# QoS-Driven Resource Optimization for Intelligent Fog Radio Access Network: A Dynamic Power Allocation Perspective

Jun Yu, Rui Wang, *Senior Member, IEEE*, Jun Wu, *Senior Member, IEEE*

*Abstract*—The fog radio access network (Fog-RAN) has been considered a promising wireless access architecture to help shorten the communication delay and relieve the large data delivery burden over the backhaul links. However, limited by conventional inflexible communication design, Fog-RAN cannot be used in some complex communication scenarios. In this study, we focus on investigating a more intelligent Fog-RAN to assist the communication in a high-speed railway environment. Due to the train's continuously moving, the communication should be designed intelligently to adapt to channel variation. Specifically, we dynamically optimize the power allocation in the remote radio heads (RRHs) to minimize the total network power cost considering multiple quality-of-service (QoS) requirements and channel variation. The impact of caching on the power allocation is considered. The dynamic power optimization is analyzed to obtain a closed-form solution in certain cases. The inherent tradeoff among the total network cost, delay and delivery content size is further discussed. To evaluate the performance of the proposed dynamic power allocation, we present an invariant power allocation counterpart as a performance comparison benchmark. The result of our simulation reveals that dynamic power allocation can significantly outperform the invariant power allocation scheme, especially with a random caching strategy or limited caching resources at the RRHs.

*Index Terms*—Fog-RAN, caching, dynamic power allocation, quality-of-service.

## I. Introduction

As an evolution of the conventional cloud radio access network (C-RAN) [1], [2], the fog radio access network (Fog-RAN) has been regarded as a promising new wireless access network architecture, which can help relieve the large data traffic burden in the backhaul links of a cellular network and satisfy the stricter delay requirements in beyond-fifth-generation (B5G) and sixth-generation (6G) wireless communications [3], [4]. The key technology of the Fog-RAN to achieve this performance enhancement is the introduction of caching resources on the edge devices, i.e., the remote radio heads (RRHs), of the access network. With some popular contents cached at the RRHs, the frequently requested data by the terminals can be instantly acquired from the RRHs without needing to be fetched from remote servers. This scenario avoids possible traffic congestion over the backhaul links and further shortens the content delivery delay.

J. Yu, and R. Wang are with the School of Electronics and Information Engineering at Tongji University, Shanghai, 201804, P. R. China (e-mails: 1610986@tongji.edu.cn, ruiwang@tongji.edu.cn).

J. Wu is with the School of Computer Science and Technology at Fudan University, Shanghai, 201203, P. R. China (e-mail: wujun@fudan.edu.cn).

Due to the advantage of the Fog-RAN, the corresponding studies have recently received much attention. Current studies on the Fog-RAN mainly focus on two aspects: performance characterization and resource optimization. Here, the resources include caching, beamforming, power and computation resources. The authors in [5] investigated the delay and energy efficiency of the Fog-RAN by considering the hybrid caching strategy, which combines coded cached, nonpartitioned cached and uncached files. The hybrid caching strategy was further optimized to balance the delay and energy efficiency. Successful transmission probability (STP) is often used a performance metric to characterize the impact of content caching at an RRH. In [6], the authors first derived the STP of the Fog-RAN system with a proactive probabilistic caching strategy. The caching probability for different contents was further optimized to maximize the STP. In [7], Fog-RAN-assisted transmission was studied in a heterogeneous Fog-RAN wireless network. Different from [6], the authors in [7] optimized the caching strategy in both the RRHs and mobile users, aiming to optimize the STP. The power allocation of the Fog-RAN was studied in [8]–[11]. In [8], the authors proposed using the nonorthogonal multiple access (NOMA) technique to distinguish multiple users in the Fog-RAN. An improved fractional transmit power allocation algorithm was developed for the three-user NOMA-assisted Fog-RAN system. In [9], the latency minimization problem for the Fog-RAN was formulated with a dynamic user demand. To understand the demand of the user and to intelligently perform the joint optimization of the proactive cache strategy and power allocation, a deep reinforcement learning approach was used. The authors in [10] jointly optimized the power and the model selection for uplink transmission of the Fog-RAN. The reinforcement learning approach was used to solved the nonconvex mixed-integer programming problem. Subchannel assignment and power control were jointly optimized in [11] for a mmWave-based Fog-RAN. The alternative direction method of multipliers (ADMM) was utilized to solve the power control problem.

The beamforming design problem for the Fog-RAN was investigated in [12]–[14]. In particular, the uplink Fog-RAN was studied [12], aiming to maximize the delivery rate subject to the constraints, including the backhaul capacity, transmit power, and file size. A two-layer transmission scheme including the cache level and network level was applied for content transmission. Both the centralized algorithm and the decentralized algorithm were proposed for optimizing

the beamforming problem. In [13], the authors considered the Fog-RAN with multicast transmission. The beamforming design problem was constructed to minimize the network total cost subject to the power and signal-to-interference-plus-noise ratio (SINR) constraints. The sparse beamforming vector was found with the convex-concave procedure. The authors in [14] further extended the channel model in [13] to a scenario with both multicast and unicast transmission. A branch-and-bound algorithm was utilized to find the global optimal solution.

Computational resource allocation was considered in [15]–[18] for performance improvement. In [15], edge/cloud computing and edge computing task migration were included in the design problem to improve the quality-of-service (QoS) at the users. Optimizing the user association and computational offloading, the communication resources and computational resources could be balanced. A reinforcement-learning-based approach was used in [16] to optimize the computational resources with the objective of reducing the latency and energy consumption. In [17], the authors proposed using the computational resources at the edge devices of the Fog-RAN to create cooperative downlink transmission to decrease the latency. An order-optimal upload-download communication latency pair was characterized to evaluate the network performance. In [18], the authors formulated a joint decision problem for communication, caching and computing resources. The problem was modeled as a multiple-choice multidimensional knapsack problem, and was then solved by the Lagrangian dual decomposition approach.

As we summarized above, even though a large number of contributions to the Fog-RAN have been reported, the investigations are still not sufficient to address all challenges under different application scenarios. A key shortcoming of the existing work is that current studied Fog-RAN are not intelligent enough, leading to that it cannot be used in some complex dynamic communication scenarios. A typical scenario is the high-speed railway scenario [19]. The technology of the high-speed railway, identified as a typical wireless communication scenario for future cellular networks, has developed rapidly on a global scale [20]. How to use the Fog-RAN architecture to provide reliable and low-latency communication services is an interesting topic that will require in-depth studies.

To compensate for the deficiencies of the current research on the Fog-RAN, we aim to design a more intelligent Fog-RAN to assist the communication in a high-speed railway environment. In specific, we investigate the Fog-RAN-assisted downlink data transmission in the high-speed railway scenario. Because of the time-varying channel state information, the corresponding studies become more challenging. Although wireless communication in the high-speed railway scenario has received much attention [21]–[25], to our knowledge, few works have studied the Fog-RAN in the high-speed railway scenario. This lack motivates the study of this work.

In the considered high-speed railway Fog-RAN, we assume that the train is cooperatively served by multiple RRHs. By taking the time-varying channel into account, we investigate an intelligently dynamic power allocation to minimize the cost of the entire network power, which includes the transmit power at the RRHs and the power consumed over the backhaul
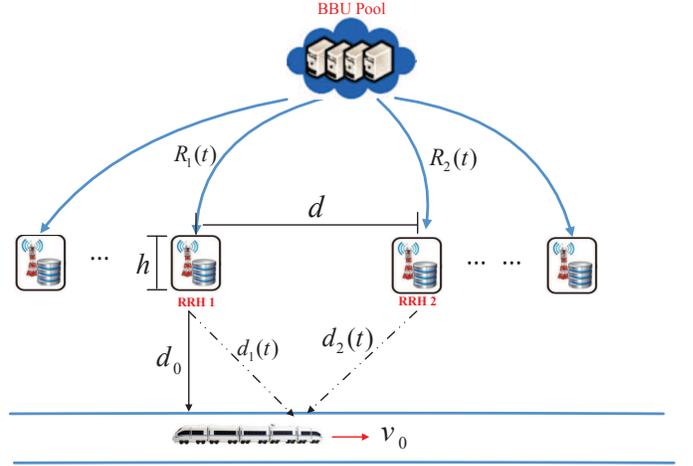


Fig. 1. Fog-RAN-assisted high-speed communication system.

links, subject to a few quality-of-service (QoS) requirements. With this goal, the caching placement on the RRHs affects the physical layer power allocation. To find a reasonable solution to the considered nonconvex problem, smoothed $l_0$-norm approximation is employed to convert the nonconvex problem into a tractable form. By analyzing the unique channel properties of the considered high-speed railway Fog-RAN, we provide a closed-form solution in certain special cases. With the obtained solution, the inherent tradeoff among the total network cost, delay and delivery content size is further discussed. Moreover, the invariant power allocation counterpart is derived to evaluate the performance of the proposed dynamic power allocation. Simulation results reveal that the intended dynamic power allocation is able to significantly outperform the invariant power allocation scheme. The performance gain is more pronounced when the random caching strategy is used or the caching resources at the RRHs are limited.

The organization of the remainder of the paper is shown below. In Section II, we present the channel model as well as problem formulation of the considered high-speed railway Fog-RAN. In Section III, optimization of the dynamic power allocation is discussed. In Section IV, we analyze the tradeoff among the total network cost, delay and delivery content size. In Section V, we present the invariant power allocation counterpart. In Section VI, extensive numerical results are stated to demonstrate the performance. Finally, we conclude the paper in Section VII.

## II. CHANNEL MODEL AND PROBLEM FORMULATION

In this section, we present the channel model of the considered Fog-RAN. After that, the dynamic optimization problem is formulated.

### A. Channel model

As illustrated in Fig. 1, we consider a dynamic wireless transmission scenario of a train running at high speed along a straight-line railway. To provide the required high data rate and low latency downlink transmission, the train is served by a Fog-RAN in which multiple base stations, i.e., RRHs, are

uniformly deployed along one side of the road at intervals of the same size. The RRHs are assumed to be connected to the same baseband unit (BBU). The BBU is responsible for the RRH resource allocation as well as the serving-RRH switch. At any time, the train is served by the two nearest RRHs. Because the service repeats periodically for different pairs of neighboring RRHs, with no loss of generality, only one time interval is considered, where the train is assisted by RRH 1 and RRH 2.

It is assumed that the contents requested by the train include $L$ different popular contents of the same size $Q$. All RRHs have a limited local storage size, and all $L$ contents cannot be cached in a single RRH simultaneously. Suppose $F_n$ as the local storage size of RRH $n$; thus, we have $F_n < LQ$. According to the requested frequency of the contents, the popularity distribution of the contents is denoted by $\mathbf{p} = [p_1, p_2, \cdots, p_L]$, where $p_l \in (0, 1)$ denotes the popularity of content $l$. With unchanged generality, in this study, we make the assumption that $p_1 \geq p_2 \geq \cdots \geq p_L$ and that the content popularity follows a Zipf distribution given by $p_l = \frac{l^{-\eta}}{\sum_{l=1}^{L} l^{-\eta}}$, where $\eta$ denotes the shaping parameter defining the skewness of the popularity distribution [7], [13], [26]. According to the content popularity distribution, the RRHs can store the requested contents with different caching strategies. Two caching strategies that widely used are the popularity-aware caching (PopC) strategy and random caching (RndC) strategy. The PopC strategy allows each RRH to cache the most popular contents until its storage size is fully utilized, while the RndC strategy asks each RRH to cache the contents randomly with identical probabilities regardless of the content popularity distribution. Therefore, to identify the used caching strategy, we define a cache placement matrix $C \in \mathbb{B}^{N \times L}$ with $c_{n,l} \in 0, 1$. Specifically, $c_{n,l} = 1$ occurs if content $l$ is cached in RRH $n$; otherwise, content $l$ is not cached in RRH $n$. To satisfy the RRH storage size constraints, we have $\sum_{l=1}^{F} c_{n,l}Q \leq F_n, \forall n$.

We consider the Fog-RAN-assisted downlink transmission described in Fig. 1. We assume $d$ equal intervals between every two RRHs that $d_0$ is the distance between each RRH and the road, and that $h$ is the height of the transmit antenna on each RRH. The train moves down the railway at a unchanging velocity $v$. To determine the coordinates of the RRHs, we assume that there is an original point $o$ and the system time when the train passes through $o$ equals to 0. The coordinates of RRH $n$ are represented as $(l_n, d_0)$. During the time interval $(0, T]$, RRH 1 and 2 serve the train. According to the geometric structure of the system, the distance between RRH $n$ and the train can be denoted by

$$d_n(t) = \sqrt{(vt - l_n)^2 + d_0^2 + h^2}, \quad t \in (0, T]. \quad (1)$$

When $t > T$, the BBU coordinates the handoff process, and the train is served by a new set of RRHs. However, because the communications over different time intervals are periodic, with no loss of generality, we next focus only on the dynamic resource allocation over time interval $t \in (0, T]$.

During the time interval $t \in (0, T]$, we assumed that a requested content has to be delivered from the RRHs to the train. Denote $x(t)$ as the modulated signal for the requested content transmitted in the downlink. Here, we assume that signals transmitted from different RRHs are sent over an orthogonal bandwidth. Moreover, assume that signal $x(t)$ is a stochastic process with mean of zero and with unit variance. Thus, the baseband signal transmitted from RRH $n$ at time $t$ can be represented by

$$y_n(t) = \sqrt{P_n(t)} h_n(t) x(t) + n_n(t), \quad (2)$$

where $P_n(t)$ denotes the instantaneous transmit power at RRH $n$, $h_n(t)$ shows the instantaneous channel coefficient, and $n_n(t)$ represents the additive complex cycle symmetric Gaussian noise at the train following $CN(0, \sigma^2)$. In this study, we assume that the train runs in an open area and that the channel coefficient is dominated by the line-of-sight (LOS) component without any scatter. In this way, the propagation attenuation model can be represented as $h_n(t) = \frac{\sqrt{G}}{d_n^\alpha(t)}$, where $G$ shows the constant channel gain and $\alpha$ represents the path-loss exponent.

At the receiver, the received signals can be combined using a maximal ratio combiner. The instantaneous achievable rate at time $t$ can be given as

$$C(t) = B \log_2 \left( 1 + \sum_{n \in \mathcal{N}} \frac{P_n(t)|h_n(t)|^2}{\sigma^2} \right), \quad (3)$$

where $\mathcal{N} = \{1, 2\}$ and $B$ denotes the frequency bandwidth that is distributed for every channel between an RRH and the train.

### B. Problem formulation

In the considered Fog-RAN-assisted downlink transmission, if the content requested by the train has been cached at the serving RRH, the serving RRH is able to instantly convey the content to the train; if not, fetching the content from the BBU via backhaul links is required for the RRH , which consumes extra backhaul transmission resources. Denote the instantaneous content delivery rate over the backhaul link between RRH $n$ and the train at time $t$ as $R_n(t)$. The target of dynamic power allocation is to reduce the entire network power cost as much as we can, which includes the backhaul power consumption and RRH convey power consumption. Specifically, assuming that content $l$ is requested by the train, the backhaul power consumption can be represented as

$$\text{Cost}_b = \int_0^T \sum_{n=1}^2 \beta || \int_0^T P_n(t)dt ||_0 (1 - c_{n,l}) R_n(t) dt. \quad (4)$$

Here, we consider that only the backhaul link associated with the active RRH at which the requested content is not cached consumes extra power. In (4), we use the term $|| \int_0^T P_n(t)dt ||_0$ to indicate the active RRH, which implies that if $P_n(t)$ is not always zero over time period $(0, T]$, RRH $n$ is an active RRH. $1 - c_{n,l}$ is used to indicate the impact of content $l$ caching on the backhaul link associated with RRH $n$. We observe that if $c_{n,l} = 1$, which means that content $l$ has been at RRH $n$, then $1 - c_{n,l} = 0$ indicates that no extra power is needed for backhaul link $n$; otherwise, extra power is required for

backhaul link $n$. Parameter $\beta$ in (4) is a constant number establishing a relationship between the rate $R_n(t)$ and the cost of power.

In addition, the entire RRH transmit power consumed over time period $(0, T]$ can be denoted as

$$\text{Cost}_p = \int_0^T \sum_{n \in \mathcal{N}} P_n(t)dt. \tag{5}$$

In this way, the network power cost in total can be written as

$$\text{Cost} = \text{Cost}_b + \text{Cost}_p. \tag{6}$$

In addition to minimizing the total network power consumption, our dynamic power allocation also considers several QoS-related constraints. The most important one is the delay requirement. Basically, for an RRH that does not cache the requested content, the delay contains two hops: one for the backhaul link and another for the wireless transmission link between this RRH and the train. However, as we here assume that $R_n(t) \geq C(t)$ always succeeds, the instantaneous delay can be written as

$$\tau(t) = \frac{1}{C(t)}. \tag{7}$$

In brief, the overall dynamic power allocation problem is formulated as

$$\min_{P_n(t)} \text{Cost} \tag{8a}$$

$$s.t. \quad \frac{1}{T}\int_0^T P_n(t)dt \leq P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{8b}$$

$$\tau(t) \leq \tau_{\max} \tag{8c}$$

$$\int_0^T C(t)dt \geq Q \tag{8d}$$

$$P_n(t) \geq 0, \tag{8e}$$

where (8b) denotes the average power constraint of every RRH, with $P_{n,\text{avg}}$ being the maximum average power at RRH $n$; (8c) denotes the instantaneous transmission delay requirement, with $\tau_{\max}$ being the maximum delay requirement; (8d) is the requested content delivery that requires to be completed through the network during the time period $(0, T]$; and (8e) shows that the instantaneous power at each time $t$ should not be negative. Our final objective is to minimize the total network cost through maximizing the instantaneous power at every RRH and the instantaneous transmission rate over each backhaul link.

## III. DYNAMIC POWER ALLOCATION FOR HIGH-SPEED RAILWAY FOG-RAN

In this section, our solution to (8) is presented. Different from the widely studied conventional static power allocation problem, our considered dynamic is more challenging. The reason includes the following aspects. First, because our optimization problem considered the backhaul consumption, a nonconvex $l_0$-norm function is introduced, which results in the nonconvexity of our problem. Moreover, our optimization variable is a function with respect to time $t$, and the constraints in problem (8) include an integration form. Basically, the

dynamic optimization problems are widely applied in the fields of smart power systems, robotics, etc. [27]. In what follows, by using certain approximations, we present several ways to find the optimal solution to the approximated problem.

To address the nonconvex objective function in (8), we utilize the continuous smooth log-function to approximate the $l_0$-norm function as [13]

$$||x||_0 \approx \frac{\log\left(\frac{x}{\theta} + 1\right)}{\log\left(\frac{1}{\theta} + 1\right)}, \tag{9}$$

where the function of $\theta$ is to control the smoothness of the approximation. A smaller value of $\theta$ results in a better approximation, while leads to a worse smooth function, vice versa. With (9), the nonconvex part $\text{Cost}_b$ in (8) can be written as

$$\text{Cost}_b \approx c\beta \sum_{n=1}^{N} (1 - c_{n,f}) \log\left(\frac{\int_0^T P_n(t)dt + \theta}{\theta}\right) \\ \times \int_0^T R_n(t)dt, \tag{10}$$

where $c = \frac{1}{\log\left(\frac{1}{\theta}+1\right)}$. With (10), the objective function in (8) can be rewritten as

$$\text{Cost}_{\text{appro1}} \approx \int_0^T \sum_{n \in \mathcal{N}} P_n(t)dt + \sum_{n \in \mathcal{N}} b_n \\ \times \log\left(\frac{\int_0^T P_n(t)dt + \theta}{\theta}\right), \tag{11}$$

where $b_n = c\beta(1 - c_{n,f})\int_0^T R_n(t)dt$.

It is observed that the approximated function $\text{Cost}_{\text{appro1}}$ is still nonconvex with respect to $P_n(t)$ as it involves the summation of concave functions $\log\left(\frac{\int_0^T P_n(t)dt+\theta}{\theta}\right)$. To address this problem, we apply the majorization-minimization (MM) theory to find the upper bound of $\text{Cost}_{\text{appro1}}$. Considering that the logarithmic function is a concave function, the MM theory applied in our case involves using its first-order Taylor expansion as the upper bound. Then, in the MM algorithm, a solution to (8) is generated through minimizing the following upper-bounded function:

$$\text{Cost}_{\text{appro2}} \approx \int_0^T \sum_{n \in \mathcal{N}} P_n(t)dt + \sum_{n \in \mathcal{N}} b_n \Bigg[ \\ \log\left(\frac{\theta + \int_0^T P_n^0(t)dt}{\theta}\right) + \frac{\int_0^T P_n(t)dt - \int_0^T P_n^0(t)dt}{\theta + \int_0^T P_n^0(t)dt} \Bigg], \tag{12}$$

where $\int_0^T P_n^0(t)dt$ denotes a basis point of the Taylor expansion of $\log\left(\frac{\int_0^T P_n(t)dt+\theta}{\theta}\right)$.

Assuming $k_n = 1 + b_n \frac{1}{\theta + \int_0^T P_n^0(t)dt}$, solving (8) finally

reduces to solving the following convex problem:

$$\min_{P_n(t)} \int_0^T \left( \sum_{n \in \mathcal{N}} k_n P_n(t) \right) dt \tag{13a}$$

$$s.t. \quad \frac{1}{T} \int_0^T P_n(t) dt \leq P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{13b}$$

$$C(t) \geq \frac{1}{\tau_{\max}} \tag{13c}$$

$$\int_0^T C(t) dt \geq Q \tag{13d}$$

$$P_n(t) \geq 0. \tag{13e}$$

For solving (13), firstly, we obtain the following theorem.

THEOREM 1. *If $\frac{T}{\tau_{\max}} \geq Q$, we have $C(t) = \frac{1}{\tau_{\max}}$ at the optimal solution, and the optimal solution of (13) can be represented as*

$$P_1(t) = \tilde{a}_0(t) - \tilde{a}_2(t) P_2(t), \tag{14}$$

*where $\tilde{a}_n(t)$ is as defined in (17). Denote the critical time points $\{t', t''\}$ and $\tilde{t}'$ as defined in (29) and (31), respectively; the optimal $P_2(t)$ is given in the following four cases:*

- *If $\int_0^T \tilde{a}_3(t) dt \leq T P_{2,\text{avg}}$ and $B \leq 0$*

$$P_2^*(t) = \begin{cases} 0 & 0 < t < t' \\ \tilde{a}_3(t) & t' \leq t \leq T \end{cases}$$

- *If $\int_0^T \tilde{a}_3(t) dt \leq T P_{2,\text{avg}}$ and $B > 0$*

$$P_2^*(t) = \begin{cases} 0 & 0 < t < \min\{t', t''\} \\ \tilde{a}_3(t) & \min\{t', t''\} \leq t \leq T \end{cases}$$

- *If $\int_0^T \tilde{a}_3(t) dt > T P_{2,\text{avg}}$ and $B \leq 0$*

$$P_2^*(t) = \begin{cases} 0 & 0 < t < \max\{t', \tilde{t}'\} \\ \tilde{a}_3(t) & \max\{t', \tilde{t}'\} \leq t \leq T \end{cases} \tag{15}$$

- *If $\int_0^T \tilde{a}_3(t) dt > T P_{2,\text{avg}}$ and $B \leq 0$*
  - *If $t' \geq \tilde{t}'$, the optimal solution is given by (15).*
  - *Otherwise, the optimal solution can be obtained by solving linear program problem (33).*

*Proof:* It is noted that if $\frac{T}{\tau_{\max}} \geq Q$ succeeds, for condition (13d), we have $\int_0^T C(t) dt \geq \int_0^T \frac{1}{\tau_{\max}} dt = \frac{T}{\tau_{\max}} \geq Q$. This implies that condition (13d) in (13) is redundant. Next, we show that at the optimal solution, (13c) must be active. We prove this result using a contradiction. Assume that at the optimal solution of (13), the optimal $P_n(t)$ makes $C(t) > \frac{1}{\tau_{\max}}$ at some time $t$. Now, we can multiply $P_n(t)$ by a coefficient $\delta_n(t) \in (0, 1)$, which can further reduce the value of the objective function while not violating the constraints. According to the above analysis, under the condition of $\frac{T}{\tau_{\max}} \geq Q$, problem (13) is equivalent to

$$\min_{P_n(t)} \int_0^T \left( \sum_{n \in \mathcal{N}} k_n P_n(t) \right) dt \tag{16a}$$

$$s.t. \quad \int_0^T P_n(t) dt \leq T P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{16b}$$

$$B \log_2 \left( 1 + \sum_{n \in \mathcal{N}} \frac{G P_n(t)}{d_n(t)^\alpha \sigma^2} \right) = \frac{1}{\tau_{\max}} \tag{16c}$$

$$P_n(t) \geq 0, \forall n. \tag{16d}$$

With the equality constraint (16c), by denoting $a_n(t) = \frac{G}{d_n(t)^\alpha \sigma^2}$, we have

$$P_1(t) = \tilde{a}_0(t) - \tilde{a}_2(t) P_2(t), \tag{17}$$

where $\tilde{a}_0(t) = \frac{2^{\frac{1}{B \tau_{\max}}} - 1}{a_1(t)}$ and $\tilde{a}_2(t) = \frac{a_2(t)}{a_1(t)}$. Problem (16) can be further transformed into

$$\min_{P_2(t)} \int_0^T \left( k_2 - k_1 \tilde{a}_2(t) \right) P_2(t) dt \tag{18a}$$

$$s.t. \quad \int_0^T P_2(t) dt \leq T P_{2,\text{avg}} \tag{18b}$$

$$B \leq \int_0^T \tilde{a}_2(t) P_2(t) dt \leq A \tag{18c}$$

$$0 \leq P_2(t) \leq \tilde{a}_3(t), \tag{18d}$$

where $A = \int_0^T \tilde{a}_0(t) dt$, $B = \int_0^T \tilde{a}_0(t) dt - T P_{1,\text{avg}}$ and $\tilde{a}_3(t) = \frac{\tilde{a}_0(t)}{\tilde{a}_2(t)}$. Constraint (18c) comes from the fact that $0 \leq \int_0^T P_1(t) dt \leq T P_{1,\text{avg}}$, while (18d) comes from the fact that $P_1(t) \geq 0$.

According to the definitions of $\tilde{a}_3(t)$ and $\tilde{a}_2(t)$, we observe that if constraint (18d) is satisfied, we have

$$\int_0^T \tilde{a}_2(t) P_2(t) dt \leq \int_0^T \tilde{a}_2(t) \tilde{a}_3(t) dt = A. \tag{19}$$

Then, problem (18) can be equivalent to

$$\min_{P_2(t)} \int_0^T \left( k_2 - k_1 \tilde{a}_2(t) \right) P_2(t) dt \tag{20a}$$

$$s.t. \quad \int_0^T P_2(t) dt \leq T P_{2,\text{avg}} \tag{20b}$$

$$\int_0^T \tilde{a}_2(t) P_2(t) dt \geq B \tag{20c}$$

$$0 \leq P_2(t) \leq \tilde{a}_3(t). \tag{20d}$$

To find the analytical solution to (20), we next discuss certain particular cases that help simplify the problem.

***Case 1)*** when $\int_0^T \tilde{a}_3(t) dt \leq T P_{2,\text{avg}}$ and $B \leq 0$ are satisfied.

In this case, constraints (33b) and (33c) are redundant. Problem (20) reduces to

$$\min_{P_n(t)} \int_0^T \left( k_2 - k_1 \tilde{a}_2(t) \right) P_2(t) dt \tag{21a}$$

$$s.t. \quad 0 \leq P_2(t) \leq \tilde{a}_3(t). \tag{21b}$$

It is noted that the optimal solution to (21) depends on the sign of $k_2 - k_1 \tilde{a}_2(t)$. To minimize the objective function, the optimal solution to (21) is given as

$$P_2^*(t) = \begin{cases} 0 & t \in \{t | k_2 - k_1 \tilde{a}_2(t) > 0\} \\ \tilde{a}_3(t) & t \in \{t | k_2 - k_1 \tilde{a}_2(t) \leq 0\} \end{cases}. \tag{22}$$

***Case 2)*** when $\int_0^T \tilde{a}_3(t) dt \leq T P_{2,\text{avg}}$ and $B > 0$ are satisfied.

In this case, constraint (33b) is redundant. Problem (20) reduces to

$$\min_{P_2(t)} \int_0^T \left( k_2 - k_1 \tilde{a}_2(t) \right) P_2(t) dt \quad (23a)$$

$$s.t. \quad \int_0^T \tilde{a}_2(t) P_2(t) dt \geq B \quad (23b)$$

$$0 \leq P_2(t) \leq \tilde{a}_3(t). \quad (23c)$$

Within the time period $(0, T]$, as the train departs from RRH 1 and move closer to RRH 2, we see that functions $k_2 - k_1 \tilde{a}_2(t)$, $\tilde{a}_2(t)$, and $\tilde{a}_3(t)$ are decreasing, increasing and decreasing functions with respect to $t$, respectively. According to this observation, the two critical time points are defined as

$$t' \triangleq \{t | k_2 - k_1 \tilde{a}_2(t) = 0\}$$
$$t'' \triangleq \left\{ t | \int_0^T \tilde{a}_2(t) \tilde{a}_3(t) dt = \int_0^T \tilde{a}_0(t) dt = B \right\}. \quad (24)$$

It is noted that for the time region $(t', T]$, $P_2(t)$ should be equal to $\tilde{a}_3(t)$ to minimize the value of the objective. Therefore, in the scenario with $t' \geq t''$, the optimal solution is presented as

$$P_2^*(t) = \begin{cases} 0 & t < t'' \\ \tilde{a}_3(t) & t \geq t'' \end{cases}. \quad (25)$$

For the scenario with $t' < t''$, it is easy to see that for the time region over $t \in (t'', T]$, the optimal $P_2(t)$ should be equal to $\tilde{a}_3(t)$ as given in (25). However, this kind of solution cannot satisfy the constraint (23b). To this end, we need to find a certain $P_2(t)$ in the time period $t \in (0, t'']$ to satisfy

$$\int_0^{t''} \tilde{a}_2(t) P_2(t) dt = B - \int_{t''}^T \tilde{a}_0(t) dt. \quad (26)$$

Next we show that the optimal solution over $t \in (0, t'']$ is

$$P_2^*(t) = \begin{cases} 0 & t < t' \\ \tilde{a}_3(t) & t \in (t', t''] \end{cases}. \quad (27)$$

We prove this result with contradiction analysis. Assume that there is a new optimal solution $P_2^{**}(t)$ over $t \in (0, t'']$ given by

$$P_2^{**}(t) = \begin{cases} 0 & t < \tilde{t} \\ < \tilde{a}_3(t) & t \in (\tilde{t}, t''] \end{cases}. \quad (28)$$

To satisfy (26), we have $\tilde{t} < t'$. However, $k_2 - k_1 \tilde{a}_2(t)$ is a decreasing function over $t$. Hence, $\int_0^{t''} \left( k_2 - k_1 \tilde{a}_2(t) \right) P_2^{**}(t) dt$ must be larger than $\int_0^{t''} \left( k_2 - k_1 \tilde{a}_2(t) \right) P_2^*(t) dt$ with $P_2^*(t)$ given in (27), which implies that $P_2^{**}(t)$ cannot be an optimal solution. According to the above analysis, the optimal solution for *Case 2* is denoted as

$$P_2^*(t) = \begin{cases} 0 & t < \min\{t', t''\} \\ \tilde{a}_3(t) & t \geq \min\{t', t''\} \end{cases}. \quad (29)$$

***Case 3)*** when $\int_0^T \tilde{a}_3(t) dt > T P_{2,\text{avg}}$ and $B \leq 0$ are satisfied.

In this case, the constraint (33c) is redundant. Problem (20) reduces to

$$\min_{P_2(t)} \int_0^T \left( k_2 - k_1 \tilde{a}_2(t) \right) P_2(t) dt \quad (30a)$$

$$s.t. \quad \int_0^T P_2(t) dt \leq T P_{2,\text{avg}} \quad (30b)$$

$$0 \leq P_2(t) \leq \tilde{a}_3(t). \quad (30c)$$

We define $\tilde{t}'$ as

$$\tilde{t}' \triangleq \left\{ t | \int_{\tilde{t}'}^T \tilde{a}_3(t) dt = T P_{2,\text{avg}} \right\}. \quad (31)$$

Again, because $k_2 - k_1 \tilde{a}_2(t)$ is a decreasing function, similar to the analysis in *Case 2*, the optimal solution to problem (30) is presented as

$$P_2^*(t) = \begin{cases} 0 & t < \max\{t', \tilde{t}'\} \\ \tilde{a}_3(t) & t \geq \max\{t', \tilde{t}'\} \end{cases}. \quad (32)$$

It is noted that in (32), the critical time point is $\max\{t', \tilde{t}'\}$, which is different from *Case 2* due to the inequality constraint (30b).

***Case 4)*** when $\int_0^T \tilde{a}_3(t) dt > T P_{2,\text{avg}}$ and $B > 0$ are satisfied.

In this case, according to Lemma 3, the feasibility of problem (20) requires $\tilde{t}' \leq t''$. Under this condition, if $t' \geq \tilde{t}'$, we see that the optimal solution is equal to (29) and that constraint (33b) is redundant. Otherwise, the optimal solution can be approximately obtained through solving the following linear programming:

$$\min_{P_2(t_m)} \sum_{m=1}^M \left( k_2 - k_1 \tilde{a}_2(t_m) \right) P_2(t_m) \triangle t \quad (33a)$$

$$s.t. \quad \sum_{m=1}^M P_2(t_m) \triangle t \leq T P_{2,\text{avg}} \quad (33b)$$

$$\sum_{m=1}^M \tilde{a}_2(t_m) P_2(t_m) \triangle t \geq B \quad (33c)$$

$$0 \leq P_2(t_m) \leq \tilde{a}_3(t_m), \forall m. \quad (33d)$$

The linear problem is obtained by sampling the time period $t \in (0, T]$ as the discrete time points $\{t_1, t_2, \cdots, t_M\}$ with the adjacent sampling point interval given by $\triangle t$. It is noted that $\triangle t$ is sufficiently small; thus, we can obtain an approximately optimal solution via (33).

Now, by combining the solutions of *Case 1-Case 4*, we obtain the final optimal solution. This completes the proof of Theorem 1.
∎

Consider a special case where the train is assisted by only RRH 1 over the time period $(0, T]$, i.e., $\mathcal{N} = \{1\}$. Theorem 1 can be reduced to the following lemma.

**LEMMA 1.** *If $\frac{T}{\tau_{\max}} \geq Q$ and $\mathcal{N} = \{1\}$, we have $C(t) = \frac{1}{\tau_{\max}}$ at the optimal solution, and the optimal solution to (13) can be denoted as*

$$P_1(t) = \left( 2^{\frac{1}{B \tau_{\max}}} - 1 \right) \frac{d_1^\alpha(t) \sigma^2}{G}. \quad (34)$$

*Proof:* When $\mathcal{N} = \{1\}$, problem (13) can be written as

$$\min_{P_1(t)} \int_0^T \left( k_1 P_1(t) \right) dt \tag{35a}$$

$$s.t. \quad \frac{1}{T} \int_0^T P_1(t) dt \leq P_{1,\text{avg}} \tag{35b}$$

$$C(t) = \frac{1}{\tau_{\max}} \tag{35c}$$

$$P_1(t) \geq 0. \tag{35d}$$

If the problem is feasible, the solution is determined by constraint (35c), which completes the proof of Lemma 1. ∎

If condition $\frac{T}{\tau_{\max}} \geq Q$ is not satisfied, the conclusions shown in Theorem 1 and Lemma 1 will not be applicable. To propose an efficient way to address this problem, we first give the following lemma.

LEMMA 2. *If $\frac{T}{\tau_{\max}} < Q$, the optimal solution to (13) can be represented as*

$$P_n(t) = P_{n,1}(t) + P_{n,2}(t), \tag{36}$$

*where $P_{n,1}(t)$ is the optimal solution obtained by solving (16) and $P_{n,2}(t)$ is the solution to*

$$\min_{P_{n,2}(t)} \int_0^T \left( \sum_{n \in \mathcal{N}} k_n P_{n,2}(t) \right) dt \tag{37a}$$

$$s.t. \quad \frac{1}{T} \int_0^T P_{n,2}(t) dt \leq b_n \quad \forall n \in \mathcal{N} \tag{37b}$$

$$\sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,2}(t) \geq 0 \tag{37c}$$

$$\int_0^T B \log_2 \left( c_n(t) + \sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,2}(t) \right) dt \geq Q \tag{37d}$$

$$P_{n,2}(t) \geq -P_{n,1}^*(t), \tag{37e}$$

*where $b_n = P_{n,\text{avg}} - \frac{1}{T} \int_0^T P_{n,1}^*(t) dt$, $\frac{G}{d_n^\alpha(t) \sigma^2} = \kappa_n(t)$ and $c_n(t) = 1 + \sum_{n=1}^N \kappa_n(t) P_{n,1}^*(t)$, with $P_{n,1}^*(t)$ being the optimal solution to $P_{n,1}(t)$.*

*Proof:* It is noted that for any feasible $P_n(t)$ of (13), we can always decompose it into a sum of $P_{n,1}(t)$ and $P_{n,2}(t)$. Basically, the choice of $P_{n,1}(t)$ and $P_{n,2}(t)$ can be arbitrary as long as their sum is equal to $P_n(t)$. Without reduction of generality, we assume that the term $P_{n,1}(t)$ is chosen as the optimal solution to (16), denoted by $P_{n,1}^*(t)$. Then, problem

(13) changes to

$$\min_{P_{n,2}(t)} \int_0^T \left( \sum_{n \in \mathcal{N}} k_n(P_{n,1}^*(t) + P_{n,2}(t)) \right) dt \tag{38a}$$

$$s.t. \quad \frac{1}{T} \int_0^T \left( P_{n,1}^*(t) + P_{n,2}(t) \right) dt \leq P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{38b}$$

$$B \log_2 \left( 1 + \sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,1}^*(t) + \right.$$
$$\left. \sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,2}(t) \right) dt \geq \frac{1}{\tau_{\max}} \tag{38c}$$

$$\int_0^T B \log_2 \left( 1 + \sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,1}^*(t) + \right.$$
$$\left. \sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,2}(t) \right) dt \geq Q \tag{38d}$$

$$P_{n,1}^*(t) + P_{n,2}(t) \geq 0. \tag{38e}$$

Problem (38) involves only solving variable $P_{n,2}(t)$. We next rewrite the constraints in (38) by analyzing the value range of $P_{n,2}(t)$. Because we set $B \log_2 \left( 1 + \sum_{n \in \mathcal{N}} \frac{G P_{n,1}^*(t)}{d_n^\alpha(t) \sigma^2} \right) = \frac{1}{\tau_{\max}}$, it is observed that $\sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,2}(t) \geq 0$ must be satisfied; otherwise, the delay constraint (38c) is violated. Then, by rewriting constraints (38c) and (38e) as (37c) and (37e), respectively, we have Problem (37). This completes the proof of Lemma 2.

∎

With Lemma 2, solving $P_n(t)$ from (13) under the condition $\frac{T}{\tau_{\max}} < Q$ reduces to finding $P_{n,2}(t)$ from (37). The result in Lemma 2 will be useful in the performance tradeoff analysis in Section IV.

We can easily observe that problem (37) is a convex problem. Then, to obtain the optimal solution, we construct an algorithm based on the Karush-Kuhn-Tucker (KKT) conditions. To proceed, firstly, the Lagrangian function is presented as

$$L = \int_0^T \left( \sum_{n \in \mathcal{N}} \left( k_n P_{n,2}(t) + \mu_{1,n}(P_{n,2}(t) - T b_n) \right) \right.$$
$$\left. - \mu_2 \left( B \log_2(c_n(t) + \sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,2}(t)) - Q \right) \right) dt$$
$$- \mu_3(t) \sum_{n=1}^N \kappa_n(t) P_{n,2}(t), \tag{39}$$

where $\mu_{1,n}$, $\mu_2$ and $\mu_3(t)$ are non-negative multipliers related to the constraints (37b), (37d) and (37c), respectively. To minimize the Lagrangian function, it is necessary to differentiate the Lagrangian function with respect to $P_{n,2}(t)$ and set the derivative to zero for each time $t$, which is,

$$\frac{\partial L}{\partial P_{n,2}(t)} = k_n + \mu_{1,n} - \mu_3(t) \kappa_n(t) -$$
$$\frac{\mu_2 B}{\log 2} \frac{\kappa_n(t)}{c_n(t) + \sum_{n \in \mathcal{N}} \kappa_n(t) P_{n,2}(t)} = 0. \tag{40}$$

Through combining with constraint (37e), the solution can be obtained given by

$$P_{n,2}(t) = \left[ \left( \frac{\mu_2 GB}{\log 2 d_n(t)^\alpha \sigma^2 (k_n + \mu_1 - \mu_3(t)\kappa_n(t))} - c_n(t) - \sum_{m \neq n} \frac{GP_{m,2}(t)}{d_m(t)^\alpha \sigma^2} \right) \times \frac{d_n(t)^\alpha \sigma^2}{G}, -P_{n,1}^*(t) \right]^+ .$$
(41)

(41) shows that $P_{n,2}(t)$ values with different $n$ are coupled with each other, and the final solution of $P_{n,2}(t)$ can be obtained by iteratively updating them until convergence. During the iteration, Lagrangian multipliers $\mu_{1,n}$, $\mu_2$ and $\mu_3(t)$ can be obtained via the subgradient technique.

Another way to solve (37) is to sample the time period to generate certain discrete time points. If the time interval between two adjacent discrete points is small enough, the obtained result can approximately be considered as a solution to (37). Denote the discrete time points as $\{t_1, t_2, \cdots, t_M\}$ and the adjacent sampling point interval as $\triangle t$; the power $P_{n,2}(t_m)$ can be efficiently obtained through solving the following problem:

$$\min_{P_{n,2}(t_m)} \sum_{m=1}^{M} \sum_{n \in \mathcal{N}} k_n P_{n,2}(t_m) \triangle t \tag{42a}$$

$$s.t. \quad \sum_{m=1}^{M} P_{n,2}(t_m) \triangle t \leq T b_n \quad \forall n \in \mathcal{N} \tag{42b}$$

$$\sum_{n=1}^{N} \kappa_n(t_m) P_{n,2}(t_m) \geq 0 \quad \forall m \tag{42c}$$

$$\sum_{m=1}^{M} B \log_2 \left( c_n(t_m) + \sum_{n=1}^{N} \kappa_n(t_m) \right.$$
$$\left. \times P_{n,2}(t_m) \right) \triangle t \geq Q \tag{42d}$$

$$P_{n,2}(t_m) \geq -P_{n,1}^*(t_m) \; \forall m. \tag{42e}$$

It is noted that problem (42) has only one nonlinear convex constraint (42d), and thus can be efficiently solved by an interior point algorithm using the software CVX [28].

To summarize, the algorithm we proposed to solve (13) is as follows:

## IV. Cost, Delay and Delivery Content Size Tradeoff Analysis

In our system design, our target is to minimize the total network cost subject to the delay constraint and total content delivery constraint. Both constraints actually determine the total cost consumed by the system. In fact, there is inherently a tradeoff between the total network cost, delay and delivery content size. In this section, we present the tradeoff analysis of the system, which may help simplify the design of the system.

We first present the tradeoff between the total network cost and the delay requirement for a given delivery content size.

PROPOSITION 1. *In the considered system design problem (13), for a given delivery content size $Q$, we have*

---

**Algorithm 1** Dynamic power optimization in problem (13)

- **Input:** BS interval $d$, speed $v_0$, BS height $h_0$, time interval $(0, T)$, BS coordinate $(l_n, d_0)$, content size $Q$, local storage size $F_n$, cache placement matrix $C$, bandwidth $B$, noise power $\sigma^2$, backhaul cost ratio $\beta$, backhaul rate $R_n(t)$, delay requirement $\tau_{\max}$, average power of BS $P_{n,\text{avg}}$.
- **Initialization:** Basis point of the Taylor expansion of $P_n(t)$, that is, $P_n^0(t)$.
- **Output:** Dynamic power allocation $P_n(t)$.
- **While** not convergence **do**
  - Update power $P_n(t)$;
    1) Update $P_n(t)$ using Theorem 1 if $\frac{T}{\tau_{\max}} \geq Q$;
    2) Update $P_n(t)$ using KKT conditions or solving (42) if $\frac{T}{\tau_{\max}} < Q$;
  - Update $P_n^0(t)$ using $P_n(t)$;
- **End while**

---

- *when $\frac{T}{\tau_{\max}} \geq Q$, the total network cost increases as the delay requirement becomes strict.*
- *when $\frac{T}{\tau_{\max}} < Q$, if we denote the optimal solution to (13) as $P_n^*(t)$, which is decomposed into the sum $P_{n,1}^*(t) + P_{n,2}^*(t)$, with $P_{n,1}^*(t)$ being the optimal power term obtained by solving (16), there exists a delay requirement region*

$$\tau_{\max}^{\text{reg}} = \left[ \frac{1}{\tau}, \tau_{\max} \right], \tag{43}$$

*where $\tau = B \log_2(1 + \sum_{n \in \mathcal{N}} \frac{GP_{n,1}^*(t)}{d_n^\alpha(t)\sigma^2} + \min_t \sum_{n \in \mathcal{N}} \frac{GP_{n,2}^*(t)}{d_n^\alpha(t)\sigma^2})$ such that the total network cost does not increase with the decrease in delay.*

*Proof:* Under the condition $\frac{T}{\tau_{\max}} \geq Q$, the original system design problem (13) is equivalent to problem (16), where we have only the power constraint and the delay constraint. At the optimal solution, the delay constraint is always active. Hence, if we decrease the value of $\tau_{\max}$ to achieve a strict delay requirement, the total network cost must be increased.

If under the condition $\frac{T}{\tau_{\max}} < Q$, for a given delay requirement $\tau_{\max}$, the optimal solution is $P_n^*(t)$, which is decomposed into the sum $P_{n,1}^*(t) + P_{n,2}^*(t)$, with $P_{n,1}^*(t)$. We first prove that the solution $P_n^*(t)$ with its decomposition $P_{n,1}^*(t)$ and $P_{n,2}^*(t)$ is also the solution to (13) for a smaller given delay requirement $\tau_{\max}' \in \tau_{\max}^{\text{reg}}$. Then, we prove that this solution is the optimal solution.

It is noted that as $B \log_2 \left(1 + \sum_{n \in \mathcal{N}} \frac{P_{n,1}^*(t)|h_n(t)|^2}{\sigma^2}\right) = \tau_{\max}$, we have $\sum_{n \in \mathcal{N}} \frac{P_{n,2}^*(t)|h_n(t)|^2}{\sigma^2} \geq 0$, as claimed in Lemma 2. This further produces

$$B \log_2 \left(1 + \sum_{n \in \mathcal{N}} \frac{P_n^*(t)|h_n(t)|^2}{\sigma^2}\right)$$
$$\geq B \log_2 \left(1 + \sum_{n \in \mathcal{N}} \frac{GP_{n,1}^*(t)}{d_n^\alpha(t)\sigma^2} + \min_t \sum_{n \in \mathcal{N}} \frac{GP_{n,2}^*(t)}{d_n^\alpha(t)\sigma^2}\right) \tag{44}$$
$$\geq \frac{1}{\tau_{\max}}.$$

Thus, if we change the delay requirement $\tau_{\max}$ to any smaller value $\tau'_{\max}$ in region $\boldsymbol{\tau}^{\text{reg}}_{\max}$, $P^*_n(t)$ can still be a solution to the design problem (13), and they have the same total network cost.

In the following, we prove that $P^*_n(t)$ is the optimal solution to the design problem (13) with a smaller delay requirement $\tau'_{\max} \in \boldsymbol{\tau}^{\text{reg}}_{\max}$, $P^*_n(t)$. It is noted that if we decrease the delay requirement $\tau_{\max}$ to $\tau'_{\max}$, the design problem (13) with delay requirement $\tau_{\max}$ has a smaller feasible region than that of the design problem with delay requirement $\tau'_{\max}$. Then, the value of the objective function, i.e., the total network cost, of the former cannot be smaller than that of the latter. This indicates that $P^*_n(t)$ is still the optimal solution to the design problem with delay requirement $\tau'_{\max}$, which completes the proof of Proposition . ∎

Next, we discuss the tradeoff between the total network cost and delivery content size for a given delay requirement.

PROPOSITION 2. *In the considered system design problem (13), for a given delay requirement $\tau_{\max}$, we have*

- *when $Q \leq \frac{T}{\tau_{\max}}$, increasing the delivery content size will not lead to an increase in the total network cost.*
- *when $Q > \frac{T}{\tau_{\max}}$, increasing the delivery content size must also increase the total network cost.*

*Proof:* Under the condition $Q \leq \frac{T}{\tau_{\max}}$, because constraint (13d) is the same as (13), the system cost is determined only by the delay requirement. Therefore, increasing the delivery content size will not lead the total network cost to increase.

Under the condition $Q > \frac{T}{\tau_{\max}}$, based on the result presented in Lemma 2, with a constant power $P^*_{n,1}(t)$, we have to increase the value of $\sum_{n \in \mathcal{N}} \frac{P^*_{n,2}(t)|h_n(t)|^2}{\sigma^2}$ to meet constraint (13d) with larger $Q$. This leads to a larger total network cost. It is noted that when increasing the size of delivery content in (13), we cannot find a solution to $P^*_{n,2}(t)$ that satisfies constraint (13d) while not increasing the value of the objective function in (13). We prove this conclusion by using a contradiction statement. Assume that with a delivery content size $Q'$, the optimal solution to (13) is $P'^*_n(t) = P'^*_{n,1}(t) + P'^*_{n,2}(t)$. Now, assume that with a larger delivery content size $Q''$, the optimal solution to (13) is $P''^*_n(t) = P'^*_{n,1}(t) + P''^*_{n,2}(t)$. If $P'^*_n(t)$ and $P''^*_n(t)$ have the same objective function value in (13), we can always obtain a new solution to (13) with a delivery content size $Q'$ as $\tilde{P}'^*_n(t) = P'^*_{n,1}(t) + \alpha P''^*_{n,2}(t)$, where $\alpha$ is a positive value smaller than 1. This new solution $\tilde{P}'^*_n(t)$ produces a smaller objective function value than $P'^*_n(t)$. This contradicts the fact that $P'^*_n(t)$ is the optimal solution, which completes the proof of Proposition 2. ∎

## V. INVARIANT POWER OPTIMIZATION WITH QOS CONSTRAINTS

As another simple power allocation scheme, we consider a constant power optimization design where power does not vary with the channel. Under this situation, the overall optimization problem can be modified as

$$\min_{P_n} \quad T \sum_{n \in \mathcal{N}} P_n + \tag{45a}$$

$$\int_0^T \sum_{n \in \mathcal{N}} \beta ||P_n||_0 (1 - c_{n,f}) R_n(t) dt$$

$$s.t. \quad 0 \leq P_n \leq P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{45b}$$

$$\frac{1}{C(t)} \leq \tau_{\max} \tag{45c}$$

$$\int_0^T C(t) dt \geq Q \tag{45d}$$

$$P_n(t) \geq 0 \quad \forall n \in \mathcal{N}, \tag{45e}$$

where $C(t) = B \log_2 \left( 1 + \sum_{n \in \mathcal{N}} \frac{G P_n}{d_n(t)^\alpha \sigma^2} \right)$. Similar to the solution we presented in Section III, we determine $P_n$ by solving the following problem:

$$\min_{P_n} \sum_{n \in \mathcal{N}} k'_n P_n \tag{46a}$$

$$s.t. \quad 0 \leq P_n \leq P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{46b}$$

$$\frac{1}{C(t)} \leq \tau_{\max} \tag{46c}$$

$$\int_0^T C(t) dt \geq Q, \tag{46d}$$

where $k'_n = T + \frac{1}{\log(1/\theta+1)} \frac{\beta(1-c_{n,f}) \int_0^T R_n(t) dt}{\theta + P_n^0}$, with $P_n^0$ being a basis point of the Taylor expansion.

To solve (46), two specific cases are considered. If $\frac{T}{\tau_{\max}} \geq Q$, we have $C(t) = \frac{1}{\tau_{\max}}$. Then, constraint (46d) is redundant. $P_n$ can be found by solving

$$\min_{P_n} \sum_{n \in \mathcal{N}} k'_n P_n \tag{47a}$$

$$s.t. \quad 0 \leq P_n \leq P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{47b}$$

$$\sum_{n \in \mathcal{N}} \frac{G P_n}{d_n(t)^\alpha \sigma^2} \geq 2^{\frac{1}{\tau_{\max} B}} - 1. \tag{47c}$$

It is noted in (47) that constraint (47c) should be satisfied for any arbitrary $t$, and thus causes (47) to contain infinite constraints. We next solve (47) as problem (42) by sampling the time period at certain discrete time points. Denote the discrete time points $\{t_1, t_2, \cdots, t_M\}$ and the adjacent sampling point interval as $\triangle t$; the power $P_n$ can be efficiently obtained through solving the following linear programming problem:

$$\min_{P_{n,1}} \sum_{n \in \mathcal{N}} k'_n P_n \tag{48a}$$

$$s.t. \quad 0 \leq P_n \leq P_{n,\text{avg}} \quad \forall n \in \mathcal{N} \tag{48b}$$

$$\sum_{n \in \mathcal{N}} \frac{G P_n}{d_n(t_m)^\alpha \sigma^2} \geq 2^{\frac{1}{\tau_{\max} B}} - 1, \forall i. \tag{48c}$$

If $\frac{T}{\tau_{\max}} < Q$, similar to its dynamic counterpart, the power can be denoted as $P_n = P_{n,1} + P_{n,2}$ where $P_{n,1}$ is used to activate the constraint (48c) in (48). Then, $P_{n,2}$ can be obtained by using the Lagrangian method or time period sampling approach.

## VI. NUMERICAL RESULTS

In the following, we provide numerical results to show the superiority of using dynamic power allocation in a high-speed moving transmission scenario. In particular, we compare the performance of dynamic power allocation and invariant power optimization with respect to different caching schemes. The parameter settings for our simulation are summarized in TABLE I. Moreover, we consider three caching strategies to illustrate the effect of caching on the performance, that is, PopC, RndC and the noncaching scheme (NonC). Because the PopC caching strategy asks each RRH to cache the most popular contents until its storage is full, based on our parameter settings, RRH 1 and RRH 2 both cache contents $\{1, 2, 3, 4, 5\}$. NonC assumes that the RRH has no storage resources and that no content is cached at the RRH.

TABLE I
PARAMETER SETTINGS IN SIMULATIONS

| Parameter | Notation | Value |
|---|---|---|
| Height of RRH antennas | $h$ | 20 m |
| Interval between two RRHs | $d$ | 1000 m |
| Distance between RRH and road | $d_0$ | 100 m |
| Coordinates of RRH 1 | $(l_1, d_0)$ | $(-200$ m, $100$ m$)$ |
| Coordinates of RRH 2 | $(l_2, d_0)$ | $(800$ m, $100$ m$)$ |
| Path-loss exponent | $\alpha$ | 0.8 |
| Channel gain | $G$ | 2 |
| Train speed | $v_0$ | 200 Km/h |
| Ratio of backhaul power cost and rate | $\beta$ | 2.8 |
| Noise power | $\sigma^2$ | 1 |
| Content number | $L$ | 15 |
| Content size | $Q$ | 1 |
| RRH storage size | $F_n$ | 5 |
| Shaping parameter of content popularity | $\eta$ | 1 |

In Fig. 2, the convergence behavior of the proposed algorithm is illustrated. We can observe that the proposed algorithm converges fast in no more than five iterations. Three curves with different delay requirements also demonstrate that a decrease in delay can significantly increase the total network power cost.

In Fig. 3, the dynamic power is illustrated with the time-varying channel. As the train departs from RRH 1 and moves closer to RRH 2, we see that the channel gain of $h_1(t)$ decreases as time passes, while the channel gain of $h_2(t)$ increases with time. To satisfy the QoS requirements, the power at RRH 1 gradually increases to compensate for the channel gain loss of $h_1(t)$ while the power at RRH 2 remains zero. As the train moves closer to RRH 2, RRH 2 begins to serve it, and the power at RRH 1 becomes zero to maintain a lower network power cost. In particular, when the train gets close to RRH 2, the required power at RRH 2 gradually decreases as the quality of channel $h_1(t)$ improves.

In Fig. 4, we illustrate the effect of the delay requirements on the total network power cost for different caching strategies. It is observed that a stricter delay requirement enhances the total network power cost and that the proposed dynamic power allocation significantly outperforms the invariant power allocation. Moreover, the PopC caching strategy has the lowest total power cost. The RndC strategy performs worse than the PopC strategy. The NonC scheme has the maximum total power cost. The result implies that caching at the RRH is
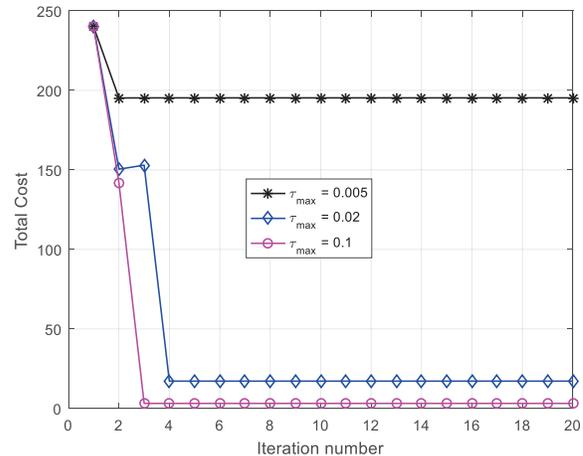


Fig. 2. Convergence behavior of the proposed Algorithm 1 at an average SNR $\frac{P_{1,\text{avg}}}{\sigma^2} = \frac{P_{2,\text{avg}}}{\sigma^2} = 10$ dB.
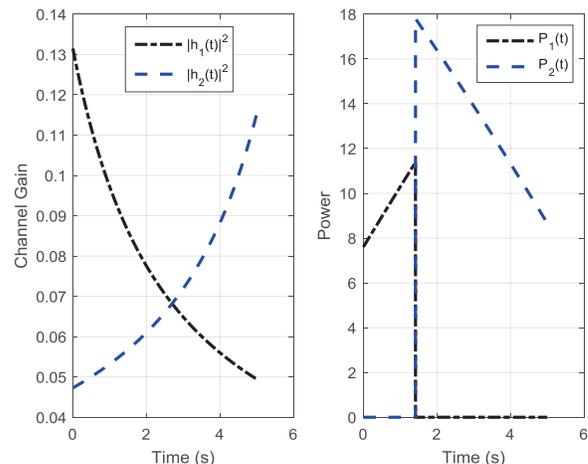


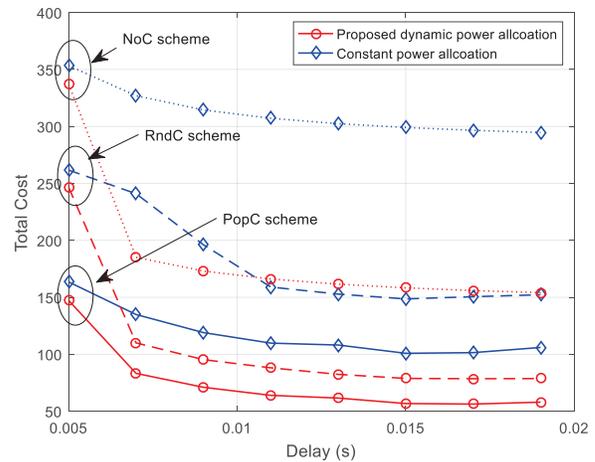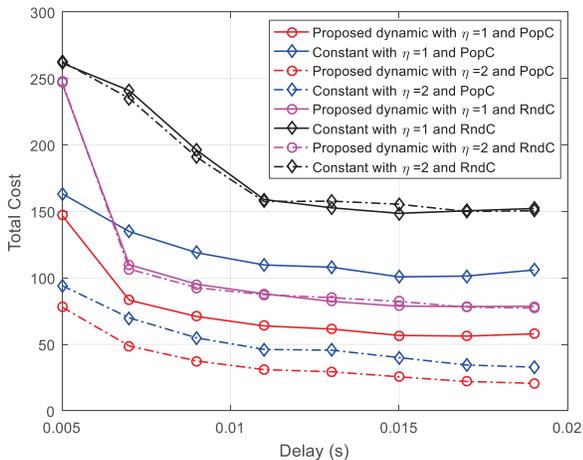Fig. 3. Power variance with a dynamic channel.



Fig. 4. Total power cost for different delay requirements at an average SNR $\frac{P_{1,\text{avg}}}{\sigma^2} = \frac{P_{2,\text{avg}}}{\sigma^2} = 10$ dB.
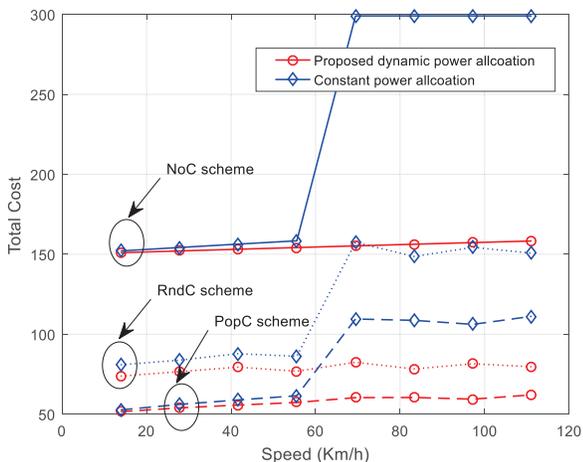
Fig. 5. Total power cost for different delay requirements and $\eta$ at an average SNR $\frac{P_{1,\mathrm{avg}}}{\sigma^2} = \frac{P_{2,\mathrm{avg}}}{\sigma^2} = 10$ dB.



Fig. 6. Total power cost for different train speeds at an average SNR $\frac{P_{1,\mathrm{avg}}}{\sigma^2} = \frac{P_{2,\mathrm{avg}}}{\sigma^2} = 10$ dB.

beneficial for reducing the network power consumption and that caching popular contents is more helpful.

Fig. 5 further illustrates the effect of the shaping parameter of popularity $\eta$. In general, a larger $\eta$ implies a larger popularity difference among the contents. We observe that a larger $\eta$ produces a lower total network power for the PopC caching strategy. This observation is reasonable as the request contents are very likely to be cached at the RRHs, which can reduce the backhaul power consumption. However, for the RndC strategy, we see that the change in $\eta$ has little effect on the total network power cost. The main reason is that the RndC strategy does not consider the content popularity and caches all contents with equal probabilities.

In Fig. 6, we show the effect of train speed on the total network power consumption. It is observed that the total network power cost increases with the speed, especially for dynamic power allocation. In the given time period $(0, T]$, a higher speed implies that a longer distance will be covered by the train. This increases the total power consumption.

## VII. CONCLUSIONS

In this paper, we studied the dynamic power allocation of the Fog-RAN-assisted high-speed railway system. For a given caching strategy, we optimized the instantaneous power allocation at the RRHs with the aim to minimize the network power consumption subject in total to several QoS constraints. By analyzing the dynamic power optimization problem, we derived the analytical power solution. Our results showed that caching at the RRHs can significantly reduce the total network power consumption. More so, the dynamic power allocation is significantly superior to the invariant one, as it takes the time-varying characteristic of the channel into consideration.

## APPENDIX A
### FEASIBILITY ANALYSIS OF PROBLEM (20)

LEMMA 3. *For problem* (20), *if* $\tilde{t}' > t''$, *the optimization problem is infeasible.*

*Proof:* With an assumption $\tilde{t}' > t''$, we have $\int_{t''}^{T} \tilde{a}_3(t)dt > TP_{2,\mathrm{avg}}$. Next, we prove the infeasibility of problem (20) by contradiction analysis. Assume that we have a feasible solution $P_2'(t)$ that satisfies $P_2'(t) < \tilde{a}_3(t)$ for $t \in (t'', T]$ and

$$\int_0^{t''} \tilde{a}_2(t)P_2'(t)dt + \int_{t''}^{T} \tilde{a}_2(t)P_2'(t)dt = \int_{t''}^{T} \tilde{a}_2(t)\tilde{a}_3(t)dt = B. \tag{49}$$

Because $\tilde{a}_2(t)$ is an increasing function, when (49) is satisfied, although $\int_{t''}^{T} P_2'(t)dt < TP_{2,\mathrm{avg}}$, we have

$$\int_0^{t''} P_2'(t)dt > \int_{t''}^{T} \tilde{a}_3(t)dt - \int_{t''}^{T} P_2'(t)dt$$
$$> TP_{2,\mathrm{avg}} - \int_{t''}^{T} P_2'(t)dt, \tag{50}$$

which implies that $P_2'(t)$ cannot be a feasible solution. We thus complete the proof of Lemma 3. ∎

## REFERENCES

[1] Z. Zhao, M. Peng, Z. Ding, W. Wang, and H. V. Poor, "Cluster content caching: An energy-efficient approach to improve quality of service in cloud radio access networks," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 5, pp. 1207–1221, 2016.

[2] D. Yan, R. Wang, E. Liu, and Q. Hou, "Admm-based robust beamforming design for downlink cloud radio access networks," *IEEE Access*, vol. 6, pp. 27912–27922, 2018.

[3] M. Peng and K. Zhang, "Recent advances in fog radio access networks: Performance analysis and radio resource allocation," *IEEE Access*, vol. 4, pp. 5003–5009, 2016.

[4] Z. Zhao, C. Feng, H. H. Yang, and X. Luo, "Federated-learning-enabled intelligent fog radio access networks: Fundamental theory, key techniques, and future trends," *IEEE Wireless Communications*, vol. 27, no. 2, pp. 22–28, 2020.

[5] Y. Jiang, C. Wan, M. Tao, F. C. Zheng, P. Zhu, X. Gao, and X. You, "Analysis and optimization of fog radio access networks with hybrid caching: Delay and energy efficiency," *IEEE Transactions on Wireless Communications*, pp. 1–1, 2020.

[6] R. Wang, R. Li, P. Wang, and E. Liu, "Analysis and optimization of caching in fog radio access networks," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 8, pp. 8279–8283, 2019.

[7] R. Wang, R. Li, E. Liu, and P. Wang, "Performance analysis and optimization of caching placement in heterogeneous wireless networks," *IEEE Communications Letters*, vol. 23, no. 10, pp. 1883–1887, 2019.

[8] W. Bai, T. Yao, H. Zhang, and V. C. M. Leung, "Research on channel power allocation of fog wireless access network based on noma," *IEEE Access*, vol. 7, pp. 32 867–32 873, 2019.

[9] G. M. S. Rahman, M. Peng, S. Yan, and T. Dang, "Learning based joint cache and power allocation in fog radio access networks," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 4, pp. 4401–4411, 2020.

[10] H. Xiang, M. Peng, Y. Sun, and S. Yan, "Mode selection and resource allocation in sliced fog radio access networks: A reinforcement learning approach," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 4, pp. 4271–4284, 2020.

[11] H. Zhang, L. Zhu, K. Long, and X. Li, "Energy efficient resource allocation in millimeter-wave-based fog radio access networks," in *2018 2nd URSI Atlantic Radio Science Meeting (AT-RASC)*, 2018, pp. 1–4.

[12] S. He, C. Qi, Y. Huang, Q. Hou, and A. Nallanathan, "Two-level transmission scheme for cache-enabled fog radio access networks," *IEEE Transactions on Communications*, vol. 67, no. 1, pp. 445–456, 2019.

[13] M. Tao, E. Chen, H. Zhou, and W. Yu, "Content-centric sparse multicast beamforming for cache-enabled cloud RAN," *IEEE Transactions on Wireless Communications*, vol. 15, no. 9, pp. 6118–6131, 2016.

[14] E. Chen, M. Tao, and Y. Liu, "Joint base station clustering and beamforming for non-orthogonal multicast and unicast transmission with backhaul constraints," *IEEE Transactions on Wireless Communications*, vol. 17, no. 9, pp. 6265–6279, 2018.

[15] Y. Ma, H. Wang, J. Xiong, J. Diao, and D. Ma, "Joint allocation on communication and computing resources for fog radio access networks," *IEEE Access*, vol. 8, pp. 108 310–108 323, 2020.

[16] N. Khumalo, O. Oyerinde, and L. Mfupe, "Reinforcement learning-based computation resource allocation scheme for 5g fog-radio access network," in *2020 Fifth International Conference on Fog and Mobile Edge Computing (FMEC)*, 2020, pp. 353–355.

[17] K. Li, M. Tao, and Z. Chen, "Exploiting computation replication for mobile edge computing: A fundamental computation-communication tradeoff study," *IEEE Transactions on Wireless Communications*, vol. 19, no. 7, pp. 4563–4578, 2020.

[18] T. Dang and M. Peng, "Joint radio communication, caching, and computing design for mobile virtual reality delivery in fog radio access networks," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 7, pp. 1594–1607, 2019.

[19] B. Ai, K. Guan, M. Rupp, T. Kurner, X. Cheng, X. Yin, Q. Wang, G. Ma, Y. Li, L. Xiong, and J. Ding, "Future railway services-oriented mobile communications network," *IEEE Communications Magazine*, vol. 53, no. 10, pp. 78–85, 2015.

[20] J. Wu and P. Fan, "A survey on high mobility wireless communications: Challenges, opportunities and solutions," *IEEE Access*, vol. 4, pp. 450–476, 2016.

[21] P. Muneer and S. M. Sameer, "Joint ML estimation of CFO and channel, and a low complexity turbo equalization technique for high mobility OFDMA uplinks," *IEEE Transactions on Wireless Communications*, vol. 14, no. 7, pp. 3642–3654, 2015.

[22] J. Wang, H. Zhu, and N. J. Gomes, "Distributed antenna systems for mobile communications in high speed trains," *IEEE Journal on Selected Areas in Communications*, vol. 30, no. 4, pp. 675–683, 2012.

[23] T. Li, K. Xiong, P. Fan, and K. B. Letaief, "Service-oriented power allocation for high-speed railway wireless communications," *IEEE Access*, vol. 5, pp. 8343–8356, 2017.

[24] C. Zhang, P. Fan, K. Xiong, and P. Fan, "Optimal power allocation with delay constraint for signal transmission from a moving train to base stations in high-speed railway scenarios," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 12, pp. 5775–5788, 2015.

[25] X. Liu and D. Qiao, "Location-fair beamforming for high speed railway communication systems," *IEEE Access*, vol. 6, pp. 28 632–28 642, 2018.

[26] R. Wang, R. Li, P. Wang, and E. Liu, "Analysis and optimization of caching in fog radio access networks," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 8, pp. 8279–8283, 2019.

[27] A. Simonetto, E. Dall'Anese, S. Paternain, G. Leus, and G. B. Giannakis, "Time-varying convex optimization: Time-structured algorithms and applications," *Proceedings of the IEEE*, pp. 1–17, 2020.

[28] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.1," http://cvxr.com/cvx, Mar. 2014.