

# Model Draft

Yu Wu\*, Chuan Wu\*, Francis C.M. Lau\*

\*Department of Computer Science, The University of Hong Kong,

Email: {ywu,cwu,fcmlau}@cs.hku.hk

## Abstract

### I. MODEL OVERVIEW

The model consists of a set of geographically diverse cloud clusters  $F$ , a set of videos  $O$  and a set of clients  $D$ . ( $D_f$  denotes the consolidated viewer group within domain of cloud cluster  $f$  ( $f \in F$ ).) Without loss of generality, I will first assume all the videos have unit length. (Extend this when the model is done)

#### A. Alphabet Soup

- 1) Each cloud cluster is assigned a storage capacity,  $S_f$ .
- 2) Each cloud cluster has a bandwidth capacity,  $\mu_f$ . Here I plan to borrow the idea from our ICDCS paper's model *i.e.*, bandwidth will be abstracted into a VM instance, which will actually provide the bandwidth.
- 3)  $x_{jf}^{(o)}$  denotes the variable indicating whether the request for video  $o$  issued from viewer  $j$  will be directed to cloud cluster  $f$ .
- 4)  $y_f^{(o)}$  denotes the variable indicating whether to store a copy of video  $o$  at the cloud cluster  $f$ .
- 5)  $c_f$  denotes the storage cost of cloud cluster  $f$ .
- 6)  $v_f$  denotes the transferring cost from viewer  $j$ 's location to cloud cluster  $f$ .
- 7)  $R_{jf}$  denotes the transferring latency from viewer  $j$ 's location to cloud cluster  $f$ . It can be assigned with value of RTT between these two geographical regions.

There should be a mapping function to map user  $j$  to a location  $f$ . For simplicity, we can denote it as  $D^{-1}(j)$ . In that way, we can represent  $x_{jf}^{(o)}$ ,  $c_{jf}$  and  $R_{jf}$  as  $x_{D^{-1}(j)f}^{(o)}$ ,  $c_{D^{-1}(j)f}$  and  $R_{D^{-1}(j)f}$  respectively. For clearness, I will use  $j$  afterwards.

#### B. Objective function

To minimize the operation cost, based on the premise that the expected average global latency should be bounded below some tolerant value.

$$\min \sum_{o \in O} \sum_{f \in F} y_f^{(o)} \times c_f + \sum_{o \in O} \sum_{f \in F} \sum_{j \in D_f} x_{jf}^{(o)} \times v_f$$

#### C. Constraints

- 1) *Storage*  
 $\sum_{o \in O} y_f^{(o)} \leq S_f, \forall f \in F$
- 2) *VM capacity*  
 $\sum_{o \in O} \sum_{j \in D} x_{jf}^{(o)} \leq \mu_f, \forall f \in F$

3) *Placement*

$$\sum_{f \in F} y_f^{(o)} \geq 1, \forall o \in O$$

4) *Latency guarantee*

$$\frac{\sum_{o \in O} \sum_{f \in F} \sum_{j \in D_f} x_{jf}^{(o)} \times R_{jf}}{|D|} \leq R_{threshold}, \text{ where } R_{threshold} \text{ is an input into the system.}$$

5) *Variable constraint*

$$y_f^{(o)} \in \{0, 1\}, x_{jf}^{(o)} \in \{0, 1\}$$

6) *Variable constraint*

$$x_{jf}^{(o)} \leq y_f^{(o)}, \forall j \in D$$

## II. ALTERNATIVE LP (RELAXATION)

Obviously, the optimization in Sec. I is an integer problem. Here we want to make an intuitive relaxation. The reason is two-fold. First, the number of users makes the optimization problem too large to solve. Second, we want to transform the original problem into a more tractable one.

### A. Consolidate users

As what we have assumed, at any time, each user can at most view one video. So we can treat the users within one specific region  $f$  ( $f \in F$ ) as one, which will make our optimization much slimmer. Based on that, we are able to eliminate all the variables  $x_{jf}^{(o)}$  and consolidates them as one viewer. Suppose the total user set at time slot  $T$  is represented as  $D(T)$ , so the viewer set in region  $f$  at that time slot is  $D_f(T)$ . We introduce a new variable  $\alpha_{jf}^{(o)}$ , which denote the portion of request for video  $o$  issued from the aggregate user  $j$  to cloud cluster  $f$ . To note that,  $\alpha_{jf}^{(o)}$  is a fractional variable. Due to our charging mode, the storage deployment is scheduled at a larger time scale while the VM rental is done at a smaller one, i.e.  $T$ . From above, the original ILP has only one type of integer variable  $y_f^{(o)}$  and the original optimization problem is changed to,

$$\min \sum_{o \in O} \sum_{f \in F} y_f^{(o)} \times c_f + \sum_{o \in O} \sum_{f \in F} \sum_{j \in F} \alpha_{jf}^{(o)} \times D_f(T) \times v_f \quad (1)$$

### B. Constraints

1) *Storage*

$$\sum_{o \in O} y_f^{(o)} \leq S_f, \forall f \in F$$

2) *VM capacity*

$$\sum_{o \in O} \sum_{j \in F} \alpha_{jf}^{(o)} \times D_f(T) \leq \mu_f, \forall f \in F$$

3) *Placement*

$$\sum_{f \in F} y_f^{(o)} \geq 1, \forall o \in O$$

4) *Latency guarantee*

$$\frac{\sum_{o \in O} \sum_{f \in F} \sum_{j \in F} \alpha_{jf}^{(o)} \times D_f(T) \times R_{jf}}{|D|} \leq R_{threshold}, \text{ where } R_{threshold} \text{ is an input into the system.}$$

5) *Variable constraint*

$$y_f^{(o)} \in \{0, 1\}, \alpha_{jf}^{(o)} \in [0, 1]$$

6) *Variable constraint*

$$\alpha_{jf}^{(o)} \leq y_f^{(o)}, \forall j, f \in F$$

7) *Variable constraint*

$$\sum_{f \in F} \alpha_{jf}^{(o)} = 1$$

### C. How to solve?

If we relax the variable  $y_f^{(o)}$ , the optimization problem in Sec. II-A is a problem with complicating constraints and we can utilize an efficient dual decomposition to solve it. More specifically, the original constraint (5) is relaxed into  $y_f^{(o)} \in [0, 1], \alpha_{jf}^{(o)} \in [0, 1]$ . The original problem can be formulated as,

$$\begin{aligned} \min & \sum_{o \in O} \sum_{f \in F} y_f^{(o)} \times c_f + \sum_{o \in O} \sum_{f \in F} \sum_{j \in F} \alpha_{jf}^{(o)} \times D_f(T) \times v_f \\ \text{s.t.} & \begin{cases} y_f^{(o)} \in \mathbb{C}_1 (\forall o \in O, f \in F) \\ \alpha_{jf}^{(o)} \in \mathbb{C}_2 (\forall o \in O, j, f \in F) \\ \alpha_{jf}^{(o)} - y_f^{(o)} \leq o (\forall o \in O, j, f \in F) \end{cases} \end{aligned} \quad (2)$$

$\mathbb{C}_1$  is the convex set defined by linear constraints (1), (3) and (5), while  $\mathbb{C}_2$  is the convex set defined by linear constraints (2), (4), (5) and (7).

1) *Dual Decomposition:* Let  $\mathbb{L}(y, \alpha, L)$  be the Lagrangian problem of Eqn. 2, where  $y$  and  $\alpha$  is column-formatted primal variables,

i.e.,  $y = (y_1^1, y_1^2, \dots, y_1^{|O|}, y_2^1, y_2^2, \dots, y_{|F|}^{|O|})^T$ .

$$\begin{aligned} \mathbb{L}(y, \alpha, L) &= \sum_{o \in O, f \in F} y_f^{(o)} \times c_f + \sum_{o \in O, f \in F, j \in F} \alpha_{jf}^{(o)} \times D_f(T) \times v_f + \sum_{o \in O, f \in F, j \in F} \lambda_{jf}^{(o)} \times (\alpha_{jf}^{(o)} - y_f^{(o)}) \\ &= [\sum_{o \in O, f \in F} y_f^{(o)} \times c_f - \sum_{o \in O, f \in F, j \in F} \lambda_{jf}^{(o)} \times y_f^{(o)}] + [\sum_{o \in O, f \in F, j \in F} \alpha_{jf}^{(o)} \times (D_f(T) \times v_f + \lambda_{jf}^{(o)})] \end{aligned} \quad (3)$$

Obviously, Eqn. 3 can be easily decomposed into two sub problems, i.e., shown as Eqn. 4

$$\begin{aligned} g_1 \lambda &= \sum_{o \in O, f \in F} y_f^{(o)} \times c_f - \sum_{o \in O, f \in F, j \in F} \lambda_{jf}^{(o)} \times y_f^{(o)} \quad (A) \\ \text{s.t. } y_f^{(o)} &\in \mathbb{C}_1 (\forall o \in O, f \in F) \\ g_2(\lambda) &= \sum_{o \in O, f \in F, j \in F} \alpha_{jf}^{(o)} \times (D_f(T) \times v_f + \lambda_{jf}^{(o)}) \quad (B) \\ \text{s.t. } \alpha_{jf}^{(o)} &\in \mathbb{C}_2 (\forall o \in O, j \in F, f \in F) \end{aligned} \quad (4)$$

Both sub problems are standard linear optimizations with efficient polynomial-time solutions. However, a subsequent rounding for sub problem (A) seems essential since only integrals (0, 1) are feasible in the original problem (Eqn. 1). We will prove that the optimal solution of (A) are integrals.

**Theorem 1.** *The optimal solutions for sub problem (A) are integrals.*

*Proof:* Since (A) is LP, the optimal solution should be a vertex. We have known that If  $A$  is *totally unimodular*, then every vertex solution of  $Ax \leq b$  is integral. So we prove the Theorem. 1 iff we can show the constraint matrix is *totally unimodular*. As mentioned, all the constraints of sub problem (A) is,

$$\begin{cases} \sum_{o \in O} y_f^{(o)} \leq S_f, \forall f \in F \\ -\sum_{f \in F} y_f^{(o)} \leq -1, \forall o \in O \\ y_f^{(o)} \leq 1, \forall o \in O, f \in F \\ -y_f^{(o)} \leq 0, \forall o \in O, f \in F \end{cases} \quad (5)$$

Thus, if we present Eqn. 5 into  $Ax \leq b$ .  $A =$

$$\left( \begin{array}{cccc|cccc|cccc}
 1 & \dots & 1 & 0 & \dots & 0 & \dots & & 0 & \dots & 0 \\
 0 & \dots & 0 & 1 & \dots & 1 & \dots & & 0 & \dots & 0 \\
 & & \vdots & & & \vdots & & & & & \\
 0 & \dots & 0 & 0 & \dots & 0 & \dots & & 1 & \dots & 1 \\
 -1 & \dots & 0 & -1 & \dots & 0 & \dots & & -1 & \dots & 0 \\
 & & \ddots & & & \ddots & & & & \ddots & \\
 & & & -1 & & -1 & & & -1 & & -1
 \end{array} \right) \left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} \begin{array}{l} |F| \\ |O| \\ |F| \times |O| \\ |F| \times |O| \end{array}$$

$$\left( \begin{array}{cccc|cccc|cccc}
 \mathbf{E} & & & & 0 & & & & \dots & & 0 \\
 & & & & & & & & & & \\
 0 & & & & \mathbf{E} & & & & \dots & & 0 \\
 & & & & & & & & & & \\
 & & & & & & & & \ddots & & \\
 & & & & & & & & & & \mathbf{E} \\
 -\mathbf{E} & & & & 0 & & & & \dots & & 0 \\
 & & & & & & & & & & \\
 0 & & & & -\mathbf{E} & & & & \dots & & 0 \\
 & & & & & & & & & & \\
 & & & & & & & & \ddots & & \\
 & & & & & & & & & & -\mathbf{E}
 \end{array} \right) \left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} \begin{array}{l} |O| \\ |O| \\ |O| \times (|F|-3) \\ |O| \end{array}$$

$$x = (y_1^1, y_1^2, \dots, y_1^{|O|}, y_2^1, y_2^2, \dots, y_2^{|O|}, \dots, y_{|F|}^{|O|})^T$$

$$b = (S_1, S_2, \dots, S_F, -1, -1, \dots, -1, 1, 1, \dots, 1, 0, 0, \dots, 0)^T$$

Since every  $2 \times 2$  sub matrix must at least contain one 0, we can easily know that the determinant of any  $2 \times 2$  sub matrix is  $\{0, -1, +1\}$ . In other words, the constraint matrix  $A$  is *totally unimodular*. ■

2) *Master Algorithm*: Like ordinary sub-gradient algorithms, an iterative algorithm (Table. I) is applied. In each cycle, with fixed Lagrangian variables  $\lambda$ , we solve those two subproblems derived in Sec. II-C1, together with sub gradients  $\partial g_1(\lambda)$  and  $\partial g_2(\lambda)$ . ( $\partial g_1(\lambda_{jf}^{(o)}) = -y_f^{(o)}$ ,  $\partial g_2(\lambda_{jf}^{(o)}) = \alpha_{jf}^{(o)}$ )

TABLE I  
MASTER ALGORITHM

Repeat
Solve subproblem $g_1(\lambda)$ over $y$
Solve subproblem $g_2(\lambda)$ over $\alpha$
Update dual variables $\lambda_{jf}^{(o)} := \lambda_{jf}^{(o)} + \beta_k \times (\alpha_{jf}^{(o)} - y_f^{(o)})$