



SG6: High Performance Computing

Distributed algorithms on multi-dimensional topologies

Stéphane Vialle



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

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

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

2 - Step circulation:
A slice: $\frac{\sqrt{N}}{P} \sqrt{N}$ elements

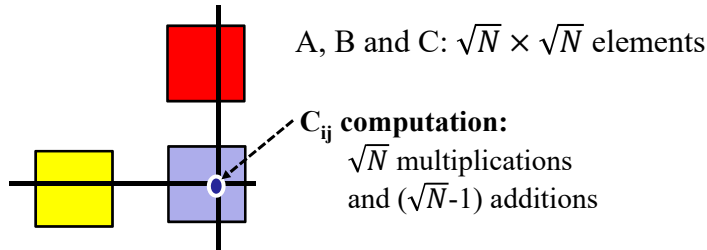
1 - Step computation:
C block: $\frac{\sqrt{N}}{P} \cdot \frac{\sqrt{N}}{P}$ elements

Static partition
B slice: $\sqrt{N} \cdot \frac{\sqrt{N}}{P}$ elements

P_{i-1} P_i P_{i+1}

Performance model

Sequential computation of C matrix:



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

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

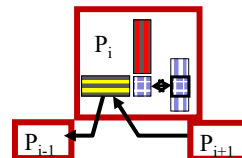
➔ $T_{C \text{ comput}}^{seq} \approx N \cdot (2 \cdot \sqrt{N} \cdot t_{flop})$

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

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

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

Performance model

Distributed computation of C matrix:

During one step:

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

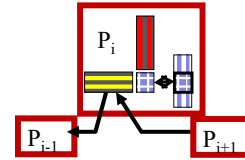
- Time of one communication: $P_{i-1} \leftarrow P_i$

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

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

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

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



Performance model

Distributed computation of C matrix:

During one step:

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

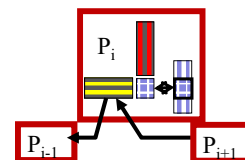
$$t_{1-step}^{P-process} \approx \frac{N}{P^2} \cdot (2 \cdot \sqrt{N} \cdot t_{flop}) + \frac{N}{P} \cdot t_w$$

During the P steps of C distributed computation:

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

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

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



Performance model

Distributed computation of C matrix:

Execution times:

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

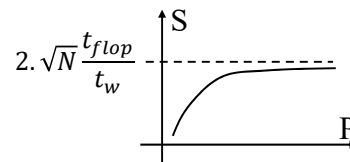
$$T_{C\text{ comput}}^{distributed} \approx \frac{N}{P} \cdot (2 \cdot \sqrt{N} \cdot t_{flop}) + N \cdot t_w \quad \text{with: } P > 1$$

Speedup:

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

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

With $N = N_0$

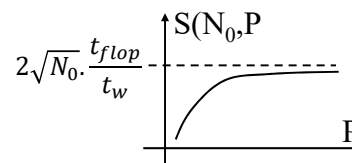


Performance model

Distributed computation of C matrix:

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

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



- When $P = P_0$ is fixed, and data size increases:

$$\lim_{N \rightarrow \infty} (S(N, P_0)) = P_0 \quad \text{Speedup increases and tends to } P_0 \text{ (ideal } S)$$

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

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

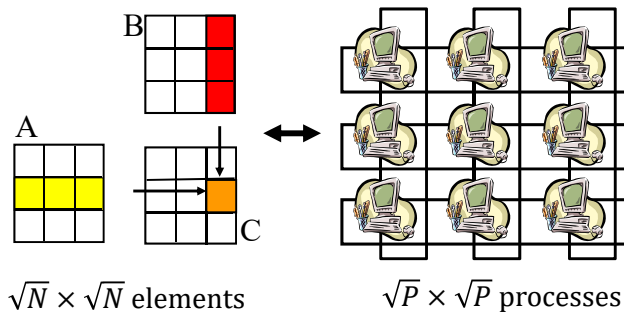
Dense matrix product on a 2D torus

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

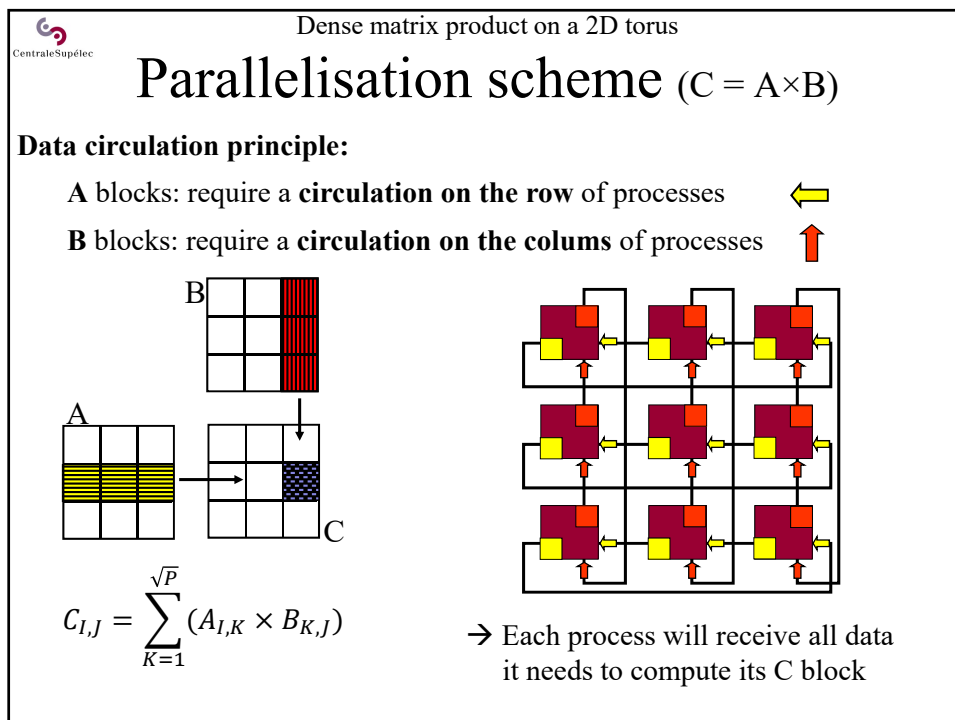
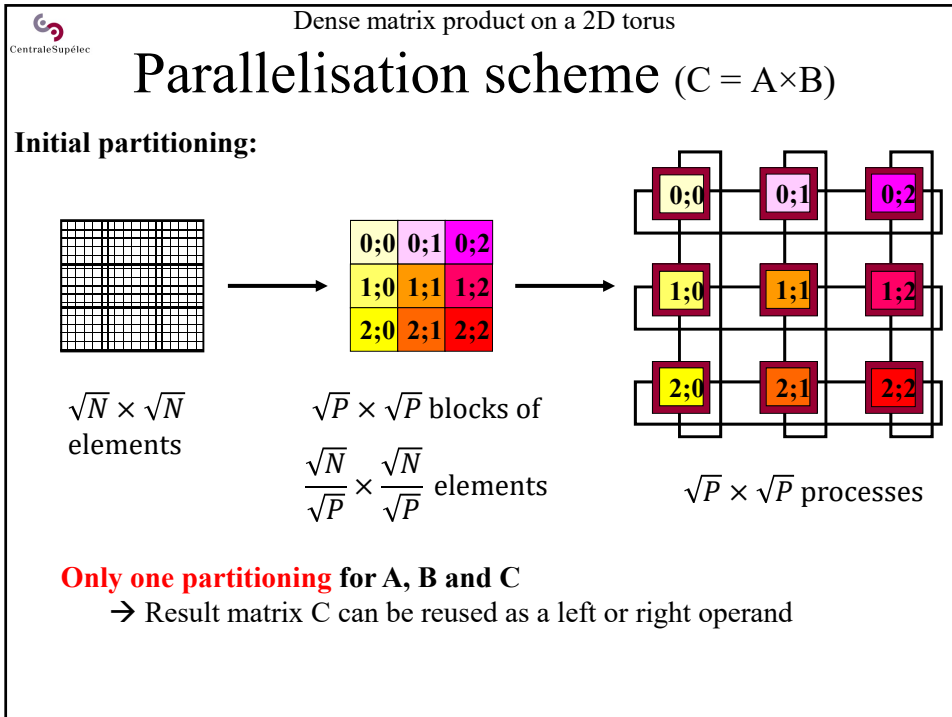
A, B, C : $n \times n = N$ elements

$$C = A \times B$$

$$c_{ij} = \sum_{k=1}^n (a_{ik} \cdot b_{kj})$$

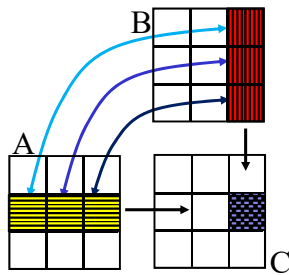


→ Which data partitioning & circulation ?



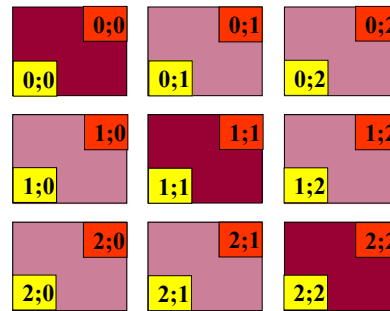
Parallelisation scheme ($C = A \times B$)

But process $P_{I,J}$ requires **simultaneously** the blocks $A_{I,K}$ and $B_{K,J}$:



$$C_{I,J} = \sum_{K=1}^{\sqrt{P}} (A_{I,K} \times B_{K,J})$$

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

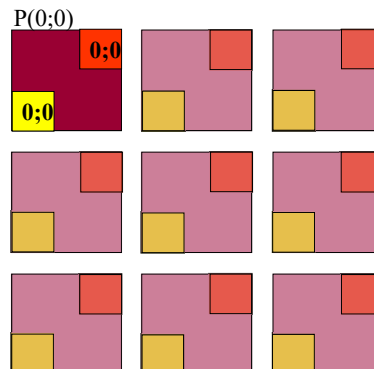


→ We need to improve the initial partitioning

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

Approach:

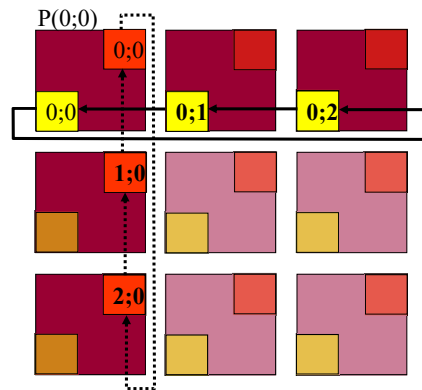
- We keep the 2D partitioning principle
- We keep the circulation scheme (A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
 - we install the first blocks required by $P(0;0)$
 - and we **propagate the constraints...**



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

Approach:

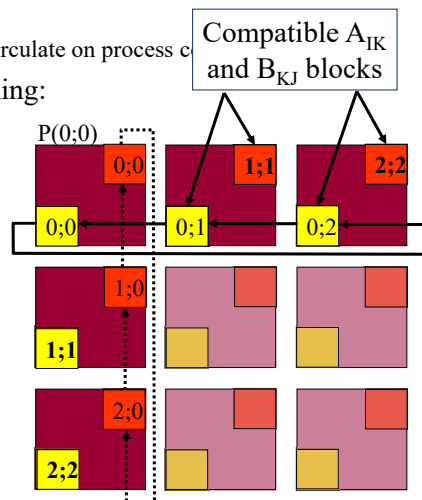
- We keep the 2D partitioning principle
- We keep the circulation scheme
(A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
 - we install the first blocks required by $P(0;0)$
 - **the circulation scheme impose:**
 $A_{0,j}$ on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



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

Approach:

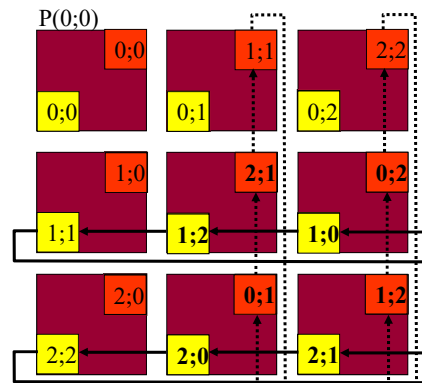
- We keep the 2D partitioning principle
- We keep the circulation scheme
(A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
 - we install the first blocks required by $P(0;0)$
 - the circulation scheme impose
 $A_{0,j}$ on the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



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

Approach:

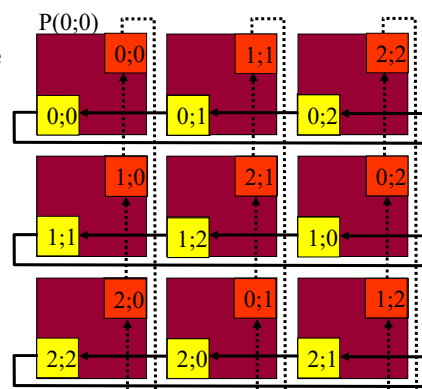
- We keep the 2D partitioning principle
- We keep the circulation scheme
(A blocks circulate on process rows, B blocks circulate on process columns)
- But we look for a better initial partitioning:
 - we install the first blocks required by $P(0;0)$
 - the circulation scheme impose $A_{0,j}$ on the first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



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

Approach:

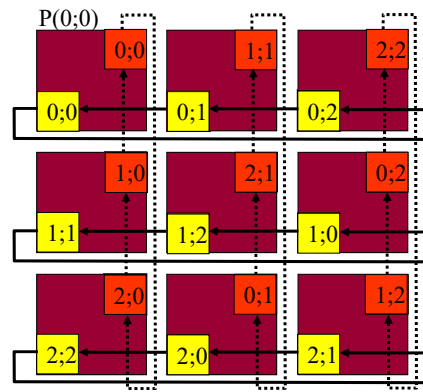
- We keep the 2D partitioning principle
- We keep the circulation scheme
(A blocks circulate on process rows, B blocks circulate on process columns)
- **We get an initial partitioning where all A_{IK} and B_{KJ} are compatible**
- **Couples of blocks remain compatible after each circulation step**



Parallelisation 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:

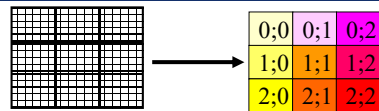
- Row I has been **pre-shifted** of I steps to the left
- Column J has been **pre-shifted** of J steps to the top

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

Canon algorithm: $C = A \times B$

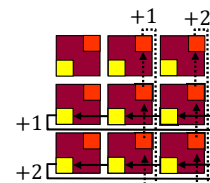
Product of dense matrixes of $\sqrt{N} \times \sqrt{N}$ elements
on a 2D torus of $\sqrt{P} \times \sqrt{P}$ processes

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: $A_{IK} \leftarrow, B_{KJ} \uparrow$ }

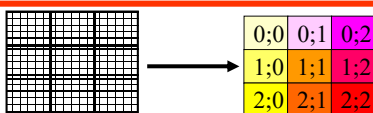
4 : Post-shifting of A and B blocks to restore input matrix partitioning

Dense matrix product on a 2D torus

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

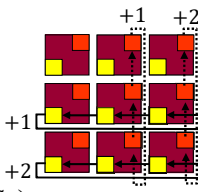
Canon algorithm: $C = A \times B$
 Product of dense matrixes of $\sqrt{N} \times \sqrt{N}$ elements
 on a 2D torus of $\sqrt{P} \times \sqrt{P}$ processes

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



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

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

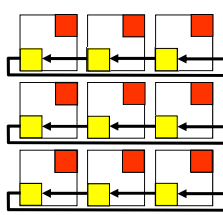
Dense matrix product on a 2D torus

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

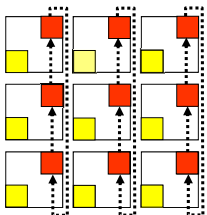
1 – Implementation with blocking communications (no overlapping)

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 rcv 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

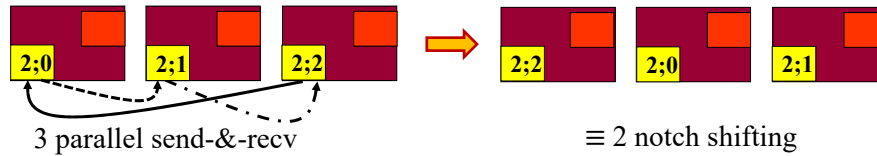
- 4 - All communications cross the same « network »
 (depends on the deployment !)

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

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)



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

Considering blocking communications:

$$\Rightarrow t_{pre-shifting-A-B} \approx 2 \cdot \frac{N}{P} \cdot t_w \quad \text{Same result for post-shifting}$$

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

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:

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

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

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:

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

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

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:

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

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

1 – Implementation with blocking communications (no overlapping)

Complete distributed computation:

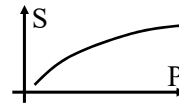
$$T_{C \text{ comput}}^{\text{distributed}} \approx 2 \cdot \frac{N^{3/2}}{P} t_{flop} + 2 \cdot (\sqrt{P} + 1) \cdot \frac{N}{P} t_w$$

Speedup:

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

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

With $N = N_0$



$$\text{Again: } \lim_{N \rightarrow \infty} S(N, P_0) = P_0$$

When pb size increase, speedup increases and tends to ideal speedup

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

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:

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

$$\rightarrow 2 \cdot \frac{N^{3/2}}{P} t_{flop} + (\sqrt{P} + 1) \cdot \frac{N}{P} t_w$$

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

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:

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

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

Speedup increases faster

$$\text{Again: } \lim_{N \rightarrow \infty} S(N, P_0) = P_0$$

When pb size increase, speedup increases and tends to ideal speedup

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

Comparison 1D Ring / 2D Torus solutions:

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

2D torus solution *should be* faster when

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

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

Distributed algorithms on multi-dimensional topologies

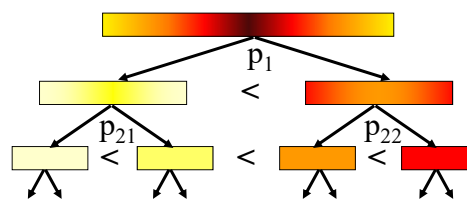
- Dense matrix product on a 1D ring
- Dense matrix product on a 2D torus
- Hyper-quicksort on a kD hypercube**
 - **Parallelisation scheme**

Hyper-quicksort on a kD hypercube

Principe séquentiel

Principe :

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



Performances :

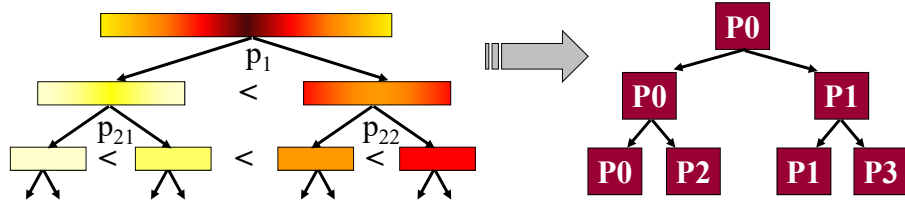
pire cas : $\left(\frac{N \cdot (N-1)}{2}\right) t_{comp} \rightarrow O(N^2)$

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

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

Parallélisation naïve

Principe :



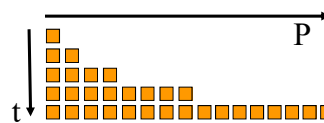
On crée un nouveau processus et on utilise un nouveau processeur à chaque appel récursif (en réutilisant les anciens).

Simple **mais inefficace !**

Parallélisation naïve

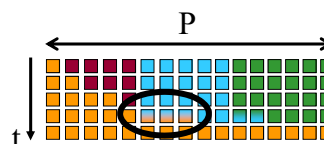
1 - Faiblesse conceptuelle :

Algorithme pas pleinement parallèle
 → Trouver plus parallèle !

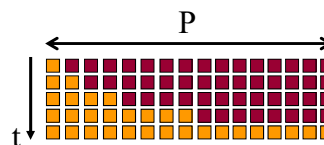


2 - Faiblesse de mise en œuvre :

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 $n-1$:

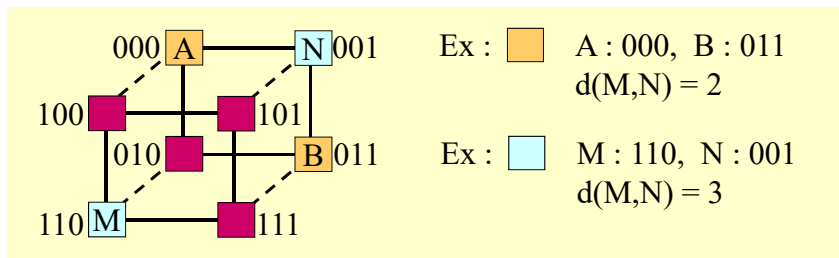
1 x dim 3 2 x dim 2 ou 2 x dim 2 ou 2 x dim 2

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

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)



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

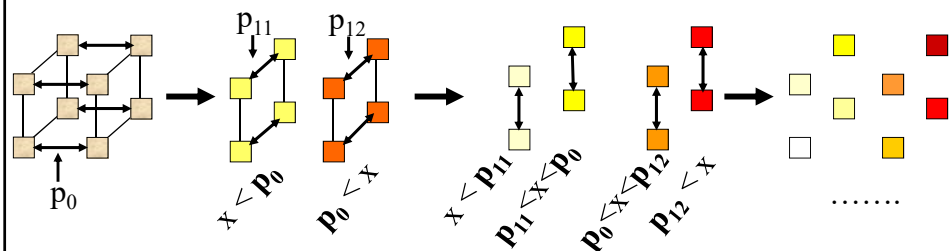
Principe de parallélisation sur Hyper-Cube

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

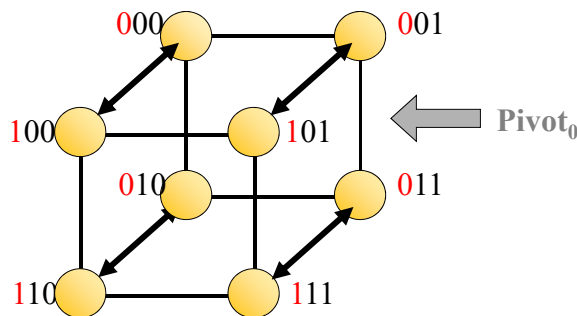


Principe de parallélisation sur Hyper-Cube

Init : Chaque nœud charge N/P données en local

Étape 0 :

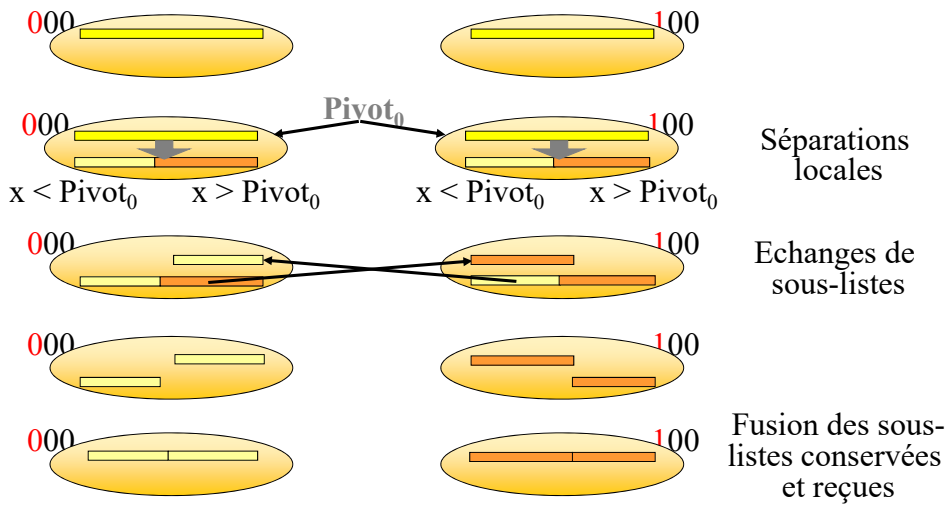
- Choix de 2^0 pivots pour les 2^0 hypercubes de dimension d-0
- Séparation selon le pivot sur chaque nœud de l'hypercube
- Echange des listes inférieures et supérieures en dimension d-0
- Fusions locales des listes conservées et des listes reçues



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

Principe de parallélisation sur Hyper-Cube

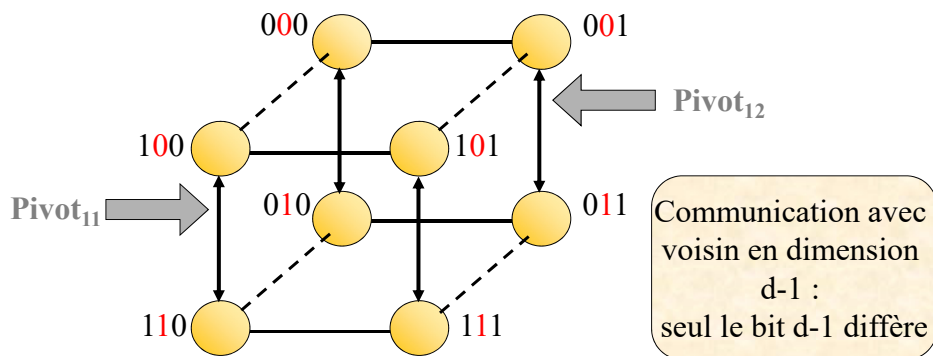
Détails de l'étape 0 sur PE-000 et PE-100 :



Principe de parallélisation sur Hyper-Cube

Etape 1 :

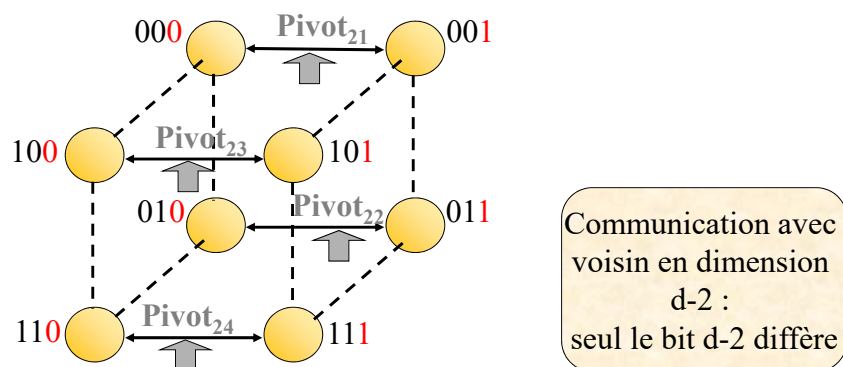
- Choix de 2^1 pivots pour les 2^1 hypercubes de dimension d-1
- Séparation selon un pivot sur chaque nœud des 2^1 hypercubes
- Echange des listes inférieures et supérieures en dimension d-1
- Fusions locales des listes conservées et des listes reçues



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

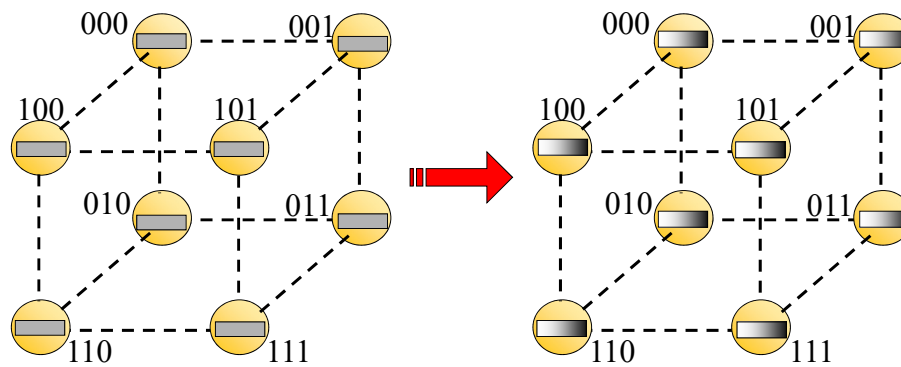


Principe de parallélisation sur Hyper-Cube

Etape finale :

Chaque processeur tri en local sa liste finale

→ ex : quick-sort local et séquentiel sur chaque processeur



1^{ère} implantation sur Hyper-Cube

```

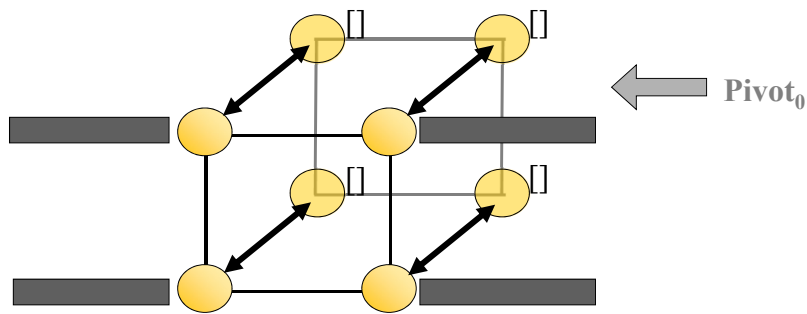
HcubeQuickSort(double *A, int d)
{
  int i; /* Compteur de dimension. */
  double *Ainf, *Asup, *Abuf; /* Tables de données. */
  double x; /* Pivot. */
  for (i = d-1; i >= 0; i--) { /* Pour chaque dimension: */
    x = choix_pivot(me,i); /* - Partitionnement local*/
    partitioner(A,x,Ainf,Asup);
    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 */

```

1^{ère} implantation sur Hyper-Cube

Faiblesse de ce premier hyper-quicksort :

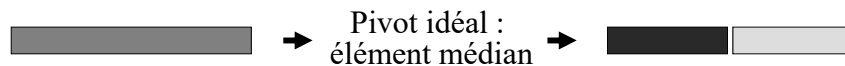
Si mauvais choix de pivot : déséquilibre définitif !
 → processeurs surchargés et processeurs à vide



→ Soigner le choix du pivot.

Parallélisation optimisée sur H-Cube

Choix de pivot optimisé :



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

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

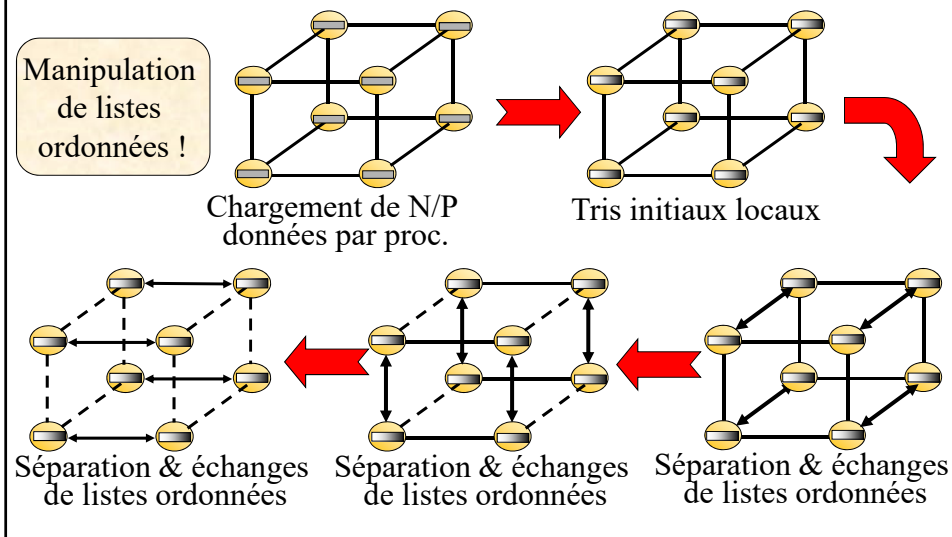
En parallèle :

- 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 !



Parallélisation optimisée sur H-Cube

Principe de l'hyper-quicksort optimisé :



2^{ème} implantation sur Hyper-Cube

Implantation de l'hyper-quicksort optimisé :

```

HcubeQuickSort(double *A, int d)
{
  int i;
  double *Ainf, *Asup, *Abuf;
  double x;
  QuickSortSequentiel(A);
  for (i = d-1; i >= 0; i--) {
    x = choix_pivot(me, i);
    partitioner(A, x, Ainf, Asup);
    if (me & exp2(i) == 0) {
      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