
Replication of Experience Replay for Continual Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we investigate the results reported in the paper “Experience Replay
2 for Continual Learning” [Rolnick et al., 2018] through a replication study of the
3 CLEAR method. More specifically, we review the contributions of this paper
4 and report a detailed specification of our replication process for sequential task
5 learning and implementing the CLEAR method for experience replay. Through
6 our investigation, we found that the authors did not release a codebase of their
7 contributions, and only the baseline, IMPALA [Espenholt et al., 2018], was available.
8 We also found that the CLEAR method required many modifications to the baseline
9 code with vague details not well described in the paper. To this end, we include our
10 solutions to fixing the IMPALA codebase and how we adapted it for the CLEAR
11 method. Lastly, we describe our attempts to implement the CLEAR method with
12 and without behavioral cloning (Figure 2 in the original paper) and report our
13 replicated graph of IMPALA for sequential task learning.

14 1 Introduction

15 When interacting with a complex environment, agents must be able to continually learn and adapt
16 as task specifications change. One desiderata of an intelligent agent is that after it has learned how
17 to perform one task, it should be able to perform more effectively in the future when exposed to
18 similar tasks. However, one major issue in Deep Reinforcement Learning (RL) is the problem of
19 catastrophic forgetting, in which an agent overrides previously learned information after being trained
20 on a new task. Catastrophic forgetting has plagued Deep RL agents from transferring previously
21 learned policies to new scenarios, making continual learning in RL a difficult hurdle to overcome.
22 Simply put, catastrophic forgetting is a symptom of the agent adapting too quickly to new experiences,
23 which destabilizes the learning process.

24 One common method of overcoming catastrophic forgetting is to teach the agent tasks in a simultane-
25 ous manner, which prevents the agent from forgetting critical information for any given task since the
26 agent will be exposed to all the tasks consistently throughout training. However, while this method
27 is feasible when access to simulated training data is abundant, it does not apply well in domains
28 where computational resources for training data is limited and scarce, or when not all tasks are known
29 beforehand, such as robotics. Therefore, the motivation of this paper is to enable Deep RL agents
30 to leverage past experiences (experience replay) in order to overcome catastrophic forgetting when
31 learning sequential tasks.

32 Rolnick et al. [2018] contributes a method for Continual Learning with Experience And Replay,
33 termed the CLEAR method. The CLEAR method enables Deep RL agents to utilize both on-line
34 learning methods to adapt new experiences into the learning process, while also using off-line
35 learning through replay to stabilize learning and prevent catastrophic forgetting in sequential task
36 learning scenarios. The authors also compare their method of implementing CLEAR with and without

37 behavioral cloning, and report the results in Figure 2 of the original paper. The reason why this figure
 38 is of interest is because it demonstrates that even without behavioral cloning, the CLEAR method is
 39 able to do very well in overcoming catastrophic forgetting, making it an attractive method for enabling
 40 continual learning in Deep RL agents. It was this figure that we originally sought out to replicate
 41 and verify that the CLEAR method truly resolves the issue of catastrophic forgetting for continual
 42 learning. However, as noted in the following article, we encountered many issues in both fixing the
 43 original code base and implementing the CLEAR method as described by the authors. Therefore, our
 44 final results include our replication of Figure 1 in the original paper, where the baseline was ran on
 45 several different domains in a sequential manner. By replicating this figure, we demonstrate that we
 46 were able to fix the codebase to get the baseline code working, which is useful and necessary for any
 47 future researchers interested in replicating the CLEAR results.

48 2 The CLEAR Method

49 The CLEAR method [Rolnick et al., 2018] uses actor-critic training on a mixture of new and replayed
 50 experiences. In this actor-critic training method, there are many distributed actors each using
 51 different parameterized policies to choose actions for their given task. Each of these actors feed their
 52 experiences into a global buffer, which a single learner uses to update global parameters that the
 53 actors pull from. Because of this global buffer, the learner uses delayed experiences in the learning
 54 process, requiring an off-policy algorithm for integrating in the experiences from the actors.

55 For this actor-critic method, the value function is updated with an L2 loss, and the policy is regularized
 56 by reducing the entropy loss:

$$L_{value} := (V_{\theta}(h_s) - v_s)^2$$

$$L_{entropy} := \sum_a \pi_{\theta}(a|h_s) \log \pi_{\theta}(a|h_s)$$

57 CLEAR uses the V-Trace off-policy learning algorithm [Espeholt et al., 2018] and adds two loss
 58 terms to induce behavioral cloning between the network and its past self.

$$L_{policy-cloning} := \sum_a \mu(a|h_s) \log \frac{\mu(a|h_s)}{\pi_{\theta}(a|h_s)}$$

$$L_{value-cloning} := \|V_{\theta}(h_s) - V_{replay}(h_s)\|_2^2$$

59 Note that the policy-cloning loss function reflects a minimization of the KL divergence between μ
 60 and π_{θ} , while the value-cloning loss attempts to minimize the L2 norm between the replay value
 61 function and the value function induced by the current policy θ .

62 3 Replication Procedures

63 For the DMLab experiments, Rolnick et al. [2018] modified the network in IMPALA [Espeholt et al.,
 64 2018], the source code for which is available at github.com/deepmind/scalable_agent. We
 65 also started from that codebase and made modifications to it in our attempt to replicate CLEAR.

66 3.1 Setup

67 3.1.1 Using the provided docker file

68 We started with building the `scalable_agent` repository using the docker file provided by the developers
 69 of the repositories. Details of our efforts and the errors faced for the first attempt are as follow:

- 70 • One of the pre-requisites for setting up the `scalable_agent` repository was to set up the
 71 `deepmind-lab` repository (<https://github.com/deepmind/lab>). The first round of er-
 72 rors were related to building up the `deepmind-lab` repo and revolved around **pip issues**, and
 73 **linking .so files using the Bazel compiler**

74
75
76
77

- We used our college's jupyter-repo2docker open source project to attempt to resolve the dependency issues inherent in their docker file as suggested. This route did not resolve all the issues. Even the solutions provided on the closed issues of the github repositories were not useful



Figure 1: Closed issue regarding the dockerfile not working

78
79
80

- We then switched to making our own virtual environment with manual installations of the dependencies to build the repository from scratch
- Figuring this out took us around a day

81 3.1.2 Manual Installation

82 Details of our efforts and the errors faced while manually installing the dependencies and building
83 the scalable_agent repository are as follow:

84
85
86
87

- The deepmind-lab repository with the tasks/environments to run by the agent is compiled with a Bazel compiler. The instructions to install Bazel were wrongly provided in the readme of repository. Closed issues of the github repo were able to better guide us in making the Bazel compiler work

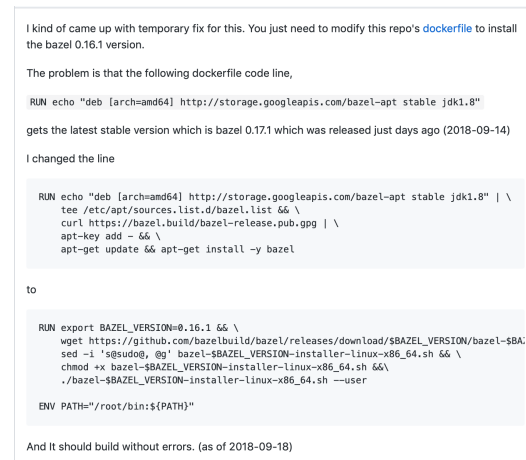


Figure 2: Closed issue regarding solution to making the Bazel compiler work

88
89
90
91
92
93
94

- Instead of doing deb and apt-get we had to wget a stable script release and update the environment path. This solved the Bazel error we were getting
- For the pip issues in the docker file, we had to search and look at the Python.BUILD script. Although the script suggested that the Scalable_agent repository supported both python2 and python3, the python2 functions had function incompatibility issues, leading us to reinstall all the previous dependencies for python3. We had to install a specific version of Tensorflow==1.15.0 and dm-sonnet==1.23 that made the repository work

- 95 • For resolving the linking issues; the dynamic batching module had to be re-run with the
96 right version of the GCC library (gcc=4.8)

97 After spending around a week purely on installation, we were able to get the `scalable_agent` repository
98 up and running, and have a `requirements.txt` file detailing the correct versions of all the necessary
99 dependencies for the `scalable_agent` codebase. We are happy to release it for people who want to
100 work on this project in the future.

101 3.2 Implementation

102 3.2.1 Replay Buffer

103 Rolnick et al. [2018] described the implementation of a replay buffer to store necessary information
104 for the off-policy V-trace algorithm. Each actor needed a separate replay buffer of $capacity =$
105 $number_of_frames / (2 * number_of_actors)$. Every actor forwarded a tuple to the learner,
106 containing the current unroll (trajectory) and the replay history. The learner selected samples from
107 the replay buffer and the current unroll using (50%, 50%) probability. Once the buffers were filled
108 up to the capacity, the new unrolls were added using reservoir sampling, as mentioned by Isele and
109 Cosgun [2018], so that the buffer contains uniform random samples up to the present point and not
110 just recent unrolls. In reservoir sampling, each experience is assigned a random value, which serves
111 as the key in priority queue where the experiences are preserved with the highest key values.

112 For the implementation of the replay buffer, we had to figure out where in the code should the
113 buffer be populated (whether it was inside the `build_actor` function or the `train` function). The
114 `build_actor` function generates the output containing the agent state, agents output and environments
115 output, while in the `train` function the `build_actor`, `build_learner` and `run.session` loop is called.

116 We made multiples attempts at designing the data structure for the replay buffer since the authors
117 provided no information on this. When asked about the implementation of the replay buffer as a
118 GitHub issue, the authors provided a vague answer that they used a custom C++ implementation with
119 a promise of release around a year back. The only information provided was to use `tf_py` function to
120 make wrappers for the replay buffer to increase the speed. A screenshot of the issue is provided in
121 Figure 3

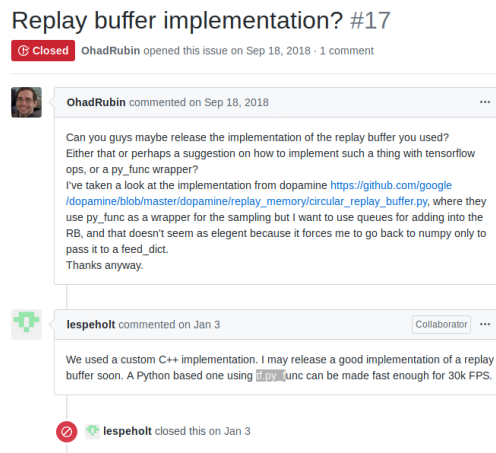


Figure 3: Github issue - Replay buffer implementation

122 The order in which the data structures were implemented for the replay buffer and the errors we ran
123 into implementing each of them are as follows:

124 3.2.1.1 Replay buffer as list per actor

125 The replay buffers were implemented as a list of lists data structure where each replay buffer was
126 implemented as a list per actor, and all the replay buffers per actor were put into one list. This did not

127 work since the samples of unrolls had to be reorganized into a specific internal structure initialized by
 128 the developers of the authors, and the samples from the replay buffers could not be packed into this
 129 specific internal structure to be fed to the learner.
 130 The exact error we faced is shown in Figure 4.

```

Traceback (most recent call last):
  File "experiment_queues.py", line 791, in <module>
    if __app__run()
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/platform/app.py", line 125, in run
    _sys_exit(main(argv))
  File "experiment_queues.py", line 785, in main
    train(action_set, level_names)
  File "experiment_queues.py", line 626, in train
    sample_replay_buffer = nest.pack_sequence_as(structure, sample_replay_buffer)
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/util/nest.py", line 316, in pack_sequence_as
    % (len(flat_structure), len(flat_sequence), structure, flat_sequence))
ValueError: Could not pack sequence. Structure had 12 elements, but flat_sequence had 24 elements. Structure: ActorOutput(level_name=cf.Tensor 'StopGradient:0' shape=() dtype=string, agent_state=SIMSta
dtype=float32), info=StepOutputInfo(episode_return=cf.Tensor 'StopGradient:4:0' shape=(101,) dtype=float32), episode_step=cf.Tensor 'StopGradient:3:0' shape=(101,) dtype=uint32), done=cf.Tensor 'StopGr
adient:6:0' shape=(101,) dtype=bool), observations=cf.Tensor 'StopGradient:7:0' shape=(101, 72, 89, 3) dtype=uint8), <cf.Tensor 'StopGradient:8:0' shape=(101,) dtype=string>), agent_output=AgentOutput(

```

Figure 4: Error - Replay Buffer as a list data structure per actor

131 **3.2.1.2 Replay buffer as queue per actor**

132 Implementing each replay buffer as a queue per actor led to issues with TensorFlow 's serializability.
 133 Initializing multiple queue runners to enqueue the unrolls in the replay buffers caused errors, as
 134 shown in Figure 5.

```

INFO:tensorflow:Error reported to Coordinator: <class 'tensorflow.python.framework.errors_impl.UnknownError'>, EOFError: Ran out of input
Traceback (most recent call last):
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/ops/script_ops.py", line 158, in __call__
    ret = func(*args)
  File "/home/lfrah/lab/scalable_agent/py_process.py", line 86, in py_call
    result = self.out.recv()
  File "/usr/lib/python3.6/multiprocessing/connection.py", line 251, in recv
    return _ForkingPickleler.loads(buf.getbuffer())
EOFError: Ran out of input

[[Node: scan/while/flow_environment_step/step = PyFunc[In=[DT_STRING, DT_INT32], Tout=[DT_FLOAT, DT_BOOL, DT_UINT8, DT_STRING], class=[*loc@scan/while/TensorArrayWrite_8/TensorArrayWriteV3],
tokens=pyfunc_3*, _device=/job:local/replica:0/task:0/device:CPU:0][scan/while/flow_environment_step/step/input_0, scan/while/GatherV2]]2019-11-30 18:57:16.577932: W tensorflow/core/framework/op_kernel
.cc:138] Unknown: UnpicklingError: Invalid load key, '\x00'.
Traceback (most recent call last):
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/ops/script_ops.py", line 158, in __call__
    ret = func(*args)
  File "/home/lfrah/lab/scalable_agent/py_process.py", line 86, in py_call
    result = self.out.recv()
  File "/usr/lib/python3.6/multiprocessing/connection.py", line 251, in recv
    return _ForkingPickleler.loads(buf.getbuffer())
pickle.UnpicklingError: Invalid load key, '\x00'.

```

Figure 5: Error - Replay Buffer as a queue data structure per actor

135 Asynchronicity among multiple threads lead to multiple blocking statements, and the above errors of
 136 running out of input to fill the receiving buffer. As stated by the documentation of TensorFlow, using a
 137 coordinator object with the queue runners could potentially resolve the issues, but the scalable_agent
 138 codebase was written to optimize for speed concerns specifically, and these optimizations were so
 139 interwoven in the code that it made it hard for us to incorporate the coordinator object.

140 **3.2.1.3 Replay buffer as list per actor - second attempt**

141 We went back to the implementation of replay buffer as a list, with the goal of overcoming the
 142 structure packing issues described in Section 3.2.1.1. We were able to overcome the structure issue
 143 by separating the replay buffer into separate replay buffers for each actor correctly, but the code
 144 processed this structure before sending it to the learner. This structure processing can be found in
 145 lines 551 – 571 of the original code in the scalable_agent repository, and we started facing a new
 146 error during this processing, as shown in Figure 6.

```

INFO:tensorflow:Error reported to Coordinator: <class 'tensorflow.python.framework.errors_impl.UnknownError'>, EOFError: Ran out of input
Traceback (most recent call last):
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/ops/script_ops.py", line 158, in __call__
    ret = func(*args)
  File "/home/lfrah/lab/scalable_agent/py_process.py", line 86, in py_call
    result = self_out.recv()
  File "/usr/lib/python3.6/multiprocessing/connection.py", line 251, in recv
    return _forkingPickle.loads(buf.getbuffer())
EOFError: Ran out of input

[[None scan/while/flow_environment_step/step = PyFunc[[In=[DT_STRING, DT_INT32], Tout=[DT_FLOAT, DT_BOOL, DT_UINT8, DT_STRING], class=[loc:@scan/while/TensorArrayWrite_8/TensorArrayWriteV3]],
tokens=[pyrun_3_1, _devices/job-local/replicas/task:0/device:CPU:0]]scan/while/flow_environment_step/step/input_0_scan/while/GatherV2]]2019-11-30 18:57:16.577932: W tensorflow/core/framework/op_kernel.cc:1385] unknown: UnpicklingError: Invalid load key, '\x00'.
Traceback (most recent call last):
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/ops/script_ops.py", line 158, in __call__
    ret = func(*args)
  File "/home/lfrah/lab/scalable_agent/py_process.py", line 86, in py_call
    result = self_out.recv()
  File "/usr/lib/python3.6/multiprocessing/connection.py", line 251, in recv
    return _forkingPickle.loads(buf.getbuffer())
pickle.UnpicklingError: Invalid load key, '\x00'.

```

Figure 6: Error - Replay Buffer as a list data structure per actor - transpose error

147 We believe that this error came up because our replay buffer was a list of lists, while the way the
 148 unroll was processed by the developers of the code was a list of a custom tensor object class. We
 149 tried converting our replay buffer list into a tensor object but ran into the error of incompatible types.

150 3.2.1.4 Replay buffer as list per task

151 We went back to the implementation of replay buffer as a list, but this time implemented it as a list
 152 per task to overcome the type error, as mentioned in 3.2.1.3. This implementation led us to the error
 153 of not being able to pack the sequence into the specific internal structure, as shown in Figure 7.

```

Traceback (most recent call last):
  File "experiment_queues.py", line 791, in <module>
    if __app__run()
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/platform/app.py", line 125, in run
    _sys_exit(main(argv))
  File "experiment_queues.py", line 785, in main
    train(action_set, level_names)
  File "experiment_queues.py", line 626, in train
    sample_replay_buffer = nest.pack_sequence_as(structure, sample_replay_buffer)
  File "/home/lfrah/.local/lib/python3.6/site-packages/tensorflow/python/ops/nest.py", line 316, in pack_sequence_as
    N(len(flat_sequence), len(flat_sequence))
ValueError: Could not pack sequence. Structure had 12 elements, but flat_sequence had 4 elements. Structure: ActorOutput(level_name=tf.Tensor 'StopGradient:0' shape=() dtype=string, agent_state=STM5ta
tuple=(<tf.Tensor 'StopGradient_130' shape=(256,) dtype=float32>, hxxtf.Tensor 'StopGradient_210' shape=(256,) dtype=float32>), env_outputs=StepOutput(reward=tf.Tensor 'StopGradient_310' shape=(101,)
dtype=float32>, info=StepOutputInfo(episode_return=tf.Tensor 'StopGradient_410' shape=(101,) dtype=float32>, episode_step=tf.Tensor 'StopGradient_510' shape=(101,) dtype=int32>), done=<tf.Tensor 'StopGra
dient_610' shape=(101,) dtype=bool>, observations=<tf.Tensor 'StopGradient_710' shape=(101, 72, 96, 3) dtype=uint8>, <tf.Tensor 'StopGradient_810' shape=(101, 9) dtype=float32>, agent_outputs=AgentOutput(<tf.Ten
sor 'StopGradient_910' shape=(101,) dtype=int32>, policy_logits=<tf.Tensor 'StopGradient_1010' shape=(101, 9) dtype=float32>, baseline=<tf.Tensor 'StopGradient_1110' shape=(101,) dtype=float
32>)), flat_sequences=[(<tf.Tensor 'StopGradient_1210' shape=() dtype=string>, <tf.Tensor 'StopGradient_1310' shape=(256,) dtype=float32>, <tf.Tensor 'StopGradient_210' shape=(256,) dtype=float32>, <tf.Ten
sor 'StopGradient_310' shape=(101,) dtype=float32>, <tf.Tensor 'StopGradient_410' shape=(101,) dtype=float32>, <tf.Tensor 'StopGradient_510' shape=(101,) dtype=int32>, <tf.Tensor 'StopGradient_610' shape=(10
1,) dtype=bool>, <tf.Tensor 'StopGradient_710' shape=(101, 72, 96, 3) dtype=uint8>, <tf.Tensor 'StopGradient_810' shape=(101,) dtype=string>, <tf.Tensor 'StopGradient_910' shape=(101,) dtype=int32>, <tf.Ten
sor 'StopGradient_1010' shape=(101, 9) dtype=float32>, <tf.Tensor 'StopGradient_1110' shape=(101,) dtype=float32>], (<tf.Tensor 'StopGradient_1210' shape=() dtype=string>, <tf.Tensor 'StopGradient_1310' sh
ape=(256,) dtype=float32>, <tf.Tensor 'StopGradient_1410' shape=(256,) dtype=float32>, <tf.Tensor 'StopGradient_1510' shape=(101,) dtype=string>, <tf.Tensor 'StopGradient_1610' shape=(101,) dtype=floa
t32>, <tf.Tensor 'StopGradient_1710' shape=(101,) dtype=int32>, <tf.Tensor 'StopGradient_1810' shape=(101, 72, 96, 3) dtype=uint8>, <tf.Tensor 'StopGradient_1910' shape=(101, 9) dtype=float32>, <tf.Ten
sor 'StopGradient_2010' shape=(101,) dtype=string>, <tf.Tensor 'StopGradient_2110' shape=(101,) dtype=int32>, <tf.Tensor 'StopGradient_2210' shape=(101, 9) dtype=float32>, <tf.Ten
sor 'StopGradient_2310' shape=(101, 9) dtype=float32>, <tf.Tensor 'StopGradient_2410' shape=(101,) dtype=float32>, <tf.Tensor 'StopGradient_2510' dtype=string>, <tf.Tensor 'StopGra
dient_2610' shape=(256,) dtype=float32>, <tf.Tensor 'StopGradient_2710' shape=(256,) dtype=float32>, <tf.Tensor 'StopGradient_2810' shape=(101,) dtype=string>, <tf.Tensor 'StopGradient_2910' shape=(101,
72, 96, 3) dtype=uint8>, <tf.Tensor 'StopGradient_3010' shape=(101,) dtype=string>, <tf.Tensor 'StopGradient_3110' shape=(101,) dtype=int32>, <tf.Tensor 'StopGradient_3210' shape=(101, 9) dtype=fl
oat32>, <tf.Tensor 'StopGradient_3310' shape=(101, 9) dtype=float32>, <tf.Tensor 'StopGradient_3410' shape=(101,) dtype=string>, <tf.Tensor 'StopGradient_3510' shape=(101, 72, 96, 3) dtype=uint8>
, <tf.Tensor 'StopGradient_3610' shape=(101,) dtype=string>, <tf.Tensor 'StopGradient_3710' shape=(101,) dtype=int32>, <tf.Tensor 'StopGradient_3810' shape=(101, 9) dtype=float32>, <tf.Tensor 'StopGra
dient_3910' shape=(101,) dtype=float32>]].

```

Figure 7: Error - Replay Buffer as a list data structure per task

154 3.2.1.5 Replay buffer as queue per task

155 The methods of the tensor class Queue were able to handle the type errors and sequence packing
 156 errors by themselves. The reason we were more inclined to a list structure than a queue structure as
 157 used by developers of the Scalable agent code was that we didn't want to dequeue from the replay
 158 buffer; instead, we wanted to keep appending to it to maintain the history of trajectories. Nevertheless,
 159 we again attempted to use the queue data structure, but this time per task. This implementation again
 160 did not work and gave us the same error as Section 3.2.1.2, shown in Figure 5.

161 We had to try all these attempts because the authors were not clear about the details for the implemen-
 162 tation of the replay buffer. We imagine that for researchers wanting to use this codebase as a baseline
 163 for their research, spending so much in implementing the original code is not an effort well-spent.
 164 Replication of the replay buffer would have been more feasible if we at least had an idea of what
 165 underlying data structures were used and a better understanding of the optimizations of the code. The
 166 latter could have been achieved by providing better documentation of the scalable_agent repository.

167 **3.2.2 Sequential Training**

168 The IMPALA (scalable_agent) codebase has settings for single-task and multi-task training, cor-
169 responding to the separate and simultaneous training shown in Figure 1 in Rolnick et al. [2018].
170 However, the setting used in Rolnick et al. [2018] is sequential training (the right-most image in
171 Figure 1), to demonstrate the network’s performance for continual learning. Rolnick et al. [2018] did
172 not specify how they modified the scalable_agent codebase to have the network trained on sequences
173 of tasks, therefore, we operated solely on our guesses of how sequential training was implemented.

174 Sequential training is essentially single-task training where the task switches from time to time.
175 Therefore, we initially sought to modify the single-task training setting to switch the training task
176 after a number of frames. In the single-task setting, all actors are initialized with the same task. We
177 modified the code to reinitialize the actors with a new task after a number of frames, and in effect
178 switching the training task. However, trying to reinitialize the actors gave us an error stating that the
179 Tensorflow graph is finalized and cannot be modified.

180 As re-initializing the actors was not feasible, we decided to initialize all actors at the beginning with
181 all the tasks we wanted to train the network on, and have separate queues to hold the output from the
182 actors for a specific task. However, the way the learner was accessing the actors’ output was unclear
183 and not well-documented in the codebase, and we were unsuccessful in modifying the code to switch
184 the source of training data (and thus the training task) for the network.

185 **3.2.3 Additional Loss Functions**

186 In order to implement policy-cloning and value-cloning as detailed in the CLEAR method, we wrote
187 two different python functions for each of the loss functions. The python code for these two loss
188 functions can be found in Figure 8. Although we were unable to apply these loss functions in the
189 training of the network, as the functions rely on the replay buffer being correctly implemented, we
190 were able to write the loss functions based on the details in the paper, and include them here for
191 posterity.

```
#####  
# New Functions:  
# Applied only for replay experiences  
# The weights were added as described in Section A.4.  
# The motivation for behavioral cloning is to prevent network output on  
# replayed tasks from drifting while learning new tasks.  
#####  
  
def compute_policy_cloning_loss(logits, observed_exp_logits):  
    policy = tf.nn.softmax(logits)  
    observed_policy = tf.nn.softmax(observed_exp_logits)  
    log_frac = tf.math.log(tf.math.divide(policy, observed_policy), axis=1)  
    loss = tf.reduce_sum(- observed_policy * log_frac, axis=-1)  
    return - 0.01 * loss  
  
def compute_value_cloning_loss(replay_advantages):  
    return 0.005 * tf.reduce_sum(tf.square(tf.l2_normalization(replay_advantages)))
```

Figure 8: Additional loss functions

192 **4 Results**

193 With our implemented corrections to the IMPALA codebase, we were able to replicate the first figure
194 in the CLEAR paper where actors were trained on separate tasks. Our replicated figure and the figure
195 from the original paper are included in Figure 9. This figure is important because it helps differentiate
196 between the effects of catastrophic forgetting and interference. Destructive interference is when two
197 tasks that are being learned are in conflict with one another in terms of learned behavior. Catastrophic
198 forgetting is when newly learned experiences override previous experiences. While these two are not
199 mutually exclusive, interference can either be destructive (hurtful) or constructive (helpful), while
200 catastrophic forgetting is often present whenever the agent switches between learning tasks. Figure 1

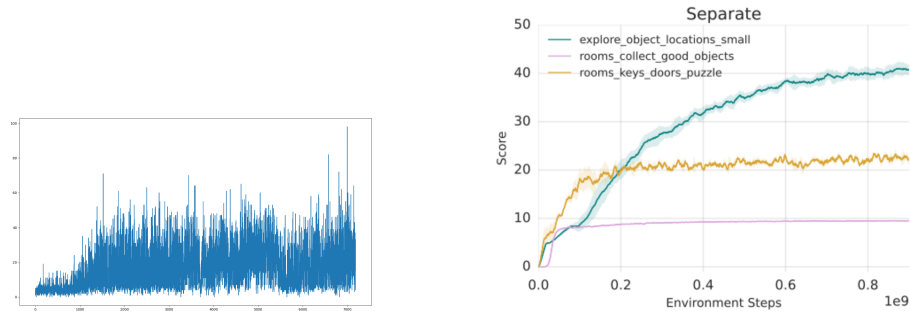


Figure 9: (Left) Our replicated graph. (Right) Figure 1 from the CLEAR paper, separate tasks. Note that we only replicated the `explore_object_location_small` curve

201 attempts to address the differences in these two failure modes by comparing the baseline method,
 202 IMPALA, on separate, simultaneous, and sequential task learning paradigms.

203 We note that the original figure ran for one billion environment steps, but we were unable to run our
 204 model to completion under the same conditions due to our limited compute capability. We also note
 205 that in our figure, we only report the training curve for the `explore_object_location_small` curve.

206 Overall, our graph has more variance than the figure in the original paper, but this is accounted for by
 207 the fact that we only report a single training instance of our model, while the graph from the paper
 208 appears to have averaged results over many runs (However, there is no mention in the original paper
 209 of how many times they ran their models, nor do they specify whether the graph includes standard
 210 deviation or confidence intervals). It is promising to note that our model starts off poorly (as expected
 211 when starting from scratch), and begins to learn such that the score begins to rise over time. It appears
 212 that our replicated graph has an average score that is around 20, which is similar to the actual figure
 213 at around the same amount of environment steps. Therefore, we find our replicated figure to positively
 214 suggest that our implementation results similarly reflects those found in the original paper.

215 5 Conclusion

216 This article describes our process of replicating the paper “Experience Replay for Continual Learning”
 217 and what results we gathered. We found issues with the original codebase and report detailed
 218 descriptions of how we approached these problems. In addition, many important implementation
 219 details for the model and figures were not included in the paper, and we therefore include our methods
 220 of implementing the CLEAR method based on our best guesses. In the end, while we were not able
 221 to successfully implement the replay buffer and sequential task training, which was necessary for
 222 replicating Figure 2 in the original paper, we were able to get the original baseline working on the
 223 separate tasks and write the loss functions for policy cloning and value cloning. Therefore, while we
 224 were not able to confirm that results of this paper, we contributed many corrections that are useful
 225 for any other researchers who wish to implement the CLEAR method.

226 Acknowledgments

227 We would like to give special thanks to Neev Parikh, our assigned Teacher Assistant mentor, for
 228 discussions regarding our replication. We would also like to give thanks to Deniz Bayazit, who was
 229 originally on our team and contributed to our coding replication.

230 References

231 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymir Mnih, Tom Ward, Yotam
 232 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with
 233 importance weighted actor-learner architectures. In *Proceedings of the International Conference
 234 on Machine Learning (ICML)*, 2018.

- 235 David Isele and Akansel Cosgun. Selective experience replay for lifelong learning. *CoRR*,
236 abs/1802.10269, 2018. URL <http://arxiv.org/abs/1802.10269>.
- 237 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P Lillicrap, and Greg Wayne. Experience
238 replay for continual learning. *arXiv preprint arXiv:1811.11682*, 2018.