



# DRLearner

Open Source Deep Reinforcement Learning  
AGI-22  
Chris Poulin

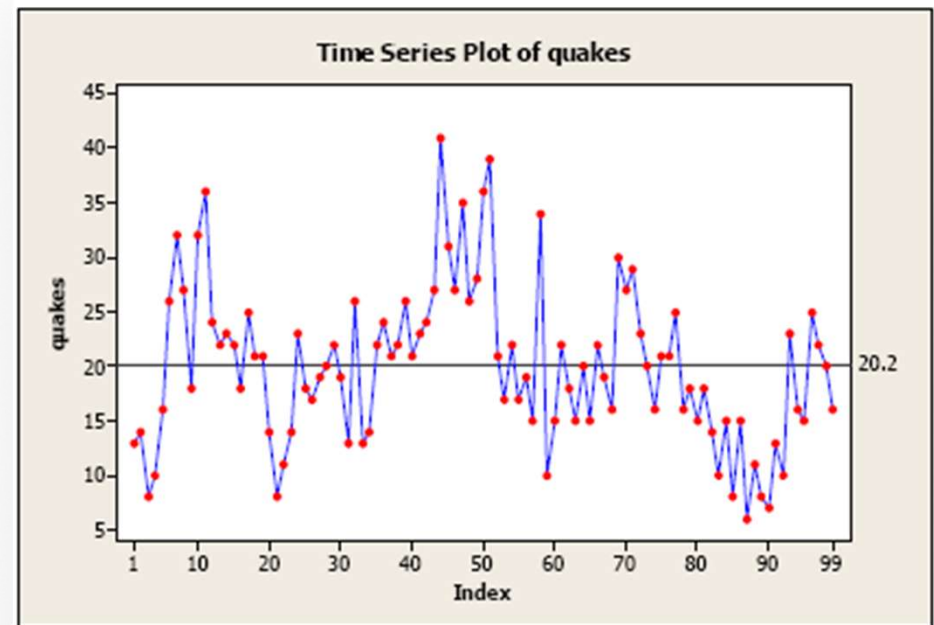
[DRLearner.org](https://DRLearner.org)

# Time Series, RL, and Temporal Difference Learning

## Time Series, and advanced analysis of time-based decisions

- A time series is a series of data points in time order. Most commonly [1].
- A time series is usually a sequence (equally spaced points in time) [2]
- Reinforcement learning (RL) is an area of machine learning concerned with how agents take actions in an environment to maximize reward [3][4]
- Temporal difference (TD) learning refers to model-free reinforcement learning methods which learn by bootstrapping from the current value function.[5] Was discovered that the firing rate of dopamine neurons appear to mimic the error function in the algorithm. [4]
- And many people were working intensely on the math around these problems (e.g. Schmidhuber et al) [6]

- [1] [https://en.wikipedia.org/wiki/Time\\_series](https://en.wikipedia.org/wiki/Time_series)
- [2] <https://online.stat.psu.edu/stat510/book/export/html/661>
- [3] [https://en.wikipedia.org/wiki/Reinforcement\\_learning](https://en.wikipedia.org/wiki/Reinforcement_learning)
- [4] Richard Sutton & Andrew Barto (1998). Reinforcement Learning. MIT Press. ISBN 978-0-585-02445-5
- [5] [https://en.wikipedia.org/wiki/Temporal\\_difference\\_learning](https://en.wikipedia.org/wiki/Temporal_difference_learning)
- [6] <https://people.idsia.ch/~juergen/naturedeepmind.html>



Source: "Time Series Basics", Penn State [2]

# DeepMind's Breakthrough

However, there were a high-profile series of papers that...

- Demonstrated human-level game play on Atari games

---

- *"Playing Atari with Deep Reinforcement Learning"* (arXiv, 2013) [7]

---

- *"Human-level control through deep reinforcement learning"* (Nature, 2015) [8]

---

- This "deep Q network" started a frenzy in machine Learning around the area of "Deep Reinforcement Learning" or "DRL"

[7] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, M. Riedmiller. Playing Atari with Deep Reinforcement Learning. Tech Report, 19 Dec. 2013, <http://arxiv.org/abs/1312.5602>

[8] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, D. Hassabis. Human-level control through deep reinforcement learning. Nature, vol. 518, p 1529, 26 Feb. 2015. <http://www.nature.com/nature/journal/v518/n7540/full/nature14236.html>

©2022 Patterns and Predictions, Phil Tabor, and SoftServe, Inc



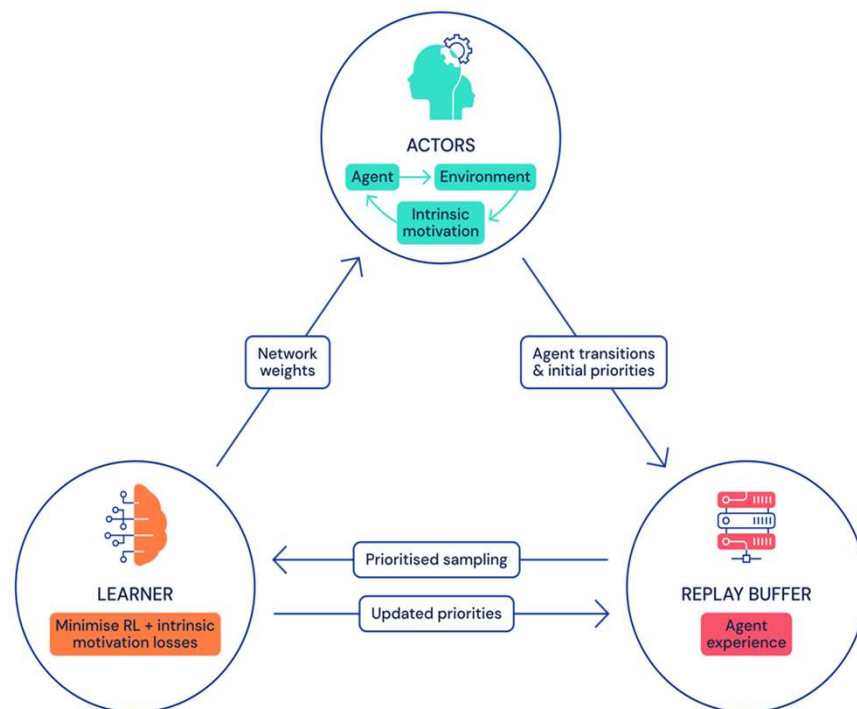
Source: DeepMind

# Agent 57 (from DeepMind)

[Agent57: Outperforming the Atari human benchmark](#)

[By Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, and Charles Blundell - 2020](#)

- **Actors**, which include Environments and Agents
- \* The Agents in this case, have a different scoring paradigm (intrinsic motivation) in **CURIOSITY**
- A **Learner**, which is a Reinforcement Learning front end to a (any) Neural Network.
- Using a **Replay Buffer** (R2D2) to give the system episodic memories to pull from.
- **Massively parallel**, requiring thousands of GPU hours



Source: Agent 57, Badia et. al.

# Agent 57: Performance

## Agent 57 shows better than human performance

- **Mastery level game play** (compared to human average)

---

- Compared to state-of-the-art ML (2020), **best performance by some metrics** (capped mean)

---

- Only **second place in overall average for DeepMind's "best"** (2020)

---

- Relatively **simple architecture**, given the high performance.

---

- Implementation is still relatively new (non-public)

Statistics	Agent57	NGU	R2D2 (Retrace)	R2D2	MuZero
Capped mean	<b>100.00</b>	95.07	94.20	94.33	89.92
Number of games > human	<b>57</b>	51	52	52	51
Mean	4766.25	3421.80	3518.36	4622.09	<b>5661.84</b>
Median	1933.49	1359.78	1457.63	1935.86	<b>2381.51</b>
40th Percentile	1091.07	610.44	817.77	1176.05	1172.90
30th Percentile	614.65	267.10	420.67	529.23	503.05
20th Percentile	<b>324.78</b>	226.43	267.25	215.31	171.39
10th Percentile	<b>184.35</b>	107.78	116.03	115.33	75.74
5th Percentile	<b>116.67</b>	64.10	48.32	50.27	0.03

Source: Agent 57, Badia et. al.

# Agent 57: Demo

## [Here Agent 57 is playing Atari Games](#)

- The Arcade Learning Environment (ALE) is a standard tool for time dependent agents to compare performance

---

- **Zaxxon** is a “medium” difficulty game

---

- The classic **PitFall** gets into “hard” category. As it requires a large amount of understanding of the ‘game map’.

---

- **Montezuma’s Revenge** is considered one of the “hardest” Atari games, in terms of game space complexity.

---

- *In all three cases, Agent 57 is better than an average human player.*

### **Zaxxon**

[https://www.youtube.com/watch?v=FhJ2yLzi4Kk&list=PL2D\\_SqpHWZGgQu4glUARUWk3rrwBrBEgg](https://www.youtube.com/watch?v=FhJ2yLzi4Kk&list=PL2D_SqpHWZGgQu4glUARUWk3rrwBrBEgg)

### **Pitfall**

[https://www.youtube.com/watch?v=96UNx8SvU\\_U](https://www.youtube.com/watch?v=96UNx8SvU_U)

### **Montezuma’s Revenge**

<https://www.youtube.com/watch?v=A32KP0DCbaE>

Source: Agent 57, Badia et. al.

# Our journey: LIT Review

1

We wanted to duplicate the capability of DeepMind's Agent 57, announced 2020, here: **Agent 57: Outperforming the human Atari benchmark**

<https://deepmind.com/blog/article/Agent57-Outperforming-the-human-Atari-benchmark>

2

This was apparently an extension of prior work (Badia et al) on a system called NGU, here: **Never Give Up: Learning Directed Exploration Strategies**

<https://arxiv.org/abs/2002.06038>

3

Once we started communication with Adria Badia, he recommended that I start my efforts based on another effort called R2D2, here: **Recurrent Experience Replay in Distributed Reinforcement Learning (R2D2)**

<https://deepmind.com/research/publications/2019/recurrent-experience-replay-distributed-reinforcement-learning>

4

Adria and Chris then discussed various implementations of this available online, and we thought that the best baseline would be the ACME framework from DeepMind, here: **Acme: A Research Framework for Distributed Reinforcement Learning**

<https://arxiv.org/abs/2006.00979>  
<https://github.com/deepmind/acme>

5

Experimented with ACME, and was able to successfully implement/test (on GCP) most of these examples, including R2D2, here:

<https://github.com/deepmind/acme/tree/master/examples>

6

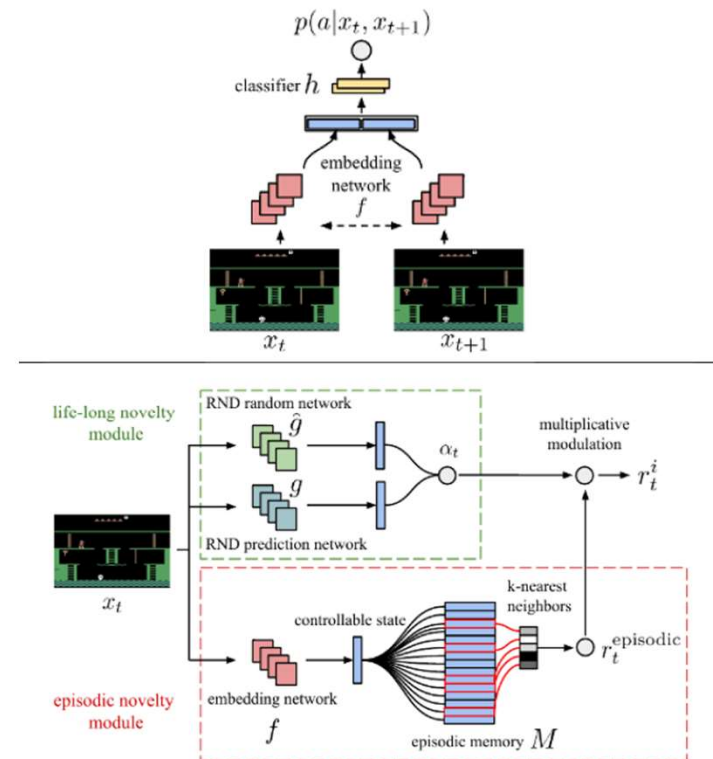
Where we left off was that Adria recommended combining the functionality of NGU with the functionality of R2D2 (on Acme), as he thought this would be a lesser clone of Agent 57 that he could make further suggestions on. He has provided a few code samples, but we (Chris) got stuck on the TF matrix transformations.



# Never Give Up (NGU)

## NGU IS A BOREDOM FREE CURIOSITY-BASED RL METHOD

- Converts screen images to abstract feature representations
- Features used to learn controllable states and to generate intrinsic rewards
- Needed a modification that was applied to open-source ACME R2D2
- DeepMind tools used to design embedding networks, episodic novelty module, and calculate intrinsic rewards
- Single threaded and lacking the life-long novelty module



Credit: Never Give Up: Learning Directed Exploration Strategies, Badia et. al.



# NGU: Implementation

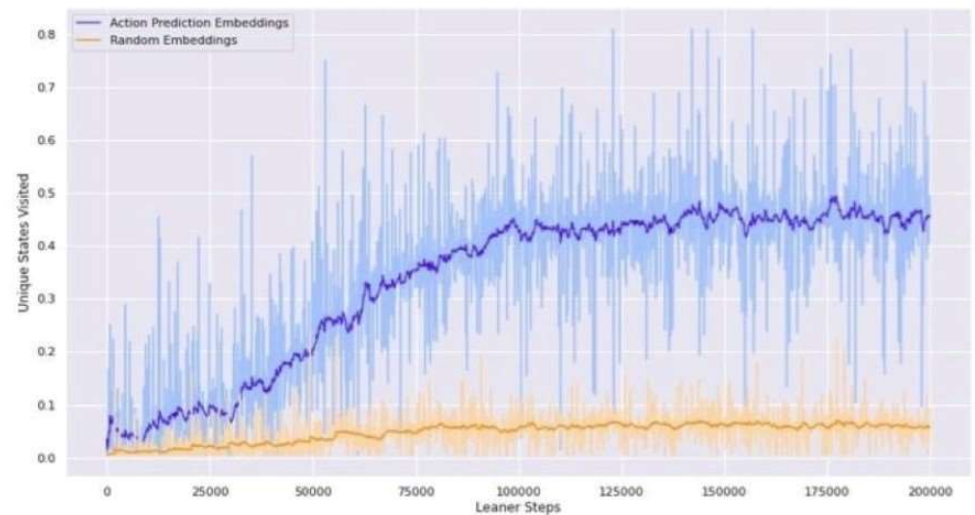
With advice from Badia et al, we first created:

- Implementation of episodic memory using reverb data table extra spec
- Curiosity rewards implemented using k-nearest neighbors in embedding space
- Implement embedding network architecture using sonnet
- Store state embedding representations in reverb data table

```
File Edit View Bookmarks Settings Help
20 class Embedding:
21     def __init__(self, network: network.EmbeddingNetwork,
22                 environment_spec: specs.EnvironmentSpec,
23                 address: Any,
24                 ):
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
File Edit View Bookmarks Settings Help
10 class EmbeddingNetwork(base.Module):
11     def __init__(self, environment_spec: specs.EnvironmentSpec,
12                 n_outputs: int = 18):
13         super(EmbeddingNetwork, self).__init__()
14         self.n_outputs = n_outputs
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
File Edit View Bookmarks Settings Help
15 class EmbeddingLearner:
16     def __init__(self, network: network.EmbeddingNetwork,
17                 environment_spec: specs.EnvironmentSpec,
18                 dataset: tf.data.Dataset,
19                 reward_dataset: tf.data.Dataset,
20                 reverb_client: Optional[reverb.TFClient] = None,
21                 ):
22         self._network = network
23         self._client = reverb_client
24         self._dataset = dataset
25         self._rew_dataset = reward_dataset
26         self._iterator: Iterator[reverb.ReplaySample] = iter(dataset)
27         self._r_iterator: Iterator[reverb.ReplaySample] = iter(reward_dataset)
28
29         self.distance_sum = np.zeros(1)
30         self.distance_counts = np.zeros(1)
31
32         tf2_utils.create_variables(network=self._network,
33                                   input_spec=[
34                                       environment_spec.observations.observation,
35                                       environment_spec.observations.observation]
36                                   )
37
38     def train(self):
39         data = next(self._iterator)
40         data = data.data
41         _, actions, _, _, extra = \
15,1 20%
```

# TensorFlow to JAX

- Switched from TensorFlow to JAX as backend framework to enable *distributed implementation*
- Inverse dynamics model for intrinsic rewards computations – verified by DiscoMaze experiment
- Random Network Distillation for intrinsic reward modulation
- Replace n-step bootstrapping with retrace learning algorithm
- Replace Q-network with Universal Value Function Approximator (UVFA)



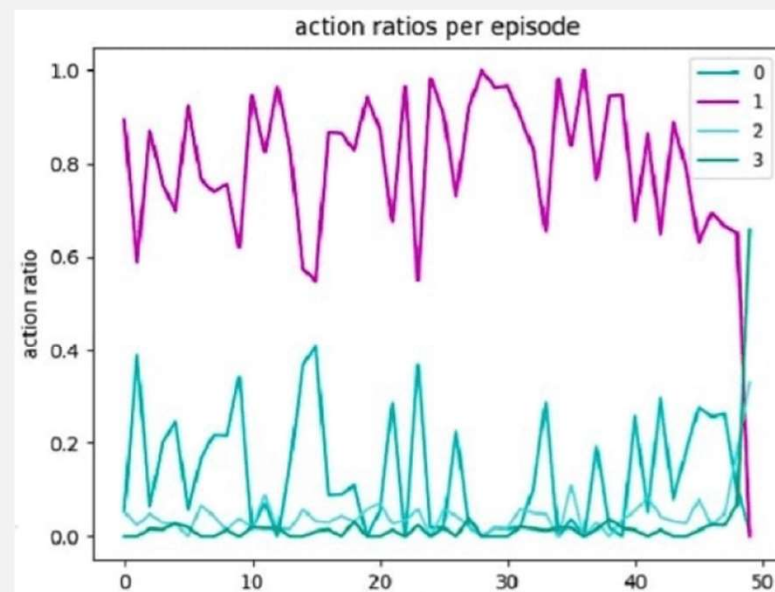
# Initial Testing

## Restoring model from checkpoint

- Implemented script for restoring model from checkpoint, because it wasn't initially supported by ACME
- Note, that optimizer state and training metadata are preserved, so training can be continued at any point

## Metrics customization

- Implemented flexible logging module which can be extended with any custom metric
- Added new metric (actions ratio per episode), which was not previously available in the framework



Actions ratios per episode plot

# Testing and Logging

## LOG EPISODE EXAMPLES

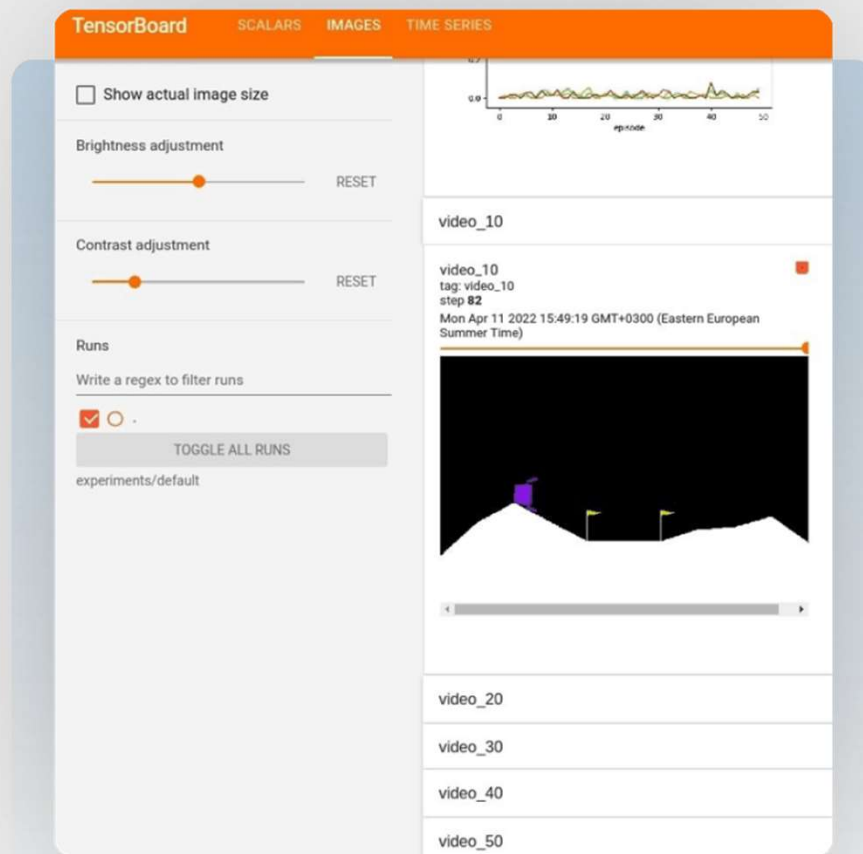
Also, saving episode examples was implemented for monitoring model performance over time.



Episode examples are logged each n-th episode



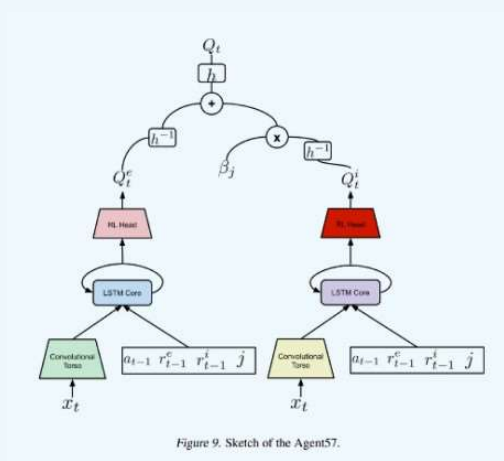
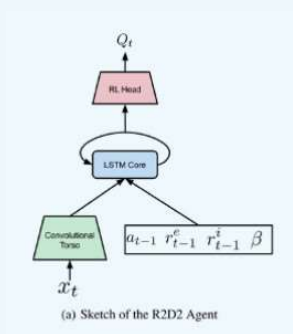
Videos can be logged into TensorBoard or locally



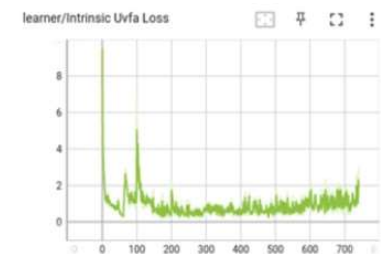
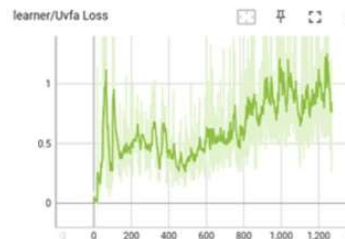
Lunar Lander episode stored in TensorBoard

# Rewards Modeling

Separate parametrization for extrinsic (points) and intrinsic (curiosity) rewards



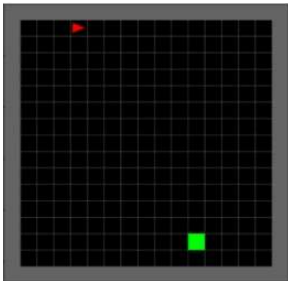
The implementation of the split of UVFA training mostly includes changes in learner, and the inference in the actor.



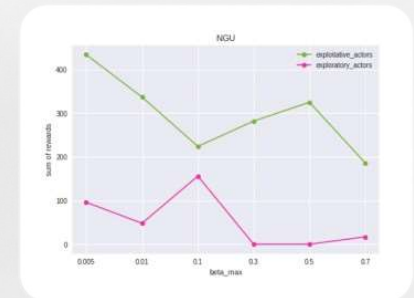
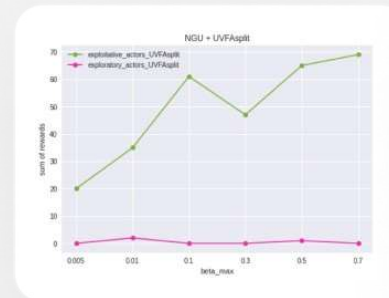
# Experiments testing Rewards

## Random coin environment

- 4 actions: move right, left, up, down. Maximum 200 moves.
- an agent(a red arrow) and a coin(a green square) are randomly placed in 15x15 grid
- when agent steps over a coin a reward of 1 is given and the episode terminates
- Gym-minigrid environment here: [Random Coin](#)



- Further comparison shows that indeed NGU is unstable as the network tries to learn both exploratory and exploitative policies jointly.
- The bigger the beta\_max parameter the more sum of returns for exploitative NGU agents decreases on random coin environment as opposed to the split version



# META Controller/Meta-learning

- **Selects which policy to use at training and evaluation time**

---

- Policies are represented by  $[\beta, \gamma]$  pair (also called mixtures – larger mixture index means more exploratory behavior)

---

- Results in agent learning *when it's better to explore and when to exploit*

---

- Implemented as UCB multi-armed bandit

---

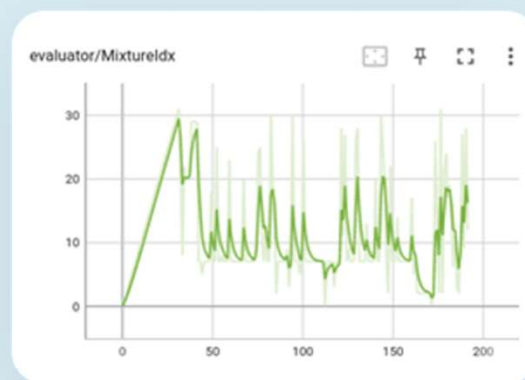
- Extrinsic episode returns are used as rewards for the bandit

---

- Each actor has its own meta controller

---

- Implemented in `actor_core.py`



# Fully Distributed Training

1

Distributed agents can be trained in multiple processes or on multiple machines.

2

Multiple machines training is executed on Vertex AI – a GCP for building and deploying AI models.

3

In a distributed setup, communication between nodes is handled by launchpad package.

## VERTEX AI TRAINING PIPELINE:

- Package the code for every type of launchpad node (Actor, Learner, Replay Buffer) into a Docker container (using launchpad)
- Build the Docker images locally or on Cloud Build
- Specify the hardware requirements for each node type
- Create a custom job on Vertex AI to train the agent
- All training artifacts are saved into a Cloud Storage bucket.

← montezuma\_128\_actors\_agent57\_208g\_mem\_replay\_1652688960365\_1

**Custom job failed with error message: CANCELED**

Status	Stopped
Custom job ID	7765795101644685312
Created	May 16, 2022, 11:38:51 AM
Start time	May 16, 2022, 11:51:06 AM
Elapsed time	7 days 13 sec
Region	us-central1
Encryption type	Google-managed key

---

Machine type (Worker pool 0)	n1-highmem-32
Machine count (Worker pool 0)	1
Container Location (Worker pool 0)	gcr.io/gcp101494-agent57/tmpb2s63d95:20220516-111601-499499
Machine type (Worker pool 1)	e2-highmem-2
Machine count (Worker pool 1)	1
Container Location (Worker pool 1)	gcr.io/gcp101494-agent57/tmpzkb8p4x:20220516-112052-316547
Machine type (Worker pool 2)	n1-highmem-16
Machine count (Worker pool 2)	1
Accelerator (Worker pool 2)	NVIDIA_TESLA_P100
Accelerator count (Worker pool 2)	1
Container Location (Worker pool 2)	gcr.io/gcp101494-agent57/tmp0k6fjk1f:20220516-112556-189971
Machine type (Worker pool 3)	e2-standard-4
Machine count (Worker pool 3)	129
Container Location (Worker pool 3)	gcr.io/gcp101494-agent57/tmp18p9len8:20220516-113212-691428

---

Dataset	No managed dataset
---------	--------------------

---

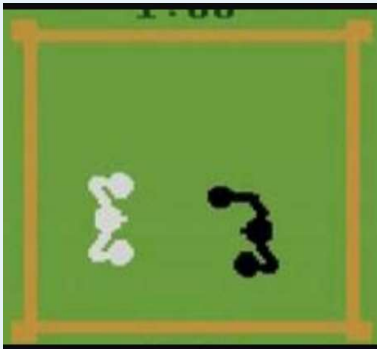
Algorithm	Custom training
Objective	Custom
Container (Training)	Custom
Logs	<a href="#">View logs</a>

Screenshot of Vertex AI training job configuration



# Game Environment Validation

Arcade Learning Environment (ALE)



*Easy: Boxing*



*Medium: Zaxxon*



*Hard: Montezuma Revenge*

Other Modality Testing: Procgen



*Coinrun*

# Code



Source Code:

[https://github.com/PatternsandPredictions/DRLearner\\_beta](https://github.com/PatternsandPredictions/DRLearner_beta)

Dev Mailing List:

<https://groups.google.com/g/drlearner/>



Contact:

**Chris Poulin, Project Lead**

[chris@patternsandpredictions.com](mailto:chris@patternsandpredictions.com)



Sponsored by Patterns and Predictions

[patternsandpredictions.com](https://patternsandpredictions.com)



# Lessons Learned & Next Steps

- **Computational cost (cloud) is non-trivial**

---
- **Noise reduction needed (e.g. contrastive learning?)**

---
- **Refocus on Representations-AGI**

---
- **Continuous environments (e.g. robotics) combinatorially explosive**

---
- **API Exploration (other modalities)**

---
- **Better documenting/on boarding (e.g. instructional videos)**

---
- **Come join us!**

# DRLearner

Thank you

Chris Poulin (Project Lead-US)

Phil Tabor (Co-Lead-US)

Dzvinka Yarish (Ukraine)

Ostap Viniavskyi (Ukraine)

Oleksandr Buiko (Ukraine)

Yuriy Pryyma (Ukraine)

Mariana Temnyk (Ukraine)

Volodymyr Karpiv (Ukraine)

Mykola Maksymenko (Advisor-Ukraine)

Iurii Milovanov (Advisor-Ukraine)

[DRLearner.org](https://DRLearner.org)