

ON LINEAR REPRESENTATIONS AND PRETRAINING DATA FREQUENCY IN LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Pretraining data has a direct impact on the behaviors and quality of language models (LMs), but we only understand the most basic principles of this relationship. While most work focuses on pretraining data’s effect on downstream task behavior, we investigate its relationship to LM representations. Previous work has discovered that, in language models, some concepts are encoded as “linear representations”, but what factors cause these representations to form (or not)? We study the connection between differences in pretraining data frequency and differences in trained models’ linear representations of factual recall relations. We find evidence that the two are linked, with the formation of linear representations strongly connected to pretraining term frequencies. First, we establish that the presence of linear representations for subject-relation-object (s-r-o) fact triplets is highly correlated with both subject-object co-occurrence frequency and in-context learning accuracy. This is the case across all phases of pretraining, i.e., it is not affected by the model’s underlying capability. In OLMo 7B and GPT-J (6B), we discover that a linear representation consistently (but not exclusively) forms when the subjects and objects within a relation co-occur at least 1-2k times, regardless of when these occurrences happen during pretraining. In the OLMo 1B model, consistent linearity only occurs after 4.4k occurrences, suggesting a connection to scale. Finally, we train a regression model on measurements of linear representation quality that can predict how often a term was seen in pretraining. We show such model achieves low error even for a different model and pretraining dataset, providing a new unsupervised method for exploring possible data sources of closed-source models. We conclude that the presence or absence of linear representations in LMs contains signal about their pretraining corpora that may provide new avenues for controlling and improving model behavior. We release our code to support future work¹

1 INTRODUCTION

Understanding how the content of pretraining data affects language model (LM) behaviors and performance is an active area of research (Ma et al., 2024; Xie et al., 2024; Aryabumi et al., 2024; Longpre et al., 2024; Antoniadou et al., 2024; Seshadri et al., 2024; Razeghi et al., 2023; Wang et al., 2024). For instance, it has been shown that for specific tasks, models perform better on instances containing higher frequency terms than lower frequency ones (Razeghi et al., 2022; Mallen et al., 2023a). The ways in which frequency affects the internal representations of LMs to cause this difference in performance remain unclear. We connect dataset statistics to recent work in interpretability, which focuses on the emergence of simple linear representations of factual relations in LMs. Our findings demonstrate a strong correlation between these features and the frequency of terms in the pretraining corpus.

Linear representations in LMs have become central to interpretability research in recent years (Ravfogel et al., 2020; Elazar et al., 2021; Elhage et al., 2021; Slobodkin et al., 2023; Olah et al., 2020; Park et al., 2024; Jiang et al., 2024; Black et al., 2022; Chanin et al., 2024). Linear representations are essentially linear approximations (linear transforms, directions in space) that are simple to understand, and strongly approximate the complex non-linear transformations that networks are

¹Anonymized

054 implementing. These representations are crucial because they allow us to localize much of the be-
 055 havior and capabilities of LMs to specific directions in activation space. This means that certain
 056 behaviors can be activated or modulated by intervening on these directions with linear projections at
 057 inference time, a process also known as steering (Todd et al., 2024; Subramani et al., 2022; Hendl
 058 et al., 2023; Rimsky et al., 2023).

059 Recent work by Hernandez et al. (2024) and Chanin et al. (2024) highlight how the linearity of dif-
 060 ferent types of relations varies greatly depending on the specific relationships being depicted. For
 061 example, over 80% of “country largest city” relations can be approximated by a single linear trans-
 062 formation on the contextual embedding of the country, but less than 30% of “star in constellation”
 063 can be. Their methods for identifying representations with linear structure do not offer an explan-
 064 ation for this. Such findings complicate the understanding of the Linear Representation Hypothesis,
 065 which proposes that LMs will represent features linearly (Park et al., 2024). While Jiang et al. (2024)
 066 provide both theoretical and empirical evidence that the training objectives of LMs implicitly en-
 067 courage linear representations, it remains unclear why some features are represented this way while
 068 others are not. This open question is a central focus of our investigation.

069 Whether linear representations for “common” concepts are actually more prevalent in models or
 070 simply easier to identify (using current methods) than those for less common concepts remains un-
 071 clear. We hypothesize that factual relations exhibiting linear representations are correlated with
 072 higher mention frequencies in the pretraining data (as has been shown with static embeddings, see
 073 Ethayarajh et al., 2019), which we confirm in Section 4. Our results also indicate that this can occur
 074 at any point in pretraining, as long as a certain average frequency is reached across subject-object
 075 pairs in a relation. In order to count the appearance of terms in data corpora throughout training,
 076 we develop an efficient tool for counting tokens in tokenized batches of text, which we release to
 077 support future work in this area. We also explore whether the presence of linear representations can
 078 provide insights into relation frequency. In Section 5, we fit a regression model to predict the fre-
 079 quency of individual terms (such as “The Beatles”) in pretraining data, based on metrics measuring
 080 the presence of a linear feature for some relation. For example, how well a linear transformation
 081 approximates the internal computation of the “lead singer of” relation mapping “John Lennon” to
 082 “The Beatles” can tell us about the frequency of those terms in the pretraining corpus.

083 Our findings indicate that the predictive signal, although approximate, is much stronger than that
 084 encoded in log probabilities and task accuracies alone, allowing us to estimate the frequencies of
 085 held-out relations and terms within approximate ranges. Importantly, this regression model gen-
 086 eralizes beyond the specific LM it was trained on without additional supervision. This provides a
 087 valuable foundation for analyzing the pretraining corpora of closed-data models with open weights.

088 To summarize, in this paper we show that:

- 089 1. The development of linear representations for factual recall relations in LMs is related to
 090 frequency as well as model size.
- 091 2. Linear representations form at predictable frequency thresholds during training, regardless
 092 of when this frequency threshold is met for the nouns in the relation. The formation of
 093 these also correlates strongly with recall accuracy.
- 094 3. Measuring the extent to which a relation is represented linearly in a model allows us to
 095 predict the approximate frequencies of individual terms in the pretraining corpus of that
 096 model, even when we do not have access to the model’s training data.
- 097 4. We release a tool for accurately and efficiently searching through tokenized text to support
 098 future research on training data.

101 2 BACKGROUND

103 2.1 LINEAR REPRESENTATIONS

104
 105 Representing information in distributed vector spaces has a long history in language processing,
 106 where geometric properties of these spaces were used to encode semantic information (Salton et al.,
 107 1975; Paccanaro & Hinton, 2001). When and why linear structure emerges without explicit bias to
 do so has been of considerable interest since the era of static word embeddings. Work on skipgram

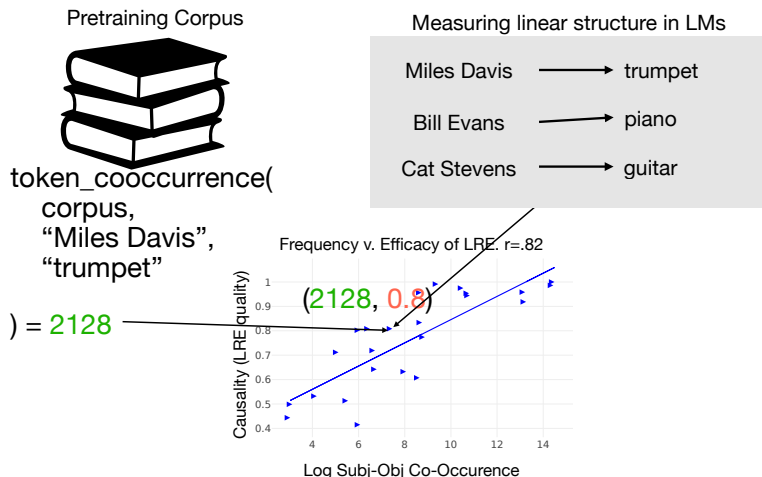


Figure 1: Overview of this work. Given a dataset of subject-relation-object factual relation triplets, we count subject-object co-occurrences throughout pretraining batches. We then measure how well the corresponding relations are represented within an LM across pretraining steps, using the Linear Relational Embeddings (LRE) method from Hernandez et al. (2024). We establish a strong relationship between average co-occurrence frequency and a model’s tendency to form linear representations for relations. From this, we show that we can predict frequencies in the pretraining corpus

models (Mikolov, 2013) found that vector space models of language learn regularities which allow performing vector arithmetic between word embeddings to calculate semantic relationships (e.g., France-Paris+Spain=Madrid) (Mikolov et al., 2013; Pennington et al., 2014). This property was subject to much debate, as it was not clear why word analogies would appear for some relations and not others (Köper et al., 2015; Karpinska et al., 2018; Gladkova et al., 2016). Followup work showed that linguistic regularities form in static embeddings for relations under specific dataset frequency constraints for relevant terms (Ethayarajh et al., 2019), but does not clearly relate to how modern LMs learn. More recently, there has been renewed interest in the presence of similar linear structure in models with contextual embeddings like transformer language models (Park et al., 2024; Jiang et al., 2024; Merullo et al., 2024). As a result, there are many ways to find and test for linear representations in modern LMs, though the relationship to pretraing data is not addressed (Huben et al., 2023; Gao et al., 2024; Templeton et al., 2024; Rimsky et al., 2023; Todd et al., 2024; Hendel et al., 2023; Hernandez et al., 2024; Chanin et al., 2024). Many of these share similarities in how they compute and test for the linear representations, typically through counterfactuals. We focus on a particular class of linear representations called Linear Relational Embeddings (LREs) (Paccanaro & Hinton, 2001).

Linear Relational Embeddings (LREs) Hernandez et al. (2024) use a particular class of linear representation called a Linear Relational Embedding (Paccanaro & Hinton, 2001) to approximate the computation performed by a model to predict the objects that complete common subject-relation-object triplets as an affine transformation. This transform is calculated from hidden state s , the subject token representation at some middle layer of the model, to o , the hidden state at the last token position and layer of the model (i.e., the final hidden state that decodes a token in an autoregressive transformer). For example, given the input sequence “Miles Davis (subject) plays the (relation)”, the goal is to approximate the computation of the object “trumpet”, assuming the model predicts the object correctly. It was found that this transformation holds for nearly every subject and object in the relation set (such as “Cat Stevens plays the guitar”) for some relations. This is surprising because, despite the non-linearities within the many layers and token positions separating s and o , a simple structure within the representation space well approximates the model’s prediction process for a number of factual relations. In this work we study LREs under the same definition and experimental setup, because it allows us to predefine the concepts we want to search for (e.g., factual relations), as well as use a handful of representations to relate thousands of terms in the dataset by learning linear representations on a per-relation level.

Hernandez et al. calculate LREs to approximate an LM’s computation as a first-order Taylor Series approximation. Let $F(\mathbf{s}, c) = \mathbf{o}$ be the forward pass through a model that produces object representation \mathbf{o} given subject representation \mathbf{s} and a few-shot context c , this computation is approximated as $F(\mathbf{s}, c) \approx W\mathbf{s} + b = F(\mathbf{s}_i, c) + W(\mathbf{s} - \mathbf{s}_i)$ where we approximate the relation about a specific subject \mathbf{s}_i . Hernandez et al. propose to compute W and b using the average of n examples from the relation ($n=8$ here) with $\frac{\partial F}{\partial \mathbf{s}}$ representing the Jacobian Matrix of F :

$$W = \mathbb{E}_{\mathbf{s}_i, c_i} \left[\frac{\partial F}{\partial \mathbf{s}} \Big|_{(\mathbf{s}_i, c_i)} \right] \quad \text{and} \quad b = \mathbb{E}_{\mathbf{s}_i, c_i} \left[F(\mathbf{s}, c) - \frac{\partial F}{\partial \mathbf{s}} \mathbf{s} \Big|_{(\mathbf{s}_i, c_i)} \right] \quad (1)$$

In practice, LREs are estimated using hidden states from LMs during the processing of the test example in a few-shot setup. For a relation like “instrument-played-by-musician”, the model may see four examples (in the form “X plays the Y”) and on the fifth example, when predicting e.g., “trumpet” from “Miles Davis plays the”, the subject representation \mathbf{s} and object representation \mathbf{o} are extracted.

2.2 INFERRING TRAINING DATA FROM MODELS

There has been significant interest recently in understanding the extent to which it is possible to infer the training data of a fully trained neural network, including LMs, predominantly by performing membership inference attacks (Shokri et al., 2017; Carlini et al., 2022), judging memorization of text (Carlini et al., 2023; Oren et al., 2024; Shi et al., 2024), or inferring the distribution of data sources (Hayase et al., 2024b; Ateniese et al., 2015; Suri & Evans, 2022). Our work is related in that we find hints of the pretraining data distribution in model itself, but focus on how linear structure in the representations relates. Carlini et al. (2024); Finlayson et al. (2024) do not focus on extracting dataset information, but on inferring information architectural information about a black-box model behind an API.

3 METHODS

Our analysis is twofold: counts of terms in the pretraining corpus of LMs, and measurements of how well factual relations are approximated by affine transformations. We use the OLMo model v1.7 (0424 7B and 0724 1B) (Groeneveld et al., 2024) and GPT-J (6B) (Wang & Komatsuzaki, 2021) and their corresponding datasets: Dolma (Soldaini et al., 2024) and the Pile (Gao et al., 2020), respectively. To understand how these features form over training time, we test 8 model checkpoints throughout training in the OLMo family of models (Groeneveld et al., 2024).

3.1 LINEAR RELATIONAL EMBEDDINGS (LREs)

The original Relations dataset includes factual, commonsense, gender bias, and linguistic relations, but we reduce this set to the 25 factual relations used by Hernandez et al. (2024)². These are relations such as capital-city and person-mother (full list in Appendix A). The reason for this is due to the way we count occurrences of a relation in training data not being accurate for non-factual relations (see §3.2). Across these relations there are 10,488 unique subjects and objects. Following Hernandez et al. (2024), we fit an LRE for each relation on 8 examples from that relation, each with a 5 shot prompt. We use the approach from this work as described in Section 2.1.

Fitting LREs Hernandez et al. (2024) find that Equation 1 underestimates the optimal slope of the linear transformation, so they scale each relation’s W by a scalar hyperparameter β . Unlike the original work, which finds one β per model, we use one β per relation, as this avoids disadvantaging specific relations. Another difference in our calculation of LREs is that we do not impose the constraint that the model has to predict the answer correctly to be used as one of the 8 examples used to approximate the Jacobian Matrix. Interestingly, using examples that models predict incorrectly to fit Equation 1 works as well as using only correct examples. We opt to use this variant as it allows us to compare different checkpoints and models (§4) with linear transformations trained on the same 8

²For the analysis, we drop “Landmark on Continent” because 74% of the answers are Antarctica, making it uninteresting for studying true relational knowledge.

examples, despite the fact that the models make different predictions on these instances. We explore the effect of example choice in Appendix A.

Metrics To evaluate the quality of LREs, (Hernandez et al., 2024) introduce two metrics that measure the quality of the learned transformations. **Faithfulness** measures whether the transformation learned by the LRE produces the same object token prediction as the original LM. **Causality** measures the proportion of the time a prediction of an object can be changed to the output of a different example from the relation (e.g., editing the *Miles Davis* subject representation so that the LM predicts he plays the *guitar*, instead of the *trumpet*). For specifics on implementation we refer the reader to Hernandez et al. (2024). We consider an LRE to be high ‘quality’ when it scores highly on these metrics, as this measures when an LRE works across subject-object pairs within the relation. In general, we prefer to use causality in our analysis, as faithfulness can be high when LMs predict the same token very often (like in early checkpoints).

3.2 COUNTING FREQUENCIES THROUGHOUT TRAINING

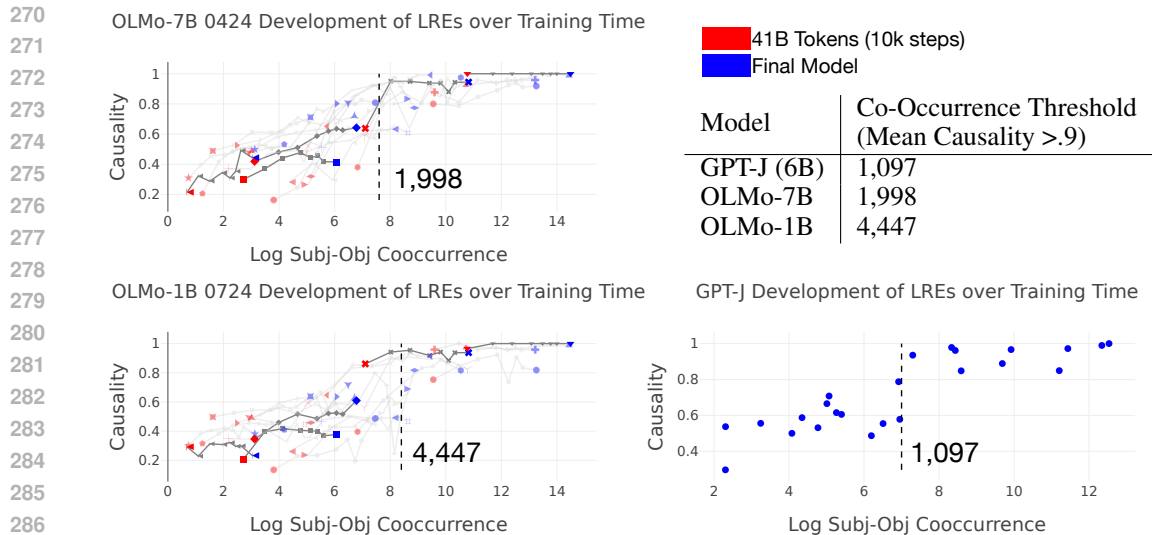
A key question we explore is how term frequencies affect the formation of linear representations. We hypothesize that more commonly occurring relations will lead to higher quality LREs for those relations. Following Elsayar et al. (2018); Elazar et al. (2022), we count an occurrence of a relation when a subject and object co-occur together. While term co-occurrence is used as a proxy for the frequency of the entire triplet mentioned in text, Elsayar et al. (2018) show that this approximation is quite accurate. We now discuss how to compute these co-occurrence counts.

What’s in My Big Data? (WIMBD) Elazar et al. (2024) index many popular pretraining datasets, including Dolma and the Pile, and provide search tools that allows for counting individual terms and co-occurrences within documents. However, this only gives us counts for the full dataset. Since we are interested in counting term frequencies throughout pretraining, we count these within training batches of OLMo instead. When per-batch counts are not available, WIMBD offers a good approximation for final checkpoints, which is what we do in the case of GPT-J. We compare WIMBD co-occurrence counts to the Batch Search method (described below) for the final checkpoint of OLMo in Appendix C, and find that the counts are extremely close.

Batch Search Data counting tools can not typically provide accurate counts for model checkpoints at arbitrary training steps. Thus, we design a tool to efficiently count exact co-occurrences within sequences of tokenized batches. This also gives us the advantage of counting in a way that is highly accurate to how LMs are trained; since LMs are trained on batches of fixed lengths which often split documents into multiple sequences, miscounts may occur unless using tokenized sequences. Using this method, we note every time one of our 10k terms appears throughout a dataset used to pretrain an LM. We count a co-occurrence as any time two terms appear in the same sequence within a batch (a (batch-size, sequence-length) array). We search 10k terms in the approximately 2T tokens of the Dolma dataset (Soldaini et al., 2024) this way. Using our implementation we are able to complete this on 900 CPUs in about a day. To support future work, we release our code as Cython bindings that integrate out of the box with existing libraries.

4 FREQUENCY OF SUBJECT-OBJECT CO-OCCURRENCES ALIGNS WITH EMERGENCE OF LINEAR REPRESENTATIONS

In this section we explore when LREs begin to appear in training time, and how these are related to pretraining term frequencies. Our main findings are that 1) average co-occurrence frequency within a relation strongly correlates with whether an LRE will form; 2) the frequency effect is independent of the pretraining stage; if the average subject-object co-occurrence within a relation surpasses some threshold it is very likely to have a high-quality LRE, even for early pretraining steps. This finding is exclusive to *co-occurrences* rather than individual subject or object occurrences. In addition to confirming dataset frequencies strongly align with LREs forming, we aim to confirm that this relationship is strongest with subject-object **co-occurrences** rather than just mentions of relevant subjects or objects.



288 Figure 2: We find that LREs have consistently high causality scores across relations after some
 289 average frequency threshold is reached (table, top right). In OLMo models, red dots show the
 290 model’s LRE performance at 41B tokens, and blue dots show the final checkpoint performance
 291 (550k steps in 7B). Gray dots show intermediate checkpoints. We highlight Even at very early
 292 training steps, if the average subject-object cooc. count is high enough, the models are very likely to
 293 already have robust LREs formed in the representation space. Symbols represent different relations.
 294 Highlighted relations are shown in darker lines.⁵

295 4.1 SETUP

297 Using the factual recall relations from the Hernandez et al. (2024) dataset, we use the Batch Search
 298 method (§3.2) to count subject and object co-occurrences within sequences in Dolma (Soldaini
 299 et al., 2024) used to train the OLMo 1B (v. 0724) and 7B (v. 0424) models (Groeneveld et al.,
 300 2024). The OLMo family of models provide tools for accurately recreating the batches from Dolma,
 301 which allow us to reconstruct the data the way the model was trained. We also use GPT-J (Wang
 302 & Komatsuzaki, 2021) and the Pile (Gao et al., 2020) as its training data, but since we do not have
 303 access to accurate batches used to train it, we use WIMBD (Elazar et al., 2024) to count s-o counts
 304 in the entire data. We fit LREs on each relation and model separately. Hyperparameter sweeps are
 305 in Appendix B. OLMo also releases intermediate checkpoints, which we use to track development
 306 over pretraining time. We use checkpoints that have seen {41B, 104B, 209B, 419B, 628B, 838B, 1T,
 307 and 2T} tokens³. We use the Pearson coefficient for measuring correlation unless other specified.

309 4.2 RESULTS

311 Our results are summarized in Figure 2. We report training tokens because the step count differs
 312 between 7B and 1B. Co-occurrence frequencies highly correlate with causality ($r=.82$). This is
 313 notably higher than the correlations with subject frequencies: $r=.66$, and object frequencies: $.59$ for
 314 both OLMo 7B and OLMo 1B, respectively.

315 We consider a causality score above $.9$ to be nearly perfectly linear. The table in Figure 2 shows the
 316 co-occurrence counts above which the average causality is above $.9$ and is shown by dashed black
 317 lines on the scatterplots. Regardless of pretraining step, models that surpass this threshold have very
 318 high causality scores. Although we can not draw conclusions from only three models, it is possible
 319 that scale also affects this threshold: OLMo 7B and GPT-J (6B params) require far less exposure
 320 than OLMo 1B.

322 ³In OLMo 7B 0424, this corresponds to 10k, 25k, 50k, 100k, 150k, 200k, 250k, 409k pretraining steps

323 ⁵These are: ‘country largest city’, ‘country currency’, ‘company hq’, ‘company CEO’, and ‘star constella-
 tion name’ in order from best to worst performing final checkpoints.

4.3 RELATIONSHIP TO ACCURACY

Increased frequency (or a proxy for it) is shown to lead to better factual recall in LMs (Chang et al., 2024; Mallen et al., 2023b). However, it remains unknown whether high accuracy entails the existence of a linear relationship. Such a finding would inform when we expect an LM to achieve high accuracy on a task. We find that the correlation between causality and subject-object frequency is higher than with 5-shot accuracy (.82 v.s. .74 in OLMo 7B), though both are clearly high. In addition, there are a few examples of high accuracy relations that do not form single consistent LREs. These relations are typically low frequency, such as star constellation name, which has 84% 5-shot accuracy but only 44% causality (OLMo 7B), with subjects and objects only co-occurring about 21 times on average across the full dataset. In general, few-shot accuracy closely tracks causality, consistent with arguments that in-context learning allows models to identify linear mappings between input-output pairs (Hendel et al., 2023; Garg et al., 2022). We find that causality increases first in some cases, like “food from country” having a causality of 65% but a 5-shot accuracy of only 42%. This gap is consistently closed through training. In the final model, causality and 5-shot accuracy is within 11% on average. We report the relationship between every relation, zero-shot, and few-shot accuracy for OLMo models across training in Appendix E.

A fundamental question in the interpretability community is why linear structures form. While previous work has claimed that the training objective encourages this type of representation (Jiang et al., 2024), our results suggest that the reason why some concepts form a linear representation while others do not, is strongly related to the pretraining frequency.

5 LINEAR REPRESENTATIONS HELP PREDICT PRETRAINING CORPUS FREQUENCIES

In this section, we aim to understand this relationship further by exploring what we can understand about pretraining term frequency from linearity of LM representations. We target the challenging problem of predicting how often a term, or co-occurrence of terms, appears in an LM’s training data from the representations alone. Such prediction model can be useful, if it generalizes, when applied to other models whose weights are open, but the data is closed. For instance, such predictive model could tell us whether a model was trained on specific domains (e.g., Java code) by measuring the presence of relevant LREs. First, we show that LRE features encode information about frequency that is not present using probabilities alone. Then, we show how a regression fit on one model generalizes to the features extracted from another without any information about the new model’s counts.

5.1 EXPERIMENTAL SETUP

We train a random forest regression model with 100 decision tree estimators to predict the frequency of terms (either the subject-object frequency, or the object frequency alone; e.g., predicting “John Lennon” and “The Beatles” or just “The Beatles”) from one of two sets of features. Our baseline set of features is based on likelihood of recalling a fact. Given some few-shot context from the relations dataset (“John Lennon is a lead singer of”) we extract the log probability of the correct answer, as well as the average accuracy on this prompt across 5 trials. The intuition is that models will be more confident about highly frequent terms. The other set of features include the first, as well as faithfulness and causality measurement.

We use Faithfulness and Causality as defined in Hernandez et al. (2024) as well as two other metrics: **Faith Prob.**, which is the log probability of the correct answer as produced by an LRE, and **Hard Causality**, which is the same as the “soft” variant, but only counts the proportion of times the causality edit produces the target answer as the number one prediction. We use every example from the Relations for which there are more than 1 object occurrence or subject-object co-occurrence. We drop the “Landmark in Continent” relation because it is too imbalanced.⁶ We do not provide an explicit signal for which relation an example comes from, but due to the bias of subjects/objects having similar frequencies within a relation, we train multiple models and evaluate on held out

⁶Most answers are “Antarctica” which was artificially inflating our results.

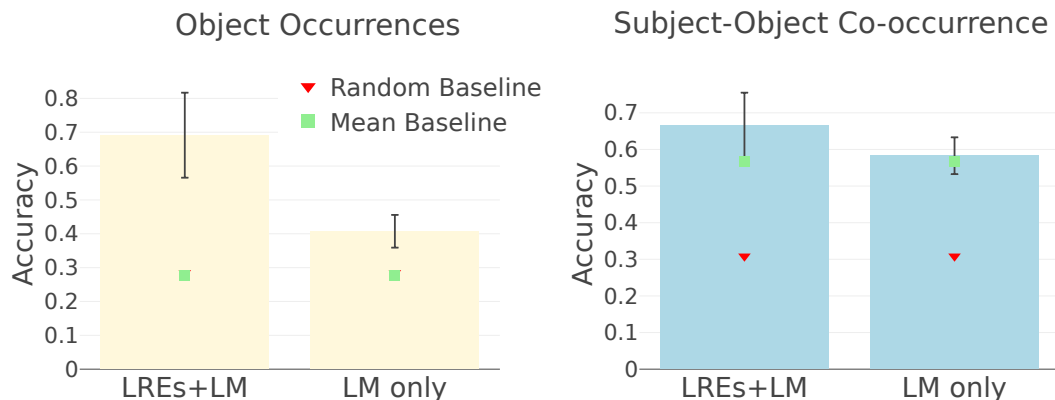


Figure 3: Within-Magnitude accuracy (aka the proportion of predictions within one order of magnitude of ground truth) for models predicting object and subject-object co-occurrences in heldout relations. Using LRE features outperforms LM only features by about 30%. We find that it is much easier to predict object frequencies; the subj-obj. prediction models with LRE features only marginally outperform baseline performance.

relations and average performance. In all settings, the held out set objects are guaranteed to not have been in the training set.

5.2 LRE METRICS ENCODE FINE-GRAINED FREQUENCY INFORMATION

We fit a regression to the Relations dataset (Hernandez et al., 2024) using OLMo-7B LRE features and log probabilities. We fit 24 models such that each relation is held out once per random seed across 4 seeds. Because of the difficulty of predicting the exact number of occurrences, we report accuracy within one order of magnitude of the ground truth. This measures whether the predicted value is within a reasonable range of the actual value. Results are shown in Figure 3. We find that language modeling features do not provide any meaningful signal towards predicting object or subject-object frequencies, and are only marginally above the baseline of predicting the average or random frequencies from the training data. On object frequency predictions, we find that LRE features encode a strong signal allowing for accurate predictions about 70% of the time. Mean absolute error of the predictions (in natural log space) for LRE features (LM-only features) are 2.1, (4.2) and 1.9, (2.3) on object prediction and subject-object predictions tasks, respectively. We find that subject-object co-occurrence frequency is likely too difficult to predict given the signals that we have here, as our predictions are higher than, but within one standard deviation of the mean baseline.

Feature Importance: How important are LRE features for predicting the frequency of an item? We perform feature permutation tests to see how much each feature (LRE features and log probs) contributes to the final answer. First, we check to see which features used to fit the regression are correlated, as if they are, then perturbing one will leave the signal present in another. In Appendix D, we show that only faithfulness and faith probability are strongly correlated, so for this test only, we train models with a single PCA component representing 89% of the variance of those two features. We find that hard causality is by far the most important feature for generalization performance, causing a difference of about 15% accuracy, followed by faithfulness measures with 5% accuracy, providing evidence that the LRE features are encoding an important signal.

5.3 GENERALIZATION TO A NEW LM

In this section, we test the ability to generalize the regression fit on one LM to another for which we do not have access to pretraining term counts, without requiring further supervision. We keep the objective the same and apply the regression model, fit for example on OLMo (“Train OLMo” setting), to features extracted from GPT-J, using ground truth counts from The Pile (or vice versa, i.e., the “Train GPT-J” setting).

	Predicting Object Occs.		Predicting Subject-Object Co-Occs.	
	Eval. on GPT-J	Eval. on OLMo	Eval. on GPT-J	Eval. on OLMo
LRE Features	0.65±0.12	0.49±0.12	0.76±0.12	0.68±0.08
LogProb Features	0.42±0.10	0.41±0.09	0.66±0.09	0.60±0.07
Mean Freq. Baseline	0.31±0.15	0.41±0.17	0.57±0.15	0.67±0.16

Table 1: Within-Magnitude accuracy for different settings of train and test models. Overall, we find that fitting a regression on one model’s LREs and evaluating on the other provides a meaningful signal compared to fitting using only log probability and task performance, or predicting the average training data frequency. The metric here is proportion of predictions within one order of 10x the ground truth. Here, Eval. on GPT-J means the regression is fit on OLMo and evaluated on GPT-J.

Predicting Object Frequency in GPT-J, Regression fit on OLMo					
Relation	Subject	Object	Prediction	Ground Truth	Error
landmark-in-country	Menangle Park	Australia	2,986,989	3,582,602	1.2x
country-language	Brazil	Portuguese	845,406	561,005	1x
star-constellation name	Arcturus	Boötes	974,550	2,817	346x
person-mother	Prince William	Princess Diana	5,826	27,094	4.6x
person-mother	Prince Harry	Princess Diana	131	27,094	207x

Table 2: Examples of a regression fit on OLMo LRE metrics and evaluated on GPT-J on heldout relations, demonstrating common error patterns: 1. Predictions are better for relations that are closer to those found in fitting the relation (country related relations), 2. Some relations, like star-constellation perform very poorly, possibly due to low frequency, 3. the regression model can be sensitive to the choice of subject (e.g., William vs. Harry), telling us the choice of data to measure LREs for is important for predictions.

We again train a random forest regression model to predict the frequency of terms (either the subject-object frequency, or the object frequency alone; e.g., predicting “John Lennon” and “The Beatles” or just “The Beatles”) on features from one of two models: either OLMo 7B (final checkpoint) or GPT-J, treating the other as the ‘closed’ model. We test the hypothesis that LRE features (**faithfulness**, **causality**) are useful in predicting term frequencies across different models, with the hope that this could be applied to dataset inference methods in the future, where access to the ground truth pretraining data counts is limited or unavailable.

Results Our results are presented in Table 1. First, we find that there is a signal in the LRE features that does not exist in the log probability features: We are able to fit a much better generalizable model when using LRE features as opposed to the LM probabilities alone. Second, evaluating on the LRE features of a heldout model (scaled by the ratio of total tokens trained between the two models) maintains around the same accuracy when fit on exact counts from OLMo, allowing us to predict occurrences without access to the GPT-J pretraining data. We find that predicting either the subject-object co-occurrences or object frequencies using LREs alone is barely better than the baseline. This task is much more difficult than predicting the frequency of the object alone, but our model may just also be unable to account for outliers in the data, which is tightly clustered around the mean (thus giving the high mean baseline performance of between approx. 60-70%). Nevertheless, we show that linearity of features within LM representations encode a rich signal representing dataset frequency.

5.4 ERROR ANALYSIS

In Table 2 we show example predictions from a regression model fit on OLMo evaluated on heldout relations with LREs measured on GPT-J. We find that some relations transfer more easily than others, with the star constellation name transferring especially poorly. In general, the regression transfers well, without performance deteriorating much (about 5% accuracy: see Figure 3 compared to the evaluation of GPT-J in Table 1), suggesting LREs are encoding information in a consistent way across models. We also find that the regression makes use of the full prediction range, producing values in the millions (see Table 2) or in the tens: The same regression shown in the table also predicts 59 occurrences for “Caroline Bright” (Will Smith’s mother) where the ground truth is 48.

6 DISCUSSION

Connection to Factual Recall Work in interpretability has focused largely around linear representations in recent years, and our work aims to address the open question of the conditions in which they form. We find that coherent linear representations form when the relevant terms (in this case subject-object co-occurrences) appear in pretraining at a consistent enough rate. Analogously, Chang et al. (2024) show that repeated exposure encourages higher retention of facts. It isn't clear whether accuracy on factual recall entails that a linear representation exists (at least for some cases) from our work, however future research could study this connection more closely.

Linear Representations in LMs The difficulty of disentangling the formation of linear representations from increases in relation accuracy, especially in the few-shot case, is interesting. Across 24 relations, only the "star constellation name" and "product by company" relations have few shot accuracies that far exceed their causality scores (and both are low frequency). Thus, it is still an open question how LMs are able to recall these tasks. While there is not a single LRE that can solve the relation, it is not necessarily true that the model is preferring a non-linear solution, as multiple incomplete LREs could account for different parts of the data. The fact that few-shot accuracy and causality seem so closely linked is consistent with findings that ICL involves locating the right task (Min et al., 2022) and applying a 'function' to map input examples to outputs (Hendel et al., 2023; Todd et al., 2024). That frequency controls this ability is perhaps unsurprising, as frequency also controls this linear structure emerging in static embeddings (Ethayarajh et al., 2019). Jiang et al. (2024) prove a strong frequency-based condition (based on matched log-odds between subjects and objects) and an implicit bias of gradient descent (when the frequency condition is not met) encourage linearity in LLMs; our work empirically shows conditions where linear representations tend to form in more realistic settings. If LMs are 'only' solving factual recall or performing ICL through linear structures, it is surprising how well this works at scale, but the simplicity also provides a promising way to understand LMs and ICL in general. An interesting avenue for future work would be to understand if and when LMs use a method that is not well approximated linearly to solve these types of tasks, as non-linear representations, as recent work has shown non-linearity can be preferred for some tasks in recurrent networks (Csordás et al., 2024).

Future Work in Predicting Dataset Frequency The ability to predict the contents of pretraining data is an important area for investigating memorization, contamination, and privacy of information used to train models. In our approach, we show it's possible to extract signal without supervision by first fitting on an opens source model. The fact that there is some transferable signal between models is indicative that this relationship between pretraining frequency and linearity is consistent between models. Mosbach et al. (2024) discuss the role of interpretability on the broader field of NLP. Without interpretability work on the nature of representations in LMs, we would not know of this implicit dataset signal, and we argue that interpretability can generate useful insights more broadly as well. Extensions on this work could include more information to tighten the prediction bounds on frequency, such as extracting additional features from the tokenizer (Hayase et al., 2024a). A likely candidate task that could integrate our method is for predicting whether a certain domain of data (e.g., code) was included in pretraining, since extensive exposure would lead to LREs forming. Regardless, we hope this work encourages future research in other ways properties of pretraining data affect LM representations for both improving and better understanding these models.

7 CONCLUSION

We find a connection between linear representations of subject-relation-object factual triplets in LMs and the pretraining frequencies of the subjects and objects in those relations. This finding can guide future interpretability work in deciphering whether a linear representation for a given concept will exist in a model, since we observe that frequencies below a certain threshold for a given model will not yield LREs (a particular class of linear representation). From there we show that we can use the presence of linear representations to predict with some accuracy, the frequency of terms in the pretraining corpus of a closed-data model without supervision. Future work could aim to improve on our bounds of predicted frequencies. Overall, our work presents a meaningful step towards understanding the interactions between pretraining data and internal LM representations.

8 LIMITATIONS

While our approach thoroughly tracks exposure to individual terms and formation of LRE features across pretraining, we can not draw causal claims about how exposure affects individual representations, due to the cost of counterfactual pretraining. We try to address this by showing the frequency of individual terms can be predicted with some accuracy from measurements of LRE presence. We motivate this approach as a possible way to detect the training data of closed-data LMs, however, we are not able to make any guarantees on its efficacy in settings not shown here, and would caution drawing strong conclusions without additional information. Furthermore, we find that our method is relatively worse at predicting subject-object co-occurrences than object occurrences, and our method fails to account for the harder task. Future work could expand on this tool by incorporating it with other data inference methods for greater confidence. We also do not discuss the role of the presentation of facts on the formation of LRE features, but following Elsahar et al. (2018) and the strength of the relationship we find, we speculate this has minimal impact. Note that the BatchSearch tool we release tracks the exact position index of the searched terms, thus facilitating future work on questions about templates/presentation of information.

REFERENCES

- Antonis Antoniadou, Xinyi Wang, Yanai Elazar, Alfonso Amayuelas, Alon Albalak, Kexun Zhang, and William Yang Wang. Generalization v.s. memorization: Tracing language models’ capabilities back to pretraining data. *ArXiv*, abs/2407.14985, 2024.
- Viraat Aryabumi, Yixuan Su, Raymond Ma, Adrien Morisot, Ivan Zhang, Acyr Locatelli, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. To code, or not to code? exploring impact of code in pre-training. *arXiv preprint arXiv:2408.10914*, 2024.
- Giuseppe Ateniese, Luigi V Mancini, Angelo Spognardi, Antonio Villani, Domenico Vitali, and Giovanni Felici. Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers. *International Journal of Security and Networks*, 10(3):137–150, 2015.
- Sid Black, Lee Sharkey, Leo Grinsztajn, Eric Winsor, Dan Braun, Jacob Merizian, Kip Parker, Carlos Ramón Guevara, Beren Millidge, Gabriel Alfour, and Connor Leahy. Interpreting neural networks through the polytope lens, 2022. URL <https://arxiv.org/abs/2211.12312>.
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr. Membership inference attacks from first principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1897–1914, 2022. doi: 10.1109/SP46214.2022.9833649.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=TatRHT_1cK.
- Nicholas Carlini, Daniel Paleka, Krishnamurthy Dj Dvijotham, Thomas Steinke, Jonathan Hayase, A Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conmy, et al. Stealing part of a production language model. In *Forty-first International Conference on Machine Learning*, 2024.
- Hoyeon Chang, Jinho Park, Seonghyeon Ye, Sohee Yang, Youngkyung Seo, Du-Seong Chang, and Minjoon Seo. How Do Large Language Models Acquire Factual Knowledge During Pretraining? 2024. URL <http://arxiv.org/abs/2406.11813>.
- David Chanin, Anthony Hunter, and Oana-Maria Camburu. Identifying Linear Relational Concepts in Large Language Models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 1524–1535. Association for Computational Linguistics, 2024. doi: 10.18653/v1/2024.naacl-long.85. URL <https://aclanthology.org/2024.naacl-long.85>.

- 594 Róbert Csordás, Christopher Potts, Christopher D. Manning, and Atticus Geiger. Recurrent neural
595 networks learn to store and generate sequences using non-linear representations, 2024. URL
596 <https://arxiv.org/abs/2408.10920>.
597
- 598 Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. Amnesic Probing: Behavioral
599 Explanation with Amnesic Counterfactuals. *Transactions of the Association for Computational*
600 *Linguistics*, 9:160–175, 03 2021. URL https://doi.org/10.1162/tacl_a_00359.
- 601 Yanai Elazar, Nora Kassner, Shauli Ravfogel, Amir Feder, Abhilasha Ravichander, Marius Mos-
602 bach, Yonatan Belinkov, Hinrich Schütze, and Yoav Goldberg. Measuring causal effects of data
603 statistics on language model’s factual’ predictions. *arXiv preprint arXiv:2207.14251*, 2022.
604
- 605 Yanai Elazar, Akshita Bhagia, Ian Helgi Magnusson, Abhilasha Ravichander, Dustin Schwenk,
606 Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hannaneh Ha-
607 jishirzi, Noah A. Smith, and Jesse Dodge. What’s in my big data? In *The Twelfth International*
608 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=RvfPnOkPV4)
609 [id=RvfPnOkPV4](https://openreview.net/forum?id=RvfPnOkPV4).
- 610 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,
611 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep
612 Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt,
613 Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and
614 Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*,
615 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- 616 Hady Elsahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon Hare, Fred-
617 erique Laforest, and Elena Simperl. T-REx: A large scale alignment of natural language with
618 knowledge base triples. In Nicoletta Calzolari, Khalid Choukri, Christopher Cieri, Thierry De-
619 clerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène
620 Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga (eds.), *Proceed-*
621 *ings of the Eleventh International Conference on Language Resources and Evaluation (LREC*
622 *2018)*, Miyazaki, Japan, May 2018. European Language Resources Association (ELRA). URL
623 <https://aclanthology.org/L18-1544>.
- 624 Kawin Ethayarajh, David Duvenaud, and Graeme Hirst. Towards Understanding Linear Word
625 Analogies. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the*
626 *57th Annual Meeting of the Association for Computational Linguistics*, pp. 3253–3262. As-
627 sociation for Computational Linguistics, 2019. doi: 10.18653/v1/P19-1315. URL [https:](https://aclanthology.org/P19-1315)
628 [//aclanthology.org/P19-1315](https://aclanthology.org/P19-1315).
- 629 Matthew Finlayson, Swabha Swayamdipta, and Xiang Ren. Logits of api-protected llms leak pro-
630 prietary information. *arXiv preprint arXiv:2403.09539*, 2024.
631
- 632 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
633 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text
634 for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 635 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
636 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint*
637 *arXiv:2406.04093*, 2024.
638
- 639 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn
640 in-context? a case study of simple function classes. *Advances in Neural Information Processing*
641 *Systems*, 35:30583–30598, 2022.
- 642 Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuoka. Analogy-based detection of morphologi-
643 cal and semantic relations with word embeddings: what works and what doesn’t. In *Proceedings*
644 *of the NAACL Student Research Workshop*, pp. 8–15, 2016.
645
- 646 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,
647 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the
science of language models. *arXiv preprint arXiv:2402.00838*, 2024.

- 648 Jonathan Hayase, Alisa Liu, Yejin Choi, Sewoong Oh, and Noah A. Smith. Data mixture inference:
649 What do bpe tokenizers reveal about their training data?, 2024a. URL [https://arxiv.org/
650 abs/2407.16607](https://arxiv.org/abs/2407.16607).
651
- 652 Jonathan Hayase, Alisa Liu, Yejin Choi, Sewoong Oh, and Noah A. Smith. Data mixture inference:
653 What do bpe tokenizers reveal about their training data?, 2024b. URL [https://arxiv.org/
654 abs/2407.16607](https://arxiv.org/abs/2407.16607).
655
- 656 Roe Hendel, Mor Geva, and Amir Globerson. In-Context Learning Creates Task Vectors. In
657 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational
658 Linguistics: EMNLP 2023*, pp. 9318–9333. Association for Computational Linguistics, 2023.
659 doi: 10.18653/v1/2023.findings-emnlp.624. URL [https://aclanthology.org/2023.
660 findings-emnlp.624](https://aclanthology.org/2023.findings-emnlp.624).
661
- 662 Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas,
663 Yonatan Belinkov, and David Bau. Linearity of Relation Decoding in Transformer Language
664 Models. 2024. URL <https://openreview.net/forum?id=w7LU2sl4kE>.
665
- 666 Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse
667 Autoencoders Find Highly Interpretable Features in Language Models. 2023. URL <https://openreview.net/forum?id=F76bWRSLeK>.
668
- 669 Yibo Jiang, Goutham Rajendran, Pradeep Kumar Ravikumar, Bryon Aragam, and Victor
670 Veitch. On the Origins of Linear Representations in Large Language Models. 2024. URL [https://openreview.net/forum?id=otuTw4Mghk&referrer=
671 %5Bthe%20profile%20of%20Goutham%20Rajendran%5D\(%2Fprofile%3Fid%
672 3D~Goutham_Rajendran1\)](https://openreview.net/forum?id=otuTw4Mghk&referrer=%5Bthe%20profile%20of%20Goutham%20Rajendran%5D(%2Fprofile%3Fid%3D~Goutham_Rajendran1)).
673
- 674 Marzena Karpinska, Bofang Li, Anna Rogers, and Aleksandr Drozd. Subcharacter information in
675 japanese embeddings: When is it worth it? In *Proceedings of the Workshop on the Relevance of
676 Linguistic Structure in Neural Architectures for NLP*, pp. 28–37, 2018.
677
- 678 Maximilian Köper, Christian Scheible, and Sabine Schulte im Walde. Multilingual reliability and
679 “semantic” structure of continuous word spaces. In *Proceedings of the 11th international confer-
680 ence on computational semantics*, pp. 40–45, 2015.
681
- 682 Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny
683 Zhou, Jason Wei, Kevin Robinson, David Mimno, et al. A pretrainer’s guide to training data:
684 Measuring the effects of data age, domain coverage, quality, & toxicity. In *Proceedings of the
685 2024 Conference of the North American Chapter of the Association for Computational Linguis-
686 tics: Human Language Technologies (Volume 1: Long Papers)*, pp. 3245–3276, 2024.
687
- 688 Yingwei Ma, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li.
689 At which training stage does code data help LLMs reasoning? In *The Twelfth International
690 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum/
691 id=KIPJKST4gw](https://openreview.net/forum?id=KIPJKST4gw).
692
- 693 Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi.
694 When not to trust language models: Investigating effectiveness of parametric and non-parametric
695 memories. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of
696 the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long
697 Papers)*, pp. 9802–9822, Toronto, Canada, July 2023a. Association for Computational Linguis-
698 tics. doi: 10.18653/v1/2023.acl-long.546. URL [https://aclanthology.org/2023.
699 acl-long.546](https://aclanthology.org/2023.acl-long.546).
700
- 701 Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi.
702 When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-
703 Parametric Memories. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Pro-
704 ceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1:
705 Long Papers)*, pp. 9802–9822. Association for Computational Linguistics, 2023b. doi: 10.18653/
706 v1/2023.acl-long.546. URL <https://aclanthology.org/2023.acl-long.546>.

- 702 Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Language models implement simple Word2Vec-
703 style vector arithmetic. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings*
704 *of the 2024 Conference of the North American Chapter of the Association for Computational*
705 *Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 5030–5047, Mexico
706 City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.
707 naacl-long.281. URL <https://aclanthology.org/2024.naacl-long.281>.
- 708 Tomas Mikolov. Efficient estimation of word representations in vector space. *arXiv preprint*
709 *arXiv:1301.3781*, 2013.
- 710
711 Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Dis-
712 tributed representations of words and phrases and their compositionality. In C.J.
713 Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger (eds.), *Ad-*
714 *vances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.,
715 2013. URL [https://proceedings.neurips.cc/paper_files/paper/2013/](https://proceedings.neurips.cc/paper_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf)
716 [file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf).
- 717 Sewon Min, Xinxin Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
718 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In
719 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.
720 11048–11064, 2022.
- 721 Marius Mosbach, Vagrant Gautam, Tomás Vergara-Browne, Dietrich Klakow, and Mor Geva. From
722 insights to actions: The impact of interpretability and analysis research on nlp. *arXiv preprint*
723 *arXiv:2406.12618*, 2024.
- 724
725 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.
726 Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- 727
728 Yonatan Oren, Nicole Meister, Niladri S. Chatterji, Faisal Ladhak, and Tatsunori Hashimoto. Prov-
729 ing test set contamination in black-box language models. In *The Twelfth International Confer-*
730 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=KS8mIvetg2)
731 [KS8mIvetg2](https://openreview.net/forum?id=KS8mIvetg2).
- 732 Alberto Paccanaro and Geoffrey E Hinton. Learning Hierarchical Structures with Linear Rela-
733 tional Embedding. In *Advances in Neural Information Processing Systems*, volume 14. MIT
734 Press, 2001. URL [https://papers.nips.cc/paper_files/paper/2001/hash/](https://papers.nips.cc/paper_files/paper/2001/hash/814a9c18f5abff398787c9cfcbf3d80c-Abstract.html)
735 [814a9c18f5abff398787c9cfcbf3d80c-Abstract.html](https://papers.nips.cc/paper_files/paper/2001/hash/814a9c18f5abff398787c9cfcbf3d80c-Abstract.html).
- 736 Kiho Park, Yo Joong Choe, and Victor Veitch. The Linear Representation Hypothesis and the Ge-
737 ometry of Large Language Models. 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=UGpGkLzwpP)
738 [UGpGkLzwpP](https://openreview.net/forum?id=UGpGkLzwpP).
- 739
740 Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word
741 representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), *Proceedings*
742 *of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp.
743 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.
744 3115/v1/D14-1162. URL <https://aclanthology.org/D14-1162>.
- 745 Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. Null it out: Guard-
746 ing protected attributes by iterative nullspace projection. In *Proceedings of the 58th Annual Meet-*
747 *ing of the Association for Computational Linguistics*, pp. 7237–7256, Online, July 2020. As-
748 sociation for Computational Linguistics. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/2020.acl-main.647)
749 [2020.acl-main.647](https://www.aclweb.org/anthology/2020.acl-main.647).
- 750 Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. Impact of pretraining term
751 frequencies on few-shot numerical reasoning. In *Findings of the Association for Computational*
752 *Linguistics: EMNLP 2022*, pp. 840–854, 2022.
- 753
754 Yasaman Razeghi, Hamish Ivison, Sameer Singh, and Yanai Elazar. Backtracking mathematical
755 reasoning of language models to the pretraining data. In *NeurIPS Workshop on Attributing Model*
Behavior at Scale, 2023. URL <https://openreview.net/forum?id=EKvqw9k3lC>.

- 756 Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner.
757 Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023.
758
- 759 G. Salton, A. Wong, and C. S. Yang. A vector space model for automatic indexing. *Commun.*
760 *ACM*, 18(11):613–620, November 1975. ISSN 0001-0782. doi: 10.1145/361219.361220. URL
761 <https://doi.org/10.1145/361219.361220>.
- 762 Preethi Seshadri, Sameer Singh, and Yanai Elazar. The bias amplification paradox in text-to-image
763 generation. In *Proceedings of the 2024 Conference of the North American Chapter of the Associ-*
764 *ation for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*,
765 pp. 6367–6384, 2024.
- 766 Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi
767 Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. In
768 *The Twelfth International Conference on Learning Representations*, 2024. URL [https://](https://openreview.net/forum?id=zWqr3MQUNs)
769 openreview.net/forum?id=zWqr3MQUNs.
- 771 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference at-
772 tacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*,
773 pp. 3–18. IEEE, 2017.
- 774 Aviv Slobodkin, Omer Goldman, Avi Caciularu, Ido Dagan, and Shauli Ravfogel. The curious
775 case of hallucinatory (un) answerability: Finding truths in the hidden states of over-confident
776 large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural*
777 *Language Processing*, pp. 3607–3625, 2023.
- 778 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur,
779 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of
780 three trillion tokens for language model pretraining research. *arXiv preprint arXiv:2402.00159*,
781 2024.
- 782 Nishant Subramani, Nivedita Suresh, and Matthew Peters. Extracting Latent Steering Vectors from
783 Pretrained Language Models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio
784 (eds.), *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 566–581. As-
785 sociation for Computational Linguistics, 2022. doi: 10.18653/v1/2022.findings-acl.48. URL
786 <https://aclanthology.org/2022.findings-acl.48>.
- 787 Anshuman Suri and David Evans. Formalizing and estimating distribution inference risks. *Proceed-*
788 *ings on Privacy Enhancing Technologies*, 2022, 2022.
- 789 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian
790 Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham,
791 Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman,
792 Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris
793 Olah, and Tom Henighan. Scaling Monosemanticity: Extracting Interpretable Features
794 from Claude 3 Sonnet. 2024. URL [https://transformer-circuits.pub/2024/](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html)
795 [scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 796 Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau.
797 Function Vectors in Large Language Models. 2024. URL [https://openreview.net/](https://openreview.net/forum?id=AwyxtyMwaG¬eId=6Qv7kx00La)
798 [forum?id=AwyxtyMwaG¬eId=6Qv7kx00La](https://openreview.net/forum?id=AwyxtyMwaG¬eId=6Qv7kx00La).
- 801 Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language
802 Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
803
- 804 Xinyi Wang, Alfonso Amayuelas, Kexun Zhang, Liangming Pan, Wenhui Chen, and William Yang
805 Wang. Understanding reasoning ability of language models from the perspective of reasoning
806 paths aggregation. In *Forty-first International Conference on Machine Learning*, 2024.
- 807 Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang,
808 Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up
809 language model pretraining. *Advances in Neural Information Processing Systems*, 36, 2024.



Figure 4: Average Causality and Faithfulness results across relations depending on if the LRE was fit with correct or incorrect samples. We find no notable difference in the choice of examples.

A EFFECT OF TRAINING ON INCORRECT EXAMPLES

In Hernandez et al. (2024), examples are filtered to ones in which the LM gets correct, assuming that an LRE will only exist once a model has attained the knowledge to answer the relation accuracy (e.g., knowing many country capitals). We find that the choice of examples for fitting LREs is not entirely dependent on the model ‘knowing’ that relation perfectly (i.e., attains high accuracy). This is convenient for our study, where we test early checkpoint models, that do not necessarily have all of the information that they will have seen later in training. In Figure 5, we show faithfulness on relations where the LRE was fit with all, half, or zero correct examples. We omit data for which the model did not get enough incorrect examples. Averages across relations for which we have enough data are shown in Figure 4, which shows that there is not a considerable difference in the choice of LRE samples to train with.

B LRE HYPERPARAMETER TUNING

There are three hyperparameters for fitting LREs: **layer** at which to edit the subject, the **beta** term used to scale the LRE weight matrix, and the **rank** of the pseudoinverse matrix used to make edits for measuring causality. Beta is exclusive to measuring faithfulness and rank is exclusive to causality. We test the same ranges for each as in Hernandez et al. (2024): $[0, 5]$ beta and $[0, \text{full_rank}]$ in for causality at varying intervals. Those intervals are every 2 from $[0, 100]$, every 5 from $[100, 200]$, every 25 from $[200, 500]$, every 50 from $[500, 1000]$, every 250 from $[1000, \text{hidden_size}]$. We perform the hyperparameter sweeps across faithfulness and causality, but we choose the layer to edit based on the causality score. In cases where this is not the same layer as what faithfulness would decide, we use the layer causality chooses, as it would not make sense to train one LRE for each metric. We refer the reader to Hernandez et al. (2024) for more details on the interactions between hyperparameters and the choice of layer. The results of our sweeps on OLMo 7B across layers in Figures 6 and 7 and across beta and rank choices in Figures 8 and 9.

C BATCH SEARCH COUNTS COMPARED TO WIMBD

In Figure 10, we find that What’s in My Big Data (Elazar et al., 2024) match very well to batch search co-occurrences, however, WIMBD tends to overpredict co-occurrences (slope less than 1), due to the sequence length being shorter than many documents, as discussed in the main paper.

D FEATURE CORRELATIONS AND IMPORTANCES

Our feature importance test is shown in Figure 12. This permutation test was done on the heldout data to show which features contribute the most to generalization performance. We use PCA to reduce the faithfulness features to one feature for the purposes of this test. Correlations are shown in Figure 11

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

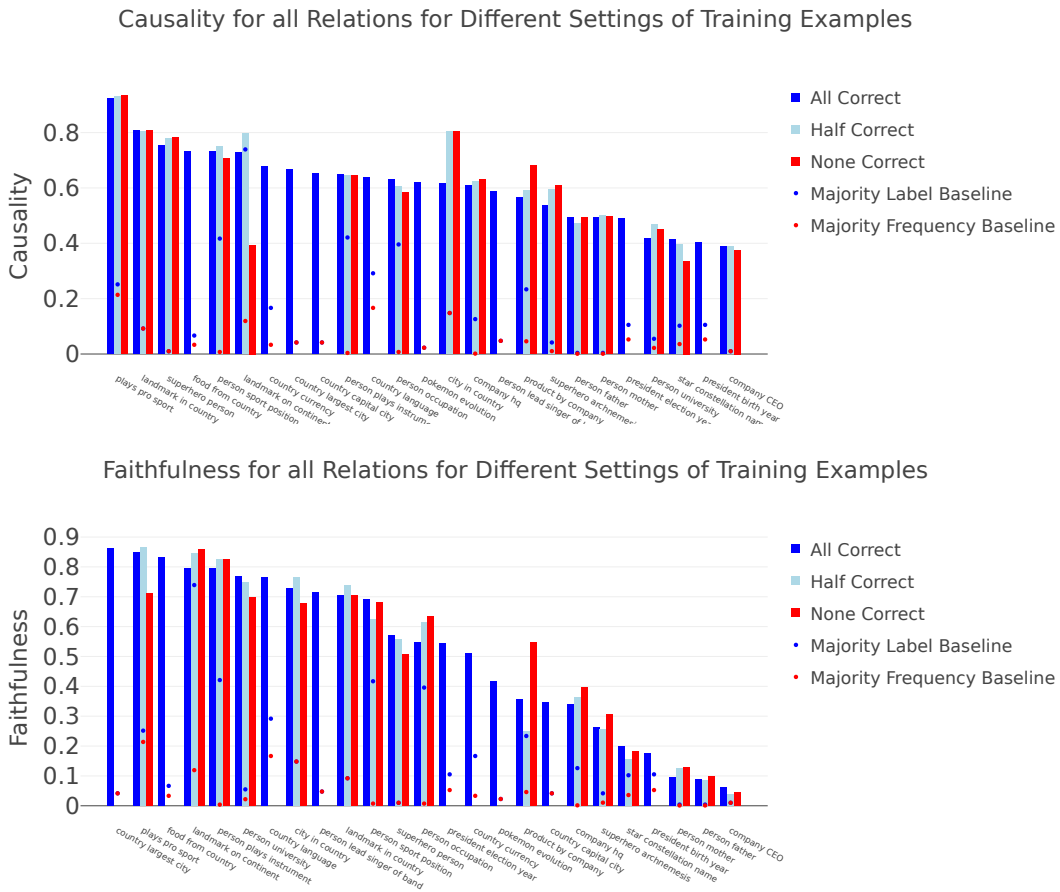


Figure 5: Causality and Faithfulness results for each relation depending on if the LRE was fit with correct or incorrect samples. Note that relations with only one bar do not have zeros in the other categories. It means that there was not enough data that the model (OLMo 7B) got wrong to have enough examples to fit.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

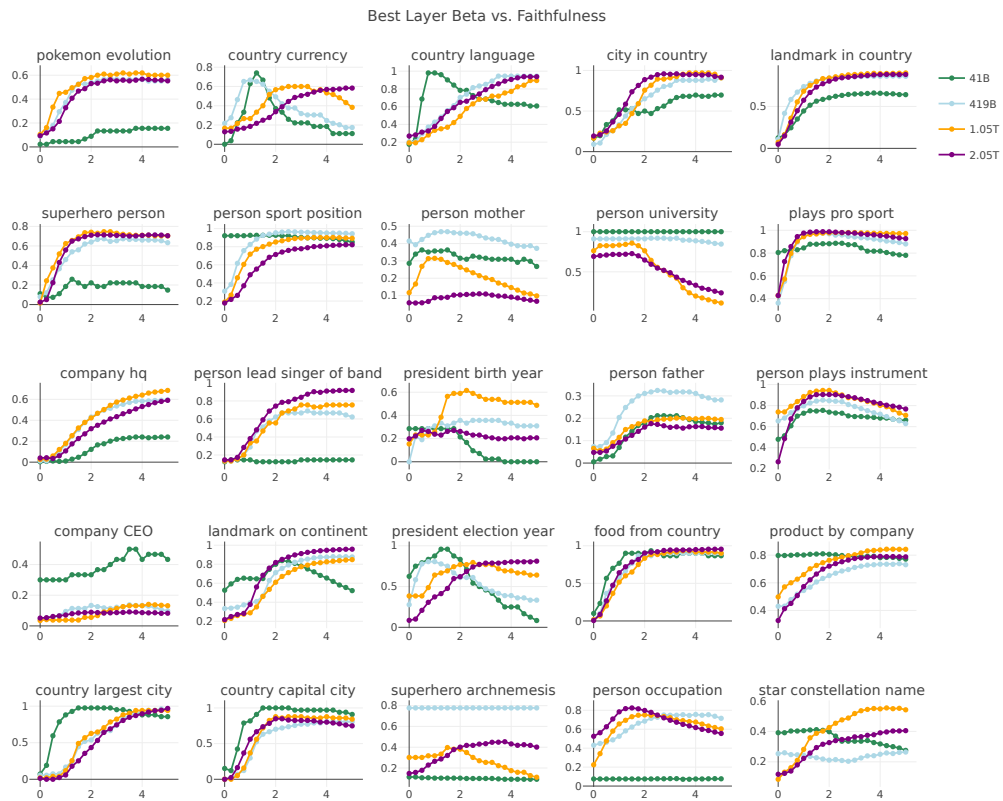


Figure 6: OLMo 0424 7B per layer faithfulness scores as a function of the choice of layer at which to fit the LRE. Note we do not use these results to choose the layer for the LRE, instead preferring the results from the causality sweep.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

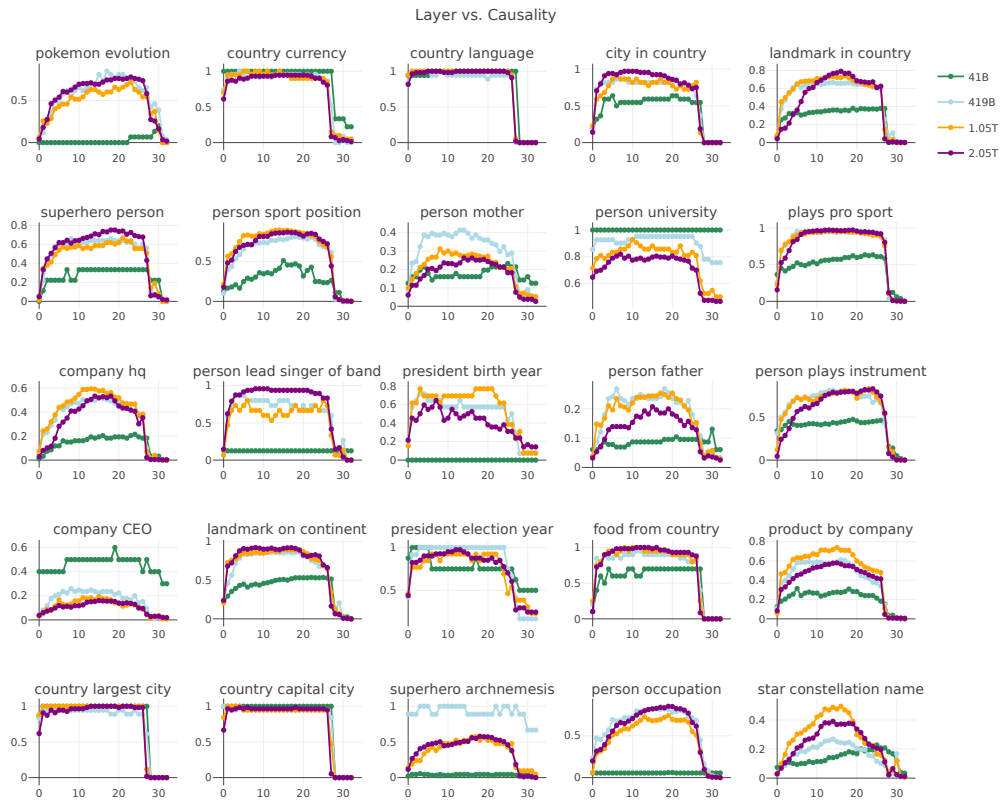


Figure 7: OLMo 0424 7B per layer causality scores as a function of the choice of layer at which to fit the LRE.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

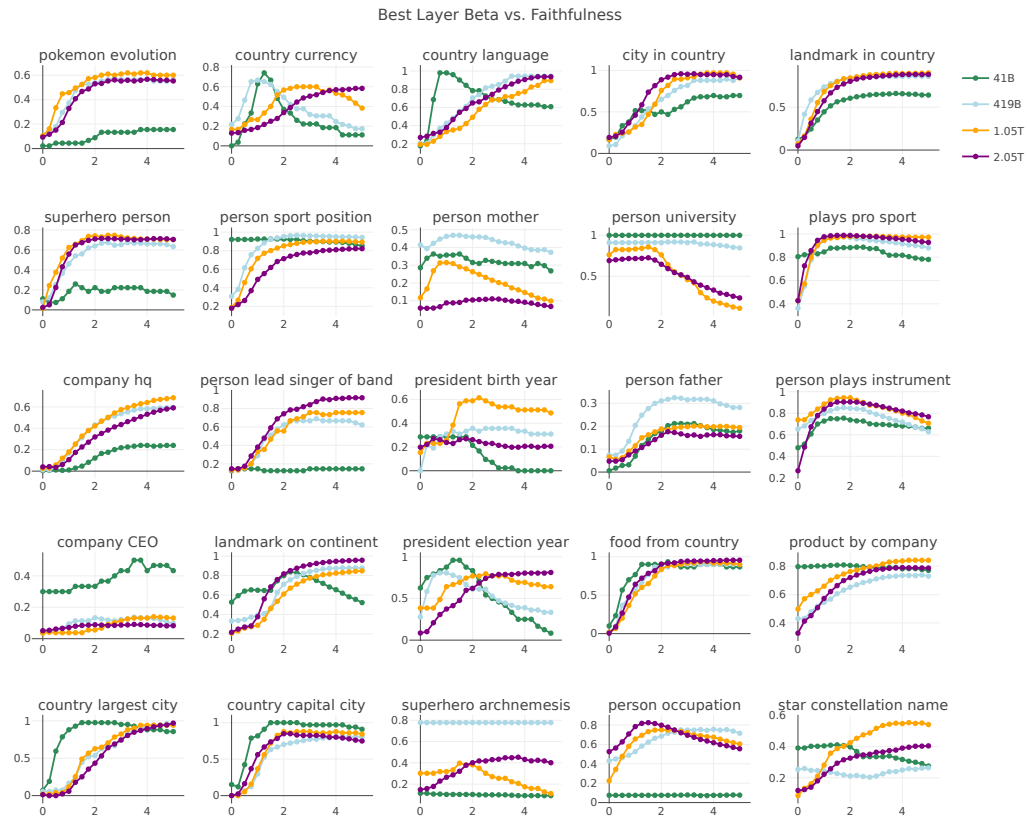


Figure 8: OLMo 0424 7B LRE Beta hyperparameter sweep at highest performing layer.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

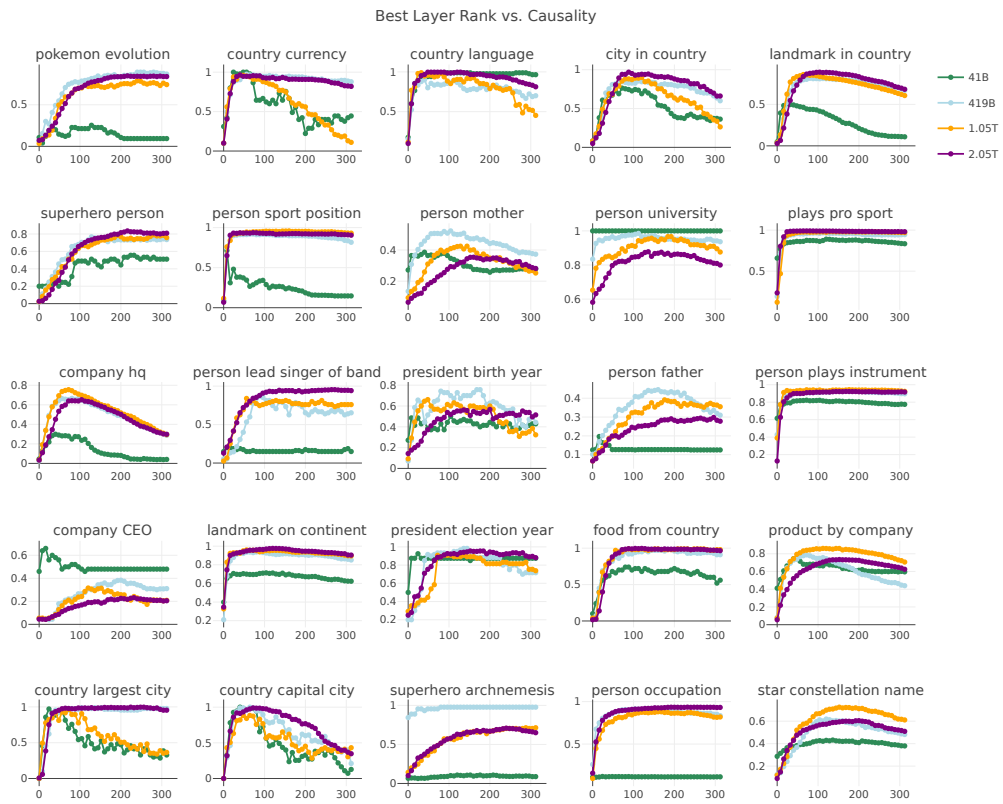


Figure 9: OLMo 0424 7B LRE Rank hyperparameter sweep at highest performing layer.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

WIMBD vs Batch Cooccurrence. slope=0.94, r=0.99

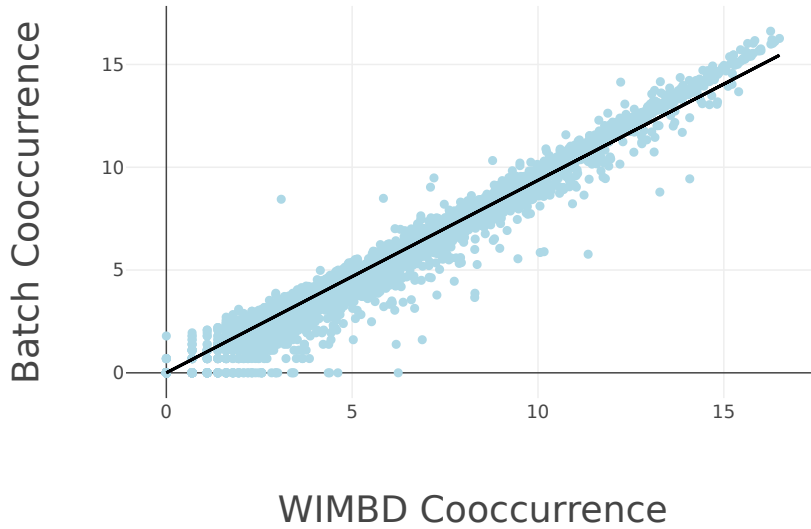


Figure 10: Comparison between WIMBD and Batch Search subject-object co-occurrences

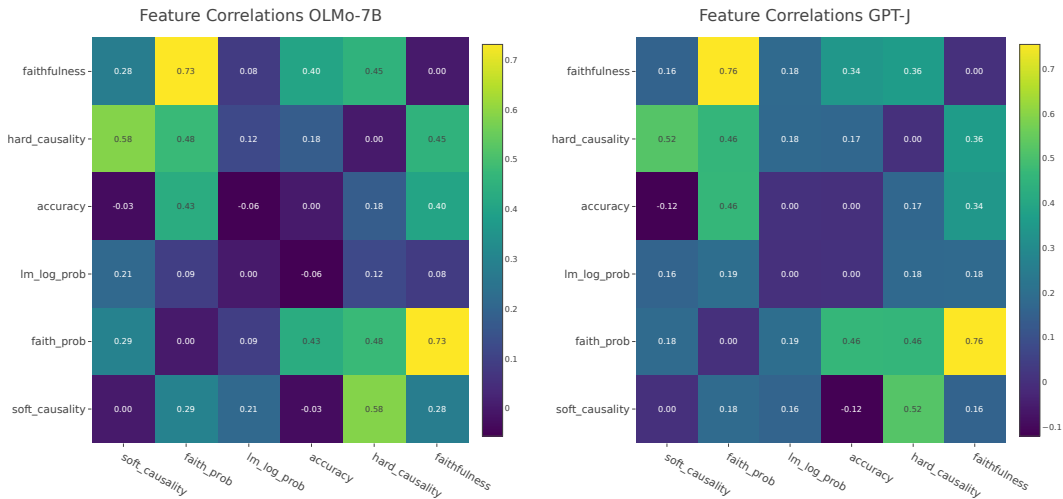
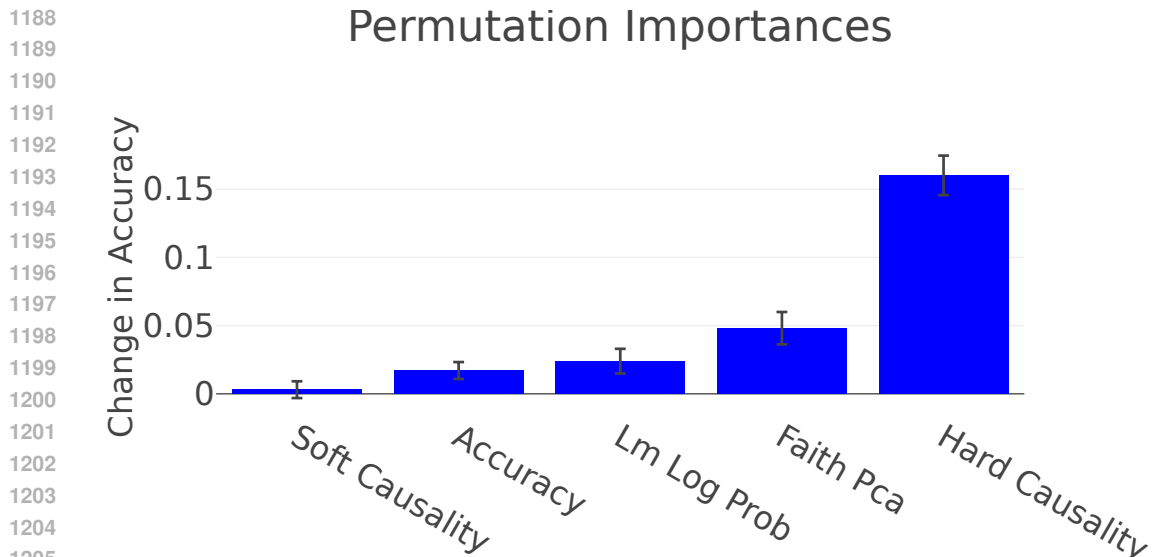


Figure 11: Correlations between each feature in our regression analysis. Because of the high correlation between faithfulness metrics, we use a single dimensional PCA to attain one feature that captures 89% of the variance of both for the purposes of doing feature importance tests. Note that we zero out the diagonal (which has values of 1) for readability.



1207 Figure 12: Hard causality is by far the most important feature for generalizing to new relations when
1208 predicting Object frequencies, causing a change in about 15% accuracy.

1209

1210 E RELATIONSHIP BETWEEN CAUSALITY AND ACCURACY

1211

1212 In this section we provide more detail on the relationship between the formation of linear represen-
1213 tations and accuracy on in-context learning tasks. Although the two are very highly correlated, we
1214 argue that accuracy and LRE formation are somewhat independent.

1215

1216 We show this relationship across training For OLMo 1B in Figure 13 and 7B in Figure 14.

1217

1218 F EXTENDING TO COMMONSENSE RELATIONS

1219

1220 Following Elsahar et al. (2018), we focus on factual relations because subject-object co-occurrences
1221 are shown to be a good proxy for mentions fo the fact. For completeness, we consider 8 additional
1222 commonsense relations here. Results for OLMo 7B are shown in Figure 15. We show that fre-
1223 quency is correlated with causality score (.42) in these cases as well, but it is possible subject-object
1224 frequencies do not accurately track occurrences of the relation being mentioned. For example, in
1225 the “task person type” relation, the co-occurrence count of the subject ”researching history” and
1226 the object “historian” does not convincingly describe all instances where the historian concept is
1227 defined during pretraining. Co-occurrences are perhaps more convincingly related to how a model
1228 learns that the outside of a coconut is brown, however (the fruit outside color relation). Therefore,
1229 we caution treating these under the same lens as the factual relations. Nevertheless, we believe these
1230 results are an interesting perspective on how a different relation family compares to factual relations.

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

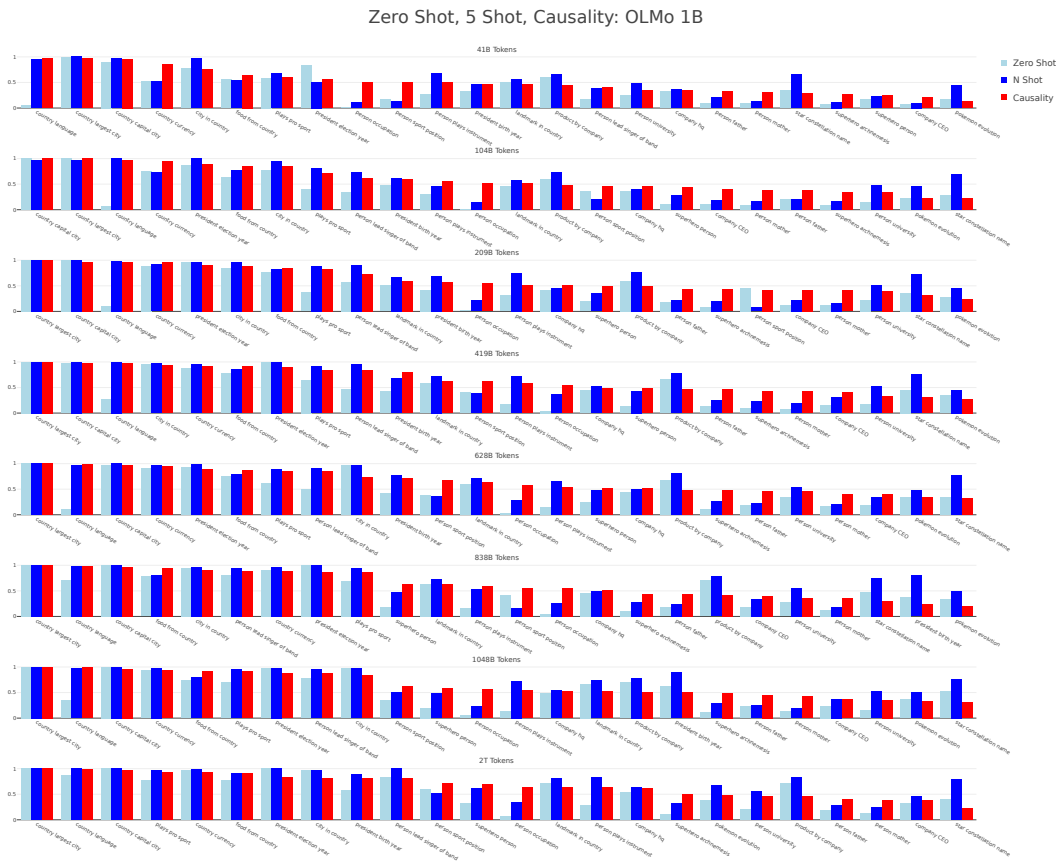


Figure 13: Zero shot, 5-shot accuracies against causality for each relation across training time in OLMo-1B

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

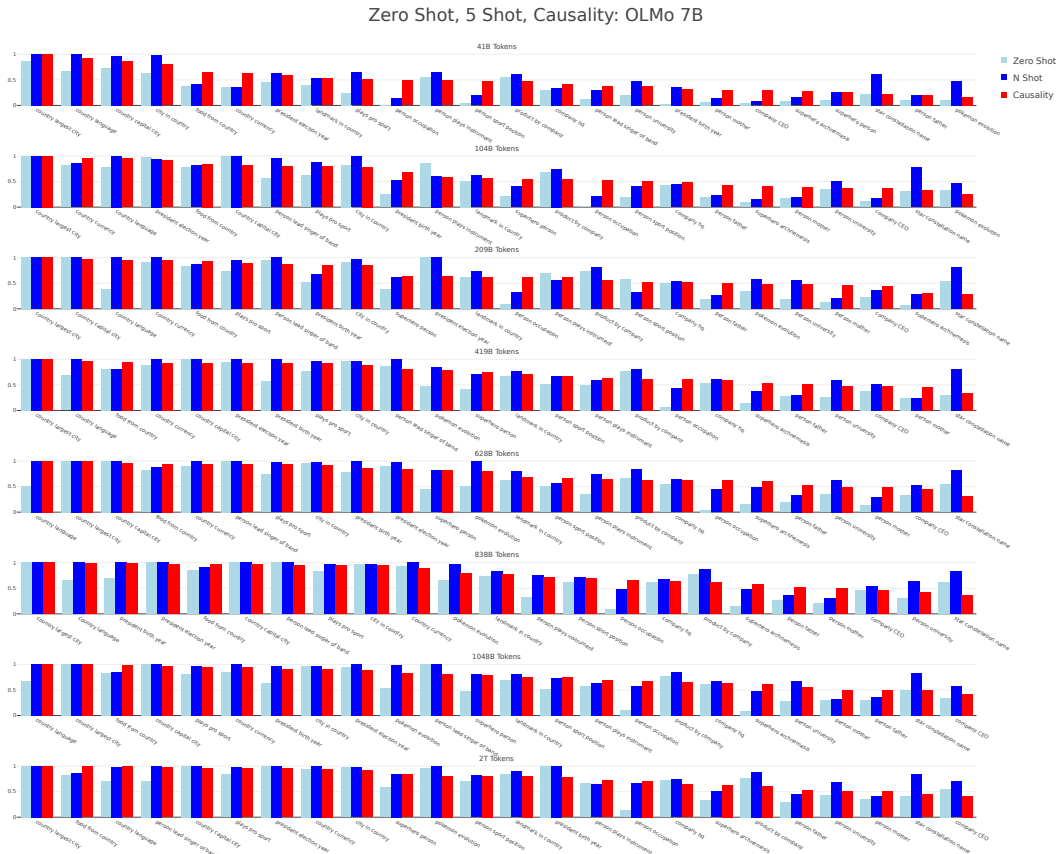


Figure 14: Zero shot, 5-shot accuracies against causality for each relation across training time in OLMo-7B

OLMo-7B 0424 Development of Commonsense LREs over Training Time

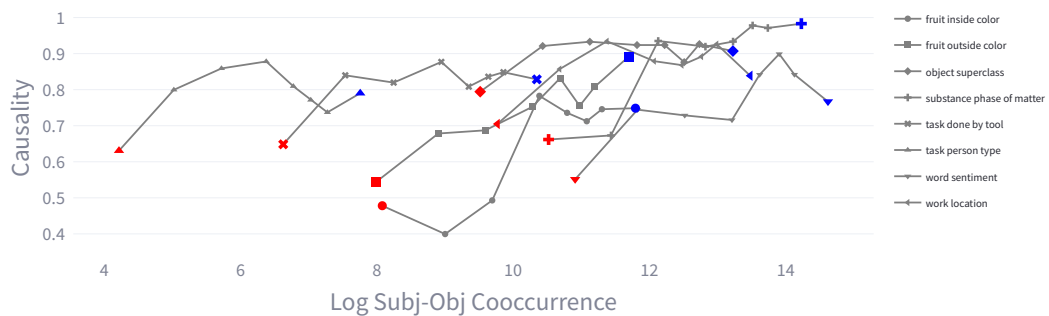


Figure 15: Commonsense relations compared to pretraining time in OLMo-7B.