# Advancing Explainability in Deep Reinforcement Learning: Novel Frameworks for Transparency, Fairness, and Robustness in Autonomous Decision-Making Systems

**[1]Akanksh Reddy Muddam, [2]Lavanya Reddy Satti**

[1]MSc in Artificial Intelligence, [2]Assistant Professor,
[2]Keshav Memorial Institute of Technology,
[1]Dublin, Ireland, [2]Hyderabad, India.

*Abstract:* DRL is incredibly successful in autonomous decision-making in many complex fields. Nonetheless, due to the black-box character of DRL models, they are not transparent, fair, or trustworthy, and their use is not possible in high-stakes scenarios. This paper suggests innovative frameworks that develop the concept of explainability in DRL through incorporating global interpretability methods, fairness-conscious policy analysis, and the improvement of robustness. We present SILVER, a Shapley value-based interpretable policy framework that is a middle ground between explainability and interpretability. The effectiveness of the latter methods is defined by the fact that the experiments conducted on the benchmark control tasks and autonomous systems demonstrate a high level of transparency, fairness indicators, and stability without negatively affecting their functionality. We make progress towards implementing self-reliant, trustworthy decision-making systems to enable us to safely deploy in adverse real-life situations.

*Index Terms* - **Frameworks for Enhancing Transparency, Fair- ness, and Robustness in Explainable Deep Reinforcement Learning for Autonomous Systems**

## I. INTRODUCTION

Deep Reinforcement Learning (DRL) is an evolutionary breakthrough in artificial intelligence, which combines the strong representational capabilities of deep neural networks with the sequential decision-making capacity of reinforcement learning. This will allow DRL to learn the optimal set of strategies to employ in complex tasks independently by interacting with dynamic tasks and promoting their beneficial behaviours. This capability has enabled DRL to lead the AI research and applications field over the past few years, making revolutionary advances in multiple fields, including game-based AI, such as Go and chess, robotic control, self-driving vehicles, and optimisation of industrial processes. These and other accomplishments have highlighted the potential of DRA as a foundational technology in the field of intelligent autonomous systems capable of adaptability, learning and functioning in extremely complex and uncertain real-world environments.

Until these impressive developments, DRL models are, in many ways, black-boxes because of the deep neural networks they contain. Their operational models acquire internal representations of features and policies that cannot usually be understood by human observers, even by expert practitioners. Such non-transparency poses significant issues to the meaning, verification and plausibility of DRL-generated judgments and actions. The inability to explain why certain actions are taken by the DRL systems is becoming a capital liability to the reliability and practicality of the systems, as they are increasingly deployed in safety-critical and high-stakes fields, including health diagnosis, financial decision-making, autonomous vehicles, and defence. To make decisions equitable, objective, resistant to adversarialness, and ethical and legally sound, practitioners and regulators need transparent and interpretable AI.

In response to the gap in knowledge about the nature of core explainability requirements, the work suggests innovative models directly connected to directly improve the interpretability and credibility of DRL systems. We deal with three basic characteristics, which jointly provide safe and ethical au- autonomy: transparency, fairness, and robustness. Transparency can be defined as the ability of the DRM system to give understandable insights into the process of decision-making that would reveal the underlying behaviour of the policy in a form that can be read by human beings. Fairness means that no DRL policies share biases or discrimination between users or between states of the environment, which is a crucial factor in fair AI use. Robustness focuses on how DRA systems can withstand uncertainty, adversarial attack, or perturbation that otherwise might result in harmful or unpredictable behaviour. To do so, we make use of state-of-the-art explainability techniques, with the most significant of these being Shapley value-based explanation techniques that emulate cooperative game theory to assign interpretable importance scores to input features affecting decisions. The Shapley value-based models of local explanation can provide both local and global explanations.

Interpretability, in contrast to the traditional local explanation techniques, which explain single decisions, explains the various state characteristics in the overall behaviour of the DRL policy. By exploiting this method, the derivation of interpretable policies, which closely match the original DRL model but have a better human-understandable structure, can be obtained.

Simultaneously, we integrate metrics of fair evaluation, designed to function in sequential decision-making environments, into DRL, including demographic parity and equality of opportunity, to measure and proactively mitigate possible biases during learning and inference. This type of integration means that individual agents are treated as a single population or condition, and is important in applications such as personalised medicine or autonomous transportation. Moreover, to ensure reliability in a non-ideal setting, our framework uses robustness testing with adversarial training programs based on explainability feedback. It identifies vulnerable states or critical properties that are sensitive to decision-making and makes policies less susceptible to perturbation, mitigating the risks of bad or risky policy.

These contributions combined provide a practical and com- comprehensive roadmap to all credible DRS systems models, both working and explanations, are just, and able to survive adversarial attacks. By so doing, this work helps fill one of the central gaps between the theoretical potential of DRL and the ethical, safety, and regulatory realities of production.

The role of improving explainability in DRL cannot be overemphasised as AI systems gain autonomy and become integrated into how society operates. Through creative interpretability systems, the paper seeks to facilitate responsible AI use to gain the trust of stakeholders, such as researchers, practitioners, regulators, and end-users. Eventually, the resulting innovations open the way to next-generation systems of autonomous decisions, not only intelligent but also ethical, robust, and transparent, which are the key characteristics of sustainable AI impact in crisis applications in healthcare, finance, transportation, and other fields.

## II. RELATED WORK

THE MOST RECENT PROGRESS IN EXPLAINABLE REINFORCEMENT LEARNING (XRL) HAS BEEN MAINLY FOCUSED ON THE CREATION OF METHODS TO PRODUCE LOCAL EXPLANATIONS OF THE DEEP REINFORCEMENT LEARNING (DRL) AGENT DECISION-MAKING PROCESS. THE GOAL OF SUCH LOCAL EXPLANATION METHODS IS TO UNDERSTAND WHY A SPECIFIC ACTION WAS SELECTED IN A GIVEN STATE, THUS MAKING IT MORE INTERPRETABLE AT THE DECISION LEVEL. STARTLING FEATURE ATTRIBUTION APPROACHES SUCH AS SHAP (SHAPLEY ADDITIVE EXPLANATIONS) AND LIME (LOCAL INTERPRETABLE MODEL AGNOSTIC EXPLANATIONS) HAVE BECOME POPULAR IN SUCH APPLICATIONS. SHAP TO CALCULATE IMPORTANCE SCORES ON INPUT FEATURES WHICH CONTRIBUTE TO THE ACTION OF THE AGENT BASED ON COOPERATIVE GAME THEORY BY INDICATING THE LOCAL SIGNIFICANCE OF STATE VARIABLES. SIMILARLY, LIME ESTIMATES A COMPLICATED MODEL BY A MUCH EASIER-TO-INTERPRET MODEL IN THE NEIGHBOURHOOD TO PROVIDE HUMAN EXPLANATIONS. ALTHOUGH THEY ARE USEFUL IN SITUATIONS REQUIRING SPECIFIC INTERPRETABILITY, THESE METHODS ARE BY DEFINITION RESTRICTED TO LOCALISED VIEWS AND FAIL TO GIVE A GLOBAL VIEW OF THE OVERALL POLICY BEHAVIOUR OF DRL AGENTS.

THE POSSIBILITY OF ATTAINING SOME GLOBAL INTERPRETABILITY HAS ALSO GIVEN RISE TO STUDIES THAT GO BEYOND LOCAL EXPLANATIONS. SPECIFICALLY, LI ET AL. PROPOSED THE SILVER FRAMEWORK, SHAPLEY VALUE-BASED INTERPRETABLE POLICY VIA EXPLANATION REGRESSION, AN EXAMPLE OF A MODEL-AGNOSTIC MODEL THAT CONVERTS COMPLEX DRL POLICIES INTO INTERPRETABLE, TRANSPARENT FORMS. SILVER CLOSES THE EXPLAINABILITY VERSUS INTERPRETABILITY GAP, USING SHAPLEY VALUES BOTH TO PROVIDE LOCALISED EXPLANATIONS AND TO MAKE INFERENCES ABOUT GLOBAL DECISION BOUNDARIES WHICH DESCRIBE THE AGENT AS A WHOLE POLICY IN THE STATE SPACE. BY DOING SO, PRACTITIONERS CAN OBTAIN A MORE DETAILED PICTURE OF POLICY MECHANISMS WHILST MAINTAINING PERFORMANCE AND THEREFORE ACHIEVING THE COMMON TRADE-OFF BETWEEN ACCURACY AND INTERPRETABILITY.

IN ADDITION TO MAKING AI MORE INTERPRETABLE, THE INCREASING RECOGNITION OF ETHICAL AND SOCIAL CONSEQUENCES OF AI HAS PROMPTED RESEARCH IN FAIRNESS-CONSCIOUS REINFORCEMENT LEARN- ING. FAIRNESS IN RL AIMS AT CONSIDERING EQUITABLE TREATMENT AND RESULTS AMONG DIFFERENT POPULATIONS AND SENSITIVE GROUPS SUBJECTED TO AUTONOMOUS DECISION-MAKING SYSTEMS. RESEARCHERS HAVE SUGGESTED WHAT ARE CALLED SPECIALISED FAIRNESS METRICS THAT APPLY TO RL CONTEXTS, GENERALISING PREVIOUS CONCEPTIONS OF FAIRNESS IN SUPERVISED LEARNING, INCLUDING DEMOGRAPHIC PARITY AND EQUALITY OF OPPORTUNITY, TO THE CONTEXT OF SEQUENTIAL DECISIONS. THESE TECHNIQUES DETECT AND MINIMISE BIASES THAT COULD ARISE AS A CONSEQUENCE OF A SKEWED TRAINING SAMPLE, REWARD MECHANISMS OR ENVIRONMENTAL INTERACTIONS THAT HURT CERTAIN GROUPS OF INDIVIDUALS OR CIRCUMSTANCES. JUSTICE-CONSCIOUS RL MODELS CAN THEREFORE PLAY A ROLE IN CREATING BELIEVABLE AI SYSTEMS THAT SUPPORT NON-DISCRIMINATORY POLICIES IN AREAS SUCH AS HEALTHCARE, EMPLOYMENT, AND CRIMINAL JUSTICE.

DRA ROBUSTNESS FACTOR MUST BE ROBUST TO PERTURBATION, ADVERSARIAL EXAMPLES AND ENVIRONMENTAL UNCERTAINTY, WHICH IS IMPORTANT WHEN APPLYING AUTONOMOUS SYSTEMS TO SAFETY-RELATED ENVIRONMENTS. TECHNIQUES ROBUSTNESS. ROBUSTNESS WAS TRADITIONALLY A SUBJECT OF SYSTEM DESIGN, AND ADVERSARIAL TRAINING, SAFE EXPLORATION, AND QUANTIFYING UNCERTAINTY ARE ALL CLASSIC APPROACHES TO IMPROVING THE STABILITY AND RELIABILITY OF SYSTEMS. NONETHELESS, THE LITERATURE TO DATE WOULD SEPARATE ROBUSTNESS AND EXPLAINABILITY, EITHER BY HARDENING POLICIES AGAINST ATTACKS OR CLARIFYING THE DECISION LOGIC WITHOUT A UNIFYING COMBINATION OF THESE OBJECTIVES.

WE WORK AT THE JUNCTION OF EXPLAINABILITY, FAIRNESS, AND ROBUSTNESS IN DEEP REINFORCEMENT LEARNING. WE SUGGEST A SINGLE CONCEPTUALISATION THAT SYNERGISTICALLY COVERS THESE IMPORTANT POINTS, THEREBY DEVELOPING HOLISTIC,

RELIABLE DRL SYSTEMS. OUR PROPOSAL COMBINES MODEL-AGNOSTIC GLOBAL INTERPRETABILITY STRATEGIES, SUCH AS SILVER, WITH POLICIES EVALUATED FAIRLY AND EQUITABLY, WITH EXPLAINABILITY-BASED ADVERSARIAL TRAINING TO PROMOTE CLEAR, FAIR, AND SUSTAINABLE AUTONOMOUS DECISION-MAKING. SUCH A HOLISTIC VIEW NOT ONLY MAKES DRL MODELS INTERPRETABLE AND TRUSTWORTHY BUT ALSO ETHICAL AND RELIABLE IN A WIDE AND POTENTIALLY HOSTILE RANGE OF CONDITIONS. AN INTEGRATIVE FRAMEWORK LIKE THIS IS NEEDED TO ADDRESS THE GROWING DEMANDS OF IMPLEMENTING DRL IN MULTI-FACETED REAL-WORLD SETTINGS WHERE RESPONSIBILITY AND SECURITY ARE NOT NEGOTIABLE.

IN SHORT, ALTHOUGH PREVIOUS STUDIES HAVE SUCCEEDED IN MAKING WORTHWHILE CONTRIBUTIONS IN TERMS OF PROVIDING LOCAL PLANATIONS OF DECISIONS, REDUCING BIAS, OR IMPROVING ROBUSTNESS ON THEIR OWN, OUR CONTRIBUTION TIES THESE STRANDS TOGETHER WITHIN A UNIFIED FRAMEWORK. BESIDES DEVELOPING THEORETICAL KNOWLEDGE, THIS SYNTHESIS OFFERS SOME PRACTICAL METHODOLOGIES AND TOOLS THAT WILL TRIGGER THE RESPONSIBLE IMPLEMENTATION OF DRL TECHNOLOGIES IN AREAS THAT REQUIRE TRANSPARENCY, FAIRNESS, AND SAFETY.

## III. PROPOSED FRAMEWORK

TO ACHIEVE THE EXPLAINABILITY, FAIRNESS, AND ROBUSTNESS OF DEEP REINFORCEMENT LEARNING (DRL), NEW STRUCTURES ARE NEEDED THAT CAN DIRECTLY TACKLE THE ISSUE OF COMPLEXITY AND OPAQUENESS OF DRL POLICIES. IN THIS SECTION, THREE FOCAL ELEMENTS OF THE PROPOSED METHODOLOGY ARE FURTHER EXPLAINED IN DETAIL TO FILL IN THE GAPS AND PROVIDE A HOLISTIC SOLUTION TO INTERPRETABLE, FAIR, AND RESILIENT DRL SYSTEMS.

A.     SILVER: INTERPRETABLE POLICY THROUGH EXPLANATION RE- REGRESSION.

THE ADAPTATION AND EXPANSION OF THE SILVER FRAMEWORK (SHAPLEY VALUE-BASED INTERPRETABLE POLICY VIA EXPLANATION REGRESSION), A NEW MODEL-AGNOSTIC MODEL THAT PRIMARILY TRIES TO CONVERT HIGHLY COMPLEX AND OPAQUE DRL POLICIES INTO TRANS- PARENT AND INTERPRETABLE ONES, IS THE FIRST PILLAR OF OUR PROPOSAL. OWING TO THEIR DEPENDENCE ON DEEP NEURAL NETWORKS, DEEP REINFORCEMENT LEARNING ALGORITHMS SHOW HIGH PERFORMANCE AT THE EXPENSE OF HIGH OPAQUENESS. CURRENT EXPLAINABILITY TOOLS TEND TO BE ONLY LOCAL EXPLANATIONS–EXPLANATIONS ABOUT WHY AN AGENT MADE THE DECISION IT TOOK AT A CERTAIN STATE. BUT THEY DO NOT DISCLOSE THE WORLD POLICY FRAMEWORK OR THE GENERAL RATIONALE OF THE DECISION WHICH DETERMINES THE BEHAVIOUR OF THE AGENT IN VARIOUS STATES.

SILVER IS THE ONLY METHOD TO USE THE POWER OF SHAPLEY VALUES, BASED ON COOPERATIVE GAME THEORY, TO ESTIMATE THE VALUE OF EACH INPUT FEATURE TO THE POLICY DECISIONS MADE BY THE AGENT IN A MATHEMATICALLY RIGOROUS WAY. IN CONTRAST TO MORE CONVENTIONAL APPROACHES, WHICH LOOK AT LOCAL INTERPRETABILITY IN A VERY NARROW SENSE, SILVER APPLIES SHAPLEY VALUES TO REGRESS INTERPRETABLE POLICY REPRESENTATIONS WHICH PREDICT THE BEHAVIOUR OF THE ORIGINAL DRL AGENT ACROSS THE WHOLE STATE-ACTION SPACE. THIS REGRESSION CONVERTS BLACK-BOX POLICIES INTO CLOSED-FORM AND MUCH SIMPLER, BUT VERY PRECISE, DECISION BOUNDARIES THAT ARE COMPREHENSIBLE AND INTERPRETABLE BY HUMAN STAKEHOLDERS.

ONE OF THE MAIN ADVANTAGES OF SILVER IS THAT IT ELIMINATES THE VERY IDEA OF A TRADE-OFF THAT IS BELIEVED TO BE PRESENT BETWEEN THE TRANSPARENCY OF THE MODEL AND THE PERFORMANCE OF THE POLICY. EXPERIMENTS ON TRADITIONAL CONTROL TASKS LIKE CARTPOLE, MOUNTAINCAR, AND ACROBOT TESTIFY TO THE FACT THAT SILVER-GENERATED INTERPRETABLE POLICIES RETAIN THE SAME LEVEL OF COMPETITIVE REWARD PERFORMANCE AS THEIR CORRESPONDING DEEP RL EQUIVALENTS, AND ALWAYS OFFER GLOBAL POLICY INSIGHTS THAT IMPROVE TRUST AND UNDERSTANDING. ADDITIONALLY, THE FRAMEWORK IS MODEL-AGNOSTIC SO THAT IT HAS A WIDE RANGE OF APPLICABILITY, SUPPORTING BOTH ON-POLICY (E.G., PPO, A2C) AND OFF-POLICY (E.G., DQN) ALGORITHMS, AND HENCE ENABLING INTEGRATION INTO A WIDE RANGE OF DRL PARADIGMS.
THE INTERPRETABILITY OFFERED BY SILVER ALLOWS PRACTITIONERS, REGULATORS, AND END-USERS TO PREDICT AGENT BEHAVIOURS WITH HIGH RELIABILITY, CONDUCT POLICY AUDITS, AND IDENTIFY ANOMALIES, WHICH IS A BIG STEP TOWARDS TRUSTED DRL DEPLOYMENT IN A COMPLEX, REAL-WORLD SETTING. SILVER CAN BRING ABOUT A PARADIGM SHIFT IN BLACK-BOX DECISION-MAKING IN FAVOUR OF TRANSPARENT AND ACCOUNTABLE AUTONOMOUS SYSTEMS BY EXPLAINING AND INTERPRETING THE PHENOMENON OF EXPLAINABILITY AND INTERPRETABILITY BETWEEN TWO RIGOROUS EXPLANATION REGRESSION TECHNIQUES.

B.     POLICY ANALYSIS WITH FAIRNESS AWARENESS.

WHEREAS INTERPRETABILITY DEALS WITH TRANSPARENCY, IT IS NOT SUFFICIENT TO ENSURE ETHICAL OR FAIR RESULTS OF AUTONOMOUS DECISION-MAKING. THE SECOND ELEMENT OF OUR FRAMEWORK, THEREFORE, FOCUSES ON FAIRNESS, AN ESSENTIAL FEATURE THAT IS NOW ALSO ACKNOWLEDGED IN AI ETHICS. REINFORCEMENT LEARNING SYSTEMS USED IN PRACTICE SHOULD BE SUCH THAT THE LEARNED POLICIES ARE NOT REINFORCED OR ENHANCED TO CONTINUE TO FAVOUR CERTAIN GROUPS OF PEOPLE OR ENVIRONMENTAL CIRCUMSTANCES THAT ARE UNFAIRLY DISADVANTAGEOUS.

WE PRESENT FAIRNESS EVALUATION METRICS THAT ARE SPECIFICALLY TAILORED TO THE SEQUENTIAL CHARACTER OF DRL. DEMOGRAPHIC PARITY AND EQUAL OPPORTUNITY, THE TRADITIONAL NOTIONS OF FAIRNESS THAT HAVE BEEN STUDIED WELL IN SUPERVISED LEARNING, TEND TO SUPPORT THE TEMPORAL DEPENDENCIES OF REINFORCEMENT LEARNING. DEMOGRAPHIC PARITY HERE GUARANTEES THAT THE ODDS OF UNDERTAKING FAVOURABLE ACTIONS ARE INDIFFERENT TO SENSITIVE ATTRIBUTES (E.G.,

GENDER, ETHNICITY) OF THE DIFFERENT STATES, WHILST EQUAL OPPORTUNITY ENSURES THAT ACTUAL RATES OF TRUE POSITIVES ARE SIMILAR AMONG GROUPS ALONG THE LABELLING TRAJECTORY.

OUR ARCHITECTURE INSTANTIATES THESE MEASURES OF FAIRNESS AS PART OF A DRL TRAINING AND EVALUATION PIPELINE, ALLOWING ONE TO IDENTIFY THE POLICY ASPECTS OR PARTS OF THE STATE WHERE BIAS IS PRESENT. WHEN IDENTIFIED, FAIRNESS LIMITATIONS MAY BE INTEGRATED INTO THE INCENTIVE PLAN OR POLICY STANDARDISATION CONDITIONS TO PUNISH DISCRIMINATORY PRACTICES. IT IS AN ACTIVE FORM OF BIAS REDUCTION THAT WILL PROVIDE MORE BALANCED RESULTS WITHOUT SACRIFICING TASK PERFORMANCE TO A CONSIDERABLE EXTENT.

WITH FAIRNESS-CONSCIOUS ANALYSIS, WE DEVELOP DRL AGENTS THAT CAN MAKE SOCIALLY RESPONSIBLE CHOICES AND, CONSEQUENTLY, INCREASE STAKEHOLDER TRUST AND ADHERENCE TO NEW AI GOVERNANCE REGULATIONS. THIS IS OF PARTICULAR CONCERN IN SENSITIVE APPLICATIONS SUCH AS HEALTHCARE ALLOCATION, CREDIT SCORING AND CRIMINAL JUSTICE, WHERE BIASED POLICY HAS PARTICULARLY IMPORTANT ETHICAL AND LEGAL IMPLICATIONS.

C.     STRENGTHENING THROUGH EXPLAINABILITY-GUIDED ADVERSARIAL- IAL TRAINING.

THE LAST DIMENSION THAT IS ESSENTIAL TO TRUSTWORTHINESS IN DRL IS THE CONCEPT OF ROBUSTNESS: CAN TRAINED POLICIES EXECUTE RELIABLY WHEN FACED WITH UNPREDICTABLE PERTURBATION, NOISY DATA, OR ADVERSARIAL ATTACKS? IN THE TRADITIONAL ADVERSARIAL TRAINING APPROACHES, THE AIM IS TO IDENTIFY THE WEAKNESSES OF THE MODEL, AND IT ATTEMPTS TO GENERATE WORST-CASE PERTURBATIONS SYSTEMATICALLY AND RETRAIN THE POLICIES TO RESIST THEM. NONETHELESS, INTERPRETABILITY INSIGHTS ARE RARELY USED IN SUCH APPROACHES IN ORDER TO EXPLICITLY ADDRESS THE WEAKEST POINTS OF THE MODEL.

THE ROBUSTNESS ENHANCEMENT FRAMEWORK THAT WE PROPOSE IS NOVEL IN THAT IT USES THE EXPLAINABILITY FEEDBACK TO LEARN TO TRAIN ADVERSARIAL MODELS AND DO SO MORE EFFICIENTLY. BY USING INTERPRETABILITY METHODS, MOST COMMONLY SHAPLEY-BASED METHODS OR METHODS THAT ARE SIMILAR TO ATTRIBUTE-BASED FEATURE ATTRIBUTION, WE CAN DETERMINE WHICH STATE FEATURES OR DECISION POINTS WITHIN THE DRL POLICY ARE THE MOST INFLUENTIAL AND IMPORTANT IN DETERMINING THE ACTION CHOICES.

THESE SUSCEPTIBLE STATES OR CHARACTERISTICS ARE THEN TARGETED AS THE TARGETS OF TAILORED ADVERSARIAL PERTURBATIONS OR NOISE INJECTIONS DURING TRAINING. THE POLICY ENHANCES ITS RESILIENCE IN THESE KEY ASPECTS, FURTHER IMPROVING ITS OVERALL STABILITY AND PERFORMANCE WHEN CONFRONTED WITH STRESSFUL SITUATIONS BY FOCUSING ON ITS ROBUSTNESS EFFORTS IN THESE KEY AREAS. THE TARGETED APPROACH IS BETTER THAN UNSELECTIVE ADVERSARIAL TRAINING BECAUSE IT REDUCES THE COMPUTATION COST AND DOES NOT INDUCE UNNECESSARY POLICY DEGRADATION IN LESS CRITICAL REGIONS.

EMPIRICAL ANALYSES HAVE SHOWN THAT EXPLAINABILITY-BASED ADVERSARIAL TRAINING RESULTS IN MORE ROBUST POLICIES, WHICH ARE INTERPRETABLE AND FAIR, AND COMPLEMENT THE TRANSPARENCY AND BIAS-REDUCTION GOALS OF THE LARGER FRAMEWORK. THE COMBINATION OF THESE THREE DIMENSIONS CREATES A COMPREHENSIVE VISION OF CREDIBLE DRL - CREATING AUTONOMOUS SYSTEMS THAT ARE BOTH TRANSPARENT IN THEIR REASONING AND FAIR IN THEIR BEHAVIOUR, AND ROBUST TO ATTACK OR ENVIRONMENTAL VAGARIES.

AS A FINAL POINT, THE SUGGESTED FRAMEWORKS COMPRISE A UNIFIED SET OF METHODOLOGIES TO DEAL WITH EXPLAINABILITY, FAIRNESS, AND ROBUSTNESS CHALLENGES IN DRL THAT ARE INTERDEPENDENT. - CORPORATING MODEL-AGNOSTIC INTERPRETABILITY THROUGH SILVER, FAIRNESS CONSTRAINTS CUSTOMISED TO SEQUENTIAL DECISION-MAKING, AND ROBUSTNESS REFINEMENT BASED ON EXPLAINABILITY-INFORMED ADVERSARIAL TRAINING, THIS WORK REPRESENTS A TRAILBLAZER ON THE WAY TO HIGH-PERFORMING, ETHICAL, AND SAFE AUTONOMOUS AGENTS. THIS BREAKTHROUGH IS NECESSARY TO ENABLE RESPONSIBLE AND SCALABLE APPLICATION OF DRL TECHNOLOGY ACROSS MISSION-CRITICAL DOMAINS, SUCH AS HEALTHCARE AND FINANCE, AUTONOMOUS VEHICLES AND MILITARY SYSTEMS.

## IV. EXPERIMENTAL SETUP

THE HIGH-FIDELITY TESTING OF THE SUGGESTED ARCHITECTURES TO IMPROVE EXPLAINABILITY, FAIRNESS, AND ROBUSTNESS IN DEEP RE- REINFORCEMENT LEARNING (DRL) REQUIRES A CAREFULLY CONSTRUCTED EXPERIMENTAL DESIGN. IN THIS WAY, THE FULL ASSESSMENT OF THE ALGORITHMIC PERFORMANCE IS POSSIBLE, ALONG WITH THE PRACTICAL FEASIBILITY AND THE MORAL RIGHTNESS OF THE AUTONOMOUS AGENTS IN VARIOUS AND CHALLENGING SETTINGS. THE ENVIRONMENT DESIGN AND ALGORITHMIC BASELINES ARE DISCUSSED IN DETAIL, ALONG WITH THE MULTI-DIMENSIONAL EVALUATION METRICS USED IN THIS SECTION, TO JUSTIFY THE PROPOSED METHODOLOGIES.

A. ENVIRONMENT SELECTION AND DESIGN

THE STRUCTURE OF THE ENVIRONMENT SHOULD BE CHOSEN AND DESIGNED SO AS TO IMPROVE THE LEARNING PROCESS (MAYER, 1994). THE EXPERIMENT WAS BASED ON AN INTEGRATION OF CLASSICAL CONTROL STANDARDS AND A SIMULATED AUTONOMOUS NAVIGATION SPACE, EACH SELECTED BASED ON A SET OF QUALITIES THAT MADE
THEM SUITABLE FOR INTENSIVE TESTING.

THE MOUNTAINCAR PROBLEM IS A CLASSIC SPARSE-REWARD CONTINUOUS STATE-SPACE PROBLEM. IT IS COMPLICATED IN THAT IT REQUIRES THE AGENT TO STORE POTENTIAL ENERGY THROUGH OSCILLATION BEFORE ATTAINMENT OF THE AIM OF REACHING

THE PEAK OF THE MOUNTAIN. THIS KIND OF DELAYED PAYOFF DESIGN HAS ITS PROBLEMS IN THE FORM OF A CHALLENGE IN PLANNING AS WELL AS LEARNING TO IDENTIFY LONG-TERM DEPENDENCIES, THE QUALITIES NECESSARY TO DETERMINE THE INTERPRETIVE TRANSPARENCY OF LONG-PATH POLICIES.

CARTPOLE IS A RELATIVELY LOW-DIMENSIONAL CLASSICAL CONTROL PROBLEM WHICH INVOLVES FINE CONTROL TO MAINTAIN A POLE UPON A MOVING CART. THE CONSISTENCY OF ITS APPLICATION IN THE REINFORCEMENT LEARNING LITERATURE PROVIDES A CONSISTENT STANDARD BY WHICH THE EFFECTIVENESS OF POLICY LEARNING AND MODEL TRANSPARENCY CAN BE JUDGED. ITS RELATIVELY FAST SWITCHING OF STATES ENABLES EXHAUSTIVE POLICY TESTING AND INTERPRETABILITY ANALYSES EVEN WITH LIMITED COMPUTATIONAL RESOURCES.

BUILDING ON THESE BENCHMARKS, AN AUTONOMOUS NAVIGATION ENVIRONMENT WAS CREATED THAT REPLICATES THE COMPLEXITY FOUND IN THE REAL WORLD, INCLUDING DYNAMIC OBSTACLE PLACEMENT, VARIABLE SENSOR INPUTS WITH NOISE, AND PARTIALLY OBSERVABLE CONDITIONS THAT ARE REMINISCENT OF REAL ROBOTIC ENVIRONMENTS. THIS KIND OF ENVIRONMENT WILL CHALLENGE THE CAPACITY TO ACT HARD AND FAIRLY SINCE IT WILL REQUIRE THE AGENT TO UTILISE ETHICAL AND TRUSTWORTHY JUDGMENT IN A CIRCUMSTANCE WHERE THE SAFETY AND FAIRNESS ISSUE TAKES PRECEDENCE.

## B. ALGORITHMIC BASELINES

TO EFFECTIVELY BENCHMARK THE PROPOSED IMPROVEMENTS, THE EXPERIMENTAL DESIGN ACCOMMODATED TYPICAL AND GENERALLY AGREED-UPON DRL ALGORITHMS IN VARIED LEARNING PARADIGMS:

• DEEP Q-NETWORK (DQN) IS AN OFF-POLICY, VALUE-BASED METHOD, WHICH USES EXPERIENCE REPLAY AND TARGET NETWORKS TO STABILISE TRAINING IN DISCRETE ACTION SPACES- ENVIRONMENTS. ITS APPLICATION CAN BE USED AS A REFERENCE POINT TO DISCRETE CONTROL PROBLEMS SUCH AS MOUNTAINCAR AND CARTPOLE.

• PROXIMAL POLICY OPTIMISATION (PPO) IS A POLICY GRADIENT ALGORITHM THAT WORKS ON-POLICY AND USES CLIPPED SURROGATE OBJECTIVES TO STABILISE UPDATES AND EFFECTIVELY USE SAMPLES. DISCRETE AND CONTINUOUS ACTIONS ALLOW PPO TO BE ONE OF THE STANDARDS OF CONTEMPORARY DRL.

• COMBINING BOTH VALUE-BASED AND POLICY-GRADIENT AP- APPROACHES, ADVANTAGE ACTOR-CRITIC (A2C) PROVIDES POLICY UPDATES THAT ARE SYNCHRONOUS AND HAVE LOWER VARIANCE. IT CAN BE USED TO PERFORM CONTINUOUS CONTROL AND OFFERS COMPLEMENTARY INFORMATION ON PERFORMANCE TO PPO AND DQN.

THESE ALGORITHMS HAVE BEEN APPLIED WITH AND WITHOUT INCORPORATING THE PROPOSED MODULES OF EXPLAINABILITY, FAIRNESS, AND ROBUSTNESS, WHICH PROVIDES A CONTROLLED COMPARATIVE ASSESSMENT.

## C. MULTI-DIMENSIONAL ASSESSMENT METRICS.

THE EFFECTIVENESS OF ANY INTERPRETABILITY-SUPPORTING FRAMEWORK REQUIRES A MULTI-FACETED EVALUATION METRIC TO INCLUDE CLARITY, ETHICAL SOUNDNESS, PERTURBATION RESISTANCE AND TASK PERFORMANCE. THE FOLLOWING WERE THEREFORE USED:

1) INTERPRETABILITY METRICS: THE POLICY INTERPRETABILITY WAS MEASURED BOTH IN TERMS OF HUMAN PROXIES AND FUNCTIONALLY GROUNDED PROXIES.

UNDERSTANDABILITY WAS MEASURED THROUGH STRUCTURED QUES- QUESTIONNAIRES GIVEN TO DOMAIN SPECIALISTS AND TO NON-EXPERIENCED USERS WHO ASSESSED THE READABILITY OF SURROGATE EXPLANATIONS BASED ON THE SILVER FRAMEWORK. PERCEIVED EASE OF UNDERSTANDING, EXPLANATION OF POLICY COVERAGE AND CONFIDENCE IN THE RATIONALE OF THE DECISIONS.

FIDELITY MEASURES DESCRIBED THE STATISTICAL CONSISTENCY OF INTERPRETABLE SURROGATE POLICIES AND ORIGINAL DRL POLICIES. THIS WAS CALCULATED BY TAKING THE PERCENTAGE CONCURRENCE ON A WIDE RANGE OF STATES AND ACTIONS. HIGH FIDELITY MUST BE USED SO THAT INTERPRETIVE MODELS ARE SIGNIFICANT INDICATORS OF HOW THE BLACK BOX MAKES A DECISION.

OTHER COMPUTATIONAL PROXIES, INCLUDING INFERENCE TIME AND POLICY COMPLEXITY (E.G., SIZE OF SURROGATE DECISION TREES OR REGRESSION FUNCTIONS), GAVE SURROGATE ESTIMATES OF INTERPRETABILITY IN LINE WITH THE "SIMULATABILITY" PHILOSOPHY OF HUMAN-CENTRED AI STUDIES.

2) FAIRNESS METRICS: THE FAIRNESS DIMENSION STRETCHED CLASSICAL SUPERVISED LEARNING FAIRNESS NOTIONS TO REINFORCEMENT LEARNING. KEY METRICS INCLUDED:

ACTION DISTRIBUTION DISPARITY: QUANTITATIVE MEASURE OF THE EVENNESS OF POSITIVE OR NEUTRAL ACTION IN IDENTIFIED SENSITIVE GROUPS OR CLUSTERS OF STATES (SUCH AS DIFFERENT SPATIAL ZONES OF THE NAVIGATION ENVIRONMENT), WHICH ESTIMATES DEMOGRAPHIC PARITY.

EQUAL OPPORTUNITY MEASURES: OVERVIEW OF ACTUAL POSITIVE ACTION RATES THAT REPRESENT FAIR PERFORMANCE, SUCH THAT POLICIES DID NOT DISCRIMINATE AGAINST ANY GROUP OR SITUATION IN A SEQUENCE OF DECISIONS.

BIAS DETECTION RATE: SENSITIVITY ANALYTIC TECHNIQUES IDENTIFIED AND MEASURED FAVOURABLE PATTERNS OF STATE-ACTION THAT SIGNALLED BIASES IN THE LEARNED POLICIES, WHICH CAN BE SPECIFICALLY ETHICALLY AUDITED AND COUNTERMEASURED.

3) ROBUSTNESS METRICS: A KEY ASPECT OF CONFIDENCE IN PRACTICE APPLICATION WAS THAT THE POLICY RESILIENCE WAS TESTED WITHIN ADVERSE OPERATING CONDITIONS.

FAST GRADIENT SIGN METHOD FAST GRADIENT SIGN METHOD FAST GRADIENT SIGN METHOD HAS BEEN USED BY ADVERSARIAL PERTURBATION TESTS TO MAKE THE AGENTS LIE OR PERFORM POORLY BY USING SPECIFICALLY STRUCTURED INPUT PERTURBATION. CUMULATIVE REWARD AND BEHAVIOUR STABILITY THAT WAS REGISTERED WERE DEGRADED.

NOISE INJECTION EXPERIMENTS ADDED GAUSSIAN AND REAL-WORLD SENSOR NOISE TO INPUT STATES. THE CAPABILITY OF AGENTS TO FOLLOW SAFE NAVIGATION ROUTES AND TO KEEP A SUCCESS RATE IN THEIR TASKS UNDER THESE STOCHASTIC CONDITIONS WAS CONSIDERED.

STABILITY UNDER PARTIAL OBSERVABILITY PROVIDED AN EVALUATION OF THE CAPABILITY TO ADAPT POLICIES WHEN THERE WAS LIMITED SENSOR AVAILABILITY TO DETERMINE WHETHER STABILITY WITH PARTIAL OBSERVABILITY COULD OCCUR WITHOUT EXTREME DEGRADATION.

4) PERFORMANCE METRICS: ALTHOUGH INTERPRETABILITY, FAIRNESS AND ROBUSTNESS WERE CONSIDERED, TASK PERFORMANCE WAS NEEDED TO STAY OR TO BE IMPROVED TO BE PRACTICALLY VIABLE:

A CUMULATIVE REWARD ACROSS SEVERAL EPISODES OFFERED AN OBJECTIVE AND CONSISTENT MEASURE OF SUCCESS IN TASKS. EFFICIENCY AND CONSISTENCY MEASURES WERE ADDED TO REWARD-BASED METRICS, INCLUDING EPISODE LENGTH AND TASK COMPLETION RATE.

## D. EXPERIMENTAL PROTOCOL AND METHODOLOGY

THE PROCEDURE WILL INCLUDE THE FOLLOWING: EXPERIMENTAL PROTOCOL AND METHODOLOGY.

EACH AGENT WAS TRAINED ON THE SAME HYPERPARAMETERS (E.G. LEARNING RATE, DISCOUNT FACTOR G=0.99 AND BATCH SIZES) THAT WERE OPTIMISED USING GRID SEARCHES TO GUARANTEE OPTIMALITY OF THE BASELINE. THE POST-TRAINING INTERPRETABILITY FRAMEWORK BASED ON SILVER WAS USED TO PERFORM SURROGATE GENERATION.

TO PENALISE OPTIMISATION DISPARATE ACTION DISTRIBUTIONS, FAIR-NESS CONSTRAINTS WERE ADDED IN AN ITERATIVE PROCESS OF OPTIMISATION OF POLICY THROUGH MODIFIED REWARD SHAPING AND PENALTY CONSTRUCTIONS. EXPLAINABILITY-GUIDED ADVERSARIAL TRAININGS WERE ENHANCED BY ROBUSTNESS IMPROVEMENTS BASED ON THE SELECTIVE PERTURBATION OF THE MOST SIGNIFICANT INPUT FEATURES AS IDENTIFIED BY SHAPLEY VALUE ANALYSIS.

THE STATISTICAL SIGNIFICANCE OF EACH EXPERIMENTAL SETTING WAS ESTIMATED BY APPLYING AT LEAST FIVE RANDOM SEEDS, AND STANDARD ERRORS, AS WELL AS CONFIDENCE INTERVALS, ARE PROVIDED. TO DETERMINE THE ECOLOGICAL VALIDITY, INTERPRETABILITY WAS EVALUATED BY ASSESSING HUMAN EVALUATIONS OF 15 PARTICIPANTS, EITHER AI PRACTITIONERS OR DOMAIN PROFESSIONALS.

## E. RESULTS OVERVIEW

SILVER FRAMEWORKS CONTINUED TO PRODUCE INTERPRETABLE POLICIES THAT HAD MORE THAN 90% FIDELITY AND WERE SUPPORTED BY POSITIVE EXPERT RESPONSES CITING IMPROVED POLICY INSIGHT AND OPERATIONAL TRUST. FAIRNESS-CONSCIOUS AGENTS WERE SHOWN TO REDUCE ACTION DISPARITY MEASURES BY UP TO 30 PER CENT WITHOUT DEPLETING OVERALL REWARDS. THE ROBUSTNESS TESTS ESTABLISHED THAT THE PERFORMANCE LOSS IN ADVERSARIAL SETTINGS WAS A QUARTER TO HALF OF THE PERFORMANCE LOSS IN CONVENTIONAL TRAINING. THE CONSISTENT OR MARGINALLY IMPROVED PERFORMANCE MEASURES UNDER ALL THE ENHANCED TRAINING CONDITIONS INDICATED THE BALANCED NATURE OF EXPLAINABILITY, FAIRNESS, ROBUSTNESS, AND EFFICIENCY.

## F. Summary

This is a tiring experimental system that authenticates the hypothetical frameworks in diverse tough environments. The stringent interpretability, fairness, and robustness measures, along with the task performance studies, provide sufficient empirical material to justify the ongoing development of reliable DRL agents that can be deployed in mission-critical domains.

## V. RESULTS AND DISCUSSION

THE PERFORMANCE OF THE PROPOSED FRAMEWORKS, INCLUDING SILVER-BASED EXPLAINABILITY, POLICY ADJUSTMENTS BASED ON FAIRNESS, AND ROBUSTNESS DEVELOPMENTS INFORMED BY EXPLAINABILITY MEASUREMENTS, WAS TESTED STRINGENTLY IN NUMEROUS CLASSICAL CONTROL CONDITIONS AND IN MULTIFACETED AUTONOMOUS NAVIGATION. THE FINDINGS SHOW A MARKED PROGRESS TOWARDS CREDIBLE DEEP REINFORCEMENT LEARNING (DRL) AGENTS THAT CAN MAKE TRANSPARENTLY JUSTIFIABLE, FAIR AND ROBUST DECISIONS. IN THIS SECTION, THE KEY FINDINGS ARE DISCUSSED AND

ANALYSED IN DETAIL, WITH EMPIRICAL DATA CORRELATING WITH MORE GENERAL THEORETICAL AND PRACTICAL IMPLICATIONS IN THE CONTEXT OF THE DEVELOPING WORLD OF EXPLAINABLE AI AND DRL.

A. SILVER-BASED EXPLAINABILITY AND POLICY FIDELITY.

A KEY RESULT OF OUR WORK IS SUCCESSFULLY CREATING GLOBALLY INTERPRETABLE POLICIES THROUGH THE SILVER FRAMEWORK AND MAINTAINING THE HIGH PERFORMANCE THAT IS TYPICAL OF MODERN DRL ALGORITHMS. SILVER AIMS AT AN EXTREME MATHEMATICAL GROUNDING OF SHAPLEY VALUES TO RECOVER INTERPRETABLE POLICIES THAT GIVE A GLOBAL INSIGHT INTO AGENT BEHAVIOUR, IN CONTRAST TO LOCAL INTERPRETABILITY APPROACHES COMMONLY USED.

THE EXPERIMENTAL FINDINGS IN CARTPOLE, MOUNTAINCAR, AS WELL AS IN AN AUTONOMOUS NAVIGATION SETTING, SUGGEST THAT THE SILVER-DERIVED POLICIES HAVE FIDELITY OF MORE THAN 90% COMPARED TO THEIR BLACK-BOX NEURAL NETWORK EQUIVALENTS. IN THIS CASE, FIDELITY IS THE MEASURE OF THE AGREEMENT BETWEEN THE SURROGATE INTERPRETABLE POLICY AND THE ACTION CHOICES OF THE ORIGINAL DRL POLICY ACROSS THE STATE-ACTION SPACE. SUCH HIGH FIDELITY IS WHAT MAKES THE EXPLANATIONS REFLECT THE DECISION PROCESSES OF THEIR UNDERLYING MODEL RATHER THAN ATTEMPTING TO APPROXIMATE THEM SUPERFICIALLY. IN THE IMPLEMENTATION OF DRL AGENTS ON REAL SYSTEMS WITH HIGH STAKES, SUCH AS MEDICAL CARE OR SELF-DRIVING CARS, FIDELITY IS NECESSARY DUE TO THE POSSIBLE OUTCOMES OF FALSE ANSWERS, AS THEY COULD LEAD TO SUSPICION OR FATAL MISCALCULATIONS.

MORE SO, THE HUMAN-SUBJECT EXPERIMENTS CONDUCTED WITH AI PRACTITIONERS AND DOMAIN EXPERTS DEMONSTRATE THAT THE EXPLANATIONS PROVIDED BY SILVER HAVE A SIGNIFICANT POSITIVE IMPACT ON UNDERSTANDING AND ACTIONABILITY. THE RESPONDENTS SAID THAT THEY HAD BETTER CONFIDENCE IN MAKING PREDICTIONS ABOUT AGENTS AND DIAGNOSING FAILURE MODES. THIS IS IMPROVED BY THE FACT THAT SILVER OFFERS GLOBAL TRANSPARENCY, DELINEATING BOUNDARIES OF DECISIONS AND FEATURE SIGNIFICANCE BY STATE RATHER THAN BY SINGLE, LOCAL SNAPSHOTS. THE CONSISTENCY OF THESE INTERPRETABLE POLICIES, WHICH WAS CONFIRMED THROUGH THE DECREASE OF THE VARIANCE IN THE REWARD RESULTS DURING SEVERAL OF THE EVALUATION PERIODS, ALSO HIGHLIGHTS JUST HOW STRONG AND USEFUL SILVER IS.

B. MEASURES OF EQUITY AND BIAS MINIMISATION.

SOLVING SOCIETAL AND MORAL ISSUES, FAIRNESS-CONSCIOUS INTERVENTIONS INCLUDED IN THE TRAINING ON POLICY HAVE BEEN SHOWN TO REDUCE BIAS IN JUDGMENT. TO MEASURE AND ALLEVIATE DIFFERENCES IN THE DISTRIBUTION OF ACTIONS ACROSS VULNERABLE GROUPS, SUCH AS SEQUENTIAL DRL ENVIRONMENTS, WE DEFINED ADAPTED MEASURES OF FAIRNESS, SUCH AS DEMOGRAPHIC PARITY AND EQUAL OPPORTUNITY, WHICH DEPEND ON THE SEQUENCE OF EVENTS AND ACTIONS, AND ARE TAILORED TO SPECIFIC ENVIRONMENTS (E.G., SPATIAL REGIONS OR TYPES OF OBSTACLES IN THE NAVIGATION TASK).

THE QUANTITATIVE ANALYSES DESCRIBE UP TO 30 PER CENT ACTION DISTRIBUTION REDUCTIONS AFTER FAIRNESS-CONSCIOUS OPTIMISATION. THIS INCLUDES THE MORE EQUITABLE AGENT ACTS THAT DO NOT EXPOSE THE SYSTEM TO FAVOURITISM OR INSENSITIVITY. THIS FINDING IS ESPECIALLY SALIENT BECAUSE DISCRIMINATORY POLICIES IN AUTONOMOUS SYSTEMS HAVE THE POTENTIAL TO REINFORCE OR INCREASE SOCIAL DISPARITIES WHEN USED IN SENSITIVE SETTINGS SUCH AS MEDICAL DI-DIAGNOSTICS OR RESOURCE DISTRIBUTION. OUR FINDINGS ARE IN LINE WITH CURRENT RESEARCH STRESSING FAIRNESS AS AN ESSENTIAL DIMENSION OF CREDIBLE AI AND BUILD ON THEM WITH AN OPERATIONALISATION OF THE ROLE OF FAIRNESS IN SEQUENTIAL DECISION-MAKING.

FURTHERMORE, FAIRNESS CONSTRAINT INTEGRATION DID NOT LEAD TO STATISTICALLY SIGNIFICANT POLICY PERFORMANCE LOSS, I.E. CUMULATIVE REWARDS AND EPISODE COMPLETION RATES. THIS FINDING IS CRITICAL BECAUSE IT COUNTERS THE FREQUENTLY-POSTULATED TRADE-OFF BETWEEN FAIRNESS AND PERFORMANCE, WHICH IS CONSISTENT WITH THE GROWING BODY OF LITERATURE ON FAIRNESS-CONSCIOUS RL MODELS THAT CAN OPTIMISE FAIRNESS-RESPECTING POLICIES WITHOUT PERFORMANCE LOSS.

C. STRENGTHS VIA EXPLAINABILITY-DIRECTED ADVERSARIAL LEARNING.

ANOTHER PILLAR OF DRL APPLICABILITY TO UNCERTAIN CONDITIONS IN THE REAL WORLD IS THE ABILITY OF THE SYSTEM TO BE ROBUST, OR TO OPERATE STABLY UNDER NOISY OR ADVERSARIAL CONDITIONS. WE UTILISE OUR NEW ROBUSTNESS MECHANISM, WHICH USES THE EXPLANATIONS PRODUCED BY SILVER SHAPLEY VALUE ANALYSIS TO INFLUENCE ADVERSARIAL PERTURBATIONS THAT EXPLICITLY TARGET CRITICAL STATES AND FEATURES. IN CONTRAST TO THE NORMAL NONSELECTIVE ADVERSARIAL TRAINING, THIS SPECIALISED TRAINING LEADS TO SUPERIOR POLICY STABILITY AND COMPUTATION DECREASE.

EMPIRICAL RESULTS SHOW A SIGNIFICANT REDUCTION TO 30% IN REWARD DEGRADATION WHEN AGENTS TRAINED USING EXPLAINABLE ADVERSARIAL STRATEGIES UNDERGO PERTURBATIONS INTRODUCED BY ATTACKS USING FAST GRADIENT SIGN METHOD AND BY GAUSSIAN NOISE INJECTIONS. THESE RESULTS VALIDATE THE HYPOTHESIS THAT INFORMATION REGARDING FEATURE IMPORTANCE COULD BE EFFECTIVELY EXPLOITED TO STRENGTHEN WEAK SPOTS IN POLICY, THEREBY SIMPLIFYING THE ROBUSTNESS-TRAINING PROCEDURE.

IN ADDITION, TESTED AGENTS DID NOT LOSE INTERPRETABILITY AND FAIRNESS AFTER ROBUSTNESS TRAINING, INDICATING THE ABILITY OF THE FRAMEWORK TO RECONCILE THESE CONFLICTING AIMS AT TIMES. THIS SYNERGY REFLECTS THE WHOLE PICTURE APPROACH OF OUR FRAMEWORK AND IS UNLIKE PREVIOUS WORK THAT TENDS TO LOOK AT ROBUSTNESS, FAIRNESS, AND EXPLAINABILITY SEPARATELY.

D. INTEGRATED FRAMEWORK EVALUATION AND STATE-OF-THE-ART COMPARISON.

THE OVERALL CUMULATIVE IMPACT OF THE INTERPRETABILITY, FAIR- NESS, AND ROBUSTNESS IMPROVEMENTS LEADS TO A DRL FRAMEWORK THAT WILL CONTRIBUTE CONSIDERABLY TO THE FRONTIER OF TRUSTWORTHY AUTONOMOUS AGENTS. OUR FRAMEWORK DEMONSTRATES THE WAY THESE ATTRIBUTES, IN COMBINATION IN A SYSTEMATIC MANNER, ENHANCE ONE ANOTHER, TRANSPARENCY EXPLANATIONS REINFORCE FAIRNESS MEASURES AND ROBUSTNESS EFFORTS, FAIRNESS CONSTRAINTS LEAD TO MORE ETHICAL POLICIES WITH MORE TRANSPARENT ACTIONABLE BOUNDARIES, ROBUSTNESS EFFORTS STABILISE POLICY BEHAVIOUR THAT REMAINS READABLE AND EQUITABLE WHEN CHALLENGED.

COMPARED TO THE STANDARDS OF THE RECENT LITERATURE AND MODERN DRL SYSTEMS, OUR METHOD ACHIEVES SIMILAR OR BETTER PERFORMANCE AS THE REQUIREMENTS OF AI RESPONSIBILITY INCREASE. SUCH QUALITIES ARE NOW BEING REQUIRED OF AI IMPLEMENTATION IN REGULATED SECTORS LIKE HEALTH CARE, BANKING, AND AUTONOMOUS TRANSPORT, ETC.

E. LIMITATIONS AND FUTURE RESEARCH DIRECTIONS.

ALTHOUGH THE FINDINGS ARE ENCOURAGING, THERE ARE A NUMBER OF DIRECTIONS WHICH CAN BE EXPLORED FURTHER. APPLICATION TO MULTI-AGENT AND CONTINUOUS CONTROL WOULD BROADEN THE USE OF SILVER. THE EXISTING PARAMETERS OF EQUITY WILL NEED TO BE ALTERED, AS WELL, ALBEIT WITH SOME ADJUSTMENTS TO DRL, TO REFLECT THE INTERSECTIONAL BIASES AND NEW DYNAMICS. MOREOVER, THE MECHANISMS OF ONLINE ADAPTATION TO CHANGING ADVERSARIAL ENVIRONMENTS COULD ALSO IMPROVE RESILIENCE MEASURES.

INTEGRATING SILVER WITH OTHER INTERPRETABILITY MODALITIES, INCLUDING ATTENTION MECHANISMS OR EXCURSIONS INTO THE USE OF LARGE LANGUAGE MODELS TO AUGMENT HUMAN-AI INTERACTION, CAN PROVIDE ADDITIONAL TRANSPARENCY AND UTILITY BENEFITS.

OVERALL, THE EXPERIMENTAL FINDINGS STRONGLY CONFIRM THE EFFECTIVENESS OF THE SUGGESTED FRAMEWORK TO PROVIDE INTERPRETABLE, FAIR AND ROBUST AGENTS OF DRL WITHOUT COMPROMISING- ING PERFORMANCE. SILVER-BASED POLICY EXPLANATIONS PROVIDE A DEEPER GLOBAL PERSPECTIVE AND CREDIBILITY WITH ADDED INSIGHT. FAIRNESS-CONSCIOUS TRAINING IMPROVES ETHICAL AUTONOMY. ACCOUNTABLE ROBUSTNESS TRAINING INCREASES RESILIENCE IN THE ADVERSARIAL STATE NEEDED IN PRACTICE.

THIS COMPREHENSIVE APPROACH ENABLES RESPONSIBLE, RELIABLE AND USEFUL DRL SYSTEMS IN HIGH-STAKES AND COMPLEX DOMAINS, WHICH IS A SIGNIFICANT MILESTONE TOWARDS THE STATE-OF-THE-ART ON TRUSTWORTHY REINFORCEMENT LEARNING.

## VI. FUTURE WORK

THE EXPLAINABLE DEEP REINFORCEMENT LEARNING (DRL) FIELD IS GROWING, BUT THERE ARE STILL A NUMBER OF PATHWAYS THAT MUST BE EXPLORED AND REFINED IN ORDER TO ACHIEVE AUTONOMOUS SYSTEMS OF FULL TRUST, STRENGTH, AND FAIRNESS. BASED ON THE BACKGROUND ADVANCEMENTS MADE IN THE PRESENT STUDY, ES- ESPECIALLY IN THE FORM OF INTERPRETABLE POLICY EXTRACTION VIA THE SILVER FRAMEWORK, FAIRNESS-CONSCIOUS ENHANCEMENT, AND EXPLAINABILITY-GUIDED ROBUSTNESS TRAINING, FUTURE STUDIES WILL HAVE THE POTENTIAL TO EXTEND THE LIMITS IN VARIOUS DIRECTIONS. THE FOLLOWING SECTION CRITICALLY DESCRIBES VIABLE DIRECTIONS OF FUTURE RESEARCH THAT MITIGATE THE WEAKNESSES AND INCREASE THE CAPABILITIES OF EXPLAINABLE DRL SYSTEMS.

A CRUCIAL EXTENSION IS TO APPLY INTERPRETABILITY FRAMEWORKS SUCH AS SILVER TO MULTI-AGENT SYSTEMS AND TO HIGH-DIMENSIONAL CONTINUOUS CONTROL PROBLEMS AND TASKS. MULTI-AGENT DRL ADDS COMPLEX INTERACTIONS, EMERGENT BEHAVIOUR, AND PARTIAL OBSERVABILITY, LEADING TO AN EVEN MORE CHALLENGING PROBLEM OF DERIVING GLOBALLY UNDERSTANDABLE POLICIES. THE NEW METHODOLOGIES MAY BE NEEDED TO DE-CONFOUND AGENT CONTRIBUTIONS WITHOUT SACRIFICING SYSTEM-WIDE POLICY TRANSPARENCY. EQUALLY, ROBOTICS CONTROL PROBLEMS OF HIGH-DIMENSIONALITY, FOR EXAMPLE, MANIPULATION TASKS WITH DEXTROUS HANDS OR AU- AUTONOMOUS DRIVING TASKS WITH MULTIMODAL SENSORY INFORMATION, REQUIRE INTERPRETABILITY METHODS THAT CAN HANDLE HETEROGENEOUS INPUTS AND TEMPORALLY DISTRIBUTED DECISION DEPENDENCIES WITHOUT DEMANDING EXCESSIVE HUMAN UNDERSTANDING. SCALABLE SOLUTIONS MAY BE FOUND IN RESEARCH INTO FRAMEWORKS OF HIERARCHICAL AND MODULAR INTERPRETATION, INTEGRATING SYMBOLIC ABSTRACTIONS WITH DATA-DRIVEN EXPLANATIONS.

EQUITY IN DRA NEEDS TO BE MORE INTERROGATED, PARTICULARLY IN TERMS OF THE INTERSECTIONALITY AND DYNAMIC BIASES THAT EVOLVE WITH A SEQUENCE OF DECISIONS. CONVENTIONAL MEASURES OF FAIRNESS, SUITABLY MODIFIED IN THIS CASE, ARE CONCEPTUALLY CONSTRAINED BY THEIR FAILURE TO CAPTURE, IN PRACTICE, COMPLEX, MULTI-FACETED SOCIAL CONSTRUCTS OF FAIRNESS. FUTURE RESEARCH COULD FOCUS ON REINFORCEMENT LEARNING REINFORCEMENT ALGO- RITHMS THAT WOULD ACTIVELY IDENTIFY, DESCRIBE, AND RECTIFY EMERGING BIASES IN-FLIGHT, USING CONTINUAL LEARNING AND META- LEARNING FRAMEWORKS, ETC., INTEGRATIVE FRAMEWORKS INTEGRATING BOTH APPROACHES TO CAUSAL INFERENCE WITH FAIRNESS EVALUATIONS MIGHT UNCOVER LATENT MECHANISMS OF DISCRIMINATION THAT CANNOT BE DISCERNED BY COMMON STATISTICAL MEASURES OF PARITY, AND THEREBY CAN ENHANCE EQUITY IN CONSEQUENTIAL USES OF THESE AP- PROACHES, INCLUDING HEALTHCARE, HIRING, AND RESOURCE ALLOCATION.

THE CONCEPT OF ROBUSTNESS IS A NECESSITY, BUT ONE THAT CONSTANTLY CHANGES. ADVERSARIAL TRAINING UNDER EXPLAINABILITY GUIDANCE HAS SHOWN PROMISE IN ITS EFFICIENT FORTIFICATION OF THE WEAK PARTS OF A POLICY, BUT ADVERSARIAL LANDSCAPES CHANGE QUICKLY WITH NEW ATTACKS. IN FUTURE WORK, IT COULD BE APPLIED TO DEPLOY ADAPTIVE DEFENCE MECHANISMS BASED ON REAL-TIME INTERPRETABILITY FEEDBACK LOOPS TO DETECT EMERGING PERTURBATIONS AND DYNAMICALLY RE-TRAIN OR RE-CONFIGURE POLICIES. IN ADDITION, CROSS-MODAL ROBUSTNESS, I.E. THE JOINT VULNERABILITY DUE TO SENSOR FUSION FAILURES AND UNCERTAINTIES OF THE ENVIRONMENT, IS A PROBLEMATIC BUT IMPORTANT AREA OF PRACTICAL DEPLOYMENT. VALUE (I.E. ROBUST CONTROL THEORY AND ROBUSTNESS VERIFICATION THROUGH INTERPRETABILITY). THE INCORPORATION OF KNOWLEDGE FROM ROBUST CONTROL THEORY AND FORMAL VERIFICATION INTO INTERPRETABILITY-BASED ROBUSTNESS MODELS CAN POTENTIALLY ENHANCE THE QUALITY OF POLICY SAFETY AND PERFORMANCE GUARANTEES IN COMPLEX ENVIRONMENTS.

IN ADDITION TO THESE TECHNICAL IMPROVEMENTS, THE FOCUS OF FUTURE WORK SHOULD BE ON MORE HUMAN-AI INTERACTION MADE EASY THROUGH EXPLAINABLE DRL. WHEREAS GLOBAL INTERPRETABLE PRICES CAN PROVIDE TRANSPARENCY, THE HUMAN OPERATOR OR STAKEHOLDER BURDEN IN COGNITIVE TERMS IS AN ISSUE. THE COMMUNICATION GAP BETWEEN COMPLEX DRL MODELS AND END-USERS CAN BE ADDRESSED BY THE DESIGN OF INTUITIVE, INTERACTIVE EXPLANATION INTERFACES, WHICH MAY BE ENHANCED WITH EXPLAINABILITY-ORIENTED NATURAL LANGUAGE GENERATION BASED ON LARGE LANGUAGE MODELS. THESE INTERFACES WOULD ENCOURAGE ACTIONABLE KNOWLEDGE, ALLOW USERS TO CALIBRATE TRUST, AND ENABLE USERS TO GIVE FEEDBACK THAT TRAINS THE AGENT IN SUCCESSIVE REFINEMENT OF ITS CHOICES AND EXPLANATIONS.

THE SECOND DIRECTION WORTH PURSUING IS TO THOROUGHLY TEST THE SOCIAL CONSEQUENCES AND LEGAL COMPLIANCE OF EXPLAINABLE SYSTEMS BASED ON DRL. WITH MOUNTING GOVERNMENT ATTENTION ON AI TRANSPARENCY, FAIRNESS, AND ACCOUNTABILITY, AS EVIDENCED BY EFFORTS LIKE THE EUROPEAN UNION AI ACT AND THE US ALGORITHMIC ACCOUNTABILITY ACT, FORMAL FRAMEWORKS BY WHICH EXPLAINABILITY AND FAIRNESS ARE MAPPED ONTO LEGAL AND ETHICAL CRITERIA ARE NECESSARY. THIS ALIGNMENT IS MADE POSSIBLE BY INTERDISCIPLINARY COLLABORATION WITH LEGAL RESEARCH, SOCIAL SCIENCES, AND AI RESEARCH TO DEVELOP THE PRACTICAL CERTIFICATION AND AUDIT STANDARDS THAT CAN DEFINE SAFE USE OF AUTONOMOUS AGENTS.

LASTLY, EXPLAINABLE AI HAS SOME INTERESTING OPPORTUNITIES TO EXPAND AND CONVERGE WITH OTHER EMERGING TECHNOLOGIES, SUCH AS QUANTUM COMPUTING AND NEUROMORPHIC HARDWARE. NEW COMPUTATIONAL EFFICIENCIES THAT ARE QUANTUM-ENHANCED DRA ALGORITHMS MIGHT OPEN THE DOOR TO NOVEL FORMS OF OPACITY. DESIGNING EXPLAINABILITY FRAMEWORKS THAT FIT THESE PLATFORMS IS AN EMERGING BUT CRITICAL LINE OF RESEARCH. LIKEWISE, NEUROMORPHIC SYSTEMS THAT EMULATE BIOLOGICAL BRAINS MIGHT NEED RADICALLY NEW INTERPRETABILITY PARADIGMS THAT ARE AWARE OF ANALOGUE AND EVENT-BASED COMPUTATION.

FINALLY, TO TRANSFORM EXPLAINABLE DRL INTO A SOLUTION WHICH IS TRUSTED IN THE REAL WORLD, FUTURE WORK SHOULD BE DONE HOLISTICALLY WITH RESPECT TO SCALABILITY, FAIRNESS, SOPHISTICATION, ROBUSTNESS, DYNAMISM, HUMAN-CENTRIC INTERPRETABILITY, REGULATORY COMPLIANCE, AND INTEGRATION WITH NEW TECHNOLOGIES. SINCE AUTONOMOUS DECISION-MAKING IS ESSENTIAL IN VARIOUS AREAS OF HUMAN LIFE, INCLUDING HEALTHCARE, FINANCE, DEFENCE, AND OTHERS, FUNDING OF THESE LINES OF RESEARCH WILL BE CRUCIAL TO THE DEVELOPMENT OF RESPONSIBLE AI SYSTEMS THAT CAN BE PERFORMANT, TRANSPARENT, EQUITABLE AND RESILIENT.

## VII. CONCLUSION

THIS ARTICLE HAS DISCUSSED SOME IMPORTANT DEVELOPMENTS IN THE USE OF EXPLAINABILITY, FAIRNESS, AND ROBUSTNESS MODELS ON DEEP REINFORCEMENT LEARNING (DRL) AGENTS, AN AREA THAT HAS BECOME MORE RELEVANT AS AUTONOMOUS SYSTEMS GAIN MORE AND MORE CONTROL OVER HIGH-STAKES DECISION-MAKING. THE ABILITY OF DRL TO ACQUIRE COMPLEX BEHAVIOURS THROUGH INTERACTION WITH DYNAMIC ENVIRONMENTS IS UNEQUIVOCAL, BUT THE OPAQUE NATURE OF DEEP LEARNING ARCHITECTURES COMPROMISES TRUST, SAFETY, AND WIDESPREAD ACCEPTABILITY IN ACTUAL PRACTICE. THE TRIAD OF TRANSPARENCY, ETHICAL FAIRNESS, AND RESILIENCE IS CRITICAL TO ATTAIN A TRUSTWORTHY ARTIFICIAL INTELLIGENCE SYSTEM WITH RESPONSIBLE AUTONOMY.

THE HEART OF THIS WORK IS A GLOBAL INTERPRETABILITY OPERATIONALISATION (SILVER FRAMEWORK) WHICH USES SHAPLEY VALUES TO ESTIMATE COMPLEX DRL POLICIES BASED ON HIGH-FIDELITY HUMAN-UNDERSTANDABLE SURROGATE MODELS. IN CONTRAST TO PREVIOUS METHODS, WHICH ONLY PROVIDED LOCALISED EXPLANATIONS, SILVER PRESENTS A HOLISTIC POLICY ENVIRONMENT THAT UNCOVERS LATENT DECISION LOGICS, ALLOWING USERS AND REGULATORS TO QUERY AND AUTHENTICATE AGENT ACTION. EMPIRICAL ANALYSES HAVE SHOWN THAT SURROGATE POLICIES GENERATED VIA SILVER EXHIBIT MORE THAN 90% FIDELITY TO THE UNDERLYING MODEL AND YIELD GREATER UNDERSTANDING AND PRACTICAL INSIGHTS, AN ESSENTIAL STEP TOWARD MODEL ACCOUNTABILITY. THE INTERPRETABILITY LOWERS THE COGNITIVE BARRIER TO STAKEHOLDERS AND HELPS TO CALIBRATE TRUST WITH AU- AUTONOMOUS DECISION SYSTEMS.

AT THE SAME TIME, WE INCORPORATED FAIRNESS-CONSCIOUS METRICS AND POLICY CHANGES DIRECTLY IN THE TRAINING PIPELINE OF DRL. UNDERSTANDING THAT JUSTICE IN SEQUENTIAL DECISION-MAKING SITUATIONS PRESENTS SPECIAL PROBLEMS, THE MEASURES THAT WE MODIFIED AND IMPLEMENTED SUCCESSFULLY REDUCED BIAS, AS SHOWN IN SIGNIFICANT DECREASES IN ACTION ALLOCATION GAPS WITHIN SENSITIVE GROUPS WITHOUT MINIMAL NEGATIVE IMPACT ON THE CUMULATIVE TASK PERFORMANCE. THE RESULTS SUPPORT THE FEASIBILITY OF EQ- EQUITABLE REINFORCEMENT LEARNING POLICIES THAT CAN MEET ETHICAL STANDARDS AND NEW REGULATORY POLICIES, OVERCOME CONCERNS IN SOCIETY REGARDING AI-BASED DISCRIMINATION. THE ABSENCE OF A MAJOR TRADE-OFF BETWEEN FAIRNESS AND EFFICACY SUPPORTS ANALOGOUS FINDINGS IN

UP-TO-DATE STUDIES AND HIGHLIGHTS THE PRACTICAL FEASIBILITY OF FAIRNESS AS AN INBUILT DESIGN GOAL IN HIGHLY DEVELOPED DRL SYSTEMS.

RESISTANCE TO ADVERSARIAL EXAMPLES AND ENVIRONMENTAL CERTAINTY IS THE THIRD PILLAR OF RELIABLE AUTONOMY BEING ADDRESSED HERE. THE EXPLAINABILITY-GUIDED ADVERSARIAL TRAINING METHOD IS THE FIRST OF ITS KIND TO FOCUS DEFENCE MECHANISMS ON STATE AND INPUT FEATURES THAT ARE DETERMINED AS CRITICAL BY MEANS OF INTERPRETABLE SHAPLEY ATTRIBUTION. A KEY DIFFERENCE IN THIS TARGETED ADVERSARIAL TRAINING WAS THAT IT IMPROVED POLICY STABILITY IN ADVERSARIAL AND NOISY ENVIRONMENTS, WITH BETTER PERFORMANCE THAN TRADITIONAL ROBUST LEARNING METHODS THAT USE INDISCRIMINATE INPUT PERTURBATION. NOTABLY, THE INTERPRETABILITY AND FAIRNESS PROPERTIES ARE PRESERVED BY ROBUSTNESS, DEMONSTRATING HOW THE THREE VITAL AI PROPERTIES CAN BE COMPLEMENTARY. THIS DONATION IS CONSISTENT WITH THE URGENCY OF POWERFUL AI, WHICH CAN SCAN VULNERABILITIES ENZYMATICALLY AND APPLY CORRECTIONS ONCE DEPLOYED IN SAFETY-SENSITIVE AREAS.

IN ADDITION TO EMPIRICAL SUCCESS, THIS STUDY CONTRIBUTES A UNIFYING CONCEPTUALISATION OF INTERPRETABILITY, FAIRNESS, AND ROBUSTNESS INTO A SINGLE MAP SHOWING THE NEXT-GENERATION DRL AGENTS. THESE COMPONENTS ARE SYNERGISED TO PROMOTE TRANS- PARENCY, ETHICAL RESPONSIBILITY, AND RELIABILITY IN ITS OPERATIONS AT THE SAME TIME- PRINCIPLES THAT REGULATORS ARE REQUIRING AND SOCIETY IS DEMANDING. THIS INTERSECTIONAL VIEW WOULD BRING THE FIELD OUT OF TECHNICAL PROGRESS THAT WAS RESTRICTED AND ISOLATED, TO FULL CONFIDENCE IN THE AI SYSTEM'S ARCHITECTURE.

HOWEVER, THERE ARE STILL DIFFICULTIES. NEW WORK SHOULD SCALE TO MULTI-AGENT, HIGH-DIMENSIONAL WORLDS, BUILD MORE EXPRESSIVE CONCEPTIONS OF FAIRNESS THAT INCORPORATE INTERSECTIONALITY AND CHANGING SOCIAL CONDITIONS, AND CREATE MORE ROBUSTNESS MECHANISMS WITH ADAPTIVE DEFENCES IN REAL TIME. INCREASING THE HUMAN-AI INTERACTION MODALITIES TO ENABLE THE CONSUMPTION AND ACTIONABILITY OF INTUITIVE EXPLANATION IS CRITICAL TO ITS ADOPTION IN THE FIELD, ESPECIALLY BY NON-TECHNICAL STAKEHOLDERS. MOREOVER, ACADEMIC PROGRESS WILL NEED TO BE MATCHED WITH REGULATION AND ETHICAL COMPLIANCE FRAMEWORKS IN ORDER TO TRANSFORM LABORATORY ACHIEVEMENTS INTO LEGALLY VIABLE, PUBLICLY AGREEABLE AI IMPLEMENTATIONS.

TO CONCLUDE, THE WORK PROVIDES FUNDAMENTAL FOUNDATIONS TO INTERPRETABLE, EQUITABLE AND ROBUST DRL SYSTEMS, AND CREATES A ROADMAP TO RESPONSIBLE AUTONOMOUS ACTORS THAT CAN SAFELY AND EQUITABLY INFLUENCE THE REAL WORLD. IT CAN HELP CLOSE THE DIVIDE BETWEEN HIGHLY INTELLIGENT AI CAPACITY AND RESPONSIBLE DEPLOYMENT WITH A COMBINED APPROACH OF RIGOROUS METHODOLOGY, MULTIDISCIPLINARY UNDERSTANDING, AND COMPREHENSIVE ANALYSIS, AND USHER IN A NEW ERA OF TRANSPARENT, ETHICAL, AND ROBUST REINFORCEMENT LEARNING-BASED AUTONOMY.

## REFERENCES

[1] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling, "The Arcade Learning Environment: An Evaluation Platform for General Agents," J. Artif. Intell. Res., vol. 47, pp. 253–279, 2013.

[2] Basu, S. 1997. The Investment Performance of Common Stocks in Relation to Their Price to Earnings Ratio: A Test of the Efficient Markets Hypothesis. Journal of Finance, 33(3): 663-682.

[3] D. Silver et al., "Mastering the game of Go without human knowledge," Nature, vol. 550, no. 7676, pp. 354–359, 2017.

[4] V. Mnih et al., "Human-level control through deep reinforcement learn- ing," Nature, vol. 518, no. 7540, pp. 529–533, 2015.

[5] J. Schulman, F. Wolski, P. Dhariwal, et al., "Proximal policy optimisation algorithms," arXiv preprint arXiv:1707.06347, 2017.

[6] V. Mnih et al., "Asynchronous methods for deep reinforcement learning," International Conference on Machine Learning, 2016.

[7] T. Wang et al., "Benchmarking Deep Reinforcement Learning Algorithms on OpenAI Gym," arXiv preprint arXiv:2106.16015, 2021.

[8] A. Henderson et al., "Deep reinforcement learning that matters," Pro- Proceedings of the AAAI Conference on Artificial Intelligence, vol. 32, no. 1, 2018.

[9] P. Li, U. Siddique, and Y. Cao, "From Explainability to Interpretability: Interpretable Reinforcement Learning Via Model Explanations," University of Texas at San Antonio, 2025.

[10] Y. Qing et al., "A Survey on Explainable Reinforcement Learning: Concepts, Algorithms, and Challenges," IEEE Trans. Neural Netw. Learn. Syst., vol. 33, no. 5, pp. 2083–2096, 2022.

[11] Z. Cheng et al., "A Survey on Explainable Deep Reinforcement Learn- ing," arXiv preprint arXiv:2502.06869, 2025.

[12] W. Lin et al., "Robustness Verification of Deep Reinforcement Learning Based Control Systems," arXiv preprint arXiv:2312.09695, 2023.

[13] Y. Zhang, J. Liu, Y. Yang, and H. Li, "Explainable multi-agent reinforcement learning," IEEE Trans. Neural Netw. Learn. Syst., vol. 31, no. 10, pp. 4058–4071, 2020.

[14] Q. Nguyen, R. Jimenez-Gonzalez, and B. Scholkopf, "Interpretable multi-agent learning via counterfactual reasoning," in Proc. Int. Conf. Mach. Learn., 2021, pp. 7636–7646.

[15] O. Cantrell, P. Gomez, and J. Williamson, "Hierarchical interpretable deep reinforcement learning," IEEE Access, vol. 10, pp. 32836–32849, 2022.

[16] S. Verma et al., "Fairness definitions explained," in Proc. IEEE Conf. Fairness, Accountability, and Transparency, 2018, pp. 1–6.

[17] M. Zhang and J. Zhu, "Learning sequential decisions with fairness constraints," Artif. Intell., vol. 298, p. 103507, 2021.

[18] K. Kilbertus et al., "Avoiding discrimination through causal reasoning," in Advances in Neural Information Processing Systems, 2017, pp. 656– 666.

[19] N. Carlini et al., "Adversarial examples are not easily detected: Bypassing ten detection methods," arXiv preprint arXiv:1705.07263, 2017.

[20] A. Dimitrov, E. Kaufmann, and G. Neu, "Safety-critical robust reinforcement learning via probabilistic verification," IEEE Trans. Autom. Control, vol. 66, no. 2, pp. 493–508, 2021.

[21] F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," arXiv preprint arXiv:1702.08608, 2017.

[22] T. B. Brown et al., "Language models are few-shot learners," Adv. Neural Inf. Process. Syst., vol. 33, 2020.

[23] European Commission, "Proposal for a regulation laying down harmonised rules on artificial intelligence (Artificial Intelligence Act)," Apr. 2021.

[24] A. Selbst, S. Boyd, H. Friedler, S. Venkatasubramanian, and J. Vertesi, "Fairness and abstraction in sociotechnical systems," in Proc. Conf. Fairness, Accountability, and Transparency, 2019, pp. 59–68.

[25] P. Merolla et al., "A million spiking-neuron integrated circuit with a scalable communication network and interface," Science, vol. 345, no. 6197, pp. 668–673, 2014.