Making the most of small Software Engineering datasets with modern machine learning

Julian Aron Prenner, Romain Robbes

Abstract—This paper provides a starting point for Software Engineering (SE) researchers and practitioners faced with the problem of training machine learning models on small datasets. Due to the high costs associated with labeling data, in Software Engineering, there exist many small (< 1000 samples) and medium-sized (< 100 000 samples) datasets. While deep learning has set the state of the art in many machine learning tasks, it is only recently that it has proven effective on small-sized datasets, primarily thanks to pre-training, a semi-supervised learning technique that leverages abundant unlabelled data alongside scarce labelled data. In this work, we evaluate pre-trained Transformer models on a selection of 13 smaller datasets from the SE literature, covering both, source code and natural language. Our results suggest that pre-trained Transformers are competitive and in some cases superior to previous models, especially for tasks involving natural language; whereas for source code tasks, in particular for very small datasets, traditional machine learning methods often has the edge.

In addition, we experiment with several techniques that ought to aid training on small datasets, including active learning, data augmentation, soft labels, self-training and intermediate-task fine-tuning, and issue recommendations on when they are effective. We also release all the data, scripts, and most importantly pre-trained models for the community to reuse on their own datasets.

Index Terms—Small Datasets, Transformer, BERT, RoBERTA, Pre-training, Fine-Tuning, Data Augmentation, Back Translation, Soft Labels, Active Learning.

1 Introduction

MALL datasets are commonplace for many Software Engineering problems. While the creation of a labelled dataset is always a significant undertaking, this is even more the case for Software Engineering. In many cases, significant expert knowledge is required to label Software Engineering data, making it difficult to use crowd-sourcing techniques, as is often done in other fields such as in computer vision [1]. Moreover, some labelling tasks involve detailed (text or source code) understanding, making the labelling of a single example time consuming. Dataset size may be further reduced by the need to label the same examples multiple times and to compute inter-rater agreement. Due to all these factors, it is thus not uncommon for hand-labelled Software Engineering datasets to number only a few thousands or even hundreds of samples. For instance, the 13 datasets used in this work (described in Section 2) range from 200 to 62,275 samples, with three datasets having more than 5,000 samples, and four having less than 1,000.

Historically, small Software Engineering datasets were used with traditional machine learning algorithms, such as Support Vector Machines (SVMs), Logistic Regression or Random Forests, often combined with manual feature engineering. In recent years, early experiments with deep learning architectures [2]–[5], such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNSs), showed mixed results, suggesting that for many tasks where training data is scarce, deep learning does not provide a clear benefit, especially in light of its considerably higher computational costs.

Whether deep learning is in fact not well suited for these small datasets, and if so, for which kind of tasks and dataset sizes is the central question of this paper. The motivation to take a second look at this problem is the recent advent of semi-supervised learning [6]. Semi-supervised learning is a machine learning paradigm in which both labelled data and unlabelled data are leveraged in the learning process, with the latter being much cheaper to acquire. While prevalent in computer vision, it is only since 2018 that semi-supervised learning has become viable in the NLP domain [7], [8], in the form of pre-training. Since 2019, pre-trained Transformerbased models such as BERT [9] or RoBERTa [10] have set many records in NLP and related fields, and considerably improved the state of the art on important benchmarks such as GLUE [11] and SQuAD [12], in which there are small datasets. In addition to pre-training, several additional techniques have the potential to benefit small datasets, including Domain-specific Fine-Tuning, Intermediate-Task fine-tuning, Active Learning, Self-Training, Data Augmentation, and Soft Labels. Thus, a better understanding of when to combine these techniques is necessary. We provide background on the Transformer Architecture, Semi-supervised learning via pre-training and other techniques in Section 3.

While previous work showed promising results in applying pre-training to SE problems [13]–[17], this work examines this phenomenon in more depth, by applying the pre-training paradigm on thirteen different small and medium-sized Software Engineering datasets selected from the literature. These datasets span natural language, source code, and source code comments, in a variety of domains (several sentiment analysis tasks, several app review classification tasks, technical debt detection, comment classification, code comment coherence, code smell detection, code readability, code complexity). In addition, and unlike previous work, we also investigate the impact of the additional techniques mentioned above, when they are rele-

vant. Section 4 presents methodological details such as preprocessing, baselines, and training, testing and validation modalities, for all the scenarios we consider. This section also presents the pre-trained and fine-tuned models we use in this work, including StackOBERTflow, a Transformer model pre-trained on 26 million Stack Overflow comments.

Section 5 presents the results of the paper, answering the following research questions:

- RQ1: For which domains and tasks does the pretraining paradigm outperform the baselines? We find that pre-training is effective for tasks working on natural language and source code comments, but is not as effective for tasks working on source code yet.
- RQ2: Which additional techniques are effective, and if so in which circumstances? We find that some techniques, such as domain-specific pre-training, and data augmentation are effective in some (but not all) settings, while we find limited evidence for the effectiveness of others, such as active learning.

Finally, we close the paper by documenting the limitations of our study, and the opportunities for additional studies in Section 6. We conclude the work in Section 7, summarizing initial recommendations on the effectiveness of pre-training and the additional techniques. Additional material can be found in three appendices: Appendix A provides additional information on datasets; Appendix B provides results; last but not least, Appendix C provides instructions on how to access the data, scripts, and pre-trained models we used in our experiments. These models can be fine-tuned for a wide range of tasks relating to software artifacts; we hope that they will prove useful to other researchers in the field.

2 DATASETS AND RELATED WORK

We selected thirteen datasets introduced in the Software Engineering literature in recent years, aiming for both variety in terms of artifacts, dataset size, and classification tasks [2], [3], [18]-[26]. Another selection factor was the availability of a comparable baseline or a way to reproduce the initial experiment. Seven datasets involve natural language, two code comments, and the remaining four source code (one with comments). They vary from 341 to 62,275 examples, and from 2 to 16 classes. The datasets cover nine different tasks: a) sentiment classification of software artifacts, such as Stack Overflow comments, app and code reviews b) detection of informative app reviews c) classification of app reviews d) detection of self-admitted technical debt through code comments e) classification of code comments f) prediction of code-comment coherence g) detection of linguistic code smells h) prediction of code runtime complexity and i) prediction of code readability.

Next, we discuss each of these tasks in greater detail. For a more concise overview refer to Table 1.

2.1 Natural Language Datasets

These datasets contain mainly natural language, but may occasionally contain some source code identifiers. They are the closest to the original setting for models pre-trained on a generic English corpus, although they come from very specific domains.

Name	Size	# Cl.	Type	Usage
Sentiment Classification (Stack Overflow) [18]	4 423	3	Natural language	Train, Test, Valid.
Sentiment Classification (Stack Overflow) [19]	1 500	3	Natural language	Test
Sentiment Classification (JIRA Issues) [3]	926	2	Natural language	Test
Sentiment Classification (App Reviews) [19]	341	3	Natural language	Test
Informative App Review Detection [20]	12 000	2	Natural language	Train, Test, Valid.
App Review Classification [21]	3 691	4	Natural language	Train, Valid.
App Review Classification [27]	3 000	7	Natural language	Train, Valid.
Self-Admitted Technical Debt Detection [22]	62 275	2	Comments	Train, Test, Valid.
Comment Classification [23]	11 232	16	Comments	Train, Valid.
Code-Comment Coherence Prediction [24]	2 881	2	Code w/ Comments	Train, Test, Valid.
Linguistic Smell Detection [2]	1753	2	Code	Train, Valid.
Code Runtime Complexity Classification [25]	933	5	Code	Train, Valid.
Code Readability Prediction [26]	200	2	Code	Train, Valid.

TABLE 1
Datasets considered in this work, along with their size, number of classes and usage.

Sentiment Classification

In sentiment classification, a model assigns a sentiment class (e.g., one of *positive*, *negative*, *neutral*) to a sentence or short piece of text; in our case the text's domain is related to software development.

The dataset compiled by Calefato *et al.* (Senti4SD) contains 3 097 training and 1 326 test samples each labeled as either *positive*, *negative* or *neutral*; all samples were extracted from questions, answers and comments of Stack Overflow posts. Along with their dataset, the authors also released individual rater annotations, i.e., three rater labels per sample, from which the final labels were obtained by applying a majority vote rule. The dataset is used to evaluate an SVM classifier using word embeddings trained on a Stack Overflow corpus as features; we include this SVM as a baseline (see Table 8).

Lin et al. created a sentiment classification datasets consisting of 1,500 sentences (from 178 text fragments) extracted from Stack Overflow discussions. They also adapted two previous datasets of 636 Jira issues (926 sentences), and 130 app reviews (341 sentences) [28], [29]. The JIRA issues dataset has only two sentiment classes (positive and negative) while the other two have an additional neutral class. All three datasets were used in a study of several sentiment analysis tools, and a novel model introduced in the same work [3]; we include all of them as baselines in our compari-

son (Table 8). To allow for a comparison with these baselines, all three datasets are *only used as test sets* in this work.

Informative App Review Detection

An app review is considered informative if it contains valuable information for the application developer, such as feature suggestions or bug reports. The dataset presented by Chen *et al.* [20] contains 12 000 app reviews belonging to four popular mobile apps from the Google Play Store. The dataset is partitioned into predefined test (2000 samples per app) and train sets (1000 samples per app). Three raters annotated each review as either *informative* or *non-informative*, with a majority vote to determine the final label. This work also presented AR-MINER, a tool based on an expectation maximization with Naive Bayes (EMNB) classifier to detect such informative app reviews; this tool is our baseline.

App Review Classification

Maalej et al. [21] compiled a dataset of 4 400 reviews crawled from Apple's App Store and Google Play. Each review was categorized by two raters into one of four classes (bug report, feature request, user experience or rating); reviews with rater disagreement were discarded. The authors experimented with various types of classifiers, finding that an ensemble of binary classifiers performs considerably better than a single multiclass classifier. We only study multiclass classification, for comparability reasons, use the author's multiclass classifier as a baseline. Scalabrino et al. [27] and Villarroel et al. [29] introduce CLAP, a tool for automatic classification and clustering of app reviews that is evaluated on a dataset of 3000 app reviews created by the same authors. We use the CLAP dataset in a data augmentation experiment.

Related Work

There is a vast literature on sentiment analysis and classification. For a general overview see e.g., Mäntylä et al. [30]. Zhang et al. [15] provide an in-depth comparison of pre-trained Transformers with six different sentiment analysis tools, including both, tools for general sentiment analysis and tools specifically targeted towards Software Engineering. The study uses some datasets that we also use (e.g., from Lin et al. [3]), they defined custom training and testing sets while we used them solely for testing, which makes comparisons difficult. While they explored additional Transformer architectures (XLNet [31] and ALBERT [32]), their study was limited to only fine-tuning: they did not investigate the use of task-specific pre-training, nor any other of the additional methods that we study. Ahmed et al. [33], introduced SentiCR, a sentiment analysis tool that uses of Part Of Speech (POS) tags and Gradient Boosting Trees, while Chen et al. [34] present SEntiMoji, a model that leverages emojis to improve SE sentiment classification.

Dhinakaran *et al.* [35] investigated the application of *active learning* for app review analysis, using the dataset by Maalej *et al.* They employ traditional machine learning algorithms (naive Bayes, logistic regression, and SVM). A similar experiment, carried out on the same dataset can also be found in this work.

2.2 Datasets of Code Comments

Source code comments somewhat differ from natural language: they may often contain source code identifiers, code annotations, and specific idioms common in source code documentation.

Self-Admitted Technical Debt Detection

Self-admitted technical debt (SATD) is technical debt known to and acknowledged by the author. It is often expressed in code comments with a short description of a flaw or shortcoming and sometimes, but not always, marked with specific keywords, such as FIXME, TODO or HACK. Detection of such comments can be useful to assess software quality, aid decision-making or direct further development.

The dataset by S. Maldonado et al. contains SATD comments extracted from 10 prominent Java projects (for more details, see Table 15 in Appendix A). With over 60 000 samples, it is by far the largest dataset used in this work. Each dataset sample is assigned to one of five SATD categories informed by an established ontology of technical debt [36] (design debt, requirement debt, defect debt, documentation debt or test debt) or labeled as not containing any SATD at all. In this work, we concentrate on the binary version of the problem (detecting presence of SATD), for two reasons: (1) because there are large class imbalances (i.e., the document debt class makes up less than 0.1% of the total data), and (2) to compare performance to binary classifiers from previous work. In addition to introducing this dataset, the authors also perform various SATD detection and classification experiments using traditional NLP and machine learning methods; we use their SATD detection model as a baseline. We also include the CNN from Ren et al. [4] as a second baseline.

Comment Classification

Pascarella and Bacchelli [23] released a dataset of over 11 000 comments from open-source Java projects classified according to a taxonomy of 16 different comment categories. They evaluate their dataset on a multinomial Naive Bayes classifier which serves as our baseline for this task.

Related Work

For a more general survey on self-admitted technical debt see e.g., Sierra *et al.* [37]. Santos *et al.* [5] use a Long short-term memory network (LSTM) to *classify* SATD, also making use of the dataset by S. Maldonado *et al.*; since we confined ourselves to SATD *detection* this work was not included. A text-mining based approach to SATD detection can be found in [38], [39]; in this work, unfortunately, only a subset of the dataset was used and results are thus not comparable. In a closely related comment classification work, Pascarella [40] focus on comment classification in mobile applications.

2.3 Datasets of Source Code

These datasets significantly differ from natural language: on the one hand, code has a very specific and unambiguous syntax and are much more repetitive than natural language [41]; on the other hand, code has very complex semantics, and has many identifiers, leading to vocabulary issues [42].

Code-Comment Coherence Prediction

Corazza *et al.* [24], [43] and Cimasa *et al.* [44] examine the concept of code-comment coherence, i.e., the "relatedness" of a method's code and its lead comment. The authors introduce a dataset of 2 883 *Java* methods along with their leading comment and a binary label indicating whether coherence exists between the two or not [24]. In follow-up work, the authors trained an SVM classifier on their dataset using features based on tf-idf [43] and later word embedings [44]. We include both of these models as baselines.

Linguistic Smell Detection

Fakhoury *et al.* [2] examine the automatic detection of *linguistic* code smells (also known as linguistic antipatterns), that is, code smells emerging from the use of misleading identifier names or the violation of common naming conventions. Examples of this are variable names as if they were lists or arrays when in fact they have a scalar type, or getter methods with side effects.

They labeled a dataset of roughly 1700 code snippets, following a taxonomy of linguistic smells [45], which comprises 18 different types of linguistic antipatterns. They then trained a number of models on this dataset and compared their performances. These models include CNNs in various configurations, SVMs with different kernels and a Random Forest classifier. All models are binary, that is, they only determine whether a given sample is "smelly" or not and do distinguish between different types of antipatterns. Interestingly, the authors found that a thoroughly tuned SVM model outstrips the CNN in all its configurations. This not only in terms of performance metrics but also in terms of memory consumption, training time and ease of use. Their best-performing models were selected as baselines.

Code Runtime Complexity Classification

Sikka *et al.* [25] investigate the use of machine learning to automatically predict the code runtime complexity class (e.g., $\mathcal{O}(n^2)$) of short programs. To this end, they collected 933 Java implementations of various algorithms from a competitive programming platform and annotated each with the corresponding complexity class (i.e., one of $\mathcal{O}(1)$, $\mathcal{O}(\log n)$, $\mathcal{O}(n)$, $\mathcal{O}(n\log n)$, $\mathcal{O}(n^2)$). They experiment with various traditional machine learning approaches, training, such as Random Forests and SVMs on manually engineered features (such as numbers of loops, numbers of variables etc.) and code embeddings obtained from the programs' abstract syntax tree through graph2vec [46].

Code Readability Prediction

What constitutes readable code and what does not, seems to be largely a matter of personal taste. Notwithstanding this, research by Buse and Weimer [47] suggests that code readability can, at least in part, be measured objectively.

A relatively small number of papers [26], [47]–[50] examine models for automatic code readability estimation. Most recently, Scalabrino *et al.* [26], [50] compiled a dataset by letting 30 Computer Science students rate the readability of 200 methods, previously selected from well-known Java projects. Each method received 9 readability ratings: these ratings were then averaged and compared against a

threshold value to assign a single binary readability label. Further, they developed a logistic regression model for code readability estimation by combining structural readability features proposed in Buse and Weimer [47] and Dorn [49] with novel textual features. We base all of our experiments on above dataset; previous datasets by Buse and Weimer [47] and Dorn [49] were not included due to lack of baselines suitable for comparison.

Related Work

Wang et al. [51] analyzed code-comment coherence by means of a Bi-LSTM model which was evaluated also on the above dataset. However, because a different evaluation methodology was used, a direct comparison was not possible. Arnaoudova et al. [45], in addition to introducing the already mentioned code smell taxonomy, provide an exhaustive treatment of this subject, including also an in-depth empirical study of how developers perceive such smells. A system to detect linguistic code smells in infrastructure as code scripts was developed by Borovits et al. [52]. There exists further literature on smell detection in a broader context: for instance, Fontana et al. [53] presented a machine learning approach to code smell detection whose results, however, were called into question in a later replication study [54]. More recently, Sharma et al. [55] used deep learning models such as CNNs and LSTM networks to detect code smells in C# and Java code. Also related, is work done by Arcelli Fontana and Zanoni [56], who experiment with machine learning models for code smell *severity* prediction, which can be considered an extension of the simpler detection task.

3 BACKGROUND

This section provides an overview of the various machine learning techniques that we investigate in order to evaluate their effectiveness starting with the pre-training paradigm, then covering self-training, data augmentation, active learning, and soft labels. We also highlight the uses of pre-training in software engineering.

3.1 Pre-training with The Transformer and BERT

The Transformer architecture

Transformer [57] networks are a relatively recent architecture, particularly popular in the domain of natural language processing (NLP). They have replaced Long short-term memory (LSTM) networks as the prevailing architecture for text-based data. Transformers are based on attention, a mechanism previously used in LSTM networks to align the information flow between the encoder and decoder part of the network [58]. Attention allows a model to connect related parts of a sentence and form complex structures of interdependence between them. Unlike LSTMs, Transformer networks are not recurrent and instead have a fixed-size input window of tokens (typically 512 tokens): this allows for more efficient training and avoids vanishing-gradient or long-dependency problems extant in recurrent architectures. To summarize, each layer of the Transformer uses the attention mechanism to learn relationships between its inputs, which, in the case of the first layer are the input tokens; when used for classification, a final fully-connected layer is used as output layer.

Pre-training via Language Modelling

BERT [9] (Bi-Directional Encoder Representations from Transformers) is an extension of the Transformer architecture and comes with a specific *semi-supervised learning* training regimen: BERT heavily relies on *pre-training*, a form of *unsupervised learning*, before being *fine-tuned* on a downstream task in a classical *supervised* fashion.

During pre-training, BERT is trained on large amounts of unlabeled data via Mask Language Modelling (MLM). MLM is a prediction task where some of the input tokens are randomly replaced by blanks ("masked") and the model is trained to predict the tokens behind these blanks, taking into account the textual context on both sides of the blank (see the BERT paper for more details on the pre-training itself [9]). Intuitively, this general task is supposed to initialize the weights to a state in which certain general concepts and relationships useful for a large number of downstream tasks are already present: BERT learns a *Representation* of the tokens. Unlike word embeddings [59], these are contextual representations: they depend both on the token, and its surrounding tokens.

Of note, earlier work also used Language Modelling as a pre-training task (ELMo and ULMFit [7], [8]) with LSTMs, and were used with some varying amount of success in Software Engineering [13], [14]. BERT's pre-training is more efficient for two reasons: BERT's bidirectional architecture uses the context before and after the token, whereas LSTMs use only the context before the token; and BERT uses Byte-Pair Encoding (BPE) [60] to tokenise text in subwords rather than entire words, leading to better modelling of the vocabulary (see previous work by Karampatsis *et al.* for an extended discussion of this aspect for source code [42]).

RoBERTa [10] is a refinement of BERT, in particular relating to its pre-training regimen (e.g., RoBERTa uses a larger pre-training corpus, dynamic masking, and a variation of the pre-training task) and with only minor architectural changes (RoBERTa uses Byte-level BPE tokenization, rather than character-level BPE).

Fine-tuning

Both BERT and RoBERTa are hardly ever trained from scratch. Instead, starting from a pre-trained model with pre-initialized weights, the model weights are further *fine-tuned* by training on task-specific labeled data (called a downstream task). This involves replacing the last layer of the model (useful for the pre-training task), with a task-specific layer, and resuming training. The model can leverage the pre-trained representations to be able to learn the downstream task effectively, even with a limited amount of data, allowing BERT and RoBERTa to set the state of the art on NLP benchmarks, even on tasks with limited data (the GLUE benchmark [11] includes several task with less than 10,000 examples).

Impact of the Pre-training corpora

The standard BERT and RoBERTa models have both been pre-trained on a large English natural language corpus, with several models available in various sizes. There exist pretrained BERT models for many other natural languages and even programming languages [61]. Intuitively, one would

EN	Leppie, that's great news! I look forward to trying IronScheme!
EN →DE	Leppie, das sind großartige Neuigkeiten! Ich freue mich darauf, IronScheme auszuprobieren!
DE →EN	leppie, those are great news! I am looking forward to try out IronScheme!
EN →FR	Leppie, c'est une excellente nouvelle! J'ai hâte d'essayer IronScheme!
FR →EN	leppie, this is great news! I can't wait to try Iron-Scheme!

Fig. 1. Example of back-translation. The original English sentence is first translated to German and French, then translated back into English; resulting variation underlined. Google Translate was used for the translation.

expect a generic pre-training corpus to be a "jack of all trades, master of none", with a more specific pre-training corpus to be more suited for more specific domains (such as software engineering). There is evidence of this for word embeddings in Software Engineering [62], but how much of an impact a domain-specific pre-training corpus has for a BERT or RoBERTa model is still an open question, which we investigate. Of note, the ULMFit approach [8] continues the pre-training task on the task-specific data (without using labels), before the actual fine-tuning, finding that it does improve performance.

3.2 Additional Techniques

Intermediate-Task Fine-Tuning

Intermediate-task fine-tuning (ITT), also known as two stage fine-tuning, STILTs [63], or TANDA [64] is a technique whereby the model is fine-tuned twice (with labeled data): first on an *intermediate task*, a task different from but closely related to the target task, and finally on the actual target task (e.g., training for sentiment analysis on movies, before switching to sentiment analysis on books). This is particularly attractive whenever only little data is available for the target task whilst large amounts of data are available for a similar, possibly slightly simpler, but different intermediate task. The idea is that the target task might benefit from "knowledge" that the model acquired during intermediate-task training. Pruksachatkun *et al.* [65] presents a survey on when this method offers good prospects in NLP.

Self-Training

Self-training (also known as self-labelling or self-learning) [66], [67], is a very simple semi-supervised learning method. It can be explained as follows: A model is first trained on a (possibly too small) labeled dataset. Next, this model is used to evaluate a number of additional unlabeled samples. The model's predictions for these unlabeled samples are then simply used as their gold labels. We now have additional labeled data, albeit noisier ones; after adding it to the original dataset we retrain the model. Predictions can be filtered by confidence to reduce the probability of introducing noise into the training set.

Data Augmentation and Back-Translation

Data augmentation is a well-known technique to increase the amount of labeled data without any human labeling effort, which is especially valuable in cases where training data is in short supply. It works by adding slightly varied copies of already existing, labeled samples to the dataset, assuming the variations do not affect the label. The technique was first used in computer vision, where data augmentations are easier to define, such as flipping images horizontally (a dog looking left instead of right), or cropping images randomly (a closeup of the dog's head should still be classified as a dog). For text data, several such methods for augmentation have been proposed in recent years, among others: a) replacing words with synonyms [68], [69], b) replacing, adding or deleting words randomly [69], c) replacing words with the nearest neighbor in an embedding space [70], [71], d) replacing words with predictions from a masked language model such as BERT [72], e) translating into an intermediate language and then back into the source language (back-translation) [73].

Augmentation is typically applied at training time by simply adding the augmented samples to the training set and then proceeding as usual. Alternatively, augmentation can also be carried out at test time by aggregating (e.g., averaging) the prediction for an original test sample with the predictions for its augmented copies, thus obtaining potentially more stable or more accurate predictions.

Figure 1 shows an example of augmentation through back-translation: a sample in the dataset (here an English sentence) is translated into German and French, then back into English, causing slight variations.

Active Learning

The goal of active learning is to make the process of manual data labeling more efficient. Active learning avoids presenting samples to the rater that the model is likely to classify correctly and thus provide little new information.

Initially, a human rater labels a small number of samples, called the seed. There is also a second, larger set of yet unlabeled samples, called the *pool*. A model is first trained on the seed. In a next step, this model is used to select those samples from the pool that the model found most "difficult" to classify. "Difficulty" is measured by means of a confidence or acquisition function which calculates a confidence score from the model's prediction. In classification, this is usually a distribution over the target classes and acquisition functions are thus applied to class probabilities. Selecting samples by confidence score is called *confidence sampling*. The rater then labels the selected samples, which are then removed from the pool and added to the model's training set. This process is repeated until a satisfying number of samples have been labeled or the model reaches a particular target accuracy. A possible problem with confidence sampling is that selected samples, albeit being difficult for the model, might all be very similar, reducing the efficiency of the process. Confidence sampling is often paired with diversity sampling: selected samples are subsequently filtered for diversity, for instance using a clustering algorithm such as *k*-means. A common way to evaluate and compare active learning approaches is a simulation with an already labeled dataset. See Settles [74] for an overview of variants and extensions of active learning.

Soft Labels

In classification, usually, every sample is associated with a single target class. For many machine learning algorithms, in particular for neural networks, the target label of a sample is represented as a probability distribution over classes. While optimizations exist for handling the common single-label case, conceptually we can say that the target label is denoted by a distribution vector which assigns probability one to the class it belongs to and probability zero to all other classes. Take, for example, a classification problem where each sample belongs to either class A, B or C. A sample will have the target vector (A: 0, B: 1, C: 0) if it belongs to class B and (A: 0, B: 0, C: 1) if it belongs to class C.

The term *soft label* is used when this distribution vector is fuzzy, i.e., is not comprised of a single one and many zeros. Intuitively, this means that a sample can belong to multiple classes, with a degree expressed by the class probabilities: a target distribution such as (*A*: 0.4, *B*: 0.6, *C*:0) belongs to both, class *B* and class *A*. Most datasets do not come with soft-labels. In cases where each sample in a dataset was classified by multiple raters (which is common in order to compute inter-rater agreement), instead of using a majority vote, the rater's votes can be converted into a soft-label. Intuitively, an example in which raters disagree can be seen as more ambiguous. Providing this information to the model can help it differentiate between "easy" examples and "hard" examples.

3.3 Uses of Pre-training and Transformers in Software Engineering

In previous work, we investigated the usefulness of the pre-training paradigm (using the earlier ULMFit approach [8]), finding it promising in limited sentiment analysis experiments [13]. Mahadi *et al.* [14] experimented with cross-dataset classification of design discussions, but had mixed results.

Zhang *et al.* [15] provide a detailed study on the use of Transformers for sentiment analysis of Software Engineering artifacts, comparing existing sentiment analysis tools with Transformer models (BERT, RoBERTa, XLNet). Biswas *et al.* [17] pursue a similar avenue, training BERT on a newly compiled dataset of 4000 sentences from Stack Overflow discussions and comparing results with recurrent models.

Hey *et al.* [16] present NoRBERT, a BERT model finetuned to classify functional and non-functional requirements that achieves results competitive to state-of-the-art models on the PROMISE NFR dataset.

Keim *et al.* [75] attempt to use a standard BERT (i.e., pretrained on English natural language) for the detection of architectural tactics in Java code and report mixed results that lag behind state-of-the-art approaches in one case study.

Svyatkovskiy *et al.* [76] introduce IntelliCode compose, a system for intelligent code completion based on the GPT-2 Transformer language model.

Finally, Feng *et al.* [61] present CodeBERT, a RoBERTa-based Transformer that was trained on natural language and code (bimodal), allowing for code-related tasks that also involve natural language, such as code search or documentation generation.

4 METHODOLOGY

This section covers general aspects of the methodology, that apply to all the experiments. To ease readability, methodological details that refer to a specific technique (e.g., active learning) are described jointly with the results of this technique in Section 5.

4.1 Pre-Trained models

Off-the-shelf models

We use several "off-the-shelf" pre-trained model, which were trained on a corpus of generic English text.

- BERT-base, a 12 layer Transformer model (110 million parameters), pre-trained on a 3.3 billion words from books and Wikipedia [9].
- BERT-large, a 24 layer Transformer model (340 million parameters), pre-trained on the same corpus [9].
- DistillRoBERTa, a 6 layer Transformer model (82 million parameters), a compressed version of the larger RoBERTa model.

Domain-specific models

We pre-train a range of Transformers on data from different domains and of different sizes, to gain insights on the effectiveness of pre-training on domain-specific data. One model (StackOBERTflow) is entirely trained on domain-specific data, while the others are "off-the-shelf" models pre-trained on English, that are further pre-trained on some domain-specific data.

- BERT-reviews, a 12-layer (base) BERT model trained on 169 097 (8.4 MB) unlabeled app reviews from the AR-MINER dataset.
- BERT-comments, a 12-layer (base) BERT model trained on 487 693 (48 MB) comments extracted from well-known Java projects.
- BERT-SO-1M and BERT-SO-2M, two 12-layer (base) BERT models trained on one (147 MB) and two millions (304 MB) of Stack Overflow comments, respectively, taken from the *Stack Exchange Data Dump*¹.
- BERT-SO-1M-large, a 24-layer (large) BERT model trained on one million Stack Overflow comments (as above), used, however, only in the sentiment classification experiments.
- StackOBERTflow, a 6-layer (small) RoBERTa model trained from scratch on 26.2 million Stack Overflow comments (3.6 GB). The Stack Overflow corpus was tokenized using byte pair encoding (BPE) subwords and a large vocabulary size of 52 000.

Source code models

For the code tasks, we use the following two models:

- CodeBERTa², a small (6-layer) RoBERTa model trained on the polyglot *CodeSearchNet* [77] source code corpus and released by Hugging Face. The model supports multiple programming languages: Go, Java, Javascript, PHP, Python and Ruby.
- CodeBERT [61], a larger 12-layer model trained on the same corpus, but with a bimodal training regimen: the

model takes as input pairs of natural language and code and primarily targets code-related tasks that also involve natural language (e.g., code search or summarization). Since none of our code experiments involves natural language, we use an empty string as natural language input, except for the code-comment coherence task, where, after stripping comment markers, comments are treated as natural language.

Scratch model

To get a rough estimate of the effect of pre-training, we also train, for the sentiment classification task, a randomly initialized small RoBERTa model solely on the training set.

4.2 Preprocessing

Pre-trained models do not require extensive pre-processing, such as stemming or removing stop words. In fact, these may be harmful to performance, as the models were pretrained on data that was not pre-processed. In addition, large neural networks have enough parameter capacity to pick up on subtleties such as word order and negation. Thus for most of the datasets, in line with the practice, we did very little preprocessing. For the sentiment analysis and app review datasets, we used raw, unpreprocessed input. For the app review classification tasks we concatenated the review title and body and prepended the review's rating (a number in the range 1-5). We applied heavier preprocessing for tasks with code comment input. Here, similarly to [22], we removed newlines, comment delimiters (such as //, /*, */), stripped HTML tags and removed all punctuation except periods and question marks as well as repeated whitespace characters. Our goal was to reduce the length of the comments to fit in the Transformer's input window: each punctuation mark is treated as an additional token, taking away a spot in the window.

Preprocessing was also necessary for some source code tasks. In the code-comment coherence task, we simply took the concatenation of the lead comment and the method body as input. Moreover, we reformated all code files in the complexity prediction dataset using Google's Java code formatter³ such that, for instance, all of them use the same indentation width. For the remaining dataset and tasks no preprocessing was done.

After this, we used each pre-trained model's tokenizer to properly segment the data in the subword units specific to this model, as each model may have a different vocabulary.

4.3 Dataset Partitioning and Evaluation

We tried to replicate the baseline models' training and evaluation methodology as closely as possible. Independent of the different evaluation strategies, we repeat all of our experiments at least three times with varying random seeds and average results to reduce noise.

Sentiment classification. We trained our models on the predefined training set of the Senti4SD dataset [18], 30% of which we use for validation; the corresponding test set was used for testing. All the remaining sentiment analysis

^{1.} https://archive.org/details/stackexchange

^{2.} https://huggingface.co/huggingface/CodeBERTa-small-v1

^{3.} https://github.com/google/google-java-format

datasets were solely used as test sets (as was done in previous work). Whenever a test set lacked a neutral sentiment class, as was the case for the JIRA issues dataset, we treated *neutral* predictions from the model as negative.

Informative app reviews. A predefined train-test split is also given for the informative app reviews detection task. Here, the test set is actually larger than the training set (2000 and 1000 samples, respectively). We used 15% of the training set for validation.

App review classification. We used Monte Carlo cross-validation: we split the dataset in 10 random training and validation partitions with a ratio of 70:30. Reported results are averages over 10 runs.

SATD. We use cross-validation, with a $9 \rightarrow 1$ cross-project setting: we train on 9 out of the 10 total projects; the remaining project acts a test set.

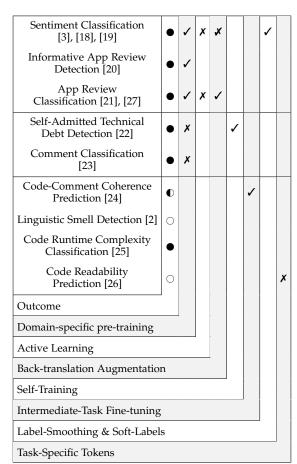
Code-comment coherence. The model in [43] was trained on 75% of the dataset while the remaining 25% were used for testing. Because this train-test split was chosen randomly, an exact comparison is not possible. To obtain more stable performance metrics, we re-evaluated this baseline model with three random train-splits and averaged the results. We train our own model in the same way, and use 10% of the training data for validation.

Linguistic code smells. As pointed out in the initial work, the leave-one-out cross validation strategy used for training the linguistic smell detection baselines is prohibitively expensive for a deep neural network. We resorted to 15-fold cross validation, putting the baselines at a slight advantage: they were trained on over 99% of the entire dataset, while our model uses only 93% of the data.

Other datasets. Finally, the code runtime complexity prediction, comment classification and code readability prediction datasets were evaluated using k-fold cross-validation with k equal to 5, 10 and 10 respectively.

4.4 Fine-Tuning and Testing

We tried several different hyper-parameters on the Senti4SD dataset, varying learning rate (2e-5, 4e-5, 5e-5), batch size (8, 12, 16, 32), and drop-out rate (0.05, 0.07, 0.1). We found that general recommendations give good results in most cases (e.g., a learning rate of 5e-5 for fine-tuning BERT). Interestingly, on Senti4SD a configuration with a relatively small batch size of 12 worked best. On the remaining tasks batch size was selected so as to fill the available GPU memory (16-48 depending on GPU and task). This means that for medium-sized models, the batch size was halved with respect to small models, such as CodeBERTa or StackOBERTflow. We also reduced the Transformers input window size from 512 to 256 tokens on datasets where input sentences where so short that the bulk of them fit this narrower window; a smaller window reduces memory consumption (and thus allows for a larger batch size) and speeds up training. Exceptionally long samples that occurred occasionally were truncated to the used window size. We stopped training after validation performance converged, which usually happened after 4-6 epochs.



○= below baseline ◆= close to baseline ◆= above baseline **X**= no clear benefit **X**= little benefit **V**= likely benefit

TABLE 2
Overview of tasks, experiments, and results in this work.

5 RESULTS AND DISCUSSION

We first start by giving an overview of the performance on each task compared to the available baseline, before diving into the details of the impact on the performance of each of the techniques that we investigate. Of note, due to limitations in the datasets and the run-time needed for each experiment, we were limited in the number of experiments we could run for each additional technique. For an overview of our results, refer to Table 2.

5.1 Comparison with Baselines

Sentiment Analysis. On the sentiment analysis datasets, Transformers are ahead of previous methods on most datasets (Table 8). In particular, this is also true for the slightly out-of-domain datasets, such as the JIRA dataset, which the Transformers were not directly trained on, with one exception: Transformers lag behind SentiStrength on the second Stack Overflow test set, but only in terms of F1 (by less than 1%), not accuracy.

App Review Analysis. In both, app review classification (Table 10) and informative app review detection (Table 3) Transformer models clearly outperform baselines. The BERT model that was further pre-trained on app reviews is

in the lead, but all of the Transformers manage to improve upon baselines. Finally, also on the CLAP dataset StackOBERTflow-comments was able to achieve a 5% higher macro F1 score over the previous random forest model proposed by Scalabrino *et al.* [27] (Table 11).

SATD. The Transformer models are able to outperform the *CNN* model by Ren *et al.* [4], which, to the best of our knowledge, represents the current state of the art (see Table 6). The dataset contains a considerable number of exact duplicates and near duplicates (those arising after preprocessing): we report results with and without removal of such duplicates; we do not know whether baseline have been trained with or without such duplicates.

Comment Classification. The Naive Bayes baseline lags behind all three Transformer models (StackOBERTflow-comments, standard BERT, and a domain-specific pretrained BERT) by a margin of 4% (Table 7). A BERT model further pre-trained on task-specific data (i.e., Java comments) performed slightly worse than standard BERT.

Code-Comment Coherence. The SVM baseline by Corazza *et al.* [43] performs better than the CodeBERTa Transformer, even when employing intermediate-task training (+1% accuracy, Table 13). In a later work, Cimasa *et al.* [44] experiment with word embeddings: the resulting baseline is weaker than their first and outperformed by the Transformer.

Linguistic Smell Detection. We compare with the baselines established by Fakhoury *et al.* [2] in Table 4. CodeBERTa is able to outperform the manually tuned SVM and the also CNN but clearly remains behind the SMO (sequential minimal optimization) model that was automatically tuned using Bayesian optimization (through Auto-Weka [78]).

Runtime Complexity Classification. The complexity classification task was the only code task where the Transformer exceeded all baselines (Table 5), including the Random Forest classifier and the SVM trained on AST embeddings.

Code Readability Prediction. The logistic regression baseline trained on manually engineered features by Scalabrino *et al.* [50] is out of reach for the Transformer: the accuracy achieved by the baseline is over 10% higher (Table 17). From all the selected tasks, the readability prediction task was the hardest for the Transformer.

Model	Face- book		Temple Run2	Swift- Key	Avg.
StackOBERTflow	90.9	89.4	88.6	85.1	88.5
BERT (base)	90.6	88.2	89.2	85.4	88.4
BERT-SO-1M	92.1	90.0	91.1	87.5	90.2
BERT-reviews	93.3	91.4	91.3	89.9	91.5
EMNB [20]	87.7	76.1	79.7	76.4	80.0

TABLE 3

Results (macro F1) for four apps in the AR-MINER dataset. Our numbers are averages over five runs with different seeds.

Pre-trained transformers were able to outperform baselines on domains closer to natural language; for source code, results were mixed.

Model	F1	Prec.	Rec.
CodeBERTa ¹	81.1	81.5	81.1
CodeBERT ¹	71.2	73.7	71.5
SMO Poly ² [2]	88.8	91.8	86.0
SVM RBF ² [2]	74.8	76.2	73.4
$CNN^2[2]$	74.5	75.6	73.5

¹ 15-fold cross validation

TABLE 4

Macro F1, precision and recall for the linguistic smell detection task. Our numbers are averages over three runs with different seeds.

Model	Acc.
CodeBERTa	78.2
CodeBERT	78.4
Random Forest [25]	74.3
Logistic Regression [25]	73.2
SVM [25]	73.0
SVM+graph2vec [25]	73.9

TABLE 5

Results for the code runtime complexity prediction task. Our numbers are means over five runs with different seeds.

5.2 Pre-Training

Our results suggest that, in particular for natural language tasks, the most promising approach seems to be to further pre-train models already pre-trained on general English. When available, in-domain data should be used for pretraining, but even close-to-domain data can yield good improvements. For instance, pre-training BERT on Stack Overflow comments helped to improve accuracy also on the app review dataset (Table 8). Further pre-training an already pre-trained model should also be preferred over pre-training from scratch: The further pre-trained models outperformed our model pre-trained from scratch for most tasks and metrics even though pre-training from scratch required considerably more training time and (unlabeled) training data. This comparison comes with a grain of salt: our further pre-trained models have twice as many layers as our pre-trained-from-scratch model, which, in turn, has a much larger vocabulary (52 000 vs 30 522). While not having a larger model pre-trained from scratch is a limitation of this work, it also highlights how expensive it is.

What speaks for our small model, and for small models in general, is of course their size: with only half the layers, training and evaluation is roughly twice as fast, the memory footprint is much smaller and, depending on the task, the performance hit may be acceptable.

Our experiments also indicate that further pre-training is effective even with relatively small amounts data. In the sentiment classification task as little as 150 MB (1 million samples) of pre-training data seems to be sufficient and able to "saturate" the model. Doubling the amount of pre-training data resulted in virtually negligible improvements (see BERT-SO-1M versus BERT-SO-2M in Table 8). Similarly, for our large model (BERT-SO-1M-large) improvements are marginal: on the Senti4SD test set, i.e., the test set that "matches" the training set, it outperforms the base-sized models by only 0.6%, while on the other test tests it lags behind them.

² leave-one-out cross validation

		without duplicates			with duplicates					
	BERT- SO-1M	BERT- comments	BERT (base)	Stack- OBERTflow	BERT- SO-1M	BERT- comments	BERT (base)	Stack- OBERTflow	CNN	NLP
Apache Ant	70.2	66.9	65.4	67.5	70.3	68.8	67.2	69.0	66.0	51.2
ArgoUML	89.8	90.0	89.7	89.3	89.4	90.0	89.7	89.0	87.8	81.9
Columba	90.9	90.4	91.0	90.0	91.4	91.4	91.9	90.6	85.2	75.0
EMF	73.5	72.4	73.2	69.2	72.5	69.5	74.4	71.0	67.9	46.2
Hibernate	88.6	87.5	88.2	87.4	88.7	88.3	88.9	87.7	82.6	76.3
JEdit	72.7	70.6	71.9	73.9	72.0	70.8	70.5	72.2	59.9	46.1
JFreeChart	77.7	80.7	79.2	77.6	64.1	64.7	63.8	63.0	73.9	51.3
JMeter	87.0	87.5	87.4	86.6	86.4	85.8	86.2	84.3	82.8	71.5
JRuby	91.1	91.2	90.8	90.5	92.1	92.4	92.3	91.4	86.3	77.3
SQuirrel	79.1	78.2	78.8	78.6	80.5	79.5	80.5	78.4	73.9	59.3
Average	82.1	81.5	81.6	81.1	80.7	80.1	80.5	79.7	76.6	63.6

TABLE 6

Macro F1 scores for the SATD detection task. As far as our results are concerned, numbers are means over five runs, each with different seed.

Model	F1	Prec.	Rec.
StackOBERTflow	88.4	89.6	90.1
BERT-comments	88.3	90.6	89.0
BERT	88.4	90.5	89.1
Naive Bayes Multinomial [23]	84.3	82.0	87.2

TABLE 7

Results for the comment classification task. Task-specific pre-training failed to improve performance. As far as our models are concerned, results are averages over three runs with different seeds.

Interestingly, our BERT-comments model, a general English model further pre-trained on Java comments, performs slightly worse than the same model without this task-specific pre-training (i.e., a standard BERT) on both datasets it was applied to (Tables 7 and 6). As to why this is the case we can only *speculate*: A possible explanation is that the comments in our pre-training dataset are very repetitive and have low linguistic diversity (e.g., Java docstrings). Thus, the model might have unlearned some of its general language capabilities during task-specific pre-training.

Figure 2 demonstrates that pre-training is essential: a randomly initialized model not only converges much slower, it also has higher variance and typically reaches much lower peak performance. In sum, our experiments show that there is very little reason not to use an already pre-trained, general natural language model as the basis for further domain-specific pre-training and should in most cases be preferred over pre-training from scratch, which, in relation to training time, hardly seems worth the effort.

Pre-training is essential. Further pre-training a generic model on domain-specific data is often beneficial, and is much more effective than pre-training from scratch.

5.3 Soft Labels

Since Calefato *et al.* [18] released multi-rater labels (three per sample) along with the majority label, we conducted a soft label experiment. For instance, if one voter assigned the *positive* label to a sentence, while two raters assigned *neutral*, majority voting would label it as *neutral*. Instead,

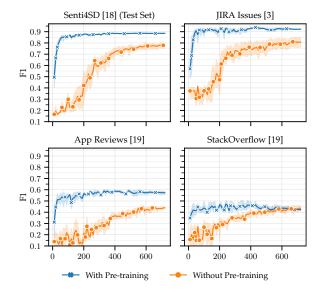


Fig. 2. F1 score on different sentiment classification datasets with and without pre-training. Number of optimization steps is shown on the x-axis; error bands are 95% confidence intervals.

the soft label captures all three rater labels, assigning the distribution: (positive: $0.\overline{33}$, negative: 0, neutral: $0.\overline{66}$) to the sentence.

We train a subset of the models with hard-labels and softlabels on Senti4SD, and evaluate on all sentiment analysis datasets (see Table 9). The results look rather promising: on the Senti4SD test set, training with all three rater labels resulted in an increase of 0.5%. On the JIRA test set, soft labels yielded an improvement of 1%. On the other hand, performance dropped on the second Stack Overflow test set. While these improvements are not certain, and might seem modest, they come almost for free. Whenever multirater labels are available, we recommend to tentatively use them in this way and encourage creators of datasets to also release labels of individual raters.

When available, individual rater labels may improve performance at very low cost.

Dataset	Model	Acc.	F1
	BERT (base)	87.9	87.8
	BERT-SO-1M-large	89.4	89.4
	BERT-SO-2M	88.8	88.7
	BERT-SO-1M	88.8	88.7
Senti4SD [18]	DistilRoBERTa (small)	87.4	87.3
(Test Set)	RoBERTa (small, no pretr.)	78.4	77.8
(lest set)	Senti4SD [18]	-	86.0
	SentiCR [18]	-	82.0
	SentiStrengthSE [18]	-	80.0
	SentiStrength [18]	-	84.0
	StackOBERTflow	88.7	88.6
	BERT (base)	64.9	50.2
	BERT-SO-1M-large	66.0	51.9
	BERT-SO-2M	67.7	53.1
	BERT-SO-1M	68.3	53.7
	DistilRoBERTa (small)	60.9	48.5
. D : [40]	NLTK [3]	54.0	40.8
App Reviews [19]	RoBERTa (small, no pretr.)	48.0	42.2
	SentiStrength+AC0-SE [3]	58.9	46.6
	SentiStrength [3]	62.5	48.2
	StackOBERTflow	72.0	57.4
	Stanford CoreNLP SO [3]	41.6	35.5
	Stanford CoreNLP [3]	69.5	56.0
	BERT (base)	95.7	95.0
	BERT-SO-1M-large	94.8	93.8
	BERT-SO-2M	95.2	94.5
	BERT-SO-1M	95.1	94.2
	DistilRoBERTa (small)	93.7	92.7
	NLTK [3]	29.8	46.5
JIRA Issues [3]	RoBERTa (small, no pretr.)	84.0	80.1
	SentiStrength+AC0-SE [3]	76.0	87.0
	SentiStrength [3]	77.1	85.4
	StackOBERTflow	93.5	92.5
	Stanford CoreNLP SO [3]	36.0	44.2
	Stanford CoreNLP [3]	67.6	73.7
	BERT (base)	79.9	47.6
	BERT-SO-1M-large	79.2	43.3
	BERT-SO-2M	80.3	48.6
	BERT-SO-1M		47.9
		80.1	
	DistilRoBERTa (small) NLTK [3]	78.8 77.9	44.3 43.2
StackOverflow [19]			43.2
omeno (emon [15]	RoBERTa (small, no pretr.)	75.5	
	SentiStrength+AC0-SE [3]	78.0	46.8
	SentiStrength [3]	69.5	49.5
	StackOBERTflow	79.3	44.1
	Stanford CoreNLP SO [3]	75.9	47.5
	Stanford CoreNLP [3]	40.3	35.5

TABLE 8

Accuracy, macro F1 and per-class precision and recall for different models and datasets. Values reported are means over five runs, each with different seed (only our models). All models were trained on the Senti4SD [18] Stack Overflow dataset.

Dataset	Label Type	F1		
Dataset	Laber Type	$\mu \pm \sigma$	max	
Amm Darriarus [10]	All Label Votes	57.5 ± 2.2	59.1	
App Reviews [19]	Majority Label		59.7	
[IRA Issues [3]	All Label Votes	93.5 ± 1.0	94.7	
JIKA Issues [5]	Majority Label		93.2	
StackOverflow [19]	All Label Votes		43.8	
StackOvernow [19]	Majority Label		46.4	
Senti4SD [18]	All Label Votes	89.1 \pm 0.5	89.6	
Setti 43D [16]	Majority Label	88.6 ± 0.5	89.1	

TABLE 9

Macro F1 for different label types and datasets: mean, maximum and standard deviation over five runs, each with different seed. All models were trained on the Senti4SD [18] Stack Overflow dataset.

5.4 Back-Translation Augmentation

We performed back-translation experiments on the sentiment and app review classification tasks by translating the entire datasets into French, German and Russian using Google Translate and from these languages back into English (see Figure 1 for an example). For the CLAP and Senti4SD datasets we do training-time and test-time augmentation, both, separately and combined, using the Stack-OBERTflow model (Table 11). On the other app review dataset [21] we do training-time augmentation alone, and combine it with test-time augmentation: here the the experiment is carried out on several different models (Table 10).

	Bug reports	Feature request	Ratings	User experience	Avg.
StackOBERTflow	61.1	42.0	79.4	49.8	58.1
StackOBERTflow+BT	62.5	43.3	79.0	50.2	58.7
StackOBERTflow+BTT	62.7	44.3	79.6	50.6	59.3
BERT (base)	61.8	39.7	79.4	49.4	57.6
BERT (base)+BT	62.0	43.1	79.4	50.2	58.7
BERT (base)+BTT	62.5	42.6	80.0	50.1	58.8
BERT-reviews	64.1	44.7	80.2	51.0	60.0
BERT-reviews+BT	63.2	43.9	79.5	51.0	59.4
BERT-reviews+BTT	63.7	46.7	80.1	51.2	60.4
Decision Tree ¹ [21]	62.0	42.0	54.0	50.0	52.0
Naive Bayes ¹ [21]	62.0	47.0	54.0	53.0	54.0

¹ multiclass, bag of words + metadata

TABLE 10

F1 scores for app review classification, reported with training-time back-translation augmentation (+BT), with training-time and test-time back-translation augmentation (+BTT) and without any augmentation. Results are averages over three runs with different seeds.

Dataset	Augmentation	Acc.	F1	Prec.	Rec.
	Train+Test	89.2	81.4	84.2	82.3
CLAP	None	88.1	79.7	80.6	81.7
CLAP	Test	88.3	80.2	81.8	81.5
	Train	88.9	81.1	83.7	82.0
	Train+Test	88.7	88.6	88.7	88.5
Senti4SD	None	88.4	88.2	88.2	88.4
	Test	88.2	88.0	88.1	88.0
	Train	88.4	88.3	88.3	88.5

TABLE 11

Macro F1, precision and recall for back-translation augmentation at training and test time on the Senti4SD (sentiment classification) and CLAP (app review classification) datasets using the StackOBERTflow model. Results are averages over three runs with different seeds.

Back-translation augmentation led to a clear increase in F1 and accuracy on the CLAP app review dataset, in particular when training and testing time augmentation were combined (+1.1% accuracy, Table 11). On Senti4SD, data augmentation yields modest improvements (+0.3%); in fact, augmenting at test time only caused a slight drop in performance. Table 10 suggests that the effect of back-translation augmentation depends on the model and pre-training choice. With train-time augmentation only, we see a modest increase of 0.6% in F1 for our small StackOBERTflow model and 1.1% on a general BERT model, while BERT-reviews shows better performance without augmentation. The latter does however benefit from combined augmentation (+0.4%). However, the question of whether in general

task-specific pre-training diminishes the effects of data augmentation cannot be answered given this limited data and would require further experiments.

When possible, back translation yields improvements, particularly if used at both training and test time.

5.5 Active Learning

$$C_{LC}(x) = \frac{n}{n-1} (1 - P_{\theta}(y_1^*|x))$$

$$C_{MC}(x) = 1 - (P_{\theta}(y_1^*|x) - P_{\theta}(y_2^*|x))$$

$$C_{RC}(x) = \frac{P_{\theta}(y_1^*|x)}{P_{\theta}(y_2^*|x)}$$

$$C_E(x) = -\frac{1}{\log_2(n)} \sum_{i=1}^n P_{\theta}(y_i|x) \log_2(P_{\theta}(y_i|x))$$

$$C_{rand}(x) = rand([0, 1])$$

TABLE 12

Acquisition functions used in our active learning experiment, adapted from [79]: least confidence (C_{LC}) , margin of confidence (C_{MC}) , ratio of confidence (C_{RC}) , entropy (C_E) and random confidence (C_{rand}) . y_1^* and y_2^* are the classes with highest and second highest probabilities, respectively; n is the number of classes. All functions have range [0,1].

We try active learning on the Senti4SD sentiment analysis dataset and an app review dataset. In both cases we compare several acquisition functions (Table 12). We carry out the experiment as follows: initially we split the *training* set into the seed set \mathcal{D}_{seed} containing 5% of all samples and the pool set \mathcal{D}_{pool} , containing the remaining samples. We let $\mathcal{D}_{train} := \mathcal{D}_{seed}$ and train the model. Then we evaluate \mathcal{D}_{pool} as well as the test set on this model. Next, for each $x \in \mathcal{D}_{pool}$ we calculate a confidence score by applying the acquisition function (from Table 12). After that, we let \mathcal{D}_{top} be the k=180 samples with the highest confidence score. We remove these samples from \mathcal{D}_{pool} and add them to \mathcal{D}_{train} This procedure is repeated until \mathcal{D}_{pool} is empty.

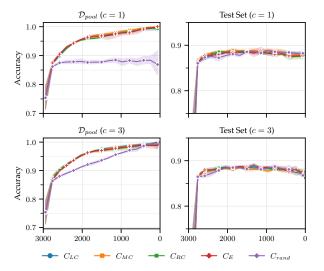


Fig. 3. Accuracy of \mathcal{D}_{pool} and the test set at each iteration of the active learning process for different acquisition functions with (c=3) and without (c=1) diversity sampling for the sentiment classification task (Senti4SD). Error bands are 95% confidence intervals. Same plot for the review classification task can be found in the appendix (Figure 5)

We also combine confidence sampling with diversity sampling to avoid introducing similar samples into the pool. Instead of k samples, we select $c \cdot k$ samples from the pool, where c determines the cluster size. We use the k-means algorithm to cluster the $c \cdot k$ samples into k clusters, each of size c. We then select a single sample from each cluster, for a total of k samples. Then, we proceed as above.

Figure 3 shows evaluation results at each iteration step with a cluster size of c=3 for the sentiment classification task. The outcome of our active learning experiments remained behind expectations: in both tasks, neither confidence sampling alone nor confidence and diversity sampling combined showed an appreciable advantage over the random baseline. The choice of acquisition function did not seem crucial, but a more systematic study would be needed to draw more solid conclusions. On the other hand, the plot of pool accuracy (top left) indicates that the active learning process worked as expected: a random acquisition function without diversity sampling had constant performance.

Active learning simulations did not improve results.

5.6 Self-Training

We investigate the use of self-training for the SATD detection task. We extract 350000 comments from various popular Java libraries and frameworks. Then we train a classifier model on the entire *original* dataset, which we use to classify the 350000 comments as either *technical debt* or *not technical debt*. We only keep the 7904 positive comments, i.e., those classified as *technical debt*, and discard all other samples to avoid increasing the class imbalance already extant in the original dataset. For each positive sample we calculate a confidence score using C_{LC} (Table 12). We take the top 5%, and 80% most confidently classified comments, equal to 6092 and 7880 additional comments, respectively, and add them to the original training set. Finally, the model is trained and evaluated on this extended training set.

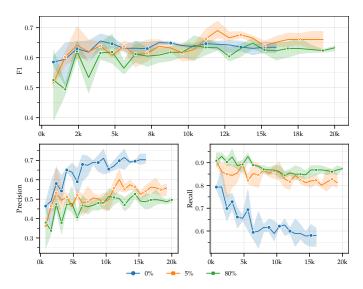


Fig. 4. Precision, recall and F1 under different self-training settings for the *Apache Ant* project. Error bands are 95% confidence intervals.

Figure 4 shows that self-training increases recall (and possibly F1) but causes precision to drop. This is, of course, not surprising: the added samples increase dataset variance which likely explains an increased recall. Similarly, the precision drop can be explained by the lower quality self-training labels. Thus, one can tune the precision-recall trade-off according to task-specific needs, such as when a recall is more important than precision, or the model's precision is high enough to be partly sacrificed for better recall. Figure 4 shows different self-training settings for *Apache Ant*: F1 score went up 4% (precision: -1%, recall +10%), when using a 5% confidence threshold. The change in F1 strongly depends on the confidence threshold and varies across projects: *EMF* sees an 8% F1 drop (precision: -28%, recall: +15%); other projects range from -1% to +4%.

Self-training increases recall at the expense of precision, but the confidence threshold should be tuned carefully.

5.7 Intermediate-Task Training

We evaluate Intermediate-Task Training (ITT) on the code-comment coherence task, as it is the only setting for which we could define such an intermediate task. We use 38 000 Java methods along with their lead comments from the *CodeSearchNet* [77] dataset. We assign half of the methods to their actual lead comments (assumed to be coherent) and shuffle the other half randomly (thus assumed to be incoherent). The model is then fine-tuned on the intermediate task of detecting whether a method was paired with its true lead comment or a random one. Finally, we fine-tune on the code-comment coherence dataset as usual.

Acc.	AUC
81.4	80.1
82.5	80.9
81.6	80.0
83.5	81.3
83.3	82.2
80.5	-
	81.4 82.5 81.6 83.5 83.3

¹ Single seed; value reported in [43]

TABLE 13

Results on the comment-code coherence dataset. CodeBERTa was evaluated with and without intermediate-task training (ITT). Our numbers are averages over five runs with different seeds.

ITT improved the performance of the Transformer model on the code-comment coherence task (Table 13): we observed a modest rise in AUC of 0.8%, and a slightly higher increase in accuracy (+1.1%). If an appropriate intermediate task for the task at hand can be found, ITT can be done with relatively little effort: in our case, the training procedure for the intermediate-task was mostly identical to the one for the target task; the bulk of the work consisted in generating the intermediate-task dataset (e.g. selecting, shuffling and preprocessing the data from *CodeSearchNet*).

When applicable, ITT may improve performance; however finding a suitable intermediate task may be difficult.

6 Discussion

6.1 Summary and implications of the results

Types of datasets. Overall, we see that Transformers work very well for natural language datasets, but that performance on source code is "hit or miss". This comes with caveats: the models used so far are multilingual, which might reduce performance. They are also trained with less computational resources, and on an order of magnitude less data: While CodeSearchNet is around 1.7 GB, the training data for BERT is around 16GB, while it is 160GB for RoBERTa (CodeSearchNet). Models where also small (6 layers), or had a dual input (text and code). We would expect a Java-specific model trained on a similar size corpus as BERT to perform better. Moreover, source code is quite different from natural language: code snippets are often larger than sentences, and much more structured. This might pose limits on what an unstructured model might achieve. Recent adaptations of the Transformer architecture (e.g. [80], [81]) allow the model to better make use of the tree-like structure of code. Investigating such code-specific architectures in connection with small datasets remains an issue for future research.

Domain-specific pre-training. proved effective in natural language settings, improving performance at a moderate cost in terms of computation and data. The only case where it did not work well was for code comments. While we are not sure why, one reason could be that code comments are too far way from regular English (needing a specific model instead), or that careful curation of the data set (avoiding too many duplicates) is needed. Both cases could lead to catastrophic forgetting [82] of the initial pretraining. Leveraging the resources that were used to train BERT and fine-tuning it further proved much more effective than training a model from scratch. We have not evaluated domain-specific training from English to source code, as we hypothesized that the two domains are very different the tokenization alone might differ significantly [42]. This intuition is supported by the literature, which reports an example where an English BERT was applied to source code, with underwhelming results [75].

Back-translation. While data augmentation is effective for natural language, it is not immediately applicable to source code. Source code can not be "back translated" easily. Specific data augmentations for code should be investigated, but may not be trivial (e.g., renaming identifiers could be investigated, but what should be done with API methods?).

Intermediate Task Training. Another alternative is to define suitable intermediate training tasks. We have found initial evidence of this, and a recent paper adds further evidence, in the context of traceability [83]. However, it used a very similar task and dataset. Thus, the challenge here is not whether intermediate task training helps, but rather *whether a suitable task exists* for a given problem.

Soft Labels. Soft labels that reflect the uncertainty of raters (and thus the difficulty of the samples) can be useful as well, and at a minimal cost. However these are not common, as of now. We call on dataset builders to release them alongside the majority label, as was done by Calefato [18].

² Average over 5 seeds; reproduction

6.2 Limitations of this work

Limited Number of experiments. While we try to report results as extensively as possible to increase their generalisability, we are limited for two main reasons: 1) we have limited computational resources, and 2) some techniques are specific to some settings.

Limited resources. Deep learning is famously resource intensive. While fine-tuning is less resource intensive than training models from scratch, it still requires significant time on one or more dedicated GPUs, particularly for larger models. A single run is measured in hours. This limits the number of experiments, particularly as we repeat experiments several times with different random seeds.

Results in specific settings. While resources are limited, we still wanted to try each technique on at least two datasets. However, some techniques were applied to a single dataset. For soft labels, we needed multiple ratings: only a single sentiment analysis dataset had the required three ratings per sample. We could only define a reasonable intermediate task for code-comment readability prediction. We considered using self-training for comment classification, but did not, due to the large number of imbalanced classes.

Hyper-parameters. Limited resources also impact the extent to which we perform hyper-parameter optimisation, as thorough parameter searches (whether by grid, random or bayesian methods) would be prohibitively expensive. A second limit is that some hyper-parameters are fixed by the usage of a pre-trained model (e.g. number of layers, number of attention heads, embedding size, vocabulary size). A silver lining is that, given the interest in pre-trained models, general recommendations for hyper-parameters exist and are broadly applicable. Thus, we started with these recommendations, and investigated some variations of the hyper-parameters on the Senti4SD dataset, confirming that the recommendations worked well. We then applied those hyper-parameters on other experiments, varying only the most important ones in some cases (learning rate, batch size). Cross-validation also makes evaluation and hyperparameter tuning more complex and resource intensive. Since we limited hyper-parameter tuning, we are not at risk of overfitting to the test fold when doing cross validation. An alternative would be to use doubly nested cross validation, but this further increases the resource needed. We note that dedicated test sets ease this considerably.

Comparisons with previous work. We do our best to provide a fair comparison with previous work, while avoiding methodological issue (e.g. averaging seeds). We do not always exactly know how previous work was evaluated evaluation (e.g., hyper-parameter selection strategy, whether simple or nested cross-validation was used, or whether some data points were excluded) as code is not always released. In some other cases, other factor presents us to make an exact comparison (e.g. use of leave-one-out cross validation is not practical for our setting). To alleviate this in the future, we release our source code (see Appendix C).

Active Learning. While we could not see an advantage to active learning, this is not in line with previous work by Dhinakaran *et al.* [35] and Tu *et al.* [84].Of note, our results are obtained through simulation based on an existing labelled

dataset. While this is a practice often used to evaluate active learning methods, a realistic application of active learning on a larger set of unlabelled data would lead to a different training set, which may be substantially more varied, and thus more effective. But it is also possible that the impact of active learning is less visible when pre-training is used.

Random Seeds. Dodge *et al.* [85] found that the choice of the random seed can have a substantial impact on performance, especially for small datasets. We ran most of our experiments five times and all of them at least three times, with different seeds. While this surely mitigates the problem, it might not fully clear it up.

Implementation Bugs. Our implementations are based on Hugging Face's transformers Python package [86], a high quality implementation of common Transformer models. However, despite careful reviews we cannot fully preclude errors in our own code and adaptations.

7 CONCLUSIONS

Software Engineering datasets are often small, by necessity. In this work, we trained various Transformer models on 13 small and medium-sized dataset selected from the recent Software Engineering literature. We not only compared Transformers of different size and different pre-training regimes but also applied several machine learning techniques that promised a possible benefit for small datasets. These techniques were data augmentation, self-training, intermediate-task training, active learning and soft labels.

Overall, we found that on natural language tasks, Transformers usually outperform existing baselines. On source code tasks, however, results were mixed. Significant work lies ahead to define effective pre-trained source code models either by training larger models on more data, or by incorporating more structural information during training.

In general, we advise *against* pre-training a new model from scratch as it is extremely resource intensive, for mixed results. Instead, an already pre-trained model can be *further pre-trained* on task-specific data. If such task-specific data is unavailable, training on close-to-domain data is worth a try. We provide several such pre-trained models in Appendix C.

Several additional techniques were useful at a relatively low cost. We particularly recommend the use of soft labels derived from multi-rater labels if available, and call on dataset authors to release these multi-rater labels. Backtranslation is similarly useful, if more expensive. It is unfortunately not easily applicable to source code.

Other techniques were less applicable. We find that self-training is advisable only in cases where the user wants to boost recall and is willing to sacrifice precision. If circumstances allow it, intermediate task training seems promising, but it seems rarely applicable, and has a much higher cost. Finally, our active learning experiments were inconclusive; a wider study on a larger set of dataset might be required to draw a clearer picture.

While these general guidelines are useful on their own, their applicability is limited. To this extent, we release all the scripts and pretrained models that were built as part of this work, so that the community can easily fine-tune the models on their own Software Engineering datasets, and apply additional techniques as they see fit (see Appendix C).

REFERENCES

- [1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in 2009 IEEE conference on computer vision and pattern recognition, Ieee, 2009, pp. 248–255.
- [2] S. Fakhoury, V. Arnaoudova, C. Noiseux, F. Khomh, and G. Antoniol, "Keep it simple: Is deep learning good for linguistic smell detection?" In 2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER), 2018, pp. 602–611. DOI: 10. 1109/SANER.2018.8330265.
- [3] B. Lin, F. Zampetti, G. Bavota, M. Di Penta, M. Lanza, and R. Oliveto, "Sentiment analysis for software engineering: How far can we go?" In 2018 IEEE/ACM 40th International Conference on Software Engineering (ICSE), 2018, pp. 94–104.
- [4] X. Ren, Z. Xing, X. Xia, D. Lo, X. Wang, and J. Grundy, "Neural network-based detection of self-admitted technical debt: From performance to explainability," ACM Transactions on Software Engineering and Methodology (TOSEM), vol. 28, no. 3, pp. 1–45, 2019.
- [5] R. M. Santos, M. C. R. Junior, and M. G. de Mendonça Neto, "Self-admitted technical debt classification using lstm neural network," in 17th International Conference on Information Technology—New Generations (ITNG 2020), Springer, 2020, pp. 679–685.
- [6] X. Zhu and A. B. Goldberg, "Introduction to semisupervised learning," Synthesis lectures on artificial intelligence and machine learning, vol. 3, no. 1, pp. 1–130, 2009.
- [7] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," *arXiv preprint* arXiv:1802.05365, 2018.
- [8] J. Howard and S. Ruder, "Universal language model fine-tuning for text classification," arXiv preprint arXiv:1801.06146, 2018.
- [9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, 2018. arXiv: 1810.04805 [cs.CL].
- [10] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, *Roberta: A robustly optimized bert pretraining approach*, 2019. arXiv: 1907.11692 [cs.CL].
- [11] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, *Glue: A multi-task benchmark and analysis platform for natural language understanding*, 2019. arXiv: 1804.07461 [cs.CL].
- [12] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, Squad: 100,000+ questions for machine comprehension of text, 2016. arXiv: 1606.05250 [cs.CL].
- [13] R. Robbes and A. Janes, "Leveraging small software engineering data sets with pre-trained neural networks," in 2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER), IEEE, 2019, pp. 29–32.
- [14] A. Mahadi, K. Tongay, and N. A. Ernst, "Cross-dataset design discussion mining," in 2020 IEEE 27th Inter-

- national Conference on Software Analysis, Evolution and Reengineering (SANER), IEEE, 2020, pp. 149–160.
- [15] T. Zhang, B. Xu, F. Thung, S. A. Haryono, D. Lo, and L. Jiang, "Sentiment analysis for software engineering: How far can pre-trained transformer models go?" In 2020 IEEE International Conference on Software Maintenance and Evolution (ICSME), IEEE, pp. 70–80.
- [16] T. Hey, J. Keim, A. Koziolek, and W. F. Tichy, "Norbert: Transfer learning for requirements classification," in 2020 IEEE 28th International Requirements Engineering Conference (RE), IEEE, pp. 169–179.
- [17] E. Biswas, M. E. Karabulut, L. Pollock, and K. Vijay-Shanker, "Achieving reliable sentiment analysis in the software engineering domain using bert," in 2020 IEEE International Conference on Software Maintenance and Evolution (ICSME), 2020, pp. 162–173. DOI: 10.1109/ICSME46990.2020.00025.
- [18] F. Calefato, F. Lanubile, F. Maiorano, and N. Novielli, "Sentiment polarity detection for software development," *Empirical Software Engineering*, vol. 23, no. 3, pp. 1352–1382, 2018.
- [19] B. Lin, F. Zampetti, R. Oliveto, M. Di Penta, M. Lanza, and G. Bavota, "Two datasets for sentiment analysis in software engineering," in 2018 IEEE International Conference on Software Maintenance and Evolution (ICSME), 2018, pp. 712–712.
- [20] N. Chen, J. Lin, S. C. Hoi, X. Xiao, and B. Zhang, "Ar-miner: Mining informative reviews for developers from mobile app marketplace," in *Proceedings of* the 36th international conference on software engineering, 2014, pp. 767–778.
- [21] W. Maalej, Z. Kurtanović, H. Nabil, and C. Stanik, "On the automatic classification of app reviews," *Requirements Engineering*, vol. 21, no. 3, pp. 311–331, 2016.
- [22] E. d. S. Maldonado, E. Shihab, and N. Tsantalis, "Using natural language processing to automatically detect self-admitted technical debt," *IEEE Transactions* on Software Engineering, vol. 43, no. 11, pp. 1044–1062, 2017.
- [23] L. Pascarella and A. Bacchelli, "Classifying code comments in java open-source software systems," in 2017 IEEE/ACM 14th International Conference on Mining Software Repositories (MSR), IEEE, 2017, pp. 227–237.
- [24] A. Corazza, V. Maggio, B. Kessler, and G. Scanniello, "A new dataset for source code comment coherence," *CLiC it*, p. 100, 2016.
- [25] J. Sikka, K. Satya, Y. Kumar, S. Uppal, R. R. Shah, and R. Zimmermann, "Learning based methods for code runtime complexity prediction," in *Advances in Information Retrieval*, J. M. Jose, E. Yilmaz, J. Magalhães, P. Castells, N. Ferro, M. J. Silva, and F. Martins, Eds., Cham: Springer International Publishing, 2020, pp. 313–325, ISBN: 978-3-030-45439-5.
- [26] S. Scalabrino, M. Linares-Vasquez, D. Poshyvanyk, and R. Oliveto, "Improving code readability models with textual features," in 2016 IEEE 24th International Conference on Program Comprehension (ICPC), IEEE, 2016, pp. 1–10.
- [27] S. Scalabrino, G. Bavota, B. Russo, M. D. Penta, and R. Oliveto, "Listening to the crowd for the release planning of mobile apps," *IEEE Transactions on Soft-*

- ware Engineering, vol. 45, no. 1, pp. 68–86, 2019. DOI: 10.1109/TSE.2017.2759112.
- [28] M. Ortu, A. Murgia, G. Destefanis, P. Tourani, R. Tonelli, M. Marchesi, and B. Adams, "The emotional side of software developers in jira," in *Proceedings of the 13th International Conference on Mining Software Repositories*, ser. MSR '16, Austin, Texas: Association for Computing Machinery, 2016, pp. 480–483, ISBN: 9781450341868. DOI: 10.1145/2901739.2903505. [Online]. Available: https://doi.org/10.1145/2901739.2903505.
- [29] L. Villarroel, G. Bavota, B. Russo, R. Oliveto, and M. Di Penta, "Release planning of mobile apps based on user reviews," in *Proceedings of the 38th International Conference on Software Engineering*, ser. ICSE '16, Austin, Texas: Association for Computing Machinery, 2016, pp. 14–24, ISBN: 9781450339001. DOI: 10.1145/ 2884781.2884818. [Online]. Available: https://doiorg.libproxy.unibz.it/10.1145/2884781.2884818.
- [30] M. V. Mäntylä, D. Graziotin, and M. Kuutila, "The evolution of sentiment analysis—a review of research topics, venues, and top cited papers," *Computer Science Review*, vol. 27, pp. 16–32, 2018, ISSN: 1574-0137. DOI: https://doi.org/10.1016/j.cosrev.2017.10.002. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S1574013717300606.
- [31] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, "Xlnet: Generalized autoregressive pretraining for language understanding," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32, Curran Associates, Inc., 2019, pp. 5753–5763. [Online]. Available: https://proceedings.neurips.cc/paper/2019/file/dc6a7e655d7e5840e66733e9ee67cc69-Paper.pdf.
- [32] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, *Albert: A lite bert for self-supervised learning of language representations*, 2020. arXiv: 1909. 11942 [cs.CL].
- [33] T. Ahmed, A. Bosu, A. Iqbal, and S. Rahimi, "Senticr: A customized sentiment analysis tool for code review interactions," in 32nd IEEE/ACM International Conference on Automated Software Engineering (NIER track), ser. ASE '17, 2017.
- [34] Z. Chen, Y. Cao, H. Yao, X. Lu, X. Peng, H. Mei, and X. Liu, "Emoji-powered sentiment and emotion detection from software developers' communication data," *ACM Trans. Softw. Eng. Methodol.*, vol. 30, no. 2, Jan. 2021, ISSN: 1049-331X. DOI: 10.1145/3424308. [Online]. Available: https://doi-org.libproxy.unibz.it/10.1145/3424308.
- [35] V. T. Dhinakaran, R. Pulle, N. Ajmeri, and P. K. Murukannaiah, "App review analysis via active learning: Reducing supervision effort without compromising classification accuracy," in 2018 IEEE 26th International Requirements Engineering Conference (RE), 2018, pp. 170–181. DOI: 10.1109/RE.2018.00026.
- [36] N. S. R. Alves, L. F. Ribeiro, V. Caires, T. S. Mendes, and R. O. Spínola, "Towards an ontology of terms on technical debt," in 2014 Sixth International Workshop on Managing Technical Debt, 2014, pp. 1–7.

- [37] G. Sierra, E. Shihab, and Y. Kamei, "A survey of self-admitted technical debt," *Journal of Systems and Software*, vol. 152, pp. 70–82, 2019, ISSN: 0164-1212. DOI: https://doi.org/10.1016/j.jss.2019.02.056. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S0164121219300457.
- [38] Q. Huang, E. Shihab, X. Xia, D. Lo, and S. Li, "Identifying self-admitted technical debt in open source projects using text mining," *Empirical Software Engineering*, vol. 23, no. 1, pp. 418–451, 2018.
- [39] Z. Liu, Q. Huang, X. Xia, E. Shihab, D. Lo, and S. Li, "Satd detector: A text-mining-based self-admitted technical debt detection tool," in 2018 IEEE/ACM 40th International Conference on Software Engineering: Companion (ICSE-Companion), 2018, pp. 9–12.
- [40] L. Pascarella, "Classifying code comments in java mobile applications," in 2018 IEEE/ACM 5th International Conference on Mobile Software Engineering and Systems (MOBILESoft), 2018, pp. 39–40.
- [41] A. Hindle, E. T. Barr, M. Gabel, Z. Su, and P. Devanbu, "On the naturalness of software," *Communications of the ACM*, vol. 59, no. 5, pp. 122–131, 2016.
- [42] R.-M. Karampatsis, H. Babii, R. Robbes, C. Sutton, and A. Janes, "Big code!= big vocabulary: Open-vocabulary models for source code," in 2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE), IEEE, 2020, pp. 1073–1085.
- [43] A. Corazza, V. Maggio, and G. Scanniello, "Coherence of comments and method implementations: A dataset and an empirical investigation," *Software Quality Journal*, vol. 26, no. 2, pp. 751–777, Jun. 2018, ISSN: 0963-9314. DOI: 10.1007 / s11219 016 9347 1. [Online]. Available: https://doi.org/10.1007/s11219 016-9347-1.
- [44] A. Cimasa, A. Corazza, C. Coviello, and G. Scanniello, "Word embeddings for comment coherence," in 2019 45th Euromicro Conference on Software Engineering and Advanced Applications (SEAA), 2019, pp. 244–251.
- [45] V. Arnaoudova, M. Di Penta, and G. Antoniol, "Linguistic antipatterns: What they are and how developers perceive them," *Empirical Software Engineering*, vol. 21, no. 1, pp. 104–158, 2016.
- [46] A. Narayanan, M. Chandramohan, R. Venkatesan, L. Chen, Y. Liu, and S. Jaiswal, *Graph2vec: Learning distributed representations of graphs*, 2017. arXiv: 1707. 05005 [cs.AI].
- [47] R. P. Buse and W. R. Weimer, "Learning a metric for code readability," *IEEE Transactions on Software Engineering*, vol. 36, no. 4, pp. 546–558, 2009.
- [48] D. Posnett, A. Hindle, and P. Devanbu, "A simpler model of software readability," in *Proceedings of the 8th Working Conference on Mining Software Repositories*, ser. MSR '11, Waikiki, Honolulu, HI, USA: Association for Computing Machinery, 2011, pp. 73–82, ISBN: 9781450305747. DOI: 10.1145/1985441.1985454. [Online]. Available: https://doi-org.libproxy.unibz.it/10. 1145/1985441.1985454.
- [49] J. Dorn, "A general software readability model," 2012.
- [50] S. Scalabrino, M. Linares-Vásquez, R. Oliveto, and D. Poshyvanyk, "A comprehensive model for code

- readability," Journal of Software: Evolution and Process, vol. 30, no. 6, e1958, 2018.
- [51] D. Wang, Y. Guo, W. Dong, Z. Wang, H. Liu, and S. Li, "Deep code-comment understanding and assessment," *IEEE Access*, vol. 7, pp. 174200–174209, 2019. DOI: 10.1109/ACCESS.2019.2957424.
- [52] N. Borovits, I. Kumara, P. Krishnan, S. D. Palma, D. Di Nucci, F. Palomba, D. A. Tamburri, and W.-J. van den Heuvel, "Deepiac: Deep learning-based linguistic anti-pattern detection in iac," in *Proceedings of the 4th ACM SIGSOFT International Workshop on Machine-Learning Techniques for Software-Quality Evaluation*, ser. MaLTeSQuE 2020, Virtual, USA: Association for Computing Machinery, 2020, pp. 7–12, ISBN: 9781450381246. DOI: 10.1145/3416505.3423564. [Online]. Available: https://doi.org/10.1145/3416505.3423564.
- [53] F. A. Fontana, M. V. Mäntylä, M. Zanoni, and A. Marino, "Comparing and experimenting machine learning techniques for code smell detection," *Empirical Software Engineering*, vol. 21, no. 3, pp. 1143–1191, 2016.
- [54] D. Di Nucci, F. Palomba, D. A. Tamburri, A. Serebrenik, and A. De Lucia, "Detecting code smells using machine learning techniques: Are we there yet?" In 2018 ieee 25th international conference on software analysis, evolution and reengineering (saner), IEEE, 2018, pp. 612–621.
- [55] T. Sharma, V. Efstathiou, P. Louridas, and D. Spinellis, "On the feasibility of transfer-learning code smells using deep learning," *arXiv preprint arXiv:1904.03031*, 2019.
- [56] F. Arcelli Fontana and M. Zanoni, "Code smell severity classification using machine learning techniques," *Knowledge-Based Systems*, vol. 128, pp. 43–58, 2017, ISSN: 0950-7051. DOI: https://doi.org/10.1016/j.knosys.2017.04.014. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S0950705117301880.
- [57] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [58] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," English (US), arXiv, 2014.
- [59] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," *arXiv* preprint *arXiv*:1310.4546, 2013.
- [60] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016, pp. 1715–1725.
- [61] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, and M. Zhou, Codebert: A pre-trained model for programming and natural languages, 2020. arXiv: 2002.08155 [cs.CL].
- [62] V. Efstathiou, C. Chatzilenas, and D. Spinellis, "Word embeddings for the software engineering domain," in

- Proceedings of the 15th international conference on mining software repositories, 2018, pp. 38–41.
- [63] J. Phang, T. Févry, and S. R. Bowman, "Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks," *arXiv* preprint *arXiv*:1811.01088, 2018.
- [64] S. Garg, T. Vu, and A. Moschitti, "Tanda: Transfer and adapt pre-trained transformer models for answer sentence selection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 2020, pp. 7780–7788.
- [65] Y. Pruksachatkun, J. Phang, H. Liu, P. M. Htut, X. Zhang, R. Y. Pang, C. Vania, K. Kann, and S. R. Bowman, Intermediate-task transfer learning with pretrained models for natural language understanding: When and why does it work? 2020. arXiv: 2005.00628 [cs.CL].
- [66] H. Scudder, "Probability of error of some adaptive pattern-recognition machines," *IEEE Transactions on Information Theory*, vol. 11, no. 3, pp. 363–371, 1965.
- [67] D. Yarowsky, "Unsupervised word sense disambiguation rivaling supervised methods," in 33rd Annual Meeting of the Association for Computational Linguistics, Cambridge, Massachusetts, USA: Association for Computational Linguistics, Jun. 1995, pp. 189–196. DOI: 10.3115/981658.981684. [Online]. Available: https://www.aclweb.org/anthology/P95-1026.
- [68] X. Zhang, J. Zhao, and Y. LeCun, Character-level convolutional networks for text classification, 2016. arXiv: 1509. 01626 [cs.LG].
- [69] J. Wei and K. Zou, Eda: Easy data augmentation techniques for boosting performance on text classification tasks, 2019. arXiv: 1901.11196 [cs.CL].
- [70] X. Jiao, Y. Yin, L. Shang, X. Jiang, X. Chen, L. Li, F. Wang, and Q. Liu, *Tinybert: Distilling bert for natural language understanding*, 2020. arXiv: 1909.10351 [cs.CL].
- [71] W. Y. Wang and D. Yang, "That's so annoying!!!: A lexical and frame-semantic embedding based data augmentation approach to automatic categorization of annoying behaviors using #petpeeve tweets," in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, Lisbon, Portugal: Association for Computational Linguistics, Sep. 2015, pp. 2557–2563. DOI: 10.18653/v1/D15-1306. [Online]. Available: https://www.aclweb.org/anthology/D15-1306.
- [72] S. Garg and G. Ramakrishnan, Bae: Bert-based adversarial examples for text classification, 2020. arXiv: 2004. 01970 [cs.CL].
- [73] R. Sennrich, B. Haddow, and A. Birch, *Improving neural machine translation models with monolingual data*, 2016. arXiv: 1511.06709 [cs.CL].
- [74] B. Settles, "Active learning literature survey," University of Wisconsin–Madison, Computer Sciences Technical Report 1648, 2009.
- [75] J. Keim, A. Kaplan, A. Koziolek, and M. Mirakhorli, "Does bert understand code? – an exploratory study on the detection of architectural tactics in code," in Software Architecture, A. Jansen, I. Malavolta, H. Muccini, I. Ozkaya, and O. Zimmermann, Eds., Cham: Springer International Publishing, 2020, pp. 220–228, ISBN: 978-3-030-58923-3.

- [76] A. Svyatkovskiy, S. K. Deng, S. Fu, and N. Sundaresan, "Intellicode compose: Code generation using transformer," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. New York, NY, USA: Association for Computing Machinery, 2020, pp. 1433–1443, ISBN: 9781450370431. [Online]. Available: https://doi-org.libproxy.unibz.it/10.1145/3368089.3417058.
- [77] H. Husain, H.-H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "CodeSearchNet Challenge: Evaluating the State of Semantic Code Search," arXiv:1909.09436 [cs, stat], Sep. 2019, arXiv: 1909.09436. [Online]. Available: http://arxiv.org/abs/1909.09436 (visited on 03/12/2020).
- [78] L. Kotthoff, C. Thornton, H. H. Hoos, F. Hutter, and K. Leyton-Brown, "Auto-weka: Automatic model selection and hyperparameter optimization in weka," in *Automated Machine Learning*, Springer, Cham, 2019, pp. 81–95.
- [79] R. Munro, *Human-in-the-Loop Machine Learning*. Manning Publications Co., 2021, ISBN: 9781617296741.
- [80] S. Kim, J. Zhao, Y. Tian, and S. Chandra, "Code prediction by feeding trees to transformers," arXiv preprint arXiv:2003.13848, 2020.
- [81] V. Shiv and C. Quirk, "Novel positional encodings to enable tree-based transformers," in *Advances in Neural Information Processing Systems*, 2019, pp. 12 081–12 091.
- [82] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio, *An empirical investigation of catastrophic forgetting in gradient-based neural networks*, 2015. arXiv: 1312.6211 [stat.ML].
- [83] J. Lin, Y. Liu, Q. Zeng, M. Jiang, and J. Cleland-Huang, "Traceability transformed: Generating more accurate links with pre-trained bert models," in 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), 2021, pp. 324–335. DOI: 10.1109/ICSE43902.2021.00040.
- [84] H. Tu, Z. Yu, and T. Menzies, "Better data labelling with emblem (and how that impacts defect prediction)," *IEEE Transactions on Software Engineering*, 2020.
- [85] J. Dodge, G. Ilharco, R. Schwartz, A. Farhadi, H. Hajishirzi, and N. Smith, "Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping," arXiv preprint arXiv:2002.06305, 2020.
- [86] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush, *Huggingface's transformers: State-of-the-art natural language processing*, 2020. arXiv: 1910.03771 [cs.CL].
- [87] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
- [88] R. Müller, S. Kornblith, and G. E. Hinton, "When does label smoothing help?" In *Advances in Neural Information Processing Systems*, 2019, pp. 4694–4703.

- [89] E. Frank, M. A. Hall, G. Holmes, R. Kirkby, B. Pfahringer, and I. H. Witten, "Weka: A machine learning workbench for data mining.," in *Data Mining and Knowledge Discovery Handbook: A Complete Guide for Practitioners and Researchers*, O. Maimon and L. Rokach, Eds. Berlin: Springer, 2005, pp. 1305–1314. [Online]. Available: http://researchcommons.waikato.ac.nz/handle/10289/1497.
- [90] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [91] T. T. Le, W. Fu, and J. H. Moore, "Scaling tree-based automated machine learning to biomedical big data with a feature set selector," *Bioinformatics*, vol. 36, no. 1, pp. 250–256, 2020.



Julian Aron Prenner received his M.Sc. in Computer Science from the Free University of Bozen-Bolzano, where he is also currently pursuing his Ph.D. His research interests include program repair, automatic test case generation and applications of machine learning in Software Engineering.



Romain Robbes is an Associate Professor at the Free University of Bozen-Bolzano, in the SwSE research group. Before that, he was an Assistant, then Associate Professor at the University of Chile. He earned his PhD in 2008 from the University of Lugano, Switzerland and his Master's degree from the University of Caen, France. His research interests include Empirical Software Engineering, Software Evolution, Mining Software Repositories, and Machine Learning for Software Engineering.

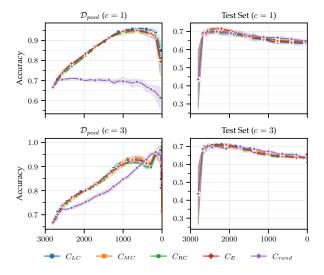


Fig. 5. Accuracy of \mathcal{D}_{pool} and the test set at each iteration of the active learning process for different acquisition functions with (c=3) and without (c=1) diversity sampling for the review classification task. Error bands are 95% confidence intervals.

APPENDIX A ADDITIONAL INFORMATION ON THE DATASETS

Data Quality. On the SATD dataset we noticed that false positives often contain keywords such as FIXME, TODO, or HACK; while this might raise questions about quality, we found that only 1.2% of the negative instances contain such keywords, compared to 6% among positives. Similarly, in the app review dataset by Maalej *et al.* [21] we found that 9% of the samples appear twice, with different labels. We left these duplicates in the dataset as we did not know how this issue was handled by previous work, nor which duplicates to remove and which to keep.

Additional examples and Statistics. Table 14 shows a representative example of each dataset, alongside with its class. Table 15 shows detailed statistics (per project) on the Self-Admited Technical Debt Dataset.

APPENDIX B ADDITIONAL RESULTS

B.1 Active Learning

Figure 5 shows the results of active learning on the app review classification dataset. On this dataset, all pool curves show a noticeable drop at the end. A possible explanation for this might be data quality: we observed that the app review dataset by Maalej *et al.* [21] contains a number of duplicates with conflicting labels. Having learned one of the duplicate samples, all its copies will be considered very "easy" and not be selected until the very end, at which point the model will predict the label of the duplicate selected first, which will, as labels are conflicting, be wrong, causing accuracy to drop towards the end of the active learning process.

B.2 Label Smoothing

A method to obtain soft-labels that does not require any additional information is *label smoothing* [87], [88]. In label smoothing, the original target distribution is mixed with the uniform distribution over all classes: For a given target vector y, its smoothed version is calculated as: $\mathbf{y}_{smooth} = (1-\alpha)\cdot\mathbf{y} + \alpha\cdot\frac{1}{K}$, where K is the number of classes and α controls the smoothing strength. As an example, smoothing the target vector (A: 0, B: 1, C: 0) with $\alpha = 0.2$ results in (A: 0.06, B: 0.86, C: 0.06), now a soft-label. Label smoothing is a form of regulation: intuitively, it dampens the model's prediction confidence, forcing it to make more "cautious" predictions.

We carried out a label-smoothing experiment on the Senti4SD dataset in addition to soft labels. We train a subset of our Transformer models with hard-labels, soft-labels and different degrees of label-smoothing ($\alpha=0.1,0.05$ and 0.03) respectively, and evaluate on all sentiment analysis datasets; refer to Table 16 for a comparison of the results. While label smoothing can occasionally improve performance (e.g. on Jira issues), it is more likely to either degrade performance, or not affecting it significantly.

B.3 Task-Specific Tokens

In the code readability task, our Transformer models cannot compete with manually engineered features used by the baseline. Since Buse and Weimer [47] found that line length is one of the most important features for predicting code readability, we attempt to provide this information explicitly to our model in form the of special line length tokens, added to line ends. These tokens range from <11>, indicating a short line, up to <110> for very long lines and are inserted before newline tokens. We fine-tuned the same model with and without these special tokens.

Line length tokens failed to improve the performance. In fact, they seem to hurt performance (see Table 17). The logistic regression model by Scalabrino *et al.* outperforms our Transformer model by a wide margin. We were able to successfully reproduce the results of Scalabrino's model, which was implemented using Weka [89]. We found that for their model, attribute selection is crucial; without it, in our experiments, accuracy dropped significantly (below 60%).

A simple logistic regression model implemented using scikit-learn [90], even with attribute selection, was similarly unable to beat their Weka model; neither was TPOT [91], a framework for automated machine learning, that automatically evaluates a large number of combinations of different machine learning algorithms. By first using Weka's attribute selection algorithm and feeding selected attribute to a scikit-learn logistic regression model we were eventually able to obtain results close to the Weka-only model.

B.4 Sentiment Analysis

Full results, including per-class precision and recall for sentiment classification (Table 18). Finally, as can be seen in table 19, sentences that are most confusing to the model are hard to classify even for humans. Similar observations can be made in other datasets (e.g., SATD, not shown here).

Name	Sample	Class
Sentiment Classification (Stack Overflow) [18]	I want them to resize based on the length of the data they're showing.	neutral
Sentiment Classification (Stack Overflow) [19]	When I run my client, it throws the following exception.	negative
Sentiment Classification (JIRA Issues) [3]	This is always a really bad way to design software.	negative
Sentiment Classification (App Reviews) [19]	amazing! a must have app	positive
Informative App Review Detection [20]	not able to download any pictures please fix these bugs immediately	informative
App Review Classification [21]	Best game I've played on Android	rating
App Review Classification [27]	good but it has adsplease remove ads from this	usability
Self-Admitted Technical Debt Detection [22]	// FIXME: Is "No Namespace is Empty Namespace" really OK?	SATD
Comment Classification [23]	@return a string for throwing	usage
Code-Comment Coherence Prediction [24]	<pre>/** * Returns the current number of milk units in * the inventory. * @return int */ public int getMilk() { return milk; }</pre>	coherent
Linguistic Smell Detection [2]	<pre>public void ToSource(StringBuilder sb) { sb.append(";"); this.NewLine(sb); }</pre>	smelly (transform method does not return)
Code Runtime Complexity Classification [25]	<pre>class GFG { static int minJumps(int arr[], int n) { int[] jumps = new int[n]; int min; jumps[n - 1] = 0; for (int i = n - 2; i >= 0; i) { if (arr[i] == 0) jumps[i] = Integer.MAX_VALUE; else if (arr[i] >= n - i - 1) jumps[i] = 1; else { } } return jumps[0]; } public static void main(String[] args) {}</pre>	$\mathcal{O}(n\log n)$
Code Readability Prediction [26]	<pre>@Override public void configure(Configuration cfg) { super.configure(cfg); cfg.setProperty(Environment.USE_SECOND_LEVEL_CACHE,) cfg.setProperty(Environment.GENERATE_STATISTICS,); cfg.setProperty(Environment.USE_QUERY_CACHE, "false"); // more cfg.setProperty calls }</pre>	; readable

	SATD	not SATD	Total
Apache Ant	131	3 967	4 098
ÂrgoUML	1 413	8 039	9 452
Columba	204	6 2 6 4	6468
EMF	104	4 286	4390
Hibernate	472	2 496	2968
JEdit	256	10 066	10322
JFreeChart	209	4 199	4408
JMeter	374	7 683	8 057
JRuby	622	4 275	4897
SQuirrel	286	6 9 2 9	7 2 1 5
Total	4071	58 204	62 275

TABLE 15
Projects in the SATD dataset along with the number of samples classified as self-admitted technical debt.

Dataset	Label Type	F1		
Dataset	Laber Type	$\mu \pm \sigma$	max	
	All Label Votes ($\alpha = 0.03$)	56.9 ± 3.3	61.6	
	All Label Votes ($\alpha = 0.05$)	56.2 ± 2.9	58.8	
	All Label Votes ($\alpha = 0.1$)	56.4 ± 2.6	59.0	
A D [10]	All Label Votes (no smooth.)	$\textbf{57.5} \pm 2.2$	59.1	
App Reviews [19]	Majority Label ($\alpha = 0.03$)	55.7 ± 3.8	60.4	
	Majority Label ($\alpha = 0.05$)	55.9 ± 1.5	57.6	
	Majority Label ($\alpha = 0.1$)	56.7 ± 1.6	58.8	
	Majority Label (no smooth.)	57.4 ± 1.6	59.7	
	All Label Votes ($\alpha = 0.03$)	92.7 ± 1.4	94.2	
	All Label Votes ($\alpha = 0.05$)	93.2 ± 0.7	94.0	
	All Label Votes ($\alpha = 0.1$)	93.1 ± 1.1	94.2	
IID A Jaguage [2]	All Label Votes (no smooth.)	93.5 ± 1.0	94.7	
JIRA Issues [3]	Majority Label ($\alpha = 0.03$)	91.6 ± 1.9	94.3	
	Majority Label ($\alpha = 0.05$)	92.4 ± 1.3	93.6	
	Majority Label ($\alpha = 0.1$)	$\textbf{93.5} \pm 0.7$	94.3	
	Majority Label (no smooth.)	92.5 ± 0.5	93.2	
	All Label Votes ($\alpha = 0.03$)	41.6 ± 1.3	43.6	
	All Label Votes ($\alpha = 0.05$)	42.5 ± 2.4	46.1	
	All Label Votes ($\alpha = 0.1$)	41.0 ± 1.4	42.7	
StackOverflow [19]	All Label Votes (no smooth.)	42.4 ± 1.0	43.8	
StackOvernow [19]	Majority Label ($\alpha = 0.03$)	42.8 ± 3.5	47.4	
	Majority Label ($\alpha = 0.05$)	42.7 ± 2.3	45.6	
	Majority Label ($\alpha = 0.1$)	43.1 ± 1.4	45.2	
	Majority Label (no smooth.)	44.1 \pm 2.9	46.4	
	All Label Votes ($\alpha = 0.03$)	88.9 ± 0.2	89.2	
	All Label Votes ($\alpha = 0.05$)	89.0 ± 0.6	89.9	
	All Label Votes ($\alpha = 0.1$)	89.0 ± 0.4	89.5	
Senti4SD [18]	All Label Votes (no smooth.)	$\textbf{89.1} \pm 0.5$	89.6	
(Test Set)	Majority Label ($\alpha = 0.03$)	88.4 ± 0.5	89.0	
	Majority Label ($\alpha = 0.05$)	88.4 ± 0.4	89.1	
	Majority Label ($\alpha = 0.1$)	88.5 ± 0.3	88.8	
	Majority Label (no smooth.)	88.6 ± 0.5	89.1	

TABLE 16

Macro F1 for different label types and datasets: mean, maximum and standard deviation over five runs, each with different seed. All models were trained on the training set of the Stack Overflow dataset from Calefato *et al.* [18].

Model	Acc.
CodeBERTa	73.1
CodeBERTa+LLT	72.5
CodeBERT	69.3
Logistic Regression [50]	84.0

TABLE 17

Results for the code readability prediction task. As far as our number are concerned, values are means over five runs with different seeds.

APPENDIX C FURTHER DETAILS ON IMPLEMENTATION, RUNTIME, REPLICATION

Obtaining back-translation data. While the original works introducing back-translation used an ad-hoc neural translation model, we found that the most efficient way to obtain back-translations is to load the dataset into Google Sheet and use the GOOGLETRANSLATE macro. An example is available online⁴.

Model implementations. For all our experiments we use HuggingFace's transformers package [86], a Python library based on pytorch that implements many different Transformer architectures, including BERT and RoBERTa.

Runtime considerations. All of our experiments were carried out either on an NVIDIA V100 GPU with 32 GB of memory or on up to three NVIDIA RTX 2080TIs with 10 GB memory each.

Our pre-training regimes are generally affordable even with relatively modest computational budget, although an extensive hyper-parameter search is hardly feasible. We thus followed common recommendations and only tried a few parameter combinations. With a training time of two weeks, pre-training StackOBERTflow from scratch was by far the most expensive (especially considering that this was clearly not enough, as it ended up being out-performed by the further pre-trained models). Further pre-training the 12-layer models required considerably less training time, usually below 24 hours (on an NVIDIA V100 GPU).

Using the pre-trained models. Our models are publicly available: the StackOBERTflow model can be obtained through the Huggingface Model Hub⁵; our fine-tuned BERT and RoBERTa models can be downloaded from GitHub⁶.

You can download our pre-trained models and use them for your own experiments. Our StackOBERTflow model can be automatically downloaded using the transformers library. You can instantiate a classification model using model = AutoModel.from_pretrained ('giganticode/StackOBERTflow-comments-small-v1'), and then fine-tune on your task-specific data. You can also use our other models: first download the model as ZIP archive from our GitHub page and unpack it; then, likewise, load them as follows: model = AutoModel.from_pretrained('/path/to/model')

Rerunning First experiments. clone our repository⁷; GitHub then run python -mdl4se.experiments.<experiment>.default -seeds 100 200 300 400 500 -out_file=result_file.csv, where experiment is one of the experiments listed in Table 20. Configuration options and default hyper-parameters can for each experiment be found in /dl4se/config/<experiment>.py,

- 4. https://docs.google.com/spreadsheets/d/ 19X8vvV3LF9m2fqUwS9L9yEOPxD-RzshC-fCBTGjkD8I/edit?usp=sharing
- 5. https://huggingface.co/giganticode/ StackOBERTflow-comments-small-v1
- 6. https://github.com/giganticode/small-datasets-ml-resources/releases/tag/0.1
 - 7. https://github.com/giganticode/small-datasets-ml-resources

Dataset	Model	Acc. F1		Precision			Recall		
Dataset	Wiodei	Acc.	LI	Pos.	Neg.	Neu.	Pos.	Neg.	Neu.
	BERT (base)	64.9	50.2	76.1	84.7	9.8	86.6	41.4	25.6
	BERT-SO-1M-large	66.0	51.9	77.2	83.7	10.6	84.7	46.8	27.2
	BERT-SO-2M	67.7	53.1	77.2	85.4	11.7	87.4	47.4	27.2
	BERT-SO-1M	68.3	53.7	77.2	86.9	12.7	88.0	47.8	28.8
	DistilRoBERTa (small)	60.9	48.5	75.3	74.7	10.1	77.3	43.5	28.8
App Reviews [19]	NLTK [3]	54.0	40.8	75.1	100.0	9.3	81.2	16.9	44.0
App Reviews [19]	RoBERTa (small, no pretr.)	48.0	42.2	73.8	68.9	8.5	53.0	42.0	42.4
	SentiStrength+AC0-SE [3]	58.9	46.6	74.1	92.9	10.6	81.7	30.0	40.0
	SentiStrength [3]	62.5	48.2	74.5	81.5	11.3	86.6	33.8	32.0
	StackOBERTflow	72.0	57.4	79.7	86.4	15.1	88.4	56.6	29.6
	Stanford CoreNLP SO [3]	41.6	35.5	77.0	47.0	8.4	25.3	66.9	32.0
	Stanford CoreNLP [3]	69.5	56.0	83.1	66.7	17.6	71.5	75.4	24.0
	BERT (base)	95.7	95.0	92.4	97.2	-	94.0	96.4	-
	BERT-SO-1M-large	94.8	93.8	93.3	95.4	-	89.7	97.0	-
	BERT-SO-2M	95.2	94.5	92.2	96.8	-	92.9	96.3	-
	BERT-SO-1M	95.1	94.2	93.1	96.0	-	91.0	96.9	-
	DistilRoBERTa (small)	93.7	92.7	89.9	95.5	-	90.1	95.3	-
IID A I [2]	NLTK [3]	29.8	46.5	84.0	100.0	-	36.2	26.9	-
JIRA Issues [3]	RoBERTa (small, no pretr.)	84.0	80.1	78.6	86.3	-	67.0	91.7	-
	SentiStrength+AC0-SE [3]	76.0	87.0	94.8	99.6	-	88.3	70.4	-
	SentiStrength [3]	77.1	85.4	85.0	99.3	-	92.1	70.3	-
	StackOBERTflow	93.5	92.5	87.8	96.4	-	92.1	94.1	-
	Stanford CoreNLP SO [3]	36.0	44.2	63.5	72.4	-	25.2	40.9	-
	Stanford CoreNLP [3]	67.6	73.7	72.6	94.5	-	62.1	70.1	-
	BERT (base)	79.9	47.6	32.1	79.9	82.0	15.9	21.1	95.7
	BERT-SO-1M-large	79.2	43.3	29.8	78.2	81.2	13.4	13.8	96.2
	BERT-SO-2M	80.3	48.6	33.7	79.8	82.4	16.2	23.0	95.9
	BERT-SO-1M	80.1	47.9	34.9	80.0	82.1	16.8	20.8	95.9
	DistilRoBERTa (small)	78.8	44.3	31.5	61.3	81.5	14.8	16.2	95.2
CharleOrrondlarer [10]	NLTK [3]	77.9	43.2	31.7	62.5	81.5	24.4	8.4	94.1
StackOverflow [19]	RoBERTa (small, no pretr.)	75.5	42.2	21.1	37.2	81.6	10.8	21.1	90.7
	SentiStrength+AC0-SE [3]	78.0	46.8	31.2	50.0	82.6	22.1	18.5	93.0
	SentiStrength [3]	69.5	49.5	20.0	39.7	85.8	35.9	43.3	77.2
	StackOBERTflow	79.3	44.1	32.1	78.1	81.5	15.7	13.7	96.1
	Stanford CoreNLP SO [3]	75.9	47.5	31.7	36.5	83.6	14.5	36.5	88.6
	Stanford CoreNLP [3]	40.3	35.5	23.1	17.7	88.4	34.4	83.7	34.4
	BERT (base)	87.9	87.8	92.8	88.0	83.7	93.0	83.1	86.8
	BERT-SO-1M-large	89.4	89.4	93.4	88.2	86.8	93.6	86.8	87.4
	BERT-SO-2M	88.8	88.7	93.2	85.6	87.1	94.9	86.7	84.7
	BERT-SO-1M	88.8	88.7	92.7	86.9	86.7	94.5	85.9	85.7
C ('4CD [10]	DistilRoBERTa (small)	87.4	87.3	92.3	81.2	87.9	92.7	88.6	81.8
Senti4SD [18]	RoBERTa (small, no pretr.)	78.4	77.8	84.0	72.4	77.6	86.1	69.1	78.1
(Test Set)	Senti4SD [18]	-	86.0	92.0	80.0	87.0	92.0	89.0	80.0
	SentiCR [18]	-	82.0	88.0	79.0	79.0	90.0	73.0	82.0
	SentiStrengthSE [18]	-	80.0	89.0	75.0	75.0	83.0	79.0	77.0
	SentiStrength [18]	-	84.0	89.0	67.0	95.0	92.0	96.0	64.0
	StackOBERTflow	88.7	88.6	93.0	85.6	86.9	94.4	87.0	84.6

TABLE 18

Accuracy, macro F1 and per-class precision and recall for different models and datasets. Values reported are means over five runs, each with different seed (only our models). All models were trained on the training set of the Stack Overflow sentiment dataset from Calefato et al. [18].

dataset loading and pre-processing code lies in /dl4se/datasets/<experiment>.py Note that you cannot use the original datasets, as datasets need to adhere to a specific format. We will provide all of the datasets upon request.

Sample	Act. Label	Pred. Label	Agr.
This worked for me.	positive	neutral	no
In addition to firebug (which should be your first port of call), the will also tell you where a given style is sourced from, just in case IE - shock, horror - should be different.	negative	neutral	no
The named scopes already proposed are pretty fine. The clasic way to do it would be:	positive	neutral	yes
I have a basic app written with ATL, using the wizard with VS2008. I have a treeview in the left side of the app. I see how to (painfully) add tree items. Question is how do I show a menu when the mouse is right clicked? How do I trap any click events on each item that could be selected?	negative	neutral	yes
Your implementation looks absolutely fine to me! A range-based away subscript is a type for performance reasons. It does not copy the indicated sub-array, instead it just points to the range defined by the you provide to the subscript.	positive	neutral	no
Is it possible (or desirable?!) to set up to behave more like ? For example, instead of writing why can't I just write Similarly, instead of why not just	negative	neutral	no
This whole DB is almost entirely read only so I'm not too worried about it changing.	positive	neutral	no
I would agree except they are related and I really hated the idea of writing 4 separate questions since they seemed to close.	negative	positive	yes
I understand its not a desirable circumstance, however if I NEEDED to have some kind of HTML within JSON tags, e.g.: is this possible to do in Python without requiring to to be escaped beforehand? It will be a string initially so I was thinking about writing a regular expression to attempt to match and escape these prior to processing, but I just want to make sure there isn't an easier way.	negative	neutral	no
You can easily define a comparator for a one-level , so that lookup becomes way less cumbersome. There is no reason of being afraid of that. The comparator defines an ordering of the _Key template argument of the map. It can then also be used for the multimap and set collections. An example:	positive	neutral	no
I'm new to Flash but want to create a nice video for a product. It takes a long time to make a nice looking presentation , and I'm hoping for a jump start. Are there any good templates which are free on the internet where I can quickly change the text in ,for example, to make my video? I've tried looking in google, and there are too many websites, many of which look gimmicky. Any recommendations? (A video like this one would be amazing!)	neutral	positive	no

TABLE 19
Sentences from the test set of the Stack Overflow [18] with highest loss along with predicted and actual labels and whether all raters agreed on the actual label.

Experiment	Description	Original Dataset URL
ar_miner	Informative app reviews	https://github.com/jinyyy666/AR_ Miner/tree/master/datasets
coherence	Code-comment coherence	http://www2.unibas.it/gscanniello/coherence/
comment_classification	Comment classification	https://zenodo.org/record/2628361
corcod	Runtime complexity classification	https://github.com/midas-research/ corcod-dataset
readability	Code readability classification	https://dibt.unimol.it/report/ readability/
review_classification	Review classification	https://mast.informatik.uni-hamburg. de/wp-content/uploads/2014/03/ REJ_data.zip, CLAP was requested from the respective authors
satd	Self-admitted debt detection	https://github.com/maldonado/tse. satd.data
senti4sd	Sentiment analysis on Stack Overflow comments and JIRA issues	https://github.com/collab-uniba/ Senti4SD, https://sentidata.github.io/
smell_detection	Linguistic smell detection	https://github.com/Smfakhoury/ SANER-2018-KeepItSimple-

TABLE 20 Different code modules along with the source of the used datasets.