

Adaptive Event-Triggered Policy Gradient for Multi-Agent Reinforcement Learning

Umer Siddique, Abhinav Sinha, *Senior Member, IEEE*, and Yongcan Cao, *Senior Member, IEEE*

Abstract—Conventional multi-agent reinforcement learning (MARL) methods rely on time-triggered execution, where agents sample and communicate actions at fixed intervals. This approach is often computationally expensive and communication-intensive. To address this limitation, we propose ET-MAPG (Event-Triggered Multi-Agent Policy Gradient reinforcement learning), a framework that jointly learns an agent’s control policy and its event-triggering policy. Unlike prior work that decouples these mechanisms, ET-MAPG integrates them into a unified learning process, enabling agents to learn not only what action to take but also when to execute it. For scenarios with inter-agent communication, we introduce AET-MAPG, an attention-based variant that leverages a self-attention mechanism to learn selective communication patterns. AET-MAPG empowers agents to determine not only when to trigger an action but also with whom to communicate and what information to exchange, thereby optimizing coordination. Both methods can be integrated with any policy gradient MARL algorithm. Extensive experiments across diverse MARL benchmarks demonstrate that our approaches achieve performance comparable to state-of-the-art, time-triggered baselines while significantly reducing both computational load and communication overhead.

Index Terms—Multi-agent reinforcement learning, event-triggered learning and control, self-attention, data-driven control.

I. INTRODUCTION

Event-triggered control (ETC) [1] is an approach in which signals are updated or exchanged only when a certain state or output condition is met, rather than at fixed, periodic intervals. The main goal of ETC is to reduce communication and computation while maintaining closed-loop performance, in contrast to time-triggered control (TTC). ETC has been widely studied [2]–[5], where most of these methods assume access to an accurate system dynamics or model. Although this may be possible in small-scale simulated environments, it’s often unrealistic or nearly impossible in complex real-world applications.

To mitigate this model dependence, data-driven ETC methods have been recently proposed. These methods include discrete-time formulations [3], [5]–[8] but continuous-time approaches are scarce, e.g., [4]. For linear time-invariant systems, several works simplify learning by ignoring disturbances during offline data collection [3], [5], [6]. In

contrast, the work in [8] proposes a more realistic method that includes designing the controller and trigger from a single batch of noisy data, thereby accounting for disturbances and measurement errors. Moreover, the work in [4] incorporates disturbances both during learning and in closed-loop operation via a dynamic triggering strategy that guarantees \mathcal{L}_2 stability.

Reinforcement learning (RL) has achieved strong empirical results in sequential decision-making and control, including robotics [9]–[12]. Yet, most of the RL works focus on designing a time-triggered control policy, often overlooking communication cost. While a few model-free RL-based ETC methods exist [7], [13]–[16], they are developed for a single agent. In practice, many systems are inherently multi-agent systems with tight bandwidth constraints. In MARL, multiple agents are interacting, learning, and coordinating with each other to solve a shared task, making event-triggered learning even more important, where multiple agents should act and communicate only when necessary.

Communication in MARL is essential for coordination and efficient problem-solving, especially under partial observability. Early work in MARL introduced deep distributed recurrent Q-networks (DDRQN), e.g., [17], which demonstrate that agents can learn communication protocols for coordination. Building on this, the work in [18] proposed Reinforced Inter-Agent Learning (RIAL) and Differentiable Inter-Agent Learning (DIAL) to learn communication end to end, and the authors in [19] investigated communication scheduling with relational inductive biases. However, these approaches often rely on specialized communication networks or extensive parameter sharing, which can be costly at scale. Inspired by the success of self-attention [20], [21], we instead learn compact, attention-based messages where agents, during the learning phase, compute attention scores over shared representations and exchange only the relevant information.

ETC has also been explored in MARL to some extent. ETCNet [22] reduces bandwidth by sending messages only when necessary, but all agents still interact with the environment at every time step, and triggering applies only to inter-agent communication. ETMAPPO [23] integrates ETC with multi-agent proximal policy optimization algorithms via a Beta strategy to compress transmitted information and accelerate convergence in specific UAV environments. Although their model-free multi-agent PPO method performs well in the anti-UAV jamming scenario, their approach also

U. Siddique and Y. Cao are with the Unmanned Systems Lab, Department of Electrical and Computer Engineering, The University of Texas at San Antonio, San Antonio, TX 78249, USA. (e-mails: muhammad-umer.siddique@my.utsa.edu, yongcan.cao@utsa.edu). A. Sinha is with the GALACxIS Lab, Department of Aerospace Engineering and Engineering Mechanics, University of Cincinnati, OH 45221, USA. (email: abhinav.sinha@uc.edu).

applies ETC only to communication among agents.

To address these limitations, we propose Event-Triggered Multi-Agent Policy Gradient reinforcement learning (ET-MAPG), which jointly learns both the control action head and an event-trigger head for each agent and decides when to update the action and/or communicate. Unlike approaches that learn triggering conditions and control actions with separate policies [14], [24], [25], ET-MAPG employs a single shared network with two heads. Due to this, ET-MAPG reduces the policy network parameters, latency, and improves efficiency, which is crucial when scaling many agents in MARL. When inter-agent communication is allowed, we further introduce AET-MAPG, an attention-based variant of ET-MAPG that leverages self-attention to facilitate selective, learned message passing during training. As a consequence of the proposed approach, the communication graph in AET-MAPG is inherently sparse since an agent resamples an action or transmits messages only when its triggering condition is satisfied. Otherwise, it reuses its previous action and suppresses messaging. This design reduces communication and computation while supporting stable and efficient deployment. Our main contributions are summarized as follows:

- We propose ET-MAPG, a method that jointly learns a control action head and an event-trigger head for each agent. The trigger head determines when a new action should be sampled, thereby providing an improvement over approaches that learn triggering and control with separate policies.
- We further propose AET-MAPG, an attention-based variant of ET-MAPG that uses self-attention as a communication mechanism during training, which can improve the coordination efficiency.
- We demonstrate the generality of ET-MAPG and AET-MAPG by integrating them with three state-of-the-art multi-agent policy gradient algorithms, including IPPO [26], MAPPO [27], and IA2C [28].
- Through our extensive experiments, we show that ET-MAPG and AET-MAPG match the performance of standard MARL algorithms while reducing computation cost by up to 50%.

II. BACKGROUND AND PRELIMINARIES

We consider a multi-agent system of N agents, indexed by $i \in \mathcal{I} = \{1, \dots, N\}$, communicating over a fully connected undirected graph $\mathcal{G} = (\mathcal{I}, \mathcal{E})$. The dynamics of each agent are governed by a discrete-time nonlinear equation given by

$$\mathbf{x}_{i,k+1} = \mathbf{f}_i(\mathbf{x}_{i,k}, \mathbf{u}_{i,k}, \{\mathbf{x}_{j,k}\}_{j \in \mathcal{N}_i}), \quad (1)$$

where $\mathbf{x}_{i,k}$ is the state of agent i , $\mathbf{u}_{i,k}$ is its control input, and \mathcal{N}_i is the set of its neighbors. We model the problem of multi-agent event-triggered learning as a decentralized partially observable Markov decision process (Dec-POMDP) which is defined by a tuple $\langle \mathcal{I}, \mathcal{X}, \{\mathcal{U}_i\}, \mathcal{P}, r, \{O_i\}, \mathcal{Z}, \gamma \rangle$, where $\mathcal{I} = \{1, \dots, N\}$ is the set of N agents. The true state of the environment is $\mathbf{x} \in \mathcal{X}$. At each timestep k , every

agent $i \in \mathcal{I}$ selects an action $\mathbf{u}_{i,k} \in \mathcal{U}_i$ from its individual action set. This forms a joint action $\mathbf{u}_k = (\mathbf{u}_{1,k}, \dots, \mathbf{u}_{N,k}) \in \mathcal{U}$, where $\mathcal{U} = \times_{i \in \mathcal{I}} \mathcal{U}_i$ is the joint action space. The joint action governs the state transition according to the probability function $\mathcal{P}(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k)$. In this cooperative setting, the team of agents receives a single shared reward, $\mathbf{r}_k = \mathbf{r}(\mathbf{x}_k, \mathbf{u}_k)$. The agents, however, do not observe the true state \mathbf{x}_k . Instead, after the transition to state \mathbf{x}_{k+1} , each agent i receives a private observation $\mathbf{o}_{i,k+1} \in \mathcal{O}_i$. The joint observation $\mathbf{o}_{k+1} = (\mathbf{o}_{1,k+1}, \dots, \mathbf{o}_{N,k+1})$ is determined by the observation function $\mathcal{Z}(\mathbf{o}_{k+1} \mid \mathbf{x}_{k+1}, \mathbf{u}_k)$. Finally, $\gamma \in [0, 1)$ is the discount factor.

Each agent maintains a local action-observation history, $\tau_{i,k} = (\mathbf{o}_{i,0}, \mathbf{u}_{i,0}, \dots, \mathbf{u}_{i,k-1}, \mathbf{o}_{i,k})$. Actions are chosen according to a local, stochastic policy, $\mathbf{u}_{i,k} \sim \pi_i(\cdot \mid \tau_{i,k})$. The team of agents aims to learn a joint policy π that factorizes into the local policies

$$\pi(\mathbf{u}_k \mid \tau_k) = \prod_{i=1}^N \pi_i(\mathbf{u}_{i,k} \mid \tau_{i,k}), \quad (2)$$

where $\tau_k = (\tau_{1,k}, \dots, \tau_{N,k})$ is the joint history. The objective is to find a joint policy that maximizes the expected discounted return

$$J(\pi) = \mathbb{E}_{\pi, \mathcal{P}, \mathcal{Z}} \left[\sum_{k=0}^{\infty} \gamma^k \mathbf{r}(\mathbf{x}_k, \mathbf{u}_k) \right]. \quad (3)$$

Remark 1. In this model, we adopt the centralized training with decentralized execution (CTDE) paradigm (see [17], [18], [26]–[28]), which combines the benefits of having global information during training and decentralized scalability at execution. In CTDE, agents are assumed to have access to the full state of the system during training, which helps in mitigating non-stationarity in dynamic environments (a common challenge that arises when multiple agents interact with a shared environment, causing the dynamics to shift from the perspective of each agent). However, during execution, agents' actions only depend on their local observations, which is crucial for real-world deployment where global state information is usually unavailable or may not be possible to get due to bandwidth, latency, or privacy constraints.

To conserve computational and communication resources, we depart from the standard time-triggered paradigm. Instead, we employ an event-triggered scheme where each agent decides independently when to compute and broadcast a new action. Each agent i maintains its own sequence of event times $\{t_j^i\}_{j \in \mathbb{N}}$. At an event time t_j^i , it updates its history τ_{i,t_j^i} and computes a new action by sampling from its policy, $\mathbf{u}_{i,t_j^i} \sim \pi_i(\cdot \mid \tau_{i,t_j^i})$. This action is then held constant until the next event

$$\mathbf{u}_{i,k} = \mathbf{u}_{i,t_j^i}, \quad \forall k \in [t_j^i, t_{j+1}^i). \quad (4)$$

The next event time, t_{j+1}^i , is determined by a local triggering rule based on an error signal, $\mathbf{e}_{i,k} = \mathbf{x}_{i,k} - \mathbf{x}_{i,t_j^i}$, such that

$$t_{j+1}^i = \inf\{k > t_j^i \mid \mathcal{F}_i(\mathbf{x}_{i,k}, \mathbf{e}_{i,k}) \geq 0\}. \quad (5)$$

The goal is to co-design the policies $\{\pi_i\}$ and triggering functions $\{\mathcal{T}_i\}$ to maximize the expected return $J(\pi)$ while significantly reducing the frequency of policy evaluations and network communication.

Previously, the approaches often relied on parameter sharing (e.g., latent information) or specialized architectures for message exchange, both of which typically assume access to full state information or sufficient bandwidth (see e.g., [17], [18] and references therein). To address these limitations, we leverage *attention mechanisms* as a communication tool. Self-attention computes queries Q , keys K , and values V , and evaluates attention

$$\mathcal{A}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V, \quad (6)$$

where d_k is the dimensionality of K . Multi-head self-attention extends this by projecting Q, K, V into multiple subspaces, applying attention in parallel, and concatenating the outputs. This allows agents to selectively focus on relevant information and capture diverse interaction patterns. By integrating this multi-head attention as a communication mechanism with event-triggered learning, we let MARL agents decide *what*, *how*, and *when* to communicate. Note that we assume that the agents communicate over a complete undirected graph at the triggering instants, which indicates a bidirectional messaging sharing at the triggering instants only.

III. PROPOSED APPROACH

We now present the proposed event-triggered framework in multi-agent settings. Since our objective is to address large-scale multi-agent problems with high-dimensional (or possibly even continuous state-action spaces in general), we focus on multi-agent deep policy gradient methods such as IPPO [26], MAPPO [27], and IA2C [28]. Among these methods, IPPO and IA2C learn independent actors and critics for each agent, whereas MAPPO employs a centralized critic with decentralized actors. Since all these methods share the same actor-critic architecture, the proposed framework can be easily extended to any of them.

In a cooperative setting, the goal of agents is to learn a joint policy $\pi_\theta = (\pi_{1,\theta}, \dots, \pi_{n,\theta})$ that maximizes the discounted sum of rewards given by

$$\max_{\pi_\theta} J(\pi_\theta) = \max_{\pi_\theta} \left(\mathbb{E}_{\pi, \mathcal{P}, \mathcal{Z}} \left[\sum_{k=0}^{\infty} \gamma^k \mathbf{r}(\mathbf{x}_k, \mathbf{u}_k) \right] \right), \quad (7)$$

where \mathcal{P} denotes the environment transition dynamics.

Remark 2. In independent learning settings, where agents are learning independently by observing their local observations and performing actions individually, this decomposes into maximizing each agent's expected return in the form of

$$J_i(\pi_{i,\theta}) = \mathbb{E}_{\pi_{i,\theta}} \left[\sum_{k=0}^{\infty} \gamma_i^k r_{i,k} \right], \quad \forall i = 1, \dots, n. \quad (8)$$

To maximize (7), standard MARL algorithms update policies at every time step, which is equivalent to TTC. Traditional ETC designs have attempted to subdue the limitations of time-triggered execution, although they design triggering conditions manually (e.g., based on state deviations) and treat them separately from the control policy [7], [14], [24]. To address this issue, we propose ET-MAPG, an independent learning MARL algorithm in which each agent i jointly learns both its control action $\mathbf{u}_{i,k}$ and its triggering condition \mathcal{T}_i such that

$$(\mathcal{T}_i, \mathbf{u}_{i,k}) = \pi_{i,\theta}(\mathbf{z}_{i,k}, \tau_{i,k}), \quad \forall \mathbf{z}_{i,k} \in \mathcal{Z}. \quad (9)$$

Proposition 1. *Let each agent i in a multi-agent system, at a given timestep k , maintain a local observation $\mathbf{z}_{i,k}$ and a state-action history $\tau_{i,k}$. Consider a parametric policy $\pi_{i,\theta}$ that jointly outputs the control action $\mathbf{u}_{i,k}$ and the triggering decision \mathcal{T}_i as given in (9). If agent i maximizes the expected cumulative reward with a triggering regularization*

$$J_i(\pi_{i,\theta}) = \mathbb{E}_{\pi_{i,\theta}} \left[\sum_{k=0}^{\infty} \gamma_i^k r_i(\mathbf{z}_{i,k}, \mathbf{u}_{i,k}) \right] - \Psi \cdot \mathbb{I}(\mathcal{T}_i = 1), \quad (10)$$

where Ψ penalizes frequent triggering, and \mathbb{I} is the indicator function, then the optimal policy $\pi_{i,\theta}^* = \arg \max_{\pi_{i,\theta}} J_i(\pi_{i,\theta})$ simultaneously learns what action to take and when to update.

Remark 3. This unified approach reduces the model complexity by avoiding the need for hand-crafted triggers and improves sample efficiency by allowing agents to dynamically decide both what action to take and when to update. In fact, if $\mathcal{T}_i = 1$ at every timestep, then a higher penalty discourages frequent policy updates and incentivizes efficient communication and computation by trading off performance with triggering frequency.

Following the policy gradient theorem [29], the gradient for each agent i is computed as

$$\nabla_{i,\theta} J(\pi_{i,\theta}) = \mathbb{E}_{\pi_{i,\theta}} \left[\mathbf{A}_{i,k}(\mathbf{z}_{i,k}, \mathbf{u}_{i,k}) \nabla_{i,\theta} \log \pi_{i,\theta}(\mathbf{u}_{i,k} | \mathbf{z}_{i,k}) \right] - \Psi \cdot \mathbb{I}(\mathcal{T}_i = 1), \quad (11)$$

where $\mathbf{A}_{i,k}$ is the vector advantage function operated componentwise at a given timestep. Positive components indicate improvement for that particular agent, and negative components indicate a worse policy than the baseline.

Remark 4. Our proposed method ET-MAPG is model-agnostic and is sufficiently general to be extended to a large class of MARL algorithms.

As an illustrative example, we first extend IPPO [26] to incorporate the proposed framework by employing the clipped surrogate PPO objective augmented with the triggering penalty

$$\mathbb{E}_{\pi_{i,\theta}} \left[\min \left(\rho_{i,\theta} \mathbf{A}_{i,k}^{\text{IPPO}}, \bar{\rho}_{i,\theta} \mathbf{A}_{i,k}^{\text{IPPO}} \right) \right] - \Psi \cdot \mathbb{I}(\mathcal{T}_i = 1), \quad (12)$$

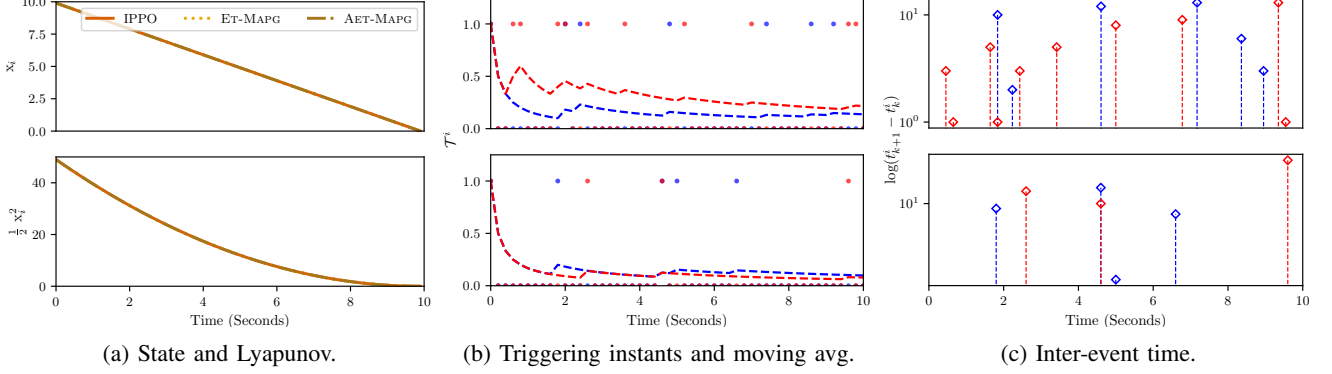


Fig. 1: Performance comparison between IPPO and our proposed methods in the perturbed multi-agent single integrators. Subfigures (b) and (c) show ET-MAPG (top) and AET-MAPG (bottom) results, respectively.

where $\rho_{i,\theta}$ and $\bar{\rho}_{i,\theta}$ denote PPO's ratio and clipping terms for agent i . The advantage $\mathbf{A}_{i,k}^{\text{IPPO}}$ is estimated using TD(λ) given by

$$\mathbf{A}_{i,k}^{\text{IPPO}} = \sum_k (\gamma_i \lambda)^{k-1} (r_i(\mathbf{z}_{i,k}, \mathbf{u}_{i,k}) + \gamma \mathbf{V}_{i,\theta}(\mathbf{z}_{i,k+1}) - \mathbf{V}_{i,\theta}(\mathbf{z}_{i,k})),$$

where $\mathbf{V}_{i,\theta}$ is the vector value function for agent i . As a consequence of this modular design, ET-MAPG can be extended to other policy gradient MARL algorithms.

For instance, with MAPPO [27], the centralized critic enables advantage estimation as

$$\mathbf{A}_{i,k}^{\text{MAPPO}} = \sum_k (\gamma_i \lambda)^{k-1} (r_i(\mathbf{z}_{i,k}, \mathbf{u}_{i,k}) + \gamma \mathbf{V}_{i,\theta}(\mathbf{x}_{k+1}) - \mathbf{V}_{i,\theta}(\mathbf{x}_k)),$$

where \mathbf{x}_k denotes the global state. Similarly, in IA2C, the advantage reduces to

$$\mathbf{A}_{i,k}^{\text{IA2C}} = r_i(\mathbf{z}_{i,k}, \mathbf{u}_{i,k}) - \mathbf{V}_{i,\theta}(\mathbf{z}_{i,k}).$$

Therefore, our approach provides a general event-triggered extension to any policy gradient MARL algorithms.

While ET-MAPG efficiently learn both the control action and policy and triggering condition, it assumes agents act independently without explicit communication. However, many cooperative tasks require coordination, where agents can communicate to reach a consensus or mitigate non-stationarity [30]. To this end, we propose AET-MAPG, a variant of ET-MAPG which integrates event-triggered communication with self-attention.

Proposition 2. *Let ET-MAPG be an independent learning framework in which each agent i jointly learns its control action $\mathbf{u}_{i,k}$ and the triggering decision \mathcal{T}_i . AET-MAPG is a variant of ET-MAPG that integrates event-triggered communication with self-attention, such that when $\mathcal{T}_i = 1$, agent i broadcasts its learned message to all other agents over*

\mathcal{G} . Each agent aggregates received messages via multi-head self-attention

$$b_i = \sum_j \alpha_{ij} v_j; \quad \alpha_{ij} = \text{softmax} \left(\frac{[\mathbf{Q}]_i [\mathbf{K}_j]^\top}{\sqrt{d_k}} \right), \quad (13)$$

where α_{ij} denotes the attention weight of agent i attending to agent j . This mechanism allows agents to selectively exchange information only when triggering conditions are met, ensuring both coordination and efficiency.

The aggregated message b_i is then fused into the policy network. With multiple attention heads ($h = 4$ in our experiments), agents capture diverse interaction patterns, improving robustness in cooperative settings.

IV. EXPERIMENTAL RESULTS

To demonstrate the effectiveness of our proposed methods, we evaluate them across diverse environments with different levels of complexity and communication demands. Specifically, we consider (i) perturbed chain of single integrators, (ii) the repeated penalty matrix game [31], and (iii) multi-agent particle environments (MPE) [32], which are a collection of 2D simulated physics environments that require cooperation or competition among agents. Via these experiments, we evaluate both control and coordination abilities of our methods. For all experiments, we perform hyperparameter optimization and report results for the best-performing configurations. Furthermore, for the sake of reproducibility, we run each experiment with five different seeds and report the mean performance. Unless otherwise specified, we benchmark ET-MAPG and AET-MAPG against IPPO [26] in the main experiments and show the generality of our methods with other policy gradient methods in ablation studies.

In the first case, the agents seek stabilization of the origin from different initial conditions. At each time step k , agent i updates according to $\mathbf{x}_{i,k+1} = \mathbf{x}_{i,k} + \mathbf{u}_{i,k} T_s$, where T_s is the sampling time. The reward function for agent i is given by the improvement in state convergence, which means encouraging an agent to reduce the absolute value of the state (i.e., getting closer to the equilibrium point) while penalizing large or

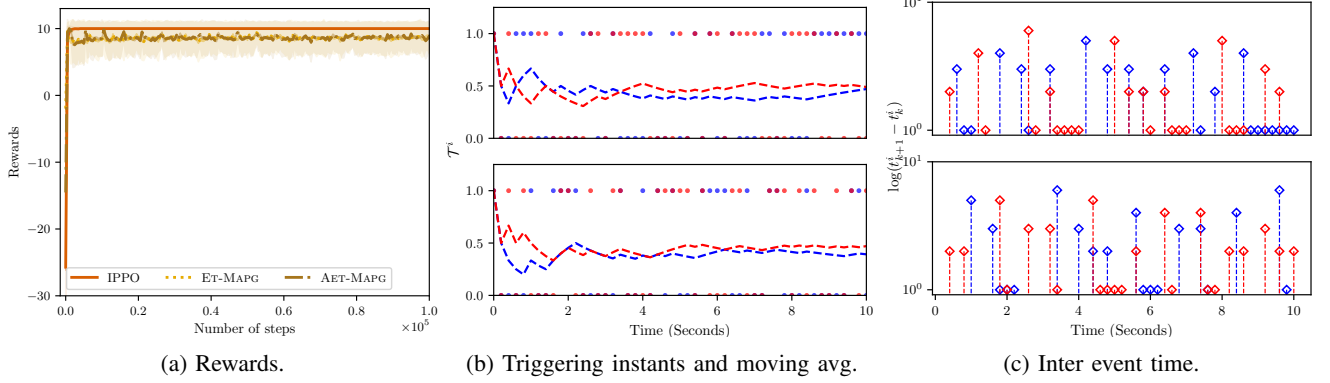


Fig. 2: Performance comparison between IPPO and our proposed methods in the repeated matrix game. Subfigures (b) and (c) show TIGER (top) and ATT-TIGER (bottom) results, respectively.

unnecessary control actions. In a cooperative setting, the goal of all agents is to reach the equilibrium point.

The results in Fig. 1 show that both proposed methods successfully stabilize the system while requiring substantially fewer action updates. The initial state value, which is 10, is driven to equilibrium in each case. However, ET-MAPG and AET-MAPG achieve this while reducing triggering events by at least 60%. Fig. 1a (top) illustrates the state trajectories of the system, while the bottom plot shows the decay of the standard quadratic Lyapunov function over time for a representative agent i . In both cases, ET-MAPG and AET-MAPG match IPPO performance. Fig. 1b demonstrates the communication frequency with the policy, where the triggering frequencies for ET-MAPG (top) and AET-MAPG (bottom) are significantly lower than those of the IPPO baseline, which triggers policy interaction at each time step. Here, blue and red triggering instants and moving average curves correspond to the two agents in the environment. Finally, Fig. 1c shows the inter-event times across agents, demonstrating that they remain strictly positive and adaptive, avoiding Zeno behavior while ensuring stable convergence.

We now test the proposed algorithms on the repeated penalty matrix game by Claus and Boutilier [31], which is a two-agent cooperative setting defined by the payoff matrix

$$\begin{bmatrix} \ell & 0 & 10 \\ 0 & 2 & 0 \\ 10 & 0 & \ell \end{bmatrix},$$

with $\ell \leq 0$ (set to $\ell = -100$ for more complexity). At each timestep, agents observe their local state, perform their own actions, and receive rewards. Specifically, agent 1 selects a row as its action and agent 2 selects a column, producing a joint payoff from the corresponding matrix entry (the entry where their choices intersect). To earn the highest reward of 10, the agents must select the correct combination of actions simultaneously. Alternatively, certain actions are considered safe, which guarantee a smaller but reliable reward of 2, regardless of what the other agent chooses. At the same time, failure to cooperate may incur the penalty of

$\ell = -100$. Each episode length is 25, where each position in the matrix remains the same and agents only rely on a constant observation, which makes the environment stateless aside from episode progression. As the penalty is strongly negative, even small deviations from cooperative behavior lead to catastrophic consequences, and agents fall into the local Nash equilibria (i.e., agents keep choosing the safe actions that yield rewards equal to 2).

Fig. 2a demonstrates the rewards achieved by our proposed methods compared to the IPPO, indicating that our methods perform comparably to the IPPO. Although our methods have slightly lower rewards, their triggering frequency is significantly lower than the standard IPPO, which triggers at every time step (see Fig. 2b). The inter-event time in Fig. 2c confirms that our methods maintain strictly positive intervals, adapting the triggering schedule dynamically and avoiding Zeno behavior. These results show that both ET-MAPG and AET-MAPG preserve performance while reducing communication overhead, even in sparse environments with high-risk coordination.

Our third evaluation domain is the Multi-Agent Particle Environments (MPE), a widely used suite of continuous 2D tasks where agents with simple dynamics must solve cooperative or competitive problems under partial observability. Since our focus is on cooperative MARL, we consider two tasks: *Simple Reference* and *Simple Spread*. These tasks require agents to balance control efficiency with inter-agent communication and make them suitable for testing event-triggered methods. In these environments, each agent observes its local position and velocity, relative positions of landmarks and other agents, and optional communication inputs. The action space of each agent includes *no-action*, *move-left*, *move-right*, *move-down*, and *move-up*. In *Simple Reference*, the environment consists of two agents and three landmarks. Each landmark is a fixed circular location, and each agent is assigned a private target landmark known only to the other agent. In this environment, agents act both as speakers and listeners, and their goal is to navigate to their assigned targets. The reward function in this environment

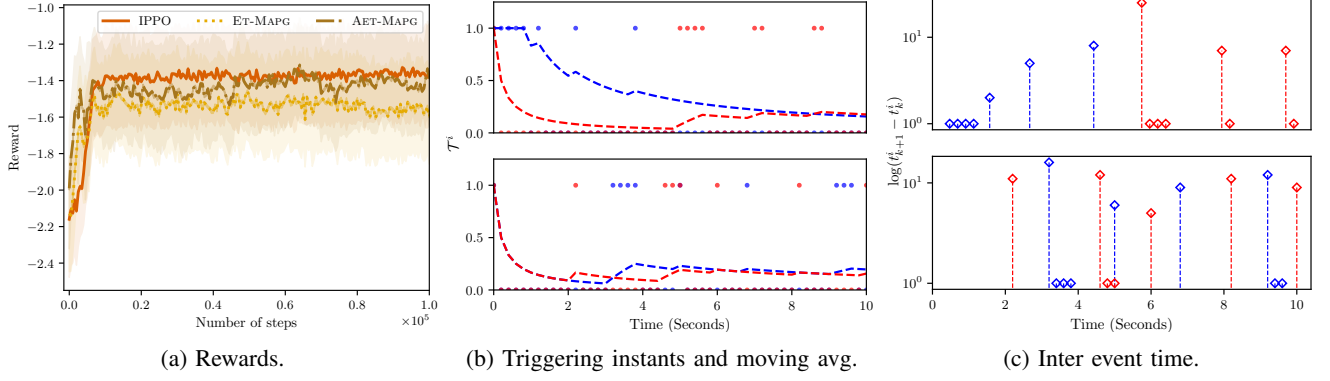


Fig. 3: Performance comparison between IPPO and our proposed methods in the Simple Reference MPE. Subfigures (b) and (c) show TIGER (top) and ATT-TIGER (bottom) results, respectively.

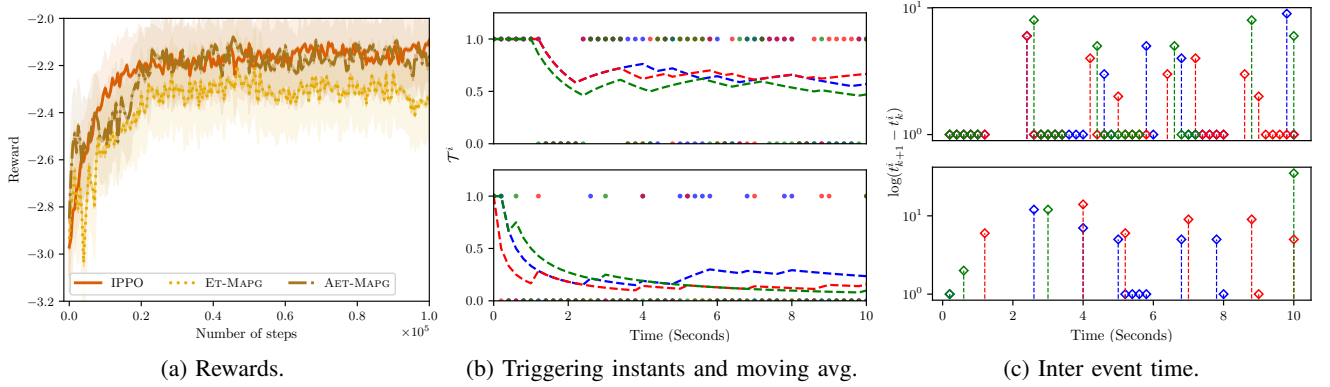


Fig. 4: Performance comparison between IPPO and our proposed methods in the Simple Spread MPE. Subfigures (b) and (c) show TIGER (top) and ATT-TIGER (bottom) results, respectively.

consists of local and global rewards, where a local reward for each agent is the negative distance to its own target, and the global reward is defined as the average distance of all agents to their respective targets. In *Simple Spread*, there are three agents and three landmarks. The goal in this environment is to cover all landmarks while avoiding collisions. Again, the reward function consists of both local and global components, where a local reward is a penalty of -1 for each collision and a global reward is the negative sum of the minimum distances from each landmark to the closest agent, encouraging agents to spread out and efficiently cover all landmarks.

Figs. 3a and 4a present the learning curves in both environments and show that while ET-MAPG significantly reduces communication frequency, it may converge to slightly lower rewards compared to the baseline. This is because, in an event-triggered setting, and especially in environments where communication is necessary, the proposed methods allow an efficient communication mechanism. Instead of constant updates, agents learn to share information only at the most critical moments. This targeted communication helps them build a better model of the environment’s dynamics and achieve tighter coordination as a team. This is evident as AET-MAPG leverages selective communication to match

the reward performance of IPPO, while also achieving the most resource-efficient behavior. Communication frequencies shown in Figs. 3b and 4b demonstrate that both methods substantially reduce communication frequencies, where AET-MAPG performs the best. Figs. 3c and 4c present the inter-event time. Once again, these results show that triggering intervals remain strictly positive and dynamically adjusted over time, which avoids Zeno behavior. These results demonstrate that our proposed event-triggered methods also generalize to high-dimensional, partially observable MPE domains, achieving both performance and resource-efficiency.

To further evaluate the generality of our proposed framework, we conduct ablation studies by integrating the joint learning of event-triggered policies with control policies in IA2C [28] and MAPPO [27], which are two other state-of-the-art MARL algorithms. We evaluate these methods in the multi-agent single integrator environment first, owing to the fact that greater insights could be obtained from a simple environment. Since the objective in this environment is to drive the system toward the origin (the equilibrium point), event-triggered methods achieve this efficiently by reusing previously sampled actions rather than resampling at every timestep, as required in standard MARL algorithms. Fig. 5

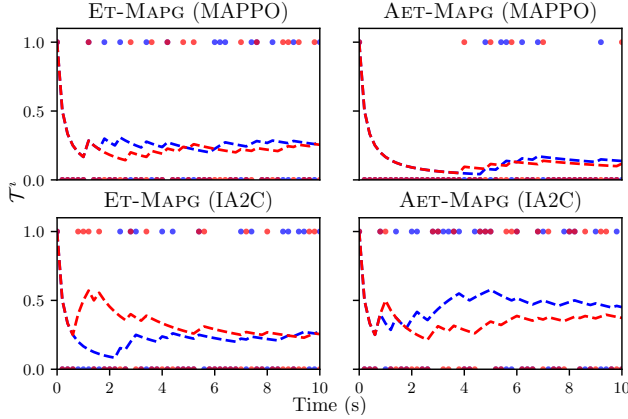


Fig. 5: Communication frequencies of event-triggered MAPPO and IA2C, along with their attention-based variants in multi-agent single-integrator environment.

presents the communication frequencies of event-triggered MAPPO and IA2C, along with their attention-based counterparts. The results show that both event-triggered MAPPO and IA2C significantly reduce the number of action samples compared to their standard MARL baselines that are time-triggered. Although IA2C demonstrates weaker performance relative to MAPPO and IPPO, it still achieves over 50% reduction in communication. Fig. 6 shows the inter-event times, verifying that triggering intervals remain strictly positive and adapt dynamically, thereby avoiding Zeno behavior. Overall, these results demonstrate that our framework is both effective and generalizable across different MARL paradigms. In particular, while MAPPO employs a centralized critic with decentralized actors and IPPO/IA2C adopts a fully independent actor-critic architecture, our event-triggered methods integrate seamlessly with these architectures while retaining their performance.

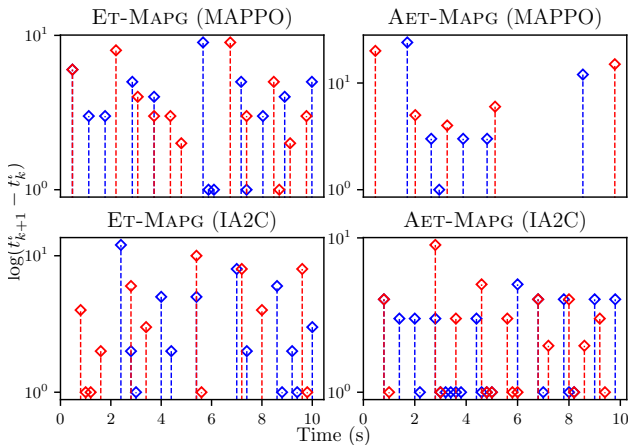


Fig. 6: Inter event time of event-triggered MAPPO and IA2C, along with their attention-based variants in multi-agent single-integrator environment.

V. CONCLUSIONS AND FUTURE WORK

In this paper, we introduced a novel framework for event-triggered MARL. We proposed ET-MAPG, a method that jointly learns a control action head and an event-trigger head for each agent within a single policy network, where the triggering head dynamically determines when a new action should be sampled from the action head, in contrast to prior approaches that rely on separate policies for control and triggering. Building on this, we further proposed AET-MAPG, an attention-based variant of ET-MAPG that incorporates self-attention as a communication mechanism during training. By enabling agents to share messages only when triggering conditions are satisfied selectively, AET-MAPG achieves efficient coordination while maintaining a sparse communication graph. Through extensive experiments across both control and standard MARL benchmarks, we demonstrated that ET-MAPG and AET-MAPG achieve performance comparable to MARL baselines while reducing communication and computation costs by up to 50%. Finally, we show that our proposed methods are general and can be seamlessly integrated with any multi-agent policy gradient methods, including IPPO, MAPPO, and IA2C.

Our work also has some limitations. First, our current framework is restricted to discrete action spaces. Second, communication in AET-MAPG is assumed to occur over a complete undirected graph, where all agents share messages bidirectionally. Third, our framework is currently suitable only for policy gradient MARL methods. As future work, we plan to address these limitations by extending our methods to continuous action spaces and more complex nonlinear systems, developing techniques for communication over dynamic graphs, and extending our framework to value-based MARL approaches.

REFERENCES

- [1] M. Miskowicz, *Event-Based Control and Signal Processing*, 1st ed. CRC Press, 2015.
- [2] A. Selivanov and E. Fridman, “Event-triggered h_∞ control: A switching approach,” *IEEE Transactions on Automatic Control*, vol. 61, no. 10, pp. 3221–3226, 2015.
- [3] V. Digge and R. Pasumathy, “Data-driven event-triggered control for discrete-time lti systems,” in *2022 European Control Conference (ECC)*, 2022, pp. 1355–1360.
- [4] W.-L. Qi, K.-Z. Liu, R. Wang, and X.-M. Sun, “Data-driven \mathcal{L}_2 -stability analysis for dynamic event-triggered networked control systems: A hybrid system approach,” *IEEE Transactions on Industrial Electronics*, vol. 70, no. 6, pp. 6151–6158, 2023.
- [5] L. A. Q. Cordovil Jr, P. H. S. Coutinho, I. Bessa, M. L. C. Peixoto, and R. M. Palhares, “Learning event-triggered control based on evolving data-driven fuzzy granular models,” *International Journal of Robust and Nonlinear Control*, vol. 32, no. 5, pp. 2805–2827, 2022.
- [6] W. Liu, J. Sun, G. Wang, F. Bullo, and J. Chen, “Data-driven self-triggered control via trajectory prediction,” *IEEE Transactions on Automatic Control*, vol. 68, no. 11, pp. 6951–6958, 2023.
- [7] D. Baumann, J.-J. Zhu, G. Martius, and S. Trimpe, “Deep reinforcement learning for event-triggered control,” in *2018 IEEE Conference on Decision and Control (CDC)*, 2018, pp. 943–950.
- [8] X. Wang, J. Berberich, J. Sun, G. Wang, F. Allgöwer, and J. Chen, “Model-based and data-driven control of event- and self-triggered discrete-time linear systems,” *IEEE Transactions on Cybernetics*, vol. 53, no. 9, pp. 6066–6079, 2023.

- [9] L. Buşoniu, T. De Bruin, D. Tolić, J. Kober, and I. Palunko, “Reinforcement learning for control: Performance, stability, and deep approximators,” *Annual Reviews in Control*, vol. 46, pp. 8–28, 2018.
- [10] D. Bertsekas, *Reinforcement learning and optimal control*. Athena Scientific, 2019, vol. 1.
- [11] J. Ibarz, J. Tan, C. Finn, M. Kalakrishnan, P. Pastor, and S. Levine, “How to train your robot with deep reinforcement learning: lessons we have learned,” *The International Journal of Robotics Research*, vol. 40, no. 4-5, pp. 698–721, 2021.
- [12] U. Siddique, A. Sinha, and Y. Cao, “On deep reinforcement learning for target capture autonomous guidance,” in *AIAA SCITECH*, 2024.
- [13] K. G. Vamvoudakis and H. Ferraz, “Model-free event-triggered control algorithm for continuous-time linear systems with optimal performance,” *Automatica*, vol. 87, pp. 412–420, 2018.
- [14] X. Zhong, Z. Ni, H. He, X. Xu, and D. Zhao, “Event-triggered reinforcement learning approach for unknown nonlinear continuous-time system,” in *2014 International Joint Conference on Neural Networks (IJCNN)*, 2014, pp. 3677–3684.
- [15] X. Yang, H. He, and D. Liu, “Event-triggered optimal neuro-controller design with reinforcement learning for unknown nonlinear systems,” *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 49, no. 9, pp. 1866–1878, 2019.
- [16] U. Siddique, A. Sinha, and Y. Cao, “Adaptive event-triggered reinforcement learning control for complex nonlinear systems,” in *2025 American Control Conference (ACC)*. IEEE, 2025, pp. 212–217.
- [17] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson, “Learning to communicate to solve riddles with deep distributed recurrent q-networks,” *arXiv preprint arXiv:1602.02672*, 2016.
- [18] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, “Learning to communicate with deep multi-agent reinforcement learning,” *Advances in neural information processing systems*, vol. 29, 2016.
- [19] D. Kim, S. Moon, D. Hostallero, W. J. Kang, T. Lee, K. Son, and Y. Yi, “Learning to schedule communication in multi-agent reinforcement learning,” *arXiv preprint arXiv:1902.01554*, 2019.
- [20] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [22] G. Hu, Y. Zhu, D. Zhao, M. Zhao, and J. Hao, “Event-triggered communication network with limited-bandwidth constraint for multi-agent reinforcement learning,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 8, pp. 3966–3978, 2021.
- [23] Z. Feng, M. Huang, Y. Wu, D. Wu, J. Cao, I. Korovin, S. Gorbachev, and N. Gorbacheva, “Approximating nash equilibrium for anti-uav jamming markov game using a novel event-triggered multi-agent reinforcement learning,” *Neural Networks*, vol. 161, pp. 330–342, 2023.
- [24] J. Lu, L. Han, Q. Wei, X. Wang, X. Dai, and F.-Y. Wang, “Event-triggered deep reinforcement learning using parallel control: A case study in autonomous driving,” *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 4, pp. 2821–2831, 2023.
- [25] J. Chen, X. Meng, and Z. Li, “Reinforcement learning-based event-triggered model predictive control for autonomous vehicle path following,” in *ACC*. IEEE, 2022, pp. 3342–3347.
- [26] C. S. De Witt, T. Gupta, D. Makoviichuk, V. Makoviychuk, P. H. Torr, M. Sun, and S. Whiteson, “Is independent learning all you need in the starcraft multi-agent challenge?” *arXiv preprint arXiv:2011.09533*, 2020.
- [27] C. Yu, A. Velu, E. Vinitzky, J. Gao, Y. Wang, A. Bayen, and Y. Wu, “The surprising effectiveness of ppo in cooperative multi-agent games,” *Advances in neural information processing systems*, vol. 35, pp. 24 611–24 624, 2022.
- [28] G. Papoudakis, F. Christianos, L. Schäfer, and S. V. Albrecht, “Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks,” in *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS)*, 2021. [Online]. Available: <http://arxiv.org/abs/2006.07869>
- [29] 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.
- [30] G. Papoudakis, F. Christianos, A. Rahman, and S. V. Albrecht, “Dealing with non-stationarity in multi-agent deep reinforcement learning,” *arXiv preprint arXiv:1906.04737*, 2019.
- [31] C. Claus and C. Boutilier, “The dynamics of reinforcement learning in cooperative multiagent systems,” *AAAI/IAAI*, vol. 1998, no. 746-752, p. 2, 1998.
- [32] I. Mordatch and P. Abbeel, “Emergence of grounded compositional language in multi-agent populations,” *arXiv preprint arXiv:1703.04908*, 2017.