



DRLearner

Open Source Reinforcement Learning: Deep Dive
AGI-22
Phil Tabor

DRLearner.org

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

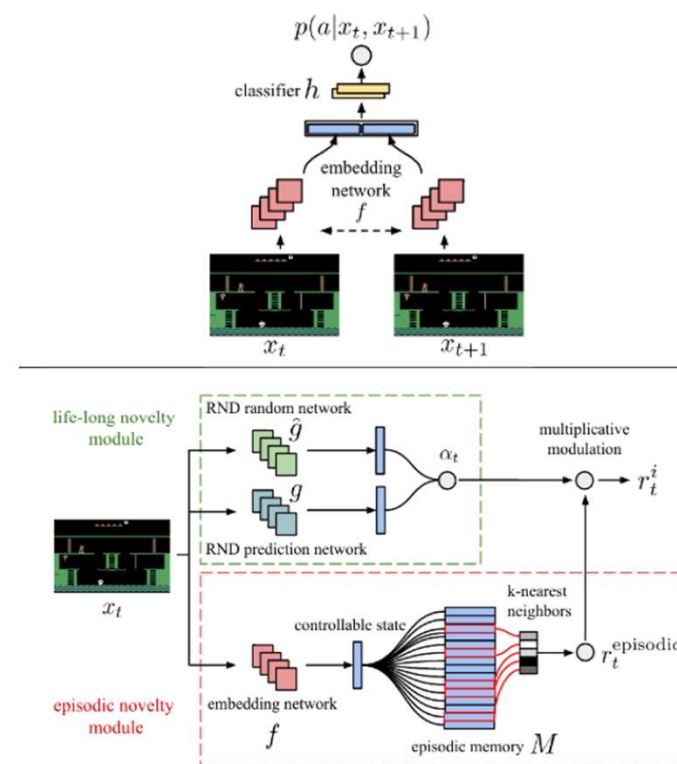
Where we left off was that Adria recommended combining the functionality of NGU with the functionality of R2D2 (on Acme), as he thought this would be a lesser clone of Agent 57 that he could make further suggestions on. He has provided a few code samples, but we (Chris) got stuck on the TF matrix transformations.



Curiosity Learning and NGU

NGU IS A BOREDOM FREE CURIOSITY-BASED RL METHOD

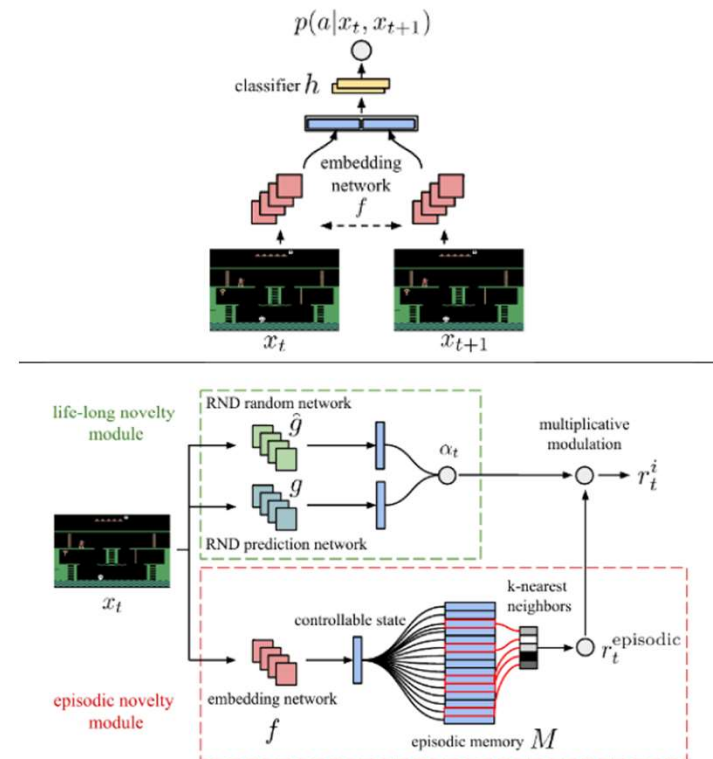
- Curiosity: learning in environments with sparse rewards
- Give intrinsic rewards to the agent based on its inability to predict actions generating successive states
- As the agent visits the same (or similar states) it gets better at predicting transitions, thus intrinsic rewards go to 0 (boredom)
- Another pitfall is rewarding passive observation (agent is rewarded for observing unpredictable noise without acting)
- Never Give Up solves these problems using concepts of lifelong/episodic curiosity and controllable states



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

NGU Architecture

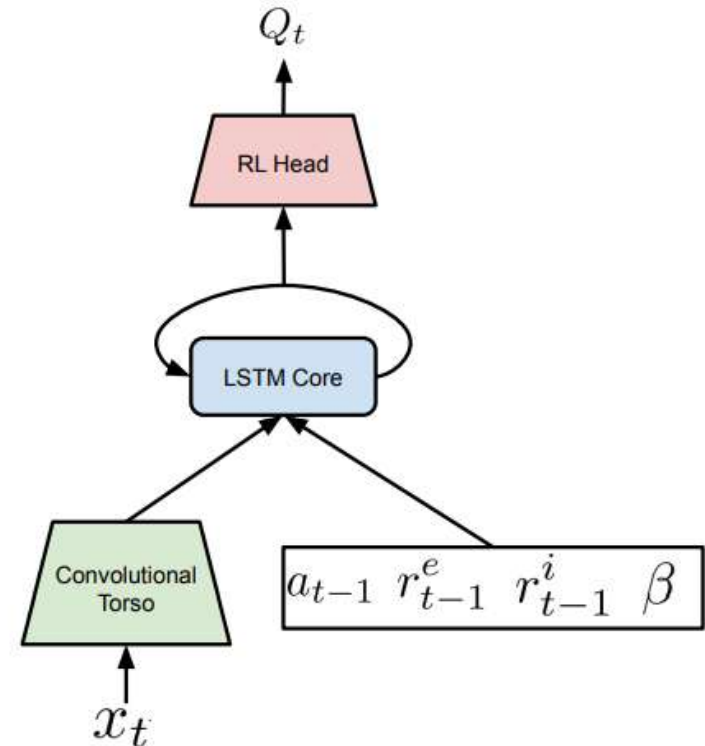
- CNN converts screen images to abstract feature representations
- Embedding features used to learn controllable states and to generate intrinsic rewards
- Distance between inter-episode controllable states is used to generate intrinsic rewards (episodic curiosity)
- RND network generates multiplicative constant for the episodic reward (life-long curiosity)
- What about the explore-exploit dilemma?



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

NGU Architecture

- Use UVFA to approximate $Q(x, a, \theta, \beta_i) \rightarrow r_t = r^e + \beta_i r^i$
- Discrete number of β between β_{\min} and β_{\max} where we include 0 and 1.
- Turn off exploratory policy by acting greedily with respect to $Q(x, a, \theta, 0)$
- Learn to exploit without seeing any extrinsic reward
- Concatenate one hot encoding of beta to action and both rewards and feed into LSTM core for the agent



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

NGU Implementation

V0.1 Implementation

- Start with minimum viable implementation of never give up
- Use DeepMind's tools (ACME) to the greatest extent possible
- Sonnet for the embedding network
- Reverb for data storage and retrieval

```
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
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._reward_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%
```

NGU Implementation

V0.1 Implementation Drawbacks

- Lack of multithreading (TF2)
- No life-long curiosity module
- Used the n-step loss instead of retrace loss
- Didn't use UVFA, so no β

```
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._reward_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%
```

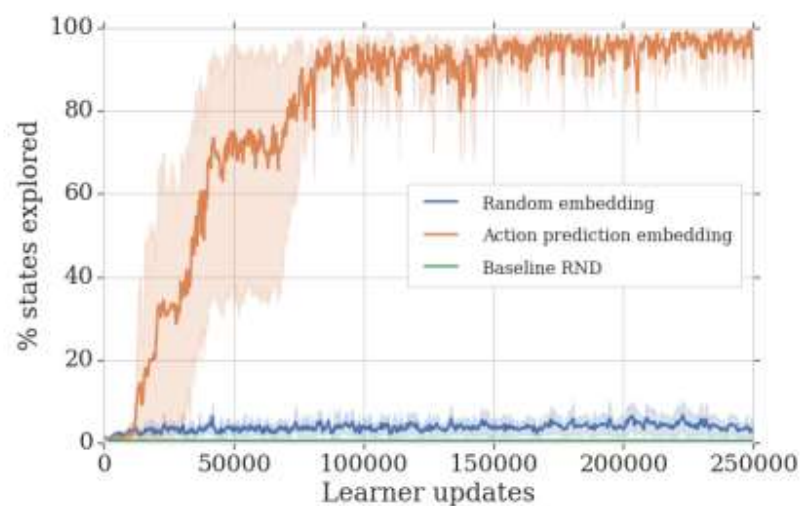
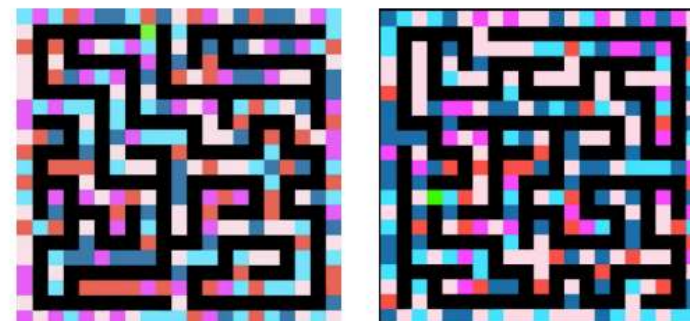
From R2D2 to Full NGU

- 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)

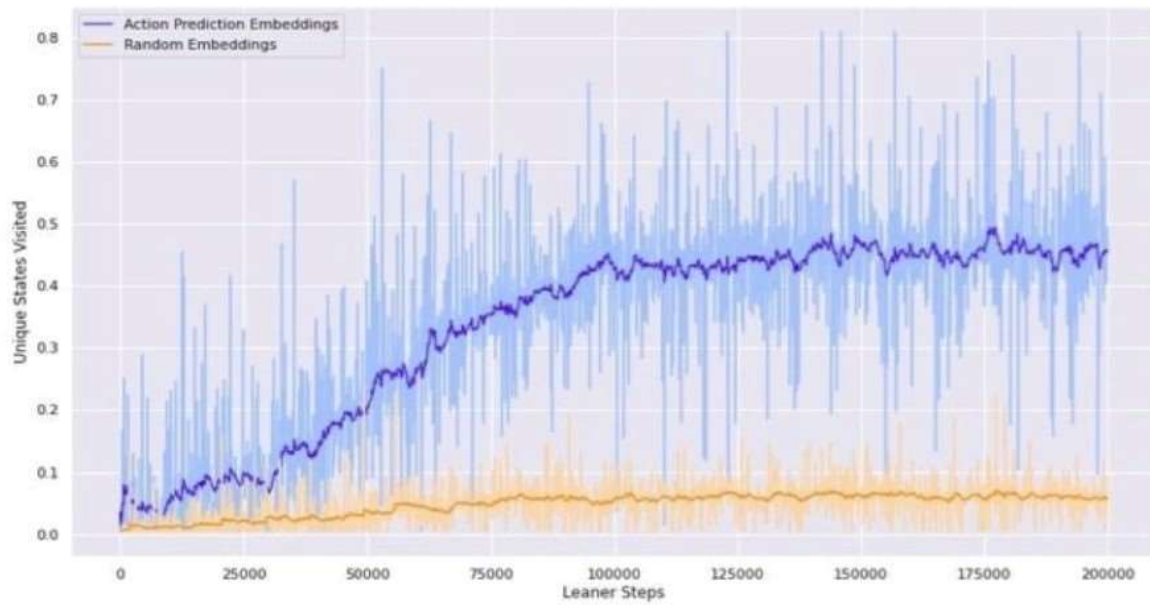
```
1 """Defines local DRLearner agent, using JAX."""
2
3 from typing import Optional
4
5 from acme import specs
6 from acme.utils import counting
7
8 from ..core import local_layout
9 from .builder import DRLearnerBuilder
10 from .config import DRLearnerConfig
11 from .networks import make_policy_networks, DRLearnerNetworks
12
13
14 class DRLearner(local_layout.LocalLayout):
15     """Local agent for DRLearner.
16
17     This implements a single-process DRLearner agent.
18     """
19
20     def __init__(
21         self,
22         spec: specs.EnvironmentSpec,
23         networks: DRLearnerNetworks,
24         config: DRLearnerConfig,
25         seed: int,
26         workdir: Optional[str] = '~/acme',
27         counter: Optional[counting.Counter] = None,
28     ):
29         ngu_builder = DRLearnerBuilder(networks, config, num_actors_per_mixture=1)
30         super().__init__(
31             seed=seed,
32             environment_spec=spec,
33             builder=ngu_builder,
34             networks=networks,
35             policy_network=make_policy_networks(networks, config),
36             workdir=workdir,
37             min_replay_size=config.min_replay_size,
38             samples_per_insert=config.samples_per_insert if config.samples_per_insert \
39                 else 10 / (config.burn_in_length + config.trace_length),
40             batch_size=config.batch_size,
41             num_sgd_steps_per_step=config.num_sgd_steps_per_step,
42             counter=counter,
43         )
44
```


Disco Maze Test

- Test single actor implementation in disco maze environment
- 21x21 randomly generated grid; no extrinsic rewards
- At each time step the blocks change color, which tempts pathological behavior
- NGU agent manages to explore a significant fraction of the available states
- RND and random embeddings flop in this environment



Our Disco Maze Results



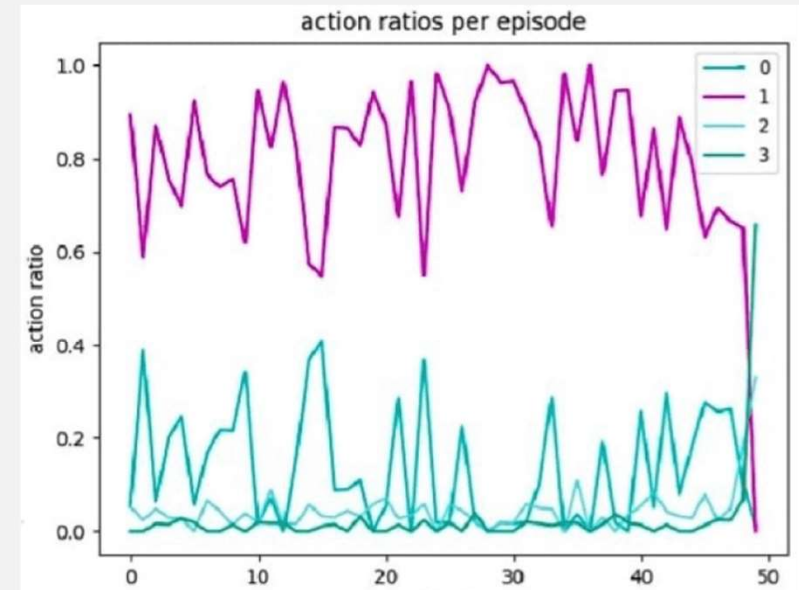
Expanding the ACME Feature Set

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



Action ratios per episode plot

Expanding the ACME Feature Set

LOG EPISODE EXAMPLES

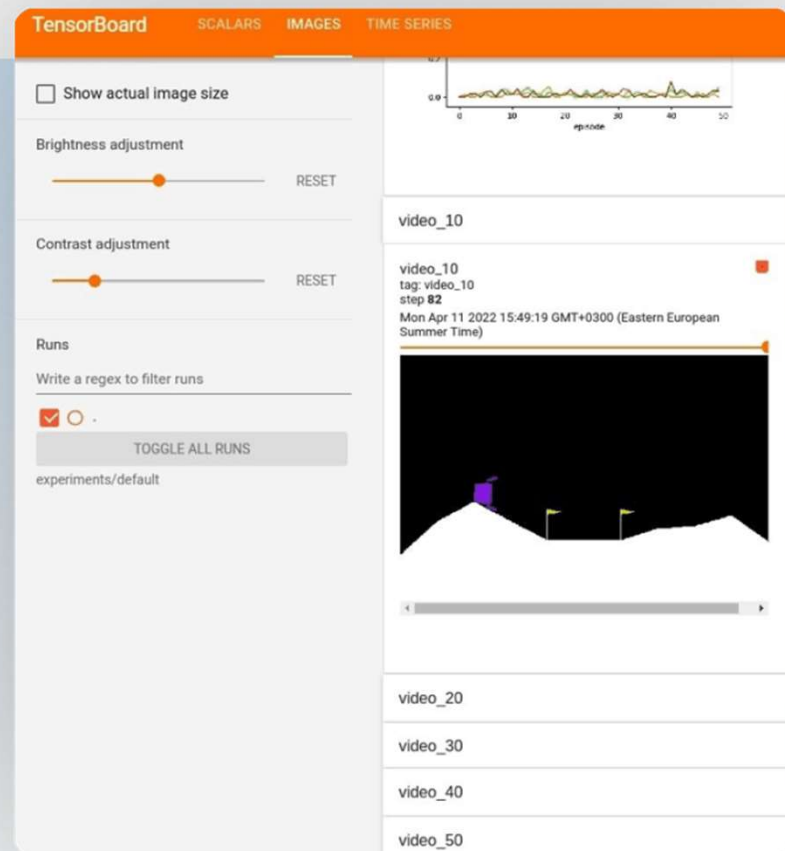
Also, saving episode examples was implemented for monitoring model performance over time.



Episode examples are logged each n-th episode



Videos can be logged into TensorBoard or locally



Lunar Lander episode stored in TensorBoard

Beyond NGU

- DeepMind eventually improved upon NGU with Agent57

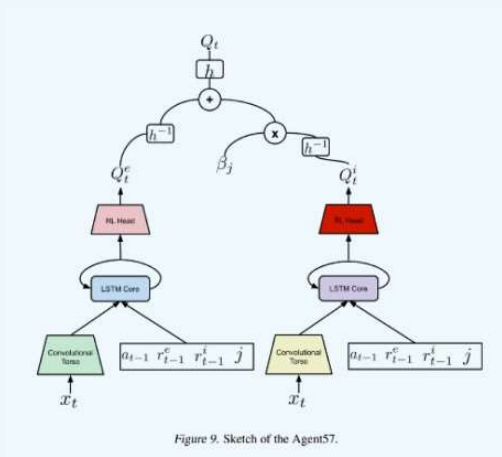
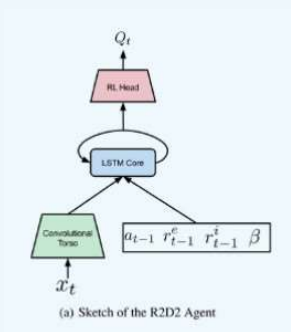
- Split parameterization for UVFA
 $Q(x, a, \theta, j) = Q(x, a, \theta^e, j) + \beta_j Q(x, a, \theta^i, j)$

- Meta-controller for automatically selecting which policy to use

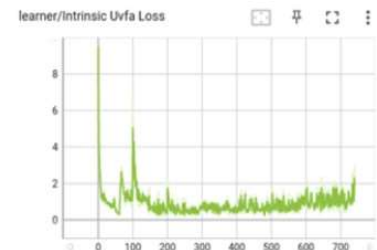
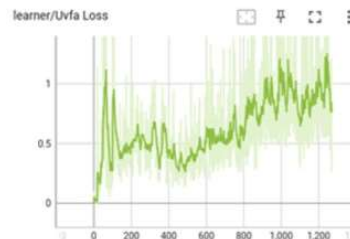
```
28 class DRLearnerLearner(acme.Learner):
29     """DRLearnerLearner."""
30
31     def __init__(self,
32                 uvfa_unroll: networks_lib.FeedForwardNetwork,
33                 uvfa_initial_state: networks_lib.FeedForwardNetwork,
34                 idm_action_pred: networks_lib.FeedForwardNetwork,
35                 distillation_embed: networks_lib.FeedForwardNetwork,
36                 batch_size: int,
37                 beta_min: float,
38                 beta_max: float,
39                 gamma_min: float,
40                 gamma_max: float,
41                 num_mixtures: int,
42                 random_key: networks_lib.PRNGKey,
43                 burn_in_length: int,
44                 target_epsilon: float,
45                 importance_sampling_exponent: float,
46                 max_priority_weight: float,
47                 target_update_period: int,
48                 iterator: Iterator[reverb.ReplaySample],
49                 uvfa_optimizer: optax.GradientTransformation,
50                 idm_optimizer: optax.GradientTransformation,
51                 distillation_optimizer: optax.GradientTransformation,
52                 idm_clip_steps: int,
53                 distillation_clip_steps: int,
54                 retrace_lambda: float,
55                 tx_pair: flax.TxPair = flax.SIGNED_HYPERBOLIC_PAIR,
56                 clip_rewards: bool = False,
57                 max_abs_reward: float = 1.,
58                 use_core_state: bool = True,
59                 prefetch_size: int = 2,
60                 replay_client: Optional[reverb.Client] = None,
61                 counter: Optional[Counting.Counter] = None,
62                 logger: Optional[Loggers.Logger] = None):
63         """Initializes the learner."""
64
65         batched_epsilon_greedy_prob = jax.vmap(
66             jax.vmap(epsilon_greedy_prob, in_axes=(0, None), out_axes=0),
67             in_axes=(0, None),
68             out_axes=0
69         )
70
71         def uvfa_loss(
72             uvfa_params: networks_lib.Params,
73             uvfa_target_params: networks_lib.Params,
74             key_grad: networks_lib.PRNGKey,
75             sample: reverb.ReplaySample,
76             rewards_t: jnp.ndarray,
77             core_state_extraction_name: str = 'extrinsic_core_state'
78         ) -> Tuple[jnp.ndarray, jnp.ndarray]:
79             """
80             Computes mean transformed N-step loss for a batch of sequences.
81             """
```

UVFA Parametrization

Separate parametrization for extrinsic and intrinsic rewards



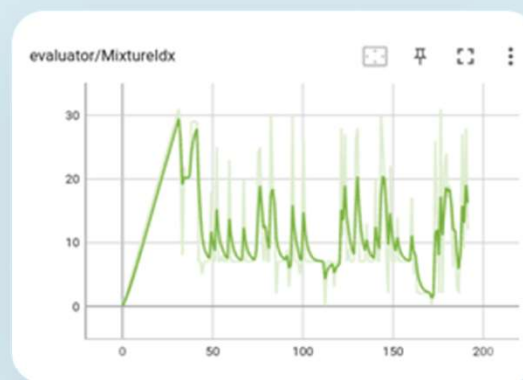
The implementation of the split in UVFA training mostly includes changes in learner, and the inference in the actor.



META Controller

- 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`

©2022 Patterns and Predictions, Phil Tabor, and SoftServe, Inc. All rights reserved.



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

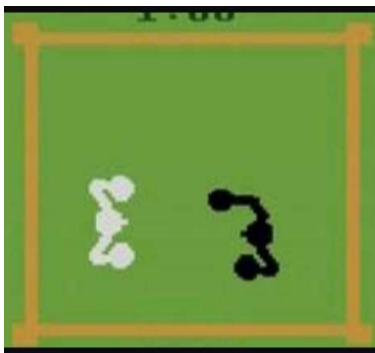
Dataset	No managed dataset
---------	--------------------

Algorithm	Custom training
Objective	Custom
Container (Training)	Custom
Logs	View logs

Screenshot of Vertex AI training job configuration

Training Environments

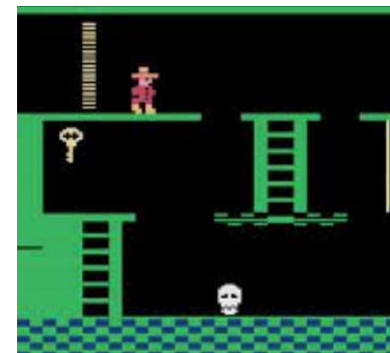
Arcade Learning Environment (ALE)



Easy: Boxing



Medium: Zaxxon



Hard: Montezuma Revenge

Code



Apache License

Code Mailing List:

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



Contact:

Chris Poulin, Project Lead

chris@patternsandpredictions.com



Sponsored by Patterns and Predictions

www.patternsandpredictions.com

