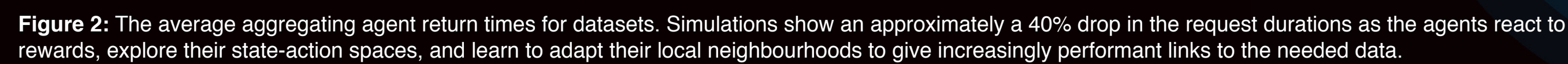
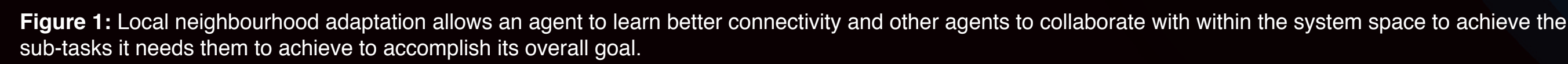
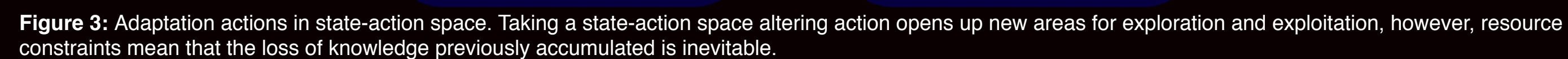


*NIALL CREECH, Kings College London*

The Local Neighbourhood Learning Algorithm (Local-NL) tackles the need for agents to be able to adapt their connectivity and knowledge of other agents within the system to find sub-spaces that allow completion of a composite task. Using an extension of Q-learning, agents learn utilities of the range of actions available to them. They then have the choice of taking an action that will exploit their current local neighbourhood of other agents, bringing them closer to task completion, explore the neighbourhood to optimise their requests of other agents, or reshape the neighbourhood altogether. Altering the neighbourhood topology will bring in connections to new agents that may open up more optimal actions. However, due to resource constraints this comes at the cost of the loss of knowledge of some of the previous sub-space.

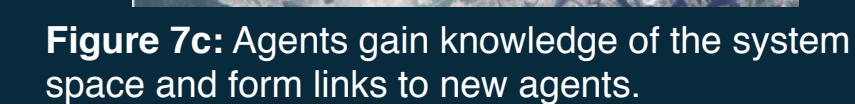
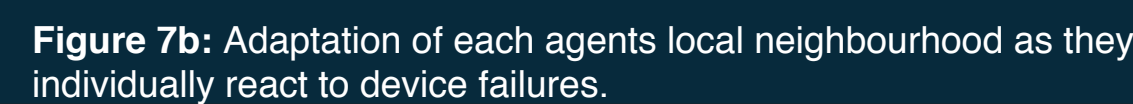
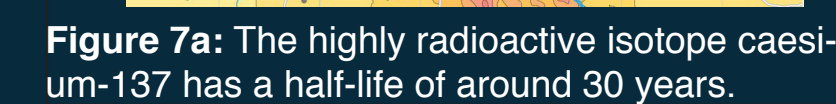
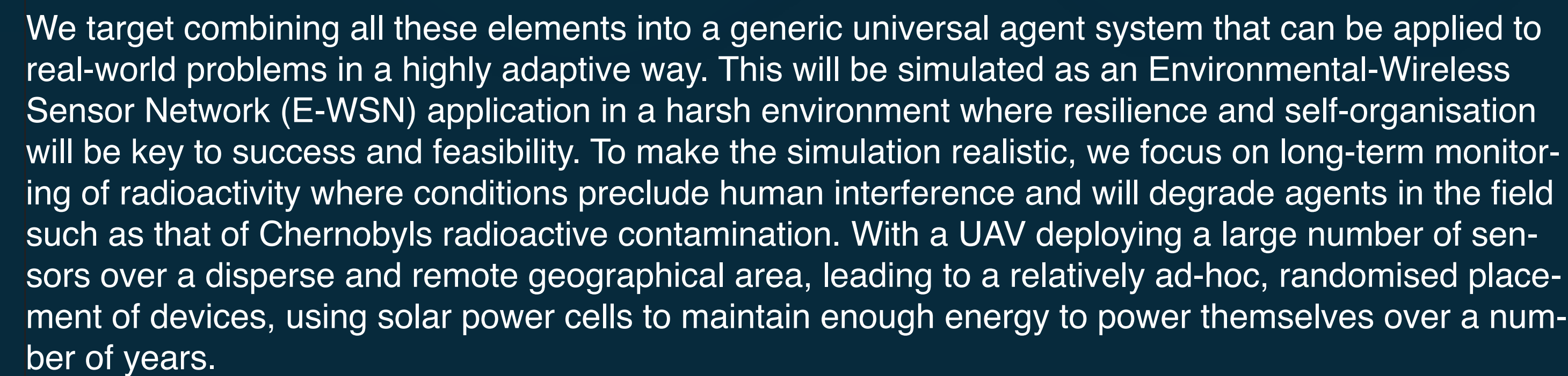


The Reward Trends for Risk-Action Probabilities (RT-RAP) algorithm we introduce combines a relative performance metric for an agent and a transformation function for its role to generate behaviours that dynamically alter its exploration and optimisation strategy. This allows an agent to use a comparison of its current learning policies performance against its historical reward trends to optimise and exploit subspaces of the systems state-action space without losing their flexibility to adapt in the face of variations and system disruption. This functionality is crucial to the agents ability to find an optimisation solution within the system, without it the agent will make too many changes to its connectivity and knowledge of other agents to realistically find exploitable knowledge within the system.

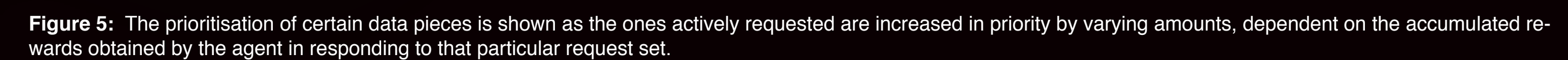
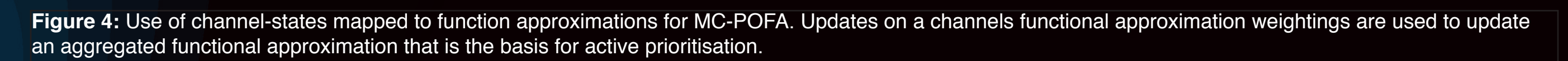


In systems with a large number of agents there are fundamental pressures on the centralised coordination techniques used to provide inter-communication, task orchestration, and routing of messages. As the scale of interacting components expands, we reach resource constraint plateaus, where computation, storage, or communication pathways become saturated. At these points we must decompose each agents functionality into a number of specialisms that can then be taken up by other agents, at the cost of even more orchestration communications and synchronisation to provide this distributed functionality. To provide solutions to tackle these issues we focus on distributed agent systems where reinforcement learning behaviours are constrained by resource usage limits and hence by local neighbourhood awareness rather than global system knowledge.

In this work we develop algorithms for autonomous intelligent agents within a distributed multi-agent system that enable agents to learn and repeatedly adapt a subset of state-action space while also exploiting it to achieve a goal. Through the investigation of these systems we also provide some insight and define concepts that illustrate the behaviours of agents under these conditions that prove useful to build further contributions, including the use of constrained local neighbourhoods as units of scalability in large distributed systems.



Approaching the communication challenge from the opposite perspective, here agents learn to prioritise their responses to requests to optimise sub-task completion. The Multi-Channel Priority Optimisation through Function Approximation Algorithm (MC-POFA) allows agents to flexibly prioritise actions to satisfy responses in a way that preserves knowledge over incoming channels and shrink or grow the function approximations capacity to handle a broad spectrum of request demand. The adaptive capacity functionality is a cornerstone of the agents learning scalability, ensuring it can focus on a manageable subset of prioritisations when placed in a highly communicative environment.



We tackle the problem of driving behavioural change of agents dependent on how well its performing through the Relative Environmental Signals Q-Transformation Algorithm (RES-QT). At a high level this algorithm generates an agent-specific internal metric based on a combination of its rewards and the entropy of its learned knowledge over its current state-action space. The internal reward signal is compared against similar metrics collected from the agents local neighbourhood through informational communications. This then allows the transformation of an agents state-action space Q-values based on the agents view of its relative performance, driving risk-taking or conservative behaviours in response its belief in its success.

