# Closing the teacher-learner loop: the role of affective signals in interactive reinforcement learning

Bernhard Hilpert, Prof. Joost Broekens, Prof. Kim Baraka, Prof. Aske Plaat
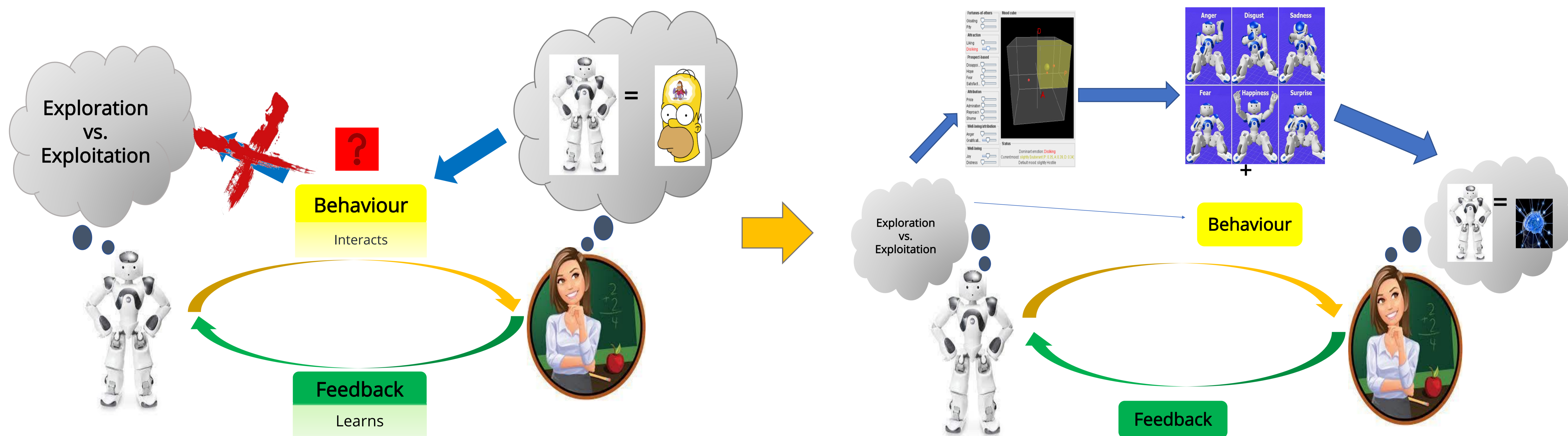
## Problem

Robot behavior can be preprogrammed (scripted) and learned, which is important for personalization of behavior towards humans and for adapting behavior to particular contexts. In Reinforcement Learning, learning is based on exploration and the optimization of feedback. During this process the agent makes action selection decisions, many of which are not obvious to the observing human. For example, is a suboptimal action of a collaborative robot an exploratory move, or, does the robot really think that move is optimal?
Transparency of the learning process is needed [1].

## Approach

We aim to investigate if and how robot affective expression, grounded in RL-based simulation (joy/distress/hope/worry) , can help to make the learning process more transparent.
In humans affective signals reflect the internal state of the person, in particular the appraisal of the situation. This natural way to express one's interpretation of a situation, when simulated by an agent in a human-agent interaction context, could positively influence human teaching performance, but structured user experiments to test this hypothesis are still missing.



## First experiments

One of the goals of this project is to assess the applicability of biologically-inspired affective signals to agent communication during learning, and potentially investigate novel agent-specific affective signals that are perhaps better at communicating the agent's learning process. To tackle this problem generally, we first need to define the type of tasks in which such RL-based emotions are useful generically and subsequently the type of tasks suitable to evaluate the approach. Further experimental questions include the timing and occurrence of the affective signals, the techniques suitable for expressing the signals. For the simulation of the affective signals we will start with model-based RL, as this enables the simulation of signals reflecting the current situation  but also the anticipated near future .

The project brings expertise from different research areas including human-agent interaction, affective computing, emotion psychology, and interactive machine learning, towards richer ways of tackling nontrivial HI problems.

## Project relevance

This is a prime example of a collaboration between an AI agent and a human where we aim to create synergy between a human teacher and an agent learner through real-time information sharing. Specifically, in our project the agent has to learn and the human has to help the agent by showing efficient teacher behavior.

In terms of use cases, this project is relevant to personalisation of behavior of e.g. an educational robot or agent.

## References

1. Broekens, J. and M. Chetouani, Towards Transparent Robot Learning through TDRL-based Emotional Expressions. IEEE Transactions on Affective Computing, 2021. 12(2): p. 352-362.
2. Broekens, J. and L. Dai, A TDRL Model for the Emotion of Regret, in 2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII). 2019. IEEE.
3. Broekens, J., A Temporal Difference Reinforcement Learning Theory of Emotion. arXiv preprint arXiv1807.08941, 2018.
4. Moerland, T., J. Broekens, and C.M. Jonker, Fear and Hope Emerge from Anticipation in Model-Based Reinforcement Learning. In Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI-16), S. Kambhampati, Editor. 2016, AAAI Press. p. 848-854.
5. Broekens, J., E. Jacobs, and C.M. Jonker, A reinforcement learning model of joy, distress, hope and fear. Connection Science, 2015: p. 1-19.