
Competence-Based Analysis of Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite the recent successes of large, pretrained neural language models (LLMs),
2 comparatively little is known about the representations of linguistic structure they
3 learn during pretraining, which can lead to unexpected behaviors in response to
4 prompt variation or distribution shift. To better understand these models and be-
5 haviors, we introduce a general model analysis framework to study LLMs with
6 respect to their representation and use of human-interpretable linguistic properties.
7 Our framework, CALM (Competence-based Analysis of Language Models), is
8 designed to investigate LLM competence in the context of specific tasks by inter-
9 venting on models' internal representations of different linguistic properties using
10 causal probing, and measuring models' alignment under these interventions with a
11 given ground-truth causal model of the task. We also develop a new approach for
12 performing causal probing interventions using gradient-based adversarial attacks,
13 which can target a broader range of properties and representations than prior tech-
14 niques. Finally, we carry out a case study of CALM using these interventions to
15 analyze and compare LLM competence across a variety of lexical inference tasks,
16 showing that CALM can be used to explain and predict behaviors across these
17 tasks.

18 1 Introduction

19 The rise of large, pretrained neural language models (LLMs) has led to rapid progress in a wide
20 variety of natural language processing tasks [7, 12, 17]. However, these models can also be quite
21 sensitive to minor changes in input prompts [19, 42, 41] and fail to generalize outside their training or
22 fine-tuning distribution [68, 73, 60]. Understanding the means by which these models can perform as
23 well as they do while exhibiting such limitations is a key question in the science of LLM interpretation
24 and analysis [3], and is likely necessary in enabling robust, trustworthy, and socially-responsible
25 LLM-enabled applications [59, 32, 78, 3].

26 We approach this question in terms of *competence*, drawing on the traditional competence-
27 performance distinction in linguistic theory (see Section 2) to motivate the study of LLMs in terms
28 of their underlying representation of language. We define LLM competence in the context a given
29 linguistic task as the alignment between the ground-truth causal structure of the task and the LLM's
30 latent representation of the task's structure, measured by intervening on the LLM's representation of
31 task-causal or non-causal properties and observing how its behavior changes in response. While such
32 representations are not directly observable, we take inspiration from *causal probing*, which damages
33 LLMs' latent representations of linguistic properties using causal interventions to study how these
34 representations contribute to their behavior [18, 31]. We introduce a general framework, CALM (for
35 Competence-based Analysis of Language Models), to study the competence of LLMs using causal
36 probing and define the first quantitative measure of LLM competence.

37 While CALM can be instantiated using a variety of existing causal probing interventions (e.g., [52,
38 50, 51, 58, 2]), we develop a new methodology, gradient-based interventions (GBIs), to intervene on
39 LLM representations using gradient-based adversarial attacks against structural probes, extending
40 causal probing to arbitrarily-encoded representations of relational properties and thereby enabling the
Submitted to NeurIPS 2024 Workshop on Causality and Large Models (CaLM). Do not distribute.

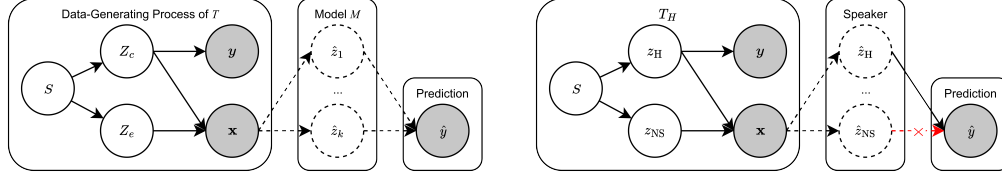


Figure 1: **Structural causal models (SCMs)** of task \mathcal{T} 's data-generating process and how it may be performed by model M (left), or hypernym prediction task (\mathcal{T}_H) and how it is performed by a competent English speaker (right). Shaded and white nodes denote observed and unobserved variables, respectively. In CALM, the goal is to determine which representations $Z_j = z_j$ are causally implicated in M 's predictions \hat{y} .

41 investigation of new questions in language model interpretation. Finally, we carry out a case study of
 42 CALM using two well-studied LLMs by implementing interventions as GBIs in order to measure
 43 and compare these LLMs' competence across 14 lexical inference tasks, showing that CALM can
 44 indeed explain and predict important patterns in behavior across these tasks by distinguishing between
 45 models' use of causal and spurious properties.

46 2 Competence-based Analysis of Language Models

47 **Linguistic Competence** Linguistic competence is generally understood as the ability to utilize one's
 48 knowledge of a language in producing and understanding utterances in that language, and is typically
 49 defined in contrast with linguistic performance, which is speakers' actual use of their language
 50 in practice, considered independently of the underlying knowledge that supports it [40]. Given a
 51 linguistic task, we may understand competence in terms of the underlying linguistic knowledge
 52 that one draws upon to perform the task. If fluent human speakers rely on (implicit or explicit)
 53 knowledge of the same set of linguistic properties to perform a given task, then we may understand
 54 their performance of this task as being causally determined by these properties, and invariant to other
 55 properties. (See Appendix A.1 for further discussion of linguistic competence.)

56 While the study of human competence has a rich history in linguistics [10, 36, 43, 57, 39, 40], there is
 57 currently no generally accepted framework for studying LLM competence [38, 45]. In order to make
 58 the study of competence tractable in the context of LLMs, we introduce the CALM (Competence-
 59 based Analysis of Language Models) framework, which describes an LLM's competence with respect
 60 to a given linguistic task in terms of its latent representation of the causal structure of the task.

61 **Task Structure** Formally, given supervised task $\mathcal{T} \sim P(\mathcal{X}, \mathcal{Y})$ where the goal is to correctly predict
 62 $y \in \mathcal{Y}$ given $\mathbf{x} \in \mathcal{X}$, and a collection of latent properties $\mathbf{Z} = \{Z_j\}_{j=1}^m$ that are (potentially) involved
 63 in generating \mathbf{x} , we formulate the causal structure of \mathcal{T} in terms of the data-generating process

$$\mathbf{x} \sim \Pr(\mathbf{x}|\mathbf{Z}_c, \mathbf{Z}_e), \quad \mathbf{y} \sim P(\mathbf{y}|\mathbf{Z}_c) \quad (1)$$

64 where \mathbf{Z} may be decomposed into $\mathbf{Z} = \mathbf{Z}_c \cup \mathbf{Z}_e, \mathbf{Z}_c \cap \mathbf{Z}_e = \emptyset$, where \mathbf{Z}_c contains all properties
 65 that causally determine \mathbf{y} , and \mathbf{Z}_e are the remaining properties that may be involved in generating \mathbf{x}
 66 (cf. [28]). However, there may be an unobserved confounder S that produces spurious correlations
 67 between \mathbf{y} and \mathbf{Z}_e , which, if leveraged by language model M in the course of predicting \hat{y} , can lead
 68 to unexpected failures on \mathcal{T} when the spurious association is broken [46]. The structural causal model
 69 (SCM)¹ of this data-generating process is visualized on the left side of both diagrams in Figure 1.

70 **Internal Representation** Our main concern is measuring how attributable an LLM M 's behavior
 71 in a given task \mathcal{T} is to its representation of various properties $\mathbf{Z} = \{Z_1, \dots, Z_m\}$, and how these
 72 properties correspond to the causal structure of the task. If M respects the data-generating process of
 73 \mathcal{T} , then its behavior should be attributable only to causal properties $Z \in \mathbf{Z}_c$ (and not to environmental
 74 properties $Z \in \mathbf{Z}_e$), in which case we say that M is *competent* with respect to \mathcal{T} . We study model
 75 M 's use of each property $Z_j \in \mathbf{Z}$ by performing causal interventions $\text{do}(Z_j)$ on its representation of
 76 Z_j in the course of performing task \mathcal{T} , and measure the impact that these interventions have on its
 77 predictions.

78 **Measuring Competence** We evaluate the competence of M with respect to task $\mathcal{T} \sim P(\mathcal{X}, \mathcal{Y})$ by
 79 measuring its causal alignment with a *competence graph* $\mathcal{G}_{\mathcal{T}}$, which we define as a structural causal

¹An SCM is a directed acyclic graph where each node represents a variable and directed edges indicate causal dependencies (see [5] for an introduction to SCMs).

80 model (SCM) of \mathcal{T} with nodes corresponding to each latent variables $Z_j \in \mathbf{Z}$ and an additional node
 81 for outputs $\mathbf{y} \in \mathcal{Y}$ and directed edges denoting causal dependencies between these variables. That is,
 82 the set of causal properties \mathbf{Z}_c defined by $\mathcal{G}_{\mathcal{T}}$ is the set of all properties $Z_j \in \mathbf{Z}$ such that there is an
 83 edge or path from Z_j to \mathbf{y} . To determine the extent to which M 's behavior is correctly explained by
 84 the causal dependencies (and lack thereof) in $\mathcal{G}_{\mathcal{T}}$, we measure their consistency under interventions
 85 $\text{do}(\mathbf{z})$, where setting $\mathbf{z} = \{z_j\}_{j=1}^m \sim \text{val}(\mathbf{Z})$ is a combination of values $Z_j = z_j \in \text{val}(Z_j)$ taken by
 86 each corresponding latent variable $Z_j \in \mathbf{Z}$.

87 The alignment of M with $\mathcal{G}_{\mathcal{T}}$ is measured in terms of the similarity S of their predictions under
 88 interventions $\text{do}(\mathbf{z})$ given input $\mathbf{x} \sim P(\mathcal{X})$, and can be computed using a given similarity metric
 89 $S : \mathcal{Y}, \mathcal{Y} \rightarrow [0, 1]$ (e.g., equality, n-gram overlap, cosine similarity, etc.) depending on the SCM $\mathcal{G}_{\mathcal{T}}$
 90 and output space \mathcal{Y} . That is, we define $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$ as M 's competence with respect to task \mathcal{T} as a
 91 function of its alignment with corresponding task SCM $\mathcal{G}_{\mathcal{T}}$ under interventions $\text{do}(\mathbf{z})$ measured by
 92 similarity metric S , as follows:

$$\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}}) = \mathbb{E}_{\mathbf{x}, \mathbf{z} \sim P(\mathcal{X}, \text{val}(\mathbf{Z}))} S(M(\mathbf{x}|\text{do}(\mathbf{z})), \mathcal{G}_{\mathcal{T}}(\mathbf{x}|\text{do}(\mathbf{z}))) \quad (2)$$

93 This $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$ metric (bounded by $[0, 1]$) is an adaptation of the Interchange Intervention Accuracy
 94 (IIA) metric [23, 22] to the context of causal probing, where instance-level interventions are replaced
 95 with concept-level interventions enabled by the gradient-based intervention methodology we introduce
 96 in Section 3. (See Appendix C.1 for a detailed comparison of our competence metric with IIA.)

97 **Causal Probing** A key technical challenge in implementing CALM (and causal probing more
 98 generally) is designing an algorithm to perform causal interventions $\text{do}(Z)$ that maximally damage
 99 the representation of a property Z while otherwise minimally damaging representations of other
 100 properties Z' [50]. For example, *amnesic probing* [18] uses the INLP algorithm [52] to produce
 101 interventions g_Z that remove all information that is linearly predictive of property Z from a set
 102 of embedding representations \mathbf{H} . However, when such information removal methods are used to
 103 remove representation of these properties in early LLM layers, models are often able to “recover”
 104 this representation in later layers [18, 50], which is likely due to models encoding these properties
 105 nonlinearly. Beyond recoverability, linear information removal methods like INLP also cannot
 106 account for relational properties between multiple input embeddings (see Appendix A.1). Thus, it
 107 is important to develop interventions that do not require restrictive assumptions about the structure
 108 of LLMs’ representations such as linearity [67], a problem which we aim to solve in the following
 109 section.

110 3 Gradient-based Interventions

111 Our goal in developing gradient-based interventions (GBIs) as a causal probing technique is to enable
 112 interventions over arbitrarily-encoded LLM representations. GBIs allow users to flexibly specify
 113 the class of representations they wish to target, expanding the scope of causal probing to arbitrarily-
 114 encoded properties. We take inspiration from Kos, Fischer, and Song [29], who developed a technique
 115 to perturb latent representations using gradient-based adversarial attacks.² They begin by training
 116 probe $g_Z : \mathbf{h} \mapsto z$ to predict image class $z \in Z$ from latent representations $\mathbf{h} = f_{\text{enc}}(\mathbf{x})$ of images \mathbf{x} ,
 117 where f_{enc} is the encoder of a VAE-GAN [30] trained on an unsupervised image reconstruction task
 118 (i.e., $f_{\text{dec}}(f_{\text{enc}}(\mathbf{x})) = \hat{\mathbf{x}} \approx \mathbf{x}$, for decoder f_{dec} and reconstructed image $\hat{\mathbf{x}}$ approximating \mathbf{x}). Next,
 119 gradient-based attacks like FGSM [24] and PGD [37] are performed against g_Z in order to minimally
 120 manipulate \mathbf{h} such that it resembles encoded representations of target image class $Z = z'$ (where
 121 $z' \neq z$, the original image class), yielding perturbed representation \mathbf{h}' . Finally, \mathbf{h} and \mathbf{h}' are each fed
 122 into the VAE decoder to reconstruct corresponding output images $\hat{\mathbf{x}}$ and $\hat{\mathbf{x}}'$ (respectively), where $\hat{\mathbf{x}}$
 123 resembles input image class $Z = z$ and $\hat{\mathbf{x}}'$ resembles target class $Z = z'$.

124 We reformulate this approach in the context of causal probing as visualized in Figure 2, treating layers
 125 $L = 1, \dots, l$ as the encoder and layers $L = l + 1, \dots, |L|$ (composed with language modeling head
 126 f_{LM}) as the decoder, allowing us to target representations of property Z across embeddings \mathbf{h}_i^l of
 127 token $x_i \in \mathbf{x}$ in layer l . We train g_Z to predict Z from a set of such \mathbf{h}_i^l , then attack g_Z using FGSM
 128 and PGD to intervene on \mathbf{h}_i^l (representing the original value $Z = z$), producing $\mathbf{h}_i^{l'}$ (representing
 129 the counterfactual value $Z = z'$). Finally, we replace \mathbf{h}_i^l with $\mathbf{h}_i^{l'}$ in the LLMs’ forward pass from

²Notably, Tucker, Qian, and Levy [66] developed a similar methodology without explicit use of such attacks (see Section 7).

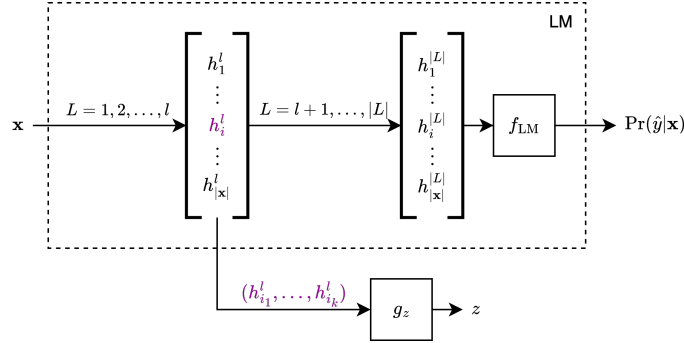


Figure 2: **Gradient-based Interventions.** Input tokens $\mathbf{x} = (x_1, \dots, x_{|\mathbf{x}|})$ are passed through layers $L = 1, \dots, l$, where embedding h_i^l (encoding the value $Z = z$) is extracted from layer l and given to g_z as input. Next, the embedding is modified by gradient-based attacks on g_z to encode the counterfactual value $Z = z'$, then fed back into subsequent layers $L = l + 1, \dots, |L|$ and language modeling head f_{LM} to obtain the intervened predictions $M(\mathbf{x} | \text{do}(Z = z'))$.

130 layers $L = l + 1, \dots, |L|$, simulating the intervention $\text{do}(Z = z')$, and observe the impact on its word
 131 predictions $M(\mathbf{x} | \text{do}(Z = z'))$.

132 There are several key benefits associated with GBIs relative to existing causal probing interventions
 133 (e.g., they can be applied to any differentiable probe), as well as some important limitations (e.g.,
 134 lack of theoretical guarantees), as we discuss in detail in Appendix A.2.

135 4 Experiments

136 In this work, we begin by examining BERT [16] and RoBERTa [34],³ two language models which
 137 have been extensively studied in the context of probing [54, 52, 35, 18, 31]. Our primary goal in the
 138 following experiments is to develop and test an experimental implementation of CALM using GBIs
 139 in the context of comparatively small, well-studied models and tasks in order to validate whether
 140 CALM can explain behavioral findings of earlier work in this simplified environment. (We motivate
 141 this choice in greater detail in Appendix B.1.)

142 **Tasks** We use the collection of 14 lexical inference tasks included in the ConceptNet [61] subset of
 143 LAMA [48], each of which are formulated as a collection of cloze prompts [33]. For example, the
 144 LAMA “IsA” task contains $\sim 2\text{K}$ hypernym prompts corresponding to the “IsA” ConceptNet relation
 145 (including, e.g., “A laser is a [MASK] which creates coherent light.”, where the task is to predict that
 146 the [MASK] token should be replaced with “device”, a hypernym of “laser”), with the remaining 13
 147 LAMA ConceptNet tasks corresponding to other lexical relations such as “PartOf”, “HasProperty”,
 148 and “CapableOf”. (See Appendix B.2 for additional details.)

149 Using these task datasets allow us to test how the representation of each relation is used across
 150 all other tasks. In the context of a single task \mathcal{T}_j , intervening on a model’s representation of the
 151 task-causal relation Z_j allows us to measure the extent to which its predictions are attributable to its
 152 representation of the causal property $\mathbf{Z}_c = \{Z_j\}$ (where a large impact indicates competence). On the
 153 other hand, intervening on the representations of the other 13 lexical relations $Z_k \in \mathbf{Z}_e$ allows us (in
 154 the aggregate) to measure how much the model is performing task \mathcal{T}_j by leveraging representations
 155 of general, non-causal lexical information (where a large impact indicates incompetence).⁴

156 **Experimentally Measuring Competence** Given LLM M and task \mathcal{T} , measuring the competence
 157 $\mathcal{C}_{\mathcal{T}}(M | \mathcal{G}_{\mathcal{T}})$ of M given $\mathcal{G}_{\mathcal{T}}$ requires us to specify an experimental model $E = (\mathbf{Z}, \mathcal{G}_{\mathcal{T}}, S)$, where \mathbf{Z}
 158 is a set of properties, $\mathcal{G}_{\mathcal{T}}$ is a competence graph for task \mathcal{T} , and S is a scoring function that compares
 159 the predictions of M and $\mathcal{G}_{\mathcal{T}}$ in order to compute the approximated $\hat{\mathcal{C}}_{\mathcal{T}}$. Given that each task \mathcal{T}_i
 160 is defined by a single causal lexical relation Z_i (i.e., $\mathbf{Z}_{c_i} = \{Z_i\}$), we model settings \mathbf{z} as a collection
 161 of values $Z_j = z_j$ taken by each property Z_j in the context of a specific task instance $(\mathbf{x}, \mathbf{y}) \sim \mathcal{T}_i$,
 162 where $Z_j = 1$ if $i = j$ (i.e., where the property Z_j is the causal property for the task \mathcal{T}_i) or $Z_j = 0$
 163 otherwise. That is, for each instance $(\mathbf{x}, \mathbf{y}) \sim \mathcal{T}_i$, the corresponding setting \mathbf{z} is a one-hot vector

³Specifically, BERT-base-uncased and RoBERTa-base [70].

⁴Note that the strictest interpretation of this formulation of competence makes the simplifying assumption that each non-causal property is equally (un)related to the target property, which is not always true; see Appendix B.2.

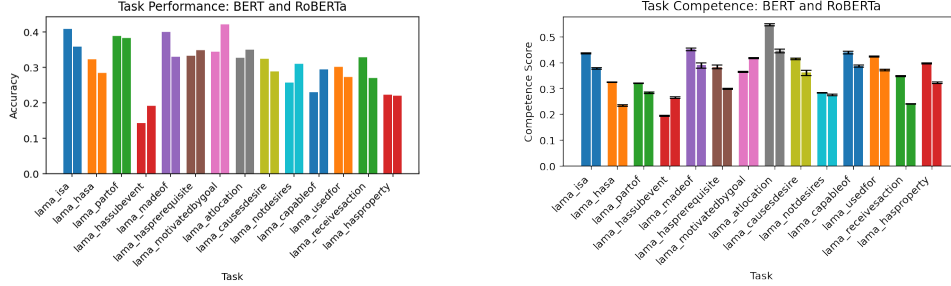


Figure 3: Performance (top) and competence (bottom) of BERT (left bars) and RoBERTa (right bars) for all tasks, using FGSM with $\epsilon = 0.1$. In the competence plot, y-values are the average competence score and error bars are the maximum and minimum competence score, as measured over 10 experimental iterations (each with a different randomly-initialized probe g_Z).

164 whose i -th element $\mathbf{z}_i = 1$. We may specify $\mathcal{G}_{\mathcal{T}_i}$ in a similar manner: for task $\mathcal{T}_i \sim P(\mathcal{X}, \mathcal{Y})$, outputs
165 $\mathbf{y} \in \mathcal{Y}$ are causally dependent on the property Z_i , and invariant to other concepts $Z_j, j \neq i$, meaning
166 that the only direct parent node of \mathbf{y} in $\mathcal{G}_{\mathcal{T}_i}$ is Z_i . Finally, as we are dealing with masked language
167 models whose output space \mathcal{Y} for each task consists only of single tokens in M 's vocabulary V_M , our
168 experimental model can define the scoring function S as the overlap $\text{overlap}(\mathbf{y}_i, \mathbf{y}_j)$ for top- k token
169 predictions $\mathbf{y}_i = \{y_1, \dots, y_k\} \subset V_M$, where $\text{overlap}(\cdot, \cdot)$ is the size of the intersection of each set of
170 predictions divided by the total number of predictions $\text{overlap}(\mathbf{y}_i, \mathbf{y}_j) = \frac{|\mathbf{y}_i \cap \mathbf{y}_j|}{k}$, and $\hat{\mathcal{C}}_{\mathcal{T}_k}$ denotes
171 $\hat{\mathcal{C}}_{\mathcal{T}}$ as measured using the top- k token predictions \mathbf{y}_i . (See Appendix C.2 for additional details on
172 how we compute competence in each experiment.)

173 **Probes** We implement probes g_Z as a 2-layer MLP over each language model's final hidden layer,
174 and train the probe on the task of classifying whether there is a particular relation Z between a
175 final-layer [MASK] token in the context of a cloze prompt and the final-layer object token from
176 the "unmasked" version of the same prompt. All reported figures are the average of 10 runs of our
177 experiment, using different randomly-initialized g_Z each time. (See Appendix B.3 for further details.)

178 **Interventions** We implement GBIs against g_Z using two gradient attack strategies, FGSM [24] and
179 PGD [37]. We bound the magnitude of each intervention as follows: where h is the input to g_Z and
180 h' is the intervened representation following a GBI, $\|h - h'\|_{\infty} \leq \epsilon$. For all experiments reported in
181 our main paper, we use FGSM with $\epsilon = 0.1$. (See Appendix B.4 for more details and PGD results.)

182 5 Results

183 In Figure 3, we visualize the performance and competence of BERT and RoBERTa across the test
184 set of each LAMA ConceptNet task. Performance is measured using (0, 1)-accuracy, competence is
185 measured using the experimental competence metric in Equation (3), and both metrics are averaged
186 across the top- k predictions of each model for $k \in [1, 10]$. That is, for ground truth (\mathbf{x}, y) and $n = 10$,
187 we compute accuracy and competence as follows:

$$\text{acc}(M) = \frac{1}{n} \sum_{k=1}^n \mathbb{1}[y \in \text{top-}k \Pr_M(\hat{y}|\mathbf{x})] \quad \text{and} \quad \hat{\mathcal{C}}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}}) = \frac{1}{n} \sum_{k=1}^n \hat{\mathcal{C}}_{\mathcal{T}_k}(M|\mathcal{G}_{\mathcal{T}})$$

188 To account for stochasticity in initializing and training probes g_Z , scores are also averaged over 10
189 randomized experiments for each target task where the probe is randomly re-initialized each time
190 (resulting in different GBIs).

191 **Performance** While their accuracies on individual tasks vary, BERT and RoBERTa have quite
192 similar aggregate performance: BERT outperforms RoBERTa on just over half (8/14) of the tasks,
193 achieving essentially equivalent performance when averaged across all tasks (0.3099 versus 0.3094).

194 **Competence** Given our experimental model E with $m = 14$ tasks, consider a random baseline
195 language model R whose predictions always change in response to each intervention, making equal
196 use of all properties in each task. R would yield a competence score of $\mathcal{C}(R|\mathcal{G}_{\mathcal{T}}) = \frac{1}{m} \approx 0.0714$
197 for each task. Both BERT and RoBERTa score above this threshold for all tasks, meaning that their
198 competence is consistently greater than that of a model (R) that does not distinguish between causal
199 and environmental properties. However, RoBERTa is consistently less competent than BERT (on

200 12/14 tasks), and also has lower competence scores averaged across all tasks (0.381 vs. 0.334).
201 Furthermore, relative performance and competence are correlated: the Spearman’s Rank correlation
202 coefficient between the average difference in accuracy and average difference in performance is a
203 moderately strong positive correlation $\rho = 0.508$ with significance $p = 0.064$.

204 6 Discussion

205 **Explananda** The performance of BERT and RoBERTa on lexical inference tasks such as hypernym
206 prediction has been shown to be highly variable under small changes to prompts [27, 53, 21, 19]. Our
207 findings offer one possible explanation for such brittle performance: BERT and RoBERTa’s partial
208 competence in hypernym prediction indicates that it should be possible to prompt these models in a
209 way that will yield high performance, but that its reliance on spurious lexical associations may lead it
210 to fail when these correlations are broken – e.g., by substituting singular terms for plurals [53] or
211 paraphrasing a prompt [19].

212 **Future Work** While the simplified experimental context considered in this work is a necessary
213 first step in empirically validating our theoretical CALM framework, competence metric, and
214 GBI methodology, we anticipate a much broader range of future research directions and potential
215 applications for CALM. We elaborate several such directions in Appendix D.

216 7 Related Work

217 **Causal Probing** Most related to our work is amnesic probing [18], as discussed in Section 2.
218 Lasri et al. [31] applied amnesic probing to study the use of grammatical number representations in
219 performing an English verb conjugation prompt task. As this experiment involves intervening on the
220 representation of a property which is causal with respect to the prompt task, it may be understood as an
221 informal instantiation of CALM (albeit without considering environmental properties or measuring
222 competence).

223 **Gradient-based Interventions** Tucker, Qian, and Levy [66] developed a similar approach to our
224 GBI causal probing methodology (as outlined in Section 2) without explicit use of gradient-based
225 adversarial attacks. Their methodology is equivalent to performing a targeted, unconstrained attack
226 using standard gradient descent.⁵ In such attacks, it is standard practice to constrain the magnitude of
227 resulting perturbations [24, 37, 29], which we do here in order to minimize the effect of “collateral
228 damage” done by such attacks (see Section 4 and Appendix B.4); so failing to impose such constraints
229 may result in indiscriminate damage to representations.

230 **Unsupervised Probing** Instead of training supervised probes to predict a pre-determined property
231 of interest (as we do here), an alternative approach is to train *unsupervised* probes such as Sparse
232 Auto-Encoders (SAEs; [62, 74, 14]) to automatically learn an overcomplete basis of features that are
233 useful for sparsely representing embeddings, which can also be used to control models’ use of these
234 learned features [6, 63]. However, as SAEs are unsupervised probes, they yield feature vectors that
235 are not inherently interpretable and must be retroactively interpreted, meaning that the task of creating
236 a supervised probe training dataset (as required for conventional causal probing) is substituted for the
237 task of interpreting learned features [15]. However, given features that can be reliably interpreted as
238 representing task-causal or -environmental features, it is also possible to implement CALM using
239 unsupervised probes like SAEs.

240 8 Conclusion

241 In this work, we introduced CALM, a general analysis framework that enables the study of LLMs’
242 linguistic competence using causal probing, including the first quantitative measure of linguistic
243 competence. We developed the gradient-based intervention (GBI) methodology, a novel approach
244 to causal probing that can target a far greater range of representations than previous techniques,
245 expanding the scope of causal probing to new questions in LLM interpretability and analysis. Finally,
246 we carried out a case study of CALM using GBIs, analyzing BERT and RoBERTa’s competence
247 across a collection of lexical inference tasks, finding that even a simple experimental model is
248 sufficient to explain and predict their behavior across a variety of lexical inference tasks.

⁵I.e., they continue running gradient updates until the targeted probe loss saturates, irrespective of resulting perturbation magnitude.

References

- 249
- 250 [1] Martin Arjovsky et al. “Invariant risk minimization”. In: *arXiv preprint arXiv:1907.02893*
251 (2019).
- 252 [2] Nora Belrose et al. “Leace: Perfect linear concept erasure in closed form”. In: *Advances in*
253 *Neural Information Processing Systems* 36 (2024).
- 254 [3] Leonard Bereska and Efstratios Gavves. “Mechanistic Interpretability for AI Safety—A Review”.
255 In: *arXiv preprint arXiv:2404.14082* (2024).
- 256 [4] BigScience et al. “Bloom: A 176b-parameter open-access multilingual language model”. In:
257 *arXiv preprint arXiv:2211.05100* (2022).
- 258 [5] Stephan Bongers et al. “Foundations of structural causal models with cycles and latent vari-
259 ables”. In: *The Annals of Statistics* 49.5 (2021), pp. 2885–2915.
- 260 [6] Trenton Bricken et al. “Towards Monosemanticity: Decomposing Language Models
261 With Dictionary Learning”. In: *Transformer Circuits Thread* (2023). [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
262 [circuits.pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 263 [7] Tom Brown et al. “Language models are few-shot learners”. In: *Advances in neural information*
264 *processing systems* 33 (2020), pp. 1877–1901.
- 265 [8] Sébastien Bubeck et al. “Convex optimization: Algorithms and complexity”. In: *Foundations*
266 *and Trends® in Machine Learning* 8.3-4 (2015), pp. 231–357.
- 267 [9] Peter Bühlmann. “Invariance, causality and robustness”. In: *Statistical Science* 35.3 (2020),
268 pp. 404–426.
- 269 [10] Noam Chomsky. “Aspects of the Theory of Syntax”. In: *MIT Press* (1965).
- 270 [11] Leshem Choshen et al. “Where to start? Analyzing the potential value of intermediate models”.
271 In: *arXiv preprint arXiv:2211.00107* (2022).
- 272 [12] Aakanksha Chowdhery et al. “Palm: Scaling language modeling with pathways”. In: *arXiv*
273 *preprint arXiv:2204.02311* (2022).
- 274 [13] Arthur Conmy et al. “Towards Automated Circuit Discovery for Mechanistic Interpretability”.
275 In: *Thirty-seventh Conference on Neural Information Processing Systems*. 2023. arXiv: 2304.
276 14997 [cs.LG].
- 277 [14] Hoagy Cunningham et al. “Sparse autoencoders find highly interpretable features in language
278 models”. In: *arXiv preprint arXiv:2309.08600* (2023).
- 279 [15] Adam Davies and Ashkan Khakzar. “The Cognitive Revolution in Interpretability: From
280 Explaining Behavior to Interpreting Representations and Algorithms”. In: *arXiv preprint*
281 *arXiv:2408.05859* (2024).
- 282 [16] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language
283 Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the*
284 *Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*
285 *and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, June
286 2019, pp. 4171–4186. DOI: 10.18653/v1/N19-1423. URL: [https://aclanthology.org/](https://aclanthology.org/N19-1423)
287 [N19-1423](https://aclanthology.org/N19-1423).
- 288 [17] Abhimanyu Dubey et al. “The llama 3 herd of models”. In: *arXiv preprint arXiv:2407.21783*
289 (2024).
- 290 [18] Yanai Elazar et al. “Amnesic Probing: Behavioral Explanation with Amnesic Counterfactuals”.
291 In: *Transactions of the Association for Computational Linguistics* 9 (2021), pp. 160–175. DOI:
292 10.1162/tacl_a_00359. URL: <https://aclanthology.org/2021.tacl-1.10>.
- 293 [19] Yanai Elazar et al. “Measuring and Improving Consistency in Pretrained Language Models”. In:
294 *Transactions of the Association for Computational Linguistics* 9 (Dec. 2021), pp. 1012–1031.
295 ISSN: 2307-387X. DOI: 10.1162/tacl_a_00410. eprint: [https://direct.mit.edu/](https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00410/1975957/tacl_a_00410.pdf)
296 [tacl/article-pdf/doi/10.1162/tacl_a_00410/1975957/tacl_a_00410.pdf](https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00410/1975957/tacl_a_00410.pdf).
297 URL: https://doi.org/10.1162/tacl%5C_a%5C_00410.
- 298 [20] Nelson Elhage et al. “A mathematical framework for transformer circuits”. In: *Transformer*
299 *Circuits Thread* 1 (2021).
- 300 [21] Allyson Ettinger. “What BERT Is Not: Lessons from a New Suite of Psycholinguistic Diagnos-
301 tics for Language Models”. In: *Transactions of the Association for Computational Linguistics*
302 8 (Jan. 2020), pp. 34–48. ISSN: 2307-387X. DOI: 10.1162/tacl_a_00298. eprint: [https://direct](https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00298/1923116/tacl_a_00298.pdf)
303 [.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00298/1923116/](https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00298/1923116/tacl_a_00298.pdf)
304 [tacl_a_00298.pdf](https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00298/1923116/tacl_a_00298.pdf). URL: https://doi.org/10.1162/tacl%5C_a%5C_00298.

-
- 305 [22] Atticus Geiger, Chris Potts, and Thomas Icard. “Causal Abstraction for Faithful Model Inter-
306 pretation”. In: *arXiv preprint arXiv:2301.04709* (2023).
- 307 [23] Atticus Geiger et al. “Inducing Causal Structure for Interpretable Neural Networks”. In:
308 *Proceedings of the 39th International Conference on Machine Learning*. Ed. by Kamalika
309 Chaudhuri et al. Vol. 162. Proceedings of Machine Learning Research. PMLR, 17–23 Jul 2022,
310 pp. 7324–7338. URL: <https://proceedings.mlr.press/v162/geiger22a.html>.
- 311 [24] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. “Explaining and Harnessing Ad-
312 versarial Examples”. In: *International Conference on Learning Representations*. 2015. URL:
313 <http://arxiv.org/abs/1412.6572>.
- 314 [25] Sven Gowal et al. “Scalable verified training for provably robust image classification”. In:
315 *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019, pp. 4842–
316 4851.
- 317 [26] Dirk Groeneveld et al. “OLMo: Accelerating the Science of Language Models”. In: *arXiv*
318 *preprint arXiv:2402.00838* (2024).
- 319 [27] Michael Hanna and David Mareček. “Analyzing BERT’s Knowledge of Hypernymy via
320 Prompting”. In: *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Inter-*
321 *preting Neural Networks for NLP*. Punta Cana, Dominican Republic: Association for Computa-
322 tional Linguistics, Nov. 2021, pp. 275–282. DOI: 10.18653/v1/2021.blackboxnlp-1.20.
323 URL: <https://aclanthology.org/2021.blackboxnlp-1.20>.
- 324 [28] Maximilian Ilse, Jakub M Tomczak, and Patrick Forré. “Selecting data augmentation for
325 simulating interventions”. In: *International Conference on Machine Learning*. PMLR. 2021,
326 pp. 4555–4562.
- 327 [29] Jernej Kos, Ian Fischer, and Dawn Song. “Adversarial examples for generative models”. In:
328 *2018 IEEE Security and Privacy Workshops (SPW)*. IEEE. 2018, pp. 36–42.
- 329 [30] Anders Boesen Lindbo Larsen et al. “Autoencoding beyond pixels using a learned similarity
330 metric”. In: *Proceedings of The 33rd International Conference on Machine Learning*. Ed. by
331 Maria Florina Balcan and Kilian Q. Weinberger. Vol. 48. Proceedings of Machine Learning
332 Research. New York, New York, USA: PMLR, 20–22 Jun 2016, pp. 1558–1566. URL: <https://proceedings.mlr.press/v48/larsen16.html>.
- 333 [31] Karim Lasri et al. “Probing for the Usage of Grammatical Number”. In: *Proceedings of the 60th*
334 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
335 Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 8818–8831. DOI:
336 10.18653/v1/2022.acl-long.603. URL: [https://aclanthology.org/2022.acl-](https://aclanthology.org/2022.acl-long.603)
337 [long.603](https://aclanthology.org/2022.acl-long.603).
- 338 [32] Q Vera Liao and Jennifer Wortman Vaughan. “AI Transparency in the Age of LLMs: A
339 Human-Centered Research Roadmap”. In: *arXiv preprint arXiv:2306.01941* (2023).
- 340 [33] Pengfei Liu et al. “Pre-Train, Prompt, and Predict: A Systematic Survey of Prompting Methods
341 in Natural Language Processing”. In: *ACM Comput. Surv.* 55.9 (Jan. 2023). ISSN: 0360-0300.
342 DOI: 10.1145/3560815. URL: <https://doi.org/10.1145/3560815>.
- 343 [34] Yinhan Liu et al. “Roberta: A robustly optimized bert pretraining approach”. In: *arXiv preprint*
344 *arXiv:1907.11692* (2019).
- 345 [35] Zeyu Liu et al. “Probing Across Time: What Does RoBERTa Know and When?” In: *Findings*
346 *of the Association for Computational Linguistics: EMNLP 2021*. Punta Cana, Dominican
347 Republic: Association for Computational Linguistics, Nov. 2021, pp. 820–842. DOI: 10.
348 18653/v1/2021.findings-emnlp.71. URL: [https://aclanthology.org/2021.](https://aclanthology.org/2021.findings-emnlp.71)
349 [findings-emnlp.71](https://aclanthology.org/2021.findings-emnlp.71).
- 350 [36] John Lyons. *Semantics: Volume 2*. Vol. 2. Cambridge university press, 1977.
- 351 [37] Aleksander Madry et al. “Towards deep learning models resistant to adversarial attacks”. In:
352 *arXiv preprint arXiv:1706.06083* (2017).
- 353 [38] Kyle Mahowald et al. “Dissociating language and thought in large language models: a cognitive
354 perspective”. In: *arXiv preprint arXiv:2301.06627* (2023). URL: [https://doi.org/10.](https://doi.org/10.48550/arXiv.2301.06627)
355 [48550/arXiv.2301.06627](https://doi.org/10.48550/arXiv.2301.06627).
- 356 [39] D. Marconi. *Lexical Competence*. A Bradford book. Bradford Book, 1997. ISBN:
357 9780262133333. URL: https://books.google.com/books?id=lcrEq%5C_7o5m0C.
- 358 [40] Diego Marconi. “Semantic competence”. In: *The Routledge Handbook of Philosophy of Skill*
359 *And Expertise*. Routledge, 2020, pp. 409–418.
- 360

-
- 361 [41] Moran Mizrahi et al. “State of what art? a call for multi-prompt llm evaluation”. In: *Transactions of the Association for Computational Linguistics* 12 (2024), pp. 933–949.
- 362
- 363 [42] Milad Moradi and Matthias Samwald. “Evaluating the robustness of neural language models to input perturbations”. In: *arXiv preprint arXiv:2108.12237* (2021).
- 364
- 365 [43] Frederick J Newmeyer. “The Prague School and North American functionalist approaches to syntax”. In: *Journal of Linguistics* 37.1 (2001), pp. 101–126.
- 366
- 367 [44] Catherine Olsson et al. “In-context learning and induction heads”. In: *arXiv preprint arXiv:2209.11895* (2022).
- 368
- 369 [45] Ellie Pavlick. “Symbols and grounding in large language models”. In: *Philosophical Transactions of the Royal Society A* 381.2251 (2023), p. 20220041.
- 370
- 371 [46] Judea Pearl. “Causal inference in statistics: An overview”. In: *Statistics Surveys* 3 (2009).
- 372 [47] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. “Causal inference by using invariant prediction: identification and confidence intervals”. In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 78.5 (2016), pp. 947–1012.
- 373
- 374
- 375 [48] Fabio Petroni et al. “Language Models as Knowledge Bases?” In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 2463–2473. DOI: 10.18653/v1/D19-1250. URL: <https://aclanthology.org/D19-1250>.
- 376
- 377
- 378
- 379
- 380 [49] Tiago Pimentel et al. “The Architectural Bottleneck Principle”. In: *arXiv preprint arXiv:2211.06420* (2022). URL: <https://doi.org/10.48550/arXiv.2211.06420>.
- 381
- 382 [50] Shauli Ravfogel et al. “Adversarial Concept Erasure in Kernel Space”. In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, Dec. 2022, pp. 6034–6055. URL: <https://aclanthology.org/2022.emnlp-main.405>.
- 383
- 384
- 385
- 386 [51] Shauli Ravfogel et al. “Linear adversarial concept erasure”. In: *International Conference on Machine Learning*. PMLR. 2022, pp. 18400–18421.
- 387
- 388 [52] Shauli Ravfogel et al. “Null It Out: Guarding Protected Attributes by Iterative Nullspace Projection”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, July 2020, pp. 7237–7256. DOI: 10.18653/v1/2020.acl-main.647. URL: <https://aclanthology.org/2020.acl-main.647>.
- 389
- 390
- 391
- 392
- 393 [53] Abhilasha Ravichander et al. “On the Systematicity of Probing Contextualized Word Representations: The Case of Hypernymy in BERT”. In: *Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics*. Barcelona, Spain (Online): Association for Computational Linguistics, Dec. 2020, pp. 88–102. URL: <https://aclanthology.org/2020.starsem-1.10>.
- 394
- 395
- 396
- 397
- 398 [54] Anna Rogers, Olga Kovaleva, and Anna Rumshisky. “A Primer in BERTology: What We Know About How BERT Works”. In: *Transactions of the Association for Computational Linguistics* 8 (2020), pp. 842–866. DOI: 10.1162/tacl_a_00349. URL: <https://aclanthology.org/2020.tacl-1.54>.
- 399
- 400
- 401
- 402 [55] Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. “The risks of invariant risk minimization”. In: *arXiv preprint arXiv:2010.05761* (2020).
- 403
- 404 [56] Sebastian Ruder. “An overview of multi-task learning in deep neural networks”. In: *arXiv preprint arXiv:1706.05098* (2017).
- 405
- 406 [57] Ivan A Sag and Thomas Wasow. “Performance-Compatible Competence Grammar”. In: *Non-Transformational Syntax: Formal and Explicit Models of Grammar* (2011), pp. 359–377.
- 407
- 408 [58] Shun Shao, Yftah Ziser, and Shay B. Cohen. *Gold Doesn’t Always Glitter: Spectral Removal of Linear and Nonlinear Guarded Attribute Information*. arXiv:2203.07893 [cs]. Mar. 2022. DOI: 10.48550/arXiv.2203.07893. URL: <http://arxiv.org/abs/2203.07893> (visited on 09/22/2022).
- 409
- 410
- 411
- 412 [59] Donghee Shin. “The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable AI”. In: *International Journal of Human-Computer Studies* 146 (2021), p. 102551.
- 413
- 414

-
- 415 [60] Charlotte Siska et al. “Examining the robustness of LLM evaluation to the distributional
416 assumptions of benchmarks”. In: *Proceedings of the 62nd Annual Meeting of the Association
417 for Computational Linguistics (Volume 1: Long Papers)*. 2024, pp. 10406–10421.
- 418 [61] Robyn Speer, Joshua Chin, and Catherine Havasi. “ConceptNet 5.5: An Open Multilingual
419 Graph of General Knowledge”. In: *Proceedings of the AAAI Conference on Artificial Intelli-
420 gence*. Vol. 31. 2017.
- 421 [62] Anant Subramanian et al. “Spine: Sparse interpretable neural embeddings”. In: *Proceedings of
422 the AAAI conference on artificial intelligence*. Vol. 32. 1. 2018.
- 423 [63] Adly Templeton et al. “Scaling Monosemanticity: Extracting Interpretable Features from
424 Claude 3 Sonnet”. In: *Transformer Circuits Thread* (2024). URL: [https://transformer-
425 circuits.pub/2024/scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 426 [64] Hugo Touvron et al. “Llama 2: Open foundation and fine-tuned chat models”. In: *arXiv preprint
427 arXiv:2307.09288* (2023).
- 428 [65] Hugo Touvron et al. “Llama: Open and efficient foundation language models”. In: *arXiv
429 preprint arXiv:2302.13971* (2023).
- 430 [66] Mycal Tucker, Peng Qian, and Roger Levy. “What if This Modified That? Syntactic Inter-
431 ventions with Counterfactual Embeddings”. In: *Findings of the Association for Computa-
432 tional Linguistics: ACL-IJCNLP 2021*. Online: Association for Computational Linguistics,
433 Aug. 2021, pp. 862–875. DOI: 10.18653/v1/2021.findings-acl.76. URL: [https://
434 aclanthology.org/2021.findings-acl.76](https://aclanthology.org/2021.findings-acl.76).
- 435 [67] Francisco Vargas and Ryan Cotterell. “Exploring the Linear Subspace Hypothesis in Gender
436 Bias Mitigation”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural
437 Language Processing (EMNLP)*. Online: Association for Computational Linguistics, Nov.
438 2020, pp. 2902–2913. DOI: 10.18653/v1/2020.emnlp-main.232. URL: [https://
439 aclanthology.org/2020.emnlp-main.232](https://aclanthology.org/2020.emnlp-main.232).
- 440 [68] Jindong Wang et al. “On the Robustness of ChatGPT: An Adversarial and Out-of-distribution
441 Perspective”. In: *ICLR 2023 Workshop on Trustworthy and Reliable Large-Scale Machine
442 Learning Models*. 2023. URL: <https://openreview.net/forum?id=uw6HSkgoM29>.
- 443 [69] Kevin Ro Wang et al. “Interpretability in the Wild: a Circuit for Indirect Object Identification
444 in GPT-2 Small”. In: *The Eleventh International Conference on Learning Representations*.
445 2023. URL: <https://openreview.net/forum?id=NpsVSN6o4u1>.
- 446 [70] Thomas Wolf et al. “Huggingface’s transformers: State-of-the-art natural language processing”.
447 In: *arXiv preprint arXiv:1910.03771* (2019).
- 448 [71] Zhengxuan Wu et al. “Interpretability at Scale: Identifying Causal Mechanisms in Alpaca”. In:
449 *arXiv preprint arXiv:2305.08809* (2023). arXiv: 2305.08809 [cs.LG].
- 450 [72] Karren Yang, Abigail Katcoff, and Caroline Uhler. “Characterizing and learning equivalence
451 classes of causal DAGs under interventions”. In: *International Conference on Machine Learn-
452 ing*. PMLR. 2018, pp. 5541–5550.
- 453 [73] Linyi Yang et al. “GLUE-X: Evaluating Natural Language Understanding Models from an Out-
454 of-Distribution Generalization Perspective”. In: *Findings of the Association for Computational
455 Linguistics: ACL 2023*. Ed. by Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki.
456 Toronto, Canada: Association for Computational Linguistics, July 2023, pp. 12731–12750.
457 DOI: 10.18653/v1/2023.findings-acl.806. URL: [https://aclanthology.org/
458 2023.findings-acl.806](https://aclanthology.org/2023.findings-acl.806).
- 459 [74] Zeyu Yun et al. “Transformer visualization via dictionary learning: contextualized embedding
460 as a linear superposition of transformer factors”. In: *Proceedings of Deep Learning Inside Out
461 (DeeLIO): The 2nd Workshop on Knowledge Extraction and Integration for Deep Learning
462 Architectures*. Ed. by Eneko Agirre, Marianna Apidianaki, and Ivan Vulić. Online: Association
463 for Computational Linguistics, June 2021, pp. 1–10. DOI: 10.18653/v1/2021.deelio-1.1.
464 URL: <https://aclanthology.org/2021.deelio-1.1>.
- 465 [75] Susan Zhang et al. “Opt: Open pre-trained transformer language models”. In: *arXiv preprint
466 arXiv:2205.01068* (2022).
- 467 [76] Han Zhao et al. “Fundamental limits and tradeoffs in invariant representation learning”. In:
468 *The Journal of Machine Learning Research* 23.1 (2022), pp. 15356–15404.
- 469 [77] Ziqian Zhong et al. “The clock and the pizza: Two stories in mechanistic explanation of neural
470 networks”. In: *arXiv preprint arXiv:2306.17844* (2023).

471 [78] Andy Zou et al. “Representation engineering: A top-down approach to ai transparency”. In:
472 *arXiv preprint arXiv:2310.01405* (2023).

473 A Additional Context

474 A.1 Background and Related Work

475 **Linguistic Competence** There has been significant debate in linguistics and the philosophy of
476 language regarding the precise definition and nature of competence [36, 43, 57, 40]. However, the
477 formalization of competence provided by CALM is sufficiently general to incorporate most notions
478 of competence, which may be flexibly specified by instantiating CALM in different ways. In this
479 work, we focus on *lexicosemantic competence*, the ability to utilize knowledge of word meaning
480 relationships in performing tasks such as lexical inference [39, 40].

481 **Relational Properties** Why is it not possible for linear information removal methods such as INLP
482 [52] to remove relational properties between multiple input embeddings? Consider a binary relational
483 property Z denoting whether a relation $Z(i, j)$ holds between multiple embeddings $\mathbf{h}_i, \mathbf{h}_j \in \mathbb{R}^d$. In
484 INLP, we begin by training linear classifier $W \in \mathbb{R}^{2d}$ to predict Z from concatenated embeddings
485 $(\mathbf{h}_i; \mathbf{h}_j)$, where the goal is to train $W : (\mathbf{h}_i; \mathbf{h}_j) \mapsto Z$. We may decompose $W(\mathbf{h}_i; \mathbf{h}_j) = W_{[0:d]}\mathbf{h}_i +$
486 $W_{[d+1:2d]}\mathbf{h}_j$, where $W_{[0:d]}$ is the first d dimensions of W and $W_{[d+1:2d]}$ are the second d dimensions.
487 In this case, there is no interaction between the two inputs $\mathbf{h}_i, \mathbf{h}_j$, meaning there is no way for W to
488 take into account any relationship Z between them.

489 A.2 Benefits and Limitations of GBIs

490 **Benefits** The key advantage of gradient-based interventions (GBIs) as a causal probing methodology
491 is that they may be applied to any differentiable probe. For example, if we are investigating the
492 hypothesis that M 's representation of Z is captured by a linear subspace of representations in a given
493 layer (see [67]), then we may train a linear probe and various nonlinear probes on representations
494 and observe whether GBIs against the linear probe have a comparable impact to those against the
495 nonlinear probes. Alternatively, if we believe that a probe's architecture should mirror the architecture
496 of the model it is probing (as argued by [49]), we may implement probes as such. Finally, where
497 previous intervention methodologies for causal probing have focused on *nullifying* interventions that
498 remove the representation of the target property Z [52, 50, 51, 58, 2], GBIs allow one to perform
499 targeted interventions that set LLMs' representations to counterfactual values $do(Z = z')$, effectively
500 simulating the model's behavior under counterfactual inputs, which may be useful for predicting
501 behaviors under various distribution shifts (see Appendix C.1).

502 **Limitations** While GBIs are applicable to a more general range of model representations than
503 other interventions, this generality comes with a lack of constraints on probes (g_Z); and as a result,
504 GBIs cannot provide the strong theoretical constraints on collateral damage as can methods like,
505 e.g., INLP [52], which provably preserves distances between embeddings as well as possible while
506 completely removing the linear representation of the target property. To minimize collateral damage
507 to representations, the magnitude of perturbations should be modulated via constraints on gradient
508 attacks against g_Z (see Section 4) and experimentally validated to control the damage done to
509 representations (see Appendix B.4). Thus, in cases where the structure of representations is believed
510 to satisfy strong assumptions (e.g., being restricted to a linear subspace; [67]) or strong upper bounds
511 on collateral damage are required, CALM interventions can be implemented with methods like INLP
512 rather than GBIs.⁶

513 B Experimental Details

514 B.1 Simplified Environment

515 As noted in Section 4, our primary goal in our experiments is to validate CALM by testing it in a
516 simplified experimental setting consisting of comparatively small, well-studied models and tasks. As
517 such, we need models that are *just complex enough* for CALM to be applicable (i.e., neural language
518 models that are capable of performing the tasks we consider at a nontrivial level of performance),
519 making BERT and RoBERTa ideal candidates; and in future work plan to scale CALM to more

⁶It may also be possible to control for collateral damage by developing GBI strategies that offer more principled protection against damage to non-targeted properties, such as adding a loss term to penalize damage to non-targeted probes or leveraging interval bound propagation [25] to place intervened embeddings inside the adversarial polytope for non-targeted properties. We leave such possibilities to future work.

520 complex contexts covering larger, more powerful models as they perform more difficult tasks (see
521 Appendix D). This is a common setting in the context of substantial recent interpretability work:
522 first, a theoretical framework is developed for interpreting an internal representation or mechanism
523 and initially tested in the context of “toy” models or tasks [20, 44, 77, 22], and subsequent work
524 scales these frameworks to the context of larger models “in the wild” [69, 13, 71]. We anticipate that
525 all of our major contributions (the CALM framework, competence metric, and GBI causal probing
526 method) will in principle be scalable to much larger, more recent LLMs (e.g., [75, 4, 65, 64, 26], etc.),
527 and predict that the main challenge will be in finding an appropriate probing architecture (see [49]).

528 B.2 Tasks

529 The full set of LAMA ConceptNet tasks is as follows: IsA, HasA, PartOf, HasSubEvent, MadeOf,
530 HasPrerequisite, MotivatedByGoal, AtLocation, CausesDesire, NotDesires, CapableOf, UsedFor,
531 ReceivesAction, and HasProperty. We split each task dataset into train, validation, and test sets with
532 a random 80%/10%/10% split. Train and validation instances are fed to each model to produce
533 embeddings used to train g_Z and select hyperparameters, respectively; and test instances are used to
534 measure LLMs’ competence with respect to each task by observing how predictions change under
535 various interventions. In all experiments, we restrict each model M ’s output space for each task
536 \mathcal{T} to the subset of vocabulary V_M that occurs as a ground-truth answer y^* for at least one instance
537 $(\mathbf{x}, y^*) \sim \mathcal{T}$ in the respective task dataset. This lowers the probability of false negatives in evaluation
538 (e.g., penalizing the model for predicting $\hat{y} = \text{“mammal”}$ for “a dog is a type of y ” instead of $y^* =$
539 “animal”).

540 **Experimental Limitation** In our experiments, we modeled the 14 LAMA ConceptNet tasks as
541 representing fully independent properties, which is not necessarily true – e.g., knowing that a tree is
542 made of bark or contains leaves tells us something about whether it is a type of plant. However, in
543 the aggregate (with impacts summed across 14 widely-varying lexical relation types in computing
544 the final competence score for each task; see Appendix C.2), it may nonetheless be appropriate to
545 treat the relations which are not causal with respect to a given task as collectively capturing spurious
546 lexical associations.

547 B.3 Probes

548 We use BERT’s final layer L to encode h_i^l embeddings for each such example, where i is the index
549 of the [MASK] token or target word in the input prompt x_i . To encode the [MASK] token, we issue
550 BERT masked prompts (as discussed above) to extract $h_{[\text{MASK}]}$, then repeat with the [MASK] token
551 filled-in with the target word to encode it as h_+ (e.g., “device” in “A laser is a device which creates
552 coherent light.”), and concatenate matching embeddings $h = (h_{[\text{MASK}]}; h_+)$ to produce positive
553 ($y = 1$) training instances. We also construct one negative ($y = 0$) instance, $h = (h_{[\text{MASK}]}; h_-)$, for
554 each $h_{[\text{MASK}]}$ by sampling an incorrect target word x_i corresponding to an answer to a random prompt
555 from the same task, feeding it into the cloze prompt in the place of the correct answer, and obtaining
556 BERT’s contextualized final-layer embedding of this token (h_-). Finally, we train g_Z on the set of all
557 such (h, y) .

558 We implement g_Z as a multi-layer perceptron with 2 hidden layers, each with a width of 768 (which
559 is one half the concatenated input dimension of 1536), using ReLU activations and dropout with
560 $p = 0.1$, training it for 32 epochs using Binary Cross Entropy with Logits Loss⁷ and the Adam
561 optimizer, saving the model from the epoch with the highest validation-set accuracy for use in all
562 experiments.

563 For all competence results reported in Section 5, we run the same experiment 10 times – each with a
564 different random initialization of g_Z and shuffled training data – and report each figure as the average
565 among all 10 runs.

⁷<https://pytorch.org/docs/stable/generated/torch.nn.BCEWithLogitsLoss.html>

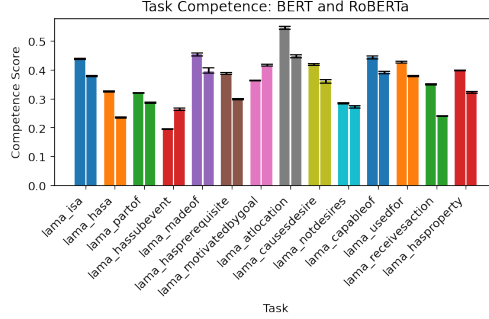


Figure 4: Competence of BERT (left bars) and RoBERTa (right bars) for all tasks, using PGD with $\epsilon = 0.1$. Y-values are the average competence score and error bars are the maximum and minimum competence score, as measured over 10 experimental iterations (each with a different randomly-initialized probe g_Z).

566 B.4 Interventions

567 For instance (h, y) , classifier g_Z , loss function \mathcal{L} , and L_∞ -bound $\epsilon \in \{0.01, 0.03, 0.1, 0.3\}$ ⁸,
 568 each intervention (gradient attack) g_z may be used to produce perturbed representations $h' =$
 569 $g_z(h, y, f_{\text{cls}}, \mathcal{L}, \epsilon)$ where $\|h - h'\|_\infty \leq \epsilon$. In particular, given $h = (h_{[\text{MASK}]}; h_\pm) \in \mathbb{R}^{2d}$, let
 570 $h'_{[\text{MASK}]}$ be the first d dimensions of h' (which also satisfies the L_∞ -bound with respect to $h_{[\text{MASK}]}$,
 571 $\|h_{[\text{MASK}]} - h'_{[\text{MASK}]}\|_\infty \leq \epsilon$). To measure BERT’s use of internal representations of Z on each
 572 prompt task, we evaluate its performance when perturbed $h'_{[\text{MASK}]}$ is used to compute masked-word
 573 predictions, compared to unperturbed $h_{[\text{MASK}]}$.

574 Our intent in intervening only on the final-layer mask embedding $h_{[\text{MASK}]}$ in our experiments is that,
 575 in the final layer of a masked language model such as BERT or RoBERTa, the only embedding which
 576 is used to compute masked-word probabilities is that of the [MASK] token. Thus, any representation
 577 of the property that is *used* by the model in its final layer must be a part of its representation of the
 578 [MASK] token, preventing “recoverability” phenomena such as those observed by Elazar et al. [18].

579 **FGSM** We implement Fast Gradient Sign Method (FGSM; [24]) interventions as

$$h' = h + \epsilon \cdot \text{sgn}(\nabla_h \mathcal{L}(f_{\text{cls}}, x, y))$$

580 **PGD** We implement Projected Gradient Descent (PGD; [8, 37]) interventions as $h' = h^T$ where

$$h^{t+1} = \Pi_{N(h)}(h^t + \alpha \cdot \text{sgn}(\nabla_h \mathcal{L}(f_{\text{cls}}, x, y)))$$

581 for iterations $t = 0, 1, \dots, T$, projection operator Π , and L_∞ -neighborhood $N(h) = \{h' : \|h -$
 582 $h'\| \leq \epsilon\}$. This method also introduces two hyperparameters, the number of PGD iterations T
 583 and step size α . We use hyperparameter grid search over $\alpha \in \{0.001, 0.003, 0.01, 0.03\}$ and
 584 $T \in \{20, 40, 60, 80, 100\}$, finding that setting $\alpha = \frac{\epsilon}{10}$ and $T = 40$ produces the most consistent
 585 impact on g_Z accuracy across all tasks; so we use these values for the results visualized in Figure 4.

586 B.5 Compute Budget

587 BERT-base-uncased has 110 million parameters, and RoBERTa-base has 125M parameters. As our
 588 goal is to study the internal representation and use of linguistic properties in existing pre-trained
 589 models, and we are not directly concerned with training or fine-tuning such models, we use these
 590 models only for inference (including encoding text inputs, using embeddings to train probes, and
 591 feeding intervened embeddings back into the language models). The only models we trained were
 592 probes g_Z , which each had 1.77M parameters.

593 Each experimental iteration (including encoding text inputs, training probes on all 14 tasks, and
 594 performing all GBIs) for either BERT or RoBERTa took less than one hour on a single NVIDIA
 595 GeForce GTX 1080 GPU, meaning that running all 10 iterations across both language models took

⁸All reported results use $\epsilon = 0.1$, as greater ϵ resulted in unacceptably high “collateral damage” across target tasks (e.g., even random perturbations of magnitude $\epsilon = 0.3$ do considerable damage), and lesser values meant that predictions changed on target tasks consisted of only a few test instances.

596 less than 20 hours on a single GPU. Each iteration, probe, and GBI can easily be parallelized across
597 GPUs: in our case, running all iterations across both models took less than 3 hours total across 8
598 GTX 1080 GPUs.

599 C Competence Metric

600 C.1 Comparison With IIA

601 As noted in Section 2, the $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$ metric defined in Equation (2) is an adaptation of the Inter-
602 change Intervention Accuracy (IIA) metric [23, 22], which evaluates the faithfulness of a causal
603 abstraction like $\mathcal{G}_{\mathcal{T}}$ as a (potential) explanation of the behavior of a “black box” system like M . In
604 our case, this is equivalent to evaluating the competence of M on task \mathcal{T} , provided that $\mathcal{G}_{\mathcal{T}}$ is the
605 appropriate SCM for \mathcal{T} , as an LLM is competent only to the extent that its behavior is determined by
606 a causally invariant representation of the task.⁹ IIA requires performing *interchange interventions*
607 $M(\mathbf{x}_i|\text{do}(\mathbf{z}_i))$, where the part of M ’s intermediate representation of input \mathbf{x}_i hypothesized to encode
608 latent variables \mathbf{Z} (taking the values \mathbf{z}_i when provided input \mathbf{x}_i) is replaced with that of \mathbf{x}_j (which,
609 in the ideal case, causes M ’s representation to encode the values \mathbf{z}_j instead of \mathbf{z}_i), and compute the
610 accuracy of $\mathcal{G}_{\mathcal{T}}(\mathbf{x}_i|\text{do}(\mathbf{z}_j))$ in predicting M ’s behavior under these interventions. Thus, given access
611 to high-quality interchange interventions over M , IIA measures the extent to which $\mathcal{G}_{\mathcal{T}}$ correctly
612 models M ’s behavior under counterfactuals, and thus its faithfulness as a causal abstraction of M .

613 To adapt IIA to the context of causal probing and define $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$, we replace instance-level
614 interchange interventions with concept-level interventions: instead of swapping M ’s representation of
615 variables \mathbf{Z} given input \mathbf{x}_i with that of \mathbf{x}_j , we intervene on representations at the level of arbitrary concept
616 settings \mathbf{z} that need not correspond to previously sampled \mathbf{x} , allowing us to simulate the behavior
617 of M under previously-unseen distribution shifts (i.e., settings \mathbf{z} representing previously-unseen
618 combinations of property values) and therefore make broader predictions about M ’s consistency
619 with a given causal model $\mathcal{G}_{\mathcal{T}}$ under such conditions. As one of the key desiderata in studying LLM
620 competence is to predict behavior under distribution shifts where spurious correlations are broken,
621 $\mathcal{C}_{\mathcal{T}}$ is more appropriate than IIA in this setting. However, it also introduces an additional challenge:
622 where interchange interventions only require localizing candidate representations – as counterfactual
623 representations are obtained merely by “plugging in” values from a different input – computing
624 $\mathcal{C}_{\mathcal{T}}$ instead requires one to both localize representations and directly intervene on them to change
625 the encoded value. Previous causal probing intervention strategies (e.g., [52, 50]) have generally
626 performed interventions by *neutralizing* concept representations, not modifying them to encode
627 specific counterfactual values; so in order to carry out our study, it is also necessary to develop a
628 novel approach to perform such interventions. We develop a solution to this problem, gradient-based
629 interventions (GBIs), in Section 3.

630 C.2 Experimental Competence Metric

631 To compute the expectation in Equation (2) for test set $\{\mathbf{x}_i, \mathbf{y}_i, \mathbf{z}_i\}_{i=1}^n \sim \mathcal{T} \times \mathbf{Z}$, we sum the
632 competence score over all samples \mathbf{x}_i and perform one intervention $\text{do}(Z_j = 0)$ corresponding to
633 each concept $Z_j \in \mathbf{Z}$.¹⁰ As our goal is to measure the extent to which M ’s behavior is attributable
634 to an underlying representation of the causal property Z_c or environmental property $Z \in \mathbf{Z}_e$, our
635 experimental model defines $\mathcal{G}_{\mathcal{T}}$ ’s predictions with reference to M ’s original predictions $M(\mathbf{x}_i) = \hat{\mathbf{y}}_i$,
636 according to the following principle: if M is competent, then its prediction $M(\mathbf{x}_i) = \hat{\mathbf{y}}_i$ is wholly
637 attributable to its representation of causal property Z_c , so its predictions $M(\mathbf{x}_i|\text{do}(Z_c)) = \hat{\mathbf{y}}_i'$ will not
638 overlap with its original predictions $\hat{\mathbf{y}}_i$ (i.e., $\text{overlap}(\hat{\mathbf{y}}_i, \hat{\mathbf{y}}_i') = 0$); and conversely, a competent M
639 will make the *same* predictions $M(\mathbf{x}_i|\text{do}(Z_j)) = \hat{\mathbf{y}}_i''$ for any $Z_j \in \mathbf{Z}_e$, because its prediction is not
640 caused by its representation of these environmental properties (i.e., $\text{overlap}(\hat{\mathbf{y}}_i, \hat{\mathbf{y}}_i'') = 1$). Motivated
641 by this reasoning, our experimental model defines $\mathcal{G}_{\mathcal{T}}(\mathbf{x}_i|\text{do}(Z_j = 0)) = M(\mathbf{x}_i)$ for environmental

⁹For many tasks, there is more than one valid $\mathcal{G}_{\mathcal{T}}$ (see, e.g., the “price tagging game” constructed by Wu et al. [71]). In such cases, $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$ should be computed with respect to each valid $\mathcal{G}_{\mathcal{T}}$ and the highest result should be selected, as conforming to any such $\mathcal{G}_{\mathcal{T}}$ carries the same implications.

¹⁰Note that this intervention changes the prediction $\mathcal{G}_{\mathcal{T}}(\mathbf{x}_i) \neq \mathcal{G}_{\mathcal{T}}(\mathbf{x}_i|\text{do}(Z_j = 0))$ if and only if $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{T}_j$ – i.e., where the corresponding $(\mathbf{z}_i)_j = 1$ – otherwise, $(\mathbf{z}_i)_j$ is already 0, so the intervention has no effect. Thus, as $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$ measures M ’s consistency with $\mathcal{G}_{\mathcal{T}}$, then to the extent that M is competent, its prediction should change under all and only the same interventions as $\mathcal{G}_{\mathcal{T}}$.

642 $Z_j \in \mathbf{Z}_c$; and for causal property Z_c , defines $\mathcal{G}_{\mathcal{T}}(\mathbf{x}_i | \text{do}(Z_c = 0)) = \{y' \in V_M : y' \notin M(\mathbf{x}_i)\}$ (i.e.,
 643 the set of all tokens y' in M 's vocabulary that were not in its original prediction $M(\mathbf{x}_i)$). Thus, under
 644 experimental model E , we approximate $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$ by computing it as follows:

$$\hat{\mathcal{C}}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}}) = \frac{1}{n \cdot m} \sum_{i=1}^n \sum_{j=1}^m \text{overlap} \left(M(\mathbf{x}_i | \text{do}(Z_j = 0)), \mathcal{G}_{\mathcal{T}}(\mathbf{x}_i | \text{do}(Z_j = 0)) \right) \quad (3)$$

645 Notably, our experimental model E only accounts for the relationship between M 's intervened and
 646 non-intervened predictions, independently of ground truth labels – instead, what is being measured is
 647 M 's consistency under meaning-preserving interventions $\text{do}(Z_{j'})$ and its mutability under meaning-
 648 altering interventions $\text{do}(Z_j)$. However, as we find in Section 5, the resulting competence metric
 649 $\mathcal{C}_{\mathcal{T}}(M|\mathcal{G}_{\mathcal{T}})$ is nonetheless useful for predicting M 's accuracy.

650 D Future Work

651 D.1 Representation Learning

652 The CALM framework, competence measure, and GBI methodology developed in Sections 2 and 3
 653 are sufficiently general to be directly applied to analyze arbitrary LLMs on any language modeling
 654 task whose causal structure is already well understood (or, for tasks where this is not the case, we
 655 may apply the causal graph discovery approach described in Appendix D.4), allowing us to study
 656 the impact of various model architectures, pre-training regimes, and fine-tuning strategies on the
 657 representations LLMs learn and use for arbitrary tasks of interest.

658 D.2 Multitask Learning

659 Are high competence scores on task \mathcal{T} correlated with an LLMs' robustness to meaning-preserving
 660 transformations (see, e.g., [19]) on tasks \mathcal{T}' that share several causal properties \mathbf{Z}_c with task \mathcal{T} .
 661 Through the lens of causally invariant prediction [47, 1, 9], this hypothesis is likely true (however,
 662 see [55] for appropriate caveats) – if so, this would make it possible to use clusters of related
 663 tasks to predict LLMs' robustness (and other behavioral patterns, such as brittleness in the face of
 664 distribution shifts introduced by spurious dependencies) between related tasks using CALM, given
 665 an appropriate experimental model. Furthermore, the ability to characterize tasks based on mutual
 666 (learned) dependency structures could be valuable in transfer learning applications such as guiding
 667 the selection of auxiliary tasks in multi-task learning [56] or predicting the impact of intermediate
 668 task fine-tuning on downstream target tasks [11].

669 D.3 Task Dependencies

670 Another possible application of CALM concerns causal invariance under multi-task applications.
 671 Existing approaches in invariant representation learning generally require task-specific training [76],
 672 as the notion of invariance is inherently task-centric (i.e., the properties which are invariant predictors
 673 of output values vary by task, and different tasks may have opposite notions of which properties are
 674 causal versus environmental), so applying such approaches to train models to be causally invariant
 675 with respect to a specific downstream task \mathcal{T} is expected to come at the cost of performance on other
 676 downstream tasks \mathcal{T}' . Therefore, considering the recent rise of open-ended, task-general LLMs [75,
 677 4, 65, 64, 26], it is important to find alternative approaches for studying models' causal dependencies
 678 in a task-general setting to account for applications involving tasks with different (and perhaps
 679 contradictory) causal structures, such as CALM.

680 D.4 Causal Competence Graph Discovery

681 One of the key benefits of CALM is that, instead of simply measuring consistency with respect to a
 682 known, static task description $\mathcal{G}_{\mathcal{T}}$, the competence metric in Equation (2) can also be used to discover
 683 a competence graph \mathcal{G} which most faithfully explains a model M 's behavior in a given task or context
 684 (see Section 2) by computing $\mathcal{C}(M|\mathcal{G})$ “in-the-loop” of existing causal graph discovery algorithms
 685 like IGSP [72]. Such algorithms can be used both to suggest likely competence graphs based on
 686 interventional data collected by running CALM experiments, to recommend the experiments that

687 would yield the most useful interventional data for the graph discovery algorithm, and to evaluate
688 candidate graphs \mathcal{G} using our competence metric, terminating the graph discovery algorithm once a
689 competence graph \mathcal{G} that offers sufficiently faithful explanations of M 's behavior has been found. In
690 this case, it is still necessary to define the set of properties \mathbf{Z} being probed and the scoring function
691 S used to compare the predictions of M and \mathcal{G} ; but no knowledge of the causal dependencies (or
692 structural functions $F : \text{pa}(Z_j) \mapsto Z_j$ mapping from causal parents $\text{pa}(Z_j)$ to causal dependents
693 Z_j ; see [5]) is required.