
From Causal to Concept-Based Representation Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To build intelligent machine learning systems, there are two broad approaches.
2 One approach is to build inherently interpretable models, as endeavored by the
3 growing field of causal representation learning. The other approach is to build
4 highly-performant foundation models and then invest efforts into understanding
5 how they work. In this work, we relate these two approaches and study how to learn
6 human-interpretable concepts from data. Weaving together ideas from both fields,
7 we formally define a notion of concepts and prove that they can be identifiably
8 recovered from diverse data. Experiments on synthetic data, CLIP models and
9 large language models show the utility of our unified approach.

10 1 Introduction

11 A key goal of modern machine learning is to learn representations of complex data that are human-
12 interpretable and can be controlled. This goal is of paramount importance given the breadth and
13 importance of ML in today’s world. There seem to be two broad approaches toward such intelligent
14 systems. The first approach is to build models that are inherently interpretable and then subsequently
15 focus on how to extract maximum performance from them; and the second approach is to build high-
16 performance neural models, and then subsequently invest efforts to understand the inner workings of
17 such models.

18 A prominent example of the first camp is the field of Causal Representation Learning (CRL) [90, 89].
19 CRL is an intricate interplay of ideas from causality, latent variable modeling and deep learning, with
20 the main goal being to reconstruct the true generative factors of data. To ensure that the true generative
21 factors can be provably identified, CRL relies on the central theme of *identifiability* which posits that a
22 unique model fits the data, which in turn implies that the problem of learning the generative factors is
23 well-posed and therefore should theoretically be amenable to modern techniques. If such a generative
24 model reconstruction can be done, the model will naturally enjoy a host of desired properties such
25 as robustness and generalization. While this endeavor has been (moderately) successful in many
26 domains such as computer vision [45, 113, 2], robotics [63, 10, 59, 126] and genomics [98, 125], it
27 is unclear how it relates to the research on foundation models.

28 The other camp is more empirical, where one tries to build a high-performance model where
29 performance is measured via various downstream tasks and then eventually invest efforts into
30 explaining or interpreting how they work. For instance, large language models and other foundation
31 models are built to be highly performant for a variety of tasks. Owing to their incredible success,
32 there is a growing but heavily-debated belief that such models are truly “intelligent” because they
33 have indeed learned the true underlying generative factors somehow, sometimes referred to as the
34 “world model”. While we are far from scientifically verifying this, the community has invested
35 tremendous efforts into interpretability research of foundation models, e.g., the field of mechanistic
36 interpretability [72] aims to reverse engineer what large language models learn.

37 In this work, we make the first step toward unifying these approaches. We focus on the goal of
38 learning identifiable human-interpretable concepts from complex high-dimensional data. Specifically,
39 we build a theory of what concepts mean for complex high-dimensional data and then study under
40 what conditions such concepts are identifiable, i.e., when can they be unambiguously recovered from
41 data. To formally define concepts, we leverage extensive empirical evidence in the foundation models
42 literature that surprisingly shows that, across multiple domains, human-interpretable concepts are
43 often *linearly* encoded in the latent space of such models (see Section 3), e.g., the sentiment of a
44 sentence is linearly represented in the activation space of large language models [105]. Motivated by
45 this rich empirical literature, we formally define concepts as affine subspaces of some underlying
46 representation space. Then we prove strong identifiability theorems for *only desired concepts* rather
47 than all possible concepts present in the true generative model. Therefore, in this work we tread
48 the fine line between the rigorous principles of causal representation learning and the empirical
49 capabilities of foundation models, effectively showing how causal representation learning ideas can
50 be applied to foundation models.

51 In CRL we generally model the input data $X = (X_1, \dots, X_{d_x})$ as $X = f(Z)$, where f is a nonlinear
52 transformation that maps structured underlying latent generative factors $Z = (Z_1, \dots, Z_{d_z})$ to X ,
53 and then to attempt to recover the model parameters Z, f from X . This is an appealing approach since
54 it implies no restrictions on the data X , and has the interpretation of recovering “ground truth” factors
55 that generated the data. It is well-known that without additional assumptions, this is impossible
56 [38, 61], a fact which has led to a long line of work on nonlinear ICA [18, 37] and unsupervised
57 disentanglement [9, 77, 52]. One approach to resolve this limitation is to assume that Z has an intrinsic
58 causal interpretation, as in CRL. Recent years have witnessed a surge of rigorous results on provably
59 learning causal representations under different assumptions [45, 28, 60, 51, 68, 128, 31, 110, 41, 102].
60 For example, as long as we have access to interventions on each latent variable Z_j (a total of at least
61 d_z interventions), under weak assumptions on Z and/or f , the causal model over Z as well as the
62 model parameters (Z, f) can be uniquely identified [98, 12].

63 While causal features are intrinsically desirable in many applications, the assumption that we can
64 feasibly perform $\Omega(d_z)$ interventions merits relaxing: Indeed, in complex models, the number of
65 true generative factors $d_z = \dim(Z)$ might be intractably large (e.g. consider all of the latent factors
66 that could be used to describe natural images, video, or text). At the same time, there are yet many
67 other applications where the strict notion of causality may not be needed, and moreover it may not be
68 necessary to learn the *full* causal model over every causal factor. Is there a middle ground where we
69 can simultaneously identify a smaller set of interpretable latent representations, without the need for
70 a huge number of interventions?

71 We study this problem in detail and provide an alternative setting under which latent representations
72 can be provably recovered. The basic idea is to recover *projections* AZ of the generative factors Z that
73 correspond to meaningful, human-interpretable concepts through *conditioning* instead of intervention.
74 The idea to model concepts as linear projections of the generative factors is derived from a growing
75 body of literature (e.g. [79, 47, 117, 67, 5, 19, 25, 15, 105, 71, 33, 65, 91], see Section 3 for even more
76 references) showing that the embeddings learned by modern, high-performant foundation models are
77 not inherently interpretable, and instead capture interpretable concepts as linear projections of the
78 (*apriori*) unintelligible embeddings. While this approach sacrifices causal semantics, it makes up for
79 this with two crucial advantages: 1) Instead of strict interventions in the latent space, it suffices to
80 *condition* on the concepts, and 2) When there are n concepts of interest to be learned, only $n+2 \ll d_z$
81 such concept conditionals are needed.

82 Furthermore, we validate and utilize our theoretical ideas via both simulations and experiments with
83 foundation models, including an effective application of our framework to large language models
84 (LLMs). First, we validate these theoretical insights on synthetic data, where we use a contrastive
85 algorithm to learn such representations for a given collection of concepts. Moving ahead to real-world
86 data, we probe our theory on embeddings learned by multimodal CLIP models [81]. The training
87 scheme for CLIP aligns with our theoretical setting and therefore, it’s reasonable to ask whether they
88 satisfy our observations. Indeed, we show that the concepts in the 3d-Shapes dataset approximately lie
89 in hyperplanes, further supporting our theoretical results. Lastly, we show an effective application of
90 our framework to large language model (LLM) alignment, where we extend the alignment technique
91 of [56] to make LLMs more truthful.

92 **Contributions** In summary, our contributions are:

- 93 1. We formalize the notion of distributions induced by abstract concepts in complex domains
94 such as images or text (see Section 2 for an overview and Section A.2 for formal defini-
95 tions). Our definition of concept conditional distributions allows both continuous and fuzzy
96 concepts.
- 97 2. We prove near-optimal identifiability results for learning a collection of concepts from
98 a diverse set of environments in Theorem 2. Thus our work can be interpreted as a new
99 direction for identifiable representation learning in order to study when interpretable concepts
100 can be recovered from data.
- 101 3. We then verify our guarantees via a contrastive learning algorithm on synthetic data. In
102 addition in Section 5, we support our geometric definition of concepts and our identifiability
103 result by analysing image embeddings of CLIP-models and we utilize our ideas to improve
104 alignment of LLMs to make them more truthful.

105 2 Overview

106 In this section, we describe our approach and put it in context of prior developments.

107 **Defining concepts geometrically** Our starting point is a geometric no-
108 tion that concepts live in linear directions in neural representation space,
109 known as linearity of representations (see extensive references in Section 3).
110 To make this precise we assume that for observed data X that has an underlying representation Z with $X = f(Z)$
111 where the latent variables Z follow an arbitrary distribution and f is a (potentially complicated) nonlinear un-
112 derlying mixing map. We do not assume that f and Z correspond to a ground truth model or that the latent vari-
113 ables Z themselves are related to a causal model or are interpretable and instead only assume linearity of repre-
114 sentations (well supported by prior works). In agreement with this hypothesis we define concepts as affine subspaces
115 $AZ = b$ of the latent space of Z s, i.e., to a concept C we
116 assign an affine hyperplane $H_C = \{Z \in \mathbb{R}^{d_z} : AZ = b\}$
117 in the embedding space and we say that $X = f(Z)$ satisfies a concept C if $Z \in H_C$. We focus on the goal of
118 identifying only a (small) set of *concepts we care about*, i.e., we want to be able to decide whether a datapoint X
119 satisfies a concept C . Our main result shows that it is possible to identify n concepts given access to
120 $n + 2$ concept conditional distributions. We now compare natural assumptions on type of data for
121 causal representation learning and the setting considered here.
122
123
124
125
126
127
128

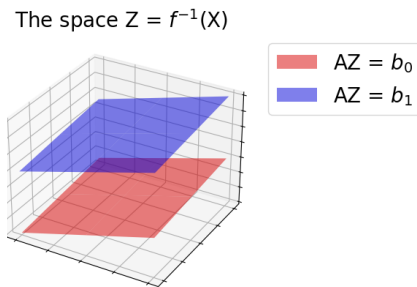


Figure 1: Concepts live in affine subspaces. The two subspaces in the figure correspond to the same concept but of different valuations.

129 **From interventions to conditioning** It is worth contrasting here the difference between viewing
130 a concept as a generic latent generative factor Z_i that non-linearly mixes together with other latent
131 factors to yield the inputs X , versus the geometric notion above, as specifying a linear subspace.
132 In the former, the natural way to provide supervision, i.e. define concept distributions, is to simply
133 intervene on a specific factor Z_i and set it to a particular value (see Section 3 for references). In
134 the latter however, it is most natural to condition on the concept, i.e., $Z \in H$. This shift is aligned
135 with the growing interest to relax the notion of interventions, and consequently dilute the notion of
136 causality [13, 88, 4], although it is still open how to properly achieve this. Two key drivers of this
137 trend are as follows. The first is that the number of additional datasets required is d_z [38, 61, 45, 12],
138 which is infeasible in many settings¹. The second is that the various assumptions that go into these
139 works are often difficult to achieve, such as requiring perfect interventions [98, 12]. Compared to
140 interventional data, *conditional* data is often easier to acquire, obtained by conditioning on particular
141 values of the latent factors (see also Appendix C.2).

142 **Concept conditional distributions** We now formalize conditioning on a concept. The obvious
143 approach to define concept conditional distributions is to simply condition on $Z \in H_C$, so $p_C(Z) =$

¹Exceptions are [49, 35], which use clever inductive biases to limit the number of environments needed.

144 $p(Z|Z \in H_C)$ where p is a base distribution of Z on \mathbb{R}^{d_z} . However, this suffers from the drawbacks
 145 that it is mathematically subtle to condition on sets of measure 0 and this does not account for inherent
 146 noise in the learned representations. Therefore we relax this strict conditioning by drawing inspiration
 147 from how data is collected in practice: We sample X from the base distribution and then keep it if
 148 it satisfies our concept C . This leads us to define $p_C(Z) \propto p(Z)q(Z|C)$ where q is defined to be
 149 the probability that Z is *perceived* to be in H by the data collector and can be chosen to incorporate
 150 noise in our data gathering scheme. Therefore, this can also be viewed from a Bayesian information
 151 gathering viewpoint, as well as a stochastic filter standpoint. This is the notion we study in this work
 152 (Definition 3) and we develop theoretical techniques to guarantee identifiability in this formulation.
 153 Depending on the specific setting other types of conditional distributions might be utilized to describe
 154 the available data and we discuss some options in Appendix D.

155 3 Related work

156 **Causal representation learning and concept discovery** Causal representation learning (CRL) [90,
 157 89] aims to learn generative factors of high-dimensional data. This exciting field has seen significant
 158 progress in the last few years [45, 10, 93, 51, 68, 49, 101, 12, 31, 1, 114, 53]. A fundamental
 159 perspective in this field is to ensure that the model parameters we attempt to recover are identifiable
 160 [45, 21, 116]. We will elaborate more on the connection of our framework to CRL in Appendix C.
 161 Concept discovery is an important sub-field of machine learning which extracts human-interpretable
 162 concepts from pre-trained models. We do not attempt to list the numerous works in this direction,
 163 see e.g., [91, 16, 122, 64, 78]. However, theoretical progress in this direction is relatively limited.
 164 The work [53] studies when concepts can be identified provided the non-linear model is known in
 165 advance, whereas we show concept identifiability for unknown non-linearity, while simultaneously
 166 allowing entangled concepts. Prior works have also attempted to formalize the notion of concepts
 167 [117, 74, 91], however their definitions seem specific to the model and domain under consideration,
 168 e.g., [74, 44] focus on binary concepts via large language model representations of counterfactual
 169 word pairs, whereas our general concept definitions are applicable to all domains.

170 **Linearity of representations** Sometimes referred to as the linear representation hypothesis,
 171 it is commonly believed that well-trained foundation models in multiple domains learn linear
 172 representations of human-interpretable concepts, with experimental evidence going back at
 173 least a decade [67, 100, 5]. This has been experimentally observed in computer vision models
 174 [79, 83, 8, 26, 47, 117, 107], language models [67, 76, 5, 19, 104, 25], large language models
 175 [15, 105, 71, 69, 56, 74, 33, 44], and other intelligent systems [65, 91]. Various works have also
 176 attempted to justify why this happens [54, 5, 30, 3, 27, 92]. We take a different angle: Given that this
 177 phenomenon has been observed for certain concepts of interest, how does this enable recovery of the
 178 concepts themselves? Consequently, our model assumptions are well-founded and our theory applies
 179 to multiple domains of wide interest.

180 4 Setup and Main Results

181 In this section, we present a brief description of our results and defer full formal details to Appendix A.
 182 For the sake of intuition, we can think of the data as images of different objects and the color of the
 183 object as a concept. We assume that the observed data X lies in a space $\mathcal{X} \subseteq \mathbb{R}^{d_x}$ of dimension d_x and
 184 has an underlying representation $X = f(Z)$ for latent variables Z that lie in a latent concept space
 185 \mathbb{R}^{d_z} of dimension d_z . We allow f to be an arbitrary nonlinearity that is injective and differentiable.

186 **Concepts** To motivate our definition, consider the color “red” as a concept. Different images have
 187 different levels of “redness” in them, so this concept is measured on a continuous scale, represented
 188 by a valuation $b \in \mathbb{R}$. We define an (atomic) concept to be represented by a vector $a \in \mathbb{R}^{d_z}$ such that
 189 $\langle a, Z \rangle = \langle a, f^{-1}(X) \rangle$ encodes the “value” of the concept in X . More precisely, for a given valuation
 190 $b \in \mathbb{R}$, the set of all observations X that satisfy this concept is given by $\{X = f(Z) | \langle a, Z \rangle = b\}$.
 191 Similarly, multi-dimensional concepts C (Appendix A) correspond to matrices A and vectors b . For
 192 a visualization, see Fig. 1.

193 **Concept conditional distributions** To define distributions of datasets over concepts, consider the
 194 case where we first collect a base dataset with some underlying distribution (e.g. a set of images

195 of all objects) and then collect concept datasets via filtering (e.g. to collect a dataset of dark red
196 colored objects, we filter them to only keep images of dark red colored objects). We call the
197 former the *base distribution* and the latter the *concept conditional distribution* corresponding to
198 our concept. Moreover, we allow for noise because humans are great at distilling concepts from
199 noisy images, e.g., we recognize cars in a misty environment. Formally, we have a noisy estimate
200 $\tilde{b} = \langle a, z \rangle + \epsilon$ where ϵ has density $q(\epsilon)$, independent of z . Then we consider the distribution
201 $p_C(z) = p(z|\tilde{b} = b) \propto p(\tilde{b} = b|z)p(z) = q(b - \langle a, z \rangle)p(z)$ where we used Bayes theorem in the
202 last step. We again extend these definitions to multi-dimensional concepts. The majority of recent
203 identifiability results relied on interventional data while we only consider conditional information here.
204 Therefore, our main problem of interest can be stated as follows: Given an observational dataset X^0
205 along with datasets X^1, \dots, X^m corresponding to concept conditional datasets for different concepts
206 C^1, \dots, C^m , under what conditions (and up to which symmetries) can we learn the concepts? This
207 is a more modest objective than learning the entire map f which is the usual goal in, say, CRL. While
208 the latter typically requires stringent assumptions, in particular $\Omega(d_z)$ environments are necessary,
209 our weaker identifiability results only need $O(d_C) \ll O(d_z)$ environments.

210 **Identifiability** Toward this end, a fundamental question is whether this problem is even possible,
211 i.e., whether it is well-defined. This is known as the question of identifiability [45, 21, 116, 49].
212 Informally, for the setting above, we say that the concepts $(C^1, A^1), \dots, (C^m, A^m)$ with associated
213 nonlinearity f are identifiable (and thus learnable) if for any other collection of different parameters
214 that fit the data, they are linearly related to the true parameters. Identifiability enables us to recover
215 the concepts of interest from our data, which is useful because they can then be used for further
216 downstream tasks such as controllable generative modeling.

217 **Main Result** To state our main result, our main assumptions are: (i) linear independence of the
218 concepts (since we want them to encode distinct concepts), (ii) Gaussianity of noise distribution
219 (conventional choice) and (iii) diversity of the environments (to motivate this, observe if two concepts
220 always occur together, it’s information-theoretically impossible to distinguish them, e.g., if an agent
221 only sees red large objects (i.e. all red objects are large and all large objects are red), it will be
222 unable to disambiguate the “red” concept from the “large” concept. Therefore, we need diversity of
223 environments to learn concepts, which we extract based on the signatures they leave on the datasets.)

224 **Theorem 1** (Informal). *Suppose we are given m context conditional datasets X^1, \dots, X^m and the*
225 *observational dataset X^0 such that the above assumptions hold. Then the concepts are identifiable.*

226 We defer formal technical details to Appendices A and B. Crucially, we only require a number of
227 datasets that depends only on the number of atoms n we wish to learn (in fact, $O(n)$ datasets), and not
228 on the underlying latent dimension d_z of the true generative process. This is a significant departure
229 from many existing works, since the true underlying generative process could have $d_z = 1000$, say,
230 whereas we may be interested to learn only $n = 5$ concepts, say. In this case, approaches based
231 on CRL necessitate at least ~ 1000 *interventional* datasets, whereas we show that $\sim n + 2 = 7$
232 *conditional* datasets are enough if we only want to learn the n atomic concepts. We will explain the
233 connection to CRL in Appendix C.

234 5 Experiments

235 In this section, we present experiments to validate and uti-
236 lize our framework. We first verify our results on synthetic
237 data, via a contrastive learning algorithm for concept learn-
238 ing. Then, we focus on experiments involving real-world
239 settings, in particular on image data using multimodal
240 CLIP models and text data using large language models
241 (LLMs).

242 **End-to-end Contrastive learning algorithm and Syn-
243 thetic experiments** We validate our framework on syn-
244 thetic data as follows. We sample the base distribution
245 from a Gaussian Mixture model and experiment with both
246 linear and nonlinear mixing functions (details deferred to

Mixing (f)	(n, d_z, d_x)	$R^2 \uparrow$	MCC \uparrow
Linear	(2, 3, 4)	0.98 ± 0.01	0.98 ± 0.03
Nonlinear	(2, 3, 4)	0.94 ± 0.06	0.96 ± 0.04
Linear	(3, 4, 6)	0.99 ± 0.01	0.86 ± 0.08
Nonlinear	(3, 4, 6)	0.97 ± 0.03	0.92 ± 0.07
Linear	(4, 8, 10)	0.97 ± 0.01	0.87 ± 0.06
Nonlinear	(4, 8, 10)	0.94 ± 0.03	0.87 ± 0.06

Table 1: Linear identifiability when number of concepts n is less than underlying latent dimension d_z with observed dimension d_x , averaged over 5 seeds.

247 Appendix H). The number of concepts n is intentionally chosen to be less than the ground truth
 248 dimension d_z and the number of concepts is $m = n + 1$ as per our theory. Inspired by [12], we
 249 use a contrastive learning algorithm to extract the concepts, with details deferred to Appendix G. In
 250 Table 1, we report the R^2 and Mean Correlation Coefficient (MCC) metrics [45, 46] with respect to
 251 the ground truth concept valuations. There are no baselines since we are in a novel setting, but our
 252 metrics are comparable to and often surpass what’s usually reported in such highly nonlinear settings
 253 [119, 12].

254 **Probing the theory on multimodal CLIP models** A real world example that approximately
 255 matches the setting considered in this paper is the training of the multimodal CLIP models [81]. They
 256 are trained by aligning the embeddings of images and their captions. We can view the caption as
 257 an indicator of the concepts present in the image. Thus the data provides access to several concept
 258 conditional distributions such as the collection of all images having the label ‘A dog’, but also to
 259 more complex distributions consisting of more than one atomic concept such as images labeled ‘A
 260 red flower’. We embed images from the 3d-Shapes Dataset [14] with known factors of variation
 261 into the latent space of two different pretrained CLIP models. Using logistic regression we learn
 262 atomic concepts for each of the factors of variations (see Appendix E.1 for details) and then evaluate
 263 the concept valuations of the learned atomic concept on held out images. We show the results for
 264 the shape attribute in Figure 2 (further results are in Appendix E.2). The results show that there
 265 are indeed linear subspaces of the embeddings space that represent certain concepts. Moreover, the
 learned valuations for different models are approximately linearly related as predicted by Theorem 2.

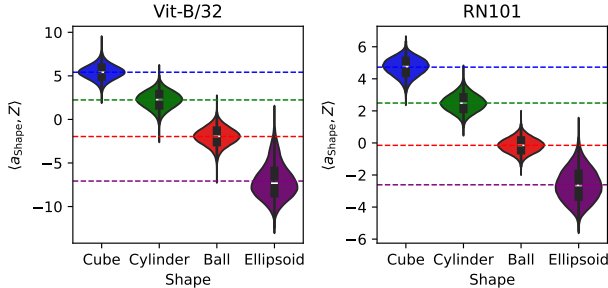


Figure 2: Violin plot of the concept valuations $\langle a_{\text{Shape}}, Z \rangle$ for the different shapes and a vision transformer CLIP embedding (left) and a residual network CLIP embedding (right). Results show concentration of the concept valuations around the concept planes indicated by the horizontal lines.

266

267 **Alignment of LLMs** Finally, we show an application of our framework to interpret representations
 268 of LLMs and improve alignment techniques. In particular, we exploit our ideas to improve the
 269 Inference-Time Intervention technique [56] to promote LLMs to be more truthful, i.e. the downstream
 270 task is to take pre-trained LLMs and during inference, change the valuation of the truthfulness concept
 271 from *false* to *true*, without affecting any other orthogonal concepts. Motivated by our framework,
 272 we propose to replace steering vectors by steering matrices for better alignment. Experiments on
 273 LLaMA [106] show an improvement of the TruthfulQA dataset [58] accuracy. Additional details,
 274 including a self-contained introduction to large language models (LLMs) and the Inference-Time
 275 Intervention (ITI) technique are deferred to Appendix F.

276 6 Conclusion

277 In this work, we study the problem of extracting concepts from data, inspired by techniques from
 278 causal representation learning. For this, we geometrically define concepts as linear subspaces, well-
 279 supported via extensive empirical literature. With this formal definition of concepts, we study under
 280 what conditions they can be provably recovered from data. Our rigorous results show that this
 281 is possible under the presence of only conditional data, requiring far fewer distributions than the
 282 underlying latent dimension. Finally, synthetic experiments, multimodal CLIP experiments and LLM
 283 alignment experiments verify and showcase the utility of our ideas.

284 **References**

- 285 [1] K. Ahuja, D. Mahajan, Y. Wang, and Y. Bengio. Interventional causal representation learning. In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*.
286 JMLR.org, 2023.
287
- 288 [2] K. Ahuja, A. Mansouri, and Y. Wang. Multi-domain causal representation learning via weak
289 distributional invariances. *arXiv preprint arXiv:2310.02854*, 2023.
- 290 [3] C. Allen and T. Hospedales. Analogies explained: Towards understanding word embeddings.
291 In *International Conference on Machine Learning*, pages 223–231. PMLR, 2019.
- 292 [4] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. Invariant risk minimization. *arXiv*
293 *preprint arXiv:1907.02893*, 2019.
- 294 [5] S. Arora, Y. Li, Y. Liang, T. Ma, and A. Risteski. A latent variable model approach to pmi-
295 based word embeddings. *Transactions of the Association for Computational Linguistics*, 4:
296 385–399, 2016.
- 297 [6] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli,
298 T. Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from
299 human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 300 [7] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirho-
301 seini, C. McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint*
302 *arXiv:2212.08073*, 2022.
- 303 [8] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba. Network dissection: Quantifying
304 interpretability of deep visual representations. In *Proceedings of the IEEE conference on*
305 *computer vision and pattern recognition*, pages 6541–6549, 2017.
- 306 [9] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new per-
307 spectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828,
308 2013.
- 309 [10] J. Brehmer, P. De Haan, P. Lippe, and T. S. Cohen. Weakly supervised causal representation
310 learning. *Advances in Neural Information Processing Systems*, 35:38319–38331, 2022.
- 311 [11] S. Buchholz, M. Besserve, and B. Schölkopf. Function classes for identifiable nonlinear
312 independent component analysis. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors,
313 *Advances in Neural Information Processing Systems*, 2022. URL [https://openreview.](https://openreview.net/forum?id=DpKaP-PY8bK)
314 [net/forum?id=DpKaP-PY8bK](https://openreview.net/forum?id=DpKaP-PY8bK).
- 315 [12] S. Buchholz, G. Rajendran, E. Rosenfeld, B. Aragam, B. Schölkopf, and P. Ravikumar.
316 Learning linear causal representations from interventions under general nonlinear mixing.
317 *arXiv preprint arXiv:2306.02235*, 2023.
- 318 [13] P. Bühlmann. Invariance, causality and robustness. *Statistical Science*, 35(3):404–426, 2020.
- 319 [14] C. Burgess and H. Kim. 3d shapes dataset. <https://github.com/deepmind/3dshapes-dataset/>,
320 2018.
- 321 [15] C. Burns, H. Ye, D. Klein, and J. Steinhardt. Discovering latent knowledge in language models
322 without supervision. *arXiv preprint arXiv:2212.03827*, 2022.
- 323 [16] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso. Machine learning interpretability: A survey
324 on methods and metrics. *Electronics*, 8(8):832, 2019.
- 325 [17] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang,
326 J. E. Gonzalez, I. Stoica, and E. P. Xing. Vicuna: An open-source chatbot impress-
327 ing gpt-4 with 90%* chatgpt quality, March 2023. URL [https://lmsys.org/blog/](https://lmsys.org/blog/2023-03-30-vicuna/)
328 [2023-03-30-vicuna/](https://lmsys.org/blog/2023-03-30-vicuna/).
- 329 [18] P. Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):
330 287–314, 1994.

- 331 [19] A. Conneau, G. Kruszewski, G. Lample, L. Barrault, and M. Baroni. What you can cram
332 into a single vector: Probing sentence embeddings for linguistic properties. *arXiv preprint*
333 *arXiv:1805.01070*, 2018.
- 334 [20] J. Cui, W. Huang, Y. Wang, and Y. Wang. Aggnce: Asymptotically identifiable contrastive
335 learning. In *NeurIPS Workshop*, 2022.
- 336 [21] A. D’Amour, K. Heller, D. Moldovan, B. Adlam, B. Alipanahi, A. Beutel, C. Chen, J. Deaton,
337 J. Eisenstein, M. D. Hoffman, et al. Underspecification presents challenges for credibility in
338 modern machine learning. *The Journal of Machine Learning Research*, 23(1):10237–10297,
339 2022.
- 340 [22] N. Dilokthanakul, P. A. Mediano, M. Garnelo, M. C. Lee, H. Salimbeni, K. Arulkumaran, and
341 M. Shanahan. Deep unsupervised clustering with gaussian mixture variational autoencoders.
342 *arXiv preprint arXiv:1611.02648*, 2016.
- 343 [23] M. Drton, B. Sturmfels, and S. Sullivant. *Lectures on Algebraic Statistics*, volume 39 of
344 *Oberwolfach Seminars*. Springer, 2009. doi: 10.1007/978-3-7643-8905-5.
- 345 [24] N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen,
346 T. Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits*
347 *Thread*, 1, 2021.
- 348 [25] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds,
349 R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wat-
350 tenberg, and C. Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022.
351 https://transformer-circuits.pub/2022/toy_model/index.html.
- 352 [26] J. Engel, M. Hoffman, and A. Roberts. Latent constraints: Learning to generate conditionally
353 from unconditional generative models. *arXiv preprint arXiv:1711.05772*, 2017.
- 354 [27] K. Ethayarajh, D. Duvenaud, and G. Hirst. Towards understanding linear word analogies.
355 *arXiv preprint arXiv:1810.04882*, 2018.
- 356 [28] F. Falck, H. Zhang, M. Willetts, G. Nicholson, C. Yau, and C. C. Holmes. Multi-facet clustering
357 variational autoencoders. *Advances in Neural Information Processing Systems*, 34, 2021.
- 358 [29] D. Ganguli, L. Lovitt, J. Kernion, A. Askell, Y. Bai, S. Kadavath, B. Mann, E. Perez,
359 N. Schiefer, K. Ndousse, et al. Red teaming language models to reduce harms: Methods,
360 scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*, 2022.
- 361 [30] A. Gittens, D. Achlioptas, and M. W. Mahoney. Skip-gram- zipf+ uniform= vector additivity.
362 In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*
363 *(Volume 1: Long Papers)*, pages 69–76, 2017.
- 364 [31] L. Gresele, J. Von Kügelgen, V. Stimper, B. Schölkopf, and M. Besserve. Independent
365 mechanism analysis, a new concept? *Advances in Neural Information Processing Systems*, 34,
366 2021.
- 367 [32] S. Gupta, S. Jegelka, D. Lopez-Paz, and K. Ahuja. Context is environment. *arXiv e-prints*,
368 pages arXiv–2309, 2023.
- 369 [33] W. Gurnee, N. Nanda, M. Pauly, K. Harvey, D. Troitskii, and D. Bertsimas. Finding neurons
370 in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*, 2023.
- 371 [34] E. Hernandez, B. Z. Li, and J. Andreas. Measuring and manipulating knowledge representations
372 in language models. *arXiv preprint arXiv:2304.00740*, 2023.
- 373 [35] D. Horan, E. Richardson, and Y. Weiss. When is unsupervised disentanglement possible?
374 *Advances in Neural Information Processing Systems*, 34:5150–5161, 2021.
- 375 [36] A. Hyvarinen and H. Morioka. Unsupervised feature extraction by time-contrastive learning
376 and nonlinear ica. *Advances in neural information processing systems*, 29, 2016.

- 377 [37] A. Hyvärinen and E. Oja. Independent component analysis: algorithms and applications.
378 *Neural networks*, 13(4-5):411–430, 2000.
- 379 [38] A. Hyvärinen and P. Pajunen. Nonlinear independent component analysis: Existence and
380 uniqueness results. *Neural networks*, 12(3):429–439, 1999.
- 381 [39] A. Hyvarinen, J. Karhunen, and E. Oja. Independent component analysis. *Studies in informatics
382 and control*, 11(2):205–207, 2002.
- 383 [40] A. Hyvarinen, H. Sasaki, and R. Turner. Nonlinear ica using auxiliary variables and generalized
384 contrastive learning. In *The 22nd International Conference on Artificial Intelligence and
385 Statistics*, pages 859–868. PMLR, 2019.
- 386 [41] A. Hyvärinen, I. Khemakhem, and R. Monti. Identifiability of latent-variable and structural-
387 equation models: from linear to nonlinear. *arXiv preprint arXiv:2302.02672*, 2023.
- 388 [42] Y. Jiang and B. Aragam. Learning latent causal graphs with unknown interventions. In
389 *Advances in Neural Information Processing Systems*, 2023.
- 390 [43] Y. Jiang, B. Aragam, and V. Veitch. Uncovering meanings of embeddings via partial orthogo-
391 nality. *Advances in Neural Information Processing Systems*, 2023.
- 392 [44] Y. Jiang, G. Rajendran, P. Ravikumar, B. Aragam, and V. Veitch. On the origins of linear
393 representations in large language models. *arXiv preprint*, 2024.
- 394 [45] I. Khemakhem, D. Kingma, R. Monti, and A. Hyvarinen. Variational autoencoders and
395 nonlinear ica: A unifying framework. In *International Conference on Artificial Intelligence
396 and Statistics*, pages 2207–2217. PMLR, 2020.
- 397 [46] I. Khemakhem, R. Monti, D. Kingma, and A. Hyvarinen. Ice-beem: Identifiable conditional
398 energy-based deep models based on nonlinear ica. *Advances in Neural Information Processing
399 Systems*, 33:12768–12778, 2020.
- 400 [47] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viegas, et al. Interpretability beyond
401 feature attribution: Quantitative testing with concept activation vectors (tcav). In *International
402 conference on machine learning*, pages 2668–2677. PMLR, 2018.
- 403 [48] B. Kivva, G. Rajendran, P. Ravikumar, and B. Aragam. Learning latent causal graphs via
404 mixture oracles. *Advances in Neural Information Processing Systems*, 34:18087–18101, 2021.
- 405 [49] B. Kivva, G. Rajendran, P. Ravikumar, and B. Aragam. Identifiability of deep generative
406 models without auxiliary information. *Advances in Neural Information Processing Systems*,
407 35:15687–15701, 2022.
- 408 [50] L. Kong, S. Xie, W. Yao, Y. Zheng, G. Chen, P. Stojanov, V. Akinwande, and K. Zhang. Partial
409 identifiability for domain adaptation. *arXiv preprint arXiv:2306.06510*, 2023.
- 410 [51] S. Lachapelle, P. Rodríguez, Y. Sharma, K. Everett, R. L. Priol, A. Lacoste, and S. Lacoste-
411 Julien. Disentanglement via mechanism sparsity regularization: A new principle for nonlinear
412 ICA. In B. Schölkopf, C. Uhler, and K. Zhang, editors, *1st Conference on Causal Learning
413 and Reasoning, CLear 2022, Sequoia Conference Center, Eureka, CA, USA, 11-13 April,
414 2022*, volume 177 of *Proceedings of Machine Learning Research*, pages 428–484. PMLR,
415 2022. URL <https://proceedings.mlr.press/v177/lachapelle22a.html>.
- 416 [52] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- 417 [53] T. Leemann, M. Kirchhof, Y. Rong, E. Kasneci, and G. Kasneci. When are post-hoc conceptual
418 explanations identifiable? In *Uncertainty in Artificial Intelligence*, pages 1207–1218. PMLR,
419 2023.
- 420 [54] O. Levy and Y. Goldberg. Neural word embedding as implicit matrix factorization. *Advances
421 in neural information processing systems*, 27, 2014.

- 422 [55] K. Li, A. K. Hopkins, D. Bau, F. Viégas, H. Pfister, and M. Wattenberg. Emergent world
423 representations: Exploring a sequence model trained on a synthetic task. *arXiv preprint*
424 *arXiv:2210.13382*, 2022.
- 425 [56] K. Li, O. Patel, F. Viégas, H. Pfister, and M. Wattenberg. Inference-time intervention: Eliciting
426 truthful answers from a language model. *arXiv preprint arXiv:2306.03341*, 2023.
- 427 [57] S. Li, B. Hooi, and G. H. Lee. Identifying through flows for recovering latent representations. In
428 *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia,*
429 *April 26-30, 2020*. OpenReview.net, 2020. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Sk10UpEYvB)
430 [Sk10UpEYvB](https://openreview.net/forum?id=Sk10UpEYvB).
- 431 [58] S. Lin, J. Hilton, and O. Evans. Truthfulqa: Measuring how models mimic human falsehoods.
432 *arXiv preprint arXiv:2109.07958*, 2021.
- 433 [59] P. Lippe, S. Magliacane, S. Löwe, Y. M. Asano, T. Cohen, and E. Gavves. Biscuit: Causal
434 representation learning from binary interactions. *arXiv preprint arXiv:2306.09643*, 2023.
- 435 [60] Y. Liu, Z. Zhang, D. Gong, M. Gong, B. Huang, A. v. d. Hengel, K. Zhang, and J. Q. Shi.
436 Identifying weight-variant latent causal models. *arXiv preprint arXiv:2208.14153*, 2022.
- 437 [61] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Schölkopf, and O. Bachem. Chal-
438 lenging common assumptions in the unsupervised learning of disentangled representations. In
439 *international conference on machine learning*, pages 4114–4124. PMLR, 2019.
- 440 [62] I. Loshchilov and F. Hutter. SGDR: stochastic gradient descent with warm restarts. In
441 *5th International Conference on Learning Representations, ICLR 2017, Toulon, France,*
442 *April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Skq89Scxx)
443 [Skq89Scxx](https://openreview.net/forum?id=Skq89Scxx).
- 444 [63] C. Lu, Y. Wu, J. M. Hernández-Lobato, and B. Schölkopf. Invariant causal representation
445 learning for out-of-distribution generalization. In *International Conference on Learning*
446 *Representations*, 2021.
- 447 [64] E. Marconato, A. Passerini, and S. Teso. Interpretability is in the mind of the beholder: A
448 causal framework for human-interpretable representation learning. *Entropy*, 25(12):1574,
449 2023.
- 450 [65] T. McGrath, A. Kapishnikov, N. Tomašev, A. Pearce, M. Wattenberg, D. Hassabis, B. Kim,
451 U. Paquet, and V. Kramnik. Acquisition of chess knowledge in alphazero. *Proceedings of the*
452 *National Academy of Sciences*, 119(47):e2206625119, 2022.
- 453 [66] K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in
454 gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- 455 [67] T. Mikolov, W.-t. Yih, and G. Zweig. Linguistic regularities in continuous space word
456 representations. In *Proceedings of the 2013 conference of the north american chapter of the*
457 *association for computational linguistics: Human language technologies*, pages 746–751,
458 2013.
- 459 [68] G. E. Moran, D. Sridhar, Y. Wang, and D. Blei. Identifiable deep generative models via sparse
460 decoding. *Transactions on Machine Learning Research*, 2022.
- 461 [69] L. Moschella, V. Maiorca, M. Fumero, A. Norelli, F. Locatello, and E. Rodola. Relative
462 representations enable zero-shot latent space communication. *arXiv preprint arXiv:2209.15430*,
463 2022.
- 464 [70] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju,
465 W. Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv*
466 *preprint arXiv:2112.09332*, 2021.
- 467 [71] N. Nanda, A. Lee, and M. Wattenberg. Emergent linear representations in world models of
468 self-supervised sequence models. *arXiv preprint arXiv:2309.00941*, 2023.

- 469 [72] C. Olah. Mechanistic interpretability, variables, and the importance of interpretable
470 bases. <https://transformer-circuits.pub/2022/mech-interp-essay/index.html>, 2022.
471
- 472 [73] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
473 K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback.
474 *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- 475 [74] K. Park, Y. J. Choe, and V. Veitch. The linear representation hypothesis and the geometry of
476 large language models. *arXiv preprint arXiv:2311.03658*, 2023.
- 477 [75] J. Pearl. *Causality*. Cambridge university press, 2009.
- 478 [76] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation.
479 In *Proceedings of the 2014 conference on empirical methods in natural language processing*
480 (*EMNLP*), pages 1532–1543, 2014.
- 481 [77] J. Peters, D. Janzing, and B. Schölkopf. *Elements of causal inference: foundations and*
482 *learning algorithms*. The MIT Press, 2017.
- 483 [78] E. Poeta, G. Ciravegna, E. Pastor, T. Cerquitelli, and E. Baralis. Concept-based explainable
484 artificial intelligence: A survey. *arXiv preprint arXiv:2312.12936*, 2023.
- 485 [79] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep
486 convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- 487 [80] A. Radford, R. Jozefowicz, and I. Sutskever. Learning to generate reviews and discovering
488 sentiment. *arXiv preprint arXiv:1704.01444*, 2017.
- 489 [81] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
490 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervi-
491 sion. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 492 [82] R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. Direct prefer-
493 ence optimization: Your language model is secretly a reward model. *arXiv preprint*
494 *arXiv:2305.18290*, 2023.
- 495 [83] M. Raghu, J. Gilmer, J. Yosinski, and J. Sohl-Dickstein. Svcca: Singular vector canonical
496 correlation analysis for deep understanding and improvement. *stat*, 1050:19, 2017.
- 497 [84] G. Rajendran, B. Kivva, M. Gao, and B. Aragam. Structure learning in polynomial time:
498 Greedy algorithms, bregman information, and exponential families. *Advances in Neural*
499 *Information Processing Systems*, 34:18660–18672, 2021.
- 500 [85] G. Rajendran, P. Reizinger, W. Brendel, and P. Ravikumar. An interventional perspective on
501 identifiability in gaussian lti systems with independent component analysis. *arXiv preprint*
502 *arXiv:2311.18048*, 2023.
- 503 [86] N. Reimers and I. Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks.
504 In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*.
505 Association for Computational Linguistics, 11 2019. URL [https://arxiv.org/abs/1908.](https://arxiv.org/abs/1908.10084)
506 10084.
- 507 [87] N. Rimsy, N. Gabrieli, J. Schulz, M. Tong, E. Hubinger, and A. M. Turner. Steering llama 2
508 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023.
- 509 [88] D. Rothenhäusler, N. Meinshausen, P. Bühlmann, and J. Peters. Anchor regression: Hetero-
510 geneous data meet causality. *Journal of the Royal Statistical Society Series B: Statistical*
511 *Methodology*, 83(2):215–246, 2021.
- 512 [89] B. Schölkopf and J. von Kügelgen. From statistical to causal learning. In *Proceedings of the*
513 *International Congress of Mathematicians (ICM)*. EMS Press, July 2022.

- 514 [90] B. Schölkopf, F. Locatello, S. Bauer, N. R. Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio.
515 Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
516 arXiv:2102.11107.
- 517 [91] L. Schut, N. Tomasev, T. McGrath, D. Hassabis, U. Paquet, and B. Kim. Bridging the human-ai
518 knowledge gap: Concept discovery and transfer in alphazero. *arXiv preprint arXiv:2310.16410*,
519 2023.
- 520 [92] Y. Seonwoo, S. Park, D. Kim, and A. Oh. Additive compositionality of word vectors. In
521 *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 387–396,
522 2019.
- 523 [93] X. Shen, F. Liu, H. Dong, Q. Lian, Z. Chen, and T. Zhang. Weakly supervised disentangled
524 generative causal representation learning. *Journal of Machine Learning Research*, 23:1–55,
525 2022.
- 526 [94] K. Shuster, S. Poff, M. Chen, D. Kiela, and J. Weston. Retrieval augmentation reduces
527 hallucination in conversation. *arXiv preprint arXiv:2104.07567*, 2021.
- 528 [95] P. Sorrenson, C. Rother, and U. Köthe. Disentanglement by nonlinear ica with general
529 incompressible-flow networks (gin). *arXiv preprint arXiv:2001.04872*, 2020.
- 530 [96] P. Spirtes, C. N. Glymour, and R. Scheines. *Causation, prediction, and search*. MIT press,
531 2000.
- 532 [97] C. Squires and C. Uhler. Causal structure learning: a combinatorial perspective. *Foundations*
533 *of Computational Mathematics*, pages 1–35, 2022.
- 534 [98] C. Squires, A. Seigal, S. S. Bhate, and C. Uhler. Linear causal disentanglement via interven-
535 tions. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors,
536 *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu,*
537 *Hawaii, USA, volume 202 of Proceedings of Machine Learning Research*, pages 32540–32560.
538 PMLR, 2023. URL <https://proceedings.mlr.press/v202/squires23a.html>.
- 539 [99] N. Subramani, N. Suresh, and M. E. Peters. Extracting latent steering vectors from pretrained
540 language models. *arXiv preprint arXiv:2205.05124*, 2022.
- 541 [100] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus.
542 Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- 543 [101] A. Taeb, N. Ruggeri, C. Schnuck, and F. Yang. Provable concept learning for interpretable
544 predictions using variational autoencoders. In *ICML 2022 2nd AI for Science Workshop*, 2022.
- 545 [102] D. Talon, P. Lippe, S. James, A. Del Bue, and S. Magliacane. Towards the reusability and
546 compositionality of causal representations. In *Causal Representation Learning Workshop at*
547 *NeurIPS 2023*, 2023.
- 548 [103] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto.
549 Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on*
550 *Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7, 2023.
- 551 [104] I. Tenney, D. Das, and E. Pavlick. Bert rediscovers the classical nlp pipeline. *arXiv preprint*
552 *arXiv:1905.05950*, 2019.
- 553 [105] C. Tigges, O. J. Hollinsworth, A. Geiger, and N. Nanda. Linear representations of sentiment
554 in large language models. *arXiv preprint arXiv:2310.15154*, 2023.
- 555 [106] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière,
556 N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models.
557 *arXiv preprint arXiv:2302.13971*, 2023.
- 558 [107] M. Trager, P. Perera, L. Zancato, A. Achille, P. Bhatia, and S. Soatto. Linear spaces of mean-
559 ings: compositional structures in vision-language models. In *Proceedings of the IEEE/CVF*
560 *International Conference on Computer Vision*, pages 15395–15404, 2023.

- 561 [108] A. Turner, L. Thiergart, D. Udell, G. Leech, U. Mini, and M. MacDiarmid. Activation addition:
562 Steering language models without optimization. *arXiv preprint arXiv:2308.10248*, 2023.
- 563 [109] B. Varici, K. Shanmugam, P. Sattigeri, and A. Tajer. Intervention target estimation in the
564 presence of latent variables. In *Uncertainty in Artificial Intelligence*, pages 2013–2023. PMLR,
565 2022.
- 566 [110] B. Varici, E. Acarturk, K. Shanmugam, A. Kumar, and A. Tajer. Score-based causal represen-
567 tation learning with interventions. *arXiv preprint arXiv:2301.08230*, 2023.
- 568 [111] B. Varıcı, E. Acartürk, K. Shanmugam, and A. Tajer. Score-based causal representation
569 learning: Linear and general transformations. *arXiv preprint arXiv:2402.00849*, 2024.
- 570 [112] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
571 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
572 30, 2017.
- 573 [113] J. Von Kügelgen, Y. Sharma, L. Gresele, W. Brendel, B. Schölkopf, M. Besserve, and F. Lo-
574 catello. Self-supervised learning with data augmentations provably isolates content from style.
575 *Advances in Neural Information Processing Systems*, 34, 2021.
- 576 [114] J. von Kügelgen, M. Besserve, W. Liang, L. Gresele, A. Kekić, E. Bareinboim, D. M. Blei,
577 and B. Schölkopf. Nonparametric identifiability of causal representations from unknown
578 interventions. In *Advances in Neural Information Processing Systems*, 2023.
- 579 [115] K. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt. Interpretability in the wild:
580 a circuit for indirect object identification in gpt-2 small. *arXiv preprint arXiv:2211.00593*,
581 2022.
- 582 [116] Y. Wang, D. Blei, and J. P. Cunningham. Posterior collapse and latent variable non-
583 identifiability. *Advances in Neural Information Processing Systems*, 34:5443–5455, 2021.
- 584 [117] Z. Wang, L. Gui, J. Negrea, and V. Veitch. Concept algebra for score-based conditional model.
585 In *ICML 2023 Workshop on Structured Probabilistic Inference & Generative Modeling*,
586 2023.
- 587 [118] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-
588 thought prompting elicits reasoning in large language models. *Advances in Neural Information*
589 *Processing Systems*, 35:24824–24837, 2022.
- 590 [119] M. Willetts and B. Paige. I don’t need u: Identifiable non-linear ica without side information.
591 *arXiv preprint arXiv:2106.05238*, 2021.
- 592 [120] D. Xu, D. Yao, S. Lachapelle, P. Taslakian, J. von Kügelgen, F. Locatello, and S. Magliacane.
593 A sparsity principle for partially observable causal representation learning. *arXiv preprint*
594 *arXiv:2403.08335*, 2024.
- 595 [121] M. Yang, F. Liu, Z. Chen, X. Shen, J. Hao, and J. Wang. Causalvae: Disentangled representa-
596 tion learning via neural structural causal models. In *Proceedings of the IEEE/CVF Conference*
597 *on Computer Vision and Pattern Recognition (CVPR)*, pages 9593–9602, June 2021.
- 598 [122] Y. Yang, A. Panagopoulou, S. Zhou, D. Jin, C. Callison-Burch, and M. Yatskar. Language in a
599 bottle: Language model guided concept bottlenecks for interpretable image classification. In
600 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
601 19187–19197, 2023.
- 602 [123] D. Yao, D. Xu, S. Lachapelle, S. Magliacane, P. Taslakian, G. Martius, J. von Kügelgen,
603 and F. Locatello. Multi-view causal representation learning with partial observability. *arXiv*
604 *preprint arXiv:2311.04056*, 2023.
- 605 [124] F. Zhang and N. Nanda. Towards best practices of activation patching in language models:
606 Metrics and methods. *arXiv preprint arXiv:2309.16042*, 2023.

- 607 [125] J. Zhang, C. Squires, K. Greenewald, A. Srivastava, K. Shanmugam, and C. Uhler. Identifiability guarantees for causal disentanglement from soft interventions. *arXiv preprint arXiv:2307.06250*, 2023.
- 608
609
- 610 [126] Y. Zhang, Y. Du, B. Huang, Z. Wang, J. Wang, M. Fang, and M. Pechenizkiy. Interpretable reward redistribution in reinforcement learning: A causal approach. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- 611
612
- 613 [127] Y. Zheng, I. Ng, and K. Zhang. On the identifiability of nonlinear ica: Sparsity and beyond. *Advances in Neural Information Processing Systems*, 35:16411–16422, 2022.
- 614
- 615 [128] R. S. Zimmermann, Y. Sharma, S. Schneider, M. Bethge, and W. Brendel. Contrastive learning inverts the data generating process. In *International Conference on Machine Learning*, pages 12979–12990. PMLR, 2021.
- 616
617
- 618 [129] A. Zou, L. Phan, S. Chen, J. Campbell, P. Guo, R. Ren, A. Pan, X. Yin, M. Mazeika, A.-K. Dombrowski, et al. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.
- 619
620

621 A Setup and Main Results

622 In this section, we provide a formal definition of concepts, which are high-level abstractions present
623 in data. This allows us to develop a theoretical framework for associated data distributions and
624 identifiability theory. For the sake of intuition, we can think of the data as images of different objects
625 and the color of the object as a concept.

626 A.1 Generative model

627 We assume that the observed data X lies in a space $\mathcal{X} \subseteq \mathbb{R}^{d_x}$ of dimension d_x and has an underlying
628 representation $X = f(Z)$ for latent variables Z that lie in a latent concept space \mathbb{R}^{d_z} of dimension d_z .
629 In contrast to most prior works we do not necessarily assume that Z represents the true underlying
630 mechanism that generated the data. Instead we simply assume that the latent representation has the
631 geometric property that it maps certain regions of the observation space to linear subspaces of the
632 latent space (motivated by previous work; see Section 3). Our first assumption is standard:

633 **Assumption 1** (Mixing function). *The non-linear f is injective and differentiable.*

634 We make no additional assumptions on f : The map from $Z \rightarrow X$ can be arbitrarily non-linear.

635 We now define concepts living in the latent space \mathbb{R}^{d_z} . Before presenting the general definition of
636 multidimensional concepts, we outline the basic ideas in the simplified setting of a one-dimensional
637 concept. Consider the color “red” as a concept. Different images have different levels of “redness”
638 in them, so this concept is measured on a continuous scale, represented by a valuation $b \in \mathbb{R}$. An
639 (atomic) concept is then represented by a vector $a \in \mathbb{R}^{d_z}$ such that $\langle a, Z \rangle = \langle a, f^{-1}(X) \rangle$ encodes
640 the “value” of the concept in X , as measured in the latent space. More precisely, for a given valuation
641 $b \in \mathbb{R}$, the set of all observations X that satisfy this concept is given by $\{X = f(Z) | \langle a, Z \rangle = b\}$.
642 For instance, for an object in an image X , if $a \in \mathbb{R}^{d_z}$ is the concept of red color, $b \in \mathbb{R}$ could indicate
643 the intensity; then all datapoints X satisfying this concept, i.e., all images with an object that has
644 color red with intensity b , can be characterized as $X = f(Z)$ where Z satisfies $\langle a, Z \rangle = b$. For a 3D
645 visualization, see Fig. 1. We make this intuition formal below.

646 **Definition 1** (Concepts). *A concept C is a linear transformation $A : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_C}$. The dimension of
647 the concept will be denoted by $\dim(C) = d_C$. A valuation is a vector $b \in \mathbb{R}^{d_C}$ and we say that a
648 datapoint X satisfies the concept C with valuation b if $AZ = b$ where $Z = f^{-1}(X)$.*

649 In this work, we are interested in learning a collection of m concepts C^1, \dots, C^m from observed
650 data. By left multiplying by the pseudo-inverse A^+ , we can equivalently assume A is a projector
651 matrix. However, the current definition is more suitable for embeddings of real models.

652 When we talk of learning concepts C , we are in particular interested in learning the evaluation map
653 $Af^{-1}(x)$. This is a more modest objective than learning the entire map f which is the usual goal in,
654 say, CRL. While the latter typically requires stringent assumptions, in particular $\Omega(d_z)$ environments
655 are necessary, our weaker identifiability results only need $O(d_C) \ll O(d_z)$ environments. To simplify
656 our analysis, we make use of the following definition:

657 **Definition 2** (Atoms). *An atom (short for atomic concept) is any concept C with $\dim(C) = 1$.*

658 The idea is that we can view each concept as being composed of atomic concepts in the following
659 sense: Atomic concepts are fundamental concepts that live in a space of co-dimension 1 in latent
660 space, and thus are equivalently defined by vectors $a \in \mathbb{R}^{d_z}$. For example, concepts such red color,
661 size of object, etc., may be atomic concepts. Any generic concept is then composed of a collection of
662 atomic concepts, e.g., the concept C of all small dark red objects will correspond to $\dim(C) = 2$
663 with row 1 corresponding to the atomic concept of red color with large valuation (dark red objects)
664 and row 2 corresponding to the atomic concept of object size with low valuation (small objects).

665 A.2 Data distributions

666 We now define the distributions of datasets over concepts. We will predominantly work with
667 distributions of Z over \mathbb{R}^{d_z} , as the resulting distribution of $X = f(Z)$ over \mathbb{R}^{d_x} can be obtained via
668 a simple change of variables.

669 To build intuition, consider the case where we first collect a base dataset with some underlying
670 distribution and then collect concept datasets via filtering. For instance, we could first collect a set of

671 images of all objects and then, to collect a dataset of dark red colored objects, we filter them to only
 672 keep images of dark red colored objects. We call the former the *base distribution* and the latter the
 673 *concept conditional distribution* corresponding to our concept.

674 Fix a nonlinearity f . We assume that the base data distribution is the distribution of $X = f(Z)$ with
 675 $Z \sim p$, where p is the underlying distribution on \mathbb{R}^{d_z} . In what follows, we will abuse notation and
 676 use p for both the distribution and the corresponding probability density which we assume exists. We
 677 make no further assumptions on p since we do not wish to model the collection of real-life datasets
 678 that have been collected from nature and which could be very arbitrary.

679 We now define the concept conditional distribution, which is a distribution over X that is induced
 680 by noisy observations of a particular concept at a particular valuation. Formally, assume we want
 681 to condition on some atomic concept $a \in \mathbb{R}^{d_z}$ with valuation b . It is reasonable to assume that this
 682 conditioning is a noisy operation. For instance, humans are great at distilling concepts from noisy
 683 images, e.g., they recognize cars in a misty environment. We formalize this by assuming that data
 684 collection is based on a noisy estimate $\tilde{b} = \langle a, z \rangle + \epsilon$ where ϵ is independent of z and its density is a
 685 symmetric distribution with density $q(\epsilon)$. Then we consider the distribution

$$p_C(z) = p(z|\tilde{b} = b) \propto p(\tilde{b} = b|z)p(z) = q(b - \langle a, z \rangle)p(z) \quad (1)$$

686 where we used Bayes theorem in the last step. This definition directly extends to higher dimensional
 687 concepts which are concisely defined as follows.

688 **Definition 3** (Concept conditional distribution). *For a concept C with associated linear map A and
 689 an arbitrary valuation $b \in \mathbb{R}^{\dim(C)}$, we define the concept conditional distribution to be the set of
 690 observations X respecting this concept, which is defined as the distribution of $X = f(Z)$ where
 691 $Z \sim p_C$ with*

$$p_C(Z) \propto p(Z) \prod_{k \leq \dim(C)} q((AZ - b)_k). \quad (2)$$

692 This is by no means the only possible definition, and we present feasible alternate definitions in
 693 Appendix D. We remark that our formulation is related to the iVAE setting [45] and the auxiliary
 694 variable setting for identifiable ICA in Hyvarinen et al. [40] and we discuss the relation later.
 695 The majority of recent identifiability results relied on interventional data while we only consider
 696 conditional information here.

697 A.3 Concept learning and identifiability

698 We are ready to define our main problem of interest.

699 **Problem 1.** *We are given an observational dataset $X^0 = f(Z^0)$ corresponding to the latent base
 700 distribution p along with datasets X^1, \dots, X^m corresponding to concept conditional datasets for
 701 different concepts C^1, \dots, C^m and corresponding valuations b^1, \dots, b^m over the same latent space
 702 \mathbb{R}^{d_z} with the same mixing f . Under what conditions (and up to which symmetries) can we learn
 703 the concepts C^1, \dots, C^m , which includes the linear maps A^1, \dots, A^m , and the concept valuations
 704 $A^1 f^{-1}(x), \dots, A^m f^{-1}(x)$?*

705 Toward this end, a fundamental question is whether this problem is even possible, i.e., whether
 706 it is well-defined. This is known as the question of identifiability [45, 21, 116, 49]. Therefore,
 707 we make the following definition. Informally, for the setting above, we say that the concepts
 708 $(C^1, A^1), \dots, (C^m, A^m)$ with associated nonlinearity f are identifiable (and thus learnable) if for
 709 any other collection of different parameters that fit the data, they are linearly related to the true
 710 parameters.

711 **Definition 4** (Identifiability). *Given datasets X^0, X^1, \dots, X^m corresponding to the observa-
 712 tional distribution and m concepts C^1, \dots, C^m with underlying latent base distribution p on
 713 \mathbb{R}^{d_z} , nonlinearity f , linear maps A^1, \dots, A^m and valuations b^1, \dots, b^m , we say the concepts
 714 are identifiable if the following holds: Consider any different collection of parameters $\tilde{f}, \tilde{d}_z, \tilde{p}$,
 715 concepts $(\tilde{C}^1, \tilde{A}^1), \dots, (\tilde{C}^m, \tilde{A}^m)$ and valuations $\tilde{b}^1, \dots, \tilde{b}^m$ that also generate the same observa-
 716 tions X^0, X^1, \dots, X^m . Then there exists a shift $w \in \mathbb{R}^{d_z}$, permutation matrices P^e and invertible*

717 diagonal matrices Λ^e such that for all e and x ,

$$\tilde{A}^e \tilde{f}^{-1}(x) = \Lambda^e P^e A^e (f^{-1}(x) + w), \quad (3)$$

718 i.e., we can evaluate the concept evaluations on the data up to linear reparametrizations. Moreover,
719 there exists a linear map $T : \mathbb{R}^{\tilde{d}_z} \rightarrow \mathbb{R}^{d_z}$ such that the concepts and their evaluations satisfy

$$\tilde{A}^e = P^e A^e T^{-1}, \quad \tilde{b}^e = \Lambda^e P^e (b^e - A^e w). \quad (4)$$

720 Identifiability implies we can identify the nonlinear map f^{-1} within the span of the subspace of the
721 concepts of interest, and therefore we can recover the concepts of interest from our data. That is, if
722 certain concepts are identifiable, then we will be able to learn these concept representations up to
723 linearity, even if they can be highly nonlinear functions of our data. Such concept discovery is useful
724 because they can then be used for further downstream tasks such as controllable generative modeling.

725 We emphasize that in contrast to previous work we are not aiming to identify f completely and
726 indeed, no stronger identifiability results on f can be expected. First, we cannot hope to resolve the
727 linear transformation ambiguity because the latent space is not directly observed. In other words, a
728 concept evaluation can be defined either as $\langle a, Z \rangle$ or as $\langle Ta, T^{-\top} Z \rangle$ for an invertible linear map T .
729 For the purposes of downstream tasks, however, this is fine since the learned concepts will still be
730 the same. Second, we cannot expect to recover f^{-1} outside the span of the concepts because we do
731 not manipulate the linear spaces outside the span therefore we do not learn this information from
732 our observed data so this is also tight. The permutation matrix captures the fact that the ordering
733 of the concepts does not matter. Therefore, this definition captures the most general identifiability
734 guarantee that we can hope for in our setting and furthermore, this suffices for downstream tasks such
735 as controllable data generation.

736 Because we will only be interested in recovering the set of concepts up to linear transformations,
737 without loss of generality, we will fix the base collection of atomic concepts. That is, we assume
738 that each concept C^e corresponds to a linear map A^e whose rows are a subset of \mathcal{C} , where $\mathcal{C} =$
739 $\{a_1, \dots, a_n\}$ is a set of atomic concepts that we wish to learn. Moreover, we assume that they are
740 linearly independent, since we want them to encode distinct concepts. This is formalized as follows.

741 **Assumption 2.** *There exists a set of atomic concepts $\mathcal{C} = \{a_1, \dots, a_n\}$ of linearly independent*
742 *vectors such that for each concept C^e under consideration the rows of the concept matrix A^e are*
743 *contained in \mathcal{C} , i.e., $(A^e)^t e_i \in \mathcal{C}$. We denote the indices of the subset of \mathcal{C} that appear as rows of A^e*
744 *by S^e and we assume that all concepts in \mathcal{C} appear in some environment e , i.e., $\bigcup_e S^e = [n]$.*

745 **Remark 1.** *Definition 4 implies that the atoms can be identified in the sense that there is a permutation*
746 *$\pi \in S_n$ and $\lambda_i \neq 0$ such that for T as in Definition 4 and some λ_i*

$$\tilde{a}_{\pi(i)}^\top = a_i^\top T^{-1} \quad (5)$$

$$\langle \tilde{a}_{\pi(i)}, \tilde{f}^{-1}(x) \rangle = \lambda_i (\langle a_i, f^{-1}(x) \rangle + \langle a_i, w \rangle), \quad (6)$$

747 i.e., we can evaluate the valuations of the atomic concepts up to linear reparametrization.

748 A.4 Main Result

749 In this section, we present our main result on identifying concepts from data. The punchline is that
750 when we have rich datasets, i.e., sufficiently rich concept conditional datasets, then we can recover
751 the concepts. Crucially, we only require a number of datasets that depends only on the number of
752 atoms n we wish to learn (in fact, $O(n)$ datasets), and not on the underlying latent dimension d_z
753 of the true generative process. This is a significant departure from many existing works, since the
754 true underlying generative process could have $d_z = 1000$, say, whereas we may be interested to
755 learn only $n = 5$ concepts, say. In this case, approaches based on CRL necessitate at least ~ 1000
756 *interventional* datasets, whereas we show that $\sim n + 2 = 7$ *conditional* datasets are enough if we
757 only want to learn the n atomic concepts. We will explain the connection to CRL in Appendix C. Let
758 us now discuss our main assumptions.

759 **Assumption 3.** *The noise distribution q is Gaussian, i.e. $q \sim N(0, \sigma^2)$ for some $\sigma^2 > 0$.*

760 We choose Gaussian noise since it is a conventional modeling choice. However, it would be feasible
761 to consider other noise families and we expect similar results to hold (albeit with modified proof

762 techniques). We now relate the concepts C^e to the atoms. Recall that we defined the index sets
 763 $S^e = \{i \in [n] : a_i \in \mathcal{C} \text{ is a row of } A^e\}$ of atomic concepts in environment e .

764 We define the environment-concept matrix $M \in \mathbb{R}^{m \times n}$ indexed by environments and atoms by

$$M_{ei} = \begin{cases} \frac{1}{\sigma^2} & \text{if } i \in S^e \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

765 Similarly, we consider the environment-valuation matrix $B \in \mathbb{R}^{m \times n}$ given by

$$B_{ei} = \begin{cases} \frac{b_k^e}{\sigma^2} & \text{if } i \in S^e \text{ and row } k \text{ of } A^e \text{ is } a_i, \\ 0 & \text{otherwise.} \end{cases} \quad (8)$$

766 Our first assumption ensures that the concept conditional distributions are sufficiently diverse.

767 **Assumption 4** (Environment diversity I). *The environment-concept matrix $M \in \mathbb{R}^{m \times n}$ has rank n
 768 and there is a vector $v \in \mathbb{R}^m$ such that $v^\top M = 0$ and all entries of $v^\top B$ are non-zero (B denotes
 769 that environment-valuation matrix).*

770 We remark that this assumption can only hold for $m \geq n + 1$ and indeed is satisfied under mild
 771 assumptions on the environments if $m = n + 1$, as the following lemma shows.

772 **Lemma 1.** *Assumption 4 is satisfied almost-surely if there are $n + 1$ concept conditional distributions
 773 such that every n rows of the environment-concept matrix are linearly independent and the b^e are
 774 drawn independently according to a continuous distribution.*

775 We also assume one additional diversity condition. To motivate this, observe if two concepts always
 776 occur together, it’s information-theoretically impossible to distinguish them, e.g., if an agent only
 777 sees red large objects (i.e. all red objects are large and all large objects are red), it will be unable
 778 to disambiguate the “red” concept from the “large” concept. Therefore, we make the following
 779 assumption.

780 **Assumption 5** (Environment diversity II). *For every pair of atoms a_i and a_j with $i \neq j$ there is an
 781 environment e such that $i \in S^e$ and $j \notin S^e$.*

782 We remark that these are the only assumptions about the sets S^e . In particular, we do not need to
 783 know the sets S^e . In the proof, we will extract these sets based on a the signatures they leave on the
 784 datasets. We can now state our main result.

785 **Theorem 2.** *Suppose we are given m context conditional datasets X^1, \dots, X^m and the observational
 786 dataset X^0 such that Assumptions 1-5 hold. Then the concepts are identifiable as in Definition 4.*

787 **Remark 2.** *Assumption 4 can only be satisfied for $m \geq n + 1$, i.e., the result requires at least $n + 2$
 788 environments. On the other hand, Lemma 1 assures that $n + 2$ environments are typically sufficient.
 789 We expect that the result could be slightly improved by showing identifiability for $n + 1$ environments
 790 under suitable assumptions. However, this would probably require more advanced techniques from
 791 algebraic statistics [23] compared to the techniques we employ here.*

792 As mentioned before, our setting somewhat resembles the iVAE setting in Khemakhem et al. [45]
 793 and therefore, their proof techniques can also be applied, with several modifications, to derive
 794 identifiability results in our setting (however our formulation and application are very different).
 795 However, this approach will require more environments because their main assumption is that the
 796 matrix $\Lambda = (M, B) \in \mathbb{R}^{m \times 2n}$ has rank $2n$ so that $2n + 1$ environments are necessary. Moreover,
 797 this rank condition is much stronger than Assumption 4. For completeness and as a warm-up we
 798 prove this result in Appendix B. The full proof of Theorem 2 is fairly involved and is deferred to
 799 Appendix B.

800 B Proofs of the main results

801 In this appendix we provide the proofs of our results, in particular the proof of our main result,
 802 Theorem 2. However, as a warm-up we first start in Appendix B.1 with a proof of the simpler
 803 result that can be shown based on the iVAE approach. In Appendix B.2 we prove Theorem 2 and in
 804 Appendix B.3 we prove the additional lemmas that appear in the paper.

805 **B.1 Proof of identifiability with $2n + 1$ environments**

806 As a warm-up and to provide a connection to earlier results we show here how to obtain identifiability
 807 by adapting the iVAE framework to our context. Indeed, our mathematical setting is related to the
 808 setting used in [45] in the sense that the environments are generated by modulation with certain
 809 exponential families. Therefore, we can essentially apply their proof techniques to prove identifiability
 810 (with some modifications), albeit this requires the suboptimal number of $2m + 1$ environments (there
 811 are two sufficient statistics for the Gaussian distribution).

812 **Theorem 3.** *Suppose data satisfies Assumption 1, 2, and 3 and the environment statistics matrix Λ*
 813 *has rank $2n$. Assume we know the number of atoms n . Then identifiability in the sense of Definition 4*
 814 *holds.*

815 We remark that the rank condition can only be satisfied for $2n + 1$ environments (observational
 816 distribution and $2n$ concept conditional distributions. For this theorem the assumption that the
 817 filtering distribution is always the same is not necessary. Instead we could consider variances $(\sigma_k^e)^2$
 818 depending on environment e and row k , i.e., the filtering distribution $q_{(\sigma_k^e)^2}$ is Gaussian with varying
 819 variance. The generalization of the environment-concept matrix $M \in \mathbb{R}^{m \times n}$ is given by

$$M_{ei} = \begin{cases} \frac{1}{(\sigma_k^e)^2} & \text{if } i \in S^e \text{ and row } k \text{ of } A^e \text{ is } a_i \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

820 Similarly the generalization of the environment-valuation matrix $B \in \mathbb{R}^{m \times n}$ is given by

$$B_{ei} = \begin{cases} \frac{b_k^e}{(\sigma_k^e)^2} & \text{if } i \in S^e \text{ and row } k \text{ of } A^e \text{ is } a_i, \\ 0 & \text{otherwise.} \end{cases} \quad (10)$$

821 We now prove Theorem 3. We use essentially the same ideas as in the proof of Theorem 1 in
 822 Khemakhem et al. [45] (followed by the same reasoning as in Sorrenson et al. [95], Kivva et al. [49]
 823 but since our concepts are not axis aligned and we only extract some information about the mixing
 824 we give a complete proof.

825 *Proof of Theorem 3.* Suppose there are 2 sets of parameters that generate the same data
 826 X^0, X^1, \dots, X^m . Denote by $\tilde{\cdot}$ the latter set of parameters, e.g., X^e is distributed as $\tilde{f}(\tilde{Z}^e)$ where
 827 $\tilde{Z}^e \in \mathbb{R}^{\tilde{d}_z}$ corresponds to the concept class \tilde{C}^e with distribution $\tilde{Z}^e \sim \tilde{p}^e$ and the same distribution is
 828 generated by $f(Z^e)$ where f and \tilde{f} are injective and differentiable. Let $\mathcal{C} = \{a_1, \dots, a_n\}$ be the set
 829 of atomic concepts in the first setting and let $\tilde{\mathcal{C}} = \{\tilde{a}_1, \dots, \tilde{a}_n\}$ be the set of atomic concepts in the
 830 second setting (here we use that n is assumed to be known). We also consider the transition function
 831 $\varphi = \tilde{f}^{-1}f$ and in the following we always write $\tilde{Z} = \varphi(Z)$. The equality $f(Z^e) \stackrel{\mathcal{D}}{=} X^e \stackrel{\mathcal{D}}{=} \tilde{f}(\tilde{Z}^e)$
 832 implies $\varphi(Z^e) \stackrel{\mathcal{D}}{=} \tilde{Z}^e$. This implies that for all environments e

$$p^e(Z) = |\det J_{\varphi^{-1}}| \cdot \tilde{p}^e(\tilde{Z}) \quad (11)$$

833 Taking the logarithm and subtracting this for some $e = 1, \dots, m$ from the base distribution we obtain

$$\ln(p(Z)) - \ln(p^e(Z)) = \ln(\tilde{p}(\tilde{Z})) - \ln(\tilde{p}^e(\tilde{Z})). \quad (12)$$

834 Using the definition (2) we can rewrite for some constants c_e and c'_e

$$\begin{aligned} \ln(p(Z)) - \ln(p^e(Z)) &= \sum_{k=1}^{\dim(C_e)} \frac{(A^e Z^e - b^e)_k^2}{2(\sigma_k^e)^2} - c'_e \\ &= \sum_{i=1}^n \left(\frac{1}{2} M_{ei} \langle a_i, Z^e \rangle^2 - B_{ei} \langle a_i, Z^e \rangle \right) - c_e. \end{aligned} \quad (13)$$

835 Here we used the environment-concept matrix and the environment-valuation matrix in the second
 836 step which were defined in (7) and (8) (in (9) and (10) for varying variance). We define the vector
 837 $\mathbf{p}(Z)$ with components $p_e(Z) = \ln(p(Z)) - \ln(p^e(Z))$. Then we find the relation

$$\mathbf{p}(Z) = \frac{1}{2} M \begin{pmatrix} \langle a_1, Z \rangle^2 \\ \vdots \\ \langle a_n, Z \rangle^2 \end{pmatrix} - B \begin{pmatrix} \langle a_1, Z \rangle \\ \vdots \\ \langle a_n, Z \rangle \end{pmatrix}. \quad (14)$$

838 Together with (12) we conclude that

$$\frac{1}{2}M \begin{pmatrix} \langle a_1, Z \rangle^2 \\ \vdots \\ \langle a_n, Z \rangle^2 \end{pmatrix} - B \begin{pmatrix} \langle a_1, Z \rangle \\ \vdots \\ \langle a_n, Z \rangle \end{pmatrix} = \frac{1}{2}\tilde{M} \begin{pmatrix} \langle \tilde{a}_1, \tilde{Z} \rangle^2 \\ \vdots \\ \langle \tilde{a}_n, \tilde{Z} \rangle^2 \end{pmatrix} - \tilde{B} \begin{pmatrix} \langle \tilde{a}_1, \tilde{Z} \rangle \\ \vdots \\ \langle \tilde{a}_n, \tilde{Z} \rangle \end{pmatrix} \quad (15)$$

839 Since by assumption $\tilde{\Lambda} = (\tilde{M}, \tilde{B}) \in \mathbb{R}^{m \times 2n}$ has rank $2n$ there is a vector v such that $v^\top \tilde{M} = 0$ and
840 $v^\top \tilde{B} = -e_i$ ($e_i \in \mathbb{R}^{d_z}$ denotes the i -th standard basis vector). Thus we find that

$$\langle \tilde{a}_i, \tilde{Z} \rangle = \frac{1}{2}v^\top M \begin{pmatrix} \langle a_1, Z \rangle^2 \\ \vdots \\ \langle a_n, Z \rangle^2 \end{pmatrix} - v^\top B \begin{pmatrix} \langle a_1, Z \rangle \\ \vdots \\ \langle a_n, Z \rangle \end{pmatrix}. \quad (16)$$

841 In other words $\langle \tilde{a}_i, \tilde{Z} \rangle$ can be expressed as a quadratic polynomial in Z . We apply the same reasoning
842 for $\langle \tilde{a}_i, \tilde{Z} \rangle^2$, i.e., pick a vector v' such that $\frac{1}{2}v'^\top \tilde{M} = e_i$ and $v'^\top \tilde{B} = 0$ to obtain a relation

$$\langle \tilde{a}_i, \tilde{Z} \rangle^2 = \sum_j \eta_j \langle a_j, Z \rangle^2 + \ell(Z) \quad (17)$$

843 for some coefficients η_j and some affine function ℓ of Z . The following reasoning is now the same
844 as in Kivva et al. [49], Sorrenson et al. [95]. We thus find that $\langle \tilde{a}_i, \tilde{Z} \rangle$ and its square can be written
845 as polynomials of degree at most 2 in Z . This implies that in fact $\langle \tilde{a}_i, \tilde{Z} \rangle$ is an affine function of Z
846 (otherwise its square would be a quartic polynomial), i.e., we can write

$$\langle \tilde{a}_i, \tilde{Z} \rangle = \sum_j \lambda_j \langle a_j, Z \rangle + C_i = \langle \sum_j \lambda_j a_j, Z \rangle + C_i. \quad (18)$$

847 Equating the square of this relation with (17) and taking the gradient with respect to Z (as a polynomial
848 the function is differentiable) we find

$$2 \sum_j \eta_j a_j \langle a_j, Z \rangle + w = 2 \sum_j \lambda_j a_j \langle \sum_j \lambda_j a_j, Z \rangle + w' \quad (19)$$

849 for two vectors w and w' . The equality (for $Z = 0$) implies $w = w'$. Now linear independence of a_j
850 implies that for each r

$$\eta_r a_r = \lambda_r \sum_j \lambda_j a_j. \quad (20)$$

851 Applying linear independence again we conclude that either $\lambda_r = 0$ or $\lambda_j = 0$ for all $j \neq r$. This
852 implies that there is at most one r such that $\lambda_r \neq 0$. The relation (18) and the bijectivity of φ implies
853 that there is exactly one $r(i)$ such that $\lambda_{r(i)} \neq 0$ and therefore

$$\langle \tilde{a}_i, \tilde{Z} \rangle = \lambda_{r(i)} \langle a_{r(i)}, Z \rangle + C_i. \quad (21)$$

854 Applying the same argument in the reverse direction we conclude that there is a permutation $\pi \in S_n$
855 such that

$$\langle \tilde{a}_{\pi(i)}, \tilde{Z} \rangle = \lambda_i \langle a_i, Z \rangle + C_i. \quad (22)$$

856 By linear independence we can find an invertible linear map T such that

$$\tilde{a}_{\pi(i)}^\top = a_i^\top T^{-1} \quad (23)$$

857 (i.e., $T^\top \tilde{a}_{\pi(i)} = a_i$) and a vector $w \in \mathbb{R}^{d_z}$ (the a_i are linearly independent) such that

$$\langle \tilde{a}_{\pi(i)}, \tilde{Z} \rangle = \lambda_i (\langle a_i, Z \rangle + \langle a_i, w \rangle). \quad (24)$$

858 In particular the relations (5) and (6) hold. Now it is straightforward to see that if $i \in S^e$, i.e., a_i is a
859 row of A^e then $\tilde{a}_{\pi(i)}$ is a row of \tilde{A}^e and vice versa. Indeed, this follows from (15) for environment e
860 together with (24) and linear independence of the atoms. Therefore we conclude from (23) that there
861 is a permutation P^e such that

$$\tilde{A}^e = P^e A^e T^{-1}. \quad (25)$$

862 Moreover, (24) then implies setting $Z = f^{-1}(x)$, $\tilde{Z} = \tilde{f}^{-1}(x)$

$$\tilde{A}^e \tilde{f}^{-1}(x) = \Lambda^e P^e A^e (f^{-1}(x) + w) \quad (26)$$

863 holds for the same permutation matrix P^e and a diagonal matrix Λ^e whose diagonal entries can be
 864 related to (24). Let us assume now that row k of A^e is a_i and row k' of \tilde{A}^e is $\tilde{a}_{\pi(i)}$. Now we consider
 865 the subspace $H \subset \mathbb{R}^{d_z}$ containing all Z such that $\langle Z, a_j \rangle = 0$ for $j \neq i$. Via (24) this implies that
 866 $\langle \tilde{a}_j, \tilde{Z} \rangle$ is constant for $j \neq \pi(i)$. Then we conclude from (15) that for $Z \in H$

$$\frac{(\langle a_i, Z \rangle - b_k^e)^2}{2(\sigma_k^e)^2} = \frac{(\langle \tilde{a}_{\pi(i)}, \tilde{Z} \rangle - \tilde{b}_{k'}^e)^2}{2(\tilde{\sigma}_{k'}^e)^2} + c_k^e \quad (27)$$

867 for some constant c_k^e . Using (24) this implies that

$$\frac{(\langle a_i, Z \rangle - b_k^e)^2}{2(\sigma_k^e)^2} = \frac{(\lambda_i(\langle a_i, Z \rangle + \langle a_i, w \rangle) - \tilde{b}_{k'}^e)^2}{2(\tilde{\sigma}_{k'}^e)^2} + c_k^e. \quad (28)$$

868 Comparing the quadratic term and the linear term (note that $\langle a_i, Z \rangle$ can take any value on H) we find

$$\frac{1}{2(\sigma_k^e)^2} = \frac{\lambda_i^2}{2(\tilde{\sigma}_{k'}^e)^2} \quad (29)$$

$$-\frac{b_k^e}{2(\sigma_k^e)^2} = -\frac{\lambda_i \tilde{b}_{k'}^e - \lambda_i^2 \langle a_i, w \rangle}{2(\tilde{\sigma}_{k'}^e)^2} \quad (30)$$

869 Combining the equation we obtain

$$\tilde{b}_{k'}^e = \lambda_i(b_k^e - \langle a_i, w \rangle) \quad (31)$$

870 This implies then the relation

$$\tilde{b} = \Lambda^e P^e (b + A^e w). \quad (32)$$

871

□

872 B.2 Proof of Theorem 2

873 In this section we prove our main Theorem 2. The proof is structured in several steps: First we remove
 874 the symmetries of the representation and derive the key relations underlying the proof. Then we show
 875 that we can identify the environment-concept matrix M and then also the valuations collected in B .
 876 Once this is done we can complete the proof. We will need the following lemma to conclude the
 877 proof.

878 **Lemma 2.** *The relations (3) and (6) in Definition 4 define an equivalence relation of representations*
 879 *if we assume that the underlying atoms form a linearly independent set.*

880 The proof of this lemma can be found in Appendix B.3.

881 **Remark 3.** *Without the assumption on the underlying atoms the lemma is not true. In this case*
 882 *a slightly different scaling must be chosen (e.g., $(\Lambda^e)^{-1} \tilde{b}^e = \Lambda^e P^e b^e - P^e A^e w$ instead of $\tilde{b}^e =$*
 883 *$\Lambda^e P^e (b^e - A^e w)$). Since our results address the case of atoms we used the simpler definition in the*
 884 *main paper.*

885 We can allow slightly more general filtering distributions where q is Gaussian with variance σ_i^2 if we
 886 filter on concept i , i.e., the variance needs to be constant for different environments and the same
 887 atom but might depend on the atom. The proof will cover this case, the simple case stated in the main
 888 paper is obtained by setting $\sigma_i^2 = \sigma^2$. Some steps of the proof (e.g., the expressions for the difference
 889 of the log-densities) agree with the proof of Theorem 3. To keep the proof self contained we repeat a
 890 few equations.

891 *Proof of Theorem 2.* We proceed in several steps.

892 **Step 1: Reduction to standard form.** Let us first transform every possible data representation into
 893 a standard form. Recall that we have the set of atomic concepts $\mathcal{C} = \{a_1, \dots, a_n\}$. Recall that we
 894 defined the environment-concept matrix $M \in \mathbb{R}^{m \times n}$ in (7) and note that the natural generalisation
 895 reads

$$M_{ei} = \begin{cases} \frac{1}{\sigma_i^2} & \text{if } a_i \text{ is a row of } A^e, \\ 0 & \text{otherwise.} \end{cases} \quad (33)$$

896 We say that concept a_n is conditioned on the environment e . Note that the nonzero entries of row e
 897 of M encode the set S^e . To pass from A^e to its rows a_i we assume that the e -th row of A^e is $a_{i_j^e}$, i.e.,
 898 $a_{i_j^e} = (A^e)^\top e_j$. Recall also consider the environment-valuation matrix B which is given by

$$B_{ei} = \begin{cases} \frac{b_k^e}{\sigma_i^2} & \text{if } a_i \text{ is row } k \text{ of } A^e, \\ 0 & \text{otherwise.} \end{cases} \quad (34)$$

899 Denoting by q_{σ^2} the centered Gaussian distribution with variance σ^2 we find in environment e

$$\begin{aligned} \ln(p(Z)) - \ln(p^e(Z)) &= - \sum_{k=1}^{\dim(C_e)} \ln q_{(\sigma_k^e)^2}((A^e Z^e - b^e)_k) = \sum_{k=1}^{\dim(C_e)} \frac{(A^e Z^e - b^e)_k^2}{2(\sigma_k^e)^2} - c'_e \\ &= \sum_{i=1}^n \frac{1}{2} M_{ei} \langle a_i, Z^e \rangle^2 - B_{ei} \langle a_i, Z^e \rangle - c_e. \end{aligned} \quad (35)$$

900 Now we consider an invertible linear map $T: \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_z}$ such that $T^{-\top} a_i = e_i$ for all $1 \leq i \leq n$.
 901 Such a map exists because we assume that the a_i are linearly independent. Moreover, we consider
 902 a shift vector $\lambda \in \mathbb{R}^{d_z}$ with $\lambda_i = 0$ for $i > n$ which we fix later. We define $\Sigma \in \mathbb{R}^{d_z \times d_z}$ to be the
 903 diagonal matrix with entries $\Sigma_{ii} = \sigma_i$ for $1 \leq i \leq n$ and $\Sigma_{ii} = 1$ for $i > n$. Now we consider the
 904 linear map $L(z) = \Sigma^{-1} T z - \lambda$ and a new representation given by

$$\bar{z} = L(z), \quad \bar{f} = f \circ L^{-1}, \quad \bar{\mathcal{C}} = \{e_1, \dots, e_n\}, \quad \bar{\sigma}_i = 1, \quad \bar{A}^e = A^e T^{-1}, \quad \bar{p}(\bar{z}) = p(L^{-1} \bar{z}) |\det T^{-1}|. \quad (36)$$

905 We also define

$$\bar{b}_k^e = \frac{b_k^e}{\sigma_i} - \lambda_i \quad \text{if row } k \text{ of } A^e \text{ is } a_i. \quad (37)$$

906 Define \bar{M} and \bar{B} in terms of \bar{A}^e , \bar{b}^e and $\bar{\sigma}_i^2$ as before. We remark that all entries of \bar{M} are either 0 or
 907 1 and note that

$$\bar{M} = M \text{Diag}(\sigma_1^2, \dots, \sigma_n^2) \quad (38)$$

$$\bar{B} = B \text{Diag}(\sigma_1^{-1}, \dots, \sigma_n^{-1}) - M \text{Diag}(\lambda_1, \dots, \lambda_n). \quad (39)$$

908 We claim that this model generates the same observations as the original model. By definition
 909 $L_* p = \bar{p}$ (as mentioned before, we slightly abuse notation and here refer to the distributions). Next,
 910 we calculate for any δ

$$\begin{aligned} -2 \ln q_1(\langle e_i, L(z) \rangle - \delta) &= (\langle e_i, L(z) \rangle - \delta)^2 \\ &= (\langle e_i, \Sigma T z - \lambda \rangle - \delta)^2 \\ &= (\sigma_i^{-1} \langle T^\top e_i, z \rangle - \lambda_i - \delta)^2 \\ &= \frac{(\langle a_i, z \rangle - \sigma_i \lambda_i - \sigma_i \delta)^2}{\sigma_i^2} \\ &= -2 \ln q_{\sigma_i^2}(\langle a_i, z \rangle - \sigma_i \lambda_i - \sigma_i \delta). \end{aligned} \quad (40)$$

911 Using this for $\delta = \bar{b}_k^e$ and some k such that row k of A^e is a_i we find

$$-2 \ln q_1(\langle e_i, L(z) \rangle - \bar{b}_k^e) = -2 \ln q_{\sigma_i^2}(\langle a_i, z \rangle - \sigma_i \lambda_i - \sigma_i \bar{b}_k^e) = -2 \ln q_{\sigma_i^2}(\langle a_i, z \rangle - b_k^e). \quad (41)$$

912 This then implies that for $\tilde{z} = L(z)$

$$\prod_k q_1((\tilde{A}^e \tilde{z} - \tilde{b}^e)_k) \propto \prod_k q_{\sigma_k^e}((A^e z - b^e)_k). \quad (42)$$

913 Combining this with the definition (2) and the definition $\bar{p}(\tilde{z}) = p(L^{-1}\tilde{z})|\det T^{-1}|$ we find that for
914 $\bar{z} = L(z)$

$$\bar{p}^e(\tilde{z}) \propto p^e(z) \quad (43)$$

915 and thus $\bar{f}(\bar{Z}^e) \stackrel{\mathcal{D}}{=} f(Z^e) \stackrel{\mathcal{D}}{=} X^e$. Moreover, one directly sees that the two representations are also
916 equivalent in the sense of Definition 4. We now fix the vector λ such that each row of \bar{B} has mean zero.
917 Finally, by changing the sign of \tilde{z}_i we can in addition assume that for every i the first non-zero \bar{B}_{ei}
918 is positive. Finally we remark that Assumption 4 is still satisfied for \bar{M} and \bar{B} . Indeed, $w^\top M = 0$
919 implies $w^\top \bar{M} = 0$ by (38). But then $w^\top \bar{B} = w^\top B \text{Diag}(\sigma_1^{-1}, \dots, \sigma_n^{-1})$ by (39) which has all
920 entries different from zero if this holds for $w^\top B$. In the following we will therefore always assume
921 that the representation satisfies the properties of the \bar{Z} variables and we remove the modifier in the
922 following. The plan is now to show that M and B can be identified up to permutations of the rows
923 (under the fixed normalization we derived in this step) and then show that every two representations
924 with the same M and B can be identified.

925 **Step 2: The key identity** Let us here restate the key identity based on the difference of the log-
926 densities. As is common in identifiability results for multi-environment data with general mixing we
927 consider the difference in log densities. Consider

$$\begin{aligned} \ln p^0(z) - \ln p^e(z) &= \sum_{i=1}^n \frac{1}{2} M_{ei} \langle e_i, z \rangle^2 - B_{ei} \langle e_i, z \rangle - c'_e \\ &= \sum_{i=1}^n \frac{1}{2} M_{ei} z_i^2 - B_{ei} z_i - c'_e \end{aligned} \quad (44)$$

928 for some constant c'_e . Those functions will play a crucial role in the following and we will denote

$$g^e(z) = \ln p^0(z) - \ln p^e(z) \quad (45)$$

929 Note that since the log-density changes only by the Jacobian for pushforward measures we find that

$$g^e(z) = \ln p^0(z) - \ln p^e(z) = \ln p_X^0(f(z)) - \ln p_X^e(f(z)) = G^e(f(z)) = G^e(x). \quad (46)$$

930 Note that the functions $G^e(x)$ can be estimated from the distributions of X^e . We remark X might be
931 supported on a submanifold if d_z and d_x do not agree making the definition of the density subtle. But
932 we can just consider any chart locally and consider the density of the pushforward with respect to the
933 Lebesgue measure. The resulting difference expressed in G^e will be independent of the chart as the
934 determinant cancels thus G^e is a well defined function. The relation

$$g^e(z) = G^e(f(z)) = G^e(x) \quad (47)$$

935 will be crucial in the following because it shows that properties of g^e are closely linked to the
936 identifiable functions G^e .

937 **Step 3: Identifiability of environment-concept matrix** Let us now show that we can identify
938 which concepts are contained in which environment (up to relabeling of the concepts). Recall that
939 $S^e = \{i \in [n] : a_i \text{ is a row of } A^e\}$ and we similarly define $S_T = \bigcup_{e \in T} S^e$ for all subsets $T \subset [m]$.
940 The main observation is that we can identify $|S_T| = |\bigcup_{e \in T} S^e|$ for all subsets $T \subset [m]$. To show
941 this we consider the set

$$I_T = \underset{z}{\operatorname{argmin}} \sum_{e \in T} g^e(z). \quad (48)$$

942 Note that the function g^e are convex functions, and they can be decomposed as sums of functions in
943 z_i , i.e., for some functions h_i^T

$$\sum_{e \in T} g^e(z) = \sum_{i=1}^n h_i^T(z_i). \quad (49)$$

944 Now if $i \in S_T$ then $i \in S^e$ for some e and thus $M_{ei} \neq 0$ for the e and h_i^T is the sum of quadratic
 945 function in x_i which as a strictly convex function has a unique minimum z_i^T . On the other hand, if
 946 $i \notin S_T$ then $i \notin S^e$ for $e \in T$ and thus $M_{ei} = 0$ for all $e \in T$ and $h_i^T(z_i) = 0$. Thus we conclude
 947 that

$$I_T = \{z \in \mathbb{R}^{d_z} : z_i = z_i^T \text{ for } i \in S_T\}. \quad (50)$$

948 This is an affine subspace of dimension $d_z - |S_T|$. The relations $G^e(f(z)) = g^e(z)$ imply that

$$f(I_T) = \operatorname{argmin}_x \sum_{e \in T} G^e(x). \quad (51)$$

949 Note that $G^e(x)$ is identifiable from the datasets X^e and thus the submanifold (by assumption on f)
 950 $f(I_T)$ is identifiable and by finding its dimension we obtain $d_z - |S_T|$. Since d_z is the dimension of
 951 the data manifold $f(X)$ we can indeed identify $|S_T|$ for all $T \subset [m]$. In particular, the total number
 952 of atomic concepts $n = |S_{[m]}|$ is identifiable (assuming that all atomic concepts are filtered upon at
 953 least once). Now, it is a standard result that we can identify the matrix M up to permutation of the
 954 atomic concepts. Indeed, we can argue by induction in m to show this. For $m = 1$ we just have $|S^1|$
 955 atomic concepts appearing in environment 1 and $n - |S^1|$ concepts not appearing. For the induction
 956 step $m \rightarrow m + 1$ we consider the sizes $|S_{T \cup \{m+1\}}|$ for $T \subset [m]$. Applying the induction hypothesis
 957 we can complete M_{ei} for all columns such that $M_{m+1,i} = 1$. Similarly, we can consider the sizes
 958 $|S_T| - |S_{T \cup \{m+1\}}|$ to identify the matrix M for concepts not used in environment $m + 1$.

959 Thus, we can and will assume after permuting the atomic concepts that M is some fixed matrix.

960 **Step 4: Identifiability of concept valuations** Next, we show that we can also identify the matrix
 961 B . We do this column by column, i.e., for one atomic concept after another. Assume we consider
 962 atomic concept i . Then we consider the set $T_i = \{e : M_{ei} = 0\}$ of concepts that not filter on atomic
 963 concept i . By Assumption 5 there is for every $i' \neq i$ an environment e such that i' is filtered on, i.e.,
 964 $M_{ei'} \neq 0$. This implies $S_{T_i} = [n] \setminus \{i\}$. Then we consider as in (50) the set I_{T_i} given by

$$I_{T_i} = \{z \in \mathbb{R}^{d_z} : z_{i'} = z_{i'}^T \text{ for } i' \in [n] \setminus \{i\}\}. \quad (52)$$

965 Note that all $z_{i'}$ for $i \neq i'$ are constant on I_{T_i} . Thus we find for any environment e such that $i \in S^e$.

$$\begin{aligned} g^e(z) &= \sum_{j=1}^n \frac{1}{2} M_{ej} z_j^2 - B_{ej} z_j - c'_e \\ &= \sum_{j \neq i} \frac{1}{2} M_{ej} z_j^2 - B_{ej} z_j - c'_e + \frac{1}{2} z_i^2 - B_{ei} z_i \\ &= c_{T_i, e} + \frac{1}{2} z_i^2 - B_{ei} z_i \end{aligned} \quad (53)$$

966 on I_{T_i} for some constant c_{T_i} .

967 Now we consider two concepts $e_1 \neq e_2$ such that atomic concept i is contained in these two
 968 environments. Then we consider the set

$$I_{T_i}^{e_1} = \operatorname{argmin}_{z \in I_{T_i}} g^{e_1}(z) = \{z \in \mathbb{R}^{d_z} : z_{i'} = z_{i'}^T \text{ for } i' \in [n] \setminus \{i\}, z_i = B_{e_1 i}\}. \quad (54)$$

969 Note that in the second equality we used that $g^{e_1}(z)$ depends on z_i through $z_i^2/2 - B_{e_1 i} z_i$ so it is
 970 minimized at $B_{e_1 i}$. Now we find using (53)

$$\begin{aligned} \min_{z \in I_{T_i}^{e_1}} g^{e_2}(z) - \min_{I_{T_i}} g^{e_2}(z) &= \min_{z \in I_{T_i}^{e_1}} c_{T_i, e_2} + \frac{1}{2} z_i^2 - B_{e_2 i} z_i - \min_{I_{T_i}} \left(c_{T_i, e_2} + \frac{1}{2} z_i^2 - B_{e_2 i} z_i \right) \\ &= c_{T_i, e_2} + \frac{1}{2} B_{e_1 i}^2 - B_{e_1 i} B_{e_2 i} - \left(c_{T_i, e_2} + \frac{1}{2} B_{e_2 i}^2 - B_{e_2 i}^2 \right) \\ &= \frac{(B_{e_1 i} - B_{e_2 i})^2}{2}. \end{aligned} \quad (55)$$

971 As before, this quantity is identifiable from observations because $f(T_i)$ can be identified and we can
 972 minimize $G^{e_2}(x)$ over $f(T_i)$.

973 This allows us to identify $B_{e_1i} - B_{e_2i}$ up to a sign. However, we can evaluate this expression over
 974 all pairs e_1 and e_2 and pick the one with the maximal difference. Then all remaining values B_{e_i}
 975 for e such that i is filtered on in e must satisfy $B_{e_i} \in [B_{e_1i}, B_{e_2i}]$. Together with identifiability of
 976 $|B_{e_i} - B_{e_1i}|$ this allows us to identify all B_{e_i} up to one sign indeterminacy and a constant shift.
 977 However, in the first step we ensured that $\sum_e B_{e_i} = 0$ for all i which determines the shift and the
 978 sign is fixed by our choice of making the first non-zero entry positive. Thus, we can assume that our
 979 two representations have the same M and B .

980 **Step 5: Identifiability of concepts** We are now ready to prove our identifiability result.

981 Assume we have two representations Z^e, f, p and \tilde{Z}^e, \tilde{f} , and \tilde{p} such that the corresponding
 982 environment-concept and environment-valuation matrices agree, i.e., $M = \tilde{M}$ and $B = \tilde{B}$. We
 983 consider the transition function $\varphi = \tilde{f}^{-1} \circ f$ which is by assumption differentiable. What we want to
 984 show is that $\varphi(z)_i = z_i$ for all $z \in \mathbb{R}^{d_z}$ and $1 \leq i \leq n$. We now decompose $z = (z^c, z^o)$ into the
 985 concept part and the orthogonal part. We fix $z^o \in \mathbb{R}^{d_z - n}$ and define the function $\iota^o(z^c) = (z^c, z^o)$,
 986 the projection $\pi^c((z^c, z^o)) = z^c$, and $\varphi^o : \mathbb{R}^n \rightarrow \mathbb{R}^n$ given by $\varphi^o(z^c)_i = \varphi(\iota^o(z^c))_i$.
 987 Note that φ^o is differentiable but not necessarily injective. Let us denote by $\mathbf{g} : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^m$ the
 988 function with coordinates $\mathbf{g}_e = g^e$ and similarly we define $\mathbf{G} : M \rightarrow \mathbb{R}^d$. Identifiability will be
 989 based on the crucial relation

$$\mathbf{g}(\iota^o(z^c)) = \mathbf{G}(f(\iota^o(z^c))) = \mathbf{G}(\tilde{f}(\varphi^o(z^c))) = \mathbf{g}(\varphi^o(z^c)). \quad (56)$$

990 Here we used in the last step that g^e is defined in terms of M and B and thus agrees for both
 991 representations. Note that \mathbf{g} is just a quadratic function. Differentiating we obtain

$$D_i g^e(z) = M_{ei} z_i - B_{ei}. \quad (57)$$

992 Concisely this can be written as

$$D\mathbf{g} = M \text{Diag}(z_1, \dots, z_n) - B. \quad (58)$$

993 Differentiating (56) we find

$$M \text{Diag}(z_1, \dots, z_n) - B = (M \text{Diag}(\tilde{z}_1, \dots, \tilde{z}_n) - B) D\varphi^o(z^c). \quad (59)$$

994 Let v be a vector as in Assumption 4. Denote by $M^+ \in \mathbb{R}^{n \times m}$ the pseudoinverse of M which has
 995 rank n because M has. We consider the matrix $\widetilde{M}^+ \in \mathbb{R}^{n+1 \times m}$ given by

$$\widetilde{M}^+ = \begin{pmatrix} M^+ \\ v^\top \end{pmatrix} \quad (60)$$

996 Let us multiply the relation (59) by \widetilde{M}^+ and find that

$$\begin{pmatrix} z_1 & & 0 \\ & \ddots & \\ 0 & & z_n \\ 0 & \dots & 0 \end{pmatrix} - \widetilde{M}^+ B = \left(\begin{pmatrix} \tilde{z}_1 & & 0 \\ & \ddots & \\ 0 & & \tilde{z}_n \\ 0 & \dots & 0 \end{pmatrix} - \widetilde{M}^+ B \right) D\varphi^o(z^c) \quad (61)$$

997 Note that the first n rows of the left hand side are $\text{Diag}(z_1, \dots, z_n) - M^+ B$. This matrix is invertible
 998 for almost all values of $z^c = (z_1, \dots, z_n)^\top$ because its determinant is a non-zero polynomial (the
 999 coefficient of the term $z_1 \dots z_n$ is 1) which vanishes only on a set of measure zero. Outside of this
 1000 set the left hand side of has rank n . Then the equality (61) implies that also the right hand side has
 1001 rank n and thus $D\varphi^o(z^c)$ has rank n and thus is invertible. For z^c outside of this set there is up to
 1002 scaling a unique vector $w \neq 0$ (depending on z_1, \dots, z_n such that

$$w^\top \left(\begin{pmatrix} z_1 & & 0 \\ & \ddots & \\ 0 & & z_n \\ 0 & \dots & 0 \end{pmatrix} - \widetilde{M}^+ B \right) = 0 \quad (62)$$

1003 From (61) we conclude using the invertibility of $D\varphi^o(z^c)$ that

$$w^\top \left(\begin{pmatrix} \tilde{z}_1 & & 0 \\ & \ddots & \\ 0 & & \tilde{z}_n \\ 0 & \dots & 0 \end{pmatrix} - \widetilde{M}^+ B \right) = 0. \quad (63)$$

1004 Next, we claim that for almost all values of z^c the vector w has all entries different from 0 (this
 1005 property is invariant under rescaling). Actually we need this only for entries 1 to n but the case $n + 1$
 1006 is a bit simpler so we show it first. We show this by proving that for each entry w_i there is only a null
 1007 set of z^c such that $w_i = 0$. Let $w = (w', 0)$ for some $w' \in \mathbb{R}^n$ and $w' \neq 0$, i.e., $w_{n+1} = 0$. Then

$$0 = w^\top \left(\begin{pmatrix} z_1 & & 0 \\ & \ddots & \\ 0 & & z_n \\ 0 & \dots & 0 \end{pmatrix} - \widetilde{M^+B} \right) = w'^\top (\text{Diag}(z_1, \dots, z_n) - M^+B) \quad (64)$$

1008 But this implies that $\text{Diag}(z_1, \dots, z_n) - M^+B$ has non-trivial kernel, i.e., does not have full rank
 1009 and we have seen above that this happens only for a subset of measure 0 of all z^c . Next we show that
 1010 the same is true if $w_1 = 0$. Decompose $0 \neq w = (0, w')$. Then we find

$$0 = w^\top \left(\begin{pmatrix} z_1 & & 0 \\ & \ddots & \\ 0 & & z_n \\ 0 & \dots & 0 \end{pmatrix} - \widetilde{M^+B} \right) = w'^\top \left(\begin{pmatrix} 0 & z_2 & 0 & 0 \\ \dots & & \ddots & \\ 0 & & & z_n \\ 0 & \dots & \dots & 0 \end{pmatrix} - (\widetilde{M^+B})_{2:(n+1)} \right) \quad (65)$$

1011 Thus we conclude that the matrix on the right hand side is not invertible. Its determinant is a
 1012 polynomial in z_2, \dots, z_n and its highest degree term is $\pm z_2 \cdot \dots \cdot z_n \cdot (\widetilde{M^+B})_{(n+1),1}$. By definition
 1013 of $\widetilde{M^+B}$ we find $(\widetilde{M^+B})_{(n+1),1} = (v^\top B)_1 \neq 0$ by Assumption 4 (recall that we showed invariance
 1014 of the assumption under the transformation of M and B). We find that the determinant is a non-zero
 1015 polynomial and the set of its zeros is a set of measure 0 of all z_2, \dots, z_n but since it does not depend
 1016 on z_1 this holds true for almost all z^c . The same reasoning for $i = 2, \dots, n$ implies that for every
 1017 i the set of z^c such that $w_i = 0$ is a set of measure zero. We have therefore shown that for almost
 1018 all z^c the rank of the left hand side of (61) is n and the corresponding vector $w \neq 0$ has all entries
 1019 different from zero. Subtracting (62) and (63) we obtain

$$0 = w^\top \begin{pmatrix} z_1 & & 0 \\ & \ddots & \\ 0 & & z_n \\ 0 & \dots & 0 \end{pmatrix} - w^\top \begin{pmatrix} \tilde{z}_1 & & 0 \\ & \ddots & \\ 0 & & \tilde{z}_n \\ 0 & \dots & 0 \end{pmatrix} = (w_1(z_1 - \tilde{z}_1), \dots, w_n(z_n - \tilde{z}_n), 0). \quad (66)$$

1020 Now $w_i \neq 0$ implies $z_i = \tilde{z}_i$. We conclude that for almost all z^c the relation $\varphi^o(z^c) = z^c$
 1021 holds. By continuity this implies that the relation actually holds everywhere. We conclude that
 1022 $\pi^c \tilde{f}^{-1} f((z^c, z^o)) = z^c$ for a fixed z^o but since z^o was arbitrary the relation holds for all z^o and all
 1023 z^c . Thus we conclude that for $1 \leq i \leq n$

$$\langle e_i, \tilde{f}^{-1}(x) \rangle = \langle e_i, \varphi(f^{-1}(x)) \rangle = \langle e_i, f^{-1}(x) \rangle \quad (67)$$

1024 holds. This implies that those two representations satisfy (3) and (4) (with $P^e = \Lambda^e = \text{Id}$ and
 1025 $T = \text{Id}$). But since this relation is an equivalence relation in our setting by Lemma 2 and since we
 1026 showed equivalence to a representation in standard form in the first step we conclude that also any
 1027 two representations are related through (3) and (4) thus finishing the proof. \square

1028 B.3 Remaining proofs

1029 Here we prove the remaining auxiliary results.

1030 *Proof of Lemma 1.* Since $M \in \mathbb{R}^{m \times n}$ has rank n and $m = n + 1$ there is exactly one vector $v \in \mathbb{R}^m$
 1031 such that $v^\top M = 0$ and $v \neq 0$. We claim that this vector has all entries different from zero.
 1032 Indeed suppose $v_m = 0$ which then implies $v_{1:(m-1)}^\top M_{1:(m-1)} = 0$. But by assumption every $n \times n$
 1033 submatrix of M is invertible (this is equivalent to the rows being linearly independent) so we conclude
 1034 that $v_{1:(m-1)} = 0$ which is a contradiction to $v \neq 0$. The same reasoning applies to every entry.
 1035 Note that the assumption on M implies that every column has at least one non-zero entry, i.e., every
 1036 column of B has one entry sampled from a continuous distribution. But then the probability that v is
 1037 orthogonal to a column is zero because this is a codimension 1 hyperplane of all valuations of this
 1038 row (since all entries of v are non-zero). \square

1039 *Proof of Lemma 2.* Reflexivity is obvious, just pick $T = \text{Id}$, $w = 0$, $\Lambda^e = P^e = \text{Id}_{\dim(C^e)}$. To show
 1040 symmetry we first consider the atoms. Let $\tilde{T} = T^{-1}$ and $\tilde{\pi} = \pi^{-1}$. Then

$$a_{\tilde{\pi}(i)}^\top = a_{\pi^{-1}(i)}^\top T^{-1} T = \tilde{a}_{\pi \circ \pi^{-1}(i)} \tilde{T}^{-1} = \tilde{a}_i \tilde{T}^{-1}. \quad (68)$$

1041 Let \tilde{w} be a vector such that for all $1 \leq i \leq n$

$$\langle a_i, w \rangle = -\frac{1}{\lambda_i} \langle \tilde{a}_{\pi(i)}, \tilde{w} \rangle. \quad (69)$$

1042 Such a vector exists by linear independence of \tilde{a}_i . Let $\tilde{\lambda}_i = \lambda_{\tilde{\pi}(i)}^{-1}$. Then we find that the relation (6),
 1043 namely

$$\langle \tilde{a}_{\pi(i)}, \tilde{f}^{-1}(x) \rangle = \lambda_i (\langle a_i, f^{-1}(x) \rangle + \langle a_i, w \rangle) \quad (70)$$

1044 implies

$$\begin{aligned} \langle a_{\tilde{\pi}(i)}, f^{-1}(x) \rangle &= \frac{1}{\lambda_{\tilde{\pi}(i)}} \langle \tilde{a}_{\pi \circ \tilde{\pi}(i)}, \tilde{f}^{-1}(x) \rangle - \langle a_{\tilde{\pi}(i)}, w \rangle = \frac{1}{\lambda_{\tilde{\pi}(i)}} \langle \tilde{a}_i, \tilde{f}^{-1}(x) \rangle + \frac{1}{\lambda_{\tilde{\pi}(i)}} \langle \tilde{a}_{\pi \circ \tilde{\pi}(i)}, \tilde{w} \rangle \\ &= \tilde{\lambda}_i (\langle \tilde{a}_i, \tilde{f}^{-1}(x) \rangle + \langle \tilde{a}_i, \tilde{w} \rangle). \end{aligned} \quad (71)$$

1045 It remains to be shown that this lifts to the concepts C^e . We first note that the relation (6) together
 1046 with (69) and (3) implies that

$$\Lambda^e P^e A^e w = -\tilde{A}^e \tilde{w}. \quad (72)$$

1047 Let $\tilde{P}^e = (P^e)^{-1}$ and $\tilde{\Lambda}^e = (P^e)^{-1} (\Lambda^e)^{-1} P^e$. Then (3) combined with the previous disply implies

$$\begin{aligned} A^e f^{-1}(x) &= (P^e)^{-1} (\Lambda^e)^{-1} \tilde{A}^e \tilde{f}^{-1}(x) - A^e w \\ &= \tilde{\Lambda}^e \tilde{P}^e \tilde{A}^e \tilde{f}^{-1}(x) + (P^e)^{-1} (\Lambda^e)^{-1} \tilde{A}^e \tilde{w} \\ &= \tilde{\Lambda}^e \tilde{P}^e \tilde{A}^e (\tilde{f}^{-1}(x) + \tilde{w}). \end{aligned} \quad (73)$$

1048 The relation

$$A^e = \tilde{P}^e \tilde{A}^e \tilde{T}^{-1} \quad (74)$$

1049 is a direct consequence of the definitions of \tilde{P}^e and \tilde{T} and (4) and the relation

$$b^e = \tilde{\Lambda}^e \tilde{P}^e (\tilde{b}^e - \tilde{A}^e w) \quad (75)$$

1050 follows exactly as in (73). The proof of transitivity is similar (first establish the relations on the
 1051 atomic concepts then lift it to C^e). \square

1052 C Comparison to Causal Representation Learning

1053 In this appendix we describe causal representation learning and discuss the similarities and differences
 1054 between the viewpoint taken in this paper and the standard setting in causal representation learning.

1055 Causal Representation Learning (CRL) [90, 89] aims to learn representations of data that correspond to
 1056 true causal generative processes. More precisely, if we assume that data X is generated as $X = f(Z)$
 1057 where Z are latent causal factors and f is some arbitrary nonlinearity, the goal is to learn f as well as
 1058 the distribution of Z . Since the latent variables Z are assumed to have causal relationships among
 1059 them, many works exploit the presence of interventional data to learn the generative model. CRL
 1060 incorporates ideas from the field of causality [96, 75, 77, 84, 97] into the field of latent variable models
 1061 and is a generalization of nonlinear independent component analysis [18, 37, 39] and disentangled
 1062 representation learning [9, 77, 52]. The field has seen a surge of advances in the last few years, e.g.,
 1063 [45, 48, 28, 60, 51, 11, 68, 128, 31, 85, 110, 42, 41, 102, 111, 123, 120]. As motivated in Schölkopf
 1064 et al. [90], CRL enables many desiderata such as robustness, out of distribution generalization, and in
 1065 addition enables planning and alignment. CRL has also been successful in many domains such as
 1066 computer vision [45, 113, 2], robotics [63, 10, 59, 126] and genomics [98, 125].

1067 In our work, we take significant inspiration from this framework of causal representation learning and
 1068 present a relaxed framework that is weaker, but more general and also importantly, aligns better with
 1069 empirical works on interpretability of large pre-trained models in the literature. We now describe
 1070 the setup of CRL more formally in Appendix C.1. Then, in Appendix C.2, we discuss conceptual
 1071 differences between causal representation learning and our framework.

1072 C.1 Formal setup

1073 We assume that we observe data $X \in \mathbb{R}^{d_x}$ with the generative model $X = f(Z)$ where $Z \in \mathbb{R}^{d_z}$
1074 are the latent variables and f is a deterministic mixing function. The dataset X is sampled from
1075 a distribution p and the goal is to recover the mixing function f as well as the distributions of
1076 the underlying latent variables Z_1, \dots, Z_{d_z} . To this end, this problem is over-parameterized since
1077 multiple pairs of Z and f could fit the dataset apriori, so the common practice in CRL is to impose
1078 various assumptions that will make this model *identifiable*. Here, identifiability is the notion that
1079 a unique set of parameters fit the model (up to trivial transformations). This makes the problem
1080 well-defined and feasible, although it could still be a hard problem to solve in practice. Below, we
1081 informally summarize two classes of prior works that enable such identifiability guarantees.

- 1082 1. Disentangled representation learning: In this setting, we assume that the distributions of
1083 Z_1, \dots, Z_{d_z} are jointly independent. Different studies constrain the distribution of the
1084 variables Z_1, \dots, Z_{d_z} , e.g., each Z_i is independently sampled from $N(0, 1)$. This is also
1085 the setting studied in nonlinear independent component analysis [18, 37].
- 1086 2. Causal Representation Learning: This setting is more general than the one above where we
1087 relax the independence assumption on the Z_i , and instead assume that they have (typically
1088 unknown) causal relationships among them. For instance, they could satisfy a linear
1089 structural causal model with Gaussian noise, i.e., $Z = AZ + \epsilon, \epsilon \sim N(0, I)$ where A
1090 encodes a weighted directed acyclic graph. This setting is generalizes the previous setting,
1091 since having no causal relationships (i.e., $A = 0$) implies joint independence.

1092 As explained earlier, in both these domains, a critical notion is that of identifiability [45, 21, 116],
1093 which posits that the given dataset(s) are diverse enough for the modeling assumptions, in order to
1094 ensure that a unique set of parameters fit the data. It’s folklore that the disentangled representation
1095 learning model is not identifiable if all Z_i are Gaussian [38, 61]. However, under appropriate as-
1096 sumptions, e.g., distributional, sparsity or observed side-information, the model becomes identifiable,
1097 see e.g., [45, 36, 10, 93, 51, 68, 127, 49, 11, 128, 31, 85]. In addition, various works have proposed
1098 methods to learn them [28, 119, 22, 121, 57, 20, 11, 53, 12].

1099 C.2 Conceptual differences

1100 In this section, we highlight the conceptual differences between causal representation learning and
1101 our framework.

1102 **Are causal generative concepts necessarily interpretable?** Moreover, we are constantly conjuring
1103 new concepts of interest since human-interpretable concepts are constantly evolving, e.g., the concept
1104 of mobile phones did not exist 100 years ago, but is a valid concept to learn now. Therefore, as
1105 opposed to working with a rigid model as in causal representation learning, we take the approach of
1106 working with a dynamic representation learning model. Finally, even if individual causal factors *are*
1107 interpretable (which may be the case in certain applications), the perspective that we take in this work
1108 is that the number of true generative factors could be prohibitively large so that attempting to extract
1109 and interpret all of them together is infeasible, whereas the number of desired human-interpretable
1110 concepts is much smaller and more manageable.

1111 **Number of environments needed** When the ground truth generative process has ambient latent
1112 dimension d_z , for causal representation learning to be feasible, we usually require d_z environments or
1113 datasets. For instance, in the iVAE setting [45] with k sufficient statistics, we require $d_z k + 1 \geq d_z + 1$
1114 environments. This is indeed necessary, as counterexamples show. However, it’s not clear what the
1115 value of d_z is for complex datasets, and it could potentially be prohibitively large.

1116 But the question remains, do we need to learn the entire generative model for solving downstream
1117 tasks? Along these lines, there is a tremendous research effort attempting to relax such requirements
1118 by imposing various inductive or domain biases and by building a theory of partial identifiability
1119 [49, 59, 50]. This is for good reason, since even though it would be ideal to learn the full ground
1120 truth generative model, it may be prohibitively large and moreover it may not be necessary for the
1121 downstream tasks we care about, therefore it suffices to learn what is necessary. On this note, the
1122 related task of learning only a subset of the generative latent variables is also not easy as the latent
1123 variables interact in potentially complicated ways.

1124 In this work, we show that if we only wish to learn $n \ll d_z$ concepts, it suffices to have $O(n)$
 1125 environments instead of $\Omega(d_z)$ environments. Therefore, our results can be viewed as a result on
 1126 partial identifiability with a sublinear number of environments.

1127 **Multi-node interventions** Multi-node interventions are an exciting area of study in CRL, since
 1128 they are a natural extension of existing works and are more useful for modeling various real-life
 1129 datasets where it can be hard to control precisely one factor of variation. This is easily incorporated in
 1130 our setting by utilizing non-atomic concepts, since each non-atomic concept is a collection of vectors
 1131 corresponding to atomic concepts and can be modified simultaneously by changing the valuation.

1132 **Conditional vs. interventional data** In this work we focus on conditional data and identification
 1133 of concept structure, while a recent trend in CRL is to focus on interventional data and identification
 1134 of the causal structure [97, 109, 12, 42, 113]. For causal models, interventions are a natural approach
 1135 to solving the identifiability problem, however, in the absence of an assumed causal model (as in
 1136 our framework), interventions may not even be formally well-defined. In our framework, we do not
 1137 think of concepts as being causal variables that are connected by a graph. (We note that an interesting
 1138 approach would be to study learning concepts over a given causal generative model, which is an
 1139 intriguing direction for future study that we do not pursue in this work).

1140 By contrast, conditional data does not require the formal framework of causal models, and is often
 1141 more frequently available in practice. Conditional data can be obtained by selection through filtering,
 1142 e.g., patients that are admitted to different hospitals based on the severity of their condition or by the
 1143 availability of label information as in the CLIP setting [81]. Thus conditional data can be obtained
 1144 by observing the system in different condtions. On the other hand interventional data requires
 1145 manipulation of the system which is more difficult to obtain in general.

1146 D Alternate definitions of concept conditional measure

1147 In this section, we present alternate feasible definitions for data distributions than the one we
 1148 introduced in Appendix A.2. While we went with the definition most suited for practice, these
 1149 alternate definitions are also justifiable in different scenarios and are exciting avenues for further
 1150 study.

1151 We want to essentially define a concept C via a conditional measure p_C where the concept C is
 1152 identified with an affine subspace $C = \{Z \in \mathbb{R}^{d_z} : A^C Z = b^C\}$ for some $A^C \in \mathbb{R}^{k \times d_z}$, $b^C \in \mathbb{R}^k$.
 1153 We consider the shifted parallel linear subspace $C_0 = \{Z : A^C Z = 0\}$ and the orthogonal splitting
 1154 $\mathbb{R}^{d_z} = C_0 \oplus V$. Suppose we have a distribution q_V on the space V which will typically be a Gaussian
 1155 centered around $v^C \in V$ which is the unique solution of $A^C v^C = b^C$. In addition we have a base
 1156 distribution p on \mathbb{R}^{d_z} . We will assume that all distributions have a smooth density so that conditional
 1157 probabilities are pointwise well defined. There are at least three ways to create the context conditional
 1158 measure p_C .

- 1159 1. The first option is to enforce that the distribution of the V marginal $p_C(v) = \int_{C_0} p_C(v, c) dc$
 1160 exactly matches $q_V(v)$ while the in-plane distribution $p_C(c|v = v_0) \propto p_C(c, v_0)$ remains
 1161 invariant, i.e., equals $p(c|v = v_0)$. Under this condition, there is a unique measure p_C given
 1162 by

$$p_C(c, v) \propto q_V(v) \frac{p(c, v)}{\int_{C_0} p(c', v) dc'}.$$

1163 In other words, to get (c, v) we sample $v \sim q_V$ and then $c \sim p(c|v)$ according to the
 1164 conditional distribution.

- 1165 2. The second option is to again enforce the V marginal but instead of keeping the in plane
 1166 distribution we average over the V space. Then we obtain

$$p_C(c, v) \propto q_V(v) \int_V p(c, v') dv'.$$

1167 This corresponds (vaguely) to a $\text{do}(v)$ operation from causal inference, i.e., we sample
 1168 according to $p(v, c)$ and then do a random intervention on v with target distribution q_V .

1169 3. The third option is to take a Bayesian standpoint. Then we view p as a prior and q_V as
 1170 the context dependent acceptance probability, i.e., we sample by p and then accept with
 1171 probability q_V . Then we find

$$p_C(c, v) = \frac{p(c, v)q_V(v)}{\int p(c, v)q_V(v) dv dc} \propto p(c, v)q_V(v). \quad (76)$$

1172 This is probably the closest aligned to practice, so this is the one we study in this work. To
 1173 justify this option, imagine the following scenario. If we wish to learn the concept of *red*
 1174 *color*, a first step would be to curate a dataset of red objects. To do this, we first consider
 1175 a collection of photos of objects of varying color and then filter out the ones that look red.
 1176 The concept conditional measure we define aligns with this process. To learn the actual
 1177 red concept accurately, our theory predicts that it is sufficient to have additional datasets
 1178 of objects that are not red, from which we can distinguish red objects, thereby learning the
 1179 concept of red color.

1180 The next question is how to define the measure q_V . When considering a single concept $A^C Z = b^C$ the
 1181 most natural option to consider $N(v^C, \sigma^2 \text{Id}_V)$ where $v^C \in V$ is the unique solution of $A^C v^C = b^C$
 1182 and $\sigma > 0$ is a positive constant. This is what we do in this work (note that σ^2 can be set to 1 by
 1183 scaling the concept and valuation accordingly).

1184 However, we can also use alternate definitions as suggested above. For instance, we can set $AZ \stackrel{\mathcal{D}}{=} N(b^C, \text{Id})$.
 1185 Then $Z \sim N(v^C, (A^T A)^{-1})$. However, this runs into some technical issues we sketch
 1186 (and leave to future work to handle this). Consider the intersection of multiple concepts C^e . In this
 1187 case the concept space is given by the intersection $C = \bigcap C^e$ and $C_0 = \bigcap (C^e)_0$ and we have the
 1188 orthogonal decomposition $\mathbb{R}^{d_z} = C_0 \oplus \sum V^e$. In general the spaces V^e are not necessarily orthogonal
 1189 but it is reasonable to assume that the non-degeneracy condition $\dim(\sum V^i) = \sum \dim(V^e)$ holds.
 1190 Now set $V = \sum V^e$. If we choose just the standard normal distribution for q_{V^e} we can define just as
 1191 in our approach

$$q_V \sim N(v^C, \sigma^2 \text{Id}_V). \quad (77)$$

1192 The second option is to enforce that the marginals of q_V agree with q_{V^e} , i.e., $q_V(\Pi_{V^e}(v) \in O) =$
 1193 $q_{V^e}(O)$ for $O \subset V^e$. This results in the set of equations for all i

$$A^e \Sigma (A^e)^T = \text{Id}_{V^e}. \quad (78)$$

1194 It is likely that this system has a unique solution when non-degeneracy holds for V^e and this is clearly
 1195 true for orthogonal spaces but it is not clear how to solve this in general.

1196 E Analysis of pretrained CLIP models

1197 In this section we provide additional experimental details and further results for the analysis of
 1198 pretrained CLIP models [81].

1199 E.1 Experimental Details

1200 We transform the images from the 3d-Shapes dataset to match the CLIP training data, i.e., reshape
 1201 to images of size 224 and match the channel distributions. Then we calculate the embeddings
 1202 for all images in the dataset using two CLIP models, a model with a vision transformer backbone
 1203 (‘ViT-B/32’) and a model with a Resnet backbone (‘RN101’)². We split the embedded images in to
 1204 training and test sets of equal size. Then for any factor of variation (orientation of the scene, shape
 1205 and scale of the object, and hue of floor, wall, and object) we perform the following procedure. For
 1206 each pair of values of a factor of variation we run logistic regression on the embeddings for those
 1207 two values of the concept to classify which value is taken for a given embedding. We average the
 1208 directions of the logistic regression vectors β_i , i.e., consider $\bar{\beta} = N^{-1} \sum_{i=1}^N \beta_i$. Since the direction
 1209 is defined only up to a sign (depending on the order of the two groups) we repeatedly replace β_i by
 1210 $-\beta_i$ if the scalar product with the current mean is negative (this is a heuristic procedure to align β_i
 1211 with $\bar{\beta}$). We then use the learned concept vectors $a = \bar{\beta}$ to evaluate the concept valuations on the

²Models are publicly available under <https://github.com/openai/CLIP>

1212 held out test data, i.e., we evaluate $\langle a, Z \rangle$ where $Z = f^{-1}(X)$ is the embedding of an image X . The
 1213 preprocessing to calculate the CLIP image embeddings required few hours on a A100-GPU. The
 1214 remaining evaluations were performed on a standard notebook.

1215 E.2 Further results

1216 Here we report the mean and standard deviations of the per-class concept valuations $\langle a, Z \rangle$ for the
 1217 concept vectors learned as described in Section E.1. The results for the six factors of variation can be
 1218 found in Tables 2, 3, and 4. We observe that shape, scale, and orientation are well aligned with linear
 1219 subspaces. For the hue variables this still holds to some degree the discrepancy might be attributed
 1220 to hue not being an atomic concept (colours are typically represented by at least two numbers).
 1221 Moreover, we consider the correlation coefficient of the valuations obtained for different embedding
 1222 models, i.e., for $\langle a^{M_1}, Z_i^{M_1} \rangle$ and $\langle a^{M_2}, Z_i^{M_2} \rangle$ where a^{M_1} and a^{M_2} are concept vectors for the same
 1223 concept and two different models and $Z_i^{M_1}$ and $Z_i^{M_2}$ denote the embeddings of the two models M_1
 1224 and M_2 of sample X_i . We report these correlation coefficients for the two CLIP models in Table 5.
 1225 The results indicate that the valuations indeed approximately agree up to a linear transformation.
 1226 Note that for the scene orientation attribute the valuation corresponds to the absolute value of the
 1227 angle.

Table 2: Mean valuations and standard deviation on the test set for the floor hue and wall hue attributes.

Floor hue	Vit-B/32	RN101	Wall hue	Vit-B/32	RN101
0.0	-1.4 ± 1.4	-0.3 ± 0.9	0.0	1.1 ± 1.3	-1.5 ± 1.4
0.1	4.5 ± 1.5	1.4 ± 0.8	0.1	2.8 ± 1.3	1.8 ± 1.0
0.2	4.3 ± 1.3	3.2 ± 0.8	0.2	3.3 ± 1.1	1.5 ± 0.9
0.3	2.2 ± 1.4	3.0 ± 0.8	0.3	1.7 ± 1.0	0.8 ± 0.8
0.4	1.2 ± 1.5	2.2 ± 0.8	0.4	0.8 ± 1.3	0.5 ± 0.9
0.5	0.0 ± 1.1	0.5 ± 0.8	0.5	-0.6 ± 1.2	-0.6 ± 1.1
0.6	-2.8 ± 1.3	-0.4 ± 0.9	0.6	-3.3 ± 1.2	-2.3 ± 1.1
0.7	-5.8 ± 1.5	-2.0 ± 1.0	0.7	-3.6 ± 1.2	-3.7 ± 1.0
0.8	-3.8 ± 1.4	-1.3 ± 0.9	0.8	-1.4 ± 1.1	-2.0 ± 1.0
0.9	-3.2 ± 1.4	-1.0 ± 0.8	0.9	-0.6 ± 1.2	-2.0 ± 1.1

Table 3: Mean valuations and standard deviation on the test set for the object hue and scene orientation attributes.

Object hue	Vit-B/32	RN101	Scene orientation ($^\circ$)	Vit-B/32	RN101
0.0	-0.3 ± 1.5	-0.1 ± 1.1	-30.0	-4.9 ± 1.4	-0.0 ± 1.1
0.1	4.8 ± 2.1	1.4 ± 1.0	-25.7	-4.0 ± 1.3	0.4 ± 1.2
0.2	6.0 ± 2.0	2.7 ± 0.8	-21.4	-2.9 ± 1.3	-0.8 ± 1.2
0.3	3.9 ± 1.7	2.6 ± 0.7	-17.1	-0.2 ± 1.4	-1.4 ± 1.1
0.4	2.3 ± 1.4	2.2 ± 0.7	-12.9	3.3 ± 1.5	-3.9 ± 1.1
0.5	-0.5 ± 1.6	0.3 ± 0.9	-8.6	7.5 ± 2.1	-6.7 ± 0.9
0.6	-4.8 ± 1.8	-1.8 ± 0.9	-4.3	7.2 ± 2.4	-7.4 ± 1.1
0.7	-5.6 ± 1.9	-2.4 ± 1.0	0.0	8.2 ± 2.7	-8.2 ± 1.2
0.8	-3.4 ± 1.4	-1.3 ± 0.9	4.3	5.8 ± 2.3	-7.6 ± 1.1
0.9	-1.9 ± 1.4	-0.6 ± 1.0	8.6	6.5 ± 1.9	-7.0 ± 1.0
			12.9	2.0 ± 1.6	-4.7 ± 0.9
			17.1	-2.9 ± 1.3	-2.2 ± 0.9
			21.4	-4.8 ± 1.3	-1.8 ± 1.1
			25.7	-5.7 ± 1.5	-0.7 ± 1.1
			30.0	-6.6 ± 1.8	-0.7 ± 1.1

Table 4: Mean valuations and standard deviation on the test set for the scale and shape attributes.

Scale	Vit-B/32	RN101	Shape	Vit-B/32	RN101
0.8	10.6 ± 2.6	7.0 ± 1.5	Cube	8.2 ± 1.4	6.9 ± 0.9
0.8	8.3 ± 2.1	5.2 ± 1.4	Cylinder	2.9 ± 1.6	2.9 ± 0.9
0.9	5.0 ± 1.9	3.6 ± 1.3	Ball	-3.6 ± 1.6	-1.2 ± 0.7
1.0	1.9 ± 1.9	1.8 ± 1.1	Ellipsoid	-11.8 ± 3.1	-5.5 ± 1.7
1.0	-1.3 ± 1.8	0.2 ± 1.1			
1.1	-4.3 ± 2.0	-1.4 ± 1.2			
1.2	-7.1 ± 2.1	-2.8 ± 1.2			
1.2	-9.3 ± 2.3	-3.9 ± 1.3			

Table 5: Correlation coefficients of the evaluations learned for two different CLIP models evaluated on the full dataset.

Concept	ρ
Floor hue	0.86
Wall hue	0.83
Object hue	0.86
Scale	0.53
Shape	0.95
Orientation	-0.70

1228 F Inference-Time Intervention of Large Language Models

1229 In this section, we first briefly describe Large Language Models and the recent Inference-Time
 1230 Intervention (ITI) technique proposed for LLM alignment, which we build on. Then, we use our
 1231 framework to provide better intuition on some intriguing observations about ITI, including why it
 1232 works. And then we exploit our ideas to improve the performance of ITI by choosing the steering
 1233 direction to be a matrix instead of a vector.

1234 F.1 Preliminaries

1235 **Large Language Models (LLMs)** LLMs are large models capable of generating meaningful
 1236 text given a context sentence. Due to large-scale training, modern LLMs have shown remarkable
 1237 capabilities and achieve expert-human-like performance in many benchmarks simultaneously. The
 1238 architecture of many generative pre-trained transformers (GPT)-style LLMs consists of several
 1239 transformer layers stacked on top of each other. Since we’ll be intervening on them during inference,
 1240 we’ll describe the transformer architecture [112, 24] briefly here. First, the sequence of input tokens
 1241 (tokens are sub-word units) are encoded into a vector x_0 using a (learned) text embedding matrix and
 1242 in many cases also a positional embedding matrix. Then, a series of transformer layers act on this
 1243 vector which passes through a residual stream, to obtain vectors x_0, x_1, \dots, x_n . The final vector x_n
 1244 is then decoded back into token probabilities with a (learned) unembedding matrix. Each transformer
 1245 layer consists of a multi-head attention mechanism and a standard multilayer perceptron, which
 1246 captures the nonlinearity.

1247 In the l th layer, each single multi-head attention mechanism can be described as

$$x_{l+1} = x_l + \sum_{h=1}^H Q_l^h x_l^h, \quad x_l^h = \text{Att}_l^h(P_l^h x_l)$$

1248 Here, P_l^h and Q_l^h are matrices that linearly map the vector to an activation space and back respectively,
 1249 and Att denotes the attention mechanism that allows communication across tokens. Here, we have
 1250 kept the notation consistent with Li et al. [56] for the sake of clarity.

1251 In our setting, we consider the entire set of activations as the learnt latent vector Z . That is, the
 1252 input is $x = x_0$ and the pre-trained model is essentially the function f such that $f(x)$ consists of
 1253 the concatenation of the vectors $\{x_l\}_{l \geq 1}$, the intermediate activations $\{x_l^h\}_{l \geq 0}$ and also the output
 1254 of the linear transformations $\{P_l^h x_l\}_{l \geq 0}, \{Q_l^h x_l^h\}_{l \geq 0}$. Our theory hinges on the assumption that
 1255 pre-trained LLMs satisfy the linear representation hypothesis, that is, various relevant concepts
 1256 can be realized via linear transformations of the latent transformation $f(x)$. Indeed, this has been
 1257 empirically observed to hold in many prior works [15, 105, 71, 69, 56, 74, 33, 44] (see also related
 1258 works on geometry of representations [43, 44] and references therein). It’s a fascinating question
 1259 why such models trained with next token prediction loss also learn linear representations of various
 1260 human-interpretable concepts such as sentiment, see Jiang et al. [44] for recent progress on this
 1261 problem.

1262 It’s well-known that despite large-scale pretraining and subsequent improvement of pre-trained
 1263 models via techniques like Reinforcement Learning with Human Feedback (RLHF) and Supervised
 1264 Fine-Tuning (SFT) [73, 6, 106], significant issues still remain [94], e.g., the model can hallucinate
 1265 or generate incorrect responses (even though the model *knows* the correct response which can be
 1266 extracted via other means, e.g., Chain-of-Thought prompting [118]). Various methods have been
 1267 proposed to fine-tune the models [73, 6, 7, 106, 82] but many of them are expensive and time-
 1268 and resource-intensive as they requires huge annotation and computation resources. Therefore,
 1269 more efficient techniques are highly desired, one of which is the category of methods known as
 1270 activation patching. activation patching (also called activation editing or activation engineering)
 1271 [34, 115, 99, 108, 129, 124, 55, 66].

1272 **Inference-Time Intervention, an activation patching method for truthfulness** Activation patch-
 1273 ing is a simple minimally invasive technique to align LLMs to human-preferences. Specifically, given
 1274 various concepts such as truthfulness, activation patching makes modifications to the model during
 1275 inference time so that the desired concepts can be aligned. This technique can be thought of as an
 1276 application of the emerging field of mechanistic interpretability [72], which aims to interpret the
 1277 learnt latent vector in terms of human-interpretable concepts, thereby allowing us to reverse-engineer
 1278 what large models learn.

1279 Activation patching has many variants [55, 34, 66], but we’ll focus on the simple technique of adding
 1280 *steering vectors* to various intermediate layers during intervention [99, 108, 56, 87]. This means that
 1281 during inference, the output activations are modified by adding a constant vector in order to promote
 1282 alignment of some concept. The vector will be learnt independently based on separate training data.

1283 In particular, a recent technique called Inference-Time Intervention (ITI) was proposed to do this
 1284 for the specific concept of truthfulness. ITI focuses on the activation heads $\{\text{Att}_l^h(P_l^h x_l)\}_{l \geq 0}$ and
 1285 add to them steering vectors in order to promote truthfulness. To learn the steering vectors, a subset
 1286 of the TruthfulQA dataset [58], namely a dataset of questions q_i with annotated true $(a_{i,j}, 0)$ and
 1287 false answers $(a_{i,j}, 1)$, are prepared as $\{q_i, a_i, y_i\}_{i=1,2,\dots}$. For each sample, the question and answer
 1288 are concatenated as a pair and the corresponding activations of the heads x_l^h (for the final token) are
 1289 computed via forward passes. Then, a linear probe sigmoid $(\langle \theta, x_l^h \rangle)$ is independently trained on each
 1290 activation head to distinguish true from false answers. Finally, the top K heads based on the accuracy
 1291 of this classification task are chosen (for a tunable hyperparameter K) and the steering vector θ_l^h for
 1292 the h -th head in layer l is chosen to be the mean difference of the activations between the true and
 1293 false inputs. The intuition is that this direction roughly captures the direction towards truthfulness.

1294 Formally, for the h th head of the l th layer, ITI adds the steering vector $\alpha \sigma_l^h \theta_l^h$ so as to get

$$x_{l+1} = x_l + \sum_{h=1}^H Q_l^h (x_l^h + \alpha \sigma_l^h \theta_l^h), \quad x_l^h = \text{Att}_l^h(P_l^h x_l)$$

1295 during inference. Here, θ_l^h is the steering vector, σ_l^h is the standard deviation of the activations of this
 1296 head along the chosen direction and α is a hyperparameter. That is, the activations are shifted along
 1297 the truthful directions by a multiple of the standard deviation, and this is repeated autoregressively.
 1298 Note that this does not depend on the specific GPT-like model being used. The intuition is that during
 1299 inference, the activations are intervened upon to shift towards the truthful direction. The top K heads
 1300 are chosen to be minimally intrusive and also a design choice based on observations of the probing
 1301 metrics.

1302 **Performance of ITI** In Li et al. [56], ITI was shown to significantly improve the truthfulness of
1303 various LLMs after having been trained on as few as a few dozen samples, compared to what’s
1304 needed for Reinforcement Learning based techniques [73, 29]. ITI was evaluated on the TruthfulQA
1305 benchmark [58], which is a hard adversarial benchmark to evaluate truthfulness of language models.
1306 In particular, it contains 817 questions with a multiple-choice and generation tracks, spanning 38
1307 categories such as logical falsehoods, conspiracies and common points of confusion. For the multiple-
1308 choice questions, the accuracy is determined by the conditional probabilities of candidate answers
1309 given the question. Evaluating the generation track questions is harder, and it is done by generating a
1310 model output and then evaluating it via a finetuned GPT-3-13B model [58, 70]. Moreover, the choice
1311 of the intervention strength α is calibrated so that it’s neither too small (to promote truthfulness)
1312 nor too large (to ensure the original capabilities of the LLM are not lost). To check if the original
1313 capabilities are preserved, [56] compute two additional quantities to measure how far the modified
1314 model deviates from the original model. These are the Cross-Entropy (CE) loss, which is standard in
1315 language modeling and the Kullback–Leibler divergence (KL div.) of the next token probabilities
1316 before and after intervention. To compute these quantities, a subset of Open Web Text is used [80].
1317 Finally, it was shown that ITI implemented on the LLaMA [106], Alpaca [103] and Vicuna [17]
1318 models significantly improved their performance on the TruthfulQA benchmark compared to the
1319 baseline models. Moreover, in many cases, it also beat other techniques such as few-shot prompting
1320 and supervised fine-tuning. Please see Li et al. [56] for additional details.

1321 **F.2 Interesting observations of ITI**

1322 While the elegant ITI technique was designed to align LLMs towards truthfulness in practice, it also
1323 raised fascinating and intriguing questions in mechanistic interpretability. In addition to improving the
1324 technique of ITI itself, our work makes progress towards some of these questions via our framework.

- 1325 1. The authors of Li et al. [56] state in section 2 that although the technique works well in
1326 practice, it’s not clear what ITI does to the model’s internal representations. In addition, prior
1327 works [15, 105, 71, 69, 74, 44] have observed empirically that the latent representations
1328 learned by LLMs seem to have interpretable linear directions, which ITI exploits. We use
1329 our framework to illustrate in more detail one possible explanation of what ITI does to the
1330 model representations and why it works, in the next section.
- 1331 2. The authors visualize the geometry of “truth” representations in section 3.2 of their work via
1332 the following experiment: For the most significant head (layer 14, head 18), after finding the
1333 first truthful direction via the linear probing technique, they remove it and attempt to find a
1334 second probe orthogonal to the first. They find surprisingly that the second probe is also
1335 very informative, leading them to predict that the concept of “truth” lies in a subspace, not
1336 a single direction. Restated in our framework, the concept of truthfulness is a non-atomic
1337 concept (as per Definition 2). This served as an inspiration for our proposed technique in
1338 the next section, where we propose to use steering matrices instead of steering vectors for
1339 LLM alignment.
- 1340 3. As α was increased, the authors observed that truthfulness of the model increased however
1341 helpfulness decreased. This suggests that the “truthfulness” and “helpfulness” concepts
1342 are not atomic (as per Definition 2) however they share certain atomic concepts. We leave
1343 to future work the exciting question of mechanistically extracting such common atomic
1344 concepts.

1345 **F.3 The choice of the steering vector**

1346 In this section, we will use our theoretical framework to get insights about the ITI technique and
1347 use it to improve alignment. First, similar to the multimodal CLIP setting, we will assume that the
1348 non-linearity has already been learned up to a linear transformation (by large-scale training of LLMs).
1349 This aligns with our theoretical insights because the training data for powerful LLMs are diverse, so
1350 they essentially satisfy our core assumptions (see also the related work [32] that proposes that context
1351 is environment in LLM training). Therefore, we simply focus on the downstream tasks, which in this
1352 section is LLM alignment. The difficulty, of course, is that we do not know the concept matrix nor
1353 the valuations.

1354 We will now analyze the truthfulness concept via our framework and give more insight on why the
 1355 mean of the differences is a reasonable choice of steering vector for ITI. Based on our theory, we
 1356 will then provide a modification to this choice that uses steering matrices instead of steering vectors.
 1357 Since this section is based on heuristics and informal assumptions, we will refrain from making any
 1358 formal claims or analyses. Indeed, a formal analysis of concepts in natural language is a hard problem
 1359 in general and we do not attempt it here. We conclude with ideas for potential extensions that’re
 1360 worth exploring in future work.

1361 Denote the function h to be the sequence of head activations $h(x) = (x_l^h)_{l,h} \in \mathbb{R}^d$. Note that while
 1362 we can study general steering vectors for the entire latent space of representations $f(x)$ learned by
 1363 LLMs as some works do, ITI focuses only on steering the head activations $h(x)$, so we will apply
 1364 our framework to this subset representation space. In addition, we will make the simplification that
 1365 we neglect the effects of the steering vector from bottom layers towards the top layers, which we do
 1366 because we are dealing with sparse steering vectors and also, each single head shift is minor and does
 1367 not in isolation change the behavior of the model as verified by experiments [56][Appendix B.1].

1368 Applying our framework, we model the concept of truth via the concept matrix $A \in \mathbb{R}^{d_C \times d}$ and two
 1369 valuations $b_0, b_1 \in \mathbb{R}^{d_C}$ corresponding to *False* and *True* respectively. In other words, the set of false
 1370 sentences and true sentences lie respectively in

$$\mathcal{S}_{false} = \{x | Ah(x) = b_0\}, \quad \mathcal{S}_{true} = \{x | Ah(x) = b_1\}$$

1371 Note that they only approximately lie in these spaces because of our notion of concept conditional
 1372 distribution. However, if we reasonably assume that the Gaussian concentration region is much
 1373 smaller than the separation between these hyperplanes, then the rest of the arguments in this section
 1374 should apply.

1375 Now, a steering vector η is a vector such that it moves the activations from the false space to the true
 1376 space, while keeping other concepts unaffected. That is, if we pick a false sentence x , i.e., $Ah(x) = b_0$,
 1377 then the steering vector $\eta \in \mathbb{R}^d$ essentially steers the activations so that $A(h(x) + \eta) = b_1$. In other
 1378 words, it moves the sentence from false to true. Indeed, many vectors η do satisfy this equality,
 1379 because we could move $h(x)$ to any point in the hyperplane $\{AZ = b_1\}$. Therefore the goal is to find
 1380 an optimal η that does not (significantly) affect other concepts of interest, i.e., $B(h(x) + \eta) \approx Bh(x)$
 1381 (equivalently $B\eta = 0$) for any other concept of interest B . Indeed, a natural choice of the steering
 1382 vector will be $A^+(b_1 - b_0)$ where A^+ is the pseudoinverse of A . This vector will precisely affect this
 1383 concept space and will not affect the concept valuations for any concept orthogonal to A . However,
 1384 there are two issues with this approach: We do not know A and therefore we will approximate this
 1385 steering vector from training samples and there is no guarantee that other concepts of interest are
 1386 orthogonal to A (note that angles between concepts are not even identifiable).

1387 Previous approaches are based on a collection of counterfactual sentence pairs c_i^F, c_i^T which corre-
 1388 spond to a false answer and a true answer for the same question q_i . Consider the i th counterfactual
 1389 pair c_i^F, c_i^T . We will assume the reasonable scenario that the only difference among their concepts is
 1390 the concept of truthfulness. That is, for any other concept of interest B_i for this sample the valuations
 1391 of B_i for these pairs c_i^F and c_i^T are identical. A common strategy is to use the mean

$$\eta = \frac{1}{n} \sum_{i=1}^n h(c_i^T) - h(c_i^F) \quad (79)$$

1392 as a steering vector. Note that if

$$A(h(c_i^T) - h(c_i^F)) \approx b_1 - b_0, \quad (80)$$

1393 i.e., the truthfulness valuation is changed as desired for all samples then

$$A\eta = b_1 - b_0. \quad (81)$$

1394 Moreover, concepts of interest are preserved in two prototypical settings. First, if concepts of interest
 1395 are the same for all samples and the new datapoint, i.e., $B = B_i = B_j$ in which case

$$B\eta = \frac{1}{n} \sum_{i=1}^n B_i(h(c_i^T) - h(c_i^F)) = 0. \quad (82)$$

1396 Similarly, if concepts of interest for a new point x are B_x and the valuations of $B_x(h(c_i^T) - h(c_i^F))$ of
 1397 the counterfactual pairs are random, independent, and centered, then we expect them to approximately
 1398 cancel and

$$B_x \eta \approx 0. \quad (83)$$

1399 Note that in this case, this is not true if just a single steering vector $h(c_i^T) - h(c_i^F)$ is used as a
 1400 steering vector.

1401 This explains why the choice of mean of the activation differences across counterfactual pairs is a
 1402 reasonable choice of steering vector. This is precisely the technique used in ITI. While they also
 1403 experiment with other steering vectors, they found that this works the best for their experiments.

1404 Now, we will continue on our insights to analyze whether we can build better steering vectors η . We
 1405 present two crucial insights based on our analysis so far.

- 1406 1. Looking at our desired equations, any *weighted combination* of $\eta_i = h(c_i^T) - h(c_i^F)$ will
 1407 satisfy $Ah(x) = b_0$, $A(h(x) + \eta) = b_1$ exactly.
- 1408 2. We could potentially choose the steering vector η to be a function of x instead of being a
 1409 constant vector, provided $\eta(x)$ is efficiently computable during inference time.

1410 Exploiting our first insight, we conclude that choosing any weighted combination of the η_i should be
 1411 a reasonable choice of steering vector provided we can control its effects on the spaces orthogonal to
 1412 A . That is, we can choose

$$\eta = \sum_i w_i \eta_i = \sum_i w_i (h(c_i^T) - h(c_i^F))$$

1413 as our steering vector. This gives us the extra freedom to tune the weights w_1, w_2, \dots based on other
 1414 heuristics. Note that this also captures the choice of the top principal component of the steering vector
 1415 as experimented in [105].

1416 Our second observation suggests that even the steering vector η could be a function of x , namely
 1417 $\eta(x)$, provided it's efficiently computable during inference. Therefore, this suggests the usage of

$$\eta(x) = \sum_i w_i(x) (h(c_i^T) - h(c_i^F))$$

1418 as our steering vector where the weights $w_i(x)$ depend on x .

1419 Based on these two observations, we propose our ITI modification. We choose the steering vector
 1420 to be dependent on the context x , with weights chosen to be $w_i = \langle \lambda(x), \lambda(c_i^F) \rangle$ for a sentence
 1421 embedding λ (such as Sentence-BERT [86]). That is,

$$\eta(x) = \sum_i \langle \lambda(x), \lambda(c_i^F) \rangle (h(c_i^T) - h(c_i^F))$$

1422 Indeed, this is reasonable as if a context x is close to c_i^F for a specific training sample i in terms of
 1423 their sentence embeddings $\lambda(x)$ and $\lambda(c_i^F)$, then this particular sample's steering vector should be
 1424 upsampled. In other words, we can think of the training sample contexts as voting on their respective
 1425 counterfactual steering vector, with weights determined by the similarity between the representation
 1426 of the test context and the representation of the sample context. A justification would be that $B(x)$
 1427 (the relevant concepts for a datapoint) depend smoothly on x (proximity is measured by similarity of
 1428 embeddings) so it makes sense to upweight close points to enforce that x preserves similar concepts.

1429 Finally, we need to argue that we can compute this efficiently during inference. For this, we exploit
 1430 the structure of our steering vector representation as follows.

$$\begin{aligned} \eta(x) &= \sum_i \langle \lambda(x), \lambda(c_i^F) \rangle (h(c_i^T) - h(c_i^F)) \\ &= \left(\sum_i (h(c_i^T) - h(c_i^F)) \lambda(c_i^F)' \right) h(x) \\ &= Mh(x) \end{aligned}$$

1431 for the matrix $M = \sum_i (h(c_i^T) - h(c_i^F)) \lambda(c_i^F)'$, where v' denotes the transposed vector. We remark
 1432 that the weights $w_i(x)$ as used could potentially be negative but this is not an issue since the
 1433 projection of the corresponding counterfactual vector in the direction of B is still random and we
 1434 finally normalize $\eta(x)$, so the magnitude doesn't matter.

1435 Therefore, this steering can be done efficiently by precomputing the *steering matrix* M and then
 1436 during inference, we simply compute the steering vector $\eta(x)$ as $\eta(x) = Mh(x)$.

1437 In Table 6, we show the results of
 1438 our experiments with steering matri-
 1439 ces. We use the open-source large lan-
 1440 guage model LLaMA [106] with 7 bil-
 1441 lion parameters (open sourced version
 1442 from Hugging Face) and the sentence
 1443 transformer SBERT [86] for the sen-
 1444 tence embedding. We report the ac-
 1445 curacy of the multiple-choice track of
 1446 TruthfulQA [56] over 3 random seeds

Technique	α	Acc.	CE loss	KL div.
Baseline	-	0.257 ± 0.00005	2.16 ± 0.02	0.0 ± 0.00
Random direction	20	0.258 ± 0.002	2.19 ± 0.02	0.02 ± 0.002
CCS direction	5	0.262	2.21	0.06
ITI: Probe weight dir.	15	0.270 ± 0.004	2.21 ± 0.02	0.06 ± 0.005
ITI: Mass mean shift	20	0.288 ± 0.004	2.41 ± 0.08	0.27 ± 0.007
Steering matrices (ours)	15	0.295 ± 0.02	2.61 ± 0.07	0.41 ± 0.04

1447 Table 6: Comparison of steering vectors for LLM alignment
 1448 and also the Cross-Entropy Loss and KL divergence of the model pre- and post-intervention. All
 1449 hyperparameters are tuned as per [56] and the experiments are performed on eight A6000 GPUs.
 1450 Higher accuracy is better and lower CE loss, and KL divergence indicate that the original model has
 1451 not been significantly modified. Here, the baselines are the unmodified model, random direction
 1452 intervention, Contrast-Consistent Search (CCS) direction [15] and two different direction choices
 using vanilla ITI; and 2-fold cross validation is used.

1453 We see that the multiple-choice accuracy improved, showcasing the potential of our steering matrices
 1454 technique which is novel in the field of LLM alignment to the best of our knowledge. This is meant
 1455 to be a proof of concept and not meant to be a comprehensive study of this specific technique.
 1456 For exploratory purposes, we outline potential modifications to our technique below, which could
 1457 potentially improve the performance, both in terms of accuracy as well as in terms of invasiveness.
 1458 These form an exciting direction for a more comprehensive study of our proposed ideas, which we
 1459 leave for future work.

1460 **Implementation considerations** We briefly note down some design choices we made in our
 1461 implementation of the above method.

- 1462 1. Since $\eta(x)$ is a function of x , the standard deviation of the activation projection on this
 1463 direction, i.e., $\sigma_i^h(x)$ cannot be precomputed (as Li et al. [56] do), therefore we compute
 1464 them dynamically during inference, which takes little overhead with fast tensorization
 1465 operations (in particular, this is not the slow step).
- 1466 2. We opted to go with evaluating the model only on the multiple-choice questions. This is
 1467 partly because to evaluate the generated text, the recommended method is to use fine-tuned
 1468 GPT-3-13B models but OpenAI have retired many of their older models as of this year,
 1469 and therefore, the entire batch of experiments would have to be rerun with their newer
 1470 models which could potentially change the baselines, and also because this work is a
 1471 proof-of-concept rather than a comprehensive evaluation.
- 1472 3. For computing the sentence embeddings, we only use the question prompts, as they contain
 1473 all relevant contexts. And we normalize $\eta(x)$ during inference time.

1474 **Additional ideas for improvement** We re-iterate that our experimental exploration is not exhaustive
 1475 and the preliminary experiments are merely meant to be a proof-of-concept. In this section, building
 1476 on our insights, we outline some further ideas to improve the performance of ITI. We leave to future
 1477 work to comprehensively explore these techniques in order to extract better performance towards
 1478 LLM alignment.

- 1479 1. Note that we opted to go with the weights $\langle \lambda(x), \lambda(c_i^F) \rangle$ where λ was chosen to be a
 1480 sentence transformer embedding [86]. While this is a reasonable choice, similarity metrics
 1481 could be measured in other ways, e.g., with other sentence embedding models.

- 1482 2. Going further, the weights do not have to be similarity scores and could be chosen via other
1483 heuristics. For instance, they could be chosen to be constants but potentially be optimized
1484 using a hold-out test set.
- 1485 3. As Li et al. [56] noted, the ITI technique could be applied on top of fine-tuned models in
1486 order to further improve their performance. Therefore, our proposed modification could also
1487 potentially be applied on top of fine-tuned models.

1488 G Contrastive algorithm for end-to-end concept learning

1489 In this section, we present an end-to-end framework based on contrastive learning to learn the
1490 nonlinearity as well as concepts from data. This is inspired by the methods of the CRL work [12].
1491 The model architecture is designed based on our concept conditional distribution parametrization.
1492 The core idea is as follows. For each concept conditional distribution X^e , we train a neural network
1493 to distinguish concept samples $x \sim X^e$ from base samples $x \sim X^0$. In Lemma 3, we derive the
1494 log-odds for this problem. Then, to learn the n atomic concepts up to linearity, we build a neural
1495 architecture for this classification problem with the final layer mimicking the log-odds expression
1496 above, which can then be trained end-to-end. Because of the careful parametrization of the last layer,
1497 this will encourage the model to learn the representations as guaranteed by our results.

1498 First, we will derive the computation of the true log-odds.

1499 **Lemma 3.** *For any concept index e , there exist some constants c_e such that*

$$\ln(p^e(Z)) - \ln(p(Z)) = \sum_{i=1}^n \left(-\frac{1}{2} M_{ei} \langle a_i, Z^e \rangle^2 + B_{ei} \langle a_i, Z^e \rangle \right) + c_e$$

1500 where M, B are the environment-concept matrix and the environment-valuation matrix defined in (7)
1501 and (8).

1502 *Proof.* This follows from Eq. (13) in the proof of Theorem 3. □

1503 From our main identifiability results, we can assume without loss of generality that the concept vectors
1504 we learn are coordinate vectors. In other words, we consider a neural network h^θ with parameters θ
1505 with output neurons $h_1^\theta, \dots, h_n^\theta$ such that the n atomic concepts will now correspond to the concept
1506 vectors e_1, \dots, e_n (which is reasonable as they are only identifiable up to linear transformations).
1507 Therefore, for each environment e , we can train classifiers of the form

$$g_e(X, \alpha^e, \beta_k^e, \gamma_k^e, \theta) = \alpha^e - \sum_{k=1}^{\dim(C_e)} (\beta_k^e h_k^\theta(X))^2 + \sum_{k=1}^{\dim(C_e)} \gamma_k^e (h_k^\theta(X))$$

1508 equipped with standard cross-entropy loss, for hyperparameters $\alpha^e, \beta_k^e, \gamma_k^e, \theta$. Indeed, this is reason-
1509 able since if the training reaches the global optima in the ideal case, then the loss function will corre-
1510 spond to the Bayes optimal classifier and therefore, $g_e(X, \alpha^e, \beta_k^e, \gamma_k^e, \theta) = \ln(p^e(Z)) - \ln(p(Z))$,
1511 which along with Lemma 3 will suggest that the learnt network h is linearly related to the function
1512 $A^e f^{-1}$, as desired. Lastly, we choose the loss function to be the aggregated CE loss and an extra
1513 regularization term. That is,

$$\mathcal{L} = \sum_e \underbrace{-\mathbb{E}_{j \sim \text{Unif}(\{0, e\})} \mathbb{E}_{X \sim X^e} \left(\ln \frac{e^{\mathbf{1}_{j=e} g_e(X)}}{1 + e^{g_e(X)}} \right)}_{\text{CE loss for environment } e} + \eta \|\beta\|_1$$

1514 for a regularization hyperparameter η .

1515 **Sampling from concept conditional distributions** A common task in controllable generative
1516 modeling is being able to generate data from a known concept. Note that this is not straightforward
1517 in our setting because the normalization term in Eq. (2) is not efficiently computable. To do this
1518 efficiently, we also outline a simple algorithm (Algorithm 1 in Appendix I) to sample from the concept
1519 conditional distribution for a known concept. Our proposed algorithm is based on rejection sampling
1520 and the algorithm as well as the complexity analysis is deferred to Appendix I.

1521 H Additional details about the synthetic setup

1522 In this section, we detail the synthetic setup in Section 5. The base distribution is sampled from a
1523 Gaussian mixture model with 3 components whose parameters are chosen randomly. The weights are
1524 randomly chosen from $\text{Unif}(0.3, 1)$ (and then normalized), the entries of the means are chosen from
1525 $\text{Unif}(-1, 1)$ and the covariance is chosen to be a diagonal matrix with entries in $\text{Unif}(0.01, 0.015)$
1526 (note that the diagonal nature doesn't really matter since a map f will be applied to this distribution).
1527 The mixing function f is chosen to be either (i) linear or (ii) nonlinear with a 1-layer MLP containing
1528 16 hidden neurons and $\text{LeakyReLU}(0.2)$ activations.

1529 The number of concepts n is intentionally chosen to be less than the ground truth dimension d_z
1530 and the number of concepts is $m = n + 1$ as per our theory. The concepts are taken to be atomic,
1531 with the concept vectors and valuations chosen randomly, where each entry of the concept vector
1532 is chosen i.i.d from $\text{Unif}(-0.3, 0.3)$, and the resampling distribution is chosen to be a Gaussian
1533 with variance 0.005. Finally, we choose 5000 samples per environment, sampled via the rejection
1534 sampling Algorithm 1. For the contrastive algorithm, we choose the architecture to either be linear or
1535 nonlinear with a 2-layer MLP with 32 hidden neurons in each layer, with the final parametric layer
1536 chosen based on the known concept, to have the form described above. We train for 100 epochs,
1537 on a single A6000 GPU, with $\eta = 0.0001$ and use Adam optimizer with learning rates 0.5 for the
1538 parametric layer and 0.005 for the non-parametric layer, with a Cosine Annealing schedule [62].

1539 I Controllable generative modeling via rejection sampling

1540 In this section, we will describe how to sample from a concept conditional distribution with a known
1541 concept. Once the concepts are learned in our framework, we can use this technique to generate new
1542 data satisfying various desired concepts, which will aid in controllable generative modeling.

1543 Consider the base distribution on $Z \in \mathbb{R}^{d_z}$ with density $p(Z)$. Suppose we wish to sample from
1544 a concept C given by $AZ = b$ and resampling distribution q . We additionally assume that q is
1545 efficiently computable and an upper bound L is known for its density, i.e., $L \geq \max(q)$.

1546 Recall that the desired density is defined as

$$p_C(Z) \propto p(Z) \prod_{i \leq \dim(C)} q((AZ - b)_i)$$

1547 Note that it's infeasible to compute the normalization constant for such complex distributions.
1548 However, we bypass this by using rejection sampling. We describe the procedure in Algorithm 1.

Algorithm 1: Rejection sampling for controllable generative modeling

Input:

- Base distribution p
- Resampling distribution q with upper bound $L \geq \max(q)$
- Concept C with transformation A and valuation C

Output: Returns a single sample from $p_C(Z)$

```
1  $M = L^{\dim(C)}$ 
  // Repeat trials until condition is met
2 while True do
3    $Z = \text{yield}(p)$ 
4    $U = \text{yield}(\text{Unif}(0, 1))$ 
5    $R = \frac{1}{M} \prod_{i \leq \dim(C)} q((AZ - b)_i)$ 
6   if  $R \geq U$  then
7     return  $Z$ 
```

1549 Informally, we first sample $Z \sim p$ (we overload notation for both density and the distribution) and an
 1550 independent variable $U \sim Unif(0, 1)$, the uniform distribution on $(0, 1)$. We accept the variable Z if

$$\frac{1}{M} \prod_{i \leq \dim(C)} q((AZ - b)_i) \geq U$$

1551 for a predetermined upper bound M on the quantity $\prod_{i \leq \dim(C)} q((AZ - b)_i)$. If the inequality is
 1552 false, we simply reject the sample and repeat.

1553 Now we will argue why this algorithm is correct, which is accomplished in Theorem 4. Let

$$N_C = \int_Z p(Z) \prod_{i \leq \dim(C)} q((AZ - b)_i)$$

1554 be the normalization constant in the definition of $p_C(Z)$. Therefore

$$p_C(Z) = \frac{1}{N_C} p(Z) \prod_{i \leq \dim(C)} q((AZ - b)_i)$$

1555 **Lemma 4.** *Let $M \geq \max(q)^{\dim(C)}$ The acceptance probability of each iteration of the while loop*
 1556 *in Algorithm 1 is $Pr[Z \text{ accepted}] = \frac{N_C}{M}$*

1557 *Proof.* We have

$$\begin{aligned} Pr[Z \text{ accepted}] &= Pr_{U,Z} \left[U \leq \frac{1}{M} \prod_{i \leq \dim(C)} q((AZ - b)_i) \right] \\ &= Pr_{U,Z} \left[U \leq \prod_{i \leq \dim(C)} \frac{q((AZ - b)_i)}{\max(q)} \right] && \text{since } M \geq \max(q)^{\dim(C)} \\ &= \int_Z Pr_U \left[U \leq \prod_{i \leq \dim(C)} \frac{q((AZ - b)_i)}{\max(q)} \right] p(Z) dZ && \text{as } U, Z \text{ are independent} \\ &= \int_Z \left[\prod_{i \leq \dim(C)} \frac{q((AZ - b)_i)}{\max(q)} \right] p(Z) dZ && \text{since } \frac{q((AZ - b)_i)}{\max(q)} \leq 1 \text{ always} \\ &= \int_Z \frac{N_C p_C(Z)}{M} dZ \\ &= \frac{N_C}{M} \end{aligned}$$

1558 □

1559 Before we prove correctness, we will remark on the expected number of trials needed for accepting
 1560 each sample.

1561 **Corollary 1.** *The expected number of trials needed to generate a single sample is $\frac{M}{N_C}$*

1562 *Proof.* Note that each iteration of the while loop is independent, therefore the number of trials until
 1563 acceptance is distributed as a geometric random variable whose expectation is the inverse of the
 1564 parameter. □

1565 This suggests that for our algorithm to be efficient in practice, M should be chosen as small as
 1566 possible, i.e., estimates of $\max(q)$ should be as tight as possible.

1567 **Theorem 4.** *Algorithm 1 yields samples from the concept conditional distribution p_C .*

1568 *Proof.* The proof is at heart the proof of correctness of rejection sampling. For arbitrary parameters
 1569 $t_1, \dots, t_{d_z} \in \mathbb{R}$, let's compute the cumulative density of the samples output by Algorithm 1 and show
 1570 that it matches the cumulative distribution function of $p_C(Z)$ evaluated at t_1, \dots, t_{d_z} , which will
 1571 complete the proof. That is, we wish to calculate

$$Pr[Z_1 \leq t_1, \dots, Z_{d_z} \leq t_{d_z} | Z \text{ accepted}] = \frac{Pr[Z_1 \leq t_1, \dots, Z_{d_z} \leq t_{d_z}, Z \text{ accepted}]}{Pr[Z \text{ accepted}]}$$

1572 We already computed the denominator in Lemma 4. Therefore,

$$\begin{aligned} & Pr[Z_1 \leq t_1, \dots, Z_{d_z} \leq t_{d_z} | Z \text{ accepted}] \\ &= \frac{M}{N_C} Pr[Z_1 \leq t_1, \dots, Z_{d_z} \leq t_{d_z}, Z \text{ accepted}] \\ &= \frac{M}{N_C} \mathbb{E}_Z [\mathbb{1}_{Z_1 \leq t_1} \dots \mathbb{1}_{Z_{d_z} \leq t_{d_z}} \cdot \mathbb{E}_U[\mathbb{1}_{Z \text{ accepted}}]] \\ &= \frac{M}{N_C} \mathbb{E}_Z \left[\mathbb{1}_{Z_1 \leq t_1} \dots \mathbb{1}_{Z_{d_z} \leq t_{d_z}} \cdot \frac{1}{M} \prod_{i \leq \dim(C)} q((AZ - b)_i) \right] \quad \text{from the proof of Lemma 4} \\ &= \int_Z \mathbb{1}_{Z_1 \leq t_1} \dots \mathbb{1}_{Z_{d_z} \leq t_{d_z}} \cdot \frac{1}{N_C} \prod_{i \leq \dim(C)} q((AZ - b)_i) p(Z) dZ \\ &= \int_Z \mathbb{1}_{Z_1 \leq t_1} \dots \mathbb{1}_{Z_{d_z} \leq t_{d_z}} \cdot p_C(Z) dZ \end{aligned}$$

1573 which is precisely the cumulative distribution function of $p_C(Z)$ evaluated at t_1, \dots, t_{d_z} . □