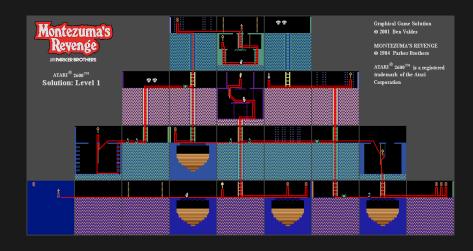


## Montezuma's Revenge



## Montezuma's Revenge



- very sparse rewards hundrets of steps
- huge state space
- hard exploration
- needs returns back

#### state of the art score

source: https://paperswithcode.com/sota/ atari-games-on-atari-2600-montezumas-revenge

year	name	score
2015	Deep Reinforcement Learning with Double Q-learning	0
2017	Curiosity-driven Exploration by Self-supervised Prediction <sup>a</sup>	0
2021	MuZero	2500
2018	Count-Based Exploration with Neural Density Models b	3705
2019	Exploration by Random Network Distillation <sup>c</sup>	8152

#### requires environment state saving/loading

GoExplore\* d

2021

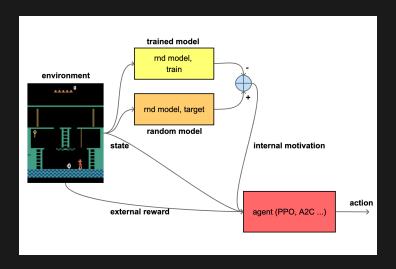
43 000

<sup>&</sup>lt;sup>a</sup>https://arxiv.org/abs/1705.05363 <sup>b</sup>https://arxiv.org/abs/1703.01310

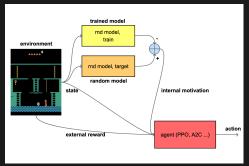
<sup>&</sup>lt;sup>c</sup>https://arxiv.org/abs/1810.12894

dhttps://arxiv.org/abs/2004.12919

#### random network distillation

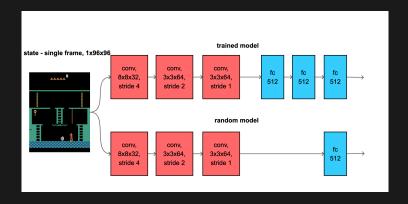


#### random network distillation

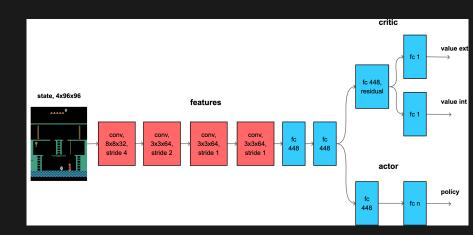


- neural network works as novelty detector
- model learns to imitate random (target) model
- less visited states produce bigger motivation signal
- orthogonal weights initialisation  $(g = 2^{0.5})$  for strong signal
- lot of fully connected layers to avoid generalisation

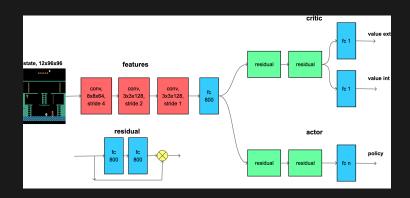
#### random network distillation architecture



## ppo model architecture



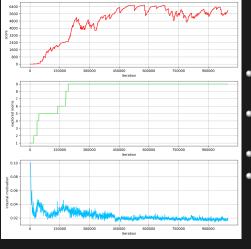
## ppo model architecture



### loss

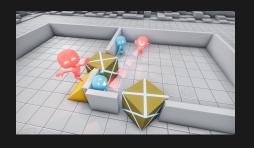
TODO

#### results



- 1M steps 20% of original paper
- 128 parallel envs = total 128M steps
- score 6400
- 9 rooms explored

## **Emergent Tool Use From Multi-Agent Autocurricula**



- multi-agent robotic environment
- hide and seek
- https:
  //openai.com/blog/
  emergent-tool-use/
- https://arxiv.org/abs/ 1909.07528

# Q&A



Michal CHOVANEC, PhD

random network distillation