







Joint Assortment and Cache Planning for Practical User Choice Model in Wireless Content Caching Networks

Yaru Fu , Member, IEEE, Xinyu Xu , Member, IEEE, Hanlin Liu , Member, IEEE, Quan Yu , Member, IEEE, Hong-Ning Dai , and Tony Q. S. Quek , Fellow, IEEE

Abstract—In wireless content caching networks (WCCNs), a user’s content consumption crucially depends on the assortment offered. Here, the assortment refers to the recommendation list. An appropriate user choice model is essential for greater revenue. Therefore, in this paper, we propose a practical multinomial logit choice model to capture users’ content requests. Based on this model, we first derive the individual demand distribution per user and then investigate the effect of the interplay between the assortment decision and cache planning on WCCNs’ achievable revenue. A revenue maximization problem is formulated while incorporating the influences of the screen size constraints of users and the cache capacity budget of the base station (BS). The formulated optimization problem is a non-convex integer programming problem. For ease of analysis, we decompose it into two folds, i.e., the personalized assortment decision problem and the cache planning problem. By using structure-oriented geometric properties, we design an iterative algorithm with examinable quadratic time complexity to solve the non-convex assortment problem in an optimal manner. The cache planning problem is proved to be a 0-1 Knapsack problem and thus can be addressed by a dynamic programming approach with pseudo-polynomial time complexity. Afterwards, an alternating optimization method is used to optimize the two types of variables until convergence. It is shown by simulations that the proposed scheme outperforms various existing benchmark schemes.

Index Terms—Cache planning, personalized assortment decision, revenue optimization, user’s choice model.

I. INTRODUCTION

WITH Internet of Everything (IoE) applications being increasingly popular, the envisioning and development for the sixth generation of wireless cellular networks (6G) have already commenced. 6G is anticipated to support various intelligent applications with diversified requirements (e.g., latency, reliability, and users’ quality of service). However, the enormous gap between the diverse requirements and what edge servers can provide is a grave challenge for 6G in terms of providing high-quality services for emerging intelligent applications [2], [3]. The challenges are also raised from the fact that popular data are frequently requested by a majority of users and thus redundantly transmitted over cellular networks, making both the fronthaul and the backhaul links congested. To counter these tremendous challenges, both the academia and the industry have shifted their attentions to exploring the edge-aware techniques. In this context, wireless content caching has been acknowledged as a successful enabler [4]. More specifically, by pre-caching reusable contents at network edge facilities during the off-peak time intervals, redundantly generated data traffic can be significantly reduced. As a consequence, highly effective networks with shortened delays and improved users’ quality of experience become a reality.

In wireless content caching networks (WCCNs), users’ content consumption is highly influenced by the assortment decision, which refers to the selection of a set of content items that are made available to users. It is widely acknowledged that having an appropriate user choice model is essential for increasing revenue, as it enables service providers to offer personalized content recommendations that are likely to be of interest to users. However, developing an effective choice model can be challenging. In addition, the assortment decision has a significant impact on the cache planning at the BS, which involves deciding which content items to store at the BS’s cache. Given the limited storage capacity of the cache, it is crucial to select the content items that are most likely to be requested by users. This requires a careful analysis of users’ content consumption patterns and preferences, which can be

Manuscript received 15 February 2023; revised 6 June 2023; accepted 10 July 2023. Date of publication 24 July 2023; date of current version 4 April 2024. This work was supported in part by the grant from the Research Grants Council (RGC) of the Hong Kong Special Administrative Region, China under Grant UGC/FDS16/E02/22, in part by the National Natural Science Foundation of China under Grant 72102098 and in part by the Guangdong Basic and Applied Basic Research Foundation under Grant 2020A1515111131, and in part by the National Research Foundation, Singapore and Infocomm Media Development Authority under its Future Communications Research & Development Programme. Recommended for acceptance by K.A. Harras. (Corresponding author: Hanlin Liu.)

Yaru Fu is with the Department of Electronic Engineering and Computer Science, Hong Kong Metropolitan University, Hong Kong 999077, China (e-mail: yfu@hkmu.edu.hk).

Xinyu Xu and Hanlin Liu are with the College of Business, Southern University of Science and Technology, Shenzhen, Guangdong 518055, China (e-mail: 12131295@mail.sustech.edu.cn; liuhl@sustech.edu.cn).

Quan Yu is with the School of Information Engineering, Wuhan University of Technology, Wuhan, Hubei 430062, China (e-mail: yuquan@whut.edu.cn).

Hong-Ning Dai is with the Department of Computer Science, Hong Kong Baptist University, Hong Kong 999077, China (e-mail: hndai@ieee.org).

Tony Q. S. Quek is with the Singapore University of Technology and Design, Singapore 487372, and also with Department of Electronic Engineering, Kyung Hee University, Yongin 17104, South Korea (e-mail: tonyquek@sutd.edu.sg).

Digital Object Identifier 10.1109/TMC.2023.3297987

supported by the user choice model. To maximize revenue, these two aspects, i.e., the assortment decision and the cache planning, should be jointly optimized. This involves developing a comprehensive optimization framework that takes into account both the user choice model and the cache planning strategy.

In this paper, we propose a revised multinomial logit (RMNL) model that aims to address the limitations of the conventional model in providing a more comprehensive perspective for content items. We derive the content request distribution per user, and formulate a revenue maximization problem by collaboratively designing the personalized assortment decision and the cache planning for WCCNs. The main contributions of our work are sorted as follows:

- We present a model called RMNL that explicitly shows how users' content request behaviors can be influenced by their personalized assortment strategies. This model is more versatile than the conventional MNL model and addresses its shortcomings. We investigate the effect of the interplay between assortment decision and cache planning on the achievable revenue of WCCNs. Based on these analyses, we formulate the revenue maximization problem for WCCNs, taking into account the cache capacity budget of the BS and the screen size requirement per user.
- The formulated optimization problem is a non-convex integer programming problem, whose optimal solution is challenging to obtain. The difficulty mainly stems from the coupling among the optimization variables. To facilitate the analysis, we decompose the original problem into two subproblems, i.e., the assortment decision subproblem and the cache placement subproblem. For the assortment problem, some structure-oriented geometric insights are provided. In particular, we show that the optimal assortment set can be searched within the intersection points among $I + 1$ straight (linear) lines, where I is the total number of contents. From this, an iterative algorithm with quadratic time complexity is developed to achieve the globally optimal solution. Meanwhile, by proceeding with some transformations, we prove that the cache planning subproblem is a 0-1 Knapsack problem, which can be optimally solved by a dynamic programming algorithm.
- With the aforementioned two-fold analysis, an alternating optimization approach is used to jointly optimize the two types of Boolean decision variables. This approach is a novel and sophisticated way to tackle the problem of optimizing WCCNs, and has shown great promise in our study. Here, the associated computational complexity analysis and the convergence analysis for the proposed joint optimization method are discussed in detail, providing a thorough understanding of the nuances of the approach. It is demonstrated that our developed joint optimization method has global convergence and polynomial time complexity, thus can meet the needs of large WCCNs.

Extensive computer simulations are performed from three aspects, including convergence performance, system revenue, and cache hit ratio, to reveal the validity and superiority of our proposed scheme when compared with multiple benchmark strategies. Moreover, to ensure that the simulation accurately

reflected real-world scenarios, we used both homogeneous and heterogeneous assortment sizes among users. Our scheme proved to be robust and efficient, making it a viable option for improving system performance and user experience. The rest of this paper is organized as follows. In Section II, the related works are reviewed. In Section III, we introduce the system model of WCCNs together with the preliminaries of the MNL choice model, followed by our developed RMNL model. In Section IV, we formulate the revenue maximization problem taking into account varied practical constraints. The designed methodology for solving the original optimization problem and the associated property analysis are elaborated in Section V. In Section VI, numerical simulations are conducted to evaluate the performance of our designed strategy compared with extensive benchmarks. At last, we summarize this paper and set forth the future research directions in Section VII.

For ease of reading, the notations used throughout this paper are summarized in Table I.

II. LITERATURE REVIEW

In recent years, enormous attention has been paid to WCCNs. For instance, the authors in [5] analyzed the files successful transmission probability (STP) of random caching-based large-scale heterogeneous wireless systems. Simulation results showed that the random caching with retransmission scheme can dramatically improve STP. Besides, it was demonstrated in [6] that with users' preference information, the optimized cache placement results in boosted cache hit ratio. In addition, the effectiveness of proactive caching in terms of improving spectrum efficiency and energy efficiency was verified by the authors in [7] and [8], respectively. Moreover, the authors in [9], [10] validated the performance of user-side caching oriented cellular networks. In these studies, the multi-winner auction approach and the two-side swapping method were used to address the cache placement problem. Furthermore, the authors in [11] jointly optimized the incentive mechanism and the price-based cache replacement to improve the profits of both Internet service providers and content providers. The joint cache rental, content placement, and user association problem for WCCNs was investigated in [12] to maximize the aggregated delay savings. The same objective was studied in [13], [14] with the consideration of recommendations. It was affirmed that recommendation-aware caching can achieve superior performance than the pure caching strategies. To address the data corruption issue, coded caching was utilized in [15], [16]. Thereof, repairable codes were designed for both the single-failure and the multi-failure cases. The applicability of content caching in other different usage scenarios, such as unmanned aerial vehicle-assisted (UAV) networks and wireless digital twin networks (WDTNs), was verified in [17] and [18], respectively. Furthermore, to break the curses of conventional optimization tools, the authors in [19], [20], [21], [22], [23] applied machine learning-based solutions, such as long-short-term-memory network, echo state network, recurrent neural network, and reinforcement learning model, to tackle the joint cache planning and transmission optimization problems for WCCNs.

TABLE I
LIST OF NOTATIONS

Notation	Definition	Notation	Definition
\mathcal{K}	Index set of users.	$\varepsilon_{k,i}$	IID random variable.
K	Total number of users.	$V_{k,i}$	Preference of user k to item i .
\mathcal{I}	Index set of content .	$p'_{k,i}$	Request probability under MNL model.
I	Total Number of contents.	\mathcal{S}_k	Assortment set of user k .
w'_i	Revenue of content i .	$p_{k,i}$	Request probability under RMNL model.
B_i	Bit size of content i .	R	System revenue.
c_i	Cache decision of content i .	\mathbf{c}	System's cache planning policy.
C_{BS}	Cache capacity of BS.	\mathbf{a}_k	User k 's assortment strategy.
α_i	Discount factor of content i .	\mathbf{a}	Assortment strategy of the system.
w_i	Effective profit of content i .	R_k	Revenue brings by user k .
$a_{k,i}$	Assortment decision of content i to user k .	$f'(\mathcal{S})$	Revenue of items in \mathcal{S} .
\bar{a}_k	Assortment set size of user k .	$f(\mathcal{S})$	Revenue under assortment strategy \mathcal{S} .
$U_{k,i}$	Random utility of user k to content i .	\mathbf{a}^t	Assortment solution in the t -th iteration.
$U_{k,0}$	Utility of outside option.	\mathbf{c}^t	Caching policy in the t -th iteration.
$\mu_{k,i}$	Mean utility of user k to content i .	N	Maximum iteration number.

Although the effectiveness of content caching has been validated by the literature, user's demand behavior therein is mainly assumed to follow Zipf's law for simplicity. Notably, Zipf's law was proposed to measure word frequencies in text [24] and the popularity of internet pages [25]. Therefore, this approximation is somewhat less accurate particularly when the heterogeneity among users' personalities is considered, i.e., each user has its own inherent content preference distribution. Moreover, a user's content selection behavior often depends on the assortment list provided by the system. Thereby, an appropriate user's choice model is pivotal for assortment decision making, which further has a dramatic impact on the system's revenue. Modeling user's (customer's) choice behavior is a compelling research area in retail and online advertisement. Thereof, user's demanding is usually characterized by MNL choice model [26], [27], [28], [29], in which the system (e.g., service, product, or content provider) shows a set of items to the users and users demand the items based on a probability distribution determined by the offered assortment. A practical application is in retail. Given a limited shelf capacity, a retailer has to determine the assortment of products that maximizes the profit. However, existing works [26], [27], [28], [29] assumed that users can only make purchase decisions among the offered assortment or leave without purchase (in the example of retail, a customer can only buy the item among the assortment on the shelves or leave without buying anything), which may not align with real-world situations. Take the streaming media platforms (such as iQIYI, YouTube or Netflix) as an example, users tend to observe the (personalized) recommended video contents (we refer as assortment throughout this work) in the home page, and they can choose the contents within the offered assortment in concordance with their attributes or choose the other options. When the latter happens (which means the offered lists are not satisfactory), the users may close (or ignore) the offered assortment list and browse the remaining contents (e.g., by clicking "More"). In this way, the users may experience another stage of choosing contents. Based on the above analysis, a more practical user's choice model is needed so as to align with practices, which motivates this work.

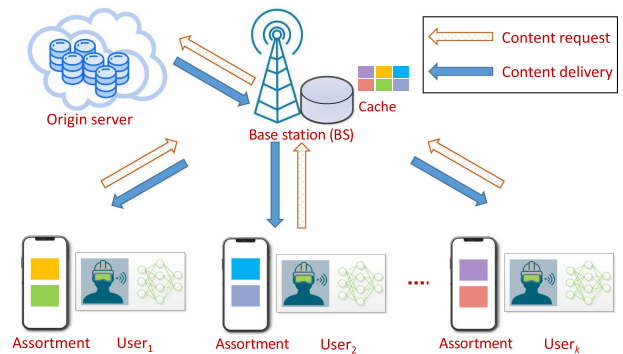


Fig. 1. Generic system model of WCCNs with personalized assortment decision.

III. SYSTEM DESCRIPTION AND THE IMPROVED USER CHOICE MODEL

In this section, we first introduce the system model of the considered WCCNs. Then, we introduce the preliminaries regarding the MNL model and outline its deficiencies. Afterwards, we elaborate on our proposed user choice model and derive the specific expression of users' content demanding distribution.

A. System Description

As illustrated in Fig. 1, we consider a generic WCCN with one BS serving a set of K users, referred to as $\mathcal{K} = \{1, 2, \dots, K\}$. Here, the BS has the capability of caching such that it can proactively store some popular contents. Therein, popular contents refer to the items that are most likely to be requested by users. To characterize this, we assume that in total there are I content items in the catalog, whose index set is referred to as $\mathcal{I} = \{1, 2, \dots, I\}$. In addition, for $i \in \mathcal{I}$, define w'_i and B_i as the associated marginal revenue¹ and bit size of content i , respectively. Besides, let $c_i \in \{0, 1\}$ be the cache decision indicator of content i . More specifically, $c_i = 1$ represents that

¹We assume that all the content items are priced as that in [30]. Details will be given in Section VI.

content i has been cached at the BS and $c_i = 0$ otherwise. As the storage size of the BS is limited, we have the following constraint:

$$\sum_{i \in \mathcal{I}} c_i B_i \leq C_{BS}, \quad (1)$$

in which C_{BS} is the cache capacity budget of the BS.

Given that the demanded content is cached at the BS, BS can delivery this content to the user directly. On this circumstance, the user can enjoy an efficient data service. As a reward, this content associated profit (e.g., w'_i) can be obtained by the system. Otherwise, the BS needs to first access the remote clouds or origin content servers to retrieve this content item and then feedback it to the corresponding user. Since additional communication cost is required, the expected profit of the un-cached item will be watered-down, which is depicted by $\alpha_i \in (0, 1)$ for the i -th content². Given that $c_i = 0$, the associated profit of content i becomes $\alpha_i w'_i$, where $i \in \mathcal{I}$. With the analysis above, the effective (actual) profit of content i , denoted by w_i , can be expressed as follows:

$$w_i = \alpha_i^{1-c_i} w'_i, \quad i \in \mathcal{I}. \quad (2)$$

For the sake of completeness, the outside option³ is denoted by 0 with $w_0 = 0$.

B. Conventional User Choice Model

In this subsection, we introduce the preliminaries of the conventional user choice model. For expression simplicity, we define a binary variable $a_{k,i} \in \{0, 1\}$ for $k \in \mathcal{K}$ and $i \in \mathcal{I}$ to depict whether content i is shown to user k or not. More precisely, $a_{k,i} = 1$ means that content item i is presented to user k and $a_{k,i} = 0$ otherwise. In addition, let \bar{a}_k be the maximal number of items that can be displayed for user k , considering the limited screen size of user k , where $k \in \mathcal{K}$. Thus, we have

$$\sum_{i \in \mathcal{I}} a_{k,i} \leq \bar{a}_k, \quad k \in \mathcal{K}. \quad (3)$$

In the MNL choice model, each user chooses the content that maximizes its own utility. Let $U_{k,i} = \mu_{k,i} + \varepsilon_{k,i}$ be the utility of user k in terms of demanding content item i , where $k \in \mathcal{K}$ and $i \in \{0\} \cup \mathcal{I}$. Specifically, $i = 0$ depicts that the user consumes nothing from the system. We assume $\mu_{k,0} = 0$, then $U_{k,0}$ represents the utility of outside option. In addition, $\mu_{k,i}$ expresses the mean utility that the user k assigns to content i , whilst $\{\varepsilon_{k,i}\}_{i \in \{0\} \cup \mathcal{I}}$ are independent and identically distributed (IID) random variables, each have a Gumbel distribution with location-scale parameters (0,1) [32].⁴ Moreover, we define $V_{k,i} = e^{\mu_{k,i}}$ as the preference of user k towards content i and let $V_{k,0} = 1$, where $k \in \mathcal{K}$ and $i \in \mathcal{I}$. Furthermore, define \mathcal{S}_k as

²It is worth mentioning that the value of α_i can be determined by the physical distance between the BS and the remote cloud or origin server [31]. Other portrayals are equally appropriated.

³In this paper, outside option means that the user purchases nothing from the content catalog.

⁴It is noteworthy that a Gumbel random variable with location-scale parameters (μ, η) has cumulative distribution function $F(x) = \exp(-e^{-(x-\mu)/\eta})$. This will be used in Lemma 1.

the content list displayed to user k . Then the probability of user k selecting content $i \in \mathcal{S}_k \cup \{0\}$ is given as follows [32]:

$$p'_{k,i} = \mathbb{P} \left(U_{k,i} \geq \max_{j \in \mathcal{S}_k \cup \{0\}} \{U_{k,j}\} \right) = \frac{V_{k,i}}{\sum_{j \in \mathcal{S}_k \cup \{0\}} V_{k,j}}. \quad (4)$$

It is worth noting that under the conventional MNL choice model, $p'_{k,i} = 0$ for $i \in \mathcal{I} \setminus \mathcal{S}_k$, where $k \in \mathcal{K}$. For completeness, we have

$$p'_{k,i} = \begin{cases} \frac{V_{k,i}}{\sum_{j \in \mathcal{S}_k \cup \{0\}} V_{k,j}}, & \text{if } i \in \mathcal{S}_k \cup \{0\}, \\ 0, & \text{if } i \in \mathcal{I} \setminus \mathcal{S}_k. \end{cases} \quad (5)$$

C. Revised User Choice Model

It can be observed from (5) that the system shows an assortment of content items to users, and a user can purchase the items among the assortment or leave without purchasing anything. In other words, the selected content items of user k can only be from the set $\mathcal{S}_k \cup \{0\}$, where $k \in \mathcal{K}$. In most cases, customers tend to observe the assorted products and select within the assortment in accordance with the individual preferences of all offered items. Given that the assortment offered is not satisfactory, the users can still choose the other items. To align with realities in practical systems, an improved choice model is designed in this paper. Therein, we assume that a second phase is available for users to select the remaining content items provided that the pre-set assortment set is unsatisfactory. Doing so, users have a wider range of choices. Based on the analysis, we explicitly derive the individual demanding distribution per user. Before giving the details, some properties of Gumbel random variables are introduced, which will be used in the derivation.

Lemma 1. Let X_1 and X_2 be two independent Gumbel random variables with location-scale parameters $(\mu_1, 1)$ and $(\mu_2, 1)$, respectively. Then we have the following properties:

- $\max(X_1, X_2)$ is Gumbel distributed with location-scale parameters $(\log(e^{\mu_1} + e^{\mu_2}), 1)$.
- $\mathbb{P}(X_1 \geq X_2) = \frac{e^{\mu_1}}{e^{\mu_1} + e^{\mu_2}}$.
- $\max(X_1, X_2)$ and $\mathbb{I}(X_1 \geq X_2)$ are independent random variables, where $\mathbb{I}(\cdot)$ represents the indicator function.

The proof of Lemma 1 can be found in [29], [33]. Details are not copied here for brevity. Thereafter, we illuminate the theorem with regard to each user's content request probability under our developed RMNL model.

Theorem 2. Under the proposed RMNL choice model, the probability of user k selecting content $i \in \mathcal{I} \cup \{0\}$, defined as $p_{k,i}$, is given below:

$$p_{k,i} = \begin{cases} \frac{V_{k,i}}{1 + V_{\mathcal{S}_k}}, & \text{if } i \in \mathcal{S}_k, \\ \frac{V_{k,i}}{(1 + V_{\mathcal{S}_k})(1 + V_{\mathcal{I}})}, & \text{if } i \in \mathcal{I} \setminus \mathcal{S}_k, \\ \frac{1}{1 + V_{\mathcal{I}}}, & \text{if } i = 0, \end{cases} \quad (6)$$

where $V_{\mathcal{S}_k} = \sum_{i \in \mathcal{S}_k} V_{k,i}$ and $V_{\mathcal{I}} = \sum_{j \in \mathcal{I}} V_{k,j}$, respectively.

Proof. We distinguish the k -th user's content request, referred to as i , into three cases. That is, $i \in \mathcal{S}_k$, $i \in \mathcal{I} \setminus \mathcal{S}_k$, and $i = 0$. The proofs are expanded from the aforementioned three aspects.

1) **Case I:** $i \in \mathcal{S}_k$

When $i \in \mathcal{S}_k$, by properties (a) and (b) of Lemma 1, we have

$$\begin{aligned} p_{k,i} &= \mathbb{P} \left(U_{k,i} \geq \max_{j \in \mathcal{S}_k \cup \{0\} \setminus \{i\}} U_{k,j} \right) \\ &= \frac{e^{\mu_{k,i}}}{\sum_{j \in \mathcal{S}_k \cup \{0\}} e^{\mu_{k,j}}} \\ &= \frac{V_{k,i}}{1 + V_{\mathcal{S}_k}}. \end{aligned} \quad (7)$$

2) **Case II:** $i \in \bar{\mathcal{S}}_k$, where $\bar{\mathcal{S}}_k = \mathcal{I} \setminus \mathcal{S}_k$.

Given that the requested content of user k falls in Case II, the utility of the outside option must be greater than the utilities of all contents in \mathcal{S}_k . Moreover, the utility of the requested content i must be greater than the utilities of all contents in $\bar{\mathcal{S}}_k \setminus \{i\}$ as well as the outside option. That is,

$$\begin{aligned} p_{k,i} &= \mathbb{P} \left(U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j}, U_{k,i} \right. \\ &\quad \left. \geq \max \left\{ U_{k,0}, \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \right\} \right) \\ &= \mathbb{P} \left(U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) \\ &\quad \times \mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \right\} \right) \\ &\quad \times \mathbb{P} \left(U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right). \end{aligned} \quad (8)$$

According to properties (a) and (b) of Lemma 1, the first item of the multiplication on the right hand side (RHS) of (8) can be expressed as follows:

$$\mathbb{P} \left(U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) = \frac{1}{\sum_{j \in \mathcal{S}_k} e^{\mu_{k,j}}} = \frac{1}{1 + V_{\mathcal{S}_k}}. \quad (9)$$

On the basis of property (a) of Lemma 1, we know that $\max_{j \in \mathcal{S}_k} U_{k,j}$ is Gumbel distributed with parameter $(\log(\sum_{j \in \mathcal{S}_k} e^{\mu_{k,j}}), 1)$. Likewise, $\max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j}$ is also Gumbel distributed, wherein the location-scale parameter is $(\log(\sum_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} e^{\mu_{k,j}}), 1)$. In addition, based on the property (c) of Lemma 1, those random variables $U_{k,0}, U_{k,i}, \max_{j \in \mathcal{S}_k} U_{k,j}$, and $\max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j}$ are mutually independent. On these grounds, the second item of the multiplication on the RHS of (8) can be expressed as (10)-(16) shown at the bottom of this page.

More specifically, (10), (11), and (12) hold because of the basic properties of conditional probability. (13) is obtained due to the fact that $\max\{U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j}\}$ and $\mathbb{I}(U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j})$ are independent in accordance with the property (c) of Lemma 1, and besides $U_{k,0}, U_{k,i}, \max_{j \in \mathcal{S}_k} U_{k,j}$ are independent. In addition, (14) holds because $\max\{U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j}\}$ and $\mathbb{I}(U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j})$ are independent, together with the fact that $U_{k,0}, U_{k,i}, \max_{j \in \mathcal{S}_k} U_{k,j}, \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j}$ are independent. Moreover, (15) again follows from the basic properties of conditional probability. Furthermore, by (a) and (b) of Lemma 1, (16) is accessible.

$$\begin{aligned} &\mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \right\} \mid U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) \\ &= \mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j}, \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \right\} \mid U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) \end{aligned} \quad (10)$$

$$= \mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j} \right\}, U_{k,i} \geq \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \mid U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) \quad (11)$$

$$\begin{aligned} &= \mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j} \right\} \mid U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) \\ &\quad \times \mathbb{P} \left(U_{k,i} \geq \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \mid U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j} \right\}, U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) \end{aligned} \quad (12)$$

$$\begin{aligned} &= \mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j} \right\} \right) \\ &\quad \times \mathbb{P} \left(U_{k,i} \geq \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \mid U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j} \right\}, U_{k,0} \geq \max_{j \in \mathcal{S}_k} U_{k,j} \right) \end{aligned} \quad (13)$$

$$= \mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j} \right\} \right) \mathbb{P} \left(U_{k,i} \geq \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \mid U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j} \right\} \right) \quad (14)$$

$$= \mathbb{P} \left(U_{k,i} \geq \max \left\{ U_{k,0}, \max_{j \in \mathcal{S}_k} U_{k,j}, \max_{j \in \bar{\mathcal{S}}_k \setminus \{i\}} U_{k,j} \right\} \right) \quad (15)$$

$$= \frac{e^{\mu_{k,i}}}{\sum_{j \in \mathcal{I} \cup \{0\}} e^{\mu_{k,j}}} = \frac{V_{k,i}}{1 + V_{\mathcal{I}}}. \quad (16)$$

By substituting (9) and (16) into (8), we derive

$$p_{k,i} = \frac{V_{k,i}}{(1 + V_{S_k})(1 + V_{\mathcal{I}})}, \text{ if } i \in \mathcal{I} \setminus S_k.$$

3) **Case III:** $i = 0$

In this case, the probability of the user in terms of purchasing nothing after foregoing two display phases is

$$\begin{aligned} p_{k,0} &= 1 - \sum_{j \in S_k} p_{k,j} - \sum_{j \in \mathcal{I} \setminus S_k} p_{k,j} \\ &= 1 - \sum_{j \in S_k} \frac{V_{k,j}}{1 + V_{S_k}} - \sum_{j \in \mathcal{I} \setminus S_k} \frac{V_{k,j}}{(1 + V_{S_k})(1 + V_{\mathcal{I}})} \\ &= \frac{1}{1 + V_{\mathcal{I}}}. \end{aligned}$$

The aforementioned three-fold analyses complete the proof.

In accordance with the foregoing discussions, the achievable revenue of the system under RMNL, which is defined as R , can be written as follows:

$$R = \sum_{k \in \mathcal{K}} \sum_{i \in \mathcal{I}} w_i p_{k,i}, \quad (17)$$

where w_i and $p_{k,i}$ are given in (2) and (6), respectively. It can be seen that the system's achievable revenue crucially depends on the assortment decision per user as well as the cache planning decision at the BS.

IV. PROBLEM FORMULATION

In this paper, we aim to maximize the WCCNs' achievable revenue by collaboratively optimizing the assortment decision (per user) and the cache planning at the BS, taking into account the constraints of users' screen size and cache capacity budget of the BS. For notation simplicity, we define $\mathbf{c} = (c_i)_{i \in \mathcal{I}}$ as the caching policy of the system. In addition, let $\mathbf{a}_k = (a_{k,i})_{i \in \mathcal{I}}$ and $\mathbf{a} = (\mathbf{a}_k)_{k \in \mathcal{K}}$ be the assortment strategy of user k and all users, respectively. With the definitions, the revenue maximization problem for WCCNs is mathematically formulated as follows:

$$\begin{aligned} & \underset{\mathbf{a}, \mathbf{c}}{\text{maximize}} R \\ \text{s.t. } C1: & \sum_{i \in \mathcal{I}} a_{k,i} \leq \bar{a}_k, \quad k \in \mathcal{K}, \\ C2: & \sum_{i \in \mathcal{I}} c_i B_i \leq C_{\text{BS}}, \\ C3: & a_{k,i} \in \{0, 1\}, \quad k \in \mathcal{K}, \quad i \in \mathcal{I}, \\ C4: & c_i \in \{0, 1\}, \quad i \in \mathcal{I}, \end{aligned} \quad (18)$$

where R in the objective function is defined in (17). $C1$ represents that the total number of content items in the assortment set of user k cannot be larger than \bar{a}_k . $C2$ shows that the cached items cannot exceed the storage capacity budget of the BS. In addition, $C3$ and $C4$ depict that all the optimization variables are binary variables. As can be seen from (6), $p_{k,i}$ is a non-convex function with respect to \mathbf{a} . Thus, the optimization problem (18) is a non-convex integer programming problem, whose globally optimal solution is difficult to obtain. The main difficulty is

originated from the coupling among the optimization variables. To make the analysis tractable, in the following, we decompose the original maximization problem into two aspects, namely, the assortment decision problem and the cache planning problem. Then we dedicate to solving each of the two subproblems optimally in an efficient manner, and finally optimize the two types of optimization variables alternately. The convergence performance with regard to the developed solution and its time complexity are analyzed as well. It is shown that our devised solution has polynomial time complexity and thus can satisfy the needs (e.g., computational latency, and complexity) of large sized WCCNs.

V. METHODOLOGY AND PROPERTY ANALYSIS

In this section, we first introduce the designed solution to problem (18). Then, we analyze its optimality, time complexity, and convergence performance.

A. Algorithm Design

To solve (18) in an efficient manner, we decompose the originally formulated problem into two subproblems, namely, the assortment decision subproblem and the cache placement subproblem. In this subsection, we first solve each of the subproblems and then use the alternating optimization approach to solve the whole problem.

1) *Assortment Decision Subproblem:* Given the caching policy, i.e., \mathbf{c} , the original optimization problem (18) degenerates to the assortment decision-making subproblem. As each user's assortment strategy is independent of the other users' policies, the assortment decision-making subproblem can be further divided into K subproblems. We take the assortment decision problem for user k as an example, and the mathematical formulation is given below:

$$\begin{aligned} & \underset{\mathbf{a}}{\text{maximize}} R_k && \mathcal{P}(1) \\ \text{s.t. } C1': & \sum_{i \in \mathcal{I}} a_{k,i} \leq \bar{a}_k, \\ C3': & a_{k,i} \in \{0, 1\}, \quad i \in \mathcal{I}, \end{aligned}$$

wherein R_k is defined as the revenue brought by user k . That is,

$$R_k = \sum_{i \in \mathcal{I}} w_i p_{k,i}, \quad k \in \mathcal{K}. \quad (19)$$

Note that $\mathcal{P}(1)$ is still a non-convex integer programming problem. We propose to attain its globally optimal solution efficiently by using its inherent structural properties. To achieve this goal, some definitions are made, which will be used in the following analyses. First of all, we remove the notation of k for expression simplicity. Then, we define $f'(\mathcal{S})$ as the achievable revenue of the content items in \mathcal{S} . According to Theorem 1, we have

$$f'(\mathcal{S}) = \frac{\sum_{i \in \mathcal{S}} V_i w_i}{1 + \sum_{i \in \mathcal{S}} V_i}. \quad (20)$$

Similarly, we denote the achievable revenue of the contents outside \mathcal{S} as $f'(\mathcal{I} \setminus \mathcal{S})$, which is obtained as follows:

$$f'(\mathcal{I} \setminus \mathcal{S}) = \frac{\sum_{i \in \mathcal{I} \setminus \mathcal{S}} V_i w_i}{(1 + \sum_{i \in \mathcal{S}} V_i)(1 + \sum_{i \in \mathcal{I}} V_i)}. \quad (21)$$

Moreover, let $f(\mathcal{S})$ be the revenue generated by user k 's content requests. Based on (20) and (21), $f(\mathcal{S})$ is expressed as follows

$$f(\mathcal{S}) = f'(\mathcal{S}) + f'(\mathcal{I} \setminus \mathcal{S}). \quad (22)$$

In accordance with (20)-(22), $f(\mathcal{S})$ can be rewritten as

$$\begin{aligned} f(\mathcal{S}) &= \frac{\sum_{i \in \mathcal{S}} V_i w_i}{1 + \sum_{i \in \mathcal{S}} V_i} + \frac{\sum_{i \in \mathcal{I} \setminus \mathcal{S}} V_i w_i}{(1 + \sum_{i \in \mathcal{S}} V_i)(1 + \sum_{i \in \mathcal{I}} V_i)}, \\ &= \frac{\sum_{i \in \mathcal{S}} V_i w_i}{1 + \sum_{i \in \mathcal{S}} V_i} + \frac{\sum_{i \in \mathcal{I}} V_i w_i - \sum_{i \in \mathcal{S}} V_i w_i}{(1 + \sum_{i \in \mathcal{S}} V_i)(1 + \sum_{i \in \mathcal{I}} V_i)}. \end{aligned}$$

Let $P = \sum_{i \in \mathcal{I}} V_i w_i$ and $Q = 1 + \sum_{i \in \mathcal{I}} V_i$, we have

$$f(\mathcal{S}) = \frac{\sum_{i \in \mathcal{S}} V_i w_i}{1 + \sum_{i \in \mathcal{S}} V_i} + \frac{P - \sum_{i \in \mathcal{S}} V_i w_i}{(1 + \sum_{i \in \mathcal{S}} V_i)Q}. \quad (23)$$

With the above definitions, $\mathcal{P}(1)$ can be equally transformed into finding Y^* as defined below:

$$Y^* = \max \{ \lambda \in \mathbb{R} : \exists \mathcal{S} \subseteq \mathcal{I}, |\mathcal{S}| \leq \bar{a}, \text{ and } f(\mathcal{S}) \geq \lambda \}, \quad (24)$$

wherein, based on (23), $f(\mathcal{S}) \geq \lambda$ can be further rephrased as follows:

$$\begin{aligned} f(\mathcal{S}) &= \frac{\sum_{i \in \mathcal{S}} V_i w_i}{1 + \sum_{i \in \mathcal{S}} V_i} + \frac{P - \sum_{i \in \mathcal{S}} V_i w_i}{(1 + \sum_{i \in \mathcal{S}} V_i)Q} \geq \lambda, \\ &\Rightarrow \sum_{i \in \mathcal{S}} V_i w_i + \frac{P}{Q} - \sum_{i \in \mathcal{S}} V_i \frac{w_i}{Q} \geq \lambda \left(1 + \sum_{i \in \mathcal{S}} V_i \right), \\ &\Rightarrow \sum_{i \in \mathcal{S}} V_i \left(w_i - \frac{w_i}{Q} - \lambda \right) + \frac{P}{Q} \geq \lambda. \end{aligned}$$

On the analysis above, Y^* in (24) is equal to the following equation, i.e.,

$$Y^* = \max \left\{ \lambda \in \mathbb{R} : \max_{\mathcal{S}: |\mathcal{S}| \leq \bar{a}} \sum_{i \in \mathcal{S}} V_i \left(w_i - \frac{w_i}{Q} - \lambda \right) + \frac{P}{Q} \geq \lambda \right\}. \quad (25)$$

To obtain Y^* , some additional definitions are made. More specifically, for each $\lambda \in \mathbb{R}$, we define functions $A: \mathbb{R} \rightarrow \{\mathcal{S} \in \mathcal{I} : |\mathcal{S}| \leq \bar{a}\}$ and $G: \mathbb{R} \rightarrow \mathbb{R}$ as

$$A(\lambda) = \arg \max_{\mathcal{S}: |\mathcal{S}| \leq \bar{a}} \sum_{i \in \mathcal{S}} V_i \left(w_i - \frac{w_i}{Q} - \lambda \right) + \frac{P}{Q}, \quad (26)$$

and

$$G(\lambda) = \sum_{i \in A(\lambda)} V_i \left(w_i - \frac{w_i}{Q} - \lambda \right) + \frac{P}{Q}, \quad (27)$$

respectively. Thereby, we have

$$Y^* = \max \{ f(A(\lambda)) : \lambda \in \mathbb{R} \}. \quad (28)$$

In other words, to find the optimal assortment set, it suffices to enumerate $A(\lambda)$ for any $\lambda \in \mathbb{R}$. In later sessions, we will show

that the time complexity of finding all the possible assortment sets, $\{A(\lambda) : \lambda \in \mathbb{R}\}$, is $\mathcal{O}(I^2)$.

Before illuminating the details of our algorithm, some geometric insights are provided at first. For any $\lambda \in \mathbb{R}$, let $g_i: \mathbb{R} \rightarrow \mathbb{R}$, $i \in \{0\} \cup \mathcal{I}$, be defined as

$$g_0(\lambda) = 0, \text{ and } g_i(\lambda) = V_i \left(w_i - \frac{w_i}{Q} - \lambda \right) \text{ for } i \in \mathcal{I}, \quad (29)$$

wherein $g_i(\lambda)$ represents a straight line in the 2-dimensional plane. In accordance with (26), for any $\lambda \in \mathbb{R}$, $A(\lambda)$ corresponds to the top- \bar{a} lines among $g_0(\lambda)$ and $g_i(\lambda)$ for $i \in \mathcal{I}$ whose values are non-negative at λ . Moreover, the intersection points among the $I + 1$ lines, i.e., $g_i(\lambda)$ for $i \in \{0\} \cup \mathcal{I}$, are at most $\frac{I^2+I}{2}$. Furthermore, for any λ falling between two consecutive intersection points, the sign of $g_i(\lambda) = V_i(w_i - \frac{w_i}{Q} - \lambda)$ will not change. In other words, for any λ exactly between two adjacent intersection points, the assortment set $A(\lambda)$ remains the same. Therefore, it is enough to enumerate all the intersections points instead of enumerating $A(\lambda)$ for any $\lambda \in \mathbb{R}$. This geometric intuition discloses the basic idea of our designed algorithm. Hereinafter, we will elaborate on the details.

For $0 \leq i < i' \leq I$ and $V_i \neq V_{i'}$, we define $\Omega(i, i')$ as the x -coordinate of the intersection point between the straight lines $g_i(\lambda)$ and $g_{i'}(\lambda)$. Based on (29), $\Omega(i, i')$ can be calculated via solving the following equation, i.e.,

$$V_i \left(w_i - \frac{w_i}{Q} - \lambda \right) = V_{i'} \left(w_{i'} - \frac{w_{i'}}{Q} - \lambda \right),$$

from which $\Omega(i, i')$ is obtained as follows:

$$\Omega(i, i') = \frac{V_{i'} \left(w_{i'} - \frac{w_{i'}}{Q} \right) - V_i \left(w_i - \frac{w_i}{Q} \right)}{V_{i'} - V_i}.$$

Let $\beta = ((i_1, i'_1), (i_2, i'_2), \dots, (i_M, i'_M))$ be the ordering of all M intersection points, wherein $i_l \leq i'_l$ for all l . Thereby, we have

$$\Omega(i_1, i'_1) \leq \Omega(i_2, i'_2) \leq \dots \leq \Omega(i_M, i'_M). \quad (30)$$

Moreover, we define $\pi^0 = (\pi_1^0, \pi_2^0, \dots, \pi_I^0)$ as the ordering of user's preference weights from the largest to the least values, i.e.,

$$V_{\pi_1^0} \geq V_{\pi_2^0} \geq \dots \geq V_{\pi_I^0}.$$

For any $\lambda \in \mathbb{R}$ that exactly falls within interval $(\Omega(i_l, i'_l), \Omega(i_{l+1}, i'_{l+1}))$, the following information will be updated accordingly in the devised assortment decision algorithm.

- 1) The ordering $\pi^l = (\pi_1^l, \pi_2^l, \dots, \pi_I^l)$ of lines $g_i(\cdot)$ for $i \in \mathcal{I}$ from the largest value to the smallest value. Namely, for any $\lambda \in (\Omega(i_l, i'_l), \Omega(i_{l+1}, i'_{l+1}))$, we have

$$g_{\pi_1^l}(\lambda) \geq g_{\pi_2^l}(\lambda) \geq \dots \geq g_{\pi_I^l}(\lambda).$$

- 2) The index set of the top- \bar{a} lines based on π^l , denoted by \mathcal{G}^l , that is,

$$\mathcal{G}^l = \{ \pi_1^l, \pi_2^l, \dots, \pi_{\bar{a}}^l \}.$$

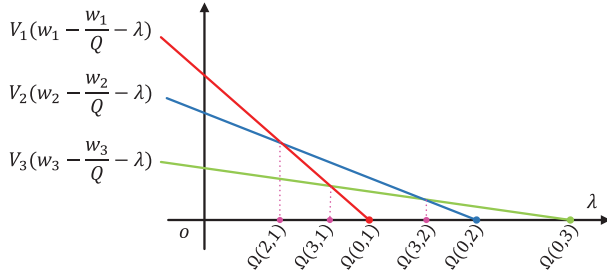


Fig. 2. Example of Algorithm 1 with $\bar{a} = 2$ and $I = 3$.

TABLE II
PARAMETERS UPDATE IN ALGORITHM 1, WHERE $\bar{a} = 2$ AND $I = 3$

l	$\Omega(i_l, i'_l)$	π^l	\mathcal{G}^l	\mathcal{H}^l	\mathcal{A}^l
0	$-\infty$	(1, 2, 3)	(1, 2)	\emptyset	{1, 2}
1	$\Omega(2, 1)$	(2, 1, 3)	{2, 1}	\emptyset	{2, 1}
2	$\Omega(3, 1)$	(2, 3, 1)	{2, 3}	\emptyset	{2, 3}
3	$\Omega(0, 1)$	(2, 3, 1)	{2, 3}	{1}	{2, 3}
4	$\Omega(3, 2)$	(3, 2, 1)	{3, 2}	{1}	{3, 2}
5	$\Omega(0, 2)$	(3, 2, 1)	{3, 2}	{1, 2}	{3}
6	$\Omega(0, 3)$	(3, 2, 1)	{3, 2}	{1, 2, 3}	\emptyset

- 3) The index set of lines whose signs become negative, referred to as \mathcal{H}^l . In other words,

$$\mathcal{H}^l = \{i : g_i(\lambda) < 0 \text{ for } \lambda \in (\lambda(i_l, i'_l), \lambda(i_{l+1}, i'_{l+1}))\}.$$

Note that the slopes of all I lines $g_i(\lambda)$ are negative, we have $\mathcal{H}^l \subseteq \mathcal{H}^{l+1}$ for all l .

- 4) The assortment set $\mathcal{A}^l = \mathcal{G}^l \setminus \mathcal{H}^l$.

With aforementioned discussions and definitions, we summarize the pseudo-code of our devised personalized assortment decision method in Algorithm 1. In each iteration of Algorithm 1, one assortment set will be generated. By comparing the achievable revenue of different sets, we can obtain the optimal assortment strategy. For ease of understanding, an illustration of Algorithm 1 is given in Fig. 2, taking $\bar{a} = 2$ and $I = 3$ as an example. Moreover, the associated parameters $\Omega(i_l, i'_l)$, π^l , \mathcal{G}^l , \mathcal{H}^l , and \mathcal{A}^l within each iteration are summed up in Table II. It is noteworthy that, in the l -th iteration, the assortment set \mathcal{A}^l corresponds to the top-2 lines that are non-negative within $\Omega(i_l, i'_l)$.

Lemma 3. The time complexity of Algorithm 1 is $\mathcal{O}(I^2)$.

Proof. Based on the foregoing analysis, it suffices to enumerate all the intersection points among $I + 1$ straight lines so as to find the optimal assortment set. The proof is completed.

2) *Cache Planning Subproblem:* Given the assortment list per user, i.e., \mathbf{a} , the resultant problem is only related to the cache planning at the BS, which is denoted by $\mathcal{P}(2)$ and stated below:

$$\begin{aligned} & \underset{\mathbf{c}}{\text{maximize}} \quad R && \mathcal{P}(2) \\ & \text{s.t. } C2 : \sum_{i \in \mathcal{I}} c_i B_i \leq C_{\text{BS}}, \\ & C4 : c_i \in \{0, 1\}, i \in \mathcal{I}. \end{aligned}$$

Algorithm 1: Assortment Decision Algorithm.

- 1: **input:** The number of all the intersection points among $I + 1$ lines, i.e., M . The ordering of these M intersection point, β , and the preference weight ordering π^0 . Let $\mathcal{G}^0 = \mathcal{A}^0 = \{\pi_1^0, \pi_2^0, \dots, \pi_{\bar{a}}^0\}$ and $\mathcal{H}^0 = \emptyset$.
 - 2: **repeat**
 - 3: **if** $i_l = 0$ **then**
 - 4: Let $\pi^l = \pi^{l-1}$ and $\mathcal{H}^l = \mathcal{H}^{l-1} \cup \{i'_l\}$
 - 5: Update $\mathcal{G}^l = \{\pi_1^l, \pi_2^l, \dots, \pi_{\bar{a}}^l\}$ and $\mathcal{A}^l = \mathcal{G}^l \setminus \mathcal{H}^l$
 - 6: **else**
 - 7: Get the permutation π^l by transposing i_l and i'_l from π^{l-1} . Let $\mathcal{H}^l = \mathcal{H}^{l-1}$
 - 8: Update $\mathcal{G}^l = \{\pi_1^l, \pi_2^l, \dots, \pi_{\bar{a}}^l\}$ and $\mathcal{A}^l = \mathcal{G}^l \setminus \mathcal{H}^l$
 - 9: **end if**
 - 10: **until** $l = M$
 - 11: **return** $\mathcal{A} = \{\mathcal{A}^1, \mathcal{A}^2, \dots, \mathcal{A}^M\}$
 - 12: Compare the achievable revenue of different sets, namely \mathcal{A}^l where $l \in \{1, 2, \dots, M\}$, and denote the set that achieves the highest revenue as \mathcal{A}^*
 - 13: **output:** the optimal assortment set \mathcal{A}^*
-

Theorem 4. $\mathcal{P}(2)$ can be regarded as a 0-1 Knapsack problem.

Proof. It is worth noting that, with \mathbf{a} , the demanding probability $p_{k,i}$ can be calculated for $k \in \mathcal{K}$ and $i \in \mathcal{I}$. Thereby, the objective function R can be rewritten as follows:

$$\begin{aligned} R &= \sum_{k \in \mathcal{K}} \sum_{i \in \mathcal{I}} p_{k,i} (c_i w'_i + (1 - c_i) \alpha_i w'_i) \\ &= \sum_{k \in \mathcal{K}} \sum_{i \in \mathcal{I}} p_{k,i} (1 - \alpha_i) w'_i c_i + \sum_{k \in \mathcal{K}} \sum_{i \in \mathcal{I}} p_{k,i} \alpha_i w'_i \\ &= \sum_{i \in \mathcal{I}} W_i c_i + Z, \end{aligned} \quad (31)$$

wherein

$$W_i = \sum_{k \in \mathcal{K}} p_{k,i} (1 - \alpha_i) w'_i \quad (32)$$

is a constant, which can be regarded as the value of content item i . Meanwhile,

$$Z = \sum_{k \in \mathcal{K}} \sum_{i \in \mathcal{I}} p_{k,i} \alpha_i w'_i \quad (33)$$

is also a constant. Thereafter, let us go back to the optimization problem $\mathcal{P}(2)$ with (31), (32), and (33), it is not difficult to check that $\mathcal{P}(2)$ is a 0-1 Knapsack problem. This completes the proof.

Given that problem $\mathcal{P}(2)$ is a 0-1 Knapsack problem, its globally optimal solution can be attained by a dynamic programming algorithm [34].

3) *Alternating Optimization Algorithm for the Joint Decision-Making Problem:* Thus far, we have explained how the two subproblems can be optimally solved. In this subsection, an alternating optimization approach is devised to jointly optimize the two types of Boolean variables. We define \mathbf{c}^t and \mathbf{a}^t as the cache placement decision as well as the assortment strategy in the t -th iteration. Moreover, we assume that \mathbf{c}^0 is set by using the

Algorithm 2: Joint Optimization Algorithm.

-
- 1: **input:** The initial cache placement strategy, i.e., c^0 . The maximum iteration number, denoted by N .
 - 2: **repeat**
 - 3: **if** $t < N$ **then**
 - 4: For $k \in \mathcal{K}$, find its optimal assortment policy, i.e., a_k^t , based on Algorithm 1 under given c^{t-1} . Determine a^t based on each user's assortment strategy
 - 5: Update the cache placement policy as c^{t+1} by using dynamic programming algorithm to solve problem $\mathcal{P}(2)$ with a^t
 - 6: **end if**
 - 7: **until** $t = N$ or the system revenue, R , cannot be further improved
 - 8: **return** (a^t, c^t)
 - 9: **output:** the assortment strategy and the cache planning of the system, i.e., (a, c)
-

top-preference scheme. With these definitions, the description of our joint optimization method is briefed in Algorithm 2.

B. Property Analysis

In this subsection, we provide the corresponding property analysis for the designed solution. To be more specific, the convergence performance regarding our joint optimization algorithm is discussed in the following lemma.

Lemma 5. The convergence of Algorithm 2 is ensured.

Proof. As each subproblem is optimally solved by our designed algorithms, the achievable revenue is non-decreasing, which is also upper bounded. The proof is completed.

Moreover, the time complexity of our developed joint optimization algorithm is analyzed in Lemma 6.

Lemma 6. The computational complexity of Algorithm 2 is $\mathcal{O}(I^2W^*)$, where $W^* = \max_{i \in \mathcal{I}} W_i$.

Proof. Based on Lemma 3, the time complexity of the designed assortment decision-making algorithm is $\mathcal{O}(I^2)$. Based on [34], the running time of the dynamic programming algorithm for solving the cache placement optimization problem is $\mathcal{O}(I^2W^*)$. The proof is completed.

Before closing this section, it is important to emphasize the significance of the developed joint optimization algorithm. While it is true that the globally optimal value of problem (18) may not be attainable, the algorithm has been designed to have only pseudo-polynomial time complexity.⁵ This makes it extremely practical for real-world applications, especially for large-scale WCCNs. The proposed algorithm can be used to optimize network performance, manage network resources, and improve overall efficiency. Furthermore, the pseudo-polynomial time complexity ensures that the algorithm can handle larger and more complex networks, making it a valuable tool for network administrators and engineers alike. In summary, the joint

⁵The time complexity of the dynamic programming for 0-1 Knapsack problem is pseudo-polynomial.

optimization algorithm is a practical and efficient solution for WCCNs, and its significance cannot be understated in the field of network optimization. Notably, in this paper, we considered a generic system model of WCCNs, wherein one BS is used to serve multiple users, as that in [35], [36]. If the scenarios with multiple BSs are taken into account, the proposed RMNL is still applicable as our assortment decision is optimally determined with the given cache placement. In regard to the cache planning problem, our developed solution can also be applied directly given that the adjacent BSs can provide the requested contents for the users whose associated BS didn't cache the corresponding items, i.e., letting the cache capacity of our BS be the summation of multiple BSs' cache size.

VI. NUMERICAL SIMULATIONS

In this section, numerical simulations are performed to validate the superiority of our developed joint optimization algorithm compared with various baseline schemes. Unless otherwise stated, the system parameters are set below. We consider a WCCN with 10 users and 100 content items, i.e., $K = 10$ and $I = 100$, respectively. The bit size per content follows a random distribution within the interval of [10,100], namely, $B_i \in [10, 100]$ for $i \in \mathcal{I}$. Following the same setting in [30], we denote $\{0.99, 1.99, 2.99, 3.99, 4.99, 5.99, 6.99, 7.99, 8.99, 9.99\}$ as the inherent marginal revenue set of all content items, i.e., the associated profit per content is randomly selected from this set. In addition, it is stipulated that the discount factor for each content item, i.e., α_i , is generated by a uniform distribution within (0,1). The mean utility of user k to item i , i.e., $\mu_{k,i}$, is randomly generated from 0 to 1. Moreover, we assume the cache capacity budget of BS to be [200, 400, 600, 800, 1000, 1200, 1400, 1600]. With the definitions, it can be seen that if the cache size is less than or equal to 1000, the storage will be fully utilized. When the cache size is larger, the cache storage may not be fully utilized, indicating that there are sufficient edge resources. The former case can be considered a system with limited storage resources. The maximum iteration number of Algorithm 2 is assumed to be 10. Furthermore, it is supposed that all users have the same capacity requirement towards the assortment set, i.e., $\bar{a}_k = \bar{a}$ for $k \in \mathcal{K}$. It is worth mentioning that the designed algorithm is also applicable when \bar{a}_k is heterogeneous. Details will be given in the following session. In this section, we assume \bar{a} is an integer that falls in [3,10].

For performance comparison, the following benchmark schemes are taken into account:

- **MNL_Top:** In this scheme, conventional MNL model is used, i.e., (5), wherein the system shows all content items to the users. It is noteworthy that in the conventional MNL choice model, only the limited number of content has been shown to users due to the screen size limitation, which is intuitively not revenue-maximized. This can also be seen from (5). To ensure the fairness of the comparison, we assume that all content items can be presented to each user. With their demand probability distributions, the BS caches the top-preferred contents based on all users' aggregated request probability to fill its cache storage. The cache

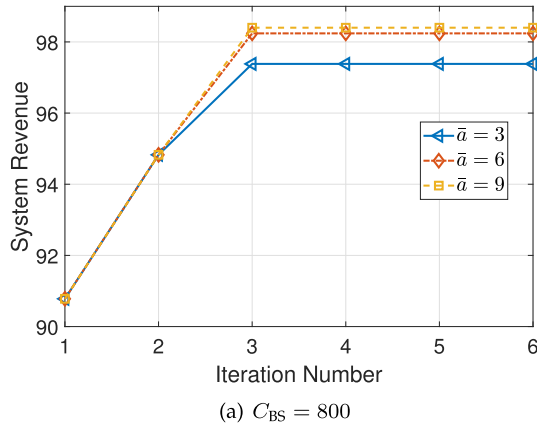
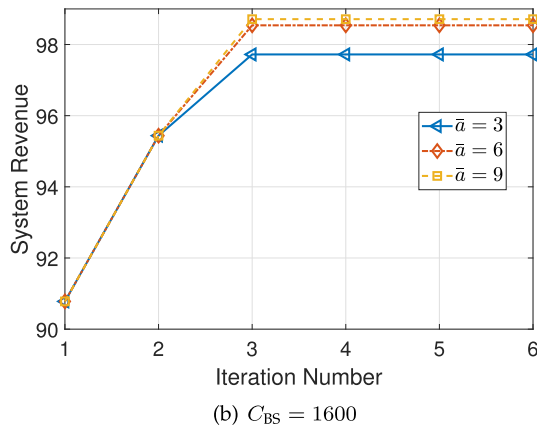
(a) $C_{BS} = 800$ (b) $C_{BS} = 1600$

Fig. 3. Convergence performance analysis.

planning scheme described in this baseline was commonly used by [10], [13], [14], [15], [16].

- **MNL_Opt**: In this scheme, conventional MNL choice model is applied as well. In addition, the BS optimally caches the contents to fully use its storage.
- **MNL_Ran**: This scheme is identical to the foregoing mentioned two methods apart from that the BS randomly caches contents. Here, the cache placement strategy was used in [5], [8], [9].
- **RMNL_Top**: The caching strategy in this baseline is the same as that in MNL_Top, whilst the optimal assortment decision is made based on our proposed RMNL model.
- **RMNL_Ran**: This scheme applies the RMNL model as proposed in this paper. Meanwhile, the BS randomly stores the contents to fill up its cache storage.

Moreover, we use “Proposed” to represent our joint optimization scheme under the designed RMNL model. In the end, we declare that to ensure fairness, all the baselines apply the same alternating optimization algorithm as our developed solution.

A. Convergence Performance

In this subsection, we study the convergence performance of our developed joint decision-making approach, i.e., Algorithm 2. The system revenue per iteration is applied to show this metric,

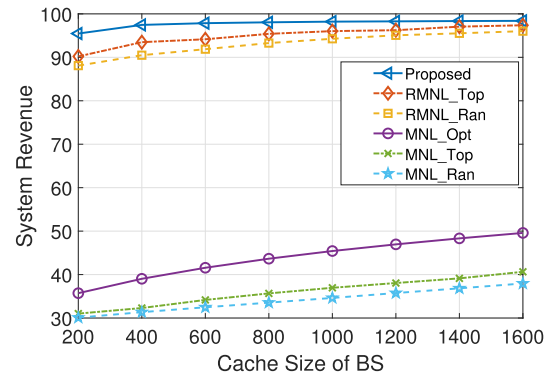
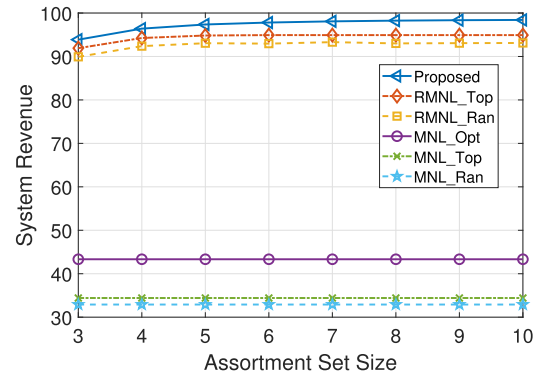
(a) System revenue vs. the cache size of BS, wherein $\bar{a} = 5$ (b) System revenue vs. the assortment size per user, $C_{BS} = 800$

Fig. 4. System revenue performance analysis under homogeneous assortment size among users.

as illustrated by Fig. 3. Thereof, the x -axis depicts the iteration number, whilst the y -axis represents the system revenue during each iteration. For Fig. 3(a) and (b), the cache capacity size at the BS, i.e., C_{BS} , is set to be 800 and 1600, respectively. It can be seen from Fig. 3 that the designed iterative scheme can converge within several steps. In addition, from Fig. 3(a), it can be observed that with the assortment size, the achievable revenue increases. This is due to the fact that increasing the assortment size, more favorable and profitable contents can be presented to the users, which motivates more content consumptions from users. Similar trends can be obtained from Fig. 3(b), which are not repeated to avoid redundancy. Moreover, comparing Fig. 3(a) and (b), one can see that under the same assortment size, a larger number of cache capacity size induces a higher system revenue, demonstrating the effectiveness of wireless content caching.

B. System Revenue

Fig. 4 illustrates the achieved system revenue of our developed joint optimization algorithm and extensive baselines under different system parameters. More specifically, Fig. 4(a) shows the system revenue versus the cache capacity size of the BS. Thereof, the assortment size per user is assumed to be 5, namely, $\bar{a} = 5$. From this, one can see that the achievable revenue of each scheme increases with the cache size at the BS, demonstrating the effectiveness of content caching in term of boosting system

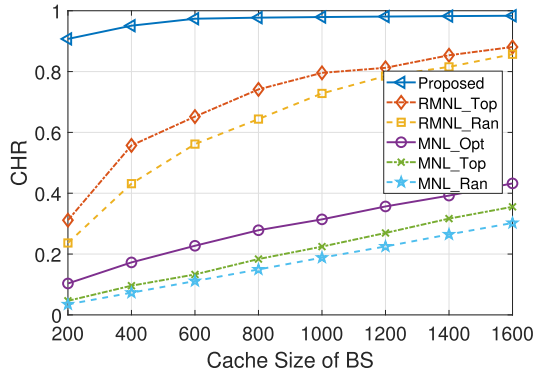
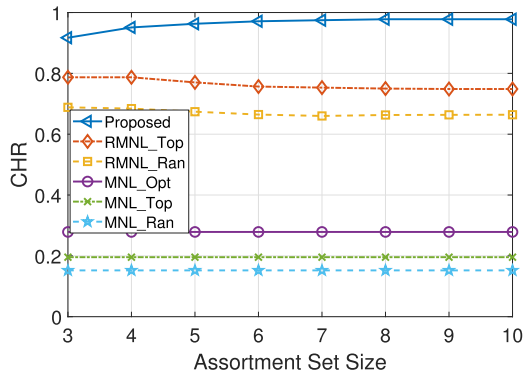
(a) CHR vs. the cache size of BS, wherein $\bar{a} = 5$ (b) CHR vs. the assortment size per user, $C_{BS} = 800$

Fig. 5. CHR performance analysis under homogeneous assortment size among users.

revenue. In addition, the proposed RMNL-oriented schemes outperform the conventional MNL-enabled solutions, showing the superiority of our devised choice model. Besides, the developed joint decision making policy achieves the highest system revenue among all the strategies. For instance, when $C_{BS} = 600$, our proposed scheme achieves 4%, 6%, 140%, 193%, and 207% more revenue than that of RMNL_Top, RMNL_Ran, MNL_Opt, MNL_Top, and MNL_Ran, respectively. Moreover, among the three MNL-enabled strategies, MNL_Opt obtains the highest revenue as the optimal cache planning has been utilized therein. Furthermore, for both MNL and RMNL-assisted schemes, Top-caching manner has superior performance in terms of system revenue compared to that of random caching scheme.

In the meantime, Fig. 4(b) depicts the system revenue versus the assortment size, in which the cache capacity at the BS is assumed to be $C_{BS} = 800$. Based on Fig. 4(b), it is not surprising to see that the achievable revenue of the MNL-oriented schemes keeps unchanged with the assortment size. This is because in MNL-enabled strategies, all content items are shown to the users at once. For the developed RMNL-assisted schemes, the system revenue increases with the assortment size and our devised scheme has the highest revenue among all strategies, demonstrating the effectiveness of our developed joint decision-making algorithm. Moreover, as expected, for both MNL and RMNL-oriented policies, optimal caching achieves the best performance and meanwhile random caching has the least total revenue.

C. Cache Hit Ratio

In content caching-oriented wireless networks, cache hit ratio (CHR) has been taken as a significant metric to evaluate the performance of caching decision strategies [37], [38], such as the dissatisfaction or satisfaction of users. Generally, CHR is characterized by the probability of the requested content that can be supplied by network edges. The detailed definition is given as follows:

$$\text{CHR} = \frac{1}{K} \sum_{k \in \mathcal{K}} \sum_{i \in \mathcal{I}} p_{k,i} c_i.$$

In this subsection, we investigate the CHR performance of each scheme with respect to the cache size at the BS and the assortment size per user, as illustrated by Fig. 5(a) and (b), respectively.

To be more specific, in Fig. 5(a), the assortment size is set to be 5. As expected, CHR per scheme increases with the cache capacity size. This is because a large size of cache entity can store more content items. It can also be seen that the developed joint optimization scheme achieves the highest CHR among all strategies. The superiority of our scheme in terms of CHR is more remarkable than that of system revenue. For example, when the cache size at the BS is 600, our developed strategy achieves 32%, 61%, 350%, 603%, and 878% higher CHR than that of RMNL_Top, RMNL_Ran, MNL_Opt, MNL_Top, and MNL_ran, respectively. In addition, all the strategies that apply our proposed RMNL model have better performance than that under MNL model, which shows the effectiveness of our devised choice model. Moreover, under both RMNL and MNL choice models, the optimal caching and the random caching schemes achieve the best and the worst performance with respect to CHR, respectively.

Meanwhile, in Fig. 5(b), we illustrate the CHR performance for each strategy with respect to the assortment size. Similar to Fig. 4(b), CHR of each MNL-oriented scheme keeps constant with \bar{a} due to the discipline set for the MNL choice model. Moreover, the CHR of RMNL_Top and RMNL_Ran decreases as \bar{a} increases. This is because in RMNL_Top and RMNL_Ran, the assortment decision is not optimized with caching decision. Due to the limited cache capacity size at the BS, the probability of the requested content that is presented to users but not cached at the BS increases under a large number of \bar{a} , resulting in a degraded CHR. Whilst, the CHR of our explored scheme increases with the assortment size as we jointly optimize the cache placement and the assortment decision. This further shows the importance of the joint optimization.

D. System Performance Comparison Under Heterogeneous Assortment Sizes

So far, we have validated the effectiveness of our designed joint optimization algorithm under the homogeneous assortment size of users. In this subsection, we compare the system performance of the proposed algorithm against the baselines under heterogeneous assortment sizes of users, as shown in Fig. 6. Specifically, we investigate the achievable revenue and the CHR performance in Fig. 6(a) and (b), respectively. Therein,

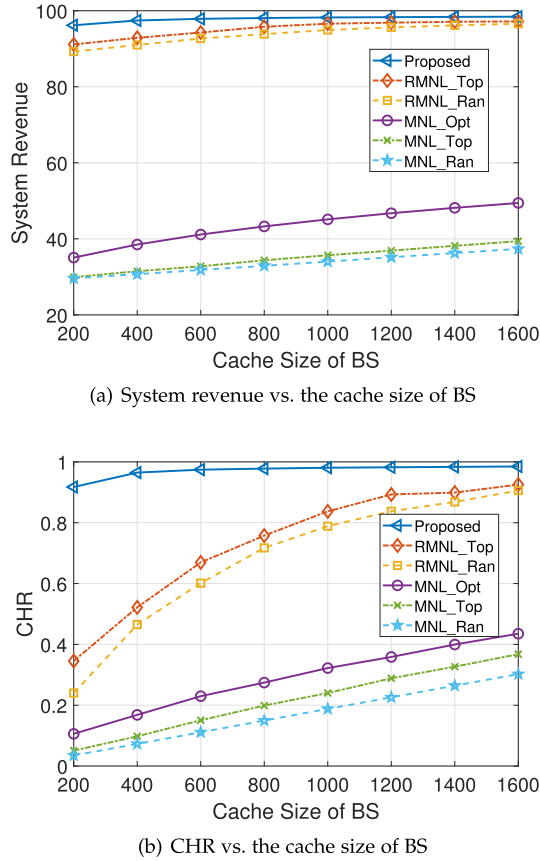


Fig. 6. Performance comparison under heterogeneous assortment sizes among users.

the assortment size of user k is randomly selected from the interval of $[3,10]$, where $k \in \mathcal{K}$. It can be seen from Fig. 6(a) that the proposed scheme achieves the highest system revenue among all schemes. In addition, all the strategies under the RMNL choice model outperforms the schemes that apply MNL model, which further demonstrates the availability of the newly designed model. Again, we can observe that the system revenue per scheme increases and eventually tends gradually towards 100 with the cache capacity at the BS. From Fig. 6(b), we can also observe the superiority of our developed scheme. Other similar trends to Fig. 5(a) are skipped here for simplicity. In summary, through Fig. 6, we demonstrate the effectiveness of our developed joint decision-making method under heterogeneous assortment size of users. These observations together with the aforementioned results show the universal applicability of the proposed algorithm.

VII. CONCLUSION

In this paper, we thoroughly investigated the choice behavior of users in WCCNs, which was characterized by a newly designed RMNL model. Based on it, a revenue maximization problem was formulated via jointly optimizing the personalized assortment decision and the cache planning, which was a non-convex integer programming problem. To render the analysis, an efficient divide-then-conquer manner was proposed.

Simulation results validated the superiority of our developed strategy compared to extensive benchmarks. In future, we plan to study the joint assortment and caching decision-making problem for WCCNs with advanced transmission technologies, such as rate-splitting multiple access (RSMA) and intelligent reflecting surface (IRS). It is noteworthy that cache planning and assortment decisions are long-term optimization variables, while radio resource management is highly dependent on time-varying channels [13], [14]. Therefore, to jointly optimize the performance of RSMA or IRS-driven WCCNs, it is imperative to solve a multi-timescale optimization problem.

ACKNOWLEDGMENTS

This paper was presented in part at the IEEE International Conference on Communications (ICC), Rome, May 2023 [1].

REFERENCES

- [1] Y. Fu, X. Xu, H. Liu, and T. Q. S. Quek, "A revised multinomial logit (RevMNL) choice model for wireless content caching networks," in *Proc. IEEE Int. Conf. Commun.*, 2023, pp. 1–7.
- [2] M. Giordani, M. Polese, M. Mezzavilla, S. Rangan, and M. Zorzi, "Toward 6G networks: Use cases and technologies," *IEEE Commun. Mag.*, vol. 58, no. 3, pp. 55–61, Mar. 2020.
- [3] W. Sun, H. Zhang, R. Wang, and Y. Zhang, "Reducing offloading latency for digital twin edge networks in 6G," *IEEE Trans. Veh. Technol.*, vol. 69, no. 10, pp. 12240–12251, Oct. 2020.
- [4] T. X. Vu, S. Chatzinotas, E. Bastug, and T. Q. S. Quek, *Wireless Edge Caching: Modeling, Analysis, and Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2020.
- [5] W. Wen, C. Liu, Y. Fu, T. Q. S. Quek, F.-C. Zheng, and S. Jin, "Enhancing physical layer security of random caching in large-scale multi-antenna heterogeneous wireless networks," *IEEE Trans. Inf. Forensics Security*, vol. 15, pp. 2840–2855, 2020.
- [6] Y. Jiang, M. Ma, M. Bennis, F.-C. Zheng, and X. You, "User preference learning-based edge caching for fog radio access network," *IEEE Trans. Commun.*, vol. 67, no. 2, pp. 1268–1283, Feb. 2019.
- [7] Z. Zhao, M. Xu, Y. Li, and M. Peng, "A non-orthogonal multiple access-based multicast scheme in wireless content caching networks," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 12, pp. 2723–2735, Dec. 2017.
- [8] W. Liu, H. Zhang, H. Ding, and D. Yuan, "Delay and energy minimization for adaptive video streaming: A joint edge caching, computing and power allocation approach," *IEEE Trans. Veh. Technol.*, vol. 71, no. 9, pp. 9602–9612, Sep. 2022.
- [9] T. Zhang, X. Fang, Y. Liu, G. Y. Li, and W. Xu, "D2D-enabled mobile user edge caching: A multi-winner auction approach," *IEEE Trans. Veh. Technol.*, vol. 68, no. 12, pp. 12314–12328, Dec. 2019.
- [10] Y. Fu, L. Salaun, X. Yang, W. Wen, and T. Q. S. Quek, "Caching efficiency maximization for device-to-device communication networks: A recommend to cache approach," *IEEE Trans. Wireless Commun.*, vol. 20, no. 10, pp. 6580–6594, Oct. 2021.
- [11] A. Ndikumana, N. H. Tran, T. M. Ho, D. Niyato, Z. Han, and C. S. Hong, "Joint incentive mechanism for paid content caching and price based cache replacement policy in named data networking," *IEEE Access*, vol. 6, pp. 33702–33717, 2018.
- [12] J. Kwak, G. Paschos, and G. Iosifidis, "Elastic FemtoCaching: Scale, cache, and route," *IEEE Trans. Wireless Commun.*, vol. 20, no. 7, pp. 4174–4189, Jul. 2021.
- [13] X. Yang, Y. Fu, W. Wen, T. Q. S. Quek, and F. Song, "Mixed-timescale caching and beamforming in content recommendation aware Fog-RAN: A latency perspective," *IEEE Trans. Wireless Commun.*, vol. 69, no. 4, pp. 2427–2440, Apr. 2021.
- [14] X. Yang, Z. Fei, B. Li, J. Zheng, and J. Guo, "Joint user association and edge caching in multi-antenna small-cell networks," *IEEE Trans. Commun.*, vol. 70, no. 6, pp. 3774–3787, Jun. 2022.
- [15] Y. Fu, Q. Yu, A. K. Y. Wong, Z. Shi, H. Wang, and T. Q. S. Quek, "Exploiting coding and recommendation to improve cache efficiency of reliability-aware wireless edge caching networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 11, pp. 7243–7256, Nov. 2021.

- [16] Y. Fu, Q. Yu, T. Q. S. Quek, and W. Wen, "Revenue maximization for content-oriented wireless caching networks (CWCNs) with repair and recommendation considerations," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 284–298, Jan. 2021.
- [17] G. Zheng, C. Xu, M. Wen, and X. Zhao, "Service caching based aerial cooperative computing and resource allocation in multi-UAV enabled MEC systems," *IEEE Trans. Veh. Technol.*, vol. 71, no. 10, pp. 10934–10947, Oct. 2022.
- [18] K. Zhang, J. Cao, S. Maharjan, and Y. Zhang, "Digital twin empowered content caching in social-aware vehicular edge networks," *IEEE Trans. Computat. Social Syst.*, vol. 9, no. 1, pp. 239–251, Feb. 2022.
- [19] M. Chen, U. Challita, W. Saad, and M. D. C. Yin, "Artificial neural networks-based machine learning for wireless networks: A tutorial," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 4, pp. 3039–3071, Fourth Quarter 2019.
- [20] Z. Yang, Y. Liu, Y. Chen, and J. T. Zhou, "Deep learning for latent events forecasting in content caching networks," *IEEE Trans. Wireless Commun.*, vol. 21, no. 1, pp. 413–428, Jan. 2022.
- [21] Z. Yang, Y. Fu, Y. Liu, Y. Chen, and J. Zhang, "A new look at AI-driven NOMA-F-RANs: Features extraction, cooperative caching, and cache-aided computing," *IEEE Wireless Commun.*, vol. 29, no. 3, pp. 123–130, Jun. 2022.
- [22] W. Jiang, G. Feng, S. Qin, T. S. P. Yum, and G. Cao, "Multi-agent reinforcement learning for efficient content caching in mobile D2D networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 3, pp. 1610–1622, Mar. 2019.
- [23] Y. Fu, Z. Yang, T. Q. S. Quek, and H. H. Yang, "Towards cost minimization for wireless caching networks with recommendation and uncharted users' feature information," *IEEE Trans. Wireless Commun.*, vol. 20, no. 10, pp. 6758–6771, Oct. 2021.
- [24] G. K. Zipf, *Selected Studies of the Principle of Relative Frequency in Language*. Cambridge, MA, USA: Harvard Univ. Press, 1932.
- [25] L. A. Adamic and B. A. Huberman, "Zipf's law and the internet," *Glottometrics*, vol. 3, pp. 143–150, 2002.
- [26] P. Rusmevichientong, Z. Shen, and B. Shmoys, "Dynamic assortment optimization with a multinomial logit choice model and capacity constraint," *Operation Res.*, vol. 58, no. 6, pp. 1666–1680, Dec. 2010.
- [27] Y. Lin, Y. Wang, D. He, and L. Lee, "Last-mile delivery: Optimal locker location under multinomial logit choice model," *Transp. Res. Part E: Logistics Transp. Rev.*, vol. 142, Oct. 2020, Art. no. 102059.
- [28] J. Feldman and P. Jiang, "Display optimization under the multinomial logit choice model: Balancing revenue and customer satisfaction," *Prod. Operations Manage.*, to be published, doi: [10.1111/poms.14040](https://doi.org/10.1111/poms.14040).
- [29] P. Gao et al., "Assortment optimization and pricing under the multinomial logit model with impatient customers: Sequential recommendation and selection," *Operations Res.*, vol. 69, no. 5, pp. 1509–1532, 2021.
- [30] A. Azaria, A. Hassidim, A. Eshkol, O. Weintraub, and I. Netanel, "Movie recommender system for profit maximization," in *Proc. 7th ACM Conf. Recommender Syst.*, Hong Kong, China, 2013, pp. 121–128.
- [31] T. Bahreini and D. Grosu, "Efficient algorithms for multi-component application placement in mobile edge computing," *IEEE Trans. Cloud Comput.*, vol. 10, no. 4, pp. 2550–2563, Fourth Quarter 2022.
- [32] G. Vulcano, G. van Ryzin, and W. Chaar, "Choice-based revenue management: An empirical study of estimation and optimization," *Manuf. Service Operations Manage.*, vol. 12, no. 3, pp. 371–392, Summer 2010.
- [33] K. T. Talluri and G. J. Ryzin, *The Theory and Practice of Revenue Management*. Boston, MA, USA: Kluwer Publishers, 2005.
- [34] S. Martello and P. Toth, *Knapsack Problems: Algorithms and Computer Implementations*, 1st ed. Hoboken, NJ, USA: Wiley, 1990.
- [35] G. Ahani and D. Yuan, "Optimal content caching and recommendation with age of information," *IEEE Trans. Mobile Comput.*, early access, Oct. 12, 2022. [Online]. Available: <https://ieeexplore.ieee.org/document/9917351>
- [36] M.-C. Lee and Y.-W. P. Hong, "Socially-aware joint recommendation and caching policy design in wireless D2D networks," in *Proc. IEEE Int. Conf. Commun.*, 2021, pp. 1–6.
- [37] D. Wei and S. Han, "An experimental study of recommendation for wireless edge caching," in *Proc. IEEE/CIC Int. Conf. Commun. China*, 2022, pp. 731–736.
- [38] S. Kastanakis, P. Sermpezis, V. Kotronis, D. Menasche, and T. Spyropoulos, "Network-aware recommendations in the wild: Methodology, realistic evaluations, experiments," *IEEE Trans. Mobile Comput.*, vol. 21, no. 7, pp. 2466–2479, Jul. 2022.



Yaru Fu (Member, IEEE) received the PhD degree in electronic engineering from the City University of Hong Kong (CityU) in 2018. She is currently an assistant professor with the School of Science and Technology, Hong Kong Metropolitan University (HKMU). She is presently serving as an associate editor for *IEEE Internet of Things Journal*, *IEEE Wireless Communications Letters*, *IEEE Networking Letters*, and *Springer Nature Computer Science*. She also serves as a review editor for the *Frontiers in Communications & Networks* and a guest editor for the *Space: Science & Technology*. She was honored with the 2021 Katie Shu Sui Pui Charitable Trust - Outstanding Research Publication Award (Gold Prize), 2022 Best Editor Award for *IEEE Wireless Communications Letters*, 2022 Katie Shu Sui Pui Charitable Trust - Excellent Research Publication Award, and 2022 Exemplary Reviewer for the *IEEE Transactions on Communications* (fewer than 5%). She was listed on the World's Top 2% Scientists 2021 ranking by Stanford University in the United States. Her research interests include intelligent wireless communications and networking, distributed storage system, and digital twin.



Xinyu Xu (Member, IEEE) is currently working toward the PhD degree with the College of Business, Southern University of Science and Technology (SUSTech), Shenzhen, China. His research interests include game theory and service operations.



Hanlin Liu (Member, IEEE) received the BS degree in mathematics from the University of Science and Technology of China (USTC) in 2014, and the PhD degree in system engineering and engineering management from the City University of Hong Kong (CityU) in 2018. Then, he worked as a postdoctoral researcher with the University of Minnesota (UMN), Twin Cities, from 2018 to 2020. He is currently an assistant professor with the School of Business, Southern University of Science and Technology (SUSTech), China. His research interests include cooperative games, queueing economics, and reliability modeling.



Quan Yu (Member, IEEE) received the BS degree in electronic information engineering from Huazhong Normal University, Wuhan, China, in 2009, and PhD degree in information engineering from the City University of Hong Kong, in 2014. From 2014 to 2016, she was a postdoctoral fellow with the Department of Electronic Engineering, City University of Hong Kong. She joined the faculty with the Wuhan University of Technology in 2017, and is now an associate professor with the School of Information Engineering. Her research interests include cloud storage, distributed storage systems, edge computing, and network coding.



and *IEEE Transactions on Industrial Cyber-Physical Systems*. He is a senior member of ACM.

Hong-Ning Dai received the PhD degree in computer science and engineering from the Department of Computer Science and Engineering, the Chinese University of Hong Kong. He is currently with the Department of Computer Science, Hong Kong Baptist University as an associate professor. His current research interests include Internet of Things and blockchain technology. He has served as associate editors *IEEE Communications Surveys & Tutorials*, *IEEE Transactions on Intelligent Transportation Systems*, *IEEE Transactions on Industrial Informatics*,



deputy director of the SUTD-ZJU IDEA. His current research topics include wireless communications and networking, network intelligence, non-terrestrial networks, open radio access network, and 6G. He has been actively involved in organizing and chairing sessions, and has served as a member of the Technical Program Committee as well as symposium chairs in a number of international conferences. He is currently serving as an area editor for the *IEEE Transactions on Wireless Communications*. He was honored with the 2008 Philip Yeo Prize for Outstanding Achievement in Research, the 2012 IEEE William R. Bennett Prize, the 2015 SUTD Outstanding Education Awards – Excellence in Research, the 2016 IEEE Signal Processing Society Young Author Best Paper Award, the 2017 CTTC Early Achievement Award, the 2017 IEEE ComSoc AP Outstanding Paper Award, the 2020 IEEE Communications Society Young Author Best Paper Award, the 2020 IEEE Stephen O. Rice Prize, the 2020 Nokia Visiting Professor, and the 2022 IEEE Signal Processing Society Best Paper Award. He is a fellow of the Academy of Engineering Singapore.

Tony Q. S. Quek (Fellow, IEEE) received the BE and ME degrees in electrical and electronics engineering from the Tokyo Institute of Technology in 1998 and 2000, respectively, and the PhD degree in electrical engineering and computer science from the Massachusetts Institute of Technology in 2008. Currently, he is the Cheng Tsang Man chair professor with Singapore University of Technology and Design (SUTD) and ST Engineering distinguished professor. He also serves as the director of the Future Communications R&D Programme, the head of ISTD Pillar, and the