

Real-time GDA

ABSTRACT

ACM Reference format:

. 2016. Real-time GDA. In *Proceedings of ACM Conference, Washington, DC, USA, July 2017 (Conference'17)*, 4 pages.
DOI: 10.475/123_4

1 SINGLE DAG QUERY

1.1 Model

Wide Area Network. Let \mathcal{S} be the set of sites that hold data and run tasks, and \mathcal{L} be the set of directed edges that represent inter-site links. For each inter-site WAN link $(i, j) \in \mathcal{L}$, let B_{ij} be the bandwidth and C_{ij} be the cost to transfer one unit of data from site i to site j .

ASSUMPTION 1. *The bandwidths are stable within the time frame of real-time data analytics.*

For the analysis we make the following assumption:

ASSUMPTION 2. *Data transfers on a particular link are non-overlapping. (To do: Show that this assumption does not add suboptimality to our objective, i.e., any overlapping schedule can be transformed into a non-overlapping schedule without any loss to the optimal value.)*

DAG of tasks. We define a stage in the DAG as a group of tasks that have the same input data dependencies. Let \mathcal{T} be the set of stages and \mathcal{D} be the directed set of edges that represent stage dependencies for the DAG, respectively. Some of the stages do not have any incoming edges and so represent the raw input data. Let $\mathcal{R} \subset \mathcal{T}$ be the set of stages for the raw input data. Each stage dependency $(k, l) \in \mathcal{D}$ also has a corresponding amount of data D^{kl} that must be transferred from stage k to stage l .

ASSUMPTION 3. *A stage must have completely received all of its input data before it can process and start transferring data to another stage.*

The DAG also has a start time T_0 and a finish time T_f for which the schedule of data transfers must respect.

1.2 Task Assignment Problem

1.2.1 Problem Statement. We model each stage as a distributable group of tasks that may be fractionally assigned to multiple sites. The fraction of tasks from stage k placed at site i is represented by the nonnegative decision variable x_i^k and proportionally determines the size of data transfers. This means that for any directed edge $(k, l) \in \mathcal{D}$, then $x_i^k D^{kl} x_j^l$ of data will be transferred across link $(i, j) \in \mathcal{L}$. Note that the set of stages \mathcal{R} which represent the raw data have a decision variable which is preset according to its site-wise distribution.

We can represent the duration of a data transfer by $x_i^k (D^{kl} / B_{ij}) x_j^l$. Each data transfer is given an upper bound on its total link duration

d_{ij}^{kl} . This allows data transfers to be more generally scheduled and not need to be uninterruptible. This is so that they can be scheduled by any general scheduling algorithm (e.g. SRPT).

The goal of the task assignment problem is to fractionally distribute the stages of tasks to minimize the WAN usage for a given set of link duration upper bounds:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \sum_{(k,l) \in \mathcal{D}} \sum_{(i,j) \in \mathcal{L}} C_{ij} D^{kl} x_i^k x_j^l \\ \text{s.t.} \quad & x_i^k (D^{kl} / B_{ij}) x_j^l \leq d_{ij}^{kl} \quad \forall (i, j) \in \mathcal{L}, \forall (k, l) \in \mathcal{D} \end{aligned} \quad (1a)$$

$$\sum_{i \in \mathcal{S}} x_i^k = 1 \quad \forall k \in \mathcal{T} \quad (1b)$$

$$x_i^k \geq 0 \quad \forall i \in \mathcal{S}, k \in \mathcal{T} \quad (1c)$$

1.2.2 Characterizing the optima. Although Problem (1) is not convex, we can still find properties of the optimal solution.

By taking the Lagrange dual and the dual variables (λ, μ, ϕ) corresponding with the respective constraints, we have the following first-order necessary stationary conditions for optimality:

$$\begin{aligned} & \sum_{k:(k,l) \in \mathcal{D}} \sum_{i:(i,j) \in \mathcal{L}} C_{ij} D^{kl} x_i^k + \sum_{k:(l,k) \in \mathcal{D}} \sum_{i:(j,i) \in \mathcal{L}} C_{ji} D^{lk} x_i^k \\ & + \sum_{k:(k,l) \in \mathcal{D}} \sum_{i:(i,j) \in \mathcal{L}} \lambda_{ij}^{kl} (D^{kl} / B_{ij}) x_i^k + \sum_{k:(l,k) \in \mathcal{D}} \sum_{i:(j,i) \in \mathcal{L}} \lambda_{ji}^{lk} (D^{lk} / B_{ji}) x_i^k \\ & - \mu^l - \phi_j^l = 0 \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \end{aligned}$$

Combine the 1st and 3rd summations, combine the 2nd and 4th summations:

$$\begin{aligned} & \sum_{k:(k,l) \in \mathcal{D}} \sum_{i:(i,j) \in \mathcal{L}} D^{kl} (C_{ij} + \lambda_{ij}^{kl} (1/B_{ij})) x_i^k \\ & + \sum_{k:(l,k) \in \mathcal{D}} \sum_{i:(j,i) \in \mathcal{L}} D^{lk} (C_{ji} + \lambda_{ji}^{lk} (1/B_{ji})) x_i^k \\ & - \mu^l - \phi_j^l = 0 \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \end{aligned}$$

Relabel the index in the 2nd summation of the second summation set, and flip the summation order:

$$\begin{aligned} & \sum_{i:(i,j) \in \mathcal{L}} \sum_{k:(k,l) \in \mathcal{D}} D^{kl} (C_{ij} + \lambda_{ij}^{kl} (1/B_{ij})) x_i^k \\ & + \sum_{i:(j,i) \in \mathcal{L}} \sum_{m:(l,m) \in \mathcal{D}} D^{lm} (C_{ji} + \lambda_{ji}^{lm} (1/B_{ji})) x_i^m \\ & - \mu^l - \phi_j^l = 0 \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \end{aligned}$$

Since we assume a complete bidirectional WAN network we have $\{i : (i, j) \in \mathcal{L}\} \equiv \{i : (j, i) \in \mathcal{L}\} \equiv \{i : i \in \mathcal{S} \setminus j\}$ and so can combine the outer summations:

$$\begin{aligned} & \sum_{i \in \mathcal{S} \setminus j} \sum_{k:(k,l) \in \mathcal{D}} D^{kl} (C_{ij} + \lambda_{ij}^{kl} (1/B_{ij})) x_i^k \\ & + \sum_{i \in \mathcal{S} \setminus j} \sum_{m:(l,m) \in \mathcal{D}} D^{lm} (C_{ji} + \lambda_{ji}^{lm} (1/B_{ji})) x_i^m \\ & = \mu^l + \phi_j^l \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \end{aligned} \quad (2)$$

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

The dual feasibility constraints give us:

$$\lambda_{ij}^{kl} \geq 0 \quad \forall (i, j) \in \mathcal{L}, \forall (k, l) \in \mathcal{D} \quad (3a)$$

$$\phi_j^l \geq 0 \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \quad (3b)$$

The complementary slackness constraints for the inequalities give us:

$$\lambda_{ij}^{kl} (x_i^k (D^{kl}/B_{ij})x_j^l - d_{ij}^{kl}) = 0 \quad \forall (i, j) \in \mathcal{L}, \forall (k, l) \in \mathcal{D} \quad (4a)$$

$$\phi_j^l x_j^l = 0 \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \quad (4b)$$

Therefore the necessary conditions for optimality are (1a) (1b) (1c) (2) (3) (4).

Inferences from the conditions

- (1) At the optimal solution, the dual variables λ correspond with the negative gradient of the objective function w.r.t. the allotted durations \mathbf{d} . (See [1] Proposition 3.3.3)
- (2) At the optimal solution, μ^l is the total WAN cost plus the product of dual prices and allotted durations for all incoming and outgoing data transfers associated with stage l . Also each fraction of stage l at site j must share that same fraction of that stage's total cost. Multiply (2) by x_j^l and substitute with (4):

$$\begin{aligned} & \sum_{i \in \mathcal{S} \setminus j} \sum_{k: (k, l) \in \mathcal{D}} (C_{ij} D^{kl} x_i^k x_j^l + \lambda_{ij}^{kl} d_{ij}^{kl}) \\ & + \sum_{i \in \mathcal{S}} \sum_{m: (l, m) \in \mathcal{D}} (C_{ji} D^{lm} x_j^l x_i^m + \lambda_{ji}^{lm} d_{ji}^{lm}) \\ & = x_j^l \mu^l \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \end{aligned} \quad (5)$$

Then sum for all $j \in \mathcal{S}$:

$$\begin{aligned} & \sum_{j \in \mathcal{S}} \sum_{i \in \mathcal{S} \setminus j} \sum_{k: (k, l) \in \mathcal{D}} (C_{ij} D^{kl} x_i^k x_j^l + \lambda_{ij}^{kl} d_{ij}^{kl}) \\ & + \sum_{j \in \mathcal{S}} \sum_{i \in \mathcal{S}} \sum_{m: (l, m) \in \mathcal{D}} (C_{ji} D^{lm} x_j^l x_i^m + \lambda_{ji}^{lm} d_{ji}^{lm}) \\ & = \sum_{j \in \mathcal{S}} x_j^l \mu^l \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \end{aligned}$$

Apply (1b) to the RHS:

$$\begin{aligned} & \sum_{j \in \mathcal{S}} \sum_{i \in \mathcal{S} \setminus j} \sum_{k: (k, l) \in \mathcal{D}} (C_{ij} D^{kl} x_i^k x_j^l + \lambda_{ij}^{kl} d_{ij}^{kl}) \\ & + \sum_{j \in \mathcal{S}} \sum_{i \in \mathcal{S}} \sum_{m: (l, m) \in \mathcal{D}} (C_{ji} D^{lm} x_j^l x_i^m + \lambda_{ji}^{lm} d_{ji}^{lm}) \\ & = \mu^l \quad \forall j \in \mathcal{S}, \forall l \in \mathcal{T} \end{aligned} \quad (6)$$

1.3 Data Transfer Scheduling Problem

Because of the dependencies specified by the DAG, each data transfer duration link allocation d_{ij}^{kl} must be scheduled after the deadline of stage k denoted t^k and before the deadline of stage l denoted t^l . Therefore the goal is to find a scheduling policy π and stage deadlines \mathbf{t} such that the query start time T_0 and deadline T_f are satisfied for the given set of durations \mathbf{d} and stage dependencies \mathcal{D} .

1.4 Jointly optimization

Both the Task Assignment and Data Transfer Scheduling problems are connected through the data transfer duration link allocations \mathbf{d} . Therefore the joint optimization problem becomes if we can find a

feasible schedule that minimizes the WAN usage cost by adjusting \mathbf{d} .

2 CLARINET

2.1 Given claims

The paper's introduction gives the following claims which were interrogated:

- (1) Formulating an optimal solution for a multi-query, network-aware, joint query planning/placement/scheduling is computationally intractable.
 - ✓ Cites two sources that even for a single query with the given scheduling assumptions is NP-hard.
- (2) CLARINET does joint multi-query planning. It picks the best WAN-aware DAG/task placement/task scheduling per query and provides "hints" (location and start times of each task) to the query execution layer. It optimizes for minimal average query completion times. Task placement and task scheduling within a query are not done jointly.
 - ✓ Joint multi-query planning although not optimal because of NP-hardness.
- (3) They show how to (heuristically) compute the WAN-optimal DAG for a single query which includes task placement and task scheduling. The solution relies on reserving WAN links. This is an effective heuristic for the best single-query DAG, that decouples placement and scheduling.
 - ✓ Task placement and scheduling are decoupled.
 - ✗ But not a jointly optimal decomposition which is NP-hard.
 - ✗ Task placement is minimizes time usage of WAN with a linear program, not completion time which would become a piecewise linear program.
 - ✓ Given a task placement for a single query, the scheduling problem gives the optimal (non-interruptible) schedule to minimize that query's completion time.
- (4) They allow for cross-query (heuristic) optimization of n queries. They order the queries by optimal DAG expected completion time. Then they choose the i th query's DAG considering the WAN impact of the $i - 1$ preceding queries.
 - ✓ This mimics the heuristic rule Shortest Job First (SJF).
- (5) They heuristically compact the schedules tightly in time by considering the above order and groups of $k \leq n$ queries to combat resource fragmentation.
 - ✓ This is done at a cost to average completion time.
- (6) They extend the above heuristic to accommodate (i) fair treatment of queries, (ii) minimize WAN bandwidth costs, and (iii) online query arrivals.
 - ✓ This is done, but only heuristically.

2.2 Model

There are n queries and each query j has a set DAG-Set QS_j from which one DAG must be chosen. Also chosen are the task locations and task start times.

DAG assumptions:

- (1) Each node in the DAG represents a stage which is a group of tasks.
- (2) For the task placement decision, each stage is assumed to be either a set of Map tasks or a set of Reduce tasks. (Josh says: This seems pretty restrictive.)
 - (i) Each Map stage is placed at the site holding its input data.
 - (ii) Each Reduce stage can be distributed to all sites and is represented by a continuous fraction r_j , i.e. the fraction of

reduce tasks for the stage at site j . (Josh says: I think this ability to geographically distribute tasks within a stage and the continuous decision simplification is the reason why they chose to model only Map and Reduce tasks.)

- (3) Each stage in the DAG has a known amount of output data D_j at each site j .

Scheduling assumptions:

- (1) Network transfers do not overlap. In fact, they have a theorem that proves that any optimal schedule has an equivalent non-overlapped schedule. This avoids having to decide concurrent flows.
- (2) Obtain non-interruptible transfer schedules.
- (3) The compute phase of a task can only start after all its inputs are available to it at its site.
- (4) Bandwidth B_{ij} between sites i and j is known and constant.

2.3 Single Query Problem

The goal is to minimize the query execution time. First, the tasks are placed across the sites. Then, the data transfers are scheduled.

2.3.1 Task placement within a (Reduce) stage. Given a DAG, each stage's task placement is decided one-by-one.

Within each stage (single node of the DAG), the decision of (Reduce) task placement is *decided* independently from other stages but *have* a dependency on each other through a reserved amount of time τ_{ij} for the link between sites i and j . This value is the length of time already reserved to use this link by other stages that do not have a DAG ordering with this stage.

The following linear program is given to decide the distribution of reduce tasks $\{r_j\}$ in a particular stage that respects the current reservation status of each link:

$$\min_r \sum_j \sum_{i \neq j} \left(\frac{D_i r_j}{B_{ij}} + \tau_{ij} \right) \quad (7a)$$

$$\text{s.t.} \quad \sum_j r_j = 1 \quad (7b)$$

$$r_j \geq 0 \quad \forall j \quad (7c)$$

After $\{r_j\}$ is decided, each τ_{ij} is incremented by $\frac{D_i r_j}{B_{ij}}$.

(Josh says: This formulation is a bit strange. First, given the above formulation τ_{ij} is a constant, and so has no effect on the linear programming problem. Second, the summation is minimizing total time usage of the WAN and *not* stage execution time. To minimize stage execution time, it should be the max operator and not the summation operator. But then the order in which these problems are solved will matter. However, there is no discussion in order in which these problems are solved since it doesn't matter with the summation operator being used.)

After deciding the task placement for each stage, task scheduling for the DAG is decided.

2.3.2 Scheduling tasks in a DAG. Each DAG gives a partial ordering of the stages. For stages that compete over the same network links but do not have an ordering, an ordering must be decided.

The DAG is augmented in two ways. First, each data transfer is made into a node that is placed between two the associated stages. Second, tasks from the same stage at the same site are coalesced into a sub-stage. From here on, only the start times s of the sub-stages will be scheduled.

For the explicitly stated order in the DAG and given two (data transfer) sub-stages u and v in a DAG, if u must finish before v then

we have the following constraint:

$$s(v) \geq s(u) + d(u) \quad (8)$$

where $s(u)$ is the start time and $d(u)$ is the duration of stage u .

To decide the order between sub-stages that don't have an ordering but need to have a non-overlapping schedule. We add the following constraints to decide whether sub-stage u should be complete before v or vice versa:

$$s(v) \geq s(u) + d(u) - N(1 - z_{uv}) \quad (9a)$$

$$s(u) \geq s(v) + d(v) - N(z_{uv}) \quad (9b)$$

where z_{uv} is a binary decision variable that indicates if v is executed after u , and N is a large constant.

With the above scheduling constraints the i -th DAG, we measure the completion time of the query as the last finished task of the query:

$$\Phi^i := \max_{u \in \text{DAG}^i} s(u) + d(u) \quad (10a)$$

$$\text{s.t. (8)(9)} \quad (10b)$$

which is a binary integer linear program. (Josh says: This binary integer linear program may not scale well for larger DAGs. Can we do better?)

2.4 Multiple Query Problem

Essentially, multiple queries are handled as processing the Shortest Job First (SJF). However, CLARINET makes some modifications (see below).

2.4.1 Handling previously scheduled queries. Suppose $\text{low}(b)$ and $\text{high}(b)$ represents the start and end times of a reservation on a link for a previously scheduled query. Then we add the following two constraints to the scheduling problem (10) when scheduling other queries:

$$s(u) \geq \text{high}(b) - N(1 - x_{ub}) \quad (11a)$$

$$\text{low}(b) \geq s(u) + d(u) - N(x_{ub}) \quad (11b)$$

where x_{ub} is a binary indicator denoting that u is scheduled after interval b . The constraints work similar to (9).

2.4.2 Handling resource fragmentation. As queries get scheduled, resources (links) become fragmented from reservations happening without accounting for other queries. They present a heuristic that optimally packs subsets of queries of size $k \leq n$.

First, they calculate the execution time for each query as if the query was scheduled alone. Based on this, the queries are put in an increasing order.

The first k queries are then scheduled. While sweeping through time, each sub-stage is scheduled from this set when it and its link first become free. When a query has been completely scheduled, it leaves the set and the next shortest query left in original ordering replaces it. This continues until all queries have been scheduled.

This heuristic allows query link usages to become better packed. However, at high values of k , the average completion time may rise which means that there is a tradeoff between efficiently using the WAN and query completion time among multiple queries.

2.4.3 Fairness. The multi-query heuristics above favor shorter queries over longer queries. To combat this, they target that each query should finish within a specific multiple m of its stand alone completion time dur_j .

At time t when a query is to be added to the query set (of size k) that is being scheduled, each of the unscheduled queries is given

the following score:

$$\text{Prox}_j(t) = 1 - \frac{n \times \text{dur}_j - t}{n \times \text{dur}_j} \quad (12)$$

and then sorted in descending order. From the first H of them in this order, the one with the shortest stand alone execution time is added to the query set being scheduled. This heuristic favors queries that are closer to their “fair” deadline.

2.4.4 WAN utilization. To limit WAN usage, they only select DAGs that have a WAN usage below a threshold β . This strictly limits each query’s WAN usage.

2.4.5 Online Arrivals. When a new query arrives, update and recompute for all queries that have yet to have any tasks start to be executed.

2.5 Possible directions

(Josh says: In thinking about our own problem to solve, I think our main focus could be on:

- (1) CLARINET assumes either a Map or Reduce structure for each stage of the DAG when deciding task placement. Can we extend this for a more general DAG stage type?
- (2) CLARINET only changes the schedules for non-executed queries when a new query arrives. Can we also change partially executed queries to make the system more flexible and in return experience a switching cost?
- (3) CLARINET focuses exclusively on minimizing the average query completion time. This may not be fair, and even with its fairness heuristic, it still favors shorter queries. Can we schedule queries’ usage of the WAN to be more fair?
- (4) CLARINET’s task scheduling problem requires the use of a binary integer linear program. This may be intractable for large DAGs. Can we design a faster algorithm or formulate a more tractable problem?

)

REFERENCES

- [1] Dimitri P Bertsekas. 1999. *Nonlinear programming*. Athena scientific Belmont.

APPENDIX