

Coping with sample inefficiency of deep-reinforcement learning (DRL) for embodied AI

WiML Un-Workshop @ ICML 2020

July 13, 2020

Session Leaders: Vidhi Jain, Simin Liu

Session Facilitators: Ganesh Iyer

Presentation:

Learning has demonstrated great potential for applications like robotics, self-driving cars and IoT [1]. Many of its notable successes have been in virtual gameplay, where thousands of hours of training is feasible. However, collecting real-world data is laborious. To this end, we discuss ways to make DRL methods more practical for embodied AI. Deep reinforcement learning has had great success in simulation, for example, AlphaGo [2] beat human experts, Deepmind's AlphaStar [3] beat top professional players at StarCraft, a challenging real-time strategy game, in 2019.. Similarly, OpenAI Five's DOTA bot [4] won the championship. But it has been much harder applying it to real world platforms (like robots, autonomous vehicles, process control).

Why learning algorithms haven't been deployed on consumer robotics/AV platforms? For example, robot vacuum cleaners use currently rely on a few simple algorithms, such as spiral cleaning (spiraling), room crossing, wall-following and random walk, angle-changing after bumping into an object or wall. There are many reasons why robotic learning is challenging in real world - lack of safety assurances, reward specification, lack of progress on the continual learning front. One of the major focus is mainly on lack of sample efficiency. While we need to have good simulators to train, we also need better ways of acquiring experience in embodied AI and robots for real platforms. We want to be efficient and reduce monotonous burden with AI tools like autonomous vehicle and home assistants.

In our breakout session, we discussed about some ways (either proven or promising) to make DRL feasible for embodied AI. We look into real embodied AI learning in order to fundamentally enhance the notion of intelligence by incorporating multi-modal interaction. We talk about two divergent approaches: algorithmic approaches to improve sample efficiency and alternatively, circumventing the sample efficiency problem by scaling up data collection for current state-of-the-art algorithms.

1 Choice of paradigm

1.1 Model-based (MB) or model-free (MF)?

Model-based methods learn an explicit model of transition function (i.e. learn the prob of moving to a state if you take an action at your current state). Model-free methods don't learn an explicit model of transition function. Instead, you might learn a value function (i.e. $V(s)$, $Q(s, a)$, $A(s, a)$) or directly learn a policy. Here the value function refers to the value of a state is the expected cumulative reward one earns by starting at that state. Model-based is more sample efficient. However, the big caveat is that its asymptotic performance levels out early (unlike model-free).

One reason for this comparative efficiency is that learning local dynamics models (i.e. models which hold in a subset of the input space) is inherently easy. Transition dynamics (especially for robots) are usually locally linear, which makes them easy to learn.

While model-based learning dominates in the low-data regime, its performance plateaus far below model-free methods in the high-data regime. In other words, there's a worrying phenomenon where adding more



Figure 1: Tradeoff between model-free and model-based algorithms¹.

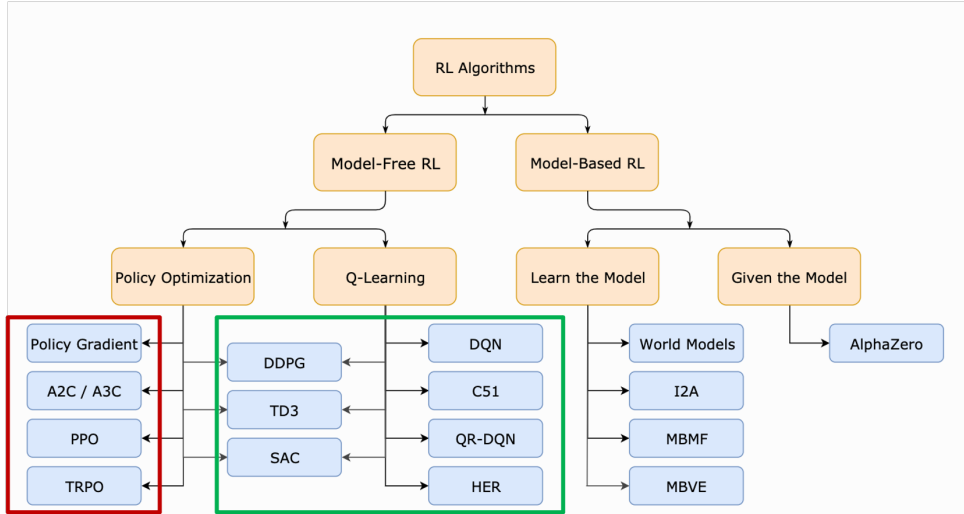


Figure 2: Hierarchy of RL algorithms, where On-policy are shown in **RED** and Off-policy are shown in **GREEN**.

data does not help [5]. This can happen because of (1) Dynamics bottleneck, and (2) Planning horizon dilemma. Dynamics bottleneck refers to the infeasibility of learning a globally accurate dynamics model. MBRL algorithms seem to only be learning locally-accurate dynamics models (i.e. accurate over a small subset of the input space). One piece of evidence for the idea that learning accurate dynamics is the bottleneck is that in some cases, if you use ground-truth dynamics with the MB controller, you can the learned-dynamics counterpart.

Planning horizon dilemma refers to the planning done by MB controllers by applying/chaining the model for multiple timesteps to see what happens when you take a certain action sequence and finding a good action sequence. However, task performance tends to degrade after 20-40 timesteps due to accumulated model error. So we can only plan short-term, which is obviously a draw-back for tasks that require longer foresight.

Some real robots that use MB approach include Millirobot path following, which combines MB with random shooting MPC controller [6] and Wayve’s AV, which uses a dynamics model where the state is the image input compressed to a latent vector by VAE, as “world” models in [7]. It uses backpropagation through time (BPTT) to apply the dynamics model for control.

1.2 Off-policy or on-policy?

Though MB excels in the low-data regime, MF methods are more popular and have been more extensively developed. So, it is also important to make MF work for real platforms. There are two choices among MF algorithms - off-policy and on-policy. Off-policy algorithm can use training samples generated by any policy. This includes Q-learning algorithms. On-policy algorithms can only use training samples generated by the current iteration of their learned policy. This includes policy gradient algorithms. Off-policy is much more data efficient because the data can be reused from different sources for example, data from the previous

iterations of the policy that is currently being updated. Off-policy updates have been shown to learn policies for robotics manipulation [8] in a sample efficient way.

Soft Actor-Critic (SAC) [9] is a max-entropy RL algorithm that learns two Q-functions. This has been successfully applied to real robotic systems, for examples, in dexterous manipulation where the goal is put blue knob on the right, and for a minitaur walking robot.

Asynchronous Q-learning is an algorithm to parallelize data collection. One learner thread is responsible for updating Q function, while multiple worker threads collect data in parallel and place it in the replay buffer. Figure ?? shows the application of this algorithm for robotic manipulation [10].

Both are off-policy MF methods, but note that while they are feasibly efficient, they are still not as efficient as MBRL. In SAC, hours of training data are required; typically, one hour of consecutive data will take many hours to collect due to intermittent hardware failures, etc. In AQL, they scale up data collection by running multiple robots in parallel.

Takeaway: While we have a few ways of making MFRL algorithms work, there is a lot to be done in the space of MBRL to improve sample efficiency, without compromising generalizability of the policy.

2 Using knowledge from other domains

Transfer learning is a kind of “Forward transfer”. Transfer learning refers to when a model trained on one task data can be applied or transferred to a new task. from mid level visual feature representations has shown effective success in embodied AI navigation tasks. Similarly, transfer approaches have been applied in language based navigation and modular approach for vision-language grounded navigation tasks.

Multi-task learning implies that the policy is trained to solve many tasks, and show good generalization to an unseen task. This approach is inspired from how human babies learn not a single but many different tasks, and has success in reinforcement learning, particularly evolutionary policy search. DeepMind’s UNREAL [11] agent learns to navigate by training on auxillary tasks in an unsupervised manner.

Meta learning refers to learning to learn by solving a variety of tasks. The meta policy learns to minimize the distribution of loss functions that allows it to generalize well to few-shot learning settings. Meta learning has been shown to enable adaptability in legged robots in case of leg joint failures.

Combining modular components of learning systems for embodied AI has been successful approach in Active Neural SLAM and object navigation challenges recently where the policy is decomposed into a mapper, global and local policy. These components have been shown to work on real LocoBot using MaskRCNN for object detection and segmentation.

3 Human Demonstrations and feedback

Copying experts reduce the number of samples needed to explore the state-action space for the optimal policy. Vanilla imitation learning involves extraction of state-action pairs from collected expert demos (trajectories), and minimization of error between policies and next actions. [12] Incorporating demonstrations can also be investigated with learning the Q-value and through inverse reinforcement learning.

4 Scaling up data collection

4.1 Sim-2-real

Can we train primarily in simulation and then transfer the policy to the real world? Sim-2-real approach allows for cheap data collection, effortless scaling and safe setting to learn policies. But we need to deal with the ‘sim-2-real gap’. There are often significant visual and physical differences between the simulation and reality.

The simplest transfer case where we can naively go from simulations to real-world, requires that we have a simulator with the perfect model of the world. Domain adaptation and domain randomization have shown empirical success with neural networks. Domain adaptation refers to approach of finding a robust mapping from simulation to reality. Domain randomization refers to perturbing the simulation to different extents in

the hope that generalization to reality becomes an interpolation task (rather than a difficult extrapolation one).

4.2 Parallelized methods

4.2.1 Parallelized asynchronous data collection

Edge workers merely send raw data to the server for training. Current attempts to parallelize data collection (Arm Farm at Google), algorithmic changes need to be made to suit in asynchronous, parallel data collection (i.e. straggler mitigation) For embodied AI navigation, [13] scaled the data collection and execution in decentralized way, while maintaining the centralized synchronized training updates.

4.2.2 Federated learning

Edge workers update personal models and asynchronously send model parameters to update the global features on the server. Federated learning allows the robot to learn local features for adaptation, leverage the personal data in privacy preserving way, and communicate the model parameters efficiently for global feature representation. Federated learning approach has been successful in improving on-device voice assistants for keyword spotting [14], and is an open research area for self-driving cars and IoT [15]. Federated deep reinforcement learning is an active area of research [16].

5 Miscellaneous

Embedding strong inductive priors. For example, reward shaping [17], crafting exploratory behavior [18], and incorporating human feedback or demonstrations [11]. Additional topics include algorithmic changes to make more efficient use of data, like experience replay [19] or curriculum learning.

6 Discussion

- When is DRL useful/necessary for embodied AI applications? (i.e. when do data-driven methods have an advantage over traditional planning control methods?)
- Is sample inefficiency a bottleneck in the progress of DRL for robotics? Pro: Sample inefficiency deters real-world applications, adaptation of the algorithms Cons: Sample inefficiency in current methods can be overcome with improved faster hardware and realistic simulation. Though sample inefficient, the current algorithms are somewhat self-reliant - perhaps we have a tradeoff here.
- Why is the sample inefficiency of current DRL algorithms not so serious? - improved hardware, simulation power, less reliance on known priors and more on rules extracted from patterns in the data i.e. data-driven?
- Does Sim2Real work? Just by publishing papers on the success of sim2real, we can not believe that sim2real works. There are no negative examples that we can discuss! Since sim-2-real somewhat works, can't we just use any of our DRL algorithms, even if data inefficient?
- Should our focus as a community be on circumventing sample efficiency (gathering data at scale, sim-2-real) or address it at algorithmic level? Should we modify our learning algorithm or not? (focus: embodied AI learning models)
- Do you see ways in which these methods can be combined? What's wrong with the way we currently measure/quantify sample efficiency?

7 Open Questions

- Should we expend more effort in improving the realism of our simulations or on algorithms to compensate for their lack of realism? If the first: how do we measure simulation realism? If the second: how do we make algorithms less ad hoc (arbitrary) than they are currently (esp. Domain Randomization - isn't the way you permute is arbitrary)?
- Federated learning has shown early promise in areas like query suggestions on mobile phones, smart speakers, etc. What other applications can you think of?
- What kind of realism is desirable - like visual realism through meshes or physical realism through CAD models and physics based game engines? For example, game engines like Unity (AI2THOR) can provide similar physical realism. Real world 3D mesh renderers like AI Habitat provide with visual realism in simulators.

References

- [1] Lei Lei, Yue Tan, Shiwen Liu, Kan Zheng, and Xuemin Shen. Deep reinforcement learning for autonomous internet of things: Model, applications and challenges. *CoRR*, abs/1907.09059, 2019.
- [2] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharmashan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- [3] Oriol Vinyals, Igor Babuschkin, Junyoung Chung, Michael Mathieu, Max Jaderberg, Wojtek Czarnecki, Andrew Dudzik, Aja Huang, Petko Georgiev, Richard Powell, Timo Ewalds, Dan Horgan, Manuel Kroiss, Ivo Danihelka, John Agapiou, Junhyuk Oh, Valentin Dalibard, David Choi, Laurent Sifre, Yury Sulsky, Sasha Vezhnevets, James Molloy, Trevor Cai, David Budden, Tom Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Toby Pohlen, Dani Yogatama, Julia Cohen, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Chris Apps, Koray Kavukcuoglu, Demis Hassabis, and David Silver. AlphaStar: Mastering the Real-Time Strategy Game StarCraft II. <https://deepmind.com/blog/alphastar-mastering-real-time-strategy-game-starcraft-ii/>, 2019.
- [4] OpenAI, :, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning, 2019.
- [5] Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking model-based reinforcement learning. *CoRR*, abs/1907.02057, 2019.
- [6] Anusha Nagabandi, Guangzhao Yang, Thomas Asmar, Gregory Kahn, Sergey Levine, and Ronald S. Fearing. Neural network dynamics models for control of under-actuated legged millirobots. *CoRR*, abs/1711.05253, 2017.
- [7] David Ha and Jürgen Schmidhuber. World models. *CoRR*, abs/1803.10122, 2018.
- [8] Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. *2017 IEEE International Conference on Robotics and Automation (ICRA)*, May 2017.
- [9] Tuomas Haarnoja, Aurick Zhou, P. Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, 2018.

- [10] S. Gu, E. Holly, T. Lillicrap, and Sergey Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3389–3396, 2017.
- [11] Max Jaderberg, V. Mnih, Wojciech M. Czarnecki, Tom Schaul, Joel Z. Leibo, D. Silver, and K. Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *International Conference on Learning Representations (ICLR)*, abs/1611.05397, 2017.
- [12] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to end learning for self-driving cars, 2016.
- [13] Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames, 2019.
- [14] Andrew Hard, Kurt Partridge, Cameron Nguyen, Niranjana Subrahmanya, Aishanee Shah, Pai Zhu, Ignacio Lopez Moreno, and Rajiv Mathews. Training keyword spotting models on non-iid data with federated learning, 2020.
- [15] Ahmed Imteaj, Urmish Thakker, Shiqiang Wang, Jian Li, and M. Hadi Amini. Federated learning for resource-constrained iot devices: Panoramas and state-of-the-art, 2020.
- [16] Hankz Hankui Zhuo, Wenfeng Feng, Yufeng Lin, Qian Xu, and Qiang Yang. Federated deep reinforcement learning, 2019.
- [17] Glen Berseth, Daniel Geng, Coline Devin, Nicholas Rhinehart, Chelsea Finn, Dinesh Jayaraman, and Sergey Levine. Smirl: Surprise minimizing reinforcement learning in dynamic environments, 2019.
- [18] Deepak Pathak, Pulkit Agrawal, Alyosha Alexei Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 488–489, 2017.
- [19] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5048–5058. Curran Associates, Inc., 2017.