

# Age-minimal Multicast by Graph Attention Reinforcement Learning

Yanning Zhang  
Zhejiang University

Guocheng Liao  
Sun Yat-Sen University

Shengbin Cao  
University of Macau

Ning Yang  
Chinese Academy of Sciences

Meng Zhang  
Zhejiang University

**Abstract**—*Age of Information (AoI)* is an emerging metric used to assess the timeliness of information, gaining research interest in real-time multicast applications such as video streaming and metaverse platforms. In this paper, we consider a dynamic multicast network with energy constraints, where our objective is to minimize the expected time-average AoI through energy-constrained multicast routing and scheduling. The inherent complexity of the problem, given the NP-hardness and intertwined scheduling and routing decisions, makes existing approaches inapplicable. To address these challenges, we decompose the original problem into two subtasks, each amenable to reinforcement learning (RL) methods. Subsequently, we propose an innovative framework based on graph attention networks (GATs) to effectively capture graph information with superior generalization capabilities. To validate our framework, we conduct experiments on three datasets including a real-world dataset called AS-733, and show that our proposed scheme reduces the average weighted AoI by 62.9% and reduces the energy consumption by at most 72.5% compared to baselines.

## I. INTRODUCTION

### A. Background and Motivations

Real-time multicast applications, such as video streaming [1] and intelligent transportation systems [2], have experienced significant growth in recent times. These applications require timely updates to ensure accurate and up-to-date information availability for critical tasks like decision-making and system control. While delay is a commonly used metric in traditional networks [3], it is now recognized that ensuring timely updates is distinct from simply minimizing delay [4]. Consequently, there is a need for a metric that captures the timeliness aspect of information dissemination. The concept of *Age of Information (AoI)* has emerged as a promising metric in various domains, including learning and network protocols [5]. AoI quantifies the freshness of information possessed by a monitor about a specific entity or process, which has been identified as a suitable metric for evaluating the performance of multicast networks [6], making it particularly relevant for real-time multicast applications. It has been shown that AoI is the most important metric for evaluating the Quality of Service (QoS) in some scenarios [7].

Multicast, a vital communication paradigm in networks, facilitates the efficient dissemination of information from a

source to multiple destinations. Recent research has extensively investigated the applications of multicast in various scenarios. For instance, optimizing the multicast Quality of Experience (QoE) is crucial for enhancing video streaming sessions [8]. One of the primary problems lies in the routing process, which entails determining the optimal paths from the source to destinations. The routing problem falls within the domain of Combinatorial Optimization (CO) problems, such as the Steiner tree problem [9], known for their NP-hardness, rendering them computationally demanding to solve in large-scale networks.

Due to the inherently distributed feature of network systems, only local information is available for a centralized controller, making it challenging to optimize the routing process. Fortunately, recent advances in Software Defined Networking (SDN) enabled the intelligent control of network devices [10]. SDN is a network architecture that separates the control plane from the data plane, where the control plane is centralized and is responsible for managing network resources and programming the network dynamically. The centralized controller has a global view of the network by monitoring and collecting the real-time network state (e.g., packets) and configurations. The above features ensure that the solutions given by intelligent algorithms (e.g., AI-based algorithms) can be implemented in real-world scenarios.

Furthermore, real-world networks often operate under energy constraints, where the overall energy consumption of the network is limited. This introduces additional complexity to the problem, as there exists a trade-off between energy consumption and AoI [11]. Consequently, certain existing algorithms (e.g., [12]) become inapplicable. Some studies have proposed scheduling algorithms to optimize AoI [6]. However, they tend to overlook the routing problem, which is also crucial for AoIs. Hence, there is a clear demand for a novel multicast scheme that offers a comprehensive solution to optimize AoI in energy-constrained multicast networks.

In this paper, we aim to answer the following main question: *How should one design multicast scheduling and routing algorithms to make the optimal tradeoffs between AoI and energy consumption?*

## B. Key Challenges

We now summarize the key challenges of answering the above problem as follows:

- 1) **Coupled Decision Variables.** In a long-term multicast process, multicast scheduling and routing are intertwined, both exerting impact on the AoIs of destinations.
- 2) **Energy Constraints.** Real-world applications often impose energy constraints on the network. This introduces a trade-off between energy consumption and AoI.
- 3) **Hidden Graph Information.** Traditional methods are inefficient when extracting relevant graph features due to the non-Euclidean nature of graphs.
- 4) **NP-hardness.** Multicast routing algorithms usually fall into the realm of CO problems, which are computationally intractable for large-scale networks.

One promising solution approach to overcoming the above challenges includes Reinforcement Learning (RL) methods, which are promising for approximating the optimal solutions of NP-hard problems via learning from environments [13]. Hence, we further prompt the following question: *How should one design an RL framework to address the problem of coupled decision variables and capture the graph information of a multicast network?*

## C. Solution Approach and Contributions

To answer the above questions, we summarize our solution approach and main contributions as follows:

- **Joint Multicast Scheduling and Routing Problem:** We tackle the complex task of joint multicast routing and scheduling, accounting for energy constraints and possible network dynamics. *To the best of our knowledge, this is the first work to consider the problem of minimizing AoI in joint multicast scheduling and routing.*
- **Hierarchical RL Framework.** To address challenges 1, 2 and 4, we decompose the original problem into two subtasks and introduce a hierarchical RL framework. The first subtask involves selecting destinations, while the second subtask generates multicast trees.
- **A Novel Graph Attention Network.** To overcome challenges 3 and 4, we propose a novel kind of Graph Attention Network (GAT) to extract graph information based on the attention mechanism.
- **Performance Evaluation.** We validate our approach on three datasets, including a real-world dataset called AS-733. TGMS outperforms other baselines and achieves an average AoI reduction of 62.9% while maintaining the energy consumption within the constraint.

## II. SYSTEM MODEL

### A. Network Description

We consider a multicast network that operates for an infinite horizon of slotted time. See an illustrative example in Fig. 1. The network topology at time  $t$  is denoted as  $\mathcal{G}_t = \{\mathcal{V}_t, \mathcal{E}_t\}$ , where  $\mathcal{V}_t$  represents the set of nodes and  $\mathcal{E}_t$  represents the

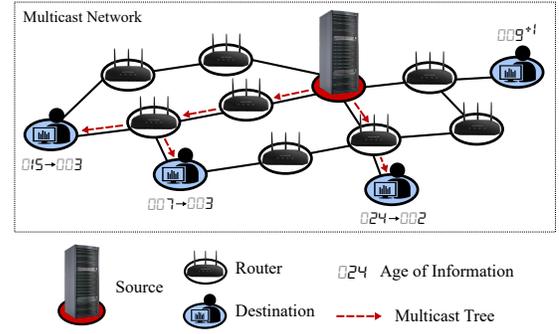


Fig. 1. An Example of a Multicast Network. The nodes are connected by links with different costs. At the beginning of each time slot, the source generates update packets, which are then forwarded to destinations by routers.

set of undirected links. The nodes in the network can be categorized into three distinct types:

- **The Source Node.** A source node generates updates for a multicast group<sup>1</sup> continuously without restrictions.
- **Router Nodes.** Router nodes are responsible for forwarding update packets to destinations.
- **Destination Nodes.** Destination nodes are expected to receive update packets from the source node. The set of destinations at time  $t$  is denoted as  $\mathcal{U}_t \subset \mathcal{V}_t$ .

A packet can be transmitted from node  $i$  to node  $j$  if the link  $(i, j)$  exists in  $\mathcal{E}_t$ , which takes one time slot with an energy cost of  $C_{i,j}$ . The multicast process can be described as follows: at the beginning of each time slot, the source node generates multiple update packets, which are then transmitted between router nodes, eventually reaching the destination nodes. Note that packets traveling through different transmission paths may not arrive at the destination simultaneously. To analyze the AoI of destinations, we initially define the AoI of update packets. Let  $\mathcal{P}_t = \{0, 1, \dots, p, \dots\}$  denotes the set of packets at time  $t$ , the AoI of packet  $p$  can be defined as:

$$\hat{A}_p(t) = t - t_p, \forall t \geq t_p, \quad (1)$$

where  $t_p$  denotes the time when packet  $p$  is generated. That is,  $\hat{A}_p(t)$  grows linearly with time. Suppose that each destination can receive only one packet during a single time slot. The AoI of a destination  $u \in \mathcal{U}_t$  is defined as:

$$A_u(t) = \begin{cases} \hat{A}_p(t), & \text{if } d_{u,p}(t) = 1, \\ A_u(t-1) + 1, & \text{otherwise,} \end{cases} \quad (2)$$

where  $d_{u,p}(t)$  is an indicator that represents whether packet  $p$  arrives at destination  $u$  at time  $t$ . If packet  $p$  arrives at destination  $u$  at time  $t$ ,  $d_{u,p}(t) = 1$ ; otherwise,  $d_{u,p}(t) = 0$ . As shown in Fig. 2, if a destination  $u$  receives a packet at time  $t$ ,  $A_u(t)$  is updated to be the AoI of the received packet at that time. Otherwise,  $A_u(t)$  grows linearly with time.

The AoIs of destinations are closely intertwined with the routing decisions. In this paper, we adopt the use of multicast

<sup>1</sup>This scenario can also be extended to accommodate multiple multicast groups, where each group has its source node.

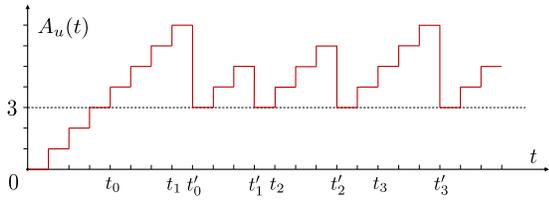


Fig. 2. An Example of the AoI.  $t_k$  denotes the time when the  $k$ -th packet is generated while  $t'_k$  denotes the time when it arrives. The AoI of destination  $u$  is updated to be the AoI of the received packet (3 in this case) at time  $t'_k$ , otherwise, it grows linearly with time.

trees to represent routing decisions as commonly done in classical multicast routing problems [9]. A tree is defined as a connected acyclic undirected graph. Accordingly, we define a multicast tree  $\mathcal{T}_t$  of network  $\mathcal{G}_t$  as  $\mathcal{T}_t = \{\mathcal{V}_t^{\mathcal{T}}, \mathcal{E}_t^{\mathcal{T}}\}$ , where  $\mathcal{V}_t^{\mathcal{T}} \subseteq \mathcal{V}_t$  denotes the set of included nodes and  $\mathcal{E}_t^{\mathcal{T}} \subseteq \mathcal{E}_t$  denotes the set of links. To establish the relation between the AoIs and multicast trees, we refer to the following lemma:

**Lemma 1** ([14]). Any two vertices in a tree can be connected by a unique simple path.

Considering the multicast process described and the store-and-forward mechanism [15], we can derive that the time required for a packet to reach destination  $u$  is equal to the number of hops between the source and destination  $u$  in a multicast tree  $\mathcal{T}_t$ , denoted as  $h_{\mathcal{T}_t}(u)$ . Then Eq. (2) can be rewritten as:

$$A_u(t + h_{\mathcal{T}_t}(u)) = \begin{cases} h_{\mathcal{T}_t}(u), & \text{if } u \in \mathcal{V}_t^{\mathcal{T}}, \\ A_u(t + h_{\mathcal{T}_t}(u) - 1) + 1, & \text{otherwise.} \end{cases} \quad (3)$$

Here,  $A_u(t + h_{\mathcal{T}_t}(u))$  will be updated to exactly  $h_{\mathcal{T}_t}(u)$  when a packet generated at  $t$  arrives at destination  $u$  at time  $t + h_{\mathcal{T}_t}(u)$ , given that  $u$  is included in the multicast tree  $\mathcal{T}_t$ . Thus,  $A_u(t)$  is influenced by two factors: (i) the frequency of which destination  $u$  is updated; (ii) the AoI of arrived packets, i.e.,  $h_{\mathcal{T}_t}(u)$ . Both factors are determined by multicast trees. Therefore, a proper multicast tree is essential for optimizing AoIs. In addition, some scenarios bring new influence factors, which are discussed in the following sections.

### B. Network Dynamics

A real network usually exhibits dynamic behavior. For instance, in wireless networks, there are topology changes from node mobility or link instability. We assume that the network topology at time  $t$  is generated by a random process, which is independent of the past. This assumption is reasonable in many scenarios (e.g., [16]). We first define a link indicator  $\sigma(e, t)$  as follows:

$$\sigma(e, t) = \begin{cases} 1, & \text{if link } e \in \mathcal{E}_t, \\ 0, & \text{otherwise.} \end{cases} \quad (4)$$

For a network, the distribution of links at time  $t$  is denoted as  $\sigma(t) \in \Xi$ , where  $\Xi$  denotes all possible network topologies.

$\sigma(t)$  evolves according to a stationary ergodic process:

$$\sum_{\sigma \in \Xi} p(\sigma) = 1, p(\sigma) > 0, \forall \sigma \in \Xi. \quad (5)$$

The statistical property of the process above depends on specific network scenarios, which are often difficult to obtain. We assume that the process is unknown in priority. In the subsequent sections, we will show how our proposed approach handles such dynamic scenarios naturally by making decisions on discrete time slots in model design and implementation.

### C. Energy-efficient Multicast Scheduling

In real networks, energy is a critical resource that needs to be managed efficiently. Here we propose a scheme to balance the energy consumption and AoIs, which we refer to as ‘‘multicast scheduling’’ in this paper. Specifically, consider a long-term energy budget denoted by  $W$  which constrains the average energy consumption. Let  $C(\mathcal{T}_t)$  denote the energy consumption of a multicast tree  $\mathcal{T}_t$ , our goal is to find multicast trees that minimize the AoIs while satisfying the energy constraint. Formally, we formulate this problem as follows:

$$\mathbf{OP:} \quad \min_{\mathcal{T}} \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\mathcal{T}} \left[ \sum_{t=0}^T \sum_{u \in \mathcal{U}_t} \omega_u A_u(t) \right], \quad (6a)$$

$$\text{s.t.} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\mathcal{T}} \left[ \sum_{t=0}^T C(\mathcal{T}_t) \right] \leq W, \quad (6b)$$

where  $\mathcal{T} = \{\mathcal{T}_t | t = 0, 1, \dots, T\}$  is a sequence of multicast trees generated over time,  $\omega_u \in (0, 1)$  represents the importance of a destination  $u \in \mathcal{U}_t$ . In other words, if we view the above problem as a sequential decision-making problem, we aim to find a policy  $\pi$  to determine  $\mathcal{T}_t$  as a function of the network state  $\{\mathcal{T}_\tau, \mathcal{G}_\tau, A(\tau) | 0 \leq \tau \leq t-1\}$ . Given the set of destinations  $\mathcal{U}_t$ , define the solution space of multicast trees including  $\mathcal{U}_t$  as  $\Omega(\mathcal{U}_t)$ , the problem above can be rewritten as:

$$\mathbf{DP:} \quad \max_{\lambda \geq 0} \quad \inf_{\mathcal{U}', \mathcal{T}} \left( \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\mathcal{T}} \left[ \sum_{t=0}^T \sum_{u \in \mathcal{U}_t} \omega_u A_u(t) + \lambda \left( \sum_{t=0}^T C(\mathcal{T}_t) - TW \right) \right] \right), \quad (7)$$

$$\text{s.t.} \quad \mathcal{T}_t \in \Omega(\mathcal{U}'_t), \forall t.$$

where  $\lambda$  is the Lagrangian multiplier corresponding to Eq. (6b). The large and discrete solution space and the mutual influence between decision variables make the problem challenging to solve. Therefore, we reformulate the original problem into two subproblems.

### III. PROBLEM REFORMULATION

Our main approach is to decompose Eq. (6) into two subproblems. The first subproblem involves identifying the set of destinations that should be updated at each time slot

and is referred to as the scheduling subproblem. The second subproblem entails finding an optimal multicast tree for the selected destinations and is known as the tree-generating subproblem. Each of them can be formulated as an MDP and solved by RL methods.

### A. Problem Decomposition

The two subproblems are formulated as follows:

**Definition 1 (Scheduling Subproblem).** Given a network  $\mathcal{G}_t$ , find a sequence  $\mathbf{U}' = \{\mathcal{U}'_t | t = 0, 1, \dots, T\}$  of destination sets  $\mathcal{U}'_t \in \mathcal{U}_t$  to:

$$\mathbf{P1:} \quad \max_{\lambda \geq 0} \inf_{\mathbf{U}'} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T g(\lambda, \mathcal{U}'_t) \quad (8)$$

where  $g(\lambda, \mathcal{U}'_t)$  is given by the optimal values of the following tree-generating subproblem.

**Definition 2 (Tree-generating Subproblem).** Given a network  $\mathcal{G}_t$  and a set of destinations  $\mathcal{U}'_t \in \mathcal{U}_t$ , find a multicast tree  $\mathcal{T}_t = \{\mathcal{V}_t^T, \mathcal{E}_t^T\}$  that:

$$\mathbf{P2:} \quad \min_{\mathcal{T}_t} \left( \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\mathcal{T}_t} \left[ \sum_{t=0}^T \sum_{u \in \mathcal{U}_t} \omega_u A_u(t) + \lambda \left( \sum_{t=0}^T C(\mathcal{T}_t) - TW \right) \right] \right), \quad (9)$$

s.t.  $\mathcal{T}_t \in \Omega(\mathcal{U}'_t)$ .

We can observe the equivalence between two subproblems and the original problem. In addition, the problem **P1** and **P2** can be both regarded as sequential decision-making problems, which can be formulated as MDPs. We first focus on problem **P1** and define MDP  $\mathcal{M}_1 = \{\mathcal{S}_1, \mathcal{A}_1, f_1, r_1\}$  as follows:

- **States:** The state  $s_t$  is defined as:

$$s_t = \{\mathcal{G}_t, \mathbf{o}_t\}, s_t \in \mathcal{S}_1, \quad (10)$$

where  $\mathbf{o}_t = \{\mathbf{x}_t, \mathbf{e}_t\}$  denotes the features at time slot  $t$ , including node features  $\mathbf{x}_t$  and link features  $\mathbf{e}_t$ .

- **Actions:** The action is a set of destinations, i.e.:

$$a_t = \{\mathcal{U}'_t | \mathcal{U}'_t \subseteq \mathcal{U}_t\}, a_t \in \mathcal{A}_1. \quad (11)$$

- The transition function  $f_1$  is unknown in priority.
- **Rewards:** To assess the network's quality at time slot  $t$ , we introduce a quality function  $q_1(s_t)$ :

$$q_1(s_t) = - \sum_{u \in \mathcal{U}_t} \omega_u A_u(t) - \lambda(C(\mathcal{G}_t) - W), \quad (12)$$

the reward function  $r_1 : \mathcal{S}_1 \times \mathcal{A}_1 \rightarrow \mathbb{R}$  is defined as:

$$r_1(s_t, a_t) = q_1(s_t) - q_1(s_{t-1}). \quad (13)$$

*Remark: By discretizing time, we can obtain the network state at each time slot and make individual decisions. Therefore, we naturally address the challenge of network dynamics.*

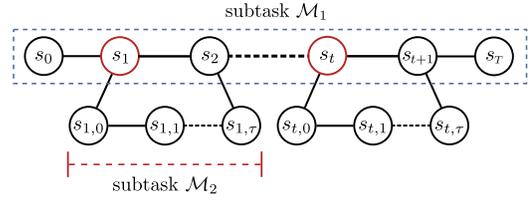


Fig. 3. Relationship between two Subtasks.  $\mathcal{M}_1$  is responsible for selecting destinations, which is utilized by  $\mathcal{M}_2$  to generate multicast trees.

### B. Tree-generating MDP

To tackle problem **P2**, notice that a multicast tree is composed of nodes and links. To approach this, we initiate with an empty set and incrementally add nodes and links in a virtual timescale. Specifically, we introduce a virtual timescale  $\tau$  and a partial solution  $\mathcal{P}_\tau = \{\mathcal{V}_\tau^P, \mathcal{E}_\tau^P\}$  with a source node at  $\tau = 0$ . An MDP  $\mathcal{M}_2 = \{\mathcal{S}_2, \mathcal{A}_2, f_2, r_2\}$  is defined as follows:

- **States:** The state  $s_\tau$  consists of network features and a partial solution  $\mathcal{P}_\tau$ , defined as <sup>2</sup>:

$$s_\tau = \{\mathcal{G}_t, \mathbf{o}_t, \mathcal{P}_\tau\}, s_\tau \in \mathcal{S}_2. \quad (14)$$

When the selected destinations  $\mathcal{U}'_t$  are covered by  $\mathcal{P}_\tau$ ,  $s_\tau$  will be a terminal state.

- **Actions:** The action  $a_\tau$  is a neighbor of  $\mathcal{P}_\tau$ , defined as:

$$a_\tau = v, v \in \mathcal{N}(\mathcal{P}_\tau), a_\tau \in \mathcal{A}_2, \quad (15)$$

where  $\mathcal{N}(\mathcal{P}_\tau)$  represents the neighbor nodes of  $\mathcal{P}_\tau$ . A link  $(v_\tau^*, a_\tau), v_\tau^* \in \mathcal{P}_\tau$  will be uniquely selected as described in section IV and added to  $\mathcal{P}_\tau$  along with node  $a_\tau$ .

- The transition function  $f_2$  is unknown in priority.
- **Rewards:** The quality function  $q_2(s_\tau)$  can be defined as:

$$q_2(s_\tau) = \sum_{u \in \mathcal{U}'_t \cap \mathcal{V}_\tau^P} \frac{\omega_u A_u(t)}{h_{\mathcal{P}_\tau}(u)} - \lambda(C(\mathcal{P}_\tau) - W), \quad (16)$$

where  $h_{\mathcal{P}_\tau}(u)$  is the number of hops between the source and destination  $u$  in  $\mathcal{P}_\tau$ . Subsequently, the reward function  $r_2 : \mathcal{S}_2 \times \mathcal{A}_2 \rightarrow \mathbb{R}$  is defined as:

$$r_2(s_\tau, a_\tau) = q_2(s_\tau) - q_2(s_{\tau-1}). \quad (17)$$

Therefore, the original problem can be solved by solving  $\mathcal{M}_1$  and  $\mathcal{M}_2$  sequentially, which is illustrated in Fig. 3. To achieve this, we propose an algorithm called Tree Generator-based Multicast Scheduling (TGMS).

## IV. TREE GENERATOR-BASED MULTICAST SCHEDULING

Our proposed approach consists of a tree generator and a scheduler, both utilizing graph embedding methods and DRL techniques.

<sup>2</sup> $s_\tau$  is actually the abbreviation of  $s_{t,\tau}$ .

### A. Graph Embedding Methods

Having a comprehensive understanding of the graph topology is essential to make informed decisions in network optimization problems. Graph embedding methods (e.g., GNNs) have demonstrated their efficacy in various CO problems [17]. Some studies focus on the attention mechanism in GNNs, which allows nodes to selectively aggregate information from neighbors (e.g., [18]). Based on these observations, we propose a novel GAT with guaranteed contraction mapping properties to extract graph information, which is defined as follows:

$$\phi(\mathbf{h}_i, \mathbf{h}_j) = \mathbf{a}^T \text{LeakyReLU}(\mathbf{W}_1(\mathbf{h}_i + \mathbf{h}_j) + \mathbf{W}_2 \mathbf{e}_{i,j}), \quad (18a)$$

$$\alpha_{ij} = \frac{\exp(\phi(\mathbf{h}_i, \mathbf{h}_j))}{\sum_{k \in \mathcal{N}_i \cup \{i\}} \exp(\phi(\mathbf{h}_i, \mathbf{h}_k))}, \quad (18b)$$

$$f_{\text{GAT}}(\mathbf{h}_i, \mathbf{x}) = \frac{1}{\|\mathbf{W}_1\|} (\alpha_{ii}(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_3 \mathbf{x}_i) + \sum_{j \in \mathcal{N}_i} \alpha_{ij}(\mathbf{W}_1 \mathbf{h}_j + \mathbf{W}_3 \mathbf{x}_j)), \quad (18c)$$

Here,  $\mathbf{h}_i \in \mathbb{R}^d$  denotes the embedding vector of node  $i$ ,  $\mathbf{x}_i$  denotes the node features of  $i$ ,  $\mathbf{e}_{i,j}$  denotes the link features of  $(i, j)$  and  $\text{LeakyReLU}(\cdot)^3$  is an activation function. The new representation of node  $i$  is obtained by a weighted sum of neighbors' node embeddings using  $\alpha_{ij}$ . Hence, we effectively disseminate information through the graph and aggregate node embeddings.

### B. Model Design

Recall that the action space of  $\mathcal{M}_1$  (see Eq. (11)) is excessively large<sup>4</sup>, making it impractical for RL algorithms to efficiently explore. To tackle this challenge, we arrange the destinations  $\mathcal{U}_t$  in descending order of their weighted AoIs, and then select a part of the destinations. However, the size of  $\mathcal{U}_t$  varies across different networks, rendering it difficult to determine a fixed value. To overcome this issue, we devise an agent with a continuous action space. Specifically, we utilize a model to predict the mean and standard deviation of a Gaussian distribution, which are then employed to sample the fraction of destinations to be selected. The forwarding process of the scheduler can be summarized as follows:

$$\mathbf{H}_t^{(l+1)} = f_{\text{GAT}}(\{\mathbf{h}_{t,i}^{(l)}\}_{i \in \mathcal{V}_t}), \quad (19a)$$

$$\tilde{\mathbf{h}}_t = \frac{1}{|\mathcal{V}|} \sum_{i=1}^{|\mathcal{V}|} \mathbf{h}_t^{(L)}, \quad (19b)$$

$$\mu = \mathbf{W}_1 \text{LeakyReLU}(\mathbf{W}_2 \tilde{\mathbf{h}}_t), \quad (19c)$$

$$\sigma = \mathbf{W}_3 \text{LeakyReLU}(\mathbf{W}_4 \tilde{\mathbf{h}}_t), \quad (19d)$$

$$\pi_1(a_t | s_t) = \mathcal{N}(\mu, \sigma), \quad (19e)$$

$$V_1(s_t) = \mathbf{W}_5 \text{LeakyReLU}(\mathbf{W}_6 \tilde{\mathbf{h}}_t). \quad (19f)$$

The initial graph embedding  $\mathbf{h}_0$  is generated by combining the graph topology and features, which is updated by the

<sup>3</sup>LeakyReLU( $x$ ) = max( $\alpha x$ ,  $x$ ) where  $\alpha$  is a hyperparameter.

<sup>4</sup>The dimension of  $|\mathcal{A}_2|$  is  $2^{|\mathcal{U}_t|}$ .

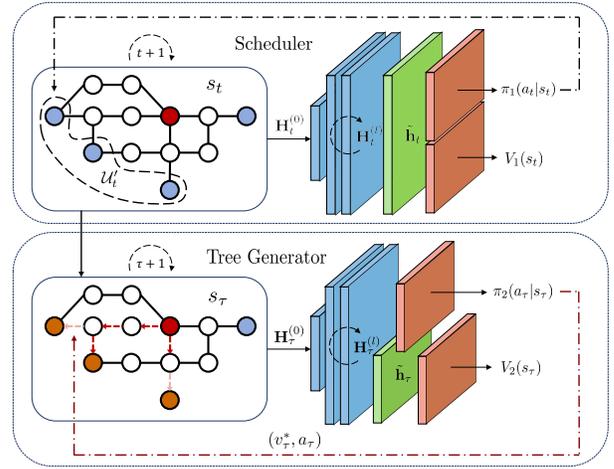


Fig. 4. The System Architecture of TGMS. A scheduler is performed to select a set of destinations, which are utilized by the tree generator to generate a multicast tree.

GAT defined by Eq. (18a)-(18c). Following this, the outputs are passed through a pooling layer serving to aggregate global information (see Eq. (19b)). Finally, the resulting graph embeddings are fed into two distinct heads: (i) a policy head responsible for predicting the mean and standard deviation of a Gaussian distribution (see Eq. (19c)-(19e)); and (ii) a critic head tasked with predicting the expected value of the current state (see Eq. (19f)). For the tree generator, we employ similar graph embedding methods as those utilized in the scheduler (see (19a)-(19b)). The difference lies in the heads as follows:

$$\pi_2(a_\tau | s_\tau) = \log \text{softmax}(\mathbf{W}'_1 \sigma(\mathbf{W}'_2 \mathbf{H}_\tau^{(L)})), \quad (20a)$$

$$V_2(s_\tau) = \mathbf{W}'_3 f_{\text{LeakyReLU}}(\mathbf{W}'_4 \tilde{\mathbf{h}}_\tau), \quad (20b)$$

where the policy  $\pi_2(a_\tau | s_\tau)$  is masked to ensure valid actions. One remaining question is that when an action  $a_\tau$  is sampled from  $\pi_2(a_\tau | s_\tau)$ , it is possible for  $a_\tau$  to have multiple links with  $\mathcal{P}_\tau$ . To address this, we select the link with the minimum cost from the set of candidate links, i.e.:

$$(v_\tau^*, a_\tau) = \arg \min_{v \in \mathcal{P}_\tau, (v, a_\tau) \in \mathcal{E}_\tau} C_{v, a_\tau}. \quad (21)$$

This approach effectively reduces the complexity of the tree generator while constraining the cost of  $\mathcal{P}_\tau$ . Importantly, it does not alter the solution space of  $\mathbf{P}_2$ . The system architecture of our approach is illustrated in Fig. 4.

## V. PERFORMANCE EVALUATION

Due to the lack of existing research on the proposed problem, we conduct extensive experiments to evaluate the performance of our approach. To validate the generalization ability of our approach, we consider three datasets of different graph topologies as follows. One of the datasets is called AS-733, which is collected from the University of Oregon Route Views Project [19]. We randomly select 240 graphs for training and 60 graphs for testing, where the testing graphs are

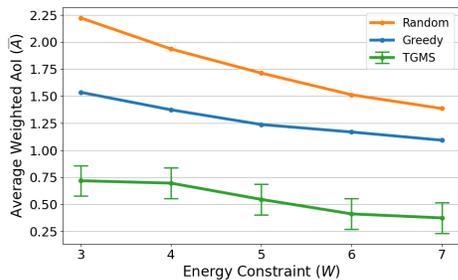


Fig. 5. Average Weighted AoI of AS-733 under  $\mathbf{W}_1$ .

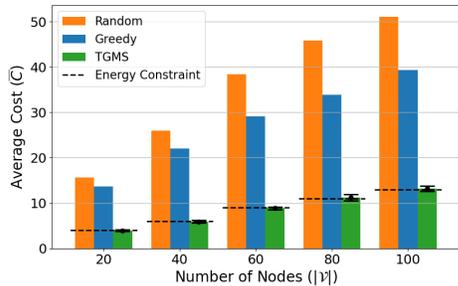


Fig. 6. Energy Consumption of AS-733 under  $\mathbf{W}_2$ .

unseen during training. Each graph will be trained or tested for 100 time slots. We randomly select 5 seeds and record the results with the mean and standard deviation. The following algorithms are considered as baselines:

- **Random:** A fraction  $m$  of nodes are randomly selected as destinations. When generating a multicast tree, we randomly select a valid node with its minimum-cost link.
- **Greedy:** A 0.5 fraction of sorted destinations are greedily selected. When generating a multicast tree, we greedily select a valid link with minimum cost.

Consider a set of energy constraints  $\mathbf{W}_1 = \{3, 4, 5, 6, 7\}$  for graph size  $|\mathcal{V}| = 60$ , we compare the average weighted AoI under different energy constraints. From Fig. 5, we conclude that TGMS has a superior performance compared to the baselines. When restricting energy consumption, TGMS can achieve a lower weighted AoI. Specifically, TGMS reduces the average weighted AoI by 57.1% and 68.7% compared to the greedy and random baselines, respectively.

Next, we evaluate the energy consumption under different energy constraints on different sizes of graphs. Consider another set of energy constraints  $\mathbf{W}_2 = \{4, 6, 9, 11, 13\}$  for graph sizes  $|\mathcal{V}| = \{20, 40, 60, 80, 100\}$ , respectively. We maintain a similar weighted AoI for all algorithms and compare the energy consumption. The results are shown in Fig. 6. We observe that TGMS achieves a lower energy consumption compared to the baselines. Specifically, TGMS reduces the energy consumption by 75.7% and 69.3% compared to the random and greedy baselines, respectively. *The additional experiments can be found in the appendix of [20].*

## VI. CONCLUSION

In this paper, we have proposed a novel hierarchical RL architecture including a GAT-based graph embedding method. It concludes with two agents: (i) a scheduler that selects a set of destinations while meeting an average power constraint, and (ii) a tree generator for generating multicast trees. We compare with several baselines and confirm that TGMS outperforms them in terms of AoI reduction and energy efficiency.

## REFERENCES

- [1] X. Jiang, F. R. Yu, T. Song, and V. C. Leung, "A survey on multi-access edge computing applied to video streaming: Some research issues and challenges," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 2, pp. 871–903, 2021.
- [2] L. Yin, J. Gui, Z. Zeng *et al.*, "Improving energy efficiency of multimedia content dissemination by adaptive clustering and d2d multicast," *Mobile Information Systems*, vol. 2019, 2019.
- [3] B. Quinn and K. Almeroth, "Ip multicast applications: Challenges and solutions," Tech. Rep., 2001.
- [4] S. Kaul, R. Yates, and M. Gruteser, "Real-time status: How often should one update?" in *2012 Proceedings IEEE INFOCOM*. IEEE, 2012, pp. 2731–2735.
- [5] R. D. Yates, Y. Sun, D. R. Brown, S. K. Kaul, E. Modiano, and S. Ulukus, "Age of information: An introduction and survey," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 5, pp. 1183–1210, 2021.
- [6] J. Li, Y. Zhou, and H. Chen, "Age of information for multicast transmission with fixed and random deadlines in iot systems," *IEEE Internet of Things Journal*, vol. 7, no. 9, pp. 8178–8191, 2020.
- [7] S. F. Lindström, M. Wetterberg, and N. Carlsson, "Cloud gaming: A qoe study of fast-paced single-player and multiplayer gaming," in *2020 IEEE/ACM 13th International Conference on Utility and Cloud Computing (UCC)*. IEEE, 2020, pp. 34–45.
- [8] A. A. Barakabitze, N. Barman, A. Ahmad, S. Zadtootaghaj, L. Sun, M. G. Martini, and L. Atzori, "Qoe management of multimedia streaming services in future networks: a tutorial and survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 1, pp. 526–565, 2019.
- [9] C. A. Oliveira and P. M. Pardalos, "A survey of combinatorial optimization problems in multicast routing," *Computers & Operations Research*, vol. 32, no. 8, pp. 1953–1981, 2005.
- [10] D. Kreutz, F. M. Ramos, P. E. Verissimo, C. E. Rothenberg, S. Azodolmoly, and S. Uhlig, "Software-defined networking: A comprehensive survey," *Proceedings of the IEEE*, vol. 103, no. 1, pp. 14–76, 2014.
- [11] M. Xie, J. Gong, X. Jia, and X. Ma, "Age and energy tradeoff for multicast networks with short packet transmissions," *IEEE Transactions on Communications*, vol. 69, no. 9, pp. 6106–6119, 2021.
- [12] I. Ljubić, "Solving steiner trees: Recent advances, challenges, and perspectives," *Networks*, vol. 77, no. 2, pp. 177–204, 2021.
- [13] M. Kim, J. Park *et al.*, "Learning collaborative policies to solve n-hard routing problems," *Advances in Neural Information Processing Systems*, vol. 34, pp. 10418–10430, 2021.
- [14] West and D. Brent, *Introduction to graph theory*. Prentice hall Upper Saddle River, 2001, vol. 2.
- [15] X. Lin and L. M. Ni, "Multicast communication in multicomputer networks," *IEEE transactions on Parallel and Distributed Systems*, vol. 4, no. 10, pp. 1105–1117, 1993.
- [16] A. Sinha, L. Tassiulas, and E. Modiano, "Throughput-optimal broadcast in wireless networks with dynamic topology," in *Proceedings of the 17th ACM International Symposium on Mobile Ad Hoc Networking and Computing*, 2016, pp. 21–30.
- [17] E. Khalil, H. Dai, Y. Zhang, B. Dilkina, and L. Song, "Learning combinatorial optimization algorithms over graphs," *Advances in neural information processing systems*, vol. 30, 2017.
- [18] S. Brody, U. Alon, and E. Yahav, "How attentive are graph attention networks?" *arXiv preprint arXiv:2105.14491*, 2021.
- [19] J. Leskovec, J. Kleinberg, and C. Faloutsos, "Graphs over time: densification laws, shrinking diameters and possible explanations," in *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, 2005, pp. 177–187.

- [20] Y. Zhang, G. Liao, S. Cao, N. Yang, and M. Zhang, "Age-minimal multicast by graph attention reinforcement learning," <https://arxiv.org/abs/2404.18084>, 2024.
- [21] G. Singal, V. Laxmi, M. S. Gaur, D. V. Rao, R. Kushwaha, D. Garg, and N. Kumar, "Qos-aware mesh-based multicast routing protocols in edge ad hoc networks: Concepts and challenges," *ACM Transactions on Internet Technology (TOIT)*, vol. 22, no. 1, pp. 1–27, 2021.
- [22] S. Kumar, A. Goswami, R. Gupta, S. P. Singh, and A. Lay-Ekuakille, "A game-theoretic approach for cost-effective multicast routing in the internet of things," *IEEE Internet of Things Journal*, vol. 9, no. 18, pp. 18 041–18 053, 2022.
- [23] D. Hatano and Y. Yoshida, "Computational aspects of the preference cores of supermodular two-scenario cooperative games." in *IJCAI*, 2018, pp. 310–316.
- [24] Z. Zhang, H. Chen, M. Hua, C. Li, Y. Huang, and L. Yang, "Double coded caching in ultra dense networks: Caching and multicast scheduling via deep reinforcement learning," *IEEE Transactions on Communications*, vol. 68, no. 2, pp. 1071–1086, 2019.
- [25] I. Kadota, A. Sinha, E. Uysal-Biyikoglu, R. Singh, and E. Modiano, "Scheduling policies for minimizing age of information in broadcast wireless networks," *IEEE/ACM Transactions on Networking*, vol. 26, no. 6, pp. 2637–2650, 2018.
- [26] J. Sun, L. Wang, Z. Jiang, S. Zhou, and Z. Niu, "Age-optimal scheduling for heterogeneous traffic with timely throughput constraints," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 5, pp. 1485–1498, 2021.
- [27] H. H. Yang, A. Arafa, T. Q. Quek, and H. V. Poor, "Optimizing information freshness in wireless networks: A stochastic geometry approach," *IEEE Transactions on Mobile Computing*, vol. 20, no. 6, pp. 2269–2280, 2020.
- [28] X. Chen, K. Gatsis, H. Hassani, and S. S. Bidokhti, "Age of information in random access channels," *IEEE Transactions on Information Theory*, vol. 68, no. 10, pp. 6548–6568, 2022.
- [29] R. D. Yates and S. K. Kaul, "The age of information: Real-time status updating by multiple sources," *IEEE Transactions on Information Theory*, vol. 65, no. 3, pp. 1807–1827, 2018.
- [30] A. M. Bedewy, Y. Sun, and N. B. Shroff, "Minimizing the age of information through queues," *IEEE Transactions on Information Theory*, vol. 65, no. 8, pp. 5215–5232, 2019.
- [31] L. Huang and E. Modiano, "Optimizing age-of-information in a multi-class queueing system," in *2015 IEEE international symposium on information theory (ISIT)*. IEEE, 2015, pp. 1681–1685.
- [32] A. M. Bedewy, Y. Sun, and N. B. Shroff, "Age-optimal information updates in multihop networks," in *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2017, pp. 576–580.
- [33] —, "The age of information in multihop networks," *IEEE/ACM Transactions on Networking*, vol. 27, no. 3, pp. 1248–1257, 2019.
- [34] B. Buyukates, A. Soysal, and S. Ulukus, "Age of information in multihop multicast networks," *Journal of Communications and Networks*, vol. 21, no. 3, pp. 256–267, 2019.
- [35] F. Wu, H. Zhang, J. Wu, Z. Han, H. V. Poor, and L. Song, "Uav-to-device underlay communications: Age of information minimization by multi-agent deep reinforcement learning," *IEEE Transactions on Communications*, vol. 69, no. 7, pp. 4461–4475, 2021.
- [36] M. Sun, X. Xu, X. Qin, and P. Zhang, "Aoi-energy-aware uav-assisted data collection for iot networks: A deep reinforcement learning method," *IEEE Internet of Things Journal*, vol. 8, no. 24, pp. 17 275–17 289, 2021.
- [37] O. S. Oubbati, M. Atiquzzaman, H. Lim, A. Rachedi, and A. Lakas, "Synchronizing uav teams for timely data collection and energy transfer by deep reinforcement learning," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 6, pp. 6682–6697, 2022.
- [38] M. A. Abd-Elmagid, H. S. Dhillon, and N. Pappas, "A reinforcement learning framework for optimizing age of information in rf-powered communication systems," *IEEE Transactions on Communications*, vol. 68, no. 8, pp. 4747–4760, 2020.
- [39] X. Wu, X. Li, J. Li, P. Ching, and H. V. Poor, "Deep reinforcement learning for iot networks: Age of information and energy cost tradeoff," in *GLOBECOM 2020-2020 IEEE Global Communications Conference*. IEEE, 2020, pp. 1–6.
- [40] K. Shi, X. Zhang, S. Zhang, and H. Li, "Time-expanded graph based energy-efficient delay-bounded multicast over satellite networks," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 9, pp. 10 380–10 384, 2020.
- [41] V. Tripathi, R. Talak, and E. Modiano, "Information freshness in multihop wireless networks," *IEEE/ACM Transactions on Networking*, 2022.
- [42] S. Mukherjee, F. Bronzino, S. Srinivasan, J. Chen, and D. Raychaudhuri, "Achieving scalable push multicast services using global name resolution," in *2016 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2016, pp. 1–6.
- [43] Q. Yu, H. Wan, X. Zhao, Y. Gao, and M. Gu, "Online scheduling for dynamic vm migration in multicast time-sensitive networks," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 6, pp. 3778–3788, 2019.
- [44] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," *Advances in neural information processing systems*, vol. 12, 1999.
- [45] K. Zhang, A. Koppel, H. Zhu, and T. Basar, "Global convergence of policy gradient methods to (almost) locally optimal policies," *SIAM Journal on Control and Optimization*, vol. 58, no. 6, pp. 3586–3612, 2020.
- [46] Y. Nesterov *et al.*, *Lectures on convex optimization*. Springer, 2018, vol. 137.

### A. Literature Review

Multicast routing poses a significant challenge in various scenarios, including ad hoc networks [21] and Internet of Things (IoTs) systems [22]. It also gains research interest in cooperative games [23]. Optimal solutions for both multicast routing and scheduling have been extensively studied and established in these contexts. However, obtaining such solutions becomes non-trivial in stochastic systems. To address the challenges stemming from unknown network dynamics, researchers have proposed deep RL algorithms to optimize multicast performance in a centralized manner (e.g., [24]). *Despite the success of these works in reducing the multicast delay, these approaches are NOT readily applied to the age-minimal multicast problem due to the aforementioned challenges of coupling decision variables and complicated graph structure.*

The concept of AoI was initially introduced by Kaul *et al.* [4] and has since garnered significant attention in wireless scenarios (e.g., [25], [26], [27], [28]) and queue theory (e.g., [29], [30], [31]). In the context of multicast networks, a substantial body of research has primarily concentrated on optimizing and analyzing stopping schemes to manage AoI. For instance, [6] considered deadlines in multicast AoI optimization, [32], [33], [34] explored AoI in multi-hop multicast networks, albeit without considering energy consumption.

On the other hand, other studies have delved into leveraging DRL methods to optimize the AoI in diverse scenarios, including UAV-assisted systems (e.g., [35], [36], [37]), RF-powered networks (e.g., [38]) and IoTs systems (e.g., [39]). This line of work mainly focuses on optimal resource allocation and trajectory design to achieve AoI optimization. [40] presented a graph-based approach in the context of multicast schemes in satellite networks. However, they did not consider the timeliness of the multicast process. On a different front, [11] conducted a study about the trade-offs between AoI and energy efficiency, while they overlooked the impact of network topology. Similarly, [41] focused on making scheduling decisions based on predetermined routings, without taking into account energy consumption. *Of utmost importance, all of the aforementioned approaches [40], [11], [41] did not address the challenge of coupled multicast routing and scheduling decisions.* Consequently, these approaches are not directly applicable to tackle our considered problem.

### B. Algorithm Details

The pseudo-code of the TGMS algorithm is described in Algorithm 1. The pseudo-code of our employed RL algorithm is described in Algorithm 2. We make a convergence analysis of our algorithm in Section F.

### C. Experiment Setting

Due to the lack of existing datasets for our problem, we built an environment for training and testing, where three datasets of different graph topologies are considered as follows:

---

### Algorithm 1 Tree Generator-based Multicast Scheduling (TGMS)

---

**Input:** A network graph  $\mathcal{G}_0$  and its features  $\mathbf{o}_0$

**Output:** Multicast trees  $\mathcal{T}$

---

```

1:  $t \leftarrow 0$ .
2: while  $t < T$  do
3:   Get state  $s_t = \{\mathcal{G}_t, \mathbf{o}_t\}$ .
4:   Compute  $\pi(a_t, s_t)$  and  $V(s_t)$ .
5:   Sort  $\mathcal{U}_t$  in descending order of  $\omega_u A_u(t)$ .
6:   Sample  $a_t \sim \pi(a_t, s_t)$  and select  $\mathcal{U}'_t \in \mathcal{U}_t$ .
7:   Initialize  $\tau \leftarrow 0$  and partial solution  $\mathcal{P}_\tau$ .
8:   while  $s_\tau$  is not a terminal state do
9:     Get state  $s_\tau = \{\mathcal{G}_t, \mathbf{o}_t, \mathcal{P}_\tau\}$ .
10:    Compute  $\pi(a_\tau, s_\tau)$  and  $V(s_\tau)$ .
11:    Sample  $a_\tau \in \mathcal{N}(\mathcal{P}_\tau)$  from  $\pi(a_\tau, s_\tau)$ .
12:     $\mathcal{P}_\tau \leftarrow$  select a link  $(v_\tau^*, a_\tau)$ .
13:     $\tau \leftarrow \tau + 1$ .
14:   end while
15:    $\mathcal{T} \leftarrow$  Multicast with  $\mathcal{P}_\tau$ .
16:    $t \leftarrow t + 1$ .
17: end while
18: return multicast trees  $\mathcal{T}$ 

```

---



---

### Algorithm 2 Advantage Actor-Critic (A2C) with Batch Normalization

---

**Input:** A batch of episodes  $B = \{\delta_0, \delta_1, \dots, \delta_b, \dots\}$  with  $\delta_b = \{s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T\}$

**Parameter:**  $\theta$  - parameter of the actor,  $\omega$  - parameter of the critic,  $\lambda$  - Lagrangian multiplier,  $\rho$  - ratio of entropy loss

$R_b \leftarrow$  Calculate returns of each episode  $\delta_b$ .

$R \leftarrow$  Combine  $R_b$  into a single batch.

$b \leftarrow 0$ .

**while**  $b < |B|$  **do**

$R_b \leftarrow (R_b - \mu_R) / \sigma_R$ .

$t \leftarrow 0$ .

**while**  $t < T$  **do**

$A_t \leftarrow R_{b,t} - V_\omega(s_t)$ .

$d\theta \leftarrow d\theta + \nabla_\theta (\log \pi_\theta(a_t | s_t) A_t + \rho H_\theta)$ .

$d\omega \leftarrow d\omega + \partial (R - V(s_t; \omega))^2 / \partial \omega$ .

$\lambda \leftarrow C(\mathcal{T}_t) - W$

$t \leftarrow t + 1$ .

**end while**

$b \leftarrow b + 1$ .

**end while**

---

- **ER-Graphs:** The Erdos-Renyi (ER) random graphs are a type of graphs where each pair of nodes is connected with a fixed probability  $p$ . It is widely used in the performance evaluation of multicast algorithms (e.g., in [42]). We set the link connection probability  $p = 6/|\mathcal{V}|$ .
- **BA-Graphs:** In Barabasi-Albert (BA) graphs, each newly introduced node connects to  $m$  pre-existing nodes. This connection process is governed by a probability

that is directly proportional to the number of links the existing nodes already possess. This characteristic endows the BA random graphs with a 'scale-free' quality. Commonly, these graphs find significant application within multicast systems. (e.g., in [43]). We set  $m = 2$  in our experiment.

- **AS-733:** The Autonomous Systems (AS)-733 dataset is a real-world dataset collected from the University of Oregon Route Views Project [19]. It contains 733 abstracted graphs of Autonomous Systems.

For each one of the three datasets we mentioned, we consider five different graph sizes  $|\mathcal{V}| = \{20, 40, 60, 80, 100\}$ . For each graph size, we randomly sample 240 graphs for training and 60 graphs for testing. For each graph, we randomly select a source node and a fraction of destinations  $\delta_d = 0.3$ . Some examples of the above datasets are shown in Fig. 7, 8 and 9, respectively.

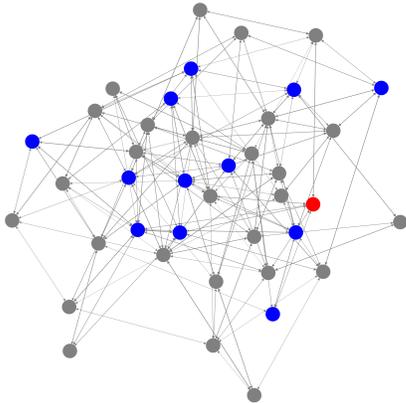


Fig. 7. An Example of ER-Graphs with  $|\mathcal{V}| = 40$ . The links are randomly connected, where the degree of each node is approximately the same.

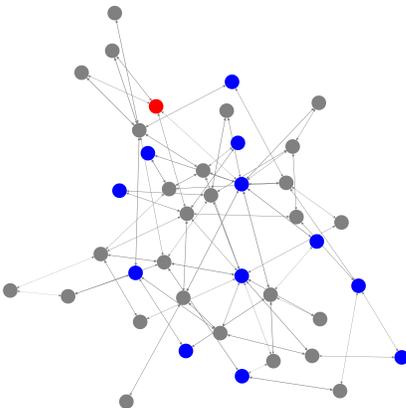


Fig. 8. An Example of BA-Graphs with  $|\mathcal{V}| = 40$ . The links are connected according to the preferential attachment mechanism. The difference in degrees is relatively small.

Module	Parameter	Value
A2C	Training Interval	20
	Discount factor ( $\gamma$ )	0.99
	Hidden dims	64
Network	Attention heads	5
	Dropout rate	0.3
Optimizer	Type	AdamW
	Actor Learning Rate	$2 \times 10^{-6}$
	Critic Learning Rate	$8 \times 10^{-7}$
	Momentum	0.9
	Centered gradient	True
	Gradient norm	0.5
	Type	CosineAnnealing
LR Scheduler	Warm restarts	True
	$T_0$	2
	$T_{mult}$	10
	$\eta_{min}$	$10^{-8}$

TABLE I  
SCHEDULER SETTING

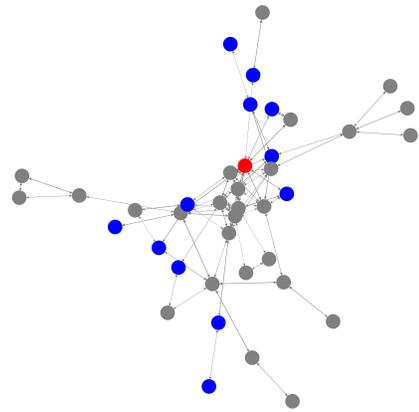


Fig. 9. An Example of AS-733 with  $|\mathcal{V}| = 40$ . The links are connected according to the real-world network topology. The difference in degrees is relatively large.

The meaning of using three types of graph topologies is to test the generalization ability and robustness of TGMS. While the EA-Graphs and BA-Graphs may not coincide with real-world networks, the experiment results still exhibit the performance of our model. Subsequently, We create some initial node features and link features on the above network topologies. We randomly assign a weight  $\omega_u \in (0, 1)$  and an initial AoI  $A_u(0) \in [1, 5]$  for every destination  $u \in \mathcal{U}_0$ . The sum of weights  $\sum_u \omega_u$  is ensured to be 1. The energy cost of each link is randomly assigned in the range  $(0, 1)$ .

#### D. Model Setting, Training and Testing

The parameters for the scheduler, tree generator and Lagrangian multiplier are shown in Table I, II and III, respectively. The unit of training interval is the time slot. All modules are trained from scratch for 100 time slots for each graph from the training dataset. The model is tested for 100 time slots per graph of the testing dataset. We use PyGraph 2.4.0 to implement TGMS, which is trained on NVIDIA GeForce RTX 3090 and tested on AMD EPYC 7763 CPU @1.50GHz with 64 cores under Ubuntu 20.04.6 LTS.

We use the following metrics to evaluate the performance of multicast algorithms:

Module	Parameter	Value
A2C	Training Interval	1
	Discount factor ( $\gamma$ )	0.99
Network	Hidden dims	64
	Attention heads	5
	Dropout rate	0.3
Optimizer	Type	AdamW
	Actor Learning Rate	$10^{-6}$
	Critic Learning Rate	$4 \times 10^{-7}$
	Momentum	0.9
	Centered gradient	True
	Gradient norm	0.5
	Type	CosineAnnealing
LR Scheduler	Warm restarts	True
	$T_0$	2
	$T_{mult}$	10
	$\eta_{min}$	$10^{-8}$

TABLE II  
TREE GENERATOR SETTING

Parameter	Value
Init Value	0.05
Training Interval	100
Learning Rate	$10^{-5}$

TABLE III  
LAGRANGIAN MULTIPLIER SETTING

- **Weighted Average AoI ( $\bar{A}$ ):** The weighted average AoI is the objective of problem **OP**, i.e.,  $\bar{A} = \frac{1}{T} \sum_{t=1}^T \sum_{u \in \mathcal{U}_t} \omega_u A_u(t)$ .
- **Average Energy Cost ( $\bar{C}$ ):** The average energy cost is mentioned in the constraint of problem **OP**, i.e.,  $\bar{C} = \frac{1}{T} \sum_{t=1}^T C(\mathcal{T}_t)$ . For an energy-efficient algorithm,  $\bar{C}$  should be kept lower than a given energy limit  $W$ .

We compare TGMS with some baselines as follows:

- **Greedy Algorithm:** We employ a greedy algorithm for both subtasks. In the scheduling subtask, we select a 0.5 fraction of destinations with the highest weighted AoI. In the tree-generating subtask, we greedily select the link with the minimum cost from the set of candidate links. When one of the selected destinations is in the candidate nodes, we select it and choose the link as described in our approach. For each type of graph topology, we choose  $\delta_{gd} \in [0.2, 0.4, 0.6, 0.8, 1.0]$ . The results are shown in Table IV, V and VI.
- **Random Algorithm:** We also adopt a random algorithm as a baseline. Specifically, the fraction of selected destinations is sampled from a uniform distribution  $\mathcal{U}(0, 1)$  in the scheduling process. The nodes of the multicast tree are also randomly selected from the neighbors of the partial solution until all selected destinations are included. The results are shown in Table VII.

From the results, we observe that the greedy algorithm performs well on ER-Graphs, but it is not robust on BA-Graphs and AS-Graphs. The reason is related to the symmetry of graph topologies. Due to the symmetry property of ER graphs, when the number of selected destinations increases, the energy costs of the greedy algorithm increase linearly and the average weighted AoI decreases almost linearly. However, the distribution of nodes' degrees in BA-Graphs is uneven.

$\delta_{gd}$  V		0.2	0.4	0.6	0.8	1.0
20	$\bar{A}$	1.005	0.625	0.514	0.461	0.372
	$\bar{C}$	1.454	2.004	2.434	2.897	3.534
40	$\bar{A}$	0.601	0.395	0.338	0.297	0.244
	$\bar{C}$	3.469	4.815	6.322	6.825	7.257
60	$\bar{A}$	0.483	0.323	0.281	0.244	0.209
	$\bar{C}$	5.808	8.592	9.987	11.229	12.165
80	$\bar{A}$	0.405	0.278	0.233	0.201	0.161
	$\bar{C}$	7.474	10.896	13.515	15.337	16.834
100	$\bar{A}$	0.298	0.222	0.199	0.179	0.144
	$\bar{C}$	12.227	15.961	18.018	20.246	21.764

TABLE IV  
GREEDY ALGORITHM ON ER-GRAPHS.

$\delta_{gd}$  V		0.2	0.4	0.6	0.8	1.0
20	$\bar{A}$	1.195	0.924	0.866	0.802	0.886
	$\bar{C}$	4.861	5.622	6.020	6.390	6.878
40	$\bar{A}$	0.532	0.447	0.367	0.339	0.340
	$\bar{C}$	9.173	11.861	13.333	13.324	13.731
60	$\bar{A}$	0.345	0.284	0.268	0.251	0.252
	$\bar{C}$	15.613	19.459	20.688	21.836	22.159
80	$\bar{A}$	0.320	0.256	0.238	0.229	0.232
	$\bar{C}$	23.589	26.396	28.790	30.246	31.030
100	$\bar{A}$	0.222	0.189	0.171	0.164	0.146
	$\bar{C}$	27.614	30.529	32.466	34.073	33.506

TABLE V  
GREEDY ALGORITHM ON BA-GRAPHS.

$\delta_{gd}$  V		0.2	0.4	0.6	0.8	1.0
20	$\bar{A}$	0.877	0.499	0.382	0.320	0.241
	$\bar{C}$	1.181	1.969	2.600	3.131	4.252
40	$\bar{A}$	0.529	0.426	0.427	0.412	0.410
	$\bar{C}$	9.553	11.355	12.869	13.377	13.858
60	$\bar{A}$	0.316	0.198	0.157	0.156	0.138
	$\bar{C}$	9.464	12.184	12.780	13.729	15.055
80	$\bar{A}$	0.212	0.150	0.134	0.126	0.114
	$\bar{C}$	20.456	23.360	25.242	26.365	27.450
100	$\bar{A}$	0.153	0.114	0.101	0.098	0.095
	$\bar{C}$	23.309	27.815	30.474	32.825	33.707

TABLE VI  
GREEDY ALGORITHM ON AS-GRAPHS.

DS  V		ER	BA	AS
20	$\bar{A}$	0.547	0.613	0.392
	$\bar{C}$	5.120	6.807	5.879
40	$\bar{A}$	0.304	0.286	0.292
	$\bar{C}$	12.933	16.355	15.439
60	$\bar{A}$	0.231	0.193	0.157
	$\bar{C}$	20.263	25.783	23.595
80	$\bar{A}$	0.177	0.166	0.101
	$\bar{C}$	28.261	35.821	33.274
100	$\bar{A}$	0.147	0.126	0.092
	$\bar{C}$	37.034	45.797	42.658

TABLE VII  
RANDOM ALGORITHM ON ALL DATASETS.

\mathcal{V}	DS	ER	BA	AS
		$\bar{A}$	0.495	0.900
20	$\bar{C}$	2.315	3.761	2.278
	$W$	2	4	2
	$\bar{A}$	0.288	0.306	0.378
40	$\bar{C}$	6.807	8.368	9.218
	$W$	6	9	9
	$\bar{A}$	0.249	0.215	0.257
60	$\bar{C}$	8.084	15.434	18.595
	$W$	8	15	18
	$\bar{A}$	0.171	0.140	0.573
80	$\bar{C}$	11.821	22.337	17.789
	$W$	12	23	18
	$\bar{A}$	0.188	0.099	0.158
100	$\bar{C}$	18.034	22.797	22.622
	$W$	18	23	20

TABLE VIII  
TGMS ALGORITHM

When the number of selected destinations increases, the energy costs of the greedy algorithm hardly increase and the decrease in the average weighted AoI is also not obvious. In a real-world network, when the number of selected destinations increases, the energy costs of the greedy algorithm increase slowly. The decrease in the average weighted AoI is obvious when the graph is small, but it is not obvious when the graph is large.

For the random algorithm, it performs poorly on all datasets. The reason is that the random algorithm does not consider the network topology and node features. When the graph is large, the random scheduler tends to select almost all links, which leads to a large energy cost.

The results of TGMS are shown in Table VIII. We find that TGMS can achieve a lower average weighted AoI than the greedy algorithm, while the energy consumption is also lower. When testing TGMS on other graph topologies, we observe similar results, which exhibits the robustness of our approach.

### E. Insights of Reward Functions

Here we provide some insights into the reward functions we adopted in our formulated MDPs. We start by calculating the cumulative reward of  $\mathcal{M}_1$ :

$$\begin{aligned}
R_1 &= \sum_{t=0}^T \gamma^{t-1} r_1(s_t, a_t) \\
&= r_1(s_0, a_0) + \sum_{t=1}^T \gamma^t (q_1(s_t) - q_1(s_{t-1})) \\
&\stackrel{\gamma=1}{=} r_1(s_0, a_0) + q_1(s_T) - q_1(s_0) \\
&= - \sum_{u \in \mathcal{U}_t} \omega_u A_u(T) - \lambda(C(\mathcal{T}_T) - W),
\end{aligned} \tag{22}$$

where we set  $q_1(s_0) = r_1(s_0, a_0)$  and assume the discount factor  $\gamma = 1$ . Comparing the RHS of Eq. (22) with problem **DP**, we observe that maximizing  $R_1$  is similar to solving problem **DP**. The difference is that problem **DP** aims to min-

imize a long-term objective, which is naturally decomposed as the reward function  $r_1$ .

To understand the reward function of  $\mathcal{M}_2$ , we first analyze how much a multicast tree  $\mathcal{T}_t$  can reduce the AoI of a destination. Denote  $\Delta A_u(t, t + h_{\mathcal{T}_t}(u))$  as the AoI reduction of destination  $u$  during time  $[t, t + h_{\mathcal{T}_t}(u)]$ , we have:

$$\begin{aligned}
&\Delta A_u(t, t + h_{\mathcal{T}_t}(u)) \\
&= A_u^+(t, t + h_{\mathcal{T}_t}(u)) - A_u^-(t, t + h_{\mathcal{T}_t}(u)) \\
&= \sum_{k=0}^{h_{\mathcal{T}_t}(u)} (A_u(t) + k) - \left( \sum_{k=0}^{h_{\mathcal{T}_t}(u)-1} (A_u(t) + k) + h_{\mathcal{T}_t}(u) \right) \\
&= A_u(t),
\end{aligned} \tag{23}$$

where  $A_u^+(t, t + h_{\mathcal{T}_t}(u))$  and  $A_u^-(t, t + h_{\mathcal{T}_t}(u))$  denotes the AoI of destination  $u$  during time  $[t, t + h_{\mathcal{T}_t}(u)]$  before and after the multicast tree  $\mathcal{T}_t$  is generated, respectively. That means whatever the multicast tree is, the AoI of a destination will be reduced by  $A_u(t)$  after  $h_{\mathcal{T}_t}(u)$  time slots. Therefore, the mean AoI reduction of a destination  $u$  during time  $[t, t + h_{\mathcal{T}_t}(u)]$  is:

$$\frac{1}{h_{\mathcal{T}_t}(u)} \Delta A_u(t, t + h_{\mathcal{T}_t}(u)) = \frac{A_u(t)}{h_{\mathcal{T}_t}(u)} \tag{24}$$

Then, we calculate the cumulative reward of  $\mathcal{M}_2$  as:

$$\begin{aligned}
R_2 &= \sum_{t=0}^T \gamma^{t-1} r_2(s_t, a_t) \\
&= r_2(s_0, a_0) + \sum_{t=1}^T \gamma^t (q_2(s_t) - q_2(s_{t-1})) \\
&\stackrel{\gamma=1}{=} r_2(s_0, a_0) + q_2(s_T) - q_2(s_0) \\
&= \sum_{u \in \mathcal{U}_t \cap \mathcal{V}_\tau^P} \frac{\omega_u A_u(t)}{h_{\mathcal{P}_\tau}(u)} - \lambda(C(\mathcal{P}_\tau) - W).
\end{aligned} \tag{25}$$

Comparing Eq. (24) and the RHS of Eq. (25), we claim that maximizing  $R_2$  is similar to maximizing the average weighted AoI reduction of destinations  $\mathcal{U}_t \cap \mathcal{V}_\tau^P$  during time  $[t, t + h_{\mathcal{P}_\tau}(u)]$ . Note that solely maximizing  $R_2$  may not equal solving problem **DP**. However, by combining two MDPs, we can obtain a solution that is close to the optimal solution of problem **DP**.

### F. Theorems and Proofs

In this section, we provide theorems and proofs for the properties of our proposed algorithms. For the MDP  $\mathcal{M}_2$ , we have the following Proposition:

**Proposition 1.** In  $\mathcal{M}_2$ , the partial solution  $\mathcal{P}_\tau$  is always a multicast tree.

*Proof.* The tree generator repeatedly adds a node that has not been chosen before, which means the node is connected to  $\mathcal{P}_\tau$  with at most one link. Therefore,  $\mathcal{P}_\tau$  is acyclic. Additionally,  $\mathcal{P}_\tau$  is connected due to the selection of neighbors, thus  $\mathcal{P}_\tau$  is a connected graph. Therefore,  $\mathcal{P}_\tau$  is a tree.  $\square$

For the proposed algorithm, we wonder how GAT influences the performance and the convergence property of our model. We observe that the tree generator and the scheduler have similar structures. Therefore, our main idea is to analyze a model consisting of GAT layers and linear layers. For the GAT defined as below:

$$\phi(\mathbf{h}_i, \mathbf{h}_j) = \mathbf{a}^\top \text{LeakyReLU}(\mathbf{W}_1(\mathbf{h}_i + \mathbf{h}_j) + \mathbf{W}_2 \mathbf{e}_{i,j}), \quad (26a)$$

$$\alpha_{ij} = \frac{\exp(\phi(\mathbf{h}_i, \mathbf{h}_j))}{\sum_{k \in \mathcal{N}_i \cup \{i\}} \exp(\phi(\mathbf{h}_i, \mathbf{h}_k))}, \quad (26b)$$

$$f_{\text{GAT}}(\mathbf{h}_i, \mathbf{x}) = \frac{1}{\|\mathbf{W}_1\|} (\alpha_{ii}(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_3 \mathbf{x}_i) + \sum_{j \in \mathcal{N}_i} \alpha_{ij}(\mathbf{W}_1 \mathbf{h}_j + \mathbf{W}_3 \mathbf{x}_j)), \quad (26c)$$

we have the following theorem:

**Theorem 1.** For any undirected graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , given a mapping  $f_{\text{GAT}}$  defined in Eq. (26c), if the attention coefficients  $\alpha_{ij}$  are symmetric, i.e.,  $\alpha_{ij} = \alpha_{ji}$ . Then  $f_{\text{GAT}}(\cdot, \mathbf{x})$  is a contraction mapping for any initial node embeddings, i.e.:

$$d(f_{\text{GAT}}(\mathbf{H}, \mathbf{x}), f_{\text{GAT}}(\mathbf{H}', \mathbf{x})) \leq d(\mathbf{H}, \mathbf{H}'), \quad (27)$$

where  $\mathbf{H} = \{\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_{|\mathcal{V}|}\}$  is the matrix of node embeddings,  $d$  is a distance metric with respect to  $\mathbf{H}$ , defined as  $d(\mathbf{H}, \mathbf{H}') = \|\sum_{i \in \mathcal{V}} (\mathbf{h}_i - \mathbf{h}'_i)\|$ .

*Proof.* Let node embeddings  $\mathbf{H} = \{\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_{|\mathcal{V}|}\}$  and  $\mathbf{H}' = \{\mathbf{h}'_0, \mathbf{h}'_1, \dots, \mathbf{h}'_{|\mathcal{V}|}\}$ . From Eq. (26c), we have:

$$\begin{aligned} & d(f_{\text{GAT}}(\mathbf{H}, \mathbf{x}), f_{\text{GAT}}(\mathbf{H}', \mathbf{x})) \\ &= \left\| \sum_{i \in \mathcal{V}} (f_{\text{GAT}}(\mathbf{h}_i, \mathbf{x}) - f_{\text{GAT}}(\mathbf{h}'_i, \mathbf{x})) \right\| \\ &= \left\| \sum_{i \in \mathcal{V}} \frac{1}{\|\mathbf{W}_1\|} (\alpha_{ii}(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_3 \mathbf{x}_i) + \sum_{j \in \mathcal{N}_i} \alpha_{ij}(\mathbf{W}_1 \mathbf{h}_j + \mathbf{W}_3 \mathbf{x}_j) - \alpha_{ii}(\mathbf{W}_1 \mathbf{h}'_i + \mathbf{W}_3 \mathbf{x}_i) - \sum_{k \in \mathcal{N}_i} \alpha_{ik}(\mathbf{W}_1 \mathbf{h}'_k + \mathbf{W}_3 \mathbf{x}_k)) \right\| \\ &= \frac{1}{\|\mathbf{W}_1\|} \left\| \sum_{i \in \mathcal{V}} (\alpha_{ii} \mathbf{W}_1 (\mathbf{h}_i - \mathbf{h}'_i) + \sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W}_1 \mathbf{h}_j - \sum_{k \in \mathcal{N}_i} \alpha_{ik} \mathbf{W}_1 \mathbf{h}'_k) \right\| \\ &= \frac{1}{\|\mathbf{W}_1\|} \left\| \mathbf{W}_1 \sum_{i \in \mathcal{V}} (\alpha_{ii} (\mathbf{h}_i - \mathbf{h}'_i) + \sum_{k \in \mathcal{D}_i} \alpha_{ki} (\mathbf{h}_i - \mathbf{h}'_i)) \right\| \\ &\stackrel{\alpha_{ij} = \alpha_{ji}}{=} \frac{1}{\|\mathbf{W}_1\|} \left\| \mathbf{W}_1 \sum_{i \in \mathcal{V}} (\alpha_{ii} + \sum_{k \in \mathcal{D}_i} \alpha_{ik}) (\mathbf{h}_i - \mathbf{h}'_i) \right\| \\ &= \frac{1}{\|\mathbf{W}_1\|} \left\| \mathbf{W}_1 \sum_{i \in \mathcal{V}} (\alpha_{ii} + \sum_{j \in \mathcal{N}_i} \alpha_{ij}) (\mathbf{h}_i - \mathbf{h}'_i) \right\| \\ &\leq \left\| \sum_{i \in \mathcal{V}} (\mathbf{h}_i - \mathbf{h}'_i) \right\| = d(\mathbf{H}, \mathbf{H}'). \end{aligned} \quad (28)$$

Therefore, the mapping  $f_{\text{GAT}}(\cdot, \mathbf{x})$  is a contraction mapping.  $\square$

We have some remarks on the above theorem: (a) According to Banach's fixed-point theorem, the mapping  $f_{\text{GAT}}(\mathbf{H}, \mathbf{x})$  has a unique fixed point for any initial node embeddings  $\mathbf{H}$ , meaning that the GAT can map the node features to a stable state. (b) In our design, the GAT is used to lower the dimension of node embeddings, which can be regarded as a dimension reduction technique with the ability to capture the importance of nodes. (c) We can view the output of GAT as a new representation of the network state with noise, which can be utilized for the subsequent linear layers. Specifically, let  $\mathbf{H}^*$  be the fixed point of  $f_{\text{GAT}}(\mathbf{H}, \mathbf{x})$ , we have  $s_{\text{Noisy}} = \mathbf{H}^*$ . For simplicity, denote  $s_N$  as the noisy state for the subsequent linear neural networks.

### G. Convergence of Algorithm 2

Based on the above analysis, we provide a convergence theorem of Algorithm 2. We first introduce some assumptions as follows.

**Assumption 1.** Suppose all reward functions  $r$  and the parameterized policy  $\pi_\theta$  satisfy the following conditions:

(a)  $r$  is bounded, i.e.:

$$-R \leq r(s, a) \leq R, \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \quad (29)$$

(b)  $\nabla \log \pi_\theta(a|s)$  is  $L_\theta$ -Lipschitz and bounded<sup>5</sup>, i.e.:

$$\|\nabla \log \pi_{\theta_1}(a|s) - \nabla \log \pi_{\theta_2}(a|s)\| \leq L_\theta \|\theta_1 - \theta_2\|, \quad (30)$$

$$\|\nabla \log \pi_\theta(a|s)\| \leq B_\theta, \forall \theta_1, \theta_2 \in \mathbb{R}^d, \quad (31)$$

where  $L_\theta, B_\theta$  are constants.

**Assumption 2.** Suppose the critic satisfies the following condition:

(a) The critic is approximated by a linear function  $V_\omega(s) = \phi(s)^\top \omega$ , where  $\phi(s)$  is the feature vector of state  $s$  and  $\omega$  is the parameter of the critic.

(b) The feature mapping  $\phi(s)$  has a bounded norm  $|\phi(s)| \leq 1$ .

(c)  $V_\omega(s)$  is bounded, i.e.:  $|V_\omega(s)| \leq B_\omega$ , where  $B_\omega$  is a constant.

*Remark:* The assumption 2 holds because we divide the critic into two parts: the graph attention network and the linear function, where the former is proved to be a contraction mapping in theorem 1.

Then we have the following convergence theorem.

**Theorem 2** (Convergence of Algorithm 2). Suppose Assumptions 1 and 2 holds, given the sequence of critic below:

$$\mathcal{E}(t) = \frac{2}{t} \sum_{k=\tau}^t \mathbb{E}[\|\omega^* - \omega_k\|^2]. \quad (32)$$

If it is bounded, then we have the following convergence:

$$\min_{\tau \leq k \leq t} \mathbb{E}[\|\nabla J(\theta_k)\|^2] = \mathcal{O}\left(\frac{\mathcal{E}(t)}{\tau}\right) + \mathcal{O}\left(\frac{1}{\tau^\sigma}\right) \quad (33)$$

<sup>5</sup>The gradient of the policy is bounded by the gradient clipping technique.

Here we provide a proof of the above theorem. Consider the following state-value function:

$$v_{\pi_{\theta}}(s) = \mathbb{E}_{\substack{a_t \sim \pi_{\theta}(\cdot|s_t), \\ s_{t+1} \sim \mathbb{P}_{\pi_{\theta}}(\cdot|s_t, a_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s \right], \quad (34)$$

where  $\mathbb{P}_{\pi_{\theta}}(\cdot|s_t, a_t)$  is the transition probability under policy  $\pi_{\theta}$ . Suppose we aim to maximize the expected state value of the initial state  $s_0$ . The objective function can be denoted as:

$$J(\theta) := v_{\pi_{\theta}}(s_0). \quad (35)$$

Obviously, we have  $|J(\theta)| < \frac{R}{1-\gamma}$ . According to the policy gradient theorem [44], the expected value of  $\nabla J(\theta)$  can be written as:

$$\nabla J(\theta) = \mathbb{E}_{s \sim \mu_{\theta}, a \sim \pi_{\theta}} [A_{\pi_{\theta}}(s, a) \nabla \log \pi_{\theta}(a|s)], \quad (36)$$

where  $A_{\pi_{\theta}}(s, a) = r(s, a) - v_{\pi_{\theta}}(s, a)$  represents the advantage value function under policy  $\pi_{\theta}$ . Notice that  $v_{\pi_{\theta}}$  is the true value function under policy  $\pi_{\theta}$ , which is unknown in practice. Therefore, as Assumption 2 states, we consider a linear-approximated critic. That means we can obtain an approximated advantage value function  $\hat{A}_t(s_t, a_t) = r(s_t, a_t) - V_{\omega}(s_t)$ , where  $\omega$  is the parameter of the critic. Then the gradient of  $J(\theta)$  can be approximated by:

$$\nabla \hat{J}(\theta) = \mathbb{E}_{s \sim \mu_{\theta}, a \sim \pi_{\theta}} [\hat{A}_{\pi_{\theta}}(s, a) \nabla \log \pi_{\theta}(a|s)], \quad (37)$$

and  $\nabla \hat{J}(\theta)$  is bounded as follows:

**Lemma 2** (Bound of Approximated Policy Gradient). Under Assumption 1 and 2, the policy gradient  $\nabla \hat{J}(\theta)$  is bounded, i.e.:

$$\|\nabla \hat{J}(\theta)\| \leq (R + B_{\omega})B_{\theta}. \quad (38)$$

*Proof:* The approximated advantage value function  $\hat{A}_t(s_t, a_t)$  can be written as:

$$\begin{aligned} \hat{A}_t(s_t, a_t) &= r(s_t, a_t) - \phi(s_t)^{\top} \omega \\ &\leq |r(s_t, a_t)| + |\phi(s_t)^{\top} \omega| \leq R + B_{\omega}, \end{aligned} \quad (39)$$

which implies:

$$\|\hat{A}_t(s_t, a_t) \nabla \log \pi_{\theta}(a_t|s_t)\| \leq (R + B_{\omega})B_{\theta}. \quad (40)$$

□

Now we analyze the approximation error of  $\nabla \hat{J}(\theta)$ . First, introduce the following lemma:

**Lemma 3** (Lipschitz-Continuity of Policy Gradient [45]). Under Assumption 1, the policy gradient  $\nabla J(\theta)$  is Lipschitz continuous with some constant  $L > 0$ , i.e.:

$$\|\nabla J(\theta_1) - \nabla J(\theta_2)\| \leq L \|\theta_1 - \theta_2\|. \quad (41)$$

According to Lemma 3 and Eq. (1.2.19) from [46], we have:

$$J(\theta_2) \geq J(\theta_1) + \langle \nabla J(\theta_1), \theta_2 - \theta_1 \rangle - \frac{L}{2} \|\theta_1 - \theta_2\|^2, \quad (42)$$

where  $\langle \cdot, \cdot \rangle$  denotes the dot product of two vectors. The actor of Algorithm 2 is updated by:

$$\theta_{t+1}^a = \theta_t^a + \alpha_t \hat{A}_t(s_t, a_t) \nabla \log \pi_{\theta_t}(a_t|s_t), \quad (43)$$

Substitute (43) into (42) yields:

$$J(\theta_{t+1}) \geq J(\theta_t) + \alpha_t \langle \nabla J(\theta_t), \nabla \hat{J}(\theta_t) \rangle - \frac{L\alpha_t^2}{2} \|\nabla \hat{J}(\theta_t)\|^2. \quad (44)$$

Here, we focus on the second term on the right-hand side of the equation above, which can be decomposed as:

$$\begin{aligned} &\langle \nabla J(\theta_t), \nabla \hat{J}(\theta_t) \rangle \\ &= \langle \nabla J(\theta_t), \phi(s_t)^{\top} (\omega^* - \omega_t) \nabla \log \pi_{\theta_t}(a_t|s_t) \rangle \\ &\quad + \langle \nabla J(\theta_t), (v(s_t) - \phi(s_t)^{\top} \omega^*) \nabla \log \pi_{\theta_t}(a_t|s_t) \rangle \\ &\quad + \langle \nabla J(\theta_t), \nabla \hat{J}(\theta_t) - \nabla J(\theta_t) \rangle + \langle \nabla J(\theta_t), \nabla J(\theta_t) \rangle \end{aligned} \quad (45)$$

We can view the first term on the RHS of (45) as the error of the critic, the second term as the error from linear function approximation of the critic and the third term as the Markovian noise. Therefore, the above equation can be rewritten as:

$$\begin{aligned} &\langle \nabla J(\theta_t), \nabla \hat{J}(\theta_t) \rangle \\ &= \Phi_{EC}(\theta_t) + \Phi_{LFA}(\theta_t) + \Phi_M(\theta_t) + \|\nabla J(\theta_t)\|^2, \end{aligned} \quad (46)$$

We will analyze the bounds of these terms in the following. First,  $\Phi_{EC}(\theta_t)$  measures the distance between the optimal parameter  $\omega^*$  and the current parameter  $\omega_t$ , which is bounded as follows:

**Lemma 4** (Error of Critic). Given Assumption 2, we have:

$$\mathbb{E}[\Phi_{EC}(\theta_t)] \geq -B_{\theta} \sqrt{\mathbb{E}[\|\nabla J(\theta_t)\|^2]} \sqrt{\mathbb{E}[\|\omega^* - \omega_t\|^2]} \quad (47)$$

*Proof:* From Cauchy inequality and Eq. (31), we have:

$$\begin{aligned} &\langle \nabla J(\theta_t), \phi(s_t)^{\top} (\omega^* - \omega_t) \nabla \log \pi_{\theta_t}(a_t|s_t) \rangle \\ &\geq -\|\nabla J(\theta_t)\| \|\phi(s_t)^{\top} (\omega^* - \omega_t)\| \|\nabla \log \pi_{\theta_t}(a_t|s_t)\| \\ &\geq -\|\nabla J(\theta_t)\| \|\omega^* - \omega_t\| \|\nabla \log \pi_{\theta_t}(a_t|s_t)\| \\ &\geq -B_{\theta} \|\nabla J(\theta_t)\| \|\omega^* - \omega_t\|. \end{aligned} \quad (48)$$

Taking the expectations of both sides yields the result. □

Then we analyze the approximation error between the optimal parameter  $\omega^*$  and the parameter  $\omega_t$ . Suppose there is an approximation error bound  $\epsilon_v > 0$  that  $\|\mathbb{E}[\phi(s)^{\top} \omega^* - v(s)]\| \leq \epsilon_v$ , we have:

**Lemma 5** (Error of Linear Function Approximation).

$$\mathbb{E}[\Phi_{LFA}(\theta_t)] \geq -\epsilon_v B_{\theta} \mathbb{E}[\|\nabla J(\theta_t)\|]. \quad (49)$$

*Proof:* From Cauchy inequality, we have:

$$\begin{aligned} &\langle \nabla J(\theta_t), (v(s_t) - \phi(s_t)^{\top} \omega^*) \nabla \log \pi_{\theta_t}(a_t|s_t) \rangle \\ &\geq -\|\nabla J(\theta_t)\| \|v(s_t) - \phi(s_t)^{\top} \omega^*\| \|\nabla \log \pi_{\theta_t}(a_t|s_t)\| \\ &\geq -\epsilon_v B_{\theta} \|\nabla J(\theta_t)\|. \end{aligned} \quad (50)$$

Taking the expectations of both sides yields the result. □

Suppose the Markovian noise is bounded, i.e., there exists  $\epsilon_m > 0$  that  $-\epsilon_m \leq \Phi_M(\boldsymbol{\theta}_t) \leq \epsilon_m$ . Substitute the inequalities above into Eq. (46) yields:

$$\begin{aligned} & \mathbb{E}[\langle \nabla J(\boldsymbol{\theta}_t), \nabla \hat{J}(\boldsymbol{\theta}_t) \rangle] \\ & \geq -B_\theta \sqrt{\mathbb{E}[\|\nabla J(\boldsymbol{\theta}_t)\|^2]} \sqrt{\mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_t\|^2]} \\ & \quad - \epsilon_v B_\theta \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_t)\|] - \epsilon_m + \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_t)\|^2]. \end{aligned} \quad (51)$$

Substitute Eq. (51) into Eq. (44) and take the expectation of both sides yields:

$$\begin{aligned} & \mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t)] \\ & \geq -\alpha_t B_\theta \sqrt{\mathbb{E}[\|\nabla J(\boldsymbol{\theta}_t)\|^2]} \sqrt{\mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_t\|^2]} \\ & \quad - \alpha_t \epsilon_v \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_t)\|] - \alpha_t \epsilon_m + \alpha_t \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_t)\|^2] - \frac{L\alpha_t^2 B_j^2}{2}, \end{aligned} \quad (52)$$

where  $B_j = (R + B_\omega)B_\theta$ . Rearrange the terms and divide both sides by  $\alpha_t$  yields:

$$\begin{aligned} & (\mathbb{E}\|\nabla J(\boldsymbol{\theta}_t)\| - \frac{\epsilon_v}{2})^2 \leq \frac{1}{\alpha_t} \mathbb{E}[J(\boldsymbol{\theta}_{t+1}) - J(\boldsymbol{\theta}_t)] \\ & \quad + B_\theta \sqrt{\mathbb{E}[\|\nabla J(\boldsymbol{\theta}_t)\|^2]} \sqrt{\mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_t\|^2]} + \frac{L\alpha_t B_j^2}{2} \\ & \quad + \epsilon_m + \frac{\epsilon_v^2}{4}. \end{aligned} \quad (53)$$

Summing both sides over  $k$  from  $\tau$  to  $t$ , we have:

$$\begin{aligned} & \sum_{k=\tau}^t (\mathbb{E}\|\nabla J(\boldsymbol{\theta}_k)\| - \frac{\epsilon_v}{2})^2 \leq \sum_{k=\tau}^t \frac{1}{\alpha_k} \mathbb{E}[J(\boldsymbol{\theta}_{k+1}) - J(\boldsymbol{\theta}_k)] \\ & \quad + B_\theta \sum_{k=\tau}^t \sqrt{\mathbb{E}[\|\nabla J(\boldsymbol{\theta}_k)\|^2]} \sqrt{\mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_k\|^2]} \\ & \quad + \frac{LB_j^2}{2} \sum_{k=\tau}^t \alpha_k + (t - \tau + 1)(\epsilon_m + \frac{\epsilon_v^2}{4}). \end{aligned} \quad (54)$$

For the first term of RHS of Eq. (54), we have:

$$\begin{aligned} & \sum_{k=\tau}^t \frac{1}{\alpha_k} \mathbb{E}[J(\boldsymbol{\theta}_{k+1}) - J(\boldsymbol{\theta}_k)] \\ & = \sum_{k=\tau}^{t-1} (\frac{1}{\alpha_k} - \frac{1}{\alpha_{k+1}}) \mathbb{E}[J(\boldsymbol{\theta}_{k+1})] - \frac{1}{\alpha_\tau} \mathbb{E}[J(\boldsymbol{\theta}_\tau)] \\ & \quad + \frac{1}{\alpha_t} \mathbb{E}[J(\boldsymbol{\theta}_{t+1})] \\ & \leq \sum_{k=\tau}^{t-1} (\frac{1}{\alpha_k} - \frac{1}{\alpha_{k+1}}) \frac{R}{1-\gamma} + \frac{1}{\alpha_\tau} \frac{R}{1-\gamma} + \frac{1}{\alpha_t} \frac{R}{1-\gamma} \\ & = \frac{2R}{\alpha_\tau(1-\gamma)} \end{aligned} \quad (55)$$

Choose  $\alpha_t = c_\alpha/(1+t)^\sigma$  and substitute it into the above equation yields:

$$\begin{aligned} & \frac{LB_j^2}{2} \sum_{k=\tau}^t \alpha_k \leq \frac{LB_j^2}{2} \sum_{k=0}^{t-\tau} \alpha_k = \frac{LB_j^2 c_\alpha}{2} \sum_{k=0}^{t-\tau} \frac{1}{(1+k)^\sigma} \\ & \leq \frac{LB_j^2 c_\alpha}{2} \int_0^{t-\tau+1} x^{-\sigma} dx \leq \frac{LB_j^2 c_\alpha}{2(1-\sigma)} (t-\tau+1)^{1-\sigma} \end{aligned} \quad (56)$$

Substitute the above two inequalities into (54) and divide both sides by  $t - \tau + 1$ , we have:

$$\begin{aligned} & \frac{1}{t-\tau+1} \sum_{k=\tau}^t (\mathbb{E}\|\nabla J(\boldsymbol{\theta}_k)\| - \frac{\epsilon_v}{2})^2 \\ & \leq \frac{B_\theta}{t-\tau+1} \sum_{k=\tau}^t \sqrt{\mathbb{E}[\|\nabla J(\boldsymbol{\theta}_k)\|^2]} \sqrt{\mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_k\|^2]} \\ & \quad + \frac{2R}{\alpha_\tau(1-\gamma)(t-\tau+1)} + \frac{LB_j^2 c_\alpha}{2(1-\sigma)} \frac{1}{(t-\tau+1)^\sigma} \\ & \quad + \epsilon_m + \frac{\epsilon_v^2}{4}. \end{aligned} \quad (57)$$

Assume  $t > 2\tau - 1$ , two terms on the RHS of Eq. (57) can be bounded as follows:

$$\frac{2R}{\alpha_\tau(1-\gamma)(t-\tau+1)} = \frac{2c_\alpha R}{(1+\tau)^\sigma(t-\tau+1)} \leq \frac{2c_\alpha R}{(1+\tau)^{\sigma\tau}}. \quad (58)$$

$$\frac{LB_j^2 c_\alpha}{2(1-\sigma)} \frac{1}{(t-\tau+1)^\sigma} \leq \frac{LB_j^2 c_\alpha}{2(1-\sigma)} \frac{1}{\tau^\sigma} \quad (59)$$

Here, we analyze the first term on the RHS of Eq. (57). By Cauchy-Schwarz inequality, we have:

$$\begin{aligned} & \frac{B_\theta}{t-\tau+1} \sum_{k=\tau}^t \sqrt{\mathbb{E}[\|\nabla J(\boldsymbol{\theta}_k)\|^2]} \sqrt{\mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_k\|^2]} \\ & \leq \frac{B_\theta}{t-\tau+1} \sqrt{\sum_{k=\tau}^t \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_k)\|^2]} \sqrt{\sum_{k=\tau}^t \mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_k\|^2]} \end{aligned} \quad (60)$$

Assume the bound of approximation error of the critic  $\epsilon_v$  satisfies  $\epsilon_v < 2\mathbb{E}\|\nabla J(\boldsymbol{\theta}_k)\|$ , substitute it into LHS of Eq. (57) yields:

$$\frac{1}{t-\tau+1} \sum_{k=\tau}^t (\mathbb{E}\|\nabla J(\boldsymbol{\theta}_k)\| - \frac{\epsilon_v}{2})^2 \quad (61)$$

$$\leq \frac{1}{t-\tau+1} \sum_{k=\tau}^t \mathbb{E}\|\nabla J(\boldsymbol{\theta}_k)\|^2$$

Denote  $X(t) := 1/(t-\tau+1) \sum_{k=\tau}^t \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_k)\|^2]$  and  $Y(t) := 1/(t-\tau+1) \sum_{k=\tau}^t \mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_k\|^2]$ , put the above inequalities into Eq. (57) yields:

$$X(t) \leq 2B_\theta \sqrt{X(t)} \sqrt{Y(t)} + \mathcal{O}(\frac{1}{\tau^\sigma}), \quad (62)$$

which is equivalent to:

$$(\sqrt{X(t)} - B_\theta \sqrt{Y(t)})^2 \leq \mathcal{O}(\frac{1}{\tau^\sigma}) + B_\theta^2 Y(t). \quad (63)$$

According to the monotonicity of the square root function, we have:

$$\sqrt{X(t)} \leq \sqrt{\mathcal{O}(\frac{1}{\tau^\sigma})} + 2B_\theta \sqrt{Y(t)}. \quad (64)$$

Note that if  $A < B + C$  and A, B, C are non-negative, then  $A^2 \leq B^2 + C^2$ . Therefore:

$$X(t) \leq \mathcal{O}(\frac{1}{\tau^\sigma}) + 4B_\theta^2 Y(t) \quad (65)$$

For  $Y(t)$ , we have:

$$\begin{aligned}
Y(t) &:= \frac{1}{(t-\tau+1)} \sum_{k=\tau}^t \mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_k\|^2] \\
&\leq \frac{2}{t} \sum_{k=\tau}^t \mathbb{E}[\|\boldsymbol{\omega}^* - \boldsymbol{\omega}_k\|^2]
\end{aligned} \tag{66}$$

Finally, substitute the above inequalities into Eq. (57) yields:

$$\begin{aligned}
\min_{\tau \leq k \leq t} \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_k)\|^2] &\leq \frac{1}{t-\tau+1} \sum_{k=\tau}^t \mathbb{E}[\|\nabla J(\boldsymbol{\theta}_k)\|^2] \\
&\leq \frac{B_{\boldsymbol{\theta}}}{t-\tau+1} 2B_{\boldsymbol{\theta}} Y(t) + \frac{2c_{\alpha} R}{(1+\tau)^{\sigma}(t-\tau+1)} + \\
&\quad \frac{LB_j^2 c_{\alpha}}{2(1-\sigma)} \frac{1}{\tau^{\sigma}} + \epsilon_m + \frac{\epsilon_v^2}{4} \\
&= \mathcal{O}\left(\frac{\mathcal{E}(t)}{\tau}\right) + \mathcal{O}\left(\frac{1}{\tau^{\sigma}}\right)
\end{aligned} \tag{67}$$

The proof is completed.  $\square$