

On the current state of reproducibility and reporting of uncertainty for Aspect-based Sentiment Analysis

Elisabeth Lebmeier^{1,2}, Matthias Aßenmacher(✉)^{1,3}[0000-0003-2154-5774], and Christian Heumann^{1,4}[0000-0002-4718-595X]

¹ Department of Statistics (LMU), Ludwigstr. 33, D-80539 Munich, Germany

² e.lebmeier@gmx.de

³ matthias@stat.uni-muenchen.de

⁴ chris@stat.uni-muenchen.de

Abstract. For the latter part of the past decade, Aspect-Based Sentiment Analysis has been a field of great interest within Natural Language Processing. Supported by the Semantic Evaluation Conferences in 2014 – 2016, a variety of methods has been developed competing in improving performances on benchmark data sets. Exploiting the transformer architecture behind BERT, results improved rapidly and efforts in this direction still continue today. Our contribution to this body of research is a holistic comparison of six different architectures which achieved (near) state-of-the-art results at some point in time. We utilize a broad spectrum of five publicly available benchmark data sets and introduce a fixed setting with respect to the pre-processing, the train/validation splits, the performance measures and the quantification of uncertainty. Overall, our findings are two-fold: First, we find that the results reported in the scientific articles are hardly reproducible, since in our experiments the observed performance most of the time fell short of the reported one. Second, the results are burdened with notable uncertainty, depending on the data splits, which is why a reporting of uncertainty measures is crucial.

Keywords: Natural Language Processing · Sentiment Analysis · Pre-trained Language Models · Reproducibility

1 Introduction

The field of Natural Language Processing (NLP) has profited a lot from technical and algorithmic improvements within the last years. Before the successful times of machine learning and deep learning, NLP was mainly based on what linguists knew about how languages work, i.e. grammar and syntax. Thus, primarily rule-based approaches were employed in the past. Nowadays, far more generalized models based on neural networks are able to learn the desired language features.

On the other hand, data in written form is available in huge amounts and thus might be an important source for valuable information. For instance, the

internet is full of comparison portals, forums, blogs and social media posts where people state their opinions on a broad range of products, companies and other people. Product developers, politicians or other persons in charge could profit from this information and improve their products, decisions and behavior.

We specifically focus on *Aspect-Based Sentiment Analysis (ABSA)* in our work. ABSA is often used as a generic umbrella term for several unique tasks, which is caused by the inconsistency of terms in literature where many different names are widely used. To be as precise as possible, we explicitly use different terms than ABSA to refer to the exact tasks. The first one (subtask 2 [14]) assumes that in each text, aspect terms are already marked and thus given exactly as written in the text (this differs from so-called aspect categories which do not necessarily appear in the text). Here, the task is to classify the sentiments for those aspect terms. This is why the term *Aspect Term Sentiment Classification (ATSC)* is most accurate.

When referring to ATSC methods, we usually think of *single-task* approaches. These methods are designed to carry out only aspect term sentiment classification as the aspect terms are already given. Whether these were identified manually or by an algorithm is not relevant in this setting. In practice, however, the aspect terms oftentimes are not already known. Thus, approaches dealing with the step of *Aspect Term Extraction (ATE)* have been developed. They can either work on their own or be combined with an ATSC method. For these combined methods, which we refer to as *ATE+ATSC*, one can further distinguish between *pipeline*, *joint* and *collapsed* models. In pipeline models, ATE and ATSC are simply stacked one after another, i.e. the output of the first model is used as input to the second model. The latter two are often also referred to as *multi-task* models, since both tasks are carried out simultaneously or in an alternating way. These models only differ in their labeling mechanisms: There are two label sets for joint models, one to indicate whether a word is part of an aspect term and the other one to state its polarity. For collapsed models, a unified labeling scheme indicates whether a word is part of a positive, negative or neutral aspect term or not.

We re-evaluate four different models for ATSC, covering a variety of different architectures. This encompasses Recurrent neural networks (RNNs), Capsule networks [6, 16], networks using a Local Context Focus (LCF [22]), BERT-based approaches [2]), as well as two different ATE+ATSC models, one of which is a pipeline approach while the other one works in a collapsed fashion. All models are re-trained five times using five different (identical) train/validation splits and tested on the respective test sets in order to (i) compare them on a common ground and (ii) quantify the epistemic uncertainty associated with the architectures and the data.

2 Related Work

Related experiments were conducted by Mukherjee et al. [11], yet with a different focus. On the one hand, the authors also try to reproduce results on the

benchmark data sets from SemEval-14 about restaurants and laptops. However, they selected six other models than we did for which the implementations are provided in one repository.⁵ For these, the authors observed a consistent drop of 1-2 % with respect to both accuracy and macro-averaged F1-Score (F_1^{macro}). Mukherjee et al. [11] report a doubling of this drop when using 15% of the training data as validation data. On the other hand, they executed additional tasks which included the creation of two new data sets about men’s t-shirts and television as well as model evaluation on these data sets. Furthermore, they also experimented with cross-domain training and testing. Yet, several important points are not addressed by their work which is why we investigate them in our work. First, while they mostly care about comparing different types of architectures (memory networks vs. BERT), we instead focus on comparing the best performing models for different *tasks* (ATSC vs. ATE+ATSC). Further, we cover a larger variety of types of architectures by selecting the best performing representatives of several different types. Second, we stick closer to the original implementations (by using them, if available) whereas they exclusively rely on community designed implementations, which adds a further potential source of errors. Third, and most important, we provide estimates for the epistemic uncertainty of performance values and are thus able to (at least tentatively) explain performance differences due to different reporting standards.

3 Materials and Methods

This section will introduce the data sets we utilized for training and evaluation as well as the selected model architectures. We start by briefly explaining the data, before the models are described, since (reported) performance values on these data sets partly motivate our choices regarding the models. We selected these data sets as they are either widely known benchmark data sets or interesting adaptations of them. We acknowledge that their sizes are not be that large, yet, the pool of available data sets for this kind of tasks is rather small. Descriptive statistics for all used data sets can be found in Table 1. Note that the data sets we eventually use for training and testing the models are all based on the *original* train/test splits. Further we apply *small* modifications (as described below) which were (a) also applied by some of the authors whose models we re-evaluate and (b) we perceive as reasonable. This allows us to evaluate all of the architectures on a common ground, which is not possible by comparing the reported values from the original publications alone. Nevertheless, we are aware of the fact that this might limit comparability of our results to the original ones to some extent.

3.1 Data Sets

SemEval-14 Restaurants This data set contains reviews about restaurants in New York. Pontiki et al. [14] chose a subset of the restaurant data from Ganu et

⁵ <https://github.com/songyouwei/ABSA-PyTorch>

al. [4] as training data⁶, while collecting test data⁷ themselves. Both were labeled for several subtasks in the same way. These data sets were designed for ATSC as well as its equivalent on category level, *Aspect Category Sentiment Classification* (ACSC), but we stick to ATSC samples only. For each identified aspect term within a sentence, the polarity is given as *positive*, *negative*, *neutral* or *conflict*. We deleted the labels of the latter category (*conflict*) from the data sets due to their rare appearance. This is similar to previous work [3, 1, 21, 8], yet, they do not all mention or explain the removing process explicitly. Rarely appearing duplicate sentences which occurred in the training set were also removed in our work. Due to their small amount, this procedure should not cause severe problems concerning the over- or underestimation of the applied metrics. We speculate that this rare appearance of duplicates also might be the reason for why a similar preprocessing step was, to the best of our knowledge, only taken in one other work [20].

MAMS A *Multi-Aspect Multi-Sentiment (MAMS)* data set for the restaurant domain was introduced by Jiang et al. [7] who criticized existing data sets for not being adequate for ATSC. Since the data sets described above mainly consist of sentences which exhibit (i) only one single aspect or (ii) several aspects with the same sentiment, they argued that the task would not be much more difficult than a sentiment prediction on the sentence-level. To circumvent this issue, they extracted sentences of Ganu et al. [4] which comprise at least two aspects with differing sentiments.⁸ The data sets have the same structure as the SemEval-14 data sets, with the difference that Jiang et al. [7] provide a fixed validation set for MAMS. The size of the validation split comprises about ten percent of the whole training set, which also inspired our choice when it comes to creating train/validation splits from the two SemEval-14 training data sets.

ARTS Xing et al. [19] questioned the suitability of existing data sets for testing the aspect robustness of a model, i.e. whether the model is able to correctly identify the words corresponding to the chosen aspect term and predict its sentiment only based on them. Thus, the authors created an automatic generation framework that takes SemEval-14 test data (restaurants and laptops) as input and creates an *Aspect Robustness Test Set (ARTS)*. They used three different strategies to enrich the existing test set: The first one, REVTGT ("*reverse target*"), aims at reversing the sentiment of the chosen aspect term (called "*target aspect*"). This is reached by flipping the opinion using antonyms or adding negation terms like "not". Additionally, conjunctions may be changed in order to make sentences sound more fluent. Another strategy to augment the test set is

⁶ <http://metashare.ilsp.gr:8080/repository/browse/semEval-2014-absa-restaurant-reviews-train-data/479d18c0625011e38685842b2b6a04d72cb57ba6c07743b9879d1a04e72185b8/>

⁷ <http://metashare.ilsp.gr:8080/repository/browse/semEval-2014-absa-test-data-gold-annotations/b98d11cec18211e38229842b2b6a04d77591d40acd7542b7af823a54fb03a155/>

⁸ <https://github.com/siat-nlp/MAMS-for-ABSA>

REVNON ("*reverse non-target*") for which the sentiment of non-target aspects are (i) changed if they have the same sentiment as the target aspect or (ii) exaggerated if the non-target aspect is of a differing polarity. The third strategy called ADDDIFF ("*add different sentiment*") adds non-target aspects with an opposite sentiment which is intended to confuse the model. These non-target aspects are selected from a set of aspects collected from the whole data set and appended to the end of the sentence. ARTS are only designed to be used as test sets after training an architecture on the respective SemEval-14 training sets. The test sets for both restaurants and laptops are publicly available.⁹ During the preparation of the ARTS data for CapsNet-BERT, we noticed that the start and end positions of some aspect terms were not correct. We changed them in order to make the code work properly and we also deleted duplicates (cf. [20]). For these specific test sets, the *Aspect Robustness Score (ARS)* was introduced by Xing et al. [19] in order to measure how well models can deal with variations of sentences. Therefore, each sentence and all its variations are regarded as one unit of observation for which the prediction is only considered to be correct if the predictions for *all* variations are correct. These units alongside with their corresponding predictions are then used to compute the regular accuracy (*ARS accuracy*) on the level of the observational unit.

SemEval-14 Laptops The second domain-specific subset of the SemEval-14 data is on laptops. The data were collected and annotated by Pontiki et al. [14] for the task of ATE and/or ATSC. The training data set is publicly available,¹⁰ just like the test data (see Footnote 7). Again, there were duplicate sentences in the training data which we deleted (cf. [20]). Unlike other benchmark data sets, both SemEval-14 data sets come without an official train/validation split.

More Data Sets Recently more data sets have been published in addition to the ones mentioned beforehand. Mukherjee et al. [11] proposed two new data sets about *men's t-Shirts* and *television*. The YASO data set [12] has a different structure as it is a multi-domain collection. This is an interesting approach, yet also the reason for not considering it for our experiments: This data set is far better suited for cross-domain analyses, which is out of the scope of this work.

3.2 Models

MGATN A *multi-grained attention network (MGATN)* was proposed by Fan et al. [3]. Its *multi-grained attention* is able to take into account the interaction between aspects. We chose MGATN since it is reported to be the best performing representative of RNN-based models on SemEval-14 data sets.

⁹ https://github.com/zhijing-jin/ARTS_TestSet

¹⁰ <http://metashare.ilsp.gr:8080/repository/browse/semEval-2014-absa-laptop-reviews-train-data/94748ff4624e11e38d18842b2b6a04d7ca9201ec33f34d74a8551626be122856>

Table 1: Descriptive Statistics for the five utilized data sets. "*Multi-Sentiment sentences*" are those with at least two different polarities after removing "conflict" polarity. "*Aspect Terms in total*" also exclude "conflict".

| Data Set | Subset | Original Sentences in total | Sentences without Duplicates | Sentences for 3-class ATSC | Multi-Sentiment Sentences | Aspect Terms in total | Positive Aspect Terms | Negative Aspect Terms | Neutral Aspect Terms | Removed Conflict Aspect Terms |
|------------------------|------------|-----------------------------|------------------------------|----------------------------|---------------------------|-----------------------|-----------------------|-----------------------|----------------------|-------------------------------|
| SemEval-14 Restaurants | Training | 3,044 | 3,038 | 1,978 | 320 | 3,605 | 2,161 | 807 | 637 | 91 |
| | Test | 800 | 800 | 600 | 80 | 1,120 | 728 | 196 | 196 | 14 |
| SemEval-14 Laptops | Training | 3,048 | 3,036 | 1,460 | 166 | 2,317 | 988 | 866 | 463 | 45 |
| | Test | 800 | 800 | 411 | 38 | 638 | 341 | 128 | 169 | 16 |
| ARTS Restaurants | Test | 2,784 | 2,784 | 2,784 | 206 | 3,528 | 1,952 | 1,103 | 473 | 0 |
| ARTS Laptops | Test | 1,576 | 1,576 | 1,576 | 74 | 1,877 | 883 | 587 | 407 | 0 |
| MAMS Restaurant | Training | 4,297 | 4,297 | 4,297 | 4,297 | 11,186 | 3,380 | 2,764 | 5,042 | 0 |
| | Validation | 500 | 500 | 500 | 500 | 1,332 | 403 | 325 | 604 | 0 |
| | Test | 500 | 500 | 500 | 500 | 1,336 | 400 | 329 | 607 | 0 |

CapsNet-BERT Capsules networks were initially proposed for the field of computer vision, with the so-called *capsules* being responsible for recognizing certain implicit entities in images. Each capsule performs internal calculations and returns a probability that the corresponding entity appears in the image. A variation of capsule networks for ATSC and its combination with BERT was introduced by Jiang et al. [7]. It was reported to outperform all other capsule networks with respect to their accuracy on the SemEval-14 restaurants data. Additionally, it performed second-best on MAMS, which is why we selected it for this study. Furthermore, we assumed their results on SemEval-14 restaurants data to be for three-class classification, as all the other results they refer to are also three-class. Yet, it is not fully clear to us which makes this experiment even more interesting.

RGAT-BERT The *Relational Graph Attention Network (RGAT)* was introduced by Bai et al. [1]. It utilizes a dependency graph representing the syntactic relationships between words of a sentence as an additional input. The RGAT encoder creates syntax-aware aspect term embeddings following the representation update procedures from *Graph Attentional Networks (GATs)* [18]. It exhibits the best performance among graph-based models and also performs best on the MAMS data in terms of both accuracy and F_1^{macro} .

LCF-ATEPC Yang et al. [21] built upon the idea of the LCF mechanism. The local context of an aspect term is defined as a fixed-size window around it, words outside this window are taken into account with lower weights or not at all. For each input token two labels, for aspect and sentiment, are assigned according to the joint labeling scheme described in Sec. 1. We chose LCF-ATEPC to be part of this meta-study since it reached the highest F_1^{macro} and accuracy on SemEval-14 data of all approaches. Yet, this only holds for the variant that is trained using additional domain adaptation.

BERT+TFM The approach described by Li et al. [9] consists of a BERT model followed by a transformer (TFM [17]) layer for classification. BERT+TFM was the best model on SemEval-14 Laptops among all collapsed models at the time point of its introduction. There were also models using other layers on top instead of the transformer layer, but our variant of choice was TFM as it produced slightly better results than the concurring models.

GRACE GRACE, a *Gradient Harmonized and Cascaded Labeling* model introduced by Luo et al. [10], belongs to the category of pipeline approaches. It includes a post-training step of the pre-trained BERT model using Yelp¹¹ and Amazon data [5]. The post-trained model then shares its first l layers between the ATE and the ATSC task. The remaining layers are only used for the former. They are followed by a classification layer for the detected aspect terms. These classification outputs are then used again as inputs for a Transformer decoder

¹¹ <https://www.yelp.com/dataset>

which performs sentiment classification. The principle of using the first set of labels as input for the second is called *Cascaded Labeling* here and is assumed to deal with interactions between different aspect terms. *Gradient Harmonization* is applied in order to cope with imbalanced labels during training. GRACE appears to be the best performing one of the pipeline models according to the literature. Furthermore, it is reported to be the best ATE+ATSC model on both SemEval-14 data sets. However, these successes have to be taken into account with care, as their results are based on four-class classification. This means that in comparison to the other authors’ settings they did not exclude conflicting reviews of SemEval-14 data. Thus, our analysis contributes to comparability even more since it has not been established yet for the model/data combinations we examine.

4 Experiments

We re-evaluate six models (cf. Sec. 3.2) on the five data sets for the English language presented in Section 3.1. Our overall goals are to establish comparability between the models, to examine whether reported performance can be reproduced and to quantify epistemic model uncertainty that might exist due to the lacking knowledge about the train/validation splits. The entire code from our experiments is publicly available on GitHub.¹²

Our proceeding is as follows: First, we re-use the implementations provided by the authors by simply cloning their git repositories and adjusting them to our setup. Subsequently we try to reproduce their results on the data sets they used. Second, we adapt their code to the remaining data sets and conduct the necessary modifications, again sticking as closely as possible to the original hyperparameter settings (cf. Table 2 in the supplementary material). The biggest change we made was increasing the number of training epochs drastically and adding an early stopping mechanism. Apart from that, we did not engage in hyperparameter tuning in order not to modify/falsify the results. For all ATSC models, we selected the optimal model during the training process based on the validation accuracy and/or F_1^{macro} . For performing the experiments, we had one *Tesla V100 PCIe 16GB* GPU at our disposal.

Data Preparation Unlike other data sets, both SemEval-14 data sets come *without* an official validation split. Thus, we created five different train/validation splits (90/10) for each of the two SemEval-14 training sets. For each split, five training runs with different random initializations were conducted per model. The resulting 25 different versions per model per data set were subsequently evaluated on the two official SemEval-14 test sets (restaurants and laptops) as well as on the ARTS test sets. In Section 5 we report overall means per model per test set as well as means and standard deviations per model and test set for each of the different splits. Since there is an official validation set for MAMS, we did not apply the splitting procedure from above when training on this data

¹² <https://github.com/el-ma-le/atsc-experiments-official>

set. Consequently, the given means and standard deviations are based on five training runs with different random initializations only.

MGATN As there exists no publicly available implementation provided by its authors, we used the one from a collection of re-implemented ABSA methods from GitHub.¹³ We slightly modified the early stopping mechanism from that repository and then implemented it also for the other re-evaluated models.

CapsNet-BERT We used the implementation of CapsNet-BERT provided by its authors.¹⁴

RGAT-BERT We relied on the implementation of RGAT-BERT provided by its authors.¹⁵ Since the authors manually created an accuracy score different to the one implemented in `scikit-learn`¹⁶ [13], we substituted their metric by the `scikit-learn` variant to ensure comparability. For data transformation, we selected the stanza tokenizer [15] over the Deep Biaffine Parser,¹⁷ which was used by Bai et al. [1], since the former provides the necessary syntactic information, whereas the latter failed to produce the syntactic dependency relation tags and head IDs the model requires.

LCF-ATEPC We were not able to run the best-performing LCF-ATEPC variant based on domain adaptation due to missing pre-trained models. Thus, we decided to go for the second best, LCF-ATEPC-Fusion, using the official implementation of LCF-ATEPC.¹⁸ During our experiments, the authors of LCF-ATEPC started building a new repository¹⁹ based on the existing code which we did not use as it was still subject to changes.

BERT+TFM We used the implementation of BERT+TFM provided by its authors.²⁰ Our model selection was based on F_1^{micro} and F_1^{macro} , which were calculated based on *(start position, end position, polarity)*-triples for each identified aspect. Due to the collapsed labeling scheme, these scores account for both ATE and ATSC.

GRACE We used the post-trained BERT model provided by Luo et al. [10].²¹ Our model selection was based on $ATSC-F_1^{micro}$ and $-F_1^{macro}$ as well as on $ATE-F_1^{micro}$, with their calculations being slightly adjusted in order to match the calculation of those from BERT+TFM.

¹³ <https://github.com/songyouwei/ABSA-PyTorch>

¹⁴ <https://github.com/siat-nlp/MAMS-for-ABSA>

¹⁵ <https://github.com/muyeby/RGAT-ABSA>

¹⁶ https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html

¹⁷ <https://github.com/yzhangcs/parser>

¹⁸ <https://github.com/yangheng95/LCF-ATEPC>

¹⁹ <https://github.com/yangheng95/pyabsa>

²⁰ <https://github.com/lixin4ever/BERT-E2E-ABSA>

²¹ <https://github.com/ArrowLuo/GRACE>

5 Results

In general, reported values were not reproducible. Fig. 1 shows a comparison of our results (averaged over all 25 runs, including 95% confidence intervals) to the reported results from the original publications on the two SemEval-14 data sets. For all architectures there exists a notable gap between the blue (reproduced) and the orange (reported) values. In general, the gap tends to be larger for the ATSC models compared to the two ATE+ATSC models, where we were even able to reach a better performance for BERT+TFM within our replication study.²²

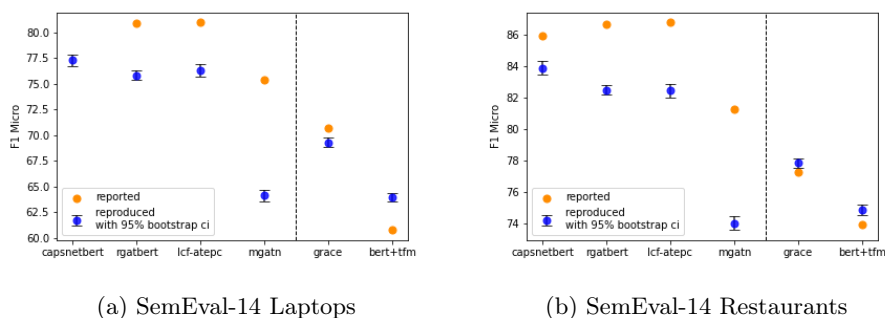


Fig. 1: Comparison of reported and reproduced performance. The reproduced value is the mean of all 25 runs per model in total. Further, 95% bootstrap ($n = 2000$) confidence intervals are displayed. Note that absolute performance of GRACE (four classes) and BERT+TFM cannot be compared to the other models due to different tasks. No F_1^{micro} was reported for CapsNet-BERT on SemEval-14 Laptops.

It is also interesting to see how different runs can lead to rather broad ranges of results, although having done only five training runs per model and data split. An example for this phenomenon is the Accuracy of MGATN on SemEval-14 Laptops (cf. Fig. 2). For the first, the fourth and fifth split, all of the values lie very close together (within $\text{mean} \pm \text{std}$), whereas the results of the other two splits show a rather high variance.

MGATN For MGATN, our reproduced results fell short of the reported values for accuracy, around five to ten percentage points for SemEval-2014 laptops and restaurants, respectively (cf. Tab. 5 and 6 in the supplementary material). Fig. 2 depicts the results on the laptops test set, the difference between reported and reproduced performance on the restaurant data (not shown) looks similar. A

²² We do not give a similar figure for MAMS or ARTS as there are not enough reported values to display the results in a meaningful way.

reason for this behavior might be that we could not use the official implementation of the authors, but had to rely on a re-implementation from the community. In terms of ARS accuracy on ARTS Restaurants, MGATN was the only model that reached only a single-digit value which means that it is not good at dealing with perturbed sentences.

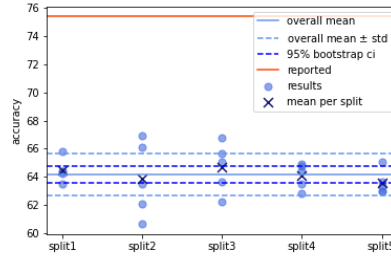


Fig. 2: Example for high differences between data splits: Accuracy of MGATN on SemEval-14 Laptops.

CapsNet-BERT Comparing all the selected models on the ATSC task, CapsNet-BERT performed best on all data sets regarding all the metrics except for ARS accuracy on the ARTS restaurant test set (cf. Tab. 5 and 6 in the supplementary material). For ARTS, it seems as if the reported ARS accuracy for laptops matched our result for restaurants, and vice versa, as Fig. 3 illustrates. As far as we can tell, we did not mix up the data sets during our calculations which makes this look quite peculiar. The difference between the reported and reproduced values on SemEval-14 Restaurants data (as shown in Fig. 1b) may be explained by the fact that we did three-class classification and we only assumed so for the reported value.

RGAT-BERT For both SemEval-14 and MAMS we missed the reported values by around five percentage points (cf. Tab. 5 and 6 in the supplementary material). ARTS restaurants is the only data set on which the best ARS accuracy was not reached by CapsNet-BERT, but RGAT-BERT. Regarding MAMS, Bai et al. [1] provided accuracy as well as F_1^{macro} , which is why we also compare these results here. Figure 4 shows the all five values of the four different measures as well as the average. For accuracy and F_1^{macro} , reported values from Bai et al. [1] were added.

LCF-ATEPC Our experiments on average resulted in about five percentage points lower accuracies for LCF-ATEPC than were reported. Yet, LCF-ATEPC reached the best ARS accuracy value on ARTS restaurant data in our analysis.

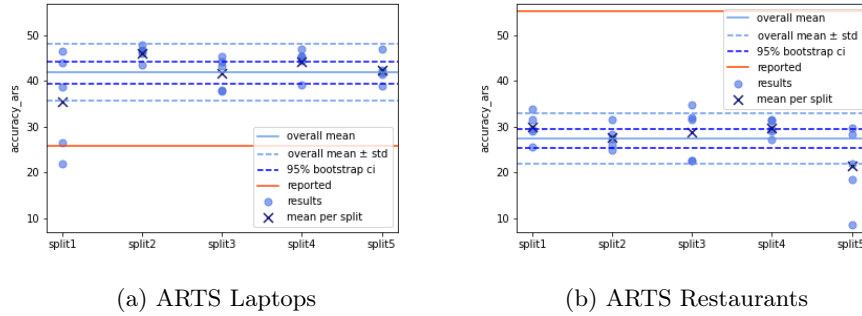


Fig. 3: Aspect Robustness Score (ARS) Accuracy of CapsNet-BERT.

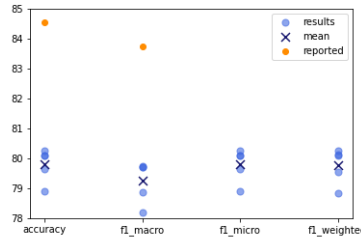
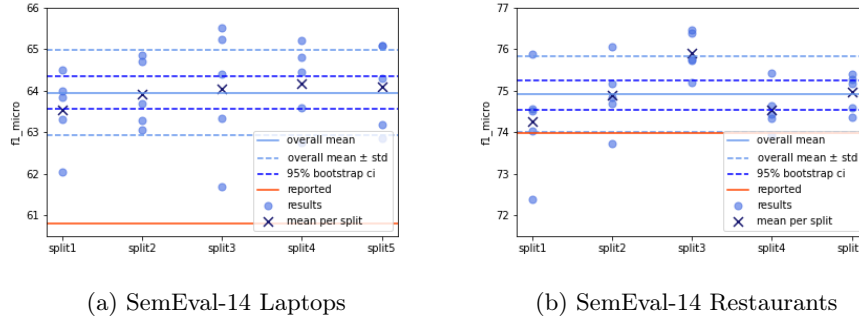
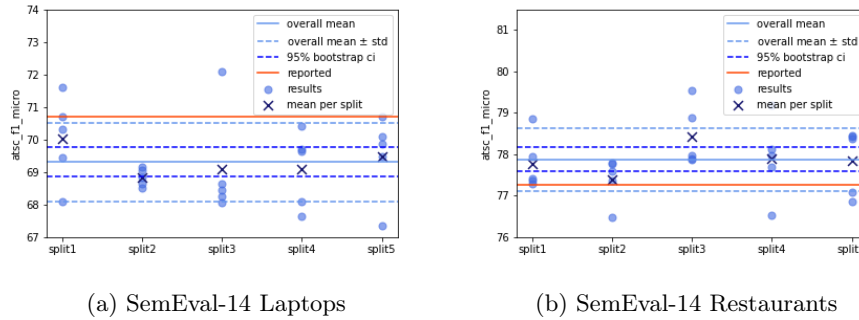


Fig. 4: Performance of RGAT-BERT on MAMS.

BERT+TFM In contrast to the majority of the other models, for BERT+TFM the (average) performance of our runs surpassed the reported performance values on the SemEval-14 data. As Fig. 5 indicates, this holds for all runs (laptop domain) and on average (restaurant domain). The reasons for our improved values may lie in the chosen hyperparameters, yet we cannot tell for sure.

GRACE During our experiments with GRACE, we were able to produce results approximately in the same range as the reported values. Regarding SemEval-14 restaurants our results on average were better than the reported ones (cf. Fig. 6b), while for laptops we could not quite reach the reported performance (cf. Fig. 6a). For the latter case, our results of single runs were better than (or at least equal to) the reported one, which is kind of a symptom of the problem. If we only reported the best of all runs, our conclusion would have been that we were able to outperform the original model. However, as we have already mentioned, reported results were based on four-class classification, whereas our results were made for three-class. This might be the reason for different results. In the ATE+ATSC task, GRACE outperformed BERT+TFM on all data sets except for MAMS (cf. Tab. 3 and 4 in the supplementary material).


 Fig. 5: F_1^{micro} of BERT+TFM.

 Fig. 6: ATSC F_1^{micro} of GRACE.

6 Discussion

6.1 General Takeaways

Results differing from the reported values can be explained by various reasons. First, we often do not know how the reported values were created, i.e. whether the authors took the best or an average value of their runs. In Fig. 6a, it becomes clearly visible that taking the best value compared to the mean over multiple runs yields a difference of about almost three percentage points. Unfortunately there are also, to the best of our knowledge, no clear guidelines for how to properly report the uncertainty resulting from different data splits. Second, our data are usually not exactly identical to the data sets used for the original papers due to the preprocessing steps we explained beforehand. Also, training and validation splits are probably different from ours. Some models required additional syntactical information which we (potentially) inferred from other packages than indicated, because either none were given or because the ones that were given did not work as stated. Third, hyperparameter configurations are often not totally clear due to a lack of concise descriptions in the original work.

In these cases we took those that were chosen by default in the implementations we used. Since those were not necessarily always provided by the authors of the models, we have no information about how close they are to the original configurations. What we could find out regarding hyperparameters can be found in Table 2 in the supplementary material. Consequently, it is not surprising that we were not able to exactly reproduce given results, since hyperparameter tuning often has a large impact on the model performance. This insight is also shared by Mukherjee et al. [11], although they tested other models in a different setup.

6.2 Possible guidelines

Taking all considerations into account, we want to tentatively propose some guidelines that might be beneficial for making NLP research reproducible and for quantifying different types of uncertainty. First, it is not enough to purely open-source your code but it also requires a thorough documentation and explanation. This should also include all the information about hyperparameters, additional training data, custom data splits (if applicable), and non-standard pre-processing, since all of this can have a (potentially) large impact on the results. Second, every information about potential randomness/variation in the results has to be acknowledged, ideally even researched further and reported/displayed properly. One potential starting point could be to *always* perform multiple runs on multiple different splits and use the results to report standard deviations between and within splits. While the former gives an impression for the uncertainty induced by data heterogeneity, the latter rather reflects the model’s share of the overall uncertainty. This would of course to some extent mean, to move away from (overly confidently) reporting single performance values. A reporting convention indicating a common procedure combined with already prepared data sets with all possible labels could improve the comparability between models a lot.

7 Conclusion & Future work

Our experiments revealed that reproducing reported results is hardly possible, given the current practice of performance reporting (at least for this subset of selected models). A tendency towards lower results is visible in our experiments, sometimes even five to ten percentage points lower than the original values. The only exception was BERT+TFM for which given values were surpassed. The reasons for these observations may lay in the data preprocessing steps, in the hyperparameters or in the absence of a convention on which values to report (best or mean of several runs). This discovery of models hardly being comparable based on their performance measures is a very important one from our point of view. When new models are proposed, one of the main aspects during their evaluation is the improvement with respect to the state-of-the-art. But when the performance of a single model can vary between single runs, the question is which results to take into account for model rankings. Also a huge practical meta-analysis of all models on several data sets would clarify the situation.

Acknowledgements

This work has been partially funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) as part of BERD@NFDI - grant number 460037581.

References

1. Bai, X., Liu, P., Zhang, Y.: Investigating typed syntactic dependencies for targeted sentiment classification using graph attention neural network. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **29**, 503–514 (2020). <https://doi.org/10.1109/taslp.2020.3042009>, <http://dx.doi.org/10.1109/TASLP.2020.3042009>
2. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>, <https://www.aclweb.org/anthology/N19-1423>
3. Fan, F., Feng, Y., Zhao, D.: Multi-grained attention network for aspect-level sentiment classification. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. pp. 3433–3442. Association for Computational Linguistics, Brussels, Belgium (Oct-Nov 2018). <https://doi.org/10.18653/v1/D18-1380>, <https://aclanthology.org/D18-1380>
4. Ganu, G., Elhadad, N., Marian, A.: Beyond the Stars: Improving Rating Predictions using Review Text Content. In: *Twelfth International Workshop on the Web and Databases (WebDB 2009)*. vol. 9, pp. 1–6. Citeseer (2009)
5. He, R., McAuley, J.: Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In: *Proceedings of the 25th International Conference on World Wide Web*. p. 507–517. WWW '16, International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE (2016). <https://doi.org/10.1145/2872427.2883037>, <https://doi.org/10.1145/2872427.2883037>
6. Hinton, G.E., Krizhevsky, A., Wang, S.D.: Transforming auto-encoders. In: *Proceedings of the 21th International Conference on Artificial Neural Networks - Volume Part I*. p. 44–51. ICANN'11, Springer-Verlag, Berlin, Heidelberg (2011)
7. Jiang, Q., Chen, L., Xu, R., Ao, X., Yang, M.: A challenge dataset and effective models for aspect-based sentiment analysis. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. pp. 6280–6285. Association for Computational Linguistics, Hong Kong, China (Nov 2019). <https://doi.org/10.18653/v1/D19-1654>, <https://aclanthology.org/D19-1654>
8. Li, X., Bing, L., Li, P., Lam, W.: A unified model for opinion target extraction and target sentiment prediction. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 33, pp. 6714–6721 (2019)
9. Li, X., Bing, L., Zhang, W., Lam, W.: Exploiting BERT for end-to-end aspect-based sentiment analysis. In: *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*. pp. 34–41. Association for Computational Lin-

- guistics, Hong Kong, China (Nov 2019). <https://doi.org/10.18653/v1/D19-5505>, <https://aclanthology.org/D19-5505>
10. Luo, H., Ji, L., Li, T., Jiang, D., Duan, N.: GRACE: Gradient harmonized and cascaded labeling for aspect-based sentiment analysis. In: Findings of the Association for Computational Linguistics: EMNLP 2020. pp. 54–64. Association for Computational Linguistics, Online (Nov 2020). <https://doi.org/10.18653/v1/2020.findings-emnlp.6>, <https://aclanthology.org/2020.findings-emnlp.6>
 11. Mukherjee, R., Shetty, S., Chattopadhyay, S., Maji, S., Datta, S., Goyal, P.: Reproducibility, replicability and beyond: Assessing production readiness of aspect based sentiment analysis in the wild. arXiv preprint arXiv:2101.09449 (2021)
 12. Orbach, M., Toledo-Ronen, O., Spector, A., Aharonov, R., Katz, Y., Slonim, N.: Yaso: A new benchmark for targeted sentiment analysis. arXiv preprint arXiv:2012.14541 (2020)
 13. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011)
 14. Pontiki, M., Galanis, D., Pavlopoulos, J., Papageorgiou, H., Androutsopoulos, I., Manandhar, S.: SemEval-2014 task 4: Aspect based sentiment analysis. In: Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014). pp. 27–35. Association for Computational Linguistics, Dublin, Ireland (Aug 2014). <https://doi.org/10.3115/v1/S14-2004>, <https://aclanthology.org/S14-2004>
 15. Qi, P., Zhang, Y., Zhang, Y., Bolton, J., Manning, C.D.: Stanza: A Python natural language processing toolkit for many human languages. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations (2020), <https://nlp.stanford.edu/pubs/qi2020stanza.pdf>
 16. Sabour, S., Frosst, N., Hinton, G.E.: Dynamic routing between capsules. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 30. Curran Associates, Inc. (2017), <https://proceedings.neurips.cc/paper/2017/file/2cad8fa47bbef282badbb8de5374b894-Paper.pdf>
 17. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *CoRR* **abs/1706.03762** (2017), <http://arxiv.org/abs/1706.03762>
 18. Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph attention networks. In: International Conference on Learning Representations (2018), <https://openreview.net/forum?id=rJXmpikCZ>
 19. Xing, X., Jin, Z., Jin, D., Wang, B., Zhang, Q., Huang, X.: Tasty burgers, soggy fries: Probing aspect robustness in aspect-based sentiment analysis. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP) pp. 3594–3605 (2020). <https://doi.org/10.18653/v1/2020.emnlp-main.292>, <https://www.aclweb.org/anthology/2020.emnlp-main.292>
 20. Xue, W., Li, T.: Aspect based sentiment analysis with gated convolutional networks. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 2514–2523. Association for Computational Linguistics, Melbourne, Australia (Jul 2018). <https://doi.org/10.18653/v1/P18-1234>, <https://aclanthology.org/P18-1234>
 21. Yang, H., Zeng, B., Yang, J., Song, Y., Xu, R.: A multi-task learning model for chinese-oriented aspect polarity classification

- and aspect term extraction. *Neurocomputing* **419**, 344–356 (2021).
<https://doi.org/https://doi.org/10.1016/j.neucom.2020.08.001>, <https://www.sciencedirect.com/science/article/pii/S0925231220312534>
22. Zeng, B., Yang, H., Xu, R., Zhou, W., Han, X.: Lcf: A local context focus mechanism for aspect-based sentiment classification. *Applied Sciences* **9**, 3389 (2019)