

---

# Detecting danger in gridworlds using Gromov’s Link Condition

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Gridworlds have been long-utilised in AI research, particularly in reinforcement  
2        learning, as they provide simple yet scalable models for many real-world applica-  
3        tions such as robot navigation, emergent behaviour, and operations research. We  
4        initiate a study of gridworlds using the mathematical framework of *reconfigurable*  
5        *systems* and *state complexes* due to Abrams, Ghrist & Peterson. State complexes  
6        represent all possible configurations of a system as a single geometric space, thus  
7        making them conducive to study using geometric, topological, or combinatorial  
8        methods. The main contribution of this work is a modification to the original  
9        Abrams, Ghrist & Peterson setup which we introduce to capture agent braiding and  
10       thereby more naturally represent the topology of gridworlds. With this modification,  
11       the state complexes may exhibit geometric defects (failure of *Gromov’s Link Condi-*  
12       *tion*). Serendipitously, we discover these failures occur exactly where undesirable  
13       or dangerous states appear in the gridworld. Our results therefore provide a novel  
14       method for seeking guaranteed safety limitations in discrete task environments  
15       with single or multiple agents, and offer useful safety information (in geometric  
16       and topological forms) for incorporation in or analysis of machine learning sys-  
17       tems. More broadly, our work introduces tools from geometric group theory and  
18       combinatorics to the AI community and demonstrates a proof-of-concept for this  
19       geometric viewpoint of the task domain through the example of simple gridworld  
20       environments.

## 21 1 Introduction

22       The notion of a state (or configuration/phase) space is commonly used in mathematics and physics to  
23       represent all the possible states of a given system as a single geometric (or topological) object. This  
24       perspective provides a bridge which allows for tools from geometry and topology to be applied to  
25       the system of concern. Moreover, certain features of a given system are reflected by some geometric  
26       aspects of the associated state space (such as gravitational force being captured by *curvature* in  
27       spacetime). Thus, insights into the structure of the original system can be gleaned by reformulating  
28       them in geometric terms.

29       In discrete settings, state spaces are typically represented by graphs or their higher dimensional  
30       analogues such as simplicial complexes or cube complexes. Abrams, Ghrist & Peterson’s *state*  
31       *complexes* [AG04, GP07] provide a general framework for representing discrete reconfigurable  
32       systems as non-positively curved (NPC) cube complexes, giving access to a wealth of mathematical  
33       and computational benefits via efficient optimisation algorithms guided by geometric insight [AOS12].  
34       These have been used to develop efficient algorithms for robotic motion planning [ABY14, ABCG17]  
35       and self-reconfiguration of modular robots [LR10]. NPC cube complexes also possess rich hyperplane  
36       structures which geometrically capture binary classification [CN05, Wis12, Sag14]. However, their  
37       broader utility to fields like artificial intelligence (AI) has until now been relatively unexplored.

38 Our main contribution is the first application of this geometric approach (of using state complexes)  
 39 to the setting of multi-agent gridworlds. We introduce a natural modification to the state complex  
 40 appropriate to the setting of gridworlds (to capture the braiding or relative movements of agents);  
 41 however, this can lead to state complexes which are no longer NPC. Nevertheless, by applying  
 42 Gromov’s Link Condition, we completely characterise when positive curvature occurs in our new  
 43 state complexes, and relate this to features of the gridworlds (see Theorem 5.2). Serendipitously,  
 44 we discover that the states where Gromov’s Link Condition fails are those in which agents can  
 45 potentially collide. In other words, collision-detection is naturally embedded into the intrinsic  
 46 geometry of the system. Current approaches to collision-detection and navigation during multi-  
 47 agent navigation often rely on modelling and predicting collisions based on large training datasets  
 48 [KFG19, FLLP20, QZC<sup>+</sup>21] or by explicitly modelling physical movements [KIU21]. However,  
 49 our approach is purely geometric, requires no training, and can accommodate many conceivable types  
 50 of actions and inter-actions, not just simple movements.

51 Our work relates to a growing body of research aimed towards understanding, from a geometric  
 52 perspective, how deep learning methods transform input data into decisions, memories, or actions  
 53 [HR17, LAG<sup>+</sup>20, SPG<sup>+</sup>21, AVBP21, SMK11]. However, such studies do not usually incorporate  
 54 the geometry of the originating domain or task in a substantial way, before applying or investigating  
 55 the performance of learning algorithms – and even fewer do so for multi-agent systems. One possible  
 56 reason for this is a lack of known suitable tools. Our experimental and theoretical results show there  
 57 is a wealth of geometric information available in (even very simple) task domains, which is accessible  
 58 using tools from geometric group theory and combinatorics.

## 59 2 State complex of a gridworld

60 A *gridworld* is a two-dimensional, flat array of *cells* arranged in a grid,  
 61 much like a chess or checker board. Each cell can be occupied or un-  
 62 occupied. A cell may be occupied, in our setting, by one and only  
 63 one freely-moving agent or movable object. Other gridworlds may in-  
 64 clude rewards, punishments, buttons, doors, locks, keys, checkpoints,  
 65 dropbears, etc., much like many basic video games. Gridworlds have  
 66 been a long-utilised setting in AI research, particularly reinforcement  
 67 learning, since they are simple yet scalable in size and sophistication  
 68 [DHLKT20, WKK20]. They also offer clear analogies to many real-  
 69 world applications or questions, such as robot navigation [HHA21], emer-  
 70 gent behaviour [KAP20], and operations research [LSS<sup>+</sup>21]. For these  
 71 reasons, gridworlds have also been developed for formally specifying  
 72 problems in AI safety [LMK<sup>+</sup>17].

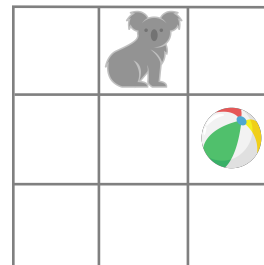


Figure 1: A  $3 \times 3$  gridworld with one agent (a koala) and one object (a beach ball).

73 A *state* of a gridworld can be encoded by assigning each cell a *label*. In  
 74 the example shown in Figure 1, these labels are shown for an agent, an  
 75 object, and empty floor. A change in the state, such as an agent moving  
 76 from one cell to an adjacent empty cell, can be encoded by *relabelling* the  
 77 cells involved. This perspective allows us to take advantage of the notion of *reconfigurable systems*  
 78 as introduced by Abrams, Ghrist & Peterson [AG04, GP07].

79 More formally, consider a graph  $G$  and a set  $\mathcal{A}$  of labels. A *state* is a function  $s : V(G) \rightarrow \mathcal{A}$ , i.e. an  
 80 assignment of a label to each vertex of  $G$ . A possible relabelling is encoded using a *generator*  $\phi$ ; this  
 81 comprises the following data:

- 82 • a subgraph  $SUP(\phi) \subseteq G$  called the *support*;
- 83 • a subgraph  $TR(\phi) \subseteq SUP(\phi)$  called the *trace*; and
- 84 • an unordered pair of *local states*

$$u_0^{loc}, u_1^{loc} : V(SUP(\phi)) \rightarrow \mathcal{A}$$

85 that agree on  $V(SUP(\phi)) - V(TR(\phi))$  but differ on  $V(TR(\phi))$ .

86 A generator  $\phi$  is *admissible* at a state  $s$  if  $s|_{SUP(\phi)} = u_0^{loc}$  (or  $u_1^{loc}$ ), in other words, if the assignment  
 87 of labels to  $V(SUP(\phi))$  given by  $s$  completely matches the labelling from (exactly) one of the two

88 local states. If this holds, we may apply  $\phi$  to the state  $s$  to obtain a new state  $\phi[s]$  given by

$$\phi[s](v) := \begin{cases} u_1^{loc}(v), & v \in V(TR(\phi)) \\ s(v), & \text{otherwise.} \end{cases}$$

89 This has the effect of relabelling the vertices in (and only in)  $TR(\phi)$  to match the other local state  
 90 of  $\phi$ . Since the local states are unordered, if  $\phi$  is admissible at  $s$  then it is also admissible at  $\phi[s]$ ;  
 91 moreover,  $\phi[\phi[s]] = s$ .

92 **Definition 2.1** (Reconfigurable system [AG04, GP07]). A *reconfigurable system* on a graph  $G$  with  
 93 a set of labels  $\mathcal{A}$  consists of a set of generators together with a set of states closed under the action of  
 94 admissible generators.

95 Configurations and their reconfigurations can be used to construct a *state graph* (or transition graph),  
 96 which represents all possible states and transitions between these states in a reconfigurable system.  
 97 More formally:

98 **Definition 2.2** (State graph). The state graph  $\mathcal{S}^{(1)}$  associated to a reconfigurable system has as its  
 99 vertices the set of all states, with edges connecting pairs of states differing by a single generator.

100 Let us now return our attention to gridworlds. We define a graph  $G$  to have vertices corresponding to  
 101 the cells of a gridworld, with two vertices declared adjacent in  $G$  exactly when they correspond to  
 102 neighbouring cells (i.e. they share a common side). Our set of labels is chosen to be

$$\mathcal{A} = \{ \text{'agent'}, \text{'object'}, \text{'floor'} \}.$$

103 We do not distinguish between multiple instances of the same label. We consider two generators:

- 104 • **Push/Pull.** An agent adjacent to an object is allowed to push/pull the object if there is an  
 105 unoccupied floor cell straight in front of the object/straight behind the agent; and
- 106 • **Move.** An agent is allowed to move to a neighbouring unoccupied floor cell.

107 These two generators have the effect of enabling agents to at any time move in any direction not  
 108 blocked by objects or other agents, and for agents to push or pull objects within the environment into  
 109 any configuration if there is sufficient room to move. For both types of generators, the trace coincides  
 110 with the support. For the Push/Pull generator, the support is a row or column of three contiguous  
 111 cells, whereas for the Move generator, the support is a pair of neighbouring cells. A simple example  
 112 of a state graph, together with the local states for the two generator types, is shown in Figure 2.

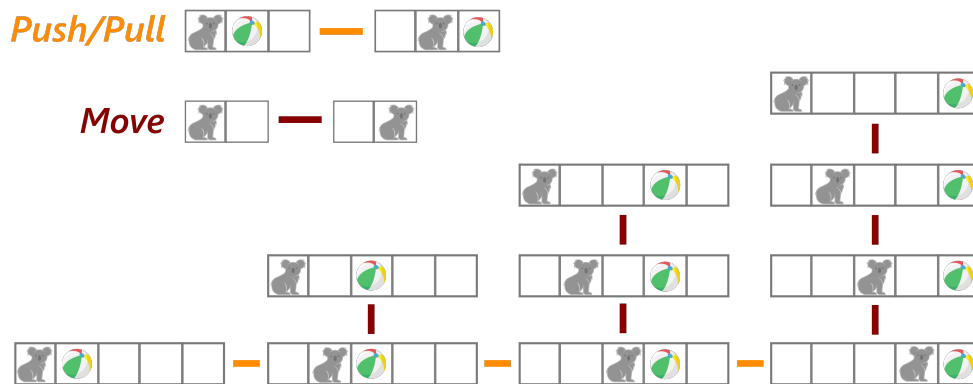


Figure 2: An example  $1 \times 5$  gridworld with one agent and one object with two generators – Push/Pull and Move – and the resulting state graph. In the state graph, edge colours indicate the generator type which relabels the gridworld.

113 In a typical reconfigurable system, there may be many admissible generators at a given state  $s$ . If the  
 114 trace of an admissible generator  $\phi_1$  is disjoint from the support of another admissible generator  $\phi_2$ ,  
 115 then  $\phi_2$  remains admissible at  $\phi_1[s]$ . This is because the relabelling by  $\phi_1$  does not interfere with

116 the labels on  $SUP(\phi_2)$ . More generally, a set of admissible generators  $\{\phi_1, \dots, \phi_n\}$  at a state  $s$   
 117 commutes if  $SUP(\phi_i) \cap TR(\phi_j) = \emptyset$  for all  $i \neq j$ . When this holds, these generators can be applied  
 118 independently of one another, and the resulting state does not depend on the order in which they are  
 119 applied. A simple example of this in the context of gridworlds is a large room with  $n$  agents spread  
 120 sufficiently far apart to allow for independent simultaneous movement.

121 Abrams, Ghrist & Peterson represent this mutual  
 122 commutativity by adding higher dimensional  
 123 cubes to the state graph to form a cube complex  
 124 called the *state complex*. We give an informal  
 125 definition here, and refer to their papers for  
 126 the precise formulation [AG04, GP07]. Further  
 127 background on cube complexes can be found  
 128 in [Wis12, Sag14]. If  $\{\phi_1, \dots, \phi_n\}$  is a set of  
 129 commuting admissible generators at a state  $s$   
 130 then there are  $2^n$  states that can be obtained  
 131 by applying any subset of these generators to  $s$ .  
 132 These  $2^n$  states form the vertices of an  $n$ -cube  
 133 in the state complex. Each  $n$ -cube is bounded  
 134 by  $2n$  faces, where each face is an  $(n-1)$ -cube:  
 135 by disallowing a generator  $\phi_i$ , we obtain a pair  
 136 of faces corresponding to those states (in the  
 137 given  $n$ -cube) that agree with one of the two  
 138 respective local states of  $\phi_i$  on  $SUP(\phi_i)$ .

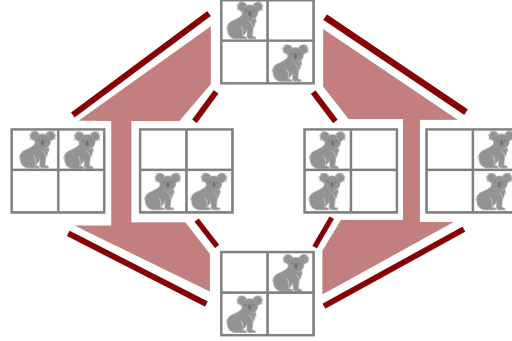


Figure 3: State complex of a  $2 \times 2$  gridworld with two agents. Shading indicates squares attached to the surrounding 4-cycles.

139 **Definition 2.3** (State complex). The *state complex*  $\mathcal{S}$  of a reconfigurable system is the cube complex  
 140 constructed from the state graph  $\mathcal{S}^{(1)}$  by inductively adding cubes as follows: whenever there is a  
 141 set of  $2^n$  states related by a set of  $n$  admissible commuting generators, we add an  $n$ -cube so that its  
 142 vertices correspond to the given states, and so that its  $2n$  boundary faces are identified with all the  
 143 possible  $(n-1)$ -cubes obtained by disallowing a generator. In particular, every cube is uniquely  
 144 determined by its vertices.

145 In our gridworlds setting, each generator involves exactly one agent. This means commuting  
 146 generators can only occur if there are multiple agents. A simple example of a state complex for  
 147 two agents in a  $2 \times 2$  room is shown in Figure 3. Note that there are six embedded 4-cycles in the  
 148 state graph, however, only two of these are filled in by squares: these correspond to independent  
 149 movements of the agents, either both horizontally or both vertically.

### 150 3 Exploring gridworlds with state complexes

151 To compute the state complex of a (finite) gridworld, we first initialise an empty graph  $\mathcal{G}$  and an  
 152 empty ‘to-do’ list  $\mathcal{L}$ . As input, we take a chosen state of the gridworld to form the first vertex of  $\mathcal{G}$   
 153 and also the first entry on  $\mathcal{L}$ . The state complex is computed according to a breadth-first search by  
 154 repeatedly applying the following:

- 155 • Let  $v$  be the first entry on  $\mathcal{L}$ . List all admissible generators at  $v$ . For each such generator  $\phi$ :
  - 156 – If  $\phi[v]$  already appears as a vertex of  $\mathcal{G}$ , add an edge between  $v$  and  $\phi[v]$  (if it does not  
 157 already exist).
  - 158 – If  $\phi[v]$  does not appear in  $\mathcal{G}$ , add it as a new vertex to  $\mathcal{G}$  and add an edge connecting it  
 159 to  $v$ . Append  $\phi[v]$  to the end of  $\mathcal{L}$ .
- 160 • Remove  $v$  from  $\mathcal{L}$ .

161 The process terminates when  $\mathcal{L}$  is empty. The output is the graph  $\mathcal{G}$ . When  $\mathcal{L}$  is empty, we have  
 162 fully explored all possible states that can be reached from the initial state. It may be possible that  
 163 the true state graph is disconnected, in which case the above algorithm will only return a connected  
 164 component  $\mathcal{G}$ . For our purposes, we shall limit our study to systems with connected state graphs.  
 165 From the state graph, we construct the state complex by first finding all 4-cycles in the state graph.  
 166 Then, by examining the states involved, we can determine whether a given 4-cycle bounds a square  
 167 representing a pair of commuting moves.

168 To visualise the state complex, we first draw the state graph using the Kamada–Kawai force-directed  
 169 algorithm [KK89] which attempts to draw edges to have similar length. We then shade the region(s)  
 170 enclosed by 4-cycles representing commuting moves. For ease of visual interpretation in our figures,  
 171 we do not also shade higher-dimensional cubes, although such cubes are noticeable and can be easily  
 172 computed and visualised if desired.

173 Constructing and analysing state complexes of gridworlds is in and of itself an inter-  
 174 esting and useful way of exploring their intrinsic geometry. For example, Figure 4  
 175 shows the state complex of a  $3 \times 3$  gridworld with one agent and one object. The  
 176 state complex reveals two scales of geometry: larger ‘blobs’ of states organised in  
 177 a  $3 \times 3$  grid, representing the location of the object; and, within each blob, copies  
 178 of the room’s remaining empty space, in which the agent may walk around and ap-  
 179 proach the object to Push/Pull. Each 12-  
 180 cycle ‘petal’ represents a 12-step choreog-  
 181 raphy wherein the agent pushes and pulls  
 182 the object around in a 4-cycle in the grid-  
 183 world. In this example, the state complex is  
 184 the state graph, since there are no possible  
 185 commuting moves.  
 186 cycle ‘petal’ represents a 12-step choreog-  
 187 raphy wherein the agent pushes and pulls  
 188 the object around in a 4-cycle in the grid-  
 189 world. In this example, the state complex is  
 190 the state graph, since there are no possible  
 191 commuting moves.

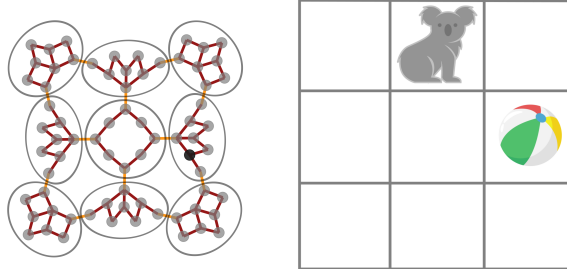


Figure 4: State complex (left) of a  $3 \times 3$  gridworld with one agent and one object (right). The darker vertex in the state complex represents the state shown in the gridworld state on the right. Edges in the state complex are coloured according to their generator – orange for Push/Pull and maroon for Move. Grey circles which group states where the ball is static have been added to illustrate the different scales of geometry.

192 The examples discussed thus far all have  
 193 planar state graphs. Planarity does not hold  
 194 in general – indeed, the  $n$ -cube graph for  
 195  $n \geq 4$  is non-planar, and a state graph can contain  $n$ -cubes if the gridworld has  $n$  agents and sufficient  
 196 space to move around. It is tempting to think that the state complex of a gridworld with more agents  
 197 should therefore look quite different to one with fewer agents. However, Figure 5 shows this may  
 198 not always be the case: there is a symmetry induced by swapping all ‘agent’ labels with ‘floor’  
 199 labels.

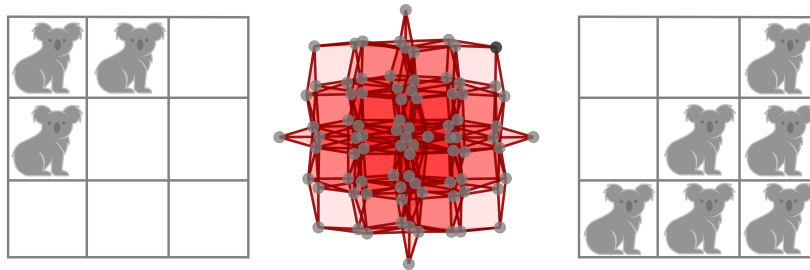


Figure 5: State complex (centre) of a  $3 \times 3$  gridworld with three agents (left) and six agents (right). They share the same state complex due to the ‘agent’  $\leftrightarrow$  ‘floor’ label inversion symmetry.

## 200 4 Dancing with myself

201 The state complex of a gridworld with  $n$  agents can be thought of as a discrete analogue of the  
 202 configuration space of  $n$  points on the 2D-plane. However, there is a problem with this analogy:  
 203 there can be ‘holes’ created by 4-cycles in the state complex where a single agent walks in a small  
 204 square-shaped dance by itself, as shown in Figure 6.

205 The presence of these holes would suggest something meaningful about the underlying gridworld’s  
 206 intrinsic topology, e.g., something obstructing the agent’s movement at that location in the gridworld  
 207 that the agent must move around. In reality, the environment is essentially a (discretised) 2D-plane  
 208 with nothing blocking the agent from traversing those locations. Indeed, these ‘holes’ are uninteresting

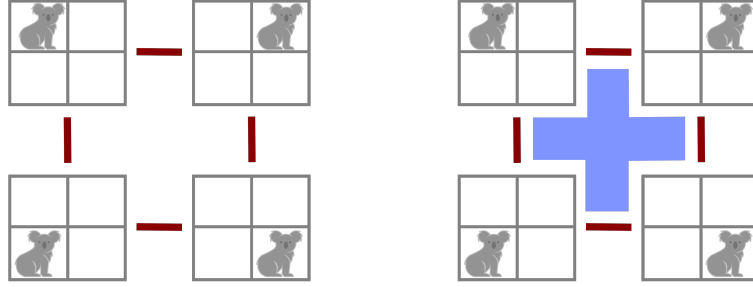


Figure 6: State complex of a  $2 \times 2$  gridworld with one agent under the original definition of Abrams, Ghrist & Peterson [AG04, GP07] (left) and with our modification (right). The blue shading is a filled in square indicating a *dance*.

209 topological quirks which arise due to the representation of the gridworld as a graph. We therefore  
 210 deviate from the original definition of state complexes by Abrams, Ghrist & Peterson [AG04, GP07]  
 211 and choose to fill in these ‘dance’ 4-cycles with squares.<sup>1</sup>

212 Formally, we define a **dance**  $\delta$  to comprise the following data:

- 213 • the support  $SUP(\delta)$  given by a  $2 \times 2$  subgrid in the gridworld,
- 214 • four local states defined on  $SUP(\delta)$ , each consisting of exactly one agent label and three  
 215 floor labels, and
- 216 • four Move generators, each of which transitions between two of the four local states (as in  
 217 Figure 6).

218 We say that  $\delta$  is *admissible* at a state  $s$  if  $s|_{SUP(\delta)}$  agrees with one of the four local states of  $\delta$ .  
 219 Moreover, these four local states are precisely the states that can be reached when we apply some  
 220 combination of the four constituent Moves. We do not define the trace of a dance, however, we may  
 221 view the trace of each of the four constituent Moves as subgraphs of  $SUP(\delta)$ .

222 The notion of commutativity can be extended to incorporate dancing. Suppose that we have a set  
 223  $\{\phi_1, \dots, \phi_l, \delta_1, \dots, \delta_m\}$  of  $l$  admissible generators and  $m$  admissible dances at a state  $s$ . We say  
 224 that this set *commutes* if the supports of its elements are pairwise disjoint. When this holds, there  
 225 are  $2^{l+2m}$  possible states that can be obtained by applying some combination of the generators and  
 226 dances to  $s$ : there are two choices of local state for each  $\phi_i$ , and four for each  $\delta_j$ . We capture this  
 227 extended notion of commutativity by attaching additional cubes to the state complex to form our  
 228 modified state complex.

229 **Definition 4.1** (Modified state complex). The *modified state complex*  $\mathcal{S}'$  of a gridworld is the cube  
 230 complex obtained by filling in the state graph  $\mathcal{S}^{(1)}$  with higher dimensional cubes whenever there  
 231 is a set of commuting moves or dances. Specifically, whenever a set of  $2^{l+2m}$  states are related by  
 232 a commuting set of  $l$  generators and  $m$  dances, we add an  $n$ -cube having the given set of states  
 233 as its vertices, where  $n = l + 2m$ . Each of the  $2n$  faces of such an  $n$ -cube is identified with an  
 234  $(n - 1)$ -cube obtained by either disallowing a generator  $\phi_i$  and choosing one of its two local states,  
 235 or replacing a dance  $\delta_j$  with one of its four constituent Moves.

236 Our modification removes uninteresting topology. This can be observed by examining 4-cycles in  $\mathcal{S}'$ .  
 237 On the one hand, some 4-cycles are trivial (they can be ‘filled in’): *dancing-with-myself* 4-cycles,  
 238 and *commuting moves* (two agents moving back and forth) 4-cycles (which were trivial under the  
 239 original definition). These represent trivial movements of agents relative to one another. On the other  
 240 hand, there is a non-trivial 4-cycle in the state complex for two agents in a  $2 \times 2$  room, as can be seen  
 241 in the centre of Figure 3 (here, no dancing is possible so the modified state complex is the same as the  
 242 original). This 4-cycle represents the two agents moving half a ‘revolution’ relative to one another –

<sup>1</sup>Ghrist and Peterson themselves ask if there could be better ways to complete the state graph to a higher-dimensional object with better properties (Question 6.4 in [GP07]).



243 indeed, performing this twice would give a full revolution. (There are three other non-trivial 4-cycles,  
 244 topologically equivalent to this central one, that also achieve the half-revolution.)

245 In a more topological sense<sup>2</sup>, by filling in such squares and higher dimensional cubes, our state  
 246 complexes capture the non-trivial, essential relative movements of the agents. This can be used  
 247 to study the braiding or mixing of agents, and also allows us to consider path-homotopic paths as  
 248 ‘essentially’ the same. One immediate difference this creates with the original state complexes is a  
 249 loss of symmetries like those shown in Figure 5, since there is no label inversion for a dance when  
 250 other agents are crowding the dance-floor.

## 251 5 Gromov’s Link Condition

252 The central geometric characteristic of Abrams, Ghrist, & Peterson’s state complexes is that they  
 253 are *non-positively curved* (NPC). Indeed, this local geometric condition is conducive for developing  
 254 efficient algorithms for computing geodesics. However, with our modified state complexes, this NPC  
 255 geometry is no longer guaranteed – we test for this on a vertex-by-vertex basis using a classical  
 256 geometric result due to Gromov (see also Theorem 5.20 of [BH99] and [Sag14]).

257 **Theorem 5.1** (Gromov’s Link Condition [Gro87]). *A finite-dimensional cube complex is NPC if and*  
 258 *only if the link of every vertex is a flag simplicial complex.*  $\square$

259 We provide a brief mathematical back-  
 260 ground on cube complexes and the finer  
 261 details of Gromov’s Link Condition in Ap-  
 262 pendix A.1. For our current purposes, it is  
 263 sufficient to know that under the Abrams,  
 264 Ghrist & Peterson setup, if  $v$  is a state in  $\mathcal{S}$   
 265 then the vertices of its link  $lk(v)$  represent  
 266 the possible admissible generators at  $v$ .  
 267 Since cubes in  $\mathcal{S}$  are associated with com-  
 268 muting sets of generators, each simplex in  
 269  $lk(v)$  represents a set of commuting gener-  
 270 ators. Gromov’s Link Condition for  $lk(v)$   
 271 can be reinterpreted as follows: whenever  
 272 a set of admissible generators is *pairwise*  
 273 commutative, then it is *setwise* commuta-  
 274 tive. Using this, it is straightforward for  
 275 Abrams, Ghrist & Peterson to verify that  
 276 this always holds for their state complexes  
 277 (see Theorem 4.4 of [GP07]).

278 For our modified states complexes, the sit-  
 279 uation is not as straightforward. The key  
 280 issue is that our cubes do not only arise  
 281 from commuting generators – we must  
 282 take dances into account. Indeed, when  
 283 attempting to prove that Gromov’s Link  
 284 Condition holds, we discovered some very  
 285 simple gridworlds where it actually fails;  
 286 see Figure 7 and Appendix A.4.

287 Failure of the Link Condition can indicate  
 288 available moves at some state that cannot be safely performed simultaneously and independently  
 289 without risking collisions between labels. Another interpretation of positive curvature in this context  
 290 is something akin to what real-time computer strategy games call ‘fog of war’ (distance-dependent  
 291 limiting of observations which extends from the player-controlled agents), and more specifically  
 292 the viewable distance from an agent’s line-of-sight. Such fog makes AI systems operating in such  
 293 environments particularly challenging, although remarkable success has been achieved in games like  
 294 StarCraft [VBC<sup>+</sup>19].

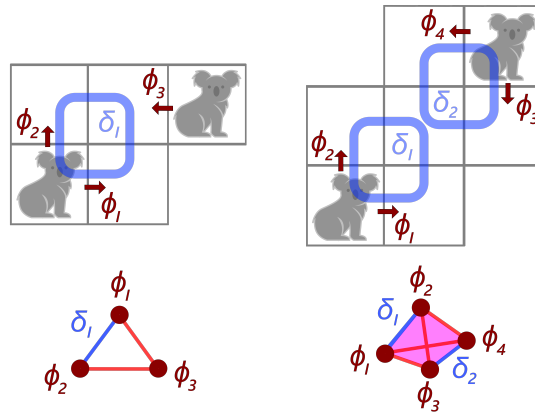


Figure 7: The two situations which lead to failure of Gromov’s Link Condition in multi-agent gridworlds. Maroon arrows indicate admissible moves and blue squares indicate admissible dances. Note that in the links (bottom row), the triangle is missing in the left example, while the (solid) tetrahedron is missing in the right (however, all 2D faces are present). This is due to the respective collections of moves and dances failing to commute – an agent interrupts the other’s dance (left) or two dances collide (right).

<sup>2</sup>By considering the fundamental group.

295 Despite this apparent drawback, we nevertheless show that Figure 7 accounts for all the possible  
296 failures of Gromov’s Link Condition in the setting of agent-only gridworlds<sup>3</sup>.

297 **Theorem 5.2** (Gromov’s Link Condition in the modified state complex). *Let  $v$  be a vertex in the*  
298 *modified state complex  $S'$  of an agent-only gridworld. Then*

- 299 •  $lk(v)$  satisfies Gromov’s Link Condition if and only if it has no empty 2–simplices nor  
300 3–simplices<sup>4</sup>, and
- 301 • if  $lk(v)$  fails Gromov’s Link Condition then there exist a pair of agents whose positions  
302 differ by either a knight move or a 2–step bishop move (as in Figure 7).

303 We provide a proof in Appendix A.2. Consequently, if the Link Condition fails at all, it must fail at  
304 dimension 2 or 3. This can be interpreted as saying that we only need a bounded amount of foresight  
305 to detect potential collisions: under fog-of-war, each agent needs a line-of-sight of only four moves.

306 Positive curvature could indicate collisions between any specified labels (e.g., objects), however, for  
307 this interpretation to be valid we would need to carefully identify which other potential cycles in the  
308 state complex ought to be filled in. Doing this in a ‘natural’ way is in itself a non-trivial task, and is  
309 the subject of further investigation.

## 310 6 Experiments and applications

311 Although our main contribution is theoretical, we conduct some small initial experiments to demon-  
312 strate the type of information which can be captured in the geometry and topology (see Appendix  
313 A.4). To run these experiments, we developed and used a custom Python-based tool (detailed in  
314 Appendix A.3). Our focus on small rooms is largely expository, i.e., they are the simplest non-trivial  
315 examples illustrating the key features we want to isolate, and naturally reoccur in all larger rooms.  
316 Our intention is also to demonstrate a combinatorial explosion in the number of states. We don’t  
317 recommend constructing the entire state complex in practical applications (indeed, to implement  
318 addition of integers on a computer, it is infeasible and unnecessary to construct *all* integers).

319 *Remark 6.1.* By a simple counting argument, one can deduce the total number of states in a gridworld.  
320 For an agent-only gridworld with  $n$  cells and  $k$  agents, there is a total of  $\binom{n}{k}$  states. If there are  $n$   
321 cells,  $k$  agents, and  $j$  objects, then there are  $\binom{n}{k} \binom{n-k}{j}$  states. Thus, even for a moderately sized  
322  $10 \times 10$  room with 50 agents, there are  $\binom{100}{50} \approx 1.008 \times 10^{29}$  vertices in the state complex.

323 By Theorem 5.2, checking if  $lk(v)$  satisfies Gromov’s Link Condition requires computing the link  
324 only up to dimension 3 and then checking whether it is a flag complex; if not, we count the number  
325 of empty simplices. Checking this for a given vertex in the state complex is not too computationally  
326 demanding, however when a state complex has many vertices it becomes more difficult. In practical  
327 applications, such as calculating collision-avoiding navigation routes, it is – again, by Theorem  
328 5.2 – only necessary to construct a small local subcomplex. But perhaps even more importantly, to  
329 detect potential collisions between agents, it is not even necessary to construct  $lk(v)$ , since Theorem  
330 5.2 provides a computational shortcut: just check for supports of knight or two-step bishop moves  
331 between agents.

332 By using Gromov’s Link Condition, we can identify a precise measure of how far ahead agents ought  
333 to look in order to safely proceed without fear of collisions. Appendix A.4 gives a summary analysis  
334 of a  $3 \times 3$  room with varying numbers of agents. We noticed several symmetries. Commuting moves  
335 and the number of states have a symmetry about 4.5 agents (due to the label-inversion symmetry as  
336 previously illustrated in Figure 5). However, curiously, the number of dances has a symmetry about  
337 3.5 agents. This difference leads to the asymmetrical distribution of positive curvature and failures  
338 of Gromov’s Link Condition – which, while maximal for 3 agents as a proportion of total states,  
339 exhibited the highest mean failure rate for 4 agents.

---

<sup>3</sup>While writing this paper, the first author was involved in two scooter accidents – collisions involving only agents (luckily without serious injury). So, while this class of gridworlds is strictly smaller than those also involving objects or other labels, it is by no means an unimportant one. If only the scooters had Gromov’s Link Condition checkers!

<sup>4</sup>In other words, if there are no “hollow” triangles or tetrahedra like those in Figure 7.



340 This shows that, heuristically, we expect most states to satisfy NPC (see Appendix A.4), and so  
341 existing greedy algorithms [AOS12] for calculating geodesics will work well in most situations.  
342 However, to implement an efficient, collision-free path-finding algorithm in our modified state  
343 complexes, we need to add an additional check. Specifically, when we are near a potentially  
344 dangerous state, we should implement a predefined ‘detour’ to avoid the collision, which can be done  
345 on a local basis using the identified supports which lead to positive curvature (as in Figure 7).

## 346 7 Conclusions and future directions

347 This study presents novel applications of tools from geometric group theory and combinatorics to the  
348 AI research community, opening new ways for recasting and analysing AI problems as geometric ones.  
349 Using these tools, we show an example of how the intrinsic geometry of a task space serendipitously  
350 embeds safety information and makes it possible to determine how far ahead in time an AI system  
351 needs to observe to be guaranteed of avoiding dangerous actions.

352 Leike et al. [LMK<sup>+</sup>17] show deep reinforcement learning agents cannot solve many AI safety prob-  
353 lems specified on gridworlds, e.g., minimising unwanted side-effects or ensuring robustness to agent  
354 self-modification. Having described the agent-only case in this study, there is now ripe opportunity to  
355 account for positive curvature or other geometric features arising due to other labels or generators  
356 (actions) present in specified AI safety problems, e.g., agents pushing/pulling objects, pressing  
357 buttons, modifying their form or behaviour, rewards/punishments, opening/unlocking doors, etc.. By  
358 considering *directed* modified state complexes, irreversible actions can be captured by “invariant  
359 subcomplexes” (i.e., you can’t escape from them), allowing geometric study of the tree/flowchart of  
360 irreversible actions and related recurrence/transience. Braiding can be used to study route planning,  
361 back-tracking, cooperation, assembly, and topological entropy in congestion [Ghr09]. Numerous  
362 extensions are possible, allowing us to study and geometrically represent further problems with a  
363 view to developing efficient, geometrically-inspired local algorithms without the need for training.

364 Do learning algorithms already implement such geometrically-inspired algorithms, the related ge-  
365 ometry, or approximations thereof? To find out, we are investigating how modified state complexes  
366 map to learned internal representations of neural networks trained to predict multi-agent gridworld  
367 dynamics. This mapping connects the geometry and topology of a task space directly to optimisation  
368 procedures and learning trajectories in latent representation spaces, highlighting unexpected topologi-  
369 cal and geometric differences and opportunities for deeper insight and improvement of optimisation  
370 procedures, in the spirit of [NZL20, ZZ22]. We can also compare biological optimisation processes  
371 and internal representations of allocentric and egocentric navigation [Bur06, GHP<sup>+</sup>22], and how this  
372 interacts with the position of other agents [DJ18, SB20].

373 From a more mathematical perspective, state complexes of gridworlds give rise to an interesting class  
374 of geometric spaces. It would be worthwhile to investigate their geometric and topological properties  
375 to more deeply understand various aspects of multi-agent gridworlds. For example, for a gridworld  
376 with  $n$  agents in a sufficiently large room, we hypothesise that the modified state complex should be  
377 a classifying space for the  $n$ -strand braid group. This is clearly false when the room is packed full of  
378 agents (in which case the state complex is a single point), so it may be fruitful to determine if there is  
379 some ‘critical’ density at which a topological transition occurs.

380 Using the *failure* of Gromov’s Link Condition in an essential way appears to be a relatively unexplored  
381 approach. Indeed, much of the mathematical literature concerning cube complexes focusses on  
382 showing that the Link Condition always holds. To our knowledge, the only other works which go  
383 against this trend are [AG04], in which failure detects global disconnection of a metamorphic system,  
384 and [BDT19], where failure detects non-trivial loops on topological surfaces. It would be interesting  
385 to explore cube complexes arising in other settings where failure captures critical information.

386 A limitation of our work is that we have so far only explored very simple AI environments. Further  
387 work is needed to expand the framework and results to more general, sophisticated, and real-world  
388 environments. For this reason, although our work provides new geometric perspectives, data,  
389 and potential algorithms for an important AI safety issue, we caution against hasty real-world  
390 implementation of the main results. To avoid potential negative societal impacts, it would still be  
391 important to perform rigorous checks and tests in application domains, since our results do not  
392 directly extend to situations beyond which the stated assumptions hold.

## 393 References

- 394 [ABCG17] Federico Ardila, Hanner Bastidas, Cesar Ceballos, and John Guo, *The configuration space of a*  
395 *robotic arm in a tunnel*, SIAM Journal on Discrete Mathematics **31** (2017), no. 4, 2675–2702.
- 396 [ABY14] Federico Ardila, Tia Baker, and Rika Yatchak, *Moving robots efficiently using the combinatorics*  
397 *of CAT(0) cubical complexes*, SIAM Journal on Discrete Mathematics **28** (2014), no. 2, 986–  
398 1007.
- 399 [AG04] Aaron Abrams and Robert Ghrist, *State complexes for metamorphic robots*, The International  
400 Journal of Robotics Research **23** (2004), no. 7-8, 811–826.
- 401 [AOS12] Federico Ardila, Megan Owen, and Seth Sullivant, *Geodesics in CAT(0) cubical complexes*,  
402 Adv. in Appl. Math. **48** (2012), no. 1, 142–163.
- 403 [AVBP21] Karen Archer, Nicola Catenacci Volpi, Franziska Bröker, and Daniel Polani, *A space of goals:*  
404 *the cognitive geometry of informationally bounded agents*, arXiv:2111.03699, 2021.
- 405 [BDT19] Mark C. Bell, Valentina Disarlo, and Robert Tang, *Cubical geometry in the polygonalisation*  
406 *complex*, Math. Proc. Cambridge Philos. Soc. **167** (2019), no. 1, 1–22.
- 407 [BH99] Martin R. Bridson and André Haefliger, *Metric spaces of non-positive curvature*, Grundlehren  
408 der mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences], vol.  
409 319, Springer-Verlag, Berlin, 1999.
- 410 [Bur06] Neil Burgess, *Spatial memory: how egocentric and allocentric combine*, Trends in Cognitive  
411 Sciences **10** (2006), no. 12, 551–557.
- 412 [CN05] Indira Chatterji and Graham Niblo, *From wall spaces to CAT(0) cube complexes*, International  
413 Journal of Algebra and Computation **15** (2005), no. 05n06, 875–885.
- 414 [DJ18] É. Duvelle and K.J. Jeffery, *Social spaces: Place cells represent the locations of others*, Current  
415 Biology **28** (2018), no. 6, R271–R273.
- 416 [DSHLKT20] Felipe Leno Da Silva, Pablo Hernandez-Leal, Bilal Kartal, and Matthew E. Taylor, *Uncertainty-*  
417 *aware action advising for deep reinforcement learning agents*, Proceedings of the AAAI Confer-  
418 ence on Artificial Intelligence **34** (2020), no. 04, 5792–5799.
- 419 [FLLP20] Tingxiang Fan, Pinxin Long, Wenxi Liu, and Jia Pan, *Distributed multi-robot collision avoidance*  
420 *via deep reinforcement learning for navigation in complex scenarios*, The International Journal  
421 of Robotics Research **39** (2020), no. 7, 856–892.
- 422 [GHP<sup>+</sup>22] Richard J. Gardner, Erik Hermansen, Marius Pachitariu, Yoram Burak, Nils A. Baas, Benjamin A.  
423 Dunn, May-Britt Moser, and Edvard I. Moser, *Toroidal topology of population activity in grid*  
424 *cells*, Nature **602** (2022), no. 7895, 123–128.
- 425 [Ghr09] Robert Ghrist, *Configuration spaces, braids, and robotics*, pp. 263–304, World Scientific Pub-  
426 lishing, 2009.
- 427 [GP07] R. Ghrist and V. Peterson, *The geometry and topology of reconfiguration*, Advances in Applied  
428 Mathematics **38** (2007), no. 3, 302–323.
- 429 [Gro87] M. Gromov, *Hyperbolic groups*, Essays in Group Theory (S. M. Gersten, ed.), Springer New  
430 York, New York, NY, 1987, pp. 75–263.
- 431 [HHA21] Victoria J. Hodge, Richard Hawkins, and Rob Alexander, *Deep reinforcement learning for drone*  
432 *navigation using sensor data*, Neural Computing and Applications **33** (2021), no. 6, 2015–2033.
- 433 [HR17] Michael Hauser and Asok Ray, *Principles of Riemannian geometry in neural networks*, Advances  
434 in Neural Information Processing Systems (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach,  
435 R. Fergus, S. Vishwanathan, and R. Garnett, eds.), vol. 30, Curran Associates, Inc., 2017.
- 436 [KAP20] Ivana Kajic, Eser Aygün, and Doina Precup, *Learning to cooperate: Emergent communication in*  
437 *multi-agent navigation*, 42nd Annual Virtual Meeting of the Cognitive Science Society (CogSci  
438 2020), 2020.
- 439 [KFGE19] Zac Kenton, Angelos Filos, Yarin Gal, and Owain Evans, *Generalizing from a few environments*  
440 *in safety-critical reinforcement learning*, Safe Machine Learning workshop at ICLR (2019), 1–9.
- 441 [KIU21] Takeshi Kano, Mayuko Iwamoto, and Daishin Ueyama, *Decentralised control of multiple*  
442 *mobile agents for quick, smooth, and safe movement*, Physica A: Statistical Mechanics and its  
443 Applications **572** (2021), 125898.
- 444 [KK89] Tomihisa Kamada and Satoru Kawai, *An algorithm for drawing general undirected graphs*,  
445 Information Processing Letters **31** (1989), no. 1, 7–15.
- 446 [LAG<sup>+</sup>20] Na Lei, Dongsheng An, Yang Guo, Kehua Su, Shixia Liu, Zhongxuan Luo, Shing-Tung Yau,  
447 and Xianfeng Gu, *A geometric understanding of deep learning*, Engineering **6** (2020), no. 3,  
448 361–374.

- 449 [LMK<sup>+</sup>17] Jan Leike, Miljan Martic, Victoria Krakovna, Pedro A. Ortega, Tom Everitt, Andrew Lefrancq,  
450 Laurent Orseau, and Shane Legg, *AI safety gridworlds*, arXiv:1711.09883, 2017.
- 451 [LR10] Tom Larkworthy and Subramanian Ramamoorthy, *An efficient algorithm for self-reconfiguration  
452 planning in a modular robot*, 2010 IEEE International Conference on Robotics and Automation,  
453 2010, pp. 5139–5146.
- 454 [LSS<sup>+</sup>21] Florian Laurent, Manuel Schneider, Christian Scheller, Jeremy Watson, Jiaoyang Li, Zhe Chen,  
455 Yi Zheng, Shao-Hung Chan, Konstantin Makhnev, Oleg Svidchenko, Vladimir Egorov, Dmitry  
456 Ivanov, Aleksei Shpilman, Evgenija Spirovska, Oliver Tanevski, Aleksandar Nikov, Ramon  
457 Grunder, David Galevski, Jakov Mitrovski, and Sharada Mohanty, *Flatland competition 2020:  
458 MAPF and MARL for efficient train coordination on a grid world*, pp. 275–301, PMLR, 08 2021.
- 459 [NZL20] Gregory Naitzat, Andrey Zhitnikov, and Lek-Heng Lim, *Topology of deep neural networks*,  
460 Journal of Machine Learning Research **21** (2020), no. 184, 1–40.
- 461 [QZC<sup>+</sup>21] Zengyi Qin, Kaiqing Zhang, Yuxiao Chen, Jingkai Chen, and Chuchu Fan, *Learning safe multi-  
462 agent control with decentralized neural barrier certificates*, International Conference on Learning  
463 Representations, 2021.
- 464 [Sag14] Michah Sageev, *CAT(0) cube complexes and groups*, Geometric group theory, IAS/Park City  
465 Math. Ser., vol. 21, Amer. Math. Soc., Providence, RI, 2014, pp. 7–54.
- 466 [SB20] Christina J. Sutherland and David K. Bilkey, *Hippocampal coding of conspecific position*, Brain  
467 Research **1745** (2020), 146920.
- 468 [SMK11] Jeremy Stober, Risto Miikkulainen, and Benjamin Kuipers, *Learning geometry from sensorimotor  
469 experience*, 2011 IEEE International Conference on Development and Learning (ICDL), vol. 2,  
470 2011, pp. 1–6.
- 471 [SPG<sup>+</sup>21] Cory Stephenson, Suchismita Padhy, Abhinav Ganesh, Yue Hui, Hanlin Tang, and SueYeon  
472 Chung, *On the geometry of generalization and memorization in deep neural networks*, Interna-  
473 tional Conference on Learning Representations, 2021.
- 474 [VBC<sup>+</sup>19] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik,  
475 Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh,  
476 Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P  
477 Agapiou, Max Jaderberg, Alexander S Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin  
478 Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L Paine, Caglar Gulcehre, Ziyu Wang,  
479 Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney,  
480 Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps,  
481 and David Silver, *Grandmaster level in StarCraft II using multi-agent reinforcement learning*,  
482 Nature **575** (2019), no. 7782, 350–354.
- 483 [Wis12] Daniel T. Wise, *From riches to raags: 3-manifolds, right-angled Artin groups, and cubical geom-  
484 etry*, CBMS Regional Conference Series in Mathematics, vol. 117, Published for the Conference  
485 Board of the Mathematical Sciences, Washington, DC; by the American Mathematical Society,  
486 Providence, RI, 2012.
- 487 [WKK20] Vikram Waradpande, Daniel Kudenko, and Megha Khosla, *Deep reinforcement learning with  
488 graph-based state representations*, arXiv:2004.13965, 2020.
- 489 [ZZ22] Yang Zhao and Hao Zhang, *Quantitative performance assessment of CNN units via topological  
490 entropy calculation*, International Conference on Learning Representations, 2022.

491 **Checklist**

- 492 1. For all authors...
- 493 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
494 contributions and scope? [Yes]
- 495 (b) Did you describe the limitations of your work? [Yes]
- 496 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 497 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
498 them? [Yes]
- 499 2. If you are including theoretical results...
- 500 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 501 (b) Did you include complete proofs of all theoretical results? [Yes] see Appendix A.2
- 502 3. If you ran experiments...
- 503 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
504 mental results (either in the supplemental material or as a URL)? [Yes] and we will  
505 provide a link to the publicly-hosted code for community usage upon acceptance (see  
506 Appendix A.3)
- 507 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
508 were chosen)? [N/A]
- 509 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
510 ments multiple times)? [N/A]
- 511 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
512 of GPUs, internal cluster, or cloud provider)? [Yes] see Appendix A.3
- 513 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 514 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 515 (b) Did you mention the license of the assets? [N/A]
- 516 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 517
- 518 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
519 using/curating? [N/A]
- 520 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
521 information or offensive content? [N/A]
- 522 5. If you used crowdsourcing or conducted research with human subjects...
- 523 (a) Did you include the full text of instructions given to participants and screenshots, if  
524 applicable? [N/A]
- 525 (b) Did you describe any potential participant risks, with links to Institutional Review  
526 Board (IRB) approvals, if applicable? [N/A]
- 527 (c) Did you include the estimated hourly wage paid to participants and the total amount  
528 spent on participant compensation? [N/A]