Original Article

## Supervised Semantic Similarity-based Conflict Detection Algorithm: S3CDA

Garima Malik<sup>1\*</sup>, Mucahit Cevik<sup>1</sup>, Devang Parikh<sup>2</sup> and Ayse Basar<sup>1</sup>

<sup>1</sup> Toronto Metropolitan University, Toronto, M5K 2B3, Ontario, Canada.

<sup>2</sup> IBM, North Carolina, USA.

\*Corresponding author(s). E-mail(s): garima.malik@torontomu.ca;

#### Abstract

In the realm of software development, the clarity, completeness, and comprehensiveness of requirements significantly impact the success of software systems. The Software Requirement Specification (SRS) document, a cornerstone of the software development life cycle, delineates both functional and nonfunctional requirements, playing a pivotal role in ensuring the quality and timely delivery of software projects. However, the inherent natural language representation of these requirements poses challenges, leading to potential misinterpretations and conflicts. This study addresses the need for conflict identification within requirements by delving into their semantic compositions and contextual meanings. Our research introduces an automated supervised conflict detection method known as the Supervised Semantic Similarity-based Conflict Detection Algorithm (S3CDA). This algorithm comprises two phases: identifying conflict candidates through textual similarity and employing semantic analysis to filter these conflicts. The similaritybased conflict detection involves leveraging sentence embeddings and cosine similarity measures to identify pertinent candidate requirements. Additionally, we present an unsupervised conflict detection algorithm, UnSupCDA, combining key components of S3CDA, tailored for unlabeled software requirements. Generalizability of our methods is tested across five SRS documents from diverse domains. Our experimental results demonstrate the efficacy of the proposed conflict detection strategy, achieving high accuracy in automated conflict identification.

**Keywords:** Software Requirement Specifications, Conflict Detection, Sentence Similarity, Sentence Embeddings, Named Entity Recognition

### 1 Introduction

Requirement Engineering (RE) is the process of defining, documenting, and maintaining the software requirements [1]. RE process involves four main activities, namely, requirements elicitation, requirements specification, requirements verification and validation, and requirements management. In the requirement specification process, the deliverable is termed as Software Requirement Specification (SRS) document which is highly important in Software Development Life Cycle (SDLC) [2]. SRS documents describe the functionality and expected performance for software products, naturally affecting all the subsequent phases in the process. The requirement set defined in SRS documents are analyzed and refined in the design phase, which results in various design documents. Then, the developers proceed with these documents to build the code for the software system [3].

SRS documents are mostly written in natural language to improve the comprehensibility of requirements. The success of any software system is largely dependent on the clarity, transparency, and comprehensibility of software requirements [4]. Conflicting and incomprehensible software requirements might lead to increased project completion times, inefficiency in software systems, and increase in the project budget. Detection of conflicts in the earlier development phase is very important, however, the manual identification of these conflicts could be tedious and time-consuming. It is necessary to develop semi-automated or automated approaches for conflict detection in SRS documents. Considering the structure of software requirements, Natural Language Processing (NLP) methods can help in analyzing and understanding the software requirements semantically. Various information extraction techniques such as Named Entity Recognition (NER), and Parts of Speech (POS) tagging can be used for this purpose, alongside the semantic similarity of the natural language text to interpret the context and syntactic nature of the software requirements.

In order to provide an automated approach for generalised conflict identification, we propose a supervised two-phase framework i.e., Supervised Semantic Similarity-based Conflict Detection Algorithm (S3CDA) which elicits the conflict criteria from the provided software requirements, and outputs the conflicting requirements. In the first phase, we convert the software requirements into high dimensional vectors using various sentence embeddings, and then identify the conflict candidates using cosine similarity. Then, in the second phase, candidate conflict set is further refined by measuring the overlapping entities in the requirement texts, with high level of overlaps pointing to a conflict. Furthermore, we formulate an unsupervised variant of our proposed algorithm, called as UnSupCDA. This algorithm seamlessly integrates the

core elements of the S3CDA approach and is adept at handling unlabeled requirements and identify the conflicts.

#### Research Contribution

The main contributions of our study can be summarized as follows:

- We introduce two novel conflict identification techniques, namely S3CDA and UnSupCDA, meticulously designed through a comprehensive analysis of software requirement structures. We assess the efficacy of these proposed methodologies through extensive numerical experiments conducted on five diverse SRS documents, providing valuable insights into their practical applicability and performance.
- We make use of information extraction techniques in detecting the requirement conflicts. Specifically, we apply software-specific NER model to extract the key entities from the software requirements which can be useful in indicating the presence of conflicts. Additionally, we conduct a thorough analysis of the correlation between semantic similarity among requirements and the overlap of entities in requirement pairs across diverse datasets.

#### Structure of the Paper

The remainder of the paper is organized as follows. Section 2 provides the background on the problem of conflict identification in software requirement datasets. Section 3 introduces our proposed method for automated conflict detection, and provides a detailed discussion over dataset characteristics, sentence embeddings, and NER. In Section 4, we present the results from our experiments and discuss the applicability of our proposed approaches. Lastly, Section 6 provides concluding remarks and future research directions.

## 2 Background

Previous studies suggest the use of NLP-based techniques to solve various software requirement related problems such as requirement classification [5, 6], ambiguity detection [7], bug report classification [8], duplicate bug report prediction [9], conflict identification [10], and mapping of natural language-based requirements into formal structures [11]. Conflict detection is one of the most difficult problems in requirement engineering [3]. Inability in identifying the conflicts in software requirements might lead to uncertainties and cost overrun in software development. Several papers discussed conflict identification in various domains, however, an autonomous, reliable and generalizable approach for detecting conflicting requirements is yet to be achieved. Below, we first introduce the basic definitions and then we review the conflict detection strategies for functional and non-functional requirements.

The terms 'Ambiguity' and 'Conflict' can be misconstrued in the requirement engineering context. Researchers have provided formal definitions for the requirement ambiguity as a requirement having more than one meaning, and

provided various techniques to detect the requirement ambiguities in SRS documents [1, 12]. On the other hand, requirement conflict detection remains as a challenging problem, lacking a well-accepted formal definition and structure. Several studies define the requirement conflict depending upon the domain of requirements. However, the term 'conflict' can be defined more broadly as the presence of interference, interdependency, or inconsistency between requirements [13]. Kim et al. [14] proposed the definition of requirement conflict as interaction and dependencies present between requirements which results into negative or undesired operation of the software systems.

Butt et al. [15] defined requirement conflicts based on the categorization of requirements to mandatory, essential, and optional requirements. Kim et al. [14] described the requirement structure as Actor (Noun) + Action (verb) + Object (object) + Resource (resource). An activity conflict can arise when two requirements achieve the same actions through different object and a resource conflict may arise when different components try to share the same resources. Moser et al. [16] categorized the conflicts as simple (if exists between two requirements) and complex (if exists between three or more requirements). Recently, Guo et al. [10] proposed a comprehensive definition for semantic conflicts amongst different functional requirements. They stated that if two requirements having inferential, interdependent, and inclusive relationship then it may lead to inconsistent behaviour in software system.

In our work, the conflicts are defined based on the premise that if the implementation of  $r_i$  and  $r_j$  cannot coexist or if the implementation of the first adversely impacts the second, then they are considered conflicts in our dataset. An example of such a conflict can be observed in the following requirements:

- 1. The UAV shall charge to 50 % in less than 3 hours.
- 2. The UAV shall fully charge in less than 3 hours.

It's not feasible to implement these requirements simultaneously as it will lead to inconsistency in the system. Notably, for the purpose of this study, requirements deemed as duplicates or paraphrased versions of each other are also considered conflicts due to their inherent redundancy.

Functional requirements specify the functionalities or features a system must possess and describe how the system should behave or what it should do. They outline the specific actions the system must be able to perform and typically address the system's core operations. Table 1 lists the studies for conflict identification in functional requirements. The table provides insights into the domain of SRS documents, the datasets employed, and the types of conflicts addressed in each study. The majority of these studies utilize rule-based and heuristic methods, facing limitations associated with the scarcity of extensive datasets and the absence of a standardized methodology. Often reliant on casestudy approaches, these investigations typically validate their methods using a limited set of requirements, posing challenges for direct comparisons with our study.

Study	Conflict Detection Method	Domain	Dataset	Type of Conflicts
[17]	Rule-based and Genetic Algorithms (GA)	Software	3 SRS documents	Functional conflicts
[10]	Rule-based and semantics	Varied	3 SRS (validation), 2 SRS (testing)	Functional conflicts
[18]	Fuzzy branching temporal logic	Internet of Things (IOT)	System of Systems (SOS) Req.	Conflicts in resource-based req.
[19]	WebSpec tool and semantics	Web software	Case study	Structural conflicts
[20]	Ontology and Semantics	Software	Casestudy	Simple and complex conflicts
[15]	Automated	Software	25 Req.	Conflicts in Mandatory, essential and optional req.
[21]	Rule-based & Tracing dependencies	Software	Voter registration system	Conflicts in Essential Use Cases (EUC)
[22]	Linear temporal logic	Bellcore industry	Telecom dataset	Functional Conflicts
[23]	Heuristics and Req. interactions	Mechanical	12 Req. Lift System	Interactions in Requirements

Table 1: Conflict detection literature for functional requirements in SRS documents.

Guo et al. [10] introduced a methodical approach, FSARC (a Finer Semantic Analysis-based Requirements Conflict Detector), aiming for a comprehensive semantic analysis of software requirements to identify conflicts. FSARC follows a seven-step procedure leveraging Stanford's CoreNLP library [24]. The initial steps involve Part-of-Speech (POS) tagging and Stanford's Dependency Parser (SDP) to transform each requirement into an eight-tuple representation (id, group\_id, event, agent, operation, input, output, restriction). Subsequent rule-based routines are applied to identify conflicts based on this tuple. While the algorithm exhibited promising results and potential for generalization, it heavily depends on the CoreNLP library's accurate generation of the eight-tuple, suggesting a reliance on a specific requirement structure for effective analysis.

Non-functional requirements define the criteria that characterize the operation of a system without detailing specific behaviors. These requirements focus on aspects such as performance, reliability, usability, scalability, and other qualities that are essential for the overall effectiveness and efficiency of the system but are not related to its specific functionalities. Table 2 presents the conflict identification studies for non-functional requirements.

Similar to the studies discussed in the previous section, the studies in Table 2 predominantly relied on rule-based methods and conducted validation on a limited number of requirements. However, the lack of substantial evidence hampers the demonstration of the generalizability of the methods.

Table 2: Conflict determine	ection literature for	r non-functional	requirements	in SRS docu-
ments.				

Study	Conflict Detection Method	Domain	Dataset	Type of Conflicts
[25]	Rule-based and Clustering	Software	15 ATM Req. and 5 SRS documents	Intra-conflicts
[26]	SVM and BiLSTM model with Word2Vec embeddings	Software	200 Req.	Software Product Quality
[27]	Rule-based and Manual	Telecom	14 Req.	Req. Conflicts in OAM&P
[28]	Quantitative analysis	Chemical	Casestudy with 2 req.	
[13]	Rule-based	Software	NFRs	Security and usability
[29]	Rule-based with quality attributes	Software	Case study	Mutually exclusive and partial
[4]	Rule-based and Manual	Software	12 Req. Video on Demand	Req. Traceability

Lastly, we note that our work differs from these existing studies in multiple ways:

- Das et al. [30] introduced the idea of sentence embeddings for similarity detection in software requirements. We extend this idea and combine the two sentence embeddings (SBERT and TFIDF) to calculate the cosine similarity between the requirements.
- Our proposed automated approaches directly works with software requirements as opposed to Guo et al. [10], which converts the software requirements into formal representations and apply rule-based procedures to detect the conflicts.
- Different from finer semantic analysis enabled by Guo et al. [10]'s rule-based approach, we define the set of software-specific entities, and train NLP-based transformer models with for software requirements. These entities provide an additional way of verifying the conflicts semantically.
- To provide the generic technique for conflict identification, we devise an unsupervised approach capable of handling raw requirements from varied domains and effectively capturing conflicts.

## 3 Methodology

In this section, we first describe the SRS datasets used in our numerical study. Then, we provide specific details of the building blocks of our conflict detection algorithms and the experimental setup.

#### 3.1 Datasets

We consider five SRS datasets that belong to various domains such as software, healthcare, transportation, and hardware. Three of these are open-source SRS datasets (OpenCoss, WorldVista, and UAV), and the other two are extracted from public SRS documents. Generally, requirements are documented in a structured format and we retain the original structure of requirements for conflict detection process. To maintain the consistency in requirement structure, we converted complex requirements (e.g., paragraphs or compound sentences) into simple sentences. Table 3 provides summary information on the SRS datasets.

Dataset	Domain	# Non-conflicts	# Known Conflicts	# Synthetic Conflicts
OpenCoss	Transportation	97	4	16
WorldVista	Medical	78	10	60
UAV	Aerospace	80	10	26
PURE	Thermodynamics	27	0	40
IBM-UAV	Hardware	75		28

**Table 3**: Dataset characteristics

We briefly describe these SRS datasets below.

- OpenCoss: OPENCOSS<sup>1</sup> refers to Open Platform for Evolutionary Certification Of Safety-critical Systems for the railway, avionics, and automotive markets. This is a challenging dataset to identify the conflicts as the samples from the OpenCoss dataset indicates a lot of similar or duplicate requirements with repeating words. Initially, this set included 110 requirements and we added 5 more synthetic conflicts.
- WorldVista: WorldVista<sup>2</sup> is a health management system that records patient information starting from the hospital admission to discharge procedures. The requirement structure is basic, and written in natural language with health care terminologies. It originally consisted of 117 requirements and we added 23 synthetic conflicts.
- UAV: The UAV (Unmanned Aerial Vehicle) [10, 31] dataset is created by the University of Notre Dame and it includes all the functional requirements which define the functions of the UAV control system. The requirement syntax is based on the template of EARS (Easy Approach to Requirements Syntax) [32]. Originally, this dataset had 99 requirements and we added 16 conflicting requirements to the set, which resulted in a conflict proportion of 30%.
- **PURE**: PURE (Public Requirements dataset), contains 79 publicly available SRS documents collected from the web [33]. We manually extracted

<sup>&</sup>lt;sup>1</sup>http://www.opencoss-project.eu

<sup>&</sup>lt;sup>2</sup>http://coest.org/datasets

set of requirements from two SRS documents, namely, THEMAS (Thermodynamic System) and Mashbot (web interface for managing a company's presence on social networks). In total, we collected 83 requirements and induced synthetic 21 conflicts to maintain consistency with the other datasets.

• IBM-UAV: This dataset is proprietary, and provided by IBM. It consists of software requirements used in various projects related to the aerospace and automobile industry. We sampled 75 requirements from the original set, and introduced 13 synthetic conflicts. The requirement text follows a certain format specified by IBM's RQA (Requirement Quality Analysis) system.

The synthetic conflicts were introduced to each of these datasets by following the standard definitions provided in the literature [3, 34]. Table 4 shows sample synthetic conflicts as indicated by requirement id, requirement text, a 'Conflict' column indicating the presence of conflict in 'Yes' or 'No' format and a 'Conflict-Label' column that indicates the pair of conflicts. For example, requirements 2 and 3 conflict with each other because of modal verb used in both requirements. Similarly, requirements 11 and 12 show mathematical operator conflict.

# 3.2 S3CDA: Supervised Semantic Similarity-based Conflict Detection Algorithm

This section is structured into two parts as depicted in Figure 1. First, we explain the similarity-based conflict detection. Second, we define the semantic-based conflict identification to validate the potential conflicts obtained in Phase I.

#### 3.2.1 Phase I: Similarity-based Conflict Detection

Algorithm 1 formalizes similarity-based conflict detection procedure (Phase I) provided in the left panel of Figure 1. The resulting set of conflicts can be used as a candidate conflict set for Phase II of our framework. We first create the sentence embedding vector for each requirement  $r \in \mathcal{R}$  using SentenceEmbedding( $\mathcal{R}$ ) procedure. It basically converts the requirements into numerical vector using sentence embeddings. Below, we describe the various sentence embedding models employed in the proposed algorithms.

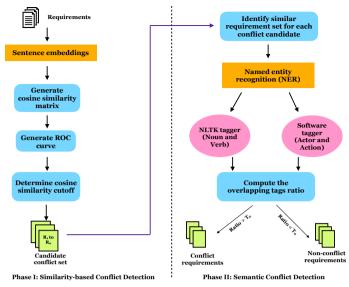
- **TFIDF**: Term Frequency Inverse Document Frequency, is employed for generating vectors from requirements [35]. Each requirement is treated as a document, and TFIDF scores are calculated for each term in the document. The TFIDF value for a term i' in a document d is given by  $TF_{i,d} \times IDF_i$ , where  $TF_{i,d}$  is the term frequency in document d', and  $IDF_i$  is the inverse document frequency.
- USE: Universal Sentence Encoder (USE) translates natural language text into high-dimensional vectors [36]. We utilize the Deep Averaging Network

Table 4: Sample conflict and non-conflict instances from each dataset with require-
ment id, text, and conflict label (Yes/No).

Dataset	Req. Id	Requirement text	Conflict	Conflict- Label
OpenCoss	1.	The OPENCOSS platform shall store the characteristics of the evidence items of an assurance project according to the CCL.	No	No
	2.	The OPENCOSS platform must provide users with the ability to specify evidence traceability links in traceability matrices.	Yes	Yes (3)
	3.	The OPENCOSS platform shall provide users with the ability to specify evidence traceability links in traceability matrices.	Yes	Yes (2)
WorldVista	4.	The system's pilot program shall use a smart card to digitally sign medication orders.	Yes	Yes (5)
	5.	The system's pilot program shall require a handwritten signature for medication orders.	Yes	Yes (4)
	6.	The system shall sort notifications based on column heading: Patient name (alphabetical or reverse alphabetical).	No	No
IBM- UAV	7.	The UAV shall charge to 50 % in less than 3 hours.	Yes	Yes (8)
	8. 9.	The UAV shall fully charge in less than 3 hours. Remote surveillance shall include video streaming for manual navigation of the surveillance platform.	Yes No	Yes (7) No
PURE	10.	If LO <= T <= UO, then this process shall output the temperature status.	Yes	Yes (11)
	11.	If LO $<=$ T $<$ UO, then this process shall output the temperature status.	Yes	Yes (12)
	12.	The THEMAS system shall maintain the ON/OFF status of each heating and cooling unit.	No	No
UAV	13.	The _InternalSimulator_ shall approximate the behavior of a UAV.	No	No
	14. 15.	The _VehicleCore_ shall support virtual UAVs. The _VehicleCore_ shall support up to three virtual UAVs.	Yes Yes	Yes (15) Yes (14)

(DAN) version of USE, a pre-trained model optimized for encoding sentences, phrases, and short paragraphs. It produces 512-dimensional vectors for each input requirement.

• SBERT-TFIDF: This method combines Sentence-BERT [37] (SBERT) and TFIDF embeddings. SBERT provides context and semantics in the vectors, while TFIDF prioritizes less frequent words. The process involves concatenating the vectors and employing Uniform Manifold Approximation and Projection (UMAP) for dimensionality reduction to ensure uniform vector size [38].



**Fig. 1:** Supervised Semantic Similarity-based Conflict Detection Algorithm (S3CDA). Proposed framework for identifying the conflicting requirements in SRS documents.

We next calculate the pairwise distance matrix  $(\Delta)$ , which measures the cosine similarity value between each pair of requirements. Below, we provide a brief example of using cosine similarity for requirements.

#### Cosine Similarity:

It is an effective measure to estimate the similarity of vectors in highdimensional space [39]. This metric models language-based input text as a vector of real-valued terms and the similarity between two texts is derived from the cosine angle between two texts term vectors as follows:

$$\cos(\mathbf{r_1}, \mathbf{r_2}) = \frac{\mathbf{r_1} \mathbf{r_2}}{\|\mathbf{r_1}\| \|\mathbf{r_2}\|} = \frac{\sum_{i=1}^{n} \mathbf{r_1}_i \mathbf{r_2}_i}{\sqrt{\sum_{i=1}^{n} (\mathbf{r_1}_i)^2} \sqrt{\sum_{i=1}^{n} (\mathbf{r_2}_i)^2}}$$
(1)

The values for cosine similarity ranges from -1 to 1 where -1 signifies dissimilarity and 1 signifies similarity. To better demonstrate how cosine similarity can be used over embedding vectors, we provide an illustrative example with three sample software requirements,  $r_1$ ,  $r_2$ , and  $r_3$ , which are defined as follows:

- $r_1$  = 'The OPENCOSS platform shall be able to export evidence traceability links of an assurance project to external tools.'
- $r_2$  = 'The OPENCOSS platform must be able to send out evidence traceability links of an assurance project to external tools and internal tools.'

•  $r_3$  = 'The OPENCOSS platform shall provide users with the ability to specify evidence traceability links in traceability matrices.'

We calculate the cosine similarity between these requirement vectors when embedded with TFIDF, SBERT, SBERT-TFIDF, and USE, which are reported in Table 5. SBERT is able to capture the semantic similarity between  $r_1$  and  $r_2$  with a cosine similarity value of 0.96. USE being the second highest with the value of 0.81. Requirement text for  $r_3$  is not similar to those of the other two requirements, and all the sentence embeddings indicate low values of cosine similarity with  $r_3$ .

**Table 5**: Cosine similarity between  $r_1$ ,  $r_2$ , and  $r_3$  with different sentence embeddings.

	TFIDF	USE	SBERT	SBERT-TFIDF
$\cos(r_1, r_2) \\ \cos(r_1, r_3) \\ \cos(r_2, r_3)$	0.44	0.81	0.96	0.72
	0.16	0.55	0.84	0.50
	0.09	0.40	0.83	0.46

Then, we use ROC curve (receiver operating characteristic curve) to identify the cosine similarity threshold  $(\delta)$ , which specifies the minimum similarity value after which requirements are labeled as conflicting. The cutoff value  $(\delta)$  is selected as the value that maximizes  $\{\text{TPR}(\Delta,k) - (1-\text{FPR}(\Delta,k))\}$  over threshold values  $k \in \{0.01,\ldots,1.00\}$  and the distance matrix  $\Delta$ . This way, we balance the false positives and true positives rates, with the conflicts having the positive labels. Lastly, we assign labels of conflict or no-conflict to the requirements using  $\delta$  as threshold value. The candidate conflict set  $(\mathcal{C})$  contains all the requirements with conflict label. Note that conflict property is symmetric, i.e., if  $r_1$  is conflicting with  $r_2$ , then  $r_2$  is also conflicting with  $r_1$ , and,  $r_1, r_2 \in \mathcal{C}$ .

#### Algorithm 1 Similarity-based Conflict Detection

#### 3.2.2 Phase II: Semantic-based Conflict Detection

Algorithm 2 describes the process of semantic conflict detection as presented in right panel of Figure 1. This algorithm serves as a second filter on the candidate conflicts generated in Phase I. Specifically, any candidate conflict  $c \in \mathcal{C}$  is semantically compared against top m most similar requirements from

 $\mathcal{R}$ . That is, by focusing on only m most similar requirements, we reduce the computational burden, and also make use of the cosine similarity between the requirements. This semantic comparison is performed based on overlap ratio between the entities present in the requirements. For a given candidate conflict  $c \in \mathcal{C}$ , overlap ratio is calculated as

$$v_c = \frac{\max_{r \in \mathcal{L}_c} \left\{ \text{Overlap}(c, r) \right\}}{\text{UniqueEntities}(c)}$$
 (2)

where  $\mathcal{L}_c$  represents the set of m most similar requirements to candidate conflict c. The function Overlap(c, r) calculates the number of overlapping entities between c and r, and function UniqueEntities(c) calculates the number of unique entities in candidate conflict (i.e., a requirement text) c. The calculated overlap ratio  $v_c$  for  $c \in \mathcal{C}$  is then compared against a pre-determined overlap threshold value,  $T_o$ , and c is added to final conflict set  $\bar{C}$  if  $v_c \geq T_o$ . In our analysis, we set m = 5 and  $T_o = 1$ , which are determined based on preliminary experiments. In Algorithm 2, GetSimilarRequirements $(c, \mathcal{R}, m)$ returns the list  $\mathcal{L}$  of m most similar requirements from  $\mathcal{R}$  for candidate conflict c, and GetMaxOverlapRatio $(c, \mathcal{L})$  returns the maximum value for the overlaps between requirements from set  $\mathcal{L}$  and candidate conflict c.

#### **Algorithm 2** Semantic Conflict Detection

```
Require:
   Candidate conflict set: C = \{c_1, c_2, \dots, c_t\}
   Requirement set: \mathcal{R} = \{r_1, r_2, \dots, r_n\}
    \# of similar requirements: m
Output: Refined conflict set \bar{C}
Initialization:
   \bar{\mathcal{C}} = \emptyset
                                                     // Initialize conflict set to an empty set
   T_o = 1
                                                              // Set overlapping threshold as 1
   m = 5
                                                 // Set number of similar requirements as 5
For c \in \mathcal{C} do:
    \mathcal{L} \leftarrow \texttt{GetSimilarRequirements}(c, \mathcal{R}, m) // Get the m similar requirements
   v \leftarrow \texttt{GetMaxOverlapRatio}(c, \mathcal{L}) // Calculate max. overlap ratio using Eqn. (2)
   If v > T_o:
                                                                      // Compare with threshold
       \bar{\mathcal{C}} \leftarrow \bar{\mathcal{C}} \cup \{c\}
                                                              // Augment the final conflict set
```

To extract the entities from the requirements, we employ two NER techniques described below.

• Part-of-Speech (POS) Tagging: [40] suggest that a software requirement should follow the structure as Actor (Noun) + Object + Action (Verb) + Resource. The generic NER method extracts 'Noun' and 'Verb' tags from the requirements based on this structure and referred as 'POS' tagging. We employ POS tagger provided in SpaCy library in Python.

• Software-specific Named Entity Recognition (S-NER): NER serves to extract relevant entities from input text, and its effectiveness can be enhanced by training machine learning models on domain-specific corpora. In our context, we leverage a software-specific NER system to extract entities crucial for understanding requirements, specifically focusing on actor, action, object, property, metric, and operator.

In the context of requirements, an "Actor" denotes an entity interacting with the software application, while an "Action" represents an operation performed by an actor within the software system. To illustrate, consider the requirement "The UAV shall charge to 75% in less than 3 hours," where UAV serves as the actor, and charge as the corresponding action. We employ transformer models trained on software-specific corpora to proficiently extract these entities from requirement pairs [41].

We also provide sample calculations in Table 6 to better illustrate the process in Algorithm 2. We show the overlapping software-specific entities present in candidate requirement (c) and similar requirement  $(r \in \mathcal{L})$  with different color codes. For instance, the entity 'UAV' is represented by blue color. The requirements  $r_1$  and  $r_2$  both return high overlap ratios, indicating a conflict with c.

**Table 6**: Sample candidate requirement (c) and set of similar requirements  $(\mathcal{L})$  to calculate the overlapping ratio (v). The maximum count of overlap is 7 which resulted in overlap ratio value v as 1.

Candidate Requirement	Similar Requirements	$\mathtt{Overlap}(c,r)$	Ratio $(v)$
	$r_1$ = 'The UAV flight range shall be no less than 20 miles.'	7	1.00
	$r_2$ = 'The UAV flight range shall be a minimum of 20 miles.'	5	0.71
c= 'The UAV flight range shall be no less than 20 kilometers.'	$r_3$ = 'The UAV shall be able to autonomously a flight plan consisting of a set of waypoints within its range and flight capabilities.'	2	0.28
	r <sub>4</sub> = 'The UAV shall be able to autonomously plan a flight consisting of a set of waypoints within its range and flight capabilities.'	2	0.28
	r <sub>5</sub> = 'The Pilot controller shall be able to download a flight plan from the Hummingbird consisting of a set of waypoints within the flight range of the Hummingbird.'	2	0.28

# 3.3 UnSupCDA: Unsupervised Conflict Detection Algorithm

We devise an unsupervised variant of the S3CDA approach to alleviate the need for labeled conflict data in the task of identifying conflicts within SRS

documents. This unsupervised version seamlessly integrates key components from both phases of the original S3CDA approach, adapting them to function without the reliance on labeled data. By combining the strengths of the two phases, our unsupervised approach retains the efficiency and effectiveness of conflict identification while eliminating the necessity for manual annotation.

In Algorithm 3, we first transform requirements into embedding vectors using SentenceEmbedding function, capturing their semantic content. Subsequently, a pairwise cosine similarity matrix  $(\Delta)$  is computed using PairwiseDistance function. For each requirement in  $\mathcal{R}$ , we identify the n similar requirements and simultaneously calculate the entity overlapping ratio using Equation 2. Then, we employs a predefined threshold  $(T_o)$  to determine the conflicts.

#### Algorithm 3 UnSupCDA

```
Require:
    Requirement set: \mathcal{R} = \{r_1, r_2, \dots, r_n\}
Output: Conflict set \bar{\mathcal{C}}
Initialization:
    \bar{C} = \emptyset
                                                          // Initialize conflict set to an empty set
    T_o = \{1.0, 0.75, 0.5\}
                                                                           // Overlapping threshold set
    n = \{1, 2, 3\}
                                                                  // n highest similarity requirement
\vec{\mathcal{R}} \leftarrow \texttt{SentenceEmbedding}(\mathcal{R})
                                                                      // Generate requirement vectors
\Delta \leftarrow \texttt{PairwiseDistance}(\vec{\mathcal{R}})
                                                                                 // Get similarity matrix
For r_i \in \mathcal{R} do:
   (r_i, r_i) \leftarrow n_HighestSimilarityPairs(\Delta, n)
    v \leftarrow \texttt{GetMaxOverlapRatio}(r_i, r_i)
                                                               // max. overlap ratio using Eqn. (2)
   If v \geq T_o:
        \bar{\mathcal{C}} \leftarrow \bar{\mathcal{C}} \cup \{r_i, r_i\}
```

## 3.3.1 Overlapping Entity Ratio v/s Cosine Similarity

We analyze the correlation between the overlapping entity ratio, as defined in Equation 2, and cosine similarity through scatter plot visualizations. The requirements from each dataset are embedded using the SBERT model, and pairwise cosine similarity is computed. Subsequently, the most similar requirement is identified for each requirement in the dataset. To calculate the entity overlapping ratio, we employ POS tagging and the S-NER method for entity extraction. The overlapping ratio is then determined based on the extracted entities. Our objective is to investigate whether there exists a relationship between the cosine similarity of requirements and their entity overlap ratio, as both metrics play crucial roles in our proposed algorithms.

Examining Figure 2, we discern an absence of a definitive correlation. Notably, in the case of the WorldVista dataset, Figures 2a and 2b reveal a

weak positive correlation, where high values of cosine similarity align with high overlapping ratios. This trend is consistent across other datasets. In the OpenCoss dataset, Figure 3b illustrates that most requirements exhibit both high cosine similarity and overlapping ratios, posing challenges for our conflict detection algorithms and elucidating the observed performance metrics. For other datasets, the scatter plots are presented in Appendix 7.1.

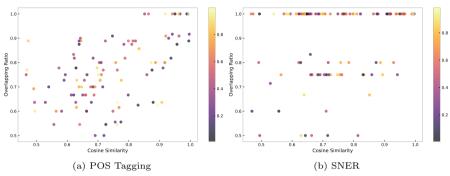


Fig. 2: Scatterplot depicting the relationship between Cosine Similarity and Overlapping Entity Ratio among WorldVista requirements.

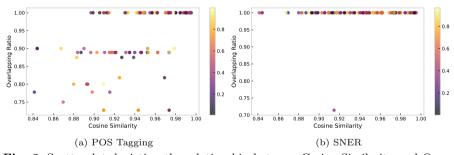


Fig. 3: Scatterplot depicting the relationship between Cosine Similarity and Overlapping Entity Ratio among OpenCoss requirements.

## 3.4 Experimental Setup

For the evaluation of our proposed approaches, we perform 3-fold cross validation over all the requirement datasets. That is, considering the distribution of conflicting requirements and the limited number of requirements, we divide each dataset into 3 different folds. Each fold includes some conflicting and non-conflicting requirements, however, we make sure that each conflict present in

the fold should have its conflict pair present in the same fold. For our techniques, we use training set to determine the cosine similarity cut-off value, and apply this value on the corresponding test set. For S3CDA, we employ standard classification metrics such as macro averaged version of F1-score, Precision, and Recall in performance evaluation. Table 7 reports the hyperparameter and model configurations for the models used in our proposed approaches.

	= -	
Model	Checkpoint/library	Model configurations/hyperparameters
SBERT	all-MiniLM-L12-V2 $^3$ all-mpnet-base-v2 $^4$	12-layer, 384-hidden, 12-heads, 33M parameters 12-layer, 768-hidden, 12-heads, 110M parameters
USE	TensorFlow 2.0 (v4)	512-dimensional vector
TFIDF	scikit-learn (sklearn) $^5$	$n\_gram = (1,4), 1024$ -dimensional vector
S-NER	Bert-base-cased [41]	12-layer, 768-hidden, 12-heads, 110M parameters
POS tagging	spaCv <sup>6</sup>	model-"en_core_web_sm"

Table 7: Model configurations and hyperparameters used in the proposed approach.

In the case of UnSupCDA, being an unsupervised task, we adopt the evaluation metrics introduced by Guo et al. [10] to assess the accuracy of conflict identification. The metrics include Precision and Recall, calculated as follows:

$$ext{Precision} = rac{Correctly\ Detected}{Overall\ Detected\ Conflicts}, \quad ext{Recall} = rac{Correctly\ Detected}{Overall\ Known\ Conflicts}$$

In these calculations, "Correctly Detected Conflicts", also known as True Positives (TP), represent the instances where the algorithm correctly identifies conflicting requirements as conflict. "Overall Detected Conflicts", also known as False Positives (FP) refer to the total number of instances where the algorithm identifies a requirement as potentially conflict, regardless of whether this identification is conflict or not. "Overall Known Conflicts" is the number of conflicts labeled by experts as true conflicts.

For entity extraction, we utilise the transformer model checkpoint trained in Malik et al. [41]. We directly employ the best model *Bert-base-cased* to our software requirements to calculate the software-specific entities.

#### 4 Results

In this section, we first assess the performance of the supervised conflict detection algorithm (S3CDA), and present the results from our comparative analysis with various sentence embeddings for each of the requirement datasets. Next, we present the performance of our unsupervised conflict detection algorithm.

#### 4.1 S3CDA Performance Evaluation

In our proposed approach S3CDA, Phase I, the step that comes after calculating the similarity matrix between the requirements is the generation of ROC curves, which are used to obtain the cosine similarity cut-off for conflict detection based on TPR and FPR values. Figure 4 shows the ROC curves for the UAV dataset with USE embedding in each fold. We generate similar ROC curves for all the other requirement sets with best sentence embeddings which are provided in the appendix (see Section 7).

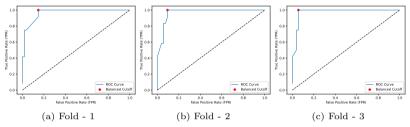


Fig. 4: ROC curves for UAV dataset with USE embedding across 3 folds

Table 8 presents a comparison of various sentence embeddings in Phase I of the S3CDA approach. The performance across different embeddings exhibits no discernible pattern, with TFIDF performing consistently at an average level. Notably, for PURE and IBM-UAV datasets, SBERT achieves F1-scores of 89.6% and 71.1%, respectively. In the case of UAV, the USE embedding attains the highest F1-score of 92.3%, while for WorldVista, SBERT-TFIDF achieves an F1-score of 87.1%. Surprisingly, the OpenCoss dataset demonstrates the highest F1-score (57.0%) with TFIDF embeddings. This anomaly may be attributed to the challenging nature of the OpenCoss dataset, characterized by a complex word structure and a substantial overlap in vocabulary among requirements. This overlap proves advantageous for frequency-based TFIDF embeddings in capturing relevant patterns.

In the subsequent step of the S3CDA approach, we validate the potential conflicts identified in Phase I. This validation involves the application of entity extraction techniques, specifically POS tagging and S-NER, to the potential conflicts associated with the best embedding identified in Table 8. Notably, Phase II of the S3CDA approach is employed to validate conflicts through semantic analysis, focusing on the overlapping entities within potential conflicts.

Table 9 summarizes the results of this validation. Notably, the S-NER method demonstrates relatively consistent performance and effectively validates Correctly Detected Conflicts. Across both techniques, there is a trend of improved Precision scores accompanied by a decrease in Recall, resulting in an overall decline in F1-scores. Additionally, the use of a hard threshold

**Table 8**: Results for various sentence embedding with requirement datasets are reported using macro-averaged metrics - Precision (P), Recall (R), and F1-score (F1) expressed as "mean  $\pm$  std" across three folds. The evaluation also includes the number of Potential and Correctly Detected conflicts in comparison to known conflicts for each dataset. Best embedding is highlighted in blue color and yellow shading indicates noteworthy results. Here,  $E_1$ : TFIDF,  $E_2$ : USE,  $E_3$ : SBERT,  $E_4$ : SBERT-TFIDF.

Dataset	Embeddings	Cosine Cutoff	P	R	F1	Potential Conflicts	Correctly Detected
OpenCoss	$\mathbf{E_1}$	0.700	$0.433\pm0.118$	$0.847 \pm 0.137$	$0.570 \pm 0.130$	41 / 20	17 / 20
	$E_2$	0.936	$0.366 \pm 0.059$	$0.652 \pm 0.272$	$0.461 \pm 0.119$	34 / 20	13 / 20
	$E_3$	0.956	$0.354 \pm 0.041$	$0.805\pm0.141$	$0.484 \pm 0.026$	46 / 20	16 / 20
	$E_4$	0.826	$0.373 \pm 0.161$	$0.722\pm0.207$	$0.487\pm0.185$	41 / 20	14 / 20
WorldVista	$E_1$	0.206	$0.812 \pm 0.093$	$0.853 \pm 0.112$	$0.821 \pm 0.036$	76 / 70	60 / 70
	$E_2$	0.663	$0.897\pm0.074$	$0.839 \pm 0.116$	$0.857 \pm 0.034$	67 / 70	59 / 70
	$E_3$	0.766	$0.856 \pm 0.112$	$0.875\pm0.102$	$0.855 \pm 0.050$	73 / 70	61 / 70
	$\mathbf{E_4}$	0.510	$0.898\pm0.087$	$0.856 \pm 0.043$	$0.871 \pm 0.022$	68 / 70	60 / 70
UAV	$E_1$	0.363	$0.852 \pm 0.051$	$0.944 \pm 0.078$	$0.894 \pm 0.051$	40 / 36	34 / 36
	$\mathbf{E_2}$	0.770	$0.860 \pm 0.050$	$1.000 \pm 0.000$	$0.923\pm0.029$	42 / 36	36 / 36
	$E_3$	0.883	$0.918\pm0.068$	$0.833 \pm 0.136$	$0.864 \pm 0.068$	31 / 36	30 / 36
	$E_4$	0.616	$0.893 \pm 0.042$	$0.944 \pm 0.048$	$0.917\pm0.059$	38 / 36	34 / 36
PURE	$E_1$	0.816	$0.910 \pm 0.063$	$0.785 \pm 0.210$	$0.819 \pm 0.112$	36 / 40	32 / 40
	$E_2$	0.816	$0.893 \pm 0.076$	$0.785 \pm 0.210$	$0.809 \pm 0.102$	37 / 40	32 / 40
	$\mathbf{E_3}$	0.883	$0.902 \pm 0.070$	$0.896\pm0.073$	$0.896\pm0.044$	40 / 40	36 / 40
	$E_4$	0.710	$0.958\pm0.058$	$0.785 \pm 0.210$	$0.841 \pm 0.123$	34 / 40	32 / 40
IBM-UAV	$E_1$	0.583	$0.708 \pm 0.212$	$0.550 \pm 0.324$	$0.557 \pm 0.178$	24 / 28	16 / 28
	$E_2$	0.836	$0.704 \pm 0.163$	$0.716 \pm 0.332$	$0.688 \pm 0.252$	28 / 28	21 / 28
	$\mathbf{E_3}$	0.899	$0.871 \pm 0.118$	$0.716\pm0.332$	$0.711\pm0.221$	26 / 28	21 / 28
	$E_4$	0.730	$0.529 \pm 0.380$	$0.566 \pm 0.418$	$0.537 \pm 0.380$	22 / 28	17 / 28

 $(T_o=1)$  leads to a reduction in Recall scores, while a more flexible threshold  $(T_o=0.75)$  shows consistent or equivalent performance to Phase I for various metrics.

## 4.2 UnSupCDA Performance Evaluation

Table 10 outlines the assessment outcomes for the UnSupCDA approach. Overall, the approach demonstrates higher Recall values across all datasets, indicating its proficiency in capturing conflicts. However, Precision struggles, leading to lower F1-scores. This is attributed to the algorithm's ability to identify conflicts, but it also includes false candidates, adversely impacting Precision.

As anticipated, employing a hard threshold  $(T_0 = 1)$  results in consistently low Precision for all n values. Conversely, a stable threshold  $(T_o = 0.75)$  shows improved Precision and Recall compared to  $T_o = 1$ . A soft threshold  $(T_o = 0.5)$  leads to the selection of numerous false conflicts, diminishing Precision. This pattern is consistent across all n values.

Contrary to expectations, increasing n from 1 to 3 does not yield enhanced performance scores. The assumption was that a higher n would improve scores,

**Table 9**: Comparison of POS tagging and S-NER entity extraction techniques with  $T_o = \{1.0, 0.75\}$ . Performance changes are highlighted in blue and red, indicating increments and decrements, respectively, with both absolute and relative percentage values. The baseline for comparison is established using the best embedding performance as identified in Table 8.

		$\mathbf{T_o} = 1$				$T_o = 0.75$			
Dataset	Method	Р	R	F1	Correctly Detected	P	R	F1	Correctly Detected
OpenCoss	POS	0.419 ± 0.097 ↓ -0.014 / 3.23%	0.791 ± 0.090 ↓ -0.056 / 6.611%	$0.543 \pm 0.092$	16 / 20	0.433 ± 0.118 ~	0.847 ± 0.137 ~	$0.570 \pm 0.130$	17 / 20
	S-NER	0.433 ± 0.118 ~	0.847 ± 0.137 ~	$0.570 \pm 0.130$	17 / 20	0.433 ± 0.118 ~	$0.847 \pm 0.137$ $\sim$	$0.570 \pm 0.130$	17 / 20
WorldVista	POS	0.925 ± 0.104 ↑ +0.027 / 3.00%	0.327 ± 0.063 ↓ -0.548 / 62.62%	$0.480 \pm 0.076$	23 / 70	0.922 ± 0.078 ↑ +0.024 / 2.67%	0.825 ± 0.077 ↓ -0.050 / 5.71%	$0.864 \pm 0.019$	58 / 70
	S-NER	$\begin{array}{c} 0.904 \pm 0.098 \\ \uparrow + 0.006 \ / \ 0.66\% \end{array}$	$\begin{array}{c} 0.756\pm0.055 \\ \downarrow -0.119 \; / \; {\color{red} 23.60\%} \end{array}$	$0.817 \pm 0.018$	53 / 70	0.898 ± 0.087 ~	$0.856 \pm 0.043$ $\sim$	$0.871 \pm 0.022$	60 / 70
UAV	POS	0.925 ± 0.104 ↑ +0.065 / 7.55%	0.361 ± 0.171 ↓ -0.639 / 63.90%	$0.484 \pm 0.155$	13 / 36	0.855 ± 0.044 ↓-0.005 / 5.81%	0.944 ± 0.078 ↓ -0.056 / 5.60%	$0.893 \pm 0.022$	34 / 36
	S-NER	$\begin{array}{c} 0.867 \pm 0.039 \\ \uparrow + 0.007 \; / \; 0.81\% \end{array}$	0.916 ± 0.068 ↓ -0.084 / 8.40%	$0.891 \pm 0.052$	33 / 36	0.860 ± 0.050 ~	1.000 ± 0.000 ~	$0.923 \pm 0.029$	36 / 36
PURE	POS	0.850 ± 0.108 ↓ -0.052 / 5.76%	0.603 ± 0.124 ↓ -0.293 / 32.70%	$0.702 \pm 0.118$	24 / 40	0.897 ± 0.075 ↑ +0.001 / 0.11%	0.869 ± 0.102 ↓ -0.027 / 3.01%	$0.879 \pm 0.068$	35 / 40
	S-NER	0.891 ± 0.078 ↓ -0.011 / 1.21%	0.750 ± 0.087 ↓ -0.146 / 16.29%	$0.807 \pm 0.035$	30 / 40	0.902 ± 0.070 ↑ +0.006 / 0.66%	$0.896 \pm 0.073$ ~	$0.896 \pm 0.044$	36 / 40
IBM-UAV	POS	0.909 ± 0.128 ↑ +0.038 / 4.36%	0.583 ± 0.239 ↓ -0.133 / 18.57%	$0.661 \pm 0.186$	17 / 28	0.864 ± 0.128 ↓-0.007 / 0.80%	0.683 ± 0.306 ↓ -0.033 / 4.60%	$0.694 \pm 0.213$	20 / 28
	S-NER	0.914 ± 0.068 ↑ +0.043 / 4.93%	$0.583 \pm 0.311$ $\downarrow -0.133 / 18.57\%$	$0.659\pm0.226$	17 / 28	0.871 ± 0.118	0.716 ± 0.332 ~	$0.711\pm0.221$	21 / 28

**Table 10**: Results with the UnSupCDA for  $n = \{1, 2, 3\}$  with various thresholds. Reported values are averaged over 3 folds and presented as "mean  $\pm$  std".

				(a) <b>n</b>	i = 1				
		$T_{o} = 1.0$			$\mathbf{T_o} = 0.75$			$T_o = 0.5$	
Dataset	P	R	F1	Р	R	F1	P	R	F1
OpenCoss	$0.178 \pm 0.023$	$1.000 \pm 0.000$	$0.302\pm0.032$	$0.170 \pm 0.020$	$1.000 \pm 0.000$	$0.290 \pm 0.030$	$0.170 \pm 0.020$	$1.000 \pm 0.000$	$0.290 \pm 0.030$
WorldVista	$0.511 \pm 0.055$	$0.656 \pm 0.136$	$0.565 \pm 0.056$	$0.510 \pm 0.029$	$0.972\pm0.039$	$0.667\pm0.015$	$0.472 \pm 0.010$	$1.000 \pm 0.000$	$0.641 \pm 0.009$
UAV	$0.330 \pm 0.019$	$0.805 \pm 0.039$	$0.468 \pm 0.024$	$0.336 \pm 0.004$	$1.000\pm0.000$	$0.503\pm0.005$	$0.310 \pm 0.003$	$1.000 \pm 0.000$	$0.473 \pm 0.004$
PURE	$0.752 \pm 0.110$	$0.821 \pm 0.050$	$0.783 \pm 0.083$	$0.741 \pm 0.126$	$1.000\pm0.000$	$0.845\pm0.086$	$0.605 \pm 0.026$	$1.000 \pm 0.000$	$0.753 \pm 0.020$
$\operatorname{IBM-UAV}$	$0.458\pm0.097$	$0.700\pm0.141$	$\bf 0.553\pm0.112$	$0.380 \pm 0.077$	$0.766 \pm 0.205$	$0.508 \pm 0.115$	$0.291 \pm 0.008$	$1.000\pm0.000$	$0.451 \pm 0.009$
				(b) n	$\mathbf{a} = 2$				
	$T_o = 1.0$			$\mathbf{T_o} = 0.75$			$T_o = 0.5$		
Dataset	P	R	F1	P	R	F1	P	R	F1
OpenCoss	$0.176 \pm 0.020$	$1.000 \pm 0.000$	$0.299 \pm 0.029$	$0.170 \pm 0.020$	$1.000 \pm 0.000$	$0.290 \pm 0.030$	$0.170 \pm 0.020$	$1.000 \pm 0.000$	$0.290 \pm 0.03$
WorldVista	$0.500 \pm 0.000$	$0.772 \pm 0.107$	$0.604 \pm 0.034$	$0.493 \pm 0.013$	$\textbf{0.972}\pm\textbf{0.039}$	$0.654\pm0.003$	$0.472 \pm 0.010$	$1.000 \pm 0.000$	$0.641 \pm 0.009$
UAV	$0.329 \pm 0.020$	$0.833 \pm 0.068$	$0.472 \pm 0.031$	$0.315 \pm 0.000$	$1.000\pm0.000$	$0.480\pm0.000$	$0.310 \pm 0.003$	$1.000 \pm 0.000$	$0.473 \pm 0.004$
PURE	$0.748 \pm 0.115$	$0.849 \pm 0.011$	$0.791 \pm 0.072$	$0.677 \pm 0.106$	$1.000\pm0.000$	$0.802\pm0.073$	$0.605 \pm 0.026$	$1.000 \pm 0.000$	$0.753 \pm 0.020$
$\operatorname{IBM-UAV}$	$0.413 \pm 0.101$	$0.700\pm0.141$	$0.518 \pm 0.115$	$0.325 \pm 0.086$	$0.766 \pm 0.205$	$0.456 \pm 0.121$	$0.279 \pm 0.022$	$1.000\pm0.000$	$0.436\pm0.027$
				(c) <b>n</b>					
		$T_o = 1.0$			$\mathbf{T_o} = 0.75$			$T_o = 0.5$	
Dataset	P	R	F1	P	R	F1	P	R	F1
OpenCoss	$0.175\pm0.021$	$1.000 \pm 0.000$	$0.297 \pm 0.031$	$0.170 \pm 0.020$	$1.000 \pm 0.000$	$0.290 \pm 0.030$	$0.170 \pm 0.020$	$1.000 \pm 0.000$	$0.290 \pm 0.030$
WorldVista	$0.510 \pm 0.014$	$0.800 \pm 0.069$	$0.621 \pm 0.011$	$0.489 \pm 0.008$	$0.972\pm0.039$	$0.650\pm0.004$	$0.472 \pm 0.010$	$1.000\pm0.000$	$0.641 \pm 0.009$
UAV	$0.318 \pm 0.020$	$0.833 \pm 0.068$	$0.461 \pm 0.030$	$0.313 \pm 0.003$	$\textbf{1.000}\pm\textbf{0.000}$	$0.476\pm0.004$	$0.310 \pm 0.003$	$1.000\pm0.000$	$0.473\pm0.004$
PURE	$0.712 \pm 0.090$	$0.849 \pm 0.011$	$0.772 \pm 0.059$	$0.661 \pm 0.086$	$\textbf{1.000}\pm\textbf{0.000}$	$0.793\pm0.061$	$0.605 \pm 0.026$	$1.000\pm0.000$	$0.753\pm0.020$
IBM-UAV	$0.413 \pm 0.101$	$0.700 \pm 0.141$	$0.518 \pm 0.115$	$0.322 \pm 0.091$	$0.799 \pm 0.216$	$0.458 \pm 0.126$	$0.279 \pm 0.022$	$1.000 \pm 0.000$	$0.436 \pm 0.027$

considering it represents the number of similar requirements considered for conflict assessment. However, the analysis reveals comparable performance across

different thresholds for each n. For instance, in WorldVista, with  $T_o=0.75$  and  $n=\{1,2,3\}$ , Recall remains nearly 100%, while F1-scores decrease (66.7%, 65.4%, and 65.0%, respectively). Similar trends are observed in UAV and PURE datasets, where Recall is consistently 100%, but F1-scores decline with increasing n. Notably, OpenCoss displays the lowest F1-score and Precision, maintaining 100% Recall, resembling the performance observed in the S3CDA approach.

#### 4.3 Discussion

Table 11 outlines key comparisons between the proposed techniques. S3CDA exhibits superior overall performance, especially evident in conflict F1-score assessment. On the other hand, UnSupCDA excels in capturing true conflicts without the need for labeled training data, showcasing higher Recall. The two techniques each have distinct advantages and are applicable in different scenarios. Notably, UnSupCDA boasts versatility, making it easily applicable to SRS documents across various domains. In contrast, S3CDA's performance is influenced by the specific characteristics, structure, and word usage within the requirements.

	S3CDA	UnSupCDA
Applicability	Well-suited for scenarios with reliable labeled data	Applicable when labeled data for conflicts is scarce or unavailable
Dataset	Require labeled data	Requires known conflicts for evaluation
Requirements	More efficient with functional requirements	Work with both functional and non-functional
Human Expert	Low false predictions, limited human expertise required	Generates false predictions, demands additional software expertise
Complexity	Higher as it involves 2 stages	Lower
Experiment Results	Balanced Precision and Recall i.e, higher F1-scores	Higher Recall and lower Precision i.e., lower F1-scores

**Table 11**: Comparison of the proposed techniques

## 5 Threats to Validity

Several potential threats to the validity of our proposed conflict detection approaches in software requirements merit consideration. Firstly, the reliance on manual labeling and the introduction of synthetic conflicts may impact the generalizability of our findings. The effectiveness of our method could be influenced by the specifics introduced during manual labeling and the nature of synthetic conflicts, potentially limiting the broader applicability of the model. Additionally, the accuracy of the software-specific NER model poses a validity threat, as its performance directly influences the identification of entities

within software requirements. Any inaccuracies in entity extraction may lead to false positives or negatives in conflict detection.

Furthermore, the validity of our approach hinges on the efficacy of the sentence embedding models in accurately embedding software requirements. Variations in the structure, vocabulary, or complexity of requirements could impact the algorithm ability to discern meaningful similarities, affecting the overall success of conflict detection.

### 6 Conclusion

This study seeks to identify conflicts within software engineering requirements, recognizing their potential to significantly impede project success by causing delays throughout the entire development process. Prior research exhibits limitations in terms of generalizability, comprehensive requirement datasets, and clearly defined NLP-based automated methodologies. We develop a two-phase supervised (S3CDA) process for automatic conflict detection from SRS documents and an unsupervised variant (UnSupCDA) which works directly on natural language-based requirements.

Our experimental design aims to evaluate the performance of the S3CDA and UnSupCDA methodologies using five distinct SRS documents. The results demonstrate the effective conflict detection capabilities of both approaches across four datasets, with OpenCoss being an exception due to elevated structural similarities among requirements. S3CDA exhibits balanced F1-scores, while UnSupCDA attains 100% Recall across all datasets, albeit encountering challenges in Precision scores.

In future, we plan to extend this study by introducing more diverse SRS documents to further validate our proposed approach. We also concur that NLP domain is highly dynamic and new methods (e.g., embeddings) are developed at a fast pace. In this regard, we aim to extend our analysis by using other transformer-based sentence embeddings. Similarly, transformer-based NER models can be explored to improve the entity extraction performance.

## Statements and Declarations

No potential conflict of interest was reported by the authors.

## Data Availability Statement

Full research data can be accessible through following link. https://gitfront.io/r/user-9946871/tiYp38xXphd7/Paper-1-Req-Conflict/

## References

[1] U. S. Shah, D. C. Jinwala, Resolving ambiguities in natural language software requirements: a comprehensive survey, ACM SIGSOFT Software Engineering Notes 40 (2015) 1–7.

- [2] M. H. Osman, M. F. Zaharin, Ambiguous software requirement specification detection: An automated approach, in: 2018 IEEE/ACM 5th International Workshop on Requirements Engineering and Testing (RET), IEEE, 2018, pp. 33–40.
- [3] M. Aldekhail, A. Chikh, D. Ziani, Software requirements conflict identification: review and recommendations, Int J Adv Comput Sci Appl (IJACSA) 7 (2016) 326.
- [4] A. Egyed, P. Grunbacher, Identifying requirements conflicts and cooperation: How quality attributes and automated traceability can help, IEEE software 21 (2004) 50–58.
- [5] R. Merugu, S. R. Chinnam, Automated cloud service based quality requirement classification for software requirement specification, Evolutionary Intelligence 14 (2021) 389–394.
- [6] L. Zhao, W. Alhoshan, A. Ferrari, K. J. Letsholo, M. A. Ajagbe, E.-V. Chioasca, R. T. Batista-Navarro, Natural language processing for requirements engineering: A systematic mapping study, ACM Computing Surveys (CSUR) 54 (2021) 1–41.
- [7] S. Ezzini, S. Abualhaija, C. Arora, M. Sabetzadeh, L. C. Briand, Using domain-specific corpora for improved handling of ambiguity in requirements, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, 2021, pp. 1485–1497.
- [8] Y. Zhou, Y. Tong, R. Gu, H. Gall, Combining text mining and data mining for bug report classification, Journal of Software: Evolution and Process 28 (2016) 150–176.
- [9] K. Aggarwal, F. Timbers, T. Rutgers, A. Hindle, E. Stroulia, R. Greiner, Detecting duplicate bug reports with software engineering domain knowledge, Journal of Software: Evolution and Process 29 (2017) e1821.
- [10] W. Guo, L. Zhang, X. Lian, Automatically detecting the conflicts between software requirements based on finer semantic analysis, arXiv preprint arXiv:2103.02255 (2021).
- [11] T. Diamantopoulos, M. Roth, A. Symeonidis, E. Klein, Software requirements as an application domain for natural language processing, Language Resources and Evaluation 51 (2017) 495–524.
- [12] A. O. J. Sabriye, W. M. N. W. Zainon, A framework for detecting ambiguity in software requirement specification, in: 2017 8th International Conference on Information Technology (ICIT), IEEE, 2017, pp. 209–213.
- [13] D. Mairiza, D. Zowghi, N. Nurmuliani, Managing conflicts among non-functional requirements, in: Australian Workshop on Requirements Engineering, University of Technology, Sydney, 2009.
- [14] M. Kim, S. Park, V. Sugumaran, H. Yang, Managing requirements conflicts in software product lines: A goal and scenario based approach, Data & Knowledge Engineering 61 (2007) 417–432.
- [15] W. H. Butt, S. Amjad, F. Azam, Requirement conflicts resolution: using requirement filtering and analysis, in: International Conference on Computational Science and Its Applications, Springer, 2011, pp. 383–397.

- [16] T. Moser, D. Winkler, M. Heindl, S. Biffl, Requirements management with semantic technology: An empirical study on automated requirements categorization and conflict analysis, in: International Conference on Advanced Information Systems Engineering, Springer, 2011, pp. 3–17.
- [17] M. Aldekhail, M. Almasri, Intelligent identification and resolution of software requirement conflicts: Assessment and evaluation., Computer Systems Science & Engineering 40 (2022).
- [18] T. Viana, A. Zisman, A. K. Bandara, Identifying conflicting requirements in systems of systems, in: 2017 IEEE 25th International Requirements Engineering Conference (RE), IEEE, 2017, pp. 436–441.
- [19] M. Urbieta, M. J. Escalona, E. Robles Luna, G. Rossi, Detecting conflicts and inconsistencies in web application requirements, in: Current Trends in Web Engineering: Workshops, Doctoral Symposium, and Tutorials, Held at ICWE 2011, Paphos, Cyprus, June 20-21, 2011. Revised Selected Papers 11, Springer, 2012, pp. 278–288.
- [20] T. Moser, D. Winkler, M. Heindl, S. Biffl, Automating the detection of complex semantic conflicts between software requirements, in: The 23rd International Conference on Software Engineering and Knowledge Engineering, Miami, 2011.
- [21] M. Kamalrudin, J. Grundy, J. Hosking, Managing consistency between textual requirements, abstract interactions and essential use cases, in: 2010 IEEE 34th Annual Computer Software and Applications Conference, IEEE, 2010, pp. 327–336.
- [22] A. P. Felty, K. S. Namjoshi, Feature specification and automated conflict detection, ACM Transactions on Software Engineering and Methodology (TOSEM) 12 (2003) 3–27.
- [23] M. Heisel, J. Souquieres, A heuristic algorithm to detect feature interactions in requirements, in: Language Constructs for Describing Features: Proceedings of the FIREworks workshop, Springer, 2001, pp. 143–162.
- [24] C. D. Manning, M. Surdeanu, J. Bauer, J. R. Finkel, S. Bethard, D. McClosky, The stanford corenlp natural language processing toolkit, in: Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations, stanford, 2014, pp. 55–60.
- [25] U. Shah, S. Patel, D. C. Jinwala, Detecting intra-conflicts in non-functional requirements, International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 29 (2021) 435–461.
- [26] G. Abeba, E. Alemneh, Identification of nonfunctional requirement conflicts: Machine learning approach, in: International Conference on Advances of Science and Technology, Springer, 2021, pp. 435–445.
- [27] Z. Chentouf, Managing oam&p requirement conflicts, Journal of King Saud University-Computer and Information Sciences 26 (2014) 296–307.
- [28] D. Mairiza, D. Zowghi, V. Gervasi, Conflict characterization and analysis of non functional requirements: An experimental approach, in: 2013 IEEE 12th International Conference on Intelligent Software Methodologies, Tools and Techniques (SoMeT), IEEE, 2013, pp. 83–91.

- [29] V. Sadana, X. F. Liu, Analysis of conflicts among non-functional requirements using integrated analysis of functional and non-functional requirements, in: 31st annual international computer software and applications conference (COMPSAC 2007), volume 1, IEEE, 2007, pp. 215–218.
- [30] S. Das, N. Deb, A. Cortesi, N. Chaki, Sentence embedding models for similarity detection of software requirements, SN Computer Science 2 (2021) 1–11.
- [31] J. Cleland-Huang, M. Vierhauser, S. Bayley, Dronology: An incubator for cyber-physical system research, arXiv preprint arXiv:1804.02423 (2018).
- [32] A. Mavin, P. Wilkinson, A. Harwood, M. Novak, Easy approach to requirements syntax (ears), in: 2009 17th IEEE International Requirements Engineering Conference, IEEE, 2009, pp. 317–322.
- [33] A. Ferrari, G. O. Spagnolo, S. Gnesi, Pure: A dataset of public requirements documents, in: 2017 IEEE 25th International Requirements Engineering Conference (RE), IEEE, 2017, pp. 502–505.
- [34] C. Haskins, K. Forsberg, M. Krueger, D. Walden, D. Hamelin, Systems engineering handbook, in: INCOSE, volume 9, International Council on Systems Engineering Seattle, 2006, pp. 13–16.
- [35] A. Aizawa, An information-theoretic perspective of tf-idf measures, Information Processing & Management 39 (2003) 45–65.
- [36] D. Cer, Y. Yang, S.-y. Kong, N. Hua, N. Limtiaco, R. S. John, N. Constant, M. Guajardo-Céspedes, S. Yuan, C. Tar, et al., Universal sentence encoder, arXiv preprint arXiv:1803.11175 (2018).
- [37] N. Reimers, I. Gurevych, Sentence-bert: Sentence embeddings using siamese bert-networks, arXiv preprint arXiv:1908.10084 (2019).
- [38] L. McInnes, J. Healy, J. Melville, Umap: Uniform manifold approximation and projection for dimension reduction, arXiv preprint arXiv:1802.03426 (2018).
- [39] X. Gao, S. Wu, Hierarchical clustering algorithm for binary data based on cosine similarity, in: 2018 8th International Conference on Logistics, Informatics and Service Sciences (LISS), IEEE, 2018, pp. 1–6.
- [40] C. Rupp, M. Simon, F. Hocker, Requirements engineering und management, HMD Praxis der Wirtschaftsinformatik 46 (2009) 94–103.
- [41] G. Malik, M. Cevik, S. Bera, S. Yildirim, D. Parikh, A. Basar, Software requirement-specific entity extraction using transformer models, in: The 35th Canadian Conference on Artificial Intelligence, 2022.

## 7 Appendix

## 7.1 Overlapping Entity Ratio v/s Cosine Similarity

This section presents additional scatter plots related to Subsection 3.3.1.

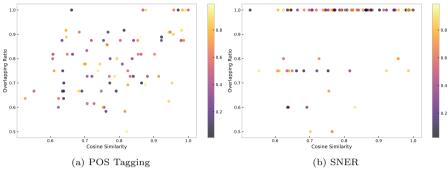


Fig. 5: Scatterplot depicting the relationship between Cosine Similarity and Overlapping Entity Ratio among UAV requirements.

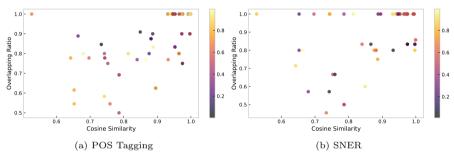


Fig. 6: Scatterplot depicting the relationship between Cosine Similarity and Overlapping Entity Ratio among PURE requirements.

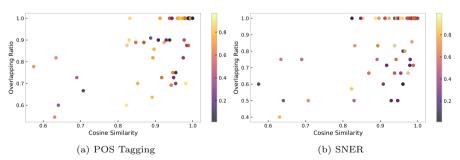
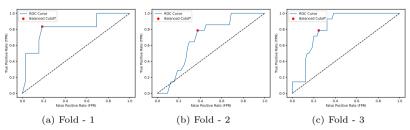


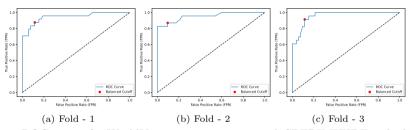
Fig. 7: Scatterplot depicting the relationship between Cosine Similarity and Overlapping Entity Ratio among IBM-UAV requirements.

#### 7.2 ROC Curves for Threshold Detection

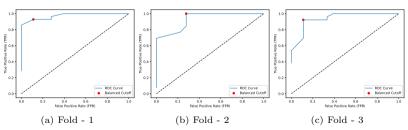
Figure 8,9,10, and 11 shows the ROC curves obtained from the 3-fold cross validation over all the requirement sets. These curves facilitate the process of finding the cosine similarity threshold in Phase I.



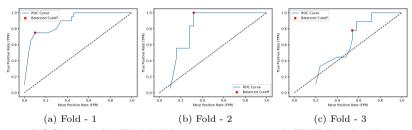
 $\textbf{Fig. 8} : \mbox{ROC curves for OpenCoss requirement set with TFIDF embedding across 3 folds}$ 



 ${\bf Fig.~9} : {\rm ROC}$  curves for WorldVista requirement set with SBERT-TFIDF embedding across 3 folds



 $\textbf{Fig. 10}: \ \, \textbf{ROC} \ \, \textbf{curves for PURE requirement set with SBERT embedding across 3} \\ \, \textbf{folds}$ 



 $\bf Fig.~11:$  ROC curves for IBM-UAV requirement set with SBERT embedding across 3 folds