

---

# Manifold Steering Reveals the Shared Geometry of Neural Network Representation and Behavior

---

Anonymous Author(s)

Affiliation

Address

email

1 Neural representations carry rich geometric structure; but does that structure  
2 causally shape behavior? To address this question, we intervene along paths  
3 through activation space defined by different geometries, and measure the behav-  
4 ioral trajectories they induce. In particular, we test whether interventions that  
5 respect the geometry of activation space will yield behaviors close to those the  
6 model exhibits naturally. Concretely, we first fit an activation manifold  $\mathcal{M}_h$  to  
7 representations and a behavior manifold  $\mathcal{M}_y$  to output probability distributions.  
8 We then test the link  $\mathcal{M}_h \leftrightarrow \mathcal{M}_y$  via interventions: we find that steering along  
9  $\mathcal{M}_h$ , which we term *manifold steering*, yields behavioral trajectories that follow  
10  $\mathcal{M}_y$ , while linear steering—which assumes a Euclidean geometry—cuts through  
11 off-manifold regions and hence produces unnatural outputs. Moreover, optimizing  
12 interventions in activation space to produce paths along  $\mathcal{M}_y$  recovers activation  
13 trajectories that trace the curvature of  $\mathcal{M}_h$ . We demonstrate this bidirectional  
14 relationship between the geometry of representation and behavior across tasks  
15 and modalities. In language models, we use reasoning tasks with cyclic and se-  
16 quential geometries as well as in-context learning tasks with more complex graph  
17 geometries. In a video world model, we use a task with geometry corresponding to  
18 physical dynamics. Overall, our work shows that geometry in neural representation  
19 is not merely incidental, but is in fact the proper object for enabling principled  
20 control via intervention on internals. This recasts the core problem of steering from  
21 finding the right *direction* to finding the right *geometry*.

## 22 1 Introduction

23 A plethora of geometric structures have been documented in neural network representations (Modell  
24 et al., 2025b; Park et al., 2025b; Kozłowski et al., 2025; Shai et al., 2024b; Pearce et al., 2025; Gurnee  
25 et al., 2026). Recent literature has begun to identify the origins of these structures by attributing  
26 them back to data statistics shaped by conceptual structure (Karkada et al., 2026; Prieto et al., 2026;  
27 Park et al., 2025b; Merullo et al., 2025). However, we have barely begun to understand what causal  
28 role these geometric structures play in a model’s computation (cf. Engels et al. 2024; Kantamneni &  
29 Tegmark 2025; Csordás et al. 2024; Sarfati et al. 2026). We address this question by intervening on  
30 model activations under different geometric assumptions and measuring the effect on behavior.

31 Currently, it is common for activation-based intervention methods to assume a Euclidean geometry  
32 for activation space, where steering is performed by adding a *steering vector* to model activations  
33 with a scalar that modulates intervention strength (Bau et al., 2019; Subramani et al., 2022; Marks &  
34 Tegmark, 2024; Panickssery et al., 2024; Turner et al., 2024; Li et al., 2023; Rinsky et al., 2024; Chen  
35 et al., 2025). This approach is motivated by the *linear representation hypothesis* (LRH), which posits  
36 that neural activations can be decomposed into atomic concepts encoded along single (approximately)  
37 orthogonal directions (Smolensky, 1986; Park et al., 2023; Elhage et al., 2022b). However, linear  
38 steering often produces degraded fluency, diversity collapse, and unstable off-target behavior (Wu  
39 et al., 2025; Da Silva et al., 2025; Bigelow et al., 2025; Tan et al., 2024; Hao et al., 2025; Bhalla et al.,  
40 2024; Pres et al., 2024), which suggests the assumed Euclidean geometry is inappropriate.

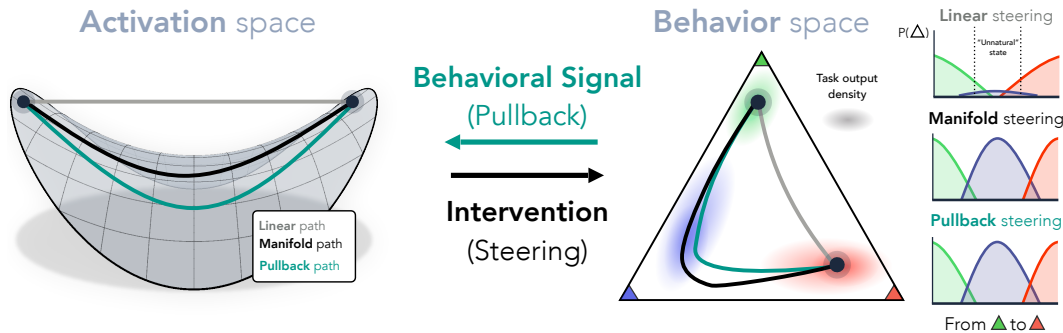


Figure 1: **How do different geometries of activation space modulate behavior?** We illustrate paths through activation space (left), each defined by a different geometry. Interventions along paths in activation space induce paths in behavior space (right, illustrated on a three-concept probability simplex). **Euclidean:** the standard approach of linear steering assumes a flat geometry and interventions follow a straight line. Such paths may cut across the activation manifold, yielding *unnatural* behavioral trajectories that pass through off-manifold regions of behavior space. **Density geometry:** a density-based metric whose geodesics follow the intrinsic geometry of a fitted activation manifold, yielding more natural transitions in behavior space. **Pullback geometry:** a behavior-aware metric obtained by “pulling back” behavior-space geometry into activation space, yielding paths that follow the manifold of natural (unintervened) output distributions. Overall, we argue that geometric structure in neural representations encodes the conceptual space a model is reasoning over, which in turn constrains its output behavior. Hence, manifolds in activation and behavior space are two images of the same underlying structure, and so we expect the density and pullback geometries to coincide.

41 In this work, we advance the hypothesis that *representation geometry* provides a blueprint for effective  
 42 steering that will overcome the limitations of the linear approach. Steering is fundamentally about  
 43 how internal representations control behavior, so to test this hypothesis we must study not only  
 44 paths through activation space, but also the behavioral trajectories induced by interventions along  
 45 these paths. Successful steering will produce trajectories that are in line with the model’s *natural*  
 46 (unintervened) output distribution. If we are right, then interventions that respect the geometry of  
 47 internal representations and interventions that respect the geometry of behavior will be one and the  
 48 same. Motivated by this, we make the following contributions in this work.

- 49 • **Uncovering isometric geometries in neural network representation and behavior.** We use tasks  
 50 where models output a distribution over a set of concepts with known structure. In each task, we  
 51 fit an *activation manifold*  $\mathcal{M}_h$  to internal representations and a *behavior manifold*  $\mathcal{M}_y$  to model  
 52 outputs (probability distributions over task-relevant concepts). We show the two geometries are  
 53 tightly interlinked via a scaled isometry relation: geodesic distances on  $\mathcal{M}_h$  align closely with  
 54 those on  $\mathcal{M}_y$ , and neither match Euclidean distances.
- 55 • **Validating the causal role of representation geometry.** We perform geometry-aware steering  
 56 experiments and compare against the baseline of linear steering (see Fig. 1). We show linear steering  
 57 cuts through low-density regions of behavior space and passes through unnatural intermediate  
 58 distributions; meanwhile, steering along the activation manifold  $\mathcal{M}_h$  yields behavioral trajectories  
 59 that follow  $\mathcal{M}_y$  closely. In fact, optimizing for paths along  $\mathcal{M}_y$  recovers activation trajectories that  
 60 trace the curvature of  $\mathcal{M}_h$ , further tightening the link between activation geometry and behavior.
- 61 • **A theoretical framework for geometry-aware steering.** Building on the results above, we  
 62 formulate steering as a problem of choosing the right *geometry* for activation space, rather than the  
 63 right *direction*. In particular, we argue steering can be defined as the problem of finding a geodesic  
 64 connecting two points under different activation-space metrics: linear steering assumes a flat metric  
 65 (Euclidean geometry), steering along the activation manifold uses a metric derived from natural  
 66 activations, and steering optimized to follow the behavior manifold uses a metric derived from  
 67 natural behaviors.

68 We demonstrate these findings hold across modalities and tasks. In large language models, we  
 69 test geometries from cyclic concepts (weekdays, months; Engels et al. 2024; Modell et al. 2025b),  
 70 sequential concepts (ages, letters), and multi-dimensional graph structures learned in context (Park  
 71 et al., 2025b). In a video world model, we test a geometry of physical position in a simulated

72 environment (mountain car; Moore 1990; Towers et al. 2024). Together, these findings provide  
73 evidence for the posited account, and support steering along neural manifolds as the principled form  
74 of activation-based intervention.

## 75 2 The Geometry of Representation and Behavior

### 76 2.1 Setup

77 **Running example.** We will explicate our framework and empirical methods using a running  
78 example where a language model is required to reason about the days of the week (Engels et al.,  
79 2024). Specifically, we consider prompts of the form: What day is  $k$  days after  $z$ ? with  
80  $z \in \mathcal{Z} = \{\text{Mon, Tue, } \dots, \text{Sun}\}$  and  $k \in \{1, \dots, 7\}$ . Given such a prompt, the LM outputs a  
81 probability distribution over all possible tokens.

82 **Concept geometry.** We draw inspiration from work on *conceptual spaces* in cognitive science,  
83 where conceptual domains, e.g., days of the week, are geometrically enriched with a metric such that  
84 distances between points encode similarity and guide patterns of inference (Shepard 1987; Gärdenfors  
85 2000; Tenenbaum & Griffiths 2001; Bellmund et al. 2018; see Fel et al. 2025b; Lubana et al. 2025;  
86 Yocum et al. 2025; Modell et al. 2025b for related work in interpretability). For example, the days of  
87 the week  $\mathcal{Z}$  may be organized in a cyclic structure that is captured by a *metric*  $d_{\mathcal{Z}}$  measuring temporal  
88 distance between days, e.g., neighboring days are closer together. Indeed, when humans mistakenly  
89 report the current day, they most often confuse it with its neighboring days (Ellis et al., 2015).

90 Karkada et al. (2026) and Prieto et al. (2026) show that similarity structure between days of the  
91 week is reflected in the statistics of training data, which in turn shape the geometry of internal  
92 representations (Engels et al., 2024; Park et al., 2025a; Modell et al., 2025a; Prieto et al., 2026). We  
93 hypothesize that a model’s output distributions over  $\mathcal{Z}$  are similarly shaped by  $d_{\mathcal{Z}}$ : e.g., when asked  
94 What day is four days after Monday?, the model concentrates mass on Friday and spreads  
95 the remainder onto nearby days like Thursday and Saturday.

96 **Notation.** We work with two spaces: the activation space  $\mathcal{A} = \mathbb{R}^n$ , and the behavior space  
97  $\mathcal{Y} = \Delta^{|\mathcal{Z}|}$ , which is the open probability simplex<sup>1</sup> over the conceptual domain  $\mathcal{Z}$ , with an additional  
98 ‘other’ class for off-concept probability mass. For an input  $x$ , let  $\mathbf{p}(x) \in \mathcal{Y}$  denote the model’s output  
99 distribution over  $\mathcal{Z}$ , given by restricting the full vocabulary distribution of the model to the tokens in  
100  $\mathcal{Z}$ , in addition to the ‘other’ class for remaining probability mass. Let  $\mathbf{h}(x) \in \mathcal{A}$  denote an activation  
101 vector of interest for input  $x$ .

102 For a class of input queries that share the same answer, e.g., What is two days after Monday?  
103 and What is three days after Sunday?, we average the hidden activations and output distri-  
104 butions to produce “activation centroids” and “behavior centroids”, respectively.

105 **Experimental tasks.** We perform language model experiments on four tasks, two with cyclic  
106 conceptual structure and two with sequential conceptual structure. The cyclic tasks require reasoning  
107 about days of the week and months of the year, e.g., What is four months after January?.  
108 The sequential tasks require reasoning about letters and ages, e.g., What is four letters after  
109 m? or Alice is 7, Bob is 5 years older. How old is Bob?.

### 110 2.2 Fitting the Manifolds

111 We fit a smooth manifold within each space to the model’s uninvolved activations or outputs for a  
112 task:  $\mathcal{M}_h \subseteq \mathcal{A}$ , the *activation manifold*, and  $\mathcal{M}_y \subseteq \mathcal{Y}$ , the *behavior manifold*. To fit the activation  
113 manifold  $\mathcal{M}_h$ , we reduce activation vectors  $\mathbf{h}(x)$  to 64 dimensions via PCA, compute “concept  
114 centroids” (e.g., averaging all activations where the correct answer is Wednesday), and fit cubic  
115 splines (Reinsch, 1967) through the centroids (see App. F.3 for further spline fitting details). To fit the  
116 behavior manifold  $\mathcal{M}_y$ , we follow a similar procedure but first map each centroid from the probability  
117 simplex onto Hellinger space via  $p \mapsto \sqrt{p}$ . This linearizes the geometry of the simplex: the Hellinger

<sup>1</sup>We require the open simplex  $\{\mathbf{p} \in \mathbb{R}_{>0}^{|\mathcal{Z}|} : \sum_i p_i = 1\}$ , i.e., strictly positive entries. The closed simplex, which includes faces where some  $p_i = 0$ , has boundary and corners and is not a smooth manifold.

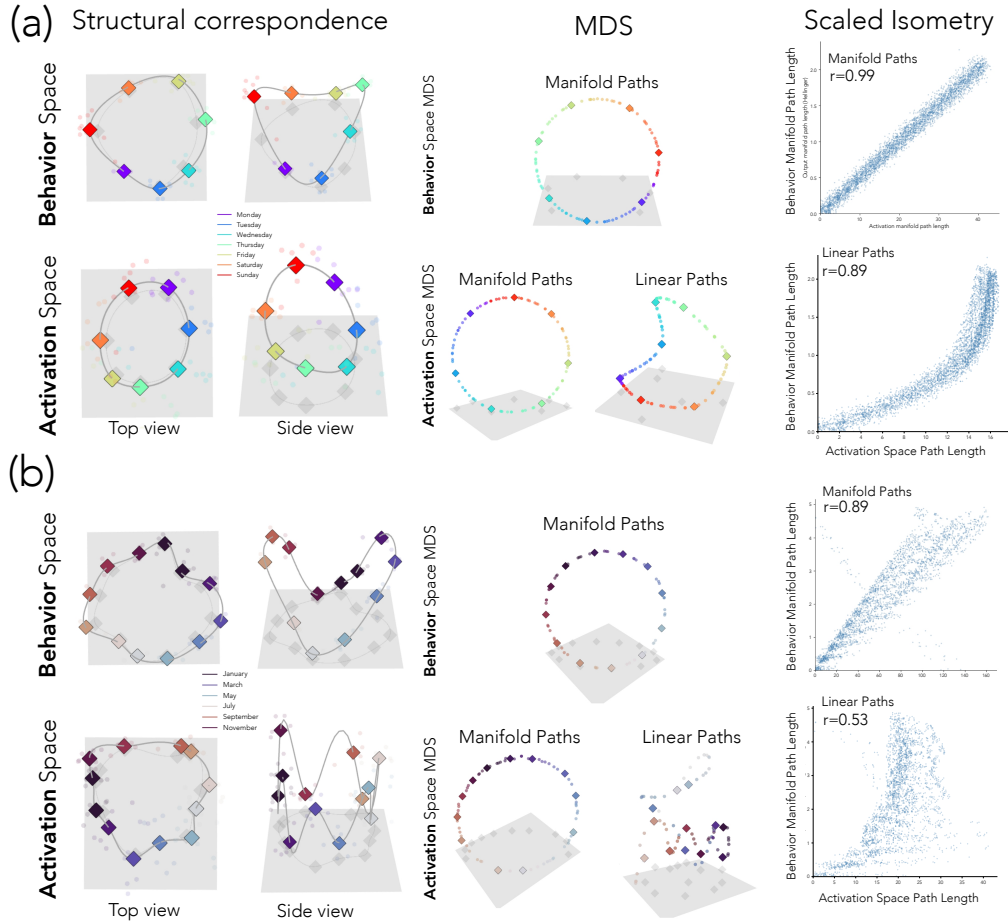


Figure 2: **Approximate isometry between activation and behavior manifolds for cyclic concepts.** Manifolds (cubic splines) fit to activation and behavior (i.e., output distributions over concept tokens) spaces of Llama 3.1 8B. The weekdays (a) and months (b) tasks consist of simple addition questions such as: What is four days after Monday?. Both activation and behavior manifolds show cyclic structure (PCA visualization shown in left column). Furthermore, on-manifold distances in activation space show strong correlation with on-manifold distances in behavior space (right column), as well as a clear structural match via a multidimensional scaling (MDS) embedding (middle column). In contrast, linear distances in activation space show weaker correlations and warped structures. These results demonstrate an approximate isometry between the activation and behavior space manifolds. See more domains in Figure 8.

118 distance between distributions becomes an ordinary Euclidean distance,  $d_H(p, q) = \frac{1}{\sqrt{2}} \|\sqrt{p} - \sqrt{q}\|$ ,  
 119 so we can fit splines and compare distributions with standard Euclidean tools while still respecting  
 120 the underlying probabilistic geometry (Amari & Nagaoka, 2000). Decoded points are squared back  
 121 to recover valid distributions (further details, including how we keep the fit on the sphere, are in  
 122 App. F.4). Unless stated otherwise, we use Llama 3.1 8B (Touvron et al., 2023) with activations from  
 123 layer 28, and visualize manifolds via 3D PCA.

### 124 2.3 Conceptual Structure Appears in Behavior and Activation Space

125 We will now begin our investigation into the connection between the activation and behavior mani-  
 126 folds. Before we perform any interventions on internal representations, we examine the structural  
 127 correspondence between these spaces, and measure whether distances along the two manifolds are  
 128 proportional (a scaled isometry).

129 We find that both the activation and behavior manifolds recapitulate conceptual structure (see Figs. 2, 8  
 130 for visualizations). For example, in the case of days of the week, the activations and output distri-

131 butions are arranged in order around a loop, with Monday adjacent to Tuesday and Sunday, and  
 132 Thursday on the opposite side (Fig. 2). The circle representing days of the week in activation space is  
 133 already known to exist (Engels et al., 2024; Modell et al., 2025b; Karkada et al., 2026). However, the  
 134 circle in behavior space is a novel discovery, and results from sharply-peaked output distributions  
 135 placing most mass on the target concept, with the remainder concentrated on its neighbors. The  
 136 correspondence is striking; output distributions and internal activations recover the same cyclic  
 137 ordering. In contrast, the conceptual structure for the ages and letters task is sequential rather than  
 138 cyclic, and so both the activations and output distributions for these tasks lie on an open curve (Fig. 8).

139 Going beyond qualitative structural correspondence, we wish to examine the mapping between  
 140 distances along each manifold. We test this by computing pairwise distances between points in  
 141 both spaces: geodesic distances  $d_{\mathcal{M}_h}(m_i, m_j)$  on the activation manifold, and geodesic distances  
 142  $d_{\mathcal{M}_y}(p_i, p_j)$  on the behavior manifold. We compute geodesic distance on  $\mathcal{M}_h$  using cumulative  
 143 Euclidean distance between points along a geodesic path, and follow the same procedure for geodesic  
 144 distance on  $\mathcal{M}_y$ , but using cumulative Hellinger distance (see App. F.5 for further details). The two  
 145 distances are highly correlated ( $r = 0.99$  weekdays,  $r = 0.89$  months,  $r = .999$  letters,  $r = .999$   
 146 ages) indicating that  $\mathcal{M}_h$  and  $\mathcal{M}_y$  are approximately isometric. Meanwhile, linear paths between the  
 147 same activation-space points correlate less well with  $\mathcal{M}_y$  geodesics, with the relationship showing  
 148 clear non-linear patterns ( $r = 0.89$  weekdays,  $r = 0.53$  months,  $r = 0.71$  letters,  $r = 0.36$  ages).

149 This correspondence leads to an intuitive hypothesis.  $\mathcal{M}_y$  was fit to unintervened task behavior, so it  
 150 traces a path through *natural* output distributions for the model. If the  $\mathcal{M}_h \leftrightarrow \mathcal{M}_y$  mapping holds,  
 151 paths along one manifold should track paths along the other. Interventions in activation space that  
 152 follow  $\mathcal{M}_h$  should produce natural trajectories along  $\mathcal{M}_y$ . Conversely, activation-space paths that  
 153 are optimized to produce trajectories on  $\mathcal{M}_y$  should recover  $\mathcal{M}_h$ . Next, we test both directions via  
 154 intervention: representation to behavior ( $\mathcal{M}_h \rightarrow \mathcal{M}_y$ ) and behavior to representation ( $\mathcal{M}_h \leftarrow \mathcal{M}_y$ ).

### 155 3 Connecting Representation and Behavior via Intervention

156 Now that we have established a correlational correspondence between the activation and behavior  
 157 manifolds across four tasks, we turn to steering interventions for causal evidence. First, we steer along  
 158  $\mathcal{M}_h$  and measure whether output trajectories follow  $\mathcal{M}_y$  (§3.2). Second, we optimize interventions  
 159 on internal representations to produce output distributions that follow  $\mathcal{M}_y$  and measure whether the  
 160 optimized activation trajectory follows  $\mathcal{M}_h$  (§3.3).

#### 161 3.1 Steering Intervention Notation

162 The basic intervention operation entails replacing the model’s activation at a chosen layer with a  
 163 target activation, and continuing the forward pass. Given a base input  $x$  and a target  $\mathbf{h}^* \in \mathcal{A}$ , we write  
 164  $\mathbf{p}_{\mathbf{h} \leftarrow \mathbf{h}^*}(x)$  for the resulting output distribution. A *steering path* is a curve  $\pi : [0, 1] \rightarrow \mathcal{A}$  between  
 165 endpoints  $\mathbf{h}_0^*$  and  $\mathbf{h}_1^*$ , inducing a trajectory  $\mathbf{p}_{\mathbf{h} \leftarrow \pi(t)}(x)$  through behavior space  $\mathcal{Y}$ . The behavioral  
 166 trajectory will be non-stationary only if the target  $\mathbf{h}$  mediates the causal effect from input to output  
 167 (Pearl, 2001; Vig et al., 2020; Mueller et al., 2024).

168 We consider two strategies, both constructed by interpolation between the endpoints; the strategies  
 169 differ only in the coordinate system in which the interpolation is taken (Fig. 1):

$$\pi_{\text{lin}}(t) = (1-t)\mathbf{h}_0^* + t\mathbf{h}_1^* \quad (\text{linear steering}); \quad (1)$$

$$\pi_{\text{m}}(t) = \mathbf{s}((1-t)\mathbf{u}_0 + t\mathbf{u}_1), \quad \mathbf{u}_i = \mathbf{s}^{-1}(\mathbf{h}_i^*) \quad (\text{manifold steering}). \quad (2)$$

170 In the above,  $\mathbf{s} : \mathbb{R}^k \rightarrow \mathcal{A}$  is a *parameterization* of  $\mathcal{M}_h$ —the map sending  $k$ -dimensional intrinsic  
 171 coordinates to the corresponding point on the manifold in the activation space  $\mathcal{A}$ . Linear steering  
 172 (also known as ‘diff-in-means steering’) (Bau et al., 2018; Subramani et al., 2022; Turner et al., 2023)  
 173 interpolates in  $\mathcal{A}$  directly—the standard additive-vector baseline. Manifold steering interpolates in  
 174 the intrinsic coordinates of  $\mathcal{M}_h$  and maps the result back through  $\mathbf{s}$ , so  $\pi_{\text{m}}$  stays on the activation  
 175 manifold  $\mathcal{M}_h$  throughout. Each strategy thus corresponds to a different choice of geometry on  
 176 activation space, which we concretize in §3.4.

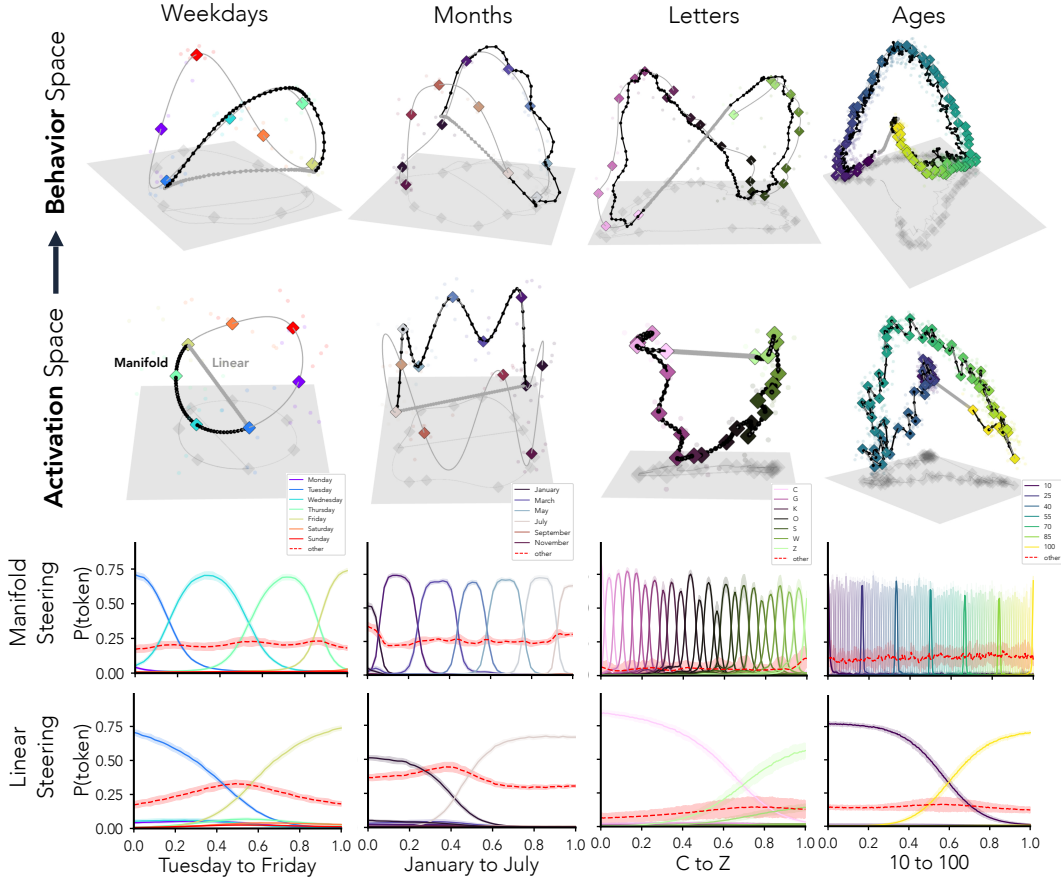


Figure 3: **Manifold steering yields smooth and ordered behavioral transitions.** Using simple addition tasks which require reasoning over structured concepts (e.g., What is four days after Monday?), we compare two steering strategies in activation space: standard linear steering, which takes direct paths, and *manifold steering*, which takes paths along a fitted activation manifold. The bottom panel shows example output paths given by each method. Across four settings, manifold steering produces smooth and ordered output transitions between adjacent concepts. In contrast, linear steering leads to ‘teleportation’ of probability between non-adjacent concepts, and at times results in probability on non-related tokens surpassing any individual concept near the path midpoint. The top panel shows output trajectories in behavior space resulting from manifold steering. We find that steering along the activation-space manifold yields paths that follow the behavior manifold, while linear steering traces paths far from the manifold. Thus, the outputs produced under manifold steering more closely resemble *natural* outputs produced without intervention.

### 177 3.2 Steering Along the Activation Manifold Follows the Behavior Manifold

178 For every pair of start and end values, e.g., from Tuesday to Friday, we steer from the start centroid  
 179 to the end centroid in activation space using manifold and linear steering with  $K = 50$  intervention  
 180 points along each path. We report the average trajectory in behavior space over a set of 16 prompts  
 181 sampled randomly from the task’s input distribution (Fig. 3, see App. F.6 for further experimental de-  
 182 tails). We find that manifold steering produces *smooth* and *ordered* behavioral transitions: probability  
 183 mass shifts steadily through adjacent values of a concept—from Monday to Tuesday to Wednesday to  
 184 Thursday—while linear steering instead exhibits ‘teleportation’: mass jumps between non-adjacent  
 185 concepts as the straight line cuts through the manifold’s interior.

186 This qualitative evidence is encouraging, but we have yet to examine our key hypothesis that  
 187 interventions along the activation manifold  $\mathcal{M}_h$  produce *natural* output trajectories that follow  $\mathcal{M}_y$ .  
 188 This would mean that outputs produced under manifold steering resemble those produced without

189 intervention. We quantify this via an “energy function”, as described next, under which a natural  
 190 trajectory is one of low cumulative energy as defined by  $\mathcal{M}_y$ .

191 **An Energy-based View of Naturalness.** Energy functions have a long history in machine learning  
 192 as a way to measure plausibility under a model (Hopfield, 1982; LeCun et al., 2006). These functions  
 193 assign low values for likely states and high values for unlikely ones, with the standard correspondence  
 194  $E(\mathbf{x}) \propto -\log p(\mathbf{x})$  giving energy the interpretation of an unnormalized log-density (Hopfield, 1982;  
 195 LeCun et al., 2006; Grathwohl et al., 2019; Song & Kingma, 2021; Béthune et al., 2025). We  
 196 adopt the same view here. The model’s output distributions on unintervened forward passes trace  
 197 out a low-energy region of behavior space (approximately captured by the manifold  $\mathcal{M}_y$ ) and a  
 198 steering trajectory is natural to the extent it stays within that region. Concretely, given a steering  
 199 path  $\pi : [0, 1] \rightarrow \mathcal{A}$ , let  $\gamma(t) = \mathbf{p}_{\mathbf{h} \leftarrow \pi(t)}(\mathbf{x})$  be the behavioral trajectory it induces. We define its  
 200 cumulative output energy:

$$E_{\text{BC}}(\gamma) = \int_0^1 d_{\text{BC}}(\gamma(t), \mathcal{M}_y) dt, \quad (3)$$

201 where  $d_{\text{BC}}(\mathbf{p}, \mathcal{M}_y) = \inf_{\sqrt{\mathbf{q}} \in \mathcal{M}_y} d_{\text{BC}}(\mathbf{p}, \mathbf{q}) = -\log(\sum_i \sqrt{\mathbf{p}_i} \sqrt{\mathbf{q}_i})$  is the Bhattacharyya distance to  
 202 the nearest point on  $\mathcal{M}_y$ , a natural choice given that it is simply the negative log of the dot product  
 203 in Hellinger space (in which  $\mathcal{M}_y$  is fit). We note that the formulation above is a tractable proxy we  
 204 use to estimate distance from the model’s natural output distribution, yet it is but one instantiation  
 205 of a more general framework we develop in §3.4. Applying this measure, we find that manifold  
 206 steering (weekdays  $E_{\text{BC}} = 0.34 \pm 0.03$ ; months  $E_{\text{BC}} = 0.36 \pm 0.01$ ; letters  $2.42 \pm 0.07$ ; ages  
 207  $E_{\text{BC}} = 5.21 \pm 0.09$ ) produces significantly more natural paths, i.e., lower cumulative energy, than  
 208 linear steering (weekdays  $E_{\text{BC}} = 0.93 \pm 0.11$ ; months  $E_{\text{BC}} = 1.09 \pm 0.06$ ; letters  $6.95 \pm 0.27$ ; ages  
 209  $E_{\text{BC}} = 13.49 \pm 0.29$ ); on average, we see an improvement of a factor of  $2.8\times$ , with all statistical  
 210 comparisons yielding  $p < 0.001$ .

211 We find further verification of the claim above by visualizing output trajectories in behavior space  
 212 (Fig. 3). Manifold steering consistently traces paths close to  $\mathcal{M}_y$ , while linear steering cuts through  
 213 regions far from the behavior manifold, yielding less natural outputs. This result provides causal  
 214 support for the correspondence between activation and output geometry, and establishes manifold  
 215 steering as a principled form of steering that yields natural behavioral trajectories ( $\mathcal{M}_h \rightarrow \mathcal{M}_y$ ).  
 216 Next, we explore whether we can find evidence from the opposite direction ( $\mathcal{M}_y \rightarrow \mathcal{M}_h$ ).

### 217 3.3 Behavior Space Geometry Recovers the Activation Manifold

218 In this section, we aim to uncover whether steering intervention paths optimized to follow the behavior  
 219 manifold  $\mathcal{M}_y$  recover the activation manifold  $\mathcal{M}_h$ . To do so, we first take a path  $\pi_y^*$  along  $\mathcal{M}_y$  in  
 220 behavior space. We work within the layer and the first 32 dimensions of the subspace in which we  
 221 fit the activation manifold (64 dimensional PCA), and optimize via L-BFGS for a path in activation  
 222 space which, upon intervention, induces the behavioral path  $\pi_y^*$  (see App. F.8 for further details  
 223 regarding the optimization procedure). We call the resulting path in activation space the *pullback*.

224 To quantify how faithfully a pullback path in activation space  $\pi_h^{\text{pullback}}$  recapitulates the manifold  
 225 steering path  $\pi_h^*$ , we report an *intrinsic*  $R^2$ . Both paths are projected into a common subspace  
 226 given by the singular directions explaining 99% of the variance in  $\pi_h^*$ , restricting the comparison  
 227 to directions where the path actually extends. We then compute the  $R^2$  in this subspace, defining  
 228 the residual at each point of  $\pi_h^{\text{pullback}}$  as its orthogonal closest-point distance to  $\pi_h^*$ . We compute this  
 229 score for each optimized pullback path  $\pi_h^{\text{pullback}}$  and compare with a linear path baseline (see App. F.9  
 230 for more details).

231 Results are shown in Fig. 4. We find that the pullback activation paths follow the activation manifold  
 232  $\mathcal{M}_h$  more closely than the linear steering path, and resemble the shape of the manifold steering  
 233 paths of §3.2 (weekdays  $R_{\text{pullback}}^2 = 0.77 \pm 0.03$  vs.  $R_{\text{linear}}^2 = 0.42 \pm 0.07$ ; months  $R_{\text{pullback}}^2 =$   
 234  $0.75 \pm 0.04$  vs.  $R_{\text{linear}}^2 = 0.32 \pm 0.05$ ; ages  $R_{\text{pullback}}^2 = 0.47 \pm 0.05$  vs.  $R_{\text{linear}}^2 = 0.24 \pm 0.01$ ; letters  
 235  $R_{\text{pullback}}^2 = 0.78 \pm 0.04$  vs.  $R_{\text{linear}}^2 = 0.23 \pm 0.03$ . All statistical comparisons yield  $p < 0.001$ ). Again  
 236 we see a striking correspondence between representation and behavior; despite being derived from  
 237 different sources—manifold steering from the density of activations and pullback from the structure  
 238 of outputs—the two geometries are tightly connected.

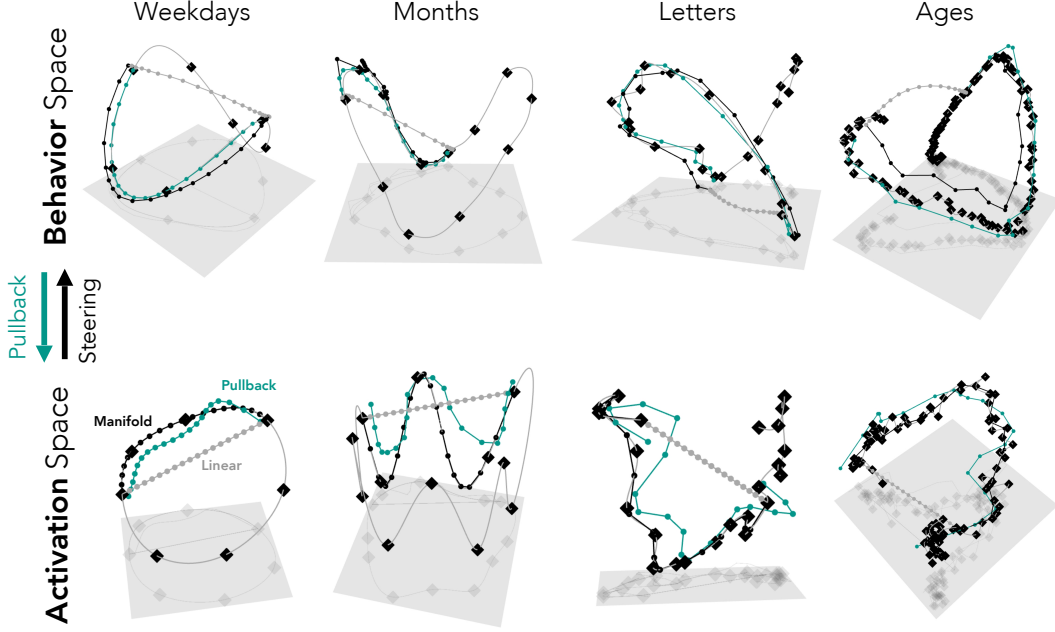


Figure 4: **Manifold steering and pullback yield coinciding trajectories in activation and behavior space.** Going in the **Activations**→**Behavior** direction, we find that steering along the activation manifold  $\mathcal{M}_h$  (black) produces paths that lie close to the behavior manifold  $\mathcal{M}_y$ . We then examine the reverse direction, **Activations**←**Behavior**: We start with paths along the behavior manifold and optimize for corresponding paths in activation space (i.e., a set of activations that yields the path on  $\mathcal{M}_y$  upon intervention). This *pullback* procedure (teal) yields trajectories that resemble the activation manifold  $\mathcal{M}_h$ . Thus, we offer bidirectional support for the connection between activation geometry and behavior, and their correspondence reflecting a shared underlying conceptual organization. Paths shown: weekdays ‘Thursday’ to ‘Sunday’; Months ‘August’ to ‘December’; Letters ‘C’ to ‘Q’; Ages 36 to 91.

239 Taken together, these results, alongside those of §3.2, provide bidirectional support for the connection  
 240 between activation geometry and behavior. This convergence indicates that  $\mathcal{M}_h$  is a core object in the  
 241 model’s representation: the geometry of activation space and the geometry of behavior are alternate  
 242 views of the same underlying conceptual organization.

### 243 3.4 Unifying Steering Strategies Through Geometry

244 In the sections above, we analyzed three methods for steering between two points in activation space  
 245 that each assume a different geometry: linear steering, which assumes a flat Euclidean geometry;  
 246 manifold steering, which derives a geometry from naturally occurring activations; and pullback  
 247 steering, which derives a geometry from naturally occurring output distributions. We provided  
 248 empirical support for our hypothesis that the geometries derived from internal activations and output  
 249 behaviors are much more similar to each other than the standard Euclidean geometry. We now  
 250 formalize the question of how to steer as *how to choose the right geometry for activation space*.

251 **The Geometry of Steering:** Consider a Riemannian metric  $\mathbf{G}$ , which assigns an inner product at  
 252 each point of  $\mathcal{A}$ ; together with a path  $\pi : [0, 1] \rightarrow \mathcal{A}$ , this defines the notion of path length as follows.

$$L_{\mathbf{G}}(\pi) = \int_0^1 \sqrt{\dot{\pi}(t)^\top \mathbf{G}(\pi(t)) \dot{\pi}(t)} dt. \quad (4)$$

253 Then, a geodesic is defined as the path of minimum length between two endpoints, and each choice of  
 254 geometry picks out a steering strategy. The strategies of linear steering and manifold steering (§3.2),  
 255 written as interpolations in two different coordinate systems (Eqs. 1, 2), are two such choices; the  
 256 pullback procedure of §3.3 is a third. Now, we make all three geometries explicit.

257 **Definition 1** (Geometries of Steering). Let  $E : \mathcal{A} \rightarrow \mathbb{R}$  be an energy function such that  $E(\mathbf{h}) \propto$   
 258  $-\log p(\mathbf{h})$ , and let  $\mathbf{g}_y$  be a chosen Riemannian metric on  $\mathcal{M}_y$ . We define:

$$\mathbf{G}_I = \mathbf{I}_n, \quad (\text{linear steering}) \quad (5)$$

$$\mathbf{G}_E(\mathbf{h}) = (\alpha e^{-E(\mathbf{h})} + \beta)^{-1} \mathbf{I}_n, \quad (\text{manifold steering}) \quad (6)$$

$$\mathbf{G}_F(\mathbf{h}) = \mathbf{J}_F(\mathbf{h})^\top \mathbf{g}_y(\mathbf{F}(\mathbf{h})) \mathbf{J}_F(\mathbf{h}) + \epsilon \mathbf{I}_n, \quad (\text{pullback}) \quad (7)$$

259 where  $\alpha, \beta > 0$  are calibration constants,  $\epsilon > 0$  regularizes the pullback,  $\mathbf{F} : \mathcal{A} \rightarrow \mathcal{Y}$  is the function  
 260 from naturally occurring activations to naturally occurring behaviors, and  $\mathbf{g}_y$  is any Riemannian  
 261 metric on  $\mathcal{M}_y$  (e.g., the induced Hellinger metric used in our experiments).

262 We discuss the intuitive interpretation of Defn. 1 below.

263 • **The Flat Geometry  $\mathbf{G}_I$ .** Linear steering treats activation space as Euclidean: all directions and  
 264 regions are equally valid, with Geodesics as straight lines  $\ell(t) = (1-t) \mathbf{h}_0 + t \mathbf{h}_1$ . This geometry  
 265 thus encodes no knowledge of naturally occurring activation or outputs.

266 • **The Density Geometry  $\mathbf{G}_E$ .** Manifold steering derives a geometry for activation space from  
 267 naturally occurring internal representations. Specifically, consider the geometry induced from an  
 268 energy function  $E(\mathbf{h}) \propto -\log p(\mathbf{h})$  by rescaling the identity according to local density. Here  
 269  $e^{-E(\mathbf{h})}$  plays the role of an unnormalized density: large where activations concentrate (on  $\mathcal{M}_h$ )  
 270 and small where they are sparse (off  $\mathcal{M}_h$ ). The inverse makes off-manifold regions expensive and  
 271 on-manifold movement cheap, with constants  $\alpha, \beta > 0$  calibrating the dynamic range (Béthune  
 272 et al., 2025). Geodesics under  $\mathbf{G}_E$  thus follow  $\mathcal{M}_h$ , recovering manifold steering.

273 • **The Pullback Geometry  $\mathbf{G}_F$ .** The steering path given by pullback derives geometric structure  
 274 from naturally occurring model outputs. Specifically,  $\mathbf{G}_F$  is the pullback of a chosen geometry  
 275 on  $\mathcal{M}_y$  through the Jacobian of the map from activation space to behavior space  $\mathbf{F} : \mathcal{A} \rightarrow \mathcal{Y}$ .  
 276 By construction, path length under  $\mathbf{G}_F$  equals path length of the induced behavioral trajectory  
 277 along  $\mathcal{M}_y$  (up to a regularization term). Geodesics under  $\mathbf{G}_F$  are therefore activation paths whose  
 278 induced behavioral trajectories are geodesics on  $\mathcal{M}_y$ —exactly the pullback construction of §3.3.  
 279 The regularization  $\epsilon \mathbf{I}_n$  ensures positive definiteness, since  $\mathbf{J}_F$  has rank at most  $|\mathcal{Z}| - 1 \ll n$ ; as  $\epsilon$   
 280 tends to 0, the geometry approaches the pure pullback in the range of  $\mathbf{J}_F$  and remains Euclidean in  
 281 its null space.

282 Overall, we claim that while the metrics  $\mathbf{G}_E$  and  $\mathbf{G}_F$  are derived from different sources (internal  
 283 activations and outputs, respectively), they converge on approximately the same paths in activation  
 284 space (§3.3). This suggests the manifolds  $\mathcal{M}_h$  and  $\mathcal{M}_y$  are two images of the same conceptual  
 285 geometry, related by an approximate Riemannian isometry. Consequently, the question of optimally  
 286 steering model behavior boils down to isolating the geometry of a concept and defining operators to  
 287 navigate it.

## 288 4 Discussion

289 **Geometry-aware steering reveals the shared structure of behavior and representation.** We  
 290 build out an empirical phenomenology that relates structure in activation space to the model output  
 291 behavior. First, we show an isometry between representations and behavior manifolds, i.e., distance  
 292 between two points on the activation manifold  $\mathcal{M}_h$  aligns with distance between the distributions  
 293 induced by those points on the behavior manifold  $\mathcal{M}_y$ . Second, we show that steering representations  
 294 along geodesics on  $\mathcal{M}_h$  induces smooth, coherent transitions in behavior that follow geodesics on  
 295  $\mathcal{M}_y$ . Third, we show that optimizing interventions to produce behaviors following geodesics on  $\mathcal{M}_y$   
 296 recover trajectories in activation space that follow  $\mathcal{M}_h$ . Thus, we establish a causal bridge between  
 297 representation and behavior that reveals shared structure reflecting underlying conceptual geometry.

298 Our results also suggest that pathologies of linear steering—brittleness, incoherence, off-target  
 299 effects (Wu et al., 2025; Bigelow et al., 2025; Da Silva et al., 2025; Bhalla et al., 2024; Tan et al.,  
 300 2024)—stem from the mismatch between assumed flat geometry and the true curved geometry of  
 301 representation space, rather than an inherent challenge with representation-based intervention. This  
 302 reframes the challenge of steering from “finding the right direction” to “finding the right geometry”.

## 303 References

- 304 Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191. American  
305 Mathematical Soc., 2000.
- 306 Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel  
307 Nanda. Refusal in language models is mediated by a single direction. *Advances in Neural  
308 Information Processing Systems*, 37:136037–136083, 2024.
- 309 Aryaman Arora, Dan Jurafsky, and Christopher Potts. CausalGym: Benchmarking causal inter-  
310 pretability methods on linguistic tasks. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.),  
311 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume  
312 1: Long Papers)*, pp. 14638–14663, Bangkok, Thailand, 2024. Association for Computational  
313 Linguistics. URL <https://aclanthology.org/2024.acl-long.785>.
- 314 Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure  
315 of word senses, with applications to polysemy. *Transactions of the Association for Computational  
316 Linguistics*, 6:483–495, 2018. doi: 10.1162/tacl.a.00034. URL [https://aclanthology.org/  
317 Q18-1034/](https://aclanthology.org/Q18-1034/).
- 318 Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass.  
319 Identifying and controlling important neurons in neural machine translation, 2018. URL <https://arxiv.org/abs/1811.01157>.
- 321 Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass.  
322 Identifying and controlling important neurons in neural machine translation. In *International  
323 Conference on Learning Representations*, 2019. URL [https://openreview.net/forum?id=  
324 H1z-PsR5KX](https://openreview.net/forum?id=H1z-PsR5KX).
- 325 Sander Beckers and Joseph Halpern. Abstracting causal models. In *AAAI Conference on Artificial  
326 Intelligence*, 2019.
- 327 Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational  
328 Linguistics*, 48(1):207–219, 2022.
- 329 Jacob LS Bellmund, Peter Gärdenfors, Edvard I Moser, and Christian F Doeller. Navigating cognition:  
330 Spatial codes for human thinking. *Science*, 362(6415):eaat6766, 2018.
- 331 Nora Belrose, David Schneider-Joseph, Shauli Ravfogel, Ryan Cotterell, Edward Raff, and  
332 Stella Biderman. LEACE: perfect linear concept erasure in closed form. In Alice Oh,  
333 Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.),  
334 *Advances in Neural Information Processing Systems 36: Annual Conference on Neural In-  
335 formation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December  
336 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
337 d066d21c619d0a78c5b557fa3291a8f4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/d066d21c619d0a78c5b557fa3291a8f4-Abstract-Conference.html).
- 338 Louis Béthune, David Vigouroux, Yilun Du, Rufin VanRullen, Thomas Serre, and Victor Boutin.  
339 Follow the energy, find the path: Riemannian metrics from energy-based models. *ArXiv e-print*,  
340 2025.
- 341 Usha Bhalla, Suraj Srinivas, Asma Ghandeharioun, and Himabindu Lakkaraju. Towards unifying  
342 interpretability and control: Evaluation via intervention. *ArXiv e-print*, 2024.
- 343 Eric Bigelow, Daniel Wurgaft, YingQiao Wang, Noah Goodman, Tomer Ullman, Hidenori Tanaka,  
344 and Ekdeep Singh Lubana. Belief dynamics reveal the dual nature of in-context learning and  
345 activation steering, 2025. URL <https://arxiv.org/abs/2511.00617>.
- 346 Eric J Bigelow, Ekdeep Singh Lubana, Robert P Dick, Hidenori Tanaka, and Tomer D Ullman.  
347 In-context learning dynamics with random binary sequences. *arXiv preprint arXiv:2310.17639*,  
348 2023.

- 349 Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Nick Fusai,  
350 Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Kay Ke, Sergey  
351 Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Shi, James Tanner,  
352 Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky.  $\pi_0$ : A vision-language-action  
353 flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- 354 Tolga Bolukbasi, Adam Pearce, Ann Yuan, Andy Coenen, Emily Reif, Fernanda Viégas, and Martin  
355 Wattenberg. An interpretability illusion for bert. *arXiv preprint arXiv:2104.07143*, 2021.
- 356 F.L. Bookstein. Principal warps: thin-plate splines and the decomposition of deformations. *IEEE*  
357 *Transactions on Pattern Analysis and Machine Intelligence*, 11(6):567–585, 1989. doi: 10.1109/  
358 34.24792.
- 359 Matthew Brand. Charting a manifold. *Advances in neural information processing systems*, 15, 2002.
- 360 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick  
361 Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,  
362 Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina  
363 Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and  
364 Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary  
365 learning. *Transformer Circuits Thread*, 2023.
- 366 Bart Bussmann, Patrick Leask, and Neel Nanda. Batchtopk sparse autoencoders. *ArXiv e-print*, 2024.
- 367 Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah, Michael Petrov, Ludwig Schubert, Chelsea  
368 Voss, Ben Egan, and Swee Kiat Lim. Thread: Circuits. *Distill*, 2020. doi: 10.23915/distill.00024.  
369 <https://distill.pub/2020/circuits>.
- 370 Runjin Chen, Andy Arditi, Henry Sleight, Owain Evans, and Jack Lindsey. Persona vectors: Mon-  
371 itoring and controlling character traits in language models. *arXiv preprint arXiv:2507.21509*,  
372 2025.
- 373 Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger  
374 Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for  
375 statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- 376 Ronald R Coifman and Stéphane Lafon. Diffusion maps. *Applied and computational harmonic*  
377 *analysis*, 21(1):5–30, 2006.
- 378 Valérie Costa, Thomas Fel, Ekdeep Singh Lubana, Bahareh Tolooshams, and Demba Ba. From  
379 flat to hierarchical: Extracting sparse representations with matching pursuit. *Advances in Neural*  
380 *Information Processing Systems (NeurIPS)*, 2025.
- 381 Thomas H Costello, Kellin Pelrine, Matthew Kowal, Antonio A Arechar, Jean-François Godbout,  
382 Adam Gleave, David Rand, and Gordon Pennycook. Large language models can effectively  
383 convince people to believe conspiracies. *arXiv preprint arXiv:2601.05050*, 2026.
- 384 Róbert Csordás, Christopher Potts, Christopher D Manning, and Atticus Geiger. Recurrent neural  
385 networks learn to store and generate sequences using non-linear representations. *arXiv preprint*  
386 *arXiv:2408.10920*, 2024.
- 387 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-  
388 coders find highly interpretable features in language models. *ArXiv e-print*, 2023.
- 389 Patrick Queiroz Da Silva, Hari Sethuraman, Dheeraj Rajagopal, Hannaneh Hajishirzi, and Sachin  
390 Kumar. Steering off course: Reliability challenges in steering language models. *arXiv preprint*  
391 *arXiv:2504.04635*, 2025.
- 392 Xander Davies, Max Nadeau, Nikhil Prakash, Tamar Rott Shaham, and David Bau. Discovering  
393 variable binding circuitry with desiderata. *CoRR*, abs/2307.03637, 2023. doi: 10.48550/ARXIV.  
394 2307.03637. URL <https://doi.org/10.48550/arXiv.2307.03637>.
- 395 Jean Duchon. Splines minimizing rotation-invariant semi-norms in sobolev spaces. In Walter Schempp  
396 and Karl Zeller (eds.), *Constructive Theory of Functions of Several Variables*, pp. 85–100, Berlin,  
397 Heidelberg, 1977. Springer Berlin Heidelberg. ISBN 978-3-540-37496-1.

- 398 Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. Amnesic probing: Behavioral expla-  
399 nation with amnesic counterfactuals. In *Proceedings of the 2020 EMNLP Workshop BlackboxNLP:  
400 Analyzing and Interpreting Neural Networks for NLP*. Association for Computational Linguistics,  
401 November 2020. doi: 10.18653/v1/W18-5426.
- 402 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,  
403 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCand-  
404 lish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of  
405 superposition. *Transformer Circuits Thread*, 2022a.
- 406 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,  
407 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition.  
408 *arXiv preprint arXiv:2209.10652*, 2022b.
- 409 David A Ellis, Richard Wiseman, and Rob Jenkins. Mental representations of weekdays. *PloS one*,  
410 10(8):e0134555, 2015.
- 411 Joshua Engels, Eric J Michaud, Isaac Liao, Wes Gurnee, and Max Tegmark. Not all language model  
412 features are one-dimensionally linear. *arXiv preprint arXiv:2405.14860*, 2024.
- 413 Thomas Fel, Ekdeep Singh Lubana, Jacob S Prince, Matthew Kowal, Victor Boutin, Isabel Papadim-  
414 itriou, Binxu Wang, Martin Wattenberg, Demba Ba, and Talia Konkle. Archetypal sae: Adaptive  
415 and stable dictionary learning for concept extraction in large vision models. *Proceedings of the  
416 International Conference on Machine Learning (ICML)*, 2025a.
- 417 Thomas Fel, Binxu Wang, Michael A Lepori, Matthew Kowal, Andrew Lee, Randall Balestriero,  
418 Sonia Joseph, Ekdeep S Lubana, Talia Konkle, Demba Ba, et al. Into the rabbit hull: From  
419 task-relevant concepts in dino to minkowski geometry. *arXiv preprint arXiv:2510.08638*, 2025b.
- 420 Jiahai Feng and Jacob Steinhardt. How do language models bind entities in context? In *The Twelfth  
421 International Conference on Learning Representations*, 2024. URL [https://openreview.net/  
422 forum?id=zb3b6oK077](https://openreview.net/forum?id=zb3b6oK077).
- 423 Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever,  
424 Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *Proceedings of the  
425 International Conference on Learning Representations (ICLR)*, 2025.
- 426 Atticus Geiger, Kyle Richardson, and Christopher Potts. Neural natural language inference models  
427 partially embed theories of lexical entailment and negation. In Afra Alishahi, Yonatan Belinkov,  
428 Grzegorz Chrupała, Dieuwke Hupkes, Yuval Pinter, and Hassan Sajjad (eds.), *Proceedings of  
429 the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pp.  
430 163–173, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/  
431 2020.blackboxnlp-1.16. URL <https://aclanthology.org/2020.blackboxnlp-1.16/>.
- 432 Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal abstractions of neural  
433 networks. In *Proceedings of the 35th International Conference on Neural Information Processing  
434 Systems, NIPS '21*, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- 435 Atticus Geiger, Duligur Ibeling, Amir Zur, Maheep Chaudhary, Sonakshi Chauhan, Jing Huang,  
436 Aryaman Arora, Zhengxuan Wu, Noah Goodman, Christopher Potts, and Thomas Icard. Causal  
437 abstraction: A theoretical foundation for mechanistic interpretability. *Journal of Machine Learning  
438 Research*, 2025a. URL <http://jmlr.org/papers/v26/23-0058.html>.
- 439 Atticus Geiger, Duligur Ibeling, Amir Zur, Maheep Chaudhary, Sonakshi Chauhan, Jing Huang,  
440 Aryaman Arora, Zhengxuan Wu, Noah Goodman, Christopher Potts, and Thomas Icard. Causal  
441 abstraction: A theoretical foundation for mechanistic interpretability. *Journal of Machine Learning  
442 Research*, 26(83):1–64, 2025b. URL <http://jmlr.org/papers/v26/23-0058.html>.
- 443 Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual  
444 associations in auto-regressive language models, 2023. URL [https://arxiv.org/abs/2304.  
445 14767](https://arxiv.org/abs/2304.14767).

- 446 Mario Giulianelli, Jack Harding, Florian Mohnert, Dieuwke Hupkes, and Willem Zuidema. Under  
447 the hood: Using diagnostic classifiers to investigate and improve how language models track  
448 agreement information. In Tal Linzen, Grzegorz Chrupała, and Afra Alishahi (eds.), *Proceedings*  
449 *of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for*  
450 *NLP*, pp. 240–248, Brussels, Belgium, November 2018. Association for Computational Linguistics.  
451 doi: 10.18653/v1/W18-5426. URL <https://aclanthology.org/W18-5426/>.
- 452 Satchel Grant, Noah D. Goodman, and James L. McClelland. Emergent symbol-like number variables  
453 in artificial neural networks, 2025. URL <https://arxiv.org/abs/2501.06141>.
- 454 Satchel Grant, Simon Jerome Han, Alexa R. Tartaglino, and Christopher Potts. Addressing divergent  
455 representations from causal interventions on neural networks, 2026. URL <https://arxiv.org/abs/2511.04638>.
- 457 Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi,  
458 and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like  
459 one. *ArXiv e-print*, 2019.
- 460 Clément Guerner, Anej Svete, Tianyu Liu, Alexander Warstadt, and Ryan Cotterell. A geometric  
461 notion of causal probing. *CoRR*, abs/2307.15054, 2023. doi: 10.48550/ARXIV.2307.15054. URL  
462 <https://doi.org/10.48550/arXiv.2307.15054>.
- 463 Yoav Gur-Arieh, Roy Mayan, Chen Agassy, Atticus Geiger, and Mor Geva. Enhancing automated  
464 interpretability with output-centric feature descriptions. In *Proceedings of the 63rd Annual*  
465 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2025. URL  
466 <https://aclanthology.org/2025.acl-long.288/>.
- 467 Wes Gurnee, Emmanuel Ameisen, Isaac Kauvar, Julius Tarnag, Adam Pearce, Chris Olah, and Joshua  
468 Batson. When models manipulate manifolds: The geometry of a counting task. *arXiv preprint*  
469 *arXiv:2601.04480*, 2026.
- 470 Peter Gärdenfors. *Conceptual Spaces: The Geometry of Thought*. The MIT Press, 03 2000. ISBN  
471 9780262273558. doi: 10.7551/mitpress/2076.001.0001. URL [https://doi.org/10.7551/](https://doi.org/10.7551/mitpress/2076.001.0001)  
472 [mitpress/2076.001.0001](https://doi.org/10.7551/mitpress/2076.001.0001).
- 473 David Ha and Jürgen Schmidhuber. World models. In *Advances in Neural Information Processing*  
474 *Systems*, 2018.
- 475 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning  
476 behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- 477 Yixiong Hao, Ayush Panda, Stepan Shabalín, and Sheikh Abdur Raheem Ali. Patterns and mecha-  
478 nisms of contrastive activation engineering. *arXiv preprint arXiv:2505.03189*, 2025.
- 479 John J Hopfield. Neural networks and physical systems with emergent collective computational  
480 abilities. *Proceedings of the national academy of sciences*, 1982.
- 481 Jing Huang, Zhengxuan Wu, Christopher Potts, Mor Geva, and Atticus Geiger. Ravel: Evaluating  
482 interpretability methods on disentangling language model representations. In *Proceedings of the*  
483 *62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
484 2024. URL <https://aclanthology.org/2024.acl-long.470>.
- 485 David Hume. *An Enquiry Concerning Human Understanding*. A. Millar, London, 1748.
- 486 Iolo Jones. Diffusion geometry. *arXiv preprint arXiv:2405.10858*, 2024.
- 487 Iolo Jones and David Lanners. Computing diffusion geometry. *arXiv preprint arXiv:2602.06006*,  
488 2026.
- 489 Subhash Kantamneni and Max Tegmark. Language models use trigonometry to do addition, 2025.  
490 URL <https://arxiv.org/abs/2502.00873>.
- 491 Dhruva Karkada, Daniel J. Korchinski, Andres Nava, Matthieu Wyart, and Yasaman Bahri. Symmetry  
492 in language statistics shapes the geometry of model representations, 2026. URL [https://arxiv.](https://arxiv.org/abs/2602.15029)  
493 [org/abs/2602.15029](https://arxiv.org/abs/2602.15029).

- 494 Daniel J Korchinski, Dhruva Karkada, Yasaman Bahri, and Matthieu Wyart. On the emergence of  
495 linear analogies in word embeddings. *arXiv preprint arXiv:2505.18651*, 2025.
- 496 Austin C. Kozlowski, Callin Dai, and Andrei Boutyline. Semantic structure in large language model  
497 embeddings, 2025. URL <https://arxiv.org/abs/2508.10003>.
- 498 Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, Fugie Huang, et al. A tutorial on energy-based  
499 learning. *Predicting structured data*, 2006.
- 500 Michael A Lepori, Tal Linzen, Ann Yuan, and Katja Filippova. Language models struggle to use  
501 representations learned in-context. *arXiv preprint arXiv:2602.04212*, 2026.
- 502 Jiwei Li, Will Monroe, and Dan Jurafsky. Understanding neural networks through representation  
503 erasure, 2017. URL <https://arxiv.org/abs/1612.08220>.
- 504 Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-  
505 time intervention: eliciting truthful answers from a language model. In *Proceedings of the 37th  
506 International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY,  
507 USA, 2023. Curran Associates Inc.
- 508 Ekdeep Singh Lubana, Can Rager, Sai Sumedh R Hindupur, Valerie Costa, Greta Tuckute, Oam Patel,  
509 Sonia Krishna Murthy, Thomas Fel, Daniel Wurgaft, Eric J Bigelow, et al. Priors in time: Missing  
510 inductive biases for language model interpretability. *arXiv preprint arXiv:2511.01836*, 2025.
- 511 Niru Maheswaranathan, Alex Williams, Matthew Golub, Surya Ganguli, and David Sussillo. Uni-  
512 versality and individuality in neural dynamics across large populations of recurrent networks.  
513 *Advances in neural information processing systems*, 32, 2019.
- 514 Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language  
515 model representations of true/false datasets, 2023.
- 516 Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large  
517 language model representations of true/false datasets, 2024. URL [https://arxiv.org/abs/  
518 2310.06824](https://arxiv.org/abs/2310.06824).
- 519 Marina Meilă and Hanyu Zhang. Manifold learning: what, how, and why, 2023. URL [https:  
520 //arxiv.org/abs/2311.03757](https://arxiv.org/abs/2311.03757).
- 521 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associ-  
522 ations in GPT. *Advances in Neural Information Processing Systems*, 36, 2022. arXiv:2202.05262.
- 523 Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. Mass-editing  
524 memory in a transformer. In *The Eleventh International Conference on Learning Representations,  
525 ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.  
526 net/pdf?id=MkbcAHYgyS](https://openreview.net/pdf?id=MkbcAHYgyS).
- 527 Jack Merullo, Noah A. Smith, Sarah Wiegrefe, and Yanai Elazar. On linear representations and  
528 pretraining data frequency in language models, 2025. URL [https://arxiv.org/abs/2504.  
529 12459](https://arxiv.org/abs/2504.12459).
- 530 Alexander Modell, Patrick Rubin-Delanchy, and Nick Whiteley. The origins of representation  
531 manifolds in large language models, 2025a. URL <https://arxiv.org/abs/2505.18235>.
- 532 Alexander Modell, Patrick Rubin-Delanchy, and Nick Whiteley. The origins of representation  
533 manifolds in large language models, 2025b. URL <https://arxiv.org/abs/2505.18235>.
- 534 Andrew William Moore. *Efficient Memory-Based Learning for Robot Control*. PhD thesis, University  
535 of Cambridge, 1990.
- 536 Depen Morwani, Benjamin L. Edelman, Cosmin Oncescu, Rosie Zhao, and Sham M. Kakade. Feature  
537 emergence via margin maximization: Case studies in algebraic tasks. In *The Twelfth International  
538 Conference on Learning Representations*, 2024.

539 Aaron Mueller, Jannik Brinkmann, Millicent L. Li, Samuel Marks, Koyena Pal, Nikhil Prakash, Can  
540 Rager, Aruna Sankaranarayanan, Arnab Sen Sharma, Jiuding Sun, Eric Todd, David Bau, and  
541 Yonatan Belinkov. The quest for the right mediator: A history, survey, and theoretical grounding of  
542 causal interpretability. *CoRR*, abs/2408.01416, 2024. doi: 10.48550/ARXIV.2408.01416. URL  
543 <https://doi.org/10.48550/arXiv.2408.01416>.

544 Aaron Mueller, Atticus Geiger, Sarah Wiegrefe, Dana Arad, Iván Arcuschin, Adam Belfki, Yik Siu  
545 Chan, Jaden Fiotto-Kaufman, Tal Haklay, Michael Hanna, Jing Huang, Rohan Gupta, Yaniv  
546 Nikankin, Hadas Orgad, Nikhil Prakash, Anja Reusch, Aruna Sankaranarayanan, Shun Shao,  
547 Alessandro Stolfo, Martin Tutek, Amir Zur, David Bau, and Yonatan Belinkov. Mib: A mechanistic  
548 interpretability benchmark, 2025. URL <https://arxiv.org/abs/2504.13151>.

549 Aaron Mueller, Jannik Brinkmann, Millicent Li, Samuel Marks, Koyena Pal, Nikhil Prakash, Can  
550 Rager, Aruna Sankaranarayanan, Arnab Sen Sharma, Jiuding Sun, Eric Todd, David Bau, and  
551 Yonatan Belinkov. The quest for the right mediator: Surveying mechanistic interpretability for  
552 nlp through the lens of causal mediation analysis. *Computational Linguistics*, pp. 1–48, 02 2026.  
553 ISSN 0891-2017. doi: 10.1162/COLI.a.572. URL <https://doi.org/10.1162/COLI.a.572>.

554 Maya Okawa, Ekdeep Singh Lubana, Robert P. Dick, and Hidenori Tanaka. Compositional abilities  
555 emerge multiplicatively: Exploring diffusion models on a synthetic task, 2024.

556 Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt  
557 Turner. Steering llama 2 via contrastive activation addition, 2024. URL <https://arxiv.org/abs/2312.06681>.

559 Core Francisco Park, Andrew Lee, Ekdeep Singh Lubana, Yongyi Yang, Maya Okawa, Kento Nishi,  
560 Martin Wattenberg, and Hidenori Tanaka. Iclr: In-context learning of representations. In Y. Yue,  
561 A. Garg, N. Peng, F. Sha, and R. Yu (eds.), *International Conference on Learning Representations*,  
562 volume 2025, pp. 53258–53284, 2025a. URL [https://proceedings.iclr.cc/paper\\_files/](https://proceedings.iclr.cc/paper_files/paper/2025/file/83fe5a77502e3d4cfab5960aed0ee6c3-Paper-Conference.pdf)  
563 [paper/2025/file/83fe5a77502e3d4cfab5960aed0ee6c3-Paper-Conference.pdf](https://proceedings.iclr.cc/paper_files/paper/2025/file/83fe5a77502e3d4cfab5960aed0ee6c3-Paper-Conference.pdf).

564 Core Francisco Park, Andrew Lee, Ekdeep Singh Lubana, Yongyi Yang, Maya Okawa, Kento Nishi,  
565 Martin Wattenberg, and Hidenori Tanaka. Iclr: In-context learning of representations. In *The*  
566 *Thirteenth International Conference on Learning Representations*, 2025b.

567 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry  
568 of large language models. *arXiv preprint arXiv:2311.03658*, 2023.

569 Kiho Park, Yo Joong Choe, Yibo Jiang, and Victor Veitch. The geometry of categorical and  
570 hierarchical concepts in large language models. *Proceedings of the International Conference on*  
571 *Learning Representations (ICLR)*, 2025c.

572 Kiho Park, Todd Nief, Yo Joong Choe, and Victor Veitch. The information geometry of softmax:  
573 Probing and steering. *arXiv preprint arXiv:2602.15293*, 2026.

574 Michael Pearce, Elana Simon, Michael Byun, and Daniel Balsam. Finding the tree of life in evo 2.  
575 *Goodfire*, August 2025. Correspondence to michael@goodfire.ai.

576 Judea Pearl. Probabilities of causation: Three counterfactual interpretations and their identification.  
577 *Synthese*, 121(1):93–149, 1999.

578 Judea Pearl. Direct and indirect effects, 2001. URL <https://arxiv.org/abs/1301.2300>.

579 Nikhil Prakash, Natalie Shapira, Arnab Sen Sharma, Christoph Riedl, Yonatan Belinkov, Tamar Rott  
580 Shaham, David Bau, and Atticus Geiger. Language models use lookbacks to track beliefs, 2025.  
581 URL <https://arxiv.org/abs/2505.14685>.

582 Itamar Pres, Laura Ruis, Ekdeep Singh Lubana, and David Krueger. Towards reliable evaluation of  
583 behavior steering interventions in llms. *arXiv preprint arXiv:2410.17245*, 2024.

584 Lucas Prieto, Edward Stevinson, Melih Barsbey, Tolga Birdal, and Pedro A. M. Mediano. Correlations  
585 in the data lead to semantically rich feature geometry under superposition. In *The Fourteenth*  
586 *International Conference on Learning Representations*, 2026. URL [https://openreview.net/](https://openreview.net/forum?id=7akSRQS5Xh)  
587 [forum?id=7akSRQS5Xh](https://openreview.net/forum?id=7akSRQS5Xh).

- 588 Goutham Rajendran, Simon Buchholz, Bryon Aragam, Bernhard Schölkopf, and Pradeep Ravikumar.  
589 From causal to concept-based representation learning. *Advances in Neural Information Processing*  
590 *Systems*, 37:101250–101296, 2024a.
- 591 Goutham Rajendran, Simon Buchholz, Bryon Aragam, Bernhard Schölkopf, and Pradeep Ravikumar.  
592 Learning interpretable concepts: Unifying causal representation learning and foundation models.  
593 *arXiv preprint arXiv:2402.09236*, 2024b.
- 594 Shauli Ravfogel, Michael Twiton, Yoav Goldberg, and Ryan Cotterell. Linear adversarial concept  
595 erasure, 2022.
- 596 Shauli Ravfogel, Yoav Goldberg, and Ryan Cotterell. Log-linear guardedness and its implications. In  
597 Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual*  
598 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023,*  
599 *Toronto, Canada, July 9-14, 2023*, pp. 9413–9431. Association for Computational Linguistics,  
600 2023a. doi: 10.18653/V1/2023.ACL-LONG.523. URL <https://doi.org/10.18653/v1/2023.acl-long.523>.
- 602 Shauli Ravfogel, Francisco Vargas, Yoav Goldberg, and Ryan Cotterell. Kernelized concept erasure,  
603 2023b.
- 604 Christian H Reinsch. Smoothing by spline functions. *Numerische mathematik*, 10(3):177–183, 1967.
- 605 Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. Steering  
606 llama 2 via contrastive activation addition. In Lun-Wei Ku, Andre Martins, and Vivek  
607 Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computa-*  
608 *tional Linguistics (Volume 1: Long Papers)*, pp. 15504–15522, Bangkok, Thailand, August  
609 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.828. URL  
610 <https://aclanthology.org/2024.acl-long.828/>.
- 611 Juan Diego Rodriguez, Aaron Mueller, and Kanishka Misra. Characterizing the role of similarity in  
612 the property inferences of language models. *CoRR*, abs/2410.22590, 2024. doi: 10.48550/ARXIV.  
613 2410.22590. URL <https://doi.org/10.48550/arXiv.2410.22590>.
- 614 Pau Rodriguez, Arno Blaas, Michal Klein, Luca Zappella, Nicholas Apostoloff, Marco Cuturi, and  
615 Xavier Suau. Controlling language and diffusion models by transporting activations. In *The*  
616 *Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-*  
617 *28, 2025*. OpenReview.net, 2025. URL <https://openreview.net/forum?id=l2zFn6TIQi>.
- 618 Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding.  
619 *science*, 290(5500):2323–2326, 2000.
- 620 Paul K. Rubenstein, Sebastian Weichwald, Stephan Bongers, Joris M. Mooij, Dominik Janzing,  
621 Moritz Grosse-Wentrup, and Bernhard Schölkopf. Causal consistency of structural equation  
622 models. In *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence (UAI)*,  
623 2017.
- 624 Naomi Saphra and Sarah Wiegrefe. Mechanistic? *CoRR*, abs/2410.09087, 2024. doi: 10.48550/  
625 ARXIV.2410.09087. URL <https://doi.org/10.48550/arXiv.2410.09087>.
- 626 Raphaël Sarfati, Eric Bigelow, Daniel Wurgaft, Jack Merullo, Atticus Geiger, Owen Lewis, Tom  
627 McGrath, and Ekdeep Singh Lubana. The shape of beliefs: Geometry, dynamics, and interventions  
628 along representation manifolds of language models’ posteriors, 2026. URL <https://arxiv.org/abs/2602.02315>.
- 630 Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic  
631 development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):  
632 11537–11546, 2019.
- 633 Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component  
634 analysis. In *International conference on artificial neural networks*, pp. 583–588. Springer, 1997.

- 635 Adam Shai, Loren Amdahl-Culleton, Casper L. Christensen, Henry R. Bigelow, Fernando E. Rosas,  
636 Alexander B. Boyd, Eric A. Alt, Kyle J. Ray, and Paul M. Riechers. Transformers learn factored  
637 representations, 2026. URL <https://arxiv.org/abs/2602.02385>.
- 638 Adam S Shai, Sarah E Marzen, Lucas Teixeira, Alexander G Oldenziel, and Paul M Riechers. Trans-  
639 formers represent belief state geometry in their residual stream. *Advances in Neural Information*  
640 *Processing Systems*, 37:75012–75034, 2024a.
- 641 Adam S. Shai, Sarah E. Marzen, Lucas Teixeira, Alexander Gietelink Oldenziel, and Paul M.  
642 Riechers. Transformers represent belief state geometry in their residual stream, 2024b. URL  
643 <https://arxiv.org/abs/2405.15943>.
- 644 Roger N Shepard. Toward a universal law of generalization for psychological science. *Science*, 237  
645 (4820):1317–1323, 1987.
- 646 P. Smolensky. *Neural and conceptual interpretation of PDP models*, pp. 390–431. MIT Press,  
647 Cambridge, MA, USA, 1986. ISBN 0262631105.
- 648 Jiajun Song and Yiqiao Zhong. Uncovering hidden geometry in transformers via disentangling  
649 position and context. *arXiv preprint arXiv:2310.04861*, 2023.
- 650 Yang Song and Diederik P Kingma. How to train your energy-based models. *ArXiv e-print*, 2021.
- 651 Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT Press,  
652 2000.
- 653 Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. A mechanistic interpretation of  
654 arithmetic reasoning in language models using causal mediation analysis. In Houda Bouamor,  
655 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in*  
656 *Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 7035–7052.  
657 Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.EMNLP-MAIN.435.  
658 URL <https://doi.org/10.18653/v1/2023.emnlp-main.435>.
- 659 Nishant Subramani, Nivedita Suresh, and Matthew Peters. Extracting latent steering vectors from  
660 pretrained language models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.),  
661 *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 566–581, Dublin, Ireland,  
662 May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.48.  
663 URL <https://aclanthology.org/2022.findings-acl.48/>.
- 664 Denis Sutter, Julian Minder, Thomas Hofmann, and Tiago Pimentel. The non-linear representation  
665 dilemma: Is causal abstraction enough for mechanistic interpretability?, 2025. URL <https://arxiv.org/abs/2507.08802>.
- 667 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2nd  
668 edition, 2018.
- 669 Daniel Tan, David Chanin, Aengus Lynch, Brooks Paige, Dimitrios Kanoulas, Adrià Garriga-Alonso,  
670 and Robert Kirk. Analysing the generalisation and reliability of steering vectors. *Advances in*  
671 *Neural Information Processing Systems*, 37:139179–139212, 2024.
- 672 Generalist AI Team. Gen-0: Embodied foundation models that scale with physical interaction.  
673 *Generalist AI Blog*, 2025. <https://generalistai.com/blog/nov-04-2025-GEN-0>.
- 674 Joshua B Tenenbaum and Thomas L Griffiths. Generalization, similarity, and bayesian inference.  
675 *Behavioral and brain sciences*, 24(4):629–640, 2001.
- 676 Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau.  
677 Function vectors in large language models. In *The Twelfth International Conference on Learning*  
678 *Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL  
679 <https://openreview.net/forum?id=AwyxtyMwaG>.
- 680 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
681 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
682 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

- 683 Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu,  
684 Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard  
685 interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- 686 Alexander Matt Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDi-  
687 armid. Activation addition: Steering language models without optimization, 2023.
- 688 Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini,  
689 and Monte MacDiarmid. Steering language models with activation engineering, 2024. URL <https://arxiv.org/abs/2308.10248>.
- 691 Daniel Vennemeyer, Phan Anh Duong, Tiffany Zhan, and Tianyu Jiang. Sycophancy is not one thing:  
692 Causal separation of sycophantic behaviors in llms. *arXiv preprint arXiv:2509.21305*, 2025.
- 693 Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and  
694 Stuart Shieber. Investigating gender bias in language models using causal mediation analy-  
695 sis. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances*  
696 *in Neural Information Processing Systems*, volume 33, pp. 12388–12401. Curran Associates,  
697 Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/](https://proceedings.neurips.cc/paper_files/paper/2020/file/92650b2e92217715fe312e6fa7b90d82-Paper.pdf)  
698 [92650b2e92217715fe312e6fa7b90d82-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/92650b2e92217715fe312e6fa7b90d82-Paper.pdf).
- 699 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.  
700 Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In  
701 *The Eleventh International Conference on Learning Representations*, 2023a. URL [https:](https://openreview.net/forum?id=NpsVSN6o4ul)  
702 [//openreview.net/forum?id=NpsVSN6o4ul](https://openreview.net/forum?id=NpsVSN6o4ul).
- 703 Zihao Wang, Lin Gui, Jeffrey Negrea, and Victor Veitch. Concept algebra for score-based conditional  
704 model. In *ICML 2023 Workshop on Structured Probabilistic Inference & Generative Modeling*,  
705 2023b.
- 706 Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christo-  
707 pher D. Manning, and Christopher Potts. Axbench: Steering llms? even simple baselines outper-  
708 form sparse autoencoders, 2025. URL <https://arxiv.org/abs/2501.17148>.
- 709 Daniel Wurgaft, Ekdeep Singh Lubana, Core Francisco Park, Hidenori Tanaka, Gautam Reddy,  
710 and Noah D. Goodman. In-context learning strategies emerge rationally, 2025. URL [https:](https://arxiv.org/abs/2506.17859)  
711 [//arxiv.org/abs/2506.17859](https://arxiv.org/abs/2506.17859).
- 712 Yongyi Yang, Hidenori Tanaka, and Wei Hu. Provable low-frequency bias of in-context learning of  
713 representations. *arXiv preprint arXiv:2507.13540*, 2025.
- 714 Julian Yocum, Cameron Allen, Bruno Olshausen, and Stuart Russell. Neural manifold geome-  
715 try encodes feature fields. In *NeurIPS 2025 Workshop on Symmetry and Geometry in Neural*  
716 *Representations*, 2025.
- 717 Carolina Zheng, Nicolas Beltran-Velez, Sweta Karlekar, Claudia Shi, Achille Nazaret, Asif Mallik,  
718 Amir Feder, and David M Blei. Model directions, not words: Mechanistic topic models using  
719 sparse autoencoders. *arXiv preprint arXiv:2507.23220*, 2025.
- 720 Tianyi Zhou, Deqing Fu, Mahdi Soltanolkotabi, Robin Jia, and Vatsal Sharan. Fone: Precise  
721 single-token number embeddings via fourier features. *arXiv preprint arXiv:2502.09741*, 2025.

## 722 A Manifold Steering Yields Factored Control in Multi-Dimensional Spaces

723 Our experiments thus far have been limited to one dimensional conceptual spaces arising from  
724 training data imbued with real-world structure, i.e., days, months, ages, and letters. In turn, the  
725 manifolds we found have been one dimensional curves with a single intrinsic coordinate. Now, we  
726 extend our results to a setting with two dimensional conceptual spaces whose geometry are defined  
727 via in-context learning. We fit manifolds and show there is a two dimensional intrinsic coordinate  
728 system for the manifold, where steering along each coordinate controls an independent dimension of  
729 the conceptual space.

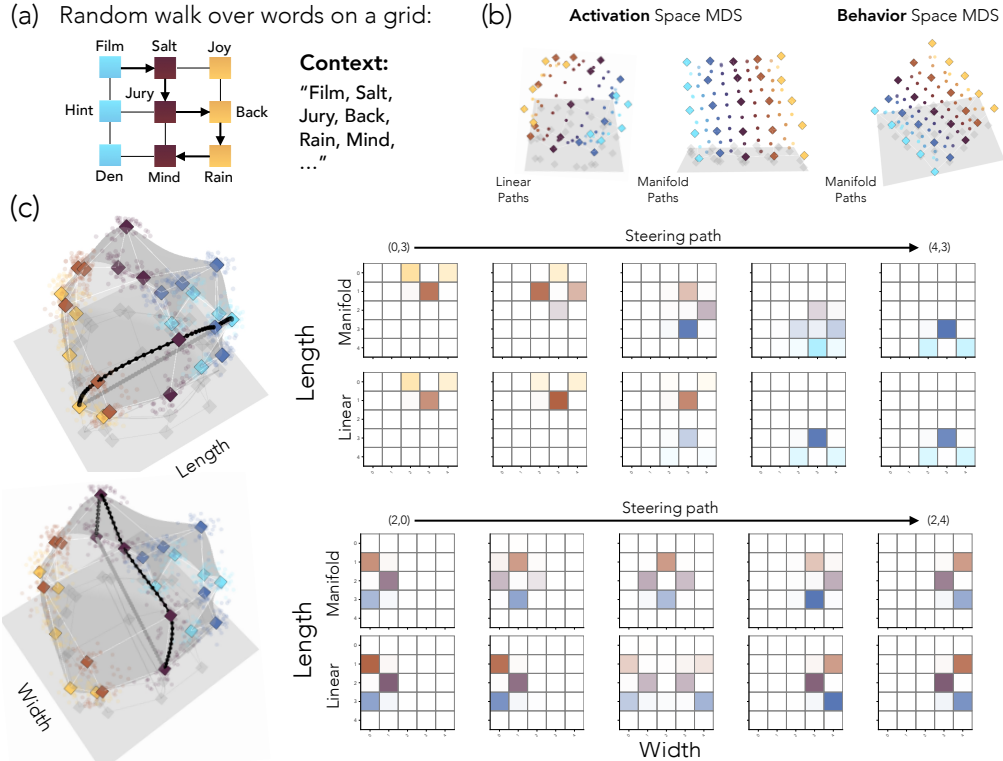
730 **In-context learning tasks with synthetic conceptual spaces.** Park et al. (2025b) introduce a family  
731 of tasks to study the in-context learning of representations (ICLR). For each task, arbitrary tokens are  
732 assigned to a discrete graphical structure and language models are supplied with sequences of tokens  
733 derived from a random walk on that graph. They show that the statistical patterns in the random  
734 walk of tokens induce a reorganization of representations that recapitulates the graphical structure  
735 used to generate data. This in turn enables the language model to match the next token distribution,  
736 i.e., predict tokens adjacent to the current location of the random walk on the grid. For our two  
737 experiments, we assign arbitrary tokens to grid and cylinder graph structures. Thus, the conceptual  
738 domain  $\mathcal{Z}$  is the set of tokens and the distance metric  $d_{\mathcal{Z}}$  is distance on the graph used to generate the  
739 random walk. Fig. 5 (a) shows an example grid and an input prompt generated by a random walk.

740 **Manifold fitting.** The ICLR grid manifold  $\mathcal{M}_h$  is topologically described by a two dimensional  
741 surface with no holes or tears. The activation geometry of  $\mathcal{M}_h$  is more complex. Its semi-spherical  
742 shape shown in Fig. 5 (c) is induced by task statistics: the random walk visits inner sites more  
743 frequently than peripheral sites, leading to slight distortions with respect to the ground truth geometry  
744 (Park et al., 2025b; Yang et al., 2025; Karkada et al., 2026). We fit two-dimensional sheets to internal  
745 activations and output distributions via thin plate splines (TPS; Duchon 1977; Bookstein 1989), which  
746 can be seen as the 2D analog of the cubic splines used previously. In this case, we use activations  
747 corresponding to the last token in the context, and compute centroids according to graph location  
748 at a given timestep. Then, TPS finds the smoothest surface interpolating through the centroids (see  
749 App. F.3 for further details).

750 **Isometry results.** For each ICLR domain, we compute pairwise distance matrices over graph-node  
751 centroids under three metrics: Euclidean (linear) distance in the activation subspace, geodesic distance  
752 along the fitted activation manifold  $\mathcal{M}_h$ , and geodesic distance along the behavior manifold  $\mathcal{M}_y$ . We  
753 find very high correlations between geodesic paths on the activation and behavior manifolds ( $r = .99$   
754 for both the  $5 \times 5$  grid and  $9 \times 9$  cylinder domains) and reduced correlations for linear paths ( $5 \times 5$   
755 grid  $r = 0.90$ ;  $9 \times 9$  cylinder  $r = 0.81$ ). To further examine these results, we embed each distance  
756 matrix with multidimensional scaling (MDS). Fig. 5(b) shows a clean grid structure in the activation  
757 manifold and behavior manifold embeddings, while the linear paths in activation space yield a warped  
758 surface. Again, we see the conceptual space recapitulated in both representation and behavior.

759 **Manifold vs. Linear steering.** We next test whether the fitted two-dimensional activation manifold  
760 affords coherent, factored control over the graph geometry used to generate the ICLR inputs. In  
761 particular, we assess whether we can control the position of the random walk input via intervention,  
762 and, moreover, whether there is an intrinsic coordinate system where steering along each coordinate  
763 independently controls the horizontal and vertical position on the grid. For each ordered pair of nodes,  
764 we use manifold steering and linear steering to interpolate between the start and end centroid and  
765 average results over 5 input prompts.

766 The top panel of Fig. 5(c) shows that manifold steering produces smooth transitions vertically along  
767 the steered graph dimension while remaining at the same horizontal position. The bottom panel of  
768 Fig. 5(c) shows similar smooth transitions but along a horizontal dimension while keeping the same  
769 vertical position. This demonstrates the manifold has an intrinsic coordinate system corresponding  
770 to the two dimensions of the grid, enabling factored control. Furthermore, this shows the smooth  
771 and ordered transitions of manifold steering generalize to multi-dimensional spaces. In contrast,  
772 linear steering again fails to provide ordered transitions through grid locations, and shows very clear  
773 ‘teleportation’ behavior between the endpoint locations along its path.



**Figure 5: Manifold steering enables factored control in multi-dimensional conceptual spaces.** (a) We examine manifold steering on multidimensional spaces using Park et al. (2025b)’s in-context learning of representations (ICLR) task. In an ICLR task, arbitrary tokens are assigned to nodes along a graph, and a language model is prompted with tokens from a random walk along the graph. Park et al. (2025b) showed that with sufficient context, models encode the structure of the latent graph in their activations. In this work, we study two graph structures learned in-context ( $5 \times 5$  grid shown above,  $9 \times 9$  cylinder in App. H). We fit manifolds to activations and output behaviors and intervene on activations using linear and manifold steering. (b) We examine the mapping between the activation and behavior manifolds by computing on-manifold and linear distances in activation space and comparing them to on-manifold distances in behavior space via a multidimensional scaling (MDS) embedding. We find a clear structural match of both the activation and output manifolds with the latent graph, providing direct evidence for these two manifolds encoding a similar underlying conceptual space. In contrast, linear distances in activation space yield a warped structure. (c) We find that manifold steering maintains the quality of smooth and ordered transitions beyond one dimension, and in conceptual spaces learned in-context. Furthermore, we find that steering along one dimension leads to minimal off-target impact, thus affording the appealing quality of factored control. In contrast, linear steering maintains its teleportation behavior.

## 774 B Manifold Steering on a Visual World Model: Mountain Car Task

775 We now ask whether the same principles of geometry-aware steering extend to the *visual* domain  
 776 of world models. This question is practically motivated: learned world models that predict future  
 777 observations from past frames and actions are central to model-based reinforcement learning and  
 778 robotic planning (Ha & Schmidhuber, 2018; Hafner et al., 2020; Team, 2025; Black et al., 2024). If  
 779 the internal representations of such models admit geometric structure, manifold-based steering could  
 780 provide a principled mechanism for intervening on a model’s behavior through changing its beliefs  
 781 about the state of the world.

782 **Environment and model architecture.** We train a recurrent world model on the Mountain Car  
 783 environment (Moore, 1990; Sutton & Barto, 2018), a classical control task in which a car must  
 784 escape a valley by building momentum. The environment has continuous position  $p \in [-1.2, 0.6]$ ,

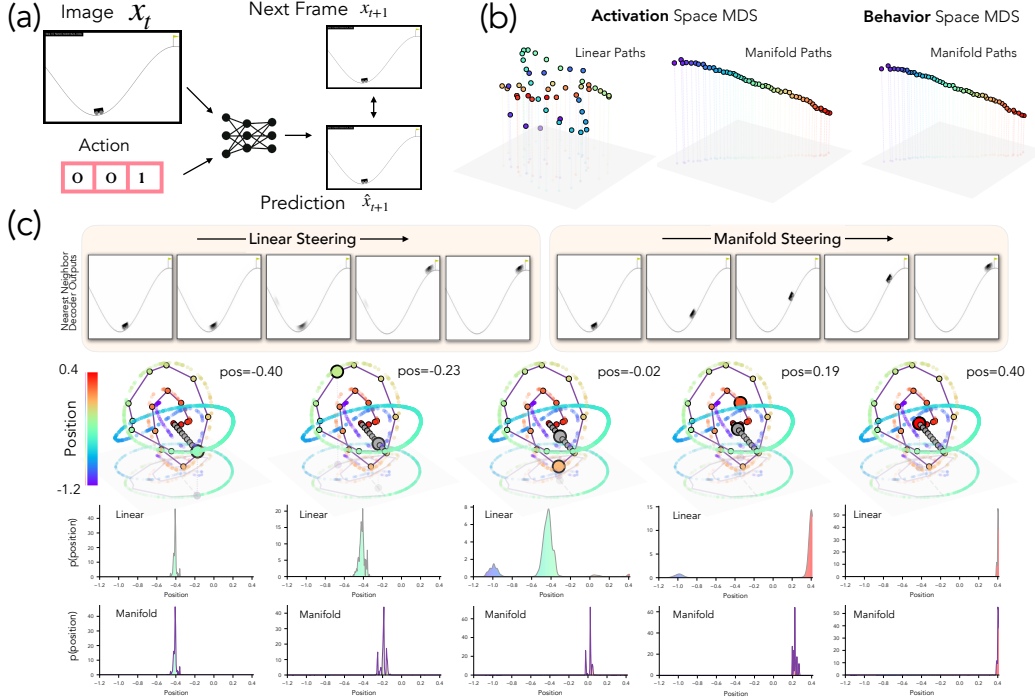


Figure 6: **Manifold steering on a visual world model produces smooth movement.** (a). We examine whether manifold steering can generalize to a visual modality by training a recurrent network on the Mountain Car environment (Moore, 1990; Sutton & Barto, 2018) to predict the next frame  $x_{t+1}$  given the previous frame  $x_t$  and an action. (b) We test the mapping between the activation and behavior manifolds by computing on-manifold and linear distances in activation space and comparing them to on-manifold distances in behavior space via an MDS embedding. On-manifold paths in activation and behavior space both recover a clean sequential ordering corresponding to location, while the embedding of linear distances scrambles it. (c) **Middle:** PCA visualization of the activation manifold and five waypoints along a path between  $p_A = -0.40$  and  $p_B = 0.40$  for both linear (red) and manifold (blue) steering. **Top:** At intermediate positions along the linear steering path, decoding shows the car as blurred or ambiguously placed, reflecting an incoherent superposition of positional beliefs (**bottom**) as the path departs from the activation manifold  $\mathcal{M}_h$ . In contrast, manifold steering along  $\mathcal{M}_h$  yields smooth movement of the car up the hill.

785 continuous velocity  $v \in [-0.07, 0.07]$ , and three discrete actions (left, no-op, right). The model  
 786 predicts the next frame  $x_{t+1}$  given the previous frame  $x_t$  and action  $a_t$  (see Fig. 6(a) for an illustration).  
 787 The full architecture is shown in Fig. 7: A convolutional encoder maps each  $128 \times 128 \times 3$  RGB  
 788 frame to a latent vector,  $v_t$ , which is concatenated with a learned action embedding  $e(a_t) \in \mathbb{R}^{16}$  and  
 789 fed to a Gated Recurrent Unit (GRU; Cho et al. 2014):

$$\mathbf{h}_t = \text{GRU}([v_t; e(a_t)], \mathbf{h}_{t-1}) \in \mathbb{R}^n; \quad v_t = \text{LayerNorm}(f_{\text{enc}}(x_t)) \in \mathbb{R}^n, \quad (8)$$

790 where  $n = 64$ . A convolutional decoder produces a residual image from the hidden state, yielding  
 791 the prediction  $\hat{x}_{t+1} = x_t + f_{\text{dec}}(\mathbf{h}_t)$ .

792 **Activation and behavior manifold fitting.** For this setting, we consider *position* to play the role of  
 793 the conceptual domain  $\mathcal{Z} = [p_{\min}, p_{\max}]$  and aim to capture the manifold structure in both activation  
 794 and behavior space of the vision encoder. To start, we first collect encoder activations from 100  
 795 rollouts in the environment (see §G.1 for details) and observe they occupy a curved, low-dimensional  
 796 manifold  $\mathcal{M}_h \subset \mathbb{R}^n$  (Fig. 6(c)). We parameterize this manifold by partitioning the position range  
 797 into bins and fit a smooth spline through the means  $\{\mu_b\}_{b=1}^B \subset \mathbb{R}^n$ . For a given input,  $x$ , we compute  
 798 the output distribution over positions,  $\mathbf{p}(x)$  using the distance of the activations  $v(x)$  to the centroid  
 799 for each position:

$$\mathbf{p}(x) = \text{softmax}\left(-\frac{\|v - \mu_b\|_2}{\tau}\right)_{b=1}^B \in \Delta^{B-1}, \quad (9)$$

800 with temperature  $\tau = 0.5$ . We follow §3.2 to parameterize the behavior manifold  $\mathcal{M}_y$  by embedding  
801 each bin in Hellinger coordinates on the unit sphere and fit a 1D smoothing spline  $\gamma_{\mathcal{M}_y}: \mathcal{Z} \rightarrow \mathbb{R}^B$   
802 parameterized by position (full details in App. G). Note that both  $\mathcal{M}_h$  and  $\mathcal{M}_y$  are 1D structures  
803 parameterized by the conceptual coordinate  $p$ ; under PCA visualization (Fig. 12), both trace closed  
804 curves in their respective ambient spaces. The curves are closed because visually distinctive states at  
805 the wall,  $p \approx -1.2$ , and goal,  $p \approx 0.4$ , are mapped to neighboring activations.

806 **Geometry-aware steering in activation space.** Fig. 6 compares linear (Eq. 1) and manifold (Eq. 2)  
807 steering through  $K = 20$  waypoints between encoder states corresponding to positions  $p_A = -0.4$   
808 and  $p_B = 0.4$ , projected into the first three principal components of encoder space. The geodesic path  
809 closely tracks  $\mathcal{M}_h$ , and the corresponding decoded frames display a smooth, coherent progression of  
810 the car through intermediate positions. The linear path departs from  $\mathcal{M}_h$  at intermediate points and  
811 decoded frames exhibit blurred or ambiguous car placement, reflecting an incoherent superposition.  
812 Then, a ‘teleportation’ to the endpoint is observed, analogous to the behavior of linear steering in  
813 the language model experiments. Moreover, linear steering causes the probability distribution over  
814 position to show greater spread compared to on-manifold paths, yielding the ambiguous car placement  
815 seen in intermediate points along the path. Finally, we reproduce the pullback procedure in §H.4 and  
816 show that optimizing paths in the output distribution over possible car positions results in activation  
817 space paths that closely track the  $\mathcal{M}_h$  (Fig. 12).

818 **The geometry assumed by linear steering is not faithful to the conceptual ordering.** To make  
819 the difference between the two steering metrics visually concrete, we apply multidimensional scaling  
820 (MDS) to the pairwise distance matrices induced by three different distance functions over  $W = 50$   
821 anchor positions evenly spaced along  $[p_{\min}, p_{\max}]$  (Fig. 6(b)). Both activation space and behavior space  
822 on-manifold distance embeddings recover a clean one-dimensional rainbow ordering of positions,  
823 while the linear distance embedding produces a scrambled three-dimensional structure whose colors  
824 are visibly out of order. This is because  $\mathcal{M}_h$  folds back on itself in the encoder’s ambient space, so  
825 two activations whose underlying positions are far apart can sit arbitrarily close in ambient space.  
826 Quantitatively, the Pearson correlation between the two arc-length distance matrices is  $r = 0.99$ ,  
827 while correlation between activation-space linear paths and behavior manifold arc-length falls to  
828  $r = 0.06$ , confirming that the linear-steering metric is not a faithful proxy for the conceptual ordering  
829 that the encoder has learned.

## 830 C Related Work

831 **Activation Steering and the Linear Representation Hypothesis.** Activation steering proto-  
832 cols (Bau et al., 2018; Subramani et al., 2022; Marks & Tegmark, 2023; Panickssery et al., 2024;  
833 Turner et al., 2024) are often motivated by the linear representation hypothesis (LRH)—a geometric  
834 assumption on model representations (Smolensky, 1986; Elhage et al., 2022a; Park et al., 2023; Costa  
835 et al., 2025; Zheng et al., 2025). In particular, LRH argues neural networks encode concepts, i.e.,  
836 latent variables underlying the data distribution (Wang et al., 2023b; Rajendran et al., 2024a,b; Okawa  
837 et al., 2024), along directions. This motivates tools like linear probing (Belinkov, 2022; Guerner  
838 et al., 2023) and sparse autoencoders (Cunningham et al., 2023; Bricken et al., 2023; Gao et al., 2025;  
839 Bussmann et al., 2024; Fel et al., 2025a). In the case where the representation geometry for a concept  
840 truly aligns with LRH, Bigelow et al. (2025) showed the effects of activation steering on model  
841 behavior can be accurately captured by a linear increase in concept log-probability. However, in the  
842 general scenario where geometry of representations does not abide by LRH, the effects of linear  
843 steering protocols are less clear. Recent work has started to fill this gap: e.g., work by Rodriguez  
844 et al. (2025); Ravfogel et al. (2022) has shown that linear steering protocols match the first moments  
845 of the current output distribution produced by the model with the target distribution; however, the  
846 effects on higher order moments can be unconstrained and adversarial (see results by Sarfati et al.  
847 (2026)), possibly explaining why it produces incoherent outputs.

848 In contrast, when the representation geometry is fully respected, our work shows that steering  
849 smoothly interpolates the source and the target distributions. Prior work has shown similar results  
850 to this effect in narrow domains, e.g., Engels et al. (2024) ablate representations and write the  
851 days-of-the-week circle directly and Kantamneni & Tegmark (2025) follow a similar protocol for a  
852 helix representing numbers, but they lack a more general account of how representation geometry  
853 and output behavior map on to each other. The closest work to ours is the contemporary paper by

854 [Park et al. \(2026\)](#), who study a toy model where the representation-to-output distribution mapping is  
855 described via a simple softmax operation.

856 **Activation Geometry and its Origins.** A large body of recent work has shown neural networks  
857 encode concepts along nonlinear, curved geometries embedded in low-dimensional subspaces across  
858 both modalities and architectures ([Fel et al., 2025b](#); [Pearce et al., 2025](#); [Yocum et al., 2025](#); [Modell  
859 et al., 2025a](#); [Lubana et al., 2025](#); [Costa et al., 2025](#); [Park et al., 2025b](#); [Engels et al., 2024](#); [Karkada  
860 et al., 2026](#); [Shai et al., 2024a, 2026](#); [Saxe et al., 2019](#); [Park et al., 2025c](#); [Morwani et al., 2024](#);  
861 [Kantamneni & Tegmark, 2025](#); [Song & Zhong, 2023](#); [Zhou et al., 2025](#); [Maheswaranathan et al.,  
862 2019](#)). While earlier work ([Saxe et al., 2019](#); [Arora et al., 2018](#); [Park et al., 2023](#); [Yocum et al.,  
863 2025](#)) concretized, in toy settings, how structure in the data-generating process imposes geometric  
864 constraints on a neural network’s representations, only recently have such accounts been extended to  
865 make predictions about the geometry of neural representations at scale (e.g., [Merullo et al. \(2025\)](#)).  
866 [Karkada et al. \(2026\)](#) and [Korchinski et al. \(2025\)](#) argue that symmetries in data statistics enforce  
867 geometries best suited for reflecting the uncertainty of the distribution in a model’s representation (cf.  
868 [Prieto et al. 2026](#)), and offer plausible accounts for the formation of representations in-context, as  
869 shown by works such as [Park et al. \(2025a\)](#) and [Lepori et al. \(2026\)](#).

870 **Causal Analysis of Neural Networks.** Several works have convincingly argued that tools like  
871 probing or visualization of representations are insufficient to make claims about model behavior, i.e.,  
872 artifacts produced via these tools can yield misleading explanations for why a model behaves the  
873 way it does ([Geiger et al., 2020](#); [Belinkov, 2022](#); [Bolukbasi et al., 2021](#); [Saphra & Wiegrefe, 2024](#)).  
874 As such, a vast array of research has used intervention on activations to study model internals ([Li  
875 et al., 2017](#); [Giulianelli et al., 2018](#); [Cammarata et al., 2020](#); [Elazar et al., 2020](#); [Ravfogel et al., 2022,  
876 2023a,b](#); [Belrose et al., 2023](#); [Geva et al., 2023](#); [Meng et al., 2022, 2023](#); [Vig et al., 2020](#); [Geiger et al.,  
877 2020](#); [Davies et al., 2023](#); [Stolfo et al., 2023](#); [Guerner et al., 2023](#); [Wang et al., 2023a](#); [Todd et al.,  
878 2024](#); [Arora et al., 2024](#); [Huang et al., 2024](#); [Feng & Steinhardt, 2024](#); [Mueller et al., 2025](#); [Prakash  
879 et al., 2025](#); [Gur-Arieh et al., 2025](#); [Grant et al., 2025](#); [Rodriguez et al., 2024](#)). This interpretability  
880 research leverages the frameworks of causal mediation ([Pearl, 2001](#); [Vig et al., 2020](#); [Mueller et al.,  
881 2026](#)) and causal abstraction ([Rubenstein et al., 2017](#); [Beckers & Halpern, 2019](#); [Geiger et al., 2021,  
882 2025a,b](#)) to ground understanding of model internals in the theory of causality ([Hume, 1748](#); [Pearl,  
883 1999](#); [Spirtes et al., 2000](#)).

## 884 D Extended Discussion

885 **Where does the shared geometry of behavior and representation come from?** While we do  
886 not study the origins of the shared geometry between behavior and representation, our experimental  
887 results are consistent with the hypothesis that conceptual structure constrains the geometry of both  
888 representation and behavior. While data statistics shape the geometry of neural representations  
889 ([Merullo et al., 2025](#); [Karkada et al., 2026](#); [Prieto et al., 2026](#)), this fails to explain how geomet-  
890 ric structure is formed for out-of-distribution inputs. For example, our in-context learning tasks  
891 (Sec. A) have synthetically defined geometries that imbue tokens with contextual meaning that is  
892 wildly different from the meaning learned during training. As such, the model must form novel  
893 representations and produce novel behaviors. The fact that we are able to establish a shared geometry  
894 for representation and behavior in these novel in-context learning tasks suggests that regardless of  
895 how training data statistics inform the geometries seen in the model, the output behavior is now  
896 computationally constrained by the activation geometry (see the contemporary work by [Yocum et al.  
897 \(2025\)](#) for a formalization of this claim).

898 **Intrinsic coordinates of representation manifolds as units of causal analysis.** [Mueller et al.  
899 \(2024\)](#) frames the field of mechanistic interpretability as being on a quest to discover a primitive unit  
900 of representation best suited for the causal analysis of neural network internals. Causal abstraction  
901 provides a theoretical framework for defining such units of analysis ([Geiger et al., 2025a,b](#)), however  
902 [Sutter et al. \(2025\)](#) point out that allowing arbitrarily complex units admits degenerate solutions.  
903 Our work suggests a path toward both answering [Mueller et al. \(2024\)](#) and addressing the problem  
904 identified by [Sutter et al. \(2025\)](#): the appropriate units of causal analysis are intrinsic coordinates on  
905 manifolds in activation space, and fitting these manifolds to naturally occurring activations provides a  
906 constraint that helps rule out degenerate solutions (cf. [Grant et al. 2026](#)).

## 907 E Future Work and Limitations

908 The goal of our paper was to understand the role of geometry in neural networks and, subsequently,  
909 use this understanding to concretize what it means to steer model behavior via representations. We  
910 have shown that the geometry of neural network representations provides a blueprint for effective  
911 control. When interventions respect the geometry of activation space, the change in behavior is  
912 smooth and coherent; when they ignore it, they risk producing states with no natural behavioral  
913 counterpart. While we believe our results have enabled significant progress towards the motivating  
914 goals, there are remaining limitations that need to be addressed in future work.

- 915 • **Expanding experimental validation to more complex domains.** To illustrate our arguments,  
916 we focused on simple settings for which the concept of interest had a well-defined domain (e.g.,  
917 weekdays), and the expected task outputs are the concepts themselves, therefore the conceptual  
918 geometry is directly displayed in the outputs. To further validate the claims posited in this  
919 paper, future work is needed to explore more abstract concepts, e.g., refusals (Arditi et al., 2024),  
920 sycophancy (Vennemeyer et al., 2025), and persuasion (Costello et al., 2026). For such concepts,  
921 output behavior will likely reflect conceptual structure in subtler ways, and it remains to be tested  
922 whether the conceptual geometry of such concepts can be inferred from behavior and related to  
923 representations as we did in this work. Moreover, in these more complex cases, it is unclear what  
924 are the right primitives for a representational account; we may need a notion of dynamics over  
925 manifolds, a view of representation as an aggregation of several geometric structures (similar to  
926 results seen by Fel et al. (2025b) in a vision context), or perhaps an altogether different object. Even  
927 if geometry is the right substrate to work with, we emphasize the simplicity of our domains allowed  
928 us to easily isolate the target concept’s geometry via synthetic, template-based text. Moving to  
929 more complex scenarios will require isolating the geometry of concepts from in-the-wild data.
- 930 • **Moving from token to sequence-level outputs.** Another way in which our tasks are simplified  
931 is our focus on the next-token distribution. This makes analysis feasible and helps avoid the  
932 combinatorial complexity involved in studying multi-token sequences. The obvious way to expand  
933 from our work’s token-level focus to sequence-level focus involves formalizing arguments in the  
934 language of “beliefs”, i.e., latent variables underlying the posterior predictive induced by a model  
935 in response to an input (Bigelow et al., 2023, 2025; Wurgaft et al., 2025). Correspondingly, what  
936 we expect to see via geometry-aware interventions is the nature of output sequences produced by  
937 a model will change as we perform steering: e.g., navigating the geometry of sycophancy (were  
938 it to exist) should allow us to alter the extent or type of sycophancy exhibited in model outputs;  
939 however, this property will be latent, rather than a concrete token-level change.
- 940 • **Fitting the geometry.** While we used a specific protocol to fit the observed geometries (Bookstein,  
941 1989), we note there is a rich literature on fitting low-dimensional manifolds (Coifman & Lafon,  
942 2006; Brand, 2002; Schölkopf et al., 1997; Roweis & Saul, 2000; Jones, 2024; Jones & Lanners,  
943 2026; Meilă & Zhang, 2023). Critically, beyond just fitting the manifold, what we seek is an  
944 operator that allows us to navigate the manifold. For the domains analyzed in this work, we have  
945 ground-truth knowledge about how different states of the concept relate to each other, which allows  
946 us to define intrinsic coordinates for spline fitting. An unsupervised protocol would however  
947 significantly broaden the applicability of our methods.
- 948 • **Manipulating intermediate algorithmic variables.** All of our experiments are about manipulating  
949 the output behavior of neural networks directly. However, the most interesting control protocols  
950 will require manipulating intermediate quantities that mediate the flow of information from input to  
951 output, e.g., an image model determining the shape of an object in service of predicting its weight.

## 952 F Experimental Details for Language Tasks

953 This Appendix describes the procedures behind the experiments of Sections §2, §3, and §A. The  
954 following section (Appendix §G) will provide details for the mountain-car experiment in §B.

### 955 F.1 Tasks and Datasets

956 **Natural domain tasks.** We use four natural-domain addition tasks: weekdays and months (cyclic),  
957 and letters and ages (sequential). The full templates and entity sets are listed in Table 1. For each

958 task, we enumerate every (entity, increment) pair whose result lies in the task’s target set, dropping  
 959 pairs whose result would fall outside it (e.g. letters past Z, or ages outside [10, 100]). The reported  
 960 activations and output distributions are computed at the answer-token position, and concept centroids  
 961 are obtained by averaging across all prompts whose ground-truth result is the same value of  $\mathcal{Z}$ .

Task	Template	Entities	Increments	Structure	$ \mathcal{D} $
Weekdays	Q: What day is {k} days after {entity}?\nA:	Monday, ..., Sunday (7)	one, ..., seven	cyclic	49
Months	Q: What month is {k} months after {entity}?\nA:	January, ..., December (12)	one, ..., seven	cyclic	84
Letters	Consider letters in the alphabet. Starting at letter {entity}, we increment by {k}. The result is letter	C, ..., Z (24)	one, two	sequential	48
Ages	Alice is {entity} years old. Bob is {k} years older than Alice. Q: How old is Bob?\nA: Bob is	1, ..., 99	1, ..., 10	sequential	909

Table 1: Natural-domain arithmetic tasks. For cyclic domains, results wrap around the modulus (7 for weekdays, 12 for months); for sequential domains, (entity, increment) pairs whose result falls outside the target set are filtered. The last column reports the dataset size  $|\mathcal{D}|$  after this filter.

962 **In-context learning of representations.** For the multi-dimensional setting of §A, we use the  
 963 in-context learning of representations (ICLR) family of Park et al. (2025b): arbitrary tokens corre-  
 964 sponding to nouns (e.g., "film", "rain") are assigned to the nodes of a graph, and prompts are random  
 965 walks on that graph. We study a  $5 \times 5$  grid and a  $9 \times 9$  cylinder, with random walks of 2048 entity  
 966 tokens. As in Park et al. (2025b)’s setup, the random walks we sample do not allow backtracking,  
 967 which we find aids models in learning the underlying structure.

## 968 F.2 Model, Intervention Site, and Output Distribution

969 We investigate Llama 3.1 8B (Touvron et al., 2023) activations at layer 28 in bfloat16 for all tasks. All  
 970 interventions are performed on the residual stream at the last-token position. We chose to examine a  
 971 late layer of the model to ensure that concept geometries are fully computed.

972 For an input  $x$ , the output distribution  $p(x) \in \mathcal{Y}$  used throughout the paper is constructed as follows:  
 973 we softmax over the full vocabulary logit distribution and aggregate probability mass over each  
 974 concept value’s variant token spellings (e.g. the tokens ‘ Monday’, ‘Monday’, and ‘monday’ are  
 975 all summed into the Monday entry). The remaining probability mass on tokens not associated with  
 976 any concept value is collected into a single ‘other’ bin, yielding a distribution on the open simplex  
 977  $\Delta^{|\mathcal{Z}|}$  over  $|\mathcal{Z}| + 1$  classes.

## 978 F.3 Fitting the Activation Manifold $\mathcal{M}_h$

979 To identify the activation manifold  $\mathcal{M}_h$ , we first obtain points in full activation space and transform  
 980 them into a 64-dimensional subspace obtained via PCA over the activations  $h(x)$  across all prompts in  
 981 the task. The manifold lives entirely in the 64-dimensional PCA subspace; the orthogonal complement  
 982 is preserved during all subsequent interventions (§F.6).

983 We compute concept centroids  $c_i$  as the mean of the projected activations across all prompts whose  
 984 ground-truth result equals the  $i$ -th concept value, and fit a smooth interpolant through them. For  
 985 the four natural-domain tasks the interpolant is a one-dimensional cubic spline (Reinsch, 1967): a  
 986 natural cubic spline (with vanishing second derivatives at the endpoints) for the sequential tasks  
 987 (letters, ages), and a periodic cubic spline for the cyclic tasks (weekdays, months) so that the curve  
 988 closes smoothly. For the sequential tasks we use the ground-truth ordinal index of each concept  
 989 as its intrinsic coordinate. For the cyclic tasks the centroids form a near-circular loop in the top  
 990 two principal components of the activation subspace, so we instead derive the intrinsic coordinate  
 991  $\theta = \text{atan2}(\text{PC}_2, \text{PC}_1)$  in an unsupervised manner.

992 The interpolant for the ICLR tasks is a thin-plate spline (TPS; Duchon (1977); Bookstein (1989)), a  
 993 multi-dimensional generalisation of the cubic spline which minimizes the bending energy  $\int \|\nabla^2 f\|^2$ .  
 994 Thin-plate splines map points in a lower-dimensional intrinsic space to the full ambient space. The  
 995 TPS parameterisation requires a choice of intrinsic coordinates for the centroids. We use the ground-  
 996 truth graph coordinates of each node in the ICLR task as intrinsic coordinates. Both the grid and  
 997 cylinder tasks use the standard TPS kernel  $r^2 \log r$ ; for the cylinder, which has both a linear and  
 998 a periodic dimension, we additionally apply a ghost-point procedure where each control point is  
 999 duplicated at one period above and below its  $\theta$  value (and we drop the linear-in- $\theta$  polynomial column)  
 1000 to enforce closure across the periodic dimension. In every case the spline interpolates the centroids  
 1001 exactly, so  $\mathcal{M}_h$  passes through every  $c_i$ .

#### 1002 F.4 Fitting the Behavior Manifold $\mathcal{M}_y$

1003 Behavior centroids  $b_i = \bar{p}_i$  are computed analogously to the activation centroids, by averaging the  
 1004 model’s output distributions across all prompts whose ground-truth result equals the  $i$ -th concept  
 1005 value. Because the probability simplex is not a proper metric space, we map each centroid into  
 1006 Hellinger coordinates via  $b_i \mapsto \sqrt{b_i}$ , placing it on the non-negative orthant of the unit  $\ell_2$  sphere in  
 1007  $\mathbb{R}^{|\mathcal{Z}|+1}$ . We then fit the same family of splines used for  $\mathcal{M}_h$  to the Hellinger-embedded centroids:  
 1008 a 1D cubic spline (natural or periodic) for the natural-domain tasks, and a thin-plate spline for the  
 1009 ICLR tasks. Analogous to the activation manifold, the fit passes exactly through every  $\sqrt{b_i}$  and we  
 1010 do not apply a smoothing penalty.

1011 The spline is fit in Euclidean space, but valid  $\sqrt{b_i}$  points lie on a curved sphere. A naive fit to their  
 1012 ambient coordinates would leave the sphere between centroids, and an off-sphere vector does not  
 1013 square to a valid distribution. We therefore fit the spline in the *tangent plane* of the sphere at a base  
 1014 point  $b_*$  – a flat space that touches the sphere at  $b_*$  – and lift back to the sphere at decode time.

1015 We take  $b_*$  to be the Euclidean mean of  $\{\sqrt{b_i}\}$ , re-normalized to unit length; because every  $\sqrt{b_i}$  lies  
 1016 in the non-negative orthant, so does  $b_*$ . The *log-map*  $t_i = \log_{b_*}(\sqrt{b_i})$  projects each centroid onto the  
 1017 tangent plane: it returns a vector whose direction points from  $b_*$  along the geodesic to  $\sqrt{b_i}$  and whose  
 1018 length equals that geodesic distance. We fit the spline to the tangent vectors  $\{t_i\}$ . To decode at a  
 1019 query coordinate  $u$ , we evaluate the spline to obtain a tangent vector  $t$  and apply the *exponential map*  
 1020  $\exp_{b_*}(t)$ , the inverse of the log-map: it walks distance  $\|t\|$  along the geodesic on the sphere starting  
 1021 at  $b_*$  in direction  $t$ . The result is unit-norm by construction, so  $\mathcal{M}_y$  stays on the sphere everywhere,  
 1022 and because  $\exp_{b_*} \circ \log_{b_*}$  is the identity on the sphere, the decoded curve passes through every  $\sqrt{b_i}$   
 1023 exactly.

#### 1024 F.5 Geodesic Distances and the Isometry Test

1025 To compute the geodesic distance  $d_{\mathcal{M}_h}(c_i, c_j)$  between two concept centroids, we discretize the  
 1026 line segment between  $s^{-1}(c_i)$  and  $s^{-1}(c_j)$  in intrinsic coordinates into 150 equal sub-intervals,  
 1027 decode each waypoint through  $s$ , and accumulate consecutive ambient distances. Each waypoint  
 1028 therefore lies on  $\mathcal{M}_h$  by construction, and the resulting arc length is measured in the 64-dimensional  
 1029 PCA subspace in which the manifold lives, with the Euclidean norm as the ambient norm. For  
 1030 the behavior manifold  $\mathcal{M}_y$  we follow the same procedure but compute distances in the full sqrt-  
 1031 probability ambient space  $\mathbb{R}^{|\mathcal{Z}|+1}$  rather than a PCA-reduced subspace, using the Hellinger distance  
 1032  $d_H(p, q) = \frac{1}{\sqrt{2}} \|\sqrt{p} - \sqrt{q}\|_2$  directly on the sqrt-embedded waypoints.

1033 The isometry score reported in §2 is the Pearson correlation between the upper-triangular entries of  
 1034 the resulting pairwise distance matrices. We augment the  $W$  centroid vertices with  $K$  interior points  
 1035 sampled at equally spaced fractions of the u-space geodesic between each centroid pair, decoded onto  
 1036 the manifold via  $s$  so that every vertex lives on  $\mathcal{M}_h$  or  $\mathcal{M}_y$ . We choose  $K$  so that the vertex set is  
 1037 dense enough to probe the geometry between centroids:  $K = 4$  for weekdays ( $W = 7$ );  $K = 1$  for  
 1038 months ( $W = 12$ ), alphabet ( $W = 24$ , letters  $C-Z$ ), and the grid  $5 \times 5$  task ( $W = 25$ ); and  $K = 0$   
 1039 for age ( $W = 91$ , ages 10–100) and the cylinder  $9 \times 9$  task ( $W = 81$ ), whose centroids are already  
 1040 dense. We then correlate every off-diagonal pair in the full vertex set except those whose two vertices  
 1041 lie on a common centroid-pair geodesic, since those distances are sub-arcs of the same geodesic  
 1042 and would inflate the correlation by construction. To visualise the resulting pairwise structure, we  
 1043 embed each distance matrix into three dimensions via classical multidimensional scaling, as shown  
 1044 in Figs. 2, 8, and 5.

## 1045 F.6 Steering Interventions

1046 For each pair of concept values  $(z_a, z_b)$ , we steer the model from the centroid  $c_a$  to the centroid  $c_b$   
1047 via a path of  $K = 50$  waypoints. We use a fixed set of base prompts sampled randomly from the  
1048 task’s input distribution (the prompts’ ground-truth results vary, and the same set is reused across all  
1049 pairs): 16 prompts for the natural-domain tasks and 5 for the ICLR tasks. At each waypoint  $\pi(t)$ ,  
1050 we intervene at the last-token residual-stream activation of the target layer and continue the forward  
1051 pass to obtain  $\mathbf{p}_{h \leftarrow \pi(t)}(x)$ . Every reported behavioral trajectory is the pointwise mean over the base  
1052 prompts. We use up to 50 randomly-sampled pairs per task. On the smaller-domain tasks where  
1053  $W \cdot (W - 1) < 50$ , all pairs are used.

1054 The two steering strategies differ both in how the waypoints  $\pi(t)$  are constructed and in what the  
1055 intervention replaces (Eqs. 1, 2). For *manifold* steering,  $c_a, c_b$  live in intrinsic coordinates and the  
1056 path is the manifold geodesic between them; at each waypoint we decode  $\pi(t)$  onto  $\mathcal{M}_h$  in the  
1057 64-dimensional PCA subspace, lift it back to the residual-stream basis via the PCA inverse, and  
1058 combine it with the prompt’s unchanged off-subspace residual – so the steered activation differs from  
1059 the base only in its top-64 PCA components. For *linear* steering,  $c_a, c_b$  are the raw activation centroids  
1060 in the full residual stream, the path is the straight line between them, and the entire residual-stream  
1061 activation is replaced by  $\pi(t)$  at each step.

## 1062 F.7 Naturalness Metric

1063 The cumulative output energy  $E_{\text{BC}}$  of §3.2 (Eq. 3) is computed as the sum of the Bhattacharyya  
1064 distances  $D_{\text{BC}}(\gamma(t), \mathcal{M}_y) = -\log \sum_i \sqrt{\gamma_i(t) q_i(t)}$  between the induced output distribution at each  
1065 of the  $K = 50$  waypoints along steering paths and the closest point  $q(t)$  on  $\mathcal{M}_y$ . We use the  
1066 Bhattacharyya distance because  $\mathcal{M}_y$  is fit in Hellinger geometry and the two are tightly related,  
1067  $D_{\text{BC}} = -\log(1 - d_H^2)$ , so  $D_{\text{BC}}$  stays inside the same geometry the manifold was constructed in. For  
1068 each of the up to 50 sampled centroid pairs, we average the per-waypoint cumulative sum across the  
1069 base prompts to obtain one scalar per pair, and report the mean and standard error of these per-pair  
1070 scalars.

## 1071 F.8 Pullback Optimization

1072 The pullback procedure of §3.3 consists of two stages. First, we fix a behavioral target by evaluating  
1073 the spline geodesic on  $\mathcal{M}_y$  between the two behavior centroids  $b_a$  and  $b_b$  at  $K = 20$  uniform fractions,  
1074 yielding a sequence of target distributions  $\hat{\mathbf{p}}_t \in \mathcal{M}_y$ . Second, we optimize an activation-space path  
1075  $\pi_h^{\text{pullback}}$  which, when used to intervene at each waypoint, induces a behavioral trajectory matching  
1076  $\hat{\mathbf{p}}_{0:K}$ .

1077 **Path Parameterization.** We parameterise  $\pi_h^{\text{pullback}}$  as a one-dimensional natural cubic spline  
1078 through 10 control vectors at uniform  $t$ -positions, all of which are optimisation variables. The  
1079 path is evaluated at the same  $K = 20$  uniform fractions used to generate the target. Each control  
1080 vector is restricted to the first 32 PCA components of the 64-dimensional subspace; the remaining 32  
1081 components, together with the orthogonal residual, are held at the base prompt’s activation values  
1082 during the intervention. We note that the linear and manifold-steering paths used as comparisons  
1083 span the full 64-dimensional subspace, so the pullback optimization is operating within a strictly  
1084 smaller search space.

1085 **Loss and optimizer.** The loss at each waypoint  $t$  is the squared Hellinger distance  
1086  $d_H^2(\mathbf{p}_{h \leftarrow \pi(t)}(x_n), \hat{\mathbf{p}}_t)$  between the induced output distribution and the target, averaged over 16  
1087 base prompts  $\{x_n\}$  sampled freshly per pair, each conditioned on ground-truth  $z_a$  (in contrast to the  
1088 steering setup of §F.6, where a fixed unfiltered set is reused across pairs). We minimize the sum of  
1089 these per-waypoint losses with L-BFGS using strong-Wolfe line search, running 50 outer steps with  
1090 up to 5 inner iterations each. The optimization is initialized by linearly interpolating between the  
1091 two centroids in the 32-dimensional subspace and then sampling the resulting line at the 10 control- $t$   
1092 positions. We stop early when the relative change in loss between two consecutive outer steps falls  
1093 below  $10^{-3}$ . We disable the path-norm regularizer for weekdays; on the other three natural-domain  
1094 tasks we add a small regularizer that penalizes deviations of  $\|\pi(t)\|$  from the linear interpolation

1095 between the endpoint centroid norms  $|c_a|, |c_b|$ . We use weight  $10^{-3}$  for age and  $5 \times 10^{-4}$  for months  
1096 and alphabet. This discourages the optimizer from drifting into a high-norm shortcut basin.

### 1097 **F.9 Pullback Recovery** $R^2$

1098 To compare the optimised pullback path  $\pi_h^{\text{pullback}}$  to the manifold-steering path  $\pi_h^*$  along  $\mathcal{M}_h$ , we  
1099 project both into the SVD basis of  $\pi_h^*$  that captures at least 99% of  $\pi_h^*$ 's variance. In this basis we  
1100 define the residual at each pullback waypoint as its orthogonal closest-point distance to  $\pi_h^*$ , and report

$$R^2 = 1 - \frac{\sum_t \|\pi_h^{\text{pullback}}(t) - \text{proj}_{\pi_h^*} \pi_h^{\text{pullback}}(t)\|^2}{\sum_t \|\pi_h^{\text{pullback}}(t) - \bar{\pi}_h^{\text{pullback}}\|^2}.$$

1101 The linear baseline used in the same comparison is the straight chord between  $c_a$  and  $c_b$  in the  
1102 64-dimensional PCA subspace—not the linear-steering trajectory after intervention. As in §F.7, the  
1103 values reported in §3.3 are mean  $\pm$  standard error across the per-pair scalars, with  $p$ -values from  
1104 paired  $t$ -tests against the linear baseline.

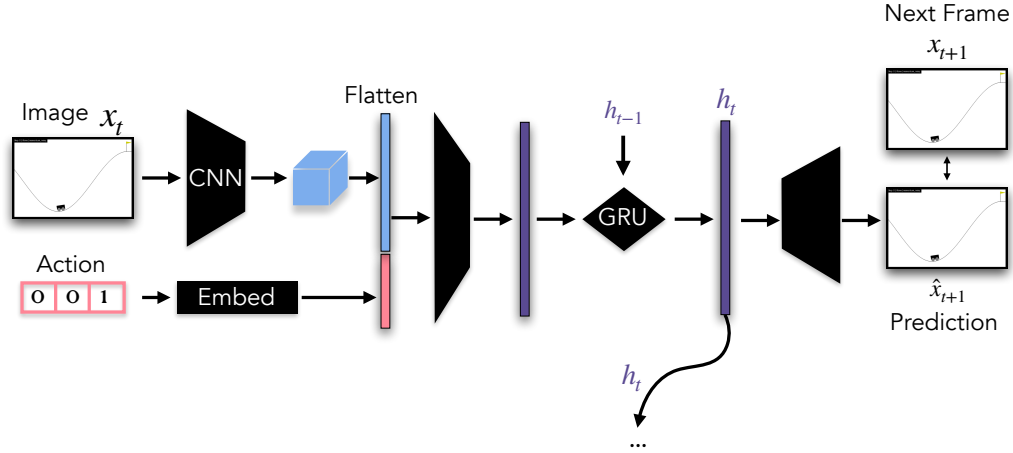


Figure 7: **Recurrent visual world-model architecture.** A convolutional encoder  $f_{\text{enc}}$  maps each  $128 \times 128 \times 3$  frame  $x_t$  to a layer-normalized latent  $z_t \in \mathbb{R}^n$  with  $n = 64$ . The discrete action  $a_t \in \{0, 1, 2\}$  is mapped to a learned embedding  $e(a_t) \in \mathbb{R}^{16}$ , concatenated with  $z_t$ , and fed to a GRU together with the previous hidden state  $\mathbf{h}_{t-1}$ . A convolutional decoder  $f_{\text{dec}}$  produces a residual image from the resulting hidden state  $\mathbf{h}_t$ , yielding the next-frame prediction  $\hat{x}_{t+1} = x_t + f_{\text{dec}}(\mathbf{h}_t)$ , supervised against the ground-truth frame  $x_{t+1}$ .

## 1105 G Experimental Details for the Vision Task

1106 This section contains additional details on the experiment from §B.

### 1107 G.1 Mountain Car.

1108 **Data collection.** To recover the encoder’s position manifold we harvest activations on 100 rollouts  
 1109 collected in MountainCar-v0 (max 200 steps per episode) under a mixed stochastic policy chosen to  
 1110 give broad coverage of the position-velocity state space. At the start of each episode we sample one of  
 1111 two policies: with probability 0.7, a *noisy momentum* policy that pushes in the direction of the current  
 1112 velocity but, at each step, replaces the action with a uniform random action with probability 0.4; with  
 1113 probability 0.3, an *oscillating square-wave* policy that alternates between full-left and full-right thrust  
 1114 on a fixed period sampled uniformly from  $\{5, \dots, 25\}$  steps. We then pass each rendered frame  
 1115 through the trained encoder, label the resulting activation with the underlying ground-truth position,  
 1116 and fit the manifold to this collection of position-labelled activations.

1117 **Manifold fitting.** To parameterize the activation manifold,  $\mathcal{M}$ , we partition the position range  
 1118 into  $B = 100$  bins, compute the mean encoder output per occupied bin, and fit a smoothing spline  
 1119  $\gamma_{\mathcal{M}}: [0, 1] \rightarrow \mathcal{M}$  through these means (one univariate spline per coordinate, weighted by the square  
 1120 root of bin counts to regularize sparse regions). We additionally verify via linear probing that the  
 1121 encoder representations  $z_t$  encode the ground-truth physics: a Ridge regression probe recovers  
 1122 position with  $R^2 \approx 0.95$  and velocity with  $R^2 \approx 0.90$ . A three-component PCA of the encoder  
 1123 outputs reveals the spline  $\gamma_{\mathcal{M}}$  as a curve that closely tracks the data manifold, while the chord  $\ell$   
 1124 between the same endpoints cuts through its interior.

1125 To parameterize the behavior manifold,  $\mathcal{M}_y$ , discretize  $\mathcal{Z}$  into  $B$  bins with centers  $\{\mu_b\}_{b=1}^B \subset \mathbb{R}^n$   
 1126 obtained by evaluating the activation-manifold spline at  $B$  evenly spaced positions,  $\mu_b = \gamma_{\mathcal{M}}(p_b)$ .  
 1127 The mapping to behavior is

$$F(z) = \text{softmax}\left(-\frac{\|z - \mu_b\|_2}{\tau}\right)_{b=1}^B \in \Delta^{B-1}, \quad (10)$$

1128 with temperature  $\tau = 0.5$ .  $F$  is a smooth, deterministic map from activations to position distributions.  
 1129 We use  $B = 128$ , which makes  $F$ ’s Jacobian full column-rank, ensuring the inverse problem has a  
 1130 locally unique solution. For each bin  $i$ , the natural centroid on the behavior manifold is the model’s

1131 average output distribution conditioned on samples in that bin:

$$b_i = \mathbb{E}[F(z) \mid \text{bin}(z) = i] \in \Delta^{B-1}. \quad (11)$$

1132 Because the bin grid is dense relative to the data manifold's intrinsic dimension,  $b_i$  is well-defined  
1133 for all bins. We embed each  $b_i$  in Hellinger coordinates  $h_i = \sqrt{b_i}$  on the unit sphere of  $\mathbb{R}^B$  and fit a  
1134 1D smoothing spline  $\gamma_{\mathcal{M}_y} : \mathcal{Z} \rightarrow \mathbb{R}^B$  through  $\{h_i\}$  parameterized by position. This is the behavior  
1135 manifold  $\mathcal{M}_y$ .

1136 **H Additional Results**

1137 **H.1 Natural Domain Arithmetic: Ages and letters**

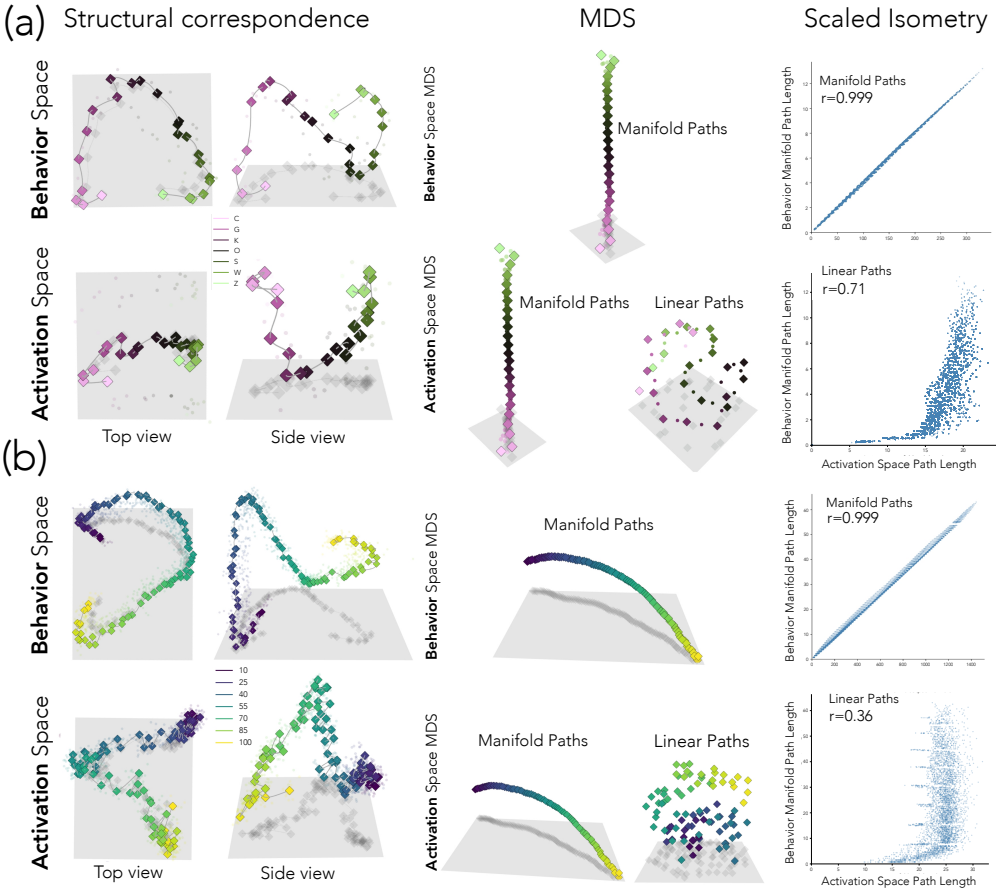


Figure 8: **Approximate isometry between activation and behavior manifolds for sequential concepts.** Manifolds (cubic splines) fit to activation and behavior (i.e., output distributions over concept tokens) spaces of Llama 3.1 8B. The letters (a) and ages (b) tasks consist of simple addition questions such as: What letter comes four letters after M?. Both activation and behavior manifolds show sequential structure (PCA visualization shown in left column). Furthermore, on-manifold distances in activation space show strong correlation with on-manifold distances in behavior space (right column), as well as a clear match via a multidimensional scaling (MDS) embedding (middle column). In contrast, linear distances in activation space show weaker correlations and warped or incoherent structure. These results demonstrate an approximate isometry between the activation and behavior space manifolds.

1138 **H.2 In-Context Learning of Representations.**

1139 In addition to results provided in the main text, we test a  $9 \times 9$  cylinder in the ICLR domains, and  
 1140 find that despite the added complexity of a periodic dimension and substantially more graph nodes,  
 1141 when Llama 3.1 8B is provided sufficient context (2048 tokens in this case), it reaches above 80%  
 1142 neighborhood accuracy (probability mass on valid neighbors). We fit a manifold and steer along this  
 1143 domain, finding that the result of factored control generalizes beyond the graph domain.

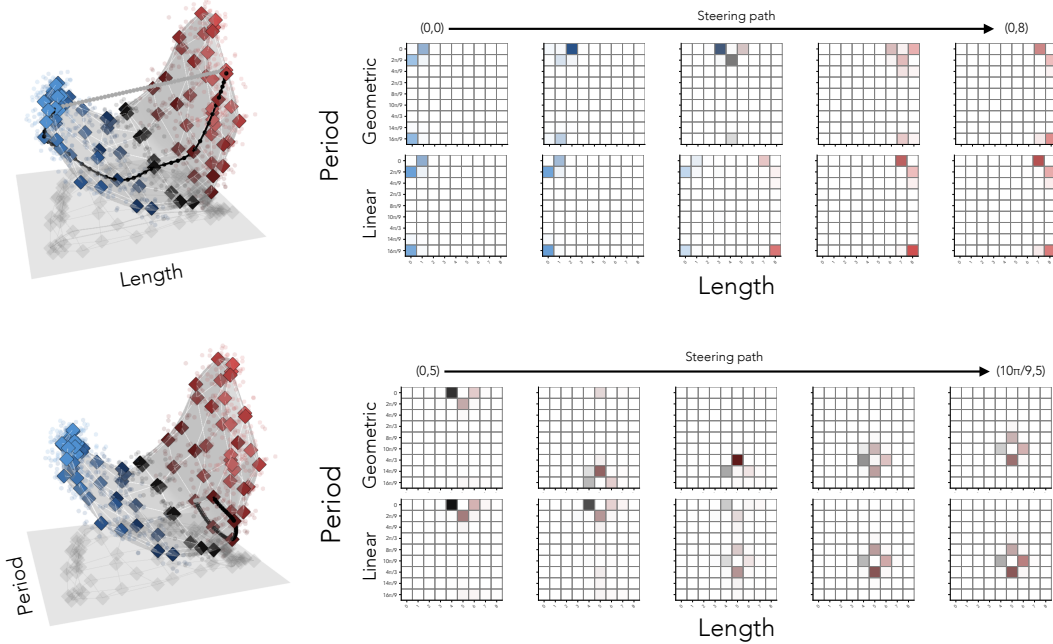


Figure 9: Results for in-context learning of representations on a  $9 \times 9$  cylinder domain. We find that, as in the grid domain, manifold steering achieves factored control: coherent steering of independent dimensions, while linear steering once again shows ‘teleportation’ behavior.

1144 **H.3 Manifold steering allows manipulation of uncertainty without loss of structure.**

1145 To examine multi-dimensional concepts in known domains, we partition weekday addition centroids  
 1146 by addition value (1–5, 6–10, 11–15, 16–20), revealing concentric circles along a second manifold  
 1147 dimension forming a cylinder-like structure (Fig. 11). Manifold steering along the circular dimension  
 1148 maintains ordered weekday transitions with increasing entropy per group. This suggests manifold  
 1149 geometry can serve as a handle for calibrating model confidence in a controlled fashion. The  
 1150 experiment was conducted with Llama 3.1 70B layer 70.

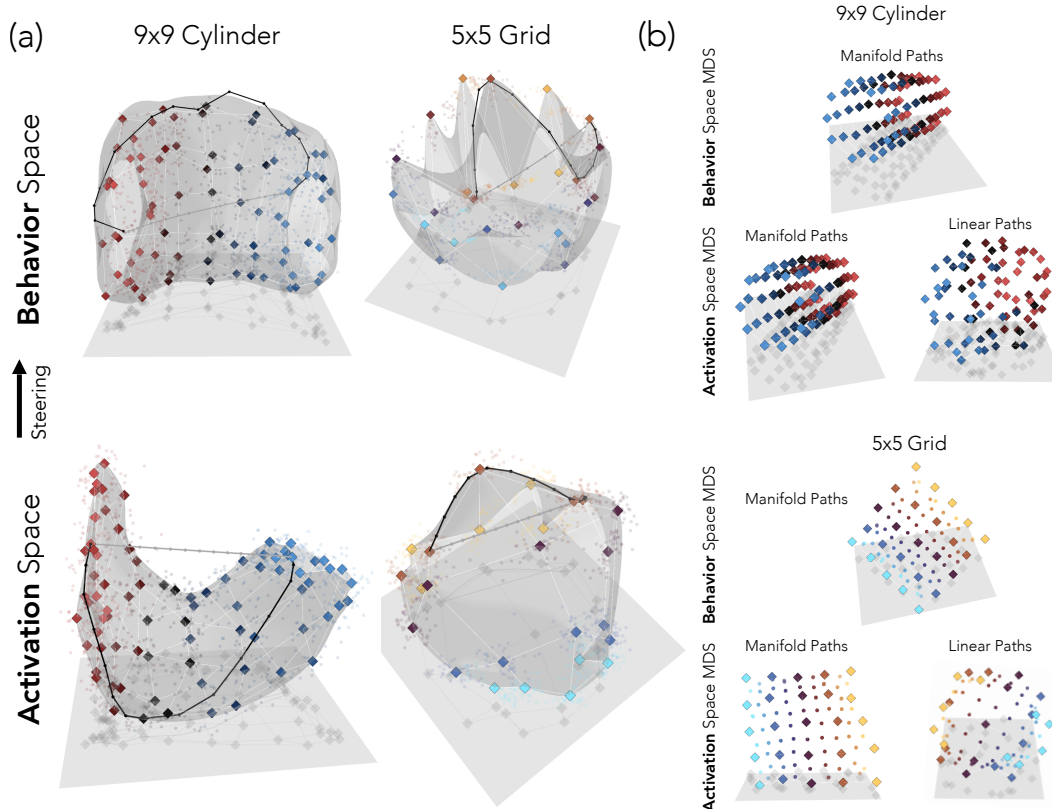


Figure 10: (a) Activation and behavior space paths for the  $5 \times 5$  Grid task and  $9 \times 9$  Cylinder. Similarly to the addition tasks with known concepts, we find that the manifold steering paths closely follow the behavior manifold  $\mathcal{M}_y$ . (b) Multidimensional scaling (MDS) embedding for linear and manifold distances in activation space and manifold distances in behavior space. As with the addition tasks with known concepts, manifold distances in activation space show a clear structural match to behavior space, whereas linear distances warp the structure.

#### 1151 H.4 Mountain Car

1152 **Structural Correspondence between  $\mathcal{M}_h$  and  $\mathcal{M}_y$ .** If  $\mathcal{M}_h$  encodes the model’s predictive distributions over  $\mathcal{Z}$ , then  $\mathcal{M}_h$  and  $\mathcal{M}_y$  should be approximately isometric—distances along one manifold  
 1153 should correlate with distances along the other. We test this by sampling  $W = 50$  anchor positions  
 1154 inside the shared parameter range and computing pairwise arc lengths along each manifold:  
 1155

$$d_{\mathcal{M}}(p_i, p_j) = \int_{p_i}^{p_j} \|\gamma'_{\mathcal{M}_h}(p)\|_2 dp, \quad d_{\mathcal{M}_y}(p_i, p_j) = \frac{1}{\sqrt{2}} \int_{p_i}^{p_j} \|\gamma'_{\mathcal{M}_y}(p)\|_2 dp, \quad (12)$$

1156 where the  $1/\sqrt{2}$  on the behavior side converts the Euclidean integral in Hellinger ambient space  
 1157 to Hellinger units. The Pearson correlation between  $\{d_{\mathcal{M}_h}(p_i, p_j)\}$  and  $\{d_{\mathcal{M}_y}(p_i, p_j)\}$  over all  
 1158  $\binom{50}{2} = 1225$  pairs is  $r = 0.996$ ; the chord distances used by linear steering correlate far less  
 1159 ( $r = 0.06$  between activation chord and behavior arc length), since chords cut across the encoder  
 1160 loop and are structurally divorced from the encoded conceptual geometry.

1161 **Pullback: Behavior space steering.** Having established the bottom-up direction in Fig. 6, i.e.,  
 1162 paths along  $\mathcal{M}_h$  produce behavior trajectories on  $\mathcal{M}_y$ , We now test the top-down direction: starting  
 1163 from a behaviorally-natural trajectory in  $\mathcal{M}_y$ , do we naturally recover an activation path that traces  
 1164  $\mathcal{M}_h$ ? For each endpoint pair  $(p_a, p_b)$  we construct the conformal behavior target  $\hat{\gamma}_\alpha$  via the procedure  
 1165 of §3.3 (geodesic on the simplex under cost  $c(p) = \exp(\alpha \cdot d_H(p, \mathcal{M}_y))$ ). We then optimize an

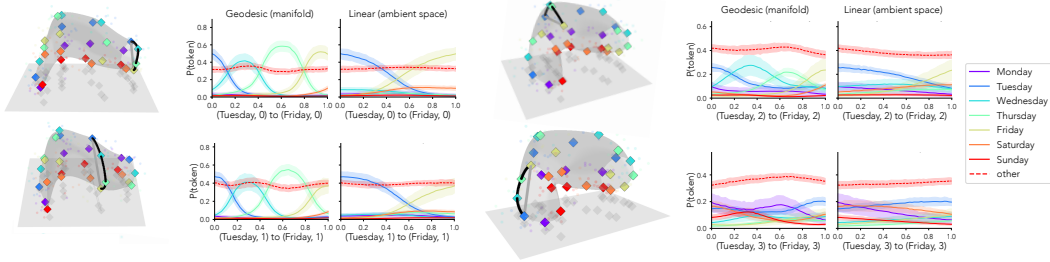


Figure 11: **Inducing greater uncertainty over a conceptual space and steering along it.** By increasing the addition value, we induce greater uncertainty in the model with respect to the right answer. Instead of grouping across addition values (as we do in Fig. 3), we visualize centroids by addition value, and find that these groups yield a series of circles organized into a curved cylinder-like shape. Manifold steering along the circle in the first three groups maintains ordered transitions, yet with increasing entropy in each group.

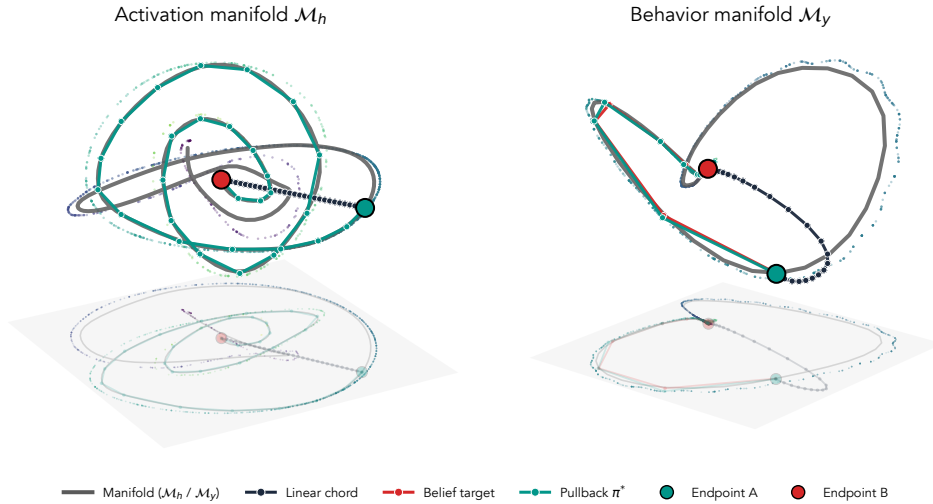


Figure 12: **Pullback from  $\mathcal{M}_y$  recovers  $\mathcal{M}_h$  in the visual world model.** **Left (Activation Space):** PCA visualization of the encoder representations, showing the geometric path along  $\mathcal{M}_h$ , the linear chord, and the pullback-optimized path  $\pi^*$  between endpoints  $p_A$  and  $p_B$ . Although initialized at the chord,  $\pi^*$  converges onto  $\mathcal{M}_h$ , closing the spiral loop traced by the encoder geometry and becoming nearly indistinguishable from the activation reference. **Right (Behavior Space):** PCA visualization of the corresponding trajectories pushed through the operator  $F$  (Eq. 10), shown in Hellinger coordinates. The conformal target  $\hat{\gamma}_\alpha$  tracks the behavior manifold  $\mathcal{M}_y$ , and the pushforward  $F(\pi^*)$  closely matches it, while the pushforward of the linear chord  $F(\ell)$  departs sharply, cutting across the simplex interior rather than following  $\mathcal{M}_y$ . Together, the two panels show that optimizing an activation path to match a behavior-manifold target recovers  $\mathcal{M}_h$  top-down: matching behavior along  $\mathcal{M}_y$  is sufficient to pull activations back onto  $\mathcal{M}_h$ , mirroring the pullback result from the language-model experiments.

1166 activation path  $\pi_\alpha = (v_0, \dots, v_K)$  in  $\mathbb{R}^n$  to minimize

$$L(\pi) = \sum_{t=0}^K \|\sqrt{F(v_t)} - \sqrt{\hat{\gamma}_\alpha(t)}\|_2^2. \quad (13)$$

1167 Following the language-model setup, all  $K + 1$  waypoints (including endpoints) are free parameters,  
 1168 initialized at the linear chord and optimized jointly via L-BFGS with strong-Wolfe line search. We  
 1169 use  $K = 30$  waypoints and run independent optimizations for each of 30 endpoint pairs.

1170 Across all 30 endpoint pairs, the pullback paths  $\pi_\alpha$  closely trace  $\mathcal{M}_h$  (Fig. 12): The mean Euclidean  
1171 distance from  $\pi$  to  $\mathcal{M}_h$ , averaged over waypoints and pairs:

linear chord: 2.22,      geometric ( $\mathcal{M}_h$ ): 0.20,      pullback: 0.29.

1172 The pullback path is at **95.4%** of the chord-to-geometric recovery and dominates the chord baseline  
1173 on 30/30 pairs. The aggregate degradation is concentrated on pairs with one endpoint at the extreme  
1174 wall position ( $p \approx -1.2$ ), where the encoder geometry has tighter curvature; on the remaining  $\sim 20$   
1175 pairs,  $\pi_\infty$  is essentially indistinguishable from  $\mathcal{M}_h$  itself. The  $\alpha$ -sweep traces the same family  
1176 of trajectories observed in the language-model experiments: at  $\alpha = 0$  the conformal target is the  
1177 unrestricted Hellinger geodesic on the simplex, and the recovered  $\pi_0$  leaves  $\mathcal{M}_h$  in order to match this  
1178 off-manifold target; as  $\alpha$  grows the target is pushed onto  $\mathcal{M}_y$  and the recovered  $\pi_\alpha$  correspondingly  
1179 tracks  $\mathcal{M}_h$ .