product  $A_{IK} \times B_{KJ}$  (partial sum of  $C_{IJ}$ );  
 post-shifting of A and B blocks to restore input matrix partitioning }

---

---

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

## Performance model ( $C = A \times B$ )

1 - Implementation with blocking communications (no overlapping)

**Hypothesis:**

- 1 - All computations are load balanced
- 2 - All computing units are identical
- 3 - All communications of one matrix shift are achieved in parallel (only one send and one recv per process: sustainable by the switch)

and then:

**No more communications at a time than with a ring of processes**  
 $\rightarrow$  Sustainable by the switch

4 - All communications cross the same « network » (depends on the deployment !)

---

---

---

---

---

---

---

---

---

---

---

---



Dense matrix product on a 2D torus

### Performance model (C = A×B)

1 – Implementation with blocking communications (no overlapping)

**Pre- and post-shifting implementation on a cluster:**  
 Each block can be moved with only one send and one recv  
 (no need to concatenate several elementary shifts on a cluster network)

3 parallel send-&-recv  $\equiv$  2 notch shifting

$\Rightarrow t_{pre-shifting-A} = t_s + \left(\frac{\sqrt{N}}{\sqrt{P}} \times \frac{\sqrt{N}}{\sqrt{P}}\right) \cdot t_w \approx \frac{N}{P} \cdot t_w$

Considering blocking communications:  
 $\Rightarrow t_{pre-shifting-A-B} \approx 2 \cdot \frac{N}{P} \cdot t_w$       **Same result for post-shifting**

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Performance model (C = A×B)

1 – Implementation with blocking communications (no overlapping)

- Pre-shifting (A & B):  $2 \cdot \frac{N}{P} \cdot t_w$
- 1 step computations:  $\frac{\sqrt{N}}{\sqrt{P}} \cdot \frac{\sqrt{N}}{\sqrt{P}} \cdot (2 \cdot \frac{\sqrt{N}}{\sqrt{P}} \cdot t_{flop}) = 2 \cdot \left(\frac{N}{P}\right)^{\frac{3}{2}} \cdot t_{flop}$   
 (P processes in parallel)
- 1 step circulation:  $2 \cdot \left(\frac{\sqrt{N}}{\sqrt{P}} \cdot \frac{\sqrt{N}}{\sqrt{P}}\right) \cdot t_w = 2 \cdot \frac{N}{P} \cdot t_w$   
 (blocking communications)
- Post-shifting (A & B):  $2 \cdot \frac{N}{P} \cdot t_w$

**Complete distributed computation:**

$$T_{C\ comput}^{distributed} \approx t_{pre-shifting} + (\sqrt{P}-1) \cdot (t_{1-step-comput} + t_{1-step-circulation}) + (t_{1-step-comput} + t_{post-shifting})$$


---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

### Performance model (C = A×B)

1 – Implementation with blocking communications (no overlapping)

- Pre-shifting (A & B):  $2 \cdot \frac{N}{P} \cdot t_w$
- 1 step computations:  $\frac{\sqrt{N}}{\sqrt{P}} \cdot \frac{\sqrt{N}}{\sqrt{P}} \cdot (2 \cdot \frac{\sqrt{N}}{\sqrt{P}} \cdot t_{flop}) = 2 \cdot \left(\frac{N}{P}\right)^{\frac{3}{2}} \cdot t_{flop}$   
 (P processes in parallel)
- 1 step circulation:  $2 \cdot \left(\frac{\sqrt{N}}{\sqrt{P}} \cdot \frac{\sqrt{N}}{\sqrt{P}}\right) \cdot t_w = 2 \cdot \frac{N}{P} \cdot t_w$   
 (blocking communications)
- Post-shifting (A & B):  $2 \cdot \frac{N}{P} \cdot t_w$

**Complete distributed computation:**

$$T_{C\ comput}^{distributed} \approx 2 \cdot \frac{N}{P} \cdot t_w + (\sqrt{P}-1) \cdot \left(2 \cdot \left(\frac{N}{P}\right)^{\frac{3}{2}} \cdot t_{flop} + 2 \cdot \frac{N}{P} \cdot t_w\right) + \left(2 \cdot \left(\frac{N}{P}\right)^{\frac{3}{2}} \cdot t_{flop} + 2 \cdot \frac{N}{P} \cdot t_w\right)$$


---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

**Performance model (C = A×B)**

**1 – Implementation with blocking communications (no overlapping)**

- Pre-shifting (A & B):  $2 \cdot \frac{N}{P} t_w$
- 1 step computations:  $\frac{\sqrt{N}}{\sqrt{P}} \cdot \frac{\sqrt{N}}{\sqrt{P}} \cdot (2 \cdot \frac{\sqrt{N}}{\sqrt{P}} \cdot t_{flop}) = 2 \cdot \left(\frac{N}{P}\right)^{\frac{3}{2}} \cdot t_{flop}$   
(P processes in parallel)
- 1 step circulation:  $2 \cdot \left(\frac{\sqrt{N}}{\sqrt{P}} \cdot \frac{\sqrt{N}}{\sqrt{P}}\right) \cdot t_w = 2 \cdot \frac{N}{P} t_w$   
(blocking communications)
- Post-shifting (A & B):  $2 \cdot \frac{N}{P} t_w$

**Complete distributed computation:**

$T_{C\text{ comput}}^{\text{distributed}} \approx 2 \cdot \frac{N^{3/2}}{P} t_{flop} + 2 \cdot (\sqrt{P}+1) \cdot \frac{N}{P} t_w$       **with: P > 1**

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

**Performance model (C = A×B)**

**1 – Implementation with blocking communications (no overlapping)**

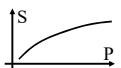
**Complete distributed computation:**

$T_{C\text{ comput}}^{\text{distributed}} \approx 2 \cdot \frac{N^{3/2}}{P} t_{flop} + 2 \cdot (\sqrt{P}+1) \cdot \frac{N}{P} t_w$

**Speedup:**

$S(N, P) = \frac{T_{C\text{ comput}}^{\text{Seq}}}{T_{C\text{ comput}}^{\text{distributed}}} \approx \frac{2 \cdot N^{3/2} \cdot t_{flop}}{2 \cdot \frac{N^{3/2}}{P} t_{flop} + 2 \cdot (\sqrt{P}+1) \cdot \frac{N}{P} t_w}$

$S(N, P) \approx P \cdot \frac{1}{1 + \frac{\sqrt{P}+1}{\sqrt{N}} \cdot \frac{t_w}{t_{flop}}}$       With  $N = N_0$



Again:  $\lim_{N \rightarrow \infty} S(N, P_0) = P_0$       When pb size increase, speedup increases and tends to ideal speedup

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

**Performance model (C = A×B)**

**2 – Implementation with non-blocking communications**

**If we assume the switch can achieve A and B communications in parallel:**

- Pre-shifting (A & B):  $2 \cdot \frac{N}{P} t_w \rightarrow \frac{N}{P} t_w$
- 1 step computations:  $2 \cdot \left(\frac{N}{P}\right)^{\frac{3}{2}} \cdot t_{flop}$  (unchanged)  
(P processes in parallel)
- 1 step circulation:  $2 \cdot \frac{N}{P} t_w \rightarrow \frac{N}{P} t_w$   
(blocking communications)
- Post-shifting (A & B):  $2 \cdot \frac{N}{P} t_w \rightarrow \frac{N}{P} t_w$

**Complete distributed computation:**

$T_{C\text{ comput}}^{\text{distributed}} \approx 2 \cdot \frac{N^{3/2}}{P} t_{flop} + 2 \cdot (\sqrt{P}+1) \cdot \frac{N}{P} t_w$   
 $\rightarrow 2 \cdot \frac{N^{3/2}}{P} t_{flop} + (\sqrt{P}+1) \cdot \frac{N}{P} t_w$

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

**Performance model (C = A×B)**

2 – Implementation with **non-blocking** communications

**Complete distributed computation:**

$$T_{C\text{ comput}}^{\text{distributed}} \approx 2 \cdot \frac{N^{3/2}}{P} t_{flop} + (\sqrt{P}+1) \cdot \frac{N}{P} t_w$$

**Speedup:**

$$S(N, P) = \frac{T_{C\text{ comput}}^{\text{Seq}}}{T_{C\text{ comput}}^{\text{distributed}}} \approx \frac{2N^{3/2} \cdot t_{flop}}{2 \cdot \frac{N^{3/2}}{P} t_{flop} + (\sqrt{P}+1) \cdot \frac{N}{P} t_w}$$

$$S(N, P) \approx P \cdot \frac{1}{1 + \frac{\sqrt{P}+1}{2\sqrt{N}} \times \frac{t_w}{t_{flop}}} \quad \text{Speedup increases faster}$$

Again:  $\lim_{N \rightarrow \infty} S(N, P_0) = P_0$       When pb size increase, speedup increases and tends to ideal speedup

---

---

---

---

---

---

---

---

---

---

Dense matrix product on a 2D torus

**Performance model (C = A×B)**

**Comparison 1D Ring / 2D Torus solutions:**

	1D ring	2D torus
Blocking comm	$T_{par}(P) \approx 2 \cdot \frac{N^{3/2}}{P} \cdot t_{flop} + N \cdot t_w$	$T_{par}(P) \approx 2 \cdot \frac{N^{3/2}}{P} \cdot t_{flop} + 2 \cdot (\sqrt{P} + 1) \frac{N}{P} \cdot t_w$
Non-blocking: A & B circ. in parallel		$T_{par}(P) \approx 2 \cdot \frac{N^{3/2}}{P} \cdot t_{flop} + (\sqrt{P} + 1) \frac{N}{P} \cdot t_w$

2D torus solution *should be* faster when

- $P \geq 3$  with blocking comm  $\rightarrow P \geq 2^2$
- $P \geq 2$  with non-blocking comm  $\rightarrow P \geq 2^2$

!! the assumptions on the network switch are right...

---

---

---

---

---

---

---

---

---

---

**Distributed algorithms on multi-dimensional topologies**

Dense matrix product on a 1D ring

Dense matrix product on a 2D torus

**Hyper-quicksort on a kD hypercube**

- Parallelisation scheme

---

---

---

---

---

---

---

---

---

---

CentralesSupélec

Hyper-quicksort on a kD hypercube

## Principe séquentiel

**Principe :**

- Tri récursif
- Diviser-pour-régner
- Séparation des données selon des valeurs *pivots*

**Performances :**

pire cas :  $\left(\frac{N \cdot (N-1)}{2}\right) t_{comp} \rightarrow O(N^2)$

moyen cas (Knuth) :  $\rightarrow O(2 \cdot N \cdot \log(N))$

meilleur cas :  $\left(N \cdot \log_2(N) - 2 \cdot N + 1\right) t_{comp} \rightarrow O(N \cdot \log(N))$

---

---

---

---

---

---

---

---

CentralesSupélec

Hyper-quicksort on a kD hypercube

## Parallélisation naïve

**Principe :**

On crée un nouveau processus et on utilise un nouveau processeur à chaque appel récursif (en réutilisant les anciens).

Simple **mais inefficace !**

---

---

---

---

---

---

---

---

CentralesSupélec

Hyper-quicksort on a kD hypercube

## Parallélisation naïve

**1 - Faiblesse conceptuelle :**  
Algorithme pas pleinement parallèle  
→ Trouver plus parallèle !

**2 - Faiblesse de mise en œuvre :**  
Ressources (CPU) partagées :  
→ Temps de création dynamique et de re-répartition de processus

Ressources (CPU) allouées au lancement de l'application :  
→ Gaspillage !

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

### Rappel des propriétés des Hyper-Cubes

**Construction récursive :**

Dim 1   Dim 2   Dim 3   Dim 4

**Numérotation récursive binaire :**

Dim 1   Dim 2   Dim 3

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

### Rappel des propriétés des Hyper-Cubes

**Décomposition en sous-hypercube :**

1 hypercube de dimension n coupé par un hyper-plan donne 2 sous-hypercubes de dimension n-1 :

1 x dim 3   2 x dim 2   ou   2 x dim 2   ou   2 x dim 2

→ Les algorithmes de routages restent les mêmes dans toutes les parties de l'hyper-cube

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

### Rappel des propriétés des Hyper-Cubes

**Distance entre nœuds :**

Le nombre de bits différents dans les adresses de deux nœuds donne leur distance (si la numérotation a suivi la construction récursive)

Ex :   A : 000, B : 011  
 $d(M,N) = 2$

Ex :   M : 110, N : 001  
 $d(M,N) = 3$

→ Les algorithmes de routages seront à base de calculs simples (et rapides)

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

## Principe de parallélisation sur Hyper-Cube

**Objectif :**

- Tous les processeurs doivent travailler tout le temps !
- S'inspirer d'un tri séquentiel rapide ( $O(N \cdot \log(N))$ ) : quick-sort

→ Trouver un schéma de parallélisation avec comm. optimisées (peu de comm, ou comm locales)

**Idée :** Quick-sort : algo récursif ↔ topologie récursive : Hyper-cube

---

---

---

---

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

## Principe de parallélisation sur Hyper-Cube

**Init :** Chaque nœud charge N/P données en local

**Etape 0 :**

- Choix de  $2^0$  pivots pour les  $2^0$  hypercubes de dimension d-0
- Séparation selon le pivot sur chaque nœud de l'hypercube
- Echange des listes inférieures et supérieures en dimension d-0
- Fusions locales des listes conservées et des listes reçues

Communication avec voisin en dimension d-0 : seul le bit d-0 diffère

---

---

---

---

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

## Principe de parallélisation sur Hyper-Cube

**Détails de l'étape 0 sur PE-000 et PE-100 :**

Séparations locales

Echanges de sous-listes

Fusion des sous-listes conservées et reçues

---

---

---

---

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

### Principe de parallélisation sur Hyper-Cube

**Etape 1 :**

- Choix de  $2^1$  pivots pour les  $2^1$  hypercubes de dimension d-1
- Séparation selon un pivot sur chaque nœud des  $2^1$  hypercubes
- Echange des listes inférieures et supérieures en dimension d-1
- Fusions locales des listes conservées et des listes reçues

Communication avec voisin en dimension d-1 : seul le bit d-1 diffère

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

### Principe de parallélisation sur Hyper-Cube

**Etape 2 :**

- Choix de  $2^2$  pivots pour les  $2^2$  hypercubes de dimension d-2
- Séparation selon un pivot sur chaque nœud des  $2^2$  hypercubes
- Echange des listes inférieures et supérieures en dimension d-2
- Fusions locales des listes conservées et des listes reçues

Communication avec voisin en dimension d-2 : seul le bit d-2 diffère

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

### Principe de parallélisation sur Hyper-Cube

**Etape finale :**

Chaque processeur tri en local sa liste finale  
 → ex : quick-sort local et séquentiel sur chaque processeur

---

---

---

---

---

---

---

---

CentralesSupélec

Hyper-quicksort on a kD hypercube

## 1ère implantation sur Hyper-Cube

```

HcubeQuickSort(double *A, int d)
{
  int i;
  double *Ainf, *Asup, *Abuf; /* Compteur de dimension. */
  double x; /* Tables de données. */
  /* Pivots. */
  for (i = d-1; i >= 0; i--) { /* Pour chaque dimension: */
    x = choix_pivot(me,i); /* - Partitionnement local*/
    partitioner(A,x,Ainf,Asup);
    if (me & exp2(i) == 0) { /* - Communications */
      asend(Asup,me | exp2(i));
      recv(Abuf,me | exp2(i));
      union(Ainf,Abuf,&A);
    } else {
      asend(Ainf,me & ~exp2(i));
      recv(Abuf,me & ~exp2(i));
      union(Abuf,Asup,&A);
    }
  }
  QuickSortSequentiel(A); /* Tri final des donnes */
} /* locales */

```

---

---

---

---

---

---

---

---

---

---

CentralesSupélec

Hyper-quicksort on a kD hypercube

## 1ère implantation sur Hyper-Cube

Faiblesse de ce premier hyper-quicksort :

**Si mauvais choix de pivot : déséquilibre définitif !**  
 → processeurs surchargés et processeurs à vide

→ Soigner le choix du pivot.

---

---

---

---

---

---

---

---

---

---

CentralesSupélec

Hyper-quicksort on a kD hypercube

## Parallélisation optimisée sur H-Cube

Choix de pivot optimisé :

Pivot idéal : élément médian

En séquentiel : on approxime l'élément médian, ex :

- l'élément médian des 5 premiers,
- l'élément médian d'un échantillonnage, ....

En parallèle :

- on trie initialement les éléments locaux  
 → on prend l'élément médian ! LocalTab[N/P/2]
- on suppose une distribution homogène sur tous les PEs  
 → le pivot idéal d'un PE est un très bon pivot pour tous !

---

---

---

---

---

---

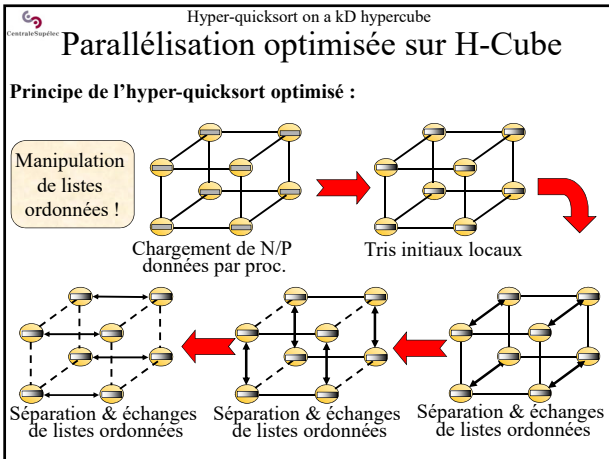
---

---

---

---






---

---

---

---

---

---

---

---

---

---

Hyper-quicksort on a kD hypercube

## 2<sup>ème</sup> implantation sur Hyper-Cube

**Implantation de l'hyper-quicksort optimisé :**

```

HcubeQuickSort(double *A, int d)
{
  int i;
  double *Ainf, *Asup, *Abuf; /* Compteur de dimension. */
  double x; /* Tables de données. */
  QuickSortSequentiel(A); /* Pivot. */
  /* Tri initial des données */
  /* locales. */
  for (i = d-1; i >= 0; i--) { /* Pour chaque dimension: */
    x = choix_pivot(me,i); /* - Partitionnement local*/
    partitioner(A,x,Ainf,Asup);
    if (me & exp2(i) == 0) { /* - Communications */
      asend(Asup,me | exp2(i));
      recv(Abuf,me | exp2(i));
      union_ordonne(Ainf,Abuf,&A);
    } else {
      asend(Ainf,me & ~exp2(i));
      recv(Abuf,me & ~exp2(i));
      union_ordonne(Abuf,Asup,&A);
    }
  }
}
  
```

---

---

---

---

---

---

---

---

---

---

Distributed algorithms on multi-dimensional topologies

End

---

---

---

---

---

---

---

---

---

---