

---

We want to demonstrate the performance of both approaches to safety. This requires a level playing field, so results can be compared. This leveling of the playing field is a recurring problem in AI research. To solve this problem, shared datasets have emerged and been adopted by researchers. One well known example of this is ImageNet, a database of images to do image processing on.

## 1 DeepMind Gridworlds

For safety in reinforcement learning, one such "dataset" being proposed is Gridworlds by Google DeepMind [1].

## 2 Approach

In order to familiarize ourselves more with the particulars of the papers, we will first explore their given implementations. Both papers include their respective approach implemented on a variety of games.

To evaluate performance of both ways to implement safety, we of course need a learner as baseline. We will look at the learners used in both papers, as well as the learners being used in research today. From those we will select one which is suitable to implement both approaches on.

Having obtained a learner, we will of course implement both approaches. Doing this we will evaluate and report on the ease of use of each approach. Then we will train them on the Gridworlds environment. Criteria for evaluation are convergence and the safety checks of Gridworlds given above.

## References

- [1] Jan Leike, Miljan Martic, Victoria Krakovna, Pedro A Ortega, Tom Everitt, Andrew Lefrancq, Laurent Orseau, and Shane Legg. Ai safety gridworlds. *arXiv preprint arXiv:1711.09883*, 2017.