

TABLE II  
SUMMARY OF NOTATIONS

Symbol	Meaning
$T$	The lifespan of a data object.
$T_t^e$	The lifespan of a cached replica determined in time slot $t$ .
$T_{t,com}^e$	The caching lifespan for the compensation mechanism.
$T_{t,pro}^e$	The caching lifespan for the prolonged caching mechanism.
$T_{t,laz}^e$	The caching lifespan for the lazy caching mechanism.
$O$	The size of a data object.
$E$	The number of edge regions.
$P_s^c$	The price of storage in the cloud (/GB/Month).
$P_b^c$	The price of bandwidth in the cloud (/GB).
$P_s^e$	The price of storage for edge region $e$ (/GB/Month).
$P_b^e$	The price of bandwidth edge region $e$ (/GB).
$g_t^e$	The number of requests in the cloud in time slot $t$ .
$g_t^e$	The number of requests in edge region $e$ in time slot $t$ .
$g_{[T_t^e]}^e$	The number of requests in region $e$ for $T_t^e$ time slots.
$G_t^e$	A request sequence from region $e$ , $G_t^e = [g_t^1, \dots, g_t^E]$ .
$G_t$	All request sequences, $G_t = [G_t^1, \dots, G_t^e, \dots, G_t^E]$ .
$C_s \rightarrow d(O)$	The replica creation cost from $s$ to $d$ for an object.
$C_t^c(O)$	The cost of not caching an object of size $O$ in time slot $t$ .
$C_t^e(O)$	The cost of caching an object for region $e$ in time slot $t$ .
$C_{[T_t^e]}^c(O)$	The cost of not caching an object for $T_t^e$ time slots.
$C_{[T_t^e]}^e(O)$	The cost of caching an object for $T_t^e$ time slots.
$C_{[T_t^e]}^{OPT}(O)$	The optimal cost of an object of size $O$ for $T_t^e$ time slots.
$C_{[T_t^e]}^{REC}(O)$	The cost of RECCaching for $T_t^e$ time slots.
$I_t^e$	The caching status indicator for region $e$ in time slot $t$ .
$J_t^e$	The violation status indicator for region $e$ in time slot $t$ .
$N_t^e$	The set of linked edge regions of region $e$ in time slot $t$ .
$d(t)$	The state for the optimal offline algorithm in time slot $t$ .
$\mathcal{D}$	The set of states for the optimal offline algorithm.
$\omega^e$	A coefficient of QoS penalty cost for edge region $e$ .
$a_t$	The action of the agent at time slot $t$ in RL.
$s_t$	The state of the environment at time slot $t$ in RL.
$r_t$	The reward of the environment at time slot $t$ in RL.
$S$	The set of states of RL.
$\mathcal{A}$	The set of actions of RL.
$\mathcal{R}$	The set of rewards of RL.
$\pi$	The policy of the agent in RL.
$h_T$	The number of consecutive uncached time slots.
$\epsilon$	The ratio of the storage price of the edge to the cloud.
$\beta$	The ratio of the bandwidth price of the edge to the cloud.

## APPENDIX PERFORMANCE ANALYSIS

For ease of reference, Table II summarizes the notions used in this paper.

We start by defining the relationship of prices between the cloud and the edge. For the price of storage in the cloud and at the edge, we have the following relationship

$$P_s^e = \epsilon P_s^c, \epsilon \in [0, 3). \quad (35)$$

Similarly, the price relationship of the bandwidth is

$$P_b^e = \beta P_b^c, \beta \in [0, 1]. \quad (36)$$

RECCaching is enhanced by three mechanisms that perform four kinds of decisions for edge caching. We perform competitive analysis for each decision separately [30]. In order to analyze the performance throughout the lifespan of a data object, we divide its lifespan into phases according to the caching decisions of RECCaching enhanced by the mechanisms. Then the performance of RECCaching for each phase is analyzed. For the caching decision given by Algorithm 2, the analysis is detailed in Proposition 1.

**Proposition 1.** *The caching decision of Algorithm 2 is  $(1 + \max\{\epsilon, \beta\})$ -competitive.*

*Proof.* For RECCaching, if it observes  $g_t^e$  requests and performs an action  $a_t = 1$ , then an object is decided to be cached at the edge for the next  $T_t^e$  time slots. In order to calculate the cost incurred by this action, or in other words, the cost of the next  $T_t^e$  time slots, we first denote the number of requests during  $T_t^e$  time slots as  $g_{[T_t^e]}^e$ . According to cost definitions of storage, bandwidth and replica creation detailed in Section III, the caching cost for this decision is

$$C_{[T_t^e]}^{REC}(O) \leq C_{[T_t^e]}^e(O) = P_s^c O T_t^e + P_s^e O T_t^e + P_b^e O g_{[T_t^e]}^e + \min_{e' \in N_t^e} \{C_t^{c \rightarrow e}(O), C_t^{e' \rightarrow e}(O)\}, \quad (37)$$

where  $C_{[T_t^e]}^e(O)$  is calculated as Eqn. (17). We use ' $\leq$ ' in the calculation of  $C_{[T_t^e]}^{REC}(O)$  because the actual lifespan of caching the object at the edge is less than or equal to  $T_t^e$  due to the proposed lazy caching mechanism.

For the offline optimal decision, we discuss its cost by case based on the comparison between  $g_t^e$  that is used to calculate the caching lifespan  $T_t^e$  and  $g_{[T_t^e]}^e$  that is the actual number of requests during the caching lifespan. For the case of  $g_{[T_t^e]}^e \geq g_t^e$ , we have  $C_{[T_t^e]}^c(O) \geq C_{[T_t^e]}^e(O)$  for  $g_{[T_t^e]}^e$  requests, which means that the offline optimal decision is to cache the object. The cost of the offline optimal decision is

$$C_{[T_t^e]}^{OPT}(O) \geq P_s^c O T_t^e + P_s^e O + P_b^e O g_{[T_t^e]}^e + \min_{e' \in N_t^e} \{C_t^{c \rightarrow e}(O), C_t^{e' \rightarrow e}(O)\}, \quad (38)$$

We use  $\geq$  in the calculation of  $C_{[T_t^e]}^{OPT}(O)$  because the minimum cost for  $T_t^e$  time slots occurs when all requests are served at the edge in one time slot. Thus, the competitive ratio in this case is

$$\frac{C_{[T_t^e]}^{REC}(O)}{C_{[T_t^e]}^{OPT}(O)} \leq \frac{P_s^c O T_t^e + P_s^e O T_t^e + P_b^e O g_{[T_t^e]}^e + \dots}{P_s^c O T_t^e + P_s^e O + P_b^e O g_{[T_t^e]}^e + \dots} \quad (39a)$$

$$\leq \frac{P_s^c T_t^e + P_s^e T_t^e + P_b^e g_{[T_t^e]}^e}{P_s^c T_t^e + P_s^e + P_b^e g_{[T_t^e]}^e} \quad (39b)$$

$$\leq \frac{P_s^c T_t^e + P_s^e T_t^e}{P_s^c T_t^e} \quad (39c)$$

$$= \frac{P_s^c + \epsilon P_s^c}{P_s^c} = 1 + \epsilon. \quad (39d)$$

In Eqn. (39a), we use ' $\dots$ ' in the numerator and denominator to indicate the replica creation cost. Eqn. (39b) holds by omitting the replica creation cost in both numerator and denominator with the fact  $C_{[T_t^e]}^{REC}(O) \geq C_{[T_t^e]}^{OPT}(O)$ . Eqn. (39c) holds due to  $P_s^e + P_b^e g_{[T_t^e]}^e > P_b^e g_{[T_t^e]}^e$ . We finally obtain Eqn. (39d) by substituting Eqn. (35) for  $P_s^e$ .

In the other case of  $g_{[T_t^e]}^e < g_t^e$ , we have  $C_{[T_t^e]}^c(O) < C_{[T_t^e]}^e(O)$  for  $g_{[T_t^e]}^e$  requests. Thus, the offline optimal decision is not to cache the object at the edge, which incurs a cost of storing the object and serving requests in the cloud only:

$$C_{[T_t^e]}^{OPT}(O) = C_{[T_t^e]}^c(O) \quad (40)$$

where  $C_{[T_t^e]}^c(O)$  is calculated as Eqn. (16). For RECCaching, we first discuss the competitive ratio when  $g_{[T_t^e]}^e = 0$ . According

to the proposed lazy caching mechanism, an object is cached at the edge when it receives an additional request after the caching decision. Because the number of requests  $g_{[T_t^e]}^e$  during the decided caching lifespan is 0, the object is not cached at the edge in fact. Thus, in this case, the cost of RECaching only includes the cost of storage for  $T_t^e$  time slots in the cloud, i.e.,

$$C_{[T_t^e]}^{REC}(O) = P_s^c OT_t^e. \quad (41)$$

Obviously, the cost of the offline optimal decision is

$$C_{[T_t^e]}^{OPT}(O) = P_s^c OT_t^e, \quad (42)$$

because there is no request, there is no QoS penalty. Thus, the competitive ratio of  $g_{[T_t^e]}^e = 0$  is

$$\frac{C_{[T_t^e]}^{REC}(O)}{C_{[T_t^e]}^{OPT}(O)} = 1. \quad (43)$$

Then we analyze the competitive ratio when  $0 < g_{[T_t^e]}^e < g_t^e$ . In this case, the offline optimal decision is not to cache the object at the edge, which incurs a cost of

$$C_{[T_t^e]}^{OPT}(O) = C_{[T_t^e]}^c(O) = P_s^c OT_t^e + P_b^c Og_{[T_t^e]}^e + \omega^e \left| \sum_{t=t+1}^{t+T_t^e} J_t^e(C_t^e(O) - C_t^c(O)) \right|. \quad (44)$$

Then the performance is analyzed as

$$\frac{C_{[T_t^e]}^{REC}(O)}{C_{[T_t^e]}^{OPT}(O)} \leq \frac{P_s^c OT_t^e + P_s^e OT_t^e + P_b^e Og_{[T_t^e]}^e + \dots}{P_s^c OT_t^e + P_b^c Og_{[T_t^e]}^e + \dots} \quad (45a)$$

$$\leq \frac{P_s^e OT_t^e + P_s^c OT_t^e + P_b^e Og_{[T_t^e]}^e + \dots}{P_s^c OT_t^e + P_b^c Og_{[T_t^e]}^e} \quad (45b)$$

$$= 1 + \frac{P_s^e T_t^e + P_b^e g_{[T_t^e]}^e - P_b^c g_{[T_t^e]}^e + \dots}{P_s^c T_t^e + P_b^c g_{[T_t^e]}^e} \quad (45c)$$

$$\leq 1 + \frac{P_s^e T_t^e + P_b^e g_{[T_t^e]}^e}{P_s^c T_t^e + P_b^c g_{[T_t^e]}^e} \quad (45d)$$

$$= 1 + \frac{\epsilon P_s^c T_t^e + \beta P_b^c g_{[T_t^e]}^e}{P_s^c T_t^e + P_b^c g_{[T_t^e]}^e} \quad (45e)$$

$$\leq 1 + \frac{\max\{\epsilon, \beta\} (P_s^c T_t^e + P_b^c g_{[T_t^e]}^e)}{P_s^c T_t^e + P_b^c g_{[T_t^e]}^e} \quad (45f)$$

$$= 1 + \max\{\epsilon, \beta\}. \quad (45g)$$

In Eqn. (45a), we use ‘ $\dots$ ’ in the numerator and denominator to indicate the replica creation cost and the QoS penalty cost, respectively. Eqn. (45b) holds because the QoS penalty cost is always greater than or equal to 0, which means that the fraction will not be smaller when ignoring the QoS penalty cost from the denominator. In Eqn. (45c), we isolate a constant and divide the numerator and denominator by  $O$ . Eqn. (45d) holds with the fact as follows

$$-P_b^c g_{[T_t^e]}^e + \min_{e' \in \mathcal{N}_t^e} \{P_b^c, P_b^{e'} u(I_t^{e'})\} \leq 0.$$

In Eqn. (45e), we use  $P_s^c$  and  $P_b^c$  to represent  $P_s^e$  and  $P_b^e$  with Eqns. (35) and (36), respectively.

To sum up, the competitive ratio for the caching decision of Algorithm 2 is  $1 + \max\{\epsilon, \beta\}$  by comparing Eqns. (39), (43) and (45).  $\square$

For the caching decision given by the compensation strategy for RECaching, we discuss its competitive analysis in Proposition 2.

**Proposition 2.** *The caching decision of the compensation strategy is  $\max\{1 + \epsilon, 2 - \beta\}$ -competitive.*

*Proof.* The compensation strategy is triggered to cache an object for  $T_t^{ecom}$  time slots when the thresholds  $h_T$  and  $h_g$  are reached in time slot  $t$ . This proof starts by counting the number of requests from time slot  $t$ . We use  $t'$  to denote the time slot when the number of requests is equal to  $h_g$ ,  $t' > t$ . The cost performance of the caching decision can be analyzed by comparing  $t'$  and  $t + T_t^{ecom}$ .

If  $t' \leq t + T_t^{ecom}$ , which means that the caching decision is cost-effective because the number of requests during the caching lifespan  $T_t^{ecom}$  is greater than or equal to the threshold  $h_g$  triggering the caching decision, the cost of caching the object at the edge for  $T_t^{ecom}$  time slots is

$$C_{[T_t^{ecom}]}^{REC}(O) \leq C_{[T_t^{ecom}]}^e(O) = P_s^c OT_t^{ecom} + P_s^e OT_t^{ecom} + P_b^e Og_{[T_t^{ecom}]}^e + \min_{e' \in \mathcal{N}_t^e} \{C_t^{c \rightarrow e}(O), C_t^{e' \rightarrow e}(O)\}. \quad (46)$$

We use ‘ $\leq$ ’ in the above equation because the actual lifespan of caching the object at the edge is less than or equal to  $T_t^{ecom}$  due to the proposed lazy caching mechanism. Similar to Eqn. (38), the cost of optimal decision is

$$C_{[T_t^{ecom}]}^{OPT}(O) \geq P_s^c OT_t^{ecom} + P_s^e O + P_b^e Og_{[T_t^{ecom}]}^e + \min_{e' \in \mathcal{N}_t^e} \{C_t^{c \rightarrow e}(O), C_t^{e' \rightarrow e}(O)\}, \quad (47)$$

Thus, similar to Eqn. (39), the competitive ratio in this case is

$$\frac{C_{[T_t^{ecom}]}^{REC}(O)}{C_{[T_t^{ecom}]}^{OPT}(O)} \leq 1 + \epsilon. \quad (48)$$

If  $t' > t + T_t^{ecom}$ , which means that the caching lifespan is not cost-effective because the number of requests during the caching lifespan is less than  $h_g$ , the cost between time slots  $t$  and  $t'$  is

$$\begin{aligned} C_{[t'-t]}^{REC}(O) &= C_{[T_t^{ecom}]}^e(O) + C_{[t'-t-T_t^{ecom}]}^c(O) \\ &\leq (P_s^c + P_s^e) OT_t^{ecom} + P_s^c O(t' - t - T_t^{ecom}) \\ &\quad + P_b^e Oh_g + \min_{e' \in \mathcal{N}_t^e} \{C_t^{c \rightarrow e}(O), C_t^{e' \rightarrow e}(O)\} \\ &\quad + \omega^e \left| \sum_{t=t+1}^{t'} J_t^e(C_t^e(O) - C_t^c(O)) \right|. \end{aligned} \quad (49)$$

We use ‘ $\leq$ ’ in the above equation because we do not know which of  $h_g$  requests are served by the cached replica at the edge and assume that  $h_g$  requests are all served by the cloud. The offline optimal decision is not to cache the object at the

$$\frac{(P_s^c + P_s^e)OT_t^{ecom} + P_s^c O(t' - t - T_t^{ecom}) + P_b^c Oh_g + \min_{e' \in \mathcal{N}_t^e} \{P_b^c O, P_b^{e'} Ou(I_t^{e'})\} + \omega^e \left| \sum_{t=t+1}^{t'} J_t^e(C_t^e(O) - C_t^c(O)) \right|}{P_s^c O(t' - t) + P_b^c Oh_g + \omega^e \left| \sum_{t=t+1}^{t'} J_t^e(C_t^e(O) - C_t^c(O)) \right|} \quad (51)$$

edge that will incur additional unnecessary cost. Its cost is calculated as

$$C_{[t'-t]}^{OPT}(O) = C_{[t'-t]}^c(O) = P_s^c O(t' - t) + P_b^c Oh_g + \omega^e \left| \sum_{t=t+1}^{t'} J_t^e(C_t^e(O) - C_t^c(O)) \right|. \quad (50)$$

Thus, the competitive ratio in this case can be obtained by simplifying Eqn. (51). By omitting the QoS penalty item in both numerator and denominator, dividing the numerator and denominator by  $O$  and re-arranging items in the numerator, we obtain Eqn. (52a) as follows

$$\frac{C_{[t'-t]}^{REC}(O)}{C_{[t'-t]}^{OPT}(O)} \leq \frac{P_s^e T_t^{ecom} + P_s^c(t' - t) + P_b^c h_g + \dots}{P_s^c(t' - t) + P_b^c h_g} \quad (52a)$$

$$= 1 + \frac{P_s^e T_t^{ecom} + \min_{e' \in \mathcal{N}_t^e} \{P_b^c, P_b^{e'} u(I_t^{e'})\}}{P_s^c(t' - t) + P_b^c h_g} \quad (52b)$$

$$= 1 + \frac{P_b^c h_g - P_b^e h_g}{P_s^c(t' - t) + P_b^c h_g} \quad (52c)$$

$$\leq 1 + \frac{P_b^c h_g - P_b^e h_g}{P_b^c h_g} = 1 + \frac{P_b^c - \beta P_b^c}{P_b^c} \quad (52d)$$

$$= 2 - \beta. \quad (52e)$$

In Eqn. (52a), ‘...’ indicates the replica creation cost. Then, we isolate a constant and replace  $T_t^{ecom}$  with Eqn. (31) in Eqns. (52b) and (52c), respectively. Eqn. (52d) holds because we omit  $P_s^c(t' - t)$  in the denominator which is greater than 0 and replace  $P_b^e$  with  $P_b^c$  using Eqn. (36).

Finally, by comparing Eqns. (48) and (52), we obtain that the competitive ratio is  $\max\{1 + \epsilon, 2 - \beta\}$ .  $\square$

Next, for the caching decision given by the prolonged caching mechanism for RECaching, the cost performance is analyzed in Proposition 3.

**Proposition 3.** *The caching decision of the prolonged caching mechanism is  $(1 + \max\{\epsilon, \beta\})$ -competitive.*

*Proof.* This proof is similar to that of Proposition 1, so we present a sketch here. If an object that has been cached at the edge is decided to continue to be cached for another  $T_t^{epro}$  time slots by the prolonged caching mechanism when it receives  $g_t$  requests in time slot  $t$ , the incurred cost for  $T_t^{epro}$  time slots is

$$C_{[T_t^{epro}]}^{REC}(O) = P_s^c OT_t^{epro} + P_s^e OT_t^{epro} + P_b^e Og_{[T_t^{epro}]}, \quad (53)$$

where  $g_{[T_t^{epro}]}$  denotes the number of requests during the cached  $T_t^{epro}$  time slots. For the offline optimal decision, similar to the proof of Proposition 1, we have two cases. If  $g_{[T_t^{epro}]} \geq g_t$ , the offline optimal decision is to cache the

object. Otherwise, it is not to cache the object. Therefore, its cost  $C_{[T_t^{epro}]}^{OPT}(O)$  has a lower bound of

$$P_s^c OT_t^{epro} + \begin{cases} P_s^e O + P_b^e Og_{[T_t^{epro}]}, & g_{[T_t^{epro}]} \geq g_t, \\ P_b^e Og_{[T_t^{epro}]}, & g_{[T_t^{epro}]} < g_t. \end{cases} \quad (54)$$

When  $g_{[T_t^{epro}]} \geq g_t$ , the competitive ratio is calculated as

$$\frac{C_{[T_t^{epro}]}^{REC}(O)}{C_{[T_t^{epro}]}^{OPT}(O)} \leq \frac{P_s^c OT_t^{epro} + P_s^e OT_t^{epro} + P_b^e Og_{[T_t^{epro}]}}{P_s^c OT_t^{epro} + P_s^e O + P_b^e Og_{[T_t^{epro}]}} \quad (55a)$$

$$\leq \frac{P_s^c T_t^{epro} + P_s^e T_t^{epro}}{P_s^c T_t^{epro}} \quad (55b)$$

$$= \frac{P_s^c + \epsilon P_s^c}{P_s^c} = 1 + \epsilon. \quad (55c)$$

When  $g_{[T_t^{epro}]} < g_t$ , the competitive ratio is calculated as

$$\frac{C_{[T_t^{epro}]}^{REC}(O)}{C_{[T_t^{epro}]}^{OPT}(O)} = \frac{P_s^c OT_t^{epro} + P_s^e OT_t^{epro} + P_b^e Og_{[T_t^{epro}]}}{P_s^c OT_t^{epro} + P_b^e Og_{[T_t^{epro}]}} \quad (56a)$$

$$= 1 + \frac{P_s^e T_t^{epro} + P_b^e g_{[T_t^{epro}]} - P_b^c g_{[T_t^{epro}]}}{P_s^c T_t^{epro} + P_b^c g_{[T_t^{epro}]}} \quad (56b)$$

$$\leq 1 + \frac{P_s^e T_t^{epro} + P_b^e g_{[T_t^{epro}]}}{P_s^c T_t^{epro} + P_b^c g_{[T_t^{epro}]}} \quad (56c)$$

$$= 1 + \frac{\epsilon P_s^c T_t^{epro} + \beta P_b^e g_{[T_t^{epro}]}}{P_s^c T_t^{epro} + P_b^c g_{[T_t^{epro}]}} \quad (56d)$$

$$= 1 + \max\{\epsilon, \beta\}. \quad (56e)$$

Eqn. (56c) holds due to  $P_b^c g_{[T_t^{epro}]} \geq 0$ . In conclusion, the competitive ratio of the caching decision of the prolonged caching mechanism is  $1 + \max\{\epsilon, \beta\}$ .  $\square$

Finally, Proposition 4 analyzes the performance of uncached periods of RECaching enhanced by the mechanisms.

**Proposition 4.** *The uncached periods of RECaching improved by the mechanisms is  $(1 + \epsilon)$ -competitive.*

*Proof.* If an object is not cached at the edge, it means that the compensation mechanism is not triggered. That is, the number of requests during the past  $h_T$  consecutive time slots without cached replicas does not reach the threshold  $h_g$ . Thus, the incurred cost of RECaching has an upper bound,

$$C_{[h_T]}^{REC}(O) \leq P_s^c Oh_T + P_b^c Oh_g + \omega^e \left| \sum_{h_T} J_t^e(C_t^e(O) - C_t^c(O)) \right|. \quad (57)$$

For the offline optimal decision, if it has a knowledge of requests during the uncached time slots in advance, it can decide to whether to cache the object at the edge or not with the goal of cost optimization. We discuss the cost of the offline optimal decision by case. If the optimal decision is not to cache the object, the incurred cost is

$$C_{[h_T]}^{OPT}(O) < P_s^c O h_T + P_b^c O h_g + \omega^e \left| \sum_{h_T} J_t^e(C_t^e(O) - C_t^c(O)) \right|. \quad (58)$$

Obviously, the competitive for the decision of uncached periods in this case is 1. If the optimal decision is to cache the object at the edge, the incurred cost has a lower bound of

$$C_{[h_T]}^{OPT}(O) \geq P_s^c O h_T + P_s^e O + P_b^e O h_g + \min_{e' \in \mathcal{N}_t^e} \{P_b^c O, P_b^{e'} O u(I_t^{e'})\}. \quad (59)$$

In this case, the competitive ratio is calculated as

$$\frac{C_{[h_T]}^{REC}(O)}{C_{[h_T]}^{OPT}(O)} \leq \frac{P_s^c h_T + P_b^c h_g + \sum_{h_T} C_t^e(1) - \sum_{h_T} C_t^c(1)}{P_s^c h_T + P_s^e + P_b^e h_g + \min_{e' \in \mathcal{N}_t^e} \{P_b^c, P_b^{e'} u(I_t^{e'})\}} \quad (60a)$$

$$\leq \frac{P_s^c h_T + P_s^e h_T + P_b^e h_g + \min_{e' \in \mathcal{N}_t^e} \{P_b^c, P_b^{e'} u(I_t^{e'})\}}{P_s^c h_T + P_s^e + P_b^e h_g + \min_{e' \in \mathcal{N}_t^e} \{P_b^c, P_b^{e'} u(I_t^{e'})\}} \quad (60b)$$

$$\leq \frac{P_s^c h_T + P_s^e h_T + P_b^e h_g}{P_s^c h_T + P_s^e + P_b^e h_g} \quad (60c)$$

$$\leq \frac{P_s^c h_T + P_s^e h_T}{P_s^c h_T} = \frac{P_s^c + \epsilon P_s^c}{P_s^c} = 1 + \epsilon. \quad (60d)$$

In conclusion, the cost performance for not caching can be bounded by the competitive ratio  $1 + \epsilon$ .  $\square$

## APPENDIX ADDITIONAL EXPERIMENTS

### A. Distance from the Optimal

This section assesses the discrepancy between RECaching and the optimal strategy *OPT* by comparing the ratios of costs obtained by these two algorithms across multiple data objects. Fig. 11 shows the distributions of the cost ratios under different pricing schemes. As we can see, the cost ratios are not large, which means that the cost performance of RECaching is close to that of *OPT*. When applying pricing schemes *UE-S1* and *AP-S1*, as depicted in Figs. 11a and 11b, RECaching demonstrates cost performance closest to that of *OPT*, which suggests that the increase in cloud pricing has minimal impact on RECaching's performance. The performance of RECaching degrades slightly when raising the storage price and lowering the bandwidth price for the edge, which is obtained by comparing Figs. 11a and 11c. Meanwhile, the performance also degrades slightly when raising the storage price for the edge, which is obtained by comparing Figs. 11b and 11d. Even so, the vast majority of cost ratios under pricing model *UE-C* are within 1.2. Hence, we can conclude that the fluctuation in edge prices has a greater impact on RECaching's performance

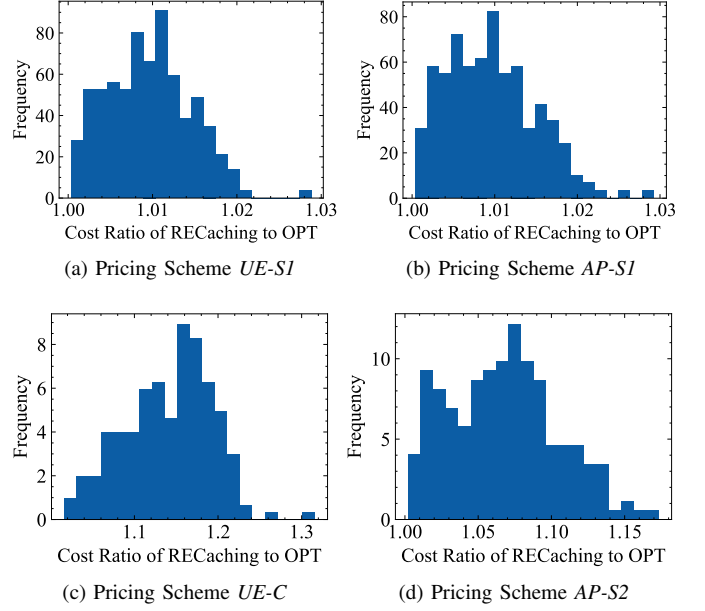


Fig. 11. Distributions of the cost ratios of RECaching to *OPT* under different pricing schemes.

than the fluctuation in cloud prices. Meanwhile, RECaching can achieve a stable performance that is close to the optimal.

At the same time, we observe that for RECaching, the distribution of cost ratios as shown in Fig. 11 and the distribution of the percentage of optimal decisions as shown in Fig. 8 are not entirely consistent. For instance, in Fig. 8a, the percentage of RECaching in agreement with the optimal decisions is about 70% consistency with the optimal decisions, while in Fig. 11a, RECaching's costs can closely approximate the optimal costs, which indicates that the decision sequence that results in a cost close to the optimal during the lifespan of an object is not unique.

In addition, we observe that the distributions of the algorithms are similar under pricing schemes *AP-S1* and *UE-C*, as shown in Figs. 8b and 8c, but the cost ratios of RECaching to *OPT* under these pricing schemes have a pronounced disparity, as shown in Figs. 11b and 11c. The reason is that the different in pricing schemes *AP-S1* and *UE-C* results in significant cost performance, which can be obtained from Figs. 6 and 7. In this case, even if the distributions are similar, pricing model *UE-C* results in higher costs for infrequently requested data objects due to its expensive storage price at the edge, thereby increasing the cost disparity between RECaching and *OPT* under the same decisions.

### B. Discussion on the Optimal Offline Algorithm

In this section, we explore the running time of the optimal offline algorithm *OPT* proposed in Section IV and discuss the result with a commercial solver Gurobi. This experiment is run on a Windows 10 machine with Intel(R) Core(TM) i7-9700F CPU @ 3.00GHz and 32 GB of RAM. *OPT* is implemented in Python 3.8.13. The version of Gurobi is v11.0.3rc0. To solve problem  $(P_1)$  with Gurobi, we encode

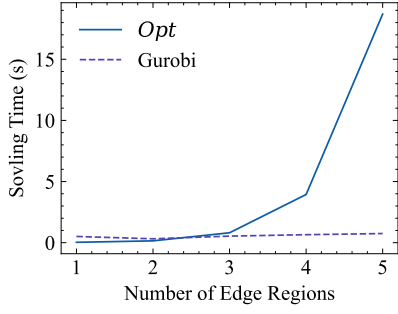


Fig. 12. Running time with respect to the number of edge regions.

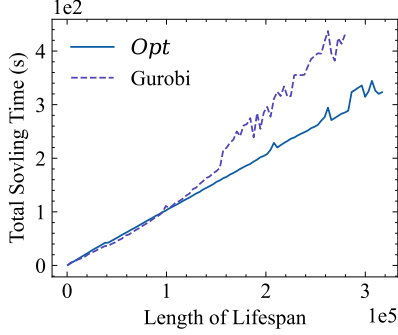


Fig. 13. Running time with respect to the length of lifespan of objects.

it using AMPL (A Mathematical Programming Language), as shown in **Algorithm 3**.

Our purpose of designing *OPT* is to show the intractability of problem  $(\mathcal{P}_1)$  and provide a principled benchmark to evaluate the distance from the optimal of RECCaching. In terms of running time, due to its implementation in pure Python and lack of low-level optimization, it is less efficient and certainly not comparable to commercial software Gurobi that is carefully optimized, integrates multiple optimization algorithms, and employs parallel computing.

As analyzed in Section IV, the time complexity of *OPT* is exponential to the number of edge regions and linear to the length of lifespan. We explore the variation of running time with the number of edge regions when the lifespan length is fixed, and the variation of running time with the lifespan length when the number of edge regions is determined, respectively. In the former experiment, the lifespan is set to 680. Fig. 12 shows the result, from which we can see that *OPT* obtains a solving time similar to Gurobi's when the number of edge regions is small. When the number of edge regions becomes large, *OPT*'s solving time increases exponentially, while Gurobi maintains a fast solving time.

However, Gurobi's advantage in running time only exists when the lifespan is short. In practice, objects stored in the cloud survive for a long time. In this experiment, the number of edge regions is set to 3. We expand the lifespan by concatenating request sequences of multiple objects. Fig. 13 shows the result, from which we can see that the solving times of *OPT* and Gurobi are similar to each other when the lifespan length is short. Gurobi's solving time is larger than that of *OPT* when the lifespan is longer. Besides, we note that there is no

---

**Algorithm 3** Encoding Problem  $(\mathcal{P}_1)$  in AMPL.

---

**Input:**

Object size  $O$ , lifespan  $T$ , edge regions  $E$ . Request sequences of all edge regions  $G_T = [G_T^1, \dots, G_T^e, \dots, G_T^E]$ .

**Output:**

Caching decisions  $[I_1^1, \dots, I_t^e, \dots, I_T^E]$ . Optimal cost.

- 1: param  $T$  integer  $> 0$ ;
  - 2: param  $E$  integer  $> 0$ ;
  - 3: param  $G_T$ ;
  - 4: param  $O$ ;
  - 5: param  $M$  default 1e6;
  - 6: var  $I_t^e$  binary; /\* Caching status for region  $e$  in time  $t$ . \*/
  - 7: var  $z_t$  binary; /\* Auxiliary variable, indicating caching status of all regions in the previous time slot of time  $t$ . \*/
  - 8: var  $y_t^e \geq 0, \leq 1$ ; /\* Auxiliary variable, used to linearize product terms in the equations such as Eqn. (7). \*/
  - 9: /\* The following 2 constraints establish the relationship between  $z_t$  and previous caching status in time  $t - 1$ . \*/
  - 10: subject to  $z_{upper}\{t \text{ in } 2..T\}$ :  
 $\text{sum}\{e \text{ in } 1..E\} I_{t-1}^e \leq M \times z_t$ ;
  - 11: subject to  $z_{lower}\{t \text{ in } 2..T\}$ :  
 $\text{sum}\{e \text{ in } 1..E\} I_{t-1}^e \geq 0.001 - M \times (1 - z_t)$ ;
  - 12: /\* The following 3 constraints linearize product terms. \*/
  - 13: subject to  $y_{upper1}\{e \text{ in } 1..E, t \text{ in } 2..T\}$ :  
 $y_t^e \leq I_t^e$ ;
  - 14: subject to  $y_{upper2}\{e \text{ in } 1..E, t \text{ in } 2..T\}$ :  
 $y_t^e \leq 1 - I_{t-1}^e$ ;
  - 15: subject to  $y_{lower}\{e \text{ in } 1..E, t \text{ in } 2..T\}$ :  
 $y_t^e \geq I_t^e + (1 - I_{t-1}^e) - 1$ ;
  - 16: /\* Objective function. \*/
  - 17: minimize the objective of problem  $(\mathcal{P}_1)$ .
- 

data for Gurobi in Fig. 13 when the lifespan length is greater than  $2.8e5$ . The reason is that the 32 GB of RAM used for our experiments can not support Gurobi to solve the problem.