

Decomposing Natural Logic Inferences in Neural NLI

Anonymous ACL submission

Abstract

In the interest of interpreting neural NLI models and their reasoning strategies, we carry out a systematic probing study which investigates whether these models capture the crucial semantic features central to natural logic: *monotonicity* and *concept inclusion*. Correctly identifying valid inferences in *downward-monotone contexts* is a known stumbling block for NLI performance, subsuming linguistic phenomena such as negation scope and generalized quantifiers. To understand this difficulty, we emphasize monotonicity as a property of a *context* and examine the extent to which models capture monotonicity information in the contextual embeddings which are intermediate to their decision making process. Drawing on the recent advancement of the probing paradigm, we compare the presence of monotonicity features across various models. We find that monotonicity information is notably weak in the representations of popular NLI models which achieve high scores on benchmarks, and observe that previous improvements to these models based on fine-tuning strategies have introduced stronger monotonicity features together with their improved performance on challenge sets.

1 Introduction

Large, black box neural models which achieve high scores on benchmark datasets designed for testing *natural language understanding* are the subject of much scrutiny and investigation.

It is often investigated whether models are able to capture specific semantic phenomena which mimic human reasoning and/or logical formalism, as there is evidence that they sometimes exploit simple heuristics and dataset artifacts instead (McCoy et al., 2019; Herlihy and Rudinger, 2021).

In this work, we consider the rigorous setting of *natural logic* (MacCartney and Manning, 2007). This is a highly systematic reasoning principle relying on only two abstract features, each of which is

in itself linguistically complex: *monotonicity* and *concept inclusion relations*. It underlies the majority of symbolic/rule-based and hybrid approaches to NLI and is an important baseline reasoning phenomenon to look for in a robust and principled NLI model.

Downward monotone operators such as negations and generalized quantifiers result in the kinds of natural logic inferences which are often known to stump neural NLI models that demonstrate high performance on large benchmark sets such as MNLI (Williams et al., 2018).

By contrast, in this work we present a **structural** study: investigating to what extent the features relevant for identifying natural logic inferences, especially context monotonicity itself, are captured in the model’s internal representations.

In this work, we carry out a systematic probing study to estimate and compare the extent to which the abstract features at the heart of monotonicity reasoning – i.e., context monotonicity and concept inclusion relations – are present in various NLI models’ representations.

Our contributions are may be summarized as follows:

1. We perform a structural investigation as to whether the behaviour of *natural logic* formalisms are mimicked within popular transformer-based NLI models.
2. For this purpose, we present a joint NLI and semantic probing dataset format (and dataset) which we call NLI-XY: it is a unique probing dataset in that the probed features relate to the NLI task output in a very systematic way.
3. We employ thorough probing techniques to determine whether the abstract semantic features of *context monotonicity* and *concept inclusion relations* are captured in the models’ internal representations.

- 082 4. We observe that some well-known NLI mod- 132
 083 els demonstrate a systematic failure to model 133
 084 context monotonicity, a behaviour we observe 134
 085 to correspond to poor performance on natu- 135
 086 ral logic reasoning in downward-monotone 136
 087 contexts. However, we show that the existing 137
 088 HELP dataset improves this behaviour.
- 089 5. We support the observations in the prob- 138
 090 ing study with several *qualitative analyses*, 139
 091 including decomposed error-breakdowns on 140
 092 the NLI-XY dataset, representation visualiza- 141
 093 tions, and evaluations on existing challenge 142
 094 sets. 143

095 2 Related Work 144

096 Natural logic dates back to the formalisms of 145
 097 Sanchez (1991), but has been received more re- 146
 098 cent treatments and reformulations in (MacCartney 147
 099 and Manning, 2007; Hu and Moss, 2018). Sym- 148
 100 bolic and hybrid neural/symbolic implementa- 149
 101 tions of the natural logic paradigm have been explored in 150
 102 (Chen et al., 2021; Kalouli et al., 2020; Abzianidze, 151
 103 2017; Hu et al., 2020). 152

104 The shortcomings of natural logic handling in 153
 105 various neural NLI models have been shown with 154
 106 several *behavioural* studies, where NLI challenge 155
 107 sets exhibiting examples of downward monotone 156
 108 reasoning are used to evaluate performance of 157
 109 models with respect to these reasoning patterns 158
 110 (Richardson et al., 2019; Yanaka et al., 2019b,a; 159
 111 Goodwin et al., 2020; Geiger et al., 2020). 160

112 In an attempt to better identify features that neu- 161
 113 ral models manage or fail to capture, researchers 162
 114 have employed *probing* strategies: namely, the *di-*
 115 *agnostic classification* (Alain and Bengio, 2018)
 116 of auxiliary feature labels from internal model
 117 representations. Most probing studies in natural
 118 language processing focus on the *syntactic* fea-
 119 tures captured in transformer-based language mod-
 120 els (Hewitt and Manning, 2019), but calls have
 121 been made for more sophisticated probing tasks
 122 which rely more on contextual information (Pi-
 123 mentel et al., 2020).

124 In the realm of semantics, probing studies have
 125 focused more on *lexical* semantics (Vulić et al.,
 126 2020): word pair relations are central to monotonic-
 127 ity reasoning, and thus form part of our probing
 128 study as well, but the novelty of our work is the task
 129 of classifying context monotonicity from contex-
 130 tual word embeddings. Due to its context-sensitive
 131 nature, it cannot be learnt by “memorizing” the

labels of specific words in the training data, a key
 shortcoming in probing studies which focus on
 tasks such as POS tagging and word-pair relation
 classification, which have much less dependency
 on context.

137 3 Problem Formulation 138

138 3.1 Decomposing Natural Logic 139

139 Natural logic inferences (as formalized in Sanchez
 140 (1991); MacCartney and Manning (2007)) are usu-
 141 ally described with respect to *substitution* opera-
 142 tions. Certain word substitutions result in either
 143 forward or reverse entailment, while others result
 144 in neither. This is the basis for a calculus of de-
 145 termining entailment from substitution sequences
 146 (MacCartney and Manning, 2007; Hu et al., 2020;
 147 Hu and Moss, 2018).

148 Broadly speaking, we wish to determine whether
 149 well-known transformer-based NLI models mimic
 150 the reasoning strategies of natural logic. However,
 151 as neural NLI models are black box classifiers that
 152 only see a premise/hypothesis sentence pair as its
 153 input, it is not immediate how to compare its pro-
 154 cess to a rule-based system.

155 To this end, we consider a formulation of natural
 156 logic which describes its rules in terms of concept
 157 pair relations and *context monotonicity* (similar to
 158 (Rožanova et al., 2021)).

159 Consider the following example of a single step
 160 natural logic inference, which we will decompose
 161 into semantic components relevant to its entailment
 label:

		NLI Label
Premise	I did not eat any fruit for breakfast.	Entailment
Hypothesis	I did not eat any raspberries for breakfast.	

162 The hyponym/hypernym pair (raspberries, fruit)
 exemplifies a more general relation which we will
 refer to as the *concept inclusion*¹ relation \sqsubset , (and
 dually, *reverse concept inclusion* \sqsupset) in reference
 to the semantic interpretation of predicates related
 with subset inclusion, as in:

$$\{x \mid \text{raspberry}(x)\} \subset \{x \mid \text{fruit}(x)\}.$$

163 In the above example, they occur in a shared
 164 **context**, namely the sentence template

¹In (MacCartney and Manning, 2007), this is treated
 as a “generalized entailment” relation which is defined on
 word/phrase pairs and extends to full sentences pairs using
 natural logic rules.

“I did not eat any _____ for breakfast”.

Such a context may be treated as a *function* f

$$f : (\mathcal{X}, \sqsubseteq) \rightarrow (\mathcal{S}, \Rightarrow)$$

between a set of concepts \mathcal{X} (ordered by the concept inclusion relation) and the set \mathcal{S} of full sentences ordered by entailment. We say that f is *upward monotone* (\uparrow) if it is order *preserving*, i.e.

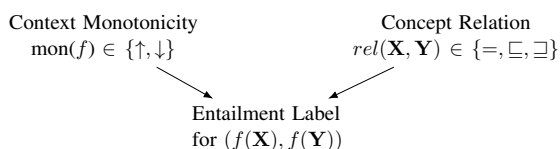
$$\forall_{X,Y}(X \sqsubseteq Y \text{ implies } f(X) \Rightarrow f(Y))$$

and that f is *downward monotone* (\downarrow) if it is order *reversing*, i.e.

$$\forall_{X,Y}(X \sqsupseteq Y \text{ implies } f(X) \Rightarrow f(Y)).$$

Given a natural language context f , any pair of grammatically valid insertions (X, Y) (e.g. ("raspberries", "fruit")) yields a sentence pair $f(X), f(Y)$. Treating $f(X)$ as a *premise* sentence and $f(Y)$ as a *hypothesis* sentence, a trained neural NLI model can provide a classification of whether $f(X)$ entails $f(Y)$.

In summary, these two abstract linguistic features, *context monotonicity* and *concept inclusion relation*, jointly determine the final gold entailment label of this type of NLI example.



3.2 NLI-XY Dataset Format

We follow this formalism as the basis for a *dataset format*, which we refer to as NLI-XY. This is the first probing dataset format (and consequently, dataset) in NLP where the auxiliary labels for intermediate semantic features influence the final task label in a rigid and determinate (yet simple) way, with these features being themselves linguistically complex. As such, it is as such a "decomposed" natural logic dataset format, where the positive entailment labels are further enriched with labels for the monotonicity and relational properties which gave rise to them. This allows for informative qualitative and structural analyses into natural logic handling strategies in neural NLI models.

The NLI-XY dataset format is comprised of the following:

			Auxilliary Label
Context	f	I did not eat any _____ for breakfast.	\downarrow (downward monotone)
Insertion Pair	(X, Y)	(fruit, raspberries)	\sqsupseteq (reverse concept inclusion)
			NLI Label
Premise	$f(X)$	I did not eat any fruit for breakfast.	Entailment
Hypothesis	$f(Y)$	I did not eat any raspberries for breakfast.	

Table 1: A typical NLI-XY example with labels for context monotonicity, lexical relation and the final entailment label.

1. A set of *contexts* f with a blank position indicated with an ‘ x ’, marked for the context monotonicity label.
2. A set of *insertion pairs* (X, Y) , which are either words or phrases, labeled with the concept inclusion relation.
3. A derived set of premise and hypothesis pairs $(f(X), f(Y))$ made up of permutations of (X, Y) insertion pairs through contexts f , controlled for grammaticality as far as possible.

The premise/hypothesis pairs may thus be used as input to any NLI model, while the context monotonicity and insertion relation information can be used as the targets of an auxiliary probing task on top of the model’s representations.

4 NLI-XY Dataset Construction

We make our NLI-XY dataset and all the experimental code used in this work is publically available². We constructed the NLI-XY dataset used here as follows:

Context Extraction We extract context examples from two NLI datasets which were designed for the behavioural analysis of NLI model performance on monotonicity reasoning. In particular, we use the manually curated evaluation set MED (Yanaka et al., 2019a) and the automatically generated HELP training set (Yanaka et al., 2019b). By design, as they are collections of NLI examples exhibiting monotonicity reasoning, these datasets mostly follow our required $(f(X), f(Y))$ structure, and are labeled as instances of upward or downward monotonicity reasoning (although the contexts are not explicitly identified).

²Anonymized github link.

We extract the common context f from these examples after manually removing a few which do not follow this structure (differing, for example, in pronoun number agreement or prepositional phrases). We choose to treat determiners and quantifiers as part of the context, as these are the kinds of closed-class linguistic operators whose monotonicity profiles we are interested in. To ensure grammatically valid insertions, we manually identify whether each context as suitable either for a singular noun, mass noun or plural noun in the blank/ x position.

Insertion Pairs Our (X, Y) insertion phrase pairs come from two sources: Firstly, the labeled word pairs from the MoNLI dataset (Geiger et al., 2020), which features only single-word noun phrases. Secondly, we include an additional hand-curated dataset which has a small number of *phrase-pair* examples, which includes intersective modifiers (e.g. (brown sugar, sugar)) and prepositional phrases (e.g. (sentence, sentence about oranges)). Several of these examples were drawn from the MED dataset. Each word in the pair is labelled as a singular, plural or mass noun, so that they may be permuted through the contexts in a reasonably grammatical way.

Premise/Hypothesis Pairs Premise/Hypothesis pairs are constructed by permuting insertion pairs through the set of contexts within the grammatical constraints. Note that the data is split into train, dev and test partitions *before* this permutation occurs, so that there are **no shared contexts or insertion pairs** between the different data partitions, in an attempt to avoid overlap issues such as those discussed in (Lewis et al., 2021)

5 Experimental Setup

Our experiments are designed to investigate the following questions: Firstly, to what extent do different NLI models differ in their encoding of context monotonicity and lexical relational knowledge? Secondly, if a model successfully captures these features, to what extent do they correspond with the model’s predicted entailment label? We investigate these questions with a detailed probing study and a supporting qualitative analysis, using decomposed error break-downs and representation visualization.

Partition	(X,Y) Relation	Context Monotonicity		
		Up ↑	Down ↓	Total
train	⊆	671	543	1214
	⊃	671	543	1214
	None	244	222	466
	Total	1586	1308	2894
dev	⊆	598	389	987
	⊃	598	389	987
	None	220	242	462
	Total	1416	1020	2436
test	⊆	1103	1066	2169
	⊃	1103	1066	2169
	None	502	516	1018
	Total	2708	2648	5356

Table 2: Dataset statistics for the NLI-XY dataset. We employ an **aggressive** 30, 20, 50 train-dev-test split for a more impactful probing result.

5.1 Models and Representations

We consider a selection of neural NLI models based on transformer language models (such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019) and BART (Lewis et al., 2020)) which are fine-tuned on one of two benchmark training sets: either SNLI (Bowman et al., 2015) or MNLI (Williams et al., 2018). Of particular interest, however, is the case where these models are trained on an additional dataset (the HELP dataset from (Yanaka et al., 2019b)) which was designed for improving the overall balance of upward and downward monotone contexts in NLI training data. We use our own random 50 – 30 – 20 train-dev-test split of the HELP dataset (ensuring unique contexts in every split), so that there is no overlap of contexts between the fine-tuning data and the few HELP-test examples we used as part of our NLI-XY dataset³.

5.2 Probing

The NLI-XY dataset is equipped with two auxiliary feature labels which are the targets of the probing task: context monotonicity and the relation of the (X, Y) phrase pair (referred to henceforth as either concept inclusion relation or lexical relation).

5.2.1 Models and Representations

For each auxiliary task, we use simple linear model architectures as the probes. We train 50 probes of varying complexities using the *probe-ably* framework (Ferreira et al., 2021). The target of the probing study is the classification token of the final

³We use the *transformers* library (Wolf et al., 2020) and their available pretrained models for this work.

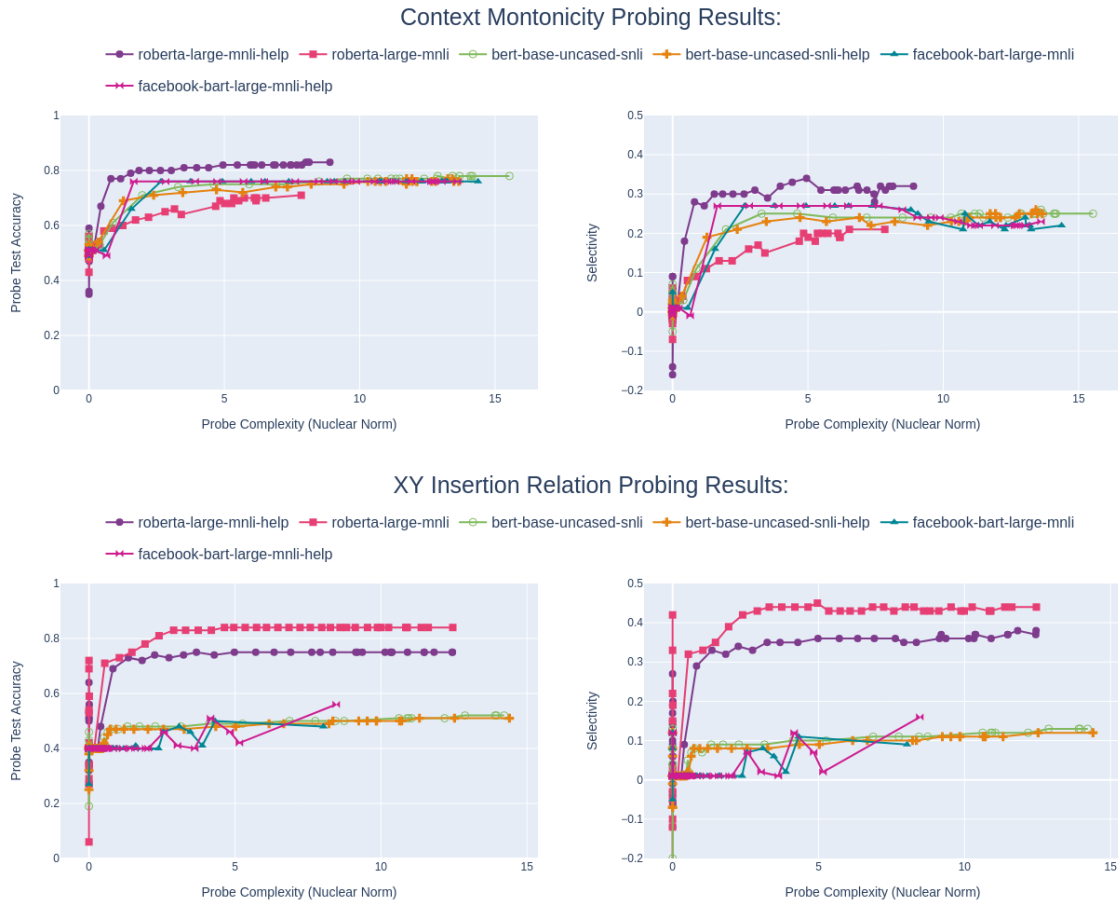


Figure 1: Probing results for all examined models.

layer of each model, as is used for the final NLI classification decision.

5.2.2 Probe Complexity Control

The complexities are represented and controlled as follows: For linear models $\hat{y} = W\mathbf{x} + \mathbf{b}$, we follow (?) in using the nuclear norm

$$\|\mathbf{W}\|_* = \sum_{i=1}^{\min(|\mathcal{T}|, d)} \sigma_i(\mathbf{W}).$$

of the matrix W as the approximate measure of complexity. In cases where the auxiliary task has a relatively large number of classes, the rank has been used as the proxy measure of model complexity (Hewitt and Manning, 2019). As the nuclear norm is a convex approximation of the *rank* of the transformation matrix, it is used in (Pimentel et al., 2020), where this allows for a larger number of informative values.

5.2.3 Metrics and Control Tasks

Accuracy and Selectivity Naively, a strong accuracy on the probing test set may be understood

to indicate strong presence of the target features within the learned representations, but there has been much discussion about whether this evidence is compelling on its own. In fact, certain probing experiments have found the same accuracy scores for random representations (Zhang and Bowman, 2018), indicating that high accuracy scores are meaningless in isolation. Hewitt and Liang (2019) describe this as a dichotomy between the representation’s encoding of the target features and the probe’s capacity for *memorization*, and propose the use of the *selectivity* measure to always place the probe accuracy in the context of a controlled probing task with shuffled labels on the same vector representations. For each fully trained probe, we report both the test accuracy and the *selectivity* measure: tracking the selectivity ensures that we are not using a probe that is complex enough to be *overly expressive* to the point of having the capacity to overfit the randomised control training set.

Control Task The *selectivity* score is calculated with respect to a *control task*. At its core, this is

NLI Models	Fine-Tuning Data	Feature Probing		NLI Monotonicity Challenge Sets		
		Context Monotonicity (%*)	XY Insertion Relation (%*)	HELP-Test (%)	MED (%)	NLI-XY (%)
roberta-large-mnli	-	71.0	84.0	36.69	46.10	59.01
roberta-large-mnli	HELP	82.0	78.0	97.63	78.22	80.68
roberta-large-mnli	HELP, HELP-Contexts	84.0	76.0	87.17	76.44	79.29
facebook/bart-large-mnli		76.0	48.0	43.61	46.54	60.59
facebook/bart-large-mnli	HELP	76.0	56.0	88.99	77.16	79.3417
bert-base-uncased-snli		77.0	50.0	63.55	0.4938	49.09
bert-base-uncased-snli	HELP	77.0	51.0	66.80	0.4613	44.79

Table 3: Summary NLI challenge test set and probing results for all considered models. *Probing results are summarized with the *accuracy at max selectivity*.

just a balanced random relabelling of the auxiliary data, but (Hewitt and Liang, 2019) advocate for more targeted control tasks with respect to the features in question and a hypothesis about the model’s possible capacity for *memorization*. For example, in their control task for POS tagging, they assign the same label to each instance of a word’s surface form (“word type”) to account for possible lexical memorization. By construction, our context monotonicity classification task is much more context-dependant and balanced: a given X insertion will occur about as often in upward and downward monotone contexts, making it harder for a probe to exploit meaningless heuristics, such as associating a given word with a context monotonicity label. For the lexical relation classification control task, we assign a shared random label for all identical insertion pairs, regardless of context.

5.3 NLI Challenge Set Evaluations

As well as the NLI-XY dataset (which can function as an ordinary NLI evaluation set), for completeness we report NLI task evaluation scores on the full MED dataset (Yanaka et al., 2019a), which was designed as a thorough stress-test of monotonicity reasoning performance. Furthermore, we report scores on the HELP-test set (from the dataset split in (Rožanova et al., 2021)): this data partition was not used in the fine-tuning of models on HELP, but we include the test scores here for insight.

5.4 Qualitative Analyses

To complement the probing and NLI results, we make two additional comparisons that may qualify the observations.

Decomposed Error Analysis The compositional structure and auxiliary labels in the NLI-XY dataset allow for extensive qualitative analysis. Firstly, we construct decomposed error analysis heatmaps which indicate whether a given premise-hypothesis data point $f(X), f(Y)$ is correctly classified by an entailment model. For brevity (and because this is representative of our observations), we include only the error breakdowns for the two subclasses of the positive entailment label: where the context monotonicity is upward and lexical relation is forward inclusion, and where the context monotonicity is downward and the lexical relation is reverse inclusion.

Representation Visualization We store the classification token ([CLS]) of the model’s last hidden layer and project it into a lower-dimensional space using the *umap* library (McInnes et al., 2018) with the default configuration. To qualify the context monotonicity probing results, we label the points according to the *gold* context monotonicity / concept relation labels.

6 Results and Discussion

6.1 Probing Results

The results for the linear probing experiments for both the *context monotonicity classification* task and the *lexical relation* classification task may be found in figure 1. The results of the control tasks are taken into account as part of the selectivity measure, which is represented on the right hand plot for each experiment. It is particularly notable that large datasets trained only on the MNLI dataset have inferior performance on context monotonicity classification. This corresponds with the further qualitative studies, suggesting that even in some of

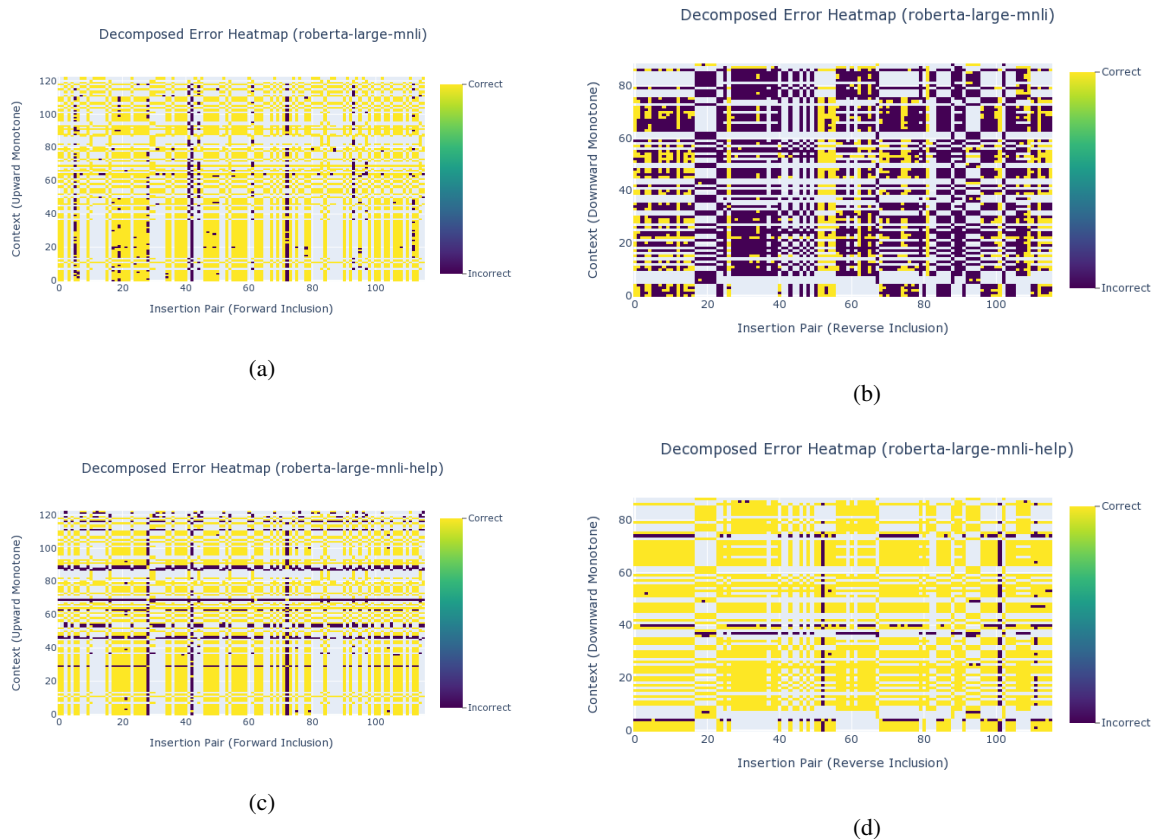


Figure 2: Decomposed error heat maps for portions of the NLI-XY dataset corresponding to the indicated context monotonicity and insertion relations (blank positions are present as only grammatical insertions were included in the dataset.)

409 the most successful transformer-based NLI models,
 410 *there is a poor “understanding” of the logical reg-*
 411 *ularities of contexts and how these are altered with*
 412 *downward monotone operators.*

413 6.2 Comparison to Challenge Set 414 Performance

415 Evaluation on the challenge test sets is relatively
 416 consistent with monotonicity probing performance,
 417 in the sense that there is a correspondence between
 418 poor/successful modeling of monotonicity features
 419 and poor/successful performance on a targeted nat-
 420 ural logic test set. As these challenge sets are fo-
 421 cused on testing monotonicity reasoning, this is a
 422 result which strongly bolsters the suggestion that
 423 explicit representation of the context monotonicity
 424 feature is crucial, especially for examples involving
 425 negation and other downward monotone operators.
 426 Furthermore, we generally confirm previous results
 427 that additional fine-tuning on the HELP data set
 428 has been helpful for these specialized test sets, and
 429 add to this that it similarly improves the explicit ex-
 430 tractability of relevant context monotonicity features

from the latent vector representations.

432 6.3 Qualitative Analyses

433 **Error Break-Downs** We are less concerned with
 434 the accuracy score (on NLI challenge sets) of a
 435 given model as with the behavioural *systematic-*
 436 *ity* visible in the errors, as we are not interested
 437 in noisy errors which may be due to words or
 438 phrases from outside the training domain. Con-
 439 sistent mis-classification for all examples derived
 440 from a fixed context or insertion pair are actually
 441 *also* strongly suggestive of a regularity in reason-
 442 ing. The decomposed error analyses paint a striking
 443 picture: we generally see that models trained on
 444 MNLI routinely fail to distinguish between the ex-
 445 pected behaviour of upward and downward mono-
 446 tone contexts, despite generally achieving high ac-
 447 curacies on large benchmark sets. This is in accor-
 448 dance with observations in Yanaka et al. (2019b)
 449 Yanaka et al. (2019a), where low accuracy on the
 450 downward-monotone reasoning sections of chal-
 451 lenge sets points to this possibility. However, they
 452 show consistently show strong behavioural regular-

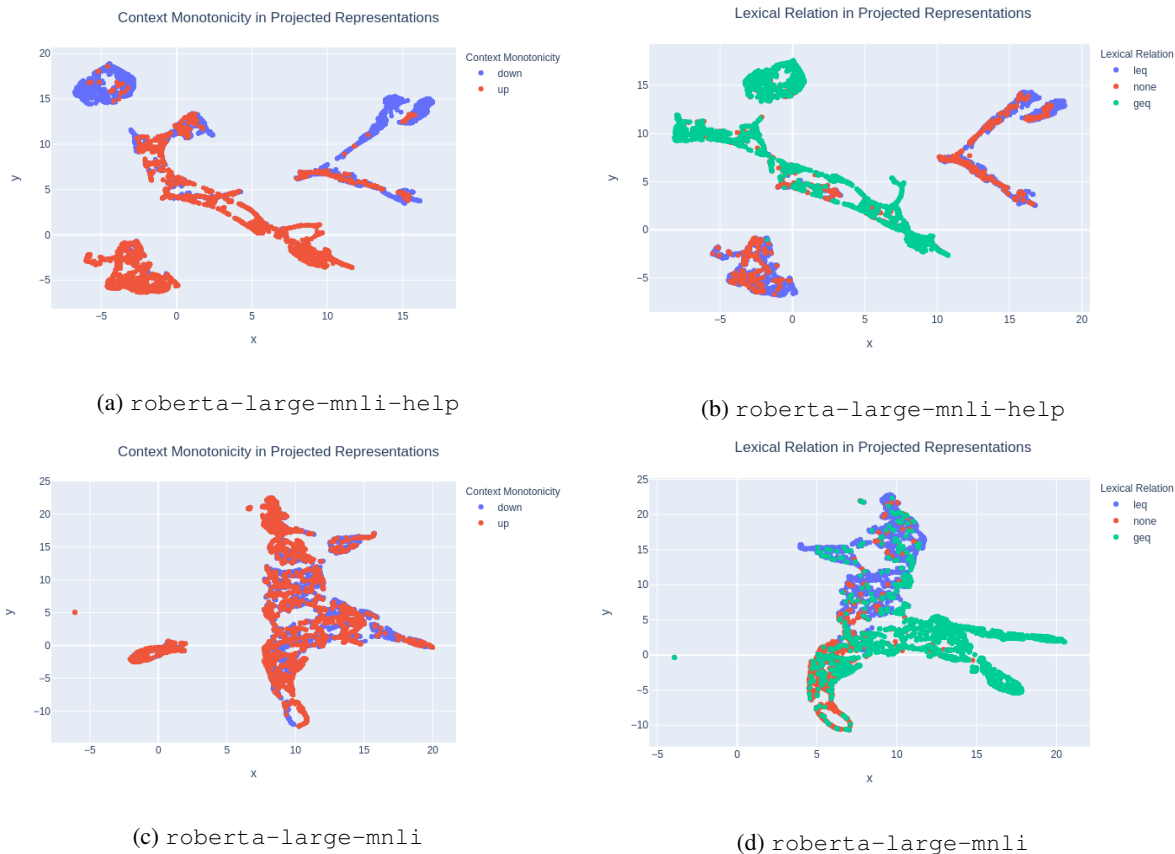


Figure 3: UMAP projections of selected classification token representations comparing `roberta-large-mnli` and the improved `roberta-large-mnli-help`, which shows greater distinction between context monotonicity features.

ity with respect to concept inclusion. Even when the contexts are downward monotone, they still treat them systematically as if they were *upward* monotone, echoing the concept insertion pair relation *only*: they completely fail to discriminate between upward/downward monotone contexts and their opposite behaviours.

Visualization Each data point corresponds to an embedded example ([CLS]) in the NLI-XY dataset, with the left and right columns colored with the *gold* auxiliary labels for context monotonicity and concept inclusion relations respectively. These illustrate the probing observations: in the well-known `roberta-large-mnli` model, concept inclusion relation features are distinguishable, whereas context monotonicity is very randomly scattered, with no emergent clustering. However, the `roberta-large-mnli-help` model shows an improvement in this behaviour, demonstrating a stronger context monotonicity distinction.

7 Conclusion

In summary, the NLI-XY has enabled us to present evidence that explicit context monotonicity feature clustering in neural model representations seems to correspond to better performance on natural logic challenge sets which test downward-monotone reasoning. In particular, many popular models trained on MNLi seem to lack this behaviour, accounting for previous observations that they systematically fail in downward-monotone contexts.

Furthermore, the probes' labels also have some explanatory value: both entailment and non-entailment labels can each further be broken down into sub-regions. This qualifies the classification with the observations that the data point occurs in a cluster of examples with a) upward (resp. downward) contexts and b) a forward (resp. backward) containment relation between the substituted noun phrases. In this sense, the analyses in this work can thus be interpreted as an explainable "decomposition" of the treatment natural logic examples in neural models.

496
497
498
499
500
501
502

503
504
505

506
507
508
509

510
511
512

513
514
515
516
517
518
519
520
521

522
523
524
525
526
527
528
529
530

531
532
533
534
535
536
537

538
539
540
541
542
543

544
545
546
547
548
549
550
551

References

Lasha Abzianidze. 2017. [LangPro: Natural language theorem prover](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 115–120, Copenhagen, Denmark. Association for Computational Linguistics.

Guillaume Alain and Yoshua Bengio. 2018. [Understanding intermediate layers using linear classifier probes](#).

Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*.

Zeming Chen, Qiyue Gao, and Lawrence S. Moss. 2021. [Neurallog: Natural language inference with joint neural and logical reasoning](#).

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.

Deborah Ferreira, Julia Rozanova, Mokanarangan Thayaparan, Marco Valentino, and André Freitas. 2021. [Does my representation capture X? probably](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 194–201, Online. Association for Computational Linguistics.

Atticus Geiger, Kyle Richardson, and Christopher Potts. 2020. [Neural natural language inference models partially embed theories of lexical entailment and negation](#). In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 163–173, Online. Association for Computational Linguistics.

Emily Goodwin, Koustuv Sinha, and Timothy J. O’Donnell. 2020. [Probing linguistic systematicity](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1958–1969, Online. Association for Computational Linguistics.

Christine Herlihy and Rachel Rudinger. 2021. [MedNLI is not immune: Natural language inference artifacts in the clinical domain](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 1020–1027, Online. Association for Computational Linguistics.

John Hewitt and Percy Liang. 2019. [Designing and interpreting probes with control tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.

John Hewitt and Christopher D. Manning. 2019. [A structural probe for finding syntax in word representations](#). 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 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.

Hai Hu, Qi Chen, Kyle Richardson, Atreyee Mukherjee, Lawrence S. Moss, and Sandra Kuebler. 2020. [MonaLog: a lightweight system for natural language inference based on monotonicity](#). In *Proceedings of the Society for Computation in Linguistics 2020*, pages 334–344, New York, New York. Association for Computational Linguistics.

Hai Hu and Larry Moss. 2018. [Polarity computations in flexible categorial grammar](#). In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 124–129, New Orleans, Louisiana. Association for Computational Linguistics.

Aikaterini-Lida Kalouli, Richard Crouch, and Valeria de Paiva. 2020. [Hy-NLI: a hybrid system for natural language inference](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5235–5249, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel. 2021. [Question and answer test-train overlap in open-domain question answering datasets](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1000–1008, Online. Association for Computational Linguistics.

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. *ArXiv*, abs/1907.11692.

Bill MacCartney and Christopher D. Manning. 2007. [Natural logic for textual inference](#). In *Proceedings of*

610	<i>the ACL-PASCAL Workshop on Textual Entailment and Paraphrasing</i> , pages 193–200, Prague. Association for Computational Linguistics.	666
611		667
612		668
613	Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3428–3448, Florence, Italy. Association for Computational Linguistics.	669
614		670
615		671
616		672
617		673
618		
619		
620	L. McInnes, J. Healy, and J. Melville. 2018. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction . <i>ArXiv e-prints</i> .	674
621		675
622		676
623	Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, Kentaro Inui, Satoshi Sekine, Lasha Abzianidze, and Johan Bos. 2019a. Can neural networks understand monotonicity reasoning? In <i>Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 31–40, Florence, Italy. Association for Computational Linguistics.	677
624		678
625		679
626		680
627		681
628		
629	Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, Kentaro Inui, Satoshi Sekine, Lasha Abzianidze, and Johan Bos. 2019b. HELP: A dataset for identifying shortcomings of neural models in monotonicity reasoning . In <i>Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)</i> , pages 250–255, Minneapolis, Minnesota. Association for Computational Linguistics.	682
630		683
631		684
632		685
633	Kelly Zhang and Samuel Bowman. 2018. Language modeling teaches you more than translation does: Lessons learned through auxiliary syntactic task analysis . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 359–361, Brussels, Belgium. Association for Computational Linguistics.	686
634		687
635		688
636		
637		
638		
639	Julia Rozanova, Deborah Ferreira, Mokbanarangan Thayaparan, Marco Valentino, and André Freitas. 2021. Supporting context monotonicity abstractions in neural nli models .	
640		
641		
642		
643		
644		
645		
646	V. Sanchez. 1991. <i>Studies on natural logic and categorial grammar</i> .	
647		
648		
649		
650		
651		
652		
653		
654	Ivan Vulić, Edoardo Maria Ponti, Robert Litschko, Goran Glavaš, and Anna Korhonen. 2020. Probing pretrained language models for lexical semantics . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7222–7240, Online. Association for Computational Linguistics.	
655		
656		
657		
658		
659		
660		
661		
662		
663		
664		
665		
666	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122. Association for Computational Linguistics.	
667		
668		
669		
670		
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		
702		
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		