# Bootstrapping Visual Understanding in Deep Reinforcement Learning Models Buck Bukaty, Amy Kanne CS 221 - Department of Computer Science, Stanford University

## Motivation

- In games, reinforcement learning algorithms must learn visual information and environmental dynamics from the ground up by taking random actions.
- Recent work shows how this lack of prior understanding puts RL algorithms at a distinct disadvantage in games intended for humans [1].
- We explore ways to use human gameplay to bootstrap the environmental understanding of deep reinforcement learning models.
- Specifically, we utilize convolutional network weights learned from human gameplay to initialize Deep Q-Network [2] and Policy Gradient [3] models.

#### Game Environment

- Our domain is the game Crypt of the NecroDancer.
- Turn based, takes place on a grid of tiles.
- Clear analogy to Markov Decision Process – take discrete action [left, right, up, down] and observe environment response.

# 1. N. 🌲 🦂

#### Infrastructure

- Created software to capture screenshots of the game window during human play and label them with keyboard inputs.
- Recorded playthroughs of 100 randomly generated levels.
- Result: labeled dataset of ~12,000 game frames and corresponding button presses.

### **RL Environment, Challenges**

No access to game state; manually scan life meter, gold meter, and level indicator for use in reward function.



×35 ZOME: | LEUEL: |

#### Architectu

Simple – Sir Deeper – Si Resnet – Sir

- Simple Fo
- Deeper Fo
- Resnet Fo
- **Optimal Hu**

#### **Policy Gradients (REINFORCE)**

We interpret our best Behavioral Cloning network as a stochastic policy network, and perform gradient updates on it according to the **REINFORCE** algorithm. [5]

We interpret our best Behavioral Cloning network as a Q-function approximator, and perform gradient updates on it according to the Q-Learning algorithm with experience replay and freezing target networks. [6]

### **Behavioral Cloning**

Trained and tested convolutional networks for predicting humanlike inputs in response to game situations.

Saw improvement when input included additional time dimension, stacked 4 images as input in the same manner as Mnih [3].

Tested network architectures, consideration given to the gamegrid-aligned nature of input images.

Transfer learning with ResNet not effective – the game's pixelated features were too dissimilar with real life curves and lines. Recorded gold acquired in 5 random test levels.

е	% Val Accuracy	1	2	3	4	5	Mean
ngle Image Input	73	16	14	0	35	0	13
ngle Image Input	77	3	0	16	35	34	17.6
ngle Image Input	78	41	31	23	18	2	23
ur Image Input	74	8	16	9	24	2	11.8
our Image Input	79	9	37	13	43	0	20.4
ur Image Input	80	9	34	39	12	46	28
man Play (Oracle)		154	185	144	200	169	170.4

Compared with a ResNet34, our architecture was much less deep with comparable performance. This was useful for intensive reinforcement learning training, when the necessary information for gradient updates for the ResNet34 overflowed GPU memory.

> Note: All reinforcement learning training seen below is performed on the same random level, functionally an MDP whose state we can only interpret visually.



— Reward Preinitialized









