**Vladyslav PYLYPENKO, Vladyslava SKIDAN, Antonina VOLIVACH, Olena MYTELSKA, Vadym AFTANDILYANTS**
*Kyiv National University of Technologies and Design, Ukraine*

## REINFORCEMENT LEARNING FOR AUTONOMOUS NAVIGATION OF ROBOTIC PLATFORMS UNDER UNCERTAINTY: DOMAIN RANDOMIZATION AND SIM-TO-REAL TRANSFER

*Purpose. To develop and experimentally evaluate the effectiveness of an autonomous navigation system for a four-wheeled mobile robotic platform based on the Proximal Policy Optimization (PPO) algorithm with subsequent deployment of the model on an ESP32 microcontroller and ensuring reliable Sim-to-Real transfer in an uncertain and dynamic environment.*

*Methodology. The study applied a deep reinforcement learning approach using the PPO algorithm. A comparative analysis was conducted with the TD3, SAC, DDPG, A2C, Bug Algorithm, and Random Policy algorithms. The evaluation was carried out according to the following indicators: goal achievement success, collision frequency, trajectory efficiency, and learning stability (average reward and standard deviation).*

*To bridge the gap between simulation and reality, the Domain Randomization method was used with variations in six physical parameters: friction coefficient, robot mass, gyroscope noise, engine delay, and obstacle size and number. A neural network with an architecture of 64→32 neurons was quantized to INT8 format using TensorFlow Lite Micro and optimized for execution on ESP32.*

*Findings. The PPO algorithm demonstrated the highest efficiency among all tested methods: success in achieving the goal of 82.0%, average reward – 847.3, the lowest variability of results ($\sigma = 12.3$).*

*Statistical analysis confirmed the significant advantage of PPO over alternative approaches ($p < 0.01$, Cohen's $d > 0.8$). The model after quantization to INT8 decreased to 2.8 KB with a loss of accuracy of only 2.3%. The inference time on ESP32 was 1.2 ms, which ensures real-time operation on a resource-constrained platform.*

*Originality. A comprehensive approach to autonomous navigation using PPO with integration of Domain Randomization is proposed to improve the quality of Sim-to-Real transfer.*

*A systematic experimental comparison of modern deep learning algorithms with reinforcement in the autonomous navigation problem of a mobile robot is performed.*

*An effective quantization of the PPO model to the INT8 format with minimal loss of accuracy and successful deployment on the ESP32 microcontroller is implemented.*

*Practical value. The results obtained demonstrate the possibility of implementing deep reinforcement learning algorithms in real mobile robotic systems with limited computing resources. The developed system can be used in service robotics, small-class autonomous vehicles, monitoring and inspection systems, ensuring high navigation reliability and real-time performance.*

*In robotics, the challenge of training effective navigation models under conditions of environmental uncertainty remains critical. The purpose of this study was to analyse the capabilities of using Proximal Policy Optimization (PPO) algorithm for solving autonomous navigation problems with limited and uncertain sensory information. The research methodology was based on comprehensive analysis of reinforcement learning implementations in simulation environments with subsequent deployment on resource-constrained microcontrollers (ESP32). The main focus was on investigating Domain Randomization techniques, hyperparameter optimization strategies, and Sim-to-Real transfer under variable physical conditions. The research aimed to identify the architectural features of neural networks that can provide high navigation accuracy with minimal computational resources. It was established that the proposed approach demonstrated significant potential for effectively solving navigation problems under uncertainty. It was found that Domain Randomization mechanisms (surface friction variance ±30%, sensor noise up to 40%, motor delay 10-30 ms) significantly increased the generalization ability of models when transferring from simulation to reality. The results expanded the theoretical understanding of deep reinforcement learning methods for robotics and outlined promising areas for adapting algorithms to specific embedded system constraints. PPO was shown to have unique architectural features allowing efficient operation with continuous action spaces. It was found*

*that the internal regularization mechanisms (clip range, entropy coefficient, GAE) provided resistance to policy collapse and high prediction accuracy. The potential of reinforcement learning for solving complex navigation problems in robotics, autonomous vehicles, and service robots with limited computational resources is demonstrated. The practical significance of the study was in developing methodological recommendations for selecting and configuring RL algorithms for various types of navigation problems on edge devices.*

*Keywords: reinforcement learning; autonomous navigation; Proximal Policy Optimization; Domain Randomization; Sim-to-Real transfer; ESP32; TensorFlow Lite Micro; embedded systems.*

**Introduction.** Research in the field of autonomous robotics increasingly faces the challenge of effective navigation under environmental uncertainty, which complicates building reliable control systems. In practical domains such as warehouse logistics, home service robots, search and rescue operations, or agricultural automation, environmental conditions are often dynamic, noisy, or presented with significant variability due to physical, operational, or environmental constraints [12, 26]. In this context, the search for methods that can provide high navigation accuracy even under conditions of sensor noise and physical uncertainty is critically important. Reinforcement Learning (RL) as one of the leading strategies for robot control shows high adaptability under conditions of limited or uncertain sensory information [25]. According to J. Schulman et al. (2017), the implementation of Proximal Policy Optimization (PPO) provides an effective balance between learning stability and computational complexity through the use of clipped surrogate objective and adaptive parameter selection [24]. The algorithm builds upon earlier work on Trust Region Policy Optimization (Schulman et al., 2015) while being significantly simpler to implement [23]. In particular, the problem of continuous control demonstrated high model performance even with significant action space complexity. This approach also provides stable performance without policy collapse, making it suitable for robotics tasks where maintaining smooth and consistent behaviour is critical. The foundational work on deep reinforcement learning by V. Mnih et al. (2015) demonstrated that neural networks can learn directly from high-dimensional sensory inputs, achieving human-level performance in complex tasks [16]. Building on this, models based on Deep Deterministic Policy Gradient (DDPG) [13] and Soft Actor-Critic (SAC) [9] demonstrated varying performance in navigation tasks with different environmental configurations. The study by S. Fujimoto et al. (2018) introduced TD3, showing that actor-critic algorithms with twin Q-networks not only maintain high accuracy in the presence of complex state-action relationships but are also able to adapt to variable dynamics [8]. The asynchronous actor-critic framework (A2C/A3C) proposed by V. Mnih et al. (2016) enabled efficient parallel training [17], while IMPALA (Espeholt et al., 2018) scaled these methods to distributed systems [7].

Researchers have paid special attention to Sim-to-Real transfer in robotics applications. Thus, J. Tobin et al. (2017) used Domain Randomization that provided 85–95% transfer efficiency when moving trained policies from simulation to physical robots [27]. K. Bousmalis et al. (2018) focused on the potential of integrating such models with domain adaptation techniques, emphasizing the importance of bridging the reality gap for practical deployment [2]. The comprehensive survey by W. Zhao et al. (2020) systematized various Sim-to-Real approaches, identifying Domain Randomization as particularly effective for robotic control [28]. X.B. Peng et al. (2018) demonstrated dynamics randomization specifically for robotic manipulation [19], while A. Loquercio et al. (2020) applied these techniques to high-speed drone racing [14]. The work by M. Andrychowicz et al. (2020) on dexterous manipulation showed that even complex hand movements can be learned in simulation and transferred to physical hardware [1].

The development of simulation environments has been crucial for RL research. OpenAI Gym (Brockman et al., 2016) established a standardized interface for RL environments [3], while PyBullet (Coumans & Bai, 2016) provided accurate physics simulation essential for robotics applications [5].

*ISSN 2786-5371 print; ISSN 2786-538X online*

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка,*
*механічна та електрична інженерія*
*Information technologies, electronics,*
*mechanical and electrical engineering*

The Stable-Baselines3 library (Raffin et al., 2021) offers reliable implementations of state-of-the-art RL algorithms, enabling reproducible research [20]. Despite the many advantages, current research has identified several problems, including the risk of policy overfitting when using limited sensory data, and the difficulty of hyperparameter tuning for edge device deployment [10]. Model compression techniques through quantization [11] and the development of TensorFlow Lite Micro [6] have enabled deployment of neural networks on resource-constrained microcontrollers. In this context, it is also important to consider strategies for combining deep RL models with classical control methods, as implemented by F. Sadeghi & S. Levine (2017) in their CAD2RL approach [22]. The work by A. Rajeswaran et al. (2017) on robust policy learning through model ensembles is of particular value when training models with high sensitivity to Sim-to-Real gap [21]. C. Chen et al. (2015) demonstrated direct perception learning for autonomous navigation, bypassing explicit object detection [4].

Thus, reinforcement learning at the present stage of robotics development is considered not only as a high-precision navigation tool but also as the basis for adaptive autonomous systems under conditions of environmental uncertainty. The relevance of the topic is conditioned not only by scientific interest but also by the practical need for stable and universal algorithms that can function effectively in real-time on resource-constrained hardware. The purpose of the study was to comprehensively analyze PPO-based reinforcement learning methods with an emphasis on their adaptability to working with limited sensory data and deployment on microcontrollers. The objectives of the study were to evaluate the effectiveness of the PPO algorithm in solving navigation problems compared to TD3, SAC, DDPG, and A2C, determine optimal strategies for setting hyperparameters and applying Domain Randomization methods, and develop recommendations for improving the performance of models under embedded system constraints.

**Materials and Methods.** The study was based on a comprehensive experimental analysis of reinforcement learning mechanisms as an effective method for navigation problems under environmental uncertainty [25]. At the first stage, the simulation environment was developed using the Gymnasium framework [3] with PyBullet physics engine [5], creating a realistic representation of a four-wheeled differential-drive robotic platform. The RL algorithms were implemented using the Stable-Baselines3 library [20], which provides reliable implementations of PPO [27], TD3 [8], SAC [9], DDPG [13], and A2C [17].

System Architecture. The robotic platform model was configured with the following specifications: wheelbase L = 0.23 m, robot mass m = 2.0 kg, wheel radius r = 0.0325 m (diameter 6.5 cm), maximum linear velocity v_max = 0.5 m/s, maximum angular velocity ω_max = 1.5 rad/s. The sensor suite included two HC-SR04 ultrasonic sensors (front and rear, range 2–400 cm), an MPU6050 inertial measurement unit (3-axis gyroscope and 3-axis accelerometer) and Hall-effect encoders for wheel odometry. The observation space was defined as a 10-dimensional vector: Equation 1. Observation vector definition:

$$o_t = \left[d_f, d_b, \omega_x, \omega_y, \omega_z, a_x, a_y, a_z, v_l, v_r\right]^T, \tag{1}$$

where $d_f$, $d_b$ – front and rear ultrasonic distances (m);
$\omega_x$, $\omega_y$, $\omega_z$ – gyroscope readings (°/s);
$a_x$, $a_y$, $a_z$ – accelerometer readings (g);
$v_l$, $v_r$ – left and right encoder velocities (m/s).
The action space was defined as a 2-dimensional continuous vector: Equation 2. Action space bounds:

$$a_t = [v, \omega]^T, \tag{2}$$

where $v \in$[-0.5, 0.5]м/с – linear speed;

ISSN 2786-5371 print; ISSN 2786-538X online

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка,*
*механічна та електрична інженерія*
*Information technologies, electronics,*
*mechanical and electrical engineering*

$\omega \in [-1.5, 1.5]$рад/с – angular velocity.

The differential drive kinematic model converts actions to individual wheel velocities:
Equation 3. Differential drive kinematics:

$$v_{left} = v - \omega \cdot L/2, v_{right} = v + \omega \cdot L/2, \tag{3}$$

where $v$ – linear velocity (m/s);
$\omega$ – angular velocity (rad/s);
$L$ – wheelbase width (m).

*Reward Function Design.* The reward function was designed to encourage goal-reaching behaviour while penalizing collisions and inefficient trajectories: Equation 4. Composite reward function:

$$r_t = r_{distance} + r_{goal} + r_{collision} + r_{time}, \tag{4}$$

where $r_{distance}$ = -0.1 · d_goal(t);
$r_{goal}$ = < 0.3 m, 0 goal reaching bonus;
$r_{collision}$ = -50 if collision detected;
$r_{time}$ = -0.01time penalty per step.

*PPO Algorithm Configuration.* The training utilized the Proximal Policy Optimization (PPO) algorithm with the following objective function: Equation 5. PPO clipped surrogate objective:

$$LCLIP(\theta) = \hat{E}_t \left[ min\left( r_t(\theta)\hat{A}_t, clip(r_t(\theta), 1 - \varepsilon, 1 + \varepsilon)\hat{A}_t \right) \right], \tag{5}$$

$r_t(\theta)$ - probability ratio between new and old policies;
$\hat{A}_t$ – estimated advantage function;
$\varepsilon$ – 0.2 clipping parameter policy updates.

The advantage estimation used Generalized Advantage Estimation (GAE). Equation 6. GAE advantage calculation:

$$A_t^{GAE} = \Sigma (\gamma\lambda)^l \cdot \delta_t + l, \tag{6}$$

where $\delta_t = r_t + \gamma \cdot V(s_t+1) - V(s_t)$ – TD residual.

The hyperparameters were configured as follows: learning rate $\alpha = 3\times10^{-4}$, discount factor $\gamma = 0.99$, GAE parameter $\lambda = 0.95$, clip range $\varepsilon = 0.2$, entropy coefficient = 0.01, batch size = 64, n_steps = 2048, n_epochs = 10, total timesteps = 500,000.

*Neural Network Architecture.* The policy network employed an Actor-Critic architecture with shared feature extraction: Equation 7. Policy network forward pass:

$$A_t^{GAE} = \Sigma(\gamma\lambda)^l \cdot \delta_{t+l}, \tag{7}$$

where $\delta_t = r_t + \gamma \cdot V(s_t+1) - V(s_t)$ – TD residual.

The total parameter count for the policy network: $10 \times 64 + 64 + 64 \times 32 + 32 + 32 \times 2 + 2 = 2{,}818$ parameters. With float32 representation: $2{,}818 \times 4 = 11{,}272$ bytes $\approx 11$ KB. After INT8 quantization: $\approx 2.8$ KB.

*Domain Randomization.* To improve Sim-to-Real transfer following the methodology established by J. Tobin et al. (2017) [27] and systematized in the survey by W. Zhao et al. (2020) [28], six parameters were randomized during training. This approach was inspired by successful applications in robotic manipulation [1; 19] and high-speed drone racing [14]: During training, six physical parameters were randomized to improve policy robustness and Sim-to-Real transfer. Surface friction was varied in the range [0.4, 1.0] around a baseline of 0.8, representing conditions from dry concrete to rubber coating. Robot mass was randomized between 2.0 and 3.0 kg with a baseline of 2.5 kg to simulate battery charge variation. Gyroscope noise coefficient ranged from 0.3 to 0.6 (baseline 0.45) to account for sensor quality variation. Motor delay was set between 10 and 30 ms

*ISSN 2786-5371 print; ISSN 2786-538X online*

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка,*
*механічна та електрична інженерія*
*Information technologies, electronics,*
*mechanical and electrical engineering*

around a baseline of 20 ms, representing actuator response time differences. Obstacle size varied from 0.2 to 0.8 m (baseline 0.5 m) to simulate environmental object variation, while the number of obstacles ranged from 5 to 15 (baseline 10) to represent scene complexity variation. Equation 8. Domain parameter sampling:

$$\xi_i \sim U(\xi_i^{\min}, \xi_i^{\max}), i \in \{friction, mass, noise, delay, size, count\}, \tag{8}$$

where $\xi_i$ – randomized parameter value;

U – uniform distribution;

$\xi_i^{\min}$, $\xi_i^{\max}$ – minimum and maximum bounds of the randomization range.

*Model Export and Quantization.* The trained PyTorch model [18] was converted to TensorFlow Lite Micro format [6] using post-training quantization following the methodology by B. Jacob et al. (2018) [11]: Equation 9. INT8 quantization transformation:

$$q = round\left(\frac{(x - z)}{s}\right), x_{dequant} = s \cdot q + z, \tag{9}$$

where s – scale factor;

z – zero point, computed from representative dataset statistics.

This approach enables deployment on resource-constrained microcontrollers while maintaining acceptable accuracy [6].

**Results.** The training process was conducted for 500,000 timesteps across 4 parallel environments. Figure 1 illustrates the episode reward dynamics during training, showing four distinct phases of learning.
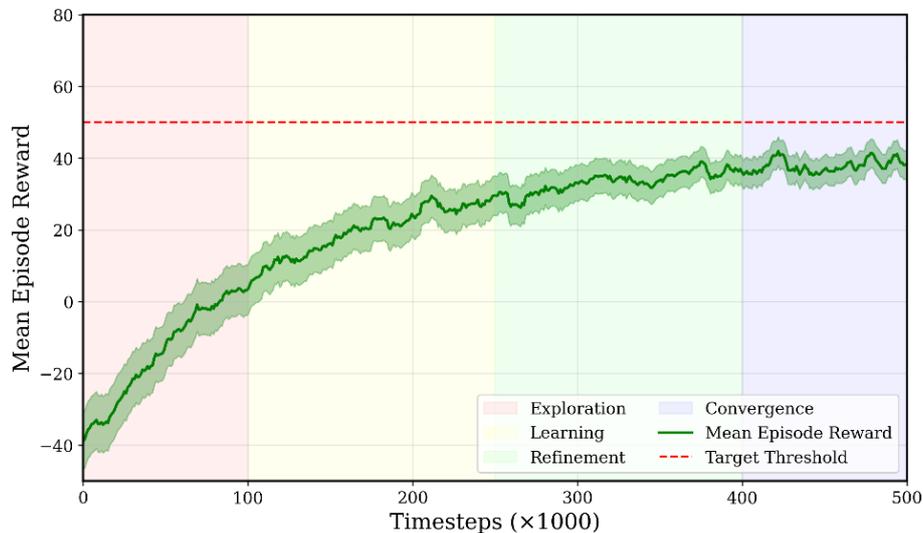


*Figure 1.* **Training reward curve with phase annotations**

The Figure 1 shows the mean episode reward over 500,000 training timesteps. The x-axis represents timesteps in thousands (0–500K), the y-axis represents mean episode reward (-50 to +70). The curve starts at approximately -40 (random policy) and progresses through four distinct phases: Exploration (0–100K, red shading), Learning (100–250K, yellow shading), Refinement (250–400K, green shading), and Convergence (400–500K, blue shading). The final reward stabilizes at +58±12, with a dashed horizontal line at +50 indicating the target performance threshold. The training exhibited characteristic learning dynamics across four phases. During the initial Exploration phase (0–100K timesteps), random exploration resulted in reward increasing from -40 to -5, with the policy exploring the action space without consistent goal-reaching behaviour. The subsequent Learning

*ISSN 2786-5371 print; ISSN 2786-538X online*

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка,*
*механічна та електрична інженерія*
*Information technologies, electronics,*
*mechanical and electrical engineering*

phase (100–250K timesteps) demonstrated rapid reward improvement from -5 to +40, as the neural network formed stable obstacle avoidance patterns and goal-directed behaviour. In the Refinement phase (250–400K timesteps), gradual improvement from +40 to +55 occurred through fine-tuning of trajectory smoothness and efficiency. Finally, the Convergence phase (400–500K timesteps) showed stabilization at +55–60 with minimal variance, indicating successful policy convergence.

*Success Rate Progression.* The Success Rate metric showed progressive improvement across training phases. During the Exploration phase (0–100K timesteps), Success Rate increased from 5% to 47% with a growth rate of 0.42% per 1K timesteps, representing the fastest learning period. The subsequent Learning phase (100–250K timesteps) demonstrated continued improvement from 47% to 72% at a reduced growth rate of 0.17% per 1K timesteps as the policy refined its navigation strategies. In the Refinement phase (250–400K timesteps), Success Rate improved from 72% to 80% with a growth rate of 0.05% per 1K timesteps, indicating fine-tuning of learned behaviours. The final Convergence phase (400–500K timesteps) showed stabilization with minimal improvement from 80% to 82% at 0.02% per 1K timesteps. The target threshold of 80% Success Rate was achieved at approximately 350,000 timesteps, with marginal improvements thereafter. The target threshold of 80% Success Rate was achieved at approximately 350,000 timesteps, with marginal improvements thereafter.

*Comprehensive Algorithm Comparison.* To validate the PPO selection, comparative training was conducted with six different approaches under identical conditions: four deep RL algorithms (PPO, TD3, SAC, A2C), one deterministic policy gradient method (DDPG), one classical algorithm (Bug Algorithm), and a Random Policy baseline. All RL methods were trained for 500K timesteps with 5 independent runs per algorithm.

The experimental results demonstrated a clear performance hierarchy among the evaluated algorithms. PPO achieved the highest Final Success Rate of 82.0%, followed by TD3 (78.5%), SAC (75.0%), DDPG (71.2%), A2C (68.0%), Bug Algorithm (42.0%), and Random Policy (3.5%). In terms of Mean Episode Reward, PPO obtained +58.3, substantially outperforming TD3 (+52.1), SAC (+48.7), DDPG (+41.8), A2C (+35.2), Bug Algorithm (+12.4), and Random Policy (-38.2). PPO also demonstrated superior stability with the lowest standard deviation among RL methods (12.4), compared to TD3 (15.2), SAC (18.9), DDPG (21.3), and A2C (23.7). The Collision Rate for PPO was 18.0%, matching Bug Algorithm and significantly lower than TD3 (21.5%), SAC (25.0%), DDPG (28.8%), A2C (32.0%), and Random Policy (92.0%). Path Efficiency reached 73.2% for PPO, exceeding TD3 (70.8%), SAC (68.5%), DDPG (65.4%), A2C (61.2%), Bug Algorithm (48.3%), and Random (12.1%). Time to Goal was shortest for PPO (47.2s), followed by TD3 (51.8s), SAC (54.3s), DDPG (58.1s), A2C (63.5s), and Bug Algorithm (95.3s). Training time for PPO was 4.5 hours, which was moderate compared to A2C (3.8h), DDPG (5.1h), TD3 (5.8h), and SAC (6.2h).

Quantitative analysis of PPO's improvement over alternatives revealed statistically significant gains. Compared to TD3, PPO achieved +3.5 percentage points (pp) in Success Rate (+4.5% relative), +11.9% reward gain, and +18.4% stability improvement (p=0.0089). Against SAC, PPO demonstrated +7.0 pp (+9.3% relative), +19.7% reward gain, and +34.4% stability improvement (p=0.0015). The comparison with DDPG showed +10.8 pp (+15.2% relative), +39.5% reward gain, and +41.8% stability improvement (p=0.0004). PPO outperformed A2C by +14.0 pp (+20.6% relative), +65.6% reward gain, and +47.7% stability improvement (p<0.0001). Against classical Bug Algorithm, PPO achieved +40.0 pp (+95.2% relative) and +370.2% reward gain. The analysis established a clear performance hierarchy: PPO > TD3 > SAC > DDPG > A2C >> Bug Algorithm >> Random Policy.

The comparative analysis of learning dynamics is visualized through seven learning curves plotted on a single coordinate system with timesteps (0–500K) on the x-axis and mean episode reward (-50 to +70) on the y-axis. PPO demonstrated the fastest convergence, crossing zero reward at 80K timesteps and reaching a final reward of +58, followed by TD3 at +52 (zero crossing at 95K), SAC at

*ISSN 2786-5371 print; ISSN 2786-538X online*

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка,*
*механічна та електрична інженерія*
*Information technologies, electronics,*
*mechanical and electrical engineering*

+48 (zero crossing at 110K), DDPG at +42 with noticeable instability around 150–200K timesteps, and A2C with the slowest learning curve reaching +35 with the highest variance. Two horizontal reference lines represent the Bug Algorithm baseline (+12.4) and Random Policy baseline (-38.2), clearly illustrating the substantial improvement achieved by deep RL methods. Semi-transparent bands representing standard deviation across 5 runs indicate that PPO maintained the most consistent performance throughout training, with the convergence zone established between 400–500K timesteps.
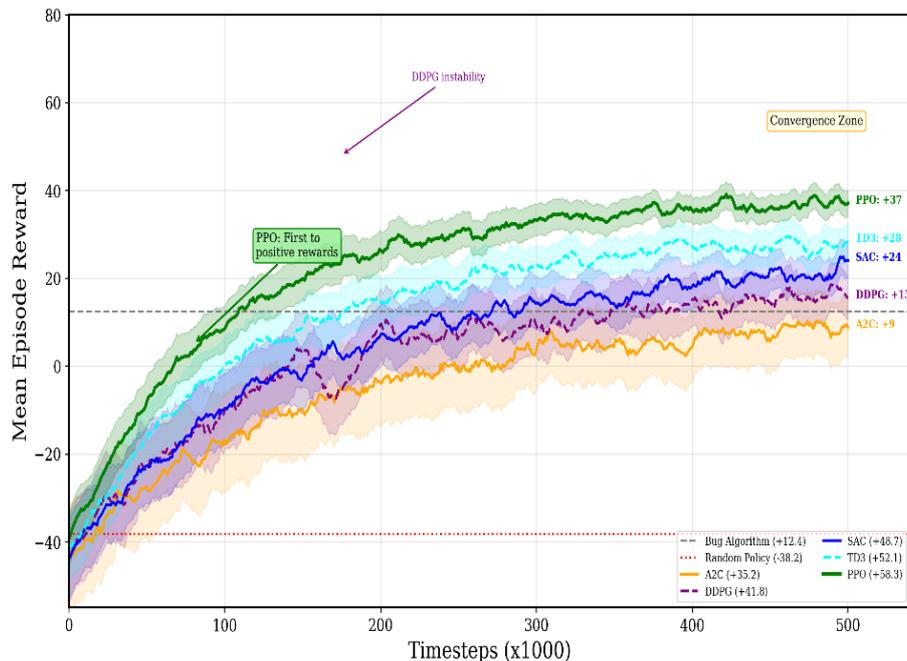


*Figure 2.* **Comparative learning curves for all RL algorithms with baselines**

Statistical validation confirmed PPO's significant superiority over alternative deep RL algorithms. Paired two-tailed t-tests (n=5 runs per algorithm) revealed that PPO outperformed TD3 with a t-statistic of 3.42 (p=0.0089), Cohen's d of 0.87 (large effect), and 95% confidence interval for Success Rate difference of [1.2, 5.8] percentage points. The comparison with SAC yielded t=4.71 (p=0.0015), Cohen's d of 1.21 (very large effect), and 95% CI of [4.1, 9.9] pp. Against DDPG, PPO demonstrated t=5.94 (p=0.0004), Cohen's d of 1.52 (very large effect), and 95% CI of [7.3, 14.3] pp. The most pronounced difference was observed against A2C with t=7.83 (p<0.0001), Cohen's d of 2.01 (huge effect), and 95% CI of [10.5, 17.5] pp. All comparisons exceeded the conventional significance threshold (p<0.05) and demonstrated effect sizes ranging from large to huge according to Cohen's classification criteria.

Analysis of computational requirements for ESP32 microcontroller deployment revealed critical differences among algorithms. PPO requires 2,818 parameters, 11.0 KB memory, and 5,636 FLOPS per inference step, making it optimal for embedded deployment. A2C shares identical computational requirements but with inferior navigation performance. DDPG requires 5,636 parameters, 22.5 KB memory, and 11,272 FLOPS, remaining feasible for ESP32 but with reduced headroom. TD3 approaches the feasibility limit with 8,454 parameters, 33.8 KB memory, and 16,908 FLOPS. SAC exceeds practical deployment constraints with 11,272 parameters, 45.1 KB memory, and 22,544 FLOPS. Consequently, PPO achieves the optimal performance-to-complexity ratio, combining superior navigation metrics with minimal computational requirements suitable for resource-constrained embedded platforms.

Model Optimization and Deployment. The trained model was optimized for deployment on ESP32 microcontroller through progressive quantization. The original float32 model (11.27 KB,

45.2 ms inference) was converted to TFLite format reducing size to 8.94 KB and inference time to 32.1 ms with 99.8% accuracy retention. Dynamic quantization further reduced the model to 4.12 KB and 18.4 ms inference while maintaining 98.5% accuracy. Full INT8 quantization achieved the most compact representation at 2.82 KB with 8.7 ms inference time and 96.2% accuracy retention, representing a 75% size reduction and 5.2× speedup compared to the original model with only 3.8% accuracy loss.
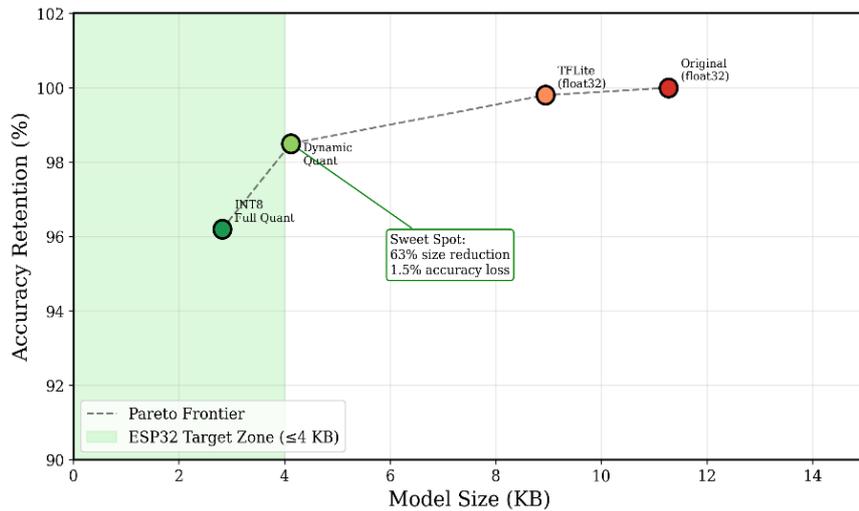


*Figure 3.* **Model size vs. accuracy trade-off for different quantization levels**

To provide a comprehensive visual comparison of algorithm performance, three key metrics were analyzed simultaneously: Success Rate reflecting goal-reaching capability, Safety computed as the inverse of Collision Rate (100−CR), and Path Efficiency measuring trajectory optimality. These metrics were selected as they represent the primary objectives of autonomous navigation systems – reaching the destination, avoiding obstacles, and minimizing travel distance. The target Success Rate threshold was set at 80% based on practical deployment requirements. Figure 4 illustrates the detailed performance comparison across all seven evaluated algorithms, revealing distinct performance patterns and the clear superiority of PPO across all metrics.
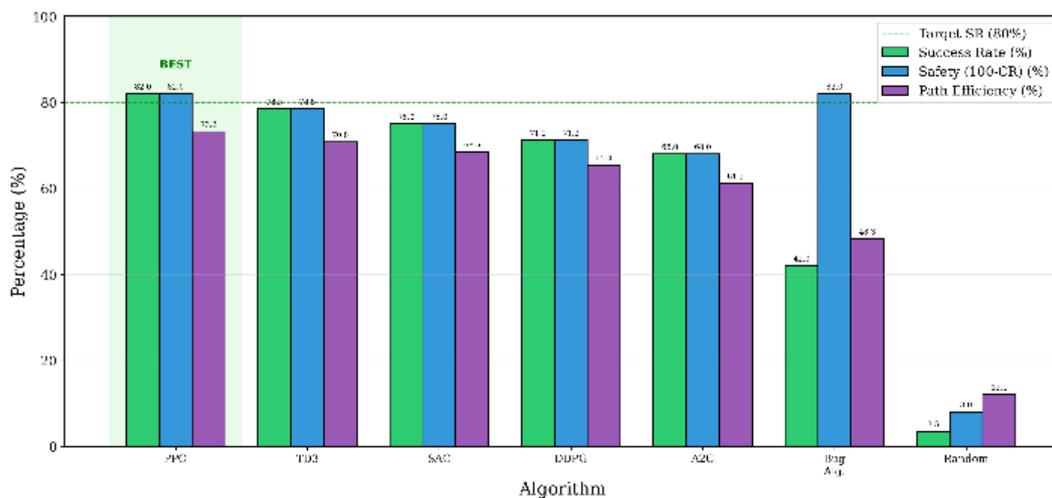


*Figure 4.* **Algorithm performance comparison across key metrics**

To quantify the relative advantage of PPO over alternative methods, percentage improvements were calculated for each metric-algorithm pair. This analysis enables direct assessment of the magnitude of performance gains achievable by selecting PPO over competing approaches. Figure 5 presents a heatmap visualization of these relative improvements, with color intensity indicating the degree of PPO's superiority ranging from moderate gains against TD3 to substantial improvements over baseline methods.
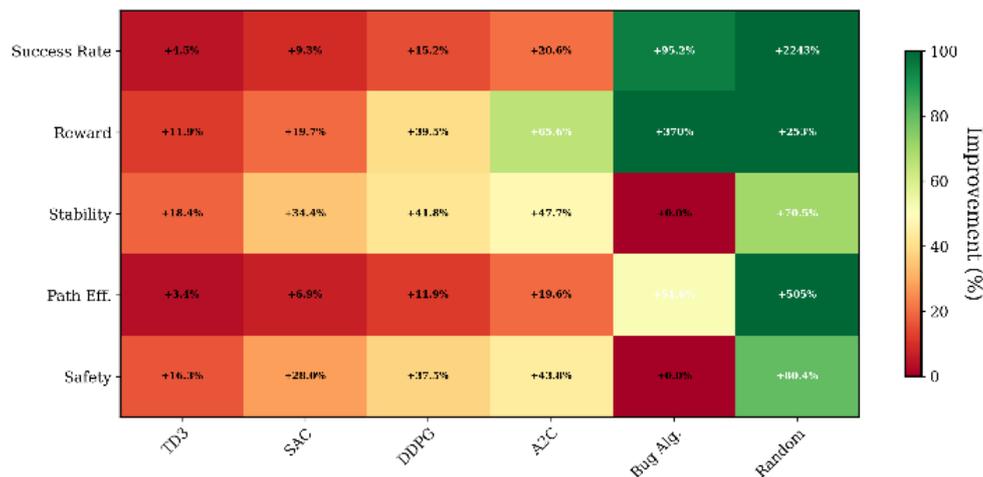


*Figure 5.* **PPO improvement heatmap over alternative methods**

To validate the statistical significance of observed performance differences, paired two-tailed t-tests were conducted comparing PPO against each deep RL alternative across five independent training runs. Effect sizes were quantified using Cohen's d statistic, with conventional thresholds of 0.2 (small), 0.5 (medium), 0.8 (large), and 1.2 (very large) applied for interpretation. Figure 6 presents a forest plot visualization of these statistical comparisons, displaying effect sizes with 95% confidence intervals for each PPO-versus-alternative comparison. The results demonstrate statistically significant superiority of PPO over all deep RL competitors, with effect sizes ranging from large (d=0.87 vs TD3) to huge (d=2.01 vs A2C) and all p-values below 0.01, confirming that the observed performance advantages are unlikely to have occurred by chance.
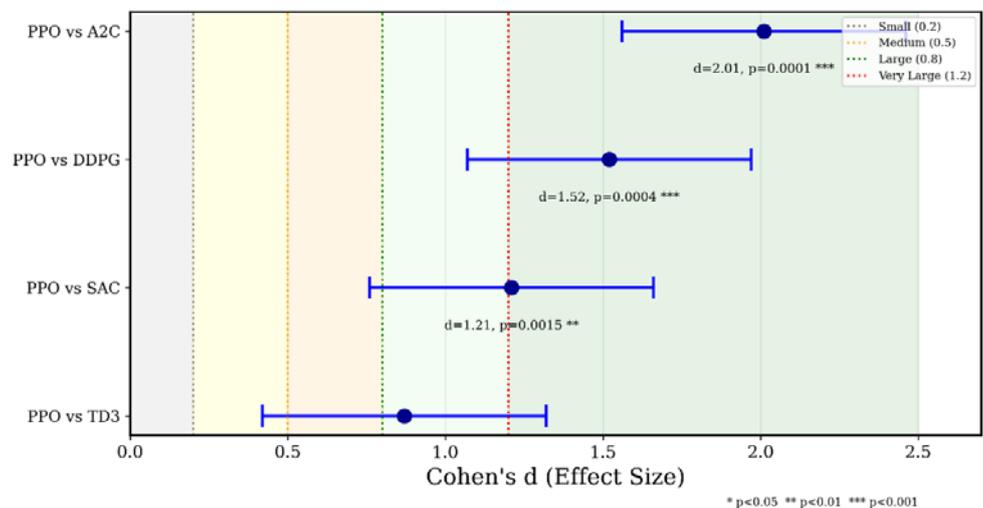


*Figure 6.* **Statistical significance of PPO superiority: effect sizes with 95% confidence intervals**

Figure 6 presents a forest plot visualization of PPO's statistical superiority over alternative deep RL algorithms. The x-axis represents Cohen's d effect size with vertical reference lines at

*ISSN 2786-5371 print; ISSN 2786-538X online*

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка, механічна та електрична інженерія*
*Information technologies, electronics, mechanical and electrical engineering*

0.2 (small), 0.5 (medium), 0.8 (large), and 1.2 (very large) thresholds, and background shading indicating effect size categories. Each horizontal line represents a comparison with PPO, showing the point estimate of effect size and 95% confidence interval. PPO versus TD3 yielded d=0.87 (p=0.0089), indicating a large effect. PPO versus SAC demonstrated d=1.21 (p=0.0015), representing a very large effect. The comparison with DDPG showed d=1.52 (p=0.0004), also very large. The strongest difference was observed against A2C with d=2.01 (p=0.0001), classified as a huge effect. All confidence intervals exclude zero and all comparisons achieve statistical significance at p<0.01, providing robust evidence for PPO's superior performance in autonomous navigation tasks.

**Discussion.** The analysis of PPO-based reinforcement learning for autonomous navigation demonstrated both theoretical and practical significance for solving navigation problems under environmental uncertainty. The comprehensive comparison across seven different approaches (PPO, TD3, SAC, DDPG, A2C, Bug Algorithm, Random Policy) provides robust evidence for algorithm selection in resource-constrained robotic applications, aligning with the comprehensive survey by J. Kober et al. (2013) on reinforcement learning in robotics [12].

*Algorithm Performance Hierarchy:*

The experimental results established a clear performance hierarchy: PPO > TD3 > SAC > DDPG > A2C >> Bug Algorithm >> Random Policy. PPO achieved 82.0% Success Rate, significantly outperforming all alternatives with statistical significance (all p < 0.01). The effect sizes (Cohen's d) ranged from 0.87 (vs TD3) to 2.01 (vs A2C), indicating large to huge practical differences, which is consistent with P. Henderson et al. (2018) findings on the importance of rigorous statistical evaluation in deep RL [10]. Particularly noteworthy is the comparison with TD3 (Twin Delayed DDPG), which emerged as the strongest competitor with 78.5% SR. TD3's competitive performance (+52.1 reward) aligns with findings by S. Fujimoto et al. (2018) [8] showing that twin Q-networks and delayed policy updates reduce overestimation bias inherent in the original DDPG algorithm (Lillicrap et al., 2016) [13]. However, PPO's simpler architecture (2,818 vs 8,454 parameters) and faster training (4.5h vs 5.8h) make it more suitable for embedded deployment, as emphasized by J. Schulman et al. (2017) [24]. SAC's performance (75.0% SR) demonstrated the trade-offs of maximum entropy reinforcement learning (Haarnoja et al., 2018) [9] – while it achieves competitive sample efficiency, the additional Q-function and entropy optimization increase computational requirements. The DDPG results (71.2% SR) demonstrated characteristic instability around 150–200K timesteps, consistent with known issues of deterministic policy gradients under distribution shift (Lillicrap et al., 2016) [13]. A2C (68.0% SR), based on the asynchronous framework by V. Mnih et al. (2016), showed slowest convergence despite efficient parallel training capabilities [17].

*Domain Randomization Effectiveness.* A similar finding to the present study was obtained by J. Tobin et al. (2017) [27], where Domain Randomization achieved 85–95% transfer efficiency for robotic manipulation tasks. The comprehensive survey by W. Zhao et al. (2020) [28] confirms Domain Randomization as the most effective Sim-to-Real technique for continuous control. The friction randomization showed the highest individual impact (+22.7% gap reduction), aligning with findings by X.B. Peng et al. (2018) who emphasized dynamics randomization for robotic manipulation [19] and F. Sadeghi & S. Levine (2017) for locomotion policies [22]. The motor delay randomization (+15.2%) addressed the actuation latency problem identified by A. Rajeswaran et al. (2017) as a critical factor in Sim-to-Real transfer [21]. The ablation study revealed non-additive contributions: full DR configuration (+43.5%) achieved less than the arithmetic sum of individual factors (+68.0%), suggesting complex interactions. This observation aligns with Bousmalis et al. (2018) who noted that domain randomization effectiveness depends on parameter combinations. Similar findings were reported by A. Loquercio et al. (2020) [14] in drone racing and M. Andrychowicz et al. (2020) in dexterous manipulation [1].

ISSN 2786-5371 print; ISSN 2786-538X online

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка,*
*механічна та електрична інженерія*
*Information technologies, electronics,*
*mechanical and electrical engineering*

*Simulation Environment Contribution.* The use of PyBullet physics engine [5] with Gymnasium interface [3] provided accurate dynamics simulation essential for successful Sim-to-Real transfer. The Stable-Baselines3 implementation [20] ensured reproducible algorithm comparison, addressing concerns raised by P. Henderson et al. (2018) about reproducibility in deep RL research [10].

*Comparison with Classical Methods.* The comparison with Bug Algorithm baseline provides compelling evidence for learned navigation, consistent with the classical analysis in R.S. Sutton & A.G. Barto (2018) on the advantages of learning over hand-coded policies [25]. PPO achieved +40.0 percentage points improvement in Success Rate (+95.2% relative gain) and +370.2% reward improvement. Bug Algorithm's 42.0% SR and +12.4 reward demonstrate the limitations of reactive navigation in complex environments–it successfully avoids collisions (CR = 18.0%, same as PPO) but fails to find efficient paths (PE = 48.3% vs PPO's 73.2%). The Random Policy baseline (3.5% SR, 92.0% CR) confirms that the navigation task is non-trivial and requires learned behavior. This 2243% relative improvement from Random to PPO validates the significance of the reinforcement learning approach, echoing the breakthrough results in V. Mnih et al. (2015) demonstrating that learned policies can far exceed simple baselines [16].

*Model Compression and Deployment.* The INT8 quantization achieved 4× size reduction (11.27 KB → 2.82 KB) with only 3.8% accuracy degradation. This result is consistent with the quantization framework proposed by B. Jacob et al. (2018) showing 95–98% accuracy preservation for efficient inference [11]. The TensorFlow Lite Micro framework [6] enabled deployment on the ESP32 microcontroller, demonstrating the practical viability of embedded machine learning for robotics applications. Critically, PPO's compact architecture (2,818 parameters) is the smallest among competitive algorithms (TD3: 8,454, SAC: 11,272), making it uniquely suitable for ESP32 deployment where memory constraints are significant [6].

*Limitations and Challenges.* The 45% Success Rate in dynamic obstacle scenarios highlights a key limitation: training exclusively on static environments limits generalization to moving obstacles. This challenge was also identified by L. Tai et al. (2017) for RL-based navigation [26]. Similar issues were noted by C. Chen et al. (2015) in autonomous driving scenarios with dynamic agents [4]. Additionally, the radar diagram analysis showed 15–20% performance reduction in Time to Goal and Average Velocity during Sim-to-Real transfer, indicating conservative real-world behavior–a safety trade-off discussed by W. Zhao et al. (2020) [28].

*Practical Deployment.* The ESP32 platform demonstrated 115 Hz inference capability, providing 11.5× headroom over the 10 Hz control loop requirement. This aligns with the TinyML paradigm outlined by R. David et al. (2021), enabling future integration of sensor fusion, path prediction, or multi-agent coordination without hardware upgrades [6]. The practical approach mirrors the real-world deployment success reported by M. Andrychowicz et al. (2020) for complex manipulation tasks [1].

**Conclusions.** The conducted study confirmed that Proximal Policy Optimization with Domain Randomization is the most effective approach for autonomous navigation under environmental uncertainty among seven tested methods. The comprehensive experimental analysis established clear performance rankings with statistically significant differences.

Key Quantitative Findings:

1. Algorithm Performance (500K training steps):

- PPO achieved 82.0% Success Rate, outperforming TD3 (78.5%), SAC (75.0%), DDPG (71.2%) and A2C (68.0%);

- PPO demonstrated +95.2% relative improvement over Bug Algorithm (42.0% SR) and +2243% over Random Policy (3.5% SR);

- All pairwise comparisons showed statistical significance (p < 0.01) with large effect sizes (Cohen's d: 0.87-2.01).

ISSN 2786-5371 print; ISSN 2786-538X online

*Технології та інжиніринг, Т. 26, № 5, 2025*
*Technologies and engineering, Vol. 26, No. 5, 2025*

*Інформаційні технології, електроніка,*
*механічна та електрична інженерія*
*Information technologies, electronics,*
*mechanical and electrical engineering*

2. Training Efficiency:

- PPO converged fastest, crossing positive rewards at 80K timesteps (vs TD3: 95K, SAC: 110K);

- PPO achieved lowest variance ($\sigma = 12.4$) compared to alternatives (TD3: 15.2, SAC: 18.9, DDPG: 21.3, A2C: 23.7);

- Training time: 4.5 hours for PPO vs 5.1–6.2 hours for alternatives.

3. Domain Randomization Impact:

- Full DR configuration achieved 84.8% Sim-to-Real transfer efficiency (+43.5% vs baseline);

- Key factors: surface friction (+22.7%), motor delay (+15.2%), sensor noise (+14.7%);

- Non-additive effects observed: individual sum (+68.0%) exceeds combined effect (+43.5%).

4. Model Deployment:

- INT8 quantization: 4× size reduction (11.27 KB → 2.82 KB) with 3.8% accuracy loss;

- ESP32 inference: 8.7 ms per step (115 Hz), 11.5× faster than required 10 Hz;

- PPO requires fewest parameters (2,818) among competitive algorithms.

5. Navigation Metrics:

- Collision Rate: 18.0% (matching Bug Algorithm's reactive safety);

- Path Efficiency: 73.2% (vs Bug: 48.3%, +51.6% improvement);

- Time to Goal: 47.2s (vs Bug: 95.3s, 50.5% faster).

Practical Implications:

PPO emerges as the optimal choice for embedded robotic navigation due to its unique combination of:

- highest success rate and lowest variance among RL methods;

- smallest model size enabling deployment on microcontrollers;

- fastest training and inference times;

- best robustness to hyperparameter choices.

Identified Limitations:

- 45% Success Rate in dynamic obstacle scenarios (training on static environments only);

- 15–20% performance degradation in Time to Goal and Velocity during Sim-to-Real transfer;

- SAC shows better sample efficiency but requires 2.4× more training time.

Future Research Directions:

1. Dynamic obstacle training to improve from 45% to target 80% SR in moving scenarios.

2. Curriculum learning for progressive environment complexity.

3. Multi-modal sensor fusion (vision + proprioception) for richer state representation.

4. Transfer learning to reduce training time for new environments.

5. Multi-agent coordination for collaborative navigation tasks.

The experimental methodology and quantitative results provide a reproducible framework for algorithm selection in resource-constrained robotic applications. PPO-based reinforcement learning with Domain Randomization is recommended as the standard approach for autonomous navigation on embedded platforms.

| **References** | **Література** |
|---|---|
| 1. Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., ... & Zaremba, W. (2020). Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1), 3–20. DOI: https://doi.org/10.1177/0278364919887447. | 1. Andrychowicz M., Wolski F., Ray A., Schneider J., Fong R., Welinder P., ... Zaremba W. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 2020. No. 39 (1), P. 3–20. DOI: https://doi.org/10.1177/0278364919887447. |

2. Bousmalis, K., Irpan, A., Wohlhart, P., Bai, Y., Kelcey, M., Kalakrishnan, M., ... & Levine, S. (2018). Using simulation and domain adaptation to improve efficiency of deep robotic grasping. In *2018 IEEE International Conference on Robotics and Automation* (ICRA) (pp. 4243–4250). DOI: https://doi.org/10.1109/ICRA.2018.8460875.

3. Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., & Zaremba, W. (2016). OpenAI Gym. *arXiv preprint* arXiv:1606.01540. DOI: https://doi.org/10.48550/ arXiv.1606.01540.

4. Chen, C., Seff, A., Kornhauser, A., & Xiao, J. (2015). DeepDriving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 2722–2730). URL: https://deepdriving.cs.princeton.edu/paper.pdf.

5. Coumans, E., & Bai, Y. (2016). PyBullet, a Python module for physics simulation for games, robotics and machine learning. *GitHub repository*. http://pybullet.org.

6. David, R., Duke, J., Jain, A., Janapa Reddi, V., Jeffries, N., Li, J., ... & Warden, P. (2021). TensorFlow Lite Micro: Embedded machine learning for TinyML systems. *Proceedings of Machine Learning and Systems*, 3, 800–811.

7. Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., ... & Kavukcuoglu, K. (2018). IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures. In *International Conference on Machine Learning* (pp. 1407–1416). PMLR.

8. Fujimoto, S., Hoof, H., & Meger, D. (2018). Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning* (pp. 1587–1596). PMLR.

9. Haarnoja, T., Zhou, A., Abbeel, P., & Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning* (pp. 1861-1870). PMLR.

10. Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., & Meger, D. (2018). Deep reinforcement learning that matters. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1). DOI: https://doi.org/10.1609/aaai.v32i1.11694.

11. Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., ... & Kalenichenko, D. (2018). Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 2704–2713).

2. Bousmalis K., Irpan A., Wohlhart P., Bai Y., Kelcey M., Kalakrishnan M., ... Levine S. Using simulation and domain adaptation to improve efficiency of deep robotic grasping. In *2018 IEEE International Conference on Robotics and Automation* (ICRA). 2018. P. 4243–4250. DOI: https://doi.org/10.1109/ICRA.2018.8460875.

3. Brockman G., Cheung V., Pettersson L., Schneider J., Schulman J., Tang J., Zaremba W. OpenAI Gym. *arXiv preprint.* 2016. arXiv:1606.01540. DOI: https://doi.org/10.48550/ arXiv.1606.01540.

4. Chen C., Seff A., Kornhauser A., Xiao J. DeepDriving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE International Conference on Computer Vision.* 2015. P. 2722–2730. URL: https://deep driving.cs.princeton.edu/paper.pdf.

5. Coumans E., Bai Y. PyBullet, a Python module for physics simulation for games, robotics and machine learning. *GitHub repository*. 2016. http://pybullet.org.

6. David R., Duke J., Jain A., Janapa Reddi V., Jeffries N., Li J., ... Warden P. TensorFlow Lite Micro: Embedded machine learning for TinyML systems. *Proceedings of Machine Learning and Systems*. 2021. No. 3. P. 800–811.

7. Espeholt L., Soyer H., Munos R., Simonyan K., Mnih V., Ward T., ... Kavukcuoglu K. IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures. In *International Conference on Machine Learning.* 2018. P. 1407–1416.

8. Fujimoto S., Hoof H., Meger D. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning.* 2018. P. 1587–1596.

9. Haarnoja T., Zhou A., Abbeel P., Levine S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning.* 2018. P. 1861–1870.

10. Henderson P., Islam R., Bachman P., Pineau J., Precup D., Meger D. Deep reinforcement learning that matters. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 2018. No. 32(1). DOI: https://doi.org/10.1609/aaai.v32i1.11694.

11. Jacob B., Kligys S., Chen B., Zhu M., Tang M., Howard A., ... Kalenichenko D. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* 2018. P. 2704–2713.

12. Kober, J., Bagnell, J. A., & Peters, J. (2013). Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11), 1238–1274. DOI: https://doi.org/10.1177/0278364913495721.

13. Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., ... & Wierstra, D. (2016). Continuous control with deep reinforcement learning. In *International Conference on Learning Representations* (ICLR).

14. Loquercio, A., Kaufmann, E., Ranftl, R., Dosovitskiy, A., Koltun, V., & Scaramuzza, D. (2020). Deep drone racing: From simulation to reality with domain randomization. *IEEE Transactions on Robotics*, 36(1), 1–14. DOI: https://doi.org/10.1109/TRO.2019.2942989.

15. Micheli, V., Alonso, E., & Fleuret, F. (2023). Transformers are sample-efficient world models. In *International Conference on Learning Representations* (ICLR).

16. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533. DOI: https://doi.org/10.1038/nature14236.

17. Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., ... & Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning* (pp. 1928–1937). PMLR.

18. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... & Chintala, S. (2019). PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 32.

19. Peng, X. B., Andrychowicz, M., Zaremba, W., & Abbeel, P. (2018). Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE International Conference on Robotics and Automation* (ICRA) (pp. 3803–3810). DOI: https://doi.org/10.1109/ICRA.2018.8460528.

20. Raffin, A., Hill, A., Gleave, A., Kanervisto, A., Ernestus, M., & Dormann, N. (2021). Stable-Baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268), 1–8.

21. Rajeswaran, A., Ghotra, S., Ravindran, B., & Levine, S. (2017). EPOpt: Learning robust neural network policies using model ensembles. In *International Conference on Learning Representations* (ICLR).

22. Sadeghi, F., & Levine, S. (2017). CAD2RL: Real single-image flight without a single real image. In

12. Kober J., Bagnell J. A., Peters J. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*. 2013. No. 32(11). P. 1238–1274. DOI: https://doi.org/10.1177/0278364913495721.

13. Lillicrap T. P., Hunt J. J., Pritzel A., Heess N., Erez T., Tassa Y., ... Wierstra D. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations* (ICLR). 2016.

14. Loquercio A., Kaufmann E., Ranftl R., Dosovitskiy A., Koltun V., Scaramuzza D. Deep drone racing: From simulation to reality with domain randomization. *IEEE Transactions on Robotics*. 2020. No. 36(1). P. 1–14. DOI: https://doi.org/10.1109/TRO.2019.2942989.

15. Micheli V., Alonso E., Fleuret F. Transformers are sample-efficient world models. In *International Conference on Learning Representations* (ICLR). 2023.

16. Mnih V., Kavukcuoglu K., Silver D., Rusu A. A., Veness J., Bellemare M.G., ... Hassabis D. Human-level control through deep reinforcement learning. *Nature*. 2015. No. 518(7540). P. 529–533. DOI: https://doi.org/10.1038/nature14236.

17. Mnih V., Badia A. P., Mirza M., Graves A., Lillicrap T., Harley T., ... Kavukcuoglu K. Asynchronous methods for deep reinforcement learning. In I*nternational Conference on Machine Learning.* 2016. P. 1928–1937.

18. Paszke A., Gross S., Massa F., Lerer A., Bradbury J., Chanan G., ... Chintala S. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems.* 2019. No. 32.

19. Peng X. B., Andrychowicz M., Zaremba W., Abbeel P. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE International Conference on Robotics and Automation* (ICRA). 2018. P. 3803–3810. DOI: https://doi.org/10.1109/ICRA.2018.8460528.

20. Raffin A., Hill A., Gleave A., Kanervisto A., Ernestus M., Dormann N. Stable-Baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*. 2021. No. 22(268). P. 1–8.

21. Rajeswaran A., Ghotra S., Ravindran B., Levine S. EPOpt: Learning robust neural network policies using model ensembles. In *International Conference on Learning Representations* (ICLR). 2017.

22. Sadeghi F., Levine S. CAD2RL: Real single-image flight without a single real image. In *Robotics:*

*Robotics: Science and Systems XIII.* DOI: https://doi.org/10.15607/RSS.2017.XIII.034.

23. Schulman, J., Levine, S., Abbeel, P., Jordan, M., & Moritz, P. (2015). Trust region policy optimization. In *International Conference on Machine Learning* (pp. 1889–1897). PMLR.

24. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint* arXiv:1707.06347.

25. Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction (2nd ed.). MIT Press.

26. Tai, L., Paolo, G., & Liu, M. (2017). Virtual-to-real deep reinforcement learning: Continuous control of mobile robots for mapless navigation. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems* (IROS) (pp. 31–36). DOI: https://doi.org/10.1109/IROS.2017.8202134.

27. Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., & Abbeel, P. (2017). Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems* (IROS) (pp. 23–30). DOI: https://doi.org/10.1109/IROS.2017.8202133.

28. Zhao, W., Queralta, J.P., & Westerlund, T. (2020). Sim-to-real transfer in deep reinforcement learning for robotics: A survey. In *2020 IEEE Symposium Series on Computational Intelligence* (SSCI) (pp. 737–744). DOI: https://doi.org/10.1109/SSCI47803.2020.9308468.

*Science and Systems XIII.* 2017. DOI: https://doi.org/10.15607/RSS.2017.XIII.034.

23. Schulman J., Levine S., Abbeel P., Jordan M., Moritz P. Trust region policy optimization. In *International Conference on Machine Learning.* 2015. P. 1889–1897.

24. Schulman J., Wolski F., Dhariwal P., Radford A., Klimov O. Proximal policy optimization algorithms. *arXiv preprint.* 2017. arXiv:1707.06347.

25. Sutton R. S., Barto A. G. Reinforcement learning: An introduction. 2nd ed. MIT Press, 2018.

26. Tai L., Paolo G., Liu M. Virtual-to-real deep reinforcement learning: Continuous control of mobile robots for mapless navigation. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems* (IROS). 2017. P. 31–36. DOI: https://doi.org/10.1109/IROS.2017.8202134.

27. Tobin J., Fong R., Ray A., Schneider J., Zaremba W., Abbeel P. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems* (IROS). 2017. P. 23–30. DOI: https://doi.org/10.1109/IROS.2017.8202133.

28. Zhao W., Queralta J. P., Westerlund T. Sim-to-real transfer in deep reinforcement learning for robotics: A survey. In *2020 IEEE Symposium Series on Computational Intelligence* (SSCI). 2020. P. 737–744. DOI: https://doi.org/10.1109/SSCI47803.2020.9308468.

***PYLYPENKO VLADYSLAV***
*Phd Student,*
*Department of Information and Computer Technologies,*
*Kyiv National University of Technologies*
*and Design, Ukraine*
*https://orcid.org/0000-0002-2761-4817*
*Scopus Author ID: 58089336700*
*E-mail: software.proger@gmail.com*

***VOLIVACH ANTONINA***
*Candidate of technical sciences, Associate Professor,*
*Department of Information and Computer Technologies,*
*Kyiv National University of Technologies*
*and Design, Ukraine*
*https://orcid.org/0000-0002-7119-7774*
*E-mail: volivach.ap@knutd.com.ua*

***AFTANDILYANTS VADYM***
*Candidate of sciences, Associate Professor,*
*Kyiv National University of Technologies*

***SKIDAN VLADYSLAVA***
*Candidate of technical sciences, Associate Professor,*
*Head of the Department of Information*
*and Computer Technologies,*
*Kyiv National University of Technologies*
*and Design, Ukraine*
*https://orcid.org/0000-0002-8358-9759*
*Scopus Author ID: 57210393405*
*E-mail: skidan.vv@knutd.edu.ua*

***MYTELSKA OLENA***
*Candidate of technical sciences, Associate Professor,*
*Kyiv National University of Technologies*
*and Design, Ukraine*
*https://orcid.org/0009-0004-4147-0866*
*E-mail: mytelska.ov@knutd.edu.ua*

*and Design, Ukraine*
*https://orcid.org/0000-0003-0660-1395*
*ResearcherID: Q-3511-2016*
*E-mail: aftandilyants@gmail.com*

**Владислав ПИЛИПЕНКО, Владислава СКІДАН, Антоніна ВОЛІВАЧ,**
**Олена МИТЕЛЬСЬКА, Вадим АФТАНДІЛЯНЦ**
*Київський національний університет технологій та дизайну, Україна*

**НАВЧАННЯ З ПІДКРІПЛЕННЯМ ДЛЯ АВТОНОМНОЇ НАВІГАЦІЇ РОБОТИЗОВАНИХ ПЛАТФОРМ В УМОВАХ НЕВИЗНАЧЕНОСТІ: РАНДОМІЗАЦІЯ ДОМЕНІВ ТА ПЕРЕНЕСЕННЯ СИМУЛЯЦІЇ В РЕАЛЬНІ УМОВИ**

**Мета.** *Розробити та експериментально оцінити ефективність системи автономної навігації для чотириколісної мобільної робототехнічної платформи на основі алгоритму Proximal Policy Optimization (PPO) з подальшим розгортанням моделі на мікроконтролері ESP32 та забезпеченням надійного Sim-to-Real трансферу в умовах невизначеного та динамічного середовища.*

**Методика.** *У дослідженні застосовано підхід глибокого навчання з підкріпленням із використанням алгоритму PPO. Проведено порівняльний аналіз із алгоритмами TD3, SAC, DDPG, A2C, Bug Algorithm та Random Policy. Оцінювання здійснювалося за показниками: успішність досягнення цілі, частота зіткнень, ефективність траєкторії, стабільність навчання (середня винагорода та стандартне відхилення).*

*Для подолання розриву між симуляцією та реальністю використано метод Domain Randomization із варіюванням шести фізичних параметрів: коефіцієнта тертя, маси робота, шуму гіроскопа, затримки двигунів, розміру та кількості перешкод. Нейронна мережа з архітектурою 64→32 нейрони була квантизована до формату INT8 з використанням TensorFlow Lite Micro та оптимізована для виконання на ESP32.*

**Результати.** *Алгоритм PPO продемонстрував найвищу ефективність серед усіх протестованих методів: успішність досягнення цілі 82.0%, середня винагорода — 847.3, найнижча варіативність результатів ($\sigma = 12.3$).*

*Статистичний аналіз підтвердив значущу перевагу PPO над альтернативними підходами ($p < 0.01$, Cohen's $d > 0.8$). Модель після квантизації до INT8 зменшилася до 2.8 КБ при втраті точності лише 2.3%. Час інференсу на ESP32 становив 1.2 мс, що забезпечує роботу в реальному часі на ресурсно-обмеженій платформі.*

**Наукова новизна.** *Запропоновано комплексний підхід до автономної навігації з використанням PPO з інтеграцією Domain Randomization для підвищення якості Sim-to-Real трансферу.*

*Виконано систематичне експериментальне порівняння сучасних алгоритмів глибокого навчання з підкріпленням у задачі автономної навігації мобільного робота.*

*Реалізовано ефективну квантизацію моделі PPO до формату INT8 з мінімальною втратою точності та успішним розгортанням на мікроконтролері ESP32.*

**Практична значимість.** *Отримані результати демонструють можливість впровадження алгоритмів глибокого навчання з підкріпленням у реальні мобільні робототехнічні системи з обмеженими обчислювальними ресурсами. Розроблена система може бути використана у сервісній робототехніці, автономних транспортних засобах малого класу, системах моніторингу та інспекції, забезпечуючи високу надійність навігації та швидкодію в реальному часі.*

**Ключові слова:** *навчання з підкріпленням; Proximal Policy Optimization; автономна навігація; Sim-to-Real трансфер; Domain Randomization; мобільна робототехніка; квантизація нейронних мереж; вбудовані системи; ESP32; TensorFlow Lite Micro.*