



# Curiosity-Driven Reinforced Learning of Undesired Actions in Autonomous Intelligent Agents

Christopher Rosser and Khalid Abed

Department of Electrical & Computer Engineering  
and Computer Science

Jackson State University

SAMI 2021

# Outline

1. Introduction
2. Related Work
3. Methodology
4. Results and Discussion
5. Conclusion and Future Work
6. References

# Introduction

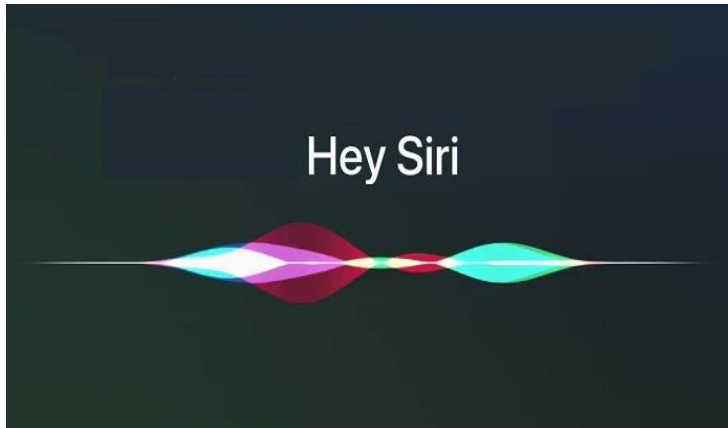
# Artificial Intelligence (AI) is everywhere!



Robotics – Warehouse, Cleaning, Delivery



Self-driving cars/Autonomous Underwater Vehicles



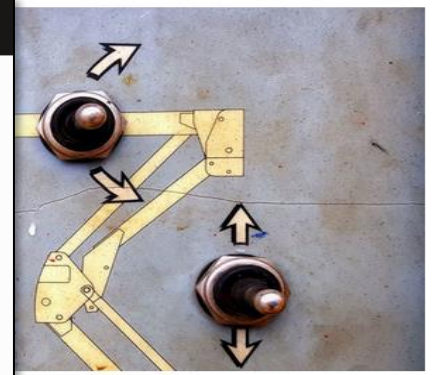
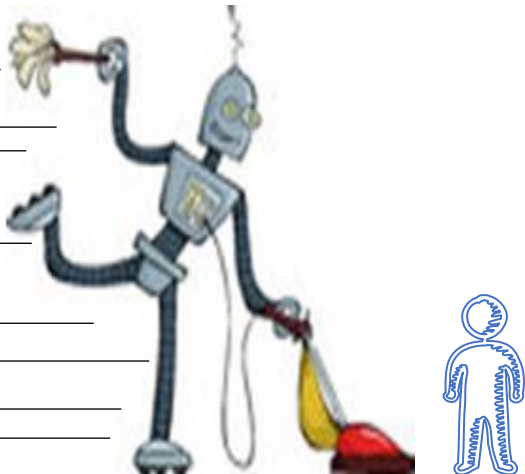
Speech Recognition (Natural Language Processing)



Video Game Playing

# Why research autonomously exploring intelligent agents?

- Increase in autonomy of AI systems has risks and can lead to accidents, posing physical threats to human AI users
- AI research is more accessible than ever
- Improvements in safe exploration may have benefits in other domains



# Why research autonomously exploring intelligent agents?

The following concrete problems in AI were provided in 2016 by researchers from Google Brain, OpenAI, UC Berkley, and Stanford in light of recent incidents

- Avoiding negative side effects
- Avoiding rewards hacking
- Scalable oversight
- Safe exploration
- Robustness to distributional shift

The same researchers assert that AI/ML accidents and risks are attributed to the following three “failures”:

- Having the wrong objective function
- Having an objective function that is too expensive to evaluate frequently
- Undesirable behavior during the learning process

# Why research autonomously exploring intelligent agents?

## Recommended Solutions for Improving Exploration Safety:

- Adversarial Blinding - preventing an agent from understanding how its reward is generated or blinding it to certain variables
- Trip Wires - introduce deliberate vulnerabilities and monitor them so that researchers are alerted and can stop the agents immediately if the vulnerabilities are exploited.
- Simulated Exploration
- Human Oversight



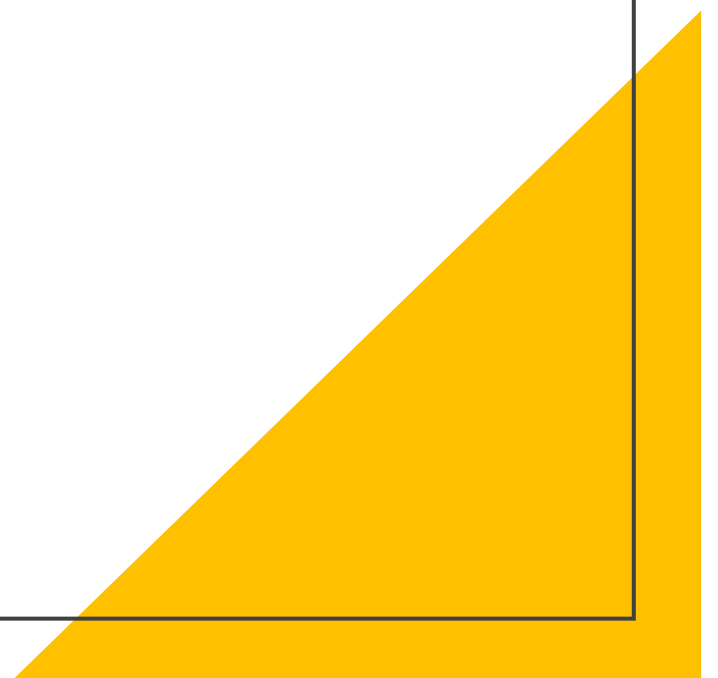
# Why research autonomously exploring intelligent agents?



- We want to improve the safety of autonomous exploring AI through simulated exploration and human oversight.
- In this research we use Unity's Machine Learning Agents Toolkit (ML-Agents) to train purposely misbehaving agents to determine when a human should intervene during the agent's learning process.

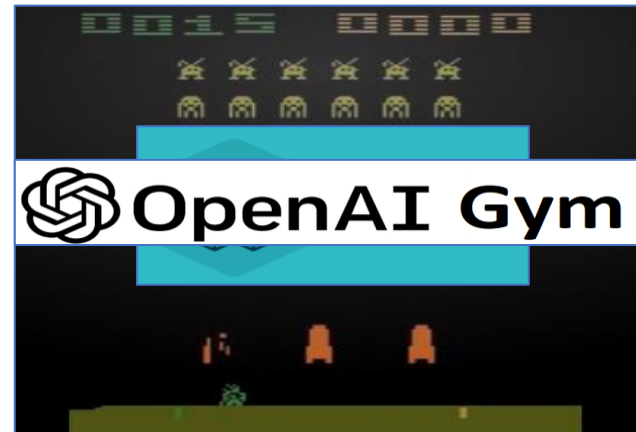


# Related Works



# Game Engines/Development Platforms as RL Environments for Investigating Autonomous Exploration

- VizDoom and Mujoco
  - Episodic curiosity-driven exploration
- OpenAI
  - Arcade Learning Environment for general intelligent agents
  - Intrinsic Curiosity Module (ICM) for curiosity-driven exploration
  - Attention-based curiosity-driven exploration



# Human-in-the-Loop Reinforcement Learning

- **Hard-coded Guidance**
  - Defining catastrophes and significant rare events before training
  - Environment-level action blockers
  - Determining where to add actions in RL
- **Actual/Learned Human Intervention**
  - Improving safety with model-based architectures and human intervention
  - Training agents to imitate human intervention
  - Runtime monitoring framework

# Methodology

A vertical line is positioned to the right of the word 'Methodology'. In the bottom right corner of the slide, there is a yellow triangle pointing upwards and to the left, partially overlapping a light gray border.

- Use Unity's Machine Learning Agents Toolkit (ML-Agents) implementation of Proximal Policy Optimization (PPO) + ICM algorithm to create purposely misbehaving autonomous exploring agents
- Identify PPO+ICM training statistics and custom environment metrics associated with agent misbehavior

# Training Agents in ML-Agents

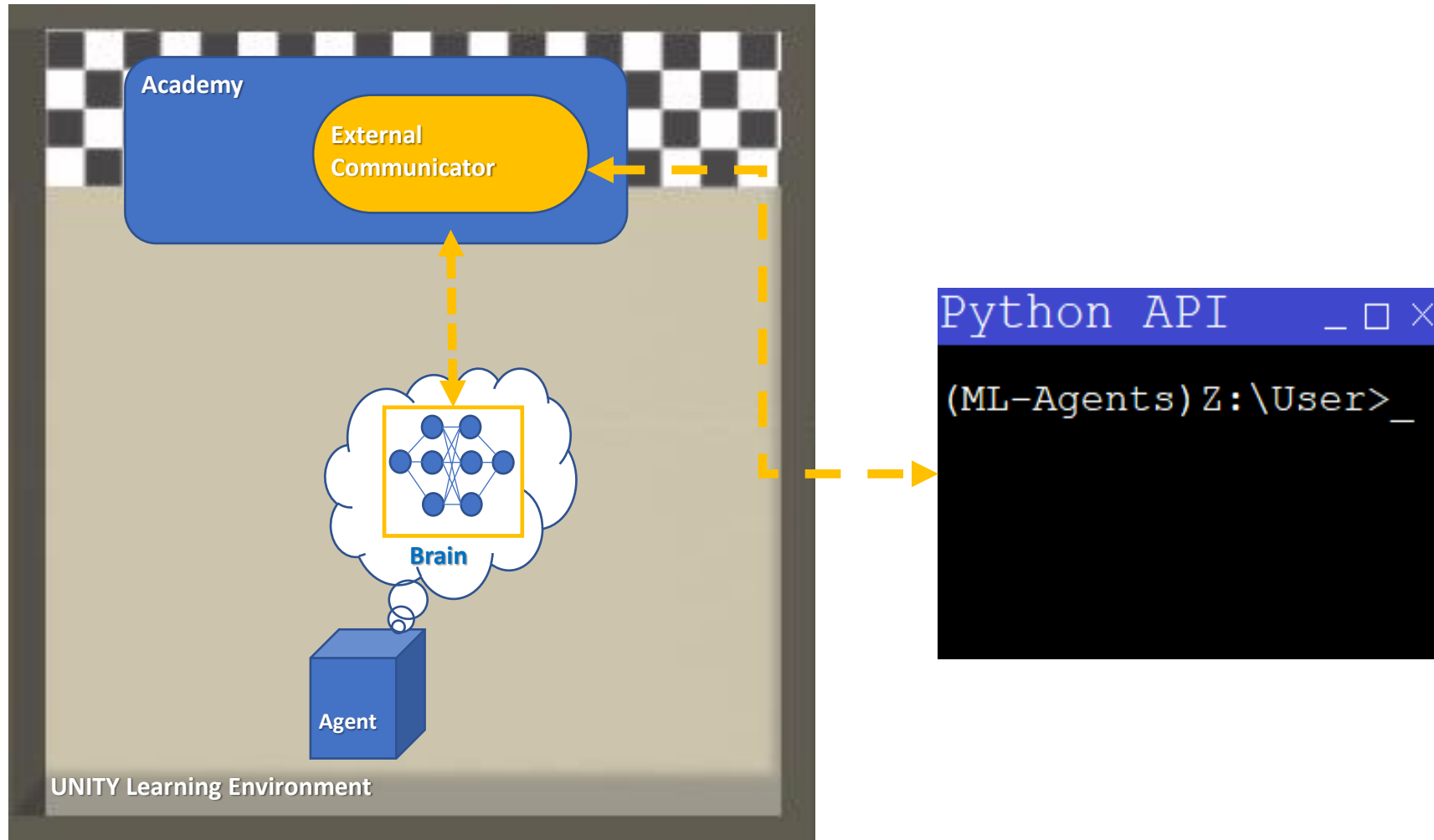


Illustration of externally training neural networks with ML-Agents

# PPO + ICM, Learning Environment, and Training Parameters

## PPO + ICM Algorithm

```
for iteration = 1, 2, ... do
  collect set of actions (a) and next states (s+1) with policy ( $\pi$ )
  encode current state (s), next state (s+1)  $\rightarrow \phi(s), \phi(s+1)$ 
  compute predicted encoded next state  $\hat{\phi}(s+1)$ 
  intrinsic reward  $\leftarrow \phi(s+1) - \hat{\phi}(s+1)$ 
  optimize  $\pi$  parameters for maximizing extrinsic + intrinsic rewards
  update  $\pi$ 
end for
```

## PPO + ICM TRAINING PARAMETERS

Parameter	Value
<i>gamma</i>	0.99
<i>lambda</i>	0.95
<i>buffer size</i>	1024
<i>batch size</i>	128
<i>epochs</i>	3
<i>learning rate</i>	0.0003
<i>time horizon</i>	64
<i>max steps</i>	150000
<i>beta</i>	0.01
<i>epsilon</i>	0.2
<i>hidden layers</i>	2
<i>hidden units</i>	128
<i>sequence length</i>	64
<i>memory size</i>	256
<i>curiosity encoding size</i>	128-256
<i>curiosity strength</i>	0.01-10.0





# Training Conditions and Experiments

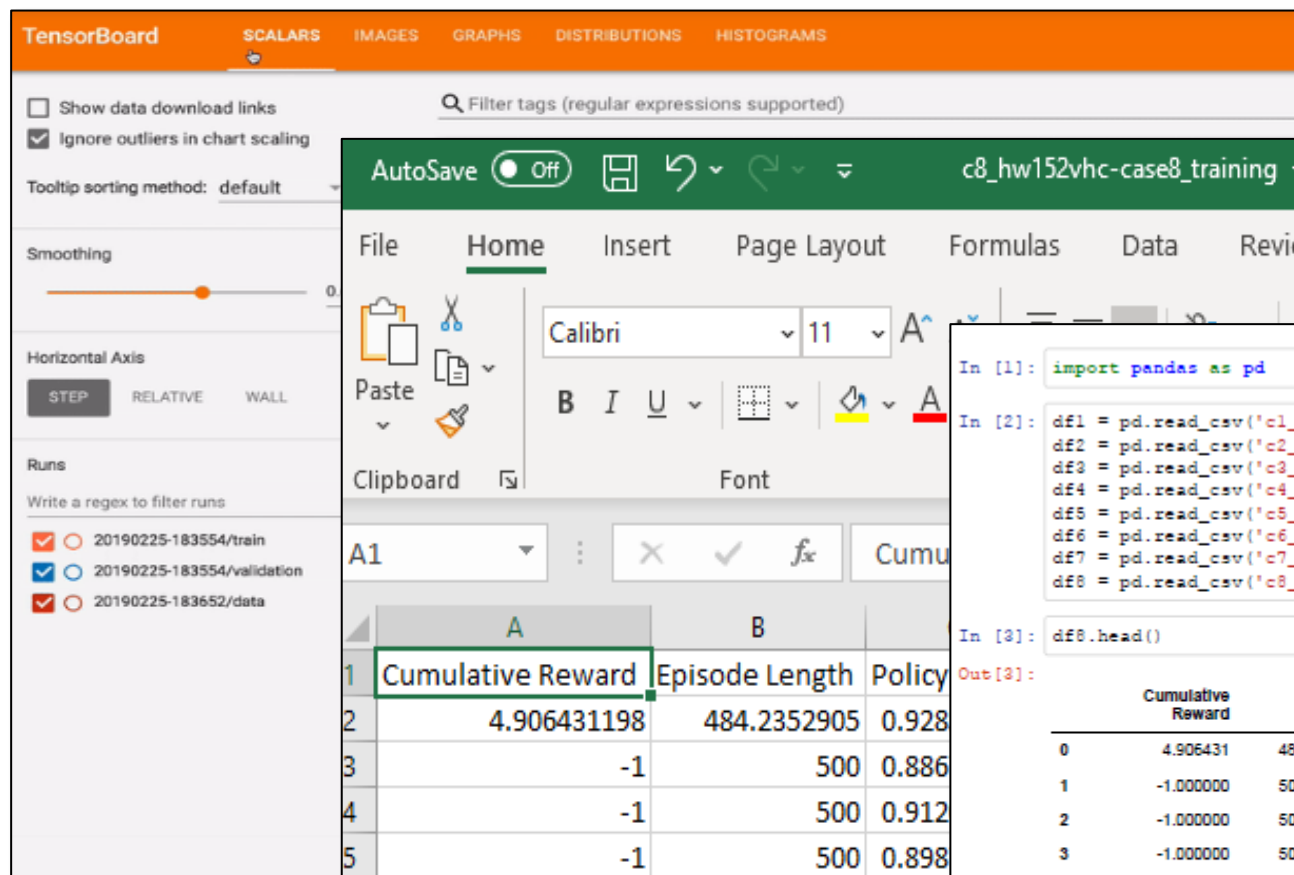
Train autonomous exploring agents with high intrinsic curiosity, large rewards, and large penalties for 150,000 timesteps under the following six conditions:

1. Low/Zero Curiosity Strength, Small Rewards, Large Penalty
2. Low/Zero Curiosity Strength, Large Rewards, Small Penalty
3. Max Recommended Curiosity Strength, Small Rewards, Large Penalty
4. Max Recommended Curiosity Strength, Large Rewards, Small Penalty
5. Very High Curiosity Strength, Small Rewards, Large Penalty
6. Very High Curiosity Strength, Large Rewards, Small Penalty

CURIOSITY STRENGTHS AND REWARDS FOR  
TEST CASES 1 - 8

Case	Curiosity Strength	Rewards	Penalties
1	Disabled	1	-100, -0.1
2	Disabled	100	-1, -0.1
3	0.1	1	-100, -0.1
4	0.1	100	-1, -0.1
5	1.0	100	-1, -0.1
6	1.0	1	-100, -0.1
7	10.0	1	-100, -0.1
8	10.0	100	-1, -0.1

# Analysis of ML-Agents Data



Microsoft Excel interface showing a table of training data. The table has columns A, B, and C. The data is as follows:

	A	B	C
1	Cumulative Reward	Episode Length	Policy
2	4.906431198	484.2352905	0.928
3	-1	500	0.886
4	-1	500	0.912
5	-1	500	0.898
6	-1	500	1.028
7	-1	500	1.002



Jupyter Notebook interface showing code and output for data analysis. The code includes importing pandas, reading CSV files, and displaying the head of a DataFrame. The output shows a table of training data.

```
In [1]: import pandas as pd
In [2]: df1 = pd.read_csv('c1_hw151_case1_training.csv')
df2 = pd.read_csv('c2_hw152_case2_training.csv')
df3 = pd.read_csv('c3_hw151_case3_training.csv')
df4 = pd.read_csv('c4_hw152_case4_training.csv')
df5 = pd.read_csv('c5_hw151hc-case5_training.csv')
df6 = pd.read_csv('c6_hw151vhc-case6_training.csv')
df7 = pd.read_csv('c7_hw152hc-case7_training.csv')
df8 = pd.read_csv('c8_hw152vhc-case8_training.csv')
In [3]: df8.head()
```

	Cumulative Reward	Episode Length	Forward Loss	Inverse Loss	Policy Loss	Value Loss	Curiosity Reward	Entropy	Learning Rate	Value Estimate
0	4.906431	484.235291	0.162945	1.530372	0.928554	104.850014	370.654266	1.520337	0.000299	10.540007
1	-1.000000	500.000000	0.061259	1.280338	0.886755	31.642797	297.494446	1.298747	0.000297	24.096262
2	-1.000000	500.000000	0.065908	0.916347	0.912093	22.592175	273.758209	1.034323	0.000295	32.691975
3	-1.000000	500.000000	0.062123	0.828616	0.898753	14.585310	277.177490	1.026057	0.000293	39.444569
4	-1.000000	500.000000	0.058550	0.737338	1.028775	10.561207	265.727387	0.893187	0.000291	43.837360

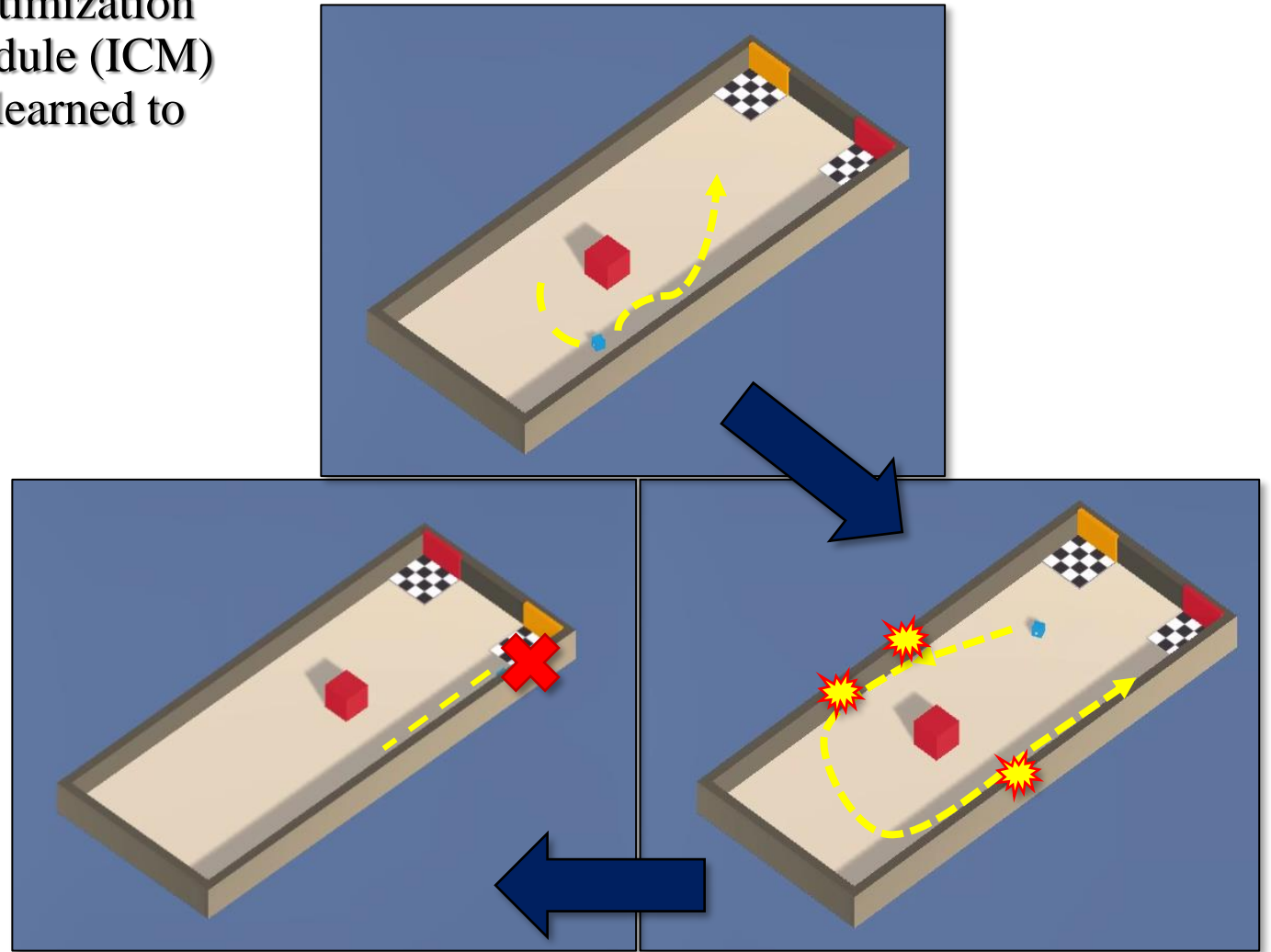
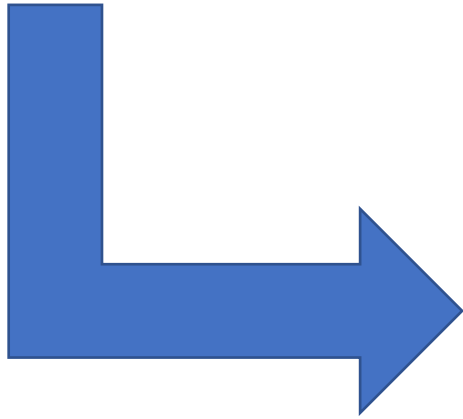
```
In [4]: #testing to make sure columns are dropped
df4.sample()
```

	Cumulative Reward	Episode Length	Forward Loss	Inverse Loss	Policy Loss	Value Loss	Curiosity Reward	Entropy	Learning Rate	Value Estimate
94	-0.982908	489.758606	0.066824	0.548569	0.570521	0.389894	3.691452	1.204852	0.000113	0.415393



# Results and Discussion

Agents trained using the Proximal Policy Optimization (PPO) algorithm with Intrinsic Curiosity Module (ICM) enabled and large rewards or large penalties learned to act undesirably within 150,000 timesteps.



- Implemented a collision counter and goal-to-collision ratio to identify three test cases for further analysis

AVERAGE VALUES FOR AGENT WITH  
DEFAULT NEURAL NETWORK

Goals	Collisions	GC-Ratio
52	2.56	0.9566

AVERAGE VALUES FOR AGENT WITH CASES 1-8  
NEURAL NETWORKS ATTACHED

Case	Goals	Collisions	GC-Ratio
1	0	12.66	0.001
2	1.333	183.33	0.008
3	0.333	140.67	0.003
4	0	418.33	0.001
5	0.333	268.67	0.006
6	0	362.33	0.001
7	0	436.33	0.001
8	0	170.67	0.001

## Case 6

Two strongest negative correlations

(1) Episode Length v. Cumulative Reward, (2) Value Loss v. Value Estimate

Third strongest

(3) Value Estima

Two strongest p

(1) Entropy v. Va

Third strongest

(3) Learning Rate

```
In [18]: df_norm6.corr()
```

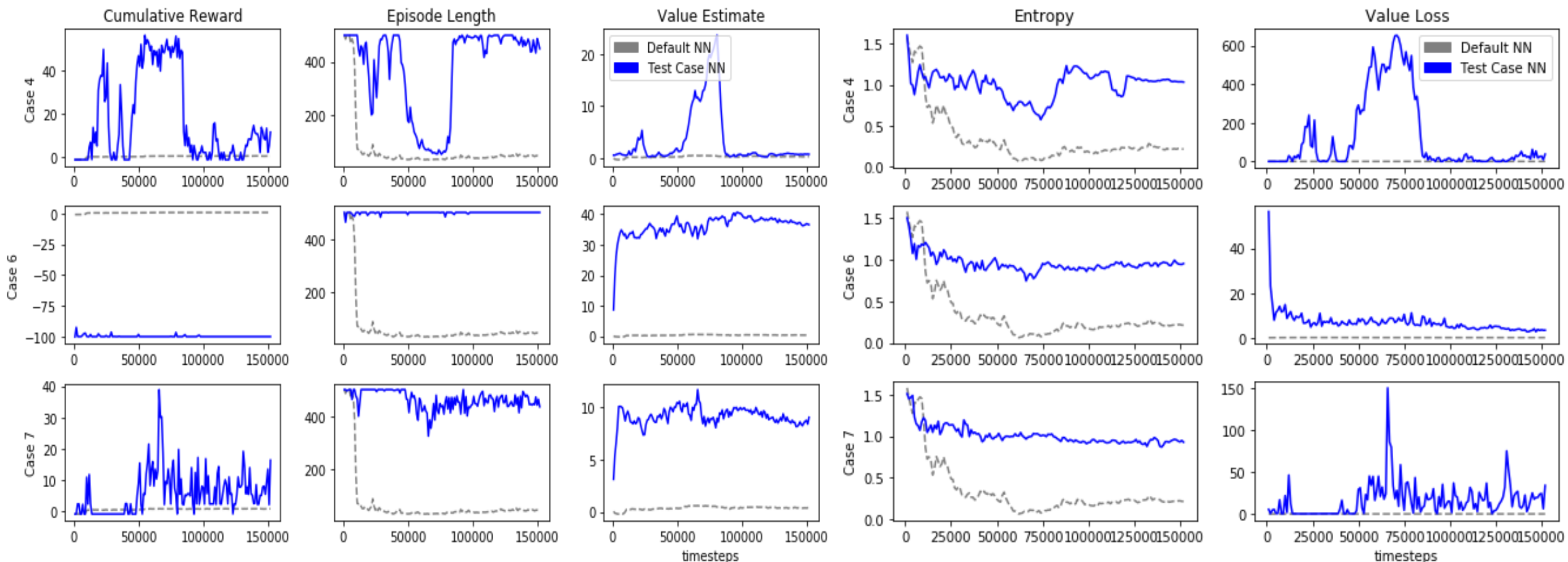
```
Out[18]:
```

	Cumulative Reward	Episode Length	Policy Loss	Value Loss	Entropy	Learning Rate	Value Estimate
Cumulative Reward	1.000000	-0.999781	0.228822	0.331685	0.434574	0.274277	-0.274521
Episode Length	-0.999781	1.000000	-0.232948	-0.346523	-0.442640	-0.276359	0.285040
Policy Loss	0.228822	-0.232948	1.000000	0.351071	0.554644	0.159015	-0.320367
Value Loss	0.331685	-0.346523	0.351071	1.000000	0.643581	0.529072	-0.801056
Entropy	0.434574	-0.442640	0.554644	0.643581	1.000000	0.470622	-0.665674
Learning Rate	0.274277	-0.276359	0.159015	0.529072	0.470622	1.000000	-0.593451
Value Estimate	-0.274521	0.285040	-0.320367	-0.801056	-0.665674	-0.593451	1.000000

Average Cumulative Reward vs. Episode Length Correlation of

**-0.96!**

# PPO+ICM for Hallway Cases 4, 6, and 7



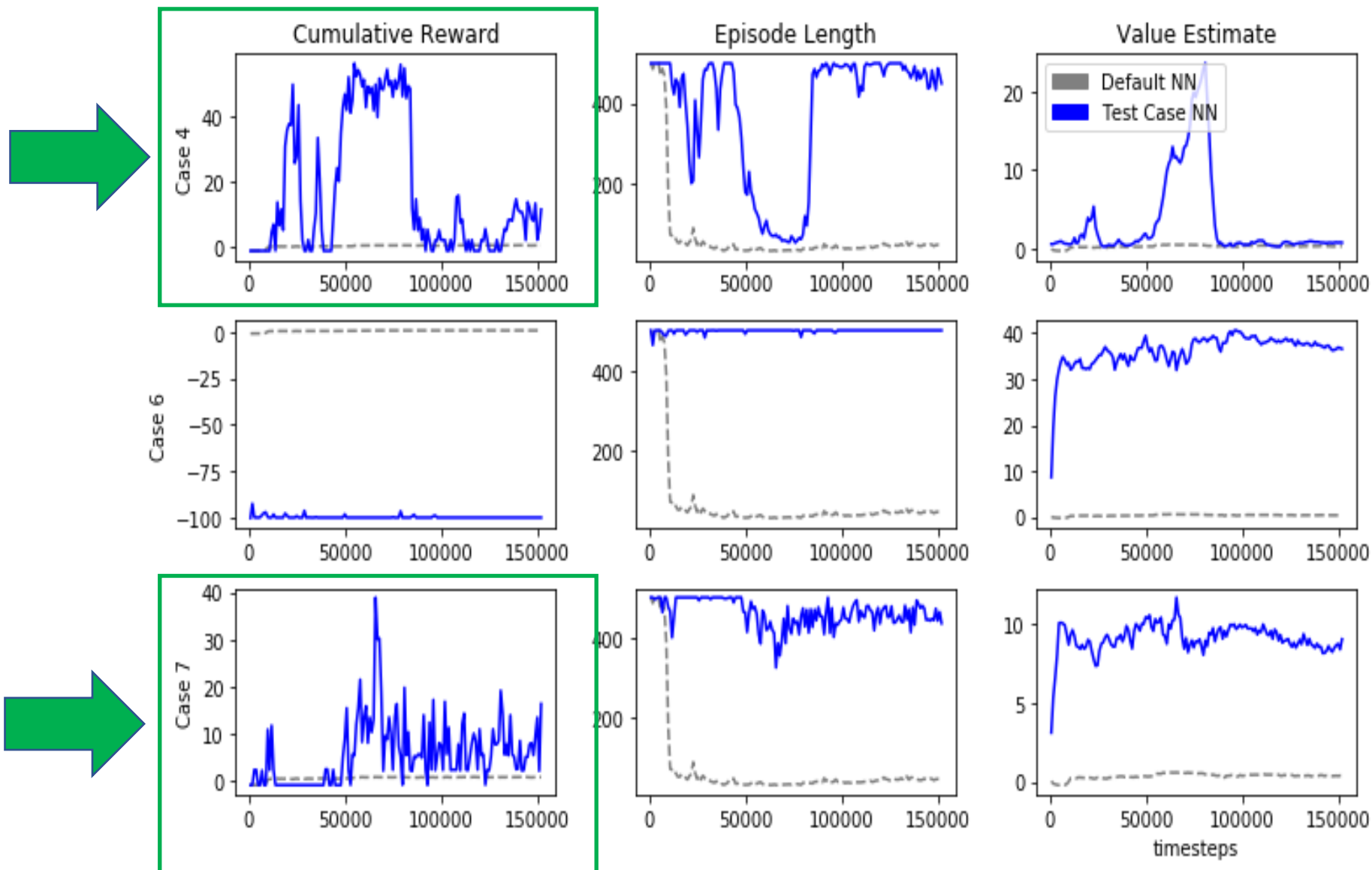




# Conclusion and Future Work

# Conclusion

High-frequency oscillation of cumulative rewards indicate agent misbehavior



# Future Work

- Expand and conduct experiments in additional learning environments
- Investigate Research Questions (RQs) and Research Objectives (ROs) 1 & 2
- RQ1: How can we modify RL algorithms to detect anomalies in training statistics?
  - RO1: Incorporate findings from experiments discussed in this research into a modified PPO + ICM algorithm.
- RQ2: How can we accommodate a distracted human's efforts to intervene during the agent's learning process in human-in-the-loop RL?
  - RO2: Design and implement a scheme for alerting and receiving input from the human during RL through an external smart device.

# The End

---

Thank you!

# References

Images used in this presentation were obtained from the following websites:

<https://www.militaryaerospace.com/unmanned/article/14069139/swarming-uavs-counter>

<https://www.wired.com/story/waymo-google-arizona-phoenix-driverless-self-driving-cars/>

<https://www.xfinity.com/tips/xfinity-x1-remote-tips-and-tricks>

<https://www.republicworld.com/technology-news/mobile/how-to-make-siri-say-things-in-ios-14-learn-in-simple-steps-here.html>

<https://en.wikipedia.org/wiki/Pac-Man>

<https://gym.openai.com/>

<http://www.mujooco.org/>

<https://github.com/Unity-Technologies/ml-agents>

[https://en.wikipedia.org/wiki/Pandas\\_\(software\)](https://en.wikipedia.org/wiki/Pandas_(software))

<https://www.libraries.rutgers.edu/node?page=23>

<https://www.slidescarnival.com/>

<https://www.fotosearch.com/CSP992/k13227884/>