



# Disentangled behavioral representations

Amir Dezfouli<sup>1\*</sup>, Hassan Ashtiani<sup>2</sup>, Omar Ghattas<sup>1,3</sup>, Richard Nock<sup>1,4,5</sup>, Peter Dayan<sup>6</sup>, Cheng Soon Ong<sup>1,4</sup>



Data61, CSIRO, Australia, 2McMaster University, 3University of Chicago, 4Australian National University, 5University of Sydney, 6Max Planck Institute for Biological Cybernetics \*akdezfuli@gmail.com

#### Introduction

- Individual characteristics in human decision-making are often quantified by fitting a parametric cognitive model to subjects' behavior and then studying differences between them in the associated parameter space. However, these models often fit behavior more poorl than recurrent neural networks (RNNs), which are more flexible and make fewer assumptions about the underlying decision-making processes.
- The parameter and latent activity spaces of RNNs are generally high-dimensional and uninterpretable, making it hard to use them to study individual differences.
- We show how to benefit from the flexibility of RNNs while representing individual differences in a low-dimensional and interpretable space. To achieve this, we propose a novel end-to-end learning framework in which an encoder is trained to map the behavior of subjects into a low-dimensional latent space. These low-dimensional representations are used to generate the parameters of individual RNNs corresponding to the decision-making process of each subject.
- We introduce terms into the loss function that ensure that the latent dimensions are informative and disentangled, i.e. encouraged to have distinct effects on behavior. This allows them to align with separate facets of individual differences.
- We illustrate the performance of our framework on synthetic data as well as a dataset including the behavior of patients with psychiatric disorders.

### **Training losses**

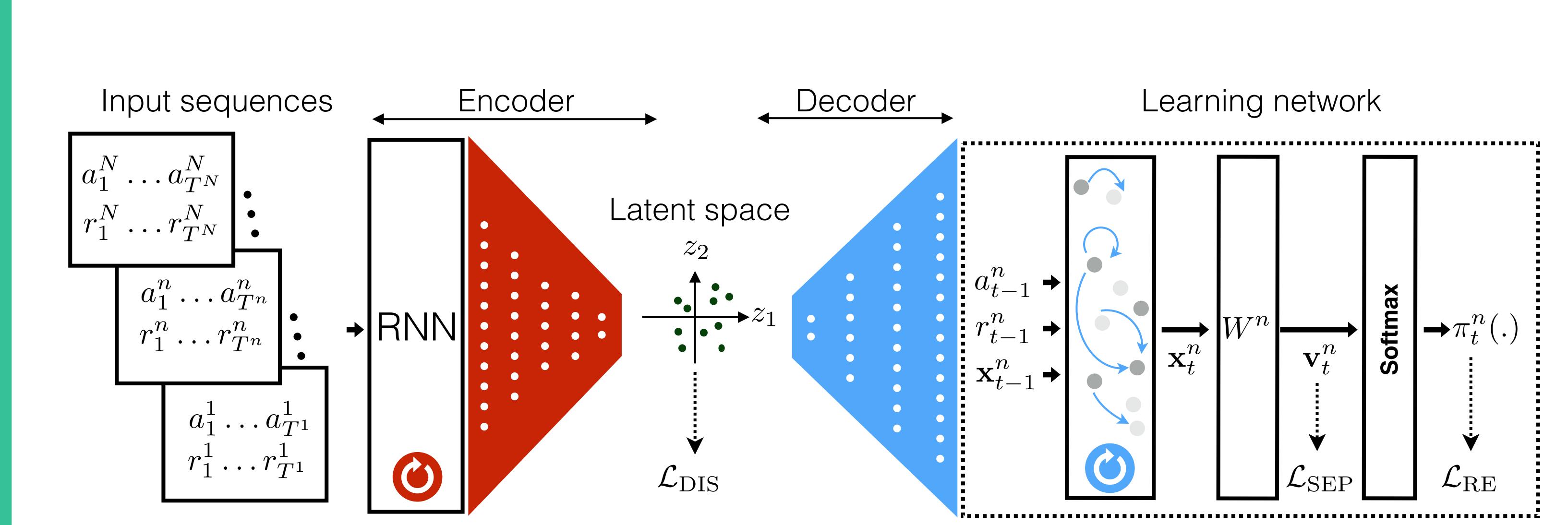
- The training loss function has three components:
- (i) a reconstruction loss which penalizes discrepancies between the predicted and actual input sequence,
- (ii) a group-level disentanglement loss which encourages sequences to spread independently across the dimensions of the input sequence,
- (iii) a separation loss which favors dimensions of the latent space that have disentangled effects on the behavior generated by the learning networks.

$$\mathcal{L}_{RE} \equiv -\frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T^n} \log \pi_t^n(a_t^n; a_{1...t-1}^n, r_{1...t-1}^n)$$

$$\mathcal{L}_{\text{DIS}} = \lambda_1 \text{MMD}(\hat{q}(\mathbf{z}), p(\mathbf{z})) + \text{KL}(g(\mathbf{z}) || p(\mathbf{z}))$$

$$\mathcal{L}_{\text{SEP}} = \frac{1}{N} \sum_{n=1}^{N} \left| \frac{\partial^2}{\partial z_1 \partial z_2} \sum_{t=1}^{T^n} u_t^n \right|$$

$$\mathcal{L} = \mathcal{L}_{RE} + \lambda_2 \mathcal{L}_{DIS} + \lambda_3 \mathcal{L}_{SEP}$$



 $\mathbf{z}_{M\times 1}^n \equiv \operatorname{enc}(a_{1\ldots T}^n, r_{1\ldots T}^n; \Theta_{\operatorname{enc}}), \quad n=1\ldots N,$ 

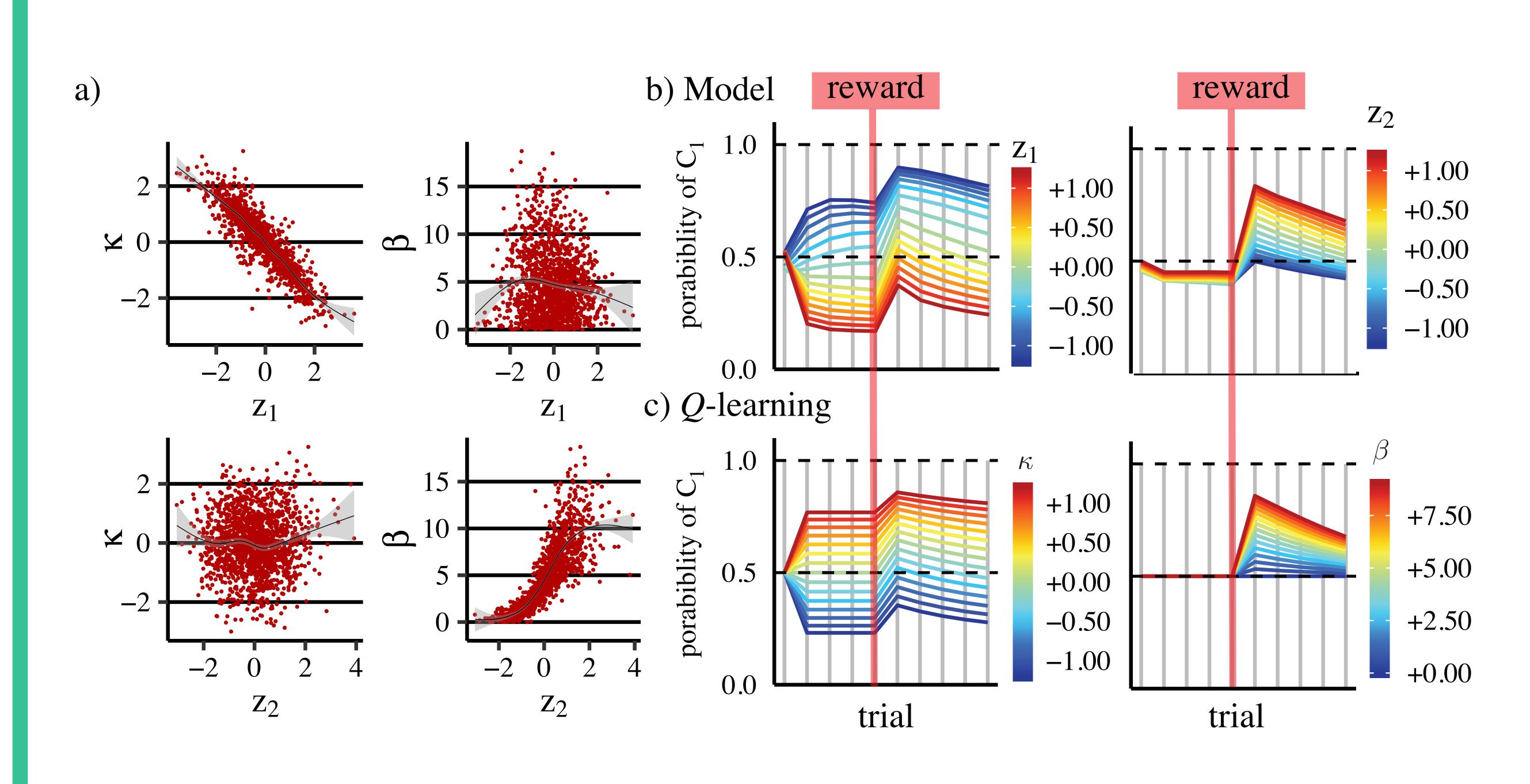
$$\Phi^n \equiv \operatorname{dec}(\mathbf{z}^n; \Theta_{\operatorname{dec}}).$$

$$\mathbf{x}_{t}^{n} \equiv f^{n}(a_{t-1}^{n}, r_{t-1}^{n}, \mathbf{x}_{t-1}^{n}; \Phi^{n})$$

$$\mathbf{v}_t^n \equiv W^n \mathbf{x}_t^n, \quad \pi_t^n(a_t \equiv C_i; a_{1...t-1}^n, r_{1...t-1}^n) = \frac{e^{\mathbf{v}_t^n[i]}}{\sum_{k=1...K} e^{\mathbf{v}_t^n[k]}},$$

### Results

Model



# Conclusion

# Summary:

We proposed a flexible autoencoder-based framework for modelling individual differences in decision-making tasks. The autoencoder maps the sequence of actions and rewards taken and received by a subject to a position in a low-dimensional latent space which is decoded to determine the weights of a 'learning network' RNN that characterizes how that subject behaves as a plastic decision-maker. The latent space was disentangled by adding a specific component to the loss function.

# References:

Amir Dezfouli, Kristi Griffiths, Fabio Ramos, Peter Dayan, Bernard W Balleine. Models that learn how humans learn: the case of decision-making and its disorders.

## **Acknowledgements:**

We are grateful to Bernard W. Balleine for sharing the dataset with us.

