# Contract theory based incentive mechanism design for buffer resource in wireless caching networks[①]

Liu Tingting(刘婷婷)[②]*, Tang Lei*, Zhu Hao*, Bao Yongqiang*, Guo Yajuan**

(*School of Information and Communication Engineering, Nanjing Institute of Technology, Nanjing 211167, P. R. China)
(**Institute of Electric Power Science, Jiangsu Electric Power Co. Ltd., Nanjing 210024, P. R. China)

## Abstract

Evidences indicate that, due to the limited caching capacity or inaccurate estimation on users' preferences, the requested files may not be fully cached in the network edge. The transmissions of the un-cached files will also lead to duplicated transmissions on backhaul channels. Buffer-aided relay has been proposed to improve the transmission performance of the un-cached files. Because of the limited buffer capacity and the information asymmetric environment, how to allocate the limited buffer capacity and how to incentivize users in participating buffer-aided relay have become critical issues. In this work, an incentive scheme based on the contract theory is proposed. Specifically, the backlog violation probability, i. e., the buffer overflow probability, is provided based on the martingale theory. Next, based on the backlog violation probability, the utility functions of the relay node and users are constructed. With the purpose to maximize the utility of the relay node, the optimal contract problem is formulated. Then, the feasibility of the contract is also demonstrated, and the optimal solution can be obtained by the interior point method. Finally, numerical results are presented to demonstrate effectiveness of the proposed contract theory scheme.

**Key words**: contract theory, buffer resource allocation, wireless caching network, incentive mechanism design

## 0 Introduction

Due to the surge of duplicated transmissions in wireless networks, wireless caching technology has been proposed in order to release the transmission redundancy on backhaul channels[14]. The advantages of wireless caching technology have been provided as releasing traffic pressures on backhaul channels, reducing the transmission latency and improving the overall network energy efficiency, etc. The recent researches on wireless caching technology can be roughly classified into 2 aspects: content placement[5,6] and content delivery[7,8]. Most of these works have assumed that the popular files can be pre-cached into the edge nodes. However, due to the limited caching capacity, or the inaccurate estimation on user preference profiles, the requested files may not be cached or only be partially cached in the edge nodes. The transmissions of the un-cached files/parts may still induce the duplicated transmissions on backhaul channels or lead to transmission latency.

In order to further enhance the performance of wireless caching, researchers begin to investigate the joint buffer and cache allocations in wireless caching networks. Cache resource can be considered as a long-term memory to pre-cache the popular files during off-peak time. The pre-cached files can be directly transmitted upon requesting without triggering traffic on backhaul channels. Buffer resource is a kind of short-term memory. It can be used to store the fetched files from the remote servers or from the cache memory. Thus, it can be utilized to enable buffer-aided relay to improve the delivery performance of the un-cached files. Buffer usually has a high input/output (I/O) speed, and therefore, it is allocated with an expensive primary memory. In practice, due to the limited battery capacity, users usually have their personal use purposes or privacy concerns[9]. It is impossible for the edge node to contribute all its buffer resource. Gener-

---

ally, it will devote a limited buffer capacity to help transmissions. It is known that, due to users' unbalanced data arrival and service rates, there may be some backlogs inside an edge or a relay node. In this paper, the edge node also acts as a relay node. The term 'relay node' is used for the ease of discussions. Backlogs may lead to queueing delay, and backlog overflow which is induced by the limited buffer capacity, and even induce re-transmissions over backhaul channels. How to allocate the limited buffer capacity to maximize relay node's utility has become a critical issue in realizing the buffer-aided relay in wireless caching networks. At the same time, users are willing to utilize the buffer-aided relay scheme to enhance its transmission performance. How to design a proper incentive scheme to incentivize users in participating buffer-aided relay communications has also become an important issue. In order to solve the above 2 issues, 2 problems need to be figured out.

First, the relationship among the arrival rate, the service rate, the backlogs, and the buffer capacity should be determined. In this way, the buffer capacity on relay node's side can be linked to each user's backlogs which are determined by this user's arrival and service rates. The effective bandwidth theory is a useful framework to analyze the backlog performance of broad classes of arrivals[10,11]. However, it may lead to loose estimation for non-Poisson processes, and it is not realistic to assume user's arrival or service data follow Poisson processes. At the same time, the martingale theory is proposed as an alternative in analyzing the backlogs, which is flexible and suitable for any kind of arrival and service processes[12]. A lot of works have demonstrated that the martingale theory can provide a tight bound compared with the real data trace. Moreover, the derived bound can be used as an objective function or as a constraint to further optimize the network performance[13-16]. Thus, in this work, the martingale theory will be used to determine the relationship among the arrival rate, the service rate, the backlogs, and the buffer capacity.

Second, in the concerned wireless caching networks, user's specific parameters cannot be fully observed by the relay node. Information asymmetry exists in this scenario. The relay node only knows the set of user's parameters, but it cannot observe the specific value of each user. Contract theory is known as a powerful tool in solving this kind of information asymmetry problem[17]. Most of the existing works on contract theory can be classified into 2 categories: one is adverse selection[18-21], and the other one is moral hazard[22,23]. The existing works have highlighted that solving the problem of asymmetric information in resource allocation is the major advantage of contract theory.

In this work, an incentive mechanism for buffer resource based on the contract theory in the wireless caching networks is designed. The proposed scheme can provide proper economic incentives to users. The main contributions of this paper are summarized as follows.

(1) A commercial wireless caching network is constructed, where user's backlog violation probability is derived based on the martingale theory. Then, the utility functions of the relay node and the users are formulated.

(2) Contract theory is proposed to design the buffer resource allocation scheme between the relay node and the users. The users are classified into different types according to their distances from the relay node. Correspondingly, the relay node divides its devoted buffer resource into different portions, defined as quality. The optimal contract problem is then constructed in maximizing the utility of the relay node.

(3) Next, the conditions of a feasible contract are developed, and they are considered as constraints in the constructed optimization problem. In order to obtain the optimal contract, the constraints need to be reduced by several steps, which are also elaborated.

(4) Numerical results are provided to demonstrate effectiveness of the proposed contract theory scheme. Also, it can be found that in terms of maximizing the relay node's utility, the proposed scheme is superior to the equally allocated scheme.

The remainders of this paper are organized as follows. The system model and the backlog violation probability derived from the martingale theory, as well as the utility functions of the relay node and users are presented in Section 1. The detailed incentive mechanism design steps based on the contract theory are elaborated in Section 2. Numerical results are presented in Section 3, and finally conclusions are drawn in Section 4.

# 1　System model

In a wireless caching network, instead of direct transmitting files from the base station (BS) to the users, as shown in Fig. 1, a nearby node will be employed as a relay to enhance the delivery performance.

## 1.1　User's data arrival model

A user $U_i$ requires a file. Data are transmitted using the buffer-aided relay scheme. Assuming that the bursty file amount $a_i(k)$ for user $U_i$ arriving at the relay
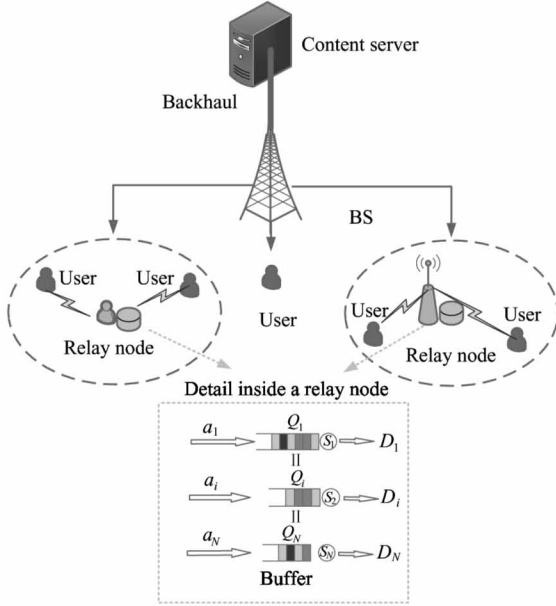
**Fig. 1**  System model of a wireless caching network and the details inside a relay node

node at time $k$ follows a Markov-modulated on off (MMOO) process which has 2 static status $[\Pi^0, \Pi^1]$. On state $\Pi^0$, there is no data arriving, i. e., $a(k) = 0$, while on state $\Pi^1$, $a_i(k) = R$ bits/s, and $R > 0$. The corresponding state transition matrix $\boldsymbol{Tr}^a$ is defined as

$$\boldsymbol{Tr}^a = \begin{bmatrix} 1 - T^\alpha & T^\alpha \\ T^\beta & 1 - T^\beta \end{bmatrix} \qquad (1)$$

where, $T^\alpha$ represents the transition probability from state $\Pi^0$ to state $\Pi^1$, $T^\beta$ represents the transition probability from $\Pi^1$ to $\Pi^0$. Accordingly, the steady state distribution of $a_i(k)$ is calculated as

$$\left[\frac{T^\beta}{T^\alpha + T^\beta}, \frac{T^\alpha}{T^\alpha + T^\beta}\right] \qquad (2)$$

Therefore, the cumulative arrival data over time interval $[m, n]$ is represented as

$$A_i(m, n) = \sum_{k=m}^{n} a_i(k) \times 1 \qquad (3)$$

where the unit of time interval is one second, and $A_i(m, n)$ can be regarded as a bivariate arrival process. If $m = 0$, $A_i(0, n) = A_i(n)$ is used for brevity.

### 1.2  Data service model

Assuming that the relay node devotes a buffer capacity $C_B$ to help the surrounding $N$ users relay their data. The transmission bandwidth is $B$. The $N$ users equally share the bandwidth. A constant transmission power $P_{tr}$ is utilized for each user. According to Shannon's channel capacity theorem[24], the service rate, i. e., the transmission rate, from the relay node to the

user can be written as

$$s_i(k) = \frac{B}{N}\log_2\left(1 + \frac{P_{tr}d_i^{-l}}{\frac{B}{N}\sigma_0^2}\right) \qquad (4)$$

where, $d_i$ is the distance between a user $U_i$ and the relay node, $l$ is the pathloss exponent, $\sigma_0^2$ is the noise density. The accumulative transmission amount $S_i$ over time interval $[m, n]$ can be written as

$$S_i(m, n) = \sum_{k=m}^{n} s_i(k) \times 1 \qquad (5)$$

Note that when the user's location is known, the service rate is irrelevant to time $k$. $s_i$ is used for brevity.

Also, the following assumption is provided to make the analysis non-trivial.

$$\mathbb{E}[a_i(k)] < s_i < \sup a_i(k), \quad \forall i \qquad (6)$$

Eq. (6) means that the service rate is larger than the average arrival rate, and it is smaller than the peak arrival rate. In this way, there will be some backlogs inside a relay node's buffer.

### 1.3  Backlog violation probability

Since there are some backlogs inside a relay node's buffer, i. e.,

$$Q_i = \sup_{n \geq 0}\{A_i(n) - S_i(n)\} \qquad (7)$$

and the relay node only devotes a limited buffer capacity $C_B$. Moreover, each user is allocated with a dedicated buffer with capacity $\alpha_i C_B$, where $\alpha_i$ represents the portion assigned to user $U_i$, $\alpha_i \geq 0$, and $\sum_{i=1}^{N} \alpha_i = 1$. The backlog violation probability is defined as follows.

**Definition 1**  Backlog violation probability is the probability that a user $U_i$'s allocated buffer $\alpha_i C_B$ is overflowed, i. e., $\Pr(Q_i \geq (\alpha_i C_B))$.

The backlogs in Eq. (7) show a stochastic characteristic, and it is infeasible to determine the precise value of the backlogs. The martingale theory is known as a power tool in handling this kind of problem. In the following, the martingale theory is utilized to analyze the property of the backlogs. First, 2 martingales should be constructed, i. e., the arrival martingale $M_i^A(n)$, i. e.,

$$M_i^A(n) = h_i^a(a_i(n))e^{\theta_i(A_i(n) - nK_i^a)}, \quad n \geq 0 \qquad (8)$$

and the service martingales $M_i^S(n)$, i. e.,

$$M_i^S(n) = h_i^s(s_i(n))e^{\theta_i(nK_i^s - S_i(n))}, \quad n \geq 0 \qquad (9)$$

The scheduling policy is first in first out (FIFO). Many works have contributed to derive the backlog violation probability. Here, it is provided directly in Theorem 1.

**Theorem 1**  The backlog violation probability is

calculated as

$$\Pr(Q_i \geqslant (\alpha_i C_B)) \leqslant \frac{\mathbb{E}[h_i^a(a_i(0))]}{H_i} e^{-\theta_i^* \alpha_i C_B}$$

(10)

where $H_i = \min\{h_i^a(a_i)h_i^s(s_i) \mid a_i - s_i > 0\}$.

**Proof**　Please refer to Refs[12,14] for the rigorous deduction.

Since all the users have the same arrival model, and the service rate $s_i$ is irrelevant to time $k$. $\dfrac{\mathbb{E}[h_i^a(a_i(0))]}{H_i}$ in Eq. (10) can be derived as a constant[12,14]. Thus, it is defined as

$$h = \frac{\mathbb{E}[h_i^a(a_i(0))]}{H_i}$$

(11)

and Eq. (10) can be rewritten as

$$\Pr(Q_i \geqslant (\alpha_i C_B)) \leqslant h e^{-\theta_i^* \alpha_i C_B}$$

(12)

### 1.4　Users' utility

Assuming that a user $U_i$ rents an $\alpha_i$ portion, it needs to pay $\pi_i$ to the relay node. If the data is overflowed, the overflowed data will be re-transmitted through the cellular link, and this will cost the user $U_i$ with some extra money. The utility of the user $U_i$ is shown as follows.

$$U_{\text{utility}}^i = \rho_0(1 - h e^{-\theta_i^* \alpha_i C_B}) - \pi_i, \quad \forall i$$

(13)

where, $\rho_0(1 - h e^{-\theta_i^* \alpha_i C_B})$ is the saving cost which is gained from renting the relay node's buffer portion, $\theta_i^*$ is a parameter that is jointly determined by the data generation rate $a_i(k)$ and the service rate $s_i(k)$. In the following, $\theta_i^*$ is defined as user's type.

**Definition 2**　User's type: $\theta_i^*$ which is jointly determined by its data generation rate and service rate, is defined as user's type. User's type is sorted by an increasing order. That is

$$\theta_1^* < \theta_2^* < \cdots < \theta_N^*$$

(14)

due to information asymmetry, the relay node only knows the set of user's parameters, but it cannot observe the specific value of each user.

### 1.5　Relay node's utility

Relay node will make profits by renting out its buffer resource to the surrounding users. The relay node rents out an $\alpha_i$ portion to the user $U_i$, and correspondingly the user $U_i$ will pay $\pi_i$ to the relay node. Thus, the utility function of the relay node is represented as

$$U_{RL}(\lambda, \alpha, \pi) = \sum_{i=1}^{N} (\pi_i + \delta_0(1 - h e^{-\theta_i^* \alpha_i C_B}))$$

(15)

where, $1 - h e^{-\theta_i^* \alpha_i C_B}$ denotes the probability that the

backlog does not overflow. In other words, it means the probability of success that the relay node can buffer the transmission data. Otherwise, the overflowed data will be transmitted through the cellular link. Thus, the corresponding saved cost is $\delta_0(1 - h e^{-\theta_i^* \alpha_i C_B})$.

## 2　Incentive mechanism design based on contract theory

The relay node designs the optimal contracts $\{\alpha_i, \pi_i\}$, $i = \{1, \cdots, N\}$ to maximize its own profits. The optimal contracts also need to comply with the following 2 constraints: individual rationality (IR) and incentive compatibility (IC) for all intended users. The definitions of the IR and IC constraints are as follows.

**Definition 3**　IR constraint: each user is assumed to be rational, and it will not accept a contract entry which produces a negative utility for its type. Therefore, IR constraint is provided as follows.

$$\rho_0(1 - h e^{-\theta_i^* \alpha_i C_B}) - \pi_i \geqslant 0, \quad \forall i$$

(16)

IR constraint means that the savings must compensate the payment. In the case of $U_i < 0$, the user will not purchase the buffer portion.

**Definition 4**　IC constraint: The incentive compatibility constraint means that $V_v$ cannot gain more utility by accepting a contract entry which is not designed for its type. That is

$$\rho_0(1 - h e^{-\theta_i^* \alpha_i C_B}) - \pi_i \geqslant \rho_0(1 - h e^{-\theta_i^* \alpha_{\tilde{i}} C_B}) - \pi_{\tilde{i}}, \quad \tilde{i} \neq i \quad (17)$$

where the contract $\{\alpha_{\tilde{i}}, \pi_{\tilde{i}}\}$ is designed for type $\tilde{i}$.

In other words, user with type $i$ should obtain the maximum utility if and only if it chooses the contract $\{\alpha_i, \pi_i\}$ designed for its type.

The optimal contract problem aims to maximize the relay node's utility, and also complies with the IR and IC constraints. Thus, the optimal contract problem is formulated as follows.

$$\{\alpha_i^*, \pi_i^*\} = \text{argmax} U_{RL}$$

(18)

$$\text{s. t. IR}(16), \text{IC}(17), \sum_{i=1}^{N} \alpha_i \leqslant 1, 0 \leqslant \alpha_i \leqslant 1, \forall i$$

Contract theory is a kind of economical tool that can provide enough incentives to the users that motivate them choose the intended contract entry. This mechanism is guaranteed by the IR and IC constraints. That is to say, if users are not interested in paying anything, they will gain nothing, and if they do not choose the intended contract entry, they will not obtain the maximum profits.

### 2.1　Constraints reduction

Since in Eq. (18), there are $N$ IR constraints and

$N \times (N-1)$ IC constraints. It is not tractable to solve an optimization problem with so many constraints. First, the number of constraints should be reduced by the following lemmas.

### 2.1.1 IR constraint reduction

Lemma 1 is provided to reduce the IR constraint.

**Lemma 1** ( IR constraint Reduction) The IR constraint for the lowest type $\theta_1^*$ is binding, i.e.,

$$\rho_0(1 - he^{-\theta_1^* \alpha_1 C_B}) - \pi_1 = 0 \qquad (19)$$

**Proof** From Definition 2, $\theta_1^* < \cdots < \theta_i^* < \cdots < \theta_N^*$. $E(\theta_i^*, \alpha_i)$ is used to represent $1 - he^{-\theta_i^* \alpha_i C_B}$ for simplicity, i.e.,

$$E(\theta_i^*, \alpha_i) = 1 - he^{-\theta_i^* \alpha_i C_B} \qquad (20)$$

Considering the properties of IC constraints and the increasing property of function $\rho_0 E(\theta_i^*, \alpha_i) - \pi_i$ over $\theta_i^*$, i.e.,

$$\rho_0 E(\theta_i^*, \alpha_i) - \pi_i \overset{(a)}{\geq} \rho_0 E(\theta_i^*, \alpha_1) - \pi_1$$
$$\overset{(b)}{\geq} \rho_0 E(\theta_1^*, \alpha_1) - \pi_1, \ \forall i \qquad (21)$$

where the property of IC constraint can assure $(a)$, while the increasing property of function $E(\theta_i^*, \alpha_i)$ over $\theta_i^*$ assures $(b)$. Therefore, if $\rho_0 E(\theta_1^*, \tau_1) - \pi_1 \geq 0$, all the computation nodes will satisfy the IR constraints. This completes the proof.

### 2.1.2 IC constraint reduction

First, in order to make the analysis consistent with the contract theory. 'Quality' is defined in this paper.

**Definition 5** Since the buffer resource is considered as trading goods, the buffer portion $\alpha_i$, $\forall i$, is defined as quality.

The IC constraints can be reduced by the following lemmas.

**Lemma 2** If the contract entries satisfy the IC constraints, the following relationship holds. Given the quality $\alpha_i > \alpha_{\bar{i}}$, if and only if the price satisfies $\pi_i > \pi_{\bar{i}}$.

**Proof** Lemma 2 is demonstrated from 2 aspects, i.e., sufficient condition and necessary condition.

(1) Sufficient condition: if the contract satisfies the IC constraints shown in Eq. (17), i.e., $\rho_0 E(\theta_i^*, \alpha_{\bar{i}}) - \pi_{\bar{i}} \geq \rho_0 E(\theta_i^*, \alpha_i) - \pi_i$, $\forall i \neq \bar{i}$, it can be rewritten as $\rho_0 E(\theta_i^*, \alpha_{\bar{i}}) - \rho_0 E(\theta_i^*, \alpha_i) \geq \pi_{\bar{i}} - \pi_i$ is an increasing function of $\alpha_i$, given $\alpha_i > \alpha_{\bar{i}} \Rightarrow E(\theta_i^*, \alpha_{\bar{i}}) - E(\theta_i^*, \alpha_i) < 0 \Rightarrow \pi_i > \pi_{\bar{i}}$.

(2) Necessary condition: if the contract satisfies the IC constraints, i.e., $\rho_0 E(\theta_i^*, \alpha_i) - \pi_i \geq \rho_0 E(\theta_i^*, \alpha_{\bar{i}}) - \pi_{\bar{i}}$, $\forall i \neq \bar{i}$, it can be rewritten as $\rho_0 E(\theta_i^*, \alpha_i) - \rho_0 E(\theta_i^*, \alpha_{\bar{i}}) \geq \pi_i - \pi_{\bar{i}}$. Given $\pi_i > \pi_{\bar{i}}$, $\rho_0 E(\theta_i^*, \alpha_i) - \rho_0 E(\theta_i^*, \alpha_{\bar{i}}) \geq 0$ holds. Since

$E(\theta_i^*, \alpha_i)$ is increasing with $\alpha_i$, then $\alpha_i > \alpha_{\bar{i}}$. This completes the proof.

**Lemma 3** If the contract entries satisfy the IC constraints, the quality $\alpha_i$ monotonically decreases with type $\theta_i^*$, i.e., if $\theta_i^* > \theta_{\bar{i}}^*$, then $\alpha_i < \alpha_{\bar{i}}$.

**Proof** Given the IC constraints, i.e., $\rho_0 E(\theta_i^*, \alpha_i) - \pi_i \geq \rho_0 E(\theta_i^*, \alpha_{\bar{i}}) - \pi_{\bar{i}}$ and $\rho_0 E(\theta_{\bar{i}}^*, \alpha_{\bar{i}}) - \pi_{\bar{i}} \geq \rho_0 E(\theta_{\bar{i}}^*, \alpha_i) - \pi_i$, $\forall i \neq \bar{i}$, then

$$\rho_0 E(\theta_i^*, \alpha_i) - \rho_0 E(\theta_i^*, \alpha_{\bar{i}}) \geq \pi_i - \pi_{\bar{i}} \qquad (22)$$

and

$$\rho_0 E(\theta_{\bar{i}}^*, \alpha_i) - \rho_0 E(\theta_{\bar{i}}^*, \alpha_{\bar{i}}) \leq \pi_i - \pi_{\bar{i}} \qquad (23)$$

Observing Eq. (22) and Eq. (23), then

$$E(\theta_i^*, \alpha_i) - E(\theta_i^*, \alpha_{\bar{i}}) \geq E(\theta_{\bar{i}}^*, \alpha_i) - E(\theta_{\bar{i}}^*, \alpha_{\bar{i}}) \qquad (24)$$

Substituting Eq. (20) into Eq. (24), and after some manipulations, then

$$- e^{-\theta_i^* \alpha_i C_B} + e^{-\theta_i^* \alpha_{\bar{i}} C_B} + e^{-\theta_{\bar{i}}^* \alpha_i C_B} - e^{-\theta_{\bar{i}}^* \alpha_{\bar{i}} C_B} \geq 0 \qquad (25)$$

Eq. (25) can be rewritten as

$$e^{-\theta_i^* \alpha_{\bar{i}} C_B} - e^{-\theta_i^* \alpha_i C_B} \geq e^{-\theta_{\bar{i}}^* \alpha_{\bar{i}} C_B} - e^{-\theta_{\bar{i}}^* \alpha_i C_B} \qquad (26)$$

Since $\theta_i^* > \theta_{\bar{i}}^*$, in order to satisfy the relationship in Eq. (26), $e^{-\theta_i^* \alpha_{\bar{i}} C_B} - e^{-\theta_i^* \alpha_i C_B}$ needs to be an increasing function over $\theta_i^*$. Then, it is defined as

$$F(\theta_i^*) = e^{-\theta_i^* \alpha_{\bar{i}} C_B} - e^{-\theta_i^* \alpha_i C_B} \qquad (27)$$

If $F(\theta_i^*)$ is an increasing function, then

$$\frac{\partial F(\theta_i^*)}{\partial \theta_i^*} = C_B (\alpha_i e^{-\theta_i^* \alpha_i C_B} - \alpha_{\bar{i}} e^{-\theta_i^* \alpha_{\bar{i}} C_B}) \geq 0 \qquad (28)$$

In order to derive the relationship between $\alpha_i$ and $\alpha_{\bar{i}}$, $\Delta_i = \theta_i^* C_B$ is defined. Then, Eq. (28) can be rewritten as

$$(\alpha_i e^{-\Delta_i \alpha_i} - \alpha_{\bar{i}} e^{-\Delta_i \alpha_{\bar{i}}}) \geq 0 \qquad (29)$$

When $\alpha_i > \frac{1}{\Delta_i}$ and $\theta_i^* > \theta_{\bar{i}}^*$, then $\alpha_i < \alpha_{\bar{i}}$. $\alpha_i$ and $\alpha_{\bar{i}}$ have the following relationship, i.e.,

$$\frac{1}{\Delta_i} \leq \alpha_i \leq \alpha_{\bar{i}} \qquad (30)$$

This completes the proof.

Moreover, if user's utility function satisfies the local downward incentive constraint (LDIC) and the local upward incentive constraint (LUIC) simultaneously, the IC constraints will be satisfied. Furthermore, the LUIC can be derived from the LDIC, and vice versa. The IC constraints can be replaced by

$$\rho_0 E(\theta_i^*, \alpha_i) - \pi_i = \rho_0 E(\theta_i^*, \alpha_{i+1}) - \pi_{i+1} \qquad (31)$$

By integrating the above lemmas, the original problem will be reduced to the following formulation.

$$\{\alpha_i^*, \pi_i^*\} = \operatorname{argmax} U_{RL} \qquad (32)$$

s. t. IR(19), IC(31), $\sum_{i=1}^{N} \alpha_i = 1, \alpha_i \geq 0, \forall i$

## 2.2 Optimal solution

By iterating IR and IC constraints shown in Eq. (32), the summation of $\pi_i$ can be written as a general equation, i.e.,

$$\sum_{i=1}^{N} \pi_i = \sum_{i=1}^{N} (N + 1 - i) \rho_0 E(\theta_i^*, \alpha_i)$$
$$- \sum_{i=1}^{N} (N - i) \rho_0 E(\theta_{i+1}^*, \alpha_i) \quad (33)$$

Substituting Eq. (33) into Eq. (32), the problem will be further reduced to the following brief formulation.

$$\{\alpha_i^*\} = \underset{i=1}{\arg\max} \sum_{i=1}^{N} C_i \quad (34)$$

s. t. $\sum_{i=1}^{N} \alpha_i = 1, \alpha_i \geq 0, \forall i$

where,

$$C_i = ((N + 1 - i) \rho_0 + \delta_0) E(\theta_i^*, \alpha_i)$$
$$- (N - i) \rho_0 E(\theta_{i+1}^*, \alpha_i) \quad (35)$$

It can be seen that each $C_i$ in Eq. (34) is only related to $\alpha_i$. In other words, it is not coupled with other $\alpha_i$. The first and second derivations of $C_i$ are

$$\frac{\partial C_i}{\partial \alpha_i} = ((N + 1 - i) \rho_0 + \delta_0) h \theta_i^* C_B e^{-\theta_i^* \alpha_i C_B}$$
$$- (N - i) \rho_0 h \theta_{i+1}^* C_B e^{-\theta_{i+1}^* \alpha_i C_B} \quad (36)$$

and

$$\frac{\partial C_i^2}{\partial^2 \alpha_i} = - ((N + 1 - i) \rho_0 + \delta_0) h (\theta_i^* C_B)^2 e^{-\theta_i^* \alpha_i C_B}$$
$$+ (N - i) \rho_0 h (\theta_{i+1}^* C_B)^2 e^{-\theta_{i+1}^* \alpha_i C_B} \quad (37)$$

respectively.

It is easy to verify that Eq. (36) $\geq 0$, and Eq. (37) $< 0$, Eq. (34) is a typical convex optimization problem. The optimal solution can be obtained by the interior point method. The numerical results are provided in the next section.

## 3 Numerical results

In this section, numerical simulation results are presented to demonstrate effectiveness of the proposed scheme. The simulation settings are listed in Table 1. Some simulation parameters are assumed, for example, $T_\alpha$, $T_\beta$, user numbers, $R_i$, total bandwidth and buffer resource capacity. Some simulation parameters are commonly used, for example, the noise density and path loss exponent. The specific simulation settings will be elaborated if the simulation scenario is changed.

Table 1　System parameters

| Parameter | Value |
|---|---|
| $T_\alpha$ | $T_\alpha = 0.4$ |
| $T_\beta$ | $T_\beta = 0.5$ |
| User number | $N = \{4, 5, 6\}$ |
| $R_i$ | $R_i \in [4,5]$ MB/s |
| Total bandwidth | $B = 0.5$ MHz |
| Transmission power | $P_{tr} = 23$ dBm[25] |
| Noise density | $\sigma_o^2 = -174$ dBm |
| Distance between relay and user | $d_i \in [20,40]$ m[26] |
| Path loss exponent | $l = 4$ |
| Buffer resource capacity | $C_B \in [60,600]$ MB |

Fig. 2 plots the variation of user's utility when the user chooses different contract entries. It can be seen that each user will obtain the maximum utility when it chooses the contract entry intended for its type. For instance, user 1 will obtain the maximum utility when it chooses the 1st contract entry. User 2 can achieve the maximum utility in the 2nd contract entry. The 3rd contract entry can make user 3 get the maximum utility. Similarly, contract entry 4th, 5th, and 6th will provide user 4, 5, and 6 with the maximum utility, respectively. The incentive compatibility can be verified from Fig. 2 that users cannot gain more utility by accepting a contract entry which is not designed for its type.
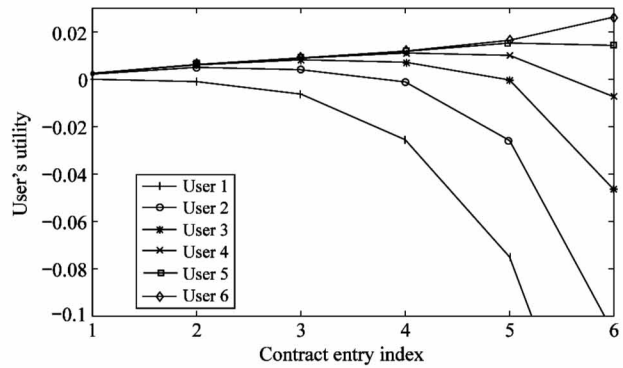


Fig. 2　The variation of users' utility values when the user chooses different contract entry

Table 2 presents the set of the parameter $\theta_i^*$, and Fig. 3 depicts the optimal contract entry which is a combination of the optimal portion $\alpha_i^*$ and the optimal price $\pi_i^*$. In this scenario, there are 6 users, and their data arrival rates $R = \{5, 4.8, 4.6, 4.4, 4.2, 4\}$ MB/s, respectively. Assuming that they are located in the same place, i.e., the distances between the relay node and the users are all 20 m. The overall buffer capacity in the relay node is 200 MB.

Table 2    The set of the parameter $\theta_i^*$

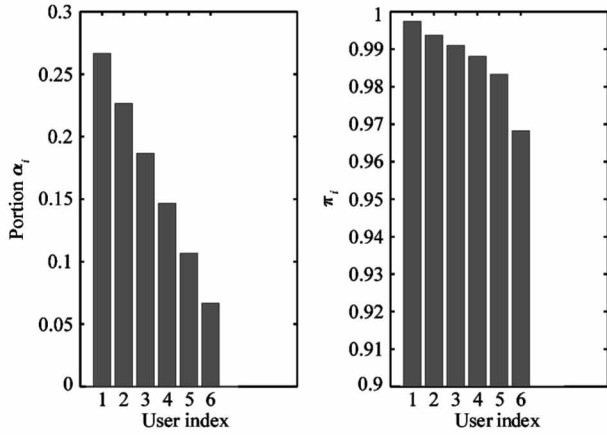| $\theta_1^*$ | $\theta_2^*$ | $\theta_3^*$ | $\theta_4^*$ | $\theta_5^*$ | $\theta_6^*$ |
|------|------|------|------|------|------|
| 0.112 | 0.148 | 0.191 | 0.243 | 0.308 | 0.39 |



Fig. 3    The optimal contract entry which is a combination of the optimal portion and the optimal price

The parameter $\theta_i^*$, i. e., user's type, is jointly determined by its data arrival rate and the service rate. The major difference among these users is their data arrival rate, and it can be noticed in Table 2 that the user with a higher data arrival rate has a smaller $\theta_i^*$, and vice versa. Parameter $\theta_i$ can reflect the inherent relationship between the arrival rate and the service rate, and further determine the backlogs inside the relay node. For example, when $\theta_i$ is small, i. e., $\theta_i = 0.112$, there are relatively more backlogs inside the relay node, while $\theta_i$ is large, i. e., $\theta_i = 0.39$, few backlogs exist inside the relay node. Accordingly, in Fig. 3, the optimal portion $\alpha_i^*$ grows along with the increase of user's index. This is consistent with Lemma 3 that if the contract entries satisfy the IC constraint, the quality $\alpha_i^*$ monotonically decreases with type $\theta_i^*$ increasing.

Moreover, in Fig. 3, the optimal portion $\alpha_i^*$ and optimal price $\pi_i^*$ designed for each user decrease as user's index increases. This is consistent with Lemma 2 that when $\alpha_i > \alpha_{\bar{i}}$, if and only if the price satisfies $\pi_i > \pi_{\bar{i}}$.

Fig. 4 compares the utility of the relay node and the users between the contract theory scheme and the equally allocated scheme. The equally allocated scheme evenly divides the buffer capacity. It can be seen from Fig. 4 that, the proposed contract theory scheme can achieve a larger utility compared to the benchmark scheme from the perspective of relay node. User's utility decreases as user's index increases when contract theory scheme is used. While this kind

of trend cannot be observed when the equally allocated scheme is employed. Moreover, the user with index 1 has the maximum utility. This observation is consistent with the intuition that user 1 purchases the maximum portion and pays the most to the relay node, it should gain the most utility. On the other hand, as shown in Fig. 4, the proposed scheme is superior to the equally allocated scheme from the perspective of relay node. The relay node who designs the mechanism will choose to use the proposed contract scheme. At the same time, users can obtain the positive utility by using the proposed contract scheme. Therefore, users would also be willing to participate in the proposed scheme.
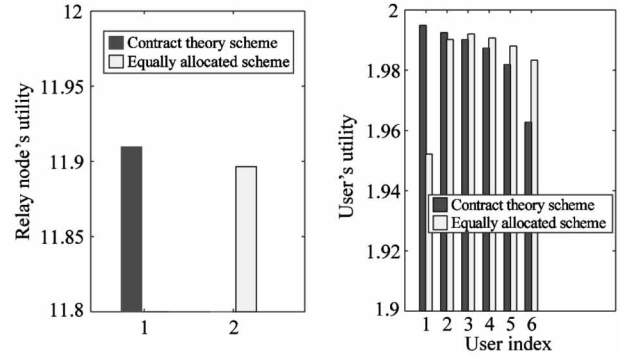


Fig. 4    Comparison between the contract theory approach and the equally allocated scheme from the perspective of utility

## 4   Conclusion

In this paper, the incentive mechanism incentivizing users in participating buffer-aided relay is investigated in wireless caching networks. Because of the information asymmetric environment, contract theory is utilized to design the incentive mechanism. Specifically, considering the limited buffer capacity, the backlog violation probability, i. e., the buffer overflow probability, is first provided based on the martingale theory. Based on the backlog violation probability, the utilities of relay node and users are formulated. Then, the optimal contract problem is modeled in order to maximize the utility of relay node, while the utilities of users are considered as IR and IC constraints. The feasibility of the contract is also demonstrated. Next, the optimal solution can be obtained by the interior point method. Numerical results are illustrated to demonstrate effectiveness of the proposed scheme. The proposed contract theory scheme has a better utility performance compared to the benchmarks.

**Reference**

[ 1] Wang X, Chen M, Taleb T, et al. Cache in the air: ex-

ploiting content caching and delivery techniques for 5G systems [J]. *IEEE Communications Magazine*, 2014 (52): 31-139

[2] Bastug E, Bennis M, Debbah M. Living on the edge: the role of proactive caching in 5G wireless networks [J]. *IEEE Communications Magazine*, 2014(52): 82-89

[3] Shanmugam K, Golrezaei N, Dimakis A G, et al. Femto-Caching: wireless contentdelivery through distributed caching helpers [J]. *IEEE Transactions on Information Theory*, 2013(59): 8402-8413

[4] Liu T, Li J, Shu F, et al. On the incentive mechanisms for commercialedge caching in 5G wireless networks [J]. *IEEE Wireless Communications*, 2018(25): 72-78

[5] Sung J, Kim M, Lim K, et al. Efficient cache placement strategy in two-tier wireless content delivery network [J]. *IEEE Transactions on Multimedia*, 2016 (18): 1163-1174

[6] Song J, Song H, Choi W. Which one is better to cache: requested contents or interfering contents? [J] *IEEE Wireless Communications Letters*, 2019(8): 861-864

[7] Jiao J, Hong X, Shi J. Proactive content delivery for vehicles over cellular networks: the fundamental benefits of computing and caching [J]. *China Communications*, 2018 (15): 88-97

[8] Kim J Y, Choi J K. Decentralized content delivery scheme using in-network caching [C]// Proceedings of the 2014 InternationalConference on Information and Communication Technology Convergence, Nanjing, China, 2014: 901-902

[9] Wang X, He J, Cheng P, et al. Privacy preserving collaborative computing: heterogeneous privacy guarantee and efficient incentive mechanism [J]. *IEEE Transactions on Signal Processing*, 2019(67): 221-233

[10] Chang C S, Thomas J A. Effective bandwidth in high-speed digital networks [J]. *IEEE Journal on Selected Areas in Communications*, 1995(13): 1091-1100

[11] Wu D P, Negi R. Effective capacity: a wireless link model for support of quality of service [J]. *IEEE Transactions on Wireless Communications*, 2003(2): 630-643

[12] Poloczek F, Ciucu F. Service-martingales: theory and applications to the delay analysis of random accessprotocols [C]//Proceedings of IEEE Conference on Computer Communications, Hong Kong, China, 2015: 945-953

[13] Hu Y, Li H, Chang Z, et al. End-to-end backlog and delay bound analysis for multi-hop vehicular Ad Hoc networks [J]. *IEEE Transactions on Wireless Communications*, 2017(16): 6808-6821

[14] Hu Y, Li H, Chang Z, et al. Scheduling strategy for multimedia heterogeneous high-speed train networks [J]. *IEEE Transactions on Vehicular Technology*, 2017 (66): 3265-3279

[15] Zhao L, Chi X, Zhu Y. Martingales-based energy-efficient D-ALOHA algorithms for MTC networks with delay-insensitive/URLLC terminals co-existence [J]. *IEEE Internet of Things Journal*, 2018(5): 1285-1298

[16] Liu T, Li J, Shu F, et al. Quality-of-service driven resource allocation based on martingale theory [C]// Proceedings of IEEE Global Communications Conference: Communication QoS, Reliability and Modeling, Abu Dhabi, UAE, 2018:1-6

[17] Bolton P, Dewatripont M. Contract Theory [M]. Massachusetts: MIT Press, 2005

[18] Gao L, Wang X, Xu Y, et al. Spectrum trading in cognitive radio networks: acontract-theoretic modeling approach [J]. *IEEE Journal on Selected Areas in Communications*, 2011(29): 843-855

[19] Li Y, Zhang J, Gan X, et al. A contract-based incentive mechanism for delayed traffic offloading in cellular networks [J]. *IEEE Transactions on Wireless Communications*, 2016(15): 5314-5327

[20] Liu T, Li J, Shu F, et al. Design of contract-based trading mechanism for a small-cell caching system [J]. *IEEE Transactions on Wireless Communications*, 2017 (16): 6602-6617

[21] Zhang Y, Song L, Saad W, et al. Contract-based incentive mechanisms for device-to-device communications in cellular networks [J]. *IEEE Journal on Selected Areas in Communications*, 2015(33): 2144-2155

[22] Zhang Y, Gu Y, Pan M, et al. Financing contract with adverse selection and moral hazard for spectrum trading in cognitive radio networks [C]// Proceedings of 2015 IEEE China Summit and InternationalConference on Signal and Information Processing, Chengdu, China, 2015: 601-605

[23] Zhang Y, Tran N H, Niyato D, et al. Multi-dimensional payment plan in fog computing with moralhazard [C]// Proceedings of IEEE International Conference on Communication System, Shenzhen, China, 2016: 1-6

[24] Shannon C E. A mathematical theory of communication [J]. *The Bell System Technical Journal*, 1948(27): 379-423

[25] Hong K, Xing D, Rai V, et al. Characterization of DSRC performance as a function of transmit power [C]// Proceedings of the 6th ACM international workshop on Vehicul ArInter NET working, Beijing, China, 2009: 63-68

[26] Morgan Y L. Notes on DSRC ampWAVE standards suite: its architecture, design, and characteristics [J]. *IEEE Communications Surveys Tutorials*, 2010(12): 504-518

**Liu Tingting**, born in 1982. She received the B. S. degree in communication engineering, and Ph. D. degree in information and communication engineering from Nanjing University of Science and Technology, Nanjing, China, in 2005 and 2011, respectively. Since 2011, she joined the School of Communication Engineering in Nanjing Institute of Technology, China. She was a post doctor in Nanjing University of Science and Technology. From 2017 to 2018, she was a visiting scholar in University of Houston, USA. Her research interests include game theory, blockchain, caching-enabled systems, mobile edge computing, network quality of service, device-to-device networks and cognitive radio networks.