

Parallélisme des opérations de haut niveau sur les fichiers

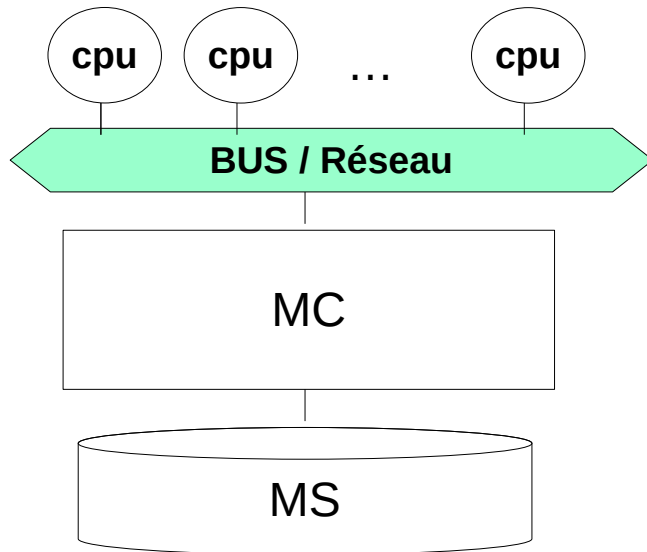
Une des manières d'améliorer les performances (en présence de volumes importants de données à traiter) est l'utilisation du parallélisme

Plan

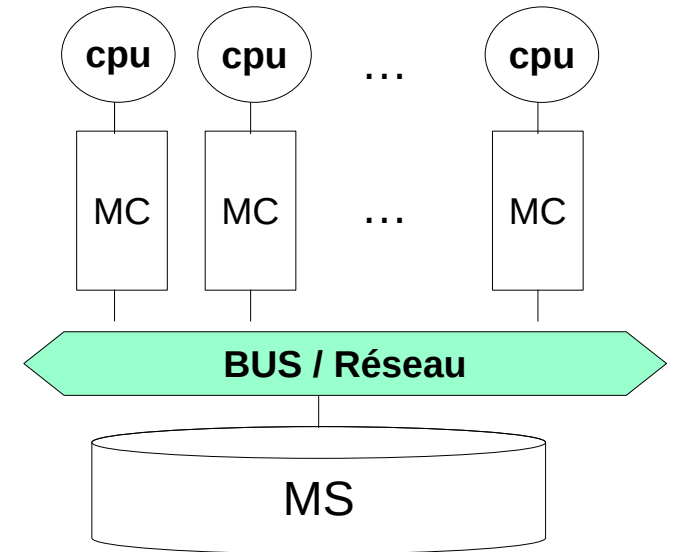
Architectures parallèles
Tri, Fusion et Jointures parallèles

Architectures Parallèles

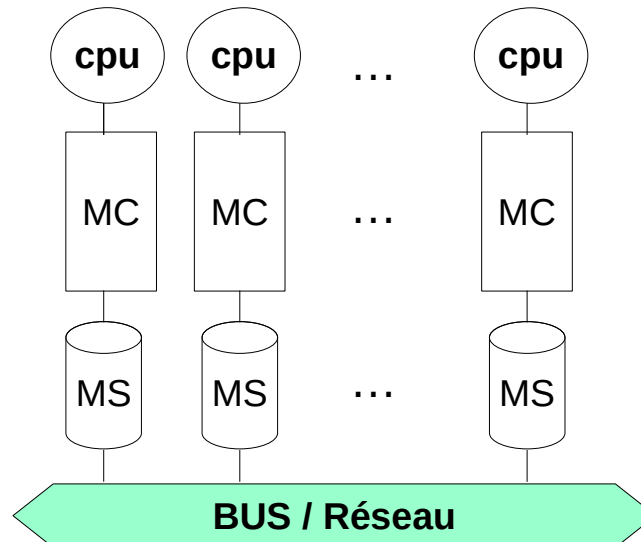
Shared Memory



Shared Disk



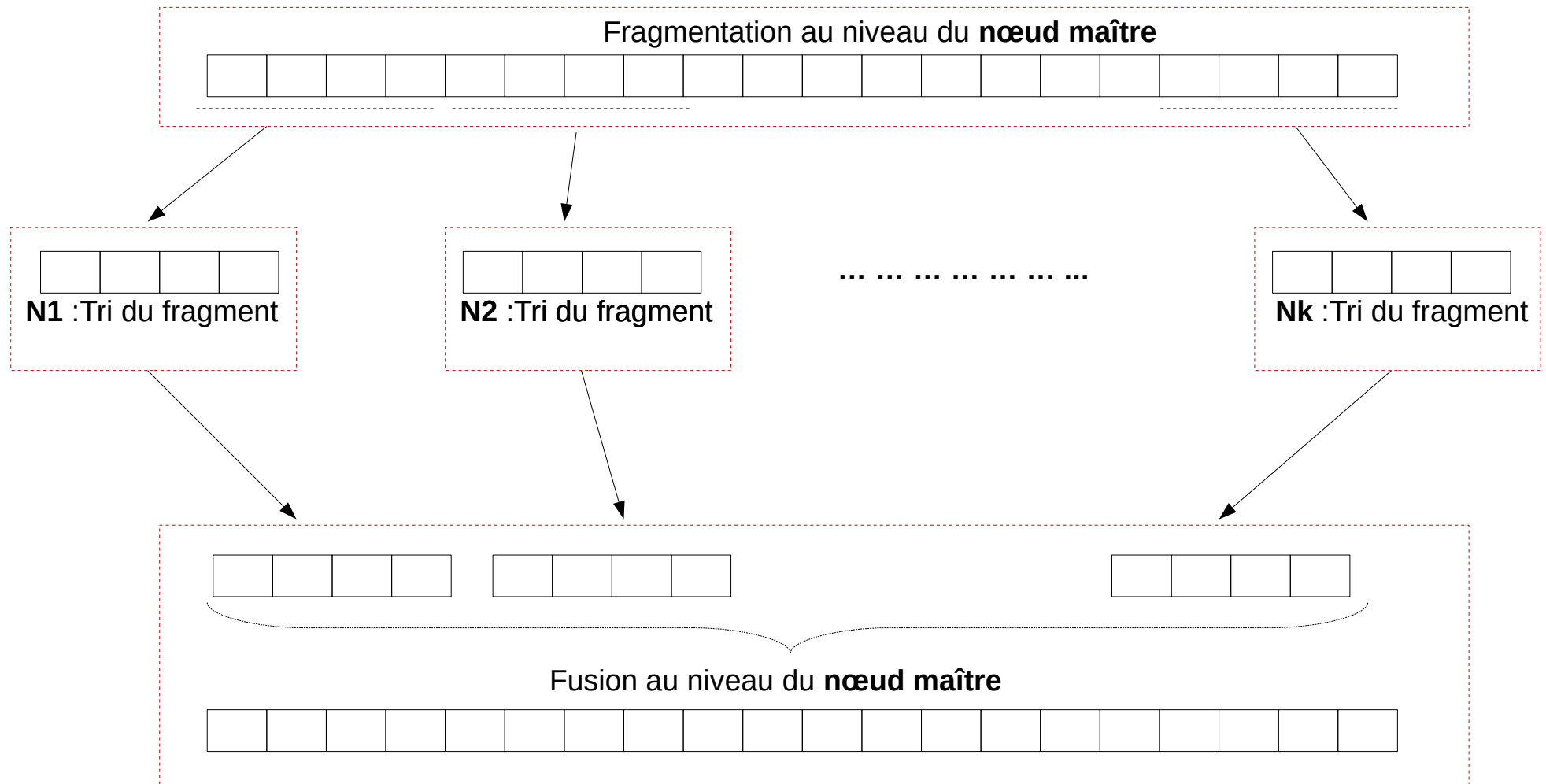
Shared Nothing (cluster)



Tri parallèle (sur cluster)

Première approche

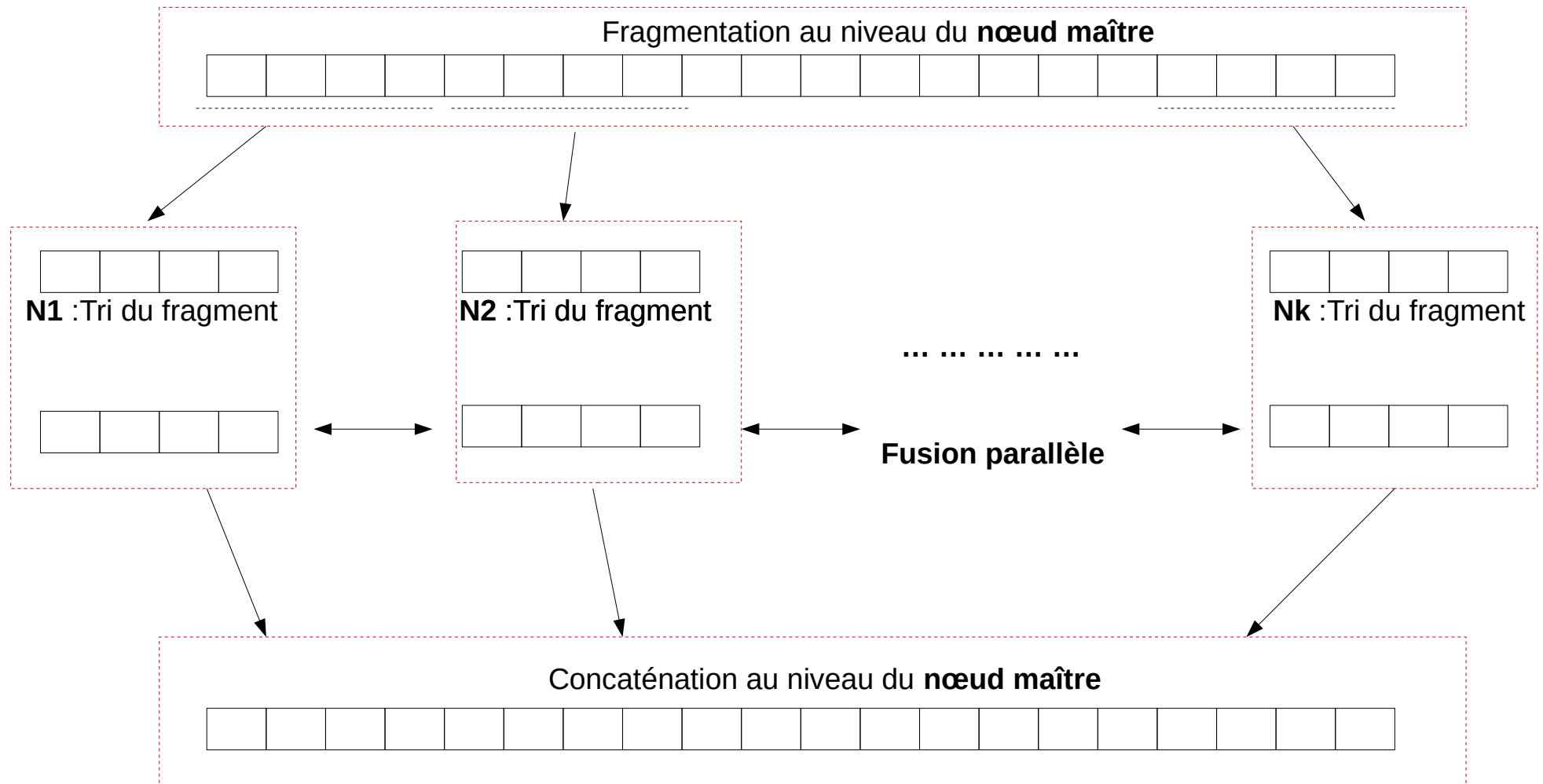
- Fragmentation par position (découpage en fragments de même taille)
- Tri des fragments (au niveau de chaque nœud)
- Fusion des résultats



Tri parallèle (sur cluster)

Première approche (améliorée)

- Fragmentation par position (découpage en fragments de même taille)
- Tri des fragments (au niveau de chaque nœud)
- Fusion des résultats

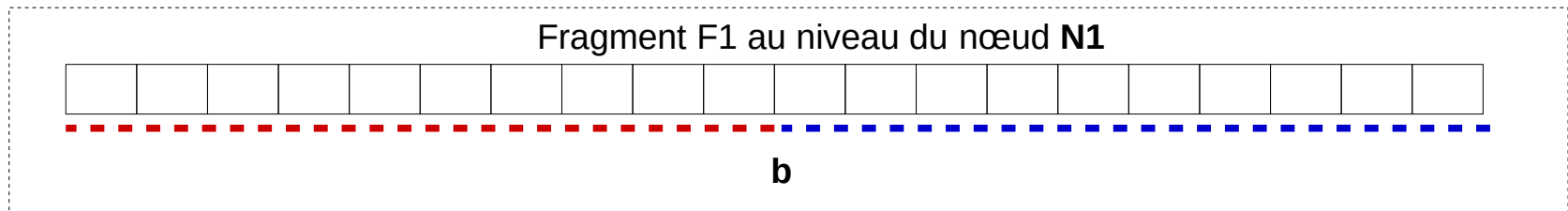


Tri parallèle (sur cluster)

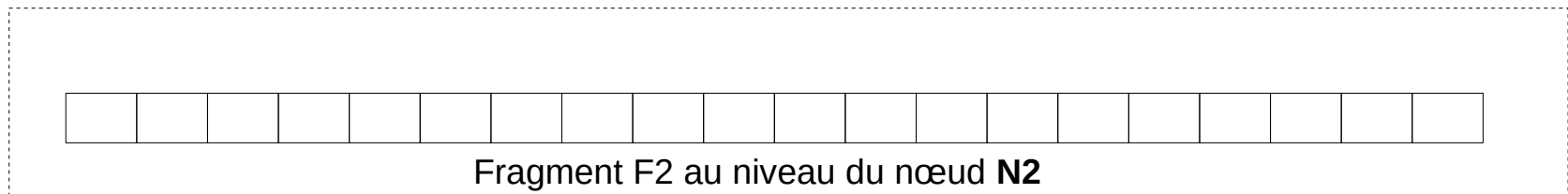
Fusion parallèle de 2 fragments (avec 2 nœuds N1 et N2)

- Le nœud N1 dispose d'un fragment trié → F1
- Le nœud N2 dispose d'un fragment trié → F2
- N1 décompose son fragment en 2 parties égales : **partie(1,1)** et **partie(1,2)**

N1



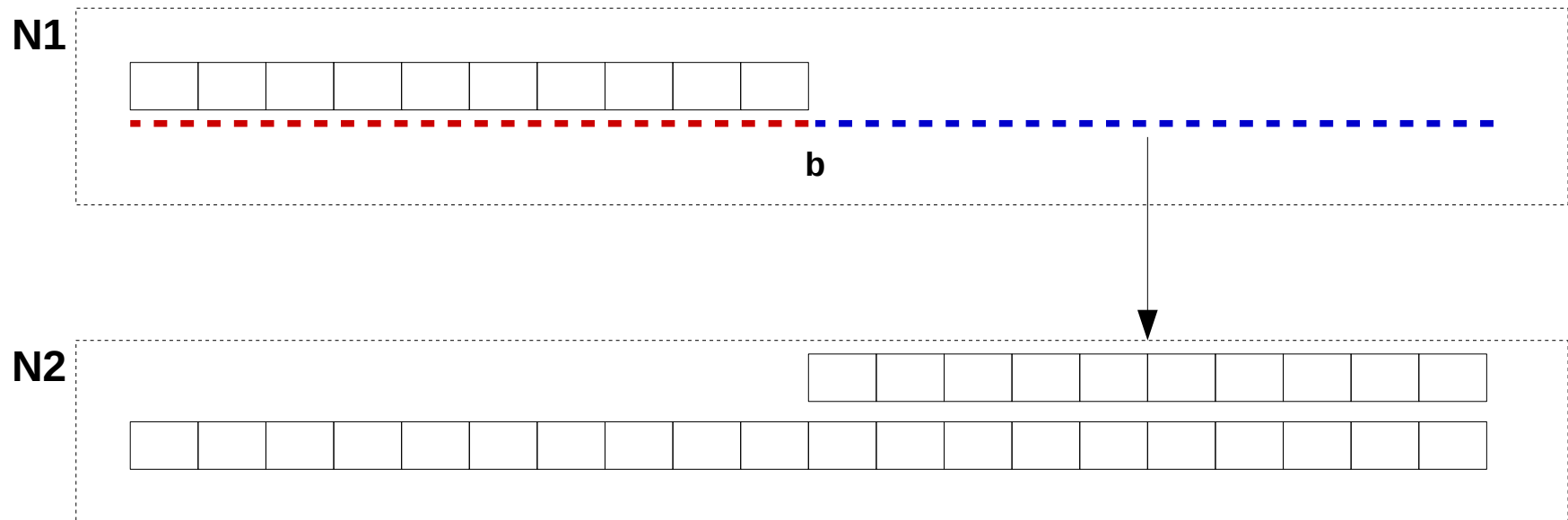
N2



Tri parallèle (sur cluster)

Fusion parallèle de 2 fragments (avec 2 nœuds N1 et N2)

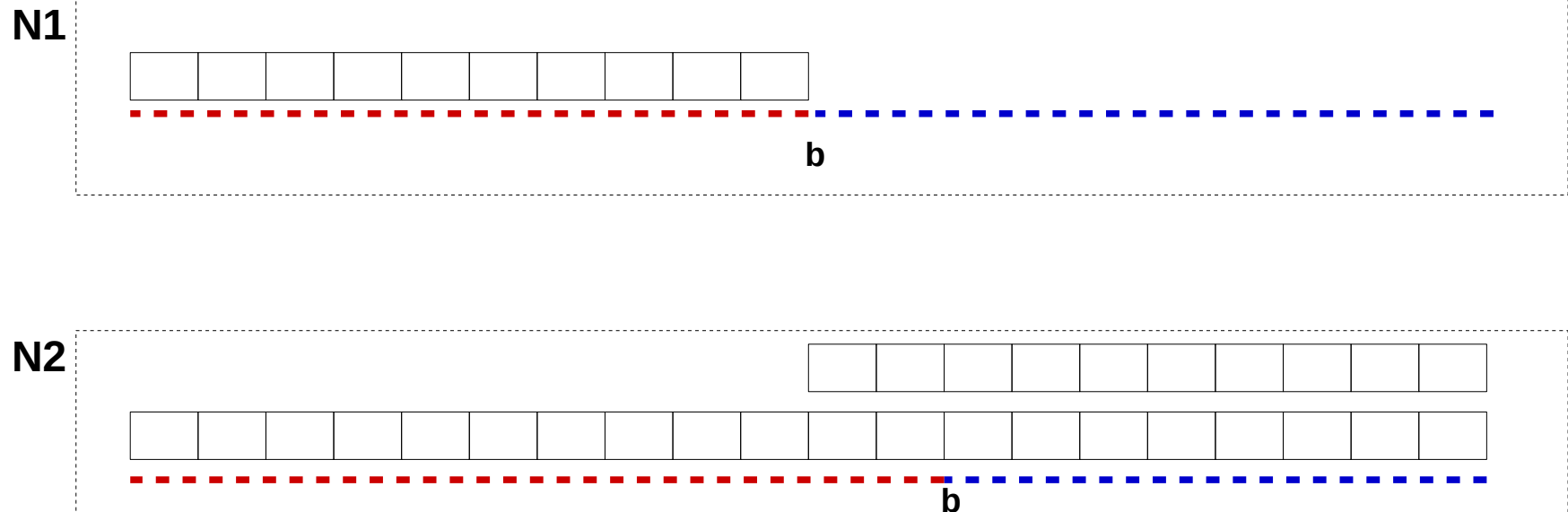
- Le nœud N1 dispose d'un fragment trié → F1
- Le nœud N2 dispose d'un fragment trié → F2
- N1 décompose son fragment en 2 parties égales : **partie(1,1)** et **partie(1,2)**
- N1 envoie une des 2 parties (avec la borne de séparation **b**) à N2



Tri parallèle (sur cluster)

Fusion parallèle de 2 fragments

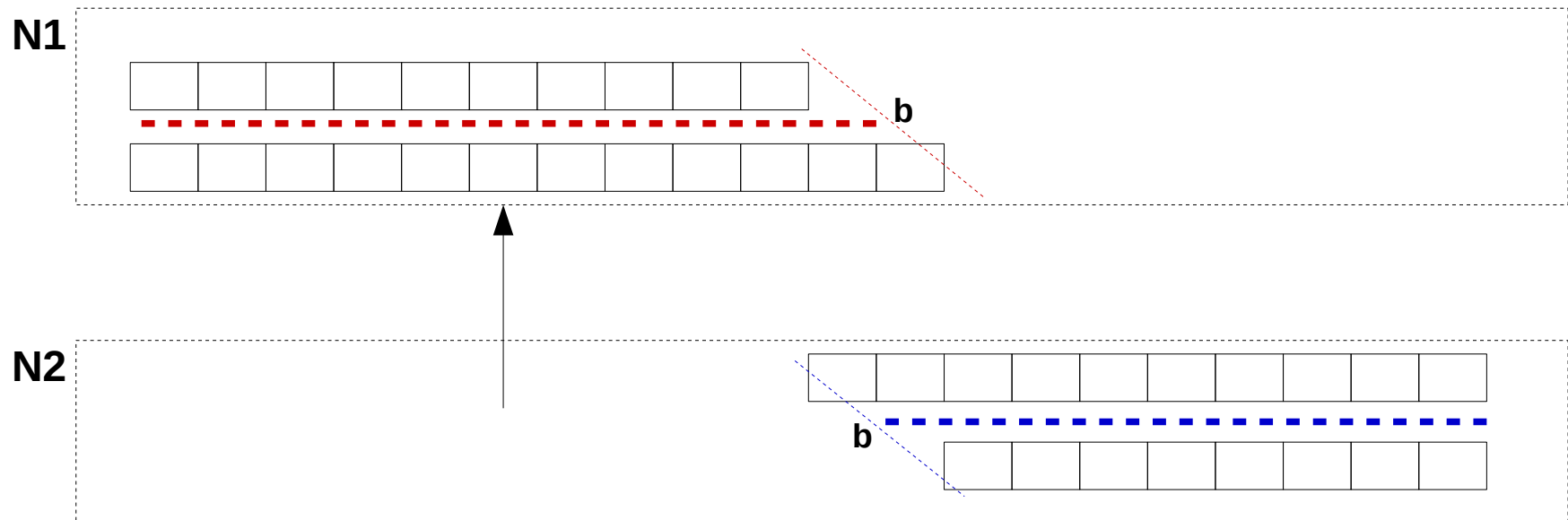
- Le nœud N1 dispose d'un fragment trié F1
- Le nœud N2 dispose d'un fragment trié F2
- N1 décompose son fragment en 2 parties égales : **partie(1,1)** et **partie(1,2)**
- N1 envoie une des 2 parties (avec la borne de séparation **b**) à N2
- N2 décompose alors son fragment en 2 parties : **partie(2,1)** et **partie(2,2)**



Tri parallèle (sur cluster)

Fusion parallèle de 2 fragments

- Le nœud N1 dispose d'un fragment trié F1
- Le nœud N2 dispose d'un fragment trié F2
- N1 décompose son fragment en 2 parties égales : **partie(1,1)** et **partie(1,2)**
- N1 envoie une des 2 parties (avec la borne de séparation **b**) à N2
- N2 décompose alors son fragment en 2 parties : **partie(2,1)** et **partie(2,2)**
- N2 envoie l'autre partie de son propre fragment à N1

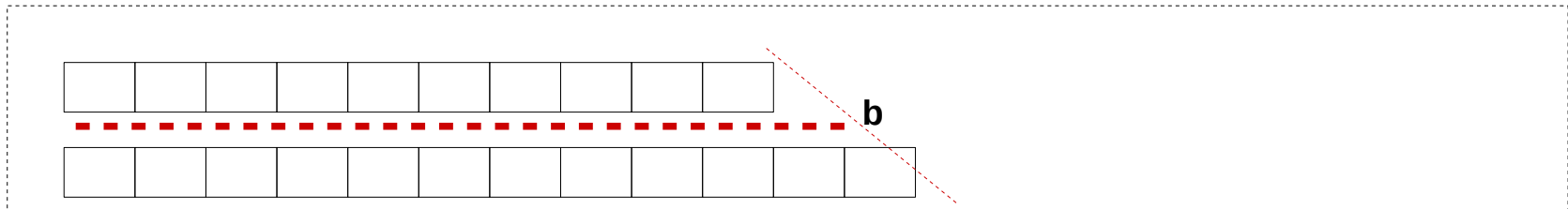


Tri parallèle (sur cluster)

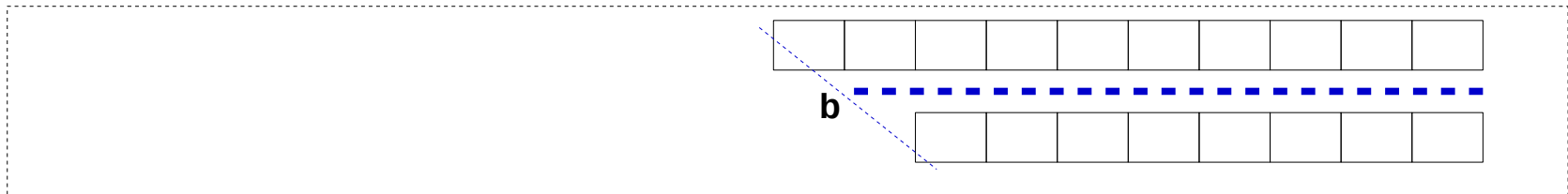
Fusion parallèle de 2 fragments

- Le nœud N1 dispose d'un fragment trié F1
- Le nœud N2 dispose d'un fragment trié F2
- N1 décompose son fragment en 2 parties égales : **partie(1,1)** et **partie(1,2)**
- N1 envoie une des 2 parties (avec la borne de séparation **b**) à N2
- N2 décompose alors son fragment en 2 parties : **partie(2,1)** et **partie(2,2)**
- N2 envoie l'autre partie de son propre fragment à N1
- N1 et N2 effectuent en parallèle la fusion des parties en leur possession
 - Dans N1 : Fusion (**partie(1,1)** , **partie(2,1)**)
 - Dans N2 : Fusion (**partie(1,2)** , **partie(2,2)**)
- Les résultats des fusions sont envoyés au nœud maître pour être concaténés.

N1



N2



Généralisation : fusion parallèle de k fragments (sur k nœuds)

Principe

Chacun des k fragments (F_i) est découpé en k parties : $[P(i,1) , P(i,2) , \dots P(i,k)]$

Chacun des k nœuds (N_i) effectue une multi-fusion des parties numéro i de chaque fragment ($P(1,i), \dots P(k,i)$)

Tâche de N_1 :

- N_1 décompose son fragment en k parties égales : $p(1,1) , p(1,2) \dots p(1,k)$
et diffuse les k-1 **bornes de séparation** aux autres nœuds ($N_2, N_3, \dots N_k$)
- N_1 garde pour lui la première partie et envoie les autres parties aux nœuds restant :
 $p(1,2) \rightarrow N_2 ; p(1,3) \rightarrow N_3 \dots p(1,k) \rightarrow N_k$
- N_1 se met en attente des parties qui le concernent (provenant des autres nœuds)
- Dès que N_1 reçoit toutes les parties $p(*,1)$ des autres nœuds, il entame la multi-fusion des parties en sa possession : $\text{MultiFusion}(p(1,1) , p(2,1), \dots p(k,1))$
- Le résultat est envoyé au nœud maître.

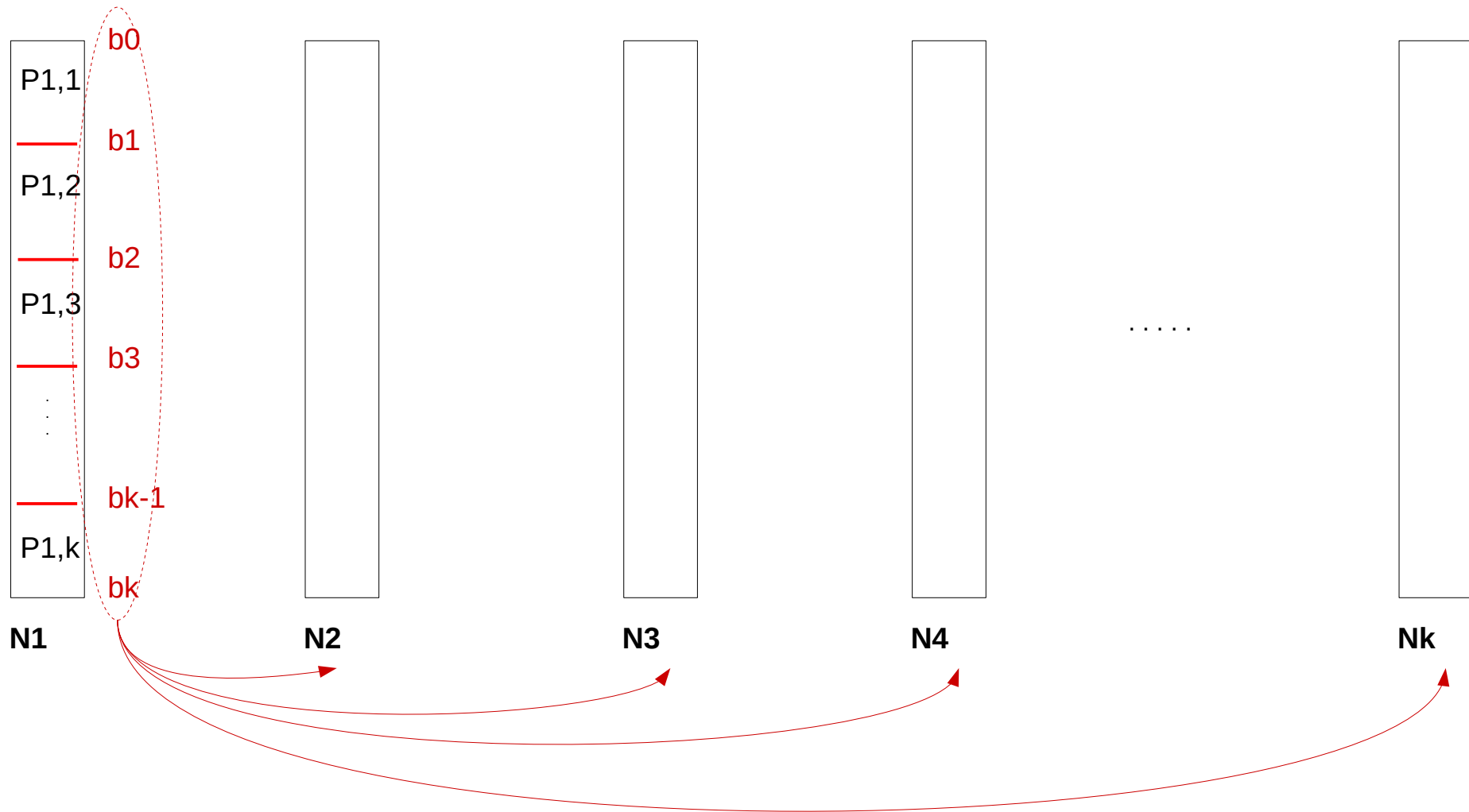
Tâche de N_i ($i = 2 \dots k$) :

- N_i se met en attente des **bornes de séparation** (provenant de N_1)
- N_i décompose son fragment local en k parties (à l'aide des bornes reçues de N_1) :
 $p(i,1) , p(i,2) \dots p(i,k)$
- N_i garde pour lui la partie $p(i,i)$ et envoie le reste aux autres nœuds :
 $p(i,1) \rightarrow N_1 ; p(i,2) \rightarrow N_2 \dots p(i,k) \rightarrow N_k$
- N_i se met en attente des parties qui le concernent (provenant des autres nœuds)
- Dès que N_i reçoit toutes les parties $p(*,i)$ des autres nœuds, il entame la multi-fusion des parties en sa possession : $\text{MultiFusion}(p(1,i) , p(2,i), \dots p(k,i))$
- Le résultat est envoyé au nœud maître.

Généralisation : fusion parallèle de k fragments (sur k nœuds)

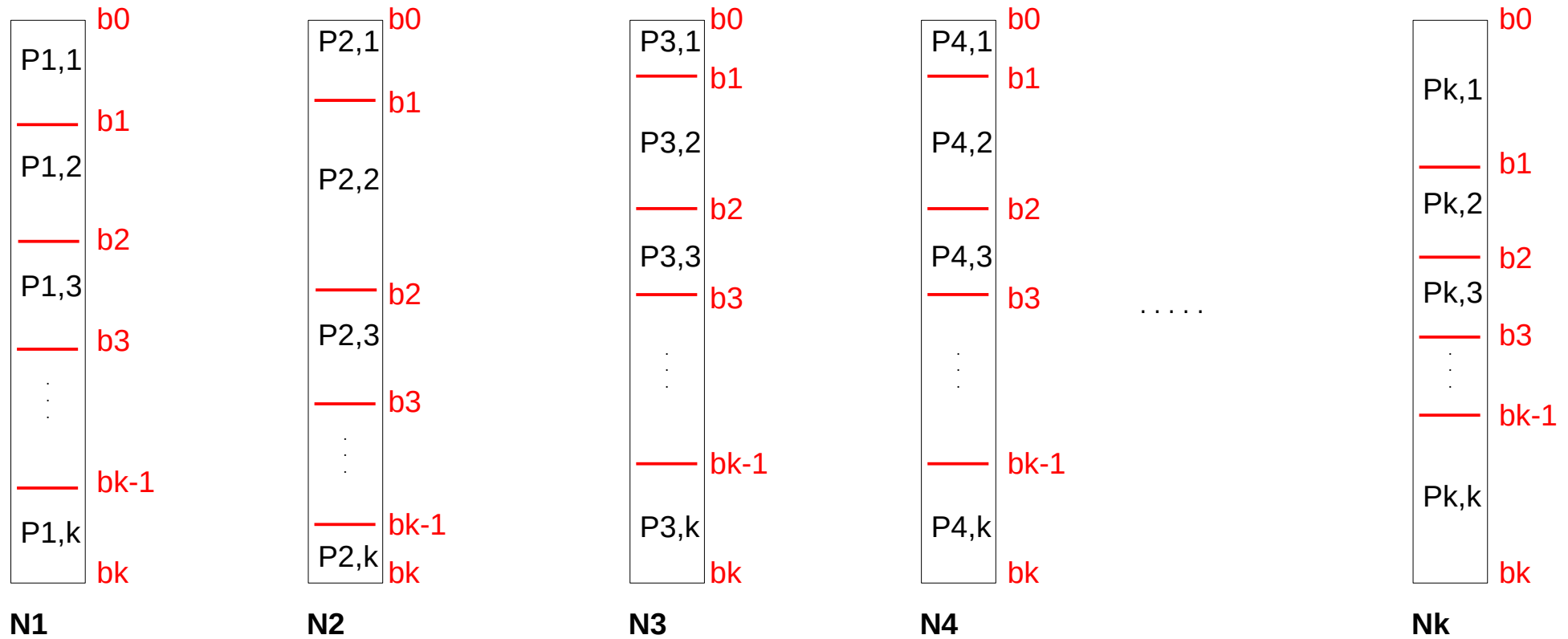
N1 découpe son fragment en k parties égales

- cela détermine les bornes de séparation : b_0, b_1, \dots, b_k
- les bornes de séparation sont envoyées aux autres nœuds

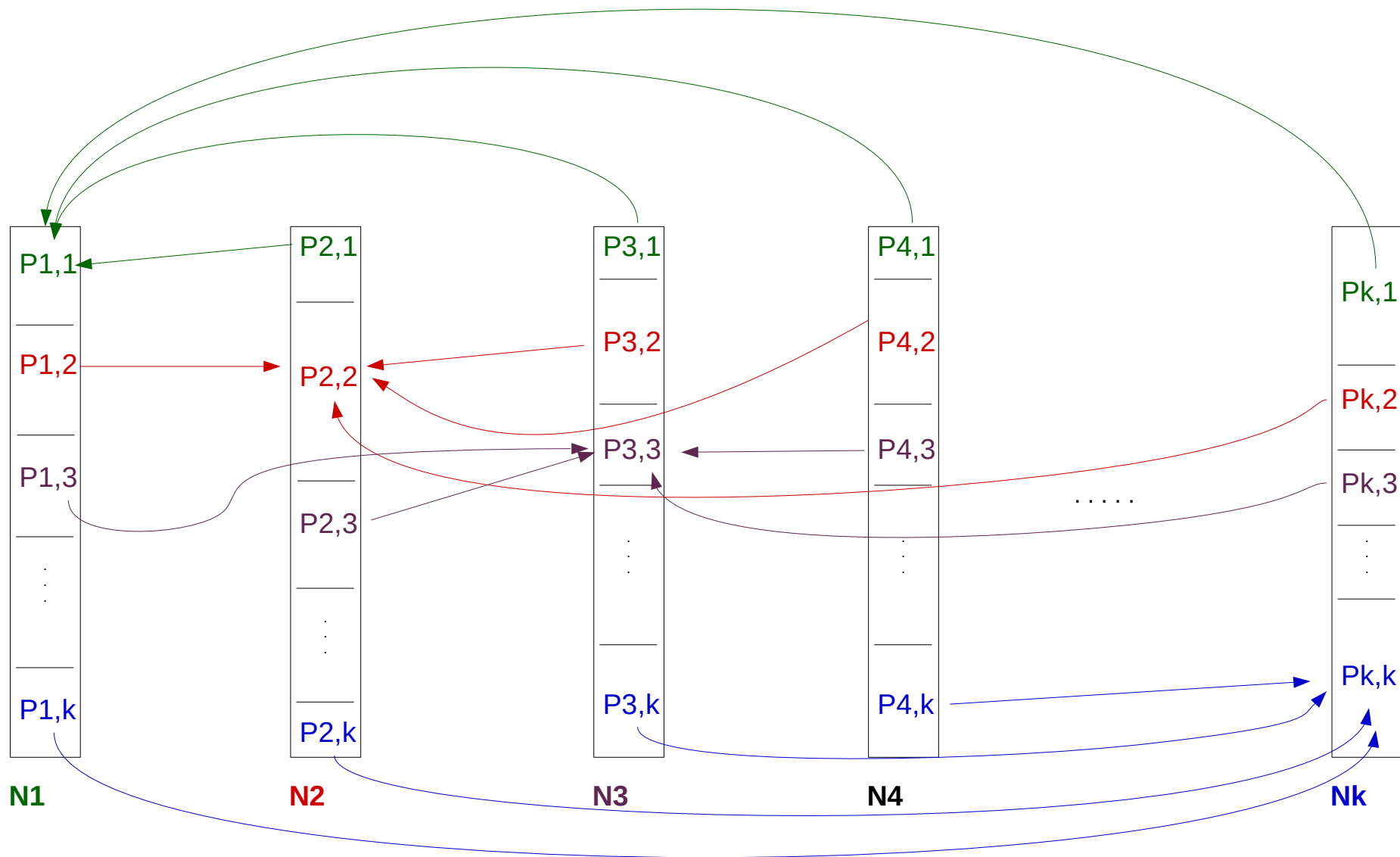


Généralisation : fusion parallèle de k fragments (sur k nœuds)

Les autres nœuds découpent leurs fragments respectifs à l'aide des bornes reçus de N1



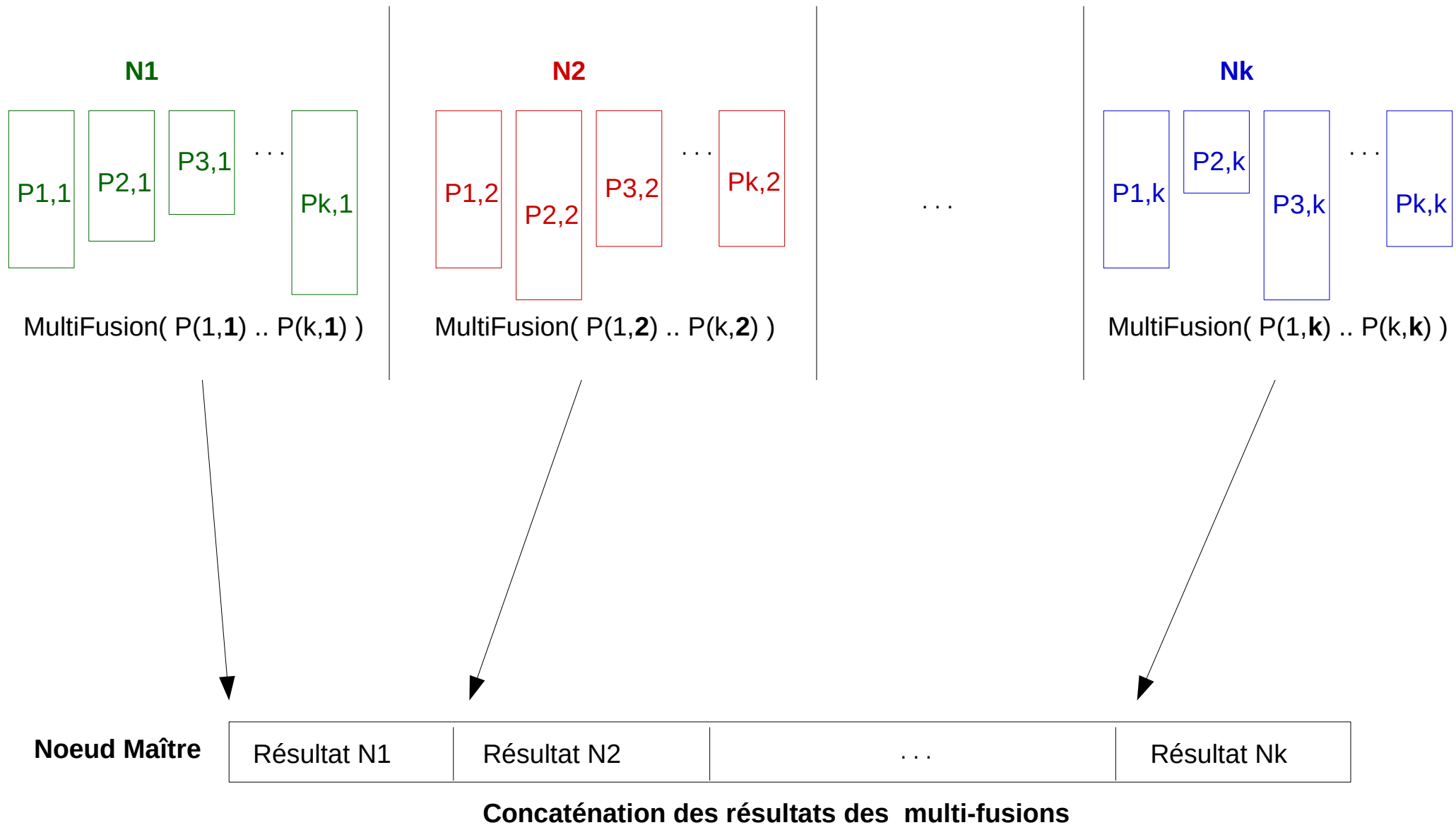
Généralisation : fusion parallèle de k fragments (sur k nœuds)



Chaque nœud Ni garde pour lui la partie $P(i,i)$ et envoie les autres parties $P(i,j)$ aux nœuds concernés : Nj (avec $j \neq i$)

Généralisation : fusion parallèle de k fragments (sur k nœuds)

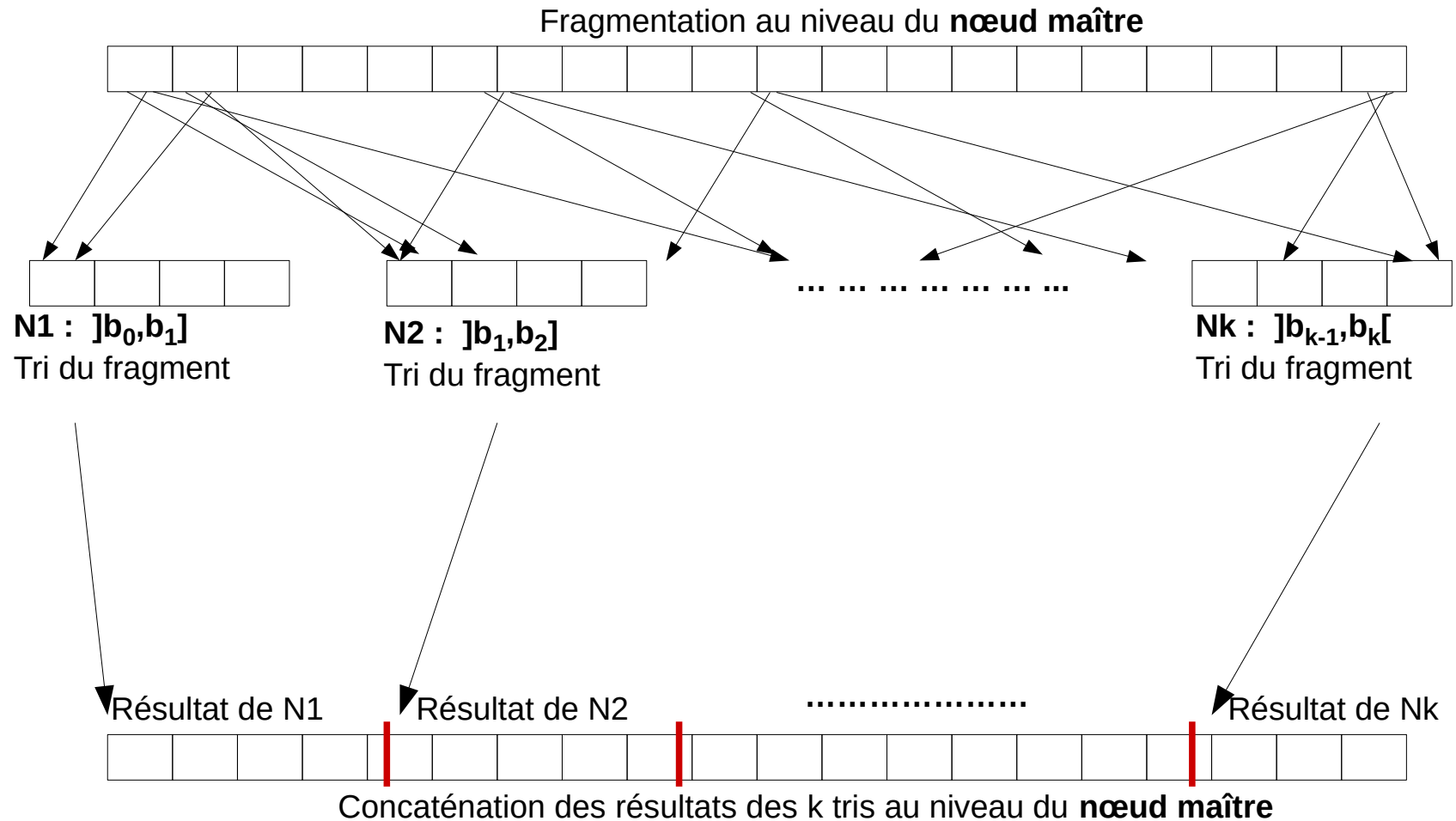
Chacun des nœuds N_i effectue une **multi-fusion** des parties qui le concernent et envoie le résultat au **nœud maître** pour être concaténé avec les autres résultats



Tri parallèle (sur cluster)

Deuxième approche

- Fragmentation par valeur (par intervalle)
- Tri des fragments (au niveau de chaque nœud)
- Concaténation des résultats



Jointure parallèle (sur cluster)

Première approche (découpage par position)

- Le nœud maître fragmente R en k parties de tailles égales : r_1, r_2, \dots, r_k
- Chaque partie r_i est envoyée avec le 2^e fichier S au nœud N_i
- Chaque Nœud N_i effectue la jointure $r_i * S$ et renvoie le résultat au nœud maître
- Le nœuds maître concatène les résultats (dans un ordre quelconque)

Deuxième approche (découpage par valeur)

- Le nœud maître fragmente R en k parties : r_1, r_2, \dots, r_k (selon les valeurs d'enreg)
- Le nœud maître fragmente S en k parties : s_1, s_2, \dots, s_k (selon les valeurs d'enreg)
- Chaque couple de parties $\langle r_i, s_i \rangle$ est envoyée au nœud N_i
- Chaque Nœud N_i effectue la jointure $r_i * s_i$ et renvoie le résultat au nœud maître
- Le nœuds maître concatène les résultats (dans un ordre quelconque)

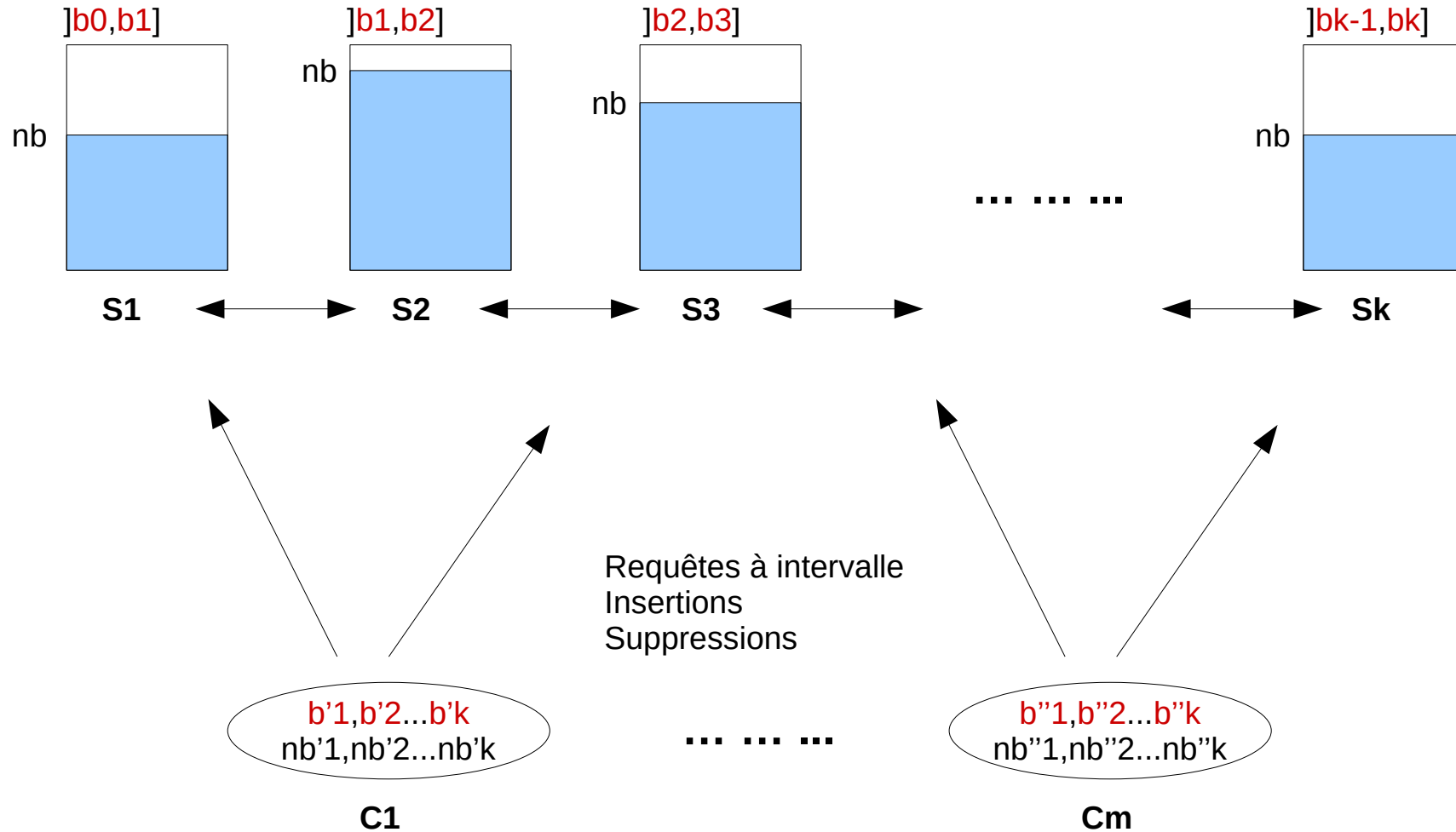
→ Le découpage par valeur peut se faire dans ce cas soit avec une fragmentation par hachage, soit avec une fragmentation par intervalle

→ risque de déséquilibre de charge causée par les données (data skew)

Exemple d'application

Système de stockage réparti, fragmenté par intervalle à **bornes dynamiques**

On a k serveurs : S_1, S_2, \dots, S_k et m clients : C_1, C_2, \dots, C_m



Chaque client possède une **image** ± correcte de la répartition des données

- les bornes (b_i) et les charges (nb) des serveurs
- en cas d'erreur d'adressage, le client est partiellement corrigé par le serveur ciblé

Chaque serveur possède aussi une **image** ± correcte de la répartition des données

- les bornes (b_i) et les charges (nb) des autres serveurs
- utilisée pour des besoins de rééquilibrage asynchrone entre serveur

