

Interventional Probing in High Dimensions: An NLI Case Study

Anonymous ACL submission

Abstract

Probing strategies have been shown to detect semantic features intermediate to certain fragments of NLI. In the case of natural logic, the relation between these features and the entailment label is explicitly known: as such, this provides a ripe setting for interventional studies on the NLI models' representations, allowing for stronger causal conjectures and a deeper critical analysis of interventional probing methods. In this work, we carry out new and existing vector-level interventions to investigate the effect of these semantic features on NLI classification: we perform *amnesic* probing (which removes features as directed by learned probes) and introduce the *mnesic* probing variation (which forgets all dimensions *except* the probe-selected ones). Furthermore, we delve into the limitations of these methods and outline pitfalls that have been obscuring the effectivity of such studies.

1 Introduction

The *probing* paradigm has emerged as a useful interpretability methodology which has been shown to have reasonable information-theoretic underpinnings (Pimentel et al., 2020; Voita and Titov, 2020; Zhu and Rudzicz, 2020), indicating whether a given feature is captured in the intermediate vector representations of neural models. It has been noted many times that this does not generally imply that the models are *using* these learnt features, and they may represent vestigial information from earlier training steps (Ravichander et al., 2021; Elazar et al., 2020).

Only through interventional analyses can we start to make claims about which modelled features are used for a given downstream task: this is the aim of works such as Elazar et al. (2020) and Geiger et al. (2021). We refer to the case where the interventions are guided by trained probes as *interventional probing*.

It has been suggested in Elazar et al. (2020) that if features are strongly detected by probes, one may use debiasing methods such as *iterative nullspace projection (INLP)* (Ravfogel et al., 2020) to intervene on the corresponding vector representations and effectively "remove" the features before reinsertion into the given classifier. This methodology is referred to as *amnesic probing* (Elazar et al., 2020). Investigating the effect of these intervention operations on the classifier performance could allow for stronger causal claims about the role of the probe-detected features.

In this work, we delve deeper into the amnesic probing methodology with an NLI case study and identify two key limitations. Firstly, there is an issue of dimensionality: when the number of dimensions is high and the number of auxiliary feature classes is low, it seems that amnesic probing is not sufficiently informative. In particular, we cannot rely on the same control baselines to reach the kind of conclusions discussed in (Elazar et al., 2020), as nulling out small numbers of random directions consistently has no impact on the downstream performance. Secondly, in the linguistic settings explored in Elazar et al. (2020), we do not have expectations for exactly *how* or even *if* the explored features should be affecting the downstream task. This makes it difficult to explore the effectivity of the methodology itself.

To this end, we propose the use of a controlled subset of NLI called *Natural Logic* (MacCartney and Manning, 2007). In this setting, the intermediate linguistic features of *context monotonicity* and *lexical relations* are already known to be highly extractable from certain NLI models' hidden layers (Rožanova et al., 2021b), allowing us a certain amount of understanding and control of these features' representations in the latent space. Using the deterministic and well-understood nature of the problem space where we have concrete *expectations* about the theoretical interaction between the

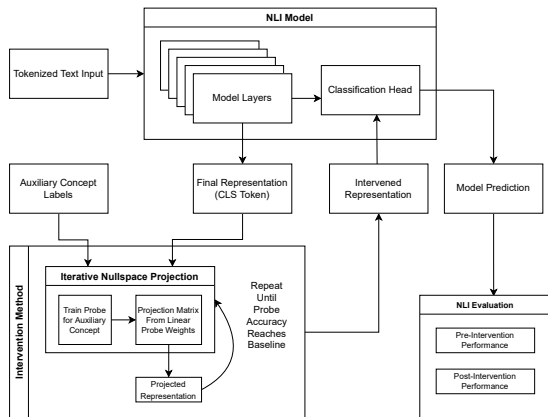


Figure 1: Workflow for Interventional Probing For NLI Models

intermediate features and the downstream label, we may critically analyse the effectivity of interventional probing.

Through the application of probe-based interventions in this setting, we show that blindly applying the amnesic probing argument structure leads to unexpected and contradictory conclusions: the two features which the final label is known to depend on are shown to have no influence on the final classification (both jointly and independently). This further calls into question the suitability of these methods for situations where a small number of feature label classes and high dimensionality of representations is concerned.

As a consequence, we introduce and study a variation which we call *mnesitic* probing, which we show to be more informative in the high-dimensional, low-class-count setting: the core idea is to *keep only* the directions identified by the iteratively trained probes. This allows us to analyse much lower dimension subspaces, and leads to more informative observations in line with expected behaviour for natural logic.

In summary, the contributions of the paper are as follows:

1. We propose the setting of *natural logic* to be ripe territory for exploration of interventional probing strategies.
2. We note two limitations of the amnesic probing methodology, demonstrating both dimensionality limitations for the control baselines

4.4 and contradictory behaviour in the NLI setting 4.2 (namely that that the expected effects of semantic features on the downstream NLI task are notably absent).

3. Building upon previous interventional methodologies, we introduce an additional *mnesitic* intervention operation based on probe outputs, which uses the outputs of the INLP process in the opposite way.

4. We contrast the *mnesitic* probing strategy with the amnesic probing results, and demonstrate it presents more informative results which are aligned with the constructed expectations in our high dimensional, low label class count setting.

2 Interventional Probing

We may summarise the general setup of interventional probing as follows: suppose we start with a classification model that may be decomposed as $f \circ g : \mathcal{X} \rightarrow \mathbb{R}^n$, where g is an encoder module which yields a representation that serves as an input to the classifier head f , and n is the number of output classes of the final classifier. We aim to intervene on the output of g and observe the change in the performance of f (usually in comparison with some kind of random control baseline intervention).

Linear probes are able to identify subspaces in which a given feature set is best represented: these may be used as a guide for vector-level intervention on the representation space. This is the class of interventions we are concerned with here: in particular, when the interventions are vector *projections* guided by the learned probes which are indicative of a given auxiliary feature.

The exact nature of this intervention is interchangeable. We consider two in particular: the *amnesic* intervention introduced in Elazar et al. (2020) (described further in section 2.2) and our *mnesitic* variation of the same INLP techniques (section 2.3).

2.1 What Should it Tell Us?

The interventional probing steps are performed on exactly the representation that would have been an input to the classifier head f . We may re-insert the intervened representations and re-calculate the classifier accuracy (note that the iterative projections

in sections 2.2 and 2.3 maintain the original dimensionality of the vector set but reduce the *rank*).

We are looking to see if the downstream performance of the classifier f drops. If it does, the interventions have removed information that was necessary for successful classification. However, as any projection would remove some information, these results must be viewed in the context of a control intervention: if the INLP process ends up removing n directions, a sample of n randomly chosen directions is selected from the original representation, Elazar et al. (2020) argue that if the amnesic downstream performance drops significantly more than the random removal control performance, we may conclude that the features were necessary for the final downstream classification. On the other hand, if the performance does not drop at all, the features were not useful for the classifier in the first place. In the ensuing sections and results, we demonstrate that this is not necessarily a valid conclusion.

2.2 The Amnesic Intervention

We follow the procedure in (Elazar et al., 2020) (in turn based on *iterative nullspace projection* (Ravfogel et al., 2020)): given a set X of encoded representations for the textual input (with dimensions `num_examples` \times `embedding_dimension`), we iteratively train linear SVM classifiers according to a set of auxiliary feature labels. For each INLP step i , This yields a linear transformation $W_i X + B$, where the vectors of W_i define directions onto which the probe projects the representations for auxiliary label classification (i.e., these are the chosen directions most aligned with auxiliary class separation). For each step i , an orthogonal basis denoted R_i is found for this rowspace. The projection to the intersection of the nullspaces is given by a matrix

$$PX = (I - (R_0 + \dots + R_n))X.$$

The matrix product PX is a matrix in the original dimensions of X , but with reduced rank by the number of iteration steps (as each projection "flattens out" the representation in these directions).

Projection to the intersection of nullspaces is thus the removal of any information pertaining to the auxiliary feature labels (or at least, the information which allows high performance for a linear probe). The training terminates these auxiliary task classifiers start consistently performing at the majority class baseline, indicating that there is no further linearly information to be extracted from the

remaining representation. As such, the resulting representation is treated as an altered representation where this feature is *removed* or forgotten.

2.3 A Variation: The Mnestic Intervention

Elazar et al. (2020) perform a series of experiments on various linguistic features which had previously been shown to be well-captured in language model representations and use the amnesic probing methodology to distinguish between features that are *used* by the model and those that are not by comparing post-intervention downstream task performance to a baseline of randomly removed directions.

Rather than projecting the embedded representations to the intersection of nullspaces of the trained probes (removing the target property), we project them to the *union of the rowspaces* with the transformation:

$$\begin{aligned}(I - P)X &= (I - (I - (R_0 + \dots + R_n)))X \\ &= (R_0 + \dots + R_n)X\end{aligned}$$

This has the opposite effect: we use projection to null out *everything except* the directions identified by the probes as indicative of the target feature. As such, we "remember" only that feature rather than forgetting it.

3 Experimental Setup

In this study, we use interventional methods¹ to study the internal behaviour of NLI models. We compare amnesic and mnestic variations of the INLP strategy, evaluating intermediate feature probing performance and downstream NLI performance after every step of the intervention process.

For each auxiliary feature label and model, we perform the *interventional probing* as outlined in figure 1.

3.1 Dataset

Our setting for this study is a fragment of NLI called *Natural Logic* (MacCartney and Manning, 2007). In particular, we focus on single-step natural logic inferences in which entailment examples are generated by replacing a noun phrase in a sentence with a hyponym, hypernym or unrelated noun phrase. The context of the substituted term

¹We reuse much of the code included with (Elazar et al., 2020), but we include our data and reproducible experimental code in *anonymized github repo*.

is either *upward* or *downward* monotone, as determined by the composition of negation markers, generalized quantifiers or determiners present in the context. The entailment label of the example is a consequence of this feature and the lexical relation between the substituted terms.

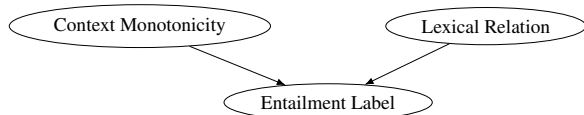


Figure 2

We use the NLI_XY dataset from (Rozanova et al., 2021b,a). By construction, the NLI_XY dataset consists of NLI examples which rely on exactly these two abstract features: context monotonicity and the lexical relation of the substituted terms.

We perform two flavours of probe-based interventions (described fully in section 2) with four feature label sets (described next).

Auxiliary Feature Labels We begin with the two relevant intermediate features (respectively, context monotonicity and lexical relation) which are already known to correlate with stronger performance on the downstream NLI_XY task (Rozanova et al., 2021b). We will refer to this as *single-feature* interventional probing, as the probing and intervention steps are only applied to one feature set at a time. Next, we combine the two features in a cross product, creating a new feature label set with all possible combinations of these intermediate features (in the dataset, they are completely independent variables by construction (Rozanova et al., 2021a)). We refer to this as the *composite feature label*.

Lastly, we also consider the *entailment label* itself (the downstream task label) as an input to the interventional probing process. The latter is particularly useful as a diagnostic sanity check, and aids the critical nature of our findings.

3.2 NLI Models and Encoding

We compare a selection of BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) models trained for NLI classification. Firstly, we include a pair of models trained respectively on the MNLI (Williams et al., 2018) and SNLI (Bowman et al., 2015) benchmark datasets. In (Rozanova et al., 2021b) and (Rozanova et al., 2021a), it is shown that when `roberta-large-mnli` (a model

which performs well on benchmarks but poorly on the targeted NLI_XY challenge set) receives additional training on the adversarial HELP dataset (Yanaka et al., 2019) it improves in NLI_XY performance and *begins to show high probing performance for the relevant intermediate features*, context monotonicity and lexical relations: this is the necessary precondition for doing interventional probing. We include two of their models with this property: `roberta-large-mnli-help` and `roberta-large-mnli-double-finetuning`, with the other models included for a contextual comparison.

We perform probing and intervention on the final representation that precedes the NLI classification head: in the case of BERT and RoBERTa, this is the [CLS] token of the final layer.

The initial input is a tokenized NLI example from the NLI_XY dataset. The findings in (Rozanova et al., 2021b) show that the intermediate feature labels (context monotonicity and lexical relations) are detectable in the concatenated tokens of the substituted noun phrases: however, for interventional purposes, we perform the probing and intervention steps on the [CLS] token which serves as an input to the NLI classifier head: we have found that the same features are detectable to a comparable standard, and this is the only position at which we are able to make a sensible intervention that would allow conclusions about the final classifier head only.

3.3 Evaluation

The significant metrics for these interventional probing paradims are the *probing accuracy* before and after the iterative nullspace projection steps (a decline to random performance indicates the feature is being “removed” from the representation in the sense that it is no longer detectable by linear probes) and the *downstream classification accuracy* on the NLI task the model’s were trained for (in our case, we report the accuracy on the NLI_XY task).

For amnesic probing, we report the performance deltas for both the probing and downstream tasks. However, for mnestic probing, a slightly more nuanced and qualitative view is helpful: it can be assumed that eventually mnestic probing will reach comparable performance to the untouched vector representations, but we are interested in the comparative rates at which this happens. As the inter-

Model	Feature	Probing Performance		NLI-XY Performance	
		Start	Intervention Δ	Start	Intervention Δ
roberta-large-mnli-help	insertion relation	80.58	-40.35	79.79	0.06
	context monotonicity	87.65	-46.22	79.79	-0.09
	composite	64.48	-43.95	79.79	0.32
	entailment label	78.05	-37.49	79.79	-1.57
roberta-large-mnli-double-finetuning	insertion relation	62.7	-36.49	80.04	0.11
	context monotonicity	89.79	-43.28	80.19	0
	composite	57.64	-49.56	80.08	-1.67
	entailment label	82.8	-24.94	80.19	-16.53
roberta-large-mnli	insertion relation	80.39	-45.59	57.22	8.99
	context monotonicity	75.44	-27.49	57.37	-0.43
	composite	72.35	-53.51	57.24	-2.27
	entailment label	73.6	-15.31	57.37	0.1
bert-base-uncased-snli-help	insertion relation	59.53	-19.1	45.95	0.28
	context monotonicity	82.72	-33.94	45.52	-2.35
	composite	37.19	-17.08	45.76	13.68
	entailment label	47.05	0.38	45.91	0
bert-base-uncased-snli	insertion relation	60.26	-35.14	48.99	1.05
	context monotonicity	81.09	-30.77	49.42	-6.25
	composite	35.37	-17.83	50.73	7.45
	entailment label	42.44	-0.24	49.42	0

Table 1: Amnesic probing performance deltas across models and target feature labels: first listed is the performance on the probing task with respect to the indicated feature, and then the accuracy on the downstream NLI-XY task. We note the results pre-intervention and the ensuing change in accuracy.

ventions are iterative, we may feed the intervened representations into the classifier head at *each step* of the intervention process - we use this to provide a step-wise presentation of results in linear plots in figure 3.

While the tabulated deltas in table 1 results are sufficient to present our observations on amnesic probing, for comparison we also include the step-wise graphical presentations in the appendix.

4 Results and Discussion

4.1 Single Feature Amnesic Probing

The results for the standard amnesic probing procedure are in table 1. In particular, the single feature results are in the rows with features labelled *insertion relation* and *context monotonicity*. The amnesic operation is successful - the respective probing accuracies approach and reach the majority class baseline. The length of this iterative process is indicative of the number of dimensions removed to reach this baseline: it can also be considered a proxy for the strength of the feature presence in the representations, or rather, the dimension of the semantic subspace corresponding to the target features.

The second phase of this process, i.e. the resub-

stitution of the modified representations as inputs to the NLI classifier head, can be seen in the right hand portion of table 1, labelled *NLI-XY Performance*. The result is unexpected: for each of these features, *the downstream task performance appears to be unaffected after their removal*. This is surprising when the dataset is explicitly controlled to rely only on these two features.

4.2 Multi Feature Amnesic Probing

The results for the amnesic probing procedure utilizing *both* auxiliary feature label sets and the entailment gold label are in the rows of table 1 with labels *composite* and *entailment label* respectively. We observe that once again, the downstream task performance is mostly unaffected. Unlike the unexpected result in the previous section, it’s difficult to argue away the fact that this is somewhat contradictory: while single feature removal may be subject to some confounding bias, the removal of both features exhausts the variables on which this classification depends. This is highly unexpected, and suggests a point of failure for the amnesic probing process. Naturally, we cannot be without doubt that despite all our best efforts to work with a controlled dataset that relies only on these two know (but still complex) features, a model may yet find

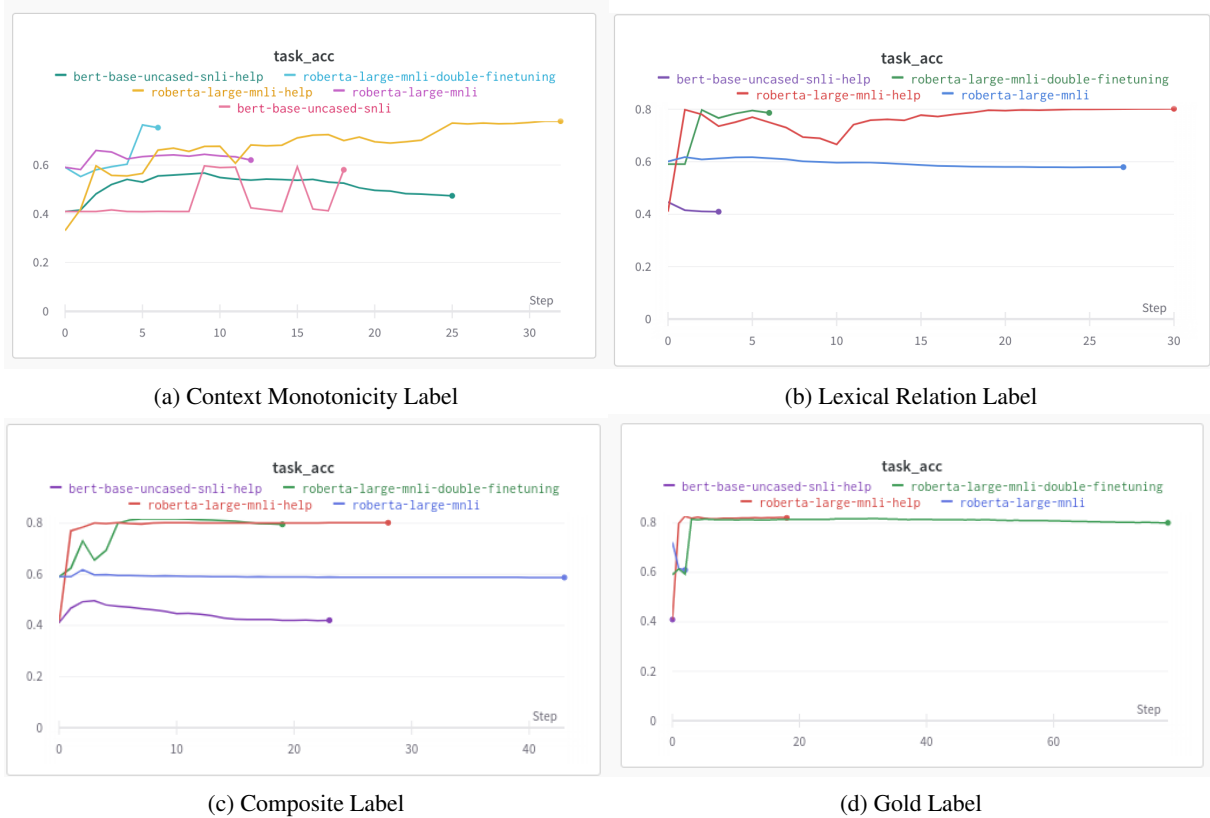


Figure 3: Downstream Task Performance After Mnestic Intervention

Figure 4: Mnestic Probing Results

381 unrelated heuristics to exploit that may correlate
 382 so strongly with the downstream task label that it
 383 may perform well without representing and using
 384 these intermediate features. However, we imagine
 385 this to be a rather low probability scenario to be
 386 that the model simultaneously learns such heuristics
 387 but simultaneously learn representations that
 388 create strong clusters for the known intermediate
 389 features *without using them at all*. The models
 390 which we have observed to perform more less well
 391 on NLI-XY (such as roberta-large-mnli) are indeed
 392 estimated to be using sub-par heuristics, but this
 393 also comes with poor probing results for the inter-
 394 mediate features - naturally, this in itself does not
 395 imply anything conclusive, but certainly adds to
 396 our convictions.

397 On a separate note, it is noted in Elazar et al.
 398 (2020) that there is no control for the number of
 399 dimensions removed, while there is a clear correla-
 400 tion between downstream task performance and the
 401 number of label classes (and thus removed probe
 402 directions) are in play. Our feature sets have only 2
 403 and 3 classes respectively. In the most analogous
 404 result in (Elazar et al., 2020) where the auxiliary

405 features had very few classes and no change on
 406 the downstream performance was observed, it was
 407 concluded that the features must have no effect on
 408 the outcome. It is very likely that *too little informa-*
 409 *tion* is being removed in this process to observe any
 410 impact on the downstream task performance. This
 411 could potentially be pointing to high redundancy in
 412 the representations which the amnesic intervention
 413 may struggle to remove appropriately.

4.3 Mnestic Probing

414 Given the possible dimensionality problem, the al-
 415 ternative method of *mnestic* probing seems promis-
 416 ing: many dimensions are removed and few remain,
 417 so it appears to be a ripe setting for observing and
 418 comparing effects on downstream NLI accuracy.
 419 The results for the *mnestic* probing procedure are
 420 in figure 3. There is a clear increase in NLI perfor-
 421 mance with subsequent addition of probe-chosen
 422 directions to the representations, but these results
 423 especially need to be viewed in the context of sec-
 424 tion 4.4, where we compare the performance to
 425 random choices of included directions.
 426

427 We observe that the *composite* label and the gold

428 *entailment* label are reflected as expected in the
 429 mnestic probing experiments: the inclusion of the
 430 probe-selected dimensions with respect to these la-
 431 bels introduces a sharp and immediate increase in
 432 the NLI classifier performance. This is significantly
 433 steeper than the baseline increase observed in ran-
 434 dom addition of representation directions. Simi-
 435 larly, the increase is nearly as sharp for the lexical
 436 relation label. However, although an increase is
 437 observed during the iterative mnestic probing in-
 438 tervention for context monotonicity, this increase
 439 is not at a dramatically higher rate than adding
 440 subsequently more directions from the original rep-
 441 resentation. For monotonicity specifically, this is
 442 not enough to conclude that the feature (or at least,
 443 the corresponding probe-selected dimensions) are
 444 critical to the final classifier.

445 Nevertheless, we have been able to make clearer
 446 observations than were possible in the amnesic
 447 probing setting.

448 4.4 Control Comparison

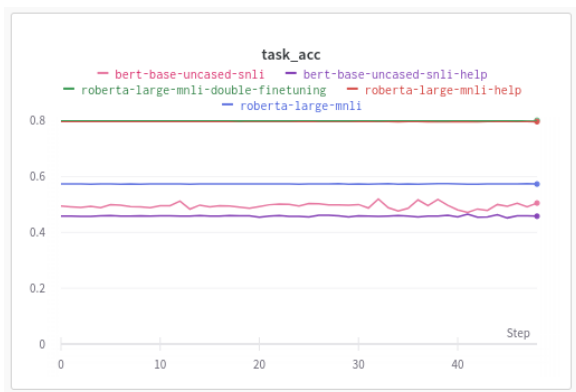


Figure 5: Amnesic control experiment: Downstream NLI accuracy upon the *removal* of n random directions of the original representation.

449 We contextualise all the preceding results with a
 450 set of control experiments both for amnesic (figure
 451 5) and mnestic (figure 6) probing. Note in partic-
 452 ular that even with very few random dimensions
 453 kept, downstream performance starts approaching
 454 comparable levels to the full representations. As
 455 such, a single random baseline as in Elazar et al.
 456 (2020) can be misleading: there is enough variabil-
 457 ity in the random direction results so as to allow
 458 for a false claim of feature irrelevance by simply
 459 getting lucky; as few as 3 dimensions can perform
 460 at the original model’s performance level or arbi-
 461 trarily lower.



Figure 6: Mnestic control experiment: Downstream NLI accuracy upon the *selection* of n random directions of the original representation.

462 Lastly, we compare to the mnestic probing re-
 463 sults in figure 3: with the probe-selected mnestic
 464 dimension choices, the increase in downstream per-
 465 formance does seem to happen faster and in a more
 466 consistent fashion, while the selection of n ran-
 467 domly chosen directions introduces very haphaz-
 468 ard performance spikes. This suggests the probe-
 469 selected dimensions are consistently adding to the
 470 model’s access to the relevant information, and
 471 this may be stronger evidence for the usefulness of
 472 the examined features for the final classification.

473 5 Related Work

474 The use of probing as an interpretability strategy
 475 dates back as far as works such as Alain and Bengio
 476 (2018) and (Conneau et al., 2018), but a core set of
 477 work on the detailed development of the method-
 478 ology includes Hewitt and Liang (2019); Belinkov
 479 and Glass (2019); Voita and Titov (2020); Pimentel
 480 et al. (2020). For a full survey, see Belinkov (2022).

481 The application of probing strategies to natural
 482 logic components has been explored in Rozanova
 483 et al. (2021b) and Geiger et al. (2020). In Rozanova
 484 et al. (2021b), probing experiments have proven
 485 effective in detecting the presence or absence of
 486 features such as *context monotonicity* and *phrase-*
 487 *pair relations* in the internal representations of NLI
 488 models.

489 Regarding interventions as interpretability tools
 490 for machine learning classifiers, there are two broad
 491 categories: those that modify the raw input (such
 492 as image or text) in a controlled way, and those that
 493 modify the hidden/latent vector representations of
 494 the data at various stages of the models’ input pro-
 495 cessing. While input-level interventions are more
 496 common as they are usually easier to control and

	Intervention	Tested Effect	Feature Characterisation	Requires Intermediate Labels	Intervention Linked to Concept Interpretation	Domain
Amnesic Probing / INLP (Elazar et al., 2020)	Debiasing / Feature Removal	Downstream Classifier Accuracy	Linear Classifier	Yes	No	Language Modelling
CausaLM: Causal Model Explanation	Re-Training Model Copy	Text representation-based individual treatment effect (TReITE)	Retrained Base Model	Yes	Yes	Sentiment Analysis
Through Counterfactual Language Models (Feder et al., 2021)	For Counterfactual Representation	Average Causal Effect Measure	VAE	Yes	Yes	Vision Classification
Causal Concept Effect	Generative Modeling	Custom Gradient Sensitivity Measure	Linear Classifier	Yes	Yes	Vision Classification
Concept Activation Vectors (TCAV) (Kim et al., 2018)	Value Shift in Vector Direction	Reconstruction Quality	VAE	No	Qualitative Judgement (Vision Only)	Vision Classification
Latent Space Explanation by Intervention	VAE Input Discretization and Reconstruction	Difference Between Concept Addition and Removal Effect	Linear Classifier	Yes	Yes	Vision Classification
Meaningfully Explaining Model Mistakes Using Conceptual Counterfactuals	Weighted Combination of Concept Vectors					

Table 2: Related Work on Latent Concept Interventions

are strongly interpretable, they don't allow us to explore and conjecture about exact high-level representational mechanisms in the latent space. We tabulate a few relevant interventional interpretability methods in table 2. Note in particular the variation in the *generation* step for the intervened input; some use generative modelling for counterfactual examples, while we use cheaper linear probes.

The only other work in which interventional methods have been applied to natural logic is Geiger et al. (2021): a similar problem setting is considered, but at a finer granularity. Our work focuses more on the summarised abstract notion of context monotonicity as a single feature, rather than the intermediate tree nodes that determine its final monotonicity profile. The interventions used in this work are vector *interchange* interventions; partial representations from transformed inputs are used, as opposed to direct manipulations of the encoded vectors.

6 Conclusion and Future Work

Our experimental setting has shown significant limitations of amnesic probing in high-dimensional settings where there are few label classes (and consequently fewer dimension modified), even if these classes are strongly detectable. Our results point out that it is misguided to conclude that a given feature is not used when post-amnesic-intervention downstream performance fails to drop, especially in our example amnesic probing studies of a) the gold downstream feature label and b) the composite of two labels that jointly determine the entailment label. This may be due to a dimension/rank confounder variable and high redundancy of information in the representations. It remains to be checked whether high performance in the random control directions corresponds to strong alignment with these probe-selected directions: we propose an analysis of the *dot products* with the fixed set of probe-selected dimensions, which indicates a shared directionality measure (0 for orthogonal vectors and 1 for codirectional ones).

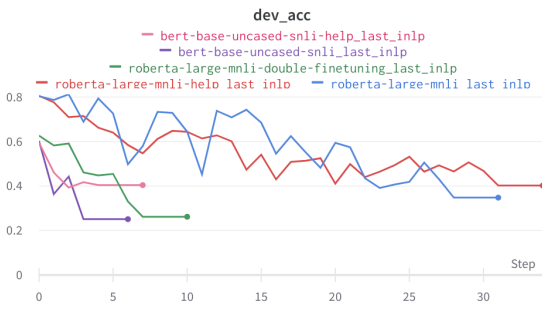
We have introduced a modification of the amnesic probing paradigm which we call *mnesic* probing which uses the same INLP process but considers the opposite intervention: using the union of projection rowspaces to keep *only* the directions the probes have identified to be modelling the target information. This strategy presents results that are more aligned with theoretical expectations, possibly because we are now able to make comparisons in a lower rank setting and also work with more useful control baselines.

References

- Guillaume Alain and Yoshua Bengio. 2018. [Understanding intermediate layers using linear classifier probes.](#)
- Yonatan Belinkov. 2022. [Probing classifiers: Promises, shortcomings, and advances.](#) *Computational Linguistics*, 48(1):207–219.
- Yonatan Belinkov and James Glass. 2019. [Analysis methods in neural language processing: A survey.](#) *Transactions of the Association for Computational Linguistics*, 7:49–72.
- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*.
- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. [What you can cram into a single vector: Probing sentence embeddings for linguistic properties.](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding.](#) In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

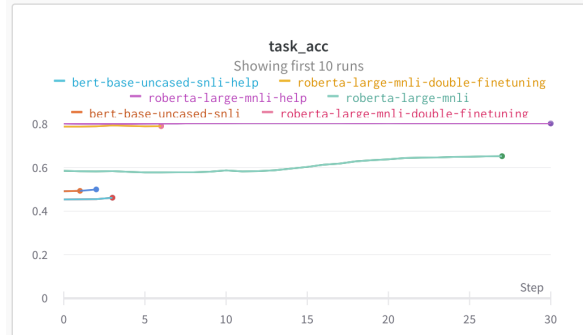
582	Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2020. Amnesic probing: Behavioral explanation with amnesic counterfactuals .	637
583		638
584		639
585	Amir Feder, Nadav Oved, Uri Shalit, and Roi Reichart. 2021. CausaLM: Causal model explanation through counterfactual language models . <i>Computational Linguistics</i> , 47(2):333–386.	640
586		641
587		642
588		643
589	Atticus Geiger, Hanson Lu, Thomas F Icard, and Christopher Potts. 2021. Causal abstractions of neural networks . In <i>Advances in Neural Information Processing Systems</i> .	644
590		645
591		646
592		647
593	Atticus Geiger, Kyle Richardson, and Christopher Potts. 2020. Neural natural language inference models partially embed theories of lexical entailment and negation . In <i>Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 163–173, Online. Association for Computational Linguistics.	648
594		649
595		650
596		651
597		652
598		653
599		654
600	John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.	655
601		656
602		657
603		658
604		659
605		660
606		661
607	Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. 2018. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav) . In <i>International conference on machine learning</i> , pages 2668–2677. PMLR.	662
608		663
609		664
610		665
611		666
612		667
613	Y. Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach . <i>ArXiv</i> , abs/1907.11692.	668
614		669
615		670
616		671
617		672
618	Bill MacCartney and Christopher D. Manning. 2007. Natural logic for textual inference . In <i>Proceedings of the ACL-PASCAL Workshop on Textual Entailment and Paraphrasing</i> , pages 193–200, Prague. Association for Computational Linguistics.	673
619		674
620		675
621		676
622		677
623	Tiago Pimentel, Josef Valvoda, Rowan Hall Maudslay, Ran Zmigrod, Adina Williams, and Ryan Cotterell. 2020. Information-theoretic probing for linguistic structure . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4609–4622, Online. Association for Computational Linguistics.	678
624		679
625		680
626		681
627		682
628		683
629		684
630	Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. Null it out: Guarding protected attributes by iterative nullspace projection . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 7237–7256. Association for Computational Linguistics.	685
631		686
632		687
633		688
634		689
635		690
636		691
	Abhilasha Ravichander, Yonatan Belinkov, and Eduard Hovy. 2021. Probing the probing paradigm: Does probing accuracy entail task relevance? In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 3363–3377, Online. Association for Computational Linguistics.	692
		693
		694
		695
		696
		697
	Julia Rozanova, Deborah Ferreira, Mokanarangan Thayaparan, Marco Valentino, and André Freitas. 2021a. Supporting context monotonicity abstractions in neural nli models .	698
		699
		700
		701
		702
	Julia Rozanova, Deborah Ferreira, Marco Valentino, Mokanarangan Thayaparan, and André Freitas. 2021b. Decomposing natural logic inferences in neural NLI . <i>CoRR</i> , abs/2112.08289.	703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		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
		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
		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

Probing Accuracy

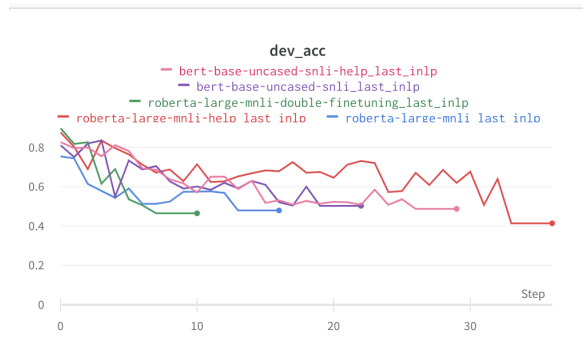


(a) Lexical Relation Probing Performance During Iterative Amnesic Intervention Process

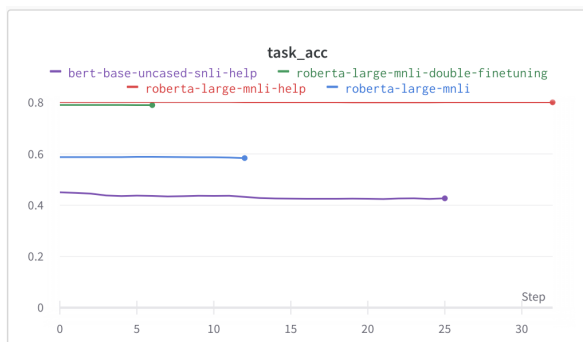
Downstream Task (NLI) Accuracy



(b) Downstream Performance On NLI_XY After Amnesic Intervention (Removing Lexical Relation Information)



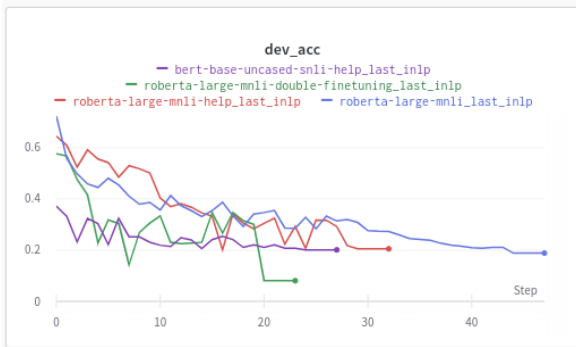
(c) Context Monotonicity Probing Performance During Iterative Amnesic Intervention Process



(d) Downstream Performance On NLI_XY After Amnesic Intervention (Removing Context Monotonicity Information)

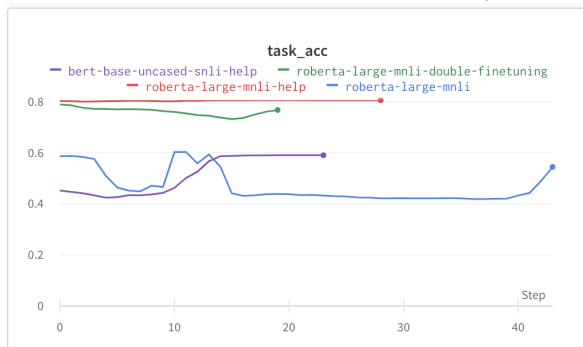
Figure 7: Single Feature Amnesic Probing

Probing Accuracy



(a) Probing Performance On NLI_XY After Composite Label Amnesic Intervention

Downstream Task (NLI) Accuracy



(b) Downstream Performance On NLI_XY After Composite Label Amnesic Intervention

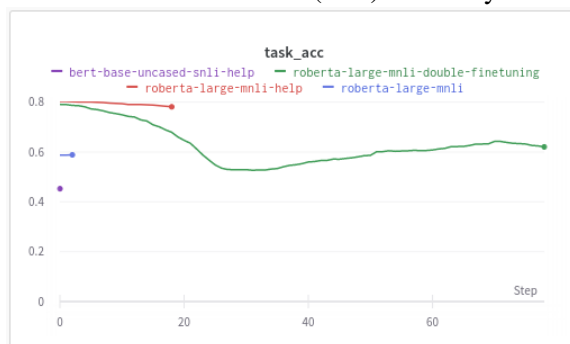
Figure 8: Composite Feature Label Amnesic Probing

Probing Accuracy



(a) Probing Performance On NLI_XY After Entailment Label Amnesic Intervention

Downstream Task (NLI) Accuracy



(b) Downstream Performance On NLI_XY After Entailment Label Amnesic Intervention

Figure 9: Sanity Check: Entailment Gold Label Amnesic Probing