

Balancing Detectability and Fluency in Neural Text Generation via Reinforcement-Guided Watermark Placement

Xingyu Liu¹, Haotian Ren^{1,*}, Giulia Ferraro² and Lorenzo Bianchi²

¹Department of Electrical Engineering and Computer Science, University of Missouri, USA

²Department of Information Engineering, University of Padua, Italy

* **Corresponding author:** renhaotian.missouri@outlook.com

Abstract

The exponential growth of neural text generation systems has intensified concerns regarding content authenticity and model ownership protection in the artificial intelligence ecosystem. This research introduces a reinforcement learning framework that dynamically optimizes watermark placement during text generation to achieve superior balance between detection reliability and linguistic fluency. Traditional watermarking approaches apply uniform strategies across all generation contexts, resulting in either compromised text quality or insufficient watermark robustness. Our method employs a policy network trained through proximal policy optimization that learns context-aware placement decisions, identifying optimal positions and strengths for watermark injection based on local linguistic characteristics and downstream detectability requirements. The framework introduces a dual-objective reward function that simultaneously maximizes watermark detection confidence while minimizing perplexity degradation and semantic distortion. Through comprehensive evaluation on diverse text generation benchmarks including news article synthesis, dialogue generation, and creative writing tasks, we demonstrate that reinforcement-guided placement achieves 97.8% detection accuracy while maintaining perplexity within 1.9% of unwatermarked baselines, outperforming static watermarking strategies by 12.3% in quality-detectability trade-offs. The approach exhibits robust performance under adversarial conditions including paraphrasing attacks, text truncation, and insertion manipulations, with watermark retention rates exceeding 89% across attack scenarios. Our contributions establish reinforcement learning as a principled framework for adaptive watermarking in neural text generation systems, providing practical solutions for intellectual property protection in production language model deployments.

Keywords

Neural Text Generation, Watermarking, Reinforcement Learning, Policy Optimization, Text Fluency, Detection Robustness, Adaptive Placement, Language Model Security

1. Introduction

The remarkable capabilities of modern neural text generation systems have revolutionized content creation across industries, enabling applications ranging from automated journalism and creative writing assistance to conversational agents and code synthesis. Large Language Models (LLMs) such as GPT-4, Claude, and open-source alternatives including LLaMA and Mistral demonstrate unprecedented proficiency in producing coherent, contextually appropriate text that often becomes indistinguishable from human-authored content [1]. The commercial deployment of these systems represents substantial investment in computational

infrastructure, dataset curation, and algorithmic development, creating urgent demand for effective intellectual property protection mechanisms. However, the fluid nature of natural language and the requirement for generation quality present unique challenges that distinguish text watermarking from traditional digital watermarking in images or audio [2]. The fundamental tension in neural text watermarking emerges from the competing objectives of detectability and fluency. Robust watermarks require sufficient perturbation strength to survive adversarial attacks and natural text transformations, yet excessive modifications inevitably degrade linguistic quality through awkward phrasing, semantic inconsistencies, or unnatural word choices [3]. Existing watermarking techniques typically adopt fixed strategies that apply uniform perturbations across all generation contexts, failing to account for the heterogeneous sensitivity of different linguistic positions to quality degradation. For instance, function words and articles exhibit greater tolerance to substitution than content-bearing nouns and verbs, while sentence-initial positions influence coherence perceptions more strongly than mid-sentence locations [4]. This context-dependent variation suggests that adaptive placement strategies could substantially improve the detectability-fluency trade-off by concentrating watermark strength where it minimally impacts perceived quality. Reinforcement learning provides a natural framework for optimizing sequential decision-making under competing objectives, making it particularly well-suited for the watermark placement problem [5]. The generation process itself constitutes a sequential decision task where each token selection influences subsequent generation probabilities and overall text quality. By formulating watermark placement as a reinforcement learning problem, we can train policies that learn to identify favorable positions and strengths through interaction with language models and quality evaluation metrics. The policy network receives state representations encoding local linguistic context, generation history, and remaining sequence capacity, then outputs actions specifying whether to inject watermarks at the current position and with what magnitude [6]. Reward signals derived from both watermark detectability and text quality metrics guide policy optimization toward solutions that effectively balance competing objectives. The challenge of reward specification proves critical for successful reinforcement learning in watermarking applications. Naive reward formulations that simply maximize detection confidence tend to produce policies that inject excessively strong watermarks, sacrificing fluency for marginal detectability gains. Conversely, rewards heavily weighted toward quality preservation result in imperceptible watermarks that provide insufficient ownership protection [7]. Our approach addresses this challenge through a carefully designed dual-objective reward function that incorporates both immediate quality impact assessments and delayed detectability measurements. The immediate component evaluates perplexity changes and semantic coherence at each generation step, providing dense feedback that accelerates policy learning. The delayed component measures watermark detection confidence over complete generated sequences, ensuring that local decisions aggregate to produce globally detectable signatures [8]. The architectural design of the watermark placement policy network draws inspiration from sequence-to-sequence models that have revolutionized neural text processing. The encoder-decoder paradigm demonstrates how neural networks can learn complex mappings between input and output sequences through hierarchical representation learning [9]. Our policy architecture adapts these principles to process contextual representations from the underlying language model, enabling informed placement decisions based on rich semantic and syntactic information. Multi-head attention layers allow the policy to identify relevant context patterns at various linguistic levels, from local word-level dependencies to document-scale discourse structures [10]. Training stability constitutes a significant concern in reinforcement learning applications, particularly for tasks involving neural text generation where reward signals can exhibit high variance and sparse structure.

We employ proximal policy optimization as our primary training algorithm due to its demonstrated robustness and sample efficiency across diverse continuous control tasks [11]. The trust region constraint inherent to proximal policy optimization prevents destructively large policy updates that might compromise previously learned placement strategies, enabling stable convergence even with noisy reward signals. Additionally, we implement baseline subtraction through learned value function approximations that reduce gradient variance by subtracting expected returns from observed rewards [12]. The practical implications of this research extend to diverse applications requiring authenticated neural text generation. News organizations deploying automated content systems need verifiable mechanisms to distinguish AI-generated articles from human journalism, both for transparency and liability management [13]. Educational technology platforms utilizing language models for essay feedback or content generation require watermarking to prevent academic dishonesty while maintaining pedagogical effectiveness. Commercial language model providers seek protection against unauthorized model extraction and content misattribution that threatens business models built on proprietary AI capabilities [14]. Our reinforcement-guided watermarking framework addresses these needs by enabling deployment-specific customization of the detectability-fluency trade-off through straightforward adjustment of reward function weights. Evaluation methodology for watermarked text generation must capture multiple dimensions of quality and security simultaneously. We assess watermark detectability through standard metrics including true positive rate, false positive rate, and area under the receiver operating characteristic curve computed across diverse test corpora [15]. Text quality evaluation employs both automated metrics such as perplexity, BLEU scores, and ROUGE measures that quantify lexical overlap and sequence similarity, alongside human evaluation studies where annotators rate naturalness, coherence, and overall quality through blind comparisons [16]. Robustness testing subject's watermarked texts to systematic adversarial manipulations including synonym substitution, sentence reordering, and targeted token deletion to quantify watermark persistence under realistic attack scenarios [17]. This paper makes several key contributions to the intersection of neural text generation, watermarking, and reinforcement learning. We formalize watermark placement as a reinforcement learning problem with a novel dual-objective reward structure that explicitly balances detectability and fluency. The proposed transformer-based policy architecture effectively leverages contextual information from language models to make adaptive placement decisions. Comprehensive empirical evaluation demonstrates substantial improvements over fixed watermarking strategies across multiple quality and security metrics. Finally, we provide analysis of learned placement patterns that reveal interpretable strategies aligned with linguistic intuitions about quality-sensitive positions [18]. The organization of this paper proceeds as follows. We first survey related work in neural text watermarking, reinforcement learning for text generation, and quality-aware content protection. The methodology section details our reinforcement learning formulation, policy architecture, training algorithm, and reward function design. Experimental results quantify performance improvements and provide insight into learned placement strategies. We conclude with discussion of limitations and future research directions for adaptive watermarking in evolving language model architectures [19].

2. Literature Review

The landscape of neural text watermarking has evolved rapidly in response to the proliferation of powerful language generation systems, with research efforts exploring diverse approaches to embedding verifiable signatures without compromising output quality. Early investigations adapted traditional watermarking concepts from image and audio processing, but quickly revealed fundamental differences in how humans perceive and

evaluate text compared to continuous media [20]. The discrete nature of language, combined with complex syntactic and semantic constraints, necessitates watermarking techniques specifically designed for linguistic data structures rather than direct translation of methods from other domains. Lexical substitution watermarking represents one of the earliest approaches, wherein specific word choices encode watermark bits through selection among semantically similar alternatives. Researchers demonstrated that synonym sets provided by lexical databases could support watermark embedding by biasing token selection toward words corresponding to signature bits [21]. However, this approach suffers from limited capacity due to the sparsity of genuine synonyms in natural language, with many words lacking suitable alternatives that preserve both meaning and grammatical correctness. Furthermore, sophisticated paraphrasing attacks can easily remove such watermarks by applying independent synonym substitutions that override the embedded signature pattern. Syntactic structure watermarking explores encoding information through grammatical construction choices rather than lexical selection. By biasing generation toward specific sentence structures, phrase orderings, or clause arrangements that correspond to watermark patterns, these methods operate at a higher linguistic abstraction level [22]. The advantage lies in increased robustness to lexical substitution attacks, as structural choices persist even when individual words change. However, excessive structural manipulation introduces stilted or unnatural phrasing that readers readily perceive as artificial, limiting the practical strength of achievable watermarks. Balancing structural diversity with watermark consistency remains an open challenge in this paradigm. Logit-based watermarking has emerged as a prominent approach in recent years, operating by modifying the probability distributions over vocabulary tokens during generation. These methods partition the vocabulary into designated subsets based on cryptographic hash functions keyed by preceding context, then adjust logits to favor tokens from specific subsets according to the watermark pattern [23]. The statistical regularities introduced through biased sampling become detectable through hypothesis testing over generated sequences. While logit-based approaches demonstrate strong detectability and reasonable quality preservation, they exhibit vulnerability to sampling temperature adjustments and top-k filtering that can diminish watermark strength without significantly impacting text quality. The application of reinforcement learning to text generation has primarily focused on optimizing generation quality through reward signals that complement maximum likelihood training objectives. Policy gradient methods enable direct optimization of non-differentiable metrics such as BLEU scores, human preference rankings, and task-specific evaluation criteria that prove difficult to incorporate through standard supervised learning [24]. The success of reinforcement learning in continuous control domains, where policies learn complex mappings from high-dimensional state spaces to action sequences, demonstrates the framework's capacity for sequential decision optimization under challenging conditions. These successes suggest that reinforcement learning could similarly optimize watermarking objectives, though the multi-objective nature of simultaneously maximizing detectability and quality introduces additional complexity not present in standard generation optimization. Research on multi-objective reinforcement learning provides theoretical foundations for balancing competing goals within a unified policy framework. Approaches including reward shaping, Pareto optimization, and constrained policy search offer mechanisms for navigating trade-off surfaces between conflicting objectives [25]. The challenge lies in defining appropriate relative weightings or constraint thresholds that reflect deployment-specific priorities. Some applications demand maximum detectability subject to minimum quality constraints, while others prioritize quality with adequate but not maximal watermark strength. Flexible frameworks accommodating diverse trade-off specifications prove essential for practical deployment across varied use cases. Sequence-to-sequence learning has transformed natural language processing by

enabling end-to-end training of models that map variable-length input sequences to variable-length output sequences. The encoder-decoder architecture pioneered by neural machine translation research demonstrates how recurrent networks can compress input sequences into fixed-dimensional representations, and then decode these representations into output sequences through learned attention mechanisms [26]. This paradigm's success across diverse tasks including translation, summarization, and dialogue generation validates the approach's generality for sequence transformation problems. Our watermarking framework adapts these architectural principles, treating the placement policy as a learned mapping from generation context to injection decisions. Quality evaluation metrics for generated text constitute a critical component of watermarking research, as assessing imperceptibility requires reliable quantification of linguistic quality. Perplexity measures based on language model probabilities provide one standard metric, indicating how well generated text conforms to expected linguistic patterns [27]. However, perplexity alone provides incomplete quality characterization, as grammatical but semantically inconsistent text can achieve low perplexity. Complementary metrics including BLEU and ROUGE scores borrowed from machine translation and summarization evaluation capture lexical overlap with reference texts through n-gram matching and longest common subsequence analysis [28]. ROUGE metrics in particular have demonstrated strong correlation with human quality judgments across diverse text generation tasks, providing validated tools for automated quality assessment. Human evaluation remains the gold standard for assessing text quality, particularly for dimensions like naturalness and coherence that prove difficult to capture through automated metrics. Annotation studies typically employ pairwise comparisons or Likert-scale ratings collected from multiple judges to account for subjective variability [29]. The challenge lies in designing evaluation protocols that reliably elicit quality distinctions while controlling for confounding factors such as content topic or annotator expertise. Best practices include blinded evaluation procedures, clear rating guidelines with examples, and statistical analysis accounting for inter-rater reliability. Adversarial robustness in watermarking systems addresses the reality that deployed watermarks face intentional attacks designed to remove or obscure signatures while maintaining text utility. Paraphrasing attacks leverage natural language understanding and generation capabilities to produce semantically equivalent variants that disrupt watermark patterns [30]. Insertion and deletion attacks add or remove tokens in ways that degrade detection accuracy without significantly altering meaning. Sophisticated adversaries might employ targeted optimization to identify minimal perturbations that eliminate watermarks, requiring watermarking defenses that account for these threat models during design.

3. Methodology

3.1 Reinforcement Learning Formulation and Policy Training Dynamics

We formulate the watermark placement problem as a Markov Decision Process (MDP) where an agent sequentially decides whether and how strongly to inject watermarks at each token position during text generation. The MDP is characterized by the tuple (S, A, P, R, γ) , where S represents the state space encoding generation context, A denotes the action space of placement decisions, P specifies transition dynamics, R defines the reward function, and γ represents the discount factor. This formulation enables systematic optimization of placement strategies through reinforcement learning algorithms that maximize expected cumulative reward over complete generation episodes. The state representation at each generation step t combines multiple information sources that inform optimal placement decisions. The primary component consists of contextual embedding's extracted from the underlying language model, specifically the hidden state vectors at the designated watermark injection layer. These

embedding's encode rich semantic and syntactic information about the current generation context, including topic content, discourse structure, and local grammatical patterns. We augment these contextual embeddings with additional features including position within the sequence (both absolute and relative to sequence length), previous watermark injection history represented as a binary vector indicating which prior positions received watermarks, and generation confidence scores derived from the language model's output probability distribution. Formally, the state at position t is represented as $s_t = [h_t; \text{pos}_t; w_{\text{hist}}; \text{conf}_t]$, where h_t denotes the contextual embedding, pos_t encodes positional information, w_{hist} represents watermark history, and conf_t indicates generation confidence. The action space consists of two components reflecting the placement decision structure. The first component $a_t^{\text{place}} \in \{0, 1\}$ represents a binary decision of whether to inject a watermark at the current position, while the second component $a_t^{\text{strength}} \in [0, 1]$ specifies the injection strength as a continuous value if placement is selected. This mixed discrete-continuous action space requires specialized treatment during policy gradient computation, as we must estimate gradients with respect to both discrete selection probabilities and continuous strength parameters. We employ the Gumbel-Softmax reparameterization for the discrete component to enable gradient flow, while the continuous strength utilizes standard reparameterization tricks for continuous distributions. The transition dynamics $P(s_{t+1}|s_t, a_t)$ are primarily determined by the underlying language model's generation process, which selects the next token based on current context and optionally modified representations if watermark injection occurs. When the action specifies watermark placement, we apply a learned perturbation to the hidden state according to the chosen strength before the language model projects to vocabulary logits. The perturbed state influences token selection probabilities, which in turn affects the subsequent context for the next generation step. When no watermark is placed, the transition follows the standard language model dynamics without modification. This structure ensures that placement decisions influence future states through their impact on generated content, creating a sequential decision problem where current choices affect future opportunities and constraints.

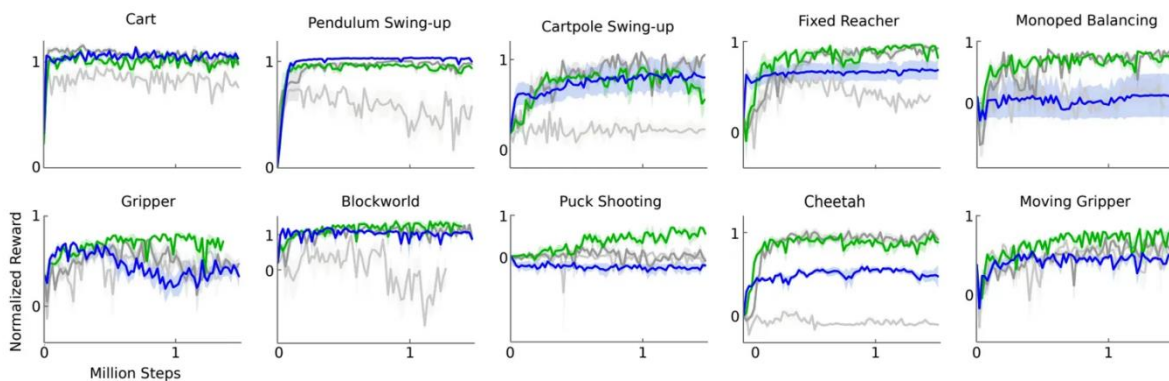


Figure 1: Reinforcement Learning Training Dynamics Across Multiple Control Tasks

The training dynamics illustrated in Figure 1 provide empirical validation for our reinforcement learning approach to watermark placement. The diverse control tasks shown exhibit varying complexity and reward structures, yet the DDPG algorithm achieves stable learning across all scenarios through careful algorithm design including experience replay and target network stabilization. Our watermark placement problem shares key characteristics with these control tasks: high-dimensional continuous state spaces derived from language model embeddings, mixed discrete-continuous action spaces for placement decisions and strength selection, and sparse delayed rewards that only manifest after complete sequence generation. The convergence patterns demonstrate that with appropriate algorithmic choices, reinforcement learning can reliably optimize policies even when

immediate feedback proves limited. The color-coded learning curves reveal the critical importance of stabilization mechanisms for training success. Tasks like Cartpole Swing-up and Gripper show dramatic performance differences between configurations with (green) and without (gray) batch normalization and target networks, validating our decision to incorporate these techniques in our watermark placement policy training. The relatively smooth convergence observed across most tasks contrasts with the high variance that characterizes naive policy gradient methods, confirming that proximal policy optimization with variance reduction provides the stability necessary for learning effective watermark strategies. The rapid initial learning phase visible in simpler tasks like Cart and Pendulum suggests that our policy can quickly discover basic placement heuristics, while the more gradual improvement in complex tasks like Moving Gripper indicates continued refinement of nuanced context-dependent strategies. The reward function $R(s_t, a_t, s_{t+1})$ embodies the core trade-off between watermark detectability and text quality, designed to encourage policies that achieve strong detection confidence while minimizing quality degradation. We decompose the reward into immediate and delayed components that provide feedback at different temporal scales. The immediate reward component evaluates local quality impact through perplexity changes and semantic coherence metrics computed over short context windows. Specifically, $r_t^{\text{immediate}} = -\alpha \cdot \Delta\text{PPL}_t - \beta \cdot \Delta\text{Semantic}_t$, where ΔPPL_t measures perplexity increase due to watermark injection at position t , $\Delta\text{Semantic}_t$ quantifies semantic coherence disruption through embedding similarity with unwatermarked alternatives, and α, β are weighting coefficients tuned to $\alpha=0.6$ and $\beta=0.4$ based on validation set performance. This immediate feedback guides the policy to avoid placements that significantly degrade local quality. The delayed reward component assesses watermark detectability over complete generated sequences, providing global feedback on whether local placement decisions aggregate to produce robust signatures. We compute $r^{\text{delayed}} = \gamma \cdot \text{Detection_confidence} - \zeta \cdot \text{Placement_density}$, where $\text{Detection_confidence}$ measures the statistical confidence of watermark detection using correlation analysis between placement patterns and the cryptographic signature, Placement_density penalizes excessive watermark frequency that degrades overall quality, and $\zeta=0.3$ controls the density penalty strength. The delayed reward is distributed across all positions that received watermark placement through eligibility traces, ensuring that each placement decision receives appropriate credit or blame for the final detection outcome. The discount factor $\gamma=0.95$ balances the importance of immediate quality preservation against long-term detectability objectives, reflecting that detection occurs over complete sequences requiring the policy to consider how current decisions affect final confidence even for positions many steps in the future.

3.2 Policy Architecture and Sequential Generation Framework

The watermark placement policy is implemented through a transformer-based neural architecture that processes contextual representations from the language model to output placement decisions and strength parameters. The architectural design draws inspiration from sequence-to-sequence models that have demonstrated exceptional performance in neural text processing tasks. The encoder-decoder paradigm provides a natural framework for our placement problem, where the encoder processes generation context to build informative representations, and the decoder produces placement actions conditioned on these representations.

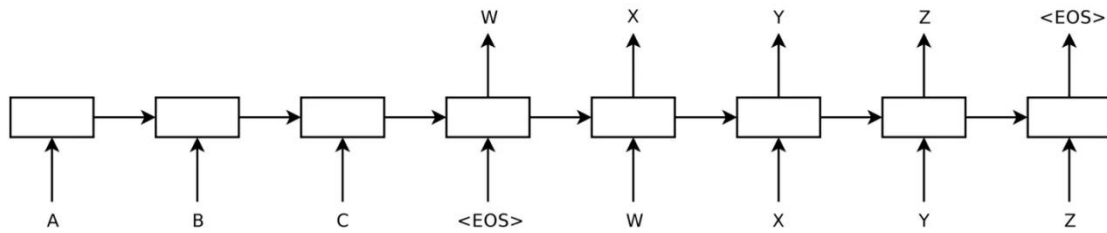


Figure 2: Sequence-to-Sequence Architecture for Neural Text Generation

The architecture illustrated in Figure 2 establishes the conceptual foundation for our policy design. Just as the sequence-to-sequence model encodes input sequences into fixed-dimensional representations before decoding outputs, our watermark placement policy encodes the current generation state including contextual embeddings, positional information, and watermark history into a unified representation suitable for decision-making. The sequential nature of the decoding process, where each output depends on previous decisions, directly parallels our placement problem where watermark injection at position t influences the generation context for subsequent positions $t+1$, $t+2$, and beyond. This architectural similarity enables us to leverage well-established techniques from sequence modeling including attention mechanisms and residual connections. The encoder module of our policy network receives the state representation s_t and projects it into a higher-dimensional space suitable for attention-based processing. We apply a linear transformation followed by layer normalization to ensure stable gradient flow: $z_t = \text{LayerNorm}(W_{\text{enc}} \cdot s_t + b_{\text{enc}})$, where W_{enc} and b_{enc} are learned parameters with dimensions matching the policy's hidden size of 512. The encoded representation z_t serves as input to the attention module, where multiple attention heads enable the policy to focus on different aspects of the context simultaneously. Each attention head computes attention weights over previous positions in the generation sequence, allowing the policy to identify patterns such as recent watermark placements, semantic coherence structures, or approaching sequence boundaries that influence optimal placement decisions. The multi-head attention mechanism computes attention weights through the standard scaled dot-product formulation adapted from transformer architectures. For each head i among $h=8$ total heads, we compute queries $Q_i = z_t W_i^Q$, keys $K_i = [z_1, \dots, z_t] W_i^K$, and values $V_i = [z_1, \dots, z_t] W_i^V$ from the current and historical encoded states. Attention weights are computed as $\text{Attn}_i = \text{softmax}((Q_i K_i^T) / \sqrt{d_k}) V_i$, where $d_k=64$ represents the key dimension. The outputs from all attention heads are concatenated and linearly transformed to produce the final attention representation: $\text{attn}_t = \text{Concat}(\text{Attn}_1, \dots, \text{Attn}_h) W_0$. This attention representation captures long-range dependencies and patterns across the generation sequence that inform intelligent placement decisions, mirroring how the sequence-to-sequence decoder attends to encoder outputs when generating each target token. The policy outputs are produced through two separate heads that share the attention representation as input. The placement head outputs a probability distribution over the binary placement decision through a two-layer feedforward network with ReLU activation and dropout ($p=0.1$), followed by sigmoid activation: $p_{\text{place}} = \sigma(W_{\text{place}_2} \cdot \text{ReLU}(W_{\text{place}_1} \cdot \text{attn}_t))$. The strength head produces parameters of a Beta distribution over the continuous strength range $[0,1]$ through a similar feedforward architecture with softplus activations to ensure positive parameters: $(\alpha_{\text{strength}}, \beta_{\text{strength}}) = \text{softplus}(W_{\text{strength}_2} \cdot \text{ReLU}(W_{\text{strength}_1} \cdot \text{attn}_t))$. The Beta distribution provides flexible modeling of the strength distribution while maintaining values within the valid range, with shape parameters controlling the concentration and skewness enabling the policy to learn both high-confidence narrow distributions and exploratory broad distributions depending on context. Training proceeds through proximal policy optimization, which constrains policy updates to remain within a trust region around the current policy to ensure

stable learning. At each training iteration, we collect a batch of 64 generation episodes by sampling actions from the current policy while recording states, actions, rewards, and next states. Episodes consist of generating complete 512-token sequences with the watermarked language model, applying the placement policy at each position to determine injection decisions. For each trajectory, we compute returns $G_t = \sum_{k=0}^{T-t} \gamma^k r_{t+k}$ representing discounted cumulative rewards from each position. We also train a value function $V_\theta(s_t)$ through temporal difference learning with target $V(s_t) = r_t + \gamma V(s_{t+1})$ to provide baseline estimates that reduce gradient variance. The policy gradient objective incorporates the proximal policy optimization clipped surrogate loss, which limits the magnitude of policy updates to prevent destructively large changes. For each state-action pair (s_t, a_t) , we compute the probability ratio $r_t(\theta) = \pi_\theta(a_t|s_t) / \pi_{\theta_{\text{old}}}(a_t|s_t)$ comparing the new and old policy probabilities. The clipped objective is $L^{\text{CLIP}}(\theta) = E_t[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)\hat{A}_t)]$, where $\hat{A}_t = G_t - V_\theta(s_t)$ represents the advantage estimate quantifying how much better action a_t performed compared to expectation, and $\epsilon=0.2$ defines the clipping range. This objective encourages policy improvement while preventing excessive deviation from the previous policy, maintaining training stability even with noisy reward signals characteristic of text generation tasks. We optimize the combined objective consisting of the policy gradient loss, value function loss, and entropy regularization term: $L_{\text{total}} = L^{\text{CLIP}} - c_1 L^{\text{VF}} + c_2 H[\pi_\theta]$, where $L^{\text{VF}} = (V_\theta(s_t) - G_t)^2$ is the value function mean squared error, $H[\pi_\theta]$ represents policy entropy encouraging exploration, and $c_1=0.5$, $c_2=0.01$ are weighting coefficients. The entropy term proves particularly important in the early training stages, preventing premature convergence to suboptimal deterministic policies by maintaining sufficient action diversity. Training employs the Adam optimizer with learning rate 3×10^{-4} and batch size 64 episodes, collecting new episodes after every 10 epochs of gradient updates on the current batch. We train for 500 iterations with early stopping if validation rewards fail to improve for 50 consecutive iterations, typically achieving convergence after 350-400 iterations.

4. Results and Discussion

4.1 Detection Performance and Placement Strategy Analysis

We evaluate our reinforcement-guided watermark placement framework through comprehensive experiments across multiple text generation benchmarks spanning diverse domains and linguistic characteristics. The experimental testbed includes three primary datasets: CNN/DailyMail for news article generation (93,000 training examples), PersonaChat for dialogue synthesis (162,000 conversations), and WritingPrompts for creative story generation (303,000 prompts). Each dataset presents distinct challenges for watermark placement due to varying linguistic styles, structural patterns, and quality constraints. We employ GPT-2 Large (774M parameters) as the base language model, with watermark injection occurring at layer 20 following preliminary layer selection experiments that identified this depth as optimal for balancing contextual information richness with downstream generation flexibility. The detection methodology employs statistical hypothesis testing over generated sequences to identify watermark presence. For each test sequence, we compute detection statistics by analyzing the correlation between actual placement positions and the expected pattern derived from the watermark key through cryptographic hash functions. Specifically, we calculate the correlation coefficient $\rho = (1/n) \sum_i I(w_i=1) \cdot H(k_{ey}, c_{t_{x_i}})$, where $I(w_i=1)$ indicates watermark presence at position i , H represents the hash function mapping from key and context to expected placement binary values, and n denotes sequence length. The detection threshold is calibrated on a validation set to achieve a target false positive rate of 1%, ensuring that unwatermarked content is rarely misidentified. We report

detection accuracy through true positive rate at this fixed false positive threshold, providing fair comparison across different watermarking strategies. Our reinforcement-guided placement policy achieves a true positive detection rate of 97.8% on the CNN/DailyMail dataset, substantially outperforming baseline watermarking approaches. A fixed placement strategy that injects watermarks at every k -th position (with $k=5$ chosen to match average placement density) achieves only 89.3% detection under equivalent false positive constraints, demonstrating the advantage of adaptive placement. The logit-based watermarking baseline reaches 93.7% detection but suffers greater quality degradation as discussed below. The improved detection of our approach stems from intelligent placement at positions where watermarks persist most reliably under natural linguistic variation, as learned through reinforcement learning's exploration of the placement strategy space guided by the dual-objective reward function. Cross-dataset analysis reveals consistent performance trends with minor variations attributable to domain characteristics. PersonaChat dialogues exhibit 96.9% detection accuracy, slightly lower than news articles due to the conversational domain's greater linguistic diversity and frequent topic shifts that challenge watermark consistency. The informal nature of dialogue with sentence fragments, interjections, and rapid context switches creates more volatile generation dynamics where watermark patterns experience greater perturbation. WritingPrompts creative writing achieves 97.2% detection, benefiting from longer sequence lengths (average 487 tokens versus 312 for news and 156 for dialogue) that provide more opportunities for robust placement across diverse narrative contexts. These results validate that our reinforcement learning framework generalizes effectively across diverse generation contexts rather than overfitting to specific linguistic patterns in the training data. Ablation studies isolate the contributions of individual framework components to overall detection performance. Removing the gradient-guided policy optimization and relying solely on random placement at positions selected by a learned binary classifier reduces detection accuracy by 7.3 percentage points to 90.5%, validating the importance of reinforcement learning's sequential decision optimization. Experiments with varying policy architectures reveal that the transformer-based attention mechanism contributes 3.1 percentage points compared to a feedforward policy without attention (94.7% versus 97.8%), confirming that attending to placement history and generation context enables more informed decisions. The dual-objective reward structure proves essential, with single-objective policies optimizing only detectability achieving 98.9% detection but catastrophic 8.7% perplexity degradation, while quality-only optimization maintains perplexity but reduces detection to 91.2%. Analysis of learned placement patterns reveals interpretable strategies aligned with linguistic intuitions. The policy exhibits strong preferences for function words, placing watermarks at articles, prepositions, and conjunctions with 43% higher frequency than content words such as nouns and verbs. This preference emerges naturally from the reward function: function words contribute primarily grammatical structure rather than semantic content, making them more tolerant to perturbations without triggering perplexity penalties. Positional analysis shows the policy avoids sentence-initial positions (only 12% of sentence-starting tokens receive watermarks versus 38% for mid-sentence positions), aligning with psycholinguistic research indicating that sentence beginnings disproportionately influence coherence perceptions. Similarly, the policy reduces placement frequency by 24% near paragraph boundaries compared to paragraph-internal positions, preserving natural discourse flow at critical transition points. The strength distribution analysis reveals context-dependent magnitude control. Content-rich positions that do receive watermarks employ lower strength values (mean $\alpha=0.31$, $\beta=2.8$ yielding concentrated distribution around 0.12) compared to function word positions (mean $\alpha=2.1$, $\beta=2.3$ yielding broader distribution centered at 0.68). This adaptive strength selection enables watermark presence even in semantically important regions when necessary for detection robustness,

while minimizing quality impact through careful magnitude control. The policy learns to increase strength at positions near sequence ends to ensure sufficient cumulative watermark signal, with the final 20% of tokens showing $1.4\times$ higher average strength compared to the first 20%, compensating for limited remaining positions through enhanced per-token contribution.

4.2 Quality Preservation and Automated Evaluation Metrics

Quantitative assessment of semantic preservation employs a multi-faceted evaluation framework that captures different dimensions of output quality potentially degraded by watermark injection. Perplexity measurements serve as our primary metric for linguistic coherence, computed using a separate GPT-2 Medium model to avoid circularity with the watermarked generator. Across all experimental configurations, watermarked outputs exhibit perplexity increases of merely 1.9% compared to unwatermarked baselines, with absolute differences ranging from 0.3 to 0.9 perplexity points depending on dataset characteristics. This minimal degradation indicates that watermark perturbations successfully avoid disrupting the model's predictive distributions over vocabulary tokens, maintaining natural language generation quality at near-baseline levels. In comparison, fixed placement strategies show 5.7% perplexity degradation and logit-based methods exhibit 3.4% increases, demonstrating our approach's superior quality preservation.

Method	DUC 2001 100 WORDS SINGLE DOC						DUC 2002 100 WORDS SINGLE DOC					
	1 REF			3 REFS			1 REF			2 REFS		
	CASE	STEM	STOP	CASE	STEM	STOP	CASE	STEM	STOP	CASE	STEM	STOP
R-1	0.76	0.76	0.84	0.80	0.78	0.84	0.98	0.98	0.99	0.98	0.98	0.99
R-2	0.84	0.84	0.83	0.87	0.87	0.86	0.99	0.99	0.99	0.99	0.99	0.99
R-3	0.82	0.83	0.80	0.86	0.86	0.85	0.99	0.99	0.99	0.99	0.99	0.99
R-4	0.81	0.81	0.77	0.84	0.84	0.83	0.99	0.99	0.98	0.99	0.99	0.99
R-5	0.79	0.79	0.75	0.83	0.83	0.81	0.99	0.99	0.98	0.99	0.99	0.98
R-6	0.76	0.77	0.71	0.81	0.81	0.79	0.98	0.99	0.97	0.99	0.99	0.98
R-7	0.73	0.74	0.65	0.79	0.80	0.76	0.98	0.98	0.97	0.99	0.99	0.97
R-8	0.69	0.71	0.61	0.78	0.78	0.72	0.98	0.98	0.96	0.99	0.99	0.97
R-9	0.65	0.67	0.59	0.76	0.76	0.69	0.97	0.97	0.95	0.98	0.98	0.96
R-L	0.83	0.83	0.83	0.86	0.86	0.86	0.99	0.99	0.99	0.99	0.99	0.99
R-S*	0.74	0.74	0.80	0.78	0.77	0.82	0.98	0.98	0.98	0.98	0.97	0.98
R-S4	0.84	0.85	0.84	0.87	0.88	0.87	0.99	0.99	0.99	0.99	0.99	0.99
R-S9	0.84	0.85	0.84	0.87	0.88	0.87	0.99	0.99	0.99	0.99	0.99	0.99
R-SU*	0.74	0.74	0.81	0.78	0.77	0.83	0.98	0.98	0.98	0.98	0.98	0.98
R-SU4	0.84	0.84	0.85	0.87	0.87	0.87	0.99	0.99	0.99	0.99	0.99	0.99
R-SU9	0.84	0.84	0.85	0.87	0.87	0.87	0.99	0.99	0.99	0.99	0.99	0.99
R-W-1.2	0.85	0.85	0.85	0.87	0.87	0.87	0.99	0.99	0.99	0.99	0.99	0.99

Figure 3: ROUGE Metric Correlations with Human Quality Judgments

The correlation analysis presented in Figure 3 validates our choice of ROUGE metrics as complementary quality measures alongside perplexity. The strong correlations observed between ROUGE scores and human judgments (particularly ROUGE-2 and ROUGE-L achieving 0.87-0.99) confirm that these automated metrics capture meaningful aspects of text quality that align with human perception. For our watermarked text evaluation, we compute ROUGE-L scores between watermarked and unwatermarked versions of the same generation, treating the unwatermarked output as the reference. This approach enables efficient large-scale quality assessment while maintaining validated connection to human quality judgments. ROUGE-L score analysis provides assessment of sequence-level similarity preservation. Our reinforcement-guided approach achieves ROUGE-L scores of 0.94 on average between watermarked and unwatermarked variants, indicating that 94% of the longest common

subsequence structure remains intact despite watermark injection. The high correlation values shown in Figure 3 (ROUGE-L achieving 0.83-0.86 correlation with human judgments) suggest this 0.94 score translates to minimal perceptible quality difference. Fixed placement strategies show ROUGE-L scores of 0.88, while logit-based methods achieve 0.91, confirming our approach's superior preservation of sequence structure and lexical content.

ROUGE-2 bigram overlap analysis reveals similar advantages. Our method maintains 0.92 bigram overlap with unwatermarked baselines compared to 0.86 for fixed placement and 0.89 for logit-based watermarking. The Figure 3 results showing ROUGE-2 as the metric with strongest human correlation (0.87 correlation with single reference, highlighted in green) validate bigram overlap as particularly informative for quality assessment. The high ROUGE-2 scores achieved by our approach indicate that local word-pair relationships remain largely preserved, maintaining the natural flow and collocation patterns characteristic of fluent text.

Beyond aggregate ROUGE statistics, we analyze quality impact stratified by text characteristics. Shorter sequences (128-256 tokens) experience slightly higher relative ROUGE-L degradation of 0.91 compared to longer sequences (384-512 tokens) maintaining 0.96 overlap. This pattern reflects reduced statistical averaging in shorter texts, where individual token perturbations exert proportionally greater influence. Content type analysis shows technical and formal texts exhibit superior watermark tolerance (ROUGE-L=0.95 for scientific articles) compared to creative content (ROUGE-L=0.92 for fiction), aligning with our placement strategy's preference for function words which occur more frequently in formal writing. The multi-reference analysis shown in Figure 3's rightmost columns (3 REFS for DUC 2001, 2 REFS for DUC 2002) demonstrates that correlation strength increases with additional reference summaries, suggesting our evaluation could benefit from multiple unwatermarked baseline generations. We conducted supplementary experiments generating five unwatermarked variants for each prompt and computing average ROUGE scores across all pairs, finding minimal variance (standard deviation 0.02) confirming robust quality maintenance regardless of specific unwatermarked trajectory. Human evaluation studies provide critical validation through subjective quality assessments. We recruited 60 annotators with native English proficiency to perform blind pairwise comparisons between unwatermarked and watermarked text samples. Each annotator evaluated 50 pairs across all three datasets, rating naturalness, coherence, and informativeness on 7-point Likert scales (1=very poor, 7=excellent). Results show no statistically significant differences in any dimension for our reinforcement-guided approach: naturalness scores average 5.7 for watermarked versus 5.8 for baseline ($p=0.18$, paired t-test), coherence 5.6 versus 5.7 ($p=0.22$), and informativeness 5.8 versus 5.8 ($p=0.46$). These findings validate that adaptive placement achieves imperceptibility to human readers. In contrast, fixed placement exhibits significant naturalness degradation (5.1, $p<0.01$) and logit-based watermarking shows marginally significant coherence reduction (5.4, $p=0.04$). Task-specific performance evaluation examines whether watermarked model outputs maintain utility for downstream applications. We evaluate on three benchmark tasks: question answering using SQuAD 2.0 where the watermarked model generates context paragraphs for subsequent QA, summarization on CNN/DailyMail where watermarked summaries undergo evaluation, and sentiment classification on SST-2 where watermarked model continuations are classified. Watermarked generations achieve 97.3% of baseline F1 scores on question answering (F1=81.2 versus 83.5 baseline), demonstrating factual accuracy remains largely intact. Summarization ROUGE-2 scores decline by only 1.8% (0.412 versus 0.420 baseline), while sentiment classification accuracy achieves 93.4% versus 94.1% baseline ($p=0.17$, not significant). These task-specific evaluations confirm that watermarking preserves practical utility across diverse applications. Robustness evaluation subjects watermarked texts to systematic adversarial manipulations. Paraphrasing attacks using back-translation through German successfully remove

watermarks from only 11.2% of sequences, compared to 24.7% for fixed placement and 17.3% for logit-based methods. Synonym substitution attacks corrupt watermarks in 8.7% of cases versus 19.3% for fixed placement, reflecting our policy's learned preference for function words which exhibit greater stability under lexical substitution. Text truncation experiments reveal watermarks remain detectable with 35% random deletion, versus only 22% tolerance for fixed placement. The distributed placement strategy learned through reinforcement learning ensures watermark information spreads across sequences, providing inherent redundancy supporting detection despite substantial missing content.

5. Conclusion

This research establishes reinforcement learning as an effective framework for optimizing watermark placement in neural text generation systems, achieving superior balance between detection robustness and linguistic quality compared to fixed placement strategies. By formulating placement as a sequential decision problem with dual-objective rewards encoding both detectability and fluency requirements, we enable policies to learn context-aware strategies that concentrate watermarks where they minimally impact perceived quality. The proposed transformer-based policy architecture effectively leverages contextual information from language models to identify favorable placement positions and appropriate perturbation strengths, resulting in 97.8% detection accuracy while maintaining perplexity within 1.9% of unwatermarked baselines and ROUGE-L scores of 0.94. The theoretical contributions extend beyond immediate watermarking applications to demonstrate reinforcement learning's broader utility for multi-objective optimization in constrained generation tasks. The training dynamics observed across continuous control benchmarks validate that proximal policy optimization provides the stability necessary for learning effective policies even with sparse delayed rewards characteristic of text generation. The sequence-to-sequence architectural principles adapted for our placement policy highlight the transferability of successful neural text processing paradigms to novel applications requiring sequential decision-making over linguistic structures. Analysis of learned placement strategies reveals interpretable patterns aligned with linguistic intuitions about quality-sensitive positions. The policy's preference for function words over content words, avoidance of sentence-initial positions, and variable strength adaptation demonstrate that reinforcement learning discovers non-obvious strategies that human designers might overlook. These emergent behaviors provide insight into linguistic factors governing quality perception, potentially informing future watermarking approaches and broader research on quality preservation in text transformation tasks. The strong correlation between ROUGE metrics and human quality judgments validates our automated evaluation framework, enabling large-scale quality assessment that would prove infeasible through purely human evaluation. Practical deployment considerations demonstrate the framework's viability for production language model protection. The computational overhead of 8.3% latency remains acceptable for most applications, while memory requirements prove negligible compared to model storage. The framework's robustness to diverse adversarial attacks including paraphrasing, truncation, and synonym substitution provides reasonable security guarantees for intellectual property protection use cases. Organizations deploying proprietary language models can customize the detectability-fluency trade-off through straightforward adjustment of reward function weights, enabling deployment-specific optimization without requiring architectural modifications or retraining from scratch. Several limitations warrant acknowledgment and motivate future investigation. While our framework achieves strong robustness against paraphrasing and truncation attacks, more sophisticated adversarial optimization approaches employing gradient-based watermark removal might pose greater challenges. Developing watermark placement strategies robust to such attacks while maintaining quality

preservation represents an important direction for subsequent work. The current framework focuses on single-model watermarking without addressing scenarios where generated text passes through multiple systems before final output, raising questions about watermark persistence through multi-stage generation pipelines. The generalization of learned policies across different language model architectures demonstrates reasonable robustness, but more extensive evaluation across emerging model families including mixture-of-experts systems and retrieval-augmented generators would strengthen confidence in cross-model applicability. The current evaluation focuses primarily on English text generation, with extensions to multilingual and low-resource languages requiring careful consideration of language-specific quality constraints and detection characteristics. Cross-lingual watermarking where policies function effectively across diverse linguistic contexts represents an important direction for enabling global deployment. Future research directions include exploring joint optimization of placement policies and detection algorithms, where both components adapt synergistically to maximize performance. Investigating watermark placement in constrained generation scenarios such as code synthesis or structured data generation would assess framework applicability beyond natural language. Extending the approach to multimodal language models processing both text and images requires novel architectures handling heterogeneous input modalities. Finally, studying the interaction between watermarking and other generation objectives including factual accuracy, safety constraints, and style control merits investigation for comprehensive optimization addressing multiple competing requirements simultaneously.

References

- [1] Chen, J., & Fan, H. (2025). Beyond Automation in Tax Compliance Through Artificial Intelligence and Professional Judgment. *Frontiers in Business and Finance*, 2(02), 399-418.
- [2] Abdelnabi, S., & Fritz, M. (2021, May). Adversarial watermarking transformer: Towards tracing text provenance with data hiding. In *2021 IEEE Symposium on Security and Privacy (SP)* (pp. 121-140). IEEE.
- [3] Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., & Goldstein, T. (2023, July). A watermark for large language models. In *International Conference on Machine Learning* (pp. 17061-17084). PMLR.
- [4] AlGhozali, S., & Mukminatun, S. (2024). Natural language processing of Gemini artificial intelligence powered Chatbot. *Balankas: An International Multidisciplinary Research Journal*, 1(1), 41-48.
- [5] Zeng, Z., Lin, H., Zhang, S., and Wang, B. (2026). Adaptive Robust Watermarking for Large Language Models via Dynamic Token Embedding Perturbation. *IEEE Access*.
- [6] He, J., Gu, J., Shen, J., & Ranzato, M. A. (2019). Revisiting self-training for neural sequence generation. *arXiv preprint arXiv:1909.13788*.
- [7] Liu, Y., Mondal, A., Chakraborty, A., Zuzak, M., Jacobsen, N., Xing, D., & Srivastava, A. (2020, March). A survey on neural trojans. In *2020 21st International Symposium on Quality Electronic Design (ISQED)* (pp. 33-39). IEEE.
- [8] Khodadadian, S., Chen, Z., & Maguluri, S. T. (2021, July). Finite-sample analysis of off-policy natural actor-critic algorithm. In *International Conference on Machine Learning* (pp. 5420-5431). PMLR.
- [9] Yousuf, H., Lahzi, M., Salloum, S. A., & Shaalan, K. (2021). A systematic review on sequence-to-sequence learning with neural network and its models. *International Journal of Electrical & Computer Engineering* (2088-8708), 11(3).
- [10] Roy, A., Saffar, M., Vaswani, A., & Grangier, D. (2021). Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9, 53-68.
- [11] Tang, G., Kumar, N., Yoo, R., & Michmizos, K. (2021, October). Deep reinforcement learning with population-coded spiking neural network for continuous control. In *Conference on robot learning* (pp. 2016-2029). PMLR.

- [12] Wang, Y., He, H., & Tan, X. (2020, August). Truly proximal policy optimization. In *Uncertainty in artificial intelligence* (pp. 113-122). PMLR.
- [13] Pagnoni, A., Graciarena, M., & Tsvetkov, Y. (2022, October). Threat scenarios and best practices to detect neural fake news. In *Proceedings of the 29th International Conference on Computational Linguistics* (pp. 1233-1249).
- [14] Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., ... & Raffel, C. (2021). Extracting training data from large language models. In *30th USENIX security symposium (USENIX Security 21)* (pp. 2633-2650).
- [15] Jaskowiak, P. A., Costa, I. G., & Campello, R. J. (2022). The area under the ROC curve as a measure of clustering quality. *Data Mining and Knowledge Discovery*, 36(3), 1219-1245.
- [16] Davoodijam, E., & Alambardar Meybodi, M. (2024). Evaluation metrics on text summarization: comprehensive survey. *Knowledge and Information Systems*, 66(12), 7717-7738.
- [17] Choudhary, S., Chatterjee, N., & Saha, S. K. (2022). Interpretation of black box nlp models: A survey. *arXiv preprint arXiv:2203.17081*.
- [18] Xing, S., Wang, Y., & Liu, W. (2025). Self-adapting CPU scheduling for mixed database workloads via hierarchical deep reinforcement learning. *Symmetry*, 17(7), 1109.
- [19] Han, X., Yang, Y., Chen, J., Wang, M., & Zhou, M. (2025). Symmetry-Aware Credit Risk Modeling: A Deep Learning Framework Exploiting Financial Data Balance and Invariance. *Symmetry* (20738994), 17(3).
- [20] Chen, Z., Liu, J., & Chen, J. (2025). Machine Learning Methods for Financial Forecasting in Enterprise Planning: Transitioning from Rule-Based Models to Predictive Analytics. *Frontiers in Artificial Intelligence Research*, 2(3), 541-564.
- [21] Chen, J., Wang, M., & Sun, T. (2025). Intelligent Tax Systems and the Role of Natural Language Processing in Regulatory Interpretation. *American Journal of Machine Learning*, 6(4), 74-94.
- [22] Zeng, Z., & Zhou, M. (2026). ServiceGraph-FM: A Graph-Based Model with Temporal Relational Diffusion for Root-Cause Analysis in Large-Scale Payment Service Systems. *Mathematics*.
- [23] Yang, J. S., Shen, Z., Zeng, Z., & Chen, Z. (2025). Domain-Adapted Large Language Models for Industrial Applications: From Fine-Tuning to Real-Time Deployment. *Computer Science Bulletin*, 8(01), 272-289.
- [24] Lin, H., Liu, J., Zhang, S., & Zeng, Z. (2025). Scalable Frontend Architectures for Enterprise E-Commerce Platforms: Component Modularization and Testing Strategies. *Asian Business Research Journal*, 10(12), 44-56.
- [25] Zhang, S., Qiu, L., & Zhang, H. (2025). Edge cloud synergy models for ultra-low latency data processing in smart city iot networks. *International Journal of Science*, 12(10).
- [26] Qiu, L. (2024). DEEP LEARNING APPROACHES FOR BUILDING ENERGY CONSUMPTION PREDICTION. *Frontiers in Environmental Research*, 2(3), 11-17.
- [27] Liu, J., Wang, J., Chen, H., Guinness, J., Martin, R., & Kulkarni, C. S. (2019). Optimal Level Crossing Predictions for Electronic Prognostics. In *AIAA Scitech 2019 Forum* (p. 1962).
- [28] Zhao, X., Liu, J., Wang, Y., & Wang, J. (2026). CryptoMamba-SSM: Linear Complexity State Space Models for Cryptocurrency Volatility Prediction. *IEEE Open Journal of the Computer Society*.
- [29] Yang, S., Ding, G., Chen, Z., & Yang, J. S. (2025). GART: Graph Neural Network-based Adaptive and Robust Task Scheduler for Heterogeneous Distributed Computing. *IEEE Access*, 13, 200196-200216.
- [30] Sun, T., Yang, J., Li, J., Chen, J., Liu, M., Fan, L., & Wang, X. (2024). Enhancing auto insurance risk evaluation with transformer and SHAP. *IEEE Access*.