FROM ISOLATION TO ENTANGLEMENT: WHEN DO INTERPRETABILITY METHODS IDENTIFY AND DISENTANGLE KNOWN CONCEPTS?

Anonymous authors

000

001

002

004

006

008 009 010

011

013

014

016

017

018

019

021

024

025

026

027

028

029

031

034

037

040

041

042

043

044

047

048

051

052

Paper under double-blind review

ABSTRACT

A central goal of interpretability is to recover representations of causally relevant concepts from the activations of neural networks. The quality of these concept representations is typically evaluated in isolation, and under implicit independence assumptions that may not hold in practice. Thus, it is unclear whether common featurization methods—including sparse autoencoders (SAEs) and sparse probes recover disentangled representations of these concepts. This study proposes a multi-concept evaluation setting where we control the correlations between textual concepts, such as sentiment, domain, and tense, and analyze performance under increasing correlations between them. We first evaluate the extent to which featurizers can learn disentangled representations of each concept under increasing correlational strengths. We then investigate whether concepts are sufficiently captured by single features or require multiple dimensions; using k-sparse probes, we find that k often needs to be much greater than 1 for optimal scores. Finally, we perform a causal investigation where we steer multiple features simultaneously and observe whether each concept is independently manipulable. Even under ideal uniform distributions of concepts, we find that unsupervised methods like SAEs struggle to learn disentangled concept representations. We then find that the feature representations we identify correspond to disjoint subspaces in activation space, but also that steering with the top feature for one concept still often affects other concepts; this suggests a fundamental entanglement of concepts in the model's representation space. These findings underscore the importance of compositional and out-of-distribution evaluations in interpretability research.

1 Introduction

Interpretability centers on understanding how and why neural networks behave how they do. This requires understanding the underlying causal variables and mechanisms that produce observed input—output behaviors; this study centers on causal variable discovery methods. To uncover causal variable representations, it is now common to deploy *featurization methods*, such as sparse autoencoders (SAEs; Olshausen & Field, 1997; Bricken et al., 2023; Huben et al., 2024) and sparse probes (Gurnee et al., 2023). These methods aim to disentangle activation vectors (wherein a dimension can have many meanings) into sparser spaces where there is a more one-to-one relationship between dimensions and concepts.

Most feature extraction studies and benchmarks focus on isolating single concepts or behaviors, such as refusal (Arditi et al., 2024) and truthfulness (Marks & Tegmark, 2024). This tells us whether the concept exists in the model, but it does not tell us to what degree the concept representation is **independent** and **disentangled** from others. How often do our feature extraction methods really recover concept representations with high precision? Answers to this question act as a ceiling for how much we can trust our steering methods to induce similar behaviors in novel contexts—i.e., to what degree we have predictive power and control over the model's future behaviors.

This is not a new idea: the fields of causal representation learning (CRL; Schölkopf et al., 2021) and disentangled representation learning (Higgins et al., 2018; Locatello et al., 2019; 2020b) have robust literatures addressing under what circumstances and under what assumptions it is possible to

055

056

057

058

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

079

081 082

083

084

085

880

089

090

091

092

094

096

098

099

100

101

102

103 104

105

106

107

identify the true latent causal variables for a task. These fields focus on learning a task from scratch, whereas the goal of interpretability is to derive a simplified causal model of a large and complex neural network that has already been trained (Geiger et al., 2024). Both lines of work are unified in asking: *in what circumstances is it possible to recover causally efficacious representations?*

Our work generalizes and extends the metrics and evaluation paradigms of CRL to mechanistic interpretability using language models. Specifically, using a probabilistic context-free grammar (PCFG), we generate sentences labeled for multiple concepts. We use this dataset to evaluate empirically successful and popular methods in interpretability, including k-sparse probes (Gurnee et al., 2023) and sparse autoencoders (Olshausen & Field, 1997; Huben et al., 2024). First, building on CRL, we use correlational evidence to understand to what degree neurons, sparse features, and probes recover disentangled representations of ground-truth concepts. Then, going beyond CRL, we develop new metrics that evaluate causal criteria we seek from disentangled representations: 1) independent manipulability: disentangled features should allow us to steer one and only one concept downstream; 2) sparse prediction: features should allow us to accurately predict the presence of a concept, ideally with a single feature (Lachapelle et al., 2023a); 3) disjointness (Zuheng et al., 2024): steering two concepts jointly should be the sum of steering each concept independently.

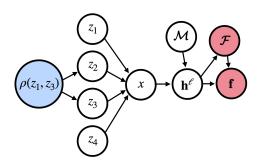


Figure 1: Causal graph of our experimental setup. The values of 4 known ground-truth concepts $\{z_i\}_{i=1}^4$ are used to generate an example x. We train a featurizer $\mathcal F$ to generate vectors $\mathbf f$ given activation vectors $\mathbf h^\ell$ from the output of layer ℓ of language model $\mathcal M$. When training $\mathcal F$ on examples with increasing correlations between pairs of concepts $\rho(z_i,z_j)$, we observe whether $\mathcal F$ learns the true latents or the correlational confound (as measured by the correlation between latents in $\mathbf f$ and the presence of the true variable z_i).

Our contributions include the following:

- Formal evaluation metrics inspired by causal representation learning. These criteria can help
 measure success in recovering precise and disentangled concept representations in interpretability
 studies.
- Experiments comparing the efficacy of common featurization methods on the proposed metrics. Top-K SAEs perform well at disentangling concepts, but do not approach the performance of supervised featurizers.
- A demonstration that even methods with strong sparsity priors still require multiple dimensions to recover known concepts.
- Causal evidence that existing methods often optimize disjointness, but not independence. That is, they succeed in recovering non-overlapping representations, but do often affect multiple unrelated concepts downstream.

2 EXPERIMENTAL SETUP

Our goal is to stress-test featurization methods by creating a dataset labeled with known concepts, but where concepts can be correlated to varying degrees. Figure 1 depicts the causal model for our experiments. We vary the correlations between concept-value pairs in the training dataset $\mathcal D$ used to train a featurizer $\mathcal F$ while holding the language model $\mathcal M$ fixed. $\mathcal F$ is trained to generate a vector $\mathbf f$ of features given activations $\mathbf h^\ell$ from layer ℓ of language model $\mathcal M$. The feature vector $\mathbf f$ should ideally encode one concept per dimension.

We fix a target correlation between two concept values—for example, positive sentiment and the science domain—and introduce an unobserved common cause (the blue node in Figure 1) to create the desired correlation. This creates a confounding variable that acts as the parent of both correlated concepts in the data generation process (DGP). Under varying correlational conditions, we observe to what extent \mathcal{F} can identify the true concepts \mathcal{Z} .

Models and featurizers. We focus primarily on unsupervised methods such as sparse autoencoders (SAEs), due to their popularity in recent unsupervised interpretability research (Costa et al., 2025; Huben et al., 2024; Mueller et al., 2025a; Marks et al., 2025). We formally define each SAE architecture we test in Appendix A. To assess how much information about the target concepts is lost relative to a supervised method, we compare to k-sparse probes, which are allowed to have non-zero weights to $\leq k$ dimensions of their inputs. Following Gurnee et al. (2023), we first train linear probes with L_1 regularization and take the top k weights to find the top k most influential neurons; then, we train logistic regression probes trained with L_2 regularization on those top k neurons.

We focus on two models: Pythia-70M (Biderman et al., 2023) and Gemma-2-2B (Team et al., 2024). We choose these because there exist publicly available SAEs trained on large natural language corpora, including the ReLU SAEs of Marks et al. (2025) and the GemmaScope SAEs (Lieberum et al., 2024).

Recent work has demonstrated the importance of the featurizer's inductive bias, especially when deploying unsupervised featurizers (Hindupur et al., 2025; Costa et al., 2025). We therefore compare SAEs that make varying geometric assumptions: ReLU SAEs (Bricken et al., 2023) assume linear separability, Top-K SAEs (Gao et al., 2025) assume angular separability, and SpADE SAEs (Costa et al., 2025) make weaker assumptions that allow for more heterogeneous concept geometries; we refer readers to Appendix A for details.

Data. Using a probabilistic context-free grammar (PCFG), we generate a training dataset \mathcal{D} containing 382,884 sentences and test dataset \mathcal{T} consisting of 1,007 sentences, where each sentence is labeled for 4 concepts $z_i \in \mathcal{Z}$: voice, tense, sentiment, and domain. In our datasets, voice (active, passive) and tense (present, past) are binary. Sentiment (positive, neutral, negative) is multinomial and ordinal, while domain (news, science, fantasy, other) is multinomial with no inherent ordering. Categorical variables will be treated as one-hot vectors of binary values—e.g., $z_i = [v_{i,0}, v_{i,1}, v_{i,2}]$ for sentiment, where $v_{i,0} = 1$ when sentiment is negative and $v_{i,0} = 0$ otherwise.

3 EVALUATING DISENTANGLEMENT

3.1 Concept identification

A key desideratum of featurizers is the ability to identify the ground-truth concepts despite potential spurious correlations between them.¹ To assess to what degree this property holds for popular featurizers, we design an identifiability evaluation.

To evaluate the ability of a featurizer to recover these concepts, we employ the **mean correlation coefficient** (MCC) metric (Hyvarinen & Morioka, 2016) common in the causal representation learning literature (Hyvarinen et al., 2019; Khemakhem et al., 2020b;a; Wendong et al., 2023; von Kügelgen et al., 2021; 2023; von Kügelgen, 2024; Reizinger et al., 2024a; 2023b;a; Gresele et al., 2021).

A featurizer consists of an encoder $\mathcal{F}:\mathbb{R}^{|\mathbf{h}|}\to\mathbb{R}^{|\mathbf{f}|}$ and optionally a decoder $\mathcal{F}^{-1}:\mathbb{R}^{|\mathbf{f}|}\to\mathbb{R}^{|\mathbf{h}|}$. The encoder \mathcal{F} maps hidden representation vector \mathbf{h}^ℓ at layer ℓ to features \mathbf{f} (where typically, $|\mathbf{f}|>|\mathbf{h}|$). Given a set of ground-truth concepts $\{z_1,\ldots,z_n\}$ that generate an input example \mathbf{x} where each concept $z_j\in\mathbb{Z}$, then $\forall i\in[1,\ldots,n]$, we compute $\hat{\mathbf{f}}_j=\arg\max_i|\rho_{\mathcal{D}}(f_i,z_j)|$, where f_i is the activation of feature \mathbf{f}_i . Intuitively, $\hat{\mathbf{f}}_j$ is the feature whose activation correlates most with the value of z_j on some training dataset \mathcal{D} . Given test set \mathcal{T} where concepts are uniformly distributed w.r.t. each other (i.e., no built-in correlations), we use $\rho_{\mathcal{T}}(\hat{f}_j,z_j)$ as a measure of how well the featurizer linearly identifies concept z_j . After locating the best features $\{\hat{\mathbf{f}}_j\}_{j=1}^n$ for each concept, we compute the MCC as the mean of their correlations with their respective concepts on \mathcal{T} . In other words:

$$MCC = \frac{1}{n} \sum_{j=1}^{n} \rho_{\mathcal{T}}(\hat{f}_j, z_j). \tag{1}$$

¹We cannot expect a model, supervised or unsupervised, to be able to disentangle two concepts if they are *completely* correlated in the data (Wiedemer et al., 2023) without making any assumptions. However, given at least a couple examples where two concepts do not covary, it is possible in theory to recover independent representations of these concepts.

²Note that this is not a literal inversion. The decoder is typically learned such that the reconstruction error is minimized, but information is lost when reconstructing **h** using the featurizer.

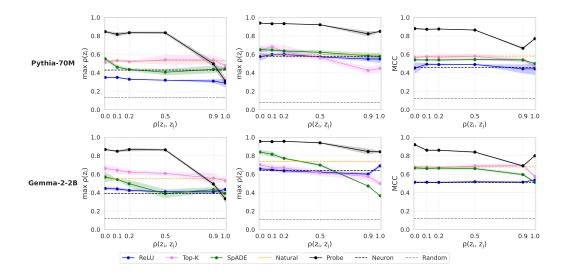


Figure 2: Maximum correlation coefficient for domain=science (left), sentiment=positive (middle), and MCC (right) under varying correlational conditions. Shaded regions represent 1 std. dev. across 3 training seeds. Ideal performance looks like a flat line at a high MCC. Probes (supervised featurizers, in black) perform best. Top-K SAEs perform best among unsupervised featurizers. SAEs trained on large-scale natural data (Natural) perform similarly to our best SAEs trained on CFG-generated data.

The MCC is measured using one-dimensional features, but multinomial concepts may not be one-dimensional in \mathbf{f} or \mathbf{h}^ℓ (Engels et al., 2025). Thus, to create a fairer evaluation, we compute the MCC over binarized concepts. That is, given a variable $z_i \in \mathbb{Z}$ with V_i possible values, we create a new binary variable $v_{i,x} \in \mathbb{B}$ for each value x corresponding to whether $z_i = v_{i,x}$. When computing the MCC, we first average the correlation coefficients for all $v_{i,x} \in V_i$ before taking the macroaverage across concepts.

A high MCC is achievable in theory only if we make the following assumption:

Assumption 1: Linear sufficiency. For each ground-truth concept z_k , there exists a linear invertible transformation T such that $z_k = T\mathbf{h}^{\ell}$ where \mathbf{h}^{ℓ} are the representations of the model \mathcal{M} .

To validate this assumption, we train linear probes for each binary concept and observe whether each probe obtains high accuracy on the concept it was trained to detect, *but also* obtains random-chance accuracy on all other concepts. Our probes satisfy these criteria and thus empirically support Assumption 1; see Figure 6 (Appendix B).

Baselines and skylines. We compare against a randomly initialized SAE (*Random*), the neurons from the residual stream whose correlations correlate most with each concept (*Neuron*, equivalent to the identity featurizer $\mathbf{f} = \mathbf{h}^{\ell}$), and publicly available SAEs trained on natural language data (*Marks* (Marks et al., 2025) and *GemmaScope* (Lieberum et al., 2024) for Pythia-70M and Gemma-2-2B, respectively).

To establish a supervised skyline (Probe), we train logistic regression probes using the binarized concept labels. We treat the probe's logit as the feature activation f_j when computing the correlation, and take the average correlation across concept-specific binary probes to compute the MCC.

Hypothesis. The ideal result is a high MCC that remains constant as the correlation between ground-truth concepts increases in the training data. We expect unsupervised featurizers, such as SAEs, to perform worse than supervised featurizers. We also expect SAEs trained on our dataset to be better able to isolate the ground-truth concepts compared to the *Natural* baselines; this is because the number of varying concepts is lower, which should make these concepts easier to isolate.

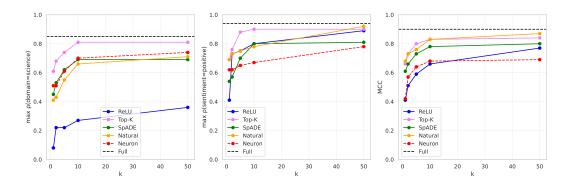


Figure 3: Correlation coefficients between probe logits and concept labels for domain=science (left), sentiment=positive (middle), and MCC (right). Results for Gemma-2-2B shown here; results for Pythia-70M are in Appendix D. We vary the number of dimensions k that the probe is allowed to have non-zero weights from. k-sparse probes trained on SAEs begin to converge around 10 dimensions for Top-K, SpADE, and Natural, and recover most of the performance of a non-sparse probe that is allowed to use the entire residual vector (Full). k-sparse probes trained on the residual stream (Neuron) require more dimensions to converge, as expected.

Results. Figure 2 shows the MCC for Pythia-70M and Gemma-2-2B for the domain and sentiment concepts as they become more correlated in the training dataset. We find that probes significantly outperform SAEs, as expected. The margin between probes and SAEs is significant; thus, if one knows *a priori* what concepts one wishes to find, then one should use supervised methods. This agrees with recommendations from Wu et al. (2025) and Mueller et al. (2025b).

The gap between SAE architectures is significant and consistent across models. Top-K in particular performs well. Our SAEs trained on synthetically generated data achieve comparable performance to SAEs trained on a much larger natural language corpus (the *Natural* SAEs in Figure 2); the best-performing methods outperform them, as hypothesized, but most methods achieve comparable or lower performance. Thus, in practice, one may not need to worry about curating concept-specific data as long as one's dataset is sufficiently large.

When do correlations between concepts start to impede concept identification? The answer depends on the method: probes and SpADE (Costa et al., 2025) maintain relatively consistent MCCs up to correlations of 0.5 between concept pairs in the training data. Beyond this, performance begins to degrade. For Top-K, MCC remains more consistent until we reach correlations of 1.0. In theory, it is always possible to disentangle concepts given at least 2 examples where those concepts do not covary. In practice, however, correlations over 0.5 cause most methods to degrade—including supervised methods. We recommend that future interpretablity studies devote effort to investigating potential correlates of the concept of focus to ensure that other concepts are not being included in learned or derived concept representations.

3.2 IS ONE DIMENSION SUFFICIENT?

In SAE-based interpretability studies, it is common to steer with a single feature, regardless of how many features receive high attributions for a given task. This corresponds to the following assumption:

Assumption 2: One feature dimension is sufficient for concept detection and control. Given binary concept z_i and feature vector \mathbf{f} , one dimension \mathbf{f}_i of \mathbf{f} is sufficient to represent and control z_i in \mathcal{M} .

To evaluate the extent to which this assumption holds in practice, we train k-sparse probes (as operationalized in Gurnee et al. (2023)) on featurized representations f. k-sparse probes are linear probes that may have non-zero weights from up to k dimensions of the representations they are trained on. Lachapelle et al. (2023a) establish a connection between disentanglement and sparse prediction: they prove that disentanglement leads to optimal loss using sparse predictors. Further, as

features become more entangled, we need to reduce sparsity regularization to maintain accuracy; this theoretical finding further motivates the following experiment.

Hypothesis. More dimensions yield monotonically increasing expressive power. Thus, performance should be non-decreasing as k increases. We care primarily about when increasing k begins to yield diminishing improvements in the MCC. Representations obtained with strong sparsity constraints, like SAEs, should reach this saturation point at smaller k than representations with no such constraints, such as residual vectors.

Results. We display the (M)CC of k-sparse probes trained on feature vectors \mathbf{f} in Figure 3. Top-K SAEs achieve the best trade-off between MCC and sparsity at all k; they also approach the MCC of training a normal probe on the full activation vector at the residual stream. ReLU SAEs do not begin saturating even at 10–50 features, whereas all other SAEs do. Top-K achieves better concept recovery at the same k as the residual neuron baseline, whereas ReLU SAEs do not.

These results suggest that SAEs do confer sparsity benefits compared to the original activation space of \mathcal{M} , but also that one-dimensionality assumptions may often be insufficient—even when the concepts are relatively simple.

4 EVALUATING COUNTERFACTUAL INDEPENDENCE AND DISJOINTNESS

4.1 Steering as a causal independence test

The identifiability evaluation above (§3.1) acts as a correlational measure of whether features identify known concepts. However, it does not necessarily provide *causal* evidence that we can independently manipulate concepts using the learned features; an SAE that achieves high MCC may not necessarily yield features that induce changes in the target concept under targeted manipulation. Thus, to measure causal efficacy, we employ steering as an independence test of the mechanisms between the features. This can be seen as testing the Independent Causal Mechanism principle prevalent in the causality literature (Pearl, 2009; Peters et al., 2018), which holds that different causal mechanisms neither influence nor inform each other.

To locate the steering feature, we could select the feature whose correlation is highest with the label, as in §3.1. However, Arad et al. (2025) has found that the features that detect the input concept (the top correlated features in our case) and the features that control the output concept are distinct. Thus, for steering experiments, we use gradient attributions (Simonyan et al., 2014) to locate the feature that should be steered. We would like features that increase the probability of some concept value $v_{i,x}$; as a proxy, we can fold the featurizer into the forward pass of the model (following Marks et al., 2025), take the logit $\Pi(\mathbf{h}^L)$ of a binary probe Π trained on the final layer L of \mathcal{M} to detect a concept value $v_{i,x}$, backpropagate from this logit to obtain its gradient with respect to a feature $\frac{\partial \Pi(\mathbf{h}^L)}{\partial f_i}$, and multiply each feature's gradient by its activation to obtain the gradient attribution $\frac{\partial \Pi(\mathbf{h}^L)}{\partial f_i} \cdot f_i$. We take the feature with the maximum average attribution across examples.

Steering of the activations of layer ℓ with the best feature $\hat{\mathbf{f}}_j$ is performed using steering function $\tilde{\mathbf{h}}^{\ell}(\mathbf{f}_i) \leftarrow \Phi(\mathbf{h}^{\ell}, \mathcal{F}, i, \alpha)$, where Φ is defined as follows:

$$\Phi(\mathbf{h}^{\ell}, \mathcal{F}, i, \alpha) = \mathcal{F}^{-1}\Big(\mathcal{F}(\mathbf{h}^{\ell})|\mathsf{do}(\mathbf{f}_i = \alpha \cdot \mathsf{max}(f_i))\Big) + \epsilon \tag{2}$$

where α controls the strength of the steering operation, $\mathcal{F}(\mathbf{h})$ corresponds to the featurized activations, and the do-operation denotes a feature intervention where feature i is set to α times its maximum on training dataset \mathcal{D}^4 $\epsilon = \mathbf{h} - \mathcal{F}^{-1}(\mathcal{F}(\mathbf{h}))$ is the reconstruction error without interventions. We set α to 5, but try different values in §4.2.

³Intuitively, this is a first-order Taylor approximation of the effect of changing feature activation f_i to 0 on $\Pi(\mathbf{h}^L)$.

⁴This is equivalent to adding the difference between the steered reconstruction and original reconstruction to the activation.

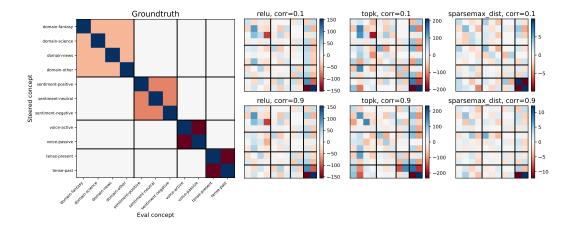


Figure 4: The effect of steering a given concept (row) on the log-odds of another (column), as measured by a probe. Results for Pythia-70M shown here; see Appendix E for Gemma-2-2B. If concept representations are causally independent, we expect a heatmap that resembles the ground-truth: Δ LOGODDS should be high on the diagonal, negative for within-concept pairs, and close to 0.0 for across-concept pairs. All SAEs demonstrate the expected diagonals, but also significant across-concept effects, indicating non-independence. Increasing correlations in the training data, even up to 0.9, do not significantly change the trends.

For all concept pairs $\{(z_i, z_j) : i, j \in [n]\}$, we steer with z_i and plot $\Delta \text{LogOdds}$ of z_j .⁵ We steer with an SAE trained on the middle layer of \mathcal{M} and then quantify $\Delta \text{LogOdds}(z_j)$ as the change in the logit of a multinomial concept probe.⁶ To validate that the concepts can be disentangled in the model, and to validate that probe logits are good proxies for concept presence, we show heatmaps of probe accuracies in Figure 7 (Appendix B). We observe that each concept probe obtains high performance on its concept's test set, and achieves random-chance performance on all other concepts. The supports the validity of the following results.

Hypothesis. If two concepts are independent, then we expect no cross-concept effects—i.e., if two features $\hat{\mathbf{f}}_i$ and $\hat{\mathbf{f}}_{j\neq i}$ correspond to independent concepts z_i and z_j , then steering z_i should not change $p(z_j)$. Note that within-concept effects are expected: for $\hat{\mathbf{f}}_i$ and $\hat{\mathbf{f}}_j$ such that i and j are really two values of the same concept z_i (e.g., positive sentiment and negative sentiment), then positive steering with one feature should necessarily decrease the probability of the other.

Results. We observe (Figure 4) that for each SAE architecture, the expected diagonal trend is present, indicating that steering is increasing the log-odds of the target concept as expected. However, in even the best architectures, steering leads to measurable impacts on many unrelated concepts, indicating widespread non-independence. This underscores the importance of both multi-concept evaluations *and* counterfactual interventions in evaluating concept representations.

4.2 DISJOINTNESS

Steering with one concept can provide causal evidence as to how disentangled two concepts are. Now, inspired by Zuheng et al. (2024), we ask whether these concept representations are **disjoint**—that is, whether they affect non-overlapping subspaces. This is non-equivalent to independence: even if two features correspond to non-overlapping subspaces (i.e., are disjoint), they could still produce non-zero effects on unrelated concepts (i.e., be entangled with other concepts). Disjointness implies that we can predict the effect of pairs of steering operations on z_i from individual steering operations,

 $^{^5\}Delta$ LOGODDS is equivalent to the logit difference.

 $^{^6}$ These are architecturally similar to the probes used in §3.1, but trained on the *final* layer of $\mathcal M$ instead of the middle layer. We use the final layer because it acts as a better proxy for the model's likely output behavior, as opposed to the model's inner representation of the input concepts. We use multinomial probes because they make the change in probabilities for within-concept pairs sum to 1.

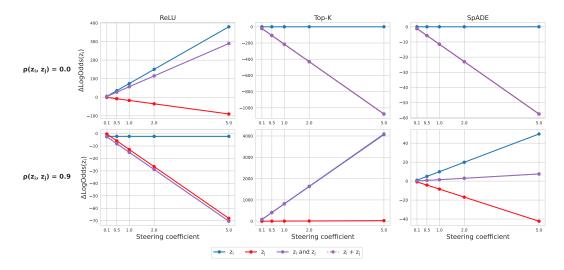


Figure 5: $\Delta \text{LogOdds}(z_i)$ under various steering coefficients α for steering feature $\hat{\mathbf{f}}_i$, and $\hat{\mathbf{f}}_j$ for a different concept z_j . Results for Pythia-70M shown here; results for Gemma-2-2B are in Appendix E. Steering only $\hat{\mathbf{f}}_i$ should increase $\text{LogOdds}(z_i)$. If $\hat{\mathbf{f}}_i$ and $\hat{\mathbf{f}}_j$ are independent, steering $\hat{\mathbf{f}}_j$ should not affect $\text{LogOdds}(z_i)$. If $\hat{\mathbf{f}}_i$ and $\hat{\mathbf{f}}_j$ are disjoint, we expect $\text{LogOdds}(z_i)$ when steering both to be equal to the sum of steering both in isolation. Here, we observe that representations of z_i and z_j are often *not* independent, which would be indicated as a flat red line. However z_i and z_j are always disjoint, as indicated by the dotted purple line and solid purple line completely overlapping such that the dotted line is not visible.

even if individual steering operations affect multiple concepts. Studying disjointness is important because its presence gives us predictive power over model behavior, even in unseen or potentially out-of-distribution scenarios. See Figure 12 for illustrations and a direct contrast of independence and disjointness. Formally,

$$p(z_i|\tilde{\mathbf{h}}^{\ell}(\hat{\mathbf{f}}_i,\hat{\mathbf{f}}_j)) - p(z_i|\mathbf{h}^{\ell}) = \left(p(z_i|\tilde{\mathbf{h}}^{\ell}(\hat{\mathbf{f}}_i)) - p(z_i|\mathbf{h}^{\ell})\right) + \left(p(z_i|\tilde{\mathbf{h}}^{\ell}(\hat{\mathbf{f}}_j)) - p(z_i|\mathbf{h}^{\ell})\right).$$
(3)

That is, the effect on $p(z_i)$ of steering both $\hat{\mathbf{f}}_i$ and $\hat{\mathbf{f}}_j$ should be equivalent to the sum of steering only $\hat{\mathbf{f}}_i$ and $\hat{\mathbf{f}}_j$ in isolation. In practice, we show LOGODDS rather than probabilities; this unbounded metric is more likely to be additive at especially high and low probabilities.

Hypothesis. Under low correlations, we expect that concepts will be disjoint, such that the effect of steering the top features for z_i and z_j on $\Delta \text{LOGODDS}(z_i)$ will be additive, regardless of their (non-)independence. Under higher correlations, we expect less disjoint representations and more non-linearly predictable interaction terms between pairs of steering operations.

Results. In Figure 5, we observe that the effect of steering with two concepts simultaneously is almost exactly equivalent to summing the impact of steering with both concepts separately. This suggests no interaction terms.

This in combination with the non-independence results of §4.1 suggests that each SAE feature is operating on a separate subspace, but also that steering with a concept representation can still affect representations of other concepts.

5 RELATED WORK

Featurization in interpretability. In interpretability, *featurization* refers to techniques that allow one to map from less interpretable and denser model representations—typically *neurons*—to more interpretable (and often sparser) representations—what are often called *features*. This has produced supervised techniques such as sparse probing (Gurnee et al., 2023), unsupervised techniques such as

sparse autoencoders (SAEs; Olshausen & Field, 1997; Bricken et al., 2023; Huben et al., 2024), and non-parametric methods such as deriving steering vectors (Subramani et al., 2022) via difference-in-means (Marks & Tegmark, 2024).

How can one evaluate the quality of a feature? Recent work has proposed standardized evaluations based on known concepts (Mueller et al., 2025b; Huang et al., 2024). These allow one to assess whether a concept discovery method discovers a concept with high recall. However, it leaves precision unexplored: how well do these concept representations disentangle the concept from others? Evaluating this requires multi-concept evaluations.

Causal representation learning. Causal representation learning (CRL; Schölkopf et al., 2021) assumes that high-dimensional observations, such as text, are generated from low-dimensional latent factors, whose relationships to other latent factors are encoded in a causal graph. Then, CRL proposes latent variable models of such observations that are **identifiable**, meaning that the recovered features (and possibly a graph over them) are related to the true factors up to permutation and element-wise transformations. Since such unsupervised learning is not identifiable without further assumptions (Hyvärinen & Pajunen, 1999; Darmois, 1951; Locatello et al., 2019), CRL methods rely on non-iid data or constraints on the decoding function (Moran et al., 2022; Gresele et al., 2021; Lachapelle et al., 2023b; Brady et al., 2025; Reizinger et al., 2023b). For example, CRL has developed identifiable models using data from sparse interventions Ahuja et al. (2022b); Zhang et al. (2023); Buchholz et al. (2023); von Kügelgen et al. (2023), contrastive pairs of samples (Ahuja et al., 2022a; Locatello et al., 2020a; Gresele et al., 2019; Brehmer et al., 2022), data from multiple environments (Ahuja et al., 2023; Layne et al., 2025; Khemakhem et al., 2020a), and temporal data with sparse or intervened mechanisms (Lachapelle et al., 2021; Lippe et al., 2023; 2022). We go further, however, and test the causal implications of disentangled features: target concept steering, accuracy with sparse probes and disjoint steering effects.

Compositional generalization. Closely related to disentanglement and the notion of disjoint effects is the ability of models to compose concepts in novel ways, called compositional generalization. Compositional generalization has a long history in the NLP literature (Ahuja & Mansouri, 2024; Han & Padó, 2024; Ramesh et al., 2024; Lake & Baroni, 2023; Nogueira et al., 2021; Dziri et al., 2023; Saparov et al., 2023; Mészáros et al., 2024; Reizinger et al., 2024b; Ujváry et al., 2025), but tends to focus on the reuse of syntactic chunks or lexemes. Some recent CRL studies investigate compositional latents; they tend to study simplified formal languages, such as regular languages or Dyck (bracketing) languages (Deletang et al., 2022; Mészáros et al., 2024; Reizinger et al., 2024b; Ujváry et al., 2025).

6 DISCUSSION AND CONCLUSIONS

Each of our experiments has revealed insufficiencies in single-concept evaluations. One may achieve far above random-chance performance under correlational evaluation methods (§3.1) and improvements in sparsity over the native residual representation space a model (§3.2). Even so, causal evidence reveals that entanglement can still be likely and widespread (§4.1,4.2) even when the aforementioned correlational metrics suggest otherwise.

Despite strong entanglement, concept pairs demonstrated very little in the way of interaction effects (§4.2). Intuitively, this implies that when features achieve the *form* of separation—that is, that the cosine similarity of the subspace on which they act is very low—it does not necessarily imply that their *functional roles* are non-interacting. This suggests that mechanistic interpretability studies aiming to establish the independence of two mechanisms cannot settle for establishing that subspaces or circuits do not overlap; one must directly establish that the functional roles on the final output are independent.

One dimension is not sufficient, even with methods with strong sparsity regularizers. This may imply that the intrinsic dimensionalities of the concepts themselves are greater than one. Given the variance of scientific domains or positive sentiment, this would not necessarily be surprising. It would be interesting for future work to investigate the relationship between causal independence metrics and the intrinsic dimensionality of feature representations—for example, using techniques like those of Engels et al. (2025). Broadly speaking, more work is needed on methods for detecting, characterizing, and steering with multi-dimensional concepts.

REPRODUCIBILITY

To ensure the robustness of our results, we average results across three random seeds and report standard deviations. For all optimization-based procedures, we fix and save these random seeds; these settings will be released alongside our code. We will release all code and data upon deanonymization.

REFERENCES

- Kartik Ahuja and Amin Mansouri. On Provable Length and Compositional Generalization, February 2024. URL http://arxiv.org/abs/2402.04875. arXiv:2402.04875 [cs, stat].
- Kartik Ahuja, Jason Hartford, and Yoshua Bengio. Weakly Supervised Representation Learning with Sparse Perturbations. October 2022a. URL https://openreview.net/forum?id=6ZI4iF_T7t.
- Kartik Ahuja, Yixin Wang, Divyat Mahajan, and Yoshua Bengio. Interventional Causal Representation Learning, September 2022b. URL http://arxiv.org/abs/2209.11924.arXiv:2209.11924 [cs, stat].
- Kartik Ahuja, Amin Mansouri, and Yixin Wang. Multi-Domain Causal Representation Learning via Weak Distributional Invariances, October 2023. URL http://arxiv.org/abs/2310.02854. arXiv:2310.02854 [cs, stat].
- Dana Arad, Aaron Mueller, and Yonatan Belinkov. Saes are good for steering if you select the right features, 2025. URL https://arxiv.org/abs/2505.20063.
- Andy Arditi, Oscar Balcells Obeso, Aaquib Syed, Daniel Paleka, Nina Rimsky, Wes Gurnee, and Neel Nanda. Refusal in language models is mediated by a single direction. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=pH3XAQME6c.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar Van Der Wal. Pythia: a suite for analyzing large language models across training and scaling. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.
- Jack Brady, Julius von Kügelgen, Sebastien Lachapelle, Simon Buchholz, Thomas Kipf, and Wieland Brendel. Interaction asymmetry: A general principle for learning composable abstractions. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Johann Brehmer, Pim de Haan, Phillip Lippe, and Taco Cohen. Weakly supervised causal representation learning, October 2022. URL http://arxiv.org/abs/2203.16437. arXiv:2203.16437 [cs, stat].
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. https://transformer-circuits.pub/2023/monosemantic-features/index.html.
- Simon Buchholz, Goutham Rajendran, Elan Rosenfeld, Bryon Aragam, Bernhard Schölkopf, and Pradeep Ravikumar. Learning Linear Causal Representations from Interventions under General Nonlinear Mixing, June 2023. URL http://arxiv.org/abs/2306.02235.arXiv:2306.02235 [cs, math, stat].
- Valérie Costa, Thomas Fel, Ekdeep Singh Lubana, Bahareh Tolooshams, and Demba Ba. From flat to hierarchical: Extracting sparse representations with matching pursuit, 2025. URL https://arxiv.org/abs/2506.03093.

George Darmois. Analyse des liaisons de probabilité. In *Proc. Int. Stat. Conferences 1947*, pp. 231, 1951

Gregoire Deletang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, and Pedro A. Ortega. Neural Networks and the Chomsky Hierarchy. September 2022. URL https://openreview.net/forum?id=WbxHAzkeOcn.

Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena D. Hwang, Soumya Sanyal, Sean Welleck, Xiang Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. Faith and fate: Limits of transformers on compositionality, 2023.

 Joshua Engels, Eric J Michaud, Isaac Liao, Wes Gurnee, and Max Tegmark. Not all language model features are one-dimensionally linear. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=d63a4AM4hb.

Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=tcsZt9ZNKD.

Atticus Geiger, Duligur Ibeling, Amir Zur, Maheep Chaudhary, Sonakshi Chauhan, Jing Huang, Aryaman Arora, Zhengxuan Wu, Noah Goodman, Christopher Potts, and Thomas Icard. Causal Abstraction: A Theoretical Foundation for Mechanistic Interpretability, August 2024. URL http://arxiv.org/abs/2301.04709.arXiv:2301.04709 [cs].

Luigi Gresele, Paul K. Rubenstein, Arash Mehrjou, Francesco Locatello, and Bernhard Schölkopf. The Incomplete Rosetta Stone Problem: Identifiability Results for Multi-View Nonlinear ICA. arXiv:1905.06642 [cs, stat], August 2019. URL http://arxiv.org/abs/1905.06642. arXiv: 1905.06642.

Luigi Gresele, Julius von Kügelgen, Vincent Stimper, Bernhard Schölkopf, and Michel Besserve. Independent mechanism analysis, a new concept? *arXiv:2106.05200 [cs, stat]*, June 2021. URL http://arxiv.org/abs/2106.05200. arXiv: 2106.05200.

Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. Finding neurons in a haystack: Case studies with sparse probing. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=JYs1R9IMJr.

Sungjun Han and Sebastian Padó. Towards Understanding the Relationship between In-context Learning and Compositional Generalization, March 2024. URL http://arxiv.org/abs/2403.11834. arXiv:2403.11834 [cs].

Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and Alexander Lerchner. Towards a Definition of Disentangled Representations. *arXiv:1812.02230 [cs, stat]*, December 2018. URL http://arxiv.org/abs/1812.02230. arXiv: 1812.02230.

Sai Sumedh R. Hindupur, Ekdeep Singh Lubana, Thomas Fel, and Demba Ba. Projecting assumptions: The duality between sparse autoencoders and concept geometry, 2025. URL https://arxiv.org/abs/2503.01822.

Jing Huang, Zhengxuan Wu, Christopher Potts, Mor Geva, and Atticus Geiger. RAVEL: Evaluating interpretability methods on disentangling language model representations. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8669–8687, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.470. URL https://aclanthology.org/2024.acl-long.470/.

Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=F76bwRSLeK.

- Aapo Hyvarinen and Hiroshi Morioka. Unsupervised Feature Extraction by Time-Contrastive Learning and Nonlinear ICA. arXiv:1605.06336 [cs, stat], May 2016. URL http://arxiv.org/abs/1605.06336. arXiv: 1605.06336.
 - Aapo Hyvarinen, Hiroaki Sasaki, and Richard E. Turner. Nonlinear ICA Using Auxiliary Variables and Generalized Contrastive Learning. *arXiv:1805.08651 [cs, stat]*, February 2019. URL http://arxiv.org/abs/1805.08651. arXiv: 1805.08651.
 - Aapo Hyvärinen and Petteri Pajunen. Nonlinear independent component analysis: Existence and uniqueness results. *Neural Networks*, 12(3):429–439, April 1999. ISSN 0893-6080. doi: 10.1016/S0893-6080(98)00140-3. URL https://www.sciencedirect.com/science/article/pii/S0893608098001403.
 - Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational Autoencoders and Nonlinear ICA: A Unifying Framework. In *International Conference on Artificial Intelligence and Statistics*, pp. 2207–2217. PMLR, June 2020a. URL http://proceedings.mlr.press/v108/khemakhem20a.html. ISSN: 2640-3498.
 - Ilyes Khemakhem, Ricardo Pio Monti, Diederik P. Kingma, and Aapo Hyvärinen. ICE-BeeM: Identifiable Conditional Energy-Based Deep Models Based on Nonlinear ICA. arXiv:2002.11537 [cs, stat], October 2020b. URL http://arxiv.org/abs/2002.11537. arXiv: 2002.11537.
 - Sebastien Lachapelle, Tristan Deleu, Divyat Mahajan, Ioannis Mitliagkas, Yoshua Bengio, Simon Lacoste-Julien, and Quentin Bertrand. Synergies between Disentanglement and Sparsity: Generalization and Identifiability in Multi-Task Learning. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 18171–18206. PMLR, July 2023a. URL https://proceedings.mlr.press/v202/lachapelle23a.html. ISSN: 2640-3498.
 - Sébastien Lachapelle, Pau Rodríguez López, Yash Sharma, Katie Everett, Rémi Le Priol, Alexandre Lacoste, and Simon Lacoste-Julien. Disentanglement via Mechanism Sparsity Regularization: A New Principle for Nonlinear ICA. arXiv:2107.10098 [cs, stat], November 2021. URL http://arxiv.org/abs/2107.10098. arXiv: 2107.10098.
 - Sébastien Lachapelle, Divyat Mahajan, Ioannis Mitliagkas, and Simon Lacoste-Julien. Additive Decoders for Latent Variables Identification and Cartesian-Product Extrapolation, July 2023b. URL http://arxiv.org/abs/2307.02598. arXiv:2307.02598 [cs, stat].
 - Brenden M. Lake and Marco Baroni. Human-like systematic generalization through a meta-learning neural network. *Nature*, 623(7985):115–121, November 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06668-3. URL https://www.nature.com/articles/s41586-023-06668-3. Publisher: Nature Publishing Group.
 - Elliot Layne, Jason Hartford, Sébastien Lachapelle, Mathieu Blanchette, and Dhanya Sridhar. Sparsity regularization via tree-structured environments for disentangled representations. *Transactions on Machine Learning Research*, 2025.
 - Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, Janos Kramar, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. In Yonatan Belinkov, Najoung Kim, Jaap Jumelet, Hosein Mohebbi, Aaron Mueller, and Hanjie Chen (eds.), *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pp. 278–300, Miami, Florida, US, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.blackboxnlp-1. 19. URL https://aclanthology.org/2024.blackboxnlp-1.19/.
 - Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M. Asano, Taco Cohen, and Efstratios Gavves. CITRIS: Causal Identifiability from Temporal Intervened Sequences, June 2022. URL http://arxiv.org/abs/2202.03169. Number: arXiv:2202.03169 arXiv:2202.03169 [cs, stat].
 - Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M. Asano, Taco Cohen, and Efstratios Gavves. BISCUIT: Causal Representation Learning from Binary Interactions, June 2023. URL http://arxiv.org/abs/2306.09643.arXiv:2306.09643 [cs, stat].

Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations. In *International Conference on Machine Learning*, pp. 4114–4124. PMLR, May 2019. URL http://proceedings.mlr.press/v97/locatello19a.html. ISSN: 2640-3498.

- Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-Supervised Disentanglement Without Compromises. *arXiv*:2002.02886 [cs, stat], October 2020a. URL http://arxiv.org/abs/2002.02886. arXiv: 2002.02886.
- Francesco Locatello, Michael Tschannen, Stefan Bauer, Gunnar Rätsch, Bernhard Schölkopf, and Olivier Bachem. Disentangling Factors of Variation Using Few Labels. *arXiv:1905.01258 [cs, stat]*, February 2020b. URL http://arxiv.org/abs/1905.01258. arXiv: 1905.01258.
- Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=aajyHYjjsk.
- Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=I4e82CIDxv.
- Gemma E. Moran, Dhanya Sridhar, Yixin Wang, and David M. Blei. Identifiable Deep Generative Models via Sparse Decoding. Technical Report arXiv:2110.10804, arXiv, February 2022. URL http://arxiv.org/abs/2110.10804. arXiv:2110.10804 [cs, stat] type: article.
- Aaron Mueller, Jannik Brinkmann, Millicent Li, Samuel Marks, Koyena Pal, Nikhil Prakash, Can Rager, Aruna Sankaranarayanan, Arnab Sen Sharma, Jiuding Sun, Eric Todd, David Bau, and Yonatan Belinkov. The quest for the right mediator: Surveying mechanistic interpretability through the lens of causal mediation analysis, 2025a. URL https://arxiv.org/abs/2408.01416.
- Aaron Mueller, Atticus Geiger, Sarah Wiegreffe, Dana Arad, Iván Arcuschin, Adam Belfki, Yik Siu Chan, Jaden Fried Fiotto-Kaufman, Tal Haklay, Michael Hanna, Jing Huang, Rohan Gupta, Yaniv Nikankin, Hadas Orgad, Nikhil Prakash, Anja Reusch, Aruna Sankaranarayanan, Shun Shao, Alessandro Stolfo, Martin Tutek, Amir Zur, David Bau, and Yonatan Belinkov. MIB: A mechanistic interpretability benchmark. In Forty-second International Conference on Machine Learning, 2025b. URL https://openreview.net/forum?id=sSrOwve6vb.
- Anna Mészáros, Szilvia Ujváry, Wieland Brendel, Patrik Reizinger, and Ferenc Huszár. Rule extrapolation in language models: A study of compositional generalization on ood prompts, 2024. URL https://arxiv.org/abs/2409.13728.
- Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. Investigating the limitations of transformers with simple arithmetic tasks, 2021.
- Bruno A. Olshausen and David J. Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision Research*, 37(23):3311–3325, 1997. ISSN 0042-6989. doi: https://doi.org/10.1016/S0042-6989(97)00169-7. URL https://www.sciencedirect.com/science/article/pii/S0042698997001697.
- Judea Pearl. Causal inference in statistics: An overview. Statistics Surveys, 3 (none), January 2009. ISSN 1935-7516. doi: 10.1214/09-SS057. URL https://projecteuclid.org/journals/statistics-surveys/volume-3/issue-none/Causal-inference-in-statistics-An-overview/10.1214/09-SS057.full.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Elements of causal inference: foundations and learning algorithms. *Journal of Statistical Computation and Simulation*, 88(16): 3248–3248, November 2018. ISSN 0094-9655, 1563-5163. doi: 10.1080/00949655.2018. 1505197. URL https://www.tandfonline.com/doi/full/10.1080/00949655. 2018.1505197.

Rahul Ramesh, Ekdeep Singh Lubana, Mikail Khona, Robert P. Dick, and Hidenori Tanaka. Compositional Capabilities of Autoregressive Transformers: A Study on Synthetic, Interpretable Tasks, February 2024. URL http://arxiv.org/abs/2311.12997. arXiv:2311.12997 [cs].

Patrik Reizinger, Luigi Gresele, Jack Brady, Julius von Kügelgen, Dominik Zietlow, Bernhard Schölkopf, Georg Martius, Wieland Brendel, and Michel Besserve. Embrace the Gap: VAEs Perform Independent Mechanism Analysis, January 2023a. URL http://arxiv.org/abs/2206.02416. arXiv:2206.02416 [cs, stat].

Patrik Reizinger, Yash Sharma, Matthias Bethge, Bernhard Schölkopf, Ferenc Huszár, and Wieland Brendel. Jacobian-based Causal Discovery with Nonlinear ICA. *Transactions on Machine Learning Research*, April 2023b. ISSN 2835-8856. URL https://openreview.net/forum?id=2Yo9xqR6Ab.

Patrik Reizinger, Siyuan Guo, Ferenc Huszár, Bernhard Schölkopf, and Wieland Brendel. Identifiable Exchangeable Mechanisms for Causal Structure and Representation Learning. October 2024a. URL https://openreview.net/forum?id=k03mB41vyM.

Patrik Reizinger, Szilvia Ujváry, Anna Mészáros, Anna Kerekes, Wieland Brendel, and Ferenc Huszár. Understanding Ilms requires more than statistical generalization, 2024b.

Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Seyed Mehran Kazemi, Najoung Kim, and He He. Testing the general deductive reasoning capacity of large language models using ood examples, 2023.

Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Towards Causal Representation Learning. *arXiv:2102.11107 [cs]*, February 2021. URL http://arxiv.org/abs/2102.11107. arXiv: 2102.11107 version: 1.

Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps, 2014. URL https://arxiv.org/abs/1312.6034.

Nishant Subramani, Nivedita Suresh, and Matthew Peters. Extracting latent steering vectors from pretrained language models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), Findings of the Association for Computational Linguistics: ACL 2022, pp. 566–581, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.48. URL https://aclanthology.org/2022.findings-acl.48/.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024. URL https://arxiv.org/abs/2403.08295.

Szilvia Ujváry, Anna Mészáros, Wieland Brendel, Patrik Reizinger, and Ferenc Huszár. Transcending bayesian inference: Transformers extrapolate rules compositionally under model misspecification. In 7th Symposium on Advances in Approximate Bayesian Inference – Workshop Track, 2025. URL https://openreview.net/forum?id=0DRAstwh5Y.

Julius von Kügelgen. Identifiable Causal Representation Learning: Unsupervised, Multi-View, and Multi-Environment. March 2024. URL https://www.repository.cam.ac.uk/handle/1810/365627.

Julius von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, and Francesco Locatello. Self-Supervised Learning with Data Augmentations Provably Isolates Content from Style, June 2021. URL http://arxiv.org/abs/2106.04619. arXiv: 2106.04619.

Julius von Kügelgen, Michel Besserve, Liang Wendong, Luigi Gresele, Armin Kekić, Elias Bareinboim, David M. Blei, and Bernhard Schölkopf. Nonparametric Identifiability of Causal Representations from Unknown Interventions, October 2023. URL http://arxiv.org/abs/2306.00542.arXiv:2306.00542 [cs, stat].

Liang Wendong, Armin Kekić, Julius von Kügelgen, Simon Buchholz, Michel Besserve, Luigi Gresele, and Bernhard Schölkopf. Causal Component Analysis, October 2023. URL http://arxiv.org/abs/2305.17225. arXiv:2305.17225 [cs, stat].

Thaddäus Wiedemer, Prasanna Mayilvahanan, Matthias Bethge, and Wieland Brendel. Compositional Generalization from First Principles, July 2023. URL http://arxiv.org/abs/2307.05596. arXiv:2307.05596 [cs, stat].

Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christopher D Manning, and Christopher Potts. Axbench: Steering LLMs? even simple baselines outperform sparse autoencoders. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=K2CckZjNy0.

Jiaqi Zhang, Chandler Squires, Kristjan Greenewald, Akash Srivastava, Karthikeyan Shanmugam, and Caroline Uhler. Identifiability Guarantees for Causal Disentanglement from Soft Interventions, July 2023. URL http://arxiv.org/abs/2307.06250. arXiv:2307.06250 [cs, math, stat].

Zuheng, Xu, Moksh Jain, Ali Denton, Shawn Whitfield, Aniket Didolkar, Berton Earnshaw, and Jason Hartford. Automated discovery of pairwise interactions from unstructured data, 2024. URL https://arxiv.org/abs/2409.07594.

A FURTHER DETAILS ON METHODS

Here, we define sparse autoencoders and describe the differences between the architectures we study.

Sparse autoencoders. The conceptually simplest architecture we deploy is the ReLU sparse autoencoder (Huben et al., 2024; Bricken et al., 2023), which learns a mapping from $\mathbf{x} = \mathbf{h}^{\ell}$ to a learned sparse feature vector \mathbf{f} , and then reconstructs the activations $\hat{\mathbf{x}}$ given \mathbf{f} . More formally:

$$\mathbf{f} = \text{ReLU}(W_{\text{enc}}\mathbf{x} + \mathbf{b}_{\text{enc}}) \tag{4}$$

$$\hat{\mathbf{x}} = W_{\text{dec}}(\mathbf{f} - \mathbf{b}_{\text{enc}}) + \mathbf{b}_{\text{dec}} \tag{5}$$

ReLU SAEs minimize $\mathcal{L} = MSE(\mathbf{x}, \hat{\mathbf{x}}) + \lambda ||\mathbf{f}||_1$.

Top-K SAEs (Gao et al., 2025) are similar to ReLU SAEs, but they strictly retain the top k activations per sample and zero out all others:

$$\mathbf{f} = \text{top-}k(W_{\text{enc}}\mathbf{x} + \mathbf{b}_{\text{enc}}) \tag{6}$$

Sparsemax distance encoders (SpADE) can capture nonlinearly separable and heterogeneous features; we refer readers to Hindupur et al. (2025) for details. In formal terms:

$$\mathbf{f} = \operatorname{Sparsemax}(-\lambda d(\mathbf{x}, W)) \tag{7}$$

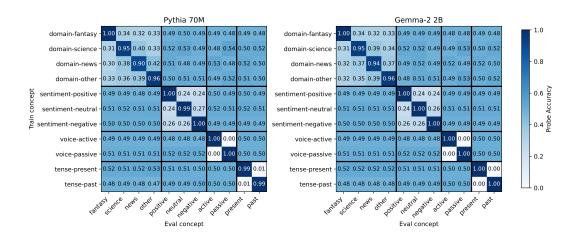


Figure 6: Accuracy of binary probes (rows) on all concept value classification tasks (columns). We expect high values on the diagonals, below random chance for within-concept value pairs, and random chance for across-concept value pairs.

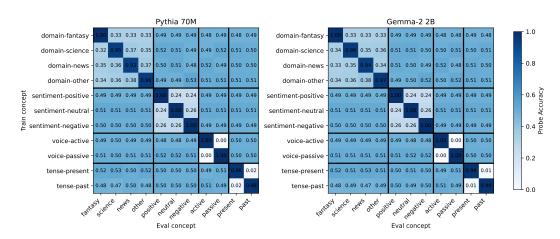


Figure 7: Accuracy of multinomial probes on all concept value classification tasks (columns). We expect high values on the diagonals, below random chance for within-concept value pairs, and random chance for across-concept value pairs.

where $d(\mathbf{x}, W)_i = \|\mathbf{x} - W_i\|_2^2$. Hindupur et al. (2025) show that this architecture can capture more irregular concept geometries, whereas ReLU SAEs assume linear separability, and Top-K SAEs assume angular separability.

B PROBE ACCURACIES

Here, we present the accuracies of each probe we use in our disentanglement experiments and evaluations. We present these as heatmaps to verify whether each probe learn an independent representation of its target concept; if it does, we expect high scores along the diagonal, lower-than-random scores for within-concept pairs, ⁷ and random-chance scores for across-concept pairs.

Binary linear probes trained on the middle layers of Pythia-70M and Gemma-2-2B (Figure 6) achieve near-perfect accuracies on their respective concepts, and achieve the expected random accuracies

⁷We expect lower-than-random scores for within-concept pairs because a classifier trained on an alternative value of a concept should be strictly worse than a random probe, as the target label will be *negatively* correlated with the target concept.

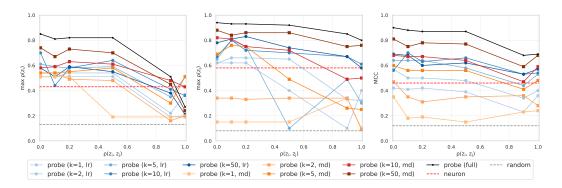


Figure 8: MCC for the two most performant sparse probing methods from Gurnee et al. (2023) at various k. The LR method achieves higher MCC at lower k, but MD overtakes LR at higher k.

on all other concepts. This empirically supports Assumption 1, and supports the idea that the MCC ceiling should be high (§3.1).

In §4.1 and §4.2, we instead use multinomial linear probes trained on the final layers of Pythia-70M and Gemma-2-2B. We find (Figure 7) that these probes also achieve the expected high accuracies on the target concepts, below-random-chance accuracies on within-concept pairs, and random-chance accuracies on across-concept pairs. This validates that the non-independence we observe in our steering experiments are not due to the probes, but rather are more likely due to the featurization methods that we use to steer.

C SPARSE PROBING

Here, we replicate the setup of Gurnee et al. (2023) in our cross-concept correlation setting. We aim to assess which k-sparse probing methods are more robust to cross-concept correlations at multiple k. We focus on the two most performant methods from Gurnee et al. (2023): max mean difference (MD), and logistic regression (LR). MD works by computing the average difference in activations between positive and negative samples, and taking the k neurons whose mean activation difference is greatest. LR works by first training a logistic regression probe with L_1 regularization on the full activation vector, and then taking the top k according to the weights of the probe.

We observe (Figure 8) that the logistic regression (LR) method of selecting neurons is more effective at lower k. Between k=5 and k=10, MD generally overtakes LR in performance. As we are more concerned with low-dimensional concept recovery, we focus on LR in the feature dimensionality experiment (§3.2).

D FURTHER DISENTANGLEMENT RESULTS

Here, we present correlation coefficients and MCCs for k-sparse probes trained with varying k on SAEs for Pythia-70M. As with Gemma-2-2B, correlation coefficients tend to converge at around 10 dimensions; this suggests that the one-dimensionality assumption may not often hold in practice, even for much smaller models. Note also that the neuron baseline is far more performant for Pythia than Gemma; perhaps this is because k=10 represents a far greater proportion of the dimensions of \mathbf{h}^ℓ for Pythia than Gemma. Other trends are largely consistent with Figure 3.

E FURTHER STEERING RESULTS

Here, we present steering heatmaps for Gemma-2-2B (Figure 10). Features appear less independent than for Pythia-70M, as indicated by more significant across-concept $\Delta LogOdds$ for many concept pairs. That said, the expected diagonal trend is still present. This is further evidence that SAE features do not often correspond to causally independent concept representations.

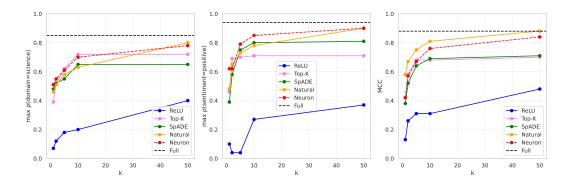


Figure 9: Correlation coefficients between probe logits and concept labels for domain=science (left), sentiment=positive (middle), and MCC (right). Results for Pythia-70M. We vary the number of dimensions k that the probe is allowed to have non-zero weights from. As with Gemma-2-2B, correlation coefficients tend to converge at around 10 dimensions. However, the neuron baseline is far more performant; perhaps this is because k = 10 represents a far greater proportion of the dimensions of h^{ℓ} for Pythia than Gemma. Other trends are largely consistent with Figure 3.

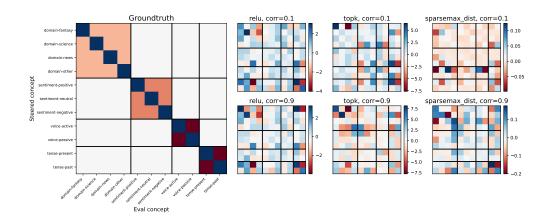


Figure 10: The effect of steering a given concept (row) on the logit of another (column), as measured by a probe. Results for Gemma-2-2B. If concept representations are causally independent, we expect a heatmap that resembles the ground-truth: $\Delta LogOddent Ddent Ddent$

We also present multi-feature steering results (Figure 11). As for the results for Pythia-70M, we observe that features are often entirely disjoint while not being independent. Here, we observe some distinction between the predicted and actual $\Delta LogOddet Ddetector Ddetector$

F FURTHER DETAILS ON METRICS

To disambiguate the conceptual distinction between independence and disjointness, we present diagrams in Figure 12. Intuitively, disjointness implies that two feature representations exist in non-overlapping subspaces of the model representations, and thus that the effect of steering of both can be predicted from the result of steering either in isolation. Independence implies that steering with one concept would not affect how the model uses other concepts. Refer to §4.2 for details.

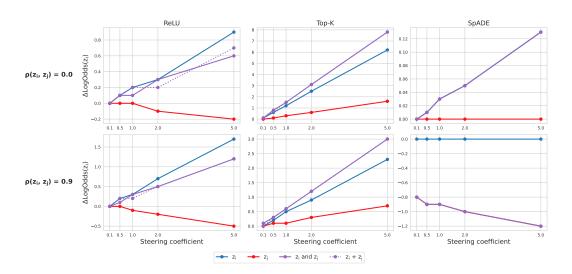


Figure 11: Δ LOGODDS(z_i) under various steering coefficients α for steering feature $\hat{\mathbf{f}}_i$, and $\hat{\mathbf{f}}_j$ for a different concept z_j . Results for Gemma-2-2B. Representations of z_i and z_j are more often independent here than for Pythia-70M, as indicated by the flat red line. z_i and z_j are typically mostly disjoint, as indicated by the dotted purple line and solid purple lines almost (but not completely) overlapping.

$$\begin{split} p(z_i|\mathbf{h}^{\ell}) & \xrightarrow{\Phi(\mathbf{h}^{\ell},\mathcal{F},i,\alpha)} p(z_i|\tilde{\mathbf{h}}^{\ell}(\hat{\mathbf{f}}_i)) \\ & \downarrow_{\Phi(\mathbf{h}^{\ell},\mathcal{F},j,\beta)} & \downarrow_{\Phi(\mathbf{h}^{\ell},\mathcal{F},i,\beta)} & p(z_j|\mathbf{h}^{\ell}) \xrightarrow{\Phi(\mathbf{h}^{\ell},\mathcal{F},i,\alpha)} p(z_j|\tilde{\mathbf{h}}^{\ell}(\hat{\mathbf{f}}_i)) \\ & p(z_i|\tilde{\mathbf{h}}^{\ell}(\hat{\mathbf{f}}_j)) \xrightarrow{\Phi(\mathbf{h}^{\ell},\mathcal{F},i,\alpha)} p(z_i|\tilde{\mathbf{h}}^{\ell}(\hat{\mathbf{f}}_i,\hat{\mathbf{f}}_j)) \end{split}$$

Figure 12: The difference between feature disjointness and independence: (Left) Two concepts z_i and z_j with feature representations $\hat{\mathbf{f}}_i$ and $\hat{\mathbf{f}}_j$, respectively, are disjoint if the left diagram commutes. (**Right**) If they are independent then there is no commutative relationship, as steering with $\hat{\mathbf{f}}_i$ should not affect $p(z_j)$. Intuitively, disjointness implies that two feature representations exist in non-overlapping subspaces of the model representations, and thus that the effect of steering of both can be predicted from the result of steering either in isolation. Independence implies that steering with one concept would not affect how the model uses other concepts. Refer to §4.2 for formulae and empirical details.

G LLM USAGE

The authors used large language models primarily as a polishing tool during writing. LLMs were not used in a significant capacity for writing experimental code nor for research ideation, although we acknowledge that libraries on which our code was based may have used LLMs.