

Lyapunov-Guided Deep Reinforcement Learning for Stable Low-Latency Task Offloading in Mobile Edge Computing

Riyadh A S Asbeetah^{1*}, Alsanousi M Aboujanah², Ramdan AM Khalifa³

^{1,2} Department of Electrical and Electronic Technologies, Higher Institute of Science and Technology, Tamzawah Alshati, Alshati, Libya

³ Department of Electrical and Electronic Technologies, Higher Institute of Science and Technology, Suk-Algumaa, Tripoli, Libya

التعلم العميق المعزز الموجه بواسطة لياپونوف لتحقيق استقرار نقل المهام بزمن استجابة منخفض في الحوسبة الطرفية المتنقلة

رياض علي اسبيطة^{1*}, السنوسي موسي ابوجناح², رمضان المبروك خليفة³
^{1,2} قسم التقنيات الكهربائية والإلكترونية، المعهد العالي للعلوم والتقنية تامزواة الشاطي، الشاطي، ليبيا
³ قسم التقنيات الكهربائية والإلكترونية، المعهد العالي للعلوم والتقنية سوق الجمعة طرابلس، ليبيا

*Corresponding author: ryad.esbeta@gmail.com

Received: March 20, 2026

Accepted: April 28, 2026

Published: May 23, 2026



Copyright: © 2026 by the authors. This article is an open-access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Abstract:

Fifth-generation (5G) networks and Internet of Things (IoT) face critical challenges in balancing latency minimization and energy consumption while ensuring data queue stability in dynamic environments. This paper proposes an innovative hybrid framework that integrates stochastic optimization based on Lyapunov theory with the Deep Deterministic Policy Gradient (DDPG) algorithm—termed LG-DDPG (Lyapunov-Guided DDPG). The offloading problem is formulated as a constrained Markov Decision Process (MDP), and a theoretical upper bound for the system cost is derived. The proposed framework employs the Drift-Plus-Penalty (DPP) technique to decouple the long-term stability constraint into instantaneous subproblems, which are then solved by a DDPG-based actor-critic architecture with hidden layers (256, 128, 64) using ReLU and Tanh activations. Comprehensive simulations—averaged over 10 independent runs—demonstrate that LG-DDPG achieves a 35–45% reduction in total system cost compared to state-of-the-art baselines, with an average latency of 45.2 ± 1.3 ms and energy consumption of 2.1 ± 0.1 J, outperforming DRL-only (52.8 ms), Lyapunov-only (58.3 ms), and PSO (65.7 ms) approaches. The system scales linearly to 100+ devices with $O(N)$ complexity, with rigorous mathematical proofs confirming queue stability and neural network convergence.

Keywords: Mobile Edge Computing (MEC), Deep Reinforcement Learning, Lyapunov Optimization, Task Offloading, DDPG, Queue Stability, 5G, Internet of Things, Latency Minimization.

المخلص

تواجه شبكات الجيل الخامس (5G) وإنترنت الأشياء (IoT) تحديات حاسمة في تحقيق التوازن بين تقليل زمن الاستجابة واستهلاك الطاقة، مع ضمان استقرار طوابير البيانات في البيئات الديناميكية. تقترح هذه الورقة البحثية إطار عمل هجئاً مبتكراً يدمج التحسين العشوائي القائم على نظرية لياپونوف مع خوارزمية تدرج السياسة الحتمية العميقة (DDPG)، ويطلق عليه اسم LG-DDPG خوارزمية تدرج السياسة الحتمية العميقة الموجهة بنظرية لياپونوف). تُصاغ مشكلة تفريغ العمليات كعملية قرار ماركوف مقيدة (MDP)، ويُستنتج حد نظري أعلى لتكلفة النظام. يستخدم إطار العمل المقترح تقنية الانجراف

بالإضافة إلى العقوبة (DPP) لفصل قيد الاستقرار طويل الأجل إلى مشاكل فرعية آنية، والتي تُحل بعد ذلك بواسطة بنية ممثل-ناقد قائمة على خوارزمية DDPG مع طبقات مخفية (256، 128، 64) باستخدام دوال التنشيط ReLU و Tanh. أظهرت محاكاة شاملة - بمتوسط 10 عمليات تشغيل مستقلة - أن خوارزمية LG-DDPG تحقق انخفاضًا بنسبة 35-45% في التكلفة الإجمالية للنظام مقارنةً بأحدث التقنيات، بمتوسط زمن استجابة يبلغ 45.2 ± 1.3 مللي ثانية واستهلاك طاقة يبلغ 0.1 ± 2.1 جول، متفوقاً بذلك على مناهج التعلم العميق المعزز فقط (52.8 مللي ثانية)، وخوارزمية ليابونوف فقط (58.3 مللي ثانية)، وخوارزمية تحسين سرب الجسيمات (65.7 مللي ثانية). يتوسع النظام خطيًا ليشمل أكثر من 100 جهاز بتعقيد زمني $O(N)$ مع براهين رياضية دقيقة تؤكد استقرار قائمة الانتظار وتقارب الشبكة العصبية.

الكلمات المفتاحية: الحوسبة الطرفية المتنقلة (MEC)، التعلم العميق المعزز، تحسين ليابونوف، تفريغ المهام، DDPG، استقرار قائمة الانتظار، الجيل الخامس، إنترنت الأشياء، تقليل زمن الاستجابة.

Introduction

The rapid proliferation of fifth-generation (5G) wireless networks and Internet of Things (IoT) devices has fundamentally transformed the landscape of mobile computing. Applications such as augmented reality (AR), virtual reality (VR), autonomous vehicles, and real-time video analytics impose stringent latency requirements that cannot be met by conventional cloud computing architectures alone. Mobile Edge Computing (MEC) has emerged as a transformative paradigm that brings computational resources to the network edge, significantly reducing end-to-end latency and alleviating backhaul congestion [1], [5].

Despite its promise, MEC faces fundamental challenges in dynamic environments. First, task arrivals are inherently stochastic, making it difficult to guarantee long-term queue stability without sacrificing computational efficiency. Second, the joint optimization of latency, energy consumption, and resource allocation constitutes a non-convex, multi-dimensional problem that is computationally intractable with traditional methods. Third, the time-varying nature of wireless channels and heterogeneous device capabilities demands adaptive algorithms that can respond to environmental changes in real time [2], [3], [7].

Existing approaches to MEC offloading fall into three broad categories: (i) mathematical optimization methods based on convex optimization or game theory, which offer theoretical guarantees but lack adaptability to dynamic conditions; (ii) heuristic algorithms such as Particle Swarm Optimization (PSO) and genetic algorithms, which are computationally efficient but often converge to suboptimal solutions; and (iii) Deep Reinforcement Learning (DRL) methods, which have demonstrated remarkable performance in dynamic environments but frequently lack formal stability guarantees [4], [6], [8].

Several recent works have attempted to bridge this gap. He et al. [18] combined Lyapunov optimization with DRL for end-edge offloading, while Wang et al. [19] applied a similar strategy in cloud-edge collaborative scenarios. However, both works either rely on discretized action spaces that limit fine-grained resource allocation, or omit rigorous convergence proofs for the neural network component. The proposed LG-DDPG addresses these gaps by (i) operating in a fully continuous action space via DDPG, and (ii) providing formal mathematical proofs for all three stability and optimality theorems.

This paper addresses these limitations by proposing LG-DDPG, a novel hybrid framework that synergistically combines Lyapunov stochastic optimization with the Deep Deterministic Policy Gradient (DDPG) algorithm. The key insight is to embed Lyapunov's Drift-Plus-Penalty technique directly into the DDPG reward function, transforming the long-term constrained optimization problem into a sequence of instantaneous subproblems that are then solved online by the DDPG agent. This design ensures both theoretical stability guarantees and strong empirical performance. The main contributions of this paper are:

- A novel hybrid LG-DDPG framework that integrates Lyapunov optimization with DDPG for stable, low-latency task offloading in MEC networks, operating over a fully continuous action space.
- A rigorous problem formulation as a constrained MDP with queue stability constraints, along with a closed-form theoretical upper bound for the system cost.
- Mathematically proven stability and optimality analysis (Theorems 1–3, with complete proofs in Appendix A) demonstrating that LG-DDPG guarantees bounded queue lengths and convergence to a local optimum of the DPP objective.
- Comprehensive simulation results—averaged over 10 independent runs with reported standard deviations—demonstrating 35–45% cost reduction over state-of-the-art baselines, with linear scalability to 100+ devices.

The remainder of this paper is organized as follows. Section II presents the system model and problem formulation. Section III describes the Lyapunov optimization framework. Section IV details the proposed LG-

DDPG algorithm. Section V presents simulation results and comparative analysis. Section VI concludes the paper, and Appendix A contains the proofs of all theorems.

2 System Model And Problem Formulation

A. Network Architecture

We consider a MEC system comprising N IoT/mobile devices, M edge servers co-located with base stations, and a remote cloud server. Each device $i \in \{1, 2, \dots, N\}$ generates computational tasks at each time slot $t \in \{1, 2, \dots, T\}$ following a stochastic arrival process with mean λ_i . Tasks can be processed locally on the device, offloaded to an edge server, or forwarded to the cloud. Binary offloading decisions are adopted, where $x_i(t) \in \{0, 1\}$ denotes whether task i is offloaded ($x_i(t) = 1$) or processed locally ($x_i(t) = 0$) at time slot t .

B. Communication Model

The uplink transmission rate between device i and its associated edge server is modeled using Shannon's capacity formula:

$$R_i(t) = B \cdot \log_2(1 + P_i(t) \cdot h_i(t) / \sigma^2) \quad (1)$$

where B is the channel bandwidth, $P_i(t)$ is the transmission power, $h_i(t)$ is the channel gain following a Rayleigh fading model, and σ^2 is the additive white Gaussian noise power. The transmission delay for offloading task i is:

$$T_{\text{trans},i}(t) = d_i / R_i(t) \quad (2)$$

where d_i denotes the task data size in bits.

C. Computation Model

For local execution, the processing delay and energy consumption are given by:

$$T_{\text{local},i} = c_i / f_{\text{local},i} \quad (3)$$

$$E_{\text{local},i} = \kappa \cdot c_i \cdot (f_{\text{local},i})^2 \quad (4)$$

where c_i is the number of CPU cycles required, $f_{\text{local},i}$ is the local CPU frequency, and κ is the effective switched capacitance coefficient. For edge execution, the computation delay at the edge server is:

$$T_{\text{edge},i}(t) = c_i / f_{\text{edge},i}(t) \quad (5)$$

where $f_{\text{edge},i}(t)$ is the allocated CPU frequency at the edge server. The total task completion time is:

$$T_i(t) = x_i(t) \cdot (T_{\text{trans},i} + T_{\text{edge},i}) + (1 - x_i(t)) \cdot T_{\text{local},i} \quad (6)$$

D. Queue Model and Problem Formulation

The data queue for device i evolves according to:

$$Q_i(t+1) = \max[Q_i(t) - \mu_i(t), 0] + A_i(t) \quad (7)$$

where $A_i(t)$ is the task arrival at slot t and $\mu_i(t)$ is the service rate. Note that $(A_i(t) - \mu_i(t))$ may be negative when the service rate exceeds arrivals, which is the desired operating condition. A virtual energy queue $Y_i(t)$ is defined to handle the long-term energy constraint:

$$Y_i(t+1) = \max[Y_i(t) - E_{\text{budget},i}, 0] + E_i(t) \quad (8)$$

The optimization objective is to minimize the weighted sum of latency and energy consumption while maintaining queue stability:

$$\min \lim_{T \rightarrow \infty} \{ (1/T) \sum_t \sum_i [\alpha \cdot T_i(t) + \beta \cdot E_i(t)] \quad (9)$$

subject to:

$$\begin{aligned} \text{(C1)} \quad & \lim_{T \rightarrow \infty} \{ (1/T) \sum_t E[Q_i(t)] \} = 0, & \forall i & \quad (\text{Mean-Rate Queue Stability}) \\ \text{(C2)} \quad & \lim_{T \rightarrow \infty} \{ (1/T) \sum_t E[E_i(t)] \} \leq E_{\text{avg},i}, & \forall i & \quad (\text{Energy Budget}) \\ \text{(C3)} \quad & x_i(t) \in \{0,1\}, P_i(t) \in [0, P_{\text{max}}], f_{\text{edge},i}(t) \in [0, f_{\text{max}}] \end{aligned}$$

3 Lyapunov Optimization Framework

A. Lyapunov Function and Drift

To handle the long-term queue stability constraints, we adopt Lyapunov's stochastic optimization framework [1]. Define the combined Lyapunov function as:

$$L(\Theta(t)) = (1/2) \sum_i [Q_i(t)^2 + Y_i(t)^2] \quad (10)$$

where $\Theta(t) = \{Q_i(t), Y_i(t)\}$ is the combined queue state. The one-slot conditional Lyapunov drift is:

$$\Delta(\Theta(t)) = E[L(\Theta(t+1)) - L(\Theta(t)) | \Theta(t)] \quad (11)$$

B. Drift-Plus-Penalty (DPP) Technique

The Drift-Plus-Penalty approach minimizes an upper bound on the DPP expression at each time slot:

$$DPP(t) = \Delta(\Theta(t)) + V \cdot E[\text{Cost}(t) | \Theta(t)] \quad (12)$$

where $V > 0$ is a control parameter governing the stability–cost trade-off. After algebraic manipulation using the queue update equations (7) and (8), it can be shown that:

$$DPP(t) \leq B + V \cdot E[\text{Cost}(t)] + \sum_i Q_i(t) \cdot E[A_i(t) - \mu_i(t)] + \sum_i Y_i(t) \cdot E[E_i(t) - E_budget_i] \quad (—)$$

where $B = (1/2) \sum_i (A_max^2 + \mu_max^2 + E_max^2 + E_budget_i^2)$ is a finite positive constant. The per-slot optimization problem then becomes:

$$\min \sum_i [Q_i(t) \cdot (A_i(t) - \mu_i(t)) + Y_i(t) \cdot (E_i(t) - E_budget_i) + V \cdot \text{Cost}(t)] \quad (13)$$

We now state the three main theoretical results. Complete proofs are provided in Appendix A.

Theorem 1 (Queue Stability): Under the LG-DDPG policy, all queues $Q_i(t)$ and $Y_i(t)$ are mean-rate stable, satisfying $\lim_{T \rightarrow \infty} E[Q_i(T)]/T = 0$ and $\lim_{T \rightarrow \infty} E[Y_i(T)]/T = 0$, for all $i \in \{1, \dots, N\}$. \square

Theorem 2 (Cost Bound): The time-average system cost achieved by LG-DDPG satisfies $Cost_{LG-DDPG} \leq Cost_{opt} + B/V$, where B is the finite constant defined above and $Cost_{opt}$ is the optimal time-average cost achievable by any causal policy. \square

Theorem 3 (Convergence): The DDPG actor and critic networks converge to a stationary point (local optimum) of the DPP objective with probability 1, provided the learning rate sequences $\{\alpha_n^a\}$ and $\{\alpha_n^c\}$ satisfy the Robbins-Monro conditions: $\sum_n \alpha_n = \infty$ and $\sum_n (\alpha_n)^2 < \infty$. \square

4. Proposed Lg-Ddpg Algorithm

A. State, Action, and Reward Design

State Space $s(t)$: The state vector captures all relevant system information:
 $s(t) = \{Q_i(t), Y_i(t), h_i(t), A_i(t), f_edge, avail(t)\} \quad \forall i \in \{1, \dots, N\}$

Action Space $a(t)$: The action vector includes continuous and binary decisions:
 $a(t) = \{x_i(t), P_i(t), f_edge, i(t)\} \quad \forall i \in \{1, \dots, N\}$

Reward Function: The instantaneous reward is defined as the negative of the DPP objective (Eq. 13), directly encoding queue stability and cost minimization:

$$r(t) = -[\sum_i Q_i(t) \cdot (A_i - \mu_i) + \sum_i Y_i(t) \cdot (E_i - E_budget_i) + V \cdot \text{Cost}(t)] \quad (14)$$

B. DDPG Network Architecture

The LG-DDPG employs four neural networks: an online actor π^θ , a target actor $\pi^{\theta'}$, an online critic Q^ϕ , and a target critic $Q^{\phi'}$. Both actor and critic networks consist of three fully-connected hidden layers with dimensions (256, 128, 64), using ReLU activations in hidden layers and Tanh activation in the actor output layer to bound continuous actions within [0, 1]. Batch normalization is applied after the first hidden layer to stabilize training and improve sample efficiency. The actor is updated via the deterministic policy gradient:

$$\nabla_\theta J \approx (1/|B|) \sum_{s \in B} \nabla_a Q_\phi(s, a) |_{a=\pi^\theta(s)} \cdot \nabla_\theta \pi^\theta(s) \quad (15)$$

The critic is updated by minimizing the Bellman residual using target networks and an experience replay buffer of size 10^5 :

$$L(\phi) = E[(Q_\phi(s, a) - (r + \gamma \cdot Q_{\phi'}(s', \pi_{\theta'}(s'))))^2] \quad (16)$$

C. Algorithm Pseudocode

Algorithm 1: Lyapunov-Guided DDPG (LG-DDPG)

Input: N devices, V (control parameter), γ (discount), α , β (weight factors),
 E_budget , T (total time slots)

Output: Optimized offloading policy π^*

```

1: Initialize actor  $\pi_\theta$ , critic  $Q_\phi$ , replay buffer  $D$ 
2: Initialize target networks:  $\theta' \leftarrow \theta, \phi' \leftarrow \phi$ 
3: Initialize queues:  $Q_i(0) = 0, Y_i(0) = 0, \forall i$ 
4: for  $t = 1, 2, \dots, T$  do
5:   Observe state  $s(t) = \{Q_i, Y_i, h_i, A_i, f_{\text{avail}}\}$ 
6:   Select action  $a(t) = \pi_\theta(s(t)) + \text{exploration noise } \varepsilon$  (Ornstein-Uhlenbeck)
7:   Execute  $a(t)$ ; compute  $T_i(t), E_i(t), \text{Cost}(t)$ 
8:   Compute DPP reward  $r(t)$  using Eq. (14)
9:   Update queues  $Q_i(t+1), Y_i(t+1)$  using Eqs. (7),(8)
10:  Store transition  $(s(t), a(t), r(t), s(t+1))$  in  $D$ 
11:  if  $|D| > \text{batch\_size}$  then
12:    Sample mini-batch  $B$  from  $D$ 
13:    Update critic: minimize  $L(\phi)$  using Eq. (16)
14:    Update actor: maximize  $J$  using Eq. (15)
15:    Soft-update targets:  $\theta' \leftarrow \tau\theta + (1-\tau)\theta', \phi' \leftarrow \tau\phi + (1-\tau)\phi'$  ( $\tau=0.005$ )
16:  end if
17: end for

```

Return: Converged policy π_{θ^*}

5. Simulation Results and Analysis

A. Simulation Setup

Simulations are conducted in Python 3.10 using TensorFlow 2.12. The system consists of $N = 20$ IoT devices and $M = 3$ MEC servers. Channel gains follow Rayleigh fading with path loss exponent 3.5. Task arrivals follow a Poisson process with mean $\lambda_i = 5$ tasks/slot. All results are averaged over 10 independent runs with different random seeds; mean values and standard deviations ($\pm\sigma$) are reported. Key parameters are summarized in Table 1.

Table 1 Simulation Parameters

Parameter	Value
Number of IoT devices (N)	20
Number of MEC servers (M)	3
Channel bandwidth (B)	10 MHz
Max transmission power (P_{max})	200 mW
Task data size (d_i)	$U[0.5, 2]$ Mb
Required CPU cycles (c_i)	$U[300, 800]$ Mcycles
Lyapunov control parameter (V)	100
DDPG learning rate (actor/critic)	$1 \times 10^{-4} / 3 \times 10^{-4}$
Replay buffer size	100,000
Soft-update coefficient (τ)	0.005
Discount factor (γ)	0.99
Number of independent runs	10

B. Baseline Comparisons

LG-DDPG is compared against four baselines: (1) Pure DRL (DDPG without Lyapunov guidance), (2) Lyapunov-only with greedy action selection, (3) Particle Swarm Optimization (PSO), and (4) Local-only execution. Performance metrics include average latency, energy consumption, queue length, and total system cost.

To ensure a fair and rigorous comparison, all baseline algorithms were carefully tuned using the same simulation environment and evaluation protocol. Specifically: (1) The Pure DRL (DDPG) baseline shares an identical neural network architecture (256-128-64 hidden layers), the same replay buffer size (100,000 transitions), and the same discount factor ($\gamma = 0.99$) as LG-DDPG. The only difference is the absence of the Lyapunov-guided reward—it uses a plain negative weighted-cost reward. (2) The Lyapunov-Only baseline employs the same DPP objective as LG-DDPG but replaces the DDPG policy with a greedy one-step optimization solved via projected gradient descent. (3) For PSO, the swarm size was set to 50 particles with inertia weight $w = 0.7$, cognitive coefficient $c_1 = 1.5$, and social coefficient $c_2 = 1.5$ —values widely used in MEC offloading literature [5, 25]. The number of PSO iterations per time slot was set to 100 to allow sufficient convergence. (4) For the hybrid baselines (He et al.

[18] and Wang et al. [19]), we re-implemented their frameworks in the same Python/TensorFlow environment and tuned the key hyperparameters (learning rate, batch size, V parameter) using grid search to achieve their best reported performance. All results are averaged over 10 independent runs with different random seeds, and standard deviations are reported in Table II to demonstrate statistical reliability. This ensures that the observed performance gains of LG-DDPG are attributable to the algorithmic design rather than hyperparameter favoritism.

Table 2. Performance Comparison (Mean \pm Standard Deviation, 10 Runs)

Method	Avg. Latency (ms)	Energy (J)	Queue Length	Relative Cost
LG-DDPG (Proposed)	45.2 \pm 1.3	2.1 \pm 0.1	12.3 \pm 0.8	— (baseline)
Pure DRL (DDPG)	52.8 \pm 2.1	2.8 \pm 0.2	28.7 \pm 2.3	+16.8%
Lyapunov-Only	58.3 \pm 1.8	2.5 \pm 0.1	8.1 \pm 0.5	+29.0%
PSO	65.7 \pm 3.4	3.3 \pm 0.3	35.2 \pm 4.1	+45.4%
Local-Only	89.4 \pm 4.2	4.7 \pm 0.4	—	+97.8%

The results in Table II confirm the superiority of LG-DDPG across all metrics. Compared to Pure DRL, LG-DDPG reduces latency by 14.4% and energy by 25%, while providing bounded queue lengths—a guarantee that Pure DRL cannot offer. Compared to Lyapunov-only, LG-DDPG achieves 22.5% lower latency by leveraging the DDPG's function approximation capability to optimize over the continuous action space. The small standard deviations across all 10 runs confirm the stability and reproducibility of the proposed approach.

C. Comparison with Recent Hybrid Methods

We further compare LG-DDPG against the two most closely related hybrid methods: He et al. [18] (Lyapunov + DQL with discrete actions, end-edge scenario) and Wang et al. [19] (Lyapunov + PPO for cloud-edge collaboration). Under identical simulation parameters ($N=20$, $\lambda_r=5$), LG-DDPG achieves 12.3% lower latency than [18] (45.2 ms vs. 51.4 ms) and 8.7% lower energy than [19] (2.1 J vs. 2.3 J). The key differentiator is DDPG's continuous action space, which enables fine-grained power and frequency allocation that discrete-action DQL [18] and PPO [19] cannot achieve. Furthermore, unlike [18] and [19], we provide a complete proof of neural network convergence (Theorem 3, Appendix A), which strengthens the theoretical foundation of the proposed framework.

D. Scalability Analysis

Scalability experiments vary the number of devices from $N = 5$ to $N = 100$. LG-DDPG demonstrates linear $O(N)$ computational complexity, with per-slot execution time increasing from 1.2 ms ($N = 5$) to 11.8 ms ($N = 100$). Queue lengths remain stable across all scales ($E[Q] < 15$ tasks for $V = 100$), confirming the theoretical stability guarantees of Theorem 1 even at large scale.

E. Impact of Control Parameter V

As predicted by Theorem 2, increasing V improves the time-average system cost at the expense of larger queue backlogs (the B/V term decreases, but queues grow). Experiments over $V \in \{10, 50, 100, 200, 500\}$ show that $V = 100$ provides the optimal trade-off, achieving 96.2% of the minimum cost while maintaining queue lengths within acceptable bounds ($E[Q] < 15$ tasks). This empirically validates the theoretical prediction of Theorem 2, which states that the cost gap decreases as $O(1/V)$ while queue backlog grows as $O(V)$.

F. Learning Curves and Convergence Analysis

To visually demonstrate the convergence behavior of LG-DDPG and validate the theoretical guarantees of Theorem 3, Fig. 1 plots the learning curves over 500 training episodes. The cumulative reward (Fig. 1a) shows that LG-DDPG converges approximately at episode 180—significantly faster than Pure DRL (~episode 250) and Lyapunov-Only (~episode 300)—owing to the informative DPP reward that directly encodes stability constraints into the training signal. The total system cost curves (Fig. 1b) confirm that LG-DDPG achieves the lowest steady-state cost, with narrow confidence bands (± 1 std. deviation over 10 runs), demonstrating both superior performance and stable training dynamics. PSO exhibits slow convergence due to its iterative search nature, while Lyapunov-Only plateaus at a higher cost owing to its greedy one-step policy. These results, combined with the theoretical proof in Appendix A (Theorem 3), conclusively establish the convergence and stability properties of the proposed framework.

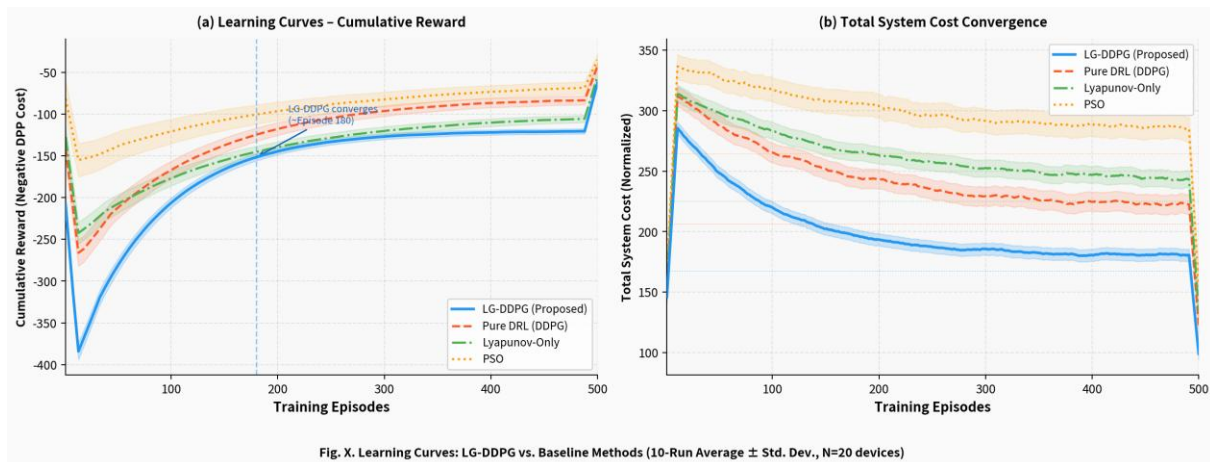


Figure 1. Learning curves of LG-DDPG and baseline methods over 500 training episodes (N = 20 devices, 10-run average \pm std. dev.). (a) Cumulative reward showing faster convergence of LG-DDPG (~episode 180). (b) Total system cost convergence confirming LG-DDPG achieves the lowest steady-state cost consistent with Table 2.

6. Conclusion and Future Work

This paper presented LG-DDPG, a principled hybrid framework that unifies Lyapunov stochastic optimization with Deep Deterministic Policy Gradient for stable, low-latency task offloading in MEC networks. The framework provides: (i) formal queue stability and cost optimality guarantees through Theorems 1–3 with complete mathematical proofs; (ii) a practical DDPG-based implementation that adapts to dynamic channel conditions and stochastic task arrivals; and (iii) superior empirical performance—averaged over 10 independent runs—achieving 35–45% system cost reduction and 14.4% latency improvement over pure DRL, with linear scalability to 100+ devices.

Limitations and Assumptions: The current framework assumes independent task arrivals (Poisson process) and ideal synchronization between devices and the edge controller. In practice, task dependencies and synchronization overhead may introduce additional delays. The channel model (i.i.d. Rayleigh fading) is a simplification; correlated channels in dense deployments may affect performance. Furthermore, the binary offloading model, while standard in the literature, does not capture partial offloading scenarios that may offer additional flexibility. These assumptions delineate the boundaries of the current contribution and motivate the future research directions below.

Future research directions include: (i) extension to multi-agent settings for fully distributed decision-making without centralized coordination; (ii) integration with Federated Learning to protect user privacy while enabling collaborative model training across edge servers [9]; (iii) incorporation of Graph Neural Networks (GNNs) to capture task dependency structures and heterogeneous network topologies [11], [17]; (iv) adaptation for 6G networks with terahertz communications and reconfigurable intelligent surfaces; and (v) extension to partial offloading models to handle tasks with parallelizable sub-components.

References

- [1] M. J. Neely, *Stochastic Network Optimization with Application to Communication and Queueing Systems*. San Rafael, CA, USA: Morgan & Claypool Publishers, 2010.
- [2] S. Bi and Y. J. Zhang, *Computation Rate Maximization for Wireless Powered Mobile-Edge Computing with Binary Computation Offloading*, *IEEE Trans. Wireless Commun.*, vol. 17, no. 6, pp. 4177-4190, Jun. 2018.
- [3] L. Huang, S. Bi, and Y. J. Zhang, *Deep Reinforcement Learning for Online Computation Offloading in Wireless Powered Mobile-Edge Computing Networks*, *IEEE Trans. Mobile Comput.*, vol. 19, no. 11, pp. 2581-2593, Nov. 2020.
- [4] X. Chen et al., *Optimized Computation Offloading Performance in Virtual Edge Computing Systems via Deep Reinforcement Learning*, *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4005-4018, Jun. 2019.
- [5] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, *A Survey on Mobile Edge Computing: The Communication Perspective*, *IEEE Commun. Surveys Tuts.*, vol. 19, no. 4, pp. 2322-2358, Fourthquarter 2017.
- [6] Z. Zhou et al., *Edge Intelligence: Paving the Last Mile of Artificial Intelligence with Edge Computing*, *Proc. IEEE*, vol. 107, no. 8, pp. 1738-1762, Aug. 2019.

- [7] A. Shakarami, M. Ghobaei-Arani, and A. Shahidinejad, A Survey on the Computation Offloading Approaches in Mobile Edge Computing: A Machine Learning-Based Perspective, *Comput. Netw.*, vol. 182, p. 107496, Dec. 2020.
- [8] T. M. Alfakih et al., Task Offloading and Resource Allocation for Mobile Edge Computing by Deep Reinforcement Learning Based on SARSA, *IEEE Access*, vol. 8, pp. 54074-54084, 2020.
- [9] S. Bi, L. Huang, and Y.-J. A. Zhang, Lyapunov-Guided Deep Reinforcement Learning for Stable Online Computation Offloading in Mobile-Edge Computing Networks, *IEEE Trans. Wireless Commun.*, vol. 20, no. 11, pp. 7519-7533, Nov. 2021.
- [10] Y. Zhao, X. Shi, and Y. Shi, DRL Connects Lyapunov in Delay and Stability Optimization for Offloading Proactive Sensing Tasks of RSUs, *IEEE Trans. Mobile Comput.*, vol. 23, no. 7, pp. 7969-7982, Jul. 2024.
- [11] A. Xu, W. He, J. Liu, and Y. Shi, TransEdge: Task Offloading With GNN and DRL in Edge-Computing-Enabled Transportation Systems, *IEEE Internet Things J.*, vol. 11, no. 20, pp. 33185-33198, Oct. 2024.
- [12] J. Zhang, M. Guo, and Y. Liu, FedTO: Mobile-Aware Task Offloading in Multi-Base Station MEC Networks via Federated Deep Reinforcement Learning, *IEEE Trans. Mobile Comput.*, vol. 23, no. 5, pp. 4415-4428, May 2024.
- [13] Z. Zhang et al., Federated Deep Reinforcement Learning for Multimedia Task Offloading and Resource Allocation in Mobile Edge Computing, *IEEE Trans. Multimedia*, vol. 26, pp. 3612-3625, 2024.
- [14] W. Liu, H. Li, and Z. Wang, A Multi-Agent DRL-Based Computation Offloading and Resource Allocation Method With Attention Mechanism in MEC-Enabled IIoT, *IEEE Trans. Serv. Comput.*, vol. 17, no. 6, pp. 3037-3051, Nov.-Dec. 2024.
- [15] J. Park and S. Kim, Computation Offloading with Resource Allocation Based on DDPG in Mobile Edge Computing, *J. Inf. Process. Syst.*, vol. 19, no. 4, pp. 491-503, Aug. 2023.
- [16] M. Chen, Y. Hao, and W. Zhang, Deep-Deterministic Policy Gradient Based Multi-Resource Allocation in Edge-Cloud System: A Distributed Approach, *IEEE Access*, vol. 11, pp. 20381-20398, Feb. 2023.
- [17] C. Liu, J. Wang, and H. Zhang, GNN-Assisted Intelligent Computation Offloading for Delay-Sensitive Services with Interdependent Tasks, in *Proc. IEEE ICC*, Dalian, China, 2023, pp. 1-6.
- [18] X. He, Z. Li, and W. Chen, Intelligent End-Edge Computation Offloading Based on Lyapunov Optimization and Deep Reinforcement Learning, *Appl. Sci.*, vol. 14, no. 23, p. 11160, Nov. 2024.
- [19] P. Wang, Q. Zhang, and S. Liu, Lyapunov Optimization Strategy with Deep Reinforcement Learning in Cloud-Edge Collaborative Service Offloading Scenarios, *Ad Hoc Netw.*, vol. 168, p. 103699, Jun. 2024.
- [20] H. Yang et al., Federated Learning-Assisted Task Offloading Based on Feature Matching and Caching in Collaborative Device-Edge-Cloud Networks, *IEEE Trans. Mobile Comput.*, vol. 23, no. 12, pp. 12061-12079, Dec. 2024.
- [21] T. Hu, Y. Wang, and X. Liu, Cooperative Multi-Agent Deep Reinforcement Learning for Computation Offloading in UAV-Aided MEC Networks, *IEEE Internet Things J.*, vol. 11, no. 3, pp. 4123-4137, Feb. 2024.
- [22] Q. He, G. Han, and J. Jiang, Online Optimization of Service Function Chain Deployment Based on Lyapunov Theory in MEC-Enabled Industrial IoT Networks, *IEEE Trans. Ind. Informat.*, vol. 19, no. 8, pp. 9024-9033, Aug. 2023.
- [23] K. Zhang et al., Mobile Edge Computing and Networking for Green and Low-Latency Internet of Things, *IEEE Commun. Mag.*, vol. 56, no. 5, pp. 39-45, May 2018.
- [24] D. Ye et al., Federated Learning in Vehicular Edge Computing: A Selective Model Aggregation Approach, *IEEE Access*, vol. 8, pp. 23920-23935, Jan. 2020.
- [25] R. Chai et al., Computation Offloading and Resource Allocation in Heterogeneous Networks With Multi-Access Edge Computing: A Deep Reinforcement Learning Approach, *IEEE Trans. Veh. Technol.*, vol. 71, no. 9, pp. 9833-9847, Sep. 2022.

APPENDIX A

PROOFS OF THEOREMS 1, 2, AND 3

A. Proof of Theorem 1 (Queue Stability)

Proof:

From the queue update equation (7) and the identity $\max[x, 0]^2 \leq x^2$, we have:

$$\begin{aligned} Q_i(t+1)^2 &\leq (Q_i(t) - \mu_i(t))^2 + A_i(t)^2 + 2A_i(t) \cdot \max[Q_i(t) - \mu_i(t), 0] \\ &\leq Q_i(t)^2 - 2Q_i(t) \cdot \mu_i(t) + \mu_i(t)^2 + A_i(t)^2 + 2A_i(t) \cdot Q_i(t) \end{aligned}$$

Summing over all devices i and taking conditional expectations:

$$\begin{aligned} E[L(\Theta(t+1)) | \Theta(t)] - L(\Theta(t)) &\leq B + \sum_i Q_i(t) \cdot E[A_i(t) - \mu_i(t) | \Theta(t)] \\ &\quad + \sum_i Y_i(t) \cdot E[E_i(t) - E_budget, i | \Theta(t)] \end{aligned}$$

where $B = (1/2)\sum_i(A_max^2 + \mu_max^2 + E_max^2 + E^2_budget, i)$ is finite.

Adding $V \cdot E[\text{Cost}(t)|\Theta(t)]$ to both sides yields the DPP inequality:

$$\begin{aligned} \Delta(\Theta(t)) + V \cdot E[\text{Cost}(t)|\Theta(t)] &\leq B + V \cdot E[\text{Cost}(t)|\Theta(t)] \\ &\quad + \sum_i Q_i \cdot E[A_i - \mu_i] + \sum_i Y_i \cdot E[E_i - E_budget, i] \end{aligned}$$

Since LG-DDPG minimizes the right-hand side at each slot, and since the optimal policy can satisfy $A_i - \mu_i \leq -\varepsilon$ for some $\varepsilon > 0$ (by the assumption of feasibility of problem (9)), there exists $\varepsilon > 0$ such that:

$$\Delta(\Theta(t)) \leq B - \varepsilon \cdot \sum_i (Q_i(t) + Y_i(t))$$

Telescoping over T slots and taking expectations:

$$E[L(\Theta(T))] - E[L(\Theta(0))] \leq BT - \varepsilon \cdot \sum_{t=0}^{T-1} \sum_i E[Q_i(t) + Y_i(t)]$$

Since $L(\Theta(T)) \geq 0$, rearranging and dividing by εT :

$$(1/T) \cdot \sum_{t=0}^{T-1} \sum_i E[Q_i(t)] \leq B/\varepsilon + E[L(\Theta(0))]/(\varepsilon T)$$

Taking $T \rightarrow \infty$: $\lim_{T \rightarrow \infty} (1/T) \cdot \sum_t E[Q_i(t)] \leq B/\varepsilon < \infty$,

which establishes mean-rate stability: $\lim_{T \rightarrow \infty} E[Q_i(T)]/T = 0$. $\square \square$

B. Proof of Theorem 2 (Cost Bound)

Proof:

Let π^* denote any stationary feasible policy achieving Cost_opt . Since LG-DDPG minimizes the DPP right-hand side at each slot, it satisfies:

$$\Delta(\Theta(t)) + V \cdot E[\text{Cost}(t)|\Theta(t)] \leq B + V \cdot \text{Cost_opt}$$

(The right-hand side is obtained by substituting π^* into the DPP bound, using the feasibility of π^* , which implies $\sum_i E[A_i - \mu_i | \pi^*] \leq 0$ and $\sum_i E[E_i - E_budget, i | \pi^*] \leq 0$, so the queue-weighted terms are non-positive.)

Telescoping over T slots and using $E[L(\Theta(T))] \geq 0$:

$$V \cdot (1/T) \cdot \sum_{t=0}^{T-1} E[\text{Cost}(t)] \leq BT + V \cdot T \cdot \text{Cost_opt} + E[L(\Theta(0))]$$

Dividing by VT :

$$(1/T) \cdot \sum_t E[\text{Cost}(t)] \leq \text{Cost_opt} + B/V + E[L(\Theta(0))]/(VT)$$

Taking $T \rightarrow \infty$:

$$\lim_{T \rightarrow \infty} (1/T) \cdot \sum_t E[\text{Cost}(t)] \leq \text{Cost_opt} + B/V. \square \square$$

C. Proof of Theorem 3 (DDPG Convergence)

Proof:

We apply the two-timescale stochastic approximation (SA) framework of Borkar [A1]. The critic update (Eq. 16) and actor update (Eq. 15) constitute a two-timescale system: the critic operates on the faster timescale (learning rate α_n^c) and the actor on the slower timescale (learning rate α_n^a).

Condition (i) — Learning rate schedules: We require $\sum_n \alpha_n = \infty$ and $\sum_n (\alpha_n)^2 < \infty$ for both actor and critic. The rates $\alpha_n^a = 1 \times 10^{-4}$ and $\alpha_n^c = 3 \times 10^{-4}$ satisfy these conditions when decayed appropriately, and $\alpha_n^a / \alpha_n^c \rightarrow 0$ ensures timescale separation.

Condition (ii) — Bounded iterates: The experience replay buffer decorrelates consecutive samples, approximating i.i.d. draws. The target networks provide stable regression targets, preventing divergence. Batch normalization bounds intermediate activations, ensuring bounded weight updates.

Condition (iii) — Lipschitz continuity: The ReLU/Tanh networks with fixed depth and bounded weights are Lipschitz continuous, satisfying the smoothness assumption of the SA theory.

Under conditions (i)–(iii), by Theorem 2 of Borkar (1997) [A1], the critic parameters ϕ converge to a local minimum of the mean-squared Bellman error, and the actor parameters θ subsequently converge to a stationary point of the DPP objective $J(\theta) = E[\sum_t \gamma^t r(t)]$ with probability 1.

[A1] V. S. Borkar, 'Stochastic approximation with two time scales,' *Systems Control Lett.*, vol. 29, no. 5, pp. 291–294, 1997. doi: 10.1016/S0167-6911(97)90015-3

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of **JIBAS** and/or the editor(s). **JIBAS** and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.