

面向元宇宙移动增强现实应用的异质内容 主动缓存与个性化交付机制

徐思雅, 付琦梦, 郭少勇

(北京邮电大学网络与交换技术全国重点实验室, 北京 100876)

摘要: 针对移动增强现实 (MAR) 应用异质数据内容传输和响应时延的问题, 提出了面向元宇宙移动增强现实应用的异质内容主动缓存与个性化交付机制, 首先综合考虑了用户特征和边缘节点服务能力, 提出了用户行为和资源感知的边缘协作服务域构建方法; 进而, 基于协作服务域设计基于存储空间划分和用户偏好预测的异质内容预缓存机制; 特别地, 针对前景内容引入了个性化推荐机制, 并设计了前景内容的差异化策略。仿真结果表明, 所提机制在缓存命中率、前景内容和异质内容的平均响应时延方面均优于 NCPCR 策略、CCS-AGP 策略和 AIEC-RSC 策略。

关键词: 元宇宙; 移动增强现实; 异质内容; 边缘缓存; 内容推荐

中图分类号: TP393

文献标志码: A

DOI: 10.11959/j.issn.1000-436x.2024187

Proactive caching and personalized delivery mechanism of heterogeneous content for MAR applications in the Metaverse

XU Siya, FU Qimeng, GUO Shaoyong

State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China

Abstract: Addressing the issues of heterogeneous data content transmission and response latency in mobile augmented reality (MAR) applications, a proactive caching and personalized delivery mechanism of heterogeneous content for MAR applications in the metaverse was proposed. Firstly, considering user characteristics and edge node service capabilities, a user behavior and resource aware edge collaboration service domain (ECSD) construction method was proposed. Then, based on the ECSD, heterogeneous content pre-caching mechanisms based on storage space division and user preference prediction were designed. Specially, by introducing the personalized recommendation mechanism of foreground content, a differentiated delivery strategy of foreground content was designed. Simulation results show that the proposed mechanism outperforms NCPCR, CCS-AGP, and AIEC-RSC strategy in terms of cache hit rate, average response latency for foreground content and heterogeneous content.

Keywords: Metaverse, mobile augmented reality, heterogeneous content, edge cache, content recommendation

0 引言

元宇宙被誉为是“下一代互联网”, 它是运用数字技术构建的、由现实相互映射或超越于现实世界且可与现实世界交互的虚拟世界^[1]。其中, 移动

增强现实 (MAR, mobile augmented reality)^[2]是支撑元宇宙融合虚拟场景和物理现实的关键技术, 广泛应用于文化旅游、在线教育和医疗培训等领域, 为元宇宙用户提供移动沉浸式交互体验。

收稿日期: 2024-08-02; 修回日期: 2024-10-11

通信作者: 徐思雅, xusiyaxsy@bupt.edu.cn

基金项目: 国家自然科学基金资助项目 (No.62201074); 北京邮电大学基本科研业务费基金资助项目 (No.2023ZCTH11)

Foundation Items: The National Natural Science Foundation of China (No.62201074), The Basic Research Fund From Beijing University of Posts and Telecommunications (No.2023ZCTH11)

在一个 MAR 场景中, 用户往往会同时请求异质内容数据, 包含静态或缓慢变化的背景内容, 如虚拟物品和场景的三维模型及其渲染资源, 以及可交互且需要频繁更新同步的前景内容, 如用户虚拟形象、用户互动数据和用户生成内容等^[3-4]。例如在 AR 游戏中, 通过摄像头和 AR 算法, 虚拟角色可以出现在现实场景中, 玩家通过摇晃手机就能与虚拟角色进行交互。针对大小和更新频率都具有显著差异的前景内容和背景内容, 使用传统的云服务器下载模式处理海量的异质内容会造成较低的数据访问速度。MAR 应用对时延非常敏感, 如果异质内容的响应时延过高, 就会引起虚拟空间匹配错位, 造成用户眩晕恶心, 破坏沉浸式体验。例如, 大型 MAR 应用的时延容忍一般不超过 50 ms^[5-7], 此外, MAR 应用是一个大型多人实时在线系统, 每个时隙产生的用户内容请求数量庞大且多变, 仅利用边缘实时缓存难以适应 MAR 用户复杂多变的决策需求^[8-10]。因此, 为了充分利用缓存资源, 需要构建有效的数据存取和缓存策略, 将前景内容和背景内容的数据及其渲染资源部署在边缘服务器上以适应 MAR 用户的移动性。同时, 利用存储空间有限的终端缓存前景内容相关数据, 可以进一步降低 MAR 应用的响应时延。并且, 在边缘缓存网络中, 缓存的有效性受到用户偏好的影响, 可以通过引入推荐机制重塑用户请求行为, 并采用终端的端到端 (D2D, device to device) 数据共享及交付模式, 提升用户体验^[11-13]。因此, 研究高效的边缘协作方法、精准的内容预缓存策略和合理的内容推荐机制, 从而提高资源利用率、降低服务响应时延, 成为业界的研究热点。

近年来, 针对 MAR 应用中的缓存和交付策略已有较多研究成果。其中, 针对边缘网络的内容缓存和推荐策略, 文献[14]面向车联网环境, 提出了一种基于服务区域划分、服务器分组和存储空间划分的协同缓存策略, 有效降低了通用数据的获取时延, 提高了车载应用的服务质量。该策略基于请求数据的相似性和边缘节点的地理位置进行协作域的划分, 使得边缘节点能够高效协作并合理利用存储空间。文献[15]提出了一种面向高速移动业务的人工智能和移动边缘计算集成的服务框架, 该框架根据边缘节点的服务能力动态构建边缘协作服务域, 并设计了用户行为感知的服务组件预缓存策略及计

算任务卸载方法, 提高了服务响应速度、资源利用率和服务可靠性。上述方案可为 MAR 应用的高效交付提供参考, 但没有考虑 MAR 应用中多源动态的异质内容、复杂多变的用户偏好和时延敏感的混合任务。

文献[16]在移动边缘云计算网络中集成了边缘缓存和推荐系统, 提出了一种联合深度强化学习和联邦学习的去中心化缓存替换算法, 以最小化综合系统成本。然而, 该算法没有考虑对终端已缓存内容的二次利用, 边缘缓存利用率有待提高。文献[17]考虑到用户的时延要求、D2D 内容交付的激励机制和链路保护机制, 提出了一种蜂窝网络和 D2D 网络配合的内容推荐和交付策略, 以提高运营商的经济效益。然而, 该方法仅考虑单一边缘节点的实时缓存与交付模式, 忽略了边缘节点间的内容共享。此外, 以上方法均没有考虑用户高速移动引起的交付中断风险, 无法保障服务的连续性和可靠性。文献[18]针对用户终端和基站之间的信息不对称问题, 提出了一种基于契约理论的激励机制, 以鼓励终端缓存 AR 应用程序并通过 D2D 通信提供缓存内容。在该机制中, 基于内容流行度、请求率等不对称信息, 设计了在个体合理性和激励相容性约束下的终端内容缓存方案, 从而提高了终端缓存的经济效益。但是, 该机制忽略了边缘服务器的计算和缓存能力, 仅依靠 D2D 方式进行内容交付, 造成了较大的 D2D 传输成本。

为解决以上问题, 本文设计了面向元宇宙移动增强现实应用的异质内容主动缓存与个性化交付机制。本文的贡献主要包括以下 3 个方面。

1) 为了提高边缘节点的协作服务能力, 本文提出了用户行为和资源感知的边缘协作服务域构建方法, 综合考虑用户的空间分布和服务内容需求、服务质量要求等用户偏好特征信息, 从全局视角层次划分最优的用户服务簇; 进而, 根据用户所属的服务簇以及边缘节点的服务能力和覆盖范围, 筛选服务能力强且用户分布密集的边缘节点, 构建边缘协作服务域, 允许偏好相似的用户共享边缘协作服务域内的缓存资源, 改善了网络的服务时延和资源效用。

2) 针对 MAR 应用中的前景内容和背景内容, 本文创新性地提出了基于存储空间划分和用户偏好预测的异质内容预缓存机制。首先, 针对前景和背

景内容在存储位置和所需存储空间差异,设计了基于域内-域间流行度的背景内容综合流行度预测方法和前景内容多因子分域流行度预测方法,并且对边缘服务器的存储空间进行精细化分区,分别存储流行缓存内容、个性化缓存内容和实时缓存内容,以均衡利用缓存空间;进而,基于构建好的边缘协作服务域,提出了异质内容的域首节点冗余预缓存策略和域成员节点精准预缓存策略,分别从全局统筹和局部优化角度提高边缘网络的缓存利用率,提高了网络对用户请求的响应速度。

3) 特别地,针对用户复杂多变的内容偏好和高速移动引起的交付中断,本文创新性地提出了前景内容个性化推荐机制和差异化交付策略。首先,综合考虑协作域中边缘节点及用户终端已缓存的前景内容的请求概率、缓存位置和传输时延,设计了前景内容个性化推荐机制,引导用户体验已有内容服务,对终端缓存进行二次利用;进而,根据用户移动性、链路保持时间等,设计了前景内容的差异化交付策略,灵活选择 D2D 直接交

付、域内边缘节点主动交付和域间边缘节点协作交付方式,进一步提高了边缘缓存网络的缓存利用率和服务质量。

1 系统模型

1.1 网络架构模型

本文采用云-边-端的3层网络架构模型,由云服务器、边缘服务器和请求MAR应用的用户终端3个部分构成,如图1所示。假设云服务器 e_0 存储了元宇宙中MAR应用所需的全部内容,总内容库用 $\mathcal{F} = \{f_i\}_{i=1}^F$ 表示,每个内容的大小不相同,表示为 s_i 。边缘节点也称为边缘服务器或基站, $e = \{e_i\}_{i=1}^E$ 代表系统中的边缘节点集合,集合大小为 E 。边缘节点之间通过光纤链路连通,可以根据用户请求进行协作交付,或从云服务器中请求缓存内容。 u 表示用户终端,使用 $u = \{u_i\}_{i=1}^U$ 来表示移动中的用户终端集合,大小为 U 。D2D技术允许通信范围内的用户终端设备之间可以不经由基站直接进行半双工数据通信。本文假设每个用户独立且随机

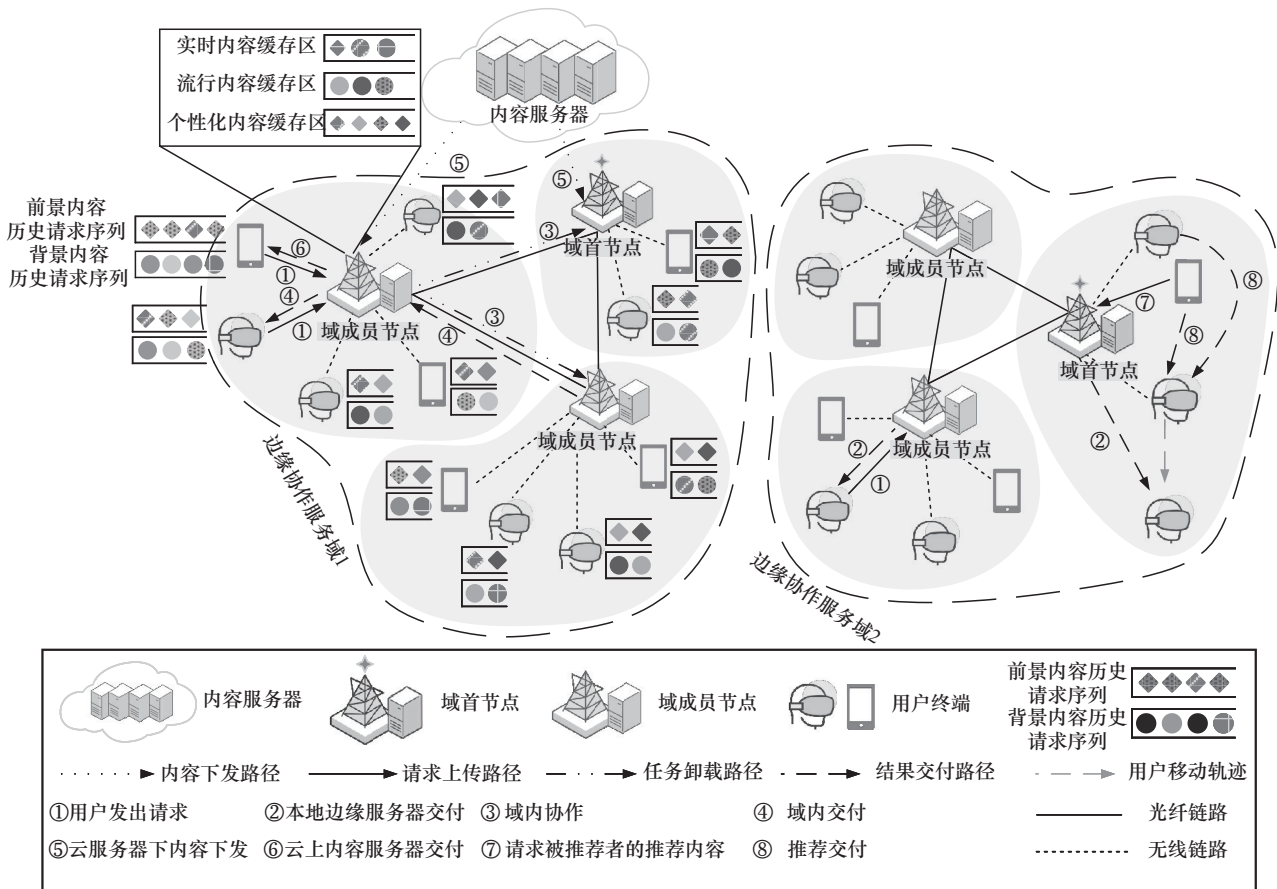


图1 面向元宇宙MAR应用的异质内容缓存服务架构

地发出请求并成为接收者，接收者在移动过程中总是向距离最近的边缘节点发出内容请求。

1.2 用户行为表征模型

1) 用户移动性模型

本文假设用户 u 的位置是随机分布的，并且在时间 t 可以被表示为 $L_u(t) = \{x_u, y_u\}$ ， x_u 和 y_u 分别代表用户在空间中所处的经度和纬度。只有在 2 个移动用户保持通信连接时，内容推荐和数据卸载等事件才会触发。

2) 用户内容请求模型

通常情况下，用户根据自身偏好来请求前景内容，在推荐机制的作用下，用户的实际请求由自身偏好和推荐引导共同影响。假设用户 u 以 a_u 的概率接受推荐，该概率是与用户的异构性相关的先验信息。用 $z_{u,i} \in \{0,1\}$ 表示内容 i 是否被推荐给用户 u ， $z_{u,i} = 1$ 表示内容 i 已经被推荐给用户 u ， $z_{u,i} = 0$ 表示内容 i 未被推荐给用户 u 。定义 $h_{u,i}^{\text{pref}}$ 为用户 u 对内容 i 的自身偏好，可以通过用户的历史请求来估计。因此，用户 u 对内容 i 的实际请求概率可表示为

$$h_{u,i}^{\text{req}} = a_u \frac{z_{u,i} h_{u,i}^{\text{pref}}}{\sum_{j \in \mathcal{F}} z_{u,j} h_{u,j}^{\text{pref}}} + (1 - a_u) \frac{(1 - z_{u,i}) h_{u,i}^{\text{pref}}}{\sum_{j \in \mathcal{F}} (1 - z_{u,j}) h_{u,j}^{\text{pref}}}, i \in \mathcal{F} \quad (1)$$

1.3 通信模型

使用香农公式计算用户终端 u 和用户终端 v 之间的数据传输时间 T_{u2v} 、用户终端 u 将任务卸载到本地边缘服务器 e_i 的传输时间 T_{u2e} 、本地边缘服务器 e 传输数据到用户终端 u 的时间 T_{e2u} 和云服务器传输数据到边缘服务器 e 的时间 T_{c2e} 分别为

$$T_{u2v} = \frac{S_i}{r_{u2v}} = \frac{S_i}{B \text{lb} \left(1 + \frac{p^{v,\text{tr}}}{\sigma^2} \right)}, \forall u, v \in \mathbf{u} \quad (2)$$

$$T_{u2e} = \frac{S_i}{r_{u2e}} = \frac{S_i}{\omega_{u,e_i} B_{e_i} \text{lb} \left(1 + \frac{p^{u,\text{tr}} g_{u,e_i}}{\sigma^2} \right)}, \forall u \in \mathbf{u}, e_i \in \mathbf{e} \quad (3)$$

$$T_{e2u} = \frac{S_i}{r_{e2u}} = \frac{S_i}{\omega_{u,e_i} B_{e_i} \text{lb} \left(1 + \frac{p^{e_i,\text{tr}} g_{u,e_i}}{\sigma^2} \right)}, \forall u \in \mathbf{u}, e_i \in \mathbf{e} \quad (4)$$

$$T_{c2e} = \frac{S_i}{r_{c2e}} = \frac{S_i}{\omega_{e_i,c} B_c \text{lb} \left(1 + \frac{p^{c,\text{tr}} g_{e_i,c}}{\sigma^2} \right)}, \forall e_i \in \mathbf{e} \quad (5)$$

其中， B 是用户终端 u 和用户终端 v 之间的通信带宽， $p^{v,\text{tr}}$ 是用户终端 v 的发射功率， σ^2 是信道的高斯白噪声功率； B_{e_i} 是边缘服务器 e_i 和所有连接的用户终端之间的总通信带宽， ω_{u,e_i} 是边缘服务器 e_i 分配给用户终端 u 的带宽资源比例， $p^{u,\text{tr}}$ 是用户终端 u 的发射功率， g_{u,e_i} 是用户终端 u 和本地边缘服务器 e_i 之间的信道条件状态， $p^{e_i,\text{tr}}$ 是边缘服务器 e_i 的发送功率； B_c 是云服务器和边缘服务器之间的总带宽， $\omega_{e_i,c}$ 是云服务器分配给边缘服务器 e_i 的带宽资源比例， $p^{c,\text{tr}}$ 是云服务器的发射功率， $g_{e_i,c}$ 是边缘服务器 e_i 和云服务器之间的信道条件状态。

1.4 时延模型

针对 MAR 应用中异质内容所需存储空间和缓存位置的差异，背景内容采用云-边交付模式，前景内容采用云-边交付和终端 D2D 推荐集成交付模式。

1) 背景内容的平均响应时延

MAR 中的背景内容请求任务需要通过蜂窝无线网络上传到边缘节点进行处理，并在 3 种可选情况下进行交付：本地边缘服务器命中，边缘协作服务域命中和云服务器命中。

具体来讲，当用户终端 u 产生背景内容请求时，需要将该请求上传至直接相连的边缘节点并检测服务器中是否有相应缓存内容，所需传输时间为 T_{u2e} 。如果直接相连的边缘服务器命中该内容，则缓存命中标识 $\theta_{u,e} = 1$ ，传输缓存内容所需的时间为 T_{e2u} 。若直接相连的边缘服务器未命中该内容，则需进一步查询所处边缘协作服务域内的其他边缘节点的缓存空间，若命中，则缓存命中标识 $\theta_{u,\text{ccsd}} = 1$ 。由于边缘服务器之间通过光纤进行数据传输且服务请求数据量较小，因此可以忽略域内传输请求的时延。边缘协作服务域内的边缘节点所组成的缓存网络将该缓存内容经由用户终端 u 直接相

连的边缘服务器传递给用户终端 u 。由于缓存数据量通常较大,通过光纤链路和无线链路传输给用户终端 u 的时延分别用 T_{c2c} 和 T_{c2u} 表示。若用户终端所处边缘协作服务域内的边缘服务器均未缓存该内容,则需从云服务器中下载所需的背景内容(即云服务器命中),缓存命中标识 $\theta_{e_0} = 1$,产生的内容下载时延和内容传输时延分别为 T_{c2c} 和 T_{c2u} 。

综上,背景内容的平均响应时延 T^{bg} 可表示为

$$T^{bg} = \begin{cases} T_{u2c} + T_{c2u}, \theta_{u,e} = 1 \\ T_{u2c} + T_{c2c} + T_{c2u}, \theta_{u,ecsd} = 1, \forall u \in \mathbf{u}, e \in \mathbf{e} \\ T_{u2c} + T_{c2c} + T_{c2u}, \theta_{e_0} = 1 \end{cases} \quad (6)$$

2) 前景内容的平均响应时延

终端 D2D 推荐集成交付模式。当用户终端 u 产生前景内容的请求任务时,首先查询用户终端 u 是否接受用户协作范围内其他空闲用户的推荐内容。如果用户终端 u 接受了推荐,则接受推荐标识符 $\rho_u = 1$ 。忽略推荐者从边缘节点获取用户终端 u 推荐列表的时间。并且,由于向用户终端发送的推荐数据仅包括内容标题等摘要信息,摘要数据量与原始数据量在数据级上具有较大差距,因此可以忽略推荐摘要信息的传输时间。如果用户终端 u 接受推荐,在交付阶段用户终端 u 优先选择通过 D2D 方式获取推荐者的缓存内容,其时延为 T_{u2v} 。如果在 D2D 通信连接时间内未获取完整的内容,用户需通过蜂窝链路从边缘节点获取剩余数据,其时延为 T_{c2u} 。

云-边交付模式。如果用户拒绝了所推荐的内容,则接受推荐标识符 $\rho_u = 0$ 。用户终端 u 首先查询其自身缓存空间是否缓存相应内容,如果终端缓存命中,则缓存命中标识 $\theta_u = 1$,用户终端 u 直接从本地缓存中获取数据。若自身设备未命中,用户终端 u 将请求上传到与用户终端 u 直接相连的边缘服务器上并检测该边缘服务器中是否缓存相应内容,所需传输时间为 T_{u2c} 。如果直接相连的边缘服务器命中,则缓存命中标识 $\theta_{u,e} = 1$,缓存内容返回所需传输时间为 T_{c2u} 。若直接相连的边缘服务器未命中该内容,则需进一步查询所处边缘协作服务域内的其他边缘节点的缓存空间,若命中,则缓存命中标识 $\theta_{u,ecsd} = 1$ 。边缘协作服务域将该缓存内容经由用户终端 u 直接相连的边缘服务器传递给用户终端 u 。边缘节点间的传输时延和边缘节点与用户终端 u 间的传输时延分别为 T_{c2c} 和 T_{c2u} 。若用户终

端所处边缘协作服务域内的边缘服务器均未缓存该内容,则需从云服务器中下载所需的前景内容,缓存命中标识 $\theta_{e_0} = 1$,产生的内容下载时延和内容传输时延分别为 T_{c2c} 和 T_{c2u} 。

综上,前景内容的平均响应时延 T^{fg} 可表示为

$$T^{fg} = \begin{cases} T_{u2v} + T_{c2u}, \rho_u = 1 \\ 0, \rho_u = 0, \theta_u = 1 \\ T_{u2c} + T_{c2u}, \rho_u = 0, \theta_{u,e} = 1 \\ T_{u2c} + T_{c2c} + T_{c2u}, \rho_u = 0, \theta_{u,ecsd} = 1 \\ T_{u2c} + T_{c2c} + T_{c2u}, \rho_u = 0, \theta_{e_0} = 1 \end{cases} \quad (7)$$

2 异质内容主动缓存与个性化交付机制

2.1 总体框架

本文的异质内容主动缓存与个性化机制包括 3 个阶段:协作域构建阶段、异质内容预缓存阶段以及前景内容推荐与交付阶段。第一个阶段构建了边缘协作服务域并进行域首节点选择,确定了多边缘节点间的协作关系;第二个阶段在已构建的协作域中分别面向前景和背景内容进行内容预缓存配置,以提高边缘网络的缓存利用率和内容命中率;第三个阶段针对前景内容引入个性化推荐机制,并融合端到端通信模式及蜂窝无线通信模式进行差异化交付,进一步提高了终端的缓存利用率和业务的响应速度。

2.2 用户行为和资源感知的边缘协作服务域构建方法

2.2.1 边缘协作服务域构建方法

本文采用用户行为和资源感知的边缘协作服务域构建方法,首先将用户聚类形成用户簇 $sucs$,并将边缘节点聚类形成边缘协作服务域 $ecsd$ 。系统中的用户簇集和边缘协作服务域集可以分别用 $\mathbf{sucs} = \{sucs_i\}_{i=1}^M$ 和 $\mathbf{ecsd} = \{ecsd_i\}_{i=1}^N$ 表示,其大小分别为 M 和 N 。边缘协作服务域 $ecsd$ 是由一组边缘服务器及其提供服务的用户组成的,这些边缘服务器协同工作并共享存储资源,为 MAR 用户提供内容缓存服务。边缘协作服务域 $ecsd$ 中的边缘服务器可以按照角色划分为 2 种类型:域首节点和域成员节点。一个边缘协作服务域包括一个域首节点和若干个域成员节点,它们之间通过光纤链路相互连接,用户终端通过无线链路连接边缘服务器,用户终端之间也可以通过无线链路实现 D2D 通信。

域首节点是边缘协作服务域的集中控制器,负

责从云服务器中获取异质内容, 维护边缘协作服务域中的节点列表、每个边缘服务器中的缓存内容列表和每个边缘服务器为所服务用户定制的前景内容推荐列表。域首节点可以通过做出内容缓存、内容交付和内容推荐决策来管理边缘协作服务域中的资源。域成员节点负责根据域首节点制定的服务策略以及周边用户的服务需求, 执行内容缓存和交付任务。

本文考虑用户的地理位置和内容偏好、响应质量、时延要求等用户特征信息, 融合密度聚类算法 DBSCAN 和 k-means++ 聚类算法划分用户簇。首先, DBSCAN 聚类算法将地理位置相近的用户合并成一个簇, 考虑到用户的高动态性和移动轨迹, 要求任一用户终端与邻域内其他用户终端的距离都不超过设定门限 D 。基于用户地理位置聚簇可以直观地显示城市人群的密集程度, 并且可以筛除地理位置上孤立的用户, 有助于后续用户簇的构建。

然后, 使用请求业务类型 $\mu_u(t)$ 、缓存质量要求 $R_u(t)$ 、时延要求 $T_u(t)$ 构建用户特征矩阵 $F_u(t) = \{\mu_u(t), R_u(t), T_u(t)\}$ 。其中, 缓存质量要求 $R_u(t)$ 和时延要求 $T_u(t)$ 表示用户 u 对同一业务的个性化需求, 如高档小区住户的缓存要求较高, 他们愿意花高价获得高质量的缓存; 偏远地区的用户可能因为昂贵的网络资费选择低质量的缓存; 校园网用户对缓存质量的要求则可能和要求的响应时间相关, 可以通过调节分辨率来调整缓存质量, 并改善响应时间。经过 DBSCAN 聚类算法去噪之后, 每个聚类结果中的用户根据用户特征矩阵使用 k-means++ 算法进行二次聚类以构建用户簇。

基于划分好的用户簇 \mathbf{sucs} , 基于算法 1 进一步划分边缘协作服务域并依据域内边缘节点的服务能力进行域首节点的选择。

算法 1 用户行为和资源感知的边缘协作服务域构建算法

输入 系统中的用户终端集合 $\mathbf{u} = \{u_i\}_{i=1}^U$, 用户地理位置集合 $\mathbf{L} = \{(x_u, y_u) | u \in \mathbf{u}\}$, 用户特征矩阵 $F_u(t) = \{\mu_u(t), R_u(t), T_u(t)\}$, 系统中的边缘服务器集合 $\mathbf{e} = \{e_i\}_{i=1}^E$

输出 用户簇集合 $\mathbf{sucs} = \{\mathbf{sucs}_i\}_{i=1}^M$, 边缘协作服务域集合 $\mathbf{ecsd} = \{\mathbf{ecsd}_i\}_{i=1}^N$

1) 根据用户地理位置集合 $\mathbf{L} = \{(x_u, y_u) | u \in \mathbf{u}\}$, 执行 DBSCAN 聚类算法, 得到聚类结果

\mathbf{sucs}' ;

2) for each $\mathbf{sucs}' \in \mathbf{sucs}'$ do:

3) 根据用户特征矩阵 $F_u(t) = \{\mu_u(t), R_u(t), T_u(t)\}$, 执行 k-means++ 聚类算法, 划分用户簇;

4) end for

5) 得到用户簇集合 $\mathbf{sucs} = \{\mathbf{sucs}_i\}_{i=1}^M$;

6) for each $e \in \mathbf{e}$ do:

7) 计算边缘服务器 e 服务范围 $R(e)$ 内的用户数量 $\text{sum}_e \leftarrow \sum u_i, \forall u_i \in R(e)$;

8) 计算边缘服务器 e 服务范围内 $R(e)$ 属于各用户簇的用户数量 $\text{sum}_{\mathbf{sucs}} \leftarrow \sum u_i, u_i \in \mathbf{sucs}_i, \forall u_i \in R(e)$;

9) if $\text{sum}_e \geq \text{NUM}$:

10) $\text{label}_e \leftarrow \text{label}_{\mathbf{sucs}}$, $\text{label}_{\mathbf{sucs}}$ 为边缘服务器 e 服务范围 $R(e)$ 内用户数量最多的用户簇的簇标;

11) end if

12) end for

13) 得到边缘协作服务域集合 $\mathbf{ecsd} = \{\mathbf{ecsd}_i\}_{i=1}^N$;

14) for each $\mathbf{ecsd} \in \mathbf{ecsd}$ do:

15) for each $e \in \mathbf{ecsd}$:

16) 根据式(9)对边缘协作服务域 \mathbf{ecsd} 的边缘节点进行评估, 选取评估值最大的节点为域首节点 e_{hd} , $e_i, \forall e_i \neq e_{hd}, e_i \in \mathbf{ecsd}$ 为域成员节点;

17) end for

18) end for

2.2.2 域首节点选择方法

在域首节点选择过程中, 本文选择缓存空间较大并靠近簇中心的边缘节点作为域首节点, 用于降低平均响应时延并缓存尽可能多的异质内容, 边缘节点的服务能力可以通过式(8)评估。

$$\kappa_n = \omega_0 \frac{m_e}{\bar{m}} + (1 - \omega_0) \frac{d_e}{\bar{d}} \quad (8)$$

其中, m_e 和 d_e 分别表示边缘节点 e 的存储空间和边缘节点 e 距离簇中心的距离, ω_0 为权重系数, 用于平衡边缘节点存储空间和地理位置对域首节点选择的影响。选择域首节点后, 簇中的剩余边缘节点为域成员节点。

2.3 基于存储空间划分和用户偏好预测的异质内容预缓存机制

2.3.1 异质内容流行度预测方法

本文将 MAR 应用中的内容分为背景内容 \mathcal{F}^{bg} (如虚拟物品和场景的三维模型及其渲染资源) 和前景内容 \mathcal{F}^{fg} (如用户虚拟形象、用户互动数据和用户生成内容)。由于背景内容和现实内容相关, 数据量大且更新频率低, 可以采用域内-域间流行度预测方法进行缓存配置并加载到边缘节点。前景内容包含个性化的数据, 更新频率较高, 可以采用多因子分域流行度预测方法进行缓存配置。并且, 为了进一步提高缓存命中率并适应用户的流动性, 可将前景内容缓存在用户终端和边缘节点的缓存空间中。

1) 域内-域间流行度预测方法

背景内容的综合流行度由域内流行度和域间流行度共同影响。在边缘协作服务域 ecsd 中, 考虑到同一用户频繁请求某一内容会造成较高的流行度, 不能反映该内容在域内的真实流行情况, 具有较大流行度偏差。因此, 通过均衡考虑内容的请求人数和频率, 背景内容 f_i^{bg} 的域内流行度 $P_{\text{bg},i}^{\text{ecsd}}$ 定义为

$$P_i^{\text{ecsd}} = \frac{\omega_1 w_i^{\text{ecsd}}}{U_{\text{ecsd}} - U_{\text{ecsd},i}} \quad (9)$$

其中, w_i^{ecsd} 表示 ecsd 中背景内容 f_i^{bg} 的请求数, U_{ecsd} 表示 ecsd 内用户总人数, $U_{\text{ecsd},i}$ 表示域内请求背景内容 f_i^{bg} 的用户总数, ω_1 是权重系数, 用于调节请求数和用户数之间的数量级差。

全局流行度 P_i^{glob} 表示背景内容 f_i^{bg} 在多个边缘协作服务域中的加权流行度

$$P_i^{\text{glob}} = \sum_{\text{ecsd} \in \text{ecsd}} \frac{E_{\text{ecsd}}}{E} P_i^{\text{ecsd}} \quad (10)$$

其中, E_{ecsd} 表示 ecsd 内的边缘节点数, E 表示系统内的边缘节点总数。

因此, 背景内容 f_i^{bg} 的综合流行度 \bar{P}_i 可以定义为

$$\bar{P}_i = \omega_2 P_i^{\text{ecsd}} + (1 - \omega_2) P_i^{\text{glob}} \quad (11)$$

其中, ω_2 为权重系数, 用于表征域内流行度和域间流行度对综合流行度的影响程度。

背景内容 f_i^{bg} 在不同边缘协作服务域中的综合流行度集合可以表示为 $\bar{\mathbf{P}}_i = \{ \bar{P}_i^{\text{ecsd}} \}, \text{ecsd} \in \text{ecsd}$ 。

2) 多因子分域流行度预测方法

前景内容的流行度由域内历史流行内容、内容请求人数和内容大小共同决定。域内历史流行内容 c_{ecsd} 由每个域内占比最多的请求内容类型决定。综上, 前景内容 f_i^{fg} 在不同边缘协作服务域中的流行度定义为

$$P_i^{\text{ecsd}} = \begin{cases} \omega_3 \frac{U_{\text{ecsd},i}}{U_{\text{ecsd}}} + \omega_4 \frac{s_i}{\bar{s}}, c_i = c_{\text{ecsd}} \\ \omega_3 \frac{U_{\text{ecsd},i}}{U_{\text{ecsd}}} - \omega_4 \frac{s_i}{\bar{s}}, c_i \neq c_{\text{ecsd}} \end{cases} \quad (12)$$

其中, $U_{\text{ecsd},i}$ 表示域内请求前景内容 f_i^{fg} 的用户总数, s_i 表示前景内容 f_i^{fg} 的存储空间需求, \bar{s} 表示前景内容占用存储空间的最大值, c_i 表示前景内容 f_i^{fg} 的类型, ω_3 、 ω_4 为权重参数且 $\omega_3 + \omega_4 = 1$ 。

3) 缓存空间划分方法

将每个边缘服务器的缓存空间划分为 3 个部分: 第一部分缓存空间 L_1 用于缓存实时数据, 第二部分缓存空间 L_2 用于预缓存个性化内容 (即前景内容数据), 第三部分缓存空间 L_3 用于预缓存流行内容 (即背景内容数据)。对于用户终端的缓存空间, 由于存储空间有限, 缓存空间被划分为 2 个部分, 分别为 L_1 和 L_2 。假设 φ_1 、 φ_2 和 φ_3 分别代表服务器缓存空间中 L_1 、 L_2 和 L_3 的比例, $\varphi_1 + \varphi_2 + \varphi_3 = 1$, 且用户终端的 φ_3 设置为 0。

2.3.2 协作域全局内容冗余预缓存策略

在边缘协作服务域构建完成后, 需要进一步制定域内的服务内容预缓存策略, 以最小化边缘节点 e 在预缓存过程中的存储成本 M_e 、时延成本 T_e 和内容缺失率 H_e , 最优化问题模型可以表示为

$$\begin{aligned} \min: & \frac{1}{U_{\text{ecsd},e}} \sum_{e \in \text{ecsd}} \omega_1 M_e + \omega_2 T_e + \omega_3 H_e \\ \text{s.t.} & \text{C1: } \sum_{i \in \mathcal{F}} o_{e,i} s_i \leq m_e^{\text{precache}}, e \in \text{ecsd} \\ & \text{C2: } \sum_{e \in \text{ecsd}} \sum_{i \in \mathcal{F}} o_{e,i} \leq \eta^e F \end{aligned} \quad (13)$$

其中, $U_{\text{ecsd},e}$ 表示 ecsd 中边缘节点 e 服务用户的数量, $o_{e,i} = 1$ 表示边缘节点 e 已缓存内容 i , s_i 表示内容 i 的存储空间需求, m_e^{precache} 表示边缘节点 e 的预缓存存储空间容量, η^e 表示冗余预缓存的比例。约束 C1 确保预缓存内容的存储空间总量不能超过预缓存空间容量。约束 C2 确保边缘协作服务域中冗余内容数量不超过规定上限。为了简化模型, 假设节点的状态、动作和奖励都是有限的, 因此该预缓

存过程可以看作一个有限的马尔可夫决策过程 (MDP, Markov decision process) [19], 上述预缓存策略 $\pi: S \times A \rightarrow [0,1]$ 可以定义为

$$\pi(a|s) = \Pr [A_t = a | S_t = s], s \in S, a \in A \quad (14)$$

由于深度 Q 学习 (DQN, deep Q-network) 可以解决具有大规模状态和动作空间的问题[19], 因此本文使用 DQN 来制定预缓存策略。边缘协作服务域服务状态、域首节点预缓存决策和域首节点预缓存奖励定义如下。

1) 边缘协作服务域服务状态

为了制定域首节点的最优预缓存决策, 需要综合考虑域内用户请求内容集合 $\mathbf{I}^{\text{ecsd}}(t)$ 、内容流行度集合 $\mathbf{P}^{\text{ecsd}}(t)$ 和域首节点的空闲存储空间 $m^{\text{hd}}(t)$ 。因此, 可定义 ecsd 在时间 t 的服务状态为 $\mathbf{S}^{\text{ecsd}}(t) = \{\mathbf{I}^{\text{ecsd}}(t), \mathbf{P}^{\text{ecsd}}(t), m^{\text{hd}}(t)\}$, 其中, $\mathbf{I}^{\text{ecsd}}(t) = \{I_{q(1)}(t), I_{q(2)}(t), \dots, I_{q(i)}(t)\}$, 如果 ecsd 内有用户请求内容 i , 则 $I_{q(i)}(t) = 1$, 否则 $I_{q(i)}(t) = 0$ 。

2) 域首节点预缓存决策

域首节点预缓存内容包含基本缓存内容集和冗余缓存内容集。定义域首节点预缓存动作集为 $\mathbf{A}^{\text{hd}} = \{\boldsymbol{\theta}^{\text{hd}}(t), \eta^e\}$, 其中 $\boldsymbol{\theta}^{\text{hd}}(t)$ 是域首节点请求内容集合, η^e 为冗余预缓存的比例。 $\boldsymbol{\theta}^{\text{hd}}(t)$ 可以表示为

$$\boldsymbol{\theta}^{\text{hd}}(t) = \{\theta_1(t), \theta_2(t), \dots, \theta_i(t)\} \quad (15)$$

其中, $\theta_i(t)$ 是内容缓存指示符, 如果 $\theta_i(t) = 1$, 则表示边缘节点缓存内容 i 。

3) 域首节点缓存奖励

为了提高缓存命中率 and 可靠性, 将收益制定为激励函数

$$R^{\text{hd}}(t) = \lg(1 + \sum_{i \in \mathcal{F}} P_i^{\text{ecsd}} r_i k_i s_i) \quad (16)$$

其中, r_i 是缓存内容 i 的命中率收益, k_i 是缓存内容 i 的可靠性收益。

本文假设 cost_i 为缓存内容 i 的时延成本, 为了降低时延, 将时延成本定义为惩罚函数

$$C^{\text{hd}}(t) = \sum_{i \in \mathcal{F}} \text{cost}_i s_i \quad (17)$$

综上, 域首节点的奖励函数可定义为

$$R^{\text{hd}}(\mathbf{S}^{\text{ecsd}}(t), \mathbf{A}^{\text{hd}}) = \lg(1 + \sum_{i \in \mathcal{F}} P_i^{\text{ecsd}} r_i k_i s_i) - \sum_{i \in \mathcal{F}} \text{cost}_i s_i \quad (18)$$

2.3.3 域成员节点的预缓存策略

1) 域成员节点服务状态

与边缘协作服务域服务状态的定义类似, 域成

员节点的服务状态需综合考虑边缘节点 e 关联用户的请求内容集合 $\mathbf{I}^e(t)$ 、边缘节点 e 服务范围内的内容流行度 $\mathbf{P}^e(t)$ 、域成员节点 e 的空闲存储空间 $m^e(t)$ 和域内缓存内容标识 $\boldsymbol{\delta}^e(t)$, 定义为 $\mathbf{S}^e(t) = \{\mathbf{I}^e(t), \mathbf{P}^e(t), m^e(t), \boldsymbol{\delta}^e(t)\}$, 如果 $\delta_i^e(t) = 1$, 表示边缘节点 e 可以从域内其他域成员节点获取请求内容 i 。

2) 域成员节点预缓存决策

将与域首节点预缓存决策的定义类似, 为了确定域成员节点的缓存内容集, 域成员节点的动作集定义为 $\mathbf{A}^e = \{\boldsymbol{\theta}^e(t)\}$, 其中 $\boldsymbol{\theta}^e(t)$ 为域成员节点请求内容的集合。

3) 域成员节点缓存奖励

与域首节点缓存奖励类似, 定义收益为激励函数, 成本为惩罚函数, 则域成员节点的预缓存奖励可表示为

$$R^e(\mathbf{S}^e(t), \mathbf{A}^e) = \lg(1 + \sum_{i \in \mathcal{F}} P_i^e r_i s_i) - \sum_{i \in \mathcal{F}} \text{cost}_i s_i \quad (19)$$

综上, 使用 DQN 求解每个边缘协作服务域中域首节点和域成员节点的预缓存决策。

2.4 前景内容个性化推荐与异质内容交付机制

2.4.1 前景内容个性化推荐机制

引入内容推荐机制可增加用户请求特定终端缓存内容的概率, 如果用户接受推荐, 则可以在 D2D 通信链路连接时间内通过 D2D 直接交付模式获取该前景内容; 如果用户拒绝了推荐并请求其他内容, 则需要向本地边缘服务器上传请求, 通过本地边缘服务器、边缘协作服务域甚至云服务器进行交付。因此, 推荐会影响终端缓存的资源利用率和数据卸载量, 从而影响平均响应时延。通过分析用户自身的历史偏好和对推荐行为的接受度, 可以为用户定制个性化的推荐列表, 引导用户请求推荐者终端已缓存的内容, 从而实现更高效的边缘缓存共享, 具体过程如下。

当推荐者 v 计划将自身已缓存内容向用户 u 推荐时, 首先向域首节点获取用户 u 的推荐内容列表, 然后从其缓存的前景内容中选择综合排名最高的内容 $f_{u,i}^{\text{fg}}$ 向用户 u 进行推荐, 如图 2 的步骤 1.1)~步骤 1.3) 所示。为了提高推荐命中率, 定义 Q_0 为推荐质量约束, 并规定只有用户 u 对前景内容 f_i^{fg} 的自身偏好大于 Q_0 时才可以推荐给用户 u 。面向用户 u 的前景内容个性化推荐排序算法

如算法 2 所示。

算法 2 前景内容个性化推荐排序算法

输入 前景内容集合 \mathcal{F}^{fg} , 用户 u 对前景内容 i 的固有偏好 $h_{u,i}^{\text{pref}}, \forall i \in \mathcal{F}^{\text{fg}}$, 推荐质量约束 Q_0

输出 用户 u 的推荐排名列表 $\mathcal{F}_u^{\text{fg}} = [\bar{f}_{u,1}^{\text{fg}}, \bar{f}_{u,2}^{\text{fg}}, \dots]$

- 1) for each $i \in \mathcal{F}^{\text{fg}}$ do:
- 2) if $h_{u,i}^{\text{pref}} \geq Q_0$:
- 3) 置 $z_{u,i} = 1$ 且 $z_{u,j} = 0, \forall j \in \mathcal{F}^{\text{fg}} \setminus i$;
- 4) 根据式 (1)、式 (2) 和式 (6), 计算 $\text{score}_i \leftarrow \frac{h_{u,i}^{\text{req}}}{T_{\text{rec},i}}$ 作为前景内容 i 的推荐综合评分;
- 5) end if
- 6) end for
- 7) 根据不同前景内容的推荐综合评分进行排序, 生成推荐排名列表 $\mathcal{F}_u^{\text{fg}} = [\bar{f}_{u,1}^{\text{fg}}, \bar{f}_{u,2}^{\text{fg}}, \dots]$

2.4.2 前景内容的差异化交付策略

用户 u 需经历推荐阶段和交付阶段获取推荐内容, 流程如图 2 所示。

1) 推荐阶段

当推荐者 v 满足以下条件, 则向用户 u 发送前景内容 $\bar{f}_{u,i}^{\text{fg}}$ 的摘要信息: 推荐者 v 空闲; 推荐者 v 与用户 u 的距离小于 R_{off} , R_{off} 是 D2D 发射器的最大卸载半径; 以推荐者 v 为中心的 R_{pro} 半径范围内没有其他接收器。在推荐阶段, 一个用户可以收到多个用户的推荐, 该用户可以接受这些推荐, 也可以选择拒绝。如果用户 u 接受了推荐者 v 的前景内容 $\bar{f}_{u,i}^{\text{fg}}$, 则请求推荐者 v 传输该前景内容, 如图 2 的步骤 1.4)~步骤 1.6) 所示。

2) 交付阶段

当用户 u 接受 v 的内容推荐之后, 可通过 D2D 直接交付、域内边缘节点主动交付和域间边缘节点协作交付模式获取请求内容。

① 域内边缘节点主动交付。如果用户 u 由于自身移动未能在 D2D 通信链路连接时间内获取完整的前景内容 $\bar{f}_{u,i}^{\text{fg}}$ 的数据, 且用户 u 未离开当前所在的边缘协作服务域, 可以从其关联的基站获得剩余的推荐数据, 如图 2 的步骤 2.5)~步骤 2.6) 所示。

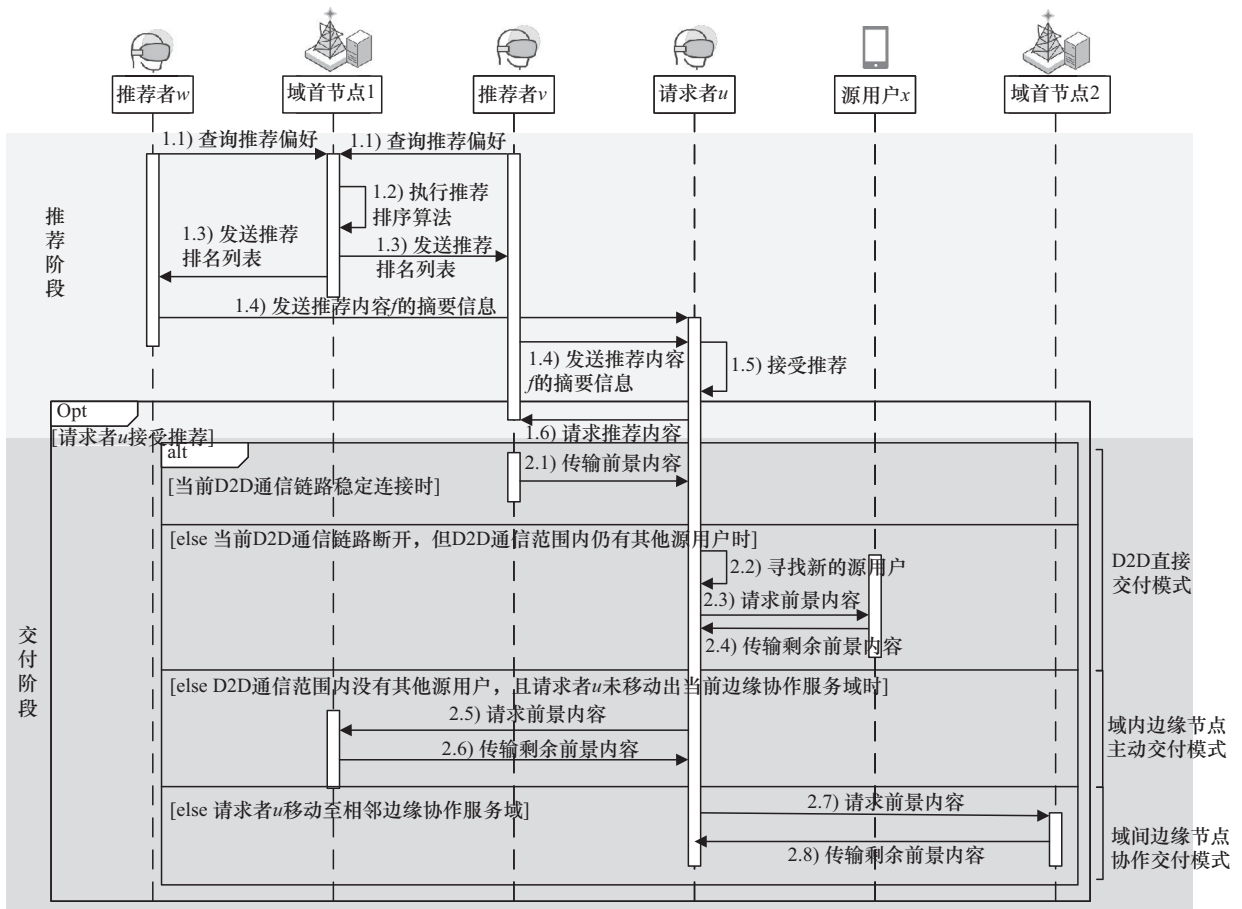


图 2 推荐交付流程

②域间边缘节点协作交付。若用户 u 高速移动, 在内容交付过程中离开其所属边缘协作服务域, 则移动前所属边缘协作服务域的域首节点通过光纤链路将前景内容 $f_{u,i}^{fg}$ 的数据发送到移动后所属边缘协作服务域的域首节点, 再经由用户 u 在移动轨迹上相关联的边缘节点发送给用户, 如图 2 的步骤 2.7)~步骤 2.8) 所示。

3 实验与分析

3.1 数据集

实验数据采用 OpenCellid 香港基站数据集, 通过数据预处理筛选出 $4\ 000\text{ m} \times 4\ 000\text{ m}$ 范围内的 100 个基站 (即边缘服务器) 的位置信息。

3.2 参数设置

假设系统中有 200 个用户终端, 用户对异质内容的请求频率服从 Zipf 分布。边缘基站的信号覆盖半径为 200 m, 边缘服务器之间的信道带宽为 10 MHz, 用户之间的最大 D2D 通信距离为 100 m, 用户的最大发射功率为 23 dBm^[20], 通用参数设置如表 1 所示。在深度强化学习模型训练过程中, 学习率为 0.05, 衰退率为 0.9, 经验池容量设置为 1 000, 批处理数量为 32, 奖励折扣因子为 0.9, 优化函数采用 Adam 优化器^[21]。本文的仿真结果为 50 次实验的平均值, 以确保实验结果的稳定性。

3.3 对比策略

本节将 DRCCS-DQN (本文提出的支持 D2D 通信的推荐感知的协同缓存策略) 与 NCPCR (无协作域、无预缓存、无推荐策略)、CCS-AGP (基于服务区划分、服务器分组和存储空间分区的协同缓存策略^[14])、AIEC-RSC (基于人工智能和边缘协作的可靠服务策略^[15]) 以及 DRCCS-Q (使用 Q-Learning 算法代替本文中的 DQN 算法) 4 种策略进

行对比, 从缓存命中率、平均响应时延等方面对划域必要性和机制先进性进行了验证。4 种策略的工作机制对比如表 2 所示。

表 1 参数设置

参数	取值
基站覆盖半径/m	200
最大 D2D 通信距离/m	100
背景内容 f_i^{bg} 占用存储空间大小 s_i /KB	[500,1 000]
前景内容 f_i^{fg} 占用存储空间大小 s_j /KB	[200,400]
带宽 B /MHz	10
终端最大发射功率 $p^{u,lr}$ /dBm	23
高斯白噪声 σ^2 /dBm	-174
边缘服务器 e_i 发射功率 $p^{e,lr}$ /dBm	43

3.4 仿真结果与分析

3.4.1 划域结果

本文使用算法 1 划分边缘协作服务域, 并展示了边缘协作服务域的一个划分实例, 如图 3 所示, $4\ 000\text{ m} \times 4\ 000\text{ m}$ 范围内的 100 个基站形成了 4 个边缘协作服务域。其中, 菱形标记代表未聚集成边缘协作服务域的孤立边缘节点, 其他相同符号的标记代表同一边缘协作服务域中的边缘节点。

3.4.2 缓存命中率

首先, 本文比较了不同策略的缓存命中率对比, 如图 4 所示。从图 4 可以看出, NCPCR 策略的缓存命中率最低, CCS-AGP 策略的缓存命中率低于 AIEC-RSC 策略、DRCCS-Q 策略和 DRCCS-DQN 策略。AIEC-RSC 策略的缓存命中率高于 NCPCR 策略和 CCS-AGP 策略, 但略低于 DRCCS-Q 策略和 DRCCS-DQN 策略, 这是由于 DRCCS-DQN 策略不仅构建了边缘协作域, 还增加了基于

表 2 4 种策略的工作机制对比

策略名称	是否划分协作域	是否进行预缓存	是否考虑异质内容	是否进行前景内容推荐	是否进行存储空间划分	算法模型
NCPCR	否	否	是	否	否	DQN
CCS-AGP	是	否	否	否	是	数值分析
AIEC-RSC	是	是	否	否	否	DQN
DRCCS-Q	是	是	是	是	是	Q-Learning
DRCCS-DQN	是	是	是	是	是	DQN

DQN 算法的内容预缓存过程,使得缓存命中率进一步升高。此外,对比 DRCCS-DQN 策略和 DRCCS-Q 策略可以得出,本文所提出的 DRCCS-DQN 策略的缓存命中率高于其他 4 种策略的原因包括分区存储异质内容以及引入前景内容个性化推荐机制。

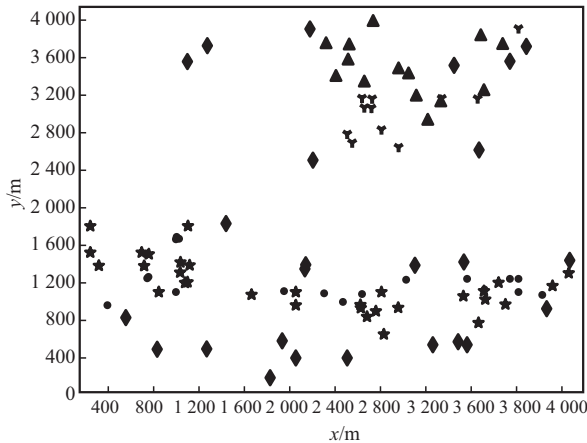


图 3 划域实例

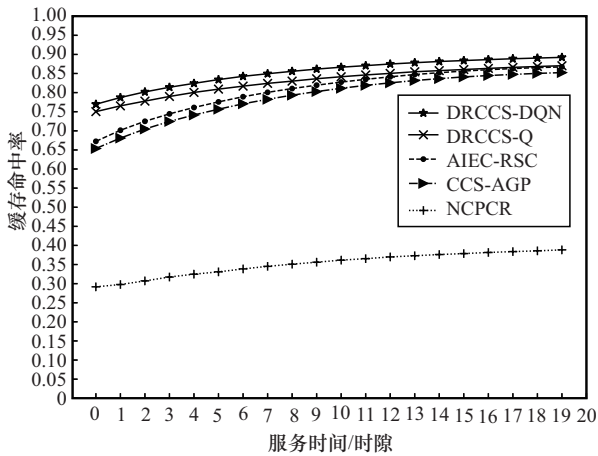


图 4 不同策略的缓存命中率对比

然后,本文比较了存储空间容量对缓存命中率的影响,如图 5 所示。从图 5 中可以看出,NCPCR 策略的缓存命中率最低且随缓存空间的增大而升高,CCS-AGP 策略和 AIEC-RSC 策略的缓存命中率并不完全随着缓存空间增大而升高,当缓存空间增大到一定程度,缓存命中率才明显升高。DRCCS-Q 策略的缓存命中率呈缓慢升高趋势,而 DRCCS-DQN 策略的缓存命中率在存储空间为 10 000 KB 时达到峰值。该仿真结果验证了预缓存策略的有效性,同时也揭示了用户偏好对预缓存策略性能的影响。

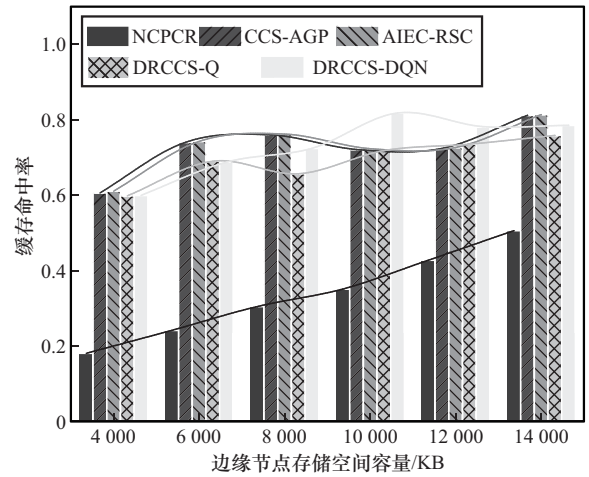
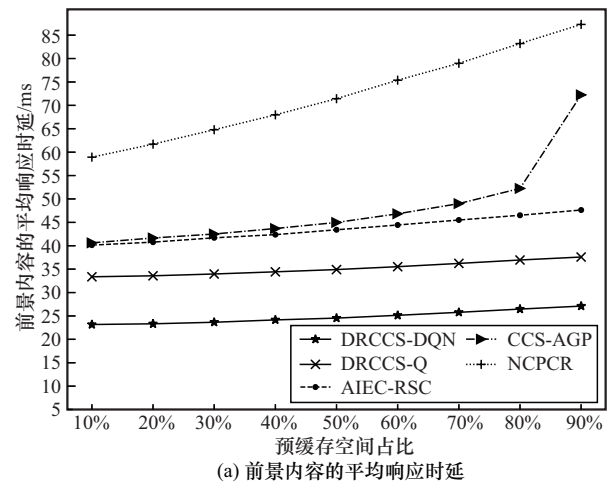


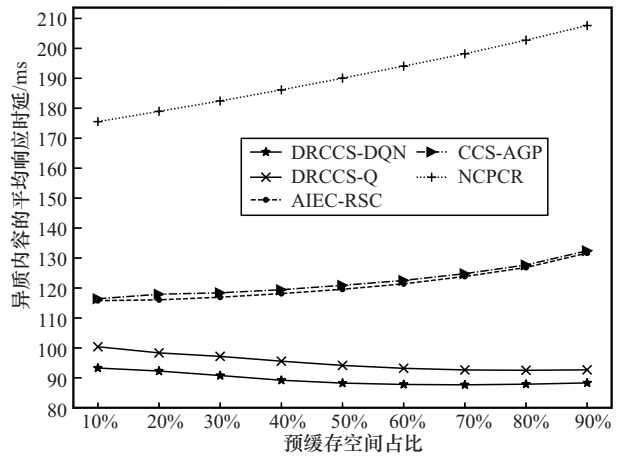
图 5 存储空间容量对缓存命中率的影响

3.4.3 预缓存/实时缓存空间占比对时延的影响

本文比较了不同策略在不同预缓存/实时缓存空间占比下前景内容和异质内容的平均响应时延,如图 6 所示,假设前景内容和背景内容的缓存空间占比为 1:1。



(a) 前景内容的平均响应时延



(b) 异质内容的平均响应时延

图 6 预缓存/实时缓存空间占比对平均响应时延的影响

如图 6(a)所示, 所有策略的前景内容的平均响应时延均随着预缓存空间的增加 (即实时缓存空间的减少) 而不断提高, 并且 NCPCR 策略的平均响应时延大于 55 ms 且最高可达 85 ms, 始终高于其他策略。NCPCR 策略仅依赖边缘节点实时缓存进行内容交付, 随着实时缓存空间的减小, 需频繁向云服务器请求所需内容, 平均响应时延不断增大。CCS-AGP 策略通过构建协作域共享边缘节点的缓存内容, 降低了从云服务器下载缓存内容的频次, 因此平均时延的增速低于 NCPCR 策略。与 CCS-AGP 策略相比, AIEC-RSC 策略的平均响应时延增速平缓, 这是由于该策略通过预缓存配置在本地快速满足了部分用户的业务需求。DRCCS-Q 策略在所有预缓存/实时缓存空间占比下的平均响应时延均低于 40 ms。本文所提出的 DRCCS-DQN 策略在所有预缓存/实时缓存空间占比下的平均响应时延均低于 30 ms。这是由于个性化的推荐机制有效利用了用户终端设备的缓存内容, 进一步提高了缓存命中率, 降低了前景内容的平均响应时延。

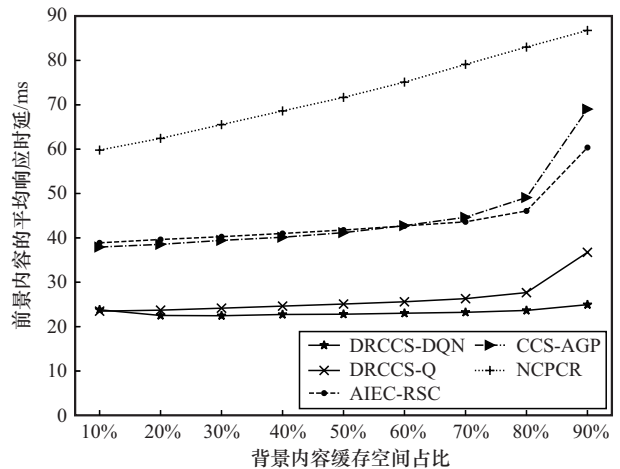
如图 6(b)所示, NCPCR 策略的异质内容的平均响应时延随着预缓存空间的增加 (即实时缓存空间的减少) 而不断增加, 同时异质内容的平均响应时延均大于 170 ms, 最高可超过 200 ms。CCS-AGP 策略和 AIEC-RSC 策略的异质内容的平均响应时延性能相近, 始终大于 110 ms 但均不超过 140 ms。DRCCS-Q 策略在所有预缓存/实时缓存空间占比下的平均响应时延低于 100 ms。本文所提出的 DRCCS-DQN 策略的异质内容的平均响应时延趋于 90 ms, 并且平均响应时延随着实时缓存空间的减少而缓慢降低, 这是由于边缘节点和用户终端通过定制化的推荐排序算法计算并缓存了用户感兴趣的部分内容数据, 并通过终端 D2D 推荐交付的方式进一步降低了异质内容的平均响应时延。

3.4.4 背景/前景内容缓存空间占比对时延的影响

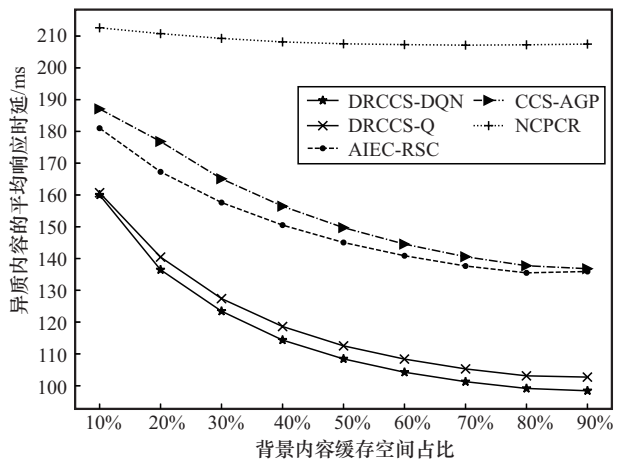
本文比较了不同策略在不同背景/前景内容缓存空间占比下前景内容和异质内容的平均响应时延, 如图 7 所示, 假设预缓存空间和实时缓存空间占比为 1:1。

如图 7(a)所示, 所有策略的前景内容的平均响应时延均随着背景内容缓存空间占比的增加 (即前

景内容缓存空间占比的减少) 而升高。本文所提出的 DRCCS-DQN 策略在所有背景内容缓存空间占比范围内变化不大且均低于 30 ms, 这是由于基于 D2D 通信范围内的用户终端缓存内容共享可满足进一步满足部分用户的服务请求。



(a) 前景内容的平均响应时延



(b) 异质内容的平均响应时延

图 7 背景/前景内容缓存空间占比对平均响应时延的影响

如图 7(b)所示, 所有策略的异质内容的平均响应时延都随着背景内容缓存空间比例的增加 (即前景内容缓存空间的减少) 而降低, NCPCR 策略的平均响应时延趋于 210 ms, CCS-AGP 策略和 AIEC-RSC 策略的异质内容的平均响应时延趋于 130 ms, DRCCS-Q 策略和本文所提出的 DRCCS-DQN 策略都趋于 100 ms, 且 DRCCS-DQN 策略在任何存储比例下都低于其他 4 种策略。这是由于随着背景内容缓存空间占比升高, CCS-AGP 策略和 AIEC-RSC 策略均缓存了较多的背景内容, 而在元宇宙应用中, 背景内容的数据量往

往大于前景内容,因此可有效降低异质内容的平均响应时延;同时,本文所提出的 DRCCS-DQN 策略通过 D2D 方式推荐和共享部分前景内容,满足了部分用户的内容需求,因此具有最低的平均响应时延。

3.4.5 前景内容数量和推荐接受率对缓存命中率的影响

本文比较了在前景内容数量分别为 40、70、100,推荐接受率分别为 30%、50%、70% 的条件下,所提出的 DRCCS-DQN 策略的前景内容缓存命中率,如表 3 所示。从表 3 可以看出,在前景内容数量不变的条件下,缓存命中率随着推荐接受率的上升而升高,这是由于推荐机制引导用户请求用户终端已缓存的前景内容,使得缓存命中率增加;在推荐接受率不变的条件下,缓存命中率随着前景内容数量的增加而升高,这说明了精准的用户偏好预测和个性化的内容推荐将提高边缘网络的缓存命中率。

表 3 缓存命中率对比

前景内容数量/个	推荐接受率		
	30%	50%	70%
40	0.938 8	0.957 9	0.962 4
70	0.947 9	0.960 6	0.974 4
100	0.960 1	0.967 3	0.975 8

4 结束语

在 MAR 场景中,应用传统的云服务器下载模式和边缘实时缓存交付模式处理多源动态的异质内容和复杂多变的用户请求会造成较低的数据访问速度和资源利用率。为了解决以上问题,本文提出了一种面向元宇宙 MAR 应用的异质内容主动缓存与个性化交付机制。仿真结果表明,所提机制在缓存命中率、前景内容和异质内容的平均响应时延方面均优于 NCPCR 策略、CCS-AGP 策略和 AIEC-RSC 策略。

本文主要关注元宇宙时延敏感业务的服务质量保障,在未来的研究工作中,笔者将继续深入研究异质内容之间的关联关系,并通过构建用户个性化的知识图谱和定价模型,提高用户服务体验和系统综合收益。

参考文献:

- [1] XU M R, NG W C, LIM W Y B, et al. A full dive into realizing the edge-enabled metaverse: visions, enabling technologies, and challenges[J]. IEEE Communications Surveys & Tutorials, 2023, 25(1): 656-700.
- [2] SIRIWARDHANA Y, PORAMBAGE P, LIYANAGE M, et al. A survey on mobile augmented reality with 5G mobile edge computing: architectures, applications, and technical aspects[J]. IEEE Communications Surveys & Tutorials, 2021, 23(2): 1160-1192.
- [3] HUANG Z H, FRIDERIKOS V. Mobility aware optimization in the metaverse[C]//Proceedings of the 2022 IEEE Globecom Workshops (GC Wkshps). Piscataway: IEEE Press, 2022: 80-86.
- [4] HUANG Z H, FRIDERIKOS V. Optimal mobility-aware wireless edge cloud support for the metaverse[J]. Future Internet, 2023, 15(2): 47.
- [5] SI P Y, ZHAO J, HAN H M, et al. Resource allocation and resolution control in the metaverse with mobile augmented reality[C]//Proceedings of the GLOBECOM 2022 - 2022 IEEE Global Communications Conference. Piscataway: IEEE Press, 2022: 3265-3271.
- [6] ZHANG L, WU X M, WANG F, et al. Edge-based video stream generation for multi-party mobile augmented reality[J]. IEEE Transactions on Mobile Computing, 2024, 23(1): 409-422.
- [7] CHEN X, LIU G Z. Energy-efficient task offloading and resource allocation via deep reinforcement learning for augmented reality in mobile edge networks[J]. IEEE Internet of Things Journal, 2021, 8(13): 10843-10856.
- [8] SOMESULA M K, ROUT R R, SOMAYAJULU D V L N. Greedy cooperative cache placement for mobile edge networks with user preferences prediction and adaptive clustering[J]. Ad Hoc Networks, 2023, 140: 103051.
- [9] WANG Y, YU T, SAKAGUCHI K. Context-based MEC platform for augmented-reality services in 5G networks[C]//Proceedings of the 2021 IEEE 94th Vehicular Technology Conference (VTC2021-Fall). Piscataway: IEEE Press, 2021: 1-5.
- [10] XU Z C, YUAN Z, LIANG W F, et al. Learning-driven algorithms for responsive AR offloading with non-deterministic rewards in metaverse-enabled MEC[J]. IEEE/ACM Transactions on Networking, 2024, 32(2): 1556-1572.
- [11] FU Y R, SALAÜN L, YANG X L, et al. Caching efficiency maximization for device-to-device communication networks: a recommend to cache approach[J]. IEEE Transactions on Wireless Communications, 2021, 20(10): 6580-6594.
- [12] 龙隆, 刘子辰, 陆在旺, 等. 移动边缘网络下服务缓存与资源分配联合优化策略[J]. 通信学报, 2023, 44(1): 64-74.
LONG L, LIU Z C, LU Z W, et al. Joint optimization strategy of service cache and resource allocation in mobile edge network[J]. Journal on Communications, 2023, 44(1): 64-74.
- [13] SEO Y J, LEE J, HWANG J, et al. A novel joint mobile cache and power management scheme for energy-efficient mobile augmented reality service in mobile edge computing[J]. IEEE Wireless Communications Letters, 2021, 10(5): 1061-1065.
- [14] ZENG F, ZHANG K W, WU L, et al. Efficient caching in vehicular edge computing based on edge-cloud collaboration[J]. IEEE Transactions on Vehicular Technology, 2023, 72(2): 2468-2481.
- [15] XU S Y, CHI J Y, WANG S, et al. AIEC-RSC: AI and edge collaboration empowered reliable service computing for high-speed mobile busi-

nesses[J]. IEEE Transactions on Services Computing, 2024, 17(1): 224-236.

- [16] CHI J Y, XU S Y, GUO S Y, et al. Federated learning empowered edge collaborative content caching mechanism for Internet of vehicles[C]// Proceedings of the NOMS 2022-2022 IEEE/IFIP Network Operations and Management Symposium. Piscataway: IEEE Press, 2022: 1-5.
- [17] SONG M Y, SHAN H G, FU Y R, et al. Joint user-side recommendation and D2D-assisted offloading for cache-enabled cellular networks with mobility consideration[J]. IEEE Transactions on Wireless Communications, 2023, 22(11): 8080-8095.
- [18] DANG T N, KIM K, KHAN L U, et al. On-device computational caching-enabled augmented reality for 5G and beyond: a contract-theory-based incentive mechanism[J]. IEEE Internet of Things Journal, 2021, 8(24): 17382-17394.
- [19] 李云, 高倩, 姚枝秀, 等. 移动边缘计算中智能服务编排和算网资源分配联合优化方法[J]. 通信学报, 2023, 44(7): 51-63.
LI Y, GAO Q, YAO Z X, et al. Joint optimization method of intelligent service arrangement and computing-networking resource allocation for MEC[J]. Journal on Communications, 2023, 44(7): 51-63.
- [20] WU Z Y, YAN D F. Deep reinforcement learning-based computation offloading for 5G vehicle-aware multi-access edge computing network[J]. China Communications, 2021, 18(11): 26-41.
- [21] CAO Z, ZHOU P, LI R, et al. Multiagent deep reinforcement learning for joint multichannel access and task offloading of mobile-edge computing in industry 4.0[J]. IEEE Internet of Things Journal, 2020, 7(7): 6201-6213.

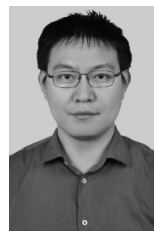
[作者简介]



徐思雅 (1988-), 女, 北京人, 博士, 北京邮电大学副教授、博士生导师, 主要研究方向为信息通信网络管理、SDN/NFV、移动边缘计算、人工智能等。



付琦梦 (2000-), 女, 河南许昌人, 北京邮电大学硕士生, 主要研究方向为移动边缘计算和人工智能。



郭少勇 (1985-), 男, 河北邢台人, 博士, 北京邮电大学教授、博士生导师, 主要研究方向为区块链、物联网等。