

An Empirical Study on the Effectiveness of Large Language Models for SATD Identification and Classification

Mohammad Sadegh Sheikhaei · Yuan Tian · Shaowei Wang · Bowen Xu

Received: date / Accepted: date

Abstract Self-Admitted Technical Debt (SATD), a concept highlighting sub-optimal choices in software development documented in code comments or other project resources, poses challenges in the maintainability and evolution of software systems. Large language models (LLMs) have demonstrated significant effectiveness across a broad range of software tasks, especially in software text generation tasks. Nonetheless, their effectiveness in tasks related to SATD is still under-researched. In this paper, we investigate the efficacy of LLMs in both identification and classification of SATD. For both tasks, we investigate the performance gain from using more recent LLMs, specifically the Flan-T5 family, across different common usage settings.

Our results demonstrate that for SATD identification, all fine-tuned LLMs outperform the best existing non-LLM baseline, i.e., the CNN model, with a 4.4% to 7.2% improvement in F1 score. In the SATD classification task, while our largest fine-tuned model, Flan-T5-XL, still led in performance, the CNN model exhibited competitive results, even surpassing four of six LLMs. We also found that the largest Flan-T5 model, i.e., Flan-T5-XXL, when used with a zero-shot in-context learning (ICL) approach for SATD identification, provides competitive results with traditional approaches but performs 6.4% to 9.2% worse than fine-tuned LLMs. For SATD classification, few-shot ICL approach, incorporating examples and category descriptions in prompts, outperforms the zero-shot approach and even surpasses the fine-tuned smaller

Mohammad Sadegh Sheikhaei and Yuan Tian
School of Computing, Queen's University, Kingston, ON, Canada
E-mail: sadegh.sheikhaei@gmail.com y.tian@queensu.ca

Shaowei Wang
Department of Computer Science, University of Manitoba, Winnipeg, MB, Canada
E-mail: Shaowei.Wang@umanitoba.ca

Bowen Xu
Department of Computer Science, North Carolina State University, Raleigh, NC, US
E-mail: bxu22@ncsu.edu

Flan-T5 models. Moreover, our experiments demonstrate that incorporating contextual information, such as surrounding code, into the SATD classification task enables larger fine-tuned LLMs to improve their performance. Our study highlights the capabilities and limitations of LLMs for SATD tasks and the role of contextual information in achieving higher performance with larger LLMs, setting a foundation for future efforts to enhance these models for more effective technical debt management.

Keywords Self-admitted technical debt (SATD) · SATD identification · SATD classification · Large language models · Fine tuning · In-context learning

1 Introduction

Technical Debt (TD) refers to the deliberate adoption of less-than-ideal solutions in software design or coding, typically to meet urgent deadlines or address immediate resource constraints (Cunningham, 1992). This metaphor portrays the trade-off between short-term expediency and long-term software quality, likening it to incurring financial debt that accrues interest in the form of increased maintenance effort over time (Buschmann, 2011). In literature, large-scale analysis on TD is often facilitated through the study of one specific type of TD, i.e., Self-Admitted Technical Debt (SATD) (Potdar and Shihab, 2014). SATDs are sub-optimal choices that developers consciously make and primarily document in code comments (Maldonado and Shihab, 2015; Guo et al., 2021; O'Brien et al., 2022), though developers may also occasionally record them in other software artifacts such as issue reports (Li et al., 2022).

Addressing SATD is crucial, as its accumulation can significantly impair the long-term maintainability and evolution of software systems Li et al. (2023b). A study by Wehaibi et al. (Wehaibi et al., 2016) indicates that source code files with SATD experience a higher frequency of bug-fixing changes than those without SATD. To tackle this challenge effectively, software development teams must manage SATD by first pinpointing the locations of existing SATDs within their software (*SATD Identification*) (Sheikhaei and Tian, 2023). Subsequently, these SATDs could be categorized into various types based on their specific impacts on software maintenance, such as design, requirement, defect, documentation, and test debt (*SATD Classification*). This categorization aids in prioritizing efforts and efficiently allocating developers to address and resolve these SATDs (Maldonado et al., 2017). However, manually performing the above two tasks, i.e., SATD identification and SATD classification, is challenging and time-consuming. For SATD identification, in large projects, only approximately 0.5 to 4% of a project's code comments are SATD (Guo et al., 2021). While certain keywords like TODO and FIXME strongly suggest a code comment as SATD (Rantala et al., 2020; Guo et al., 2021), developers often mention technical debt without using these specific keywords. Consequently, only about 20 to 90% of SATD in a project can be identified using such keywords (Yu et al., 2022). For SATD classification, there exist no specific

keyword sets that can be universally applied to effectively detect the various types of SATDs. This limitation underlines the need to develop sophisticated automated approaches to identify and categorize SATDs.

There have been numerous studies on automated SATD identification and classification (Maldonado and Shihab, 2015; Maldonado et al., 2017; Ren et al., 2019; Guo et al., 2021; Xiao et al., 2021; Cassee et al., 2022; Yu et al., 2022; O'Brien et al., 2022; Sridharan et al., 2023). These approaches are either rule-based or machine learning-based. While these models show promise, their performance on these two classification tasks could be further enhanced. Recently, machine learning has seen the emergence of Large Language Models (LLMs). LLMs have exhibited strong capabilities in text generation tasks, such as text summarization and question answering (Naveed et al., 2023). In the software engineering (SE) domain, researchers also found that LLMs can achieve state-of-the-art (SOTA) performance on generative tasks such as code generation (Wei et al., 2023), code summarization (Yuan et al., 2023), and program repair (Jin et al., 2023). However, to the best of our knowledge, no prior studies have investigated how to better use these complex models specifically for SATD identification and classification. Furthermore, LLMs have shown impressive capabilities in in-context learning (ICL) (Brown et al., 2020). With the ICL ability, an LLM can adapt to new tasks or understand new information by leveraging examples provided directly in its input (prompt) without fine-tuning - a costly and time-consuming process that adapts pre-trained models to specific tasks due to updating model parameters. This raises an intriguing question: can ICL with LLMs achieve competitive performance in SATD identification and classification compared to fine-tuned LLMs?

In this paper, we aim to further tap into the potential of LLMs in two important SATD-related tasks, i.e., SATD identification and classification. We explore the impact of various usage settings of LLMs in these tasks, encompassing fine-tuning versus ICL, model size, prompt engineering, and adaptation in model architecture. Additionally, we explore the potential of integrating contextual features beyond code comment into SATD classification. Our findings can guide future research and advance the use of LLMs in managing SATD and other classification tasks within the SE field.

To achieve our goal, we conduct an empirical study using two well-known datasets: the Maldonado-62k dataset (Maldonado et al., 2017) (we use the revised version by (Yu et al., 2022)), and the O'Brien dataset (O'Brien et al., 2022). The Maldonado-62k dataset comprises 62,275 code comments from ten popular open-source software projects across various application domains, of which 7.2% (4,497 comments) are identified as SATD. The O'Brien dataset includes 856 SATDs randomly sampled from 68,820 SATDs extracted from 2,641 popular machine-learning repositories on GitHub. These SATDs are categorized into six general types: requirement, code, test, defect, design, and documentation. Our empirical study answers the following four research questions (RQs):

RQ1 *How effective are fine-tuned LLMs in SATD identification and classification?* We evaluated the performance of six large language models: BERT-base (Devlin et al., 2019), CodeBERT (Feng et al., 2020), and four variants of Flan-T5 (small, base, large, and XL) (Chung et al., 2022), comparing them against existing baselines for the SATD identification and classification. Our analysis revealed that for SATD identification, the selected LLMs outperformed the best existing baseline, showing a 4.4% to 7.2% improvement in F1 score. The larger model, Flan-T5-XL, achieved marginally better results, with a 0.4% to 2.4% higher F1 score compared to its smaller counterparts. In the SATD classification task, while the fine-tuned Flan-T5-XL still led in performance, the CNN model exhibited competitive results, even surpassing four of six LLMs.

RQ2 *Does our proposed ICL with a larger model outperform smaller models that have been fine-tuned in identifying and classifying SATD?*

We propose an ICL approach for SATD identification and classification using the largest Flan-T5 model, i.e., Flan-T5-XXL (Chung et al., 2022). We found that a zero-shot approach with the Flan-T5-XXL model provides competitive results for SATD identification with traditional approaches but performs 6.4% to 9.2% worse than fine-tuned LLMs. For SATD classification, few-shot incorporating examples and category descriptions in prompts outperforms the zero-shot approach and even surpasses the results of fine-tuning smaller Flan-T5 models, i.e., the small and base versions.

RQ3 *What is the impact of adding the classification layer in fine-tuning LLMs?* Different LLM architectures naturally lend themselves to distinct fine-tuning strategies. We found that in the absence of sufficient training data, changing the architecture of the Flan-T5 models by substituting its original text generation layer with a classification layer is beneficial, especially when utilizing the smaller versions of the Flan-T5 models.

RQ4 *What is the impact of additional contextual features on LLM-based SATD classification?* Existing SATD identification and classification approaches only take code comments as input and overlook other contextual information. We explore the potential of LLMs to leverage this additional contextual information for improving SATD classification performance. We found that larger fine-tuned models such as Flan-T5-large and Flan-T5-XL can effectively utilize these contextual features to enhance performance. In contrast, smaller models and those employing ICL exhibit a decrease in performance when complex contextual information is included.

Our work makes the following main contributions:

- We design and perform a comprehensive evaluation to assess the performance of six popular LLMs in automated SATD identification and classification. This evaluation encompasses various aspects, including different model types, sizes, adaptation methods (fine-tuning vs. ICL), prompting techniques, and model architectures (with or without a classification layer).
- We are the first to investigate the potential of including contextual features beyond code comments in SATD classification task, leveraging LLM.

- We demonstrate that LLMs can achieve state-of-the-art performance in SATD identification and classification. Specifically, a fine-tuned Flan-T5-XL model achieved an F1 score of 0.839 in SATD identification. In SATD classification, the top-performing model was an ensemble fine-tuned Flan-T5-XL, which integrated various combinations of contextual features and code comments, achieving an overall accuracy of 0.668.
- The results and source code related to this study are available at https://github.com/RISElabQueens/SATD_LLM.

The rest of this paper is organized as follows. Section 2 summarizes research background and related work. The experimental design and result analysis are presented in Section 3 and 4, respectively. We discuss the impact of epoch numbers in Section 5 and threats to validity in Section 6. Finally, the paper concludes with the future work in Section 7.

2 Related Work

2.1 Self-Admitted Technical Debt Identification

The goal of SATD identification is to determine whether a given code comment admitted that the corresponding code is a technical debt or not (Sheikhaei and Tian, 2023). A few studies aim to detect SATD in other software artifacts, such as issues (Li et al., 2022, 2023a), leveraging comment-based SATD approaches. Existing SATD identification approaches can be classified into two categories: 1) rule-based approaches (Potdar and Shihab, 2014; Guo et al., 2021; Sridharan et al., 2023) which search for certain keywords (e.g., TODO) or phrases (e.g., “probably a bug”) in the code comment, and 2) supervised learning based approaches (Huang et al., 2017; Maldonado et al., 2017; Ren et al., 2019; Prenner and Robbes, 2022) which train a model from labeled data and evaluate it on unseen code comments.

One of the most decent rule-based approaches is Matches task Annotation Tags (MAT) (Guo et al., 2021). In this approach, the authors showed that just matching a set of popular task annotation tags, i.e., TODO, FIXME, HACK, and XXX, provides similar or even superior performance for SATD identification compared with traditional machine learning approaches such as the natural language processing (NLP) approach (Maldonado et al., 2017) which leverages maximum entropy classifier, and the text-mining based approach (Huang et al., 2017) that employs Naïve Bayes Multinomial classifier. A more recent rule-based approach is PENTACET (Sridharan et al., 2023), where the authors extend the 64 SATD identification patterns introduced by Potdar and Shihab (Potdar and Shihab, 2014), to 1,041 patterns using a tool named Sense2Vec (Trask et al., 2015). Sense2Vec captures contextually similar words with its word embedding and the authors use these extrapolated features (words) to evaluate the code comments for SATD. While it has shown that simple rule-based approaches, e.g., MAT, have a high precision, though not having a good recall, the performance of a complex rule-based approach such

as PENTACET is not well studied. In this paper, we consider both MAT and PENTACET approaches as rule-based baselines, and the NLP approach (Maldonado et al., 2017) as a traditional machine-learning-based baseline for the SATD identification task.

More recently, supervised deep learning-based methods have emerged and achieved more promising performance compared to rule-based and traditional machine-learning approaches. Ren et al. (Ren et al., 2019) proposed a Convolutional Neural Network (CNN) based approach that outperforms traditional text-mining models. Yu et al. (Yu et al., 2021) proposed a Bidirectional Long Short-Term Memory (BiLSTM) model with a balanced cross-entropy loss function to overcome the class unbalance challenge. A more recent study by Prenner and Robbes (Prenner and Robbes, 2022) showed that the BERT base models provide state-of-the-art results in different software-related tasks, including SATD identification. As the CNN and the BERT base models have achieved the state-of-the-art results in this domain, we use CNN as a strong non-LLM baseline, and BERT as a LLM baseline in our experiments for both SATD identification and classification tasks.

2.2 Self-Admitted Technical Debt Classification

Maldonado and Shihab (Maldonado and Shihab, 2015) introduced the task of SATD classification. They extracted 33,093 code comments from five Java projects and manually classified them into six categories: non-SATD, Design, Requirement, Defect, Test, and Documentation. Later, Maldonado et al. (2017) extended the dataset by adding another five Java projects with the same categories and proposed an NLP approach using the Java implementation of a maximum entropy classifier, Stanford Classifier (Manning and Klein, 2003), to identify and classify SATDs by the code comments. Later, other researchers worked on that dataset and proposed another classification from a different perspective. Fucci et al. (Fucci et al., 2021) and Cassee et al. (Cassee et al., 2022) proposed a bottom-up strategy (i.e., what do SATD comments mention?) rather than top-down (i.e., how do SATD comments map into a software development life-cycle?) which was utilized in (Maldonado and Shihab, 2015). They manually classified 1,038 SATD instances from Maldonado dataset into 41 categories and subcategories, such as *“poor implemented choices”* and *“functional issues”*. Chen et al. (Chen et al., 2022) leveraged chi-square to select representative features from the textual features, and applied the XG-Boost model for the SATD classification task on the Maldonado dataset.

In addition to SATD classification for general software systems, there are studies that focus on identifying and classifying SATDs in specific domains such as Blockchain projects (Pinna et al., 2023), machine learning (ML) software (OBrien et al., 2022; Bhatia et al., 2023), and deep learning frameworks (Liu et al., 2020). Pinna et al. (Pinna et al., 2023) employed the NLP approach (Maldonado et al., 2017) to detect Design SATDs and Requirement SATDs in Blockchain projects. OBrien et al. (OBrien et al., 2022) manu-

ally classified the extracted SATDs into six general categories, which were previously introduced by Bavota and Russo (Bavota and Russo, 2016): Requirement, Code, Test, Defect, Design, and Documentation. Then, they defined 23 specific categories for machine learning-related SATD and manually applied them to their dataset. Bhatia et al. (Bhatia et al., 2023) employed the text-mining tool (Liu et al., 2018) to identify SATDs in 318 ML and 318 non-ML open-source software projects. They then manually classified 611 randomly sampled SATDs based on the categories introduced by Bavota and Russo (Bavota and Russo, 2016).

2.3 Large Language Models

LLMs are deep learning models (Goodfellow et al., 2016), typically based on the Transformer architecture (Vaswani et al., 2017), and pre-trained on extensive corpora. One of the early milestones in LLMs was BERT (Bidirectional Encoder Representations from Transformers) (Devlin et al., 2019), which pioneered the use of bidirectional training to enhance the understanding of word context within sentences, significantly advancing performance across a range of natural language processing tasks. Following BERT’s encoder-only architecture, which excelled mainly in text classification, a new wave of LLMs emerged. These newer models adopted encoder-decoder (Raffel et al., 2019) or decoder-only architectures (Radford et al., 2019), catalyzing a revolution in text generation tasks. Concurrently, variations of the BERT model were introduced, such as CodeBERT (Feng et al., 2020), designed for understanding and generating code by blending natural language and programming language training, and RoBERTa (Liu et al., 2019), an optimized version of BERT with improved training methodology and larger datasets.

T5 (Text-to-Text Transfer Transformer) (Raffel et al., 2019) and Flan-T5 (Chung et al., 2022) are notable examples of LLMs utilizing the encoder-decoder architecture. T5, developed by Google, treats every language task as a text-to-text problem, converting tasks like translation, summarization, and question-answering into a unified framework. Flan-T5, an extension of T5, further enhances its capabilities through fine-tuning with a mixture of instruction-based tasks, improving its performance on various benchmarks. Both models are available in five different sizes, ranging from the small size with 77M parameters to the XXL size with 11.1B parameters. This variety enables researchers to study the effect of model size on their applications. Flan-T5 has been employed in various applications such as text summarization (Tam et al., 2023) and log parsing (Jiang et al., 2023) tasks and has demonstrated impressive performance, even when compared to some newer models like LLaMA (Touvron et al., 2023).

There are two general approaches for applying LLMs on downstream tasks: 1) *fine-tuning*: continuously train the LLM on the downstream task for a few more epochs, and 2) *in-context learning (ICL)*: instructing the model on what it is expected to do, and if necessary, adding a few examples in the instruction

to make the task clearer. The process of finding a good instruction to achieve good results is called *prompt engineering*. The main advantage of ICL over fine-tuning is that it doesn't require investing time and processing resources to update model parameters for the downstream task. In this study, we consider both approaches for the identification and classification of SATDs.

LLMs have been successfully employed in different Software Engineering (SE) tasks. Prenner and Robbes (Prenner and Robbes, 2022) applied BERT-based models on a selection of 13 smaller datasets from the SE literature, and achieved superior performance for tasks involving natural language. Gao et al. (Gao et al., 2023) leveraged the OpenAI's Codex and GPT-3.5 models using the ICL approach for code intelligence tasks including code summarization, bug fixing, and program synthesis. These are just two examples in LLM4SE (Large Language Models for Software Engineering) domain. Hou et al. (Hou et al., 2023) conducted a systematic literature review on LLM4SE by studying 229 research papers from 2017 to 2023. The success of LLMs in various SE-related tasks motivated us to investigate their effectiveness in SATD tasks.

3 Study Setup

3.1 Research Questions

RQ1: How effective are fine-tuned LLMs in SATD identification and classification? Previous studies have explored rule-based methods, as well as machine learning and deep learning approaches, for identifying and classifying SATD. A notable recent study by Prenner and Robbes (Prenner and Robbes, 2022) employed the BERT model for several classification tasks in software engineering, including SATD identification, and achieved significant improvements over traditional methods. The success of BERT in this domain, combined with the advent of more sophisticated, larger, open-source language models like T5 and LLaMA, has inspired our investigation of their effectiveness for SATD identification and classification. Therefore, our first research question explores whether these newer and more advanced language models, known for their proficiency in text generation tasks such as summarization and translation, can also outperform BERT in the efficient identification and classification of SATD. In our experiments for RQ1, we employed the Flan-T5 models, BERT, and CodeBERT as the chosen LLMs for fine-tuning and compared their results with non-LLM baselines. The rationale for selecting these models is presented in Section 3.3.

RQ2: Does our proposed ICL with a larger model outperform smaller models that have been fine-tuned in identifying and classifying SATD? While RQ1 focuses on fine-tuning pre-trained models, this research question explores whether ICL using a larger LLM can surpass the performance of smaller models that have been fine-tuned for SATD identification and classification. ICL is notably cost-effective during the learning phase, just requiring

prompt engineering, making it an appealing approach, especially when training data is limited. To assess the efficacy of the ICL approach, we utilize Flan-T5-XXL with 11.1B parameters, which is approximately four times larger than the biggest model fine-tuned in this study, namely, Flan-T5-XL with 2.85B parameters.

RQ3: What is the impact of adding the classification layer in fine-tuning LLMs? In order to employ the Flan-T5 models for SATD identification and classification, as a common practice we refine their architecture by substituting the original text generation layer with a classification layer to better adapt the model for classification tasks (see Section 3.6). However, the impact of this modification is unclear. Furthermore, some studies utilize the original LLM architecture and treat the classification task as a text-to-text problem (Raffel et al., 2019), or perform the classification through inference (Zhang et al., 2023), specifically by predicting the next token (which is expected to be a class name) when provided with input data as a prompt. To assess the impact of this architectural modification, we conduct fine-tuning experiments using the original Flan-T5 models, in contrast to RQ1 and RQ4, for which we use the modified architecture.

RQ4: What is the impact of additional contextual features on LLM-based SATD classification? In RQ1-RQ3, the models are given only the code comments - the default setting for existing SATD identification and classification approaches. RQ4 aims to explore whether including additional contexts, such as the file path, the containing method’s signature, and the body of the containing method, enhances the performance of large language models in the SATD classification task. A key aspect of this investigation is determining the optimal size of the LLM required to utilize this contextual data for improved results effectively. Specifically, we want to assess whether smaller LLMs can discern patterns and relationships between the code comments and the contextual features to enhance SATD classification performance. We investigate this question solely for SATD classification because the largest benchmark for SATD identification, the Maldonado dataset (Maldonado et al., 2017), does not provide the location and contextual information for each entry.

3.2 Datasets

In this study, we employ two datasets: Maldonado-62k and OBrien. The usage of these datasets is as follows: for RQ1, RQ2, and RQ3, we utilize both datasets. For RQ4, we only use the OBrien dataset, because, unlike the Maldonado-62k dataset, which only provides the comment text for each entry, the OBrien dataset includes additional features such as the file path that can be utilized by LLMs to achieve better performance. More details for these two datasets are provided below.

Table 1 Maldonado-62k dataset statistics

Project	#Comments	#SATD (original)	#SATD (Jitterbug)
ApacheAnt	4,098	131	135
ArgoUML	9,452	1,413	1,630
Columba	6,468	204	220
EMF	4,390	104	119
Hibernate	2,968	472	493
JEdit	10,322	256	259
JFreeChart	4,408	209	247
JMeter	8,057	374	416
JRuby	4,897	622	665
Squirrel	7,215	286	313
Total	62,275	4,071	4,497

Maldonado-62k: Initially, Maldonado et al. (Maldonado and Shihab, 2015) compiled a SATD dataset by manually labeling 33,093 source code comments from five well-documented open-source projects across various application domains, specifically Apache Ant, Apache JMeter, ArgoUML, Columba, and JFreeChart. This dataset was later expanded in (Maldonado et al., 2017) to include 62,275 labeled code comments from 10 projects, among which 4,071 were identified as instances of SATD. Subsequently, Yu et al. (Yu et al., 2022) conducted a meticulous review of all comments initially classified as non-SATD, particularly those containing strong SATD indicators, i.e., `todo`, `fixme`, `hack`, and `workaround`. They discovered that 426 of the 434 comments, previously labeled as non-SATD, were indeed SATD. As a result, the revised Maldonado dataset comprises 4,497 SATD and 57,778 non-SATD code comments. In our study, we utilize this modified version of the Maldonado dataset, which we refer to as the Maldonado-62k dataset. The dataset comprises three fields: “project”, “comment”, and “label”. However, it lacks contextual details like the surrounding code, commit history, or file path, as these were not included by the dataset creators. Table 1 shows the statistics of this dataset.

OBrien: OBrien et al. (OBrien et al., 2022) compiled a dataset from popular machine learning (ML) repositories written in Python, including scikit-learn and IntelPython. This dataset comprises 856 labeled instances of SATD identified within these repositories. The labeling process introduced various layers, such as “Software TD Type”, “ML TD Type”, and “ML Pipeline Stage”. However, our study focuses on the general classification of SATD, so we selected only the ‘Software TD Type’ column as our target label, omitting the others. We refined this dataset by eliminating entries with null values in the “Software TD Type” column, resulting in a total of 789 entries. The labels’ descriptions and frequencies are detailed in Table 2. A notable feature of this dataset is the inclusion of contextual information for each entry, such as the file path and the commit introducing the SATD. Leveraging this data, we added two contextual columns: “containing method signature” and “containing method body”. Among 789 entries in this dataset, 447 SATD instances are located

Table 2 Definitions and examples of the six SATD types identified in OBrien dataset (OBrien et al., 2022)

SATD Types	Description	Example Comment	#of Occ.
Requirement	Requirement debts can be functional or non-functional. In the functional case, implementations are left unfinished or in need of future feature support. In the non-functional case, the corresponding code does not meet the requirement standards (speed, memory usage, security, etc...).	TODO: handle channel modalities later, TODO: make efficient, TODO: Implement Conv Transpose.	321
Code	Bad coding practices leading to poor legibility of code, making it difficult to understand and maintain.	TODO: This next code is dense and confusing. Clean up at some point.	207
Test	Problems found in implementations involving testing or monitoring sub-components.	XXX: should we rather test if instance of estimator?	84
Defect	Identified defects in the system that should be addressed.	TODO this will fail if a parameter cant handle size=(N:)	82
Design	Areas which violate good software design practices, causing poor flexibility to evolving business needs.	TODO maybe improve this so it doesn't use a global	80
Documentation	Inadequate documentation that exists within the software system.	TODO update doc above	15

either within a method or on the very first line preceding a method, allowing us to associate each with its corresponding method.

There are two primary reasons for selecting the OBrien dataset for the SATD classification task. First, we required a dataset comprising self-admitted technical debts introduced in code comments. This is to access more advanced data, such as the surrounding code, which is crucial for delivering more accurate predictions of SATD types. Therefore, we are not considering datasets focused on other sources like issue trackers (Xavier et al., 2020) or build systems (Xiao et al., 2021). Second, among the datasets that categorize SATD in code comments, only the OBrien dataset (OBrien et al., 2022) offers additional information, namely the file path and commit hash. This information is essential for accessing the contextual details for each code comment. In contrast, the Maldonado datasets (Maldonado and Shihab, 2015; Maldonado et al., 2017) and other datasets, such as (Fucci et al., 2021; Cassee et al., 2022), which are derived from the Maldonado dataset, lack such data, making it challenging to extract the surrounding code for each code comment.

Both datasets present challenges. In the Maldonado-62k dataset, out of 62,275 code comments, only 4,497 items (7.2%) are classified as SATD. In the OBrien dataset, there are only 789 items. Therefore, these two datasets also challenge the effectiveness of large language models in dealing with imbalanced data (for SATD identification) and the scarcity of labeled data (for SATD classification).

3.3 Selected Large Language Models

For RQ1 and RQ4, we consider the following models:

- **BERT-base-uncased** (110 million parameters) (Devlin et al., 2019): This model has demonstrated superior performance in the SATD identification task (Prenner and Robbes, 2022), establishing a strong baseline for other LLMs in SATD identification and classification tasks.
- **CodeBERT** (125 million parameters) (Feng et al., 2020): Differing from BERT, which primarily focuses on natural language, CodeBERT is a bi-modal extension that utilizes both natural language and source codes. This makes it particularly suited for processing code comment SATD.
- **Flan-T5** (Chung et al., 2022): This is an enhanced version of T5 (Rafael et al., 2019), fine-tuned on over 1000 tasks. Available in various sizes, Flan-T5 allows for an exploration of the impact of model size on SATD tasks. Our infrastructure supports full fine-tuning of all Flan-T5 variants except the XXL version (11.1 billion parameters). Therefore, we will include Flan-T5 Small (77 million parameters), Base (248 million parameters), Large (783 million parameters), and XL (2.85 billion parameters) in our experiments.

For RQ2, we will utilize Flan-T5-XXL, the largest variant with 11.1 billion parameters. This model does not require fine-tuning for our purposes, as we will employ the ICL approach for the SATD identification and classification tasks. For RQ3, we evaluate the performance of the original Flan-T5 architecture against their altered version utilized in RQ1 and RQ4.

3.4 Baselines

For the task of SATD identification, we employ four baseline methods: the NLP approach by (Maldonado et al., 2017), the MAT method (Guo et al., 2021), the technique introduced in the PENTACET study (Sridharan et al., 2023), and the convolutional neural networks (CNN) model proposed by (Ren et al., 2019). For the SATD classification task, we utilize two baselines: the NLP approach (Maldonado et al., 2017) and the CNN model (Ren et al., 2019).

3.5 Evaluation metrics

Following the previous studies (Prenner and Robbes, 2022; Ren et al., 2019; Maldonado et al., 2017), we report the F1 score for the SATD class in the SATD identification task. For the SATD classification task, we report overall performance using accuracy, and per class performance using the F1 score.

3.6 Updating Flan-T5 architecture for the classification task

In this study, we propose an adaptation of the Flan-T5 architecture to address the classification task (see Figure 1). The original Flan-T5 model is inherently designed for text generation tasks, where the goal is to predict the probability

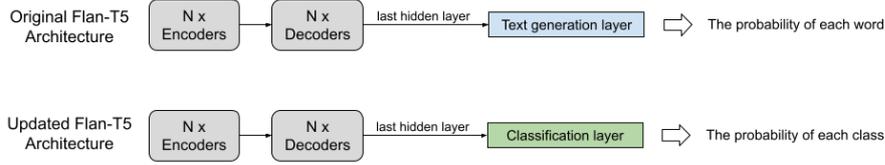


Fig. 1 Updating Flan-T5 architecture by replacing the last layer with a classification layer

distribution of subsequent words in a given text sequence. This is achieved through the model’s final layer, which is a text generation layer that outputs probabilities for each word in the vocabulary.

To tailor this architecture to classification objectives, we replace the text generation layer with a classification head, which consists of a dropout layer followed by a linear neural network layer with the size of number of classes. This new layer interprets the semantic representation of the input sequence from the last hidden state and produces a probability distribution over the predefined classes. For the loss function, we use cross-entropy which is a common choice for classification tasks. This modification leverages the pre-existing language understanding capabilities of the Flan-T5 model while aligning the output structure with the requirements of classification tasks. The updated architecture now effectively assigns input sequences to categories. We use this architecture in RQ1 and RQ4. For RQ2, we use the original architecture because for ICL approaches, there is no need to update the model parameters. In RQ3, we evaluate the effectiveness of the modified architecture against the original one.

3.7 Implementation

We begin by downloading the official checkpoints for all evaluated models from Hugging Face. Subsequently, the Torch and Transformers packages are employed to conduct fine-tuning and ICL on a single Nvidia RTX 6000 GPU with 48 GB VRAM. In all models, we use the default 32-bit precision mode, and all experiments are done in Python 3.10.

For all baseline models, we employed their default settings. For selected LLMs, we adjusted the batch size to the highest power-of-two value that our VRAM could support. Additionally, since the learning rate significantly influences LLM performance, we conducted a few experiments to determine a near-optimal learning rate for each LLM. In the CNN and all LLMs, we configured the maximum length of input data for the Maldonado-62k dataset at 128 tokens. For the OBrien dataset, we set it at 512 tokens. The rationale for the higher token limit in the OBrien dataset is related to RQ4, where we explore the impact of incorporating contextual data into the input. This necessitates additional space to include this information for the models. All large language models are trained for eight epochs using the AdamW (Loshchilov and Hut-

Table 3 The key training parameter settings for each LLM applied to each dataset. All models are trained for 8 epochs.

Model	Dataset	Maxlen	Batch size	Learning rate
BERT-base-uncased (110M)	Maldonado-62k	128	32	0.00001
CodeBERT-base (125M)	Maldonado-62k	128	32	0.00001
Flan-T5-small (77M)	Maldonado-62k	128	32	0.0001
Flan-T5-base (248M)	Maldonado-62k	128	32	0.0001
Flan-T5-large (783M)	Maldonado-62k	128	16	0.0001
Flan-T5-XL (2.85B)	Maldonado-62k	128	4	0.00002
BERT-base-uncased (110M)	OBrien	512	32	0.00005
CodeBERT-base (125M)	OBrien	512	32	0.00005
Flan-T5-small (77M)	OBrien	512	32	0.001
Flan-T5-base (248M)	OBrien	512	16	0.0005
Flan-T5-large (783M)	OBrien	512	4	0.0002
Flan-T5-XL (2.85B)	OBrien	512	1	0.00005

ter, 2019) optimizer with a linear learning rate decay schedule. As there is no validation set to choose the best model across training epochs, we take the resulting model of the last epoch and report its performance on the test data. Table 3 presents the key parameter settings for each model applied to each dataset.

4 Results and Analysis

4.1 How effective are fine-tuned LLMs in SATD identification and classification?

4.1.1 Approach

SATD Identification. The rule-based baselines, MAT and PENTACET, do not require training data. Therefore, they are run on the entire dataset at once. In contrast, other baselines and the LLMs necessitate training data. In line with previous studies (Prenner and Robbes, 2022; Ren et al., 2019), a cross-project approach is employed to evaluate these models. As the Maldonado dataset contains 10 projects, each round involves selecting one project as the test project and using the remaining nine for model training.

SATD Classification. A similar method is applied to the OBrien dataset, but with 10 randomly assigned folds instead of cross-project validation. In each round, one fold is used for testing, and the remaining nine for training. The data is split into 10 folds only once, and this division is used across all experiments on this dataset. Each fold comprises 79 data items, except the last one, which contains 78. Given the relatively small size of the OBrien dataset, LLM models may show varying performance in different runs. Therefore, each

experiment is conducted three times with different seeds, and the average accuracy or F1 score is reported.

4.1.2 Results

SATD Identification. Table 4 shows the F1 score of all selected large language models for each of the 10 Java projects against the four baselines, NLP (Maldonado et al., 2017), CNN (Ren et al., 2019), MAT (Guo et al., 2021), and PENTACET (Sridharan et al., 2023). The final row, which presents the average F1 score across the 10 projects, demonstrates that all six large language models significantly outperform the four baselines, with margins ranging from 4.4% to 11.3%. Moreover, in terms of F1 score per project, the top-performing model is an LLM for all projects except Ruby, where MAT achieves an F1 score of 95.1%, marginally higher (by 0.1%) than the scores of Flan-T5 small and large models. Among all models, Flan-T5-XL delivers the highest average score, despite not achieving the best score in every project. The results from the large language models further suggest that larger models generally surpass their smaller counterparts, supporting the notion that larger models tend to yield better results.

The Flan-T5-XL model is 37 times larger than the smallest model, Flan-T5-small, but only slightly outperforms it, with an F1 score improvement of just 2.1% (83.9% vs 81.8%). This modest improvement is surprising considering the model’s size. One significant challenge in model performance is the quality of annotations. To delve deeper, we analyzed Cohen’s Kappa coefficient, comparing Flan-T5-XL’s predictions with the ground truth. The coefficient is 0.89, higher than the 0.81 inter-rater agreement reported in the study, indicating that our Flan-T5-XL model is more accurate than human annotators. This outcome may appear counter-intuitive, as the model was trained on imperfect human labels. However, Flan-T5-XL is not a simple model; it’s extensively pre-trained on a large text corpus, enabling it to understand the meaning of words and phrases. The capability, obtained by combining its pre-training and fine-tuning, likely helps it to identify and ignore obvious mislabeling and make more accurate predictions. Indeed, improving the quality of annotated data is a prerequisite for further performance improvement.

SATD Classification. Table 5 presents the per-category F1 scores and overall accuracy for six LLMs and two baselines, NLP and CNN, for the SATD classification task. In contrast to the SATD identification task, where smaller models like Flan-T5-small and BERT-base performed competitively with larger models, a more pronounced difference is observed in this table. Here, the larger models significantly outperform the smaller ones. The Flan-T5-XL achieves the highest accuracy at 0.62, while the Flan-T5-small registers the lowest at 0.537. This disparity in performance can be attributed to two main factors. Firstly, the training dataset for SATD identification task consists of approximately 50,000 to 60,000 items, varying based on the selected project as the test data. In comparison, the training dataset for SATD classification is considerably

Table 4 Comparative F1 scores of six large language models and baselines in SATD identification on Maldonado-62k dataset. The “cross-project” row indicates if we employ cross-project prediction, i.e., using 9 projects for training and the remaining 1 for testing.

	NLP	MAT	PENT-ACET	CNN	BERT	CodeBERT	Flan-T5	Flan-T5	Flan-T5	Flan-T5
Cross-project	yes	no	no	yes	base yes	base yes	small yes	base yes	large yes	XL yes
ApacheAnt	0.532	0.654	0.505	0.625	0.629	0.691	0.643	0.679	0.689	0.710
ArgoUML	0.936	0.934	0.894	0.946	0.954	0.954	0.956	0.951	0.958	0.939
Columba	0.833	0.928	0.892	0.88	0.857	0.941	0.923	0.919	0.944	0.933
EMF	0.527	0.434	0.574	0.420	0.746	0.651	0.660	0.693	0.707	0.733
Hibernate	0.852	0.851	0.859	0.892	0.902	0.893	0.895	0.901	0.913	0.922
JEdit	0.470	0.321	0.654	0.605	0.631	0.668	0.688	0.698	0.727	0.717
JFreeChart	0.705	0.707	0.704	0.739	0.734	0.738	0.746	0.735	0.741	0.734
JMeter	0.819	0.870	0.795	0.868	0.879	0.878	0.874	0.883	0.883	0.882
JRuby	0.896	0.951	0.928	0.919	0.947	0.940	0.950	0.948	0.950	0.945
Squirrel	0.689	0.731	0.776	0.777	0.826	0.842	0.846	0.834	0.834	0.869
Average	0.726	0.738	0.758	0.767	0.811	0.820	0.818	0.824	0.835	0.839

smaller, with only 710 items. Secondly, the nature of the tasks differs: SATD identification is just a binary classification, whereas SATD classification requires a more complex 6-class classification. This additional complexity poses a greater challenge, particularly for smaller models. Consequently, the larger pre-trained models, which exhibit a more profound understanding of textual data, demonstrate superior efficacy in learning and performing the classification task compared to their smaller counterparts.

Table 5 also shows that the performance of NLP is lower than even the smallest LLM, while CNN achieves a commendable overall accuracy of 0.602, close to the best-performing LLM, Flan-T5-XL, which has an accuracy of 0.620. This effectiveness of CNN led to its selection as a strong baseline for our last research question.

Among six categories in the OBrien dataset, Flan-T5-XL scores the highest F1 in three categories and the second-highest in two of the remaining three categories. A noteworthy observation is the “Document” class, containing only 15 items, where all LLMs show low F1 scores, ranging between 0.0 and 0.23. However, Flan-T5-XL distinguishes itself with a relatively high F1 score of 0.488 in this class, indicative of its consistent performance across varied class sizes. While CNN performs best in the “Document” class, its performance is not uniformly high across all categories. For instance, it scored second-worst in the “Defect” category.

Summary: In SATD identification, the selected LLMs outperformed all existing baselines, with Flan-T5-XL achieving the highest F1 score, a 7.2% improvement compared to the best-performing baseline. In the SATD classification task, while the fine-tuned Flan-T5-XL still led in performance, the CNN model exhibited competitive results, even surpassing four of six LLMs. In both tasks, the larger models outperformed their smaller counterparts.

Table 5 Average F1 score for each category and the overall accuracy (Dataset: OBrien, Approach: 10-fold cross validation, number of runs: 3)

Class label Count	Code 207	Defect 82	Design 80	Doc. 15	M&T 84	Requirement 321	All 789
NLP	0.512	0.290	0.256	0.300	0.549	0.632	0.527
CNN	0.627	0.247	0.357	0.545	0.611	0.692	0.602
BERT-base-uncased (110M)	0.585	0.367	0.312	0.170	0.634	0.662	0.576
CodeBERT-base (125M)	0.636	0.398	0.423	0.176	0.610	0.687	0.611
Flan-T5-small (77M)	0.533	0.120	0.074	0.0	0.600	0.659	0.537
Flan-T5-base (248M)	0.600	0.378	0.299	0.0	0.593	0.654	0.565
Flan-T5-large (783M)	0.586	0.418	0.329	0.230	0.598	0.662	0.575
Flan-T5-XL (2.85B)	0.651	0.433	0.475	0.488	0.618	0.684	0.620

4.2 Does our proposed ICL with a larger model outperform smaller models that have been fine-tuned in identifying and classifying SATD?

4.2.1 Approach

In both tasks, we apply the ICL method to the largest model in the Flan-T5 family, Flan-T5-XXL, and compare its performance with that of the fine-tuned smaller versions from the same family. To perform the inference, we set *max_new_tokens* = 5 and *temperature* = 0.0. If none of the class names appear in the model’s response (indicating unknown labels), the majority class, i.e., “non-SATD” for Maldonado-62k dataset and “Requirement” for OBrien dataset, is assumed to be the predicted class.

SATD Identification. For SATD identification, we start with employing four distinct zero-shot prompts, each providing a definition of SATD and incorporating various keywords indicative of SATD, such as TODO and FIXME. This approach is based on findings by Guo et al. (2021) and Yu et al. (2022), which highlight the significance of these keywords in classifying a code comment as SATD. The prompt for each selected group of keywords is shown in Table 6. Once we find the best performing zero-shot prompt, we investigate the effect of incorporating different numbers of examples (1, 2, 3, 5, 10, 15, 20) in that prompt. Two methods are used to select samples for prompts: 1) Randomly choosing *n* examples from the training projects for each item in the test project, and 2) Using SentenceTransformer (with all-MiniLM-L6-v2 model) to embed the input data and applying cosine similarity to select the most relevant items from the training projects for each item in the test project.

SATD Classification. For SATD classification, we explore the impact of varying the number of examples (0, 1, 2, 3, 5, 10, 15, 20) in the prompt on ICL performance. Same as the SATD identification approach, two methods are used to select samples for prompts: 1) Randomly choosing *n* examples from the training folds for each item in the test fold, and 2) Using SentenceTransformer (with all-MiniLM-L6-v2 model) to embed the input data and applying cosine similarity to select the most relevant items from the training folds for each test

Table 6 The prompts applied in ICL approach on Maldonado-62k dataset

Approach	Prompt
No keywords	Self-admitted technical debt (SATD) is technical debt admitted by the developer through source code comments. Assign the label of SATD or Not-SATD for each given source code comment.
Suggested keywords by MAT	Self-admitted technical debt (SATD) are technical debt admitted by the developer through source code comments. SATD comments usually contain specific keywords: TODO, FIXME, HACK, and XXX. Assign the label of SATD or Not-SATD for each given source code comment.
Suggested keywords by Jitterbug	Self-admitted technical debt (SATD) is technical debt admitted by the developer through source code comments. SATD comments usually contain specific keywords: TODO, FIXME, HACK, and Workaround. Assign the label of SATD or Not-SATD for each given source code comment.
Suggested keywords by GPT4	Self-admitted technical debt (SATD) is technical debt admitted by the developer through source code comments. SATD comments usually contain specific keywords: TODO, FIXME, HACK, XXX, NOTE, DEBT, REFACTOR, OPTIMIZE, TEMP, WORKAROUND, KLUDGE, REVIEW, NOFIX, PENDING, and BUG. Assign the label of SATD or Not-SATD for each given source code comment.

fold item. Additionally, we explore whether including category descriptions in the prompt enhances classification performance. Table 7 presents a sample one-shot prompt for each prompting approach.

4.2.2 Results

SATD Identification. The result for each zero-shot prompting approach is shown in Table 8. The highest precision and F1 score were achieved by the prompt incorporating keywords from the MAT approach (Guo et al., 2021), while the best recall was obtained without including any keywords. The MAT keyword-based prompt attained an F1 score of 0.747, exceeding the results of the NLP and MAT approaches, but falling short of the CNN approach and significantly lagging behind all fine-tuned LLMs (see Table 4). Table 9 presents the F1 scores for the few-shot approach when the MAT keywords are incorporated into the prompt. The table also shows the count of unknown labels (instances where the model’s generated tokens do not correspond to any label), noted in parentheses beneath each F1 score. Notably, the few-shot method, whether using random samples or relevant examples by Sentence-Transformer, performs worse than the zero-shot approach. This is likely due to data imbalance and the inadequacy of methods in selecting suitable examples for the prompt in this task. According to these results, in the context of SATD identification, the most effective ICL method, which is the zero-shot prompt with MAT keywords, does not outperform even the smallest fine-tuned Flan-T5 model, despite utilizing the largest Flan-T5 model.

SATD Classification. Table 10 displays the accuracy along with the count of unknown labels noted in parentheses beneath each accuracy value. Accord-

Table 7 Sample one-shot prompts for SATD classification task

Prompt approach	Sample one-shot prompt
Describing categories + The most relevant examples using SentenceTransformer	<p>There are six types of software technical debts:</p> <p>Requirement: Requirement debts can be functional or non-functional. In the functional case, implementations are left unfinished or in need of future feature support. In the non-functional case, the corresponding code does not meet the requirement standards (speed, memory usage, security, etc...).</p> <p>Code: Bad coding practices lead to poor legibility of code, making it difficult to understand and maintain.</p> <p>M&T: Problems found in implementations involving testing or monitoring subcomponents.</p> <p>Defect: Identified defects in the system that should be addressed.</p> <p>Design: Areas which violate good software design practices, causing poor flexibility to evolving business needs.</p> <p>Documentation: Inadequate documentation that exists within the software system.</p> <p>Here are some examples:</p> <pre>### Technical debt comment: "" "" TODO normalize to make sum up to 1? "" "" ### Label: Requirement ### Technical debt comment: "" "" self.mpc.sum(3; -5) TODO: Future work: how to handle gracefully minus numbers "" "" ### Label:</pre>
Describing categories + Some random examples	<p>There are six types of software technical debts:</p> <p>[describing categories the same as in the above prompt]</p> <p>Here are some examples:</p> <pre>### Technical debt comment: "" "" TODO: delete this method when no longer needed "" "" ### Label: Code ### Technical debt comment: "" "" self.mpc.sum(3; -5) TODO: Future work: how to handle gracefully minus numbers "" "" ### Label:</pre>
relevant examples using SentenceTransformer	<pre>### Technical debt comment: "" "" TODO normalize to make sum up to 1? "" "" ### Label: Requirement ### Technical debt comment: "" "" self.mpc.sum(3; -5) TODO: Future work: how to handle gracefully minus numbers "" "" ### Label:</pre>

Table 8 Average precision, recall, and F1 score over 10 projects, obtained by the zero-shot ICL approach using Flan-T5-XXL (11.1B) on Maldonado-62k dataset

Prompt approach	P	R	F1
Zero-shot prompt with no keywords	0.212	0.889	0.329
Zero-shot prompt with suggested keywords by MAT	0.741	0.762	0.747
Zero-shot prompt with suggested keywords by Jitterbug	0.678	0.76	0.712
Zero-shot prompt with suggested keywords by GPT4	0.585	0.791	0.664

Table 9 Average F1 score over 10 projects, obtained by the few-shot ICL approach using Flan-T5-XXL (11.1B) on Maldonado-62k dataset. Note: Counts of unknown labels (instances where generated tokens do not match any label) are shown in parentheses below each F1 score.

Prompt approach	Number of examples in the prompt							
	0	1	2	3	5	10	15	20
Mentioning MAT keywords + Some examples by SentenceTransformer	0.747 (0)	0.687 (0)	0.673 (0)	0.661 (0)	0.662 (0)	0.653 (0)	0.650 (0)	0.648 (0)
Mentioning MAT keywords + Some random examples	0.747 (0)	0.707 (0)	0.679 (0)	0.665 (6)	0.660 (6)	0.649 (6)	0.651 (6)	0.649 (6)

Table 10 Average accuracy over 10 folds, obtained by the ICL approach using Flan-T5-XXL (11.1B) on the OBrien dataset. Note: Counts of unknown labels (instances where generated tokens do not match any label) are shown in parentheses below each F1 score.

Prompt approach	Number of examples in the prompt							
	0	1	2	3	5	10	15	20
Describing categories + Some examples by SentenceTransformer	0.506 (12)	0.530 (0)	0.560 (0)	0.556 (0)	0.567 (0)	0.572 (0)	0.561 (0)	0.536 (0)
Describing categories + Some random examples	0.506 (12)	0.511 (0)	0.513 (0)	<u>0.521</u> (0)	0.518 (0)	0.510 (0)	0.493 (0)	0.450 (0)
Just use SentenceTransformer to mention some examples	-	0.456 (168)	0.502 (93)	<u>0.520</u> (45)	0.516 (9)	0.477 (0)	0.496 (1)	0.470 (0)

ing to the table, unknown labels tend to occur when no examples or category definitions are provided in the prompt.

Key findings from Table 10 include:

- Providing relevant examples consistently yields better performance compared to random examples.
- Including category descriptions before examples in the prompt leads to significantly improved performance, regardless of the number of examples provided in the prompt.
- Flan-T5-XXL’s ICL approach does not outperform fine-tuning Flan-T5-XL or Flan-T5-large, though it does surpass the performance of fine-tuned Flan-T5-base and Flan-T5-small.

Summary: In the ICL approach using the Flan-T5-XXL model for the SATD identification task, the zero-shot approach provides competitive results with traditional approaches, but performs 6.4% to 9.2% worse than fine-tuned LLMs. In contrast, for the SATD classification task, incorporating relevant examples and category descriptions into the prompts enhances performance, surpassing the results achieved by fine-tuning the Flan-T5 small and base models. However, it falls behind the results obtained by fine-tuning the Flan-T5 large and XL models.

Table 11 Comparison of average F1 score: original vs. modified Flan-T5 architecture on the Maldonado-62k dataset using comment text as input

Model	Original architecture	Adding the classification layer
Flan-T5-small	0.820	0.818
Flan-T5-base	0.820	0.824
Flan-T5-large	0.835	0.835
Flan-T5-XL	0.831	0.839

4.3 What is the impact of adding the classification layer in fine-tuning LLMs?

4.3.1 Approach

SATD Identification. As described in Section 3.6, we have refined the architecture of the Flan-T5 model by substituting its original text generation layer with a classification layer to better adapt the model for classification tasks. To assess the impact of this architectural modification, we conducted some experiments using the original Flan-T5 model, and inputting only the comment text. To perform the inference, similar to RQ2 approach, we set $max_new_tokens = 5$ and $temperature = 0.0$. If none of the class names appear in the model’s response, the majority class is assumed to be the predicted class.

SATD Classification. Same as above.

4.3.2 Results

SATD Identification. Table 11 presents the comparison between the original model and the modified version with the added classification layer. We observe that the two architectures achieve nearly identical performance.

SATD Classification. Table 12 presents the comparison between the original model and the modified version for the classification task. In contrast to the outcomes observed for the identification task, the modified architecture demonstrates significant improvement over the original in the classification task, particularly when employing the smaller versions of the Flan-T5 models. We believe this enhancement can be attributed to the lesser amount of available training data for the classification task, compared to the identification task.

Summary: Substituting the word generation layer with a classification layer in Flan-T5 models, although it has no benefits for the SATD identification task, does increase their performance for the SATD classification task, especially when utilizing the smaller versions of the Flan-T5 models. This improvement is likely due to the limited available training data for the SATD classification task.

Table 12 Comparison of average accuracy: original vs. modified Flan-T5 architecture on the OBrien dataset using comment text as input

Model	Original architecture	Adding the classification layer
Flan-T5-small	0.447	0.537
Flan-T5-base	0.505	0.565
Flan-T5-large	0.561	0.575
Flan-T5-XL	0.602	0.620

4.4 What is the impact of additional contextual features on LLM-based SATD classification?

4.4.1 Approach

SATD Classification. In RQ1 to RQ3, we evaluated large language models on SATD identification and classification tasks, using code comments as input. While the Maldonado-62k dataset lacks additional features, the OBrien dataset includes informative features that can enhance model accuracy in predicting SATD categories. Undoubtedly, the code comment itself is the most significant feature in determining the category of an SATD. However, contextual features like file path or surrounding code can aid in a more comprehensive understanding of a code comment. For instance, the presence of the word ‘test’ in the file name or path may increase the likelihood of a code comment being related to a test SATD. Therefore, in this research question, we aim to investigate how effectively LLMs leverage additional contextual features for SATD classification. More specifically, RQ4 compares models using four distinct combinations of input data: 1) just the comment text (with results adopted from RQ1 and RQ2); 2) the file path along with the comment text; 3) the file path, the containing method’s signature, and the comment text; 4) the file path, comment text, and the entire containing method, which includes both the signature and body. Given that the containing method can be extensive, it is positioned at the end of the context. This arrangement ensures that in instances of long input data leading to truncation, the critical comment text is preserved. Table 13 provides an example for each input data combination.

4.4.2 Results

SATD Classification. Table 14 presents the results for the SATD classification task using three different approaches: ICL with Flan-T5-XXL¹, training the CNN and fine-tuning the selected six LLMs, and an ensemble approach by the fine-tuned Flan-T5-XL. For the fine-tune approach, when the input data includes the file path or the file path along with the method’s signature, the two largest models, Flan-T5-large and Flan-T5-XL, significantly improve

¹ To create the prompt for the ICL approach, we followed the same method as presented in RQ2 approach. We selected the best result, which was obtained when we included the category descriptions and the five most relevant examples using the SentenceTransformer.

Table 13 A sample for each combination of input data in OBrien dataset

Input data	Example
Comment text	TODO: Future work: how to handle gracefully minus numbers
File path + Comment text	file path: test/torch_test.py Technical debt comment: TODO: Future work: how to handle gracefully minus numbers
File path + Containing method’s signature + Comment text	file path: test/torch_test.py Containing method signature: “““ test_mpc_sum(self) ””” Technical debt comment: TODO: Future work: how to handle gracefully minus numbers
File path + Comment text + Containing method	file path: test/torch_test.py Technical debt comment: “““ TODO: Future work: how to handle gracefully minus numbers ””” Containing method body: “““ def test_mpc_sum(self): self.mpc_sum(3, 5) self.mpc_sum(4, 0) self.mpc_sum(5, -5) # self.mpc_sum(3, -5) TODO: Future work: how to handle gracefully minus numbers self.mpc_sum(2 ** 24, 2 ** 12) ”””

their performance compared to when we only provide the comment text. In contrast, the ICL approach and the four fine-tuned smaller LLMs and the CNN model show minimal improvement or even deteriorate. Notably, the last column demonstrates how adding both the file path and the entire containing method to the input data complicates the task for all models except Flan-T5-XL. While this input worsens the performance of all models, Flan-T5-XL utilizes it to achieve better results than when using only the comment text. This finding highlights the superior capability of fine-tuning larger models in processing and leveraging complex data for the SATD classification task. The consistent high performance of the Flan-T5-XL model (with scores ranging from 0.62 to 0.648 across different input combinations) led us to explore an ensemble approach. This method assigns the most frequent label across 12 predictions (4 different input data types \times 3 runs = 12). The ensemble approach attained an accuracy of 0.668, the highest among all experiments.

Summary: Larger models like Flan-T5-large and Flan-T5-XL significantly outperform other smaller models when contextual information are included in the input data for the SATD classification task, with Flan-T5-XL showing superior capability by utilizing complex data effectively. An ensemble approach using the Flan-T5-XL model yielded the highest accuracy, demonstrating the benefits of larger models and varied inputs.

Table 14 Average accuracy obtained by different data as input (Dataset: OBrien, Approach: 10-fold cross validation, number of runs: 3) CT: Comment Text, FP: File Path, CMS: Containing Method’s Signature, CM: Containing Method

Approach	Model	Input data			
		CT	FP+CT	FP+CMS+CT	FP+CT+CM
ICL	Flan-T5-XXL (11.1B)	0.572	0.572	0.529	0.459
Training	CNN	0.602	0.580	0.572	0.533
Fine-Tuning	BERT-base-uncased (110M)	0.576	0.578	0.567	0.496
	CodeBERT-base (125M)	0.611	0.606	0.571	0.492
	Flan-T5-small (77M)	0.537	0.520	0.506	0.436
	Flan-T5-base (248M)	0.565	0.560	0.546	0.468
	Flan-T5-large (783M)	0.575	0.599	0.590	0.545
	Flan-T5-XL (2.85B)	0.620	0.648	0.628	0.635
Ensemble	fine-tuned Flan-T5-XL			0.668	

5 Discussion

5.1 Impact of epoch number

To train the selected large language models, we set the number of training epochs to eight. After training, we used the model from the final epoch to test its performance, because there was no validation set available to select the best-performing model. In this subsection, we present the performance of the trained model on the test set across epochs for RQ1.

Figure 2 illustrates the trend of the average F1 score across Maldonado-62k projects from epochs 1 to 8. The figure reveals that all six large language models start with a high performance (approximately 0.81 F1 score) in the first epoch, with subsequent epochs showing only minor fluctuations around this score. This consistency is likely due to the binary classification nature of the SATD identification task and the substantial volume of training data in the Maldonado-62k dataset, allowing these pre-trained models to reach near-optimal performance from the first epoch.

Figure 3 presents the results for the SATD classification task using the OBrien dataset. In contrast to Figure 2, where smaller models performed competitively with larger models, a more pronounced difference is observable in this figure. Additionally, the performance of all LLMs sharply increases from epoch 1 to epoch 5. Beyond epoch 5, the performance either stabilizes with minor fluctuations or gently continues to increase, suggesting that epoch 8 is an optimal point to stop training.

6 Threats to Validity

External validity. Creating a large dataset with completely error-free labels by human effort is impractical, and the Maldonado dataset is no exception. Although some incorrect labels are not entirely wrong but rather due to dif-

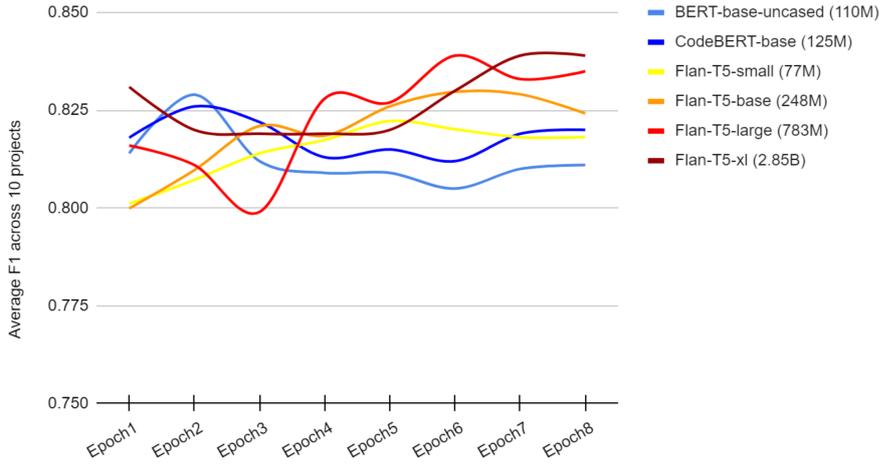


Fig. 2 Average F1 score across 10 projects in the Maldonado-62k dataset over epochs

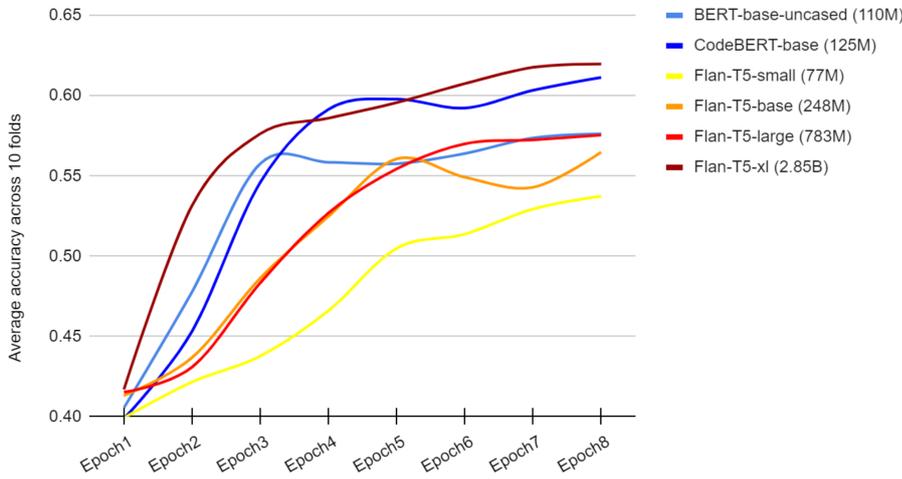


Fig. 3 Average accuracy across 10 folds over epochs (Dataset: O'Brien, Approach: 10-fold cross validation, number of runs: 3)

ferent interpretations of the SATD definition, there are numerous instances in this dataset where annotators would choose the opposite label with high agreement. We attempted to mitigate this problem by applying the label modifications proposed in Yu et al. (2022), as we found their suggestions reasonable. However, all these proposed modifications address false negative cases. There are also false positives in the dataset, a few of which are reported in Cassee et al. (2022). Therefore, we could expect better performance if we had access to a more accurate dataset.

Internal validity. In our experiments, we used only two datasets for SATD identification/classification. There are additional datasets for this purpose, but as previously mentioned, they are either created from other software sources such as issue tracking systems Li et al. (2022) and build systems Xiao et al. (2021), which are beyond the scope of this study, or they are derivatives of the Maldonado dataset Fucci et al. (2021); Cassee et al. (2022), which we have already utilized. However, due to the lack of available contextual features in this dataset, we were unable to analyze the impact of additional features on this dataset. Another concern regarding internal validity is the challenge of finding optimal parameters for fine-tuning large language models. We employed a manual investigation approach; however, considering the substantial computational resources required by larger LLMs, thoroughly exploring a wide range of parameter settings is time-consuming. Therefore, we primarily focused on identifying an effective learning rate through manual trials.

Construct validity. Regarding SATD identification with the ICL approach (RQ2), we did not find an effective method that improves performance by incorporating examples in prompts compared to simply including SATD-related keywords. There may be approaches that can efficiently select relevant examples from the training set for the specific task of SATD identification. Another potential threat to the construct validity of our study is the selection of BERT, CodeBERT, and Flan-T5 models for our experiments. There are many more open-source large language models that could have been chosen for this study. The rationale behind selecting Flan-T5 was to investigate the impact of model size on our tasks. Additionally, we faced infrastructure limitations in fine-tuning other models such as LLaMA. Although LLaMA is available in various sizes, its parameter range is significantly larger than that of Flan-T5, spanning from 7 billion to 70 billion parameters.

7 Conclusion and Future Work

In this study we investigated the effectiveness of large language models for SATD identification and classification. Our results showed that the fine-tuned selected LLMs outperform all SATD identification baselines. For the SATD classification task, while the largest fine-tuned LLM, Flan-T5-XL, still led in performance, the CNN model exhibited competitive performance, even surpassing some LLMs. In both tasks, larger LLMs outperformed smaller counterparts, particularly in SATD classification task that limited training data is available. We also applied the ICL approach on the largest Flan-T5 model, Flan-T5-XXL. For SATD identification, by including different groups of SATD related keywords in the prompts, it provided competitive results with traditional approaches, but couldn't surpass even the smallest fine-tuned LLM. In SATD classification task, incorporating examples and category descriptions in prompts outperforms the zero-shot approach and even surpasses the fine-tuned smaller Flan-T5 models.

Our experiments also showed that substituting the word generation layer with a classification layer in Flan-T5 models increases their performance for the SATD classification task likely due to the limited available training data, especially when utilizing the smaller versions of the Flan-T5 models. Finally, we discovered that large fine-tuned models, particularly Flan-T5-XL, effectively utilize additional contextual features, such as surrounding code, to enhance performance. In contrast, smaller fine-tuned models exhibit a decrease in performance when provided with complex contextual information.

One of the notable findings from this study is that the fine-tuned Flan-T5-XL model, despite being trained with human-generated labels from the Maldonado-62k dataset, managed to outperform human annotators in SATD identification, thanks to its pre-trained knowledge. We believe that the primary obstacle to achieving higher performance in large LLMs lies in the quality of the labeled data. Therefore, future research in this domain should focus on preparing data of higher quality, which includes essential details such as commit hashes and file paths. These details provides the location of code comments within repositories and facilitate access to additional information, thereby enhancing the performance of future models.

Acknowledgements We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC), [funding reference number: RGPIN-2019-05071].

Conflict of Interest

The authors declare that they have no conflict of interest.

Data Availability Statements

The results, source code, and data related to this study are available at https://github.com/RISElabQueens/SATD_LLM

References

- Bavota, G. and Russo, B. (2016). A large-scale empirical study on self-admitted technical debt. In *2016 IEEE/ACM 13th Working Conference on Mining Software Repositories (MSR)*, pages 315–326.
- Bhatia, A., Khomh, F., Adams, B., and Hassan, A. E. (2023). An empirical study of self-admitted technical debt in machine learning software.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. (2020). Language models are few-shot learners. In Larochelle,

- H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Buschmann, F. (2011). To pay or not to pay technical debt. *IEEE Software*, 28(6):29–31.
- Cassee, N., Zampetti, F., Novielli, N., Serebrenik, A., and Di Penta, M. (2022). Self-admitted technical debt and comments’ polarity: An empirical study. *Empirical Softw. Engg.*, 27(6).
- Chen, X., Yu, D., Fan, X., Wang, L., and Chen, J. (2022). Multiclass classification for self-admitted technical debt based on xgboost. *IEEE Transactions on Reliability*, 71(3):1309–1324.
- Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., Webson, A., Gu, S. S., Dai, Z., Suzgun, M., Chen, X., Chowdhery, A., Castro-Ros, A., Pellat, M., Robinson, K., Valter, D., Narang, S., Mishra, G., Yu, A., Zhao, V., Huang, Y., Dai, A., Yu, H., Petrov, S., Chi, E. H., Dean, J., Devlin, J., Roberts, A., Zhou, D., Le, Q. V., and Wei, J. (2022). Scaling instruction-finetuned language models.
- Cunningham, W. (1992). The wycash portfolio management system. In *Addendum to the Proceedings on Object-Oriented Programming Systems, Languages, and Applications (Addendum)*, OOPSLA ’92, page 29–30, New York, NY, USA. Association for Computing Machinery.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*.
- Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., and Zhou, M. (2020). CodeBERT: A pre-trained model for programming and natural languages. In Cohn, T., He, Y., and Liu, Y., editors, *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1536–1547, Online. Association for Computational Linguistics.
- Fucci, G., Cassee, N., Zampetti, F., Novielli, N., Serebrenik, A., and Di Penta, M. (2021). Waiting around or job half-done? sentiment in self-admitted technical debt. In *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*, pages 403–414.
- Gao, S., Wen, X., Gao, C., Wang, W., Zhang, H., and Lyu, M. R. (2023). What makes good in-context demonstrations for code intelligence tasks with llms? In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 761–773, Los Alamitos, CA, USA. IEEE Computer Society.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- Guo, Z., Liu, S., Liu, J., Li, Y., Chen, L., Lu, H., and Zhou, Y. (2021). How far have we progressed in identifying self-admitted technical debts? a comprehensive empirical study. *ACM Trans. Softw. Eng. Methodol.*, 30(4).
- Hou, X., Zhao, Y., Liu, Y., Yang, Z., Wang, K., Li, L., Luo, X., Lo, D., Grundy, J. C., and Wang, H. (2023). Large language models for software engineering: A systematic literature review. *ArXiv*, abs/2308.10620.

- Huang, Q., Shihab, E., Xia, X., Lo, D., and Li, S. (2017). Identifying self-admitted technical debt in open source projects using text mining. *Empirical Software Engineering*, 23:418 – 451.
- Jiang, Z., Liu, J., Chen, Z., Li, Y., Huang, J., Huo, Y., He, P., Gu, J., and Lyu, M. R. (2023). Llm-parser: A llm-based log parsing framework.
- Jin, M., Shahriar, S., Tufano, M., Shi, X., Lu, S., Sundaresan, N., and Svyatkovskiy, A. (2023). Inferfix: End-to-end program repair with llms. *arXiv preprint arXiv:2303.07263*.
- Li, Y., Soliman, M., and Avgeriou, P. (2022). Identifying self-admitted technical debt in issue tracking systems using machine learning. *Empirical Softw. Engg.*, 27(6).
- Li, Y., Soliman, M., and Avgeriou, P. (2023a). Automatic identification of self-admitted technical debt from four different sources. *Empirical Softw. Engg.*, 28(65).
- Li, Y., Soliman, M., and Avgeriou, P. (2023b). Automatically estimating the effort required to repay self-admitted technical debt.
- Liu, J., Huang, Q., Xia, X., Shihab, E., Lo, D., and Li, S. (2020). Is using deep learning frameworks free? characterizing technical debt in deep learning frameworks. In *2020 IEEE/ACM 42nd International Conference on Software Engineering: Software Engineering in Society (ICSE-SEIS)*, pages 1–10.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach.
- Liu, Z., Huang, Q., Xia, X., Shihab, E., Lo, D., and Li, S. (2018). 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)*, pages 9–12.
- Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- Maldonado, E. d. S. and Shihab, E. (2015). Detecting and quantifying different types of self-admitted technical debt. In *2015 IEEE 7th International Workshop on Managing Technical Debt (MTD)*, pages 9–15.
- Maldonado, E. d. S., Shihab, E., and Tsantalis, N. (2017). Using natural language processing to automatically detect self-admitted technical debt. *IEEE Transactions on Software Engineering*, 43(11):1044–1062.
- Manning, C. and Klein, D. (2003). Optimization, maxent models, and conditional estimation without magic. In *Companion Volume of the Proceedings of HLT-NAACL 2003 - Tutorial Abstracts*, pages 8–8.
- Naveed, H., Khan, A. U., Qiu, S., Saqib, M., Anwar, S., Usman, M., Barnes, N., and Mian, A. S. (2023). A comprehensive overview of large language models. *ArXiv*, abs/2307.06435.
- O'Brien, D., Biswas, S., Imtiaz, S., Abdalkareem, R., Shihab, E., and Rajan, H. (2022). 23 shades of self-admitted technical debt: An empirical study on machine learning software. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of*

- Software Engineering*, ESEC/FSE 2022, page 734–746, New York, NY, USA. Association for Computing Machinery.
- Pinna, A., Lunesu, M. I., Orrù, S., and Tonelli, R. (2023). Investigation on self-admitted technical debt in open-source blockchain projects. *Future Internet*, 15(7).
- Potdar, A. and Shihab, E. (2014). An exploratory study on self-admitted technical debt. In *2014 IEEE International Conference on Software Maintenance and Evolution*, pages 91–100.
- Prenner, J. and Robbes, R. (2022). Making the most of small software engineering datasets with modern machine learning. *IEEE Transactions on Software Engineering*, 48(12):5050–5067.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are unsupervised multitask learners.
- Raffel, C., Shazeer, N. M., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2019). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:140:1–140:67.
- Rantala, L., Mäntylä, M., and Lo, D. (2020). Prevalence, contents and automatic detection of kl-satd. In *2020 46th Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*, pages 385–388.
- Ren, X., Xing, Z., Xia, X., Lo, D., Wang, X., and Grundy, J. (2019). Neural network-based detection of self-admitted technical debt: From performance to explainability. 28(3).
- Sheikhaei, M. S. and Tian, Y. (2023). Automated self-admitted technical debt tracking at commit-level: A language-independent approach. In *2023 ACM/IEEE International Conference on Technical Debt (TechDebt)*, pages 22–26.
- Sridharan, M., Rantala, L., and Mäntylä, M. (2023). Pentacet data-23 million contextual code comments and 250,000 satd comments. In *2023 IEEE/ACM 20th International Conference on Mining Software Repositories (MSR)*, pages 412–416. IEEE.
- Tam, D., Mascarenhas, A., Zhang, S., Kwan, S., Bansal, M., and Raffel, C. (2023). Evaluating the factual consistency of large language models through news summarization. In Rogers, A., Boyd-Graber, J., and Okazaki, N., editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5220–5255, Toronto, Canada. Association for Computational Linguistics.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. (2023). Llama: Open and efficient foundation language models.
- Trask, A., Michalak, P., and Liu, J. C. (2015). sense2vec - a fast and accurate method for word sense disambiguation in neural word embeddings. *ArXiv*, abs/1511.06388.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In Guyon,

- I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Wehaibi, S., Shihab, E., and Guerrouj, L. (2016). Examining the impact of self-admitted technical debt on software quality. In *2016 IEEE 23rd international conference on software analysis, evolution, and reengineering (SANER)*, volume 1, pages 179–188. IEEE.
- Wei, Y., Wang, Z., Liu, J., Ding, Y., and Zhang, L. (2023). Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*.
- Xavier, L., Ferreira, F., Brito, R., and Valente, M. T. (2020). Beyond the code: Mining self-admitted technical debt in issue tracker systems. In *Proceedings of the 17th International Conference on Mining Software Repositories, MSR '20*, page 137–146, New York, NY, USA. Association for Computing Machinery.
- Xiao, T., Wang, D., McIntosh, S., Hata, H., Kula, R. G., Ishio, T., and ichi Matsumoto, K. (2021). Characterizing and mitigating self-admitted technical debt in build systems. *IEEE Transactions on Software Engineering*, 48:4214–4228.
- Yu, D., Wang, L., Chen, X., and Chen, J. (2021). Using bilstm with attention mechanism to automatically detect self-admitted technical debt. *Frontiers of Computer Science*, 15(4):154208.
- Yu, Z., Fahid, F. M., Tu, H., and Menzies, T. (2022). Identifying self-admitted technical debts with jitterbug: A two-step approach. *IEEE Trans. Softw. Eng.*, 48(5):1676–1691.
- Yuan, Z., Liu, J., Zi, Q., Liu, M., Peng, X., and Lou, Y. (2023). Evaluating instruction-tuned large language models on code comprehension and generation. *arXiv preprint arXiv:2308.01240*.
- Zhang, W., Deng, Y., Liu, B., Pan, S. J., and Bing, L. (2023). Sentiment analysis in the era of large language models: A reality check.