



# We generate a formal representation of all possible states and state transitions of AI agents interacting on a grid

Tom Burns and Robert Tang  
OIST Graduate University, Japan

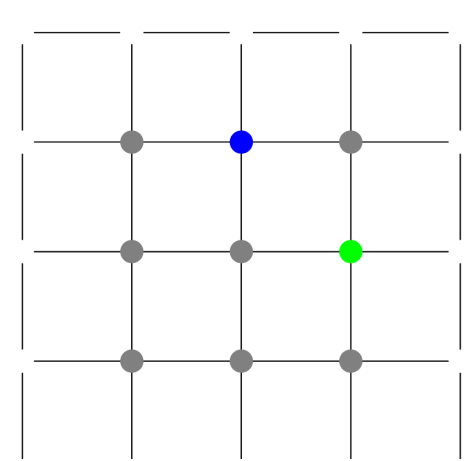


tfburns.com | thomas.burns@oist.jp | +81 (0)90 6863 7039

## ABSTRACT

Gridworlds are a popular and powerful test-bed for artificial intelligence (AI) algorithms, especially in reinforcement learning and associated AI safety problems. To describe how AI agents behave in such gridworlds, we consider gridworlds as reconfigurable systems and construct their state complexes. These state complexes reveal the underlying structures and patterns in the system's possible reconfigurations. This work incorporates the concepts of gridworlds, reconfigurable systems, and state complexes to show structures and patterns found in the state complexes of example gridworlds.

## GRIDWORLDS



- agent
- object
- floor
- wall

Gridworlds are simplified, grid-like environments in which each *cell* of the grid may be assigned a *label*. In the example above, these labels are **agent**, **object**, floor, and wall. Such environments can be used to test and develop AI algorithms, particularly in reinforcement learning (1).

## RECONFIGURABLE SYSTEMS

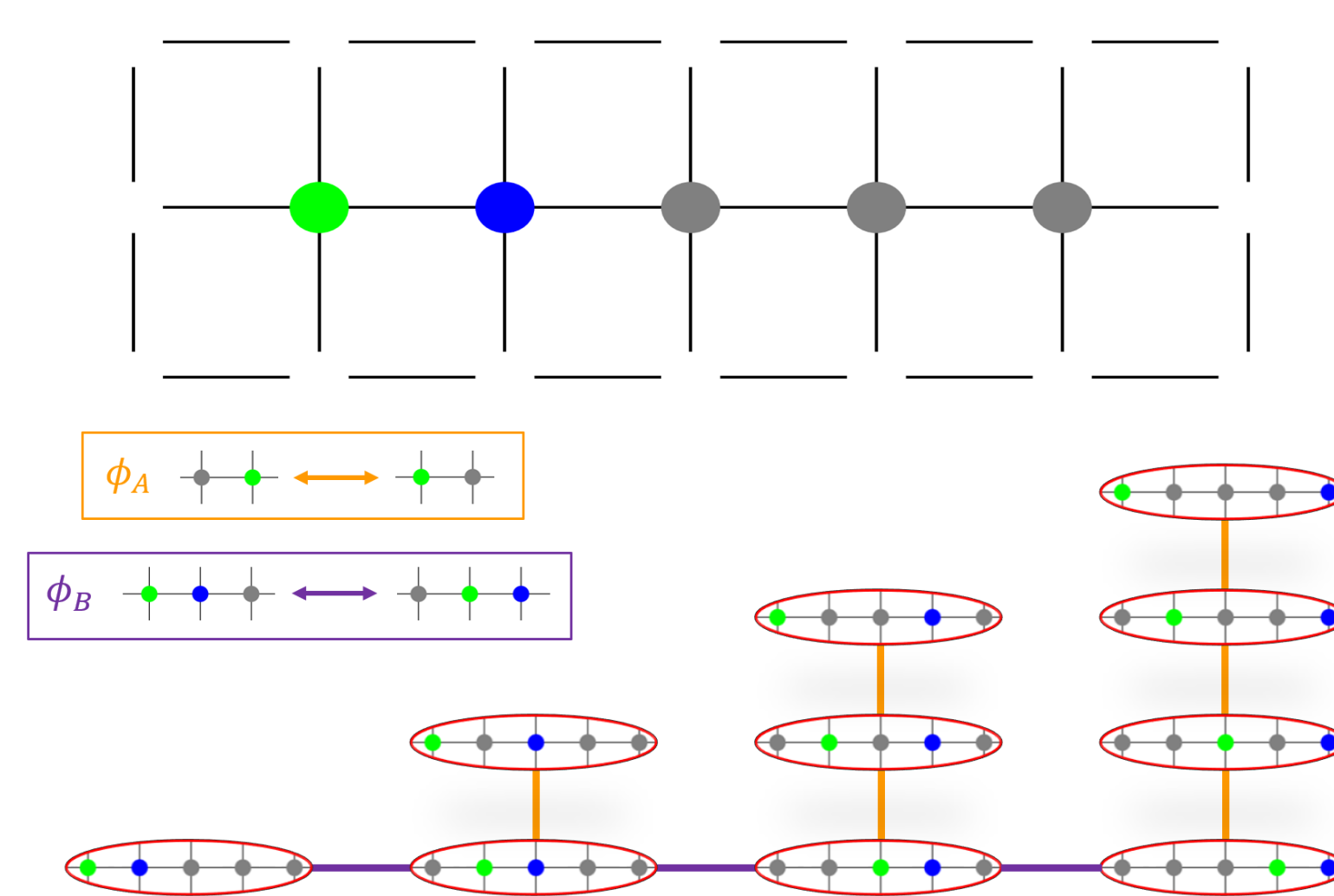
Ghrist & Peterson (2) define reconfigurable systems as a collection of labels on a graph, where local rearrangements of the labels represent reconfigurations of the system.

From (2):  $G$  is a graph.  $A$  is a set of possible labels on the vertices of  $G$ .

A *generator*  $\phi$  is a collection of three objects:

- the *support*,  $SUP(\phi) \subset G$
- the *trace*,  $TR(\phi) \subset SUP(\phi)$
- a *relabelling* for the vertex set  $TR(\phi)$

## STATE COMPLEXES



Ghrist & Peterson (2) define a **state** of a reconfigurable system as a choice of labels (chosen from  $A$ ) for every vertex of  $G$ .

$$s_i : V(G) \rightarrow A$$

The state complex  $S$  is a graph with vertices corresponding to states, with edges connecting a pair states differing by a single generator.

## REFERENCES

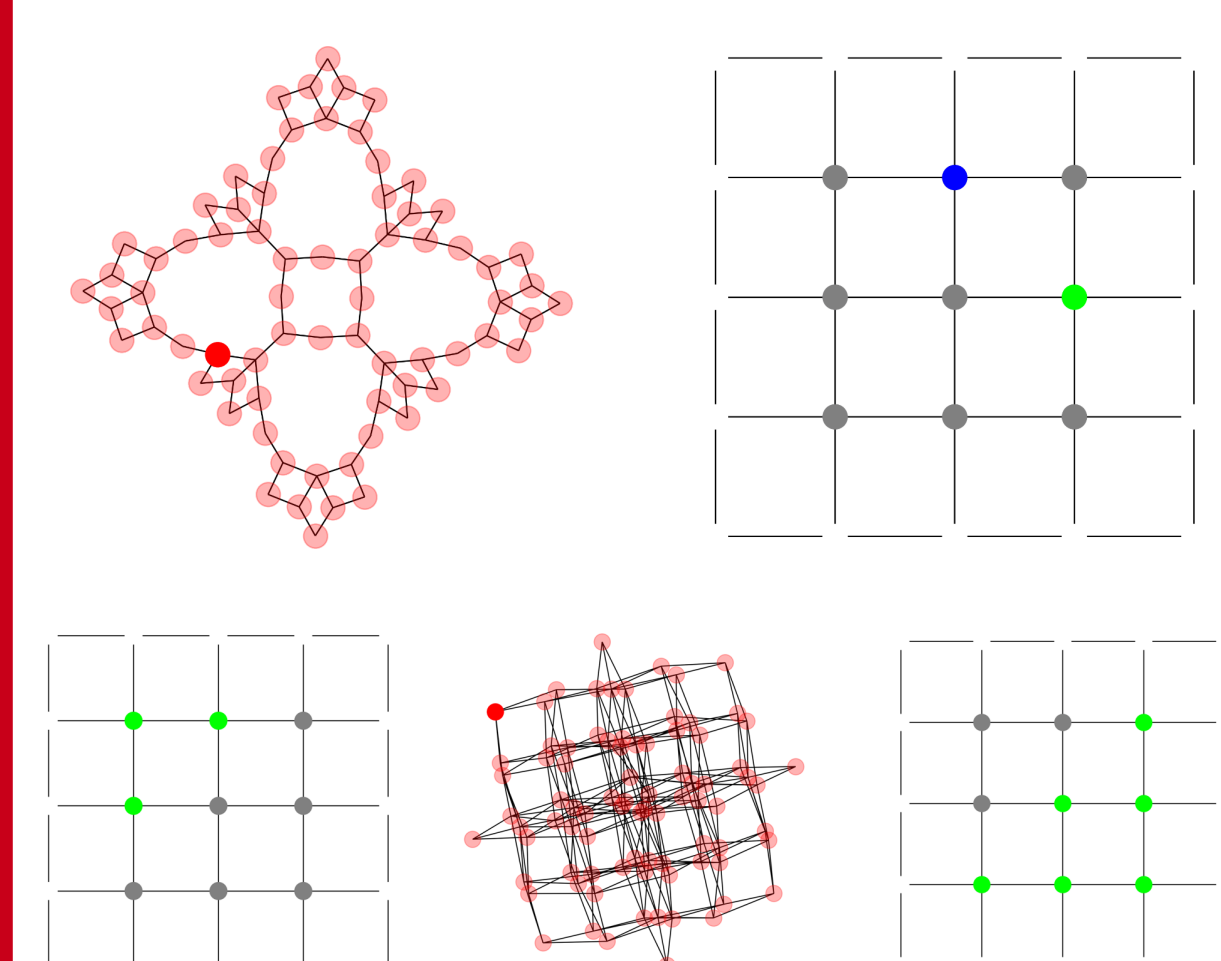
- (1) J. et al. Leike. AI safety gridworlds. *arXiv*, 1711.09883, 2017.
- (2) R. Ghrist and V. Petereson. The Geometry and Topology of Reconfiguration. *Advances in Applied Mathematics*, 38(3):302–323, 2007.

## FUTURE RESEARCH

We are currently exploring patterns and theoretical aspects of state complexes that hold across gridworlds of arbitrary size, geometry, and labelling.

For example, an upper bound for the total number of states of a gridworld without objects is  $\binom{n}{k}$  where  $n$  is the total number of non-floor and non-wall labels and  $k$  is the total number of **agent** labels. Such information may be useful to incorporate into AI algorithms or for analysis of the efficiency, accuracy, or safety of such algorithms in gridworlds.

## STATE COMPLEXES OF GRIDWORLDS



State complexes often capture symmetries in a gridworlds' geometry or labelling, as shown in these examples. The petal-like state complex above shows different scales of geometry; for each position of the **object** there exists a subgraph of 8 vertices (representing the 8 possible locations of the **agent**) which is connected to as many other subgraphs as is possible for the **agent** to push or pull the **object** to a different location from the current location represented in that subgraph.

We also find that some state complexes are subgraphs of others.

