

# Supplementary: EMOTE - An Explainable architecture for Modelling the Other Through Empathy

## 1 Experiment Layouts

Figure 5 illustrates each of the 5 environment layouts experimented on. In each game the Learning agent (red arrow) and independent agents (yellow and purple arrows) aim to collect their respective pellets.

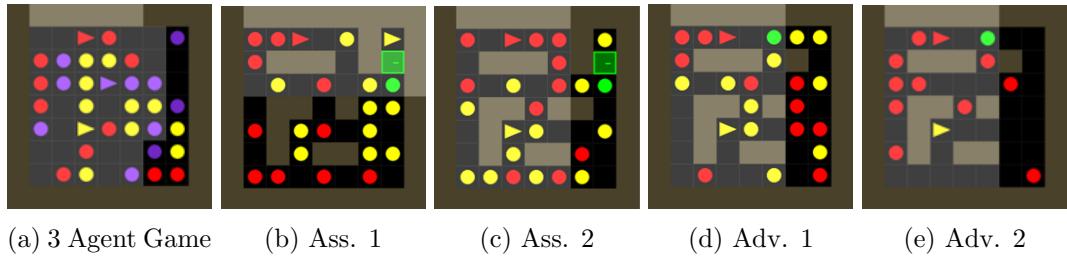


Figure 5: **3 Agents env.** (a): Learning agent can consume any coloured pellet but doing so reduces availability to other agents. **Sympathy Framework** (b)-(e) **Assistive env.** (b)-(c): Closed door (green square) only opened by learning agent **Adversarial env.** (d)-(e): Indep agent can harm learning agent. When green button pressed, learning agent can harm Indep agent, earning positive reward.

## 2 Experiment Reward Functions

Tables 3, 4 and 5 detail the reward functions for each of the environments. The reward function for the two assistive and two adversarial environments have the same reward function for the learning agent. This is to check whether the inferred reward function for the independent agent remains consistent between the two environments, even when the layout changes.

Table 3:  $R$ : 3 Agent Game

Feature	value
Consuming Learning Agent Pellet	10
Consuming Independent Agent 1's Pellet	0
Consuming Independent Agent 2's Pellet	0
Step	-1

Table 4:  $R$ : Assistive 1 and 2

Feature	value
Consuming Learning Agent Pellet	10
Consuming Independent Agent’s Pellet	0
Open Door	-1
Win	5
Step	-1

Table 5:  $R$ : Adversarial 1 and 2

Feature	value
Consuming Learning Agent Pellet	20
Consuming Independent Agent’s Pellet (Only for Adversarial 1)	0
Button Press	0
Win	30
Learning Agent Harmed	-50
Independent Agent Harmed	10
Step	-1

### 3 Experiment Hyperparameters

Tables 6, 7, 8, 9 and 10 list the experimental settings for all games. When tuning  $\delta$ , it was found that it was made easier (less sensitive) when a constant ( $> 1$ ) was multiplied with Loss term 2.  $\delta$  can be tuned without the multiplier to achieve the same performance (as with the multiplier) but it is more sensitive and as such requires finer resolution tuning. When this multiplier was set in the range [4, 8],  $\delta$  values above 0.75 had the best performance.

Table 6: 3 Agent Game

parameter	value	comment
Episodes	2000	
No. Trials	10	
State View window	5 x 5	
Batch Size	16	
$\gamma$	0.9	
optimiser	Adam	
learning rate	1e-4	$Q_{learn}$ DQN
learning rate	1e-5	Sympathy
learning rate	1e-5	E-Feature (for both Indep Agents)
learning rate	1e-5	E-Image (for both Indep Agents)
exploration initial	1	
exploration minimum	0.1	
exploration decay	0.99	
$\epsilon$ of independent agents	0.2	
target network update		every 2 episodes
E-Feature $\delta$	0.5	For both Indep Agents
E-Feature Loss 1 weight	1	For both Indep Agents
E-Feature Loss 2 weight	4	For both Indep Agents
E-Image $\delta$	0.75	For both Indep Agents
E-Image Loss 1 weight	1	For both Indep Agents
E-Image Loss 2 weight	4	For both Indep Agents
$\lambda$	0.75	
$\psi$	5e-4	

Table 7: Assistive 1

parameter	value	comment
Episodes	2000	
No. Trials	10	
State View window	5 x 5	
Batch Size	16	
$\gamma$	0.9	
optimiser	Adam	
learning rate	1e-4	$Q_{learn}$ DQN
learning rate	1e-5	Sympathy
learning rate	1e-5	E-Feature
learning rate	1e-5	E-Image
exploration initial	1	
exploration minimum	0.1	
exploration decay	0.99	
$\epsilon$ of independent agent	0.2	
target network update	every 2 episodes	
E-Feature $\delta$	0.75	
E-Feature Loss 1 weight	1	
E-Feature Loss 2 weight	4	
E-Image $\delta$	1.0	
E-Image Loss 1 weight	1	
E-Image Loss 2 weight	4	

Table 8: Assistive 2

parameter	value	comment
Episodes	2000	
No. Trials	10	
State View window	5 x 5	
Batch Size	16	
$\gamma$	0.9	
optimiser	Adam	
learning rate	1e-4	$Q_{learn}$ DQN
learning rate	1e-3	Sympathy
learning rate	1e-5	E-Feature
learning rate	1e-5	E-Image
exploration initial	1	
exploration minimum	0.1	
exploration decay	0.99	
$\epsilon$ of independent agent	0.2	
target network update	every 2 episodes	
E-Feature $\delta$	0.75	
E-Feature Loss 1 weight	1	
E-Feature Loss 2 weight	4	
E-Image $\delta$	1.0	
E-Image Loss 1 weight	1	
E-Image Loss 2 weight	4	

Table 9: Adversarial 1

parameter	value	comment
Episodes	4000	
No. Trials	10	
State View window	5 x 5	
Batch Size	16	
$\gamma$	0.9	
optimiser	Adam	
learning rate	1e-4	$Q_{learn}$ DQN
learning rate	1e-4	Sympathy
learning rate	1e-3	E-Feature
learning rate	1e-4	E-Image
exploration initial	1	
exploration minimum	0.1	
exploration decay	0.998	
$\epsilon$ of independent agent	0.2	
target network update		every 2 episodes
E-Feature $\delta$	0.95	
E-Feature Loss 1 weight	1	
E-Feature Loss 2 weight	4	
E-Image $\delta$	0.95	
E-Image Loss 1 weight	1	
E-Image Loss 2 weight	8	

Table 10: Adversarial 2

parameter	value	comment
Episodes	4000	
No. Trials	10	
State View window	5 x 5	
Batch Size	16	
$\gamma$	0.9	
optimiser	Adam	
learning rate	1e-4	$Q_{learn}$ DQN
learning rate	1e-4	Sympathy
learning rate	1e-3	E-Feature
learning rate	1e-4	E-Image
exploration initial	1	
exploration minimum	0.1	
exploration decay	0.998	
$\epsilon$ of independent agent	0.2	
target network update		every 2 episodes
E-Feature $\delta$	0.95	
E-Feature Loss 1 weight	1	
E-Feature Loss 2 weight	4	
E-Image $\delta$	0.95	
E-Image Loss 1 weight	1	
E-Image Loss 2 weight	8	

## 4 Performance

Figure 6 plots results from the Sympathy Framework, namely win rate, total rewards of the learning agent, and whether the door was opened (Assistive games) or the independent agent was harmed by the learning agent (Adversarial games) over the training period. Figure 7 presents the results from the 3 Agent Game of the learning agent's win rate, as well as the number of Independent agent 1 and 2's pellets consumed by the learning agent.

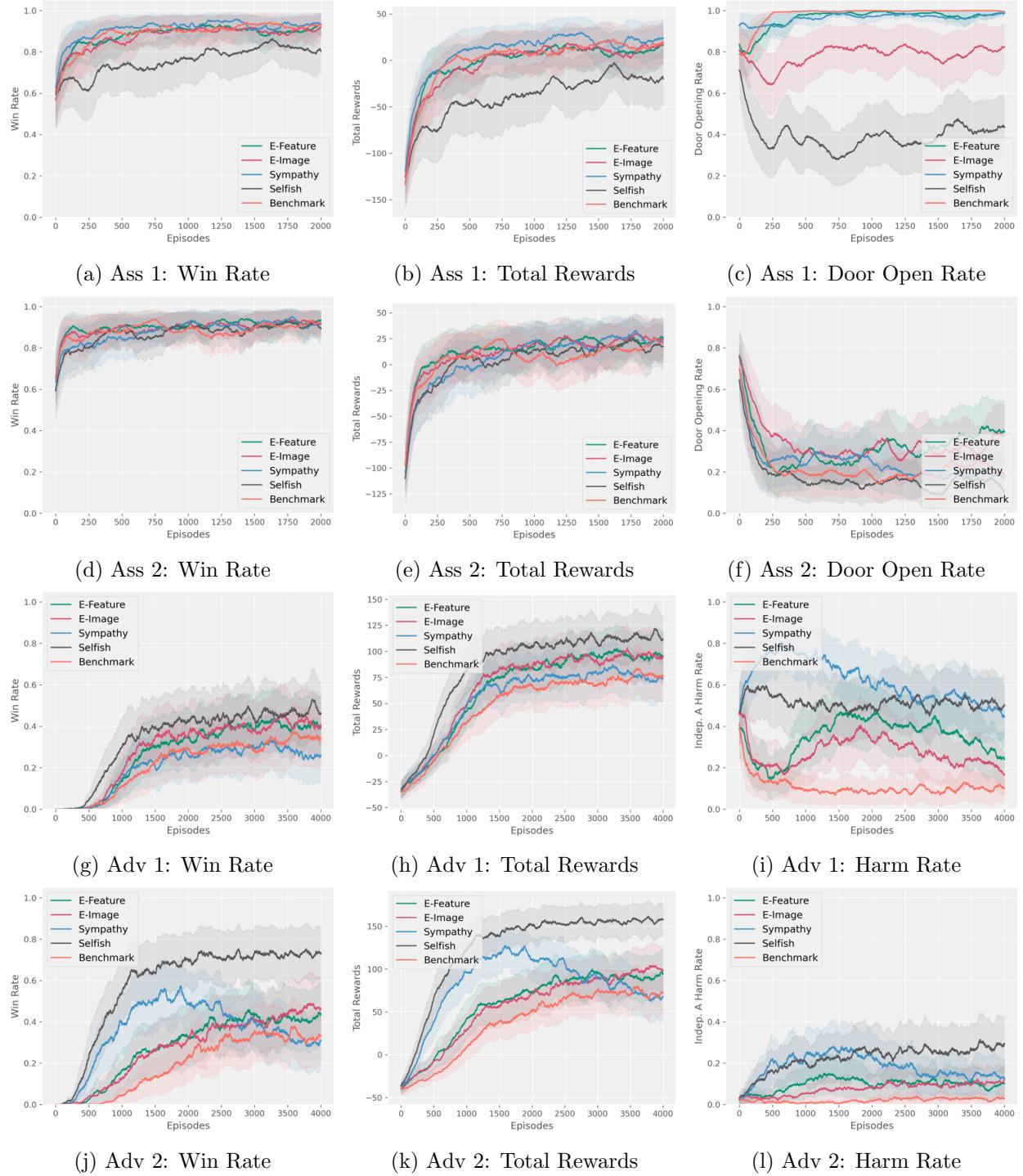


Figure 6: Sympathy Framework: Performance in each game. Assistive 1 (a-c), Assistive 2 (d-f), Adversarial 1 (g-i) and Adversarial 2 (j-l)

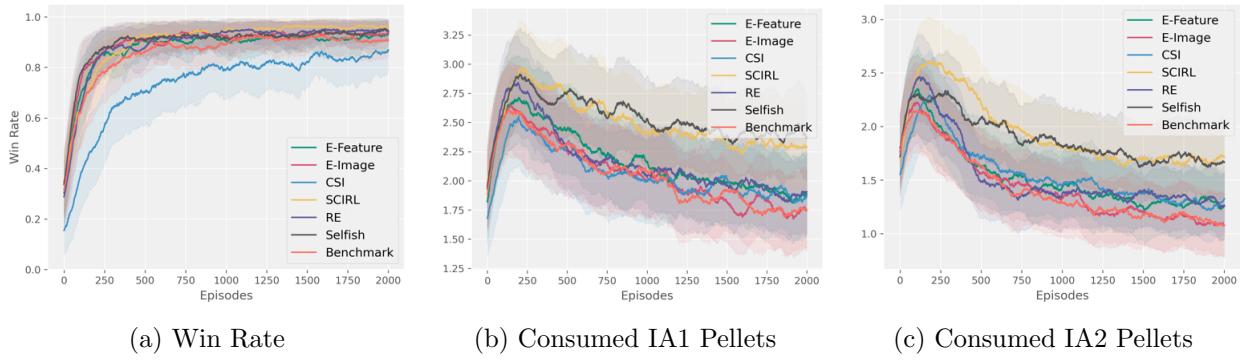


Figure 7: 3 Agent Game: Performance of learning agent. b) and c) show the number of independent agent 1 and 2's (IA1 and IA2) pellets consumed by the learning agent.

## 5 IRL Trends

Figure 8 shows the trend of the inferred rewards of the independent agent ( $\hat{R}_{indep}$ ) over the training period under the various runs of the Sympathy Framework. The weights of the Sympathy rewards have been scaled to have a  $l1$  norm equal to that of the learning agent. Figure 9 shows similar results from the 3 Agent Game. Figure 9 a) and b) shows the IRL trend for independent agent 1 and 2, respectively. The weights of the CSI, SCIRL and RE rewards have also been scaled to have a  $l1$  norm equal to that of the learning agent.

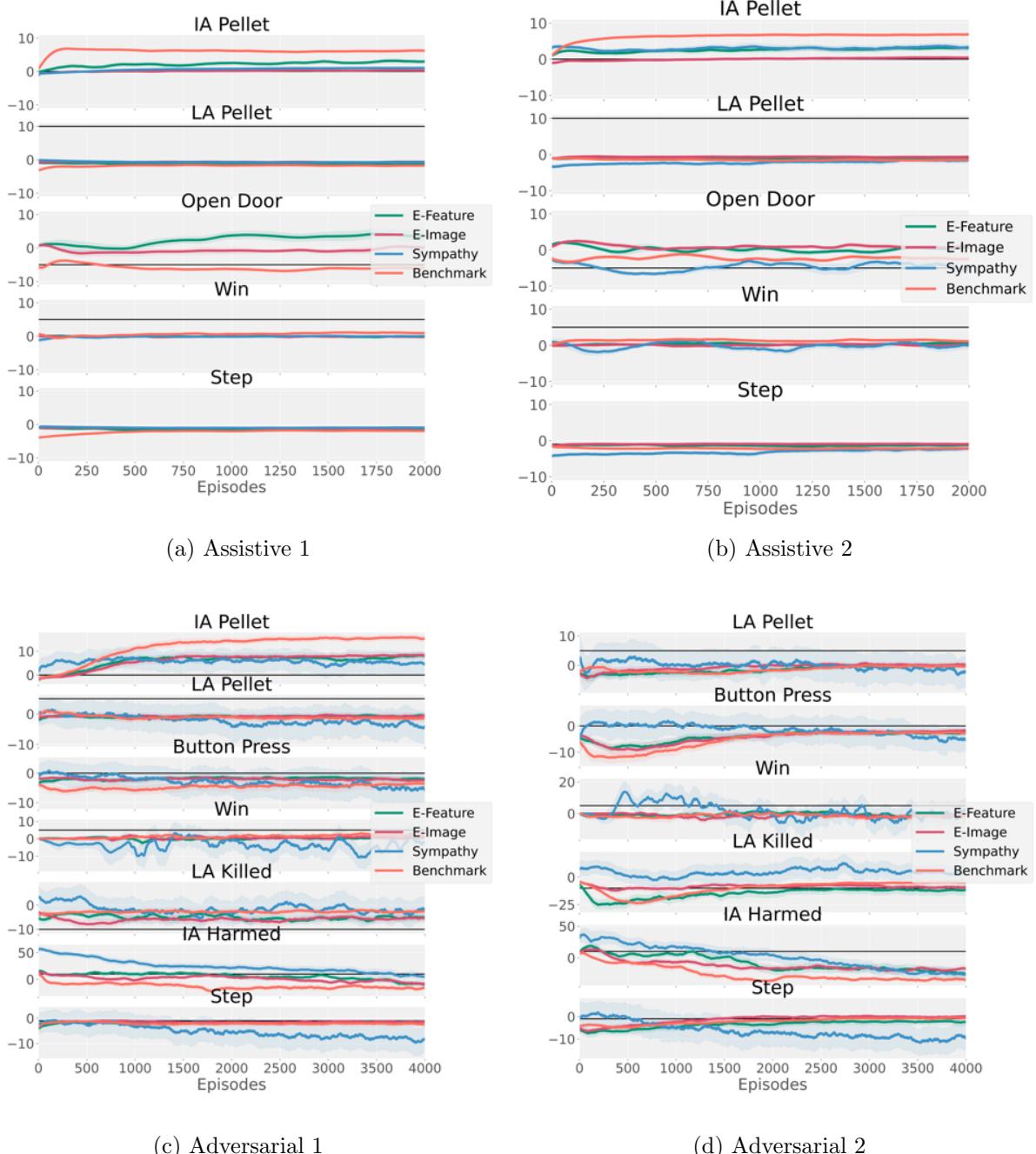
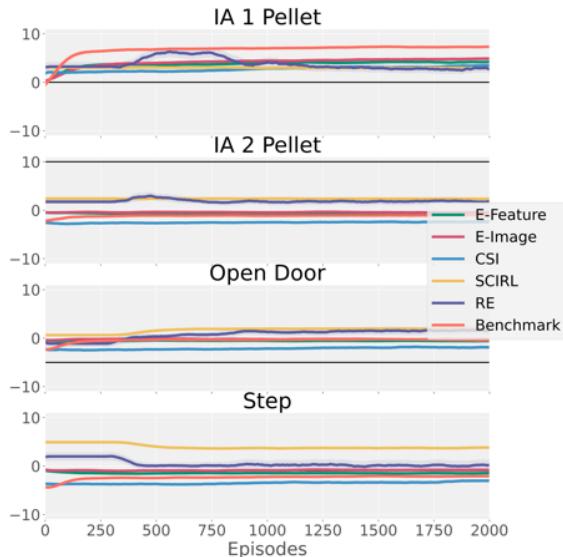
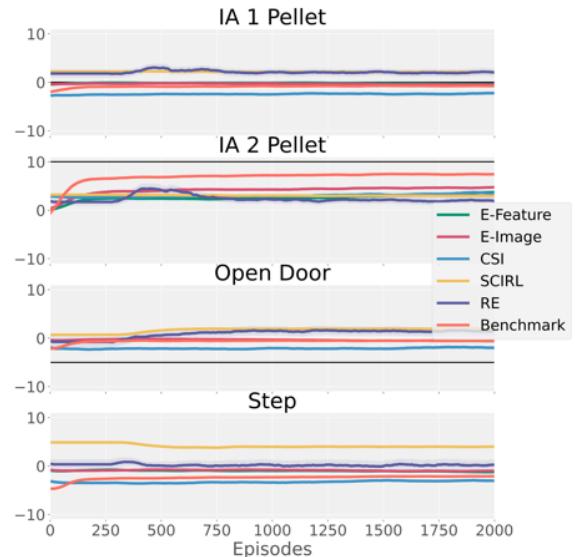


Figure 8: Sympathy Framework: Trend of  $\hat{R}_{indep}$  over training. IA - Independent Agent, LA - Learning Agent



(a) Assistive 1



(b) Assistive 2

Figure 9: 3 Agent Game: Trend of  $\hat{R}_{indep}$  over training. IA1/IA2 - Independent Agent, LA - Learning Agent

## 6 $\delta$ Sensitivity

### 6.1 $\delta$ Performance

Figures 10 and 11 show the impact of altering the hyperparameter  $\delta$  in the 3 Agent games on the win rate and the number of Independent Agent 1 and 2's pellets consumed by the learning agent under both a E-Feature and E-Image Imagination Networks, respectively. Figures 12 and 13 show the impact of altering the hyperparameter  $\delta$  in the Sympathy Framework games on the win rate, total rewards, and either door open rate or rate of harm of the independent agent by the learning agent for each of the four games under both a E-Feature and E-Image Imagination Networks, respectively.

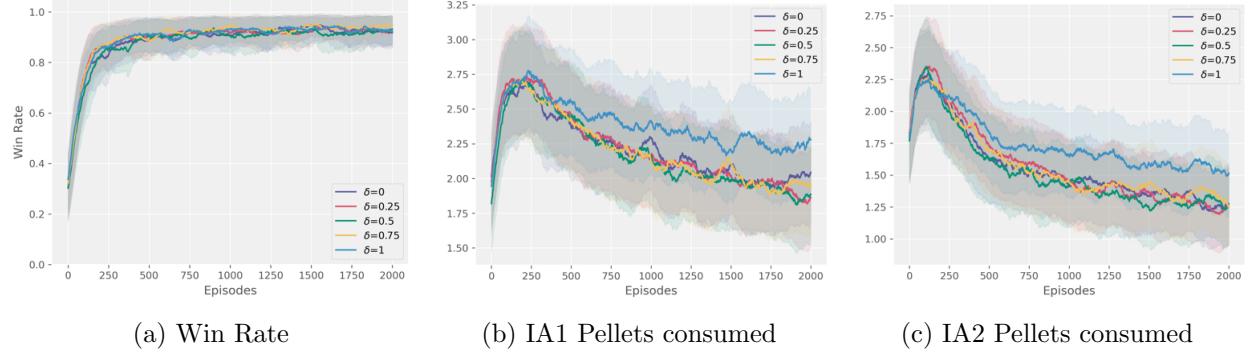


Figure 10: 3 Agent Game E-Feature: Impact of varying  $\delta$  on learning agent's performance as measured by the win rate, and number of independent agent 1 and 2's pellets consumed by the learning agent.

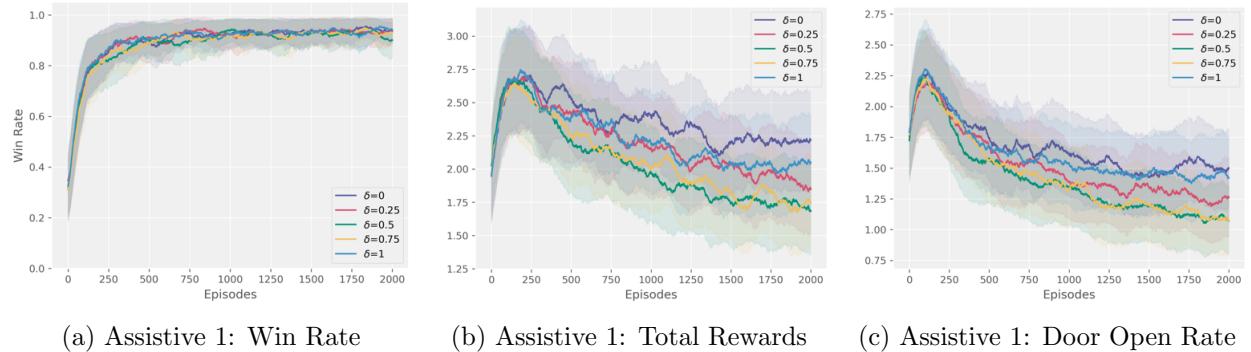


Figure 11: 3 Agent Game E-Image: Impact of varying  $\delta$  on learning agent's performance as measured by the win rate, and number of independent agent 1 and 2's pellets consumed by the learning agent.

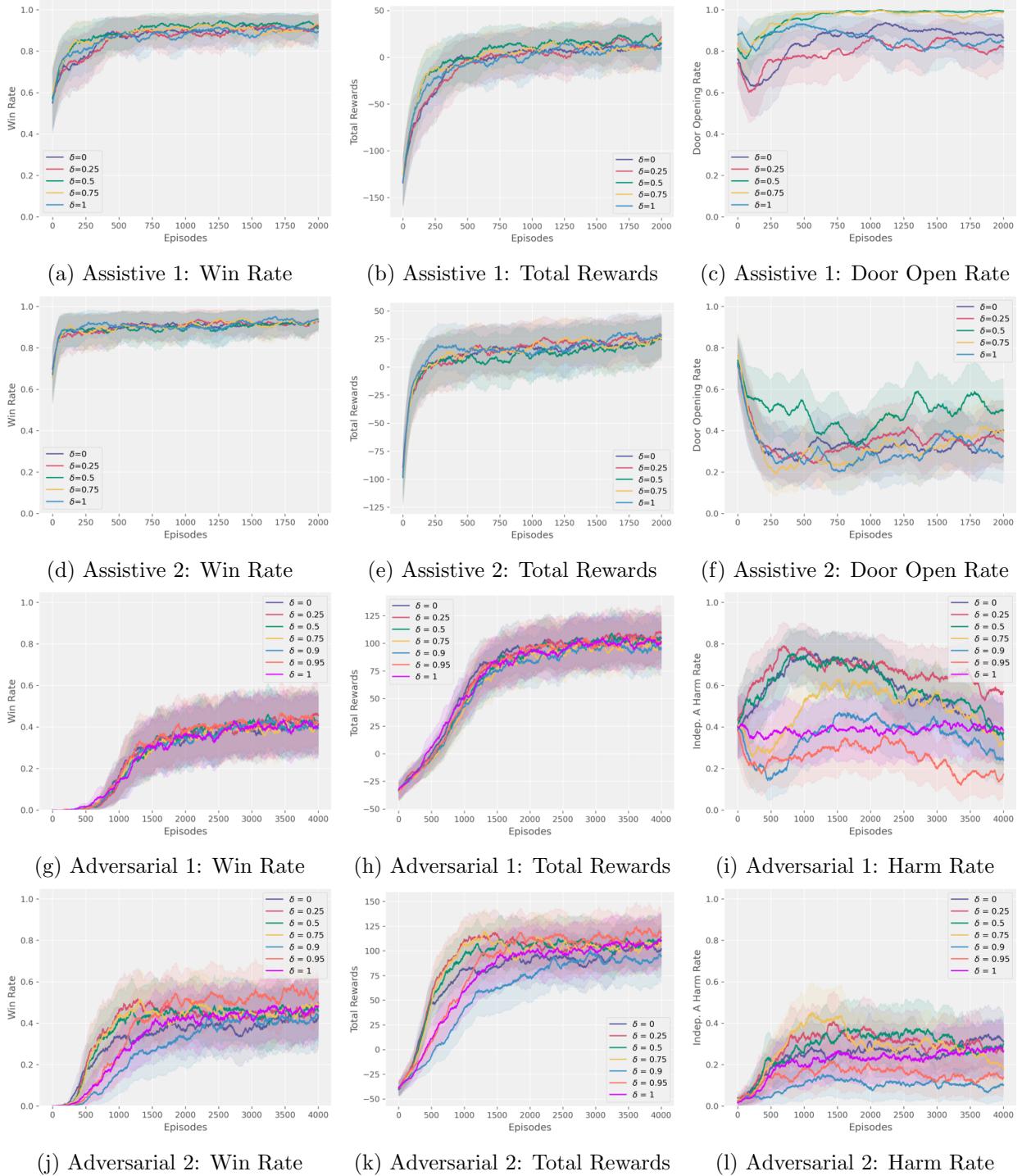


Figure 12: Sympathy Framework E-Feature: Impact of varying  $\delta$  on learning agent's performance. Assistive 1 (a-c), Assistive 2 (d-f), Adversarial 1 (g-i) and Adversarial 2 (j-l)

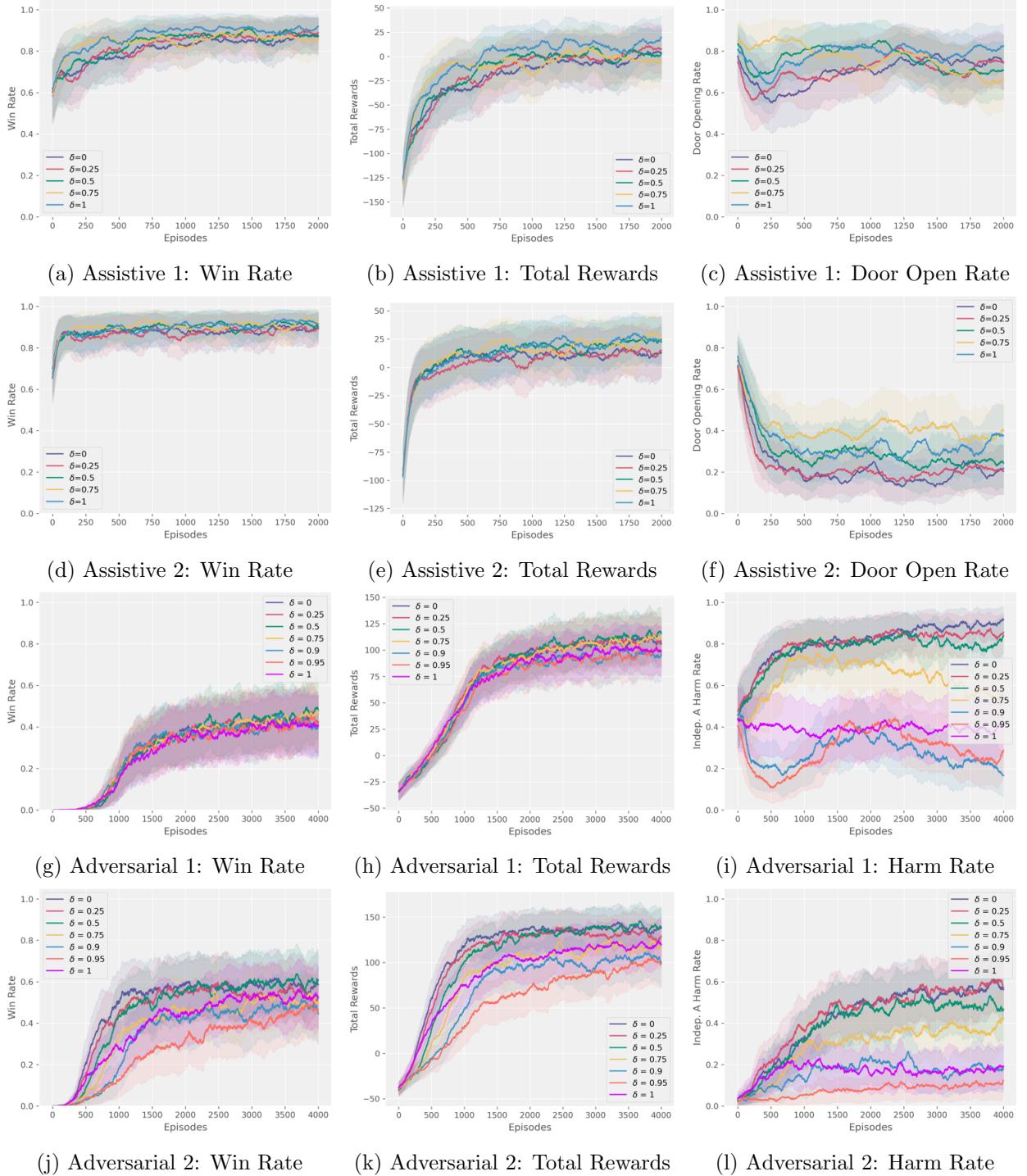


Figure 13: Sympathy Framework E-Image: Impact of varying  $\delta$  on learning agent's performance. Assistive 1 (a-c), Assistive 2 (d-f), Adversarial 1 (g-i) and Adversarial 2 (j-l)

## 6.2 $\delta$ IRL Trends

Figures 14 and 15 illustrate the impact to the resulting  $\hat{R}_{indep}$  trends by altering  $\delta$  on both the E-Feature and E-Image Imagination Networks, respectively for the Sympathy Framework. Figures 16 and 17 illustrate the impact to the resulting  $\hat{R}_{indep}$  trends by altering  $\delta$  on both the E-Feature and E-Image Imagination Networks, respectively in the 3 Agent game.

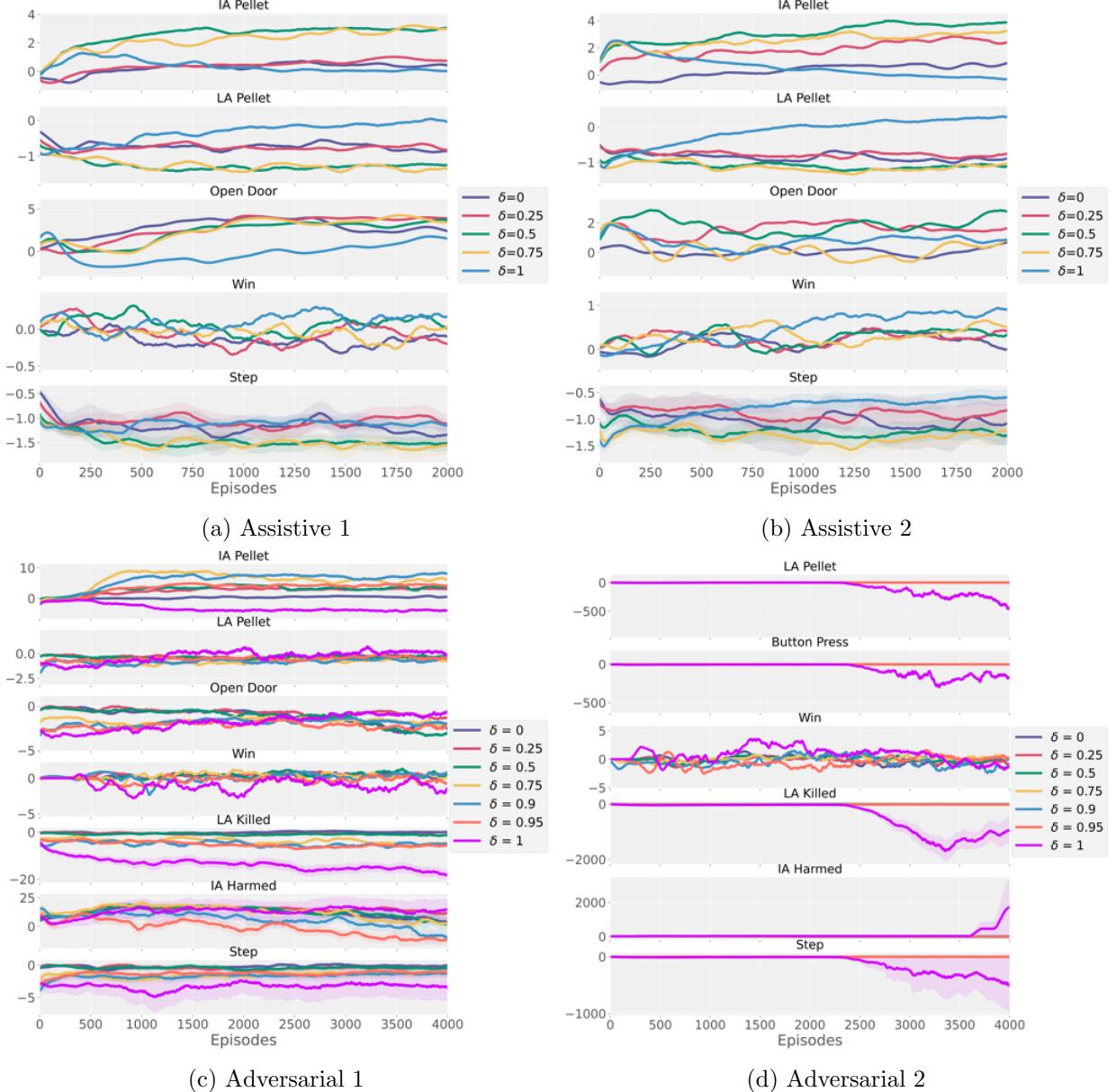


Figure 14: Sympathy Framework E-Feature: Trend of the estimated reward of the independent agent  $\hat{R}_{indep}$  as  $\delta$  is varied. IA - Independent Agent, LA - Learning Agent

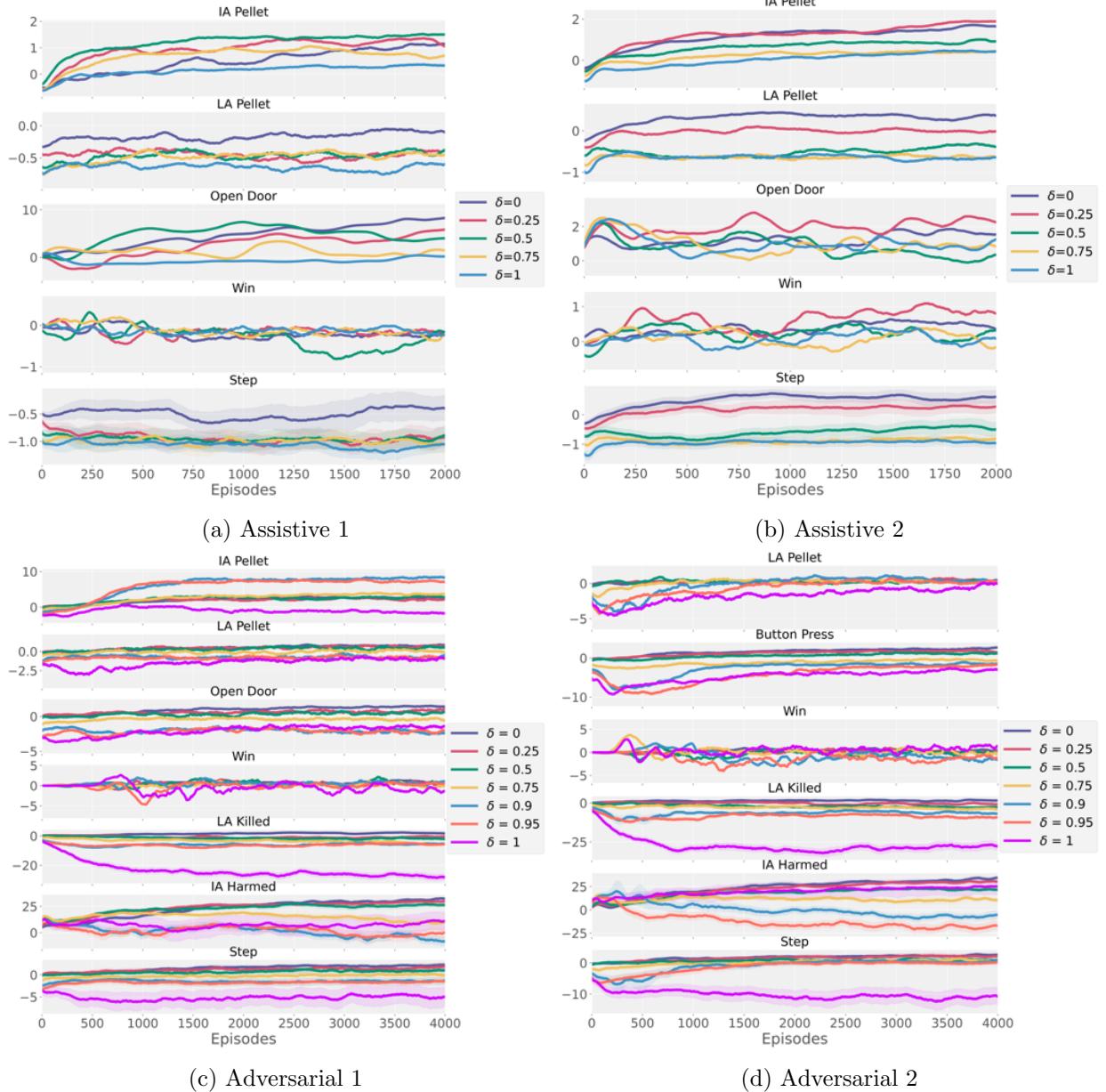


Figure 15: Sympathy Framework E-Image: Trend of the estimated reward of the independent agent  $\hat{R}_{indep}$  as  $\delta$  is varied. IA - Independent Agent, LA - Learning Agent

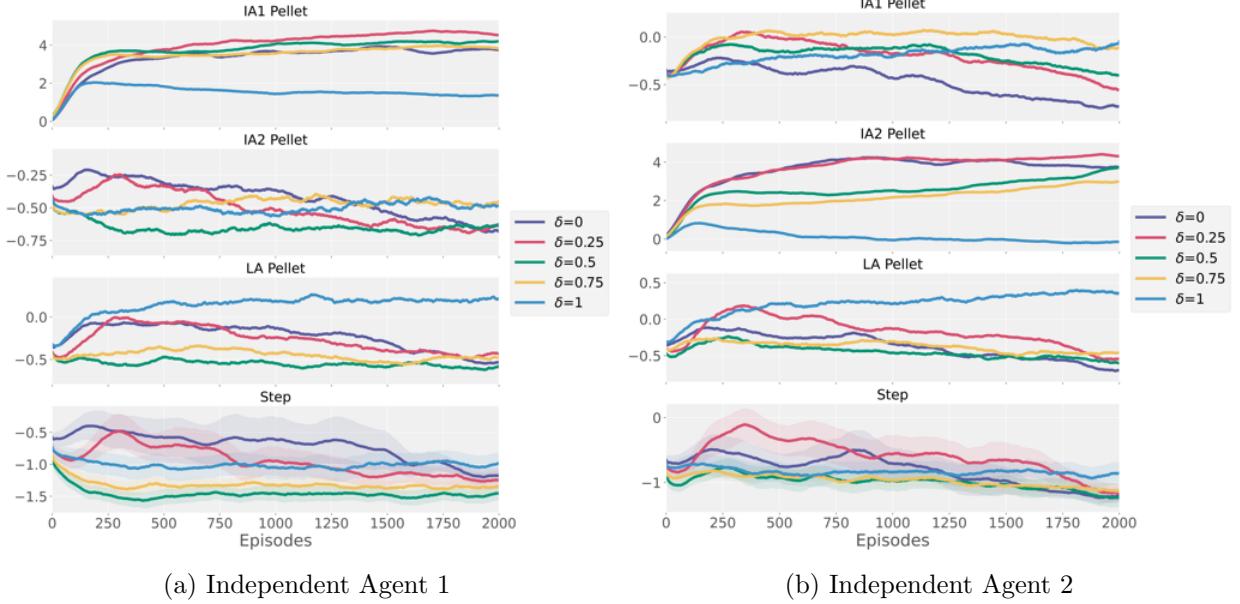


Figure 16: 3 Agent Game E-Feature: Trend of the estimated reward of the independent agent;’s  $\hat{R}_{indep}$  as  $\delta$  is varied. IA1/IA2 - Independent Agent 1 and 2, LA - Learning Agent

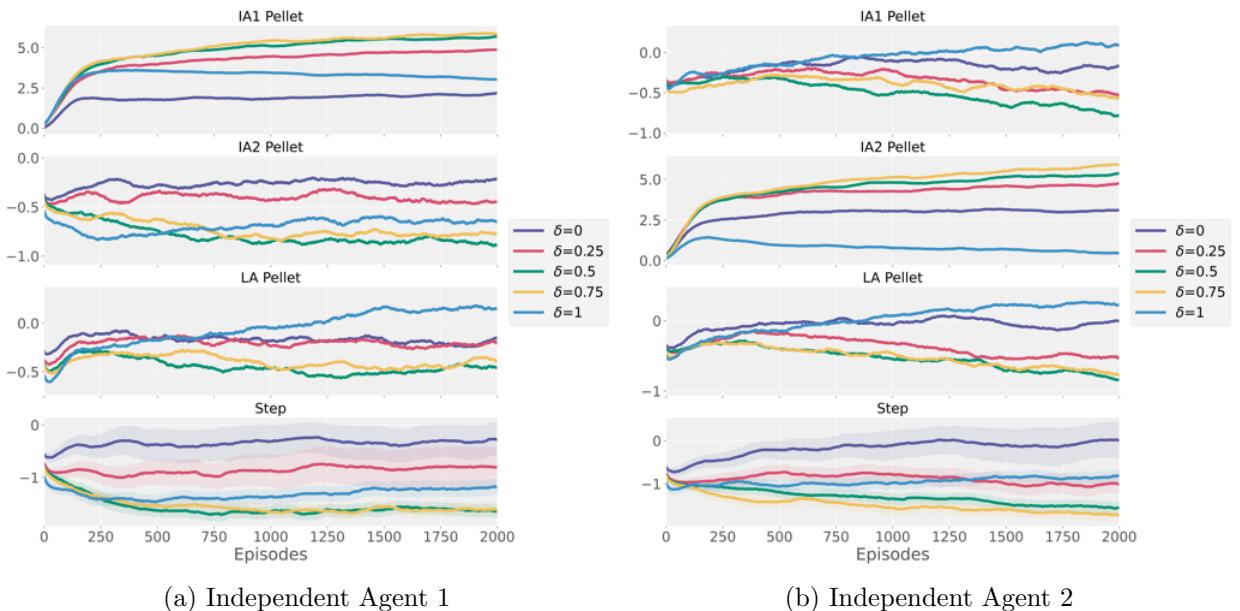


Figure 17: 3 Agents E-Image: Trend of the estimated reward of the independent agent ’s  $\hat{R}_{indep}$  as  $\delta$  is varied. IA1/IA2 - Independent Agent 1 and 2, LA - Learning Agent

### 6.3 Imagination Network Pre-training

It was found beneficial to under an initial period of training the Imagination model to take in as input  $s_i$  and produce an  $s_e$  that matched  $s_i$ . This was done to support the imagination network by quickly identifying the common features between the learning agent and the observed (independent) agent, and, after meeting a threshold  $\psi$  of average mean squared error between  $s_i$  and  $s_e$ , switch to the loss function presented in Equation 2. Figure 18 and 19 compares performance in the 3 Agent Game without pre-training and two settings of  $\psi$  for E-Feature and E-Image, respectively. With  $\psi$  better performance was achieved overall, which in our case the setting selected was  $\psi = 5e - 4$ .

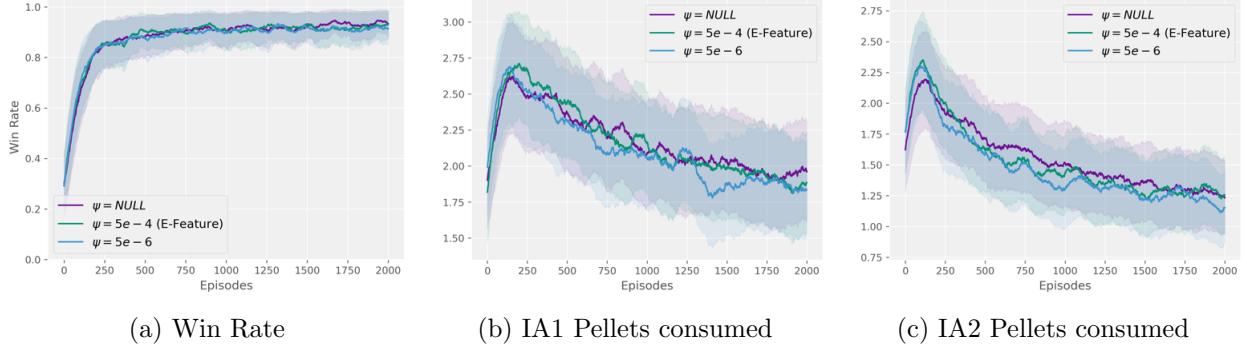


Figure 18: 3 Agent Game E-Feature: Impact of applying pre-training of the imagination network. Comparison between no pre-training ( $\psi = \text{NULL}$ ) and two  $\psi$  thresholds on learning agent’s performance as measured by the win rate, and number of independent agent 1 and 2’s pellets consumed by the learning agent.

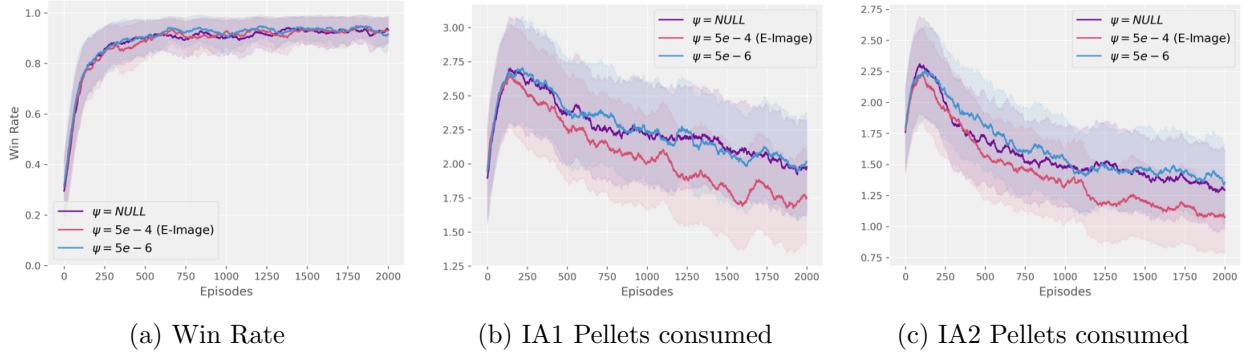


Figure 19: 3 Agent Game E-Image: Impact of applying pre-training of the imagination network. Comparison between no pre-training ( $\psi = \text{NULL}$ ) and two  $\psi$  thresholds on learning agent’s performance as measured by the win rate, and number of independent agent 1 and 2’s pellets consumed by the learning agent.

## 7 Empathetic States: 3 Agent Game

### 7.1 Generating Empathetic States $s_e$

For the original states  $s_i$  input into the Imagination Networks shown in Col 1 of Figure 5 (shown again in Figure 20), Table 11 lists the corresponding value of each of the shown state features for the 3 Agent game. In order to generate the presented empathetic state representations, we manually specified (based on the ground truth values of each feature) ranges for each feature, which was then used to reconstruct the state features. For example, for a cell with value 0.82, this corresponded to "Other Agent", and was thus rendered in orange in the reconstructed state.

Table 11: 3 Agent Game: Value of each state feature

Feature	$s_i$ value	$s_e$ range
Floor	0	0.00 - 0.05
Learning Agent Pellet (LP)	0.13	0.05 - 0.35
Indep Agent Pellet (IP1)	0.38	0.35 - 0.45
Indep Agent Pellet (IP2)	0.50	0.45 - 0.65
Other Agent (O)	0.85	0.65 - 0.90
Wall	1	0.90 - 1.00

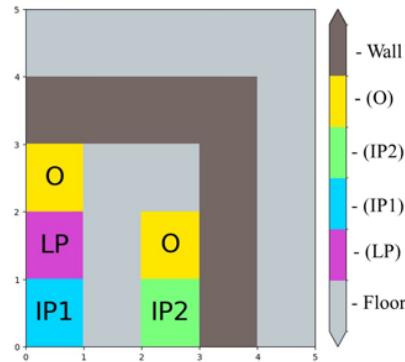
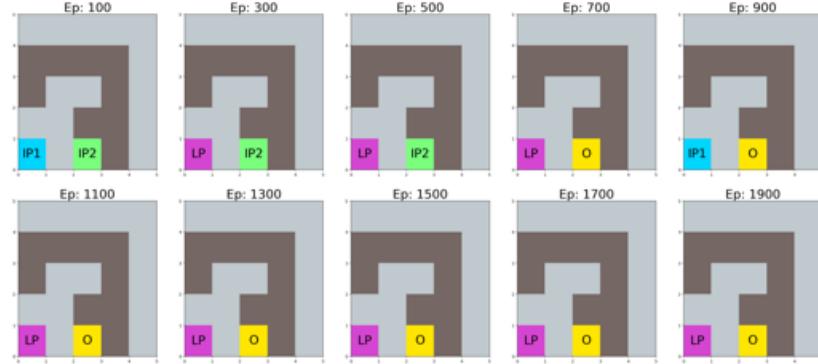


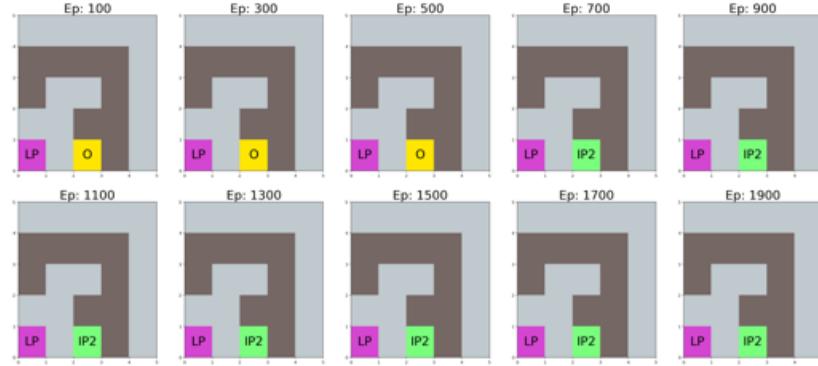
Figure 20: Original 3 Agent State  $s_i$  prior to transformation by the Imagination Networks.

### 7.2 Evolution of Empathetic State over time (during training)

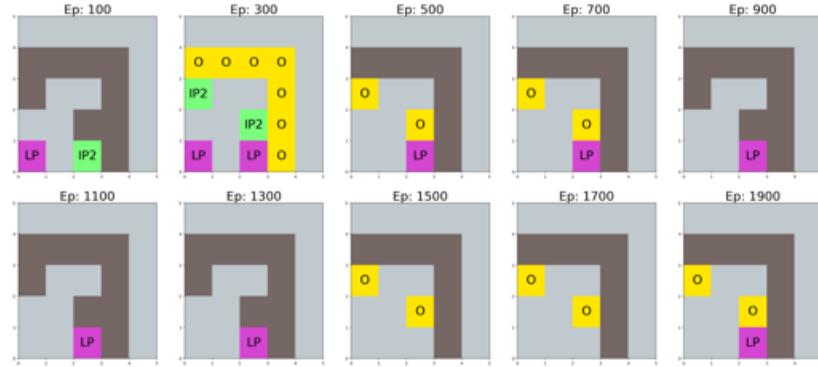
For each of the final empathetic states  $s_e$  shown in Figure 5, Figures 21 and 22 below show the evolution of the empathetic state during training. These are for both the E-Feature and E-Image Imagination Network generated  $s_e$ , shown for both Independent Agent 1 and 2.



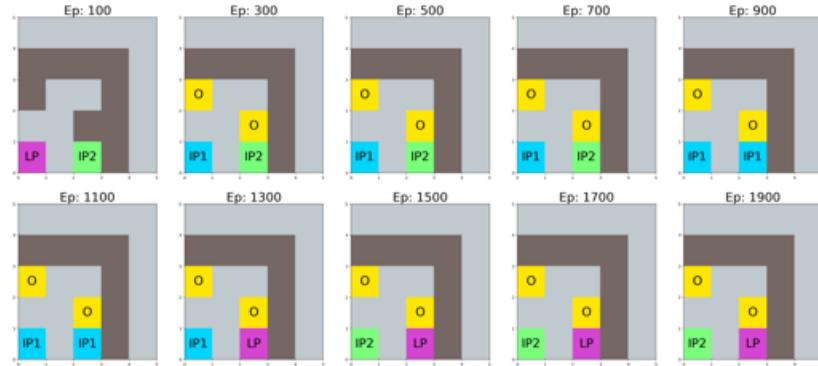
(a) Independent Agent 1 Example 1: E-Feature  $s_e$  over time



(b) Independent Agent 1 Example 2: E-Feature  $s_e$  over time

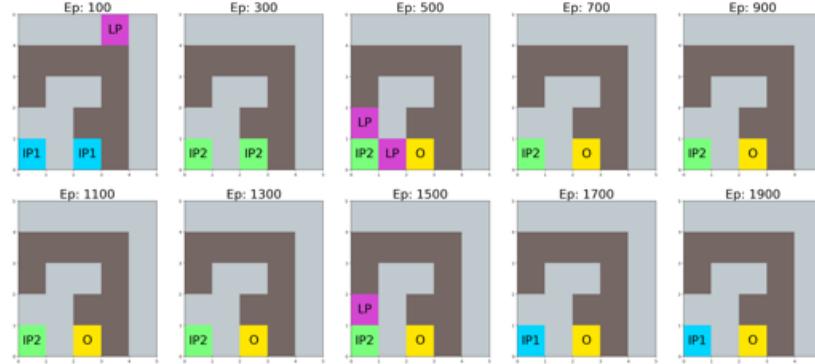


(c) Independent Agent 2 Example 1: E-Feature  $s_e$  over time

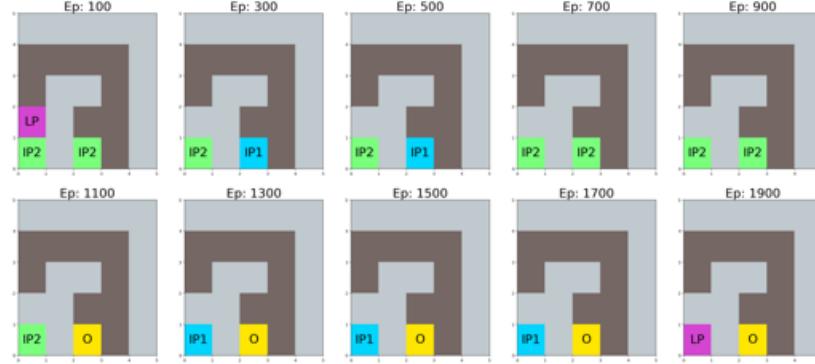


(d) Independent Agent 2 Example 2: E-Feature  $s_e$  over time

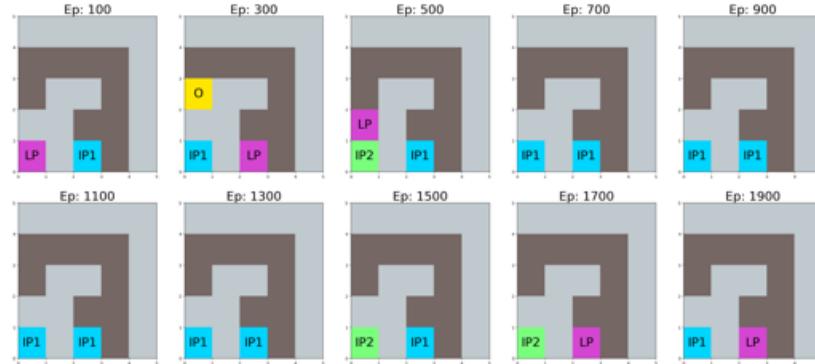
Figure 21: 3 Agent Game: Examples of the generated Empathetic State  $s_e$  over the training period under the E-Feature Imagination Networks for the final states shown in Figure 5.



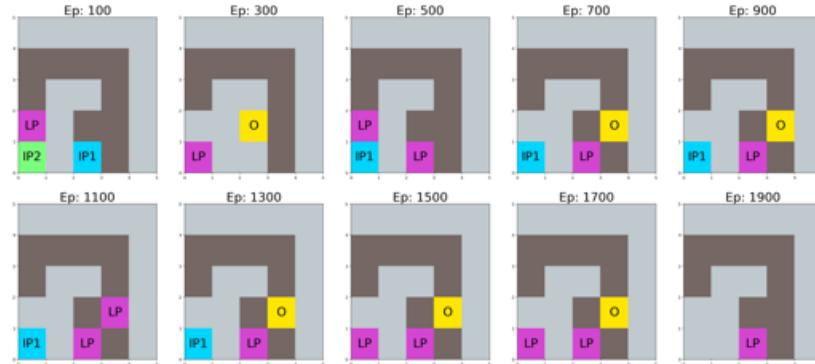
(a) Independent Agent 1 Example 1: E-Image  $s_e$  over time



(b) Independent Agent 1 Example 2: E-Image  $s_e$  over time



(c) Independent Agent 2 Example 1: E-Image  $s_e$  over time



(d) Independent Agent 2 Example 2: E-Image  $s_e$  over time

Figure 22: 3 Agent Game: Examples of the generated Empathetic State  $s_e$  over the training period under the E-Image Imagination Networks for the final states shown in Figure 5.

### 7.3 Final Empathetic States $s_e$ from random trials

Figures 23 and 24 show the final empathetic state from the 3 Agent Game for all the randomly initialised trials of each game (10 of each). Results from both E-Feature and E-Image Imagination Networks are shown for both Independent Agent 1 and 2.

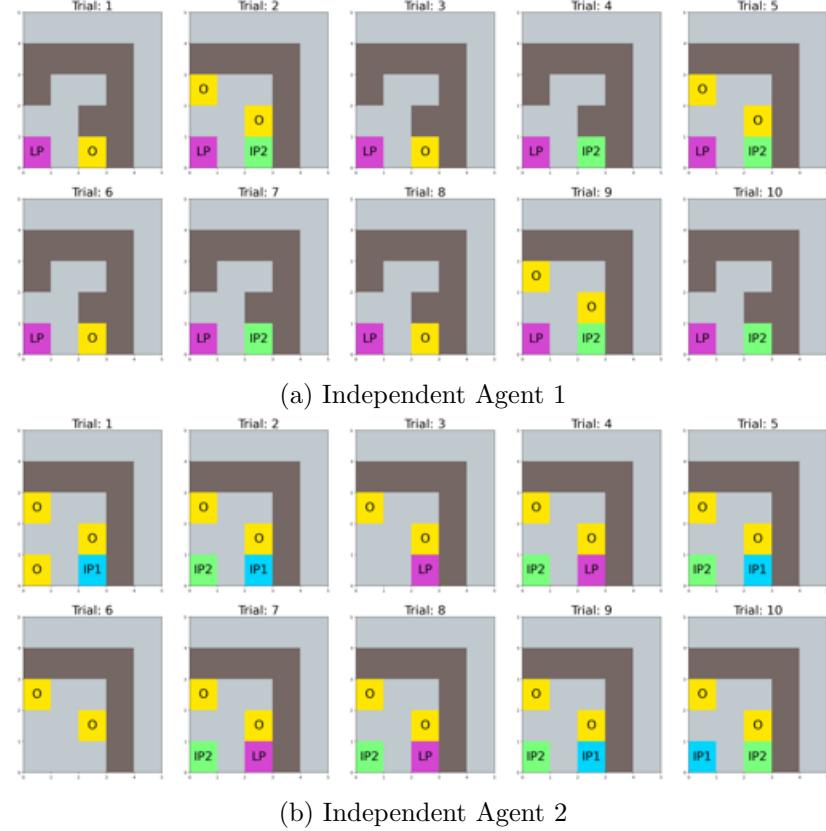
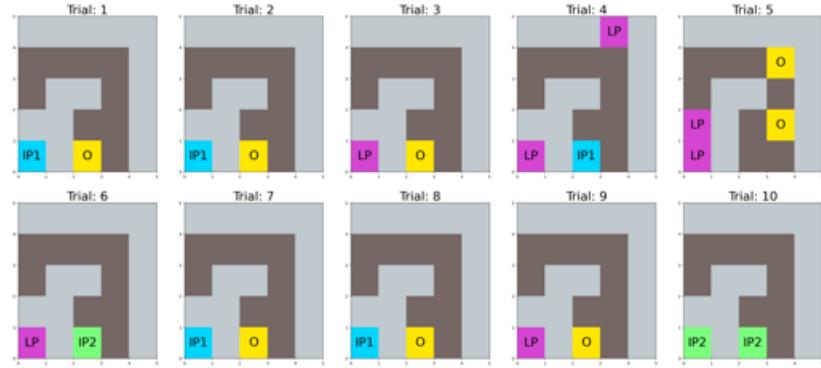
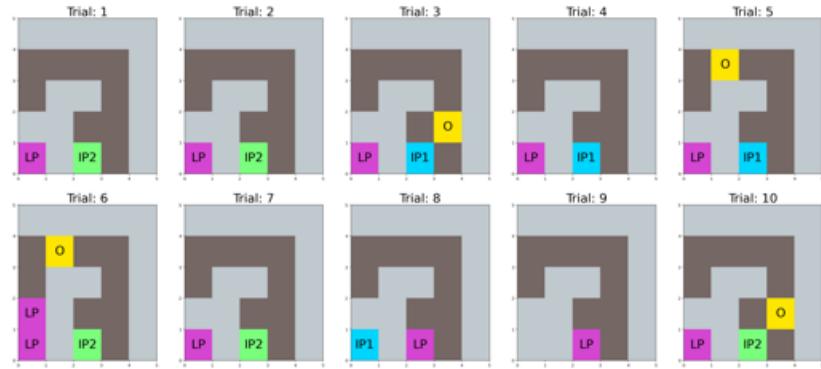


Figure 23: 3 Agent Game: Final Empathetic state  $s_e$  produced via E-Feature Imagination Network for 10 random trials of each game, for both Independent Agents.



(a) Independent Agent 1



(b) Independent Agent 2

Figure 24: 3 Agent Game: Final Empathetic state  $s_e$  produced via E-Image Imagination Network for 10 random trials of each game, for both Independent Agents.

## 8 Empathetic States: Sympathy Framework

### 8.1 Generating Empathetic States $s_e$

For the original states  $s_i$  input into the Imagination Networks shown in Col 1 of Figure 6, Table 12 and Table 13 list the corresponding value of each of the shown state features for the Assistive and Adversarial games, respectively. In order to generate the presented empathetic state representations, we manually specified (based on the ground truth values of each feature) ranges for each feature, which was then used to reconstruct the state features. For example, for a cell with value 0.82, this corresponded to "Other Agent", and was thus rendered in orange in the reconstructed state.

Table 12: Assistive Games: Value of each state feature

Feature	$s_i$ value	$s_e$ range
Floor	0	0.00 - 0.05
Learning Agent Pellet (LP)	0.13	0.05 - 0.25
Indep Agent Pellet (IP)	0.38	0.25 - 0.40
Button (B)	0.5	0.40 - 0.57
Door (D)	0.65	0.57 - 0.80
Other Agent (O)	0.85	0.80 - 0.90
Wall	1	0.90 - 1.00

Table 13: Adversarial Games: Value of each state feature

Feature	$s_i$ value	Adv 1 $s_e$ range	Adv 2 $s_e$ range
Floor	0	0.00 - 0.05	0.00 - 0.05
Learning Agent Pellet (LP)	0.13	0.05 - 0.25	0.05 - 0.32
Indep Agent Pellet (IP)	0.38	0.25 - 0.40	-
Button (B)	0.5	0.40 - 0.57	0.32 - 0.57
Other Agent (O)	0.85	0.57 - 0.90	0.57 - 0.90
Wall	1	0.90 - 1.00	0.90 - 1.00

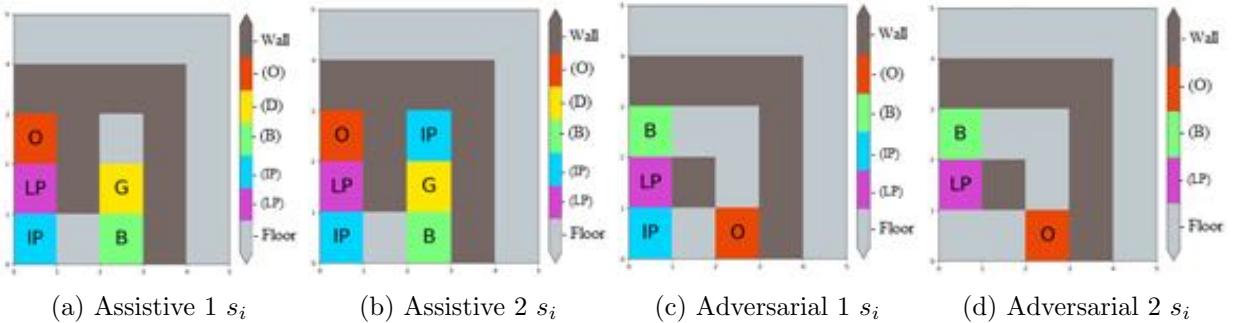


Figure 25: Original States  $s_i$  prior to transformation by the Imagination Networks.

## 8.2 Evolution of Empathetic State over time (during training)

For each of the final empathetic states  $s_e$  shown in Figure 6, the figures below show the evolution of the empathetic state during training. These are for both the E-Feature and E-Image Imagination Network generated  $s_e$ .

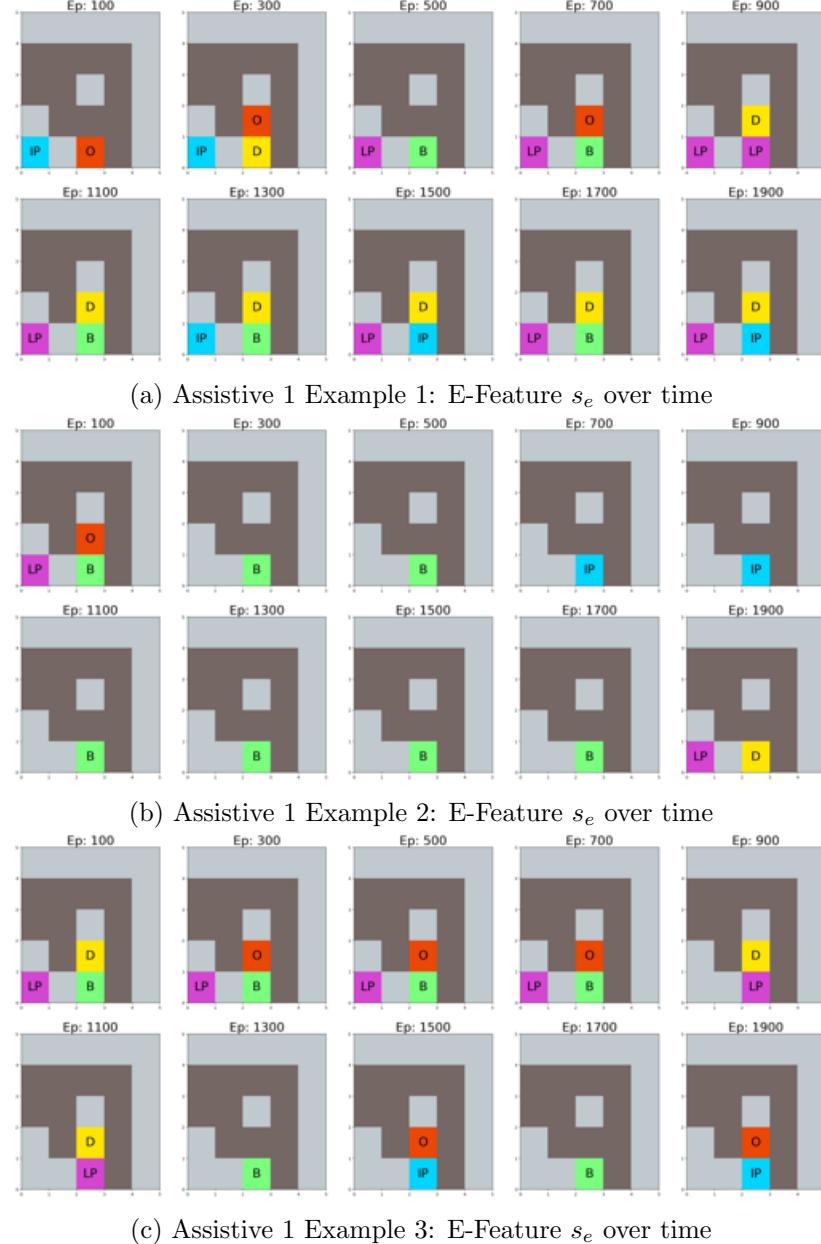
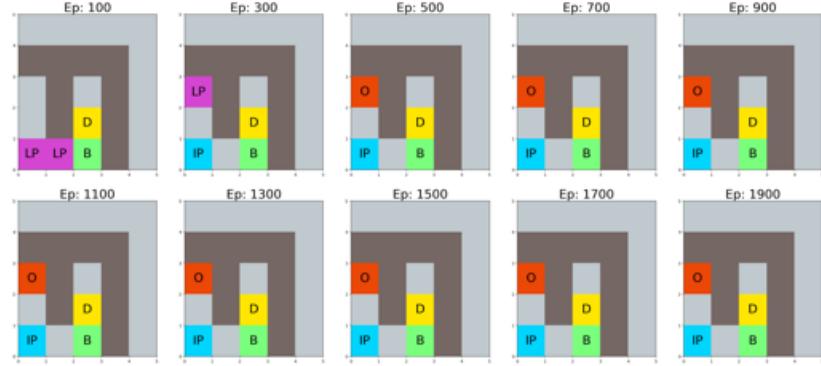
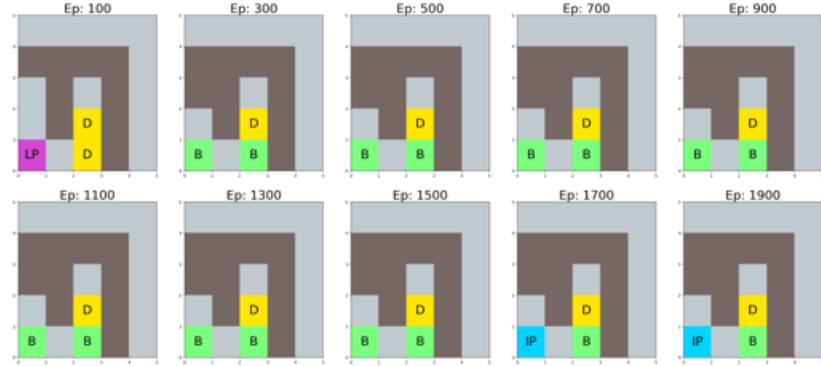


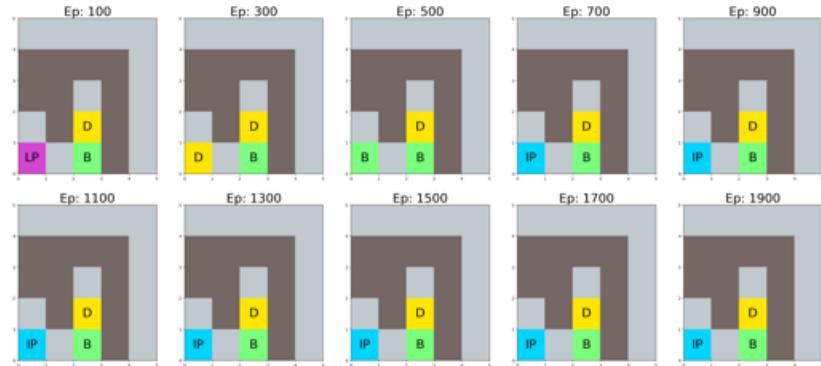
Figure 26: Assistive 1: Examples of the generated Empathetic State  $s_e$  over the training period under the E-Feature Imagination Networks for the final states shown in Figure 6.



(a) Assistive 1 Example 1: E-Image  $s_e$  over time

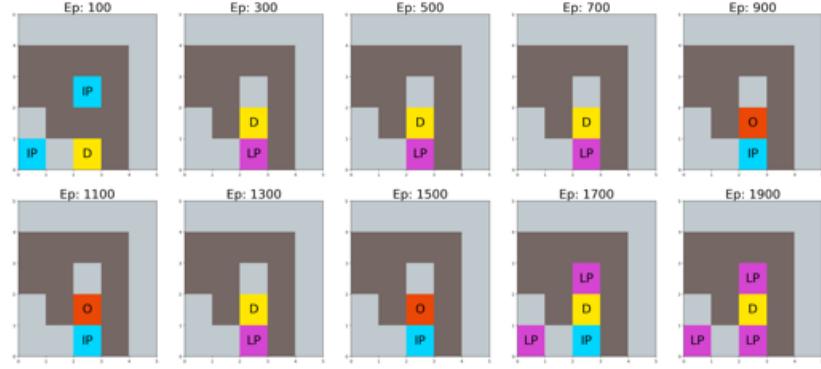


(b) Assistive 1 Example 2: E-Image  $s_e$  over time

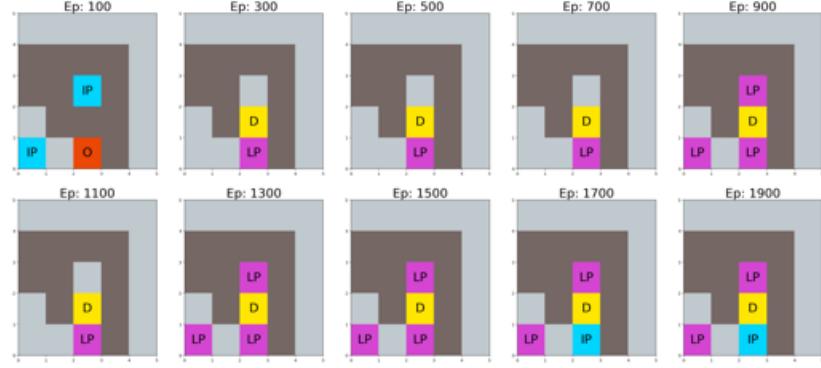


(c) Assistive 1 Example 3: E-Image  $s_e$  over time

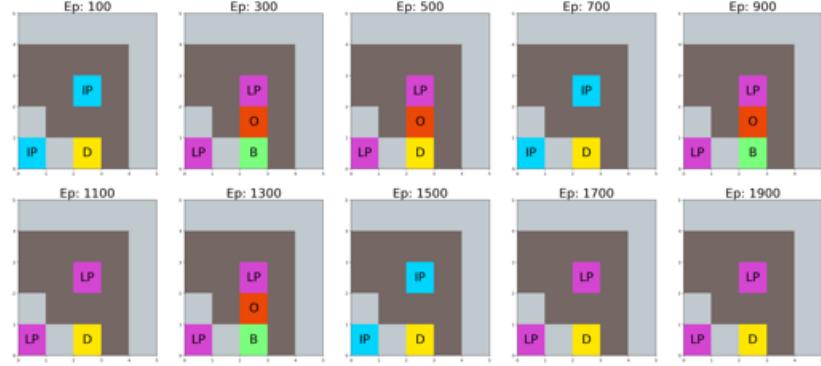
Figure 27: Assistive 1: Examples of the generated Empathetic State  $s_e$  over the training period under the E-Image Imagination Networks for the final states shown in Figure 6.



(a) Assistive 2 Example 1: E-Feature  $s_e$  over time

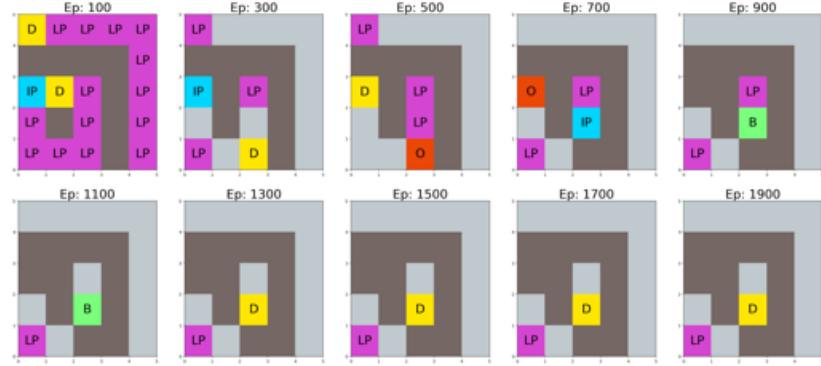


(b) Assistive 2 Example 2: E-Feature  $s_e$  over time

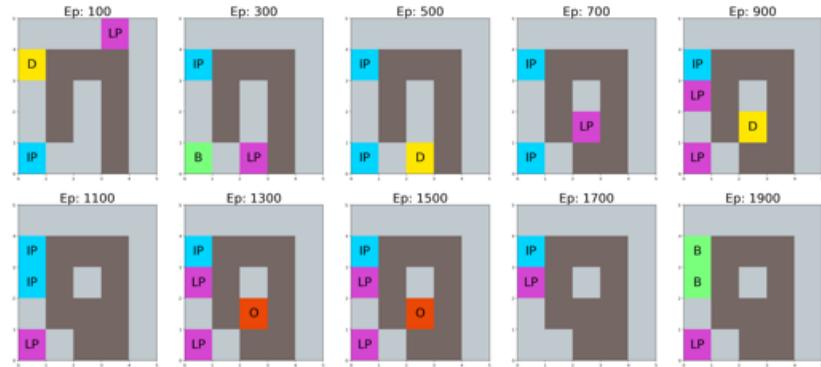


(c) Assistive 2 Example 3: E-Feature  $s_e$  over time

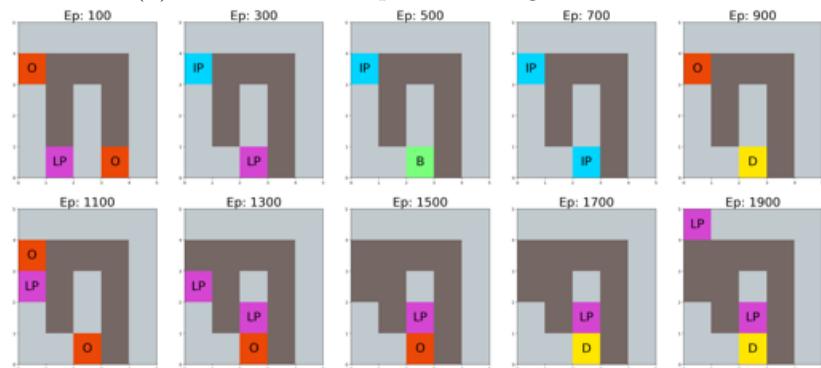
Figure 28: Assistive 2: Examples of the generated Empathetic State  $s_e$  over the training period under E-Feature Imagination Networks for the final states shown in Figure 6.



(a) Assistive 2 Example 1: E-Image  $s_e$  over time

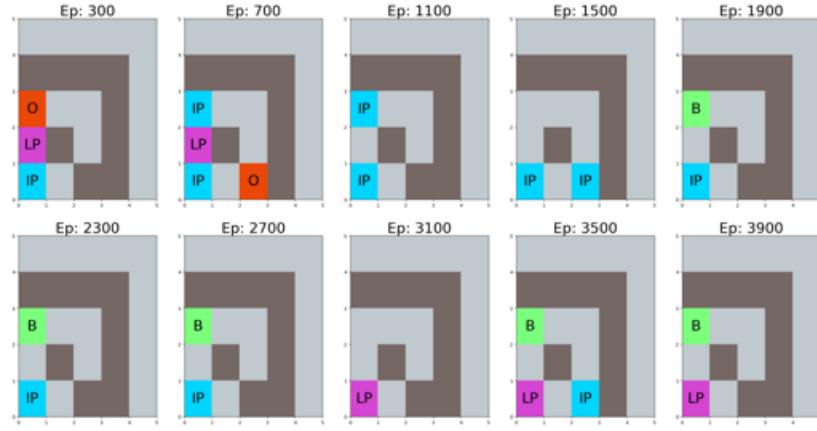


(b) Assistive 2 Example 2: E-Image  $s_e$  over time

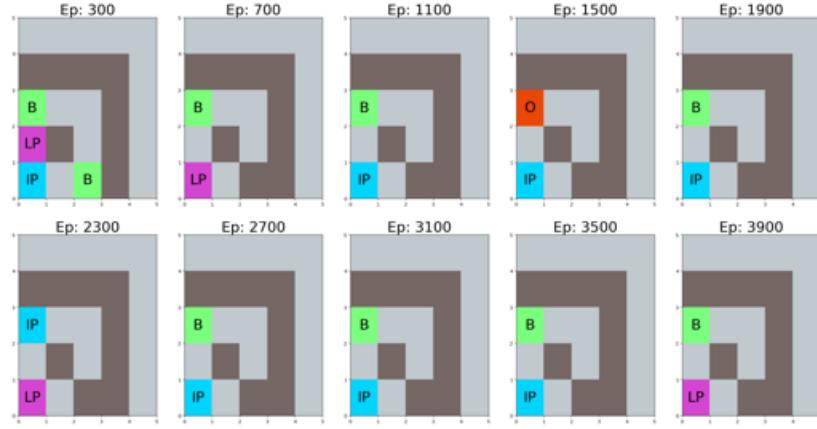


(c) Assistive 2 Example 3: E-Image  $s_e$  over time

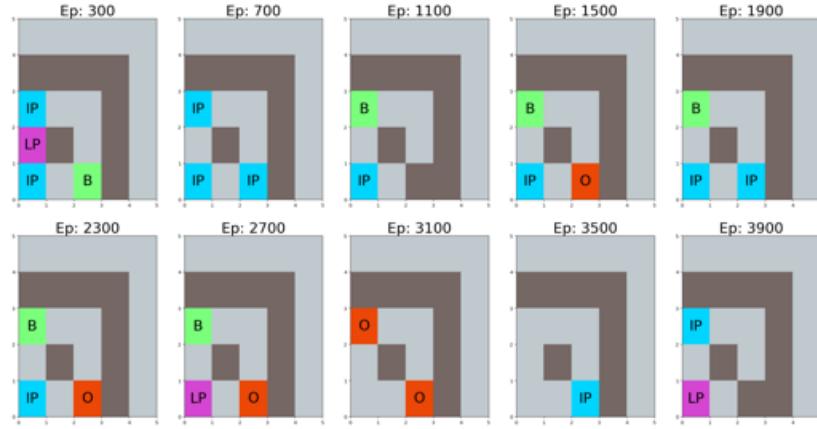
Figure 29: Assistive 2: Examples of the generated Empathetic State  $s_e$  over the training period under E-Image Imagination Networks for the final states shown in Figure 6.



(a) Adversarial 1 Example 1: E-Feature  $s_e$  over time

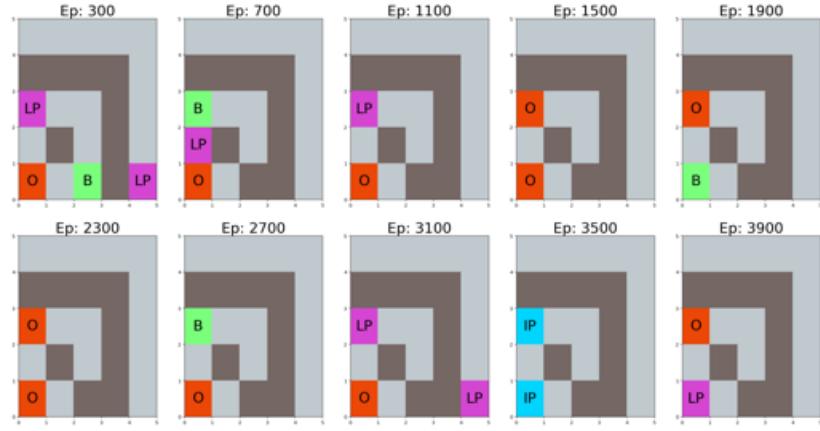


(b) Adversarial 1 Example 2: E-Feature  $s_e$  over time

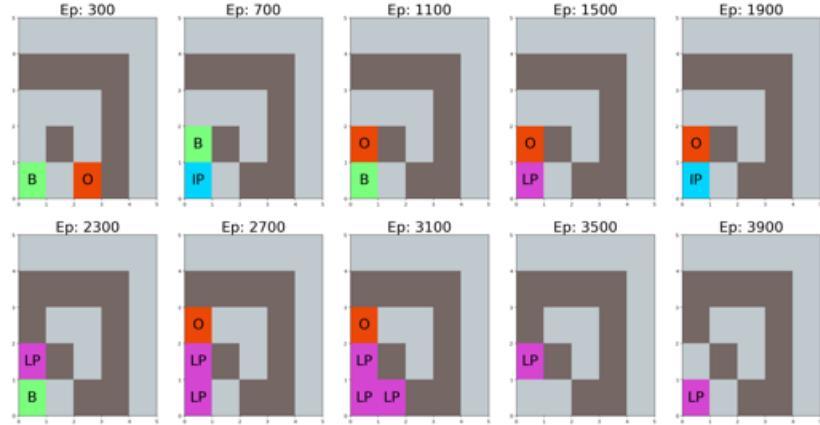


(c) Adversarial 1 Example 3: E-Image  $s_e$  over time

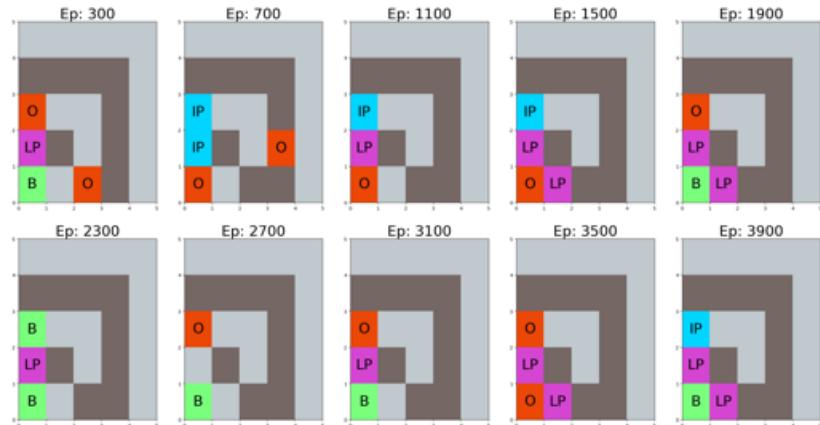
Figure 30: Adversarial 1: Examples of the generated Empathetic State  $s_e$  over the training period under E-Feature Imagination Networks for the final states shown in Figure 6.



(a) Adversarial 1 Example 1: E-Image  $s_e$  over time

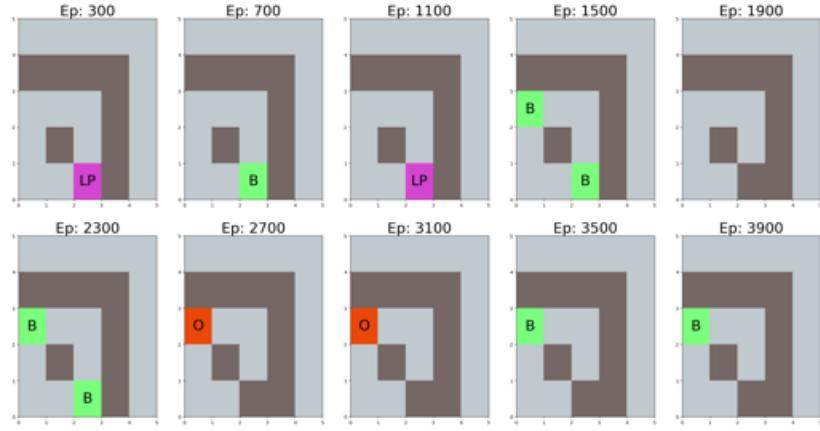


(b) Adversarial 1 Example 2: E-Image  $s_e$  over time

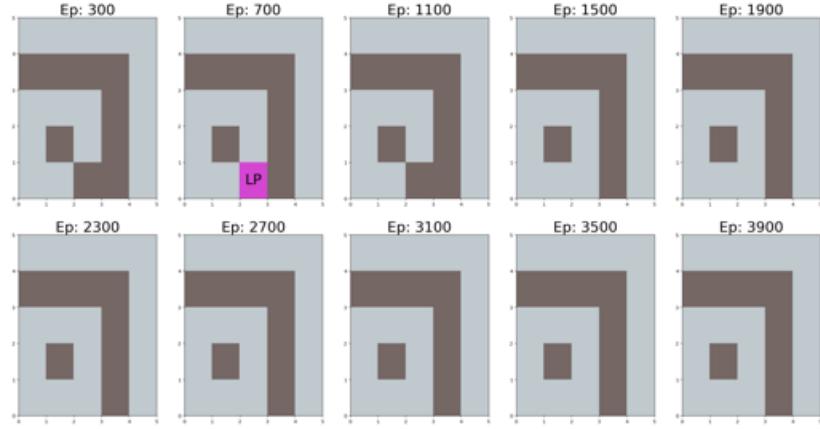


(c) Adversarial 1 Example 3: E-Image  $s_e$  over time

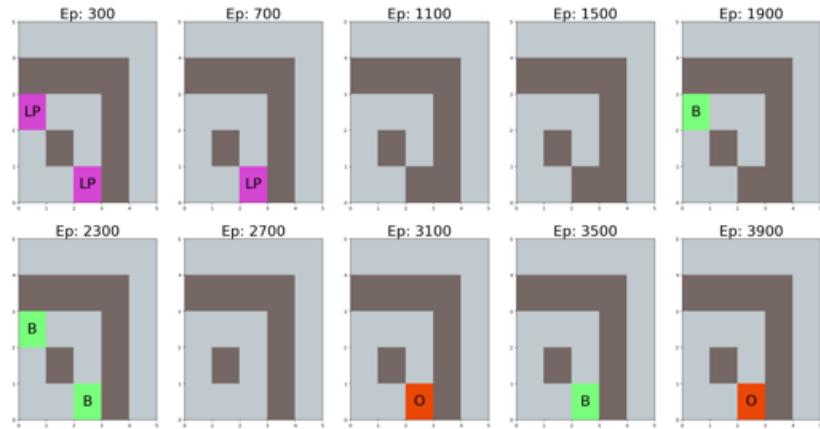
Figure 31: Adversarial 1: Examples of the generated Empathetic State  $s_e$  over the training period under E-Image Imagination Networks for the final states shown in Figure 6.



(a) Adversarial 2 Example 1: E-Feature  $s_e$  over time

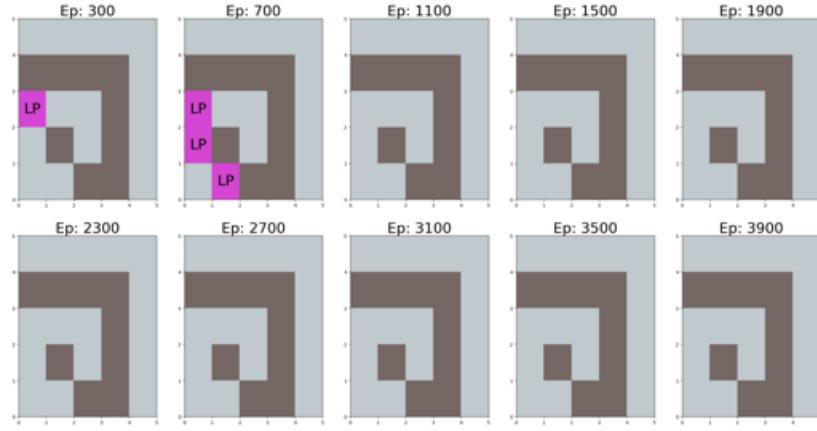


(b) Adversarial 2 Example 2: E-Feature  $s_e$  over time

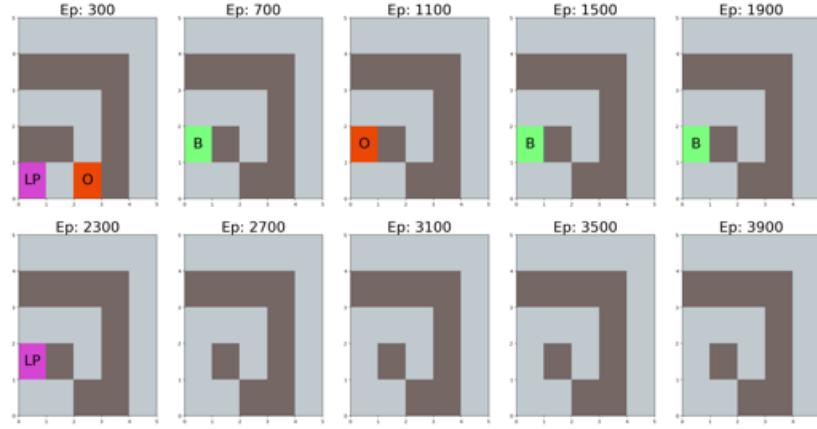


(c) Adversarial 2 Example 3: E-Image  $s_e$  over time

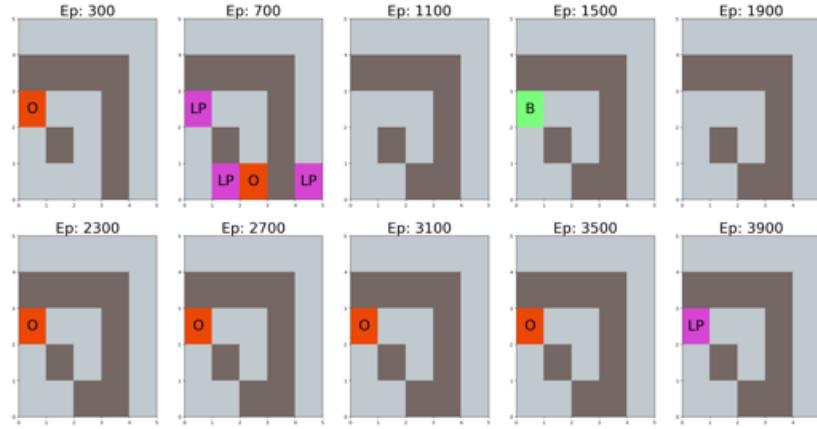
Figure 32: Adversarial 2: Examples of the generated Empathetic State  $s_e$  over the training period under E-Feature Imagination Networks for the final states shown in Figure 6.



(a) Adversarial 2 Example 1: E-Image  $s_e$  over time



(b) Adversarial 2 Example 2: E-Image  $s_e$  over time

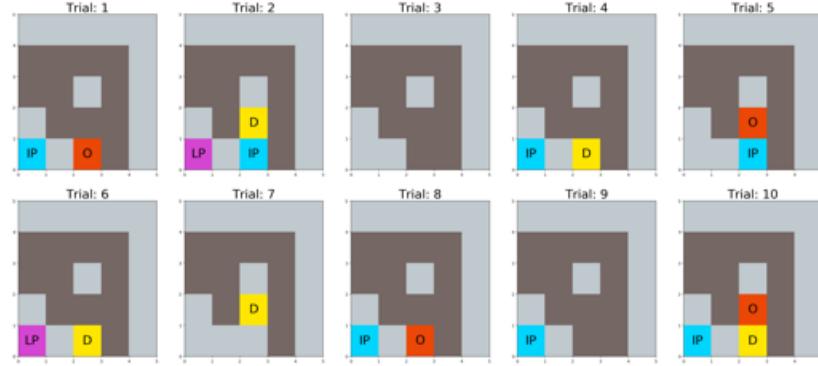


(c) Adversarial 2 Example 3: E-Image  $s_e$  over time

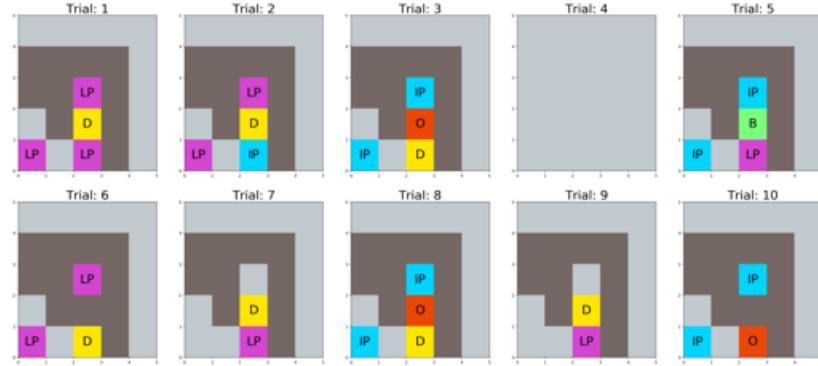
Figure 33: Adversarial 2: Examples of the generated Empathetic State  $s_e$  over the training period under E-Image Imagination Networks for the final states shown in Figure 6.

### 8.3 Final Empathetic States $s_e$ from random trials

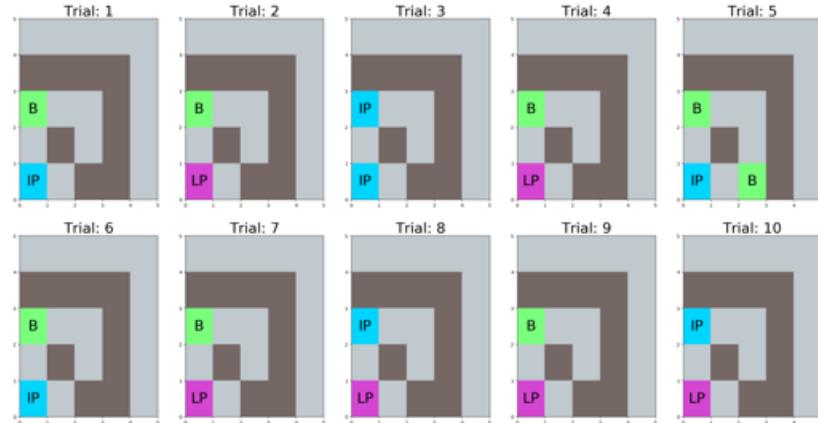
Figures 34 and 35 show the final empathetic state for all the randomly initialised trials of each game (10 of each). Results from both E-Feature and E-Image Imagination Networks are shown.



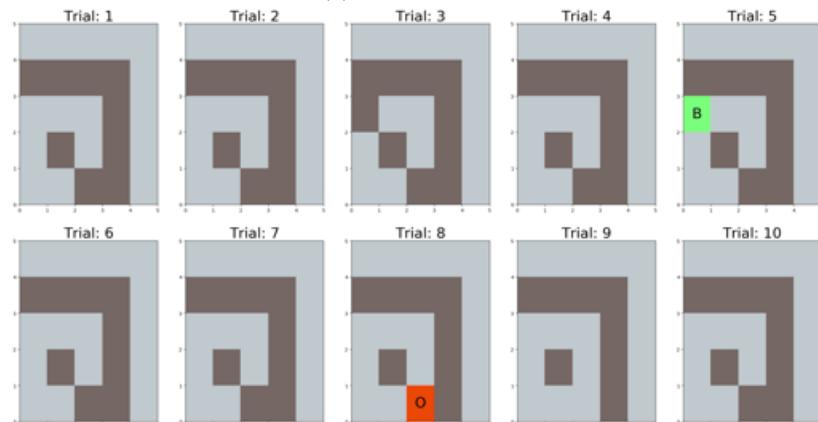
(a) Assistive 1



(b) Assistive 2

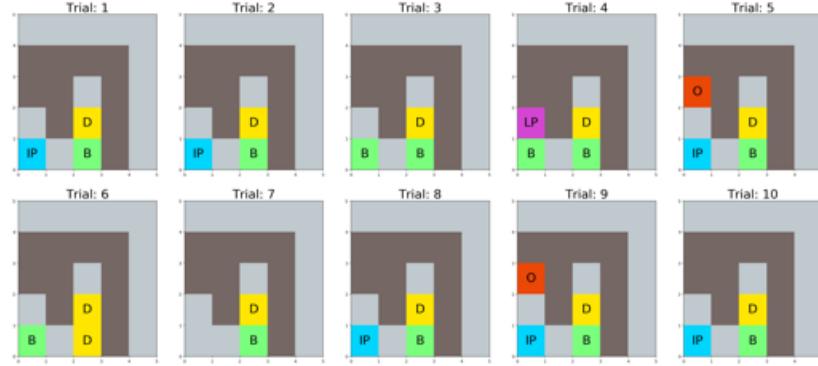


(c) Adversarial 1

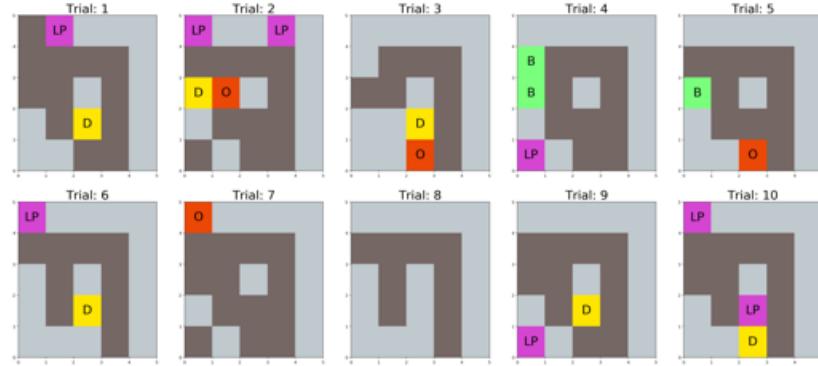


(d) Adversarial 2

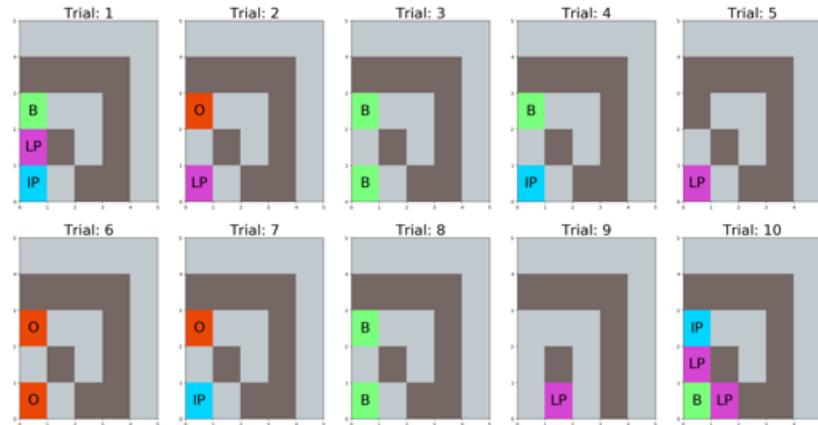
Figure 34: Final Empathetic state  $s_e$  produced via E-Feature Imagination Network for 10 random trials of each game.



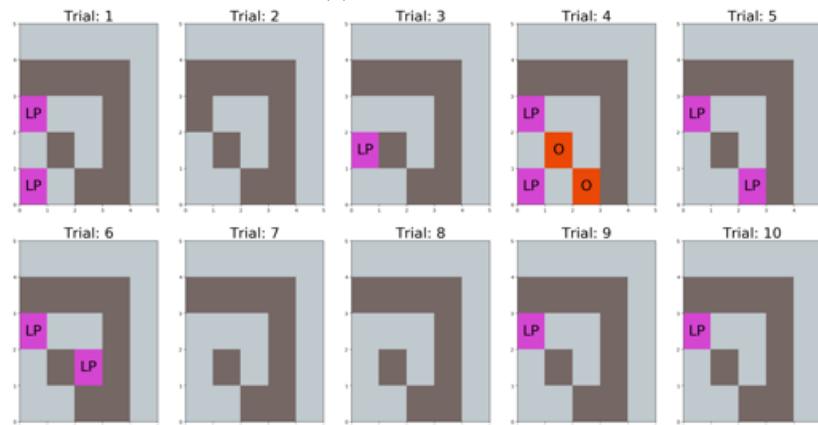
(a) Assistive 1



(b) Assistive 2



(c) Adversarial 1

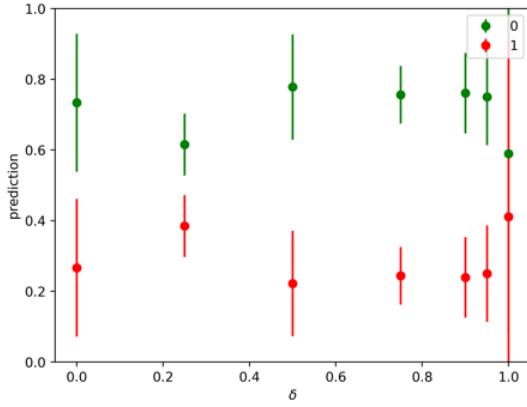


(d) Adversarial 2

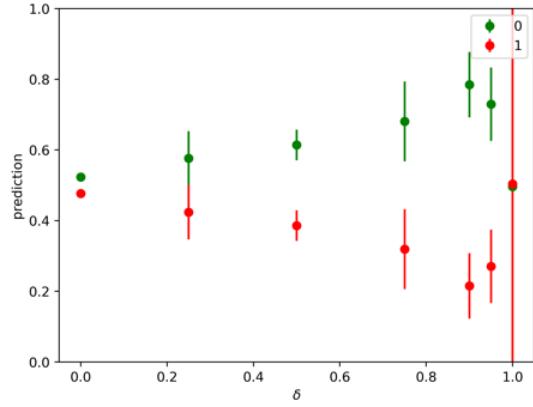
Figure 35: Final Empathetic state  $s_e$  produced via E-Image Imagination Network for 10 random trials of each game.

## 8.4 Adversarial Modelled Button

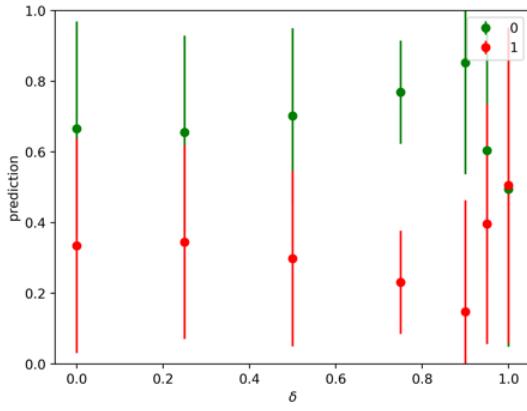
For the Adversarial games of the Sympathy Framework the button status was also modelled as part of the Imagination Model. As inputs this model took in the status of the button (either a 0 or a 1) and output a predicted empathetic button status. Figure 36 shows the predictions of this model after training for both the E-Feature and E-Image Imagination models for various  $\delta$  values.



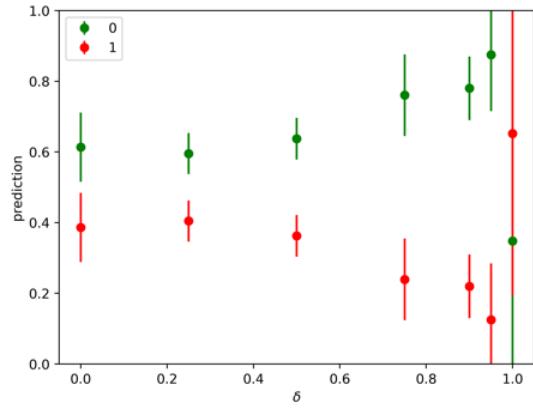
(a) Adversarial 1: E-Feature



(b) Adversarial 1: E-Image



(c) Adversarial 2: E-Feature



(d) Adversarial 2: E-Image

Figure 36: Predicted value of button status model at the end of training for various  $\delta$  values. Adversarial 1 (a-b) and Adversarial 2 (c - d).