

The Resource Allocation Problem for CERN Dataset Transfer

Qiao Xiang
qiao.xiang@cs.yale.edu

1 The Offline Model for Resource Allocation Problem

1.1 Flow prioritized by bandwidth constraint

We use a directed graph $G = (V, E)$ to denote the network. And we consider the system operating in a discrete-time, where time is divided into time slots with equal length and indexed as $t = 1, 2, \dots, T$. Given a time slot t , each link (i, j) have a link capacity denoted as $u_{ij}(t)$. At the beginning of time slot 1, we are given a set of dataset transfer requests, denoted as \mathbb{M} . Each request $m \in \mathbb{M}$ is associate with a priority level p_m . This priority p_m indicates that the transmission rate per time slot of m must be between a predefined interval $[L(p_m), U(p_m)]$ before m is accomplished, with the only exception in the final transmission slot where $L(p_m)$ may be violated. With this setting, a dataset transfer request m can be represented as a 5-tuple $\{s_m, d_m, v_m, L(p_m), U(p_m)\}$. In this tuple, s_m denotes the source location of dataset, d_m denotes where the dataset should be sent to, and v_m denotes the volume of dataset.

Decision Variables: For each request m , we have three sets of decision variables in the model:

- T_m , the number of time slots needed to finish the transfer of m ;
- $r_m(t)$, the data volume of data transfer request m that is transferred in time slot t ;
- $f_{ij}^m(t)$, the data volume of data transfer m along link $(i, j) \in E$ in time slot t .

And we formulate the resource allocation problem (**RAP**) as the following model.

$$\begin{aligned}
 & \text{subject to} & \textbf{RAP} \quad & \text{minimize} \quad T & (1) \\
 & & \sum_m f_{ij}^m(t) & \leq u_{ij}(t) & \forall (i, j) \in E \text{ and } t, & (2a) \\
 & & \sum_{j \in V} f_{ij}^m(t) - \sum_{j \in V} f_{ji}^m(t) & = 0 & \forall m, t \text{ and } i \in V / \{s_m, d_m\}, & (2b) \\
 & & \sum_{j \in V} f_{ij}^m(t) - \sum_{j \in V} f_{ji}^m(t) & = r_m(t) & \forall m, t \text{ and } i = s_m, & (2c) \\
 & & \sum_{j \in V} f_{ij}^m(t) - \sum_{j \in V} f_{ji}^m(t) & = -r_m(t) & \forall m, t \text{ and } i = d_m, & (2d) \\
 & & \sum_{t=1}^{T_m} r_m(t) & = v_m & \forall m & (2e) \\
 & & r_m(t) & \geq L(p_m) & \forall m \text{ and } t < T_m, & (2f) \\
 & & r_m(t) & = 0 & \forall m \text{ and } t > T_m, & (2g) \\
 & & T_m & \leq T & \forall m, & (2h) \\
 & & r_m(t) & \leq U(p_m), & \forall m \text{ and } t, & (2i) \\
 & & f_{ij}^m(t) & \geq 0 & \forall (i, j) \in E, m, & (2j)
 \end{aligned}$$

RAP aims to minimize the maximum transfer time slot T of all data transfer requests. In this model, Constraint (2a) ensures the total assigned data rate on each link does not exceed the link capacity. Con-

straints (2b)-(2d) represent the flow conservation. Constraints (2e)-(2i) ensure that during the data transfer, the transmission rate of each request m satisfies its predefined rate upper bound and lower bound. We see that even with equal-priority case, finding such T is very challenging. To this end, we next develop a max-min fair resource allocation algorithm for this problem.

1.2 A Max-Min Fair Resource Allocation Algorithm

To present the max-min resource allocation (MFRA) algorithm, we first define the Max-Min Fair **MMF** problem for each time slot t .

$$\mathbf{MMF}(t) \quad \text{maximize} \quad z(t) \quad (3)$$

subject to

$$\sum_m f_{ij}^m(t) \leq u_{ij}(t) \quad \forall (i, j) \in E \quad (4a)$$

$$r_m(t) \geq z(t) \cdot v_m(t) \quad \forall m \in M_{unsat}(t) \quad (4b)$$

$$r_m(t) = z_{sat}^m(t) \cdot v_m(t) \quad \forall m \in M_{sat}(t) \quad (4c)$$

$$\sum_{j \in V} f_{ij}^m(t) - \sum_{j \in V} f_{ji}^m(t) \begin{cases} = r_m(t) & \text{if } i = s_m \\ = 0 & \text{if } i \in V / \{s_m, d_m\} \\ = -r_m(t) & \text{if } i = d_m \end{cases} \quad \forall m \in M(t), \quad (4d)$$

$$r_m(t) \geq \min(L(p_m), v_m - \sum_{\tau=1}^{t-1} r_m(\tau)), \quad \forall m \in M(t), \quad (4e)$$

$$r_m(t) \leq U(p_m), \quad \forall m \in M(t), \quad (4f)$$

$$z(t) \geq 0 \quad (4g)$$

Then we present the MFRA algorithm as shown in Algorithm 1. The basic idea of MFRA is to iteratively maximize the minimum data transfer satisfaction rate until all the requests are saturated, i.e., achieving maximum $z(t)$. In each iteration, newly saturated requests are removed from the subsequent process by fixing their corresponding rate value. In particular, it first solves the **MMF** problem (Line 6). Then each unsaturated transfer request m goes through the saturation test (Line 8-15). In this test, MFRA first checks for residual paths in the network for request m (Line 8). If there is no such path, a further test is performed by solving **MMF** with a new setting, in which $M_{unsat}(t)$ only contains request m , and all other requests which were in $M_{unsat}(t)$ are moved to $M_{sat}(t)$ with a satisfaction rate $z(t)$ (Line 9). If the new $z^{temp}(t)$ is the same as $z(t)$, it means that the request m is indeed saturated. We then update both $M_{sat}(t)$ and $M_{unsat}(t)$ (Line 10-14). After all data transfers are saturated in time slot t , we can derive a feasible set of $f_{ij}^m(t)$ for each transfer request m , and the corresponding $r_m(t)$ (Line 22-23). If the total saturation rate of a transfer request m is 1, it means that this request has been fulfilled, and we can remove it from the unfinished set of requests (Line 23-25), which will be pushed for next time slots. In this way, MFRA achieves a balance between fairness of data transfers and network utilization. Furthermore, it has a low computational complexity.

Note: Here we allocate not only bandwidth, i.e., scheduling, but also routes to different transfer requests. However, MFRA can certainly handle the case where particular routing constraints are specified, e.g., all routes are fixed ahead or each transfer request only uses one single path in each time slot, by adding an additional set of linear constraints in the formulation of **MMF** problem.

Algorithm 1 MFRA: a Max-Min Fair Resource Allocation Algorithm for the RAP problem

```

1:  $t \leftarrow 0$ 
2: while  $\mathbb{M} \neq \emptyset$  do
3:    $t \leftarrow t + 1$ 
4:    $M_{sat}(t) \leftarrow \emptyset, M_{unsat}(t) \leftarrow \mathbb{M}_p$ 
5:   while  $M_{unsat}(t) \neq \emptyset$  do
6:     solve MMF( $t$ ) and get the optimal value  $z(t)$ 
7:     for each  $m \in M_{unsat}(t)$  do
8:       if there is no residual path for transfer request  $m$  to accommodate more flows then
9:         solve MMF( $t$ ) with revised inputs  $M_{unsat}^{temp}(t) \leftarrow m, M_{sat}^{temp}(t) \leftarrow M_{sat}(t) \cup M_{unsat}(t)/\{m\}$  and
           $z_{sat}^m(t) \leftarrow z(t)$  for any  $m \in M_{unsat}(t)/\{m\}$ , and get the optimal value  $z^{temp}(t)$ 
10:        if  $z^{temp}(t) == z(t)$  then
11:           $M_{sat}(t) \leftarrow M_{sat}(t) \cup \{m\}$ 
12:           $M_{unsat}(t) \leftarrow M_{unsat}(t)/\{m\}$ 
13:           $z_{sat}^m(t) \leftarrow z(t)$ 
14:        end if
15:      end if
16:    end for
17:  end while
18:  for each  $m \in M_{sat}(t)$  do
19:    if  $\sum_{t_0=1}^t z_{sat}^m(t_0) \geq 1$  then
20:       $z_{sat}^m(t) \leftarrow 1 - \sum_{t_0=1}^{t-1} z_{sat}^m(t_0)$ 
21:    end if
22:    Derive a feasible set of  $f_{ij}^m(t)$  and the corresponding  $r_m(t)$  for each data transfer request  $m \in \mathbb{M}(t)$  based
    on constraints (4a)-(4g)
23:    if  $\sum_{t_0=1}^t z_{sat}^m(t_0) == 1$  then
24:       $\mathbb{M} \leftarrow \mathbb{M}/\{m\}$ 
25:    end if
26:  end for
27: end while
28:  $T \leftarrow t$ 
29: return  $T$ 

```
