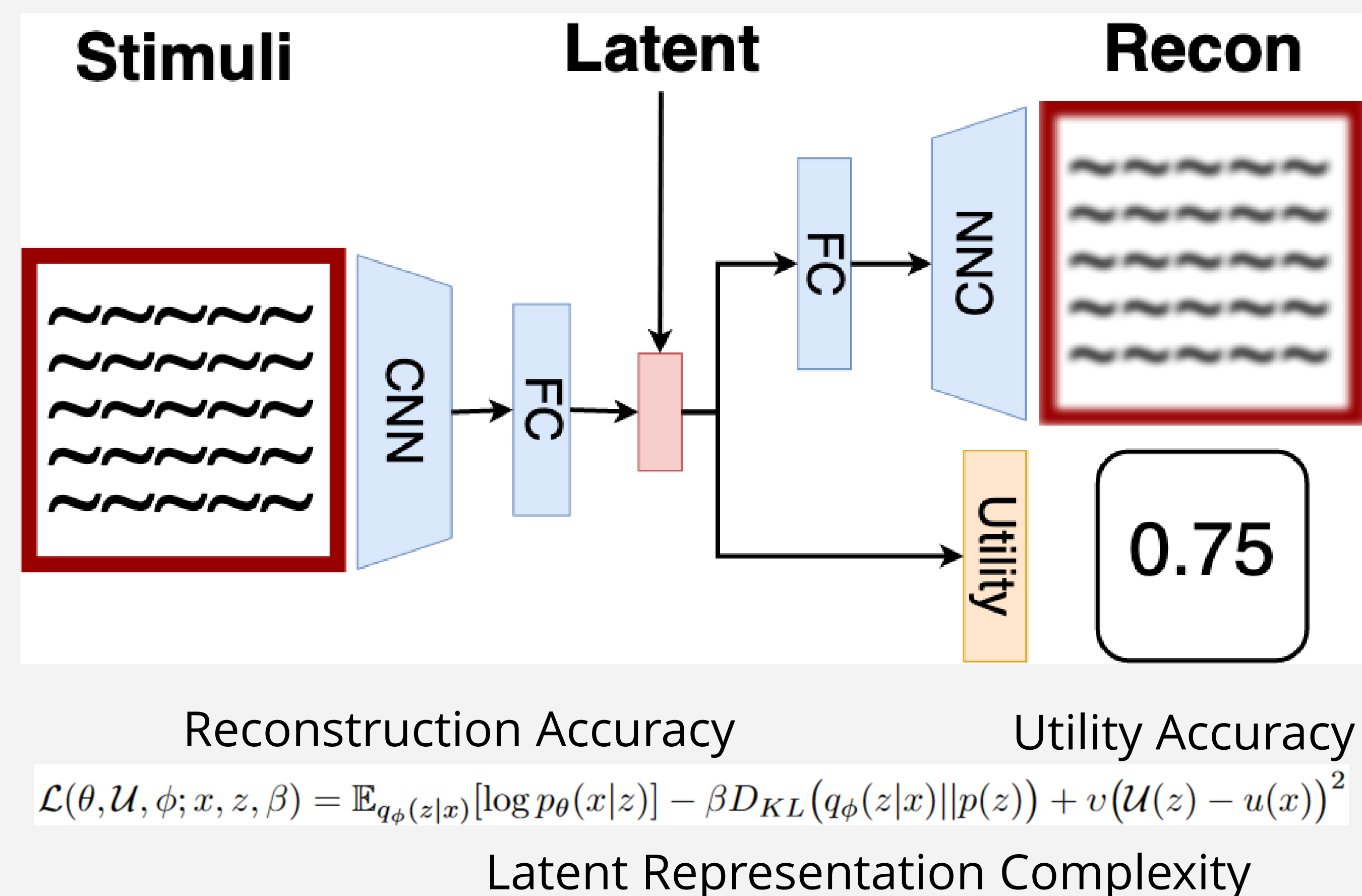


Representations of visual information used in learning can be modelled with a disentanglement objective. This model describes how changing utility of stimuli results in changes of visual representation formation.

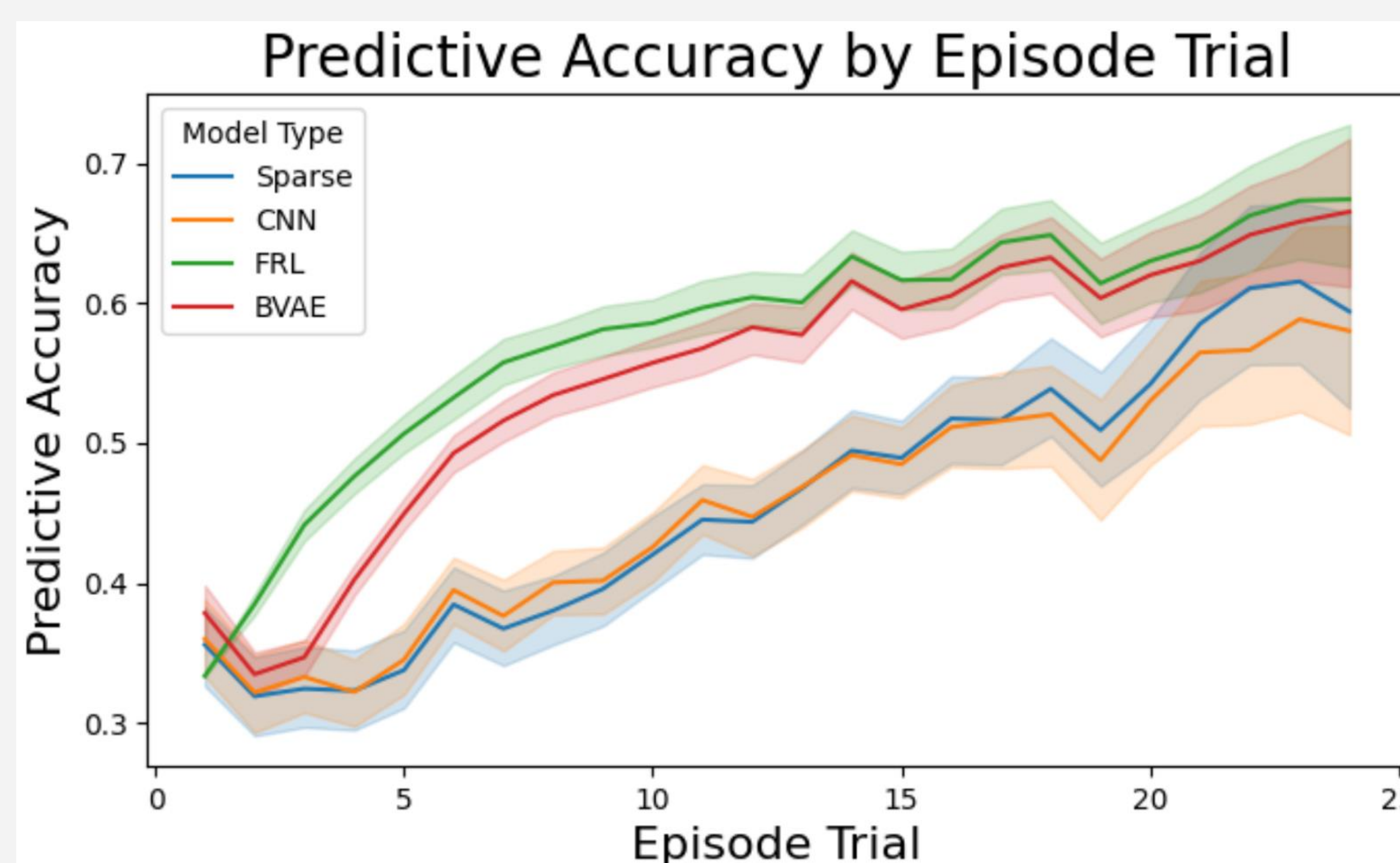
Abstract

The β Variational Autoencoder (β -VAE) learns to form representations of visual information by maximizing the disentanglement objective. These representations deviate along a single dimension as the stimuli deviates. Deep reinforcement learning (DRL) can be difficult to apply onto predicting human behaviour. This model addresses this by using representations of stimuli in learning tasks produced by a combined DRL and β -VAE model. This requires a novel model structure and training method:



Results

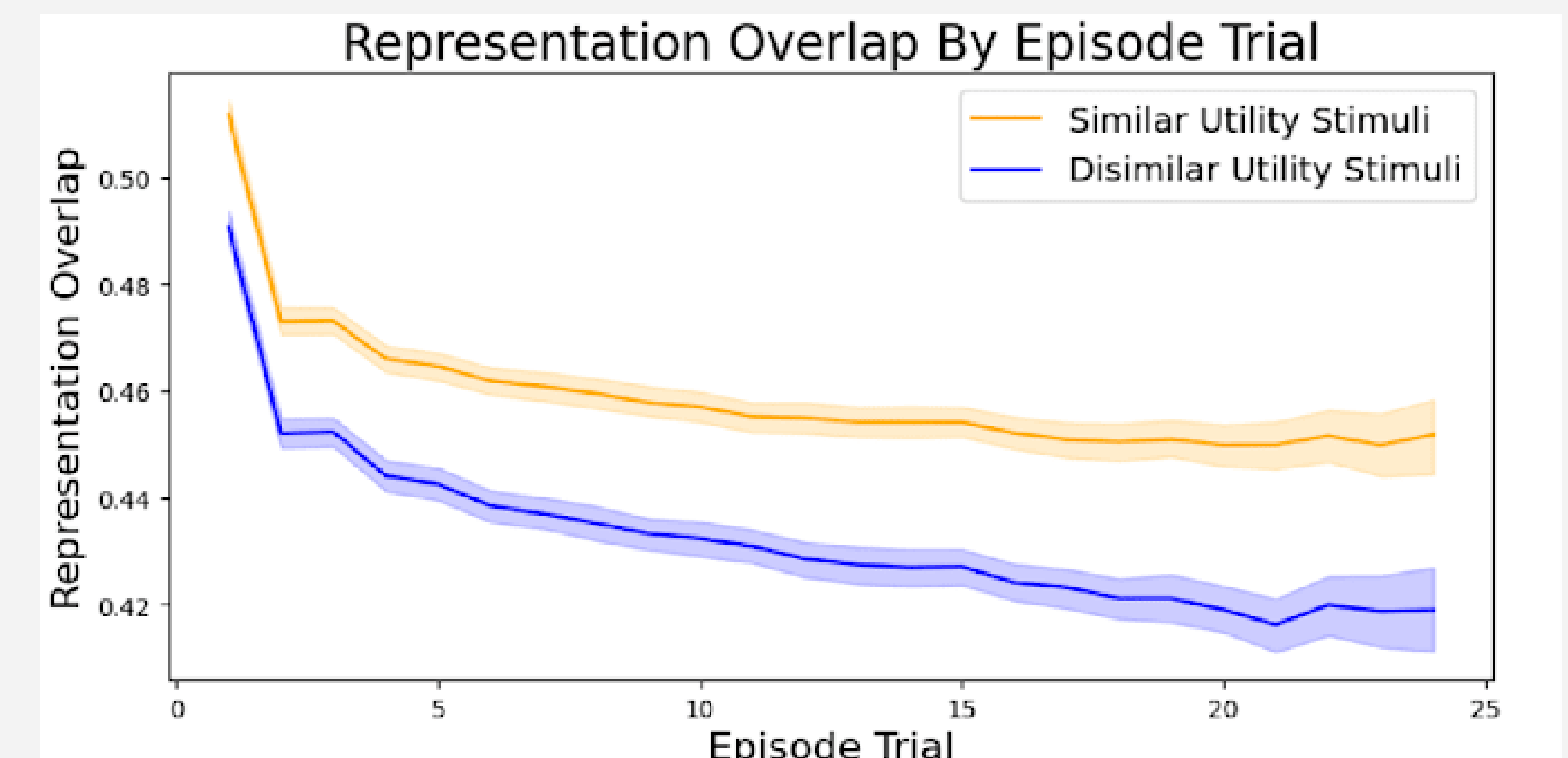
We apply the described model onto predicting human behaviour in a contextual bandit learning task. The model has a higher accuracy than alternative DRL models using traditional Convolutional Neural Networks or a sparse CNN model that mimics the information bottleneck of β -VAEs.



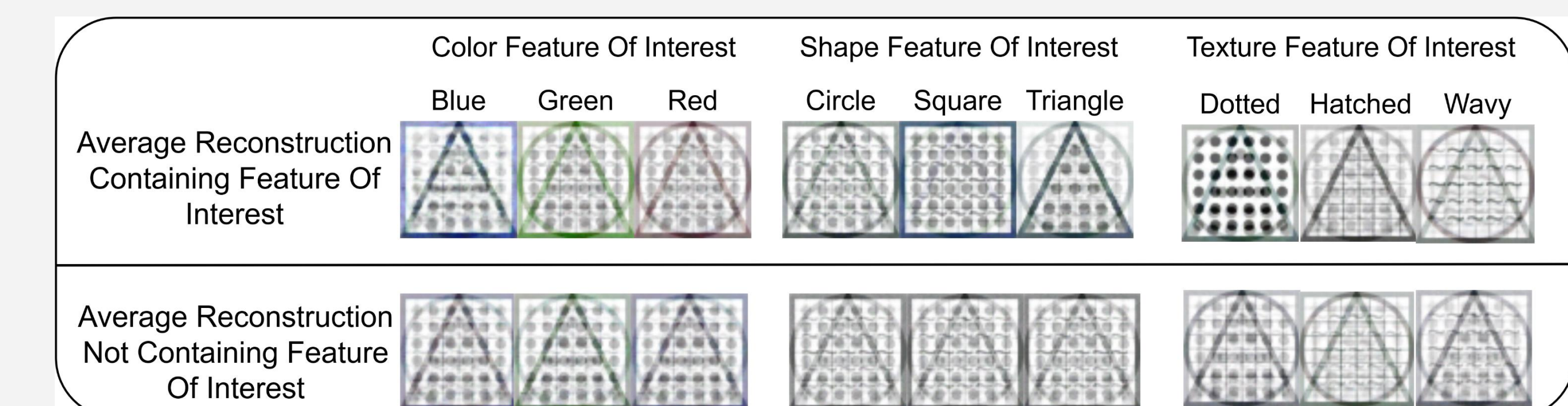
The tabular Feature Reinforcement Learning method provides utility predictions to the three deep neural network methods. This FRL model takes as input hand-crafted features instead of visual information.

Analysis

Measuring how the β -VAE model representations overlap each other gives insight into how changing utility impacts learned visual representations.



Reconstructing noisy latent samples demonstrates that high-utility representations are more resilient to noise and take up more of the latent representation space.



Conclusion

The disentanglement objective learned by β -VAE models describes how learning utility impacts visual representation formation, while also predicting behaviour. This suggests a source of acquired equivalence of visual information based on representation overlap. Future research comparing how participants detect change can further apply this model.



Modelling Human Reinforcement Learning with Disentangled Representations

Tyler Malloy, Tim Klinger, Chris R. Sims