

The mathematics of systems with goals (agents?)

Manuel Baltieri

Araya Inc., Tokyo, Japan Department of Informatics, University of Sussex, Brighton, UK

manuel_baltieri@araya.org

Areas of interest

Life sciences: agents include living organisms and possibly related systems, defining agents can help studies of the *origins of life* and *life as it* <u>could be</u>.



Agents interact with their environment





Psychology and **neuroscience**: decision making and behavioural studies rely on knowing what agents do and what drives their *actions*, intentions and plans.

Al and machine learning: explainable AI and AI safety tackle complex AI/ML models to understand how they work and *interact* with users, so to generate interpretable output.

Robotics and engineering: self-driving cars, drones and robot assistants require a certain level of *autonomy* to handle complex real-world applications.



Philosophy: studies of cognition and mind for autonomous, and perhaps conscious, agents.

Other areas: economics, policy making, sociology.



Agent

Dynamical systems theory: attractors give us a way to study invariant properties of dynamical systems and their long-term behaviour, (generalised) synchronisation gives us ways to describe particular instances of the behaviour of coupled dynamical systems.



Control theory and **reinforcement learning**: regulation condition(s) define possible goals for a system, open vs. closed loop control show us the differences between systems that are paired to their environment and systems that ignore observations from the environment.



Artificial life: adaptive fit defines a subspace of the cartesian product of agent and environment as the space of vaiability for an agent-environment system, anything outside of it cannot be an agent.



What's an "agent"? Where is *the* boundary?









