

## 基于深度强化学习的无线异构网络中继决策研究

周天玉, 官铮

引用本文

周天玉, 官铮. [基于深度强化学习的无线异构网络中继决策研究](#)[J]. 计算机科学, 2023, 50(11A): 221000088-5.

ZHOU Tianyu, GUAN Zheng. [Study on Relay Decision in Wireless Heterogeneous Networks Based on Deep Reinforcement Learning](#) [J]. Computer Science, 2023, 50(11A): 221000088-5.

---

## 相似文章推荐 (请使用火狐或 IE 浏览器查看文章)

Similar articles recommended (Please use Firefox or IE to view the article)

[一种融合CNN和Swin Transformer的医学显微图像分割模型](#)

Medical Microscopic Image Segmentation Model Based on CNN Structure and Swin Transformer  
计算机科学, 2023, 50(11A): 230200119-8. <https://doi.org/10.11896/jsjcx.230200119>

[基于GRU与自注意力网络的声源到达方向估计](#)

Sound Source Arrival Direction Estimation Based on GRU and Self-attentive Network  
计算机科学, 2023, 50(11A): 220900135-7. <https://doi.org/10.11896/jsjcx.220900135>

[一种安全高效的去中心化移动群智感知激励模型](#)

Safe Efficient and Decentralized Model for Mobile Crowdsensing Incentive  
计算机科学, 2023, 50(11A): 221000184-10. <https://doi.org/10.11896/jsjcx.221000184>

[面向工业物联网的轻量级群组密钥协商方案](#)

Lightweight Group Key Agreement for Industrial Internet of Things  
计算机科学, 2023, 50(11A): 230700075-10. <https://doi.org/10.11896/jsjcx.230700075>

[基于替代模型的批量零阶梯度符号算法](#)

Batch Zeroth Order Gradient Symbol Method Based on Substitution Model  
计算机科学, 2023, 50(11A): 230100036-6. <https://doi.org/10.11896/jsjcx.230100036>

# 基于深度强化学习的无线异构网络中继决策研究

周天玉 官 铮

云南大学信息学院 昆明 650500

(zty@mail.ynu.edu.cn)

**摘 要** 在物联网大规模多用户场景中,远端节点需通过中继接入网络。为解决中继在异构接入技术环境下的自适应接入控制问题,提出一种基于深度强化学习的智能中继接入控制策略,将中继对远端用户数据的收发过程视为一个部分可观察马尔可夫决策过程,通过动态决策中继工作状态,以实现最大化系统的总吞吐量和节点公平性目标。首先,建立具有中继的无线异构网的上行链路模型,以提高系统总吞吐量为优化目标,建立中继动态决策优化模型;其次,构建含有 LSTM 隐藏层的深度 Q 网络(DQN)作为行为状态值函数,以优化系统总吞吐量。测试结果表明深度强化学习无线异构网络中继决策方案(DRL-RAP)可在确保原有用户服务质量的前提下,为远端用户提供网络接入,系统总吞吐量在原有网络基础上显著提高,吞吐量最大可提高 30%。

**关键词:** 物联网;无线异构网络;深度强化学习;中继智能决策;神经网络

**中图法分类号** TN925

## Study on Relay Decision in Wireless Heterogeneous Networks Based on Deep Reinforcement Learning

ZHOU Tianyu and GUAN Zheng

School of Information Science & Engineering, Yunnan University, Kunming 650500, China

**Abstract** For large-scale multi-user scenarios of the Internet of Things, remote nodes need to access the network through relay. In order to solve the adaptive access control problem of relay in heterogeneous access technology environment, an intelligent relay access control strategy based on deep reinforcement learning is proposed, which regards the transmission and reception process of relay to remote user data as a partially observable Markov decision process, and dynamically decides the relay working state to maximize the total system throughput and node fairness. Firstly, the uplink model of wireless heterogeneous network with relay is established. With the goal of improving the total throughput of the system, the dynamic decision optimization model of relay is established. Secondly, a deep Q network(DQN) with LSTM hidden layer is constructed as a behavior state value function to optimize the total system throughput. Test results show that DRL-RAP can provide network access for remote users on the premise of ensuring the original user's quality of service. The total throughput of the system is significantly improved on the basis of the original network, and the maximum throughput can be increased by 30%.

**Keywords** Internet of Things, Wireless heterogeneous network, Deep reinforcement learning, Relay intelligent decision, Neural network

## 1 引言

异构性是未来无线网络的主要特征,接入技术异构是指在一个通信模型中存在两种以上通信协议。现有接入控制技术可分为两种——传统优化方法及人工智能优化方法。

在传统方法上,为解决关联调度及网络负载均衡拥塞问题,文献[1]采用提高密集无线异构网络上行链路传输可靠性的网络关联方法,双重关联大基站和小区基站,从而提高上行链路传输成功概率;文献[2]在 IEEE802.15.4 信标启用模式下控制占空比的高异构无线传感器网络中分组投递率的方法,调整 MAC 层参数以控制数据传输,降低了分组传输率;文献[3]提出可持续地、公平地提高用户总最小效用和最差最小效用的选择算法解决用户动态处理无线接入点(Access

Point, AP)决策问题;文献[4]提出了监督学习获得网络的带宽分配值的无线异构网络节点联合接入选择和带宽分配算法,用户择优选择最佳接入网络。此外考虑到基站对节点控制及资源分配文献[5]通过中央控制节点分配频谱资源的方法均衡负载,使各个无线接入技术保持均衡的负载比,高效利用无线资源。由于无线网络信道具有特殊性,因此需要良好的信道编码。文献[6]基于流水线网络信道编码的多径传输控制协议算法解决无线异构网络中更新调度及网络拥塞问题。

对于强化学习方法,文献[7]采用“无模型”的方法,用深度强化学习设计高效的异构网络协议;文献[8]针对有噪声的不完善信道提出了分布式强化学习 MAC 协议,提出反馈回复机制来恢复学习过程中丢失的信息;文献[9]针对无线异构

基金项目:国家自然科学基金(61761045);云南省科研基金资助项目(202201AT070167);云南大学科研项目(2021Y189)

This work was supported by the National Natural Science Foundation of China (61761045), Research Foundation of Yunnan Province (202201AT070167) and Research Project of Yunnan University(2021Y189).

通信作者:官铮(gz\_627@sina.com)

网络频谱利用问题,通过学习现有节点的传输方式,实现了有效利用多通道频谱;文献[10]基于多智能体决斗深度 Q 网络结合分布式协议算法,解决了两层异构无线网络设备联合关联、频谱和功率分配;文献[11]基于固定或时变异构网络提出采用 DenseNet 的双深度 Q 网络算法选择最优信道,减少碰撞率。

上述方法均针对节点直接接入网络。在实际通信中,为满足大规模网络覆盖,提高频谱利用率,远端节点可通过中继转发实现对本地 AP 的数据接入,然而,随机接入机制下,中继节点的数据收发增加了本地数据接入的碰撞概率,降低了吞吐量。

针对无线异构网络中继,中继的改进与优化受到了较多研究者的关注且他们探索了多种基于不同技术基础的方法。为降低误码率,提高计算效率<sup>[12]</sup>,采用了带有相移键控调制的异构网络编码,但对于找到相移键控调制中的 SFS 的数量和位置是较为困难的。文献[13]研究了设备到设备的中继的负载均衡策略,将负载过大的单元的负载转移到负载较小的单元,提出负载均衡算法。文献[14]针对中继节点部署问题,提出基于欧氏距离最小生成树的部署算法。在多中继协作网络中,中继节点选择是重要问题<sup>[15-17]</sup>。作者将中继节点选择问题阐述为马尔可夫决策过程,利用强化学习算法找到中继节点最优选择路径,从多个中继节点中选择最优中继。文献[18]针对中继闲置问题提出 HTT(Harvest then Transmit)功率消耗方案,提升了多中继协作网络中性能。

在大规模物联网用户场景下,为保证传输高效公平,节点接入网络需进行实时动态决策并且不具有抢占性,中继的收发决策应是实时、快速的。为解决这一问题,本文在异构网络接入中继并采用深度强化学习进行实时动态决策,将深度强化学习技术用于中继的收发决策过程。深度强化学习是一种端到端的感知控制,符合网络通信节点接入网络的端到端机制,因此本文基于深度强化学习,设计了一种无线异构网络中继决策方案,为方便阐述,本文将这种方案称为 DRL-RAP (Deep Reinforcement Learning-Relay Access Point)。

## 2 系统模型

### 2.1 网络系统模型

本文采用如图 1 所示的系统模型结构,在无线异构接入控制网络中,节点采用不同访问控制协议,通过相同公共无线频段向 AP 发送数据包。

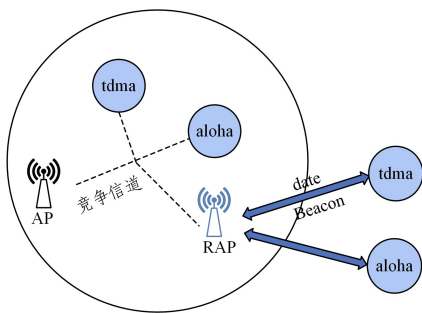


图 1 系统模型

Fig. 1 System model

本文在无线异构网络中加入深度强化学习算法驱动的中继 AP(Relay Access Point, RAP),从而向 AP 接入远端用户。

远端节点首先将数据包发送到 RAP,再由 RAP 接入 AP。每个节点以基本时隙的方式在共享频谱上传输数据包。

本文考虑了 2 个不同类型的节点:TDMA 和 ALOHA 的变体  $q$ -ALOHA。中继节点具有载波感知以及接收和转发数据包的功能,ALOHA 节点具备载波感知功能但由于没有设定退避数,所以 ALOHA 在单位时隙内以固定概率  $q$  在时隙内传输数据包,TDMA 节点在特定时隙中以固定周期的方式进行传输,远端的 TDMA 节点和 ALOHA 节点以轮流的方式向 RAP 发送数据包。RAP 的工作过程分为两个阶段,第一阶段,RAP 调用远端节点接收数据包并返回是否成功接收数据包的反馈报文;第二阶段,执行载波监听的同时接入 AP 转发数据包。同时本文考虑了原有节点传输时会对 RAP 接收数据产生干扰。

系统的目标是节点之间实现高效的数据传输、频谱资源公平共享,最大化系统的总吞吐量。针对这一目标,本文采用  $\alpha$ -公平<sup>[19]</sup>来实现。

### 2.2 目标函数

针对异构无线网络上行链路, $N$  个节点中有  $L$  类异构节点,每类有  $K$  个节点,下面给出吞吐量的定义,每个节点  $i$  的吞吐量为:

$$X(i) = \frac{N}{TU} \quad (1)$$

其中, $N$  为成功发送分组数, $TU$  代表信道占用时间。

系统信道利用率  $\eta$  定义为:

$$\eta = \frac{TU}{TU + TI} \quad (2)$$

其中, $TI$  表示信道空闲时间。

数据成功率  $S$  定义为:

$$S = \frac{N}{M} \quad (3)$$

其中, $M$  表示发送分组数量。

根据文献[19],对于  $\alpha$  公平函数, $\alpha=0$  表示最大化所有网络的总吞吐量; $\alpha=1$  表示比例公平; $\alpha \rightarrow \infty$  表示实现最大最小公平,本文采用最大最小公平性约束。所以节点  $i$  的公平效用函数定义为:

$$f_{\alpha}(X(i)) = \begin{cases} \log(X(i)), & \alpha=0 \\ \frac{(X(i))^{1-\alpha}}{1-\alpha}, & \text{其他} \end{cases} \quad (4)$$

每类节点的吞吐量为:

$$G = \sum_i \sum_j^{K_i} f_{\alpha} X(i) \quad (5)$$

$$\arg \max G; \sum_{i=1}^n G \quad (5.1)$$

$$\text{s. t. } X(i) > 0, \arg \max G < 1 \quad (5.2)$$

## 3 实现方法

### 3.1 深度强化学习

深度强化学习(DRL)致力于解决序贯决策问题,即通过连续不断地决策来实现最终的目标,解决游戏、机器人控制、无线通信、网络管理与控制等复杂决策问题。在大规模物联网用户场景下,RAP 对远端用户的数据收发的决策问题可以视作一个序贯决策问题。图 2 为深度强化学习架构。在强化学习中,神经网络通过与环境交互更新状态并以一定概率选择动作,环境给予特定奖励<sup>[20]</sup>。

在某一个时间步长  $t$ , 给予环境一个状态, 代理(agent)通过策略采取行动, 策略将环境状态映射到 agent 的动作中, 在 agent 的动作到达环境后给予 agent 一个奖励反馈, 以此评价 agent 在这一时间步长的动作表现, 同时传递下一个时间步长状态  $s_{t+1}$ <sup>[20]</sup>。DQN 在稳定环境中遵循马尔可夫决策过程, Q 值收敛于最优解<sup>[21]</sup>。中继对远端用户数据的收发过程同样可视为一个部分可观察马尔可夫决策过程, 在中继中加入深度 Q 网络进行自适应决策, 提高中继决策效率, 将空闲时隙利用起来。深度 Q 网络采用神经网络模型来逼近动作价值 Q 函数, 深度 Q 网络的输入为状态  $s$ , 网络的输出为不同动作的 Q 值,  $\{Q(s, a; \theta) | a \in A\}$ ,  $\theta$  是神经网络权重参数,  $A$  是 agent 动作的集合。同时 DQN 还结合了经验回放和目标网络两个关键技术, 以解决收敛的不稳定性问题。训练时, 参数  $\theta$  通过最小化下面的损失函数进行更新:

$$L(\theta) = \frac{1}{N_{E, \epsilon \in E}} \sum [(r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a', \theta') - Q(s_t, a_t; \theta))^2] \quad (6)$$

其中,  $E$  表示所存储的经验集合,  $r$  为 agent 学习过程所获得的奖励,  $\gamma$  为训练的折扣因子。经验回放: agent 在与环境交互过程中收集不同时间步长的经验并将其存储在经验缓冲器中通过随机采样多个经验样本放到一起进行训练更新网络参数。在训练过程中打破时间的相关性, 混合越来越多的经验, 并且一些比较好的经验能够被多次使用以此提高训练效果。目标网络: DQN 用一个单独的目标神经网络来计算式(6)中的  $r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a', \theta')$ 。

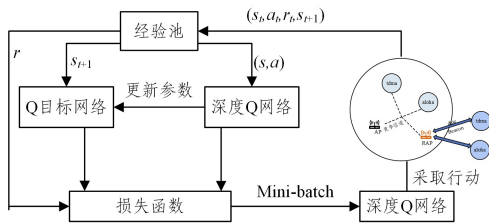


图2 实现 DRL-RAP 的架构

Fig. 2 Implement DRL-RAP architecture

### 3.2 设置 DRL-RAP

代理: 本文将 RAP 视为一个代理, RAP 通过与通信网络交互, 自己做收发决策。动作: 代理在基本时隙上所执行的动作定义为  $a_t \in \{\text{接收, 发送, 等待}\}$ 。用 0 代表等待; 1 代表发送; 2 代表接收。观察值: 代理在基本时隙  $t$  执行一个动作后会返回一个观察值, 如果  $a_t = 0$ , 返回的观察值为信道状态,  $o_t = I$  或  $o_t = B$ ,  $I$  表示信道空闲,  $B$  表示信道繁忙;  $a_t = 1$ , 返回的观察值为  $o_t = S$  或  $o_t = F$ ,  $S$  代表数据包转发成功,  $F$  代表数据包转发失败;  $a_t = 2$ , 返回的观察值为  $o_t = L$  或  $o_t = F$ ,  $L$  表示代理接收数据包成功,  $F$  表示代理接收数据包失败。奖励: 根据代理所执行的动作定义了不同的奖励, 以此评价代理动作的表现, 并假设代理知道其他所有用户的奖励, 如果返回的观察值  $o_t = L$ , 则  $r_t = 0.5$ ; 如果  $o_t = S$ , 则  $r_t = 1$ 。状态: 代理在时隙  $t$  的状态定义为  $s_t$ , 代理在时隙  $t+1$  的状态包含过去时隙的状态, 定义为  $s_{t+1} = [s_t, a_t, o_t, r_t^1, \dots, r_t^n]$ 。

具体流程如算法 1 所示。

#### 算法 1 DRL-RAP algorithm

Input: the current state  $s_t$

Output: action  $a_t$

Initialize state of the agent, the hyperparameters of DQN

Initialize the experience

1. for  $t=1, 2, \dots$  in DRL-RAP do

2. Decide to the agent action  $a_t$ ;

3. Get to know the observation  $o_t$ ;

4. If  $o_t = L$

5.  $r_t = 0.5$ ;

6. elif  $o_t = S$

7.  $r_t = 1$ ;

8. elif  $o_t = F$

9.  $r_t = 0$ ;

10. end if

11. Get  $st+1 = [s_t, a_t, r_t^1, \dots, r_t^n, st+1]$ ;

12. Store  $et = (s_t, r_t^1, \dots, r_t^n, st+1)$  into the experience buffer;

13. end for

14. Return  $o_t, r_t^1, \dots, r_t^n$

本文采用 FNN 作为神经网络模型, 网络包括两层隐藏层 LSTM 层和 FNN 层, 每层的神经元数量为 64, ReLU 作为激活函数, 采用 RMSProp 小批量梯度下降。网络的输入是当前时间步长  $t$  的状态, 输出是不同动作的近似 Q 值。

实验中的超参数如表 1 所列, 为了充分利用过去时间步长  $t$  的经验, 我们设定状态长度为 20,  $\epsilon$  初始值设为 1 直至衰减为 0.05, 折扣因子  $\gamma$  设为 0.9, 为减小学习过程中的震荡将学习率设为 0.01, 经验池大小为 1000, mini-batch 大小为 64, 训练 200 步进行一次网络更新。

表 1 DRL-RAP 超参数说明

Table 1 DRL-RAP hyperparameter description

参数	大小
状态长度	20
$\epsilon$ 策略取值	1-0.05
折扣因子 $\gamma$	0.9
经验池大小	1000
小批量大小	64
网络更新频率	1/200
学习率	0.01

图 3 为算法的收敛曲线, 收敛指标为平均奖励, 算法在训练到大概 750 步时平均奖励趋于一个平均值, 图 3 表明本文所提出的方法是收敛的。

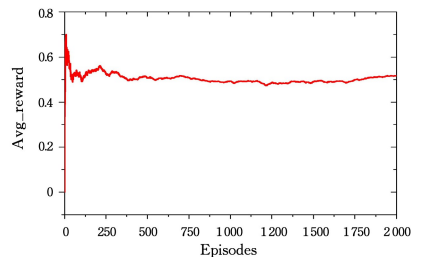


图3 DRL-RAP 收敛曲线

Fig. 3 DRL-RAP convergence curve

## 4 实验及分析

本次实验的目的是测试和验证本文的系统性能。深度强化学习仿真参数已经在表 1 中列出。实验设定近端节点两个, 远端节点两个, 近端和远端节点均包括一个 q-ALOHA 节点和一个 TDMA 节点。实验通过改变 q-ALOHA 的发送概率和 TDMA 节点的周期调整系统的总负载和远端负载。本文对比了最大最小公平性下接入 DRL-RAP 对原有用户以及

系统性能的影响。

图 4 给出了增加系统总负载对系统总平均吞吐量及原有用户平均吞吐量的影响,下面给出具体的负载数。接入中继时,总负载为 0.2 时  $q=0.05$ , TDMA 周期为 20 个时隙;总负载为 0.4 时  $q=0.1$ , TDMA 周期为 10 个时隙;总负载为 0.6 时  $q=0.2$ , TDMA 周期为 10 个时隙;总负载为 0.8 时  $q=0.2$ , TDMA 周期为 5 个时隙;总负载为 1 时  $q=0.3$ , TDMA 周期为 5 个时隙。对于原有用户,在不同的总负载下  $q$  分别为  $q=0.1, q=0.2, q=0.3, q=0.4, q=0.5$ ; TDMA 周期分别为 10 个时隙、5 个时隙、3 个时隙、2.5 个时隙、2 个时隙。

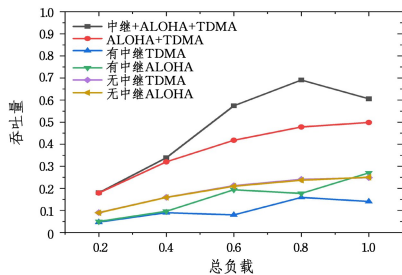


图 4 增大总负载时的吞吐量变化

Fig. 4 Throughput change when total load is increased

图 4 表明,接入中继使得远端用户得以接入 AP,在原有近端用户的基础上系统总平均吞吐量有所增大。在总负载到达 0.8 左右前虽然原有近端节点受到一定程度影响但是平均总吞吐量逐渐增大,系统总平均吞吐量与总负载呈现正比发展趋势,显然中继的接入使得原有网络的空闲时隙得到利用,表明了系统的有效性。总负载到达 0.8 左右后由于系统容量接近饱和,碰撞的概率增大,平均吞吐量开始下降,可见接入中继后系统增加的吞吐量最大可达 0.3。

图 5 为增大系统总负载对系统信道利用率的影响,接入远端用户后,信道空闲的时隙得以利用,频谱资源得到充分利用,提高了信道利用率。图 6 为增大系统总负载对系统数据发送成功率的影响,接入中继后随着信道利用率的增加,时延和碰撞急剧增大,造成数据发送成功率有所降低,但保证了 60% 左右的成功率。

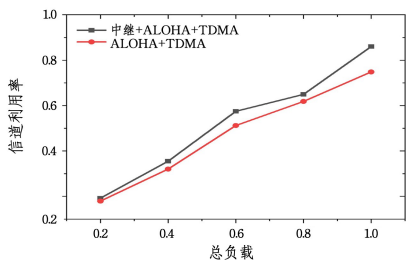


图 5 增大总负载时的信道利用率变化

Fig. 5 Change of channel utilization when the total load is increased

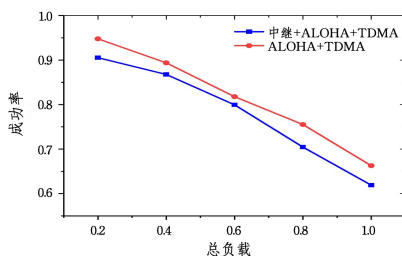


图 6 增大总负载时的成功率变化

Fig. 6 Change of success rate when total load is increased

最后为了验证增加远端负载对原有用户服务质量的影响,实验模拟了在不同远端负载下的原有用户平均吞吐量的曲线图,图 7 为在原有用户负载不变的情况下增大远端用户负载对原有用户的影响。实验设置原有用户  $q=0.2$ , TDMA 周期为 5 个时隙;在不同的远端用户负载下  $q$  分别为 0.05, 0.1, 0.2, 0.2, 0.3, 0.3; TDMA 周期分别为 20 个时隙、10 个时隙、10 个时隙、5 个时隙、5 个时隙、3 个时隙。由图 7 可见,在尽量确保原有用户服务质量的前提下,增加远端负载 RAP 后吞吐量随之逐渐增加,表明了 DRL-RAP 的有效性。远端负载增大到 0.4 左右,中继吞吐量达到最高,并且原有用户的服务质量没有太大损耗。继续增加远端负载,中继吞吐量下降后趋于稳定,原因是并发数过大,数据碰撞概率和时延急剧增大。

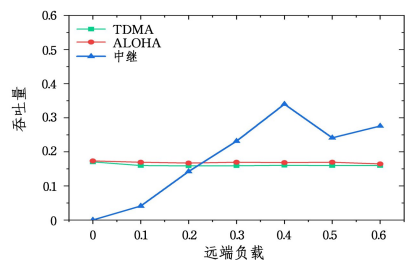


图 7 增大远端负载对原有用户的影响

Fig. 7 Effect of increasing remote load on original users

**结束语** 本文针对现有异构无线网络中继算法忽略对远端用户数据的收发过程这一问题,提出了 DRL-RAP,实现了高效公平的频谱利用,最大化系统总吞吐量。在 RAP 中加入深度强化学习 DQN 算法,将 RAP 视作一个代理,通过与系统交互快速自主决策收发,提出了 DRL-RAP 并讨论了 DRL-RAP 与其他异构节点的共存性及其产生的影响。DRL-RAP 的加入会对其他节点有所损耗但提高了系统的总体性能,提高了信道利用率,增加了用户容量,使系统总吞吐量最多提高 30%。此外,本文提出的 DRL-RAP 算法对同构网络仍然适用。在未来的研究工作中,将继续针对无线异构网络中的终端接入网络切换及选择进行研究。

## 参考文献

- [1] KIM D, POPOVSKI P. Reliable uplink communication through couple cassociation in wireless weterogeneous networks[J]. IEEE Wireless Communications Letters, 2017, 5(3): 312-315.
- [2] TOMITA T K, KOMURO N. Duty-Cycle Control Achieving High Packet Delivery Ratio in Heterogeneous Wireless Sensor Networks[C]// 2019 IEEE 8th Global Conference on Consumer Electronics(GCCE). Osaka, Japan, 2019: 1164-1167.
- [3] KOBAYASHI H, KAMEDA E, TERASHIMA Y, et al. A strategy for AP selection with mutual concessions in sustainable heterogeneous wireless networks[C]// 2016 IEEE Region 10 Conference(TENCON 2016). IEEE, 2016.
- [4] XU C Q, WANG P, XIONG C S, et al. Pipeline network coding-based multipath data transfer in heterogeneous wireless networks[J]. IEEE Transactions on Broadcasting, 2016, 63(2): 376-390.
- [5] GEN L, YU H W, GUO X X, et al. Joint access selection and bandwidth allocation algorithm supporting user requirements and preferences in heterogeneous wireless networks[J]. IEEE

- Access,2019(7):23914-23929.
- [6] ZARIN N,AGARWAL A. A centralized approach for load balancing in heterogeneous wireless access network[C]// 2018 IEEE Canadian Conference on Electrical & Computer Engineering(CCECE). IEEE,2018.
- [7] YU Y D,LIEW S C,WANG T T. Carrier-sense multiple access for heterogeneous wireless networks using deep reinforcement learning[C]// 2019 IEEE Wireless Communications and Networking Conference Workshop(WCNCW). IEEE,2019.
- [8] YU Y D,LIEW S C,WANG T T. Multi-agent deep reinforcement learning multiple access for heterogeneous wireless networks with imperfect channels[J]. IEEE Transactions on Mobile Computing,2021,21(10):3718-3730.
- [9] YE X W,YU Y D,FU L Q, et al. Multi-Channel Opportunistic Access for Heterogeneous Networks Based on Deep Reinforcement Learning[J]. IEEE Transactions on Wireless Communications,2022,21(2):794-807.
- [10] CHENG Q,WEI Z,YUAN J. Deep reinforcement learning-based spectrum allocation and power management for IAB networks[C]// 2021 IEEE International Conference on Communications Workshops(ICC Workshops). IEEE,2021.
- [11] KANG Z. Deep Reinforcement Learning-Based Dynamic Multi-Channel Access for Heterogeneous Wireless Networks with DenseNet[C]// 2021 IEEE/CIC International Conference on Communications in China(ICC Workshops). IEEE,2021.
- [12] ARUNACHALA C,BUCH S D,RAJAN S. Wireless bidirectional relaying using physical layer network coding with heterogeneous PSK modulation[J]. IEEE Transactions on Vehicular Technology,2018,67(3):2335-2344.
- [13] FAN J,YAO L,WANG B, et al. A relay-aided device-to-device-based load balancing scheme for multitier heterogeneous networks[J]. IEEE Internet of Things Journal,2017,4(5):1537-1551.
- [14] CHEN N,LI Z J,JIANG S X. Relay node deployment algorithm in heterogeneous wireless networks [J] Journal of Computer Science,2016,39(5):905-918.
- [15] KIM H,UJII T,UMEBAYASHI K. Relay nodes selection using reinforcement learning[C]// 2021 International Conference on Artificial Intelligence in Information and Communication(ICAI-IC). 2021.
- [16] HUANG C,CHEN G,GONG Y, et al. Joint buffer-aided hybrid-duplex relay selection and power allocation for secure cognitive networks with double deep Q-network[J]. IEEE Transactions on Cognitive Communications and Networking,2021,7(3):834-844.
- [17] SU Y,LU X,ZHAO Y, et al. Cooperative communications with relay selection based on deep reinforcement learning in wireless sensor networks[J]. IEEE Sensors Journal,2019,19(20):9561-9569.
- [18] SHAN Y F,JIANG R,XU Y Y, et al. A power consumption scheme for full duplex multi relay cooperative SWIPT network [J]. Computer Science,2022,49(7):280-286.
- [19] MO J,WALRAND J. Fair end-to-end window-based congestion control[C]// Performance and Control of Network Systems II, International Society for Optics and Photonics,1998.
- [20] DONG H. Deep Reinforcement Learning: Foundation, Research and Application[M]. Beijing:Electronic Industry Press,2021.
- [21] WANG H N,LIU T,ZHANG Y Y, et al. A Review of deep reinforcement learning[J]. Frontiers of Information Technology & Electronic Engineering,2020,21(12):63-82.



**ZHOU Tianyu**, born in 1999, master. Her main research interests include deep reinforcement learning and intelligent mobile communication.



**GUAN Zheng**, born in 1982, Ph.D, associate professor, master supervisor, is a member China Computer Federation. Her main research interests include wireless sensor networks, network access technology, and performance analysis and optimization of polling systems.