Scalable Reinforcement Leaning on Cray Systems

Benjamin Robbins, Director - AI & Advanced Productivity Ananda Kommaraju, Kristyn Maschhoff, Mike Ringenburg





Purpose: Demonstrate steps to reproduce a working distributed reinforcement learning configuration on a Cray system. Platform for further exploration.

Relevance: Reinforcement learning approaches have produced many of the recent state of the art results for Machine Learning. Reinforcement learning is resource and communication intensive and is, therefore, an excellent candidate to take advantage of high-performance-computing.

Results: Cray systems' ability to support mixed node types (GPU, CPUs) and resource configuration flexibility make for an ideal platform to explore and push limits of reinforcement learning.



Recent Advances in Al







An animation of the gradient descent method predicting a structure for CASP13 target T1008



Reinforcement Learning



Learning from the rewards of a given action.

PROPRIETARY & CONFIDENTIAL

Reinforcement Learning (RL)

Day 19 I have successfully conditioned him to smile and write in his book every time I drool.

- Pavlov's dog.

Image: geekshumor.com

Reinforcement Learning (RL) Basics





PROPRIETARY & CONFIDENTIAL

The Goal of RL?

Find the optimal Policy.

The policy tells the agent how to act given a particular state.



Many RL Methods



https://www.intel.ai/introducing-reinforcement-learning-coach-0-10-0/#gs.91i90q



Distributed Reinforcement Learning





PROPRIETARY & CONFIDENTIAL

© 2018 Cray Inc.

https://cdn-images-1.medium.com/max/1600/1*tYxWuyksovxA1Thu8PggPQ.jpeg

Motivation for Study of Distributed RL on Cray Systems

- Explore Reinforcement Learning algorithms
- Understand if possible to leverage existing open source frameworks
- Study nested parallelism and complex workflows
- Distributed RL is very active area of research. Participate in research and make available for other Cray end-users
- Scalability is a requirement for many real-world problems. How can distributed RL at the scale of Cray help solve these problems?

Environments for Study

CRAY



PROPRIETARY & CONFIDENTIAL



Setup and Configuration

PROPRIETARY & CONFIDENTIAL



Mapping RL to Cray XC

- Nested Parallelism
- Resource Scaling
- Multiple algorithms can be mapped



What is Ray?



- Ray is a flexible, high-performance distributed execution framework.
- Developed by RISELab at UC Berkeley
- Very active open-source project https://ray.readthedocs.io/en/latest/index.html
- Ray has libraries for tuning and reinforcement learning
- RLlib is Ray's reinforcement learning library
- RLlib supports many RL methods.



Deploy a Ray Cluster on Cray XC



- Allocate nodes through SLURM workload manager
- ccmlogin (Cluster compatible mode)

```
$salloc -N 1 -p ccm_queue -C P100 --gres=gpu --
exclusive
$module load ccm
$ccmlogin -V
```

Start a Ray head node with initial settings

```
#!/bin/bash
source activate ray
IP=$(ip -oneline -family inet addr list ipogif0 \
    | head --lines 1 | grep --perl-regexp \
    --only-matching 'inet \K[\d.]+')
echo $IP:6380 > $HOME/ray_head_node
ray start --head --node-ip-address=$IP --redis-port
    =6380
```

Deploy a Ray Cluster on Cray XC



- Independently start Ray workers with initial settings
- Resource scaling Increase or reduce workers/resources

```
$salloc -N 1 -C P100 --gres=gpu --exclusive
$srun sh worker.sh
```

done

PROPRIETARY & CONFIDENTIAL

Execute Rllib algorithms

```
import ray
from ray.rllib.agents.ppo import PPOAgent
from ray.tune import run_experiments
```

```
def train_fn(config, reporter):
  agent1 = PPOAgent(env="CartPole-v0", config=config)
  for in range(100000):
    result = agent1.train()
    result["phase"] = 1
    reporter(**result)
    phase1 time = result["timesteps total"]
  state = agent1.save()
  agent1.stop()
if name == " main ":
  ray.init(redis address="10.128.0.225:6380")
  run experiments({
     "demo": {
       "run": train fn,
       "local dir": "/lus/scratch/user/ray results/custom/",
       "config": {
         "lr": 0.01,
          "num_workers": 64,
       },
     },
  })
```





RL Methods and Results

PROPRIETARY & CONFIDENTIAL

RL Method – Actor/Critic





Distributed RL Method – A3C



Asynchronous Advantage Actor Critic



https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-8-asynchronous-actor-critic-agents-a3c-c88f72a5e9f2

Distributed RL Method – A3C vs A2C Training in parallel Training in parallel Agent 1 Agent 1 Agent 2 Agent 2 Global Global Network Coordinator Network Agent 3 Agent 3 Parameters **Parameters** Agent n Agent n

A3C (Async)

A2C (Sync)

https://lilianweng.github.io/lil-log/2018/04/08/policy-gradient-algorithms.html#a2c



Qbert_16W_PPO



Breakout_7W_PPO





PROPRIETARY & CONFIDENTIAL

© 2018 Cray Inc.

SpaceInvader_16W_A2C

BeamRider_7W_A2C

Multi-Node Scaling





PROPRIETARY & CONFIDENTIAL

Distributed RL Method – Ape-X





Distributed RL Method - IMPALA



Importance Weighted Actor Learner Architecture



Observations are trajectories of experience (sequence of states, actions, and rewards) not gradients.

https://arxiv.org/pdf/1802.01561.pdf

IMPALA and Ape-X



tag: ray/tune/episode_reward_mean

tune/episode_reward_mean



	Configuration	Reward Mean	Timestep	Walltime
r Ov	IMPALA_Beamrider	2173	10M	28m 18s
	IMPALA_Breakout	374	10M	28m 14s
	IMPALA_Qbert	1077	10M	28m 18s
	IMPALA_SpaceInvader	747	10M	28m 9s





Mapping RL to Cray CS

- 8 GPUs total
- Assign 1 GPU as Head
 Node
- 7 GPUs as Worker Nodes
- 9 CPU cores per Worker



Single GPU Node on CS Storm

XC vs Dense GPU Node on CS



tune/episode_reward_mean tag: ray/tune/episode_reward_mean



Time in hrs

	Configuration	Reward Mean	Timestep	Walltime
agoy/	Qbert_CS_A2C	1.5101e+4	10M	1h 35m 54s
	Qbert_CS_PPO	1.2775e+4	10M	2h 24m 15s
	Qbert_XC_A2C	4193	10M	1h 11m 48s
	Qbert_XC_Ppo	6998	10M	2h 1m 40s



© 2018 Cray Inc.

Key Results

- Explored Reinforcement Learning as a HPC workload
- Deploy a UC Berkeley's RISElab's Ray cluster on XC system
- Trained state-of-art RL agents
- UC Berkeley's RLlib using Ray's distributed execution
- IMPALA, Ape-X, PPO, A2C/A3C on Atari games (from gym library)
- Variety of resource configurations
- Scaled training on multi-node and single-node XC with mixed CPU and GPU





What's Next?



- Explore Optimization of Ray Libraries
- Identify problems sets that map well to Distributed RL
- Comparative studies against other published results
- Explore architecture needs for computational node support and future network requirements

SAFE HARBOR STATEMENT

This presentation may contain forward-looking statements that are based on our current expectations. Forward looking statements may include statements about our financial guidance and expected operating results, our opportunities and future potential, our product development and new product introduction plans, our ability to expand and penetrate our addressable markets and other statements that are not historical facts.

These statements are only predictions and actual results may materially vary from those projected. Please refer to Cray's documents filed with the SEC from time to time concerning factors that could affect the Company and these forward-looking statements.



PROPRIETARY & CONFIDENTIAL

THANK YOU

QUESTIONS?





cray.com

in

@cray_inc

linkedin.com/company/cray-inc-