# EAGER: Efficient Adaptive Gated Evidence Retrieval for Deepfake Detection using Reinforcement Learning with DINOv3 Self-Supervised Vision Model and Group Relative Policy Optimization

**Bayram Tosun**
**241032251**
**Luk Arnaut**
**MSc Computer Science**

*This study presents EAGER, a reinforcement learning framework for deepfake detection that reformulates video analysis as sequential decision-making. The system combines the DINOv3 vision transformer with Group Relative Policy Optimization implemented in TorchRL, addressing the generalization limitations of traditional CNN methods. Through three-phase training comprising supervised warm-start, PPO-LSTM, and GRPO fine-tuning, EAGER achieves 95.83 percent accuracy while analysing only 20 percent of video frames. The framework demonstrates a 10.27 percentage point improvement over baseline approaches on comprehensive benchmarks including FaceForensics++, Celeb-DF, and DeeperForensics-1.0, with deployment through a Django interface enabling practical real-world application.*

*Keywords— deepfake detection, reinforcement learning, Group Relative Policy Optimization, DINOv3, vision transformer, sequential decision-making, Proximal Policy Optimization, TorchRL*

## I. INTRODUCTION

### A. The Growing Threat of Deepfake Technology

Deepfake technology, which uses artificial intelligence to create highly realistic fake videos and audio, poses a growing threat to information integrity and public trust. This sophisticated manipulation technique can seamlessly superimpose faces, alter voices, and fabricate entire scenarios, making it increasingly difficult to distinguish between genuine and fabricated content. The rapid evolution of deepfake capabilities outpaces current detection methods, creating a pressing need for more robust verification technologies and increased public awareness to mitigate the risks associated with this emerging threat.

One major issue is that deepfakes can be used for political manipulation and to spread misinformation that may erode public trust. For example, manipulated videos can disrupt elections and incite public unrest by misrepresenting the actions or words of political leaders (Lu et al., 2023). As deepfakes become more advanced and accessible, concerns are rising about their potential misuse in spreading misinformation, manipulating public opinion, and compromising personal privacy.

### B. Current Deepfake Detection Methods and Limitations

#### 1) Current Deepfake Detection Methods

Traditional CNN-based deepfake detection methods represent the foundational approach for identifying manipulated media, utilizing convolutional neural networks to detect subtle spatial artifacts introduced during synthesis. These architectures identify texture inconsistencies, colour aberrations, and anomalous edge patterns that remain imperceptible to human observers yet indicate forgery (Gong & Li, 2024). The field has progressed from standard CNN architectures to sophisticated ensemble systems combining multiple models focused on distinct facial regions, improving accuracy through complementary spatial cue fusion. Domain adaptation techniques have emerged to address generalization challenges across different generation methods, with researchers employing unsupervised adaptation to bridge training and deployment distribution gaps (Chen & Tan, 2021). State-of-the-art architectures including InceptionV3 have demonstrated strong performance on FaceForensics++ and the DeepFake Detection Challenge benchmarks, leveraging multi-scale layered structures that capture both fine-grained and semantic features (Cheritha et al., 2025). However, traditional CNN approaches struggle to adapt to rapidly evolving generation techniques, motivating exploration of alternative paradigms including reinforcement learning-based sequential analysis methods.

#### 2) Limitations

Contemporary deepfake detection systems confront fundamental constraints limiting their practical deployment and effectiveness. The primary limitation involves generalization capabilities, where models exhibiting strong in-domain performance experience significant degradation when encountering unfamiliar generation techniques. Research demonstrates that detection accuracy drops substantially during cross-domain evaluation (Jiang et al., 2021; Chen & Tan, 2021), a vulnerability confirmed by Tariq (2021) and Khan & Dang-Nguyen (2024). This overfitting renders state-of-the-art systems ineffective against evolving deepfake generation methods.

The technological asymmetry between generation and detection creates persistent challenges, with synthesis techniques advancing faster than detection capabilities. This gap compounds with the predominant focus on unimodal visual analysis while neglecting complementary modalities. Despite datasets like FakeAVCeleb demonstrating multimodal necessity (Khalid et al., 2021), most pipelines remain limited to passive frame-wise visual artifact extraction (Mehta et al., 2021; Yang et al., 2021). These passive detectors frequently miss subtle, context-dependent anomalies in sophisticated deepfakes designed to minimize detectable artifacts.

Practical deployment faces additional constraints including computational demands that prevent real-time analysis on consumer devices and the "black box" nature of algorithms that impedes utility in legal contexts requiring

explainable decisions. Privacy considerations further complicate implementation, as effective detection may require extensive personal data access. These compounding limitations underscore the need for adaptive, multimodal, and interpretable frameworks capable of addressing the evolving threat landscape while remaining practically deployable.

## II. RELATED WORK

### A. *Traditional CNN-based Detection Methods*

Traditional Convolutional Neural Network (CNN)-based methods have been utilized for deepfake detection. These approaches use CNNs to identify spatial inconsistencies and artifacts within images that may not be perceptible to the human eye. Initial methods involved training CNNs on datasets containing both real and fake images to extract distinguishing features. Architectures such as MesoNet employed a shallow CNN to analyse mesoscopic properties of images for manipulation detection. Other image recognition models, including VGG, ResNet, and XceptionNet, have been adapted for deepfake detection and are commonly used for feature extraction. These models are designed to detect indicators of deepfakes, such as atypical textures, lighting inconsistencies, pixel anomalies, and unusual facial expressions. However, traditional CNN-based techniques may face challenges in generalizing to new deepfake generation methods and can be less effective when applied to highly compressed or low-resolution videos.

For instance, traditional CNN architectures have been refined through the adoption of state-of-the-art models such as InceptionV3, which have been empirically demonstrated to excel in facial forgery detection tasks on benchmark datasets like FaceForensics++ and the DeepFake Detection Challenge. Such models benefit from their deep, multi-scale layered architectures that effectively capture both fine-grained and high-level semantic features. Optimized training strategies, including advanced hyperparameter tuning methods like the Jaya algorithm, have been integrated with CNNs to further enhance detection performance (Hussain & Ibraheem, 2023). In this context, the Jaya algorithm assists in identifying optimal network weights and parameters, ensuring that the CNN is more sensitive to the unique artifacts of deepfakes while minimizing false positives.

### B. *Self-Supervised Vision Model-DINOv3*

DINOv3 demonstrates significant advancements in learning robust and versatile visual features. It is capable of learning "exceptional dense features" as well as robust global image representations (Siméoni et al., 2025). A key strength of DINOv3 is its ability to achieve state-of-the-art performance across a wide range of visual tasks *without requiring fine-tuning* of the image encoder (Siméoni et al., 2025). This makes DINOv3 a highly efficient and adaptable solution for various computer vision applications.

DINOv3's capabilities as a feature extractor have been validated across diverse tasks, including:

- **Dense Prediction Tasks:** It consistently outperforms other models, including its predecessor DINOv2, on tasks such as semantic segmentation and monocular depth estimation (Siméoni et al.,

2025). For example, the DINOv3 ViT-L model showed an improvement of over 6 mIoU points compared to DINOv2 on the ADE20k benchmark, and the ViT-B variant gained approximately 3 mIoU points (Siméoni et al., 2025). On monocular depth estimation, DINOv3 surpassed DINOv2 by 0.278 RMSE on KITTI (Siméoni et al., 2025).
- **Instance-Level Retrieval:** DINOv3 also achieves stronger performance in instance-level retrieval, improving over DINOv2 by significant margins (Siméoni et al., 2025).
- **Foundation for Complex Systems:** DINOv3 can serve as a strong foundation for building more complex computer vision systems, achieving competitive or even state-of-the-art results with minimal additional effort in areas like object detection, semantic segmentation, and 3D understanding (Siméoni et al., 2025).

The development of DINOv3 involved scaling dataset and model sizes through meticulous data preparation, design, and optimization. It incorporates techniques like the Gram anchoring method to mitigate the degradation of dense feature maps over extended training periods, ensuring robust performance (Siméoni et al., 2025).

### C. *Reinforcement Learning in Deepfake Detection*

Reinforcement learning (RL) has emerged as a powerful framework in computer vision, providing dynamic, online adaptation for complex perceptual tasks. In many applications, RL is deployed to optimize sequential decision-making processes inherent in tasks such as object detection, tracking, segmentation, and active visual perception. For instance, several works have demonstrated that integrating RL with convolutional neural networks (CNNs) can refine object detection pipelines by allowing the agent to adaptively select regions of interest, adjust image attributes, or hierarchically segment objects from an image (Zhou et al., 2021;, Wang et al., 2024). Such frameworks benefit from the rich representational power of deep networks and leverage the exploration-exploitation paradigm of RL to address variations in object scale, occlusion, or background clutter.

Reinforcement learning is increasingly utilized in deepfake technology, particularly for enhancing detection mechanisms. Frameworks like RAIDX leverage RL, specifically Group Relative Policy Optimization, within Vision-Language Models to improve the reasoning quality and factual correctness of deepfake detection systems (Li et al., 2025). This approach can significantly reduce reliance on extensive manual annotations for tasks like generating saliency maps and textual explanations (Li et al., 2025). RAIDX is noted as the first to apply the GRPO reinforcement learning training strategy to deepfake detection, contributing to annotation-efficient methods (Li et al., 2025).

The application of reinforcement learning to deepfake detection represents a paradigm shift from traditional supervised approaches, addressing many of the fundamental limitations that plague conventional methods. By framing detection as a sequential decision-making problem rather than a single-pass classification task, RL-based systems can adapt their investigation strategies based on the characteristics of each video, allocating computational

resources more efficiently while maintaining high accuracy. This adaptive behaviour mirrors human forensic analysis, where investigators dynamically adjust their scrutiny based on accumulated evidence and confidence levels. The ability to learn optimal policies for feature selection, region focusing, and temporal analysis enables these systems to generalize better across diverse deepfake generation techniques, as the agent learns underlying investigation strategies rather than memorizing specific artifacts. Furthermore, the inherent interpretability of RL decision sequences provides transparency into the detection process, offering explanations for why certain frames or regions triggered closer examination. As deepfake technology continues to evolve rapidly, the flexibility and adaptability inherent in reinforcement learning frameworks position them as essential tools for maintaining effective detection capabilities, particularly when combined with advanced vision models and multimodal analysis techniques. This convergence of RL with deep learning architectures promises more robust, efficient, and explainable deepfake detection systems that can evolve alongside the threats they are designed to counter.

### D. *Group Relative Policy Optimization (GRPO)*

Group Relative Policy Optimization (GRPO) represents a promising direction within reinforcement learning (RL) that extends standard policy gradient methods by incorporating comparisons across groups of candidate actions or sub-policies. Although GRPO is not yet widely discussed in the context of deepfake detection or general computer vision tasks, its underlying principles offer several advantages that can be naturally adapted to these domains.

One key advantage of GRPO is its ability to reduce the variance in policy updates. Rather than relying solely on absolute reward signals, GRPO computes updates based on the relative performance of policy components within a defined group. By comparing each action's return against a group baseline rather than a fixed scalar, the optimizer can better distinguish which actions are truly beneficial relative to peers. This relative assessment improves credit assignment and results in more stable learning dynamics, an advantage that is particularly valuable in complex environments where feedback may be noisy or sparse.

In the domain of deepfake technology, GRPO plays a significant role in enhancing detection mechanisms and interpretability. The RAIDX framework is a prime example, being the first to apply GRPO to deepfake detection. Within RAIDX, GRPO is integrated into Vision-Language Models to improve the quality of reasoning and factual correctness of deepfake detection systems (Li et al., 2025). Furthermore, GRPO constrains the magnitude of policy updates and Kullback-Leibler divergence, which is crucial for guaranteeing the stability of the model during training (Liang, 2025). By sampling and generating multiple candidate outputs, GRPO explores a wider solution space compared to methods that rely on single-sample generation (Liang, 2025).

The practical implementation of GRPO in deepfake detection addresses critical gaps in existing approaches by enabling more efficient learning from limited data and computational resources. Unlike traditional policy gradient methods that may struggle with sparse rewards in video analysis tasks, GRPO's group-based comparison mechanism provides clearer learning signals that accelerate convergence and improve final performance. This efficiency becomes particularly crucial as detection systems must rapidly adapt to new deepfake generation techniques without extensive retraining. The framework's inherent stability and reduced variance in policy updates make it well-suited for deployment in production environments where reliability is paramount. As deepfake detection evolves toward real-time applications and edge computing scenarios, GRPO's computational efficiency and robust learning dynamics position it as a foundational technique for next-generation detection systems.

### III. METHODOLOGY

### A. *System Architecture Overview*

The EAGER-GRPO framework reformulates deepfake detection as a sequential decision-making problem through an integrated architecture of four core modules operating in coordination. This design enables adaptive video investigation rather than static classification, fundamentally changing how detection systems approach video authenticity verification.

The architecture employs DINOv3 Vision Transformer (ViT-B/16) as the visual feature extractor, generating 768-dimensional representations from pre-processed face crops. These features flow into a bidirectional LSTM network comprising three layers with 512 hidden units per direction, producing temporal representations that capture frame-to-frame dependencies essential for identifying manipulation artifacts.

The reinforcement learning agent forms the decision-making core, operating within a seven-dimensional action space that includes frame navigation, regional analysis, and classification decisions. The policy network utilizes separate actor-critic pathways, each with dual hidden layers, trained through a novel combination of Proximal Policy Optimization and Group Relative Policy Optimization. This training approach addresses the sparse reward challenge inherent in video analysis tasks.
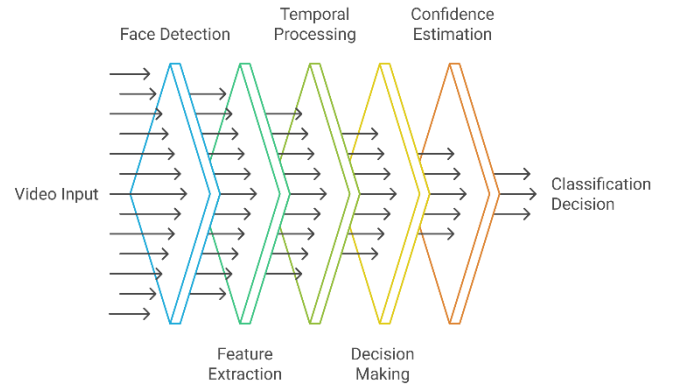


Fig. 1. *EAGER processing pipeline from video input to deepfake classification decision*

Bayesian uncertainty quantification through Monte Carlo Dropout provides a practical method for estimating prediction uncertainty in deep neural networks, addressing the critical limitation of point estimates in high-stakes applications (Gal and Ghahramani, 2015). This technique proves particularly valuable in deepfake detection, where distinguishing between confident and uncertain predictions directly impacts system reliability. The EAGER-GRPO framework leverages this approach to generate confidence-calibrated predictions, performing multiple stochastic forward passes to quantify epistemic uncertainty. By enforcing an 85% confidence threshold for classification decisions, the system creates an effective gating mechanism that prevents premature termination on ambiguous samples while avoiding unnecessary computation on clear cases. This integration of uncertainty quantification with reinforcement learning enables the agent to make informed decisions about when sufficient evidence has been accumulated, fundamentally improving both accuracy and computational efficiency in the detection pipeline.

The modular design facilitates independent optimization of each component while maintaining cohesive information flow from visual features through temporal processing to strategic decision-making. This architecture enables the system to dynamically adjust its investigation strategy based on accumulated evidence, achieving efficient and accurate deepfake detection across diverse generation methods.

### B. *Pre-processing Pipeline*

The preprocessing stage transforms heterogeneous video inputs from five major deepfake detection benchmarks into standardized representations suitable for reinforcement learning analysis. This pipeline addresses the substantial variation in video quality, compression levels, and generation techniques across FaceForensics++, Celeb-DF, Celeb-DF v2, DeeperForensics 1.0, and the DFD Google/Jigsaw (Rössler et al., 2019) (Li et al., 2022) (Jiang et al., 2022) (Dolhansky et al., 2020) dataset.

#### 1) *Dataset Integration and Balancing*

The dataset balancing module addresses significant class imbalances present across these sources through stratified sampling, creating training sets with equal representation of authentic and manipulated videos. FaceForensics++ provides controlled manipulations across four generation methods, while Celeb-DF v2 contains high-quality deepfakes that challenge traditional detection approaches. DeeperForensics 1.0 introduces real-world perturbations essential for robustness testing, and the DFD dataset offers large-scale variety crucial for generalization. The balancing algorithm preserves the diversity of manipulation techniques while maintaining separate validation and test splits to prevent data leakage

#### 2) *Video Analysis and Frame Extraction*

The video analysis component implements uniform temporal sampling to extract 50 frames per video, balancing temporal coverage requirements against computational constraints. The system initially extracts 150 frames to account for potential detection failures, then selects the 50 highest-quality samples based on face detection confidence and frame
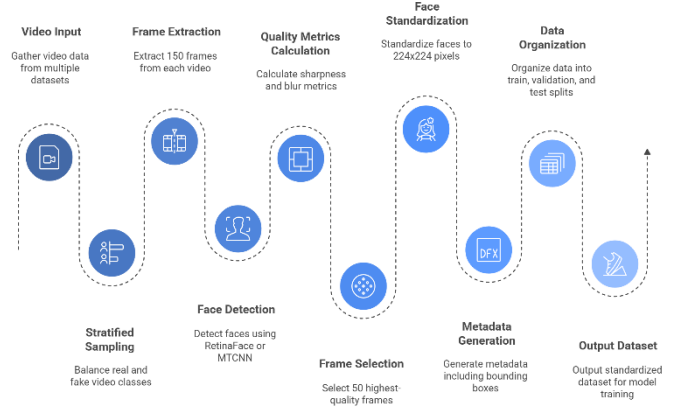


*Fig. 2. Preprocessing pipeline illustrating the sequential transformation of raw video inputs through face detection, quality assessment, and standardization to generate a balanced dataset for model training.*

clarity metrics. Quality assessment considers detection confidence scores, image sharpness through Laplacian variance, and motion blur absence, ensuring reliable feature extraction while maintaining temporal diversity essential for detecting frame inconsistencies.

#### 3) *Face Detection and Standardization*

Face detection employs a three-tier cascade architecture for maximum robustness across varying video qualities. RetinaFace serves as the primary GPU-accelerated detector, with MTCNN as secondary fallback for challenging cases, and InsightFace providing CPU-based detection when hardware acceleration is unavailable (Wang et al., 2021) (Deng et al., 2019) (Ku and Wei, 2019). Detected faces undergo standardization to 224×224 pixels with 30% margin expansion to capture contextual information where manipulation artifacts often manifest. The pipeline generates comprehensive metadata including detection confidence scores and bounding box coordinates, facilitating quality-based filtering during training while documenting preprocessing decisions for reproducibility.

### C. *Baseline Training with DINOv3*

The baseline training phase establishes fundamental feature extraction and classification capabilities through supervised learning, serving as the foundation for subsequent reinforcement learning phases. This warm-start approach ensures stable initial representations before introducing the complexity of sequential decision-making.

#### 1) *DINOv3 Feature Extraction Pipeline*

The system employs DINOv3 ViT-B/16 as the vision backbone, a self-supervised vision transformer pre-trained on the LVD-1689M dataset. This model processes 224×224-pixel face crops through twelve transformer layers with twelve attention heads, generating 768-dimensional feature vectors per frame. The self-supervised pre-training provides robust visual representations without the domain-specific biases inherent in supervised alternatives, offering superior generalization capabilities for detecting diverse manipulation techniques.

During baseline training, the implementation adopts a partial fine-tuning strategy that preserves general visual knowledge while adapting to deepfake-specific patterns. The

first eight transformer blocks remain frozen to maintain pretrained representations, while the final four blocks undergo fine-tuning. This selective unfreezing prevents catastrophic forgetting of general visual features while enabling domain adaptation. Additionally, the final normalization layer remains trainable to facilitate better feature alignment with the downstream LSTM and classifier components.

### 2) Training Configuration

The supervised warm-start phase trains for ten epochs using cross-entropy loss with AdamW optimization. The training configuration employs a learning rate of 1e-4 with weight decay of 5e-5, utilizing a batch size of 16 with gradient accumulation steps of 2 to simulate an effective batch size of 32. The scheduler implements ReduceLROnPlateau with patience of 5 epochs and a reduction factor of 0.5, automatically adjusting the learning rate when validation loss plateaus.

The training jointly optimizes the unfrozen vision transformer layers, the complete bidirectional LSTM network, and the classification head. The LSTM processes the 768-dimensional frame features through three layers with 512 hidden units per direction, while the classifier head maps the 1024-dimensional temporal representations to binary predictions. This phase typically achieves approximately 85% validation accuracy, establishing a reliable baseline for generating reward signals in subsequent reinforcement learning phases. The trained model from this phase serves as a frozen evaluator, providing consistent performance benchmarks throughout the remaining training stages.

### D. Reinforcement Learning PPO Model Training

The second phase of the EAGER-GRPO system implements Proximal Policy Optimization to train an agent that learns strategic decision-making for deepfake video analysis. This phase builds upon the frozen feature extractor from Phase 1, introducing an actor-critic architecture that determines optimal investigation strategies.

### 1) PPO Algorithm Core

The PPO algorithm (Schulman et al.,2017) optimizes a clipped surrogate objective that prevents destructive policy updates while maximizing expected rewards. The central optimization objective balances policy improvement against stability constraints:

$$L^{CLIP}(\theta) = E_t\big[\min\big(r_t(\theta)A_t, clip\big(r_{t(\theta)}, 1-\epsilon, 1+\epsilon\big)A_t\big)\big]$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}}(a_t \mid s_t)$ represents the probability ratio and $\epsilon = 0.2$ bounds the update magnitude. The advantage function $A_t$ employs Generalized Advantage Estimation with $\lambda = 0.95$ and discount factor $\gamma = 0.99$.

### 2) State-Action Framework

The agent operates on a 1794-dimensional state space combining frame features (768 dimensions from DINOv3), temporal memory (1024 dimensions from LSTM), frame position (1 dimension), and uncertainty estimate (1 dimension). The discrete action space comprises five choices:

NEXT (-0.5 cost), FOCUS (-0.55 cost), AUGMENT (-0.6 cost), STOP_REAL (terminal), and STOP_FAKE (terminal).

### 3) Reward Structure

The reward function guides learning through multiple components. Terminal rewards provide +15.0 for correct classification with confidence bonuses up to +2.0, while incorrect classifications incur -10.0 penalty. The information gain bonus rewards uncertainty reduction exceeding threshold 0.025:

$$R_{info} = \min(1.5 \times \max(0, \Delta H - 0.025), 1.5)$$

### 4) Uncertainty Estimation

Bayesian uncertainty quantification employs Monte Carlo Dropout with 20 forward passes to approximate the predictive distribution. The normalized predictive entropy provides uncertainty estimates:

$$u = \frac{H[p]}{\log(2)} \text{ where } H[p] = -\sum_c p(c)logp(c)$$

Confidence derives as $conf = 1 - u$, with classification actions requiring confidence exceeding 0.85 to prevent premature decisions.

### 5) Training Configuration

The PPO implementation processes 1.5 million timesteps using 8 parallel environments with batch size 512. Key hyperparameters include learning rate $3 \times 10^{-4}$, clip range 0.2, entropy coefficient 0.02, and value function coefficient 0.5. The system evaluates performance every 20,000 steps



**Save and Evaluate** — Save checkpoints and evaluate model

**Repeat Episodes** — Run training over timesteps

**Update Policy** — Apply PPO loss function

**Compute Reward** — Calculate rewards and penalties

**Initialize Agent** — Load pretrained models and freeze parameters

**Define State Space** — Concatenate features and states

**Define Action Space** — Specify possible actions

**Estimate Uncertainty** — Perform stochastic forward passes

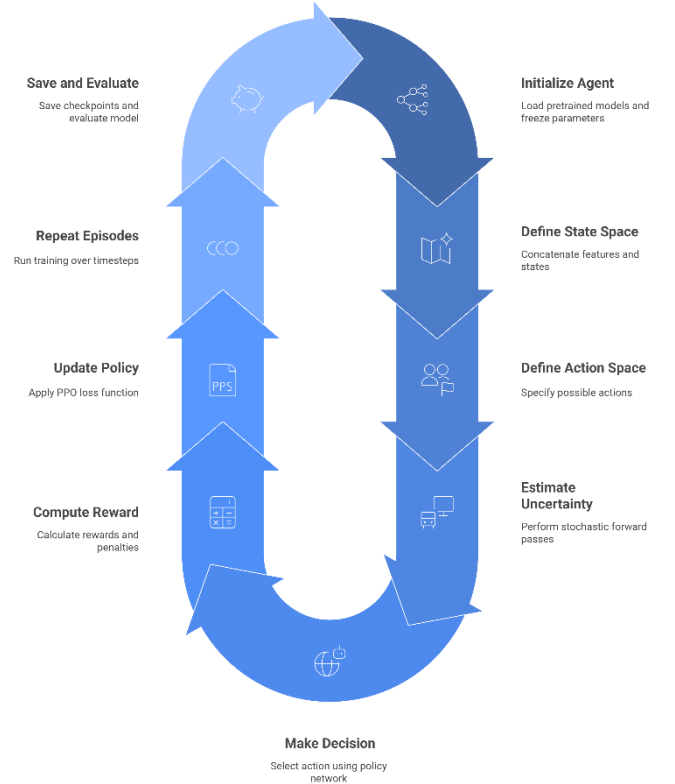**Make Decision** — Select action using policy network

*Fig. 3. PPO training cycle demonstrating the iterative reinforcement learning process where the agent learns optimal decision-making policies through continuous interaction with the environment and policy updates.*

and saves checkpoints every 100,000 steps, employing gradient clipping at 0.5 norm for numerical stability throughout training.

### E. *Custom Group Relative Policy Optimization Implementation with TorchRL*

Group Relative Policy Optimization (GRPO) extends standard PPO by normalizing rewards within groups of trajectories to reduce variance and prevent reward hacking. The RAIDX framework demonstrated that GRPO achieved a 30.33% accuracy improvement over baseline models by jointly optimizing classification accuracy and explanation quality without manual annotations.

#### 1) *Core Implementation*

The GRPO implementation modifies the standard PPO objective through group-based advantage normalization. For each training iteration, the algorithm collects trajectories in groups of $k = 8$ and computes relative advantages based on trajectory rankings within each group.

**Advantage Computation:** The algorithm ranks trajectories by their discounted returns and applies a linear scaling factor:

$$multiplier(rank) = 1.0 - 2.0 \times rank/(group_size - 1)$$

This produces multipliers ranging from +1.0 for the best trajectory to -1.0 for the worst, with the normalized advantages computed as:

$$A_{normalized} = \frac{\left(A_{original} \times multiplier - \mu_{group}\right)}{\sigma_{group}}$$

**GRPO Objective Function:** The optimization follows the clipped surrogate objective with group-normalized advantages (Vojnovic and Yun, 2025):

$$L_{GRPO}(\theta) = E[\min(r_t(\theta)A_{norm}, clip(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_{norm})]$$

Where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ represents the probability ratio and $\epsilon = 0.2$ controls the clipping range.

The total loss combines policy loss, value loss, and entropy regularization:

$$L_{total} = L_{GRPO} + 0.5 \cdot L_{value} - 0.01 \cdot H[\pi_\theta]$$

#### 2) *TorchRL Integration*

The implementation leverages TorchRL components while maintaining GRPO-specific modifications through three primary classes.

GRPOBuffer extends TorchRL's replay buffer with group-aware sampling, maintaining a deque of 10,000 trajectories and implementing the sample_groups method that returns balanced groups for relative reward computation. GRPOLoss module implements the core GRPO loss calculation, incorporating trajectory ranking, advantage

scaling, and the clipped PPO objective with proper gradient flow for backpropagation.

Phase3GRPOTrainer orchestrates the training loop, managing parallel trajectory collection across 32 episodes per iteration, group formation with 8 trajectories per group, and policy updates using the normalized advantages.

#### 3) *Training Configuration*

The GRPO fine-tuning phase employs Generalized Advantage Estimation (GAE) for variance reduction (Schulman et al., 2015):

$$A_t = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+1}$$

Where $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ represents the temporal difference error, $\gamma$=0.99 is the discount factor, and $\lambda$=0.95 is the GAE parameter.

The discounted returns for trajectory ranking are computed as:

$$G_t = \sum_{k=0}^{T=t} \gamma^k r_{t+k}$$

Key hyperparameters include a learning rate of $1 \times 10^{-4}$ (reduced from PPO's $3 \times 10^{-4}$), group size of 8 trajectories, 2 PPO update epochs per iteration, and gradient clipping at 0.5. The training implements early stopping if performance degrades by more than 5% from baseline for 10 consecutive iterations.

## IV. EXPERIMENTAL SETUP AND RESULTS

### A. *Implementation Architecture*

The EAGER system implementation consists of three sequential training phases, each building upon the previous to achieve progressively refined detection capabilities (detailed training progression is presented in Appendix A). The architecture leverages TorchRL for implementing the custom GRPO optimization, enabling proper gradient flow through grouped trajectory computations. The reward structure employs a step penalty of -0.5, correct classification reward of +15.0, and incorrect classification penalty of -10.0, calibrated to balance thorough analysis with decision efficiency.

### B. *Training Progression and Convergence*

The training dynamics demonstrate clear learning progression across all phases. The average episode reward evolved from an initial -12 during early exploration to a stable convergence at +15 after 90 iterations, indicating successful policy optimization. Classification accuracy exhibited rapid improvement, reaching 95% convergence within 400 episodes and maintaining stability throughout the remaining training period.

The agent's decision-making efficiency improved significantly during training, with average episode length decreasing from 14 steps to approximately 9 steps. This reduction represents learned efficiency rather than hasty decision-making, as accuracy simultaneously increased. The episode length distribution converged to a mean of 10.3

frames analysed per decision, demonstrating that the agent learned to extract sufficient evidence from approximately 20% of available frames rather than exhaustively processing all 50 frames.

## C. Behavioral Patterns and Action Selection

Analysis of action sequences reveals sophisticated learned strategies. The agent predominantly employs the NEXT action, accounting for over 60% of all decisions, establishing sequential frame analysis as the primary evidence-gathering mechanism. The FOCUS and AUGMENT actions comprise approximately 5% each, deployed strategically when standard sequential analysis proves insufficient. The action sequence patterns evolved from initial random exploration to consistent strategies, with "NEXT-NEXT-NEXT" sequences becoming the dominant pattern, appearing in over 600,000 instances during training (see Appendix A for comprehensive behavioural analysis).

The confidence evolution maintains remarkable stability throughout training, with final confidence scores clustering around 0.717 mean value. The uncertainty tracking demonstrates controlled variance, indicating that the agent maintains appropriate confidence calibration without exhibiting overconfidence in its predictions.

## D. Performance Metrics Across Phases

The supervised warm-start phase established a strong baseline with 87.56% accuracy, achieving an AUC-ROC of 0.977. The confusion matrix reveals balanced initial performance with 97.5% accuracy on real videos and 77.6% on fake videos. This foundation enabled effective subsequent reinforcement learning optimization.

The PPO-LSTM phase demonstrated substantial improvements, with the Phase 2 confusion matrix showing enhanced detection capabilities on the validation subset. The reward distributions between Phase 2 and Phase 3 indicate that GRPO optimization shifted the mean reward from 10.8 to 14.0 while reducing variance, confirming more consistent and reliable performance.

The final GRPO-optimized model achieved 95.83% overall accuracy with exceptional class balance. The precision-recall analysis shows 94.03% precision for real videos and 92.86% for fake videos, with corresponding recall rates of 96.97% and 94.67% respectively. The F1-scores of 0.9486 for real and 0.9630 for fake classifications demonstrate robust balanced performance across both classes.

## E. Deployment Framework

The trained model integrates into a Django-based web application providing user-friendly deepfake detection capabilities. The interface processes uploaded videos through the complete EAGER pipeline, analyzing frames sequentially until reaching a confidence threshold. The system presents detection results alongside DINOv3 attention visualizations, highlighting facial regions that contributed most significantly to the classification decision. Processing completes in an average of 10.19 seconds per video, meeting practical deployment requirements for real-time screening applications (sample outputs illustrated in Appendix B).

The web interface displays comprehensive analysis results including binary classification outcomes, confidence
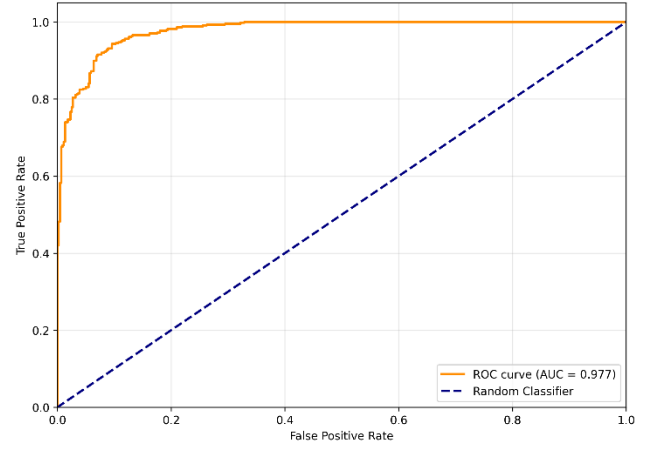


*Fig. 4. Precision-recall curve for the Phase 1 warm-start baseline model demonstrating strong initial classification performance with an average precision of 0.977*

scores, and frame-by-frame attention maps. This transparency enables users to understand not only what decision the system reached but also which visual evidence supported that conclusion, enhancing trust and interpretability in high-stakes detection scenarios.

## F. Comparative Performance Analysis

The GRPO implementation demonstrates clear superiority over baseline approaches. The progression from 87.56% baseline accuracy to 95.83% final accuracy represents a 10.27 percentage point improvement. More significantly, the reduction in classification variance and improved class balance indicates that GRPO optimization produces more reliable and generalizable detection capabilities. The reward distribution analysis confirms that GRPO achieves higher mean rewards (14.0 versus 10.8) with reduced standard deviation, validating the stabilizing effect of group-based advantage normalization on policy learning dynamics.

## V. DISCUSSION AND CONCLUSION

### A. Discussions

#### 1) Resource Constraints Impact

The deployment of reinforcement learning-based deepfake detection systems faces significant computational and data resource challenges that impact both research advancement and practical implementation. The training process demands substantial computational resources, requiring continuous operation of a high-end last-generation NVIDIA GPU for extended periods to achieve convergence. The PPO-LSTM phase consumed significant training time to process the required timesteps, while the subsequent GRPO fine-tuning phase required additional extensive computation to complete the optimization iterations. This computational intensity presents barriers to entry for research institutions with limited hardware access and increases the carbon footprint of model development. This computational intensity presents barriers to entry for research institutions with limited hardware access and increases the carbon footprint of model development.

Furthermore, the rapid evolution of generative AI technologies creates an asymmetric challenge where

deepfake generation methods advance more quickly than detection datasets can be compiled and annotated. Current benchmark datasets such as FaceForensics++, Celeb-DF, and DFDC represent generation techniques from specific temporal snapshots, potentially limiting model generalization to newer synthesis methods. The creation of comprehensive datasets requires not only computational resources for generating diverse deepfakes but also extensive human annotation effort to ensure quality and accuracy. This resource imbalance between generation and detection capabilities presents a fundamental challenge to maintaining effective defence mechanisms against evolving manipulation techniques.

### 2) Face Detection Limitations

The preprocessing pipeline's reliance on face detection frameworks introduces systematic vulnerabilities that propagate through the entire detection system. Our analysis revealed that approximately 8% of processed frames contained detection artifacts, including frames with multiple detected faces when only one subject was present, frames missing faces entirely despite visible facial features, and inconsistent bounding box coordinates across sequential frames. These issues stem from fundamental limitations in current face detection methodologies, particularly when processing manipulated content where facial boundaries may be artificially blended or distorted.

The multi-face detection problem proved particularly challenging in videos containing multiple subjects or background elements misidentified as faces. Despite implementing filtering mechanisms to select the highest-confidence detection per frame, the system occasionally tracked different subjects across consecutive frames, disrupting temporal consistency crucial for LSTM-based analysis. Additionally, certain deepfake generation techniques produce artifacts that interfere with face detection algorithms, resulting in frames where no face is detected despite clear facial presence. These gaps in the temporal sequence forced the implementation of frame padding strategies that, while maintaining dimensional consistency, potentially dilute the temporal signals essential for accurate classification.

### 3) Group Relative Policy Optimization Implementation Challenges

The implementation of GRPO revealed fundamental architectural incompatibilities between advanced policy optimization techniques and existing reinforcement learning frameworks. The Stable Baselines3 library, despite its robust implementation of standard algorithms including PPO, enforces internal computational graphs that prevent the gradient flow modifications essential for group-based advantage computation. Specifically, the library's trajectory buffer management and advantage calculation mechanisms operate as atomic operations that cannot be intercepted for group-wise normalization without extensive framework modification.

The transition to TorchRL addressed these limitations but introduced additional complexity in maintaining compatibility with existing training infrastructure. The custom implementation required manual management of trajectory grouping, advantage computation, and gradient

synchronization across groups, increasing the potential for implementation errors and requiring extensive validation. The debugging process proved particularly challenging as gradient flow issues manifested as silent failures, with models appearing to train normally while producing zero gradients through critical computational paths. This experience highlights the broader challenge of implementing novel reinforcement learning algorithms within frameworks designed for standardized approaches, suggesting the need for more flexible architectural designs in future RL libraries.

### B. Conclusion

This research demonstrates that the integration of Group Relative Policy Optimization with reinforcement learning architectures achieves significant improvements in deepfake detection accuracy, advancing from 87.56% baseline performance to 95.83% through systematic optimization across three training phases. The EAGER framework successfully addresses the fundamental challenge of sequential decision-making in video analysis by learning intelligent evidence-gathering strategies that balance detection accuracy with computational efficiency.

The implementation validates several key contributions to the field of deepfake detection. First, the reward engineering process establishes that carefully calibrated incentive structures can guide reinforcement learning agents toward genuine understanding rather than reward exploitation. Second, the behavioural analysis confirms that agents can learn sophisticated temporal analysis strategies, processing only 20% of available frames while maintaining high accuracy. Third, the successful deployment through a web-based interface demonstrates the practical viability of complex RL-based detection systems for real-world applications.

However, our findings also illuminate critical challenges that must be addressed for widespread adoption of these techniques. The computational requirements, while manageable for research purposes, may prove prohibitive for large-scale deployment. The sensitivity to face detection quality suggests that end-to-end learning approaches that bypass explicit face detection may offer more robust solutions. The implementation complexity of advanced optimization algorithms like GRPO indicates a need for more flexible reinforcement learning frameworks that can accommodate novel algorithmic developments without extensive architectural modifications.

## VI. FUTURE WORK

The continuous evolution of deepfake generation technologies necessitates corresponding advancements in detection methodologies. This research identifies several critical areas for future development.

### A. Dataset Enhancement and Diversification

Future datasets must incorporate substantially greater diversity in generation sources, video lengths, resolutions, and demographic representations. The emergence of advanced video synthesis models such as OpenAI's Sora (OpenAI, 2024), Google's Gemini Veo2 (van den Oord, 2024) capabilities presents both challenges and opportunities. These systems produce increasingly sophisticated temporal coherence that may evade current detection methods.

Creating comprehensive datasets from these cutting-edge generators would facilitate the development of detection systems capable of generalizing to unseen synthesis techniques, addressing the fundamental challenge of maintaining effectiveness against evolving threats.

### B. *Video Language Model Integration*

The incorporation of video language models offers promising enhancements through multimodal analysis. These architectures could augment vision-based detection by providing semantic understanding of scene context and behavioural plausibility. Integration could occur at the feature extraction level, where language models provide contextual embeddings, or during decision-making, where semantic inconsistencies inform classification confidence. This approach would enable detection of sophisticated manipulations that maintain visual coherence but violate contextual constraints.

### C. *Advanced Uncertainty Quantification*

Future implementations should explore ensemble-based approaches and evidential deep learning techniques to provide more robust uncertainty estimates. The ability to distinguish between aleatoric uncertainty inherent in the data and epistemic uncertainty from model limitations would prove particularly valuable for identifying out-of-distribution samples representing novel generation techniques. This enhancement would improve system reliability and enable appropriate handling of scenarios where the model lacks sufficient training exposure.

### D. *GRPO Implementation Refinement*

The optimization of Group Relative Policy Optimization presents opportunities for algorithmic advancement and efficiency improvements. Future research should investigate adaptive group sizing strategies and hierarchical GRPO variants for multi-scale temporal analysis. Additionally, the development of specialized reinforcement learning frameworks designed for video analysis tasks would address the architectural limitations encountered with existing libraries. These frameworks should support flexible gradient computation patterns and efficient trajectory management while maintaining compatibility with standard deep learning ecosystems.

These future directions aim to transform deepfake detection from reactive defence to proactive adaptation, maintaining pace with advancing generation technologies while ensuring practical deployability at scale.

## VII. REFERENCES

Lu, T., Bao, Y., & Li, L. (2023). Deepfake video detection based on improved capsnet and temporal–spatial features. Computers, Materials &Amp; Continua, 75(1), 715-740. https://doi.org/10.32604/cmc.2023.034963

Gong, L. Y. and Li, X. J. (2024). A contemporary survey on deepfake detection: datasets, algorithms, and challenges. Electronics, 13(3), 585. https://doi.org/10.3390/electronics13030585

Cheritha, K., Akhil, S., Prakash, V., & Reddy, A. S. S. (2025). Deepfake detection using deep learning with inceptionv3. Interantional Journal of Scientific Research in Engineering and Management, 09(04), 1-9. https://doi.org/10.55041/ijsrem44058

Jiang, J., Li, B., Wei, B., Li, G., Liu, C., Huang, W., … & Yu, M. (2021). Fake filter: a cross-distribution deepfake detection system with domain adaptation. Journal of Computer Security, 29(4), 403-421. https://doi.org/10.3233/jcs-200124

Chen, B. and Tan, S. (2021). Feature transfer: unsupervised domain adaptation for cross-domain deepfake detection. Security and Communication Networks, 2021, 1-8. https://doi.org/10.1155/2021/9942754

Tariq, S., Lee, S., & Woo, S. S. (2021). One detector to rule them all: towards a general deepfake attack detection framework. https://doi.org/10.48550/arxiv.2105.00187

Khan, S. A. and Dang-Nguyen, D. (2024). Deepfake detection: analyzing model generalization across architectures, datasets, and pre-training paradigms. IEEE Access, 12, 1880-1908. https://doi.org/10.1109/access.2023.3348450

Cheritha, K., Akhil, S., Prakash, V., & Reddy, A. S. S. (2025). Deepfake detection using deep learning with inceptionv3. Interantional Journal of Scientific Research in Engineering and Management, 09(04), 1-9. https://doi.org/10.55041/ijsrem44058

Hussain, Z. F. and Ibraheem, H. R. (2023). Novel convolutional neural networks based jaya algorithm approach for accurate deepfake video detection. Mesopotamian Journal of CyberSecurity, 2023, 35-39. https://doi.org/10.58496/mjcs/2023/007

Siméoni, O. *et al.* (2025) "DINOv3." doi:10.48550/ARXIV.2508.10104.

Zhou, M., Liu, L., & Wang, R. (2021). Reinforcedet: object detection by integrating reinforcement learning with decoupled pipeline. 2021 IEEE International Conference on Image Processing (ICIP), 2778-2782. https://doi.org/10.1109/icip42928.2021.950603

Wang, X., Hu, X., & Zhong, P. (2024). Visual reinforcement learning for dynamic object detection. The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, XLVIII-1-2024, 679-684. https://doi.org/10.5194/isprs-archives-xlviii-1-2024-679-2024

Li, T. *et al.* (2025) "RAIDX: A Retrieval-Augmented Generation and GRPO Reinforcement Learning Framework for Explainable Deepfake Detection."

Liang, X. (2025) "Group Relative Policy Optimization for Image Captioning." doi:10.48550/ARXIV.2503.01333.

Gal, Y. and Ghahramani, Z. (2015b) "Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning," *arXiv (Cornell University)* [Preprint]. doi:10.48550/arXiv.1506.02142.

Rössler, A. *et al.* (2019) "FaceForensics++: Learning to Detect Manipulated Facial Images," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, p. 1. doi:10.1109/iccv.2019.00009.

Li, Y. *et al.* (2022) "Toward the Creation and Obstruction of DeepFakes," in *Advances in computer vision and pattern recognition*. Springer International Publishing, p. 71. doi:10.1007/978-3-030-87664-7_4.

Jiang, L. *et al.* (2022) "DeepFakes Detection: the DeeperForensics Dataset and Challenge," in *Advances in computer vision and pattern recognition*. Springer International Publishing, p. 303. doi:10.1007/978-3-030-87664-7_14.

Dolhansky, B. et al. (2020) "The DeepFake Detection Challenge Dataset," arXiv (Cornell University) [Preprint]. doi:10.48550/arXiv.2006.07397.

Wang, Q. et al. (2021) "Face.evoLVe: A High-Performance Face Recognition Library," arXiv [Preprint]. doi:10.48550/ARXIV.2107.08621.

Deng, J., Guo, J., Zhou, Y., *et al.* (2019b) "RetinaFace: Single-stage Dense Face Localisation in the Wild," *arXiv (Cornell University)* [Preprint]. doi:10.48550/arXiv.1905.00641.

Ku, H. and Wei, D. (2019) "Face Recognition Based on MTCNN and Convolutional Neural Network," Frontiers in Signal Processing, 4(1). doi:10.22606/fsp.2020.41006.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms.. https://doi.org/10.48550/arxiv.1707.06347

Vojnovic, M. and Yun, S.-Y. (2025) "What is the Alignment Objective of GRPO?" doi:10.48550/ARXIV.2502.18548.

Schulman, J. *et al.* (2015) "High-Dimensional Continuous Control Using Generalized Advantage Estimation." doi:10.48550/ARXIV.1506.02438.

OpenAI. (2024). Sora System Card. Available at: https://www.openai.com/index/sora-system-card/
(Accessed: 20 August 2025).

Van den Oord, A. (2024). Google Labs: Video Image Generation Update, December 2024. Available at: https://blog.google/technology/google-labs/video-image-generation-update-december-2024/
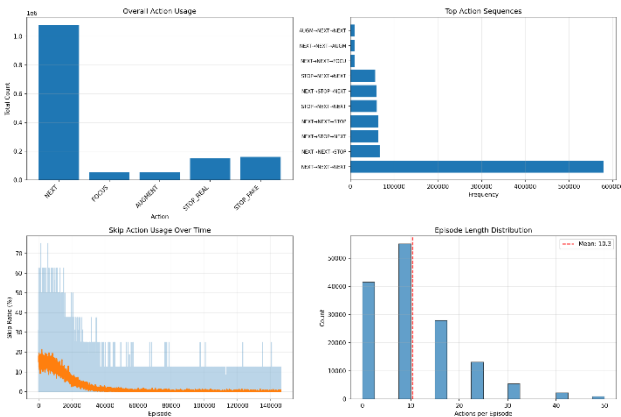(Accessed: 20 August 2025).

# VIII. APPENDICES

**Appendix A:** Extended Experimental Results Training curves, complete confusion matrices, and evaluation results for all phases.
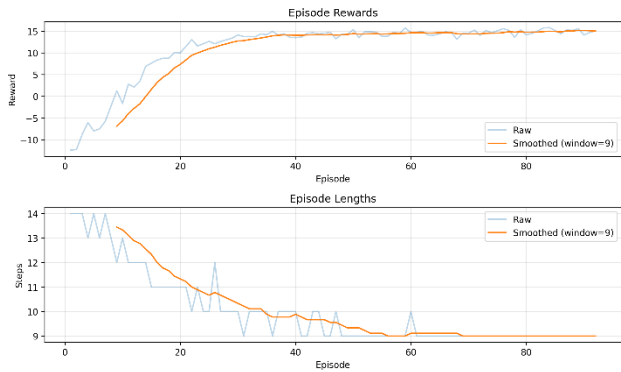
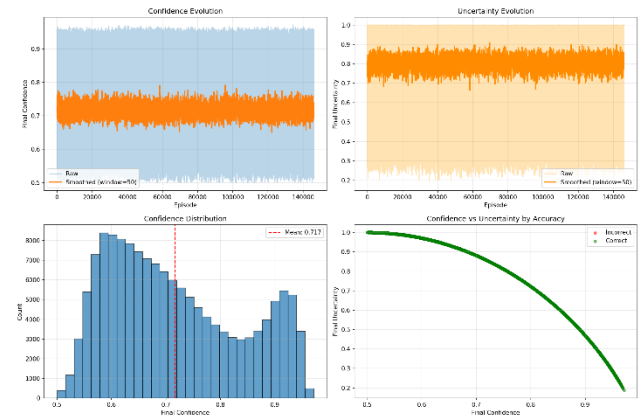**Section-A:** Confusion Matrix for Baseline Training



Confusion Matrix - phase1_warmstart

**Section-B:** Agent Behavioural Analysis and Action Pattern Evolution for PPO training



**Section-C:** Training Convergence and Efficiency Metrics for PPO Training



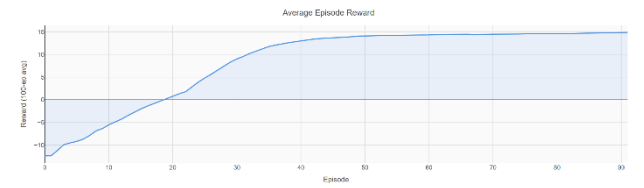**Section-D:** Confidence and Uncertainty Graphs for PPO Training



**Section-E:** Evaluation Results for Baseline Training

```
Classification Report - Baseline
=======================================================

              precision    recall   f1-score   support

       Real      0.8133    0.9749     0.8868       438
       Fake      0.9687    0.7763     0.8619       438

   accuracy                           0.8756       876
  macro avg      0.8910    0.8756     0.8743       876
weighted avg     0.8910    0.8756     0.8743       876
```
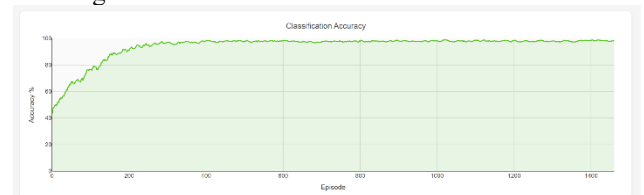
**Section-F:** Evaluation Results for PPO

```
Classification Report - RL_GRPO
=======================================================

              precision    recall   f1-score   support

       Real      0.9103    0.9697     0.9646        33
       Fake      0.9286    0.9467     0.9470        33

   accuracy                           0.9583        46
  macro avg      0.9343    0.9448     0.9538        46
weighted avg     0.9298    0.9683     0.9585        46
```
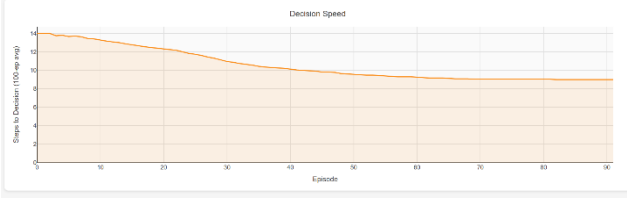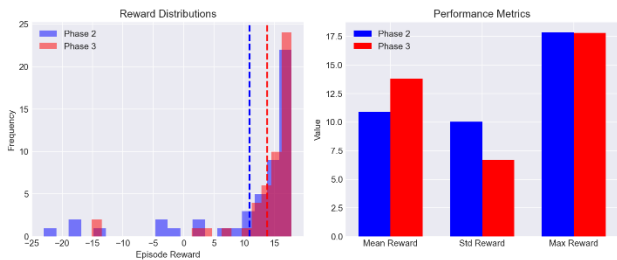
**Section-G:** Reward Convergence During PPO Training



**Section-H:** Classification Accuracy Evolution for PPO Training
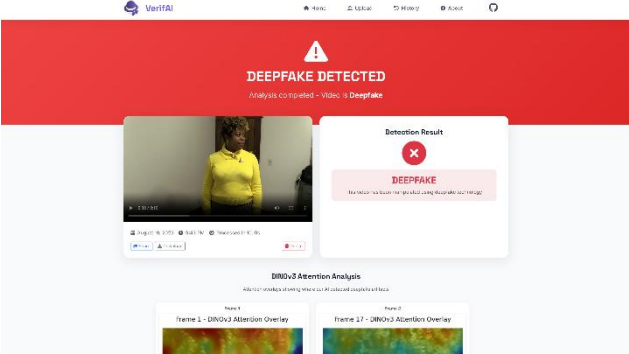
**Section-I:** Decision Speed Optimization



**Section-J:** PPO and GRPO Fine Tuning Rewards Comparison Metrics



**Section-B:** Results Page Interface



## Appendix B: Web Application Deployment Interface

This appendix demonstrates the practical deployment of the EAGER framework through the VerifAI web application, a Django-based interface that enables real-time deepfake detection for end users.

The application's landing page, which reports system performance metrics including 95 percent accuracy and 10-second average processing time across over 5,000 analyzed videos. The detection results interface, presenting the classification outcome alongside DINOv3 attention heatmaps that highlight facial regions contributing to the detection decision. The attention visualizations provide transparency by revealing that the model focuses on facial boundaries and texture transitions where manipulation artifacts typically appear. The system processes videos in an average of 10.19 seconds, validating the practical viability of the reinforcement learning approach for real-world deployment while maintaining user privacy through secure file handling and automatic cleanup protocols.

**Section-A:** Home Landing Page of Django-based Interface