#### DOI:10.20079/j.issn.1001-893x.220419009

开放科学(资源服务)标识码(OSID):

**引用格式:**陈昊宇,胡宏林.雾网络中基于统计分布的内容缓存与交付方案[J].电讯技术,2023,63(12):1902-1910.[CHEN H Y, HU H L. Content caching and delivery policy with frequency distribution in fog radio access networks[J]. Telecommunication Engineering, 2023, 63

(12):1902-1910.]

# 雾网络中基于统计分布的内容缓存与交付方案\*

## 陈昊宇1,2,胡宏林1

(1. 中国科学院上海高等研究院,上海 201210;2. 中国科学院大学 电子电气与通信工程学院,北京 100049)

摘 要:作为5G中的一种重要模型,雾无线接入网络(Fog Radio Access Network,F-RAN)通过设备 到设备通信和无线中继等技术获得了显著的性能增益,而边缘设备中合适的缓存则可以让内容缓存 用户(Caching Users,CUs)向内容请求用户(Requesting Users,RUs)直接发送缓存内容,有效减小前 传链路的负担和下载延迟。考虑一个F-RAN 模型下用户发出请求并获得交付的场景,将每个CU的 内容请求队列建模为独立的 M/D/1 模型,分析导出 CUs 缓存命中率和平均下载延迟关于内容缓存 与交付方案的表达式,证明 CUs 缓存命中率与内容统计分布之间的联系有助于实现前者的近似最 优解。针对在一段时间内的期望视角下建立的优化问题,提出了基于统计分布的算法并注意了执行 时的交付控制。仿真结果表明,相较于现有缓存策略,优化内容整体统计分布的方案能够最大化 CUs 缓存命中率,同时减小平均下载延迟。

关键词:5G;雾无线接入网络(F-RAN);雾计算;终端直通(D2D);内容缓存



中图分类号:TN915.07 文献标志码:A 文章编号:1001-893X(2023)12-1902-09

## Content Caching and Delivery Policy with Frequency Distribution in Fog Radio Access Networks

CHEN Haoyu<sup>1,2</sup>, HU Honglin<sup>1</sup>

(1. Shanghai Advanced Research Institute, Chinese Academy of Sciences, Shanghai 201210, China;
2. School of Electronic, Electrical and Communication Engineering,
Using the following of the second sec

University of Chinese Academy of Sciences, Beijing 100049, China)

**Abstract**: As an important model in 5G, the Fog Radio Access Network (F-RAN) achieves significant performance gains through technique among device-to-device (D2D) communications and wireless relays. Especially, caching appropriately in the edge devices allows content caching users (CUs) to send content to content requesting users (RUs) directly, which effectively reduces the burden of the fronthaul and download delay. The authors consider a scenario where users send content requests and get delivered under the F-RAN model. By modeling the content request queue at each CU as an independent M/D/1 queue model, the authors analyze and derive the cache hit probability of CUs and average download delay expressions under content caching and delivery policies. It is proved that the relationship between the cache hit probability of CUs and content frequency distribution can help realizing the approximate optimal solution of the former. On this basis, the authors establish an optimization problem under the expectation perspective over some time and propose a frequency distribution (FD)-based algorithm with delivery control when implementing to solve it. The simulation results show that compared with the existing caching policies, the policy optimizing the FD of all content can maximize the cache hit probability of CUs and reduce the average download delay.

Key words:5G; fog radio access network (F-RAN); fog computing; device-to-device (D2D); content caching

 <sup>\*</sup> 收稿日期:2022-04-19 修回日期:2022-06-10
 基金项目:国家重点研发计划(2020YFB1806606);国家自然科学基金资助项目(61801460)
 通信作者:胡宏林

## 0 引 言

随着智能移动设备(如智能手机、平板电脑、可穿 戴设备等)的发展和普及,用户对低延迟移动应用程 序和多媒体服务的需求大幅增加,这导致了数据流量 的爆炸性增长,对未来网络设计形成了挑战<sup>[1]</sup>。文献 [2]中预测:到 2023 年,连接到 IP 网络的设备数量将 超过全球人口的3倍,其中物联网(Internet of Things, IoT)连接将占全球在线设备和连接数量的一半,5G 设备和连接则占据10%以上。通过整合云计算和无 线接入网络, 云无线接入网络(Cloud Radio Access Network, C-RAN)为满足无线通信系统发展带来的巨 大算力与带宽需求提供了可能性<sup>[3]</sup>。然而,在 C-RAN 正面临着前传容量受限等一系列挑战的同 时<sup>[4]</sup>,集中式云计算能力的线性增速已逐渐与边缘数 据的需求拉开差距<sup>[5]</sup>。为了解决上述问题,进一步整 合云计算和雾计算则产生了新的雾无线接入网络 (Fog Radio Access Network, F-RAN)架构<sup>[6]</sup>。它将大 量存储、通信和控制功能转移,在网络边缘扩展了传 统的云计算范式<sup>[7]</sup>。在 F-RAN 中,缓存用户 (Caching Users, CUs)和雾无线接入节点(Fog Radio Access Point, F-AP)组成了容量有限的内容服务器, 直接向请求用户(Requesting Users, RUs)交付已被缓 存的内容。通过避免内容重复传输以及缩小与用户 间的物理距离,F-RAN 可以节省大量核心网络和回程 传输的资源消耗,有效减少用户的服务延迟,进而提 高用户体验质量(Quality of Experience, QoE)<sup>[1,8]</sup>。

内容分发网络(Content Delivery Network, CDN) 和信息中心网络(Information Center Network, ICN) 等架构对缓存技术的成功实践证明了在移动网络中 部署内容缓存的合理性与可行性<sup>[5]</sup>。然而,一方面 互联网与移动网络内容缓存机制的巨大差异使得传 统基于 CDN 的缓存技术在移动网络中无法直接应 用:另一方面,与一般移动边缘计算(Mobile Edge Computing, MEC)网络相比, F-RAN 能够凭借更强的 算力同时处理更多的数据<sup>[9]</sup>。因此,充分利用 F-AP 和 CUs 等边缘设备的缓存、计算和通信能力,方便 RUs 快速访问和检索内容,对于减轻前传、后传甚至 云端和主网的流量负载,有效节约网络带宽和能耗 具有重要意义<sup>[6,10]</sup>。作为提高 F-RAN 性能的关键 组件,如何缓存有使用价值的数据几乎决定了不同 终端的服务质量(Quality of Service, QoS)能否得到 保证,这导致内容缓存策略的优化成为了 F-RAN 边 缘缓存研究的热点[11-12]。

关于 F-AP 处的缓存策略,目前已经有丰富的研

究<sup>[13-15]</sup>,主要考虑了用户和缓存内容本身的各种特性。然而,这些 F-RAN 边缘缓存场景中并未考虑 D2D 技术的应用。事实上,CUs 可以与邻近的 RUs 共享缓存数据,利用 D2D 链路直接向 RUs 提供服务。该技术对于 F-RAN 在延迟和传输成本等内容交付方面的显著性能提升已经得到了研究和验证<sup>[11,16]</sup>。

上述关于 F-RAN 环境下 F-AP 与 D2D 设备合 作缓存的研究主要从各个节点的角度考虑它们与内 容集的具体映射关系,这相当程度上削弱了所有 D2D 设备形成的缓存空间的整体性。然而,若因系 统随机性而采取一段时间内的期望视角,在 CUs 具 备显著内容交付优势的前提下,直接优化 CUs 缓存 命中率是必要的。因此,本文主要以一个包含单个 F-AP 和多个 D2D 设备的 F-RAN 系统为研究对象, 通过数学证明探究有关建模与内容统计分布之间的 联系,从而导出优化目标并设计关于一般优化问题 的求解算法。针对边缘存储,本文提出了基于统计 分布的内容缓存与交付方案,实现了 CUs 缓存命中 率的近似全局最优解,同时验证了相应交付方案对 于优化平均下载延迟的有效性。

## 1 F-RAN 系统模型

## 1.1 网络模型

本文考虑一个包含多个 D2D 设备的典型 F-RAN 系统,如图 1 所示,其中共包含了一个功能完 善的云计算中心(Cloud Computing Center,CCC),一 个 F-AP,多个 CUs 和多个 RUs。在该网络结构下, F-AP 位于其覆盖区域中心,可以通过前传链路传输 并缓存来自 CCC 的内容,也可以通过无线链路直接 服务这一区域内的所有 CUs 和 RUs。



图 1 F-RAN 系统模型

假设 F-AP 覆盖区域的半径为  $R_{\rm F}$ , CUs 的最大 传输距离为  $R_{\rm C}$ , 且其分布服从密度为  $\gamma_{\rm C}$  的齐次泊 松点过程(Poisson Point Process, PPP)  $\Phi_{\rm C}$ ; RUs 的分 布服从密度为  $\gamma_{\rm R}$  的独立 PPP  $\Phi_{\rm R}^{[16]}$ , 因此 F-AP 和 CUs 的最大服务范围分别为  $A_{\rm F} = \pi R_{\rm F}^2$  和  $A_{\rm C} = \pi R_{\rm C}^2$ 。 对于不同范围下的 CUs 和 RUs 数量,本文采用 PPP 在指定区域中随机撒点个数的期望来进行表示,即 F-AP 可服务的 CUs 和 RUs 数量分别为  $K_{\rm CF} = \gamma_{\rm C} A_{\rm F}$ 和  $K_{\rm RF} = \gamma_{\rm R} A_{\rm F}$ , CUs 最大通信范围内的 CUs 和 RUs 数量分别为  $K_{\rm CC} = \gamma_{\rm C} A_{\rm C}$  和  $K_{\rm RC} = \gamma_{\rm R} A_{\rm C}$ 。

#### 1.2 内容缓存与用户请求模型

本文将内容的流行程度视为用户请求相应内容 的概率,并假定该信息可以由 CCC 端基于大量用户 请求数据定期获取。当缓存空间有限时,设 CCC 服 务器存储了一个大小为 N 的总体内容库  $C = \{c_1, c_2, \dots, c_n\}$  $c_3, \dots, c_N$ },其中 $c_i$ (*i*  $\in$  {1,2,...,N})表示所有内容 中第 i 流行的内容文件。随着 RUs 实际请求的发 生,该库及其包含的流行度信息将定期更新,下一步 F-AP 和 CUs 即可根据预先制定的缓存策略更新它 们存储的内容集,以便 RUs 在后续的一段时间内通 过不同链路直接取得已被缓存的内容。设 F-AP 本 地存储空间大小为 $N_{\rm F}$ ,执行最流行内容缓存(Most Popular Content Cache, MPC)策略<sup>[17]</sup>,缓存了 CCC 端最流行的 N<sub>F</sub> 个内容,则 F-AP 缓存的内容集为  $C_{\rm F} = \{c_1, c_2, c_3, \cdots, c_{N_{\rm F}}\}$ 。此外,假设单个 CU 本地存 储空间的大小为 $N_c$ ,其缓存的内容集 $C_c$ 为 $C_F$ 的一 个子集,由 F-AP 经无线链路直接传输获得。综上 所述,显然有 N>N<sub>F</sub>>N<sub>C</sub>。CUs 内容集的形成方法即 为本文主要讨论的内容缓存方案。

针对用户请求模型,本文规定 F-AP 覆盖区域 内所有 RUs 都首先将它们的内容请求发送给 F-AP。 不失一般性,令 CCC 服务器所存储的所有内容都有 相同的大小 *s*,且用户请求内容 *c<sub>i</sub>* 的概率服从 Zipf 分布,则 *p<sub>i</sub>* 有表达式如式(1)所示:

$$p_i = \frac{i^{-\beta}}{\sum\limits_{j=1}^{N} j^{-\beta}}$$
(1)

式中:β 是控制受欢迎内容集中程度的 Zipf 指数<sup>[18]</sup>,当β 的值逐渐增大时,越来越多的用户请求 会向越来越少的高流行度文件集中<sup>[19]</sup>。

#### 1.3 内容交付模型

当 F-AP 接收到 RUs 的内容请求以后,系统将 处理这些请求,然后交付相应内容,并于期间产生一 定的下载延迟。此处有关下载延迟的具体定义为 F-AP 接收到内容请求至请求内容完成交付之间所 经历的时间。假设每个 RU 请求到达都遵循独立泊 松过程,且有相同的强度(到达率)λ(请求数/秒),即 每个 RU 平均每秒发出 λ 个请求。对于不同 RU 的 不同请求,F-AP 接收后有 3 种服务模式可供选择。

## 1.3.1 CCC 模式

当 RUs 请求的内容  $c_i \notin C_F$  时, F-AP 就会通过前 传链路向 CCC 端转发该信息。收到指令之后, CCC 端则可通过高功率节点(High Power Nodes, HPNs)直 接向 RUs 提供此类流行度较低的突发性内容<sup>[6]</sup>。设 CCC 模式经历的下载延迟为固定常数  $t_o$ 

## 1.3.2 F-AP 模式

当 RUs 请求的内容  $c_i \in C_F$  时,若满足下列条件 之一,则请求的内容由 F-AP 直接交付给 RUs:

1) 在距离某个 RU 不超过 R<sub>c</sub> 的范围内没有 CUs 存在, 或存在若干个 CUs 但都未缓存内容 c<sub>i</sub> 时;

2) 在距离某个 RU 不超过 R<sub>c</sub> 的范围内,存在一个缓存了内容 c<sub>i</sub> 的 CU,但 F-AP 按照一定的概率判断其交付内容所需的时延过长,不适宜为该 RU 提供服务时。

本文假设 F-AP 有足够强的服务能力,可以同时 传输大量内容,因此对于随机到达的 RUs 请求,F-AP 不会形成内容请求队列。令 F-AP 到 RUs 的数据传 输速率为  $r_{\rm F}$ ,则该模式经历的下载延迟为  $c_i$  从 F-AP 传输到相应 RU 所需要的时间,即  $t_{\rm F}=s \cdot r_{\rm F}^{-1}$ 。

#### 1.3.3 CUs 模式

当 RUs 请求的内容  $c_i \in C_F$  时,若在距某个 RU 不超过  $R_c$  的范围内存在一个缓存了  $c_i$  的 CU,且 F-AP 按照一定的概率判断其适合提供服务,则  $c_i$  由 该 CU 交付给该 RU,具体分为两步:

1) F-AP 首先识别出适合的 CU,向其发送一条 包含内容请求信息的指令,通知它准备交付内容;

2) 收到指令后, CU 再通过 D2D 链路向发出请 求的 RU 交付内容。

本文假设 CUs 服务能力有限,最多同时传输一 个内容,因此对于 RUs 随机到达的请求,每个 CU 都 将形成一个内容请求队列。故该模式经历的下载延 迟包括内容请求在队列中的等待时间及其所需的服 务时间(此处服务时间由指令从 F-AP 传输到相应 CU 和  $c_i$  从指定 CU 传输到相应 RU 两部分时延组 成)。令 F-AP 发送指令部分的时延为  $t_{FC}$ , CUs 和 RUs 间的数据传输速率为  $r_c$ ,队列服务时间则可被 表示为  $t_c = t_{FC} + s \cdot r_c^{-1}$ 。本文考虑把每个 CU 的内 容请求队列都建模成独立的 M/D/1 排队模型,即请 求到达服从泊松过程且服务时间固定<sup>[20]</sup>。

综上所述,包括 F-AP 判断某个 CU 是否适合提供服务的概率的具体计算方法在内,完整的模式选择方法即为本文主要讨论的内容交付方案。

## 2 内容缓存与交付方案性能分析

#### 2.1 缓存命中率

本文选取了缓存命中率作为分析方案性能的一项指标,用 P<sub>HC</sub> 表示。它的定义为任意一个 RU 所请求的内容在其最大可通信范围内的某个 CU 处已被缓存的概率。因此,关于 CUs 缓存命中率的讨论与内容交付模型 CUs 模式下 F-AP 对相应 CU 是否适合提供服务的判断无关。同时,较大的 P<sub>HC</sub> 意味着 CUs 可以通过 D2D 链路满足更多的 RUs 内容请求。下面,本文将从内容统计分布的角度导出 P<sub>HC</sub> 的表达式。

假设缓存了内容  $c_i$  的 CUs 在所有 CUs 中的占 比为 $\delta_i$ , $\delta_i \in [0,1]$ ,且任意一个 RU 位于一个半径 为 $R_c$  的圆形区域的中心。故缓存了  $c_i$  的 CUs 在 F-AP 最大服务范围内的分布服从密度为 $\delta_i \gamma_c$  的齐次 PPP,根据二维泊松过程的性质,区域 $A_c$  内存在 m个此类 CUs 的概率如式(2)所示:

$$P_{i}(m) = \frac{(\delta_{i} \gamma_{\rm C} A_{\rm C})^{m} {\rm e}^{-\delta_{i} \gamma_{\rm C} A_{\rm C}}}{m!} {\rm o}$$
(2)

若令 m=0,则式(2)给出了区域  $A_c$  内不存在相应 CUs 的概率<sup>[21]</sup>,即缓存未命中概率,那么至少有一个缓存了  $c_i$  的 CUs 位于  $A_c$  内的概率为  $1-P_i(m=0)$ 。此外,由于 RUs 内容请求的概率分布一致且只与内容本身有关,当任意一个 RU 请求内容  $c_i$  的概率都为  $p_i$ 时,仅考虑  $c_i$  的缓存命中率如式(3)所示:

 $P_{HC,i} = p_i [1 - P_i(m=0)] = p_i (1 - e^{-\delta_i K_{CC}})$ 。 (3) 在放置 CUs 缓存的过程中,各个内容之间不涉 及相关性。结合式(3)及上述定义,若要计算 CUs 能够交付任意一个请求内容的概率,可以对所有  $c_i \in C_F$  的  $P_{HC,i}$  求和。CUs 缓存命中率如式(4) 所示:

$$P_{\rm HC} = \sum_{i=1}^{N_{\rm F}} P_{{\rm HC},i} = \sum_{i=1}^{N_{\rm F}} p_i (1 - e^{-\delta_i K_{\rm CC}})_{\circ}$$
(4)

同时易得, $\delta_i$ 关于 CUs 本地存储空间大小的约 束如式(5)所示:

$$\sum_{i=1}^{N_{\rm F}} \delta_i = N_{\rm Co} \tag{5}$$

## 2.2 平均下载延迟

本文还选取平均下载延迟作为分析内容缓存与

交付方案性能的另一项指标,用 *D*<sub>AC</sub> 表示。它的定 义为所有 RUs 长期经历的平均下载延迟。基于 F-RAN 的 3 种服务模式,本文需要对不同的内容交付 链路进行讨论。

首先是 CUs 的 M/D/1 排队模型,这里面包括两 个参数:平均请求到达率和平均服务率(队列服务 时间的倒数)。假设  $G_{CU}$  表示所有 CUs 组成的集 合, $G_{RU}$  表示所有 RUs 组成的集合。对于任意一个 RU $x \in G_{RU}$ ,到达任意一个 CU $y \in G_{CU}$  的内容请求遵 循强度为  $\lambda_y(x)$ 的泊松过程。参数  $\lambda_y(x)$ 的表达式 如式(6)所示:

$$\lambda_{y}(x) = \lambda \eta_{y} \sum_{c, \in C_{\pi}} g_{y,i}(x) p_{i} \circ$$
(6)

式中: $\eta_y \in (0,1]$ 为 CUy 的交付系数,它表示 F-AP 判断 CUy 适合为 RUs 提供服务的概率; $g_{y,i}(x)$ 为指 示函数。若  $g_{y,i}(x) = 1$ 则表示 CUy 是最接近 RUx 的缓存了内容  $c_i$ 的一个 CU,且 CUy 和 RUx 之间的 距离不大于  $R_c$ ;否则  $g_{y,i}(x) = 0$ 。

因此,对于  $CU_y$  处的泊松过程,它的内容请求 总平均到达率  $\lambda_y$  的表达式如式(7)所示:

$$\lambda_{y} = \sum_{x \in G_{\text{DV}}} \lambda_{y}(x)_{\circ} \tag{7}$$

根据 M/D/1 排队模型的性质<sup>[22]</sup>和 CUs 处的平均服务率, CUy 处的流量负载可以用系统中的平均内容请求总数  $L_y$  来表示。参数  $L_y$  的表达式如式 (8)所示:

$$L_{y} = \lambda_{y} t_{c} + \frac{(\lambda_{y} t_{c})^{2}}{2(1 - \lambda_{y} t_{c})^{\circ}}$$
(8)

其次是 F-AP 模式,它虽然并不涉及排队模型 或队列等待时间,但是不妨同样利用平均请求到达 率和平均服务率两个参数加以理解(服务时间部分 仍然成立)。对于 F-AP 处的泊松过程,内容请求到 达的总平均速率为其能直接交付的所有 RUs 请求 除去已通过 CUs 模式交付的部分。结合平均服务 率计算,F-AP 处的平均内容请求总数(即流量负 载)始终为下载延迟期间内容请求到达的总量,故 参数 L<sub>F</sub> 的表达式如式(9)所示:

$$L_{\rm F} = t_{\rm F} \left(\lambda K_{\rm RF} \sum_{i=1}^{N_{\rm F}} p_i - \sum_{y \in G_{\rm CU}} \lambda_y\right) \, . \tag{9}$$

最后,考虑到 CCC 模式的下载延迟,CCC 处平 均内容请求总数 L 的表达式同理可得如式(10) 所示:

$$L = t \left(1 - \sum_{i=1}^{N_{\rm F}} p_i\right) \lambda K_{\rm RF\,\circ} \tag{10}$$

综合上述 3 种模式下的流量负载和排队理论的 利特尔法则(Little's Law),上述模型有传输延迟方 程<sup>[23]</sup>:系统长期平均请求数=长期平均请求到达

· 1905 ·

率×系统长期平均请求等待时间。方程中有两项已 知,即3种服务系统的总平均内容请求数,以及所有 RUs请求的平均总到达率,而系统等待时间已经包 括仅CUs模式涉及的队列等待时间和3种模式都 涉及的服务时间两部分。因此,所有 RUs内容请求 经历的平均下载延迟如式(11)所示:

$$D_{\rm AC} = \frac{1}{\lambda K_{\rm RF}} \left( L + L_{\rm F} + \sum_{y \in G_{\rm CU}} L_y \right) \, \, (11)$$

## 3 基于统计分布的内容缓存与交付方案

本文对 F-RAN 中内容缓存与交付问题进行建模,导出了方案评价指标的量化表达式,目的是寻找关键影响因素以最大化 CUs 缓存命中率,并验证优化结果对降低平均下载延迟的有效性。为了求解该模型引出的一般优化问题,本文将提出基于统计分布的内容缓存与交付方案,尽可能地提升内容缓存与交付方案的性能和用户体验。

从 2.1 节中  $P_{HC}$  表达式导出的过程和结果可 知,影响其大小的因素只有各个  $c_i$  对应的  $\delta_i$ 。因此,为了最大化 CUs 缓存命中率,优化结果仅需考 虑内容的整体统计分布,而与  $C_F$  和  $G_{CU}$  之间的映射 关系无关。记缓存策略矢量  $\Delta = (\delta_1, \delta_2, \delta_3, \cdots, \delta_{N_F})$ 。本文所提方案的思路如下:通过计算求出最 优  $\Delta$  使  $P_{HC}$  达到最大,再按定义将不同内容随机放 置于一定数量的 CUs 本地存储空间以形成所有  $C_c$ 。 基于统计分布的内容缓存与交付方案最优缓存命中 率问题的最小化表述如式(12)所示:

$$\min_{\delta_i} -P_{\rm HC} \qquad (12a)$$

$$\delta_i \in [0, 1], i = 1, 2, \cdots, N_{\rm F},$$
 (12b)

$$\sum_{i=1}^{N_{\rm F}} \delta_i = N_{\rm C}, \qquad (12c)$$

$$t_{\rm C}\boldsymbol{\lambda}_{\rm y} = 1 - \boldsymbol{\xi}, \boldsymbol{y} \in \boldsymbol{G}_{\rm CU\,\circ} \tag{12d}$$

式(12b)表示内容分布的统计结果是 0~1 之间 的比例值;式(12c)表示实际缓存要完全占用 CUs 的本地存储空间;式(12d)中 $\xi$ 为一较小正数,用于 使得 M/D/1 排队模型的服务强度(即服务能力利用 率:平均到达率与平均服务率的比值)小于 1,以满 足排队系统的稳定条件,保证不会形成无限队列。 此处假定 F-AP 已经掌握了  $N_c$ , $K_{cc}$ , $\beta$ , $\lambda$ , $t_c$ 等必要 信息。记式(12a)中目标函数– $P_{HC}$ = $W(\Delta)$ ,且一般 而言  $K_{CF}$ 足够大,足以将 $\delta_i$ 视作连续变量,故容易求 得函数 W的一阶梯度和 Hessian 矩阵如式(13)和式 (14)所示:

$$\nabla_{\Delta} W = -\begin{pmatrix} p_{1} K_{CC} e^{-\delta_{1} K_{CC}} \\ p_{2} K_{CC} e^{-\delta_{2} K_{CC}} \\ \vdots \\ p_{N_{F}} K_{CC} e^{-\delta_{N_{F}} K_{CC}} \end{pmatrix}, \quad (13)$$

$$\nabla_{\Delta}^{2} W = \begin{pmatrix} \frac{p_{1} K_{CC}^{2}}{e^{\delta_{1} K_{CC}}} & 0 & \cdots & 0 \\ 0 & \frac{p_{2} K_{CC}^{2}}{e^{\delta_{2} K_{CC}}} & \cdots & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & \frac{p_{N_{F}} K_{CC}^{2}}{e^{\delta_{N_{F}} K_{CC}}} \end{pmatrix}$$

$$(14)$$

由于 $\nabla_{4}^{2}$ W为正定矩阵,且式(12b)和式(12c)均 为关于  $\Delta$ 的仿射函数,也是凸函数,故式(12)符合 凸问题的标准形式<sup>[24]</sup>。对于此类具有等式和不等 式约束的一般优化问题,本文采用带 KKT(Karush-Kuhn-Tucker)条件的拉格朗日乘数法(Lagrange Multipliers)求解,求得的  $\Delta$  即为全局最优解。定义 拉格朗日函数 Q的表达式如式(15)所示:

$$Q(\boldsymbol{\Delta},\boldsymbol{\mu},\boldsymbol{\tau},\boldsymbol{\lambda}) = W(\boldsymbol{\Delta}) - \boldsymbol{\mu} \boldsymbol{\Delta}^{\mathrm{T}} - \boldsymbol{\mu} (1 - \boldsymbol{\Delta}^{\mathrm{T}}) +$$

t

$$\Theta(\|\boldsymbol{\Delta}\|_{1} - N_{\rm C}) \circ$$
(15)

式中:不等式的约束矢量 $\mu = (\mu_1, \mu_2, \mu_3, \dots, \mu_{N_F})$ 和  $\tau = (\tau_1, \tau_2, \tau_3, \dots, \tau_{N_F})$ 各个分量均为非负,分别表示 限制 $\Delta$ 中各个分量不小于0和不大于1所对应的约 束系数;非零参数 $\theta$ 则为式(12c)对应的等式约束 系数。

进一步地,上述一般优化问题还需要的求解条 件如式(16)所示:

$$\begin{cases} \| \nabla_{\Delta} Q \|_{1} = 0 \\ \| \boldsymbol{\mu} \circ \boldsymbol{\Delta} \|_{1} = 0 \\ \| \boldsymbol{\tau} \circ (\boldsymbol{I}^{\mathrm{T}} - \boldsymbol{\Delta}) \|_{1} = 0 \end{cases}$$
(16)

式中:运算符"。"表示矩阵的哈达玛(Hadamard)积。

根据 Zipf 分布的性质,  $p_i$  的大小将随着 i 的增 大逐渐减小,因此,不难预计  $\delta_i$  的变化也遵循相似 态势,即  $\Delta$  中需要在流行度较高的内容中出现尽可 能多的 1,而在流行度较低的内容中出现尽可能少 的 0。本文由此对  $\delta_i$  进行分段讨论。

假设当  $1 \le i \le N_L$  时, $\delta_i = 1$ ,显然  $\mu_i = 0$ ,故此时 关于  $\tau_i$  的约束如式(17)所示:

$$\tau_i = p_i K_{\rm CC} e^{-\kappa_{\rm CC}} - \theta \ge 0_{\circ} \tag{17}$$

假设当  $N_{\rm L}$ +1 $\leq i \leq N_{\rm H}$  时,0< $\delta_i$ <1,显然此时 $\mu_i$ =  $\tau_i$ =0,故关于 $\delta_i$ 的约束如式(18)所示:

· 1906 ·

$$\theta = p_i K_{\rm CC} e^{-\delta_i K_{\rm CC}}$$
(18)  
假设当  $N_{\rm H} + 1 \le i \le N_{\rm F}$  时, $\delta_i = 0$ ,显然此时有 $\tau_i =$ 

0,故关于 $\mu_i$ 的约束如式(19)所示:

$$\boldsymbol{\mu}_i = \boldsymbol{\theta} - p_i \boldsymbol{K}_{\rm CC} \ge \boldsymbol{0}_{\circ} \tag{19}$$

综上,关于 $\delta_i$ 的求解方程组如式(20)所示:

$$\begin{cases} \sum_{i=N_{\rm L}+1}^{N_{\rm H}} \delta_i = N_{\rm C} - N_{\rm L} \\ p_{N_{\rm H}+1} \leq p_{N_{\rm L}} e^{-\kappa_{\rm CC}} & \circ \\ \prod_{i=N_{\rm L}+1}^{N_{\rm H}} p_i e^{-\delta_i \kappa_{\rm CC}} \leq (p_{N_{\rm L}} e^{-\kappa_{\rm CC}})^{N_{\rm H}-N_{\rm L}} \end{cases}$$
(20)

结合式(18),当 $N_{L}$ +1 $\leq i \leq \lfloor N_{L}e^{\frac{\Lambda_{CC}}{\beta}} \rfloor$ 时,可以求 得 $\delta_{i}$ 的表达式如式(21)所示:

$$\delta_{i} = \frac{\beta \sum_{j=N_{\mathrm{L}}+1}^{\lfloor N_{\mathrm{L}} \mathrm{e}^{\frac{-K_{\mathrm{CC}}}{\beta}} \rfloor} \ln\left(\frac{j}{i}\right) + N_{\mathrm{C}} - N_{\mathrm{L}}}{\lfloor N_{\mathrm{L}} \mathrm{e}^{\frac{K_{\mathrm{CC}}}{\beta}} \rfloor - N_{\mathrm{L}}} \circ$$
(21)

式中:变量 N<sub>L</sub> 由不等式确定,具体计算方式为寻找 满足式(22)的最大正整数。

$$\frac{\beta}{K_{\rm CC}} \sum_{i=N_{\rm L}+1}^{\lfloor N_{\rm L} e^{\frac{CC}{\beta}}} \ln\left(\frac{i}{N_{\rm L}}\right) \ge \lfloor N_{\rm L} e^{\frac{K_{\rm CC}}{\beta}} \rfloor - N_{\rm Co} \qquad (22)$$

另外需要说明的是,在不同内容关于 CUs 本地 存储空间的随机放置完成以后,上述算法有必要通 过交付系数保证自身的成功执行。计算后得到各个  $CUy \in G_{CU}$ 的交付系数  $\eta_x$  如式(23)所示:

$$\eta_{y} = \min\left\{1, \frac{1-\xi}{\lambda t_{C}} \left(\sum_{x \in G_{RU}^{c_{i}} \in C_{F}} g_{y,i}(x) p_{i}\right)^{-1}\right\} \circ (23)$$

至此,本文提出的基于统计分布的内容缓存与 交付方案的具体描述如下:

输入:单个 CU 的本地存储空间大小  $N_c$ , CUs 最大通信 范围内的 CUs 数量  $K_{cc}$ , Zipf 指数  $\beta$ , 每个 RU 请求遵循的相 同到达率  $\lambda$ , CUs 模式的队列服务时间  $t_c$ 

输出:缓存策略矢量  $\Delta$ ,各个 CUy  $\in G_{CU}$  处的交付系数  $\eta_{y}$ 

1 初始化  $N_{\rm L} = 1, N_{\rm L}' = \lceil N_{\rm F} e^{\frac{K_{\rm CC}}{\beta}} \rceil //确 \ge N_{\rm L}$ 的取值范围; 2 While  $N_{\rm L}' - N_{\rm L} > 1$   $\hat{N}_{\rm L} = \lceil \frac{N_{\rm L} + N_{\rm L}'}{2} \rceil$ If  $\hat{N}_{\rm L}$  满足式(22)  $N_{\rm L} = \hat{N}_{\rm L}$ Else  $N_{\rm L}' = \hat{N}_{\rm L}$ End if End while//用二分法求解  $N_{\rm L}$ ;

3 For 
$$i=1, N_{\rm L}$$
  
 $\delta_i = 1$   
End for//对  $\Delta$  赋值;  
4 For  $i=N_{\rm L}+1, \lfloor N_{\rm L} e^{\frac{K_{\rm CC}}{\beta}} \rfloor$   
按式(21)计算  $\delta_i$   
End for//对  $\Delta$  赋值;  
5 For  $i=\lceil N_{\rm L} e^{\frac{K_{\rm CC}}{\beta}} \rceil, N_{\rm F}$   
 $\delta_i = 0$   
End for//对  $\Delta$  赋值;  
6 For  $y \in G_{\rm CU}$   
按式(23)计算  $\eta_y$   
End for//计算所有的  $\eta_y$ ,共 $K_{\rm CF}$  个。

算法执行完毕以后,则需要进行内容  $c_i$  的具体 放置( $i \in \{1, 2, \dots, N_F\}$ ):考虑内容  $c_i$ ,计算出  $K_{CF}\delta_i$ 的数值并取整,代表内容  $c_i$  应当被放置在 $\lfloor K_{CF}\delta_i + \frac{1}{2} \rfloor$ 个 CUs 中,然后随机选择相应数量的 CUs,令它 们缓存指定内容即可(这里实际上取了  $K_{CF}\delta_i$  的近 似值,故得到的是近似全局最优解)。倘若 CUs 本 地存储空间已被完全利用且各个 CUy  $\in G_{CU}$  的  $C_C$ 内部没有重复元素的条件得以保证,所有内容  $c_i \in C_F$ 的放置完成之后,形成的 CUs 内容集即为本文讨 论的内容缓存方案的执行结果。当所有 RUs 都向 F-AP 发送内容请求时,F-RAN 系统最终执行的内 容交付方案的完整模式选择流程如图 2 所示。



图 2 内容交付方案模式选择流程

### 4 仿真实验与结果分析

## 4.1 参数设置

针对基于统计分布的内容缓存与交付方案的仿 真实验使用 Matlab 完成,通过执行具体 F-RAN 网络 中的重复随机试验取得 1 000 组数据,依据其均值 评估本文所提算法。仿真实验中使用的详细初始参 数如表 1 所示。

私I 内关入型生中乡奴
-------------

参数	值
独立 PPP $\Phi_{R}$ 的密度 $\gamma_{R}/(\gamma/m^{2})$	10 <sup>-3</sup>
每个 RU 请求遵循的相同到达率 $\lambda/(\gamma/s^1)$	3.5
F-AP 到 RUs 的数据传输速率 r <sub>F</sub> /(Mb/s)	$2^{[26]}$
CUs 和 RUs 间的数据传输速率 r <sub>c</sub> /(Mb/s)	20
CCC 服务器存储所有内容的大小 s/Mb	5
F-AP 发送指令部分的时延 t <sub>FC</sub> /ms	25
CCC 模式经历的固定下载延迟 t/ms	5 000
F-AP 覆盖区域的半径 R <sub>F</sub> /m	150
CUs 的最大传输距离 R <sub>c</sub> /m	30[16]
CCC 服务器存储的总体内容库的大小 N	4 000
F-AP 本地存储空间的大小 $N_{\rm F}$	1 000
单个 CU 本地存储空间的大小 $N_{\rm c}$	80
限制服务强度的较小正数 ξ	0.1

#### 4.2 实验结果分析

本文选取了文献[12]中介绍的两种内容缓存 与交付方案进行对比实验,通过仿真程序分别实现 了基于概率(Probability-based,PB)的缓存策略、基 于边缘缓存用户分类(Edge Caching Users Classification,EUC)的缓存策略以及基于统计分布 (Frequency Distribution,FD)的内容缓存与交付方 案,目的是验证本文所提方案分别在优化 CUs 缓存 命中率和平均下载延迟方面的优越性和有效性。

图 3 展示了不同内容缓存与交付方案下, CUs 缓存命中率随 CUs 分布所服从的齐次 PPP $\Phi_c$  密度 变化的情况。此处设置变量  $\gamma_c \in [2 \times 10^{-4}, 10 \times 10^{-4}]$ , Zipf 指数  $\beta = 0.5$ 。由图 3 可以看出, CUs 缓 存命中率会随着 CUs 分布密度的增加而不断增加, 但鉴于取整运算等原因,导致图像出现了数值跳变。 不难理解,当 $\gamma_c$ 数值较小时, CUs 数量少, 因而形成 的整体缓存空间较小,能够缓存的内容也很少, 这使 得  $P_{HC}$ 的数值偏低。然而当 $\gamma_c$ 的数值增大时,除了 CUs 形成的整体缓存空间增大外, 单个 RU 周边存 在 CUs 的概率也在增加。此时,本文提出的 FD 方 案尽可能地考虑了 CUs 可提供服务的多样性,避免 了单个 RU 周边存在 CUs 的数量变多时缓存内容重 复而可能出现的无效缓存现象(内容已被缓存但没 有机会进行交付)。因此,相较于其他内容缓存与 交付方案,FD 方案在 CUs 缓存命中率方面具有显 著优势。



图 4 展示了不同内容缓存与交付方案下,CUs 缓存命中率随 Zipf 指数变化的情况。该部分的其 他参数设置还包括变量  $\beta \in [0.15, 0.65]$  以及  $\gamma_c = 5 \times 10^{-4}$ 。根据 Zipf 分布的性质,当 $\beta$ 增大时,越来越 少的高流行度文件将有越来越高的概率被 RUs 请 求。此时,即使缓存内容保持不变,CUs 缓存命中率 也会提高。仿真结果表明,这一过程中 FD 方案相 较于其他缓存策略始终保持了显著的性能优势,这 是因为它有效减少了 RUs 附近的 CUs 因缓存内容 相同而造成的缓存空间冗余。从期望的视角观察, F-RAN 网络凭借 FD 方案的内容丰富性和它本身被 请求概率增大两方面因素获得了更高的 CUs 缓存 命中率。更重要的是,本文提出的算法可以达到 CUs 缓存命中率最大化问题的近似全局最优解。



图 4 不同 Zipf 指数对 CUs 缓存命中率的影响

图 5 展示了不同内容缓存与交付方案下,平均 下载延迟随 Zipf 指数变化的情况,其中补充参数的 设置与图4部分相同。这项仿真在实验设计上具备 的主要差异在于,由于 PB 和 EUC 缓存策略本身并 不包含交付系数的概念,在假定各个 CU 的交付系 数都为1的情况下,就需要舍弃会形成无限队列的 部分数据,而仅保留使平均下载延迟计算有效的数 据。图5的结果显示,当CUs缓存的内容变得更加 流行以后,更多的 RUs 请求将会通过 CUs 模式交 付,D2D 链路则会始终保持内容交付方面的性能优 势。这意味着提高 CUs 缓存命中率可以通过让更 多的流量负载向 CUs 转移从而减小平均下载延迟, 证明了本文所提算法的有效性和执行边缘缓存的价 值。除此之外,就FD方案而言,即使有队列等待时 间,相比其他两种缓存策略,它用内容缓存方案中的 随机性和内容交付方案中的交付系数对其进行控制 与均衡,最终既没有浪费缓存空间,也实现了更优越 的性能。仿真实验的设计还表明,限制单个 CU 处 的内容请求到达率是必要的,FD 方案利用交付系数 避免了无效数据,增强了算法的鲁棒性。



## 5 结束语

本文主要针对雾网络中的缓存进行性能优化, 提出了基于统计分布的内容缓存与交付方案,利用 带 KKT 条件的拉格朗日乘数法最大化了 CUs 缓存 命中率,并且通过仿真实验分析了各个变量对有关 性能指标造成的影响,验证了优化 CUs 缓存命中率 与降低平均下载延迟的相关性。结果表明,本文设 计的算法有效,实现了预期效果,达到了优化目标的 近似全局最优解,可以始终保持性能优势。因此,在 设计内容缓存与交付方案时,直接采取期望视角是 必要的,本文从节点分布和用户请求等建模开始,到 算法中考虑的内容在概率和统计分布方面的特征都 体现了这一点。值得注意的是,为了保证内容缓存 与交付方案的成功执行,需要将交付控制纳入考量, 事实上它是方案性能的重要影响因素。

下一步将考虑流行度预测、用户移动性等对模型更精确的随机表示和更多的优化方法,探讨此类研究思路与本文的差异。

## 参考文献:

- JIANG W, FENG G, QIN S, et al. Multi-agent reinforcement learning based cooperative content caching for mobile edge networks [J]. IEEE Access, 2019, 7: 61856-61867.
- [2] CISCO. CISCO annual Internet report (2018-2023) white paper[EB/OL]. (2020-03-09) [2022-03-19]. https://www.cisco.com/c/en/us/solutions/collateral/ executive-perspectives/annual-internet-report/whitepaper-c11-741490.pdf.
- [3] PENG M, LI Y, ZHAO Z, et al. System architecture and key technologies for 5G heterogeneous cloud radio access networks [J]. IEEE Network, 2015, 29(2):6-14.
- PENG M, WANG C, LAU V, et al. Fronthaul-constrained cloud radio access networks: insights and challenges [J].
   IEEE Wireless Communications, 2015, 22(2):152-160.
- [5] 李碧瑶.边缘网络下的计算卸载和边缘缓存方法研 究[D].北京:北京邮电大学,2020.
- [6] PENG M, YAN S, ZHANG K, et al. Fog-computingbased radio access networks:issues and challenges[J]. IEEE Network,2016,30(4):46-53.
- [7] BONOMI F, MILITO R, ZHU J, et al. Fog computing and its role in the Internet of Things[C]//Proceedings of the First Edition of the MCC Workshop on Mobile Cloud Computing. New York: Association for Computing Machinery, 2012:13-16.
- [8] WANG X, CHEN M, TALEB T, et al. Cache in the air: exploiting content caching and delivery techniques for 5G systems
   [J]. IEEE Communications Magazine, 2014, 52 (2):131-139.
- [9] LI Z, CHEN J, ZHANG Z. Socially aware caching in D2D enabled fog radio access networks [J]. IEEE Access, 2019, 7:84293-84303.
- [10] LI Q, NIU H, PAPATHANASSIOU A, et al. Edge cloud and underlay networks: empowering 5G cell-less wireless architecture [C]//Proceedings of 2014 20th European Wireless Conference. Barcelona: IEEE, 2014:1-6.
- [11] WANG X, LENG S, YANG K. Social-aware edge caching in fog radio access networks [J]. IEEE Access, 2017, 5: 8492-8501.
- [12] HUA H, CHU X. Content caching policy with edge caching user classification in fog radio access networks

[ C ]//Proceedings of 2021 IEEE Wireless Communications and Networking Conference. Nanjing: IEEE,2021:1-7.

- [13] LIU T, LI J, KIM B, et al. Distributed file allocation using matching game in mobile fog-caching service network[C]//Proceedings of 2018 IEEE Conference on Computer Communications Workshops. Honolulu: IEEE, 2018:499-504.
- [14] ZHOU X, LIU Z, GUO M, et al. SACC: a size adaptive content caching algorithm in fog/edge computing using deep reinforcement learning [J]. IEEE Transactions on Emerging Topics in Computing, 2022, 10(4):1810–1820.
- [15] YAN S, JIAO M, ZHOU Y, et al. Machine-learning approach for user association and content placement in fog radio access networks [J]. IEEE Internet of Things Journal, 2020, 7(10):9413-9425.
- [16] JIANG F, ZHANG X, SUN C. A D2D-enabled cooperative caching strategy for fog radio access networks
   [ C ]//Proceedings of 2020 IEEE 31st Annual International Symposium on Personal, Indoor and Mobile Radio Communications. London; IEEE, 2020; 1–6.
- [17] CHEN Z, LEE J, QUEK T, et al. Cooperative caching and transmission design in cluster-centric small cell networks
   [J]. IEEE Transactions on Wireless Communications, 2017, 16(5):3401-3415.
- [18] YU K, MA Z, NI R, et al. A caching strategy based on many-to-many matching game in D2D networks [J].

Tsinghua Science and Technology, 2021, 26(6):857-868.

- BRESLAU L, CAO P, FAN L, et al. Web caching and Zipf-like distributions: evidence and implications [C]// Proceedings of Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies. New York: IEEE, 1999:126-134.
- [20] KINGMAN J. The first Erlang century—and the next [J]. Queueing Systems, 2009, 63(1/2/3/4):3-12.
- [21] ANDREWS J, BACCELLI F, GANTI R. A tractable approach to coverage and rate in cellular networks [J]. IEEE Transactions on Communications, 2011, 59 (11): 3122-3134.
- [22] CAHN R. Wide area network design: concepts and tools for optimization [M]. San Francisco: Morgan Kaufmann Publishers Inc., 1998.
- [23] KLEINROCK L. Queueing systems [ M ]. Hoboken: Wiley-Interscience, 1975.
- [24] BOYD S, VANDENBERGHE L. Convex optimization [M]. New York:Cambridge University Press, 2004.

## 作者简介:

**陈昊宇** 男,1998年生于安徽合肥,2020年获工学学士 学位,现为硕士研究生,主要研究方向为无线通信、边缘 计算。

**胡宏林** 男,1975年生于安徽安庆,2004年获工学博士 学位,现为研究员、博士生导师,主要研究方向为移动通信技 术、脑机通信技术。