

# Real-time GDA

## ABSTRACT

### ACM Reference format:

. 2016. Real-time GDA. In *Proceedings of ACM Conference, Washington, DC, USA, July 2017 (Conference'17)*, 3 pages.  
DOI: 10.475/123.4

## 1 MODEL

Let  $\mathcal{L}$  be the set of sites that holds data and runs tasks. For each inter-site WAN link, let  $B_{ij}$  be the bandwidth from site  $i \in \mathcal{L}$  to site  $j \in \mathcal{L}$ . We assume that the bandwidths are stable within the time frame of doing real-time data analytics.

### 1.1 Single MapReduce Query

We perform the map tasks at the sites that contain the associated data and denote  $D_i$  as the output data from all of the map tasks at site  $i$ . The fraction of reduce tasks assigned to site  $j$  is denoted as  $r_j$  which is also the fraction of all other sites' data that must be transferred through the WAN to  $j$ . This means that the total WAN usage is for a given task distribution  $r$  is:

$$\sum_i \sum_{j \neq i} D_i r_j \quad (1)$$

which can be rearranged to:

$$\begin{aligned} \sum_i \sum_{j \neq i} D_i r_j &= \sum_i D_i \sum_j r_j - \sum_j D_j r_j = \sum_i D_i - \sum_j D_j r_j \\ &= \sum_j D_j (1 - r_j) \end{aligned} \quad (2) \quad (3)$$

If we are given a start time  $s$  after the map steps are all completed and a data shuffle deadline  $t$  for the reduce tasks, then the completion time is bounded as such:

$$s + \max_{(i,j): j \neq i} \left\{ \frac{D_i r_j}{B_{ij}} \right\} \leq t \quad (4)$$

which is caused by the heterogeneous WAN bandwidth and has a bottlenecking link(s). We assume that  $t > s$ .

### 1.2 Single DAG Query

Let  $\mathcal{N}$  be the set of nodes and  $\mathcal{D}$  be the directed set of edges respectively, for the DAG of a Query Execution Plan (QEP). We define  $D_i^k$  as the amount of input data for node  $k$  at site  $i$ , and  $r_j^k$  as the fraction of tasks for node  $k$  assigned to site  $j$ . We denote  $\alpha^k$  as the selectivity for node  $k$  which is the proportional data shrinkage/expansion from placed tasks at each of the sites. For each node  $l$  in the DAG that has predecessors, the sum of the predecessor nodes' output data becomes the input data:

$$\sum_{j \in \mathcal{L}} D_j^l = \sum_{k: (k,l) \in \mathcal{D}} \alpha^k \sum_{i \in \mathcal{L}} D_i^k. \quad (5)$$

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

Conference'17, Washington, DC, USA

© 2016 Copyright held by the owner/author(s). 123-4567-24-567/08/06...\$15.00  
DOI: 10.475/123.4

We assume that each node in the QEP requires a data shuffle which means that

$$D_j^l = \sum_{k: (k,l) \in \mathcal{D}} \alpha^k r_j^k \sum_{i \in \mathcal{L}} D_i^k. \quad (6)$$

We also have that the total WAN usage is:

$$\sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{L}} \sum_{j \in \mathcal{L}: j \neq i} D_i^k r_j^k \quad (7)$$

which can like previously be rearranged to

$$\sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{L}} \sum_{j \in \mathcal{L}: j \neq i} D_i^k r_j^k = \sum_{k \in \mathcal{N}} \left( \sum_{i \in \mathcal{L}} D_i^k - \sum_{j \in \mathcal{L}} D_j^k r_j^k \right) \quad (8)$$

$$= \sum_{k \in \mathcal{N}} \sum_{j \in \mathcal{L}} D_j^k (1 - r_j^k) \quad (9)$$

## 2 PROBLEM FORMULATIONS

### 2.1 Single MapReduce Query

*Minimize WAN usage.* When minimizing WAN usage for a MapReduce data shuffle we have the following optimization problem:

$$\min_r \sum_j D_j (1 - r_j) \quad (10a)$$

$$\text{s.t.} \quad \sum_j r_j = 1 \quad (10b)$$

$$r_j \geq 0 \quad \forall j \quad (10c)$$

*Feasibility to meet deadline.* We want to find a feasible  $r$  so that the data shuffle finishes at or before  $t$ :

$$\text{find } r \text{ s.t.} \quad \frac{D_i r_j}{B_{ij}} \leq t - s \quad \forall (i,j) : j \neq i \quad (11a)$$

$$\sum_j r_j = 1 \quad (11b)$$

$$r_j \geq 0 \quad \forall j \quad (11c)$$

*Minimize WAN usage given a shuffle deadline.* When minimizing WAN usage for a MapReduce data shuffle for a given deadline  $t$  we have the following optimization problem:

$$\min_r \sum_j D_j (1 - r_j) \quad (12a)$$

$$\text{s.t.} \quad \frac{D_i r_j}{B_{ij}} \leq t - s \quad \forall (i,j) : i \neq j \quad (12b)$$

$$\sum_j r_j = 1 \quad (12c)$$

$$r_j \geq 0 \quad \forall j \quad (12d)$$

## 2.2 Single DAG Query

*Minimize WAN usage.* When minimizing WAN usage for a DAG we have the following optimization problem:

$$\min_{r, D} \sum_{k \in \mathcal{N}} \sum_{j \in \mathcal{L}} D_j^k (1 - r_j^k) \quad (13a)$$

$$\text{s.t.} \quad \sum_j r_j^k = 1 \quad \forall k \in \mathcal{N} \quad (13b)$$

$$r_j^k \geq 0 \quad \forall j \in \mathcal{L}, \forall k \in \mathcal{N} \quad (13c)$$

$$D_j^l = \sum_{k: (k, l) \in \mathcal{D}} \alpha^k r_j^k \sum_{i \in \mathcal{L}} D_i^k \quad \forall l \in \mathcal{N} \quad (13d)$$

Even without adding deadlines, this problem has nonlinear terms in the objective function and the last equality constraint.

## 3 OPTIMA

### 3.1 Single MapReduce Query with a deadline

Since Problem (12) is a convex optimization problem (linear program), we look at the Karush-Kuhn-Tucker (KKT) conditions for optimality using the following dual variables  $(\mu_{ij}, \theta, \lambda_j) : \forall (i, j), i \neq j$ :

$$-D_j + \sum_{i \neq j} \frac{D_i}{B_{ij}} \mu_{ij} + \theta - \lambda_j = 0 \quad \forall j \quad (14a)$$

$$\mu_{ij} \left( \frac{D_i}{B_{ij}} r_j - (t - s) \right) = 0 \quad \forall (i, j) : i \neq j \quad (14b)$$

$$r_j \lambda_j = 0 \quad \forall j \quad (14c)$$

$$\mu_{ij} \geq 0 \quad \forall (i, j) : i \neq j \quad (14d)$$

$$\lambda_j \geq 0 \quad \forall j \quad (14e)$$

$$\frac{D_i}{B_{ij}} r_j - (t - s) \leq 0 \quad \forall (i, j) : i \neq j \quad (14f)$$

$$\sum_j r_j - 1 = 0 \quad (14g)$$

$$r_j \geq 0 \quad \forall j \quad (14h)$$

where (14a) are the first stationary conditions, (14b) (14c) are the complimentary slackness conditions, (14d) (14e) are the dual feasibility conditions, and (14f) (14g) (14h) are the primal feasibility conditions.

Notice that for every site  $j$ , (14f) implies that there could be bottlenecking link if  $\max_{i \neq j} \{D_i/B_{ij}\} > t - s$ . There would not be a bottlenecking link only if making  $r_j = 1$  meets the deadline. Let us denote

$$U_j := \max_{i \neq j} \{D_i/B_{ij}\} \quad (15)$$

as the maximum upload time for site  $j$  if  $r_j = 1$ , and denote

$$b_j := \arg \max_{i \neq j} \{D_i/B_{ij}\} \quad (16)$$

as the set of bottlenecking data sources. If for any  $i \notin b_j$ , then (14f) is satisfied if it is satisfied for  $b_j$  and also  $D_i/B_{ij} < t - s$ . This fact and (14b) implies that  $\mu_{ij} = 0$  if  $i \notin b_j$ . Now the KKT conditions (14)

have redundant conditions that can be eliminated to:

$$-D_j + U_j \sum_{i \in b_j} \mu_{ij} + \theta - \lambda_j = 0 \quad \forall j \quad (17a)$$

$$\mu_{ij} (U_j r_j - (t - s)) = 0 \quad \forall (i, j) : i \in b_j \quad (17b)$$

$$r_j \lambda_j = 0 \quad \forall j \quad (17c)$$

$$\mu_{ij} \geq 0 \quad \forall (i, j) : i \in b_j \quad (17d)$$

$$\lambda_j \geq 0 \quad \forall j \quad (17e)$$

$$U_j r_j - (t - s) \leq 0 \quad \forall (i, j) : i \in b_j \quad (17f)$$

$$\sum_j r_j - 1 = 0 \quad (17g)$$

$$r_j \geq 0 \quad \forall j \quad (17h)$$

We have three cases for each  $r_j$ :

- (1)  $r_j = (t - s)/U_j$ . Since  $t - s > 0$  and  $U_j > 0$ , then  $r_j > 0$ . Then (17c) results in  $\lambda_j = 0$ . In this case (17a) turns into:

$$\theta = D_j - U_j \sum_{i \in b_j} \mu_{ij} \quad (18)$$

Since  $\sum_{i \in b_j} \mu_{ij} \geq 0$  from (17d), then this means that  $\theta \leq D_j$ .

- (2)  $0 < r_j < (t - s)/U_j$ . Then (17c) results in  $\lambda_j = 0$  and (17b) results in  $\mu_{ij} = 0 : \forall i \in b_j$ . In this case (17a) turns into:

$$\theta = D_j \quad (19)$$

- (3)  $r_j = 0$ . Then (17b) results in  $\mu_{ij} = 0 : \forall i \in b_j$ . In this case (17a) turns into:

$$\theta = D_j + \lambda_j \quad (20)$$

Since (17e), then  $\theta \geq D_j$ .

Since  $\theta$  can only be a single value, the above cases give a natural ordering. For a particular  $\theta$ : if  $D_j > \theta$ , then Case 1 applies; if  $D_j < \theta$  then Case 3 applies; if  $D_j = \theta$  then Cases 1, 2, or 3 could apply. This also means that we can restrict  $\theta$  to the set  $\{D_j : \forall j\}$  without restricting the solution space of the task placements  $r_j : \forall j$ .

Also, if all sites are in Case 1 and we observe that from (17g)  $\sum_j \frac{1}{U_j} < \frac{1}{t-s}$ , then the deadline  $t$  is too soon and the problem is infeasible. Let us denote the lower bound on  $t$  as:

$$\underline{t} := s + \frac{1}{\sum_j \frac{1}{U_j}} \quad (21)$$

On the other hand, let  $l := \arg \max_j \{D_j\}$  and if  $U_l < t - s$ , then  $r_l = 1$  and  $r_j = 0 : \forall j \neq l$ . Let us denote the upper bound on  $t$  as:

$$\bar{t} := s + U_l \quad (22)$$

If  $t \geq \bar{t}$ , then  $r_l = 1, r_j = 0 : \forall j \neq l$ , and the WAN usage is  $\sum_i D_i - D_l$ .

If  $t = \underline{t}$ , then  $r_j = (\underline{t} - s)/U_j : \forall j$  and the WAN usage is  $\sum_i D_i - (\underline{t} - s) \sum_j (D_j/U_j)$ .

Note that if either the maximum deadline range

$$\bar{t} - \underline{t} = U_l - \frac{1}{\sum_j \frac{1}{U_j}} \quad (23)$$

$$D_l - \frac{1}{\sum_j \frac{1}{U_j}} \sum_j \frac{D_j}{U_j} \quad (24)$$

or the maximum WAN savings

$$D_l - \frac{1}{\sum_j \frac{1}{U_j}} \sum_j \frac{D_j}{U_j} \quad (25)$$

has significant cost, then developing an optimization algorithm is worthwhile.

We have the following algorithm:

- (1) **Initialize:**
  - Order  $D_j$  in descending order and relabel the indexes for this ordering.
  - Set  $y := 1$ ,  $k := 1$ , and  $r_j := 0 : \forall j$
  - For each site  $j$ , set:  
 $U_j := \max_{i \neq j} \{D_i/B_{ij}\}$   
 $b_j := \arg \max_{i \neq j} \{D_i/B_{ij}\}.$
- (2) **Test for feasibility:**
  - If  $t \leq s + \frac{1}{\sum_j \frac{1}{U_j}}$ , then stop because the problem is infeasible.
- (3) **Process:**
  - Replace  $r_k := \min\{y, (t - s)/U_k\}.$
  - Update  $y := y - r_k$  and  $k := k + 1.$
  - If  $y \leq 0$  or  $k > |\mathcal{L}|$ , then stop and output  $r_j : \forall j.$   
 Otherwise repeat Step (3).

## APPENDIX