

Dive into the Chasm: Probing the Gap between In- and Cross-Topic Generalization

Anonymous ACL submission

Abstract

Pre-trained Language Models (PLMs) excel in downstream tasks but may struggle when faced with limited training data. Cross-Topic experiments simulate these scenarios by withholding data from distinct topics and often expose significant performance gaps compared to experiments, where instances are randomly chosen for training and evaluation (In-Topic). Transferring superior performance between these two setups is not always feasible. To better understand this generalization gap, we propose a set of probing-based experiments and analyze various PLMs. We show that this generalization gap already exists after pre-training and differs amongst PLMs and tasks. We found that pre-training objectives and architectural regularization are keys for more generalizable and robust PLMs. In addition, we consider a set of Large Language Models (LLMs) to bootstrap the analysis of other training paradigms like prompt-based learning. Ultimately, our research leads to a better understanding of PLMs and guides in selecting appropriate models or building more robust models.¹

1 Introduction

Fine-tuning is a widely used approach for imparting new tasks to pre-trained language models (PLMs) (Devlin et al., 2019; Liu et al., 2019; He et al., 2021; Radford et al., 2019), resulting in a remarkable performance on a range of NLP tasks including GLUE (Wang et al., 2018) or SuperGLUE (Wang et al., 2019). However, evaluation setups of such benchmarks are not well aligned to real-world scenarios where data is limited or unavailable. At the same time, PLMs may not meet expectations of more realistic settings, such as Cross-Topic evaluation (Sapkota et al., 2014; Stab et al., 2018; Reuver et al., 2021; Ren et al., 2021), where generalization towards unseen topics is necessary. As a result, a generalization gap exists between the commonly

¹We provide data and code anonymized [online](#).

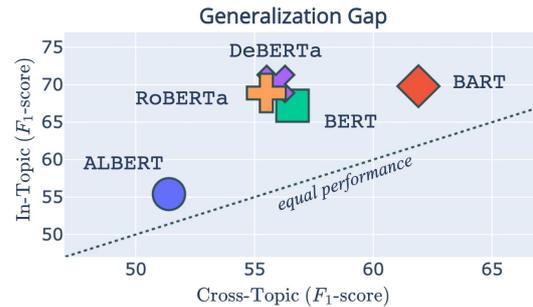


Figure 1: Generalization gap of fine-tuning PLMs on argumentative *stance detection* (Stab et al., 2018) in the In- or Cross-Topic evaluation setup. The dashed line marks the ideal case of equal performance.

used In- and the more realistic Cross-Topic evaluation setup. We see the varying semantic shift when train and test instance originating from the same (In-Topic) or different topics (Cross-Topic) as the main reason for this gap.

Exemplary, we illustrate in Figure 1 this generalization gap using our results for fine-tuning on the *UKP ArgMin* dataset (Stab et al., 2018) for In- and Cross-Topic. This dataset is part of the growing field of Argument Mining and annotates arguments being either in favor, against, or neutral towards one of eight topics like *Gun Control*. Although PLMs perform better for In-Topic, their performance differences are inconsistent across the two setups. For instance, while BART performs similarly to others for In-Topic, it outperforms them for Cross-Topic. Thus, we can not generalize findings from one setup to another. This inconsistency between these evaluation setups among different PLMs makes it challenging to draw practical conclusions - like which model to use. Therefore, we address in the paper the question *why do we see varying gaps between these two setups across PLMs?*

While probing techniques have been extensively used to study PLM behavior (Belinkov et al., 2017; Peters et al., 2018), little research has been done

on generalization gaps. Exceptions are analyzing metaphors (Aghazadeh et al., 2022) and visual probing (Zhu et al., 2022). To overcome this shortage, we suggest three probing-based experiments to study the In- vs. Cross-Topic generalization gap (§ 2) to provide ground for improving generalizability of PLMs. Firstly, we consider the prominence of differences in this generalization gap. Therefore, we evaluate this gap for three linguistic probings tasks, namely, dependency-tree parsing, part-of-speech tagging, and named-entity recognition, to measure the generalization of PLMs after pre-training and argumentative *stance detection* (UKP *ArgMin*) as a reference task (§ 4). Secondly, we investigate how PLMs depend on token-level topic information (§ 5). Finally, we re-evaluate PLMs on the probing tasks after fine-tuning on the UKP *ArgMin* dataset to better understand how this learning paradigm affects the generalization gap (§ 6). We also validate our findings with another dataset from the social media domain (Conforti et al., 2020) (see Appendix § B.1).

We show how to describe and compare PLMs in different setups and that they are more diverse than a simple evaluation could suggest. We ultimately contribute in the following way:

(1) We show that the generalization gap exists after pre-training and that diverse pre-training objectives and architectural regularization influence generalizability and robustness of PLMs. By considering Large Language Models (LLMs), we underscore the versatility and reliability of this analysis for emergent methods.

(2) We evaluate the dependence of PLMs on token-level topic information using our probing task and show that they heavily differ in how they encode and rely on token-level properties. This is crucial when we consider the susceptibility of PLMs to spurious correlations.

(3) We use the first two experiments to examine how fine-tuning in the In- and Cross-Topic setup influence PLMs and, therefore, how varying generalization gaps between them arise.

2 In- and Cross-Topic Probing

The following section formally outlines our probing setup and the used probing tasks before elaborating on the generalization gap, introducing In- and Cross-Topic probing evaluation and outline their differences.

2.1 Probing Setup and Tasks

We define in the following a probe f_p comprised of a frozen encoder h and linear classifier c without any intermediate layer. This classifier is trained to map instances $X = \{x_1, \dots, x_n\}$ to targets $Y = \{y_1, \dots, y_n\}$ for a given probing task. Using a frozen PLM as h , the probe converts x_i into a vector \mathbf{h}_i . In detail, we encode the entire sentence, which wraps x_i , and average relevant positions of x_i to find \mathbf{h}_i . Relevant positions for the considered probing task are either single tokens for *part-of-speech tagging* (POS), a span for *named entity recognition* (NER), or the concatenation of two tokens for *dependency tree parsing* (DEP). Then, the classifier c utilizes \mathbf{h}_i to generate a prediction \hat{y}_i , as shown in Equation 1.

$$\hat{y}_i = f_p(x_i) = c(h(x_i)) \quad (1)$$

2.2 Generalization Gap

Generalization gaps arise when we compare evaluation setups focusing on different capabilities for the same task. In this work, we consider evaluating a task covering various topics $T = \{t_1, \dots, t_m\}$ on instances from the same (In-Topic) or different topics (Cross-Topic) than training instances. The gap between these setups is visible when considering the semantic spread of In- and Cross-Topic test instances (in blue) in Figure 2. For Cross-Topic, these instances are less spread in the semantic space because they cover only specific topics, while In-Topic test instances are more broad spread over the space. At the same time, we note more instances with tokens never seen during training (in red) in Cross-Topic, denoted as *unseen* instances. These differences indicate more difficulties for Cross-Topic because a classifier must generalize across semantic and lexical shifts.

In an ideal case, the generalization gaps do not exist because pre-trained language models (PLMs) are robust enough to overcome such semantic shifts between different evaluation setups. However, practically, we saw in Figure 1 these gaps being pronounced on a varying scale for different models.

2.3 Difference between In- and Cross-Topic Evaluation

By evaluating probing tasks on the In- and Cross evaluation setup, we aim to understand the generalization gap better and why it differs amongst different PLMs.

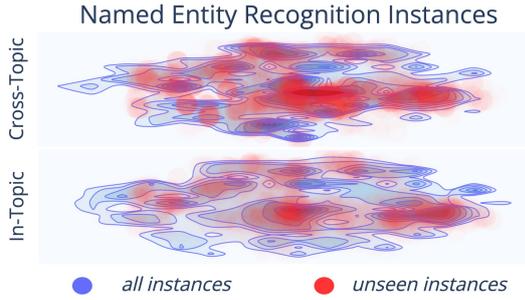


Figure 2: Density plot of the NER test split (blue) for In- and Cross-Topic, encoded with *bert-base-uncased* and reduced with the same t-SNE model (van der Maaten and Hinton, 2008). While both test splits have the same number of instances, the Cross-Topic test split has more instances (a subset of all) with *unseen* vocabulary (red) compared to In-Topic.

Cross-Topic With Cross-Topic evaluation, we investigate how well a probe generalizes when the train, dev, and test instances cover distinct sets of topics $\{T^{(train)}, T^{(dev)}, T^{(test)}\}$. In this setup, f_p must overcome a semantic shift because information about specific topics Z^T is covered only by instances from one of these splits. We formally describe this shift, denoted as ΔZ^T , as the relative complement between topic information from train and test instances - $\Delta Z^T = Z_{(train)}^T \setminus Z_{(test)}^T$. For Cross-Topic, we expect ΔZ^T to be large (Figure 2).

In-Topic In contrast, ΔZ^T is smaller for the In-Topic setup because train, dev, and test splits are randomly sampled. We expect a similar topic distribution within these splits and minor semantic differences amongst them compared to Cross-Topic (Figure 2). Thus, we expect fewer difficulties for In-Topic because a classifier does not need to generalize across semantic and lexical shifts.

Token-Level Topic Information We make differences between In- and Cross-Topic evaluation tangible using their varying divergence of topic information ΔZ^T between train and test instances - given a dataset d covering topics $T = t_1, \dots, t_m$. While being a latent property, we approximate topic information on the token-level as part of the encodings \mathbf{h}_i for a token w_i . More specifically, we adopt the approach of Kawintiranon and Singh (2021) and use the maximum log-odds-ratio $r(w_i, T)$ to capture this token-level topic information. Firstly, we calculate the odds of finding the token w_i in a topic t_j as $o(w_i, t_j) = \frac{n(w_i, t_j)}{n(\neg w_i, t_j)}$, where $n(w_i, t_j)$ is the number of occurrences of w_i in t_j , and

Model	# Params	Objectives	Data
ALBERT (Lan et al., 2020)	12M	MLM + SOP	16GB
BART (Lewis et al., 2020)	121M	DAE	160GB
BERT (Devlin et al., 2019)	110M	MLM + NSP	16GB
DeBERTa (He et al., 2021)	100M	MLM	80GB
RoBERTa (Liu et al., 2019)	110M	MLM	160GB
ELECTRA (Clark et al., 2020)	110M	MLM+DISC	16GB
GPT-2 (Radford et al., 2019)	117M	LM	40GB

Table 1: Overview of the used PLMs trained on MLM, LM, DISC, NSP, SOP, or DAE objectives.

$n(\neg w_i, t_j)$ is the number of occurrences of every other token $\neg w_i$ in t_j . We then compute r as the maximum log-odds ratio of w_i for all topics in T as $r(w_i, T) = \max_{t_j \in T} (\log(\frac{o(w_i, t_j)}{o(w_i, \neg t_j)}))$.

3 Experimental Setup

We propose three experiments to analyze how the varying generalization gap between PLMs arises by answering: *what is the generalization gap of PLMs after pre-training?* (§ 4), *what is the dependence of PLMs on topic information?* (§ 5), and *how the generalization gap evolve during fine-tuning?* (§ 6). Below, we outline relevant details about the experimental setup. Details and results are provided in the following sections.

Models We examine how various PLMs (Table 1) with varying pre-training objectives or architectural designs differ regarding our probing tasks. We cover PLMs pre-trained using masked language modeling (MLM), next sentence prediction (NSP), sentence order prediction (SOP), language modeling (LM), discriminator (DISC), and denoising autoencoder (DAE) objectives. We group them into the ones pre-trained using token- (MLM) and sentence-objectives (NSP, SOP, or DAE) and four purely token-based pre-trained (MLM, LM, DISC). We consider the base-sized variations to compare their specialties in a controlled setup. Apart from these seven contextualized PLMs, we use a static PLM with GloVe (Pennington et al., 2014).

Data We require a dataset with distinguishable topic annotations to evaluate probing tasks in the In- and Cross-Topic evaluation setup. Therefore, we mainly² rely on the *UKP ArgMin* dataset (Stab et al., 2018), which provides 25.492 arguments annotated for their argumentative stance (*pro*, *con*, or *neutral*) towards one of eight distinct topics like

²We verified our findings with another dataset in the appendix § B.1.

Nuclear Energy or *Gun Control*. Using these instances, we heuristically generate at most 40,000 instances for the three linguistic properties *dependency tree parsing (DEP)*, *part-of-speech tagging (POS)*, or *named entity recognition (NER)* using spacy.³ Additionally, we use with *stance detection (Stance)* a topic-dependent reference probe to relate the other probes and evaluate the generalization ability of PLMs on real-world tasks. We use a three-folded setup for all these four probing tasks to consider the full data variability for both In- and Cross-Topic evaluation. Details about the compositions of these folds and how we ensure a fair comparison between In- and Cross-Topic are provided in the Appendix (§ A.2) as well as examples for probing tasks (Appendix § A.1).

Evaluation We evaluate the three folds of a probing task on three random seeds to get nine measurements per task and calculate the macro averaged F_1 score to consider the variability of labels. Since recent work (Voita and Titov, 2020; Pimentel et al., 2020) questioned whether purely quantitative measures (like F_1) are enough to measure a probe’s success, we include the information compression I (Voita and Titov, 2020) for a holistic evaluation. It measures the effectiveness of a probe as the ratio ($\frac{u}{mdl}$) between uniform code length $u = n * \log_2(K)$ and minimum description length mdl , where u denotes how many bits are needed to encode n instances with label space of K . We follow *online* variation of mdl and use the same ten-time steps $t_{1:11} = \{\frac{1}{1024}, \frac{1}{512}, \dots, \frac{1}{2}\}$, where we train a probe for every t_j with a fraction of instances and evaluate with the same fraction of non-overlapping instances. Exemplary, for, t_9 we use the first fraction of $\frac{1}{4}$ instances to train and another fraction of $\frac{1}{4}$ to evaluate. We find the final mdl as the sum of the evaluation losses of every time step $t_{1:11}$. For Cross-Topic, we group training instances into two groups of distinct topics and sample the same fraction of instances to train and evaluate. Thus, we ensure the semantic shift between training and evaluation fractions.

4 The Generalization Gap of PLMs

The first experiment shows that the generalization gap already exists after pre-training and varies for the selected probing tasks and PLMs. We analyze general (Table 2 and Figure 3) and fine-grained

³We use the `en_core_web_sm` model, its evaluation metrics are available [online](#).

	DEP		POS		NER		Stance		Average		
	In	Cross	In	Cross	In	Cross	In	Cross	In	Cross	Δ
ALBERT	43.8	39.5	80.2	78.0	48.6	45.8	54.8	45.9	56.9	52.3	-4.6
BART	36.5	36.9	75.4	74.1	48.7	45.3	60.8	44.4	55.3	50.2	-5.1
BERT	25.4	25.6	68.5	67.5	45.4	41.6	56.9	43.0	49.0	44.4	-4.6
DeBERTa	32.8	29.9	73.7	74.6	48.8	42.4	59.8	45.8	53.4	48.2	-5.2
RoBERTa	25.1	23.6	64.0	65.5	48.4	42.1	51.8	40.1	47.3	42.8	-4.5
ELECTRA	33.6	33.6	75.3	75.3	41.5	41.2	46.6	43.1	49.3	48.3	-1.0
GPT-2	25.2	23.9	63.5	61.9	45.5	38.6	51.1	38.4	46.3	40.7	-5.6
GloVe	12.1	11.9	26.5	26.2	43.4	37.5	41.6	34.1	30.9	27.4	-3.5
Avg. Δ	-	-1.2	-	-0.5	-	-4.5	-	-11.0	-	-	-

Table 2: In- and Cross-Topic probing results for eight PLMs. We report the macro F_1 over three random seeds, the average difference between the two setups (last row), and their average per PLM (last two columns). Best results within a gap of 1.0 are marked by columns.

(Table 3) results and discuss them for the different evaluating setups, probing tasks, and PLMs. While we mainly focus on mid-size PLMs usable for fine-tuning, we will close this experiment by comparing them with Large Language Models (LLMs) in § 4.

Design We evaluate eight PLMs using the probe f_p (§ 2.1) on the probing tasks *DEP*, *POS*, *NER*, and *Stance*. We verified these tasks by observing significant performance drains when evaluating them on random initialized PLMs (Appendix § B.2). For a holistic evaluation, we provide results by grouping instances into two categories: *seen* and *unseen*. We define *seen* instances as already processed during training but in another context. For example, the pronoun *he* might appear in both training and test data, but in distinct sentences. By evaluating the PLMs on *seen* instances, we gain insights into the influence of token-level lexical information versus context information from surrounding tokens. In contrast, *unseen* instances were not encountered during the training of a probe. They allow assessing whether PLMs generalize to tokens that are similar to some extent (such as *Berlin* and *Washington*) but not seen during training.

Results for Evaluation Setups Upon analyzing Table 2, we observe a clear generalization gap between In- and Cross-Topic evaluation for all tasks and PLMs. As shown in Figure 3, the magnitude of this gap (ΔF_1) correlates with the degree of compression drain (ΔI). Interestingly, we find a stronger correlation between F_1 and I for Cross-Topic ($\rho = 0.72$) as compared to In-Topic ($\rho = 0.69$). A higher performance level, like for In-Topic, leaves less room for compression improvements.

Further, we examine the performance of *seen* and *unseen* instances in Table 3. It shows that *seen*

	DEP			POS			NER			
	<i>all</i>	Δ <i>seen</i>	Δ <i>unseen</i>	<i>all</i>	Δ <i>seen</i>	Δ <i>unseen</i>	<i>all</i>	Δ <i>seen</i>	Δ <i>unseen</i>	
	-	85%	15%	-	86%	14%	-	65%	35%	
In-Topic	Instance Ratio	-	-	-	-	-	-	-	-	
	ALBERT	43.8	+0.21	-3.2	80.2	+0.41	-17.7	48.6	+1.1	-5.8
	BART	36.5	+0.13	-3.0	75.4	+0.20	-16.5	48.7	+1.3	-7.0
	BERT	25.4	-0.02	-0.8	68.5	+0.20	-16.5	45.4	+1.0	-5.8
	DeBERTa	32.8	+0.07	-1.5	73.7	+0.09	-12.7	48.8	+1.0	-5.6
	RoBERTa	25.1	-0.01	-0.9	64.0	-0.04	-15.5	48.4	+1.0	-5.7
	Average	-	-0.08	-1.9	-	+0.17	-15.8	-	+1.1	-6.0
Cross-Topic	Instance Ratio	-	78%	22%	-	81%	19%	-	51%	49%
	ALBERT	39.5	+0.03	-2.3	78.0	+0.51	-12.9	45.8	+2.2	-5.3
	BART	36.9	+0.01	-4.0	74.1	+0.24	-16.5	45.3	+2.4	-5.8
	BERT	25.6	-0.09	-0.7	67.5	+0.20	-14.0	41.6	+1.9	-5.1
	DeBERTa	29.9	-0.07	-1.3	74.6	+0.14	-11.7	42.4	+2.0	-5.2
	RoBERTa	23.6	-0.22	-0.3	65.5	+0.00	-14.7	42.1	+1.9	-5.2
	Average	-	-0.08	-1.7	-	+0.22	-14.0	-	+2.1	-5.3

Table 3: Performance difference of *seen* and *unseen* instances compared to the full set (*all*). We report for ALBERT, BART, BERT, DeBERTa, & RoBERTa, and include the ratio of *seen* and *unseen* instances.

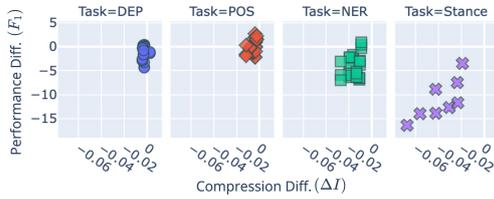


Figure 3: Comparison of the difference in ΔF_1 and ΔI between Cross-Topic and In-Topic for all eight PLMs on the four probing tasks.

performs slightly better than *all*, while *unseen* ones underperform the complete set (*all*) and *seen* instances. Considering the average over PLMs, there are fewer relative gains for *seen* for In-Topic and more loss for *unseen* instances (+1.2, -6.0 for *NER*) compared to Cross-Topic (+2.0, -5.3 for *NER*). This observation relates to the lower percentage of *unseen* instances (i.e., made of topic-specific terms) for In- compared to Cross-Topic. We see *unseen* instances on In-Topic are harder and cover rare vocabulary, and *seen* instances on Cross-Topic are easier and made of general terms. These results confirm our theoretic and semantic assumptions (§ 2).

Results for Probing Tasks Considering Table 2 and Figure 3, we note higher generalization gaps (Avg. Δ of -4.5 and -11.0) for semantic tasks (*NER* & *Stance*) than for syntactic tasks (*DEP* & *POS*) - Avg. Δ of -1.2 and -0.5. We verify this trend with results in the Appendix (§ B.5), where we observe a more pronounced gap for semantic *NER* classes (like *ORG*) than for syntactic ones - like *ORDINAL*.

Next, we separately compare tasks for *seen* and *unseen* instances. *DEP* shows the slightest performance difference compared to *all*. We assume this

is due to the pairwise task nature, which leads to a larger shared vocabulary between *unseen* and training instances. We assume frequent words (like *of*) are part of the *unseen* instances. In contrast, apparent differences between *NER* and *POS* are visible - with less performance drain on *unseen* instances for *NER* than *POS*. Therefore, we assume for *NER* a higher semantic overlap with training instances since they could include - as being an n-gram - words from the training vocabulary. In contrast, tokens of *unseen POS* instances are always single words; thus, we assume a smaller semantic overlap with the training.

Results for Encoding Models We now compare PLMs amongst themselves. The four best-performing PLMs of In-Topic differ up to 7.6 (ALBERT - BERT), while for Cross-Topic, this difference narrows to 4.1 (ALBERT - ELECTRA). These results confirm the varying generalization gap between them and, again, that we can not transfer conclusions from one evaluation setup to another. For example, the probing performance of BART for In-Topic *Stance* is the best and the third best for Cross-Topic.

Generally, we do not see a clear correlation between better average performance and a smaller generalization gap. PLMs like DeBERTa perform better for In- and Cross-Topic but show a bigger gap (-5.1) compared to lower performing PLMs like ELECTRA (-1.0), but there are also worse PLMs with a bigger gap (GPT-2, -5.6) or better ones with a smaller gap (ALBERT, -4.6). Overall, we see the generalization gap being more pronounced for better-performing PLMs.

Considering absolute performance, ALBERT & BART performs the best on average for both evaluation setups, while ELECTRA excels *POS* and *DEP*, and DeBERTa performs for *NER* & *Stance*. In contrast, BERT, RoBERTa, GPT-2, and GloVe underperform the others. Thus, PLMs with architectural regularization, such as layer-wise parameter sharing (ALBERT), encoder-decoder layers (BART), disentangled attention (DeBERTa), or discriminator (ELECTRA), tend to provide higher Cross-Topic performance. Similarly, regularized PLMs, such as ALBERT or DeBERTa, generally achieve more performance gains for *seen* instances and fewer performance drops for *unseen* ones than models without regularization such as BERT or RoBERTa. We hypothesize that architectural and regularization aspects equip PLMs with a more

	DEP		POS		NER		Stance		Average		
	In	Cross	In	Cross	In	Cross	In	Cross	In	Cross	Δ
ALBERT	43.8	39.5	80.2	78.0	48.6	45.8	54.8	45.9	56.9	52.3	-4.6
BART	36.5	36.9	75.4	74.1	48.7	45.3	60.8	44.4	55.3	50.2	-5.1
T5 (3B)	33.9	32.5	68.5	68.9	48.3	42.2	53.2	42.1	51.0	46.5	-4.5
FLAN-T5 (3B)	33.1	29.7	66.8	66.9	48.5	43.1	56.0	45.1	51.1	46.2	-4.9
GPT-Neo (2.7B)	36.4	33.1	76.4	77.1	52.9	49.6	62.4	40.5	57.0	50.1	-6.9

Table 4: Results (macro F_1) of the four probing tasks using the two overall best-performing PLMs (ALBERT and BART) in the In- and Cross-Topic setup based on the *ArgMin* dataset (Table 2) and three LLMs.

generalizable and robust encoding space.

Results for Larger Models We compare in Table 4 three relevant and open accessible LLMs with the two best performing models (ALBERT and BART) on the first experiments. In general, we see that the scaling law (Brown et al., 2020) applies to our setting for LLMs with LM-based pre-training. Specifically, GPT-Neo (2.7B) (Black et al., 2021) is more robust and outperforms GPT-2 while performing on par or slightly better than the other PLMs. In contrast, T5 (3B) (Raffel et al., 2022) or FLAN-T5 (3B) (Chung et al., 2022) underperform PLMs on syntactic tasks and perform slightly worse on semantic tasks. We hypothesize that their task-specific pre-training result in less robust and generalizable token encoding space. This is in line with the fact that amongst these two LLMs, FLAN-T5 (3B) performs worse than T5 (3B), which experienced additional instruction-based pre-training.

5 Dependence of PLMs on Topic Information

To this point, we saw that the generalization gap varies between different PLMs and probing tasks. Next, we demonstrate in Figure 4 using *Amnesic Probing* (Elazar et al., 2021) that PLMs differ on how they rely on token-level topic information - a crucial aspect to distinguish between In- and Cross-Topic evaluation. These findings show how PLMs differ beyond an evaluation score and emphasize the consideration of various PLMs when using *Amnesic Probing*. Additional insights of comparing *seen* and *unseen* instance and distinct NER classes are provided in the Appendix (§ B.4, § B.6).

Design This property \mathbf{z}_i^T influences the generalization gap since it crucially defines differences between In- and Cross-Topic evaluation (§ 2.3) To measure the influence of \mathbf{z}_i^T , we employ *Amnesic Probing* (Elazar et al., 2021) to approximate its influence by removing it from \mathbf{h}_i . More precisely,

we compare the performance of a probe f_p with an amnesic probe $f_p^{\setminus T}$ - without token-level topic information - for a given probing task like *NER*. A negative effect indicates a dependence on-topic information, while it is a hurdle when performance improves. We first train a linear model l on token-level topic information $r_{(w_i, T)}$ of a token w_i for a set of topics T . To shape it as a classification task, we categorize r into three classes (*low*, *medium*, *high*).⁴ Next, we find a projection matrix P that projects the encoded input H to the null space, where $W_l P H = 0$. This enables us to define the amnesic probe $f_p^{\setminus T}(x_i) = c(\mathbf{h}_i \setminus \mathbf{z}_i^T) = c(P\mathbf{h}_i)$, which utilizes P to neutralize the topic information T represented by \mathbf{z}_i^T . Following (Elazar et al., 2021), we verified our results by measuring less effect of removing random information from \mathbf{h}_i (see Appendix § B.3).

Results Considering Figure 4, we see ALBERT, BART, & BERT depend less on topic information. We see their diverse pre-training (token- and sentence-objectives or sentence denoising) results in a more robust embedding space. Surprisingly, they show positive effects (3.2 for *DEP* for BART) when removing topic information. This could remove potentially disturbing parts of the embedding space. Similarly, GPT-2 is less affected by removing topic information - we assume this is due to its generally lower performance level. Thus, there is less room for performance drain, and capturing topic information is less powerful.

Comparing In- and Cross-Topic setups shows a narrowing generalization gap for more affected models (like RoBERTa and GloVe on *NER* or *Stance*). Simultaneously, less affected PLMs either maintain the gap or enlarge it slightly - like BART on *DEP*, *NER*, or *Stance*. Further, DeBERTa, RoBERTa, ELECTRA, and GloVerely more on topic information since they show significant performance loss (up to 34.6 for GloVe on *POS*) when removing this information. Specifically, GloVe as a static language model, and RoBERTais affected the highest for all tasks. ELECTRA shows similar behavior, but is less pronounced for *POS*. Thus, its reconstruction pre-training objective provides a more robust embedding space than purely MLM (DeBERTa or RoBERTa). Comparing DeBERTa and RoBERTa, DeBERTa is less affected by the removal of semantic tasks (*NER* and *Stance*).

⁴Please find examples in the Appendix § A.6.

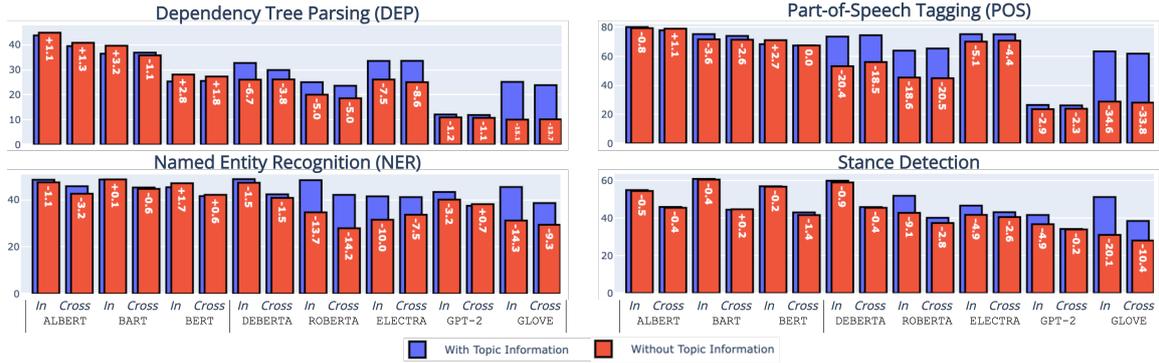


Figure 4: Comparison of the probing results with (blue bars) or without (red bars) topic information. The white text indicates the difference between these two scenarios ($\Delta F_1^{\setminus T}$).

We hypothesize that distinguishing between token content and token position via disentangled attention makes DeBERTa more robust for the semantic than for syntactic tasks (*DEP* & *POS*).

6 Evolution of the Generalization Gap during Fine-Tuning

Finally, we re-evaluate fine-tuned PLMs using the previously two probing setup and show that fine-tuning leads to a drain in probing performance. We use these results to retrace apparent differences between evaluation setups and the varying generalization gap between PLMs. This is relevant for a broader understanding of how fine-tuning affects PLMs (Mosbach et al., 2020; Kumar et al., 2022a), and what they learn during fine-tuning (Merendi et al., 2022; Ravichander et al., 2021).

Design We fine-tune the PLMs on an argumentative *stance detection* task and re-evaluate them on the probing tasks *DEP*, *POS*, and *NER*. To be consistent with our probing setup, we used the same folds for fine-tuning. Further details are in the Appendix (§ A.5). We compare these results with the probing performance of their pre-trained counterparts (§ 4 & § 5) and correlate this change with the generalization gap observed on the downstream task. We limit our analysis to ALBERT, BERT, BART, DeBERTa, and RoBERTa.

Results Table 5 shows that fine-tuning clearly boost the performance on *Stance* compared to the probing performance (§ 4) but leads to a clear performance drop (ΔF_1) for both evaluation setups and the probing tasks. Cross-Topic achieved more gains on average (+12.6) and fewer drains (-17.1) on the three linguistic properties than In-Topic (+9.5, -20.4). On average, we assume that

		<i>Stance</i>	<i>DEP</i>	<i>POS</i>	<i>NER</i>	<i>Avg.</i>	<i>DEP</i>	<i>POS</i>	<i>NER</i>
		F_1 fine-tuned	ΔF_1 probing				$\Delta F_1^{\setminus T}$		
In-Topic	ALBERT	55.4 +0.6	-27.3	-40.2	-25.0	-30.8	-0.6	-3.0	-0.1
	BART	69.8 +9.0	-17.3	-32.2	-4.0	-17.8	-0.8	-4.0	+0.3
	BERT	67.2 +10.3	-7.5	-24.8	+1.0	-10.4	+0.4	+0.7	+1.1
	DeBERTa	70.1 +10.3	-13.2	-25.3	-8.8	-15.8	-0.8	-3.8	-0.4
	RoBERTa	68.9 +17.1	-19.7	-48.6	-29.7	-27.2	-0.8	-3.0	-0.7
	Avg.	66.3 +9.5	-16.6	-32.6	-12.1	-20.4	-0.5	-2.6	+0.1
Cross-Topic	ALBERT	51.4 +5.5	-14.4	-20.3	-12.6	-15.8	+1.6	-1.3	+2.1
	BART	61.9 +17.5	-16.5	-33.9	-5.4	-18.6	-1.0	-3.5	-1.6
	BERT	56.6 +13.6	-5.7	-19.5	+0.6	-8.2	+0.7	+0.6	+1.2
	DeBERTa	55.9 +10.1	-13.4	-33.4	-11.8	-19.5	-1.2	-8.6	+1.6
	RoBERTa	55.5 +15.4	-16.6	-48.3	-23.1	-23.5	-1.9	-4.8	-0.3
	Avg.	56.3 +12.6	-13.0	-29.3	-9.1	-17.1	-0.4	-3.5	+0.6

Table 5: Results of evaluating our probing setup on fine-tuned PLMs on *Stance*. The first column shows these fine-tuned results and the gained improvement compared to probing for *Stance* on pre-trained PLMs (Table 2). Next, we show performance differences between pre-trained and fine-tuned PLMs (ΔF_1 probing) and how removing topic information affects the fine-tuned PLMs ($\Delta F_1^{\setminus T}$).

In-Topic fine-tuning affects the encoding space of PLMs more heavily than Cross-Topic. Regarding the different probing tasks, the performance drain is more pronounced for syntactic tasks (*DEP* and *POS*) than semantic tasks (*NER*). This hints that PLMs acquire competencies of semantic nature - which holds for *stance detection*. Similarly, removing topic information influences fine-tuned PLMs the least for *NER*. On the same time, this removal is more pronounced for Cross-Topic. This confirms the assumption that the Cross-Topic setup has smaller effects on PLMs internals, since we saw big impacts of this removal (§ 5).

Considering the single PLMs, we see apparent differences. For example, ALBERT, with its shared architecture and priorly best-performing PLM, experiences big probing performance drains and the smallest fine-tuning gains (+0.6, +5.5). In con-

trast, we note effective fine-tuning of BERT with +10.3 for In- and +13.6 for Cross-Topic, and that it lost the least probing performance. Comparing RoBERTa and DeBERTa reveals again the effectiveness of architectural regularization of DeBERTa. RoBERTa shows the most gains when fine-tuning on *Stance* and almost catching up with DeBERTa. However, it experiences a more clear performance drain (-27.2, -23.5) regarding the probing tasks for In- and Cross-Topic compared to DeBERTa (-15.8, -19.5). Next, we focus on BART and its superior Cross-Topic performance on *Stance*. It seems already well-equipped for this downstream task due to its high In-Topic probing performance on *Stance*. This allows it to focus on overcoming the semantic shift during fine-tuning.

7 Related Work

The rise of PLMs (Devlin et al., 2019; Liu et al., 2019; Radford et al., 2019; He et al., 2021) enabled big success on a wide range of tasks (Wang et al., 2018, 2019). Nevertheless, they still fall behind on more realistic Cross-Topic, like generalizing towards unseen topics (Stab et al., 2018; Gulrajani and Lopez-Paz, 2021; Allaway and McKeown, 2020). One primary reason is that PLMs often rely on unwanted spurious correlations. Despite PLMs saw such vocabulary during pre-training, they failed to consider test vocabulary in the required fine-grained way (Thorn Jakobsen et al., 2021; Reuver et al., 2021). Further, Kumar et al. (2022b) found linear models can outperform fine-tuning PLMs when considering out-of-distribution data. Thus, a broader understanding of PLMs in challenging evaluation setups is crucial.

Probing (Belinkov et al., 2017; Conneau et al., 2018; Peters et al., 2018) helps to analyze innards of PLMs. This includes to examine how linguistic (Tenney et al., 2019a,b), numeric (Wallace et al., 2019), reasoning (Talmor et al., 2020), or discourse (Koto et al., 2021) properties are encoded. Other works focus on specific properties used for other tasks (Elazar et al., 2021; Lasri et al., 2022), or fine-tuning dynamics (Merchant et al., 2020; Zhou and Srikumar, 2022; Kumar et al., 2022b). However, these works target the commonly used *In-Topic* setup and less work considering Cross-Topic setups. Aghazadeh et al. (2022) analyzed metaphors across domains and language, or Zhu et al. (2022) cross-distribution probing for visual tasks. They found that models generalize to some extent across dis-

tribution shifts in probing-based evaluation. Nevertheless, these works focus on specialized tasks and consider the generalizations across distributions in isolation. In contrast, we propose with our experiments a holistic evaluation of PLMs with the focus on the generalization gap between having a distribution shift (Cross-Topic) or not (In-Topic) and how this gap change after fine-tuning.

8 Conclusion

Discussion We demonstrated with our experiments how we approximate internals of PLMs using probing in different evaluation setups. We saw that varying In- vs. Cross-Topic generalization gaps already exist after pre-training and are more pronounced for semantic tasks than syntactic ones. We found PLMs with architectural regularization and diverse pre-training objectives providing more generalizable and robust encoding space. For example, they are less influenced by token-level interventions (like removing topic information). By re-evaluating PLMs after fine-tuning using this analysis, we found that generalization gaps of PLMs evolves differently during fine-tuning and that regularization contributes to narrowing this gap. Finally, we provide preliminary insights into different LLMs using our proposed setup to bootstrap further investigations. We verified our results using a second dataset from the social media domain (Conforti et al., 2020) - details in the Appendix § B.1.

To conclude, we shed light on the behavior of PLMs regarding the generalization gap and outline the multiplicity of PLMs. We suggest ways to improve their robustness and performance to close this generalization gap and identify promising PLMs. For example, how a better In-Topic probing performance on the downstream task allows PLMs (like BART) to overcome the semantic shift of Cross-Topic evaluation.

Outlook We extended probing In- and Cross-Topic to analyze the varying generalization gap of PLMs. With our findings in mind, we see extensively probing PLMs and LLMs as indispensable of a holistic evaluation of their verity. Therefore, regularly re-evaluating them on new tasks and forthcoming learning paradigms is crucial. In future work, we continue to analyze PLMs and LLMs in more depth, like how they generalize in instruction- or prompting-based scenarios or how pre-training objectives shape them differently.

Ethical Considerations and Limitations

Automatic Annotations for Linguistic Properties Our experiments require all instances origin in the same datasets with topic annotations. Thanks to this condition, we align all our experiments, like probing PLMs, with the same data as they got pre-trained. Therefore, we minimize other influences like semantic shifts of other datasets. However, there are no corresponding annotations for linguistic properties, which forces us to rely on automatically gathered annotations. This work addresses this issue by transparently stating the libraries and models we used to derive these annotations and providing the source code and the extracted labels in our repository. Further, we see the results of *DEP*, *POS*, and *NER* well aligned with previous work (Tenney et al., 2019a,b; Hewitt and Liang, 2019)

Definition of Topic Information Within this work, we consider topic information in the context of a given dataset. As previously mentioned, this focus on the context of one dataset allows in-depth analysis, like examining the change of PLMs during fine-tuning. On the other hand, we need to thoroughly re-evaluate other datasets since the semantic space and granularity of the topic are different in almost every other dataset. However, results in the Appendix (§ B.1) let us assume that our findings correlate with other datasets and domains. Further, we approximate topic information on the token-level as done previously in literature (Kawintiranon and Singh, 2021). We think that considering n-grams could give a better approximation of topic-specific terms, but we do not take them into account because *Amnesic Probing* (Elazar et al., 2021) require token-level properties to apply resulting intervention on token-level tasks like *POS*.

Impact of PLMs Design choices We analyze in this work PLMs with various properties like different pre-training objectives or architectural regularization. While we see our results influenced by these specialties, we use them to explain results on an aggregated level and not on single aspects like the usefulness of the *Sentence-Order* objective. Such experiments would require significant computational resources to pre-train models to verify such properties with full certainty. Instead, we use same-sized model variations, evaluate all probes on three folds and three random seeds to account for data variability and random processes, and verify

our results on a second dataset.

Transfer to other Learning Paradigm With this work, we mainly focus on pre-trained language models like BERT or RoBERTa in the context of the fine-tuning learning paradigm. Therefore, we just preliminary cover recently proposed Large Language Models (LLMs) like GPT-Neo (2.7B) (Black et al., 2021). Nevertheless, we see our approach as adaptable to instruction- or prompt-based learning (Liu et al., 2021).

References

- Ehsan Aghazadeh, Mohsen Fayyaz, and Yadollah Yaghoobzadeh. 2022. *Metaphors in pre-trained language models: Probing and generalization across datasets and languages*. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2037–2050, Dublin, Ireland. Association for Computational Linguistics.
- Emily Allaway and Kathleen McKeown. 2020. *Zero-Shot Stance Detection: A Dataset and Model using Generalized Topic Representations*. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8913–8931, Online. Association for Computational Linguistics.
- Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. *What do neural machine translation models learn about morphology?* In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics.
- Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. 2021. *GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow*. If you use this software, please cite it using these metadata.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. *Language models are few-shot learners*. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

738	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> .	
739		
740		
741		
742		
743	Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: Pre-training text encoders as discriminators rather than generators . In <i>ICLR</i> .	
744		
745		
746		
747	Costanza Conforti, Jakob Berndt, Mohammad Taher Pilehvar, Chryssi Giannitsarou, Flavio Toxvaerd, and Nigel Collier. 2020. Will-they-won't-they: A very large dataset for stance detection on Twitter . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1715–1724, Online. Association for Computational Linguistics.	
748		
749		
750		
751		
752		
753		
754		
755	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 <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.	
756		
757		
758		
759		
760		
761		
762		
763	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>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)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	
764		
765		
766		
767		
768		
769		
770		
771		
772	Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2021. Amnesic probing: Behavioral explanation with amnesic counterfactuals . <i>Transactions of the Association for Computational Linguistics</i> , 9:160–175.	
773		
774		
775		
776		
777	Ishaan Gulrajani and David Lopez-Paz. 2021. In search of lost domain generalization . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	
778		
779		
780		
781		
782	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention . In <i>International Conference on Learning Representations</i> .	
783		
784		
785		
786	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.	
787		
788		
789		
790		
791		
792		
	Kornrathop Kawintiranon and Lisa Singh. 2021. Knowledge enhanced masked language model for stance detection . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4725–4735, Online. Association for Computational Linguistics.	793
		794
		795
		796
		797
		798
		799
	Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2021. Discourse probing of pretrained language models . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021</i> , pages 3849–3864. Association for Computational Linguistics.	800
		801
		802
		803
		804
		805
		806
		807
	Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. 2022a. Fine-tuning can distort pretrained features and underperform out-of-distribution . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	808
		809
		810
		811
		812
		813
	Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. 2022b. Fine-tuning can distort pretrained features and underperform out-of-distribution . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	814
		815
		816
		817
		818
		819
	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A lite BERT for self-supervised learning of language representations .	820
		821
		822
		823
	Karim Lasri, Tiago Pimentel, Alessandro Lenci, Thierry Poibeau, and Ryan Cotterell. 2022. Probing for the usage of grammatical number . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8818–8831, Dublin, Ireland. Association for Computational Linguistics.	824
		825
		826
		827
		828
		829
		830
	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 <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	831
		832
		833
		834
		835
		836
		837
		838
		839
	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing . <i>CoRR</i> , abs/2107.13586.	840
		841
		842
		843
		844
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach . <i>ArXiv</i> , abs/1907.11692.	845
		846
		847
		848
		849

850	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	<i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021</i> , pages 3363–3377. Association for Computational Linguistics.	906 907 908 909 910
855	Amil Merchant, Elahe Rahimtoroghi, Ellie Pavlick, and Ian Tenney. 2020. What happens to BERT embeddings during fine-tuning? In <i>Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 33–44, Online. Association for Computational Linguistics.	Xiaoying Ren, Jing Jiang, Ling Min Serena Khoo, and Hai Leong Chieu. 2021. Cross-topic rumor detection using topic-mixtures . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1534–1538, Online. Association for Computational Linguistics.	911 912 913 914 915 916 917
861	Federica Merendi, Felice Dell’Orletta, and Giulia Venturi. 2022. On the nature of BERT: correlating fine-tuning and linguistic competence . In <i>Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022</i> , pages 3109–3119. International Committee on Computational Linguistics.	Myrthe Reuver, Suzan Verberne, Roser Morante, and Antske Fokkens. 2021. Is stance detection topic-independent and cross-topic generalizable? - a reproduction study . In <i>Proceedings of the 8th Workshop on Argument Mining</i> , pages 46–56, Punta Cana, Dominican Republic. Association for Computational Linguistics.	918 919 920 921 922 923 924
868	Marius Mosbach, Anna Khokhlova, Michael A. Hedderich, and Dietrich Klakow. 2020. On the Interplay Between Fine-tuning and Sentence-level Probing for Linguistic Knowledge in Pre-trained Transformers . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2502–2516, Online. Association for Computational Linguistics.	Upendra Sapkota, Tamar Solorio, Manuel Montes, Steven Bethard, and Paolo Rosso. 2014. Cross-topic authorship attribution: Will out-of-topic data help? In <i>Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers</i> , pages 1228–1237, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.	925 926 927 928 929 930 931 932
875	Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.	Christian Stab, Tristan Miller, Benjamin Schiller, Pranav Rai, and Iryna Gurevych. 2018. Cross-topic argument mining from heterogeneous sources . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3664–3674, Brussels, Belgium. Association for Computational Linguistics.	933 934 935 936 937 938 939
881	Matthew E. Peters, Mark Neumann, Luke Zettlemoyer, and Wen-tau Yih. 2018. Dissecting contextual word embeddings: Architecture and representation . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1499–1509, Brussels, Belgium. Association for Computational Linguistics.	Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2020. oLMpics-On What Language Model Pre-training Captures . <i>Transactions of the Association for Computational Linguistics</i> , 8:743–758.	940 941 942 943
888	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.	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019a. BERT rediscovers the classical NLP pipeline . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4593–4601, Florence, Italy. Association for Computational Linguistics.	944 945 946 947 948 949
895	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2019b. What do you learn from context? probing for sentence structure in contextualized word representations . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	950 951 952 953 954 955 956 957 958
898	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2022. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21(1).	Terne Sasha Thorn Jakobsen, Maria Barrett, and Anders Søggaard. 2021. Spurious correlations in cross-topic argument mining . In <i>Proceedings of *SEM 2021: The Tenth Joint Conference on Lexical and</i>	959 960 961 962

- 963 *Computational Semantics*, pages 263–277, Online.
964 Association for Computational Linguistics.
- 965 Laurens van der Maaten and Geoffrey E. Hinton. 2008.
966 Visualizing data using t-sne. *Journal of Machine*
967 *Learning Research*, 9:2579–2605.
- 968 Elena Voita and Ivan Titov. 2020. [Information-theoretic](#)
969 [probing with minimum description length](#). In *Pro-*
970 *ceedings of the 2020 Conference on Empirical Meth-*
971 *ods in Natural Language Processing (EMNLP)*,
972 pages 183–196, Online. Association for Computa-
973 tional Linguistics.
- 974 Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh,
975 and Matt Gardner. 2019. [Do NLP models know num-](#)
976 [bers? probing numeracy in embeddings](#). In *Proceed-*
977 *ings of the 2019 Conference on Empirical Methods*
978 *in Natural Language Processing and the 9th Inter-*
979 *national Joint Conference on Natural Language Pro-*
980 *cessing, EMNLP-IJCNLP 2019, Hong Kong, China,*
981 *November 3-7, 2019*, pages 5306–5314. Association
982 for Computational Linguistics.
- 983 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-
984 preet Singh, Julian Michael, Felix Hill, Omer Levy,
985 and Samuel Bowman. 2019. [Superglue: A stickier](#)
986 [benchmark for general-purpose language understand-](#)
987 [ing systems](#). In *Advances in Neural Information*
988 *Processing Systems*, volume 32. Curran Associates,
989 Inc.
- 990 Alex Wang, Amanpreet Singh, Julian Michael, Felix
991 Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE:](#)
992 [A multi-task benchmark and analysis platform for nat-](#)
993 [ural language understanding](#). In *Proceedings of the*
994 *2018 EMNLP Workshop BlackboxNLP: Analyzing*
995 *and Interpreting Neural Networks for NLP*, pages
996 353–355, Brussels, Belgium. Association for Com-
997 putational Linguistics.
- 998 Yichu Zhou and Vivek Srikumar. 2022. [A closer look](#)
999 [at how fine-tuning changes BERT](#). In *Proceedings*
1000 *of the 60th Annual Meeting of the Association for*
1001 *Computational Linguistics (Volume 1: Long Papers)*,
1002 pages 1046–1061, Dublin, Ireland. Association for
1003 Computational Linguistics.
- 1004 Zining Zhu, Soroosh Shahtalebi, and Frank Rudzicz.
1005 2022. [OOD-probe: A neural interpretation of out-of-](#)
1006 [domain generalization](#). In *ICML 2022: Workshop on*
1007 *Spurious Correlations, Invariance and Stability*.

A Additional Details of the Experiments

A.1 Probing Tasks

Table 14 shows examples and additional details of the different probing tasks.

A.2 Fold Composition

We rely on a three-folded evaluation for In- and Cross-Topic for a generalized performance measure. These folds cover every instance exactly once in a test split. In addition, we require that In- and Cross-Topic train/dev/test splits have the same number of instances for a fair comparison, as visualized in Figure 5. For Cross-Topic, we make sure that every topic $\{t_1, \dots, t_m\}$ is covered precisely once by one of the three test splits $X_{cross}^{(test)}$. To compose $X_{cross}^{(train)}$ and $X_{cross}^{(dev)}$, we randomly distribute the remaining topics for every fold. For In-Topic, we randomly⁵ form subsequent test splits $X_{in}^{(test)}$ for every fold from all instances $\{x_1, \dots, x_m\}$. $X_{in}^{(train)}$ and $X_{in}^{(dev)}$ are then randomly composed for every fold using the remaining instance set following the dimension of $X_{cross}^{(train)}$ and $X_{cross}^{(dev)}$.

A.3 Training Setup

For all our experiments, we use NVIDIA RTX A6000 GPUs, python (3.8.10), transformers (4.9.12), and PyTorch (1.11.0).

A.4 Probing Hyperparameters

Further, we use for the training of the probes the following fixed hyperparameters: 20 epochs, where we find the best one using dev instances; AdamW (Loshchilov and Hutter, 2019) as optimizer; a batch size of 64; a learning rate of 0.0005; a dropout rate of 0.2; a warmup rate of 10% of the steps; random seeds: [0, 1, 2]

In addition, we use the following tags from the huggingface model hub:

- [albert-base-v2](#)
- [bert-base-uncased](#)
- [facebook/bart-base](#)
- [microsoft/deberta-base](#)
- [roberta-base](#)

⁵We expect that all folds cover all topics given the small number of topics (8) and the big number of instances.

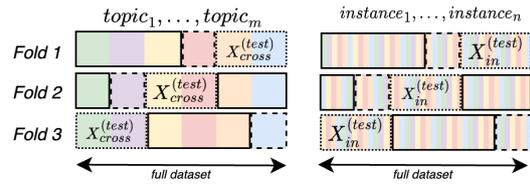


Figure 5: Overview of the In- and Cross-Topic setup using three folds. The colour indicates a topic; solid lines train-, dotted lines dev-, and dashed lines test-splits.

- [google/electra-base-discriminator](#) 1048 1049
- [gpt2](#) 1050
- [t5-3b](#) 1051
- [google/flan-t5-xl](#) 1052
- [EleutherAI/gpt-neo-2.7B](#) 1053

A.5 Fine-Tuning Hyperparameters 1054

To fine-tune on *stance detection*, we use the following setup: 5 epochs, where we find the best one using dev instances; AdamW (Loshchilov and Hutter, 2019) as optimizer; a batch size of 16; a learning rate of 0.00002; a warmup rate of 10% of the steps; random seeds: [0, 1, 2]. 1055 1056 1057 1058 1059 1060

A.6 Token-Level Examples for Topic Relevance 1061 1062

In § 5, we use the binned topic relevance (§ 5) for each token to represent its encoded topic information. We show in Table 6 examples for three bins *low*, *medium*, and *high*. The first bin (*low*) is made of tokens, which barely occur in the dataset. The second one (*medium*) consists of tokens which are part of most topics. Finally, the last bin (*high*) includes tokens with a high topic relevance for ones like *Cloning* or *Minimum Wage*. 1063 1064 1065 1066 1067 1068 1069 1070 1071

<i>low</i>	<i>medium</i>	<i>high</i>
fianc, joking, validate,	as, on, take,	cloning, uniform, wage,
latitude, poignantly, informative	some, like, how,	marijuana, minimum, gun,
ameliorate, bonding, mentors	so, one, these,	cloned, wear, clone,
brigade, emancipation, deriving,	instead, while, ago	nuclear, energy, penalty,
ignatius, 505, nominations,	where, came, still, many,	uranium, legalization, cannabis,
electorate, SWPS, 731	come, engage, seems	execution, wast, employment

Table 6: Examples of tokens with a *low*, *medium*, or *high* token relevance following § 4.

B Further Results

B.1 Generalization Across Datasets

With Table 7, Figure 7, and Table 8, we verify the results of § 4, § 5, and § 4 using another *stance detection* dataset. Namely, we use the *wttw* (*will-they-wont-they*) (Conforti et al., 2020) dataset which covers 51.284 tweets annotated either *support*, *refute*, *comment*, or *unrelated* towards five financial topics. For the overall performance comparison between In- and Cross-Topic, the results show the same trend as we already saw in § 4, but on a lower level. We assume that this is mainly due to this dataset’s more specific domain (twitter) compared to *UKP ArgMin*. Focusing on the influence of topic information verifies the previously presented results (§ 5) again. PLMs pre-trained with purely token-based objectives highly depend on topic information. Considering LLMs (Table 8), we see again similar behavior as on the *ArgMin* dataset (§ 4).

	DEP		POS		NER		Stance		Average		
	In	Cross	In	Cross	In	Cross	In	Cross	In	Cross	Δ
ALBERT	33.5	32.9	75.1	74.2	30.9	28.6	57.3	32.8	49.1	42.1	-7.0
BART	32.9	33.1	63.2	62.1	32.4	30.5	51.9	47.2	45.1	43.2	-1.9
BERT	21.6	21.2	54.8	55.9	27.2	27.8	47.4	32.1	37.8	34.2	-3.6
DeBERTa	26.9	27.6	69.6	67.9	29.4	28.5	49.5	35.7	43.9	40.0	-3.9
RoBERTa	20.4	19.9	54.7	53.5	26.1	25.5	37.0	37.8	35.6	34.2	-1.4
ELECTRA	26.6	26.6	69.6	68.6	21.7	24.1	35.1	36.7	38.2	39.0	+0.8
GPT-22	16.9	16.5	42.2	42.2	25.1	24.0	40.8	32.6	31.2	28.8	-2.4
GloVe	12.9	12.2	23.5	22.6	28.1	24.6	45.2	34.2	27.4	23.4	-4.0
Avg. Δ	-0.3		-0.7		-0.9		-9.5				

Table 7: Results of the four probing tasks using eight PLMs in the In- and Cross-Topic setup. We report the mean F_1 (macro averaged) over three random seeds, the average difference between the two evaluation setups per task (last row), and their average per PLM (last two columns). Best-performing results within a margin of 1pp are marked for every task and setup.

	DEP		POS		NER		Stance		Average		
	In	Cross	In	Cross	In	Cross	In	Cross	In	Cross	Δ
ALBERT	33.5	32.9	75.1	74.2	30.9	28.6	57.3	32.8	49.1	42.1	-7.0
BART	32.9	33.1	63.2	62.1	32.4	30.5	51.9	47.2	45.1	43.2	-1.9
T5 (3B)	25.5	26.3	59.7	59.3	34.9	36.4	53.4	38.7	43.4	40.2	-3.2
FLAN-T5 (3B)	25.5	26.3	59.7	59.3	34.9	36.4	53.4	38.7	43.4	40.2	-3.2
GPT-Neo (2.7B)	29.5	29.7	69.4	68.4	37.4	34.3	74.9	43.9	52.8	44.1	-8.7

Table 8: Results (macro F_1) of the four probing tasks using the overall best PLMs (ALBERT and BART) in the In- and Cross-Topic setup based on the *wttw* dataset (Table 7) and three LLMs.

B.2 Comparison of Probing Tasks against Random Initialized PLMs

We show in Table 9 and Table 10 the results of running the three linguistic probes on the seven contex-

tualized PLMs in their random initialized version. For In- and Cross-Topic, there is a clear performance drop of having random initialized models.

	DEP		POS		NER	
	Random	Δ	Random	Δ	Random	Δ
ALBERT	1.4	-42.4	6.8	-41.8	3.4	-76.8
BART	1.4	-35.1	5.0	-43.7	2.7	-72.7
BERT	2.7	-22.7	9.4	-36.0	4.6	-63.9
DeBERTa	7.0	-25.8	16.3	-32.5	16.1	-57.6
RoBERTa	2.2	-22.9	11.0	-37.4	4.7	-59.3
ELECTRA	1.7	-31.9	8.4	-33.1	3.8	-71.5
GPT-2	5.8	-19.4	12.3	-33.2	12.5	-51.0

Table 9: Results of evaluating *DEP*, *POS*, and *NER* using the seven contextual PLMs (random initialized) for In-Topic and the difference to their pre-trained counterparts in Table 2.

	DEP		POS		NER	
	Random	Δ	Random	Δ	Random	Δ
ALBERT	1.4	-38.1	6.2	-39.6	3.4	-74.6
BART	1.5	-35.4	5.0	-40.3	2.9	-71.2
BERT	2.1	-23.5	9.6	-32.0	4.5	-63.0
DeBERTa	6.8	-23.1	14.0	-28.4	17.2	-57.4
RoBERTa	2.6	-21.0	10.0	-32.1	5.2	-60.3
ELECTRA	3.0	-30.6	9.8	-31.4	4.1	-71.2
GPT-2	5.8	-18.1	13.6	-25.0	11.0	-50.9

Table 10: Results of evaluating *DEP*, *POS*, and *NER* using the seven contextual PLMs (random initialized) for Cross-Topic and the difference to their pre-trained counterparts in Table 2.

B.3 The Effect of Removing Random Information

We saw in § 5 that removing topic information has a big impact for some models (like RoBERTa or ELECTRA) but at the same time can even boost the performance of others like BERT. As suggested in Elazar et al. (2021), we apply a sanity check by removing random information from the encodings of PLMs. Following the results in Figure 8, removing random information (green bars) performs in between the scenarios with (blue bars) or without (red bars) topic information for cases where we see a clear negative effect when removing topic information. In contrast, removing random information can produce a more pronounced effect when we see performance improvements. This observation backs our assumption that removing information can have a regularization effect.

B.4 The Effect of Removing Topic Information on *Seen* and *Unseen* Instances

We show in Figure 6 that a performance drop affects *seen* and *unseen* instances for In- and Cross-Topic equally. Exceptionally, we see *unseen* ones are more affected on *POS* for DeBERTa and RoBERTa. This result indicates that these PLMs fall short of generalizing towards rare vocabularies - like *unseen* instances of *POS*.

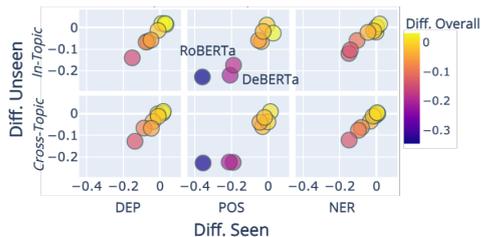


Figure 6: Performance difference for *seen* (x-axis) and *unseen* (y-axis) instances when removing topic information or not. One dot represents one PLM.

B.5 Analysis of Per-Class Results for NER

When considering the per-class results of *NER* in Table 11, we see the classes *CARDINAL*, *MONEY*, *ORG*, and *PERSON* show the biggest differences between In- and Cross-Topic. For *ORG* and *PERSON*, we see their topic-specific terms as the main reason for the performance gap. In contrast, we were surprised about the high difference for *CARDINAL*. We think this is mainly because this class embodies all numbers belonging to no other class. For *MONEY*, we see its uneven distribution over topics as the main reason for the performance difference - one topic covers more than 50% of the instances. These entities are highly topic-specific from a statistical point of view.

Despite having almost the same performance for In-Topic, BART and DeBERTa tend to outperform ALBERT on classes with more semantic complexities - like *GPE*, *ORG* or *PERSON*. For Cross-Topic, we see ALBERT performing better in classes unevenly distributed instances over topics - like *MONEY*. Further, it outperforms BART and DeBERTa on less semantical classes (*CARDINAL*, *ORDINAL*, *PERCENT*).

B.6 Effect of Removing Token-Level Topic Information of Per-Class Results for NER

Similar to the previous analysis, there are apparent effects of removing topic information when

		CARDINAL	DATE	GPE	MONEY	NORP	ORDINAL	ORG	PERCENT	PERSON
In	ALBERT	95.0	95.3	89.4	95.0	91.3	97.8	80.2	99.2	82.7
	BART	94.8	94.6	89.7	95.6	91.6	97.3	81.0	99.4	83.5
	DeBERTa	95.3	95.6	90.0	96.5	91.5	97.4	81.1	99.2	83.7
Cross	ALBERT	91.2	95.0	88.6	55.6	90.8	98.1	78.8	98.9	81.7
	BART	90.1	94.2	88.9	35.0	90.7	97.6	79.1	98.8	81.8
	DeBERTa	88.3	95.3	88.6	0.0	90.5	97.5	79.8	98.6	81.8

Table 11: Per-class results of ALBERT, BART, and DeBERTa on *NER* for In- and Cross-Topic.

considering *NER* classes separately. Table 12 shows these results for BART, BERT, DeBERTa, and RoBERTa. Like the overall result, BART, DeBERTa, and RoBERTa show lower performance when removing topic information. Whereby the effect is the most pronounced for RoBERTa with the highest performance drop for In- and Cross-Topic on classes like *NORP* or *ORDINAL*. In addition, these results show that the performance gain from removing topic information within BERT happens on *MONEY* for In-Topic and *NORP* for Cross-Topic.

		CARDINAL	DATE	GPE	MONEY	NORP	ORDINAL	ORG	PERCENT	PERSON
In	BART	-0.23	0.04	0.15	0.15	0.02	-0.04	0.08	-0.13	0.20
	BERT	1.65	-0.15	-0.04	28.00	-0.14	-0.58	0.06	0.00	0.22
	DeBERTa	-1.14	-0.13	-1.48	-7.74	-14.40	-0.30	-0.82	-0.12	-0.10
	RoBERTa	-6.00	-3.00	-7.82	-24.09	-90.61	-98.06	-2.66	-0.51	-0.58
Cross	BART	-0.48	0.01	-0.13	2.45	-0.06	-0.52	-0.38	-0.09	-0.03
	BERT	-0.05	-0.05	1.00	0.00	8.95	-0.60	0.29	0.00	0.00
	DeBERTa	-0.07	-0.16	-2.52	0.00	-21.88	-0.35	-0.91	-0.01	0.07
	RoBERTa	-9.04	-2.63	-7.45	0.00	-85.23	-98.07	-2.99	-35.97	-0.46

Table 12: Class-wise effect on the performance when removing topic information of BART, BERT, DeBERTa, and RoBERTa on *NER* for In- and Cross-Topic.

B.7 The Effect of Fine-Tuning on NER Classes

Analysing the results (Table B.7) for every *NER* class gives additional insights into where the fine-tuning had the most significant effect. We generally see the biggest effect on classes with less semantic meaning, like *ORDINAL*, *PERCENT*, or *MONEY*. At the same time, *GPE*, *PERSON*, and *ORG* are less affected as classes with more attached semantics. Regarding the different PLMs, ALBERT and DeBERTa show the most performance training, while BERT gains performance for the *MONEY* class.

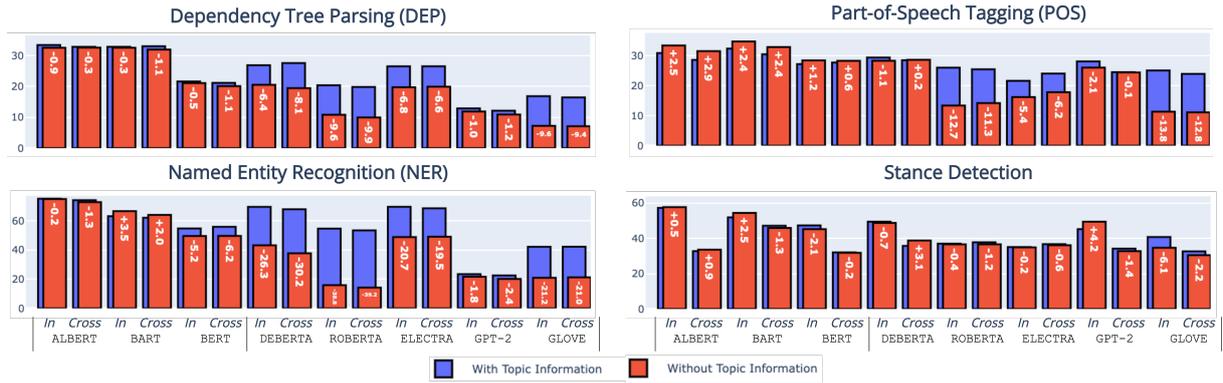


Figure 7: Comparison of the probing results with (blue bars) or without (red bars) topic information for the *will-they-wont-they* dataset (Conforti et al., 2020). The white text indicates the difference between these two scenarios.

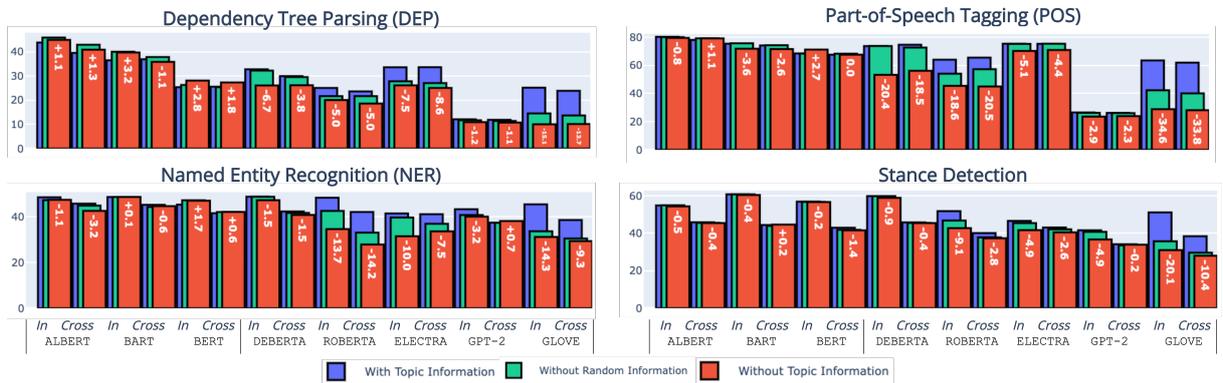


Figure 8: Comparison of the probing results with (blue bars) and without (red bars) topic information, or without random information (green bars). The white text indicates the difference between the blue and red bars.

	CARDINAL	DATE	GPE	MONEY	NORP	ORDINAL	ORG	PERCENT	PERSON	
<i>In</i>	ALBERT	-34.2	-25.4	-26.9	-95.0	-51.9	-60.3	-22.4	-99.2	-21.8
	BART	-8.5	-7.2	-7.5	-7.2	-10.4	-36.6	-4.1	-3.8	-2.7
	BERT	-1.9	-2.0	-2.0	34.8	-4.4	-17.9	-0.8	-3.9	-1.1
	DeBERTa	-15.1	-6.8	-8.7	-19.5	-43.7	-60.8	-8.8	-24.8	-8.3
<i>Cross</i>	ALBERT	-21.5	-10.4	-19.1	-55.6	-34.4	-13.1	-10.7	-81.0	-9.2
	BART	-9.2	-7.4	-7.0	-16.3	-11.2	-24.4	-3.9	-4.5	-2.1
	BERT	-2.5	-1.2	-1.2	3.6	-2.2	-9.7	-0.8	-2.6	-0.5
	DeBERTa	-18.2	-6.2	-12.7	0.0	-50.6	-76.0	-11.7	-73.5	-6.8

Table 13: Per-class difference before and after fine-tuning on *stance detection* of ALBERT, BART, BERT, and DeBERTa on NER for In- and Cross-Topic.

Task	Example	Label	# Instances	# Labels
DEP	I think there is a lot <u>we</u> can <u>learn</u> from Colorado and Washington State.	<i>nsubj</i>	40.000	41
POS	I think there is a lot <u>we</u> can learn from Colorado and Washington State.	<i>PRON</i>	40.000	17
NER	I think there is a lot we can learn from Colorado and <u>Washington State</u> .	<i>PERS</i>	25.892	17
Stance	<u>I think there is a lot we can learn from Colorado and Washington State.</u>	<i>PRO</i>	25.492	3

Table 14: Overview and examples of the different probing tasks.