

VIDEO



AN OPEN-SOURCE MODEL ZOO FOR ANALYZING, VISUALIZING, AND COMPARING DEEP REINFORCEMENT LEARNING AGENTS

Joel Lehman,
Senior Research Scientist,
Uber



VIDEO



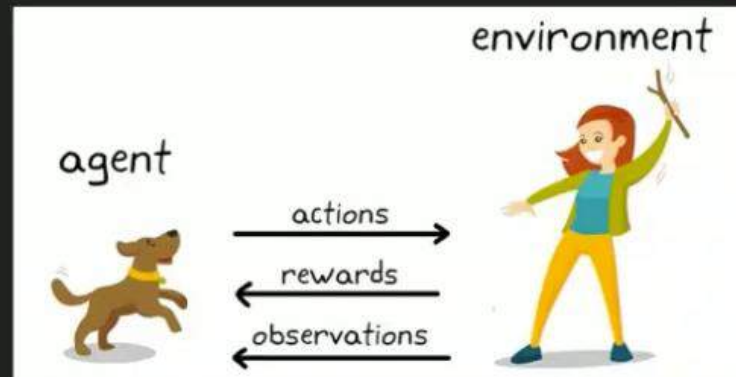
**Make it easier to research the
behavior of reinforcement learning
agents produced by different
algorithms**



VIDEO



An Open-Source Model Zoo for Analyzing, Visualizing, and Comparing **Deep Reinforcement Learning** Agents



VIDEO



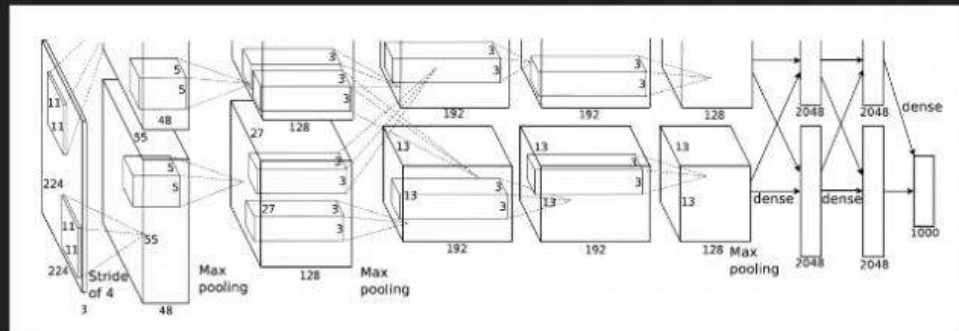
An Open-Source (**Atari**) Model Zoo for Analyzing, Visualizing, and Comparing Deep Reinforcement Learning Agents



VIDEO



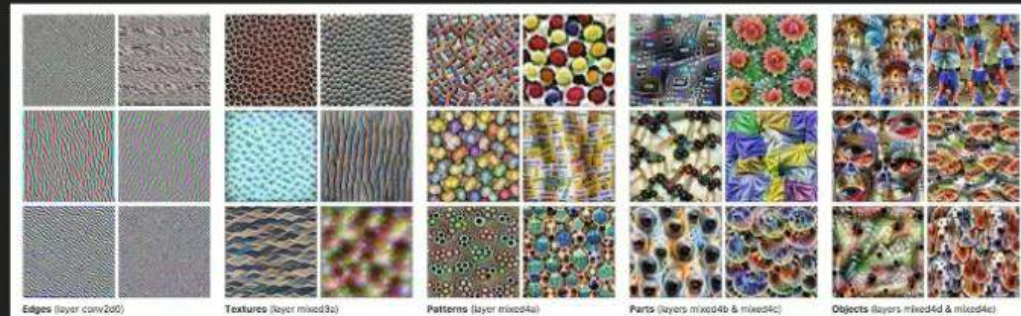
An Open-Source **Model Zoo** for Analyzing, Visualizing, and Comparing Deep Reinforcement Learning Agents



VIDEO



An Open-Source Model Zoo for **Analyzing, Visualizing, and Comparing** Deep Reinforcement Learning Agents



VIDEO



Reinforcement Learning

Learning from Rewards



VIDEO



Reinforcement Learning

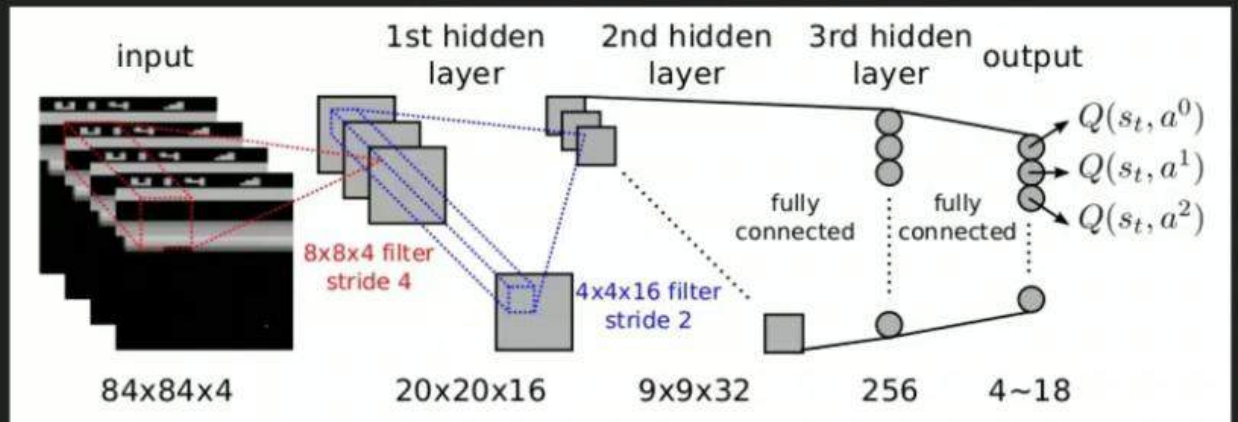
Learning from Rewards

**335**

VIDEO



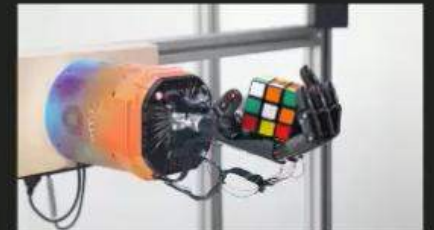
Deep Reinforcement Learning



VIDEO



Deep Reinforcement Learning



VIDEO



Deep Reinforcement Learning

- Lots of potential
- But...
 - Important in some tasks to *understand* what an agent is doing
- Just as *explainable AI* is important for supervised learning, it is also important for reinforcement learning



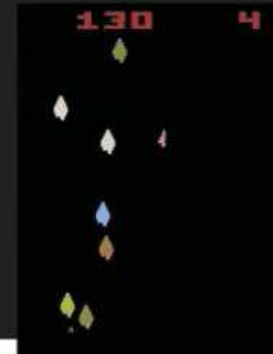
VIDEO



**Make it easier to research the
behavior of RL agents produced by
different algorithms**



VIDEO



Deep Neuroevolution: Genetic Algorithms are a Competitive Alternative for Training Deep Neural Networks for Reinforcement Learning

Felipe Petroski Such Vashisht Madhavan Edoardo Conti Joel Lehman Kenneth O. Stanley Jeff Clune

Uber AI Labs

{felipe.such, jeffclune}@uber.com

2018

Abstract

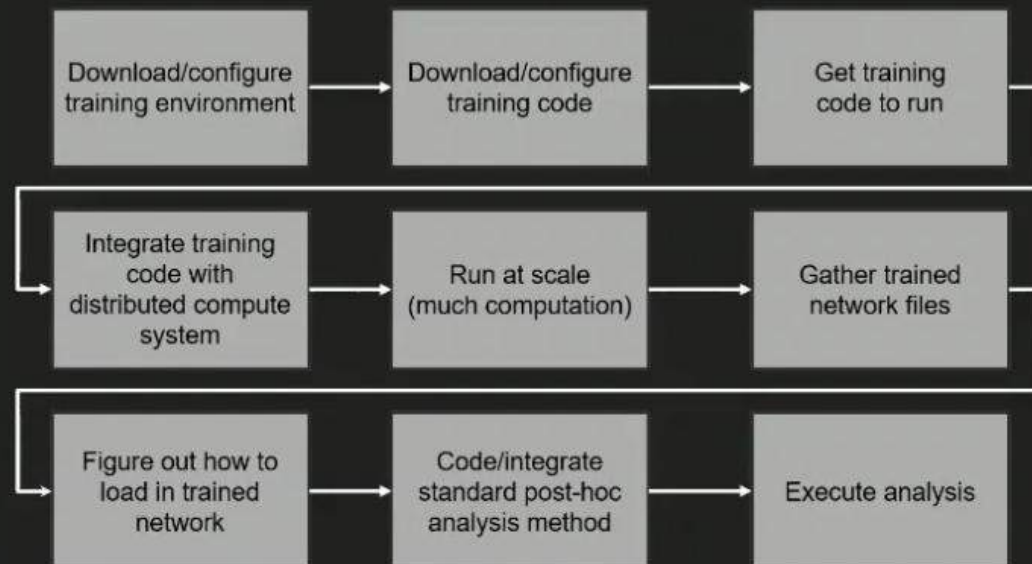
1. Introduction



VIDEO



Workflow to Analyze Agents (without Model Zoo)



VIDEO



With model zoo

```
In [15]: #load model
m = MakeAtariModel("a2c", "SeaquestNoFrameskip-v4", 1, "final")()

#grab frames and high-level neural representations
frames = m.get_frames()
rep = m.get_representation()

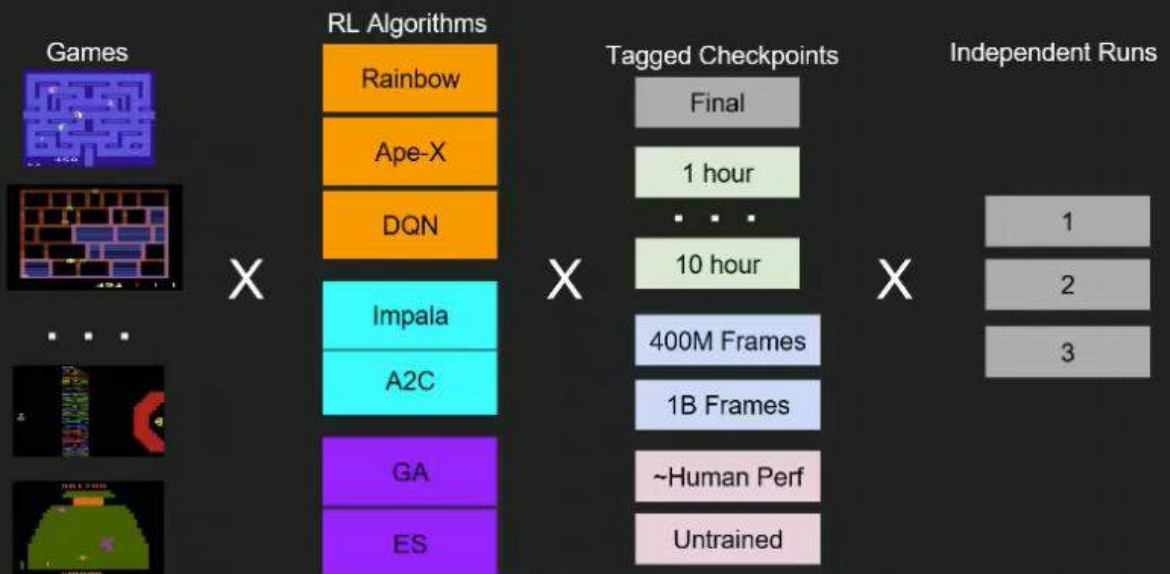
imshow(frames[10])
figure()
plot(rep[10])
```



VIDEO



An Atari Model Zoo: Lots of Trained Models

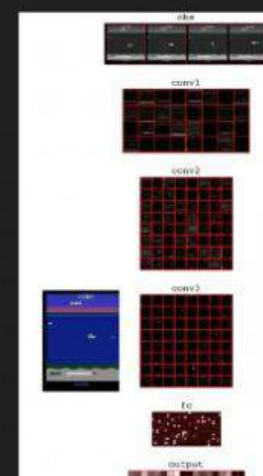
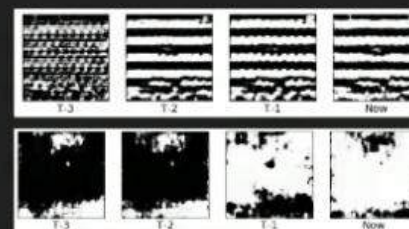
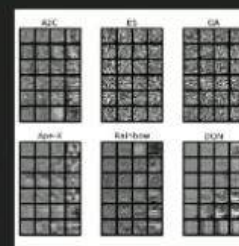
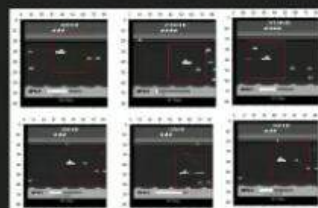
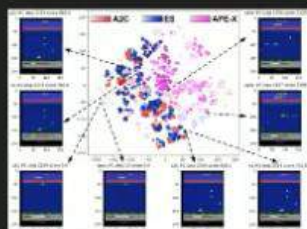


VIDEO



Accompanying Software

- Code: <http://t.uber.com/atarizoo>



VIDEO



Build on / use two previous libraries



Lucid

Lucid is a collection of infrastructure and tools for research in neural network interpretability.

<https://github.com/tensorflow/lucid>

Dopamine



Dopamine is a research framework for fast prototyping of reinforcement learning algorithms. It aims to fill the need for a small, easily grokked codebase in which users can freely experiment with wild ideas (speculative research).

Our design principles are:

- **Easy experimentation:** Make it easy for new users to run benchmark experiments.
- **Flexible development:** Make it easy for new users to try out research ideas.
- **Compact and reliable:** Provide implementations for a few, battle-tested algorithms.
- **Reproducibility:** Facilitate reproducibility in results. In particular, our setup follows the recommendations given by Machado et al. (2018).

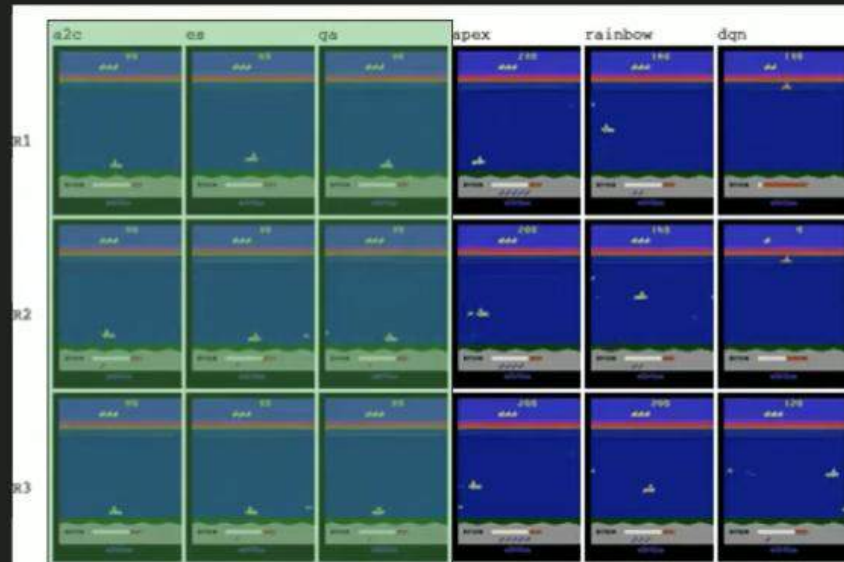
<https://github.com/google/dopamine>



VIDEO



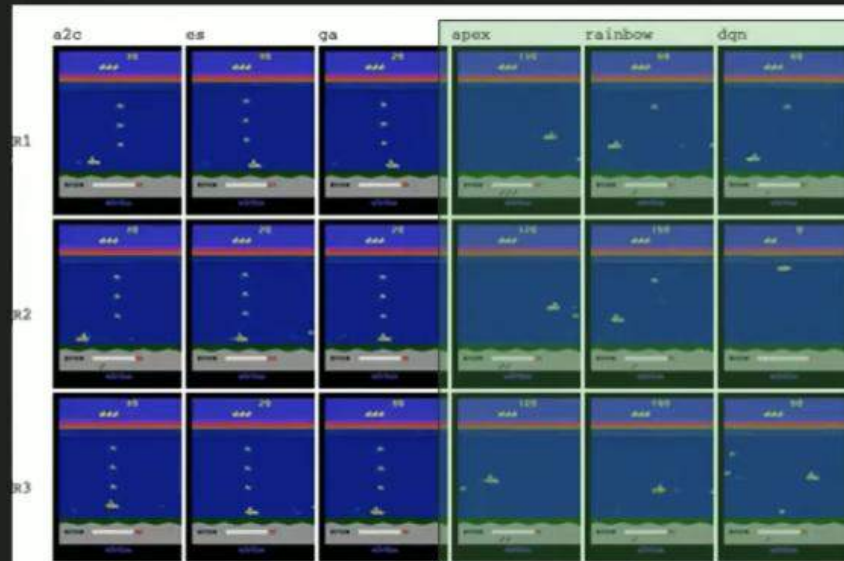
Video grids for quick qualitative comparisons



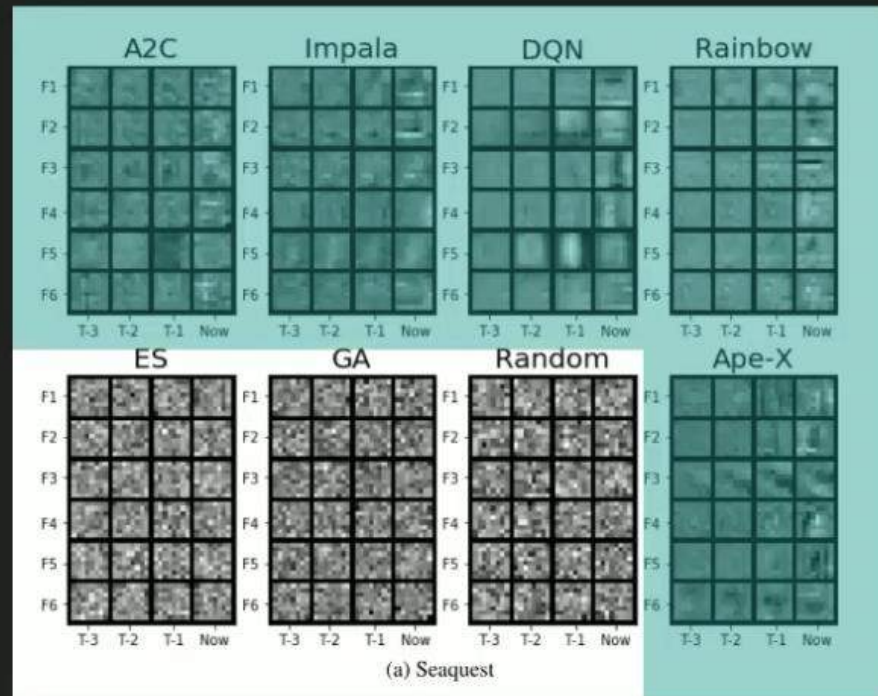
VIDEO



Video grids for quick qualitative comparisons



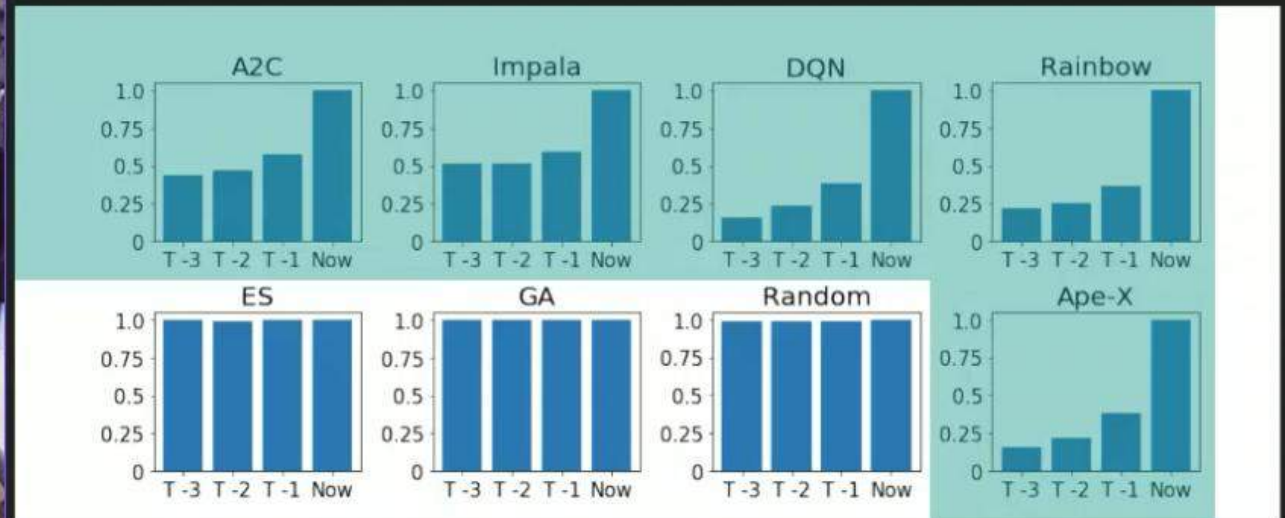
VIDEO



VIDEO



Conv1 Filters: Strength of attention to past frames

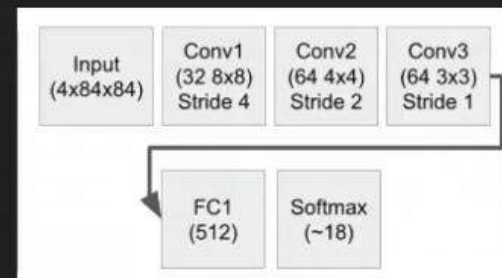


VIDEO



Robustness to Observation Noise

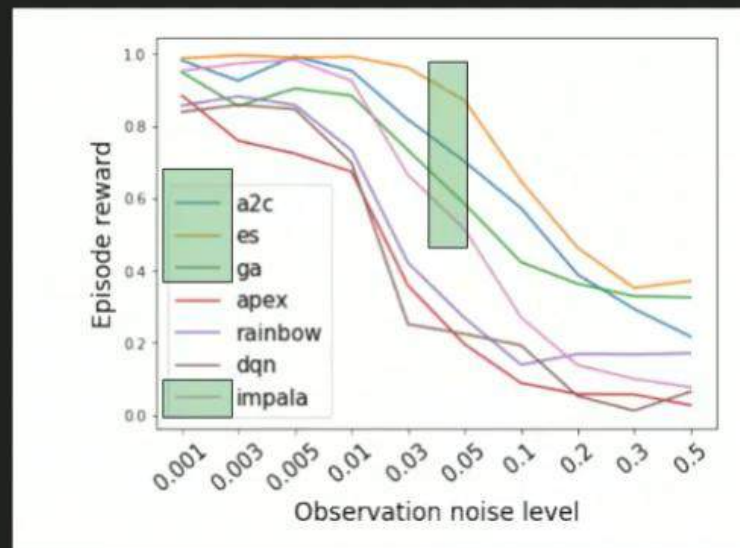
- Re-evaluate trained policy with noisy input, see how performance degrades



VIDEO



Policy Search more Robust to Observation Noise

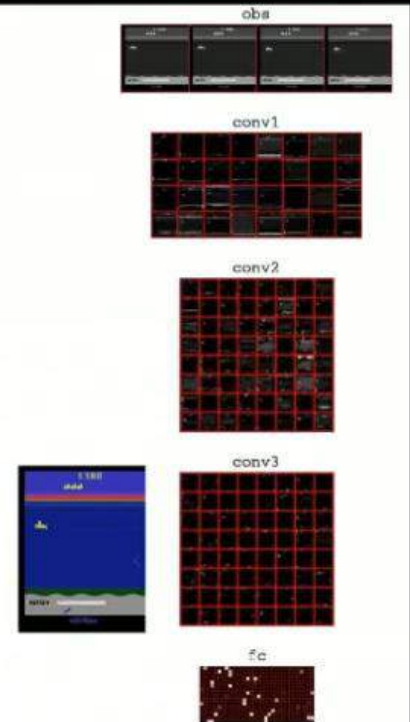


VIDEO



Visualizing activation during evaluation

In the style of DeepVis Toolbox
(Yosinski et al. 2015)



VIDEO



VIDEO



[illegible]

VIDEO



Future Zoo Directions

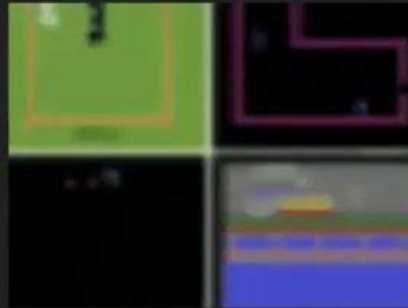
- More models (more training algorithms + different environments)
- More analysis tools / metrics
 - Implement methods from supervised learning interpretability for deep RL
- What discoveries are waiting to be found?



VIDEO



An Atari Model Zoo for Analyzing, Visualizing, and Comparing Deep Reinforcement Learning Agents



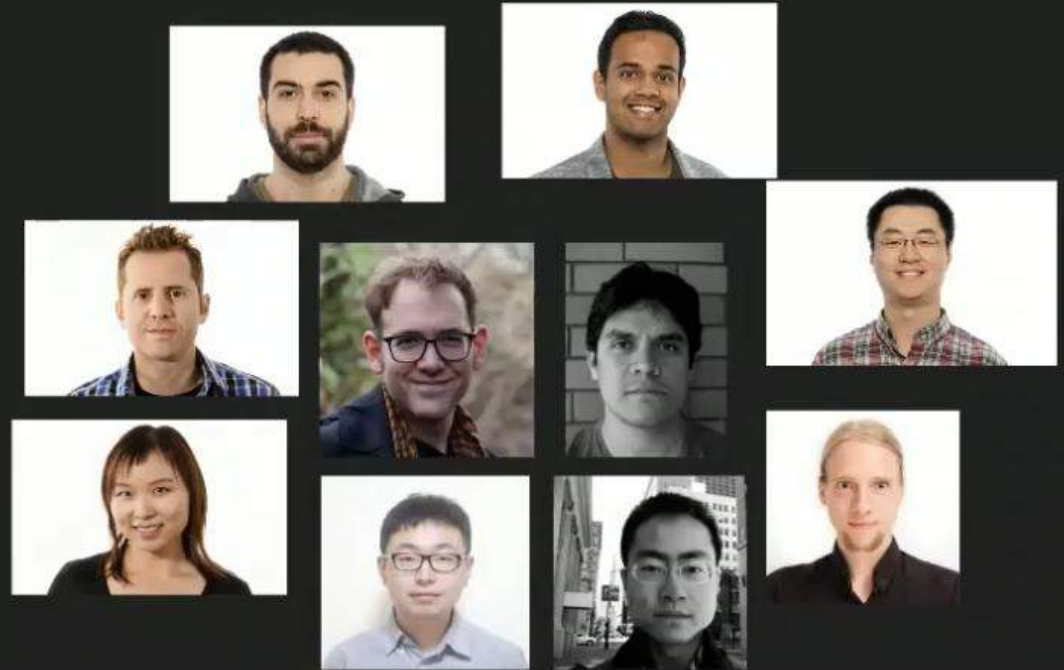
Source @ <http://t.uber.com/atarizoo>



Felipe Petroski Such, Vashisht Madhavan, Rosanne Liu,
Rui Wang, Pablo Samuel Castro, Yulun Li, Jiale Zhi,
Ludwig Schubert, Marc Bellemare, Jeff Clune, Joel Lehman



VIDEO



VIDEO



Conclusion

- We've released models trained in many Atari games across a range of RL algorithms, and software to easily load and analyze them
- Excited to see how the community uses the zoo and what research it enables
- **Questions?**
 - Lehman.154@gmail.com / @joelbot3000
(or send me a message through the coferece chat / message system)



Source (github): <http://t.uber.com/atarizoo>

Blog post: <https://eng.uber.com/atari-zoo-deep-reinforcement-learning/>

Web tool: <https://uber-research.github.io/atari-model-zoo/video2.html>

