
Control and Predictivity in Neural Interpretability

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 For the goals of mechanistic interpretability, correlational methods are typically
2 easy to scale and use, and can provide strong predictivity of Neural Network (NN)
3 representations. However, they can lack causal fidelity which can limit their rel-
4 evance to NN computation and behavior. Alternatively, causal approaches can
5 offer strong behavioral control via targeted interventions, making them superior for
6 understanding computational cause and effect. However, what if causal methods
7 use out-of-distribution representations to produce their effects? Does this raise con-
8 cerns about the faithfulness of the claims that can be made about the NN’s native
9 computations? In this work, we explore this possibility of this representational
10 divergence. We ask to what degree do causally intervened representations diverge
11 from the native distribution, and in what situations is this divergence acceptable?
12 Using Distributed Alignment Search (DAS) as a case study, we first demonstrate
13 the existence of causally intervened representational divergence in interventions
14 that provide strong behavioral control, and we show that stronger behavioral control
15 can correlate with more divergent intervened representations. We then provide
16 a theoretical discussion showing sufficient ways for this divergence to occur in
17 both innocuous and potentially pernicious ways. We then provide a theoretical
18 demonstration that causal interventions typically assume principles of additivity,
19 calling into question the use of nonlinear methods for causal manipulations. Lastly,
20 for cases in which representational divergence is undesirable, we demonstrate how
21 to incorporate a counterfactual latent loss to constrain intervened representations
22 to remain closer to the native distribution. Together, we use our results to suggest
23 that although causal methods are superior for most interpretability goals, a com-
24 plete account of NN representations balances computational control with neural
25 predictivity, with the optimal weighting depending on the goals of the research.

26 **1 Introduction**

27 In the many recent developments in mechanistic interpretability, researchers have used a variety of
28 methods to defend claims about how Neural Networks (NNs) perform their computations. These
29 methods can be broadly categorized into two groups based on their measures of success. The first
30 group observes and analyzes NN activations derived from the NN’s naturally occurring neural activity
31 without causally influencing the activity. The second camp focuses on causally manipulating the
32 neural activity in an effort to determine cause and effect relationships in the NN’s computations. We
33 will broadly refer to the former methods as correlational and the latter as causal.

34 Popular, exemplary approaches of the correlative variety are methods like Sparse Auto-Encoders
35 (SAEs) and Principal Component Analysis (PCA) [32, 6, 1, 5, 9]. Both of these methods attempt to
36 decompose NN latent activity into a linear sum of features, or vector directions, whose variability
37 corresponds to some interpretable attribute that can be used to interpretably construct the distribution
38 of native latent activity. These approaches can lead to decompositions that closely match/predict the

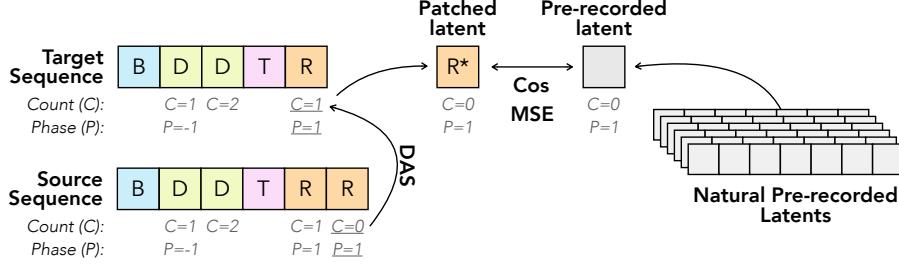


Figure 1: Figure showing counterfactual latents in counting task. The squares represent latent vectors produced from the recurrent state and the displayed input tokens. The C indicates the value of the Count of the causal abstraction, the P indicates the phase. The counterfactual latent (CL) vectors are the naturally occurring latent vectors that possess the same causal values as the post-intervention latent vector. These CL vectors can be recorded from the forward pass on curated input data.

39 native latent distribution and have been shown to provide insight into the underlying NN mechanisms
 40 and neural structures in various types of NNs [5, 9]. However, these methods often lack causal
 41 fidelity [8] and can exhibit a number of undesirable traits such as failing to identify atomic units
 42 [23], producing brittle concept representations [24], and under performing on causal mediation tasks
 43 and out-of-distribution probing [2, 20]. [22] have even gone so far as to doubt the foundations
 44 of correlational neural analyses as a means of understanding NNs due to disconnects between the
 45 strength of a feature and its importance for computation [21]. Similar doubts have surfaced in the
 46 past [19].

47 Alternatively, some popular methods of the causal variety include activation patching and Distributed
 48 Alignment Search (DAS) [12, 14, 39, 38, 28, 29, 4]. These methods use causal relevance, or behavior,
 49 as their main measures of success, where the goal is to controllably affect NN outputs by manipulating
 50 internal representations in order to make cause-and-effect claims about the internal mechanisms.
 51 These methods are, by definition, superior for determining cause and effect relationships in neural
 52 circuitry, which is perhaps of utmost importance for making mechanistic claims [31, 11, 10, 22].

53 What if these causal approaches, however, rely on intervened representations that strongly diverge
 54 from the distribution of native neural activity? Would that change what the methods tell us about
 55 the NN’s natural mechanistic components? Some activation patching examples use features with
 56 values that are multiplied by 10-15x [25], raising questions of what these experiments say about
 57 the NN’s native mechanisms. In what situations and to what degree is it okay for causal methods
 58 to deviate from the native distribution? In cases that these deviations occur, are there ways that we
 59 can mitigate the deviation? In one sense, any perspective that is critical of analyzing representations
 60 alone is suggesting that representational divergence is okay for understanding NN mechanisms, but
 61 how much deviation and in what circumstances is this deviation tolerable?

62 In this work, we provide theoretical and empirical insight on these issues. We first empirically
 63 explore representational divergence in a case study using DAS on Gated Recurrent Units (GRUs)
 64 [3]. We explore how choices of causal abstractions and alignment functions (AFs) affect behavioral
 65 control and faithfulness to native neural activity. We then theoretically show sufficient ways in which
 66 representational divergence can occur, and we provide a discussion on how these deviations can be
 67 okay for many mechanistic claims. Next, we show how causal interventions rely on principles of
 68 additivity, raising concerns of what nonlinear interpretability methods show us. Lastly, for cases in
 69 which researchers do care about representational divergence, we provide a demonstration on how to
 70 use a counterfactual auxiliary loss on intervened latent vectors to mitigate representational divergence
 71 between intervened and native latent vectors.

72 We summarize our contributions as follows:

- 73 1. We empirically explore how causally manipulated representations can control behavior,
 74 but can deviate from native neural activity. Furthermore, we show that greater behavioral
 75 accuracy can correlate with greater representational divergence.
- 76 2. We provide a theoretical discussion on two sufficient ways for representational divergence
 77 to occur, showing that it can arise from dormant-null subspace interactions introduced in

78 [26] and from innocuous co-variation of causal subspaces. We use this result to suggest that
 79 representational divergence is expected and okay for most mechanistic claims.

- 80 3. We provide a theoretical demonstration that causal interventions rely on principles of
 81 additivity, which obfuscates claims of causal interpretability made with non-linear methods.
 82 4. Lastly, we show how to use a counterfactual latent loss [17] to mitigate the divergence
 83 between the intervened and native neural distributions, bridging the gap between behavioral
 84 control and neural predictivity in causal methods.

85 We use our findings to inform a discussion on what our goals are as interpretability researchers. We
 86 suggest that a complete understanding of NN representations includes both control and predictivity of
 87 NN activity, while the notion of success with these criteria depends on the purpose of the research.

88 2 Background and Related Work

89 The majority of our empirical analyses use a specific form of activation patching [13, 37, 38, 28]
 90 known as Distributed Alignment Search (DAS) [12, 14, 39]. We provide a background on DAS here.

91 2.1 DAS Formulation

92 DAS is a framework for causally testing the degree of alignment between an NN’s latent vectors
 93 and variables from causal abstractions (CAs) (e.g. computer programs or directed acyclic graphs).
 94 DAS does this by testing the hypothesis that a latent state vector, $h \in R^{d_m}$, within an NN can be
 95 transformed into a vector $z \in R^{d_m}$ that consists of orthogonal subspaces encoding interpretable
 96 variables from CAs. The transformation is performed by a learnable, invertible *Alignment Function*
 97 (AF), \mathcal{A} , as follows: $z = \mathcal{A}(h)$ [18]. The benefit of this transformation is that it allows us to formulate
 98 the NN’s neural activity in terms of interpretable variables, and it allows us to manipulate the value of
 99 each variable without affecting the values of the others. We will refer to the space of h as the *native*
 100 *vector space* and that of z as the *aligned vector space*.

101 Concretely, for a given CA with variables $\text{var}_i \in \{\text{var}_1, \text{var}_2, \dots, \text{var}_n\}$, DAS tests the hypothesis that
 102 z is composed of subspaces $\vec{z}_{\text{var}_i} \in R^{d_{\text{var}_i}}$ corresponding to each of the variables from the CA. We
 103 include a causally irrelevant subspace, $\vec{z}_{\text{extra}} \in R^{d_{\text{extra}}}$, to encode extraneous, functionally irrelevant
 104 activity.

$$\mathcal{A}(h) = z = \begin{bmatrix} \vec{z}_{\text{var}_1} \\ \vec{z}_{\text{var}_2} \\ \vdots \\ \vec{z}_{\text{var}_n} \\ \vec{z}_{\text{extra}} \end{bmatrix} \quad (1)$$

105 Each $\vec{z}_{\text{var}_i} \in R^{d_{\text{var}_i}}$ is a column vector of potentially different lengths. We refer to d_{var_i} as the
 106 *subspace size* of var_i , and all the subspace sizes together satisfy the relation $d_{\text{extra}} + \sum_{i=1}^n d_{\text{var}_i} = d_m$.
 107 Under this assumption, the value of a single causal variable encoded in h can be freely exchanged by
 108 performing an interchange intervention defined as follows:

$$h^v = \mathcal{A}^{-1}((1 - D_{\text{var}_i})\mathcal{A}(h^{trg}) + D_{\text{var}_i}\mathcal{A}(h^{src})) \quad (2)$$

109 Where $D_{\text{var}} \in R^{d_m \times d_m}$ is a manually defined, block diagonal, binary matrix that defines the subspace
 110 size d_{var_i} . Each D_{var_i} has a set of d_{var_i} contiguous ones along its diagonal to isolate the dimensions
 111 that make up \vec{z}_{var_i} . h^{src} is the *source vector* from which the subspace activity is harvested, h^{trg}
 112 is the *target vector* into which the harvested activity is substituted/patched, and h^v is the resulting
 113 intervened vector that we use to replace h^{trg} in the model’s processing. This allows the model to
 114 make predictions using a different value of variable var_i assuming a successful intervention.

115 To train the AF, DAS uses *counterfactual behavior* from the pre-defined CA to create intervention
 116 data that can be used as training labels for the model’s processing conditioned on the intervened
 117 latent representation after an intervention. Counterfactual behavior for a given state of a CA and
 118 its context is the behavior that would have occurred had a causal variable taken a different value
 119 and everything else remained the same. We can simulate counterfactual behavior by freezing the
 120 state of the environment, changing one or more values in the CA, and using the CA to generate new

121 behavior in the same environment using the new variable values. We can then use the counterfactual
122 behavior as training labels to train the AF after each intervention while keeping the model parameters
123 frozen. We train the AF to convergence and then use fresh intervention data to evaluate the robustness
124 of the AF and to make claims about the NN’s internal mechanisms. The model’s accuracy on the
125 counterfactual behavior following each intervention is referred to as the Interchange Intervention
126 Accuracy (IIA).

127 In our experiments, we train the same AF on interventions for all causal subspaces including the
128 extraneous subspace. We consider a trial correct when the model correctly predicts all deterministic
129 tokens using the argmax over logits. We report the proportion of trials correct as the IIA. See
130 Appendix A.2 for further detail.

131 **DAS Alignment Functions:** We consider three types of AFs in this work.

- 132 1. *Orthogonal Alignment Functions (OAFs):* $\mathcal{A}(h) = Qh$ and $\mathcal{A}^{-1}(h) = Q^{-1}h$ where
133 $Q \in \mathbb{R}^{d_m \times d_m}$ is an orthogonal matrix.
- 134 2. *Linear Alignment Functions (LAFs):* $\mathcal{A}(h) = W(h+b)$ and $\mathcal{A}^{-1}(z) = W^{-1}z - b$ where W
135 is a symmetric invertible matrix [18]. See Appendix A.3 for details on how W is constructed.
- 136 3. *Reverse Resnet Alignment Functions (RRAFs):* $\mathcal{A}(h) = RevRes(h)$ and $\mathcal{A}^{-1}(z) =$
137 $RevRes^{-1}(z)$ where $RevRes$ is a reversible residual network [15, 35]. We use 3 layers
138 with no changes in dimensionality.

139 3 Methods

140 The majority of this work consists of DAS analyses performed on Gated Recurrent Unit recurrent
141 neural networks (GRUs) autoregressively trained on sequence-based tasks. We use DAS to align
142 these NNs to Causal Abstractions (CAs) by performing interchange interventions on the GRUs’
143 representations. We consider a single numeric equivalence task which has been used in prior work on
144 human cognition [16, 7] and alignment functions [18, 17].

145 3.1 Numeric Equivalence Task:

146 This task consists of a sequence of tokens produced by an environment. Each sequence starts with
147 a beginning of sequence token, B, and ends with an end of sequence token, E. After the B token,
148 the environment presents some number of demonstration (demo) tokens that are each sampled with
149 replacement from the set $\{D_a, D_b, D_c\}$. The task is to produce the same number of response (R)
150 tokens as D tokens, and end with the E token. The environment signals the end of the D tokens
151 by producing a trigger (T) token. The number of D tokens at this point is referred to as the *object*
152 *quantity* for the trial, which is uniformly sampled from 1 to 20 at the beginning. The set of possible
153 tokens includes $\{B, D_a, D_b, D_c, T, R, E\}$. An example sequence with an object quantity of 2 is: "B
154 $D_c D_a T R R E$ ". Each trial is considered correct when all deterministic tokens are correctly predicted.
155 During the model training, we include all token types in a NTP cross entropy loss, even though the *D*
156 and *T* tokens are unpredictable.

157 3.2 Model Architectures

158 In this work we consider Gated Recurrent Unit (GRU) Recurrent Neural Networks (RNNs) [3] that
159 are autoregressively trained to perform the Numeric Equivalence task. We train 3 model seeds for
160 each task variant up to $> 99.99\%$ accuracy on both training and validation data and freeze the weights
161 before analysis and interpretation. The GRUs have a dimensionality of 128. We perform all DAS
162 analyses on the output of the GRU recurrent cell, denoted h . We leave further details of the GRU
163 recurrent cell to Appendix A.1 and the referenced paper.

164 3.3 Causal Abstractions (CAs)

165 In this work, we evaluate 2 different CAs using DAS. We briefly describe them here and offer
166 Algorithms 1 and 2 in the appendix in addition to intervention data samples in Appendix A.3.1. Also,
167 refer to Figure 4 for a visual depiction of the CAs.

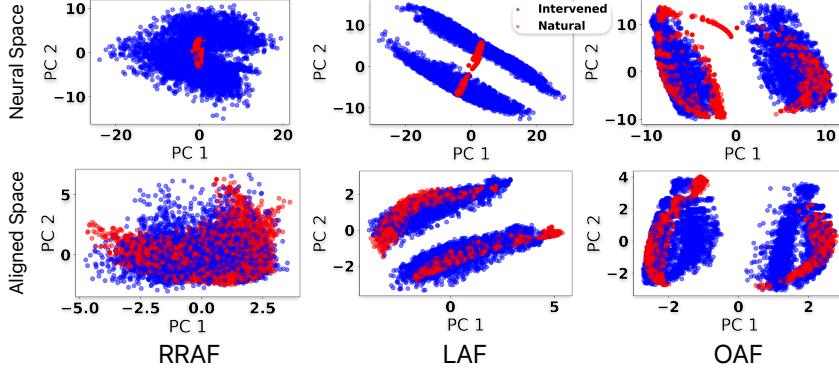


Figure 2: Principal component projections of the native latent states (in red) and the intervened latent states (blue). The top row shows the vectors in the native latent space. The bottom row shows them in the aligned space in which the extraneous subspace has been set to 0.

168 **Up-Down Program:** uses a single numeric variable, called the **Count**, to track the difference between
 169 the number of demo tokens and resp tokens at each step in the sequence. It also contains a **Phase**
 170 variable to determine whether it is in the demo phase—counting up—or response phase—counting
 171 down. The program ends when the Count is equal to 0 during the response phase.

172 **Increment-Up Program:** this program uses progress along an interval from 0 to 1 to track quantities.
 173 To do this, it first increments a Progress variable by the value of an Increment variable. The value of
 174 the Increment is initially set to $\frac{1}{\max \text{ count}}$ (in our case $\frac{1}{20}$). The value of the Progress variable is then
 175 incremented with each new demo token to track the object count (number of demo tokens) during
 176 the demo phase. Upon encountering the trigger token, a new value of the Increment is calculated as
 177 the inverse of the Progress divided by the max count: $\text{Increment} = \frac{1}{\text{Progress}} \div \max \text{ count}$. The
 178 Progress is then reset to 0. The new value of Increment is now equal to $\frac{1}{\text{obj count}}$ and is used as a step
 179 size to increment the Progress variable with each new response token. The program finishes when the
 180 Progress variable is greater than 1.

181 3.4 Counterfactual Latent Auxiliary Loss

182 To encourage intervened representations to be more similar to the native distribution of NN repre-
 183 sentations, we re-purpose the counterfactual latent auxiliary loss from [17]. This auxiliary objective
 184 relies on *Counterfactual Latent (CL) vectors* as vector objectives. CL vectors are defined as vectors
 185 that encode the causal variable values that we would expect to exist in the intervened vector, h^v .
 186 We can obtain CL vectors by searching through a pre-recorded set of h vectors for situations and
 187 behaviors that are consistent with the values of the CA to which we are aligning. See Figure 1 for a
 188 visualization.

189 As an example, assume we have a CA with variables var_y , var_w , and var_{extra} , and following a
 190 causal intervention we expect h^v to have a value of y for variable var_y and w for variable var_w . For
 191 this example, the CL vector can be obtained from a pre-recorded representation, h_{CL} , that has the
 192 same expected variable values: $var_y = y$ and $var_w = w$, as the intervened vector. The auxiliary loss
 193 $\mathcal{X}^{(k)}$ for a single intervention sample is composed of an L2 and a cosine distance using CL vectors as
 194 the labels:

$$\mathcal{X}_{L2} = \frac{1}{2} \|h^v - h_{CL}\|_2^2 \quad (3)$$

$$\mathcal{X}_{cos} = -\frac{1}{2} \frac{h^v \cdot h_{CL}}{\|h^v\|_2 \|h_{CL}\|_2} \quad (4)$$

195 where h^v is the intervened vector. We combine the CL auxiliary loss with the DAS autoregressive
 196 loss into a single loss term using a weighted sum where ϵ is a hyperparameter: $\mathcal{L}_{total} = \epsilon(\mathcal{X}_{L2} +$
 197 $\mathcal{X}_{cos}) + \mathcal{L}_{DAS}$

198 **4 Results**

199 **4.1 Intervened representations have varying degrees of divergence from the native
200 distribution**

201 We can see in Figure 2 scatter plots of NN latent vectors projected into the top two principal
202 components. The red points come from naturally occurring latent states. The blue come from
203 intervened latent states that have a new value of the Count subspace. The top row shows projections
204 from latents in the model’s native neural space; the intervened latents were transformed back from
205 the aligned space, and the native were left unchanged. The bottom row shows PCA projections of the
206 same latent vectors in the aligned space using only the Count and Phase subspaces—the extraneous
207 subspace was set to zero. All of these PCA projections together are a qualitative demonstration of the
208 possible divergence between the intervened and native representational distributions. It is important
209 to note that these AFs each exhibit a relatively high IIA as shown in Figure 3. Furthermore, we can
210 see that the divergence tends to increase for less restrictive AFs, and, from Supplemental Figure 4,
211 we can see that CAs can differ in both IIA and aligned representational divergence.

212 **4.2 How can causal interventions lead to representational divergence?**

213 In this section, we enumerate distinct cases in which systematic intervened representational divergence
214 can occur. The first draws on the work of [26] who demonstrated that it is possible for an interaction
215 between dormant and null subspaces to occur that can create interventions that produce the correct
216 counterfactual behavior without using native causal subspaces, where dormant subspaces are those
217 that do not vary between inputs, null subspaces are those that exist in the null space of the NN’s
218 next layer, and causal subspaces are those that both vary between inputs and causally affect behavior.
219 These cases of null-dormant subspace interactions can lead to representations that deviate from the
220 native distribution by definition, due to the fact that the value along the dormant direction is different
221 than the native dormant value.

222 Another way in which representational divergence can occur that ignores null and dormant subspaces
223 is in cases where segregated causal subspaces have covariance within the bounds of their behavioral
224 decision boundaries. We use a behaviorally binary subspace as a concrete example, where we define
225 a behaviorally binary subspace as one in which the behavior of the NN depends on the sign of the
226 subspace and is invariant to magnitude and angle.

227 Suppose we have an NN with two causal subspaces, $\tilde{\mathbf{z}}_{\text{var}_a}$ and $\tilde{\mathbf{z}}_{\text{var}_b}$ with values $\tilde{z}_{\text{var}_a}^{(x_i)}$ and $\tilde{z}_{\text{var}_b}^{(x_i)}$ for a
228 model input x_i , where we use the bold notation to distinguish variables from their (non-bold) values.
229 Furthermore, assume that $\tilde{\mathbf{z}}_{\text{var}_b}$ is a behaviorally binary subspace that co-varies with $\tilde{\mathbf{z}}_{\text{var}_a}$. Using
230 $h^{(x_i)}$ and $z^{(x_i)}$ from Equation 1 under a given input x_i , we use the following definition:

$$\mathcal{A}(h^{(x_i)}) = z^{(x_i)} = \begin{cases} \tilde{\mathbf{z}}_{\text{var}_a}^{(x_i)} & = \tilde{z}_{\text{var}_a}^{(x_i)} \\ \tilde{\mathbf{z}}_{\text{var}_b}^{(x_i)} & = \tilde{z}_{\text{var}_b}^{(x_i)} \end{cases} \quad (5)$$

231 Due to the assumption of covariance in $\tilde{\mathbf{z}}_{\text{var}_a}$ and $\tilde{\mathbf{z}}_{\text{var}_b}$, it is reasonable to assume that the values
232 $\tilde{z}_{\text{var}_b}^{(x_{low})}$ and $\tilde{z}_{\text{var}_b}^{(x_{high})}$ are systematically distinct for distinct values of $\tilde{\mathbf{z}}_{\text{var}_a}$ under some pedagogically
233 contrived classes of inputs x_{low} and x_{high} , while $\text{sign}(\tilde{z}_{\text{var}_b}^{(x_{low})}) = \text{sign}(\tilde{z}_{\text{var}_b}^{(x_{high})})$. Under these
234 assumptions, if we perform an interchange intervention on $\tilde{\mathbf{z}}_{\text{var}_b}$ using source representations from
235 input x_{low} and target representations from input x_{high} , the intervened representation will have values:

$$z^v = \begin{cases} \tilde{\mathbf{z}}_{\text{var}_a} & = \tilde{z}_{\text{var}_a}^{(x_{high})} \\ \tilde{\mathbf{z}}_{\text{var}_b} & = \tilde{z}_{\text{var}_b}^{(x_{low})} \end{cases} \quad (6)$$

236 Because we assumed that the value of $\tilde{z}_{\text{var}_b}^{(x_{low})}$ is systematically unique due to covariance in $\tilde{\mathbf{z}}_{\text{var}_a}$ and
237 $\tilde{\mathbf{z}}_{\text{var}_b}$, then the values $\tilde{z}_{\text{var}_a}^{(x_{high})}$ and $\tilde{z}_{\text{var}_b}^{(x_{low})}$ in Equation 6 will have never existed together in the native
238 distribution, but the behavior of the NN will remain unchanged because $\tilde{\mathbf{z}}_{\text{var}_b}$ is behaviorally binary
239 and its sign has not changed.

240 We note that this divergent effect stemming from covariance could occur in the null-space as well as
241 causal subspaces. We argue that due to the causal irrelevance of such divergence, these deviations are

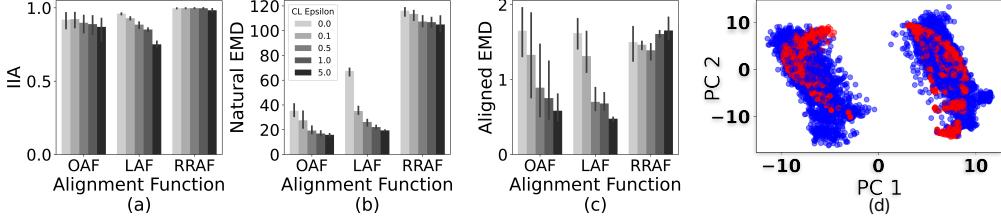


Figure 3: **(a)** The final validation interchange intervention accuracy for different weights (values of epsilon) of the Counterfactual Latent (CL) loss. **(b)** The Earth Mover’s Distance between native latent vectors and intervened latent vectors (in which only the Count subspace has been manipulated) in the NN’s neural space. **(c)** The Earth Mover’s Distance between native latent vectors and intervened latent vectors in the aligned space with non-causal subspaces set to zero. **(d)** PCA projections of the intervened and native latent states for a LAF using a CL epsilon of 5. Similar to Figure 2, blue shows projections from intervened latents whereas red shows native latents.

innocuous to causal mechanistic claims. We provide the lower half of Figure 2 to visualize how the intervened values of z can diverge in practice even in the absence of the extraneous/null subspace. These panels show the top two PCs of aligned z vectors in which the extraneous subspace has been projected out of both the intervened and native latent vectors.

If we focus only on the divergence in causal subspaces and ignore cases of null and dormant interactions, what claims can we make about the model’s native neural mechanisms using divergent, intervened representations? An interpretation consistent with the principles of superposition [34, 33, 27, 6] is that differences in the exact values of each causal subspace do not matter, *only the decision boundaries* along these subspaces matter. Under this interpretation, any divergence in the intervened distribution arising from covariance in causal subspaces is okay, because the separation of the causal subspaces and their decision boundaries—the functionally important aspects of the NN—are respected in the alignment. This interpretation allows us to ignore/abstract away functionally irrelevant covariance between the causal subspaces in our attempt to understand the NN mechanisms, allowing us to focus entirely on the functional/behavioral computations. This perspective emphasizes behavioral control over neural predictivity and is the perspective that we subscribe to.

4.3 Non-linear alignment functions violate assumptions of additivity

In this section we explore how the DAS method for computing interchange interventions (Equation 2) assumes that the AF exhibits principles of additivity, raising concerns about the validity of non-linear AFs used in [35], as additivity is only guaranteed in linear AFs.

Let $\mathcal{A} : \mathbb{R}^{d_m} \rightarrow \mathbb{R}^{d_a}$ be an AF that maps a model’s native latent representation $h \in \mathbb{R}^{d_m}$ to a shared aligned space $z = \mathcal{A}(h)$. We assume that \mathcal{A} is invertible and denote its inverse as \mathcal{A}^{-1} . In practice, \mathcal{A} is a learned transformation.

Interchange interventions are motivated by the desire to substitute values of causal subspaces between latent vectors, where the resulting intervened vector h^v can be defined equivalently to Equation 2 as follows:

$$h^v = \mathcal{A}^{-1} \left(\sum_{j \neq i} \hat{z}_{\text{var}_j}^{\text{trg}} + \hat{z}_{\text{var}_i}^{\text{src}} \right) \quad (7)$$

Each $\hat{z}_{\text{var}_j} \in \mathbb{R}^{d_m}$ corresponds to a masked subspace of the aligned representation defined as:

$$\hat{z}_{\text{var}_j} = D_{\text{var}_j} z \quad (8)$$

and $D_{\text{var}_j} \in \mathbb{R}^{d_m \times d_m}$ is a block diagonal binary matrix isolating the subspace dimensions associated with variable var_j . The D matrices form a partition of the aligned space: $\sum_{j=1}^n D_{\text{var}_j} = I$ and $D_{\text{var}_i} D_{\text{var}_j} = \mathbf{0}$ when $i \neq j$.

The standard implementation of interchange interventions assumes that one can perform the substitution at the level of aligned representations, and then apply the inverse transform to return the representation to the native space (defined in Equation 2, or equivalently in Equation 7). However,

274 this procedure implicitly assumes that the inverse transform \mathcal{A}^{-1} is *additive*, i.e., for any $x, y \in \mathbb{R}^{d_a}$:

$$\mathcal{A}^{-1}(x + y) = \mathcal{A}^{-1}(x) + \mathcal{A}^{-1}(y) \quad (9)$$

275 Without this property, the inverse of a sum is not equal to the sum of the inverses, which breaks the
276 modular interpretation of subspace substitutions.

277 To illustrate this, consider the decomposition of a target latent vector:

$$z^{\text{trg}} = \mathcal{A}(h^{\text{trg}}) = \sum_{j=1}^n \hat{z}_{\text{var}_j}^{\text{trg}} \quad (10)$$

$$h^{\text{trg}} = \mathcal{A}^{-1} \left(\sum_{j=1}^n \hat{z}_{\text{var}_j}^{\text{trg}} \right) \quad (11)$$

278 If \mathcal{A}^{-1} is additive, this becomes:

$$h^{\text{trg}} = \sum_{j=1}^n \mathcal{A}^{-1}(\hat{z}_{\text{var}_j}^{\text{trg}}) \quad (12)$$

279 Thus, we may cleanly isolate and manipulate individual subspace contributions in the aligned space
280 before inverting.

281 The same logic applies to an intervention where the i -th subspace is replaced from a source latent
282 vector:

$$h^v = \mathcal{A}^{-1} \left(\sum_{j \neq i} \hat{z}_{\text{var}_j}^{\text{trg}} + \hat{z}_{\text{var}_i}^{\text{src}} \right) \quad (13)$$

$$= \sum_{j \neq i} \mathcal{A}^{-1}(\hat{z}_{\text{var}_j}^{\text{trg}}) + \mathcal{A}^{-1}(\hat{z}_{\text{var}_i}^{\text{src}}) \quad (\text{if } \mathcal{A}^{-1} \text{ is additive}) \quad (14)$$

283 The standard interchange intervention approach to computing and isolating $\sum_{j \neq i} \hat{z}_{\text{var}_j}^{\text{trg}}$ and $\hat{z}_{\text{var}_i}^{\text{src}}$ first
284 independently computes each \hat{z} from the respective h^{trg} and h^{src} before adding them together. Thus,
285 the standard method for computing h^v relies on the principle of additivity.

286 **Proposition.** Suppose \mathcal{A}^{-1} is not additive. Then there exist vectors $x, y \in \mathbb{R}^{d_a}$ such that:

$$\mathcal{A}^{-1}(x + y) \neq \mathcal{A}^{-1}(x) + \mathcal{A}^{-1}(y). \quad (15)$$

287 Consequently, there exist source and target latent representations for which:

$$h^v = \mathcal{A}^{-1} \left(\sum_{j \neq i} \hat{z}_{\text{var}_j}^{\text{trg}} + \hat{z}_{\text{var}_i}^{\text{src}} \right)$$

288 does not equal:

$$\sum_{j \neq i} \mathcal{A}^{-1}(\hat{z}_{\text{var}_j}^{\text{trg}}) + \mathcal{A}^{-1}(\hat{z}_{\text{var}_i}^{\text{src}}).$$

289 This means that the inverse transform introduces non-linear interactions between subspaces, under-
290 mining any clean attribution of h^v to its constituent parts.

291 **Implication.** The correctness and interpretability of standard interchange interventions depend on
292 the additivity of \mathcal{A}^{-1} . In the absence of this property, it becomes unclear whether the reconstructed
293 vector h^v reflects an interpretable combination of the intended latent subspaces. Instead, \mathcal{A}^{-1} may
294 behave like an arbitrary function \mathcal{F} with no guaranteed semantic alignment to \mathcal{A} .

295 **4.4 Counterfactual Latent Vectors alleviate post-intervention divergence**

296 Although we have shown in Section 4.2 why some types of representational divergence are acceptable
297 for many mechanistic claims, there are still some cases in which it may be desirable for intervened
298 representations to be predictive of the native distribution. For these cases, we explore the use of a CL
299 auxiliary loss, which is one that minimizes the L2 and cosine distances between the intervened and
300 native representations. We can see in Figure 3 that we can successfully reduce the Earth Mover's
301 distance between the intervened and native distributions by applying the CL loss during the DAS
302 training. This is a step towards making DAS more relevant for goals of neural predictivity.

303 **5 Limitations/Future Directions**

304 The results presented in this work has been confined to synthetic GRUs and simplistic tasks. A more
305 complete demonstration of the intervened distribution shift would include evaluations performed
306 on larger, more practically oriented transformers and Large Language Models (LLMs). This is
307 particularly noteworthy for this work as the GRU architecture imposes a Tanh nonlinearity on the
308 representations analyzed in this work, whereas Transformers potentially have more linear representa-
309 tions due to the nature of their residual stream [6]. We look forward to exploring LLMs in future
310 work.

311 **6 Discussion/Conclusion**

312 In this work we examined autoregressive GRUs trained on numeric tasks to demonstrate the following:
313 we showed the degree to which intervened representations can diverge from naturally occurring; we
314 showed how intervened representations can occur in ways that are mechanically innocuous to many
315 claims; we showed how AFs that don't exhibit principles of additivity violate implicit assumptions in
316 interchange interventions; and we showed how to use a CL auxiliary loss to reduce the divergence
317 between native and intervened representations. Where does this leave us with respect to neural
318 interpretability?

319 We return to our goals underlying the notion of "understanding" neural activity. In general, it is
320 reasonable to equate the notion of "understanding neural activity" to one's ability to predict and/or
321 control the activity in ways that are deemed interpretable. Causal methods such as DAS manage to
322 control NN behavior quite well as exemplified by the strong IIA in this work, and these methods
323 do so through *interpretable* causal abstractions. However, causal methods may be limited in their
324 ability to predict native neural activity. Different research objectives will place different weights on
325 the importance of neural predictivity and behavioral control. In some cases, for example, we may
326 wish to classify native neural activity, in which case, predictivity is potentially useful. In other cases,
327 we may wish to exert influence over the computations of the NN for the purpose of AI safety or to
328 characterize the space of *potential* computations as a means of predicting NN generalization.

329 In light of classic causal mediation philosophy corroborated by more recent findings of complications
330 in non-causal representational analyses, we find ourselves favoring causal methods for their ability
331 to ignore computationally irrelevant NN details as shown in this work, and their ability to unify
332 diverse neural systems [17, 21, 22] while providing useful ways of understanding neural mechanisms
333 [11, 10]. We hesitate to diminish goals of neural predictivity, however, and we remind ourselves that
334 there are no guarantees that we find satisfying, interpretable ways of understanding the complete
335 complexities of NNs.

336 **References**

- 337 [1] Trenton Bricken, Neel Nanda, Nicholas Joseph, Arthur Conmy, Andy Jones, Anna Chen, Neal
338 DasSarma, Nelson Elhage, Ben Mann, Catherine Olsson, Kamal Ndousse, Sam Ringer, Alex
339 Tran-Johnson, Yuntao Bai, Liane Lovitt, Zac Hatfield-Dodds, Amanda Askell, Dario Amodei,
340 Tom Brown, Jack Clark, Jared Kaplan, and Sam McCandlish. Towards monosemanticity:
341 Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023.
- 342 [2] Maheep Chaudhary and Atticus Geiger. Evaluating open-source sparse autoencoders on disen-
343 tangling factual knowledge in gpt-2 small. *arXiv*, 2024. arXiv preprint.

- 344 [3] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülcehre, Fethi Bougares, Holger Schwenk,
345 and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical
346 machine translation. *CoRR*, abs/1406.1078, 2014.
- 347 [4] Róbert Csordás, Christopher Potts, Christopher D. Manning, and Atticus Geiger. Recurrent
348 neural networks learn to store and generate sequences using non-linear representations, 2024.
- 349 [5] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-
350 coders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*,
351 2023.
- 352 [6] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna
353 Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse,
354 Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah.
355 Toy models of superposition. *Transformer Circuits Thread*, 2022. https://transformer-circuits.pub/2022/toy_model/index.html.
- 356 [7] Michael C. Frank, Daniel L. Everett, Evelina Fedorenko, and Edward Gibson. Number as a
357 cognitive technology: Evidence from pirahã language and cognition. *Cognition*, 108(3):819–
358 824, 2008.
- 360 [8] Dan Friedman, Andrew Lampinen, Lucas Dixon, Danqi Chen, and Asma Ghandeharioun. Inter-
361 pretability illusions in the generalization of simplified models. *arXiv preprint arXiv:2312.03656*,
362 2023.
- 363 [9] Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
364 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv
365 preprint arXiv:2406.04093*, 2024.
- 366 [10] Atticus Geiger, Jacqueline Harding, and Thomas Icard. How causal abstraction underpins
367 computational explanation. *arXiv preprint arXiv:2508.11214*, 2025.
- 368 [11] Atticus Geiger, Duligur Ibeling, Amir Zur, Maheep Chaudhary, Sonakshi Chauhan, Jing Huang,
369 Aryaman Arora, Zhengxuan Wu, Noah Goodman, Christopher Potts, and Thomas Icard. Causal
370 abstraction: A theoretical foundation for mechanistic interpretability, 2024.
- 371 [12] Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal abstractions of neural
372 networks. *CoRR*, abs/2106.02997, 2021.
- 373 [13] Atticus Geiger, Kyle Richardson, and Christopher Potts. Neural natural language inference mod-
374 els partially embed theories of lexical entailment and negation. *arXiv preprint arXiv:2004.14623*,
375 2020.
- 376 [14] Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah D. Goodman.
377 Finding alignments between interpretable causal variables and distributed neural representations,
378 2023.
- 379 [15] Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible resid-
380 ual network: Backpropagation without storing activations. *Advances in neural information
381 processing systems*, 30, 2017.
- 382 [16] Peter Gordon. Numerical cognition without words: Evidence from Amazonia. *Science*,
383 306(5695):496–499, 2004.
- 384 [17] Satchel Grant. Model alignment search. *arXiv preprint arXiv:2501.06164*, 2025.
- 385 [18] Satchel Grant, Noah D. Goodman, and James L. McClelland. Emergent symbol-like number
386 variables in artificial neural networks. *Transactions on Machine Learning Research*, 2024.
- 387 [19] Eric Jonas and Konrad Paul Kording. Could a neuroscientist understand a microprocessor?
388 *PLOS Computational Biology*, 13(1):e1005268, 2017.
- 389 [20] Subhash Kantamneni, Joshua Engels, Senthooran Rajamanoharan, Max Tegmark, and Neel
390 Nanda. Are sparse autoencoders useful? a case study in sparse probing. *arXiv preprint
391 arXiv:2502.16681*, 2025.

- 392 [21] Andrew Kyle Lampinen, Stephanie CY Chan, and Katherine Hermann. Learned feature
 393 representations are biased by complexity, learning order, position, and more. *arXiv preprint*
 394 *arXiv:2405.05847*, 2024.
- 395 [22] Andrew Kyle Lampinen, Stephanie CY Chan, Yuxuan Li, and Katherine Hermann. Representation
 396 biases: will we achieve complete understanding by analyzing representations? *arXiv preprint*
 397 *arXiv:2507.22216*, 2025.
- 398 [23] Patrick Leask, Bart Bussmann, Michael Pearce, Joseph Bloom, Curt Tigges, Noura
 399 Al Moubayed, Lee Sharkey, and Neel Nanda. Sparse autoencoders do not find canonical
 400 units of analysis. *arXiv*, 2025. arXiv preprint.
- 401 [24] Aaron J. Li, Suraj Srinivas, Usha Bhalla, and Himabindu Lakkaraju. Interpretability illusions
 402 with sparse autoencoders: Evaluating robustness of concept representations. *arXiv*, 2025. arXiv
 403 preprint.
- 404 [25] Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner,
 405 Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar,
 406 Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan,
 407 Adam Jermyn, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman,
 408 Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. On the biology of a large
 409 language model. *Transformer Circuits Thread*, 2025.
- 410 [26] Aleksandar Makelov, Georg Lange, and Neel Nanda. Is this the subspace you are looking for?
 411 an interpretability illusion for subspace activation patching. *arXiv preprint arXiv:2311.17030*,
 412 2023.
- 413 [27] J. L. McClelland, D. E. Rumelhart, and PDP Research Group, editors. *Parallel Distributed
 414 Processing. Volume 2: Psychological and Biological Models*. MIT Press, Cambridge, MA,
 415 1986.
- 416 [28] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual
 417 associations in gpt, 2023.
- 418 [29] Neel Nanda. Attribution patching: Activation patching at industrial scale. [https://www.
 419 neelnanda.io/mechanistic-interpretability/attribution-patching](https://www.neelnanda.io/mechanistic-interpretability/attribution-patching), 2022.
- 420 [30] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
 421 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
 422 Köpf, Edward Z. Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
 423 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,
 424 high-performance deep learning library. *CoRR*, abs/1912.01703, 2019.
- 425 [31] Judea Pearl. An Introduction to Causal Inference. *The International Journal of Biostatistics*,
 426 6(2):7, February 2010.
- 427 [32] Karl Pearson. On lines and planes of closest fit to systems of points in space. *The London,
 428 Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901.
- 429 [33] D. E. Rumelhart, J. L. McClelland, and PDP Research Group, editors. *Parallel Distributed
 430 Processing. Volume 1: Foundations*. MIT Press, Cambridge, MA, 1986.
- 431 [34] Paul Smolensky. On the proper treatment of connectionism. *Behavioral and Brain Sciences*,
 432 11(1):1–23, 1988.
- 433 [35] Denis Sutter, Julian Minder, Thomas Hofmann, and Tiago Pimentel. The non-linear represen-
 434 tation dilemma: Is causal abstraction enough for mechanistic interpretability? *arXiv preprint*
 435 *arXiv:2405.05847*, 2025.
- 436 [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
 437 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- 438 [37] Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Simas Sakenis,
 439 Jason Huang, Yaron Singer, and Stuart Shieber. Causal mediation analysis for interpreting
 440 neural nlp: The case of gender bias. *arXiv preprint arXiv:2004.12265*, 2020.

⁴⁴¹ [38] Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.
⁴⁴² Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022.

⁴⁴³ [39] Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah D. Goodman.
⁴⁴⁴ Interpretability at scale: Identifying causal mechanisms in alpaca, 2024.

Input Tokens →		B	D	D	T	R	R	E
Count Up Down Variables	→	Count: 0	1	2	2	1	0	NA
Increment Up Up Variables	→	Phase: 0	0	0	1	1	1	NA
Progress:	1/21	1/21	2/21	0	1/2	1	1/2	NA
Increment:	1/21	1/21	1/21	1/2	1/2	1/2	1/2	NA

Figure 4: Visual depiction of the causal abstractions considered in this work.

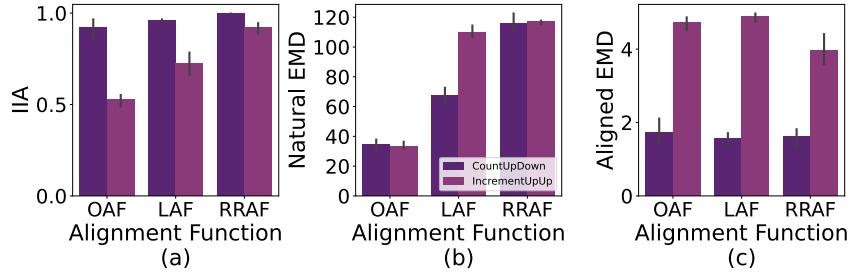


Figure 5: (a) The final validation interchange intervention accuracy for different causal abstractions aligned to the same models. (b) The Earth Mover’s Distance between naturally occurring latent vectors and intervened latent vectors (in which only the Count subspace or the Progress subspace have been manipulated depending on the causal abstraction) in the NN’s neural space. (c) The Earth Mover’s Distance between naturally occurring vectors and intervened vectors in the aligned space with non-causal subspaces set to zero.

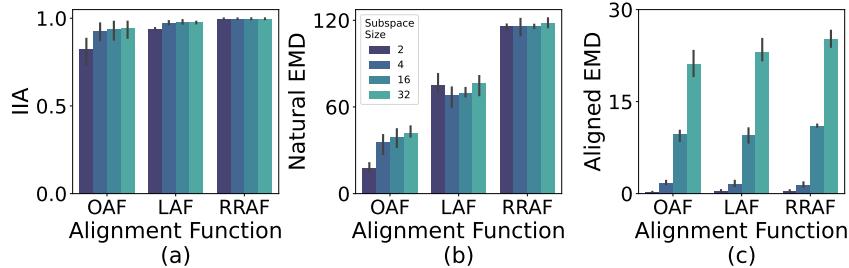


Figure 6: (a) The final validation interchange intervention accuracy for different causal subspace sizes. Each subspace size is used for both the Count and Phase variables from the Count Up Down causal abstraction in each alignment. (b) The Earth Mover’s Distance between naturally occurring latent vectors and intervened latent vectors (in which only the Count subspace has been manipulated) in the NN’s neural space. (c) The Earth Mover’s Distance between naturally occurring vectors and intervened vectors in the aligned space with non-causal subspaces set to zero.

445 A Appendix

446 A.1 Model Details

447 All artificial neural network models were implemented and trained using PyTorch [30] on Nvidia
448 Titan X GPUs. Unless otherwise stated, all models used an embedding and hidden state size of 128
449 dimensions. To make the token predictions, each model used a two layer multi-layer perceptron
450 (MLP) with GELU nonlinearities, with a hidden layer size of 4 times the hidden state dimensionality
451 with 50% dropout on the hidden layer. The model consisted of a single GRU recurrent cell followed

452 by an output multi-layer perceptron (MLP). We show the GRU model structure:

$$h_{t+1} = f(h_t, x_t) \quad (16)$$

$$\hat{x}_{t+1} = g(h_{t+1}) \quad (17)$$

453 Where h_t is the hidden state vector at step t , x_t is the input token at step t , f is the GRU cell, and g is
454 a two layer (two matrix) MLP used to make a prediction, \hat{x}_{t+1} , of the token at step $t + 1$ from the
455 updated hidden state h_{t+1} . Models were trained using a learning rate scheduler, which consisted of
456 the original transformer [36] scheduling of warmup followed by decay. We used 100 warmup steps, a
457 maximum learning rate of 0.0001, a minimum of 1e-7, and a decay rate of 0.5. We used a batch size
458 of 128, which caused each epoch to consist of 8 gradient update steps.

459 A.2 DAS Details

460 In our experiments, we perform causal interventions on individual time steps in the sequence. We
461 run the model up to a sampled timestep t in the target sequence, taking its latent representation
462 at that point as the target vector, h_t^{trg} . We do the same for the source vector, h_u^{src} , at timestep u
463 from a separate source sequence. We then construct h_t^v using Equation 2, and continue the model's
464 predictions starting from time t , using h_t^v in place of h_t^{trg} .

465 We use 10,000 intervention samples for training and 1,000 samples for validation and testing. For all
466 data, we uniformly sample trial object quantities, and unless otherwise stated, we uniformly sample
467 intervention time points, t and u , from sequence positions containing demo tokens or response tokens
468 (excluding BOS, trigger, and EOS tokens).

469 We orthogonalize the rotation matrix for OAFs using PyTorch's orthogonal parameterization with
470 default settings. We train Q with a batch size of 512 until convergence, selecting the checkpoint with
471 the best validation performance for analysis. We use a learning rate of 0.001 and an Adam optimizer.

472 To train the AFs, we sampled 10,000 sequence pairs for the intervention training datasets. See
473 Supplement A.3.1 for more details on intervention data construction and examples. We use a learning
474 rate of 0.001 and a batch size of 512. We removed models with performance below 99% to limit our
475 DAS results to perfectly performing models thus simplifying our interpretations of the results. We
476 chose 99% accuracy instead of 100% due to slight numerical underflow in accuracy calculations and
477 due to the fact that half of the Variable-Length Same-Object models would have been dropped due to
478 low performance.

479 A.3 Linear Alignment Functions

480 To construct the W matrix used in the LAF, we use the following equation: $W = (MM^\top +$
481 $\epsilon I)S$. $M \in R^{d_m \times d_m}$ is a matrix of learned parameters initially sampled from a centered gaussian
482 distribution with a standard deviation of $\frac{1}{d_m}$, $I \in R^{d_m \times d_m}$ is the identity matrix, $\epsilon = 0.1$ to prevent
483 singular values equal to 0, and $S \in R^{d_m \times d_m}$ is a diagonal matrix to learn a sign for each column of
484 X using diagonal values $s_{i,i} = \text{Tanh}(a_i) + \epsilon(\text{sign}(\text{Tanh}(a_i)))$ where each a_i is a learned parameter
485 and $\epsilon = 0.1$ to prevent 0 values.

486 A.3.1 DAS Intervention Data

487 In this section we provide intervention data examples. To construct an intervention sample, we first
488 sample a target sequence and a source sequence and a positional index from each sequence. We
489 exclude trigger tokens and beginning and end of sequence tokens from the possible positional indices.
490 We then compute the values of each of the variables at the sampled indices using the specified CA for
491 both the target and source. We then transfer the value of the variable of focus from the source into the
492 the target variable in the CA. We then continue the CA using the new variable values to produce the
493 counterfactual sequence.

494 A.3.2 Up-Down Program Examples

495 **Count Variable:** Interventions attempt to transfer the representation corresponding to the difference
496 between the number of resp tokens and demo tokens. Interventions are only performed at positional
497 indices corresponding to demo or resp tokens. The target sequence maintains its original object count

498 when the Count variable is changed in the demo phase. In cases where the new value exceeds the
 499 object count, the CA immediately produces the trigger token.

500

	<i>Multi-Object Examples</i>	1	2	3	4
Source Sequence	BOS D ₁	BOS D ₂ D ₁ D ₁	BOS D ₂ D ₁ T R	BOS D ₁ D ₃ T R R	
Target Sequence	BOS D ₃ D ₂	BOS D ₂ T R	BOS D ₁ D ₂ D ₁ T R	BOS D ₂	
Original Labels	D ₂ D ₃ T R R R R EOS	EOS	R R EOS	D ₂ T R R EOS	
Counterfactual	D ₂ D ₃ T R R R EOS	R R R EOS	R EOS	D ₂ T R EOS	
	<i>Single-Object Examples</i>	1	2	3	4
Source Sequence	BOS D	BOS D D D	BOS D D T R	BOS D D T R R	
Target Sequence	BOS D D	BOS D T R	BOS D D D T R	BOS D	
Original Labels	D D T R R R R EOS	EOS	R R EOS	D T R R EOS	
Counterfactual	D D T R R R EOS	R R R EOS	R EOS	D T R EOS	
	<i>Same-Object Examples</i>	1	2	3	4
Source Sequence	BOS C	BOS C C C	BOS C C T C	BOS C C T C C	
Target Sequence	BOS C C	BOS C T C	BOS C C C T C	BOS C	
Original Labels	C C T C C C C EOS	EOS	C C EOS	C T C C EOS	
Counterfactual	C C T C C C EOS	C C C EOS	C EOS	C T C EOS	

502 **Phase Variable:** Interventions transfer the representation corresponding to the Phase of the sequence
 503 (whether it is counting up or counting down). Interventions are only performed at positional indices
 504 corresponding to demo or resp tokens.

505

	<i>Multi-Object Examples</i>	1	2	3	4
Source Sequence	BOS D ₁	BOS D ₃ D ₁ D ₂	BOS D ₂ D ₁ T R	BOS D ₂ D ₃ T R R	
Target Sequence	BOS D ₂ D ₁	BOS D ₃ T R	BOS D ₁ D ₃ D ₁ T R	BOS D ₂	
Original Labels	D ₃ D ₁ T R R R R EOS	EOS	R R EOS	D ₁ T R R EOS	
Counterfactual	D ₃ D ₁ T R R R R EOS	D ₂ T R EOS	R R EOS	R EOS	
	<i>Single-Object Examples</i>	1	2	3	4
Source Sequence	BOS D	BOS D D D	BOS D D T R	BOS D D T R R	
Target Sequence	BOS D D	BOS D T R	BOS D D D T R	BOS D	
Original Labels	D D T R R R R EOS	EOS	R R EOS	D T R R EOS	
Counterfactual	D D T R R R R EOS	D T R EOS	R R EOS	R EOS	
	<i>Same-Object Examples</i>	1	2	3	4
Source Sequence	BOS C	BOS C C C	BOS C C T C	BOS C C T C C	
Target Sequence	BOS C C	BOS C T C	BOS C C C T C	BOS C	
Original Labels	C C T C C C C EOS	EOS	C C EOS	C T C C EOS	
Counterfactual	C C T C C C C EOS	C T C EOS	C C EOS	C EOS	

Algorithm 1 One sequence step of the Up-Down Program

```
q ← Count
p ← Phase
y ← input token
if  $y == \text{BOS}$  then
    q ← 0, p ← 0
    return sample(D)                                ▷ BOS is beginning of sequence token
else if  $y \in D$  then
    q ← q + 1
    return sample(D)                                ▷ sample a demo token
else if  $y == \text{T}$  then
    p ← 1                                            ▷ T is trigger token
else if  $y == \text{R}$  then
    q ← q - 1                                       ▷ R is response token
end if
if ( $q == 0$ ) & ( $p == 1$ ) then
    return EOS                                     ▷ EOS is end of sequence token
end if
return R
```

Algorithm 2 One sequence step of the Increment-Up Program

```
m ← Interval
q ← Progress
p ← Phase
i ← Increment
y ← input token
if  $y == \text{BOS}$  then
    q ← 0, p ← 0, i ←  $\frac{1}{m}$ 
    return sample(D)                                ▷ BOS is beginning of sequence token
else if ( $y \in D$  or  $y == \text{R}$ ) and  $q < m$  then
    q ← q + i * m                                 ▷ sample a demo token
else if  $y == \text{T}$  then
    p ← 1                                            ▷ T is trigger token
    i ←  $\frac{1}{q}$ 
    q ← 0
end if
if ( $q \geq m$ ) and ( $p == 1$ ) then
    return EOS                                     ▷ EOS is end of sequence token
else if  $q \geq m$  and  $p == 0$  then
    return T
else if  $p == 0$  then
    return sample(D)
else
    return R
end if
```
