

Investigating and Explaining Feature and Representation Learning in Translationese Classification

Anonymous ACL submission

Abstract

Recent work has shown that neural feature- and representation-learning approaches, and specifically the BERT model, demonstrates superior performance over traditional manual feature engineering and an SVM classifier for the task of translationese classification for various source and target languages. However, to date it is unclear whether the performance differences are due to better representations, better classifiers or both. Moreover, it remains unclear whether the features learnt by BERT overlap with commonly used manual features. To answer these, we exchange features between BERT-based and SVM classifiers, and show that, an SVM fed with BERT representations performs at the level of the best BERT classifiers, and BERT learning and using hand-crafted features performs at the level of traditional classifiers using hand-crafted features. Our experiments indicate that our hand-crafted feature set does not provide any additional information that BERT has not learnt already, and is likely to be a subset of features automatically learnt by BERT. Finally, we apply Integrated Gradients to examine token importance for the BERT model, and find that part of its top performance results are due to just topic differences and spurious correlations with translationese.

1 Introduction

Translationese is a descriptive (non-negative) cover term for the systematic differences between translated and originally authored text in same language (Gellerstam, 1986). Some aspects of translationese such as source interference (Toury, 1980; Teich, 2003) are language dependent, others are presumed universal, e.g. simplification, explicitation, overadherence to target language linguistic norms (Volansky et al., 2015) in the products of translations. While translationese effects can be subtle, and even human experts may not be able to reliably distinguish between original texts and professional translations (Tirkkonen-Condit, 2002),

corpus-based studies (Baker et al., 1993) and, in particular, machine-learning classifier based studies (Rabinovich and Wintner, 2015; Volansky et al., 2015; Rubino et al., 2016; Pylypenko et al., 2021) clearly reveal the differences. In this paper we focus on machine-learning classifier based research on translationese. Here, typically a classifier is trained to distinguish between original and translated texts (in the same language). Until recently, most of this research (Baroni and Bernardini, 2005; Volansky et al., 2015; Rubino et al., 2016) used manually defined, often linguistically inspired, feature-engineering based sets of features, (mostly) using support vector machines (SVM). Once a classifier is trained, feature importance and ranking methods are used to reason back to what aspects of the input is responsible for (i.e. explains) the classification. More recently, a small number of papers explored feature- and representation-learning neural network based approaches to translationese classification (Sominsky and Wintner, 2019). In a systematic study Pylypenko et al. (2021) shows that feature- and representation-learning deep neural network-based approaches (in particular BERT-based, but also other neural approaches) to translationese classification substantially outperform handcrafted feature-engineering based approaches using SVMs. However, to date, two important questions remain: (i) it is not clear whether the substantial performance differences are due to learned vs. handcrafted features, the classifiers, or the combination of both, and (ii) what the neural feature and representation learning approaches actually learn. The contributions of our paper are as follows:

1. we address (i) by feeding BERT-based learned features to SVMs and by letting BERT architectures learn handcrafted features, as well as feeding the handcrafted features into BERT as embeddings. Our experiments show that SVMs using BERT-learned features perform on a par with our best BERT-

084	translationese classifiers. Moreover, BERT	We also aim to verify if there is an overlap between	133
085	using handcrafted features only performs at	the features learnt by BERT and our handcrafted	134
086	most as good as the SVM classifier.	feature set. In some cases the classifier used for	135
087	2. we present the first steps to address (ii)	ensembling BERT and handcrafted features is a	136
088	using attribution-based explainable AI ap-	Support Vector Machine (Kazameini et al., 2020;	137
089	proaches (XAI) on our best performing full	Ray and Garain, 2020).	138
090	feature- and representation-learning BERT	Explainability methods for neural networks have	139
091	model and on BERT models that are pre-	not been widely explored for translationese clas-	140
092	trained to predict handcrafted features and	sification. Since many previous works have used	141
093	then fine-tuned for translationese classifica-	the traditional feature-engineering method, they of-	142
094	tion. We present evidence that at least part	ten quantify handcrafted feature importance. Tech-	143
095	of the high classification accuracy of BERT	niques used for that include looking at SVM feature	144
096	is due to names of places and countries, sug-	weights (Avner et al., 2016; Pylypenko et al., 2021),	145
097	gesting that part of the classification is topic-	correlation (Rubino et al., 2016), Information Gain	146
098	and not translationese-based (source texts in	(Ilisei et al., 2010), Chi-square (Ilisei et al., 2010),	147
099	Spanish translated to English e.g. may have	decision trees or random forests (Rubino et al.,	148
100	a higher likelihood of talking about Spanish	2016; Ilisei et al., 2010), ablating features and ob-	149
101	places). Moreover, some top features suggest	servating the change in accuracy (Baroni and Bernar-	150
102	that there might be certain spurious correla-	dini, 2005; Ilisei et al., 2010), training separate	151
103	tions within our dataset.	classifiers on each individual feature (or feature set)	152
104	To the best of our knowledge this is the first	and comparing accuracies (Volansky et al., 2015;	153
105	paper that shows that feature- and representation-	Avner et al., 2016). For n-grams, then difference in	154
106	learning rather than the classifier is responsible	frequencies between the original and translationese	155
107	for the substantial performance gap between deep	classes (Koppel and Ordan, 2011; van Halteren,	156
108	neural networks and machine learning approaches	2008), and the contribution to the symmetrized	157
109	using handcrafted features. It is also the first paper	Kullback-Leibler Divergence between the classes	158
110	that uses XAI methods to (begin to) explain what	(Kurokawa et al., 2009) have been used.	159
111	neural methods learn in translationese classifica-	As for looking into the neural network perfor-	160
112	tion.	mance, Pylypenko et al. (2021) quantify whether	161
113	Finally, translationese research is not just an	hand-crafted features can explain the variance in	162
114	"academic" exercise in basic research into aspects	the predictions of neural models, such as BERT,	163
115	of how translation works, but an important research	LSTM, and Simplified Transformer, by training	164
116	topic in machine translation evaluation (Stymne,	per-feature linear regression models to output the	165
117	2017; Toral et al., 2018; Freitag et al., 2019;	predicted probabilities of the neural models and	166
118	Edunov et al., 2020; Graham et al., 2020) and in	computing the R^2 measure. They find that most of	167
119	further improving machine translation (Riley et al.,	the top features are either POS-perplexity-based, or	168
120	2020).	bag-of-POS features. However, this method treats	169
121	2 Related Work	the neural network as a black-box, whereas we use	170
122	(Kaas et al., 2020; Prakash and Tayyar Madabushi,	a method that accesses the internals of the model.	171
123	2020; Lim and Madabushi, 2020) combine BERT-	In our work we use the Integrated Gradients	172
124	based and handcrafted features in an ensemble man-	method (Sundararajan et al., 2017) the method pro-	173
125	ner in order to improve over BERT's accuracy, of-	vides attribution scores for the input with respect to	174
126	ten by concatenating the pooled output of BERT	a certain class. It involves calculating the integral	175
127	with a handcrafted feature vector (sometimes addi-	of gradients with respect to the input along the path	176
128	tionally encoded by another network) and feeding	from a certain baseline (in our case, PAD tokens)	177
129	them into another classifier. They show that even	to the input.	178
130	though BERT representations are powerful, care-	3 Experimental Settings	179
131	fully picked handcrafted features may still provide	3.1 Data	180
132	additional information that aids the task in hand.	For our experiments, we use the monolingual Ger-	181
		man dataset in the Multilingual Parallel Direct Eu-	182

183 `roparl` (MPDE) (Amponsah-Kaakyire et al., 2021) 217
 184 corpus. The set contains 42k paragraphs with half 218
 185 of the texts German originals and the other half 219
 186 translations into German from Spanish. The aver- 220
 187 age length is 80 tokens per training sample. Since 221
 188 there exists a problem with pivot translations in 222
 189 Europarl (Bogaert, 2011), the DE-ES dataset con- 223
 190 tains only data from before 2004, when the pivot 224
 191 system was introduced. We additionally use a held- 225
 192 out corpus of around 30k paragraphs for estimating 226
 193 language models and n -gram quartile distributions. 227
 194 This corpus consists of originally produced texts 228
 195 only. For the heldout corpus, we sample texts from 229
 196 Europarl proceedings from 2004 onwards, since 230
 197 original data is not affected by the pivot translation 231
 198 problem. 232

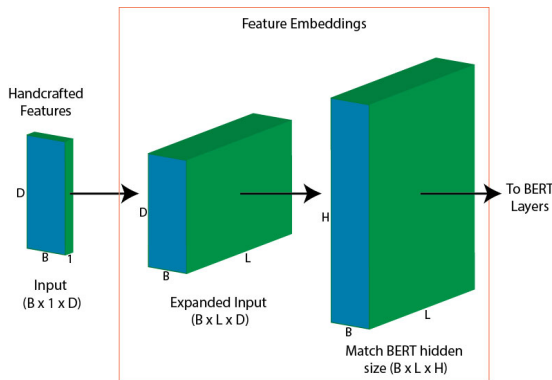


Figure 1: Mapping handcrafted features to embeddings (Section 3.4).

3.2 Base Setup

199 A fair comparison between learned and handcrafted 200
 201 features and two classifiers is non-trivial. We feed 202
 203 learned features into SVMs, and (i) let BERT learn 204
 205 handcrafted feature vectors used by SVMs and (ii) 206
 207 feed handcrafted feature vectors as embeddings 208
 209 into the BERT model. We compare this with full 210
 211 feature and representation learning settings with 212
 213 BERT, and SVMs with handcrafted features. To 214
 215 test this for translationese classification we repro- 216
 217 duce the models from Pylypenko et al. (2021):

1. a linear SVM with 108 handcrafted features (with surface, lexical, unigram bag-of-PoS, language modelling and n -gram frequency distribution features), [**handc.-features+SVM**]
2. a pretrained BERT-base model (12 layers, 768 hidden dimensions, 12 attention heads)

fine-tuned on translationese classification. 217
 [**pretrained-BERT-ft**] 218

We use multilingual BERT (Devlin et al., 2019) (BERT-base-multilingual-uncased), and fine-tuning is done with the *simpletransformers*¹ library. We use a batch size of 32, learning rate of $4 \cdot 10^{-5}$, and the Adam optimiser with epsilon $1 \cdot 10^{-8}$. We estimate n -gram language models with SRILM (Stolcke, 2002) and do POS-tagging with SpaCy.²

3.3 SVM Classifier with BERT Features

We train a SVM with linear kernel on the features learnt by the pretrained BERT model fine-tuned on the translationese classification task. We use the output of the BERT pooler, which selects the last layer [CLS] token vector, with linear projection and *tanh* activation as our feature vector. We use:

1. BERT's 768-dim pooled vector output, [**pretrained-BERT-ft-output+SVM**]
2. a 108-dim PCA projection of the vector. [**pretrained-BERT-ft-output+PCA(108 dim)+SVM**]

The later allows us to match the handcrafted feature vector dimensionality. 238
 239

3.4 BERT with Handcrafted Features

Making neural networks use handcrafted features in our comparison is non-trivial. We design two strategies. 240
 241
 242
 243

Pretraining on handcrafted feature prediction.

244 First, we train a BERT-base model from scratch, 245
 246 using it to predict the 108 dimension vectors rep- 247
 248 resenting handcrafted features originally used in 249
 250 training the SVM [**BERT-reg-full**]. The weights 251
 252 of this model encode the information of the hand- 253
 254 crafted features. With this pretrained model,

1. we freeze the weights and train a classifier on top for translationese classification, [**BERT-r2c-full-frozen**]
2. we do not freeze but fine-tune on the translationese classification task. [**BERT-r2c-full-ft**]

256 We explore the same with a smaller BERT model 257
 258 with only 6 layers instead of 12 [**BERT-reg-half**]. 259
 260 We then load the weights of the small 6 layer model into the embedding layer and the first 6 layers of a 12 layer non-pretrained BERT-base model and:

¹github.com/ThilinaRajapakse/simpletransformers
²<https://spacy.io/>

Model	Test accuracy (%)
handcr.-features + SVM	73.2±0.1
pretrained-BERT-ft-output + PCA(108 dim) + SVM	92.0±0.0
pretrained-BERT-ft-output + SVM	92.0±0.0
BERT-r2c-full-frozen-output + PCA(108 dim) + SVM	70.3±0.1
BERT-r2c-full-frozen-output + SVM	74.9±0.7
pretrained-BERT-ft	92.2±0.2
fromScratch-BERT	89.3±0.3
BERT-r2c-full-frozen	59.6±0.1
BERT-r2c-full-ft	89.3±0.4
BERT-r2c-half-frozen	67.5±0.4
BERT-r2c-half-ft	89.0±0.3
BERT-f2c $L = 1$	57.1±10.1
BERT-f2c $L = 80$	72.8±0.2
BERT-f2c $L = 256$	72.7±0.2
pretrained-BERT-f2c $L = 80$	68.0±2.1

Table 1: Translationese classification accuracy for all settings (average and standard deviation over 5 runs). All of the models were trained/fine-tuned for the translationese classification task.

3. we freeze the loaded weights in the first 6 layers and train the remaining 6 layers and classifier on the translationese classification task, [**BERT-r2c-half-frozen**]
4. we do not freeze but fine-tune it on the translationese classification task with randomly-initialised weights for the other 6. [**BERT-r2c-half-ft**]

Interestingly, according to the losses when training for predicting the handcrafted features, BERT-reg-half performs comparably to the BERT-reg-full (0.0041136 vs 0.0041148).

We also train a BERT-base model with the same settings from scratch on the translationese classification task as a baseline for all BERT models. [**fromScratch-BERT**].

Mapping handcrafted features to embeddings.

Even though the very low MSE results indicate that both versions of BERT-reg are able to learn handcrafted features well, using them in terms of frozen layers in translationese classification leads to low classification performance. This motivates us to explore an alternative way of encoding handcrafted features: we convert the single vector of handcrafted features of dimension D (108 in our experiments) into a sequence of embeddings in BERT’s layers format, that is, length of feature embedding sequence L times the dimension of the hidden states H (768), while preserving the information of the single vector. (Fig. 1)

To do this, we consider a batch of tokens with size B and take in the handcrafted features as a $B \times D$ -dimensional input to the BERT model and generate feature embeddings by passing the features through 2 linear layers as follows. We first unsqueeze the $B \times D$ input to $B \times 1 \times D$ dimensions and reshape it as $B \times D \times 1$. This is passed to the first linear layer. The resulting $B \times D \times L$ -dimensional output is reshaped as $B \times L \times D$ and fed as input to the second linear layer which outputs a $B \times L \times H$ -dimensional output as the feature embeddings.

This hand-crafted feature embedding replaces BERT’s embedding layer and serves as input to the first BERT layer. The resulting BERT model is trained on the translationese classification task. We experiment with three different values for L : 1, 80 (average length of our training samples) and 256 (maximum input for BERT). [**BERT-f2c L=1, BERT-f2c L=80, BERT-f2c L=256**]

4 Translationese Classification

Table 1 summarises results of the different translationese classification settings. As for feeding pooled output of BERT into the SVM model, we can observe that the accuracy is a lot higher comparing to feeding handcrafted features, even when the BERT vector dimensionality is reduced to match the amount of handcrafted features. This emphasizes the fact that the features learnt by BERT are superior to our current set of manual features.

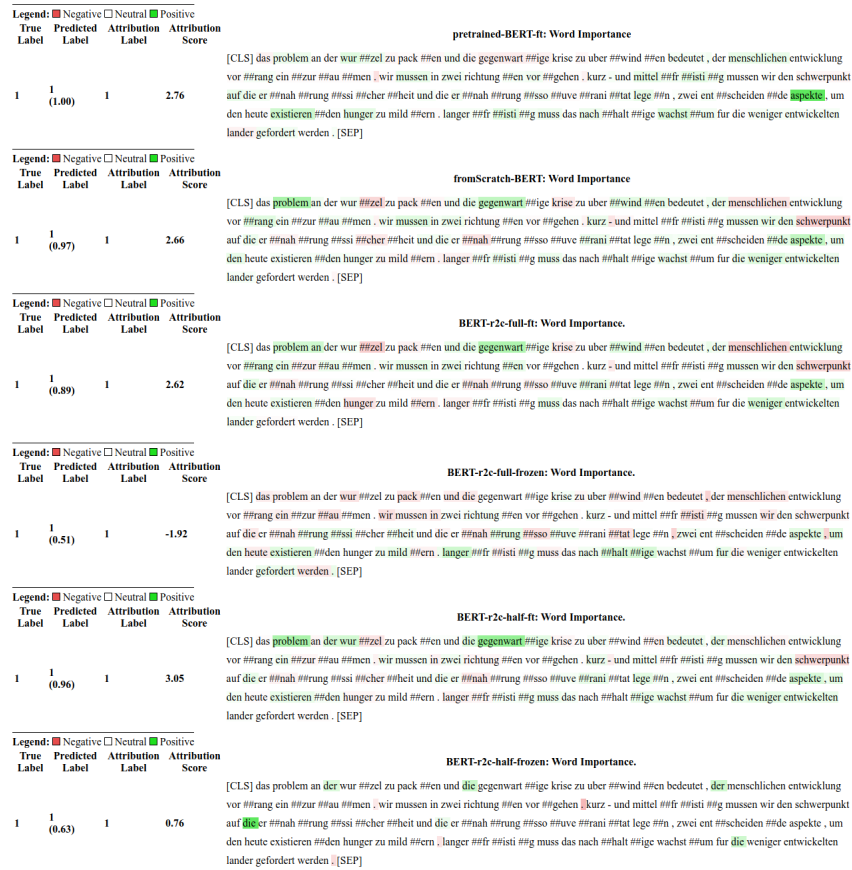


Figure 2: Layer Integrated Gradient saliency maps of input tokens contributing to the ground truth translationese label (here: translation). Comparison of different models.

When BERT is trained from scratch, translationese classification accuracy reduces by 2.93 percentage points, compared to the pretrained-BERT-ft. This suggests that pretraining on large data helps to encode additional information that turns out to be helpful in the translationese classification task.

As for BERT pretrained to predict the handcrafted features and frozen (BERT-r2c-full-frozen), it is assumed that BERT has learnt to encode the handcrafted features during pretraining. Nevertheless, its accuracy, albeit higher than a random guess, is lower by ~ 13 percentage points comparing to the SVM classifier. We perform an additional experiment, in order to check whether the difference in accuracy is due to BERT failing to sufficiently encode the handcrafted features during pretraining, or due to the SVM classifier being superior to the linear classification head of the BERT model. Namely, we train the SVM classifier on the pooled output of BERT-r2c-full-frozen model. The accuracy is around 75% (for both original and PCA-reduced dimensionality) which is as high as using SVM on handcrafted feature vectors. We conclude that

BERT encodes the handcrafted features sufficiently enough, but the linear classifier performs worse than an SVM in these conditions, possibly due to non-exhaustive hyperparameter search.

Further fine-tuning BERT, fully pretrained for handcrafted feature prediction (BERT-r2c-full-ft), for translationese classification results in accuracy comparable to BERT that was not pretrained on this task (fromScratch-BERT). This could suggest that our handcrafted feature set is either a subset of features learned by fromScratch-BERT, or that the handcrafted features are discarded during fine-tuning. The model where only the first 6 layers were pretrained (BERT-r2c-half-ft), achieves similar accuracy, likely due to the same reasons.

By contrast, freezing the 6 handcrafted feature prediction pretrained layers (BERT-r2c-half-frozen) largely reduces the accuracy, because the model only has access to the 6th layer embeddings that supposedly encode only the information about the handcrafted features, and does not have ability to extract its own features from the input text, due to its inability to tune the embeddings.

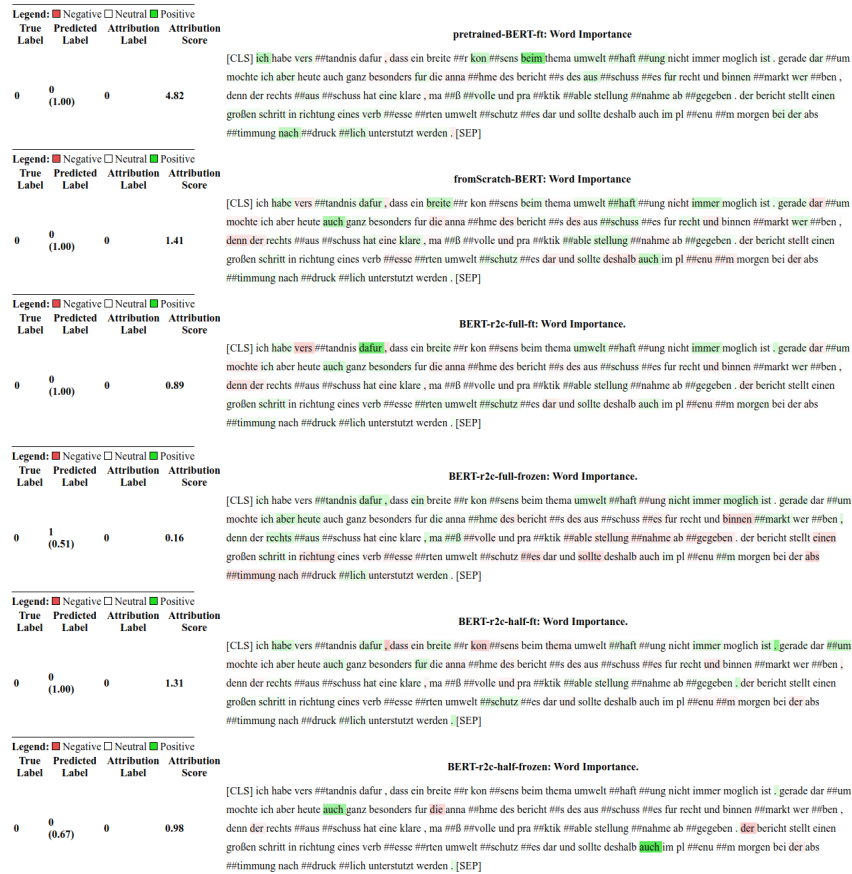


Figure 3: Layer Integrated Gradient saliency maps of input tokens contributing to the ground truth translationese label (here: original). Comparison of different models.

The results of BERT-f2c models show that BERT, when fed the handcrafted features in the form of embeddings, can reach at most the same accuracy as the handcrd.-features+SVM approach, which suggests that the BERT architecture has no advantage over the SVM classifier in utilizing the handcrafted features for classification.

5 Layer Integrated Gradients Saliency

We compare input attributions of the ground truth classification label amongst pretrained-BERT-ft, fromScratch-BERT and the different settings of the translationese classification models pretrained on the hand-crafted feature prediction task. We use Layer Integrated Gradients from the Captum library (Kokhlikyan et al., 2020), and calculate the salience score for each token by averaging the attributions over the embedding dimension.

5.1 Comparing Models

Figure 2 displays Integrated Gradient attributions for a **translated** paragraph across different BERT models. Figure 3 shows attributions for an **original**

paragraph.

Comparing the attributions of classification labels to sample inputs amongst the various settings of BERT, we observe that attributions are similar for fromScratch-BERT and the fine-tuned models: BERT-r2c-full-ft and BERT-r2c-half-ft. By contrast, freezing the weights in BERT-r2c-full-frozen and BERT-r2c-half-frozen resulted in very different attributions from the fromScratch-BERT. For BERT-r2c-half-frozen the attributions are more peaked than for other models, with only a few tokens receiving large scores, and most tokens having scores close to zero. Notably, pretrained-BERT-ft displays a pattern that is overall similar to the BERT trained from scratch, but some attributions are reversed, and the peaks are on different tokens.

For BERT-r2c-full-frozen, it can be seen that a substantial number of tokens with negative attributions have positive attributions in the model trained from scratch and also the fine-tuned models. However some attributions overlap, which suggests that BERT might be using some of the handcrafted features.

BERT-r2c-full-ft			pretrained-BERT-ft		
Rank	Token	Avg attribution score	Rank	Token	Avg attribution score
1	sagte	0.60	1	entstand	0.70
2	gebiet	0.46	2	virus	0.63
3	##dies	0.44	3	inti	0.60
4	ansicht	0.43	4	sagte	0.58
5	bezug	0.42	5	entdeckte	0.57
6	neige	0.40	6	gras	0.57
7	amt	0.40	7	nuts	0.56
8	pre	0.40	8	nicaragua	0.55
9	spanien	0.39	9	rekord	0.53
10	sprechen	0.38	10	bilbao	0.53
11	nuts	0.36	11	verfugte	0.53
12	barcelona	0.34	12	bol	0.51
13	;	0.33	13	colombia	0.51
14	##bien	0.32	14	nis	0.51
15	spanischen	0.32	15	och	0.49
16	wiederholt	0.31	16	vorkommen	0.49
17	einige	0.30	17	oecd	0.49
18	##sprache	0.29	18	;	0.46
19	weder	0.29	19	erklarte	0.45
20	territorium	0.28	20	clinton	0.45

Table 2: Top-20 tokens with highest average attribution score towards the **translationese** class in the test set. BERT-r2c-full-ft and pretrained-BERT.

5.2 Comparing Checkpoints

In Appendix A we provide additional results on examining training checkpoints for fromScratch-BERT and BERT-r2c-full-ft for an original and a translated paragraph.

Results indicate that for fromScratch-BERT some attributions change into opposite during training, whereas for BERT-r2c-full-ft the pattern appears to be already settled from the early checkpoints onwards, and does not change much over the course of fine-tuning. This may support the hypothesis that the handcrafted features are a subset of features learnt by BERT, and thus provide a favorable initialization of weights for fine-tuning for translationese classification.

5.3 Highest Average Attribution

In order to make the interpretation less local, we compute the top tokens with highest attribution on average across the test set. The results for each class for best-performing models (pretrained BERT and BERT-r2c-full-ft) are given in Tables 2 and 3.

For German translationese data translated from Spanish, some top tokens correspond to the geographical areas, where Spanish is spoken, e.g. "spanien", "barcelona", "spanischen" for BERT-r2c-full-ft; "nicaragua", "colombia", "bilbao" for

pretrained BERT. Likewise for original German data, some of the top tokens are German geographical names, e.g. "stuttgart" for pretrained BERT. The subword "##wald" also appears to be a common German toponymic suffix. This suggests that topic is one of the spurious clues that is used by BERT to determine the correct translationese class. This is also supported by the fact that some nouns that likely correspond to certain recurring discussion topics for only one class within our data sample, receive high attribution, e.g. "virus", "soja", "clinton", "orange" etc. The "ez" token, salient for the original class, appears to be a starting subword unit of the *EZB* abbreviation (Europäische Zentralbank).

The "•" token having a high attribution for the class *originals* for both models might suggest a spurious correlation within the dataset, that is apparently utilized by BERT. The ";" token is deemed important for the translationese class by both models, which might also be a spurious correlation. Conversely, this could be an indication that clauses in Spanish are more often joint with the semi-colon, than in German, which was preserved in the translation.

For both models the Präteritum forms "sagte", "erklärte" etc. are also among the top tokens impor-

BERT-r2c-full-ft			pretrained-BERT-ft		
Rank	Token	Avg attribution score	Rank	Token	Avg attribution score
1	##wegen	0.61	1	situations	0.37
2	•	0.55	2	•	0.36
3	eu	0.49	3	ria	0.34
4	daraufhin	0.49	4	##lk	0.33
5	finde	0.45	5	##iet	0.32
6	##vo	0.45	6	golden	0.32
7	gerne	0.43	7	sak	0.30
8	##abb	0.42	8	turm	0.30
9	##hrte	0.42	9	##emen	0.27
10	ausbau	0.42	10	orange	0.27
11	!	0.42	11	hang	0.26
12	bekommen	0.42	12	##wald	0.25
13	trips	0.41	13	1732	0.25
14	ez	0.41	14	dobe	0.24
15	##gemeinde	0.40	15	##pas	0.23
16	vot	0.36	16	profits	0.22
17	won	0.36	17	stuttgart	0.22
18	geplant	0.35	18	soja	0.21
19	demnach	0.35	19	r	0.21
20	ja	0.35	20	ruth	0.21

Table 3: Top-20 tokens with highest average attribution score towards the **original** class in the test set. BERT-r2c-full-ft and pretrained-BERT.

464 tant for recognizing translationese. One possible
465 explanation could be that the Perfekt form ("hat
466 gesagt") is more common in German spoken lan-
467 guage, and Präteritum is more common in writing.
468 Therefore the translators, while translating Spanish
469 speeches into German, could have preferred to use
470 the Präteritum form more common for writing.

471 6 Conclusion

472 This paper addresses two open questions in
473 classification-based translationese research: (1) are
474 the substantial performance differences between
475 feature- and representation-learning and classical
476 handcrafted feature based approaches due to (i) the
477 difference in the features, (ii) the classifiers, or (iii)
478 both, and (2) what do feature- and representation-
479 learning based approaches actually learn?

480 We address (1) by exchanging features from both
481 models examining a broad variety of settings. We
482 confirm that SVMs perform as good as BERT when
483 fed with features learnt by BERT. Likewise, BERT
484 performs at the level of traditional SVM-based clas-
485 sification with handcrafted features SVMs, when
486 fed with handcrafted features only. Our findings re-
487 veal that while pretraining on huge amount of data
488 improves the classification accuracy, pretraining
489 on handcrafted features does not guarantee an im-

490 provement on classification accuracy with respect
491 to training from scratch.

492 To address question (2), we examine BERT's in-
493 put attributions using Integrated Gradients Saliency
494 for various settings and observe that attributions are
495 indeed similar for the model trained from scratch
496 and the fine-tuned models that were pretrained on
497 handcrafted feature prediction.

498 Finally, analysis of top activated tokens in the
499 test set suggests that at least part of BERT's strong
500 translationese classification accuracy is based on
501 topical differences between the classes (rather than
502 "proper" translationese phenomena), the topical
503 differences between the classes, and spurious cor-
504 relations. The next step would be to control these
505 factors, for instance by using named entity masking
506 and cleaning/normalizing the corpus, in order to
507 investigate whether BERT would still outperform
508 the traditional approach under such conditions.

509 References

510 Kwabena Amponsah-Kaakyire, Daria Pylypenko,
511 Cristina España-Bonet, and Josef van Genabith. 2021.
512 Do not rely on relay translations: Multilingual par-
513 allel direct Europarl. In *Proceedings for the First
514 Workshop on Modelling Translation: Translatology*

515			
516		<i>in the Digital Age</i> , pages 1–7, online. Association for Computational Linguistics.	
517	Ehud Alexander Avner, Noam Ordan, and Shuly Wintner. 2016. Identifying translationese at the word and sub-word level. <i>Digit. Scholarsh. Humanit.</i> , 31:30–54.		
521	Mona Baker, Gill Francis, and Elena Tognini-Bonelli. 1993. Corpus linguistics and translation studies: Implications and applications. In <i>Text and Technology: In Honour of John Sinclair</i> , page 233–, Netherlands. John Benjamins Publishing Company.		
522			
523			
524			
525			
526	Marco Baroni and Silvia Bernardini. 2005. A New Approach to the Study of Translationese: Machine-learning the Difference between Original and Translated Text . <i>Literary and Linguistic Computing</i> , 21(3):259–274.		
527			
528			
529			
530			
531	Caroline Bogaert. 2011. Is absolute multilingualism maintainable? The language policy of the European Parliament and the threat of English as a lingua franca . Master’s thesis, UGent. Faculteit Letteren en Wijsbegeerte.		
532			
533			
534			
535			
536	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.		
537			
538			
539			
540			
541			
542			
543			
544			
545	Sergey Edunov, Myle Ott, Marc’Aurelio Ranzato, and Michael Auli. 2020. On the evaluation of machine translation systems trained with back-translation . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2836–2846, Online. Association for Computational Linguistics.		
546			
547			
548			
549			
550			
551			
552	Markus Freitag, Isaac Caswell, and Scott Roy. 2019. APE at scale and its implications on MT evaluation biases . In <i>Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)</i> , pages 34–44, Florence, Italy. Association for Computational Linguistics.		
553			
554			
555			
556			
557			
558	Martin Gellerstam. 1986. Translationese in Swedish novels translated from English. <i>Translation studies in Scandinavia</i> , 1:88–95.		
559			
560			
561	Yvette Graham, Barry Haddow, and Philipp Koehn. 2020. Statistical power and translationese in machine translation evaluation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 72–81, Online. Association for Computational Linguistics.		
562			
563			
564			
565			
566			
567	Iustina Ilisei, Diana Inkpen, Gloria Corpas Pastor, and Ruslan Mitkov. 2010. Identification of translationese: A machine learning approach. In <i>Computational Linguistics and Intelligent Text Processing</i> , pages 503–511, Berlin, Heidelberg. Springer Berlin Heidelberg.		
568			
569			
570			
571			
	Anders Kaas, Viktor Torp Thomsen, and Barbara Plank. 2020. Team DiSaster at SemEval-2020 task 11: Combining BERT and hand-crafted features for identifying propaganda techniques in news . In <i>Proceedings of the Fourteenth Workshop on Semantic Evaluation</i> , pages 1817–1822, Barcelona (online). International Committee for Computational Linguistics.		572 573 574 575 576 577 578
	Amirmohammad Kazameini, Samin Fatehi, Yash Mehta, Sauleh Eetemadi, and Erik Cambria. 2020. Personality trait detection using bagged SVM over BERT word embedding ensembles . <i>CoRR</i> , abs/2010.01309.		579 580 581 582
	Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. 2020. Captum: A unified and generic model interpretability library for pytorch .		583 584 585 586 587 588
	Moshe Koppel and Noam Ordan. 2011. Translationese and its dialects . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1318–1326, Portland, Oregon, USA. Association for Computational Linguistics.		589 590 591 592 593 594
	David Kurokawa, Cyril Goutte, and P. Isabelle. 2009. Automatic detection of translated text and its impact on machine translation. In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics</i> , Baltimore, Maryland.		595 596 597 598 599
	Wah Meng Lim and Harish Tayyar Madabushi. 2020. Uob at semeval-2020 task 12: Boosting BERT with corpus level information . <i>CoRR</i> , abs/2008.08547.		600 601 602
	Anushka Prakash and Harish Tayyar Madabushi. 2020. Incorporating count-based features into pre-trained models for improved stance detection . In <i>Proceedings of the 3rd NLP4IF Workshop on NLP for Internet Freedom: Censorship, Disinformation, and Propaganda</i> , pages 22–32, Barcelona, Spain (Online). International Committee on Computational Linguistics (ICCL).		603 604 605 606 607 608 609 610
	Daria Pylypenko, Kwabena Amponsah-Kaakyire, Koel Dutta Chowdhury, Josef van Genabith, and Cristina España-Bonet. 2021. Comparing feature-engineering and feature-learning approaches for multilingual translationese classification . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 8596–8611, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		611 612 613 614 615 616 617 618 619
	Ella Rabinovich and Shuly Wintner. 2015. Unsupervised identification of translationese . <i>Transactions of the Association for Computational Linguistics</i> , 3:419–432.		620 621 622 623
	Biswarup Ray and Avishek Garain. 2020. Factuality classification using bert embeddings and support vector machines. In <i>IberLEF@SEPLN</i> .		624 625 626

- 627 Parker Riley, Isaac Caswell, Markus Freitag, and David
628 Grangier. 2020. [Translationese as a language in “mul-](#)
629 [tilingual” NMT](#). In *Proceedings of the 58th Annual*
630 *Meeting of the Association for Computational Lin-*
631 *guistics*, pages 7737–7746, Online. Association for
632 Computational Linguistics.
- 633 Raphael Rubino, Ekaterina Lapshinova-Koltunski, and
634 Josef van Genabith. 2016. [Information density and](#)
635 [quality estimation features as translationese indica-](#)
636 [tors for human translation classification](#). In *Proceed-*
637 *ings of the 2016 Conference of the North American*
638 *Chapter of the Association for Computational Lin-*
639 *guistics: Human Language Technologies*, pages 960–
640 970, San Diego, California. Association for Compu-
641 tational Linguistics.
- 642 Iliia Sominsky and Shuly Wintner. 2019. [Automatic](#)
643 [detection of translation direction](#). In *Proceedings of*
644 *the International Conference on Recent Advances in*
645 *Natural Language Processing (RANLP 2019)*, pages
646 1131–1140, Varna, Bulgaria. INCOMA Ltd.
- 647 Andreas Stolcke. 2002. SRILM – An extensible lan-
648 guage modeling toolkit. In *Proceedings of the 7th*
649 *International Conference on Spoken Language Pro-*
650 *cessing (ICSLP 2002)*, pages 901–904.
- 651 Sara Stymne. 2017. The effect of translationese on
652 tuning for statistical machine translation. In *The*
653 *21st Nordic Conference on Computational Linguis-*
654 *tics*, pages 241–246.
- 655 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017.
656 Axiomatic attribution for deep networks. In *Proceed-*
657 *ings of the 34th International Conference on Machine*
658 *Learning - Volume 70, ICML’17*, page 3319–3328.
659 JMLR.org.
- 660 Elke Teich. 2003. *Cross-Linguistic Variation in Sys-*
661 *tem und Text. A Methodology for the Investigation*
662 *of Translations and Comparable Texts*. Mouton de
663 Gruyter, Berlin.
- 664 Sonja Tirkkonen-Condit. 2002. Translationese – a myth
665 or an empirical fact?: A study into the linguistic
666 identifiability of translated language. *Target. Interna-*
667 *tional Journal of Translation Studies*, 14(2):207–220.
- 668 Antonio Toral, Sheila Castilho, Ke Hu, and Andy Way.
669 2018. [Attaining the Unattainable? Reassessing](#)
670 [Claims of Human Parity in Neural Machine Transla-](#)
671 [tion](#). In *Proceedings of the Third Conference on Ma-*
672 *chine Translation: Research Papers*, pages 113–123,
673 Brussels, Belgium. Association for Computational
674 Linguistics.
- 675 Gideon Toury. 1980. *In Search of a Theory of Transla-*
676 *tion*. The Porter Institute for Poetics and Semiotics,
677 Tel Aviv University, Tel Aviv.
- 678 Hans van Halteren. 2008. [Source language markers in](#)
679 [EUROPARL translations](#). In *Proceedings of the 22nd*
680 *International Conference on Computational Linguis-*
681 *tics (Coling 2008)*, pages 937–944, Manchester, UK.
682 Coling 2008 Organizing Committee.
- Vered Volansky, Noam Ordan, and Shuly Wintner. 2015.
On the features of translationese. *Digital Scholarship*
in the Humanities, 30(1):98–118.

A Appendix

A.1 Extra Information on MPDE dataset

We use version 2.0.0 of the [MPDE dataset](#) licensed under CC-BY 4.0. Specifically we use the *mono_de_es* train/dev/test splits of the German-Spanish language pair. Table 4 contains summary statistics of the data.

Split	Number of Examples
Train set	29580
Validation set	6366
Test	6344

Table 4: Dataset statistics

A.2 Extra Information on BERT models

With the exception of pretrained-BERT-ft, we use the *transformers* library.³ Training is done across 4 NVIDIA GeForce GTX TITAN X GPUs with a batch size of 8 per GPU. We use a learning rate of $3 \cdot 10^{-5}$ and train or fine-tune for 5 epochs. Table 5 shows the number of parameters of the different BERT variants. Parameter counts include the embedding and respective prediction (classifier or regression) layers.

Model	Num. Params (M)
fromScratch-BERT	177.85
BERT-reg-full	177.94
BERT-reg-half	135.41
BERT-r2c-*	177.85
BERT-f2c $L = 1$	177.46
BERT-f2c $L = 80$	177.52
BERT-f2c $L = 256$	177.66
pretrained-BERT-f2c $L = 80$	177.52

Table 5: Number of parameters of the various BERT models

³https://huggingface.co/transformers/model_doc/bert.html

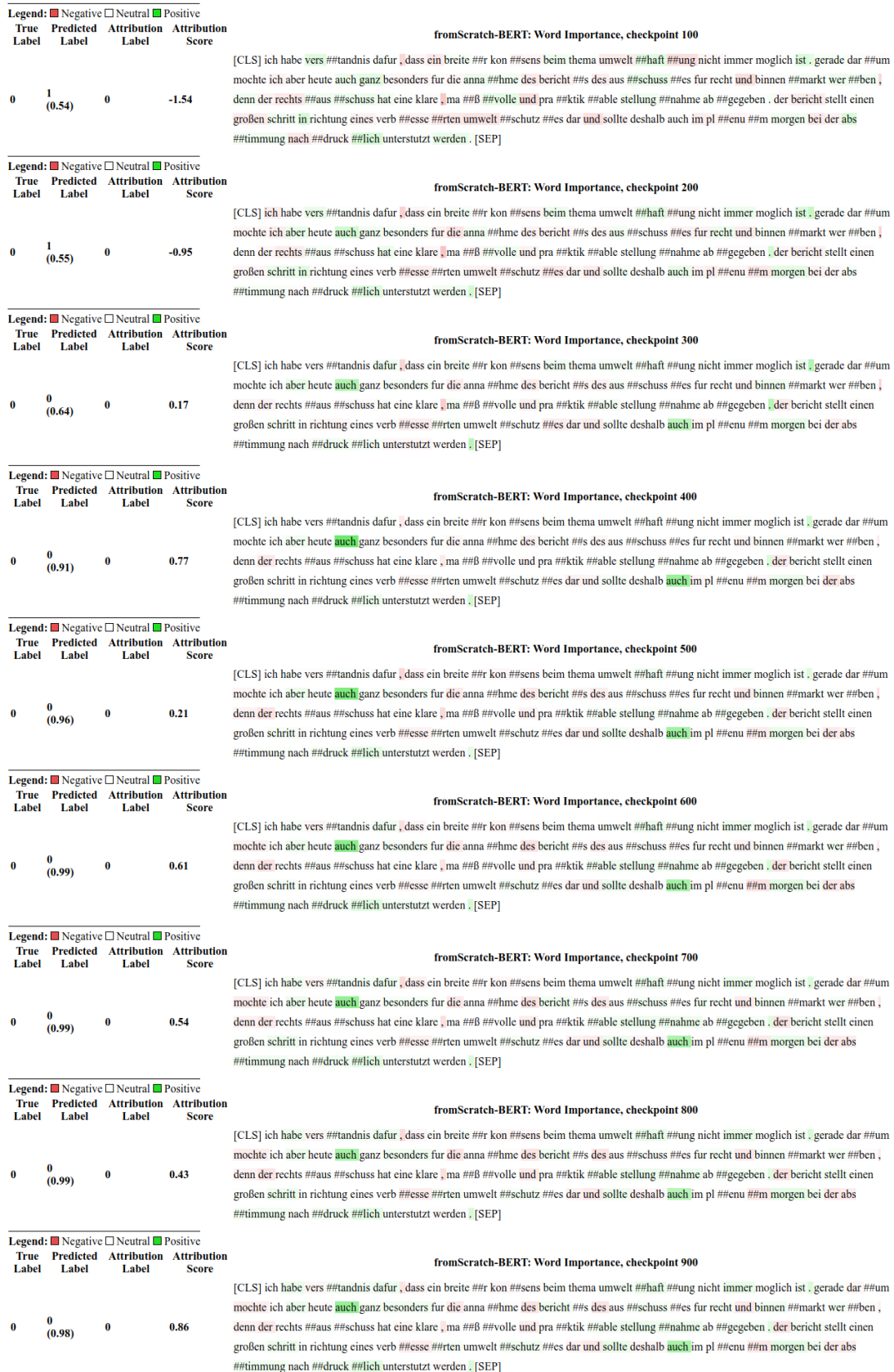


Figure 4: Layer Integrated Gradient saliency maps of input tokens contributing to the ground truth translationese label (here: original). BERT trained from scratch for translationese classification. Changes in attribution over the training checkpoints.

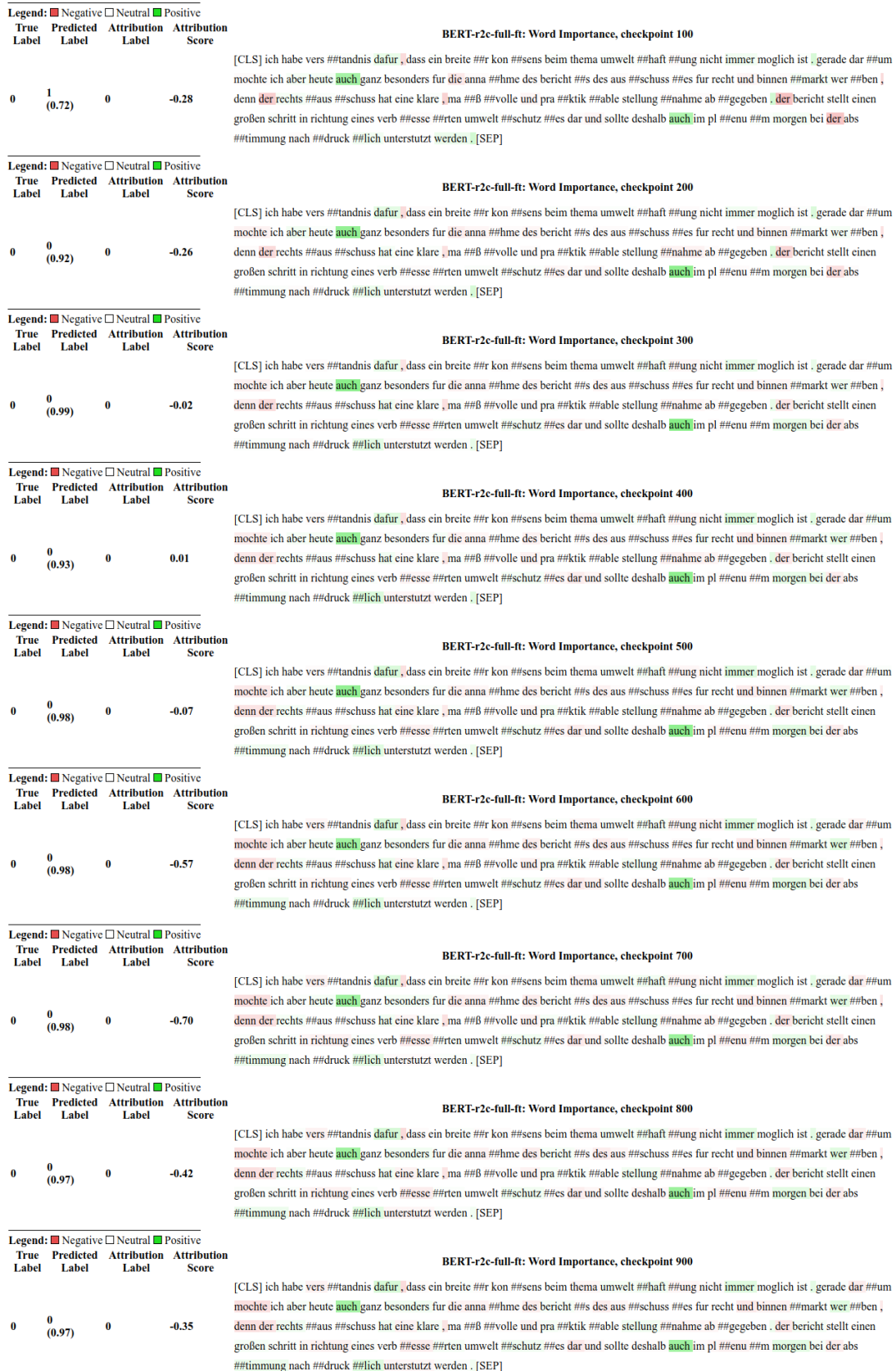


Figure 5: Layer Integrated Gradient saliency maps of input tokens contributing to the ground truth translationese label (here: original). BERT pretrained for handcrafted feature prediction, and fine-tuned for translationese classification. Changes in attribution over the training checkpoints.

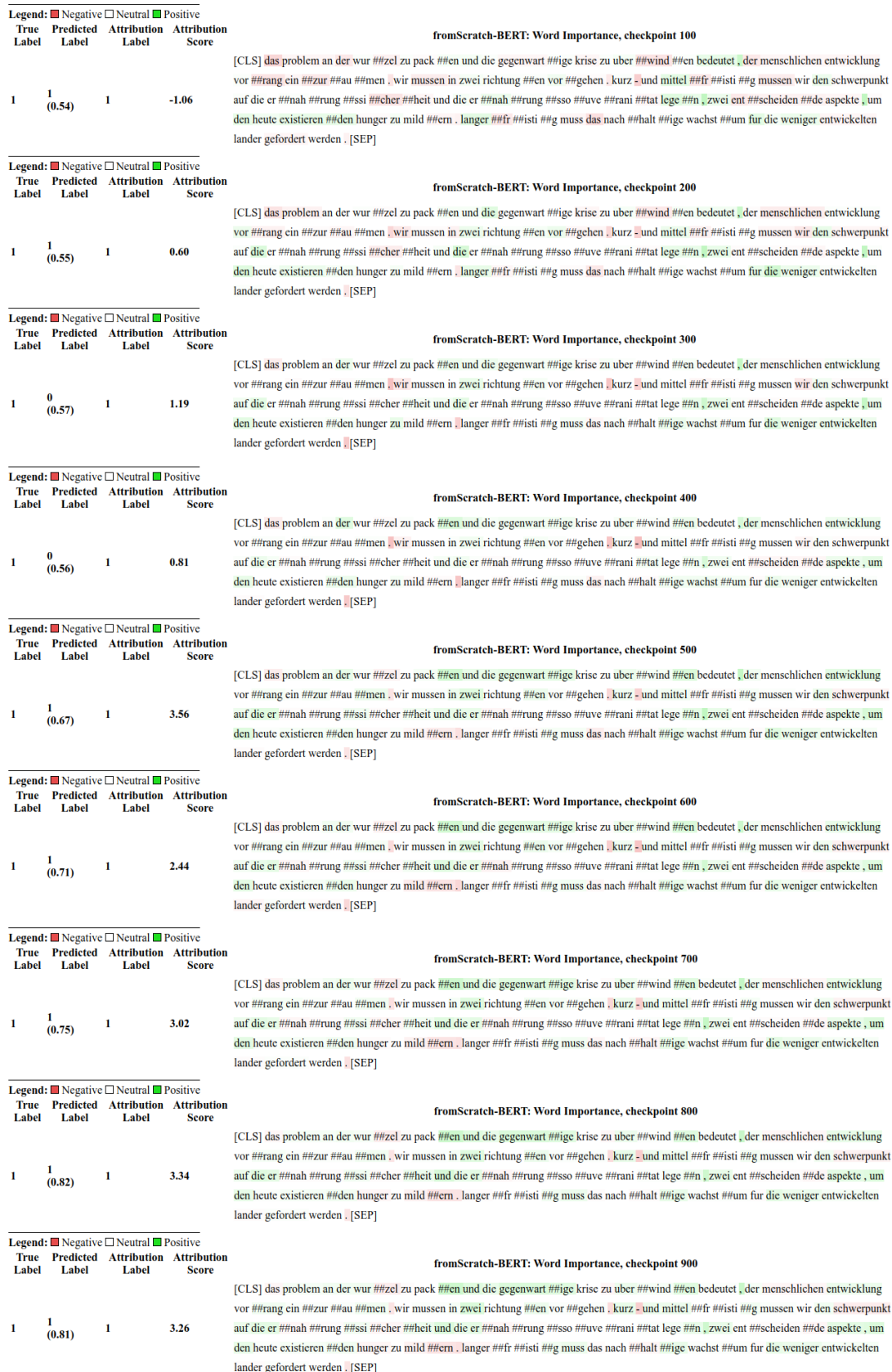


Figure 6: Layer Integrated Gradient saliency maps of input tokens contributing to the ground truth translationese label (here: translation). BERT trained from scratch for translationese classification. Changes in attribution over the training checkpoints.

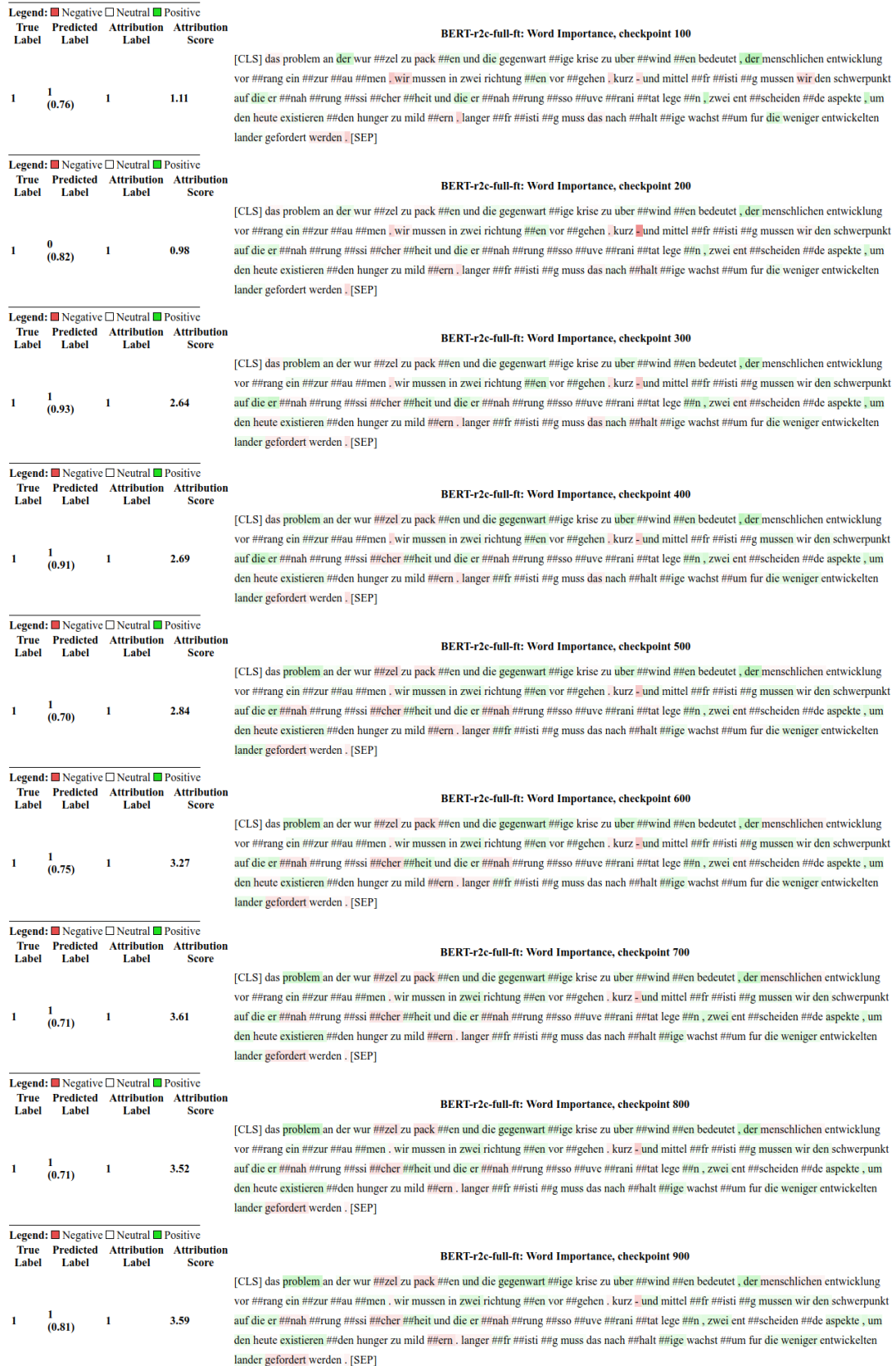


Figure 7: Layer Integrated Gradient saliency maps of input tokens contributing to the ground truth translationese label (here: translation). BERT pretrained for handcrafted feature prediction, and fine-tuned for translationese classification. Changes in attribution over the training checkpoints.