

# Behaviour Suite for Reinforcement Learning

---

**Ian Osband**, Yotam Doron, Matteo Hessel, John Aslanides,  
Eren Sezener, Andre Saraiva, Katrina McKinney, Tor Lattimore, Csaba Szepesvari,  
Satinder Singh, Benjamin Van Roy, Richard Sutton, David Silver, Hado Van Hasselt

# Reinforcement learning that matters

*“Understand agent behaviour, develop better algorithms.”*



performance



safety



exploration

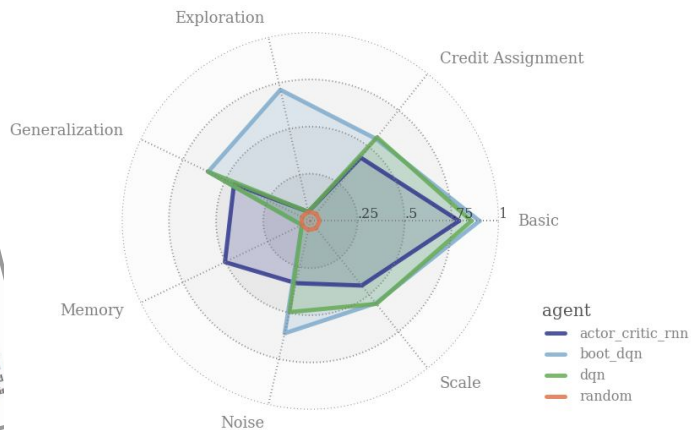
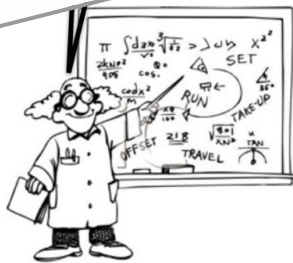


understanding

# A spectrum of agent development

**Minimax Regret Bounds for Reinforcement Learning**  
 Mohammad Gheshtlaghi Azar, Ian Osband, and Remi Munos  
 DeepMind, London, UK  
 July 4, 2017

**Abstract**  
 We consider the problem of provably optimal exploration in reinforcement learning for finite horizon MDPs. We show that an optimistic modification to value iteration achieves a regret bound of  $O(\sqrt{HSA\tau} + H^2S^2A + H\sqrt{\tau})$  where  $H$  is the time horizon,  $S$  the number of states,  $A$  the number of actions and  $\tau$  the number of time-steps. This result improves over the best previous known bound  $O(H\sqrt{SA\tau})$  achieved by the UCRL2 algorithm of Jaksch et al. (2010). The key significance of our new results is that when  $\tau \geq H^2SA$  and  $SA \geq H$ , it leads to a regret of  $O(\sqrt{HSA\tau})$  that matches the established lower bound of  $\Omega(\sqrt{HSA\tau})$  up to a logarithmic factor. Our analysis contains two key insights. We use careful applications of concentration inequalities to the optimal value function as a whole, rather than to the transition probabilities (to improve scaling in  $S$ ), and we define Bernstein-based "exploration bonuses" that use the empirical variance of the estimated values at the next states (to improve scaling in  $H$ ).



**bsuite**



# bsuite = Behaviour Suite for Reinforcement Learning

## Behaviour Suite

Ian Osband, Yotam  
Eren Sezener, Andre Saraiva, I  
Satinder Singh, Benjamin Van I

This paper introduces the **bsuite** for short. **bsuite** that investigate core capabilities of a reinforcement learning (RL) agent on these shared benchmark algorithms. Second, to study on these shared benchmark algorithms. Our research on the core issue learning algorithms. Our projects. We include examples as new reference implementations more excellent experiments from the research community, and commit to a periodic review of **bsuite** from a committee of prominent researchers.

### Trending

See what the GitHub community is most excited about today.

Repository	Stars	Forks	Build by	Stars Today
<a href="#">NVIDIA / Megatron-LM</a>	666	94		97 stars today
<a href="#">deepmind / bsuite</a>	326	21		77 stars today
<a href="#">tlbootcamp / tloadmap</a>	1,685	146		57 stars today
<a href="#">matplotlib / matplotlib</a>				

ing ( bsuite )

Assignment

1 Basic

agent

- actor\_critic\_rnn
- boot\_dqn
- dqn
- random

Scale

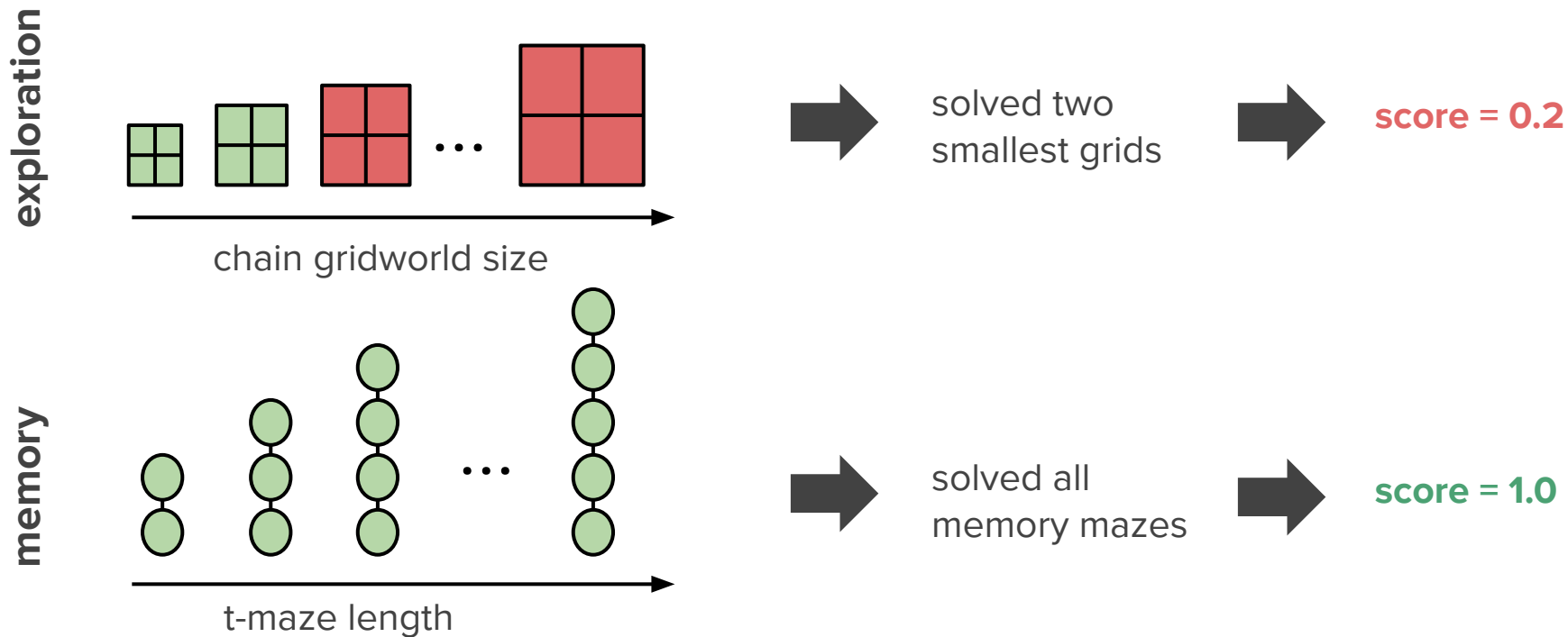
# Talk outline

1. **What is bsuite?**
2. **What type of experiments go into bsuite?**
3. **Why are these experiments useful?**
4. **How can you use bsuite going forward?**

What is bsuite?

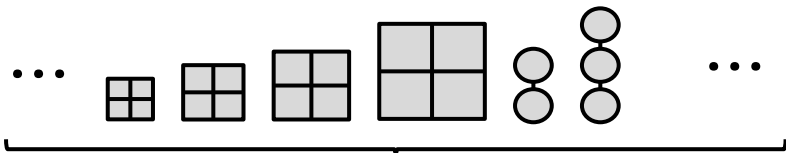
---

# bsuite is a set of *targeted* experiments



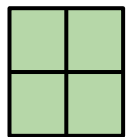
# Batteries-included sweep + analysis

```
from bsuite import sweep
for bsuite_id in sweep.SWEEP:
    run(bsuite_id=bsuite_id)
```



Flag `bsuite_id` specifies experiment.  
Easy to combine with your code.

```
import bsuite
env = bsuite.load_and_record(FLAGS.bsuite_id)
```

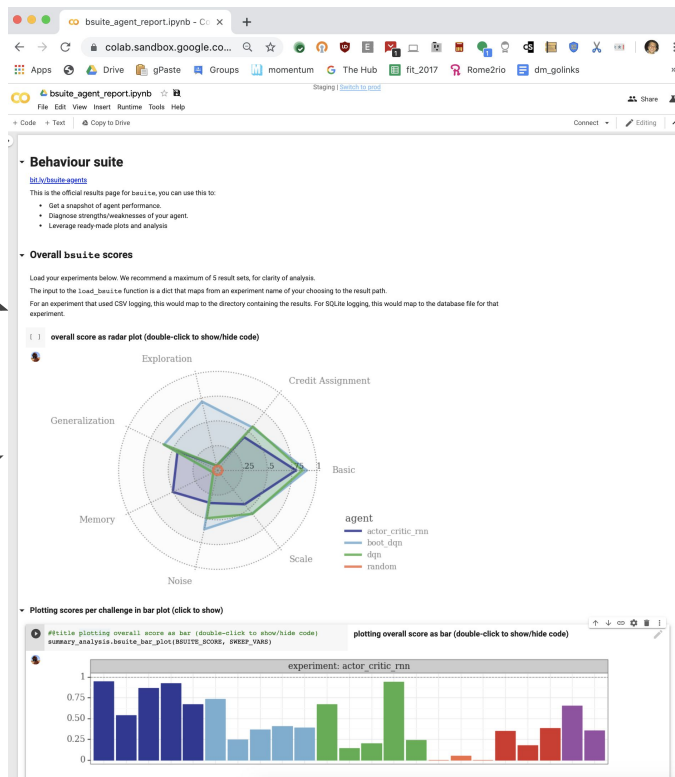


`{"reward": 0.6,`  
`"steps": 400, ...}`

environment handles  
all data logging.



Automatic analysis  
[bit.ly/bsuite-agents](https://bit.ly/bsuite-agents)





# bsuite: diagnosing agent behaviour

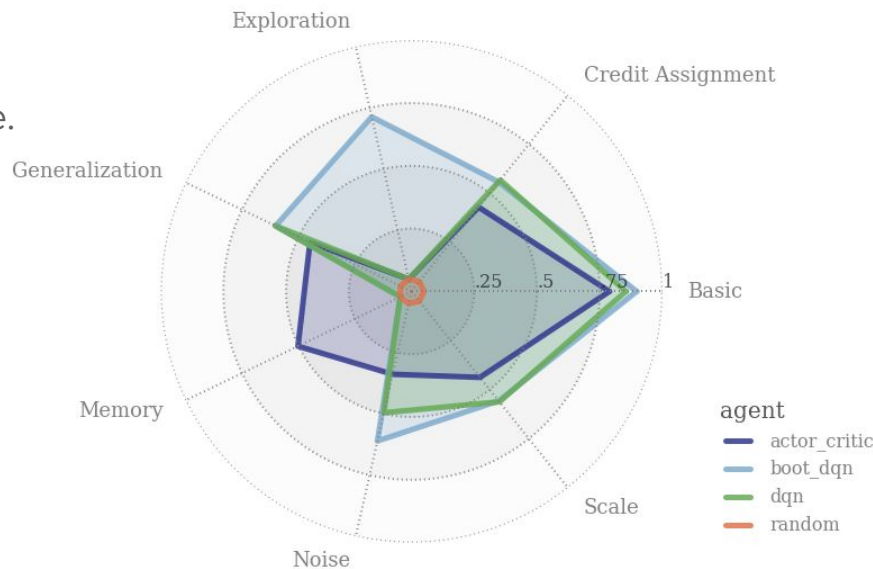
- **Quality diagnostic experiments.**

- **simple:** strip away confounding factors.
- **challenging:** push agents beyond normal range.
- **scalable:** alg scalability, not performance.
- **fast:** iterate from launch to results in < 30min.

- **Easy to use**

- run from **any agent framework.**
- integrate with < 10 lines of code.
- ready-made analysis colab.

- [github.com/deepmind/bsuite](https://github.com/deepmind/bsuite)



# It's easy to use bsuite

## 1. Swap out the environment

```
env = environments.load(name, settings)
```

```
import bsuite  
env = bsuite.load_and_record(FLAGS.bsuite_id)
```

## 2. Add one flag to your launch

```
from bsuite import sweep  
for bsuite_id in sweep.SWEEP:  
    run(bsuite_id=bsuite_id)
```

**bsuite**

environments

logging

analysis

launch

What are these experiments?

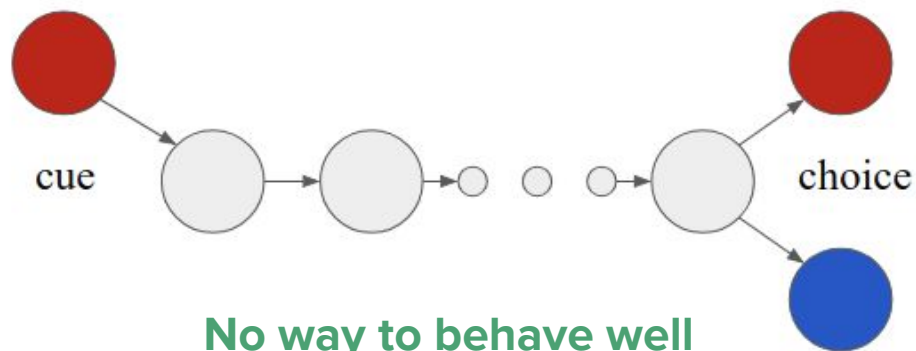
---

# What is in a bsuite experiment?

1. **Environments**: fixed set of environments that the agent interacts with.
2. **Interaction**: fixed regime of agent/environment interaction.
3. **Analysis**: fixed procedure mapping agent behaviour to results and plots.

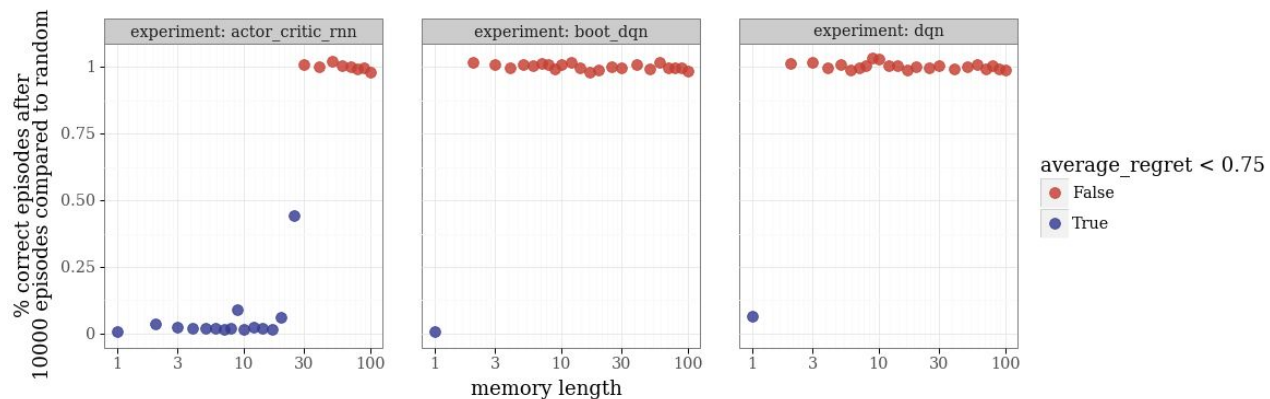
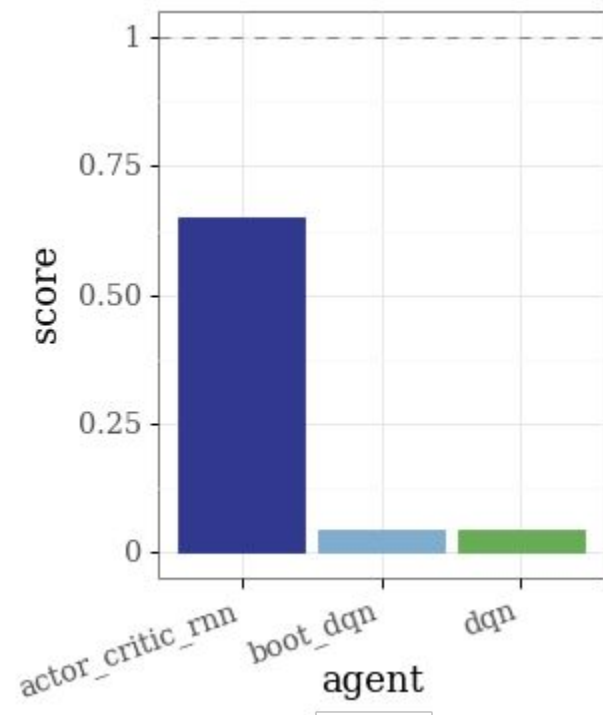
# Example 1: Memory length

- “How long can the agent remember?”
- Simple T-maze.
  - Observe context.
  - Wait N-steps.
  - Action matches context  $r=+1, -1$
- Experiment definition
  - Run for sizes  $N=1..100$
  - 10k episodes
  - Score = % of N where regret < 75 random



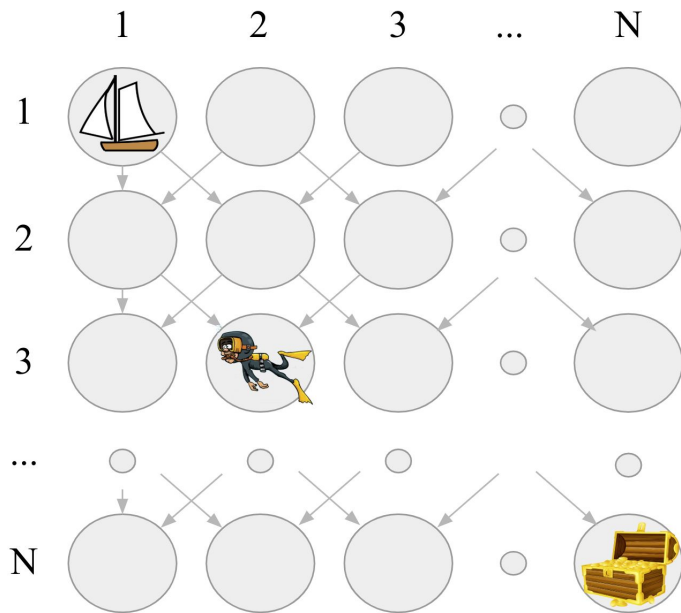
**No way to behave well  
without memory.**

# Example 1: Memory length



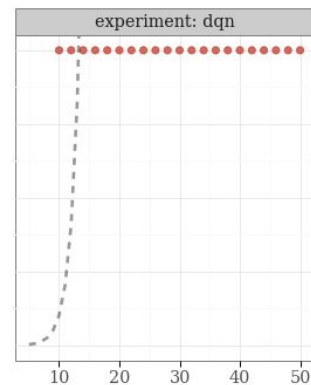
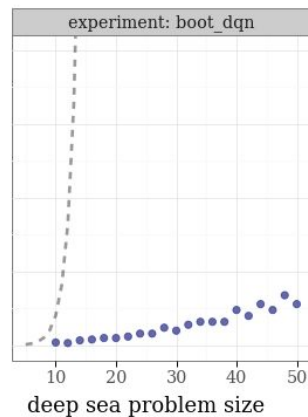
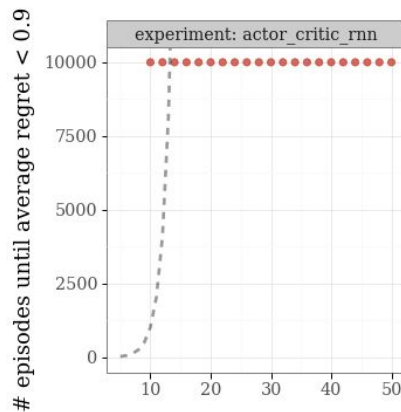
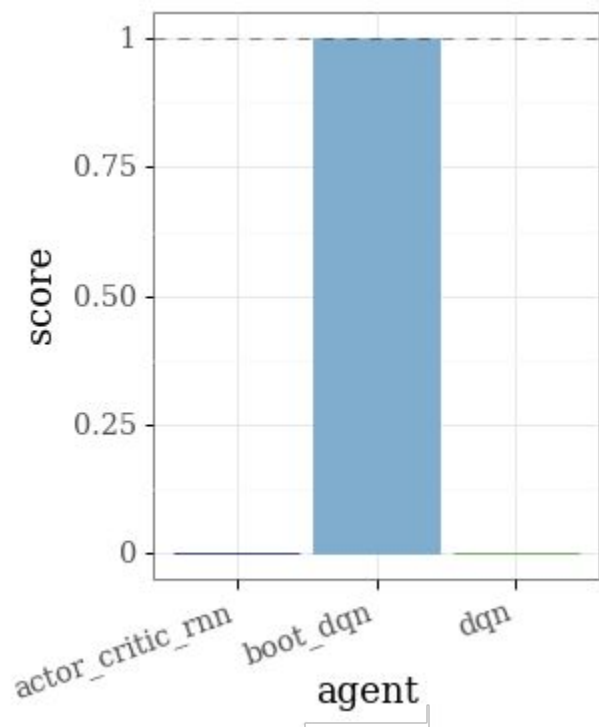
# Example 2: Deep Sea

- “Does the agent perform deep exploration?”
- Environment is an  $N \times N$  grid.
  - Start top left, fall one row per timestep.
  - Action 0 or 1 maps to “left”/“right” by state.
  - **Small cost** for going “right”.
  - **Large reward** if you make it to the treasure.
  - **Inefficient exploration  $\rightarrow 2^N$  episodes to treasure!**
- Experiment definition:
  - Run size  $N=10..50$  for 10k episodes
  - Score = % of  $N$  that have average regret  $< 0.9$  (faster than  $2^N$  baseline)



**No way to behave well  
without exploration.**

# Example 2: Deep Sea



solved  
• False  
• True



Why are these experiments useful?

---

# Why are these experiments useful?

- Practical theory often lags **practical algorithms**.
- An “MNIST” for Reinforcement Learning.
- Open source code, **reproducible research**.
- Instantiate **key issues** in RL research.
- Valuable **“unit tests”** for core RL concepts.



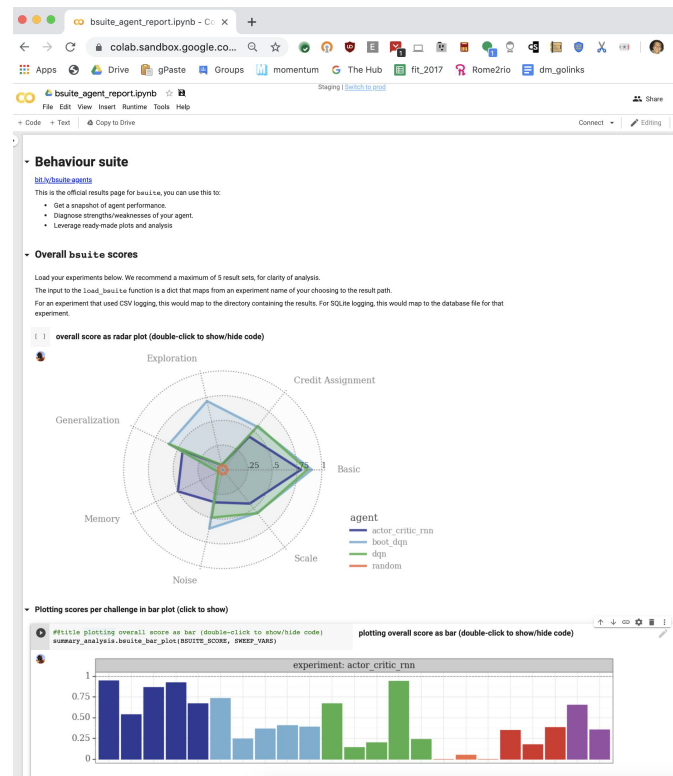
# How do you load the results?

- Pre-cooked: [bit.ly/bsuite-colab](https://bit.ly/bsuite-colab)

```
from l.d.r.bsuite import load

XIDs = {
    'name_1': xid_1,
    'name_2': xid_2,
}
DF, SWEEP_VARS = load.load_bsuite(XIDs)
```

- DF = dataframe including results.
- SWEEP\_VARS = list of hyperparams.
- **Automatically handles multiple agents and hyperparameters.**



# Interactive session

- Tutorial in colab: [bit.ly/bsuite-tutorial-colab](https://bit.ly/bsuite-tutorial-colab)
- Comparing different agents: [bit.ly/bsuite-agents](https://bit.ly/bsuite-agents)
- Changing optimizer in DQN: [bit.ly/bsuite-optimizer](https://bit.ly/bsuite-optimizer)
- Varying ensemble size bootstrapped DQN: [bit.ly/bsuite-ensemble](https://bit.ly/bsuite-ensemble)

How can you use bsuite?

---

# Automated LaTeX appendix

- Run on bsuite → 1-page summary.

- Easy to include in [conference format](#):

- ICML
- ICLR
- NeurIPS

- “Sanity check” your algorithm.

- Benchmark performance.

- Sensitivity analysis.

## C bsuite report: benchmarking baseline agents

The *Behaviour Suite for Reinforcement Learning*, or **bsuite** for short, is a collection of carefully-designed experiments that investigate core capabilities of a reinforcement learning (RL) agent. The aim of the **bsuite** project is to collect clear, informative and scalable problems that capture key issues in the design of efficient and general learning algorithms and study agent behaviour through their performance on these shared benchmarks. This report provides a snapshot of agent performance on **bsuite2019**, obtained by running the experiments from [github.com/deepmind/bsuite](https://github.com/deepmind/bsuite) [47].

### C.1 Agent definition

In this experiment all implementations are taken from **bsuite/baselines** with default configurations. We provide a brief summary of the agents run on **bsuite2019**:

- **random**: selects action uniformly at random each timestep.
- **dqn**: Deep Q-networks [38].
- **boot\_dqn**: bootstrapped DQN with prior networks [46, 45].
- **actor\_critic\_rnn**: an actor critic with recurrent neural network [37].

### C.2 Summary scores

Each **bsuite** experiment outputs a summary score in  $[0,1]$ . We aggregate these scores by according to key experiment type, according to the standard analysis notebook. A detailed analysis of each of these experiments may be found in a notebook hosted on Colaboratory [bit.ly/bsuite-agents](https://bit.ly/bsuite-agents).

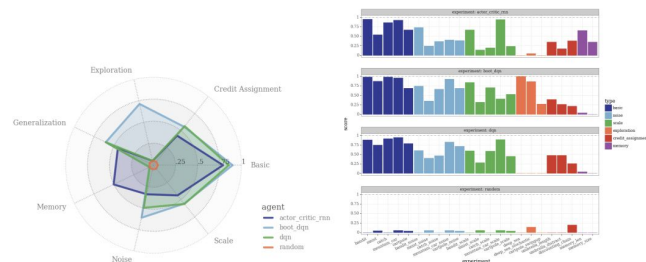


Figure 6: Radar plot gives a snapshot of agent behaviour. Figure 7: Summary score for each **bsuite** experiment.

### C.3 Results commentary

- **random** performs poorly across all aspects. This confirms that our scoring functions are working as intended.
- **dqn** performs well on basic tasks, and quite well on credit assignment, generalization, noise and scale. DQN performs extremely poorly across memory and exploration tasks. The feedforward MLP has no mechanism for memory, and  $\epsilon=5\%$ -greedy action selection is notoriously inefficient in domains that require efficient exploration.
- **boot\_dqn** performs mostly identically to DQN, except for exploration where it greatly outperforms, and also a smaller boost to performance under noise. This result matches our understanding of Bootstrapped DQN as a variant of DQN designed to estimate uncertainty and use this to guide deep exploration.
- **actor\_critic\_rnn** typically performs worse than either DQN or Bootstrapped DQN on all tasks apart from memory. This agent is the only one able to perform better than random due to its recurrent network architecture.

# We want better experiments

- This is the start of the effort, not the end.
- **Announcing bsuite committee.**
  - Joelle Pineau (Facebook)
  - John Schulmann (OpenAI)
  - ... looking for more!
- Adding a new experiment to bsuite is easy!
- Let's talk: [@lanosband](https://twitter.com/lanosband)

