

## CODED DATA REBALANCING

### 1 Data Skew and Data Rebalancing in Distributed Systems

#### Data Skew

Non-uniform distribution of data across storage nodes

#### Can arise because of

- Node additions or removals
- Load imbalance
- Behaviour of client applications
- Stragglers
- Behaviour of the file system
- Increase in task completion time

Remedy : **Data Rebalancing**

#### Data Rebalancing

Redistribute data across the available nodes to balance the distribution.

- Rebalancing may be needed at regular intervals
- Communication costs
- Reduction in performance during rebalancing.

#### Introducing Coded Data Rebalancing

- Exploit data replication for Coded transmissions during rebalancing
- Preserve database structure post rebalancing.

### 2 Coded Data Rebalancing : Formal System Model

#### System Model: Initial database

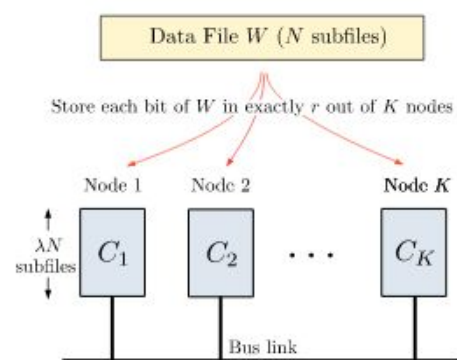
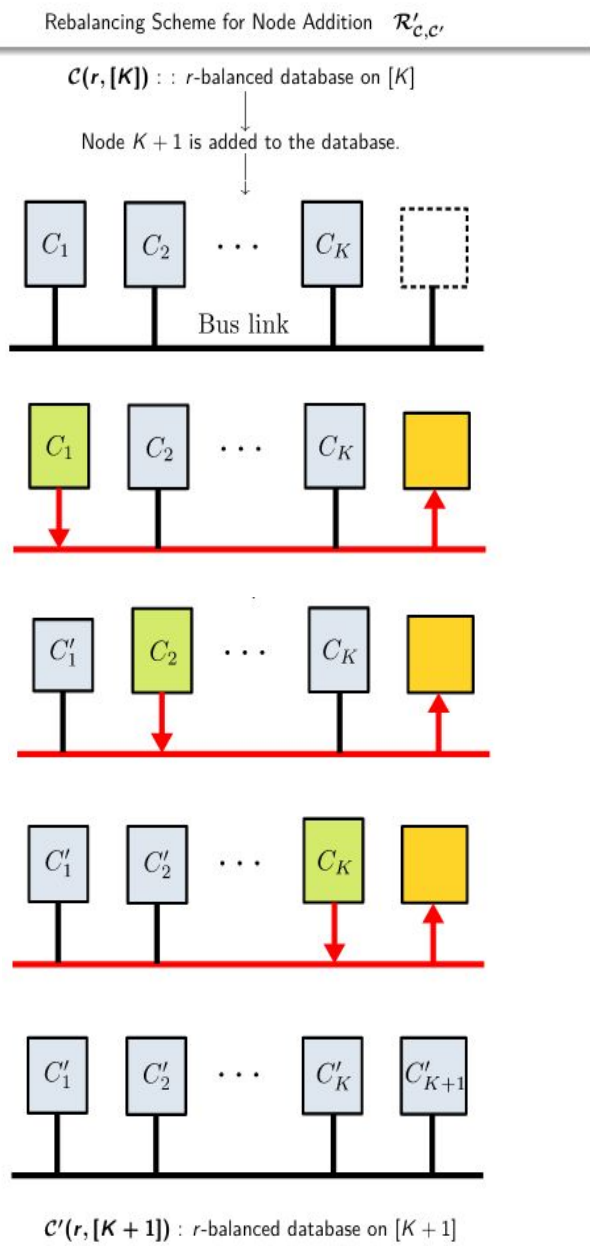
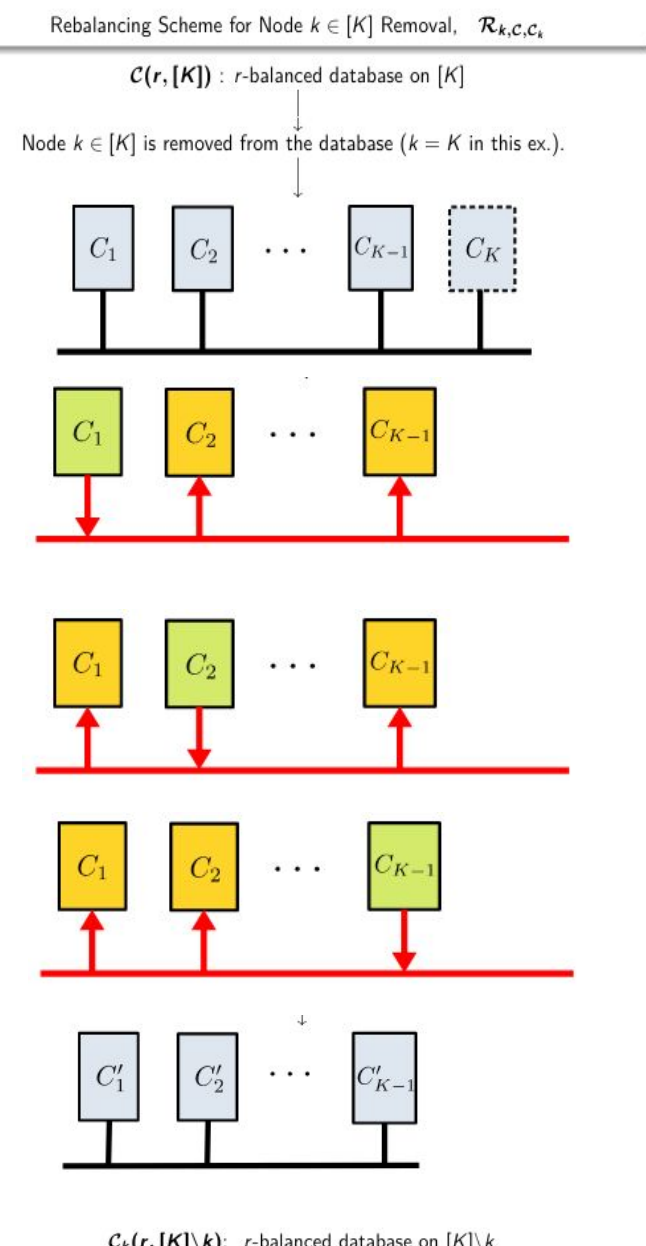


Figure: An  $r$ -balanced distributed database  $C(r, [K])$ , where  $[K] = \{1, \dots, K\}$ .

- $r$  : Replication factor
- 'Balanced': each node stores  $\lambda = \frac{r}{K}$  fraction of the data.



$C'(r, [K+1])$  :  $r$ -balanced database on  $[K+1]$



$C_k(r, [K]\kappa)$  :  $r$ -balanced database on  $[K]\kappa$

### Communication Load of Rebalancing Scheme

Let  $l_i$  be the number of bits sent by node  $i$ .

Node addition:

$$L_{add}(\mathcal{R}_{C, C'}) \triangleq \frac{\sum_{i \in [K]} l_i}{\lambda_{add} N}, \text{ where } \lambda_{add} = \frac{r}{K+1}.$$

Node Removal:

$$L_{rem}(\mathcal{R}_{k, C, C_k}) \triangleq \frac{\sum_{i \in [K]\kappa} l_i}{\lambda_{rem} N}, \text{ where } \lambda_{rem} = \frac{r}{K-1}.$$

$$L^*(r) = \inf_{C, \{C_k: k \in [K]\}, C'} \inf_{\mathcal{R}_{k, C, C_k}, \mathcal{R}_{C, C'}} \max_{k \in [K]} L_{rem}(\mathcal{R}_{k, C, C_k}) + L_{add}(\mathcal{R}_{C, C'})$$

### Main Contribution

Converse

$$L^* \geq \frac{1}{r-1} + 1 \text{ (node removal + node addition).}$$

Achievable Scheme

- Achieves **optimal** communication load (whereas  $L_{uncoded} = 1 + 1$ )
- **Optimality for any sequence** of node removals or additions
- **Structural invariance**: Maintains the structure of the database across node removal/addition.

### 3 Achievability

#### Achievable Scheme : Family of $r$ -balanced Databases

- Divide the data  $W$  into  $\frac{K!}{r!}$  subfiles.
- Index set of subfiles

$$S([K], K-r) \triangleq \{i = (i_1, \dots, i_{K-r}) : \{i_1, \dots, i_{K-r}\} \subseteq [K]\}$$

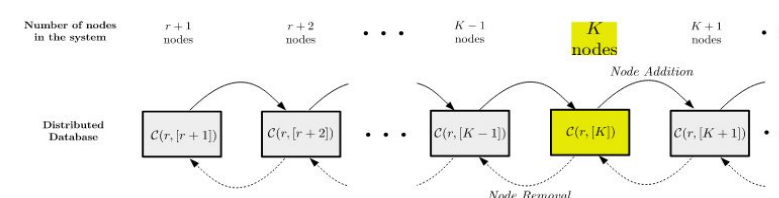
$r$ -balanced Distributed Database

For each  $i = (i_1, \dots, i_{K-r}) \in S([K], K-r)$

$W_i$  is stored in  $\{j_1, \dots, j_r\}$ , where  $\{j_1, \dots, j_r\} = [K] \setminus i$ .

Easy to check this is  $r$ -balanced database. We denote this as  $C(r, [K])$

#### Illustration of Rebalancing and Structural Invariance



#### Achievable Scheme : Node Addition

- Starting database  $C(r, [K])$
- Subfile indices:  $S([K], K-r)$
- Target Database  $C(r, [K+1])$  (keep  $r$  replication, balance node storage)
- Subfile indices  $S([K+1], K+1-r)$

For each subfile  $i = (i_1, \dots, i_{K-r}) \in S([K], K-r)$  (where  $\{i_1, \dots, i_r\} = [K] \setminus i$ )

- Split  $W_i$  into  $K+1$  subfiles, re-indexed as  $(j_1, i_1, \dots, i_{K-r}), (j_2, i_1, \dots, i_{K-r}), \dots, (j_r, i_1, \dots, i_{K-r}), (K+1, i_1, \dots, i_{K-r}), (i_1, K+1, i_2, \dots, i_{K-r}), \dots, (i_1, \dots, i_{K-r}, K+1)$

- Note that these new indices are in  $S([K+1], K+1-r)$ .
- Node  $j_i$  forwards  $(j_i, i_1, \dots, i_{K-r})$  to  $K+1$ , and then deletes it.

#### Example

$K=5, r=3$ , with database  $C([5], 3)$ . Node 6 is added.

- **Subfile Splitting**:  $W_{[2,3]}$  is split into  $K+1=6$  chunks at  $\{1, 4, 5\}$  as

$$W_{[1,2,3]}, W_{[4,2,3]}, W_{[5,2,3]}, W_{[6,2,3]}, W_{[2,6,3]}, W_{[2,3,6]}$$

- **Transmissions and Deletion**:
  - Node 1 transfers  $W_{[1,2,3]}$  to Node 6 and deletes it.
  - Node 4 transfers  $W_{[4,2,3]}$  to Node 6 and deletes it.
  - Node 5 transfers  $W_{[5,2,3]}$  to Node 6 and deletes it.

Same is done for each subfile.

Final database on 6 nodes:  $C([6], 3)$ .

#### Achievable Scheme : Node $k \in [K]$ Removal

- Starting database  $C(r, [K])$
- Subfile indices:  $S([K], K-r)$
- Target Database  $C(r, \mathcal{K}')$ , where  $\mathcal{K}'$  denotes the survivor set  $[K] \setminus k$ .
- Rebalance to reinstate replication factor of subfiles in node  $k$ , and balance storage.
- Subfile indices  $S' = S(\mathcal{K}', K-1-r)$

For subfile  $i \in S'$  (where  $\{i_1, \dots, i_r\} = \mathcal{K}' \setminus i$ )

- Consider the  $r$  subfile indices,  $(j_1, i), (j_2, i), \dots, (j_r, i)$

- Note that these indices are in  $S([K], K-r)$ , representing existing subfiles that were stored in node  $k$ .
- These have replication factor reduced to  $r-1$ , must be reinstated to  $r$ .

#### Achievable Scheme : Node $k \in [K]$ Removal (Step 1 - Grouping and Transmissions)

For each  $i \in S'$  (where  $\{i_1, \dots, i_r\} = \mathcal{K}' \setminus i$ )

- **Subfile Grouping**: Consider the  $r$  subfiles indices,  $(j_1, i), (j_2, i), \dots, (j_r, i)$

- Only subfile  $W_{(j_i, i)}$  is unavailable at  $j_i$ , available at  $\{j_1, \dots, j_r\} \setminus j_i$ .
- **Subfile Exchange via Coded Transmissions**: Each  $j_i$  does one coded transmission of size  $\frac{1}{r-1}$  of that of a subfile

$$\bigoplus_{m \in [r] \setminus i} W_{(j_m, i)}^i \text{ (for example, } W_{(j_1, i)}^1 \oplus \dots \oplus W_{(j_r, i)}^r \text{ by node } j_1)$$

- where  $W_{(j_m, i)}^i$  represents one chunk of  $(r-1)$  chunks of subfile  $W_{(j_m, i)}$ .
- After these  $r$  transmissions,  $(j_i, i)$  is available at node  $j_i$  also.

#### Achievable Scheme : Node $k \in [K]$ Removal (Step 2 - Recombining)

For each  $i \in S'$  (where  $\{i_1, \dots, i_r\} = \mathcal{K}' \setminus i$ )

- Each node  $\{j_1, \dots, j_r\}$  has following subfiles

$$(j_1, i_1, \dots, i_{K-r-1}), (j_2, i_1, \dots, i_{K-r-1}), \dots, (j_r, i_1, \dots, i_{K-r-1}), (k, i_1, \dots, i_{K-r-1}), (i_1, k, i_2, \dots, i_{K-r-1}), \dots, (i_1, \dots, i_{K-r-1}, k)$$

- Combine the above  $K$  subfiles and create a new (larger) subfile with index  $i \in S'$ .

#### Example

$K=5, r=3$ , with database  $C([5], 3)$ . Node 5 is removed.

- **Subfile Grouping**: Form a group with  $i = [3] \in S([4], 1)$

$$W_{[1,3]}, W_{[2,3]}, W_{[4,3]}$$

- **Subfile Exchange through Coded transmissions**
  - Node 1 transmits  $W_{[2,3]}^1 \oplus W_{[4,3]}^1$ .
  - Node 2 broadcasts  $W_{[1,3]}^2 \oplus W_{[4,3]}^2$ .
  - Node 4 broadcasts  $W_{[1,3]}^4 \oplus W_{[2,3]}^4$ .
- Combine the following at nodes  $\{1, 2, 4\}$ , and relabel as  $W_{[3]}$

$$W_{[1,3]}, W_{[2,3]}, W_{[4,3]}, W_{[5,3]}, W_{[3]}$$

Same is done for each  $i \in S([4], 1)$ . Final database on 4 nodes:  $C([4], 3)$ .