

LiResolver: License Incompatibility Resolution for Open Source Software

Sihan Xu

xusihan@nankai.edu.cn
DISSec, NDST, College of Cyber
Science, Nankai University
Tianjin, China

Ya Gao

gaoya_cs@mail.nankai.edu.cn
DISSec, NDST, College of Cyber
Science, Nankai University
Tianjin, China

Lingling Fan*

linglingfan@nankai.edu.cn
DISSec, NDST, College of Cyber
Science, Nankai University
Tianjin, China

Linyu Li

linyuli@mail.nankai.edu.cn
DISSec, NDST, College of Cyber
Science, Nankai University
Tianjin, China

Xiangrui Cai

caixr@nankai.edu.cn
DISSec, NDST, College of Computer
Science, Nankai University
Tianjin, China

Zheli Liu

liuzheli@nankai.edu.cn
DISSec, NDST, College of Cyber
Science, Nankai University
Tianjin, China

ABSTRACT

Open source software (OSS) licenses regulate the conditions under which OSS can be legally reused, distributed, and modified. However, a common issue arises when incorporating third-party OSS accompanied with licenses, i.e., license incompatibility, which occurs when multiple licenses exist in one project and there are conflicts between them. Despite being problematic, fixing license incompatibility issues requires substantial efforts due to the lack of license understanding and complex package dependency. In this paper, we propose LiResolver, a fine-grained, scalable, and flexible tool to resolve license incompatibility issues for open source software. Specifically, it first understands the semantics of licenses through fine-grained entity extraction and relation extraction. Then, it detects and resolves license incompatibility issues by recommending official licenses in priority. When no official licenses can satisfy the constraints, it generates a custom license as an alternative solution. Comprehensive experiments demonstrate the effectiveness of LiResolver, with 4.09% false positive (FP) rate and 0.02% false negative (FN) rate for incompatibility issue localization, and 62.61% of 230 real-world incompatible projects resolved by LiResolver. We discuss the feedback from OSS developers and the lessons learned from this work. All the datasets and the replication package of LiResolver have been made publicly available to facilitate follow-up research.

CCS CONCEPTS

• **Software and its engineering** → **Collaboration in software development.**

*Lingling Fan is the corresponding author. Email: linglingfan@nankai.edu.cn

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ISSTA '23, July 17–21, 2023, Seattle, WA, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0221-1/23/07...\$15.00

<https://doi.org/10.1145/3597926.3598085>

KEYWORDS

Open Source Software, License, License Incompatibility Resolution

ACM Reference Format:

Sihan Xu, Ya Gao, Lingling Fan, Linyu Li, Xiangrui Cai, and Zheli Liu. 2023. LiResolver: License Incompatibility Resolution for Open Source Software. In *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA '23)*, July 17–21, 2023, Seattle, WA, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3597926.3598085>

1 INTRODUCTION

Open source software (OSS) significantly facilitates software development. Developers do not need to reinvent the wheel but only focus on their unique workflows and features. Despite the benefits, misuse of OSS might also induce security issues and legal risks [31, 32, 38, 39, 43]. Many issues arise from OSS licenses, which regulate the conditions under which OSS can be *legally* reused, distributed, and modified. Generally, an OSS license plays the role as a contract signed between the software owner and the user who reuses the code or binary. Once a user of OSS does not conform to the conditions dictated by OSS licenses, license violation occurs and legal risks are induced such as copyright infringement [40].

OSS licenses, which are legally binding, protect both the rights of software owners and the freedom of OSS. Nevertheless, a common issue arises when incorporating third-party OSS accompanied with licenses, i.e., *license incompatibility*, which occurs when multiple OSS licenses exist in the same project and there are conflicts between the conditions they state [14, 27]. For example, “CANNOT sublicense” declared by the license of a third-party component but “CAN sublicense” declared by the license of the whole project are incompatible with each other. In this case, conforming to the project license cannot ensure that the requirements declared by the license of the third-party component are satisfied simultaneously. License incompatibility can also occur among the licenses of multiple third-party components, e.g., “MUST contact author” from the license of one component but “CANNOT contact author” from the license of another component. In this case, at least one license is being violated due to the license incompatibility issue.

License incompatibility represents a serious threat to all stakeholders from a legal perspective. An empirical study on 1,846 GitHub projects shows that 72.91% of the investigated projects

suffered from license incompatibility [45]. Despite being problematic, fixing license incompatibility issues requires substantial efforts to either replace the reused code or migrate licenses. Worse still, developers struggle when multiple licenses are involved in the same project [1]. To address this issue, Kapitsaki et al. [25, 27] proposed a solution based on the Software Package Data Exchange (SPDX) specification [12]. They designed a directed graph with the compatibility relationships between licenses, based on which they proposed SPDX Violation Tools (SPDX-VT) to detect license incompatibility and recommend licenses to handle incompatibility issues. Besides, Liu et al. [33] also proposed Automatic License Prediction (ALP) to prevent license incompatibility issues induced by software changes.

Despite the progress, there are mainly three problems that limit the application of previous studies. **First**, previous studies only cover a small predefined set of licenses (e.g., 20 licenses supported by SPDX-VT [25] and 25 licenses supported by ALP [33]). However, there are 489 licenses and versions in SPDX [12], not to mention potentially a considerable number of license exceptions and custom licenses. **Second**, previous works only considered the names and versions of some well-known licenses, lacking the capability to comprehend license texts and understand the inherent reasons of license incompatibility, which limits their abilities to resolve incompatibility issues in a flexible and fine-grained way. **Third**, the solutions to license incompatibility have been mainly focused on recommending an official license as the new license for the whole project (i.e., the project license). Whereas, this paper observes that there may exist other licenses whose copyright holders are the same with the project license, and thus can also be changed by the licensor for resolving license incompatibility issues.

To address the aforementioned issues, it is desirable to provide a fine-grained, scalable, and flexible solution to resolve license incompatibility issues. To this end, we propose LiResolver, an automated tool to fix license incompatibility issues for open source software. Specifically, given an OSS, it first extracts all licenses accompanied with the project, and organizes them with a hierarchical structure that represents the licensing scope. Then, it provides a fine-grained understanding of the regulations stated by each license, based on which LiResolver detects and localizes license incompatibility issues. Finally, it resolves the constraints hidden in the licensing context, and provides flexible suggestions for users to choose from. Comprehensive experiments on real-world projects demonstrate the effectiveness of LiResolver, with 4.09% false positive (FP) rate and 0.02% false negative (FN) rate for incompatibility issue localization, and 62.61% success rate for resolving 230 incompatible OSS on GitHub. We made all datasets and the replication package publicly available to facilitate follow-up research.

In summary, we made the following novel contributions:

- We propose LiResolver, a fine-grained, scalable, and flexible approach to automatically interpret licenses, localize license incompatibility issues, and provide licensors with useful suggestions for resolving issues.
- Comprehensive experiments demonstrate the effectiveness of LiResolver, with 4.09% FP rate and 0.02% FN rate for incompatibility issue localization, and 62.61% of 230 real-world incompatible projects resolved by LiResolver.

- We further investigate the impacts of license hierarchy, copyright holders, and license exceptions on the effectiveness of license incompatibility resolution. We made all the datasets and the replication package publicly available to facilitate follow-up research [6].

2 APPROACH

2.1 Overview

This section details our tool, LiResolver, a hybrid approach that automatically understands license texts, detects license incompatibility, and provides useful suggestions for license incompatibility resolution. The input of LiResolver is the path of a folder where the target OSS is stored. As shown in Figure 1, given an OSS, LiResolver first extracts all involved licenses for further analysis, including the licenses accompanied with integrated software packages. After preprocessing the extracted licenses, the main components of LiResolver include three parts: (1) License understanding. As a fundamental step, LiResolver first understands the regulations implied by licenses by fine-grained entity extraction and relation extraction, so as to provide detailed information for further incompatibility analysis and resolution. (2) Incompatibility issue localization. Based on the fine-grained understandings of licenses, LiResolver detects and localizes license incompatibility issues by considering both the hierarchical structure and copyright holders of licenses. (3) License incompatibility resolution. Given a license involved in an incompatibility issue, LiResolver fixes the issue by resolving the constraints obtained from the parent and child licenses of the target license. Finally, LiResolver gives suggestions for license incompatibility resolution by recommending official licenses in priority. When none of existing licenses are compatible with the license context, LiResolver generates a custom license as an alternative solution. LiResolver stores the license incompatibility issues, their locations, the corresponding suggestions, i.e., the name of recommended official licenses or the full-text of custom licenses in an output file.

2.2 License Hierarchy Extraction

OSS, especially in the form of frameworks or libraries, is often used as a component of software products, which can be recursively integrated into larger software. The hierarchical structure of OSS accompanied with licenses leads to the hierarchical structure of licenses. Previous studies regarded all licenses in a project independently and equally [25, 27, 45], regardless of their licensing scopes. For instance, Figure 2 illustrates a real-world example from Flask JSONDash [9], a popular OSS to create chart dashboards. In this example, the project license, which states the regulations for the whole project, is the MIT license in the main directory of the project. As shown in Figure 2(a), there are four component licenses in the modules or integrated third-party packages which are incompatible with the project license. The reason is that the component licenses are more restrictive than the project license, so that anyone who conforms to the project license (i.e., the MIT license) may violate the component licenses (e.g., Apache-2.0). To address this issue, one should either replace the third-party software packages or change the project license in the root directory. In this paper, we observe that the location of a license determines its scope of licensing, which affects license incompatibility detection and resolution. Figure 2(b)

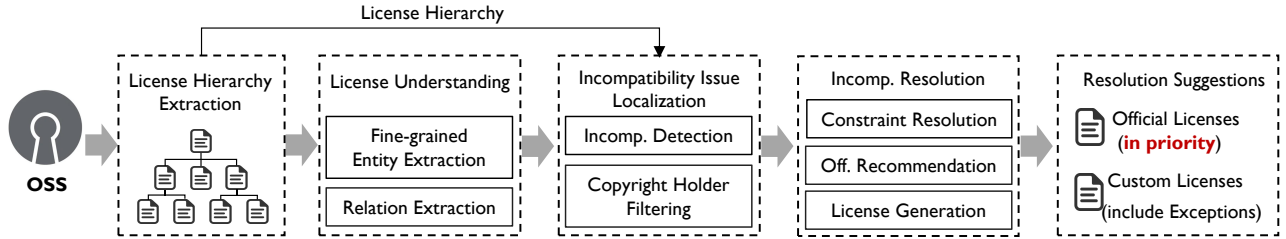
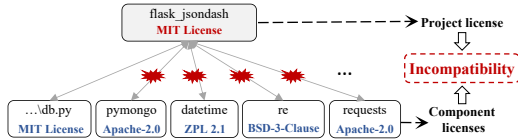
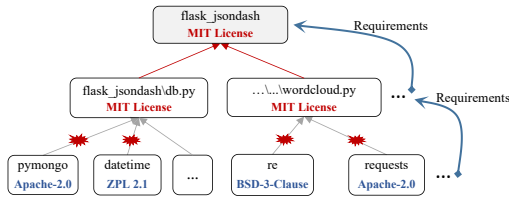


Figure 1: Overview of LiResolver



(a) License Compatibility Analysis without Hierarchy Extraction



(b) License Compatibility Analysis with Hierarchy Extraction

Figure 2: Illustration of License Hierarchy

depicts a part of the license hierarchy of JSONDash [9]. It can be seen that *pymongo* and *datetime* are two third-party packages imported by the source code file *db.py*. However, their licenses (i.e., Apache-2.0 and ZPL-2.1) are incompatible with the license of *db.py* (i.e., the MIT license). In this case, even if the project owner changes the project license, license incompatibility still remains.

To address this issue, we extract the hierarchical structure of all involved licenses in a project. Specifically, LiResolver first builds the hierarchical structure according to the file structure. Then, it extracts three types of licenses from the project, i.e., the declared licenses, the referenced licenses, and the inline licenses. The *declared licenses* state rights and obligations explicitly in license files (e.g., LICENSE). The *referenced licenses* are obtained by their names, versions, direct links, or the websites of incorporated third-party packages. The *inline licenses* typically appear at the beginning of source code files. Different licenses may have different licensing scopes. For example, an inline license is only responsible for the source code within the same source code file. The scope of the license accompanied with a third-party software package is only the package itself. After obtaining the original file structure, LiResolver removes the nodes without licenses from bottom to up and links its children nodes to its parent node if necessary. Finally, the licensing scopes of all licenses can be obtained by the hierarchy as shown in Figure 2(b), as well as the dependencies between them, which can be further used for license incompatibility resolution.

2.3 License Understanding

After obtaining the license hierarchy, LiResolver *automatically* understands the semantics of each license as a fundamental step. There

are mainly two reasons for automating license understanding. (1) There exist a variety of official licenses (i.e., 552), and the number keeps growing [13]. However, manually extracting the detailed information from long and complicated licenses is time-consuming and labor-intensive. (2) A previous study investigated 1,846 projects and obtained 5,777 unique licenses, among which 24.56% were custom licenses with flexible expressions [45]. Motivated by them, it is valuable to automatically understand licenses to facilitate further analysis. In this paper, we cast license understanding as an **information extraction** problem, which can be further decomposed into an entity extraction task and a relation extraction task.

2.3.1 Fine-grained Entity Extraction. Unlike previous works that regarded each right/obligation as a single entity [45], in this paper, we observe that each right/obligation could be decomposed into four types of entities (i.e., *action*, *object*, *attitude*, and *condition*) and construct a fine-grained structure to model it. Each entity, regardless of its type, is a flexible expression that describes a specific action, object, attitude, or condition in the form of a word, a phrase, or even a sentence in a license. For instance, from “*you are allowed to distribute modified works*”, the previous study [45] identified the term “Distribute” and inferred an attitude CAN. However, LiResolver extracts the action (“distribute”), the object (“modified works”), the attitude (“are allowed to”), and the condition if there exists. The definitions of four types of entities are as follows.

Action: the process of exercising a granted right or performing an obligation enforced by a license, e.g., *Disclose Source*.

Object: the object to be acted upon, which is typically a software artifact such as source code or binary.

Attitude: the attitudes of licensors towards licensees to grant or reserve rights, or enforce obligations.

Condition: the obligations that licensees must comply with to exercise the granted rights.

To extract the fine-grained entities from licenses, we employ sequence labelling to identify the beginning and last token of each entity. As a fundamental step, we first split each license text into sentences and preprocess them by removing non-textual parts, checking spellings, performing stemming and morphological. In the training phase, we label each token in a license as illustrated in Figure 3, where B-X and I-X represent the first token (i.e., the beginning token) and an inside token of an entity whose type is X, respectively. O denotes a token outside any named entities. Note that unlike LiDetector [45] which used X to represent the X^{th} license term, in this work, we use X to represent the type of the entity. For instance, B-object in Figure 3 indicates that the token *source* is at the beginning of an entity whose type is *object*.

Based on the labelled dataset, we trained a NER model based on spaCy [41] to identify and localize named entities. We first embed

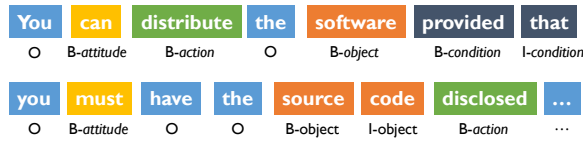


Figure 3: Illustration of the BIO Mode for Entities in Licenses

words into vector with a pre-trained *roberta-base* model [34] to employ its prior knowledge on word semantics and boost the word embedding strategy using subword features. Then, we feed the results of word embeddings into a transition-based parser, which utilizes a deep convolutional neural network with residual connections to learn the representation of each license sentence and parse named entities. In the inference phase, LiResolver sequentially predicts the label of each token in license text so as to identify and localize entities in licenses. For the example in Figure 3, LiResolver extracts seven entities, i.e., two *action* entities (*distribute* and *disclose*), two *object* entities (*software* and *source code*), two *attitude* entities (*can* and *must*), and a *condition* entity (*provide that*).

2.3.2 Relation Extraction. Relation extraction in license understanding aims to extract structured knowledge about relations between entities from unstructured license text. Unlike previous studies [28, 45] which regarded a right/obligation as a whole entity, in this paper, we decompose each regulation into four types of entities, and propose a relation extraction model to reconstruct the relations between these entities for a fine-grained and structured understanding of licenses. Since there exists no tagged dataset for entity relations in licenses and labeling is time-consuming, inspired by the promising results of prompt-tuning for few-shot tasks, we employ a prompt-based relation extraction model based on KnowPrompt [7]. The basic idea is to leverage the knowledge hidden in the pre-trained language model and formalize specific tasks as cloze tasks. With the prompt template, LiResolver predicts the mask and maps the prediction result to the classes, which are the relations between input entities. Specifically, we extract five types of relations, i.e., *action-object*, *action-attitude*, *action-condition*, *condition-action*, and *others*. *Action-object* represents the relation between an action and the object to be acted upon. *Action-attitude* represents the relation between an action and the attitude towards this action. *Action-condition* represents the relation between an action and the condition of exercising the action. *Condition-action* denotes the relation between a condition and the action required in the condition. Finally, *others* represent the relations other than the aforementioned relations. With these specific relations, LiResolver organizes the extracted entities and constructs a fine-grained understanding of licenses.

Figure 4 illustrates the model construction of relation extraction. It can be seen that given a license phrase “you can distribute the software”, LiResolver first extracts the entities *can*, *distribute*, and *software*. Then, taking the relations between *distribute* and *software* as an example, LiResolver constructs a prompt template as “[CLS] you can [E1] **distribute** [/E1] the [E2] **software** [/E2]. [SEP] [act] **distribute** [/act] [MASK] [obj] **software** [/obj] [SEP]”, and feeds it to the Mask Language Model (MLM) [10]. Here, [CLS] and [SEP] mark the input sequence and its prompt, [E1], [/E1], [E2], and [/E2] mark the beginning and end tokens of two input entities, [act], [/act], [obj], and [/obj] represent virtual types of entities learned

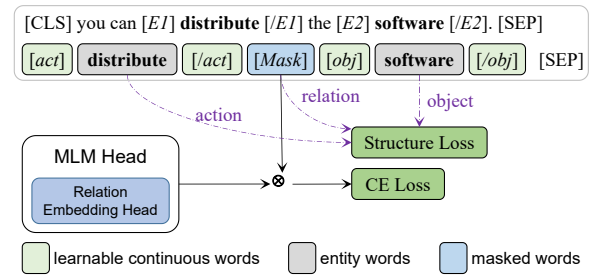


Figure 4: Illustration of Relation Extraction from Licenses

during the prompt tuning process, and [MASK] denotes the output prediction of the MLM model. After that, LiResolver obtains the probability distribution of [MASK] over the vocabulary list, and maps the prediction result to the embeddings of answer words, which are the classes of relations between input entities. Finally, LiResolver achieves the prediction result as the relation between two entities *distribute* and *software*. Note that for official licenses whose rights and obligations have already been known, we only extracted fine-grained entities and relations automatically in the preparation phase, and manually checked them to construct a database that can be directly used for further analysis. The details of the database can be seen in Section 2.5.2.

2.4 Incompatibility Issue Localization

2.4.1 Incompatibility Detection. Based on the outputs of license understanding and hierarchy extraction, LiResolver identifies incompatibility issues in a given project together with the specific localization, i.e., the packages or files involved in the incompatibility issues. Previous works only distinguished between the project and component licenses [45], in this paper, we extract license hierarchy to obtain the licensing scope of each license, as well as the relations between multiple licenses. Then, we can obtain *parent-child* pairs from the extracted license hierarchy. Given a license l and its child nodes l_c , we define that the license l is compatible with its child nodes l_c if anyone who complies with the license l will not violate any of its child nodes l_c . Otherwise, there exists license incompatibility. In other words, a parent license l is supposed to be equally or more restrictive than any of its child nodes l_c . On this basis, we propose to detect license incompatibility from bottom to up, layer by layer according to the extracted license hierarchy, so as to conduct a fine-grained license incompatibility detection.

For each *parent-child* pair of licenses, we extract the fine-grained entities and relations as mentioned in Section 2.3. Then, a fundamental step for incompatibility analysis is semantic alignment. Specifically, for each *attidue* entity (e.g., “are not allowed”), we follow previous works [29, 45] to assign it a label (i.e., CAN, CANNOT, or MUST) for further analysis. We also classify the entities of *actions* into 23 categories as previous studies [25, 28, 45]. Each group represents a type of actions that licensees may do. To this end, we trained a classification model with the Support Vector Machine algorithm [22] based on the dataset provided in [45]. To conduct an accurate analysis, we further distinguish the same *actions* with different *objects* into different groups. For example, *distribute source code* and *distribute binaries* are regarded as two different regulations.

By this means, we group the regulations from the *parent-child* pair of licenses. Each group contains the regulations about the

same *action-object* but multiple attitudes, and LiResolver detects whether there exist conflict attitudes between a parent license and its child license towards the same *action-object*. For example, if LiResolver infers from a parent license that states *can distribute source code*, while infers from one of its child licenses that claims *cannot distribute source code*, in this case, anyone who conforms to the parent license may violate at least one of its child licenses by distributing source code. As for the *conditions*, if a *condition* entity is extracted as well as its relations with other entities, we consider the conditional cases separately, i.e., separately assume two cases where the condition is *True* or *False*, respectively. By this means, we can identify incompatibility issues and the detailed information about such issues, including the incompatible parent and child nodes (i.e., licenses and the packages/files where they are located), the incompatible *action-object* pairs, the *attitudes* and *conditions* (optional) of the involved licenses, which are further utilized for license incompatibility resolution. Note that as previous studies [8, 45], the absence of a project license implies all rights are reserved, i.e., nobody can copy, distribute, or modify the work. In this case, LiResolver detects whether the licenses at the first level of the hierarchy can be incorporated into the same project without license violations, and reports incompatible licenses.

2.4.2 Copyright Holder Filtering. Previous research mainly focused on the project license whose licensing scope is the whole project, while ignoring other licenses in the project [25, 33, 45]. However, in this paper, we notice that there are often more than one license that can be changed for incompatibility resolution. The intuition behind is that if a license has the same copyright holder with the project license, the project owner ought to have the right to modify or replace it. Actually, it is a common practice for developers of large-scale OSS to place their own license claimers in the modules implemented by themselves for copyright protection and ownership assertion of these modules.

For the above reason, we propose to identify copyright holders so as to enlarge the scope of licenses that can be modified for license incompatibility resolution. Since we assume that a project owner can always modify the project license, LiResolver first identifies the copyright holder stated by the project license. Then, it examines all licenses involved in license incompatibility issues, so as to localize the module licenses (also known as component licenses) with the same copyright holder as the project license, which indicates that these licenses can also be modified by the project owner.

Specifically, we first identify the license sentences that claim copyright holders. Since the expressions of copyright holders are relatively fixed, we define a set of phrases as the signals of copyright claim, for instance, *copyright (c)* and *authored by*. The specific expressions can be seen online [6]. We detect these keywords by regular matching. After identifying the copyright-related sentences, we utilize Stanford CoreNLP [21] to recognize named entities related to *person names* and *organization names* as the copyright holder of the input license.

For the example in Figure 2(b), LiResolver first identified the copyright holder *Chris Tabor* claimed in the project license. Then, it detected 10 component licenses inside the project that have the same copyright holder with the project license, such as the inline licenses (i.e., the MIT license) at the beginning of source code files

Algorithm 1: License Incompatibility Resolution

Input: X : the target project with license incompatibility issues
Output: S : the resolution suggestions of the incompatibility issues

```

1  $LR \leftarrow detectCompt(X)$  // obtain incompatible licenses that can be
   changed.
2  $S \leftarrow \emptyset$  // resolution of the project.
3 foreach  $l \in LR$  do
   // iterate licenses in LR from bottom to top.
4    $l_p \leftarrow getParent(X, l)$ 
5    $r_p \leftarrow getReq(l_p)$ 
6    $l_c \leftarrow getChildren(X, l)$ 
7    $r_c \leftarrow \emptyset$ 
8    $f_l \leftarrow False$  // initialize the flag for exception.
9   foreach  $c \in l_c$  do
10     $r \leftarrow getReq(c)$ 
11     $r_c, f \leftarrow mergeReq(r_c, r)$ 
       // merge the requirements of child nodes.
12     $f_l \leftarrow f_l \vee f$  // if there are conflicts between child nodes.
13  $R_l \leftarrow resolve(r_p, r_c)$  // resolve the constraints for license  $l$ , i.e., no
   more restrictive than  $r_p$  but no more loose than  $r_c$ .
14 if  $R_l$  is None then
15   continue // if the constraints cannot be resolved, continue.
16 else
17    $s_{l.o} \leftarrow checkOfficial(R_l)$ 
18   if  $s_{l.o}$  is None then
19      $s_{l.c} \leftarrow generateCustom(R_l)$ 
20     if  $f_l$  is TRUE then
21        $s_{l.c} \leftarrow addException(R_l, s_{l.c})$ 
22    $S \leftarrow S \cup s_l$ 
23 return  $S$ 

```

db.py and *wordcloud.py*. In this case, there exist more than one license that can be modified to resolve license incompatibility issues.

2.5 License Incompatibility Resolution

2.5.1 License Constraint Resolution. As described in Algorithm 1, taking the project with license incompatibility issues as input, LiResolver obtains a set of problematic licenses (denoted by LR) which can be modified or replaced to handle the incompatibility issues (Line 1). Based on hierarchy extraction mentioned in Section 2.2, LiResolver iterates the licenses in LR from bottom to top for license incompatibility resolution (Lines 3–23). Specifically, for each license in LR , LiResolver first obtains its parent and child nodes, and then computes the requirements brought by them (Lines 4–12). Since there may exist multiple child nodes for a license, we merge the requirements from all child nodes by extracting the most restrictive attitude towards each *action* (Line 11). Note that to obtain constraints, we only take into account the licenses without the same copyright holder with the project license (i.e., the licenses that cannot be modified by the software owner). The intuition behind is that *licenses with the same copyright holder of the project license can be changed to address incompatibility issues, and thus should not be the constraints for the current license/node*. For the example shown in Figure 2(b), there exist three licenses that can be modified (i.e., one project license and two inline licenses stated in *db.py* and *wordcloud.py*). In this case, LiResolver resolves the incompatibility issues from bottom to top, and starts from a MIT license in a Python file. Then, it computes the requirements from two child licenses, which are imported by the third-party software packages. As for the parent license, it can be seen that although the current license should be no more restrictive than its parent node (i.e., the project

license in this case), however, the project license can be changed later to satisfy the requirements of compatibility, which should not be considered as a source of constraints for the current node.

Based on the license *actions*, *attitudes*, *conditions*, and *objects* extracted in the phase of license understanding, LiResolver resolves the constraints brought by its parent and child nodes. The basic idea is that the current license l should be no more restrictive than its parent license l_p but no more loose than its child nodes l_c . Specifically, given a license l , LiResolver merges the requirements from its child nodes by grouping the *action-object* pairs mentioned in Section 2.4.1, so that the requirements towards the same *action-object* are within the same group. Then, for each *action-object* pair, LiResolver merges the requirements from the children by obtaining the most restrictive attitude towards the same *action-object*, which serves as a constraint for this *action-object*. Similarly, LiResolver obtains the requirements from the parent license, which serve as another boundary of the constraints. By this means, LiResolver achieves the constraints from the context and stores them in R_l .

2.5.2 Official License Recommendation. After license constraint resolution, we obtain the constraints for a current node, which is a license involved in license incompatibility and can be changed by the project owner. However, we note that not all constraints from the context can be resolved. According to the “one-way-compatibility” definition, given a license, if it has a child license which is more permissive than its parent license, the requirements from the context of the current license (i.e., its parent and child licenses) cannot be satisfied simultaneously. In this case, the constraints for fixing the incompatibility issue cannot be resolved, LiResolver moves on to the next target license. Hence, as described in Algorithm 1, LiResolver only provides resolutions for licenses whose constraints from their context can be resolved by LiResolver (Lines 17–22). Specifically, after resolving the constraints, LiResolver first checks whether there exist official licenses that can satisfy the requirements. In this paper, official licenses represent the licenses from the Software Package Data Exchange (SPDX) [12] which are publicly available. Note that the high priority of recommending official licenses comes from an empirical study that reported official licenses counted for more than half of the investigated licenses rather than custom licenses [45]. However, LiResolver also provides settings for users to set priority of resolutions.

To recommend official licenses for license incompatibility resolution, we constructed a database of official licenses whose rights and obligations have already been known by previous studies [29, 45]. We manually checked the license terms labelled by previous studies (e.g., CAN *distribute*), extracted the fine-grained entities and relations mentioned in Section 2.3 using LiResolver, and verified the labels by three authors and a lawyer to construct the database. Finally, when more than one official licenses satisfy the constraints, we rank these licenses according to their similarities to the original license via cosine similarity.

2.5.3 Custom License Generation. Except official licenses, developers are also allowed to customize their own licenses, which are named by *custom licenses* in this paper. As described in Algorithm 1, given a set of constraints, it is possible that none of existing official licenses can satisfy the constraints R_l . In this case, as an alternative way, we propose to generate a custom license for the project owner

to choose (Lines 18–21). We note that although there are no conventions for custom licenses to comply with, however, as a type of contracts between the owner and users of projects, the expressions of licenses need to be precise and strict in most cases. For this reason, despite the remarkable progress of text generation techniques in natural language processing [46], we do not utilize machine learning-based methods for license text generation. Instead, in this paper, we propose a template-based approach to generate custom licenses according to the constraints and fine-grained information obtained in the aforementioned steps.

As displayed in Algorithm 1, R_l represents the results of constraint resolution for a target license l . Then, LiResolver *automatically* generates a custom license by organizing the fine-grained entities and their relations stored in R_l . Specifically, for each *action* in R_l , if there exists an *object* and an *action-object* relation between them, then LiResolver concatenates them directly or with a preposition determined by CoreNLP [21]. As for the *attitude*, we note that the results of constraint resolution may contain multiple attitudes toward the same *action-object*. For instance, given a parent node with the attitude MUST and a child node with the attitude CAN towards the same *action-object*, the result of constraint resolution for the current node contains a set of *attitudes* [CAN, MUST] for this *action-object* pair, which is no more restrictive than the parent node and no more permissive than the child node. In this case, LiResolver uses the most restrictive *attitude* in the set by default. In addition, when a *condition* entity is found along with its relation with the *action* (i.e., *action-condition*), LiResolver adds the *condition* at the end of the current sentence, equipped with the *action*, *attitude*, and *object* (optional) orderly found by the relations *condition-action*, *action-attitude*, and *action-object*. For instance, after extracting an *action* entity “modify” and an *object* entity “source code”, as well as the *action-object* relation between them, LiResolver concatenates them and generates “modify source code” as a result. If there exists a *condition* entity “as long as” in R_l and an *action-condition* relation between “modify” and “as long as”, LiResolver extends the statement and generates “modify source code as long as ...”. Finally, to generate a clear and concise license, LiResolver merges the sentences with the same *objects* (e.g., you can *distribute* and *modify* copies of software).

Exceptions. In this paper, we observe that there may exist some conflicts between the requirements of multiple child licenses which cannot be modified by the project owner. A natural solution to this problem is to replace the third-party packages involved in the conflicts. However, since the replacement of a third-party package requires substantial efforts due to the complex package dependency, we also provide an alternative way by attaching an *exception*, which grants an exception to the license or some additional permissions (e.g., LLVM Exceptions [5]). Specifically, LiResolver first localizes the files, packages or directories of the conflict licenses. Then, towards the conflict *action* and *object* (optional), it assigns different attitudes with different licensing scopes in the same license. By this means, LiResolver attaches an *exception* to the target license, so as to mitigate the incompatibility issue.

3 EVALUATION

In this section, we present the evaluation results of LiResolver on real-world datasets and answer the following research questions.

RQ1: Can LiResolver provide a fine-grained understanding of license texts by effectively extracting entities and their relations?

RQ2: Can LiResolver outperform state-of-the-art approaches in detecting and localizing license incompatibility issues?

RQ3: How effective is LiResolver for resolving license incompatibility issues in real-world OSS?

3.1 Evaluation on License Understanding

In this section, we evaluate the performance of LiResolver on fine-grained entity extraction and relation extraction. For both tasks we performed labelling, training, and testing in **sentences**. The sentence datasets were constructed based on the dataset from a previous study [45]. To reduce bias, we did not set aside any licenses when evaluating the effectiveness of license understanding provided by LiResolver. In total, there were 212 official licenses and 188 custom licenses, comprising 21,844 license sentences (i.e., 11,973 from *ilrlegal* [29] and 9,871 from GitHub). We manually labelled the entities and relations. During labeling, we took Fleiss’ Kappa to assess the reliability of agreement among the raters. The result is 0.83, indicating that the raters reach high agreement. We randomly split each dataset into the training, validation and testing datasets by 3:1:1.

3.1.1 Entity Extraction. To evaluate the performance of LiResolver for entity extraction, we compare it with two state-of-the-art tools (i.e., **FOSS-LTE** [28] and **LiDetector** [45]) and two natural language processing techniques (i.e., **regular matching** [4] and **semantic similarity** [17]). Specifically, FOSS-LTE [28] employed a topic model named Latent Dirichlet Allocation (LDA) [24] to identify license terms. LiDetector [45] utilized sequence labelling and sentiment analysis to identify license term entities (similar to the *actions* in this paper) and attitudes. Regular matching [4] was implemented by pre-defining a set of keyword patterns for entity extraction. Semantic similarity [17] utilized the word2vec pre-trained language model [35, 36] to measure the cosine similarity between the tokens and predefined keywords. Note that the outputs of FOSS-LTE are license term phrases mapped from topic sentences (e.g., MayGrant-Patents). To conduct a fair comparison of license understanding, we split the outputs of FOSS-LTE into four types of entities as other baselines. Similarly, LiDetector naturally extracts license terms (e.g., Commercial Use), attitudes (e.g., CANNOT), and conditions (e.g., *provided that*) from license texts, which can be compared with the *actions*, *attitudes*, and *conditions* in this paper. As for regular matching and semantic similarity, we analyzed 900 license sentences and summarized a set of keyword patterns for each type of entities as shown online [6].

Table 1 displays the experimental results of entity extraction for five approaches. Note that since LiDetector does not extract the *objects of actions* in license texts, the results of extracting *objects* by LiDetector are denoted by “-”. In total, LiResolver achieves 74.07% precision and 77.36% recall for entity extraction. It can be seen that the precision of LiResolver is much higher than that of the baselines by at least 26.03%, 53.47%, 3.95%, and 5.2% when extracting *actions*,

Table 1: Comparison on Fine-grained Entity Extraction

Tool	Actions		Objects		Attitudes		Conditions		Total	
	P (%)	R (%)	P (%)	R (%)	P (%)	R (%)	P (%)	R (%)	P (%)	R (%)
Reg. Match.	54.76	65.07	5.32	22.38	2.58	47.61	46.87	61.64	14.74	49.41
Similarity	50.81	59.96	3.18	12.27	2.35	40.47	46.86	61.65	14.28	42.33
FOSS-LTE	35.98	23.12	5.76	19.13	3.06	26.19	66.67	21.92	14.30	21.12
LiDetector	43.14	55.79	-	-	74.74	56.35	90.25	86.30	47.79	32.80
LiResolver	80.79	84.86	59.23	63.83	78.69	76.80	95.45	87.50	74.15	77.31

objects, *attitudes*, and *conditions*, respectively. It indicates that LiResolver can precisely extract entities in license texts without much noise. However, regular matching and similarity-based approaches achieve low precision especially for *objects* and *attitudes*. By analyzing the results, a possible explanation could be the ambiguous expressions for the two types of entities. Moreover, the recall scores achieved by LiResolver are also higher than those of four baselines by at least 19.79%, 41.45%, 29.19%, and 1.2% when extracting *actions*, *objects*, *attitudes*, and *conditions*, respectively. Among four baselines, LiDetector achieves a competitive performance with LiResolver on extracting *attitudes* and *conditions*. However, the precision and recall scores of LiDetector on extracting *actions* are lower than those of LiResolver. By analyzing the entities extracted by LiDetector, we found that the entities extracted by LiDetector are more coarse-grained, sometimes along with its *objects* or *attitudes*, which made them longer than the entities extracted by other approaches, which decreased the evaluation results of LiDetector.

3.1.2 Relation Extraction. Figure 5 shows the performance of two approaches on relation extraction. Since FOSS-LTE [28] and LiDetector [45] do not extract relations between entities, we did not compare LiResolver with them. Instead, we implemented a heuristic approach based on the positions and types of entities (denoted by **PTE**). Specifically, it determined entity relations according to their types and relative positions in the same sentence. For instance, *action-object* indicates an *action* in front of *object*. Similarly, *action-condition* and *condition-action* were also distinguished by the relative positions of the *action* and *condition*. Finally, to conduct a fair comparison, we fed both approaches the same pairs of entities as input. Since the input of entity relation is a pair of entities, there exist no false negatives and thus we evaluate the performance of two approaches only with the accuracy metric. It can be seen that LiResolver can effectively identify entity relations, achieving 95.69% accuracy in total. We can also observe that the performance of LiResolver is superior to that of the heuristic baseline for all types of relations. The baseline only achieves comparable performance with LiResolver on *action-condition* and *condition-action*. It indicates that the relation between an *action* and a *condition* can be easier inferred by their relative positions in the sentence compared to other types of entity relations. However, since PTE requires the types of entities as input, its performance relies on the accuracy of entity extraction, while LiResolver has no requirements for such information.

Answer to RQ1: LiResolver can effectively extract fine-grained entities and their relations from license texts, with 74.07% precision and 77.36% recall for entity extraction, 95.09% accuracy for relation extraction, which provides detailed information for license incompatibility localization and resolution.

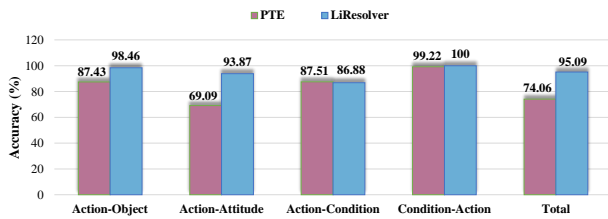


Figure 5: Accuracy of Relation Extraction

3.2 Evaluation on Incompatibility Localization

3.2.1 Setup. We evaluate the performance of LiResolver on 300 real-world projects on GitHub. Specifically, we first crawled 10,000 popular projects with more than 1,000 stars. Then, we randomly selected 300 projects and constructed a ground-truth dataset. Three authors and a lawyer were involved in labelling and verifying the dataset which has been made publicly available online [6]. In total, we extracted 81,402 licenses for incompatibility analysis. We compare LiResolver with four state-of-the-art tools, i.e., the SPDX Violation Tools (SPDX-VT) [27], Ninka [15] equipped with *tldrlegal* [29], Librariesio [30], and LiDetector [45]. SPDX-VT [27], which designed a directed graph to represent the compatibility relations between a set of licenses in Software Package Data Exchange (SPDX) [12]. Ninka [15], which identified licenses with a sentence-matching method. Combined with the summarized rights and obligations provided by *tldrlegal* [29], Ninka can be utilized to detect incompatibilities between licenses. Librariesio [30], which checked compatibility between a predefined set of SPDX licenses. LiDetector [45], which firstly introduced a learning-based method to automatically identify license terms from arbitrary licenses and detect license incompatibilities.

3.2.2 Results. Table 2 displays the experimental results of license incompatibility localization. The results show the number of incompatible projects and incompatible issues reported by each approach. Note that in this paper, an **incompatible issue** represents a pair of licenses which are incompatible with each other and should be resolved to reduce legal risks. We also report the false positive (FP) rate and the false negative (FN) rate from both the project-level and the issue-level for five approaches. From Table 2, it can be seen that LiResolver is superior to all the baselines in terms of both FP and FN rates. Specifically, LiResolver reported 231 incompatible projects with 0.43% FP rate and no false negatives, which demonstrated its ability to detect projects with license incompatibility.

It can also be seen that Ninka equipped with *tldrlegal* and SPDX-VT have low FP rates (i.e., 2.94% and 2.63%, respectively) but high FN rates (i.e., 42.61% and 67.83%). By analyzing their detection processes and results, we found that both of them complied with strict and predefined rules for license compatibility analysis, which contributed to the low FP rates. However, relying on license identification and manually defined compatibility rules, the applications of these two approaches were limited to a small set of official licenses, while ignoring other licenses in the projects, which causes their high FN rates. Similarly, the performance of Librariesio is close to that of Ninka and SPDX-VT. It categorized 35 official licenses into five classes (e.g., *permissive licenses*), and manually defined a heuristic rule based on the classification. It achieved a lower FN rate (i.e., 33.91% at the project-level) than Ninka and SPDX-VT, since it took

Table 2: Comparison on License Incompatibility Localization

Tool	#Incomp. Projects			#Incomp. Issues		
	#Report.	FP (%)	FN (%)	#Report.	FP (%)	FN (%)
Ninka	136	2.94	42.61	32,289	86.46	86.31
SPDX-VT	76	2.63	67.83	54,530	65.33	99.28
Librariesio	171	11.11	33.91	18,874	36.22	68.32
LiDetector	228	3.51	4.35	56,634	81.04	66.42
LiResolver-Flat	230	2.61	2.61	25,238	33.81	47.75
LiResolver	231	0.43	0.00	33,340	4.09	0.02

more licenses into consideration. However, due to the course-fined classification of licenses, it achieved the highest project-level FP rate among all the approaches (i.e., 11.11%).

Among four baselines, LiDetector achieved the closest performance to LiResolver for resolving license incompatibilities. Specifically, it reported 228 incompatible projects with 8 false positives and 10 false negatives, mainly caused by its coarse-grained license term extraction which did not distinguish between the same *actions* with different *objects* (e.g., *Distribute source code vs. Distribute binaries*). Moreover, like other baselines, it also ignored the license hierarchy of the projects.

From the perspective of incompatibility issues, it can be seen from Table 2 that LiResolver reported 33,340 incompatibility issues with 4.09% FP rate and 0.02% FN rate, which demonstrated the ability of LiResolver to detect license incompatibility issues. We can also observe that the FP and FN rates of four baselines, ranging from 36.22% to 99.28%, are much higher than those of LiResolver, ranging from 0.02% to 4.09%. By analyzing the results, we found that the false negatives of SPDX-VT, and Librariesio were mainly due to their limited scope of licenses that can be handled. LiDetector localized more incompatibility issues than the other three baselines, due to the flexibility adaption to arbitrary licenses. However, it treated all component licenses equally, and ignored the potential hierarchy among them, which contributed to the false negatives.

• **Effects of hierarchy extraction.** To illustrate the usefulness of the license hierarchy on incompatibility detection, we also implemented LiResolver without hierarchy extraction, denoted by LiResolver-Flat in Table 2. It can be seen that without hierarchy extraction, LiResolver reported 230 incompatible projects with 6 false positives and 6 false negatives, similar to the results of LiDetector. The reason behind is that both LiResolver-Flat and LiDetector treated all component licenses in modules or accompanied with third-party packages equally important regardless of their licensing scopes. To be specific, LiResolver-Flat reported a project with licence incompatibility issues when there existed a component license that was more restrictive than the project license. However, LiResolver reported an incompatible project when there existed a child license more restrictive than its parent license. It can be inferred that given the same project and the same results of license understanding, the true positives found by LiResolver-Flat can theoretically be found by LiResolver. Nevertheless, when the parent license is not the project license, but a license in a module of the project, the parent and child licenses are regarded as two independent licenses by LiResolver-Flat and LiDetector, which might lead to a false negative. On the other hand, it can also be seen that the false positive rate of the issues reported by LiResolver-Flat is also higher than that reported by LiResolver. By observing the results, we found that it was due to the projects without project licenses. In

Table 3: License Incompatibility Resolved by LiResolver

#Incompatible projects	230	100.00%
– #Resolved projects	144	62.61%
– #Projects that cannot be resolved	86	37.39%
#Incompatible issues	33,340	100.00%
– #Issues resolved by recommending official licenses	11,509	34.52%
– #Issues resolved by generating custom licenses	18,185	54.54%
– #Custom licenses generated w/o exceptions	13,858	41.56%
– #Custom licenses generated w/ exceptions	4,327	12.98%
– #Issues that cannot be resolved	3,646	10.94%

this case, LiResolver-Flat considered all rights reserved, and created the most strict project license for incompatibility detection, which resulted in the false positives.

Answer to RQ2: LiResolver can effectively identify OSS with license incompatibility issues, achieving 0.43% FP rate and no false negatives for 230 real-world incompatible OSS, outperforming four state-of-the-art approaches.

3.3 Evaluation on Incompatibility Resolution

3.3.1 Setup. The experiment on license incompatibility resolution was conducted based on the same dataset with Section 3.2, which consist of 230 incompatible projects detected and verified at the previous step. For each project, there exists at least one pair of licenses which are incompatible with each other. Note that a project was regarded as resolved if all incompatibility issues could be resolved by LiResolver, and no new incompatibility issues were introduced into the resolved project according to LiResolver. In addition, to verify the correctness of resolutions provided by LiResolver, we also randomly selected 5% of the licenses recommended/generated by LiResolver (i.e., 1,484 licenses) for evaluation.

3.3.2 Results. Table 3 shows the experimental results of LiResolver on license incompatibility resolution. Among the 230 incompatible projects, LiResolver resolved 144 projects, accounting for 62.61% of all incompatible projects. From Table 3 it can also be seen that 11,509 out of 33,340 incompatibility issues were resolved by LiResolver through recommending official licenses, accounting for 34.52% of all incompatibility issues. Since a majority of the incompatibility issues cannot be resolved via official licenses, as an alternative way, LiResolver also resolved issues by generating custom licenses to satisfy the constraints. In total, LiResolver generated 18,185 custom licenses, among which only 4,327 custom licenses were attached with specific exceptions to address with the problem when multiple child licenses of a target license are conflict with each other, and none of them can be modified by the software owner to resolve the conflict.

From Table 3 it can also be observed that there are also some license incompatibility issues that cannot be resolved by LiResolver. The reasons are two folds. First, as mentioned in Section 2.5.1, not all constraints from the context of a given license can be resolved; there exists some cases where the parent license of the current license is less restrictive than a child license of the current license. Second, not all licenses are allowed to be modified or replaced by the software owner. Some licenses were incorporated along with third-party software packages, and the project owner has no rights to change them. In total, 3,646 incompatible issues could

Table 4: Effects of Copyright Holder Detection

Modified	#Incomp. Projects		#Incomp. Issues	
	#Resolved Pro.	Suc. Rate	#Resolved Issues	Suc. Rate
PL	138	60.00%	26,523	79.55%
CH	76	33.04%	23,859	71.56%
PL+CH	144	62.61%	29,694	89.06%

PL: Only the project license can be changed to resolve license incompatibility.
CH: Only licenses whose copyright holder is the project owner can be changed.

not be resolved by LiResolver, accounting for 10.94% of the license incompatibility issues detected by LiResolver at the previous step.

Finally, to verify the correctness of resolutions provided by LiResolver, we randomly selected 5% of the 29,694 licenses recommended or generated by LiResolver. We manually verified and cross-validated the correctness of these licenses. By analyzing the results, we found that 1.48% of verified resolutions were incorrect, mainly caused by the incorrect predictions in the phases of fine-grained entity extraction and relation extraction, which affected the license constraints for recommending official licenses and the detailed information used for generating custom licenses.

• **Effects of copyright holder.** Table 4 displays the results of license incompatibility resolution when allowing different ranges of licenses to be modified. In this paper, we classify licenses in a project into the project license and other licenses (also known as component licenses). Unlike previous studies that only focus on the project license, this paper observed a considerable number of component licenses which were also written by the project owner and thus could be modified for resolving license incompatibility. To investigate how the ranges of the licenses allowed to be modified affects the results of license incompatibility resolution, we implemented LiResolver with different settings. The experiment was conducted on the same dataset with Table 3.

From Table 4 it can be seen that 138 projects and 26,523 incompatibility issues could be handled with only the project licenses allowed to be modified, accounting for 60.00% and 79.55% of the incompatible projects and issues, respectively. However, when only allowing the licenses with the same copyright holder with the project license to be changed, LiResolver only resolved 76 projects and 23,859 issues. By analyzing the results, we found that many project licenses did not explicitly state their copyright holders in their license texts, resulting in the low success rate of resolving incompatibility issues. In this paper, we propose to modify two types of licenses (i.e., the project licenses and the component licenses with the same copyright holder). The first intuition is that the project owner is ought to have the rights to modify or replace the project license, whose licensing scope is the whole project. The second intuition is that only a component license states the same copyright holder with as the project license states, the project owner should be allowed to modify it. From Table 4 it can be seen that LiResolver resolved 2.61% more projects and 9.51% more incompatibility issues by adding copyright holder detection.

Answer to RQ3: By recommending official licenses in priority and generating custom licenses as alternative solutions, LiResolver resolved 144 incompatible projects, among which 34.52% were resolved by official licenses, 41.56% and 12.98% were resolved by custom licenses w/o and w/ exceptions, respectively.

Table 5: Feedback from OSS Developers

Project	#Stars	Parent License	Child License	Issue State
cookiecutter	10.2k	BSD-3-clause (/LICENSE)	GPL3.0 (/project_slug/LICENSE)	Ignored
bypy	6.7k	MIT (/LICENSE)	GPL3.0 (bypy/gui.py)	Fixed
			GPL3.0 (bypy/_init_.py)	Fixed
hypnotix	798	GPL3.0 (/README.md)	AGPL3.0 (usr/lib/hypnotix/mpv.py)	Denied
mitmproxy	30.1k	MIT (/LICENSE)	LGPL2.1 (././raw_display.py)	Fixed
			Apache2.0 (././tornado/_init_.py)	Fixed
			BSD-3-clause (././click/_init_.py)	Fixed
webssh	3.4k	MIT (/LICENSE)	LGPL2.1 (tests/sshserver.py)	Fixed
OpenNRE	3.9k	MIT (/LICENSE)	Apache2.0 (./word_piece_tokenizer.py)	Confirmed
			Apache2.0 (./bert_tokenizer.py)	Confirmed
			Apache2.0 (./word_tokenizer.py)	Confirmed
			Apache2.0 (./basic_tokenizer.py)	Confirmed
websockets	879	MIT (/LICENSE.txt)	Apache2.0 (ws4redis/websocket.py)	Denied
			Apache2.0 (ws4redis/utf8validator.py)	Denied

4 DISCUSSION

4.1 Developer Feedback

To evaluate the usefulness of LiResolver, we also collected the feedback from OSS developers with regard to the incompatibility detection results and suggestions. Since there exist 33,340 license incompatibility issues in the dataset described in Section 3 and it is impractical to report all of them, we systematically selected and reported 50 issues according to the following rules. (1) We selected active and popular repositories with high stars in GitHub, which had discussions or active issues within two months before we reported the issues. (2) We only selected issues that could be fixed by LiResolver. For instance, it is difficult to resolve issues caused by licenses of two third-party components unless at least one of the components are replaced by others. (3) The selected issues covered most popular licenses such as the MIT License, Zope Public License 2.1, Apache License 2.0, GPL 3.0, and BSD 3-Clause License. The size of selected projects ranged from 100K to 1G. For each issue, we reported the licenses involved in the incompatibility issues, their locations, the inherent reason for license incompatibility, and the suggestions to resolve this issue. So far, 14 reported issues have been responded as shown in Table 5. Although the remaining issues have been confirmed by ourselves, they are still waiting for responses from OSS developers.

As shown in Table 5, 10/14 responses confirmed or fixed the reported issues. One developer ignored the reported issue and 3/14 responses denied that there existed license incompatibility in their projects. Moreover, we also investigate the reasons that prevent some developers from fixing license incompatibility issues and obtain several findings. (1) OSS developers have different awareness and knowledge about OSS licenses. For example, one developer believed that the MIT license and the Apache 2.0 license were two similar licenses and thus should be compatible with each other. Actually, the Apache 2.0 license is more restrictive than the MIT license in several aspects such as *Use Trademarks* and *State Changes*. As a result, when the project license is the MIT license and its child license is the Apache 2.0 license, there exists license incompatibility which might bring legal risks to the users of this project. (2) It is inconvenient for some large-scale projects with many contributors and users to fix license incompatibility issues, since every contributor should be contacted for the agreement to change the licenses which may also affect many users. Therefore, paying attention to the license incompatibility problem especially at an early stage of software development is necessary. (3) There exist different

opinions towards the licenses of imported packages. Some developers believed that there was no need to consider the licenses of imported packages, which was actually against the intentions of some licenses (e.g., the Apache 2.0 license, the GPL 3.0 license, and the MPL 2.0 license).

4.2 Lessons Learned

We summarize several lessons that might be learned from this work. **First**, developers should pay more attention to the licenses when reusing OSS code. One reason is that the incorporated licenses might be more restrictive than the project license, and thus developers have to modify the regulations for the whole project to comply with the incorporated licenses. The second reason is that licenses accompanied with different third-party software packages might be incompatible with each other. In this case, developers are enforced to either replace the incompatible package with an alternative one, which requires substantial efforts due to the complex package dependency, or attach an exception to the project license that states different regulations towards different source code, which might not be encouraged by the community. **Second**, not all license incompatibility issues can be resolved. Among the investigated 230 incompatible projects, 37.39% of them cannot be resolved unless developers replace the involved third-party packages. The main reason is that a project owner has no rights to change the licenses accompanied with the integrated third-party packages. Another reason is that the constraints from the license hierarchy cannot be resolved in some cases. For instance, given a license, when its parent license is more permissive than one of its child licenses, the constraints from the parent and child nodes cannot be resolved. **Third**, official licenses can only resolve a small part of license incompatibility issues. It can be seen that only 34.52% of the incompatibility issues can be addressed by recommending a new official license. Sometimes developers have to create their own licenses (i.e., custom licenses) to satisfy the requirements of all incorporated OSS licenses. Otherwise, they have to migrate the packages whose licenses are involved in incompatibility issues. **Finally**, although it is a common practice for developers of large-scale OSS to place their own licenses in the modules implemented by themselves, developers should be aware of the license hierarchy that represents the licensing scope of each license, where a module license should not be more permissive than all the licenses in this module, and should not be more restrictive than the licenses in the upper-level module that contains this module simultaneously.

4.3 Limitations and Threats to Validity

The limitations are from two aspects. **First**, given an arbitrary project, if it has no project license or its project license does not state its copyright holder, LiResolver cannot find any internal licenses (also known as component licenses) that can be changed according to the copyright holder. **Second**, LiResolver has not taken into consideration irrelevant licenses (e.g., licenses only for testing purposes). Besides, multiple project licenses in the root directory of the project are considered as the supplements of each other, which might also affect the performance of LiResolver on license incompatibility localization and resolution.

Threats that may affect the results of evaluation include the following two aspects. **First**, the quality of manual labelling could be a threat to validity. The performance of LiResolver on license understanding and incompatibility resolution was evaluated on the ground-truth datasets comprising 21,844 license sentences and 300 real-world OSS, the labelling quality of which may threaten the results. To mitigate this problem, three authors and a lawyer manually verified and cross-validated the datasets. **Second**, as previous studies [29, 45], we categorized *actions* into 23 types. However, there could be a few cases where special requirements outside the scope of these actions are claimed by the licensor, which might influence the results of LiResolver.

5 RELATED WORK

5.1 License Detection and Semantic Extraction

Much research has been done to automatically identify official licenses and extract the semantics from license texts.

License detection. Gobeille et al. [16] presented the first study on license detection. They exploited a binary Symbol Alignment Matrix algorithm to identify OSS licenses. Tuunanen et al. [42] designed license templates and identified some well-known licenses based on regular expressions. Ninka [15], a notable tool for automated license identification, was implemented based on sentence-matching. On this basis, Higashi et al. [23] employed a clustering algorithm to further identify licenses that could not be handled by Ninka.

Semantic extraction. To extract license semantics, many studies conducted the ontology study on licenses [2, 3, 11, 18–20]. For example, Alspaugh et al. [2] [3] extracted tuples from licenses to model 10 licenses. Unlike the ontology-based approaches that required much prior knowledge, Kapitsaki et al. [28] proposed FOSS-LTE to identify license terms with a topic model. Despite the progress, the topic model might induce much noise and cause low accuracy [45].

5.2 License Incompatibility Detection

A major of studies on license incompatibility detection are the graph-based approaches [26, 27, 37, 44]. Generally, these studies manually constructed a directed graph to represent the compatibility relationships between licenses. Then, they detected license incompatibility issues by examining the graph to determine whether two licenses can reach the same node. Despite being strict, there are only a small number of licenses are supported (e.g., 20 licenses and their versions supported by SPDX-VT [27]). It is difficult to manually analyze all the compatibility relations between licenses. Unlike these approaches, LiResolver provides a flexible solution that can be applied for arbitrary licenses without prior knowledge.

LiDetector [45] was the first work that proposed a machine learning-based method for license incompatibility detection. The main differences between LiResolver and LiDetector range from license understanding to incompatibility detection. **(1)** Different granularities of entities. LiDetector regarded each right/obligation as a single entity, while LiResolver splits it into four types of entities, followed by a relation extraction model to organize them. By this means, LiResolver models a regulation more accurately, which as a result benefits license incompatibility detection and resolution. **(2)** Different models. LiResolver first embeds words with a roberta-base model, and then feeds them into a transition-based parser.

In addition, LiDetector regarded each entity independently, while LiResolver employs a prompt-based model to extract the relations between entities. **(3)** Different licensing scopes. LiDetector treated all licenses except the project license equally and independently, while LiResolver observes that the location of a license determines its scope of licensing, and thus extracts license hierarchy for license incompatibility detection and resolution. **(4)** Different detection strategies. Based on license hierarchy extraction, LiResolver obtains *parent-child* pairs of licenses from the license hierarchy, and defines that a license is compatible with its child nodes if anyone who complies with the license will not violate any of its child nodes. On this basis, LiResolver detects license incompatibility from bottom to up, layer by layer according to license hierarchy. Moreover, the detection strategy of LiResolver is more fine-grained than LiDetector, which benefits from its fine-grained license understanding.

5.3 License Recommendation

Based on the aforementioned graph-based methods for license incompatibility detection, Kapitsaki et al. [25] proposed FindOSSLicense to recommend OSS licenses. It considered user requirements through a set of answers to questions, the licenses used by similar users or similar projects, and also license compatibility information according to the directed graph proposed in [27]. To this end, a hybrid approach was proposed to combine the content-based, constraint-based, and collaborative filtering techniques in the recommendation system. However, the license graph only covered a small number of licenses. Licenses not covered in the license graph were marked with “caution”. Liu et al. [33] observed that software changes might lead to license updates. Based on this observation, they proposed to predict source code file-level licenses in the presence of software changes. However, the license prediction tool is only applicable for software changes, which limits its application scope. Compared with these approaches, LiResolver automatically models license texts in a fine-grained way, analyzes the inherent reasons of license incompatibility for arbitrary licenses, and provides flexible solutions for resolving license incompatibility issues.

6 CONCLUSION

In this paper, we propose LiResolver, an automated tool to resolve license incompatibility issues for open source software. Given an OSS, it first extracts all licenses along with their licensing scopes and dependencies. Then, it conducts a fine-grained understanding of license texts, based on which it detects and localizes license incompatibility issues. Finally, for each incompatible license that can be modified by the software owner, it computes the constraints from the context of the license, and provides useful and flexible suggestions. Comprehensive experiments on 300 real-world OSS demonstrate the effectiveness of LiResolver.

ACKNOWLEDGEMENTS

This work was supported by the National Key Project of China (No. 2020YFB1005700), the National Natural Science Foundation of China (No. 62202245, 62102197, and 62002178), and the National Science Foundation of Tianjin (No. 22JCYJC01010).

REFERENCES

- [1] Daniel A Almeida, Gail C Murphy, Greg Wilson, and Mike Hoye. 2017. Do software developers understand open source licenses?. In *2017 IEEE/ACM 25th International Conference on Program Comprehension (ICPC)*. IEEE, 1–11.
- [2] Thomas A Alspaugh, Hazeline U Asuncion, and Walt Scacchi. 2009. Intellectual property rights requirements for heterogeneously-licensed systems. In *Proceedings of the 17th IEEE International Requirements Engineering Conference*. 24–33.
- [3] Thomas A Alspaugh, Walt Scacchi, and Hazeline U Asuncion. 2010. Software licenses in context: The challenge of heterogeneously-licensed systems. *Journal of the Association for Information Systems* 11, 11 (2010), 2.
- [4] Benjamin Andow, Samin Yaseer Mahmud, Wenyu Wang, Justin Whitaker, William Enck, Bradley Reaves, Kapil Singh, and Tao Xie. 2019. Policylint: investigating internal privacy policy contradictions on Google Play. In *Proceedings of the 28th USENIX Conference on Security Symposium*. 585–602.
- [5] Apache. 2004. *LLVMException*. <https://foundation.llvm.org/licensing/LICENSE.txt>
- [6] Anonymous Authors. 2022. *LiResolver*. <https://github.com/anonymous123rainy/LiResolver>
- [7] Xiang Chen, Ningyu Zhang, Xin Xie, Shumin Deng, Yunzhi Yao, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2022. KnowPrompt: Knowledge-Aware Prompt-Tuning with Synergistic Optimization for Relation Extraction. In *Proceedings of the ACM Web Conference 2022 (Virtual Event, Lyon, France) (WWW '22)*. Association for Computing Machinery, New York, NY, USA, 2778–2788. <https://doi.org/10.1145/3485447.3511998>
- [8] choosealicense. 2012. *Choose an open source license*. <https://choosealicense.com/no-permission/>
- [9] christabold. 2016. *Flask JSONDash*. https://github.com/christabor/flask_jsondashh
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186.
- [11] Gordon Thomas F. 2010. Report on prototype decision support system for oss license compatibility issues. *Qualipso* (2010), 80.
- [12] Linux Foundation and its Contributors. 2015. A common software package data exchange format, version 2.0. <https://spdx.org/sites/spdx/files/SPDX-2.0.pdf>.
- [13] Linux Foundation and its Contributors. 2023. *spdx-github*. <https://github.com/spdx/license-list-XML>.
- [14] GR Gangadharan, Vincenzo D'Andrea, Stefano De Paoli, and Michael Weiss. 2012. Managing license compliance in free and open source software development. *Information Systems Frontiers* (2012), 143–154.
- [15] Daniel M. German, Yuki Manabe, and Katsuro Inoue. 2010. A sentence-matching method for automatic license identification of source code files. In *Proceedings of the IEEE/ACM International Conference on Automated Software Engineering*. 437–446.
- [16] Robert Gobeille. 2008. The fossology project. In *Proceedings of the 2008 International Working Conference on Mining Software Repositories*. 47–50.
- [17] Yoav Goldberg and Omer Levy. 2014. word2vec Explained: deriving Mikolov et al.'s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722* (2014).
- [18] Thomas F. Gordon. 2011. Analyzing open Source license compatibility issues with Carneades. In *Proceedings of the 13th International Conference on Artificial Intelligence and Law*. 51–55.
- [19] Thomas F. Gordon. 2013. Introducing the Carneades web application. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Law (Rome, Italy)*. 243–244.
- [20] Thomas F. Gordon. 2014. A demonstration of the MARKOS license analyser. In *Proceedings of the 5th International Conference on Computational Models of Argument*. 461–462.
- [21] Stanford NLP Group. 2020. *corenlp*. <https://stanfordnlp.github.io/CoreNLP/>.
- [22] Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. 1998. Support vector machines. *IEEE Intelligent Systems and their applications* 13, 4 (1998), 18–28.
- [23] Yunosuke Higashi, Yuki Manabe, and Masao Ohira. 2016. Clustering OSS license statements toward automatic generation of license rules. In *Proceedings of the 7th International Workshop on Empirical Software Engineering in Practice*. 30–35.
- [24] Hamed Jelodari, Yongli Wang, Chi Yuan, Xia Feng, Xiahui Jiang, Yanchao Li, and Liang Zhao. 2019. Latent dirichlet allocation (LDA) and topic modeling: models, applications, a survey. *Multimed Tools Appl* (2019), 15169–15211.
- [25] Georgia M Kapitsaki and Georgia Charalambous. 2021. Modeling and recommending open source licenses with findOSSLicense. *IEEE Transactions on Software Engineering* 47, 5 (2021), 919–935.
- [26] Georgia M. Kapitsaki and Frederik Kramer. 2014. Open source license violation check for SPDX files. In *Software Reuse for Dynamic Systems in the Cloud and Beyond*. 90–105.
- [27] Georgia M. Kapitsaki, Frederik Kramer, and Nikolaos D. Tselikas. 2017. Automating the license compatibility process in open source software with SPDX. *Journal of Systems and Software* (2017), 386 – 401.
- [28] Georgia M. Kapitsaki and Demetris Paschalides. 2017. Identifying terms in open source software license texts. In *Proceedings of the 24th Asia-Pacific Software Engineering Conference*. 540–545.
- [29] kevin. 2012. *Software Licenses in Plain English*. <https://tdrlegal.com/librariesio>.
- [30] librariesio. 2015. *Check compatibility between different SPDX licenses for checking dependency license compatibility*. <https://github.com/librariesio/license-compatibility>
- [31] Open Source Licensing. 2004. *Software freedom and intellectual property law*.
- [32] Chengwei Liu, Sen Chen, Lingling Fan, Bihuan Chen, Yang Liu, and Xin Peng. 2022. Demystifying the vulnerability propagation and its evolution via dependency trees in the npm ecosystem. In *Proceedings of the 44th International Conference on Software Engineering*. 672–684.
- [33] Xiaoyu Liu, LiGuo Huang, Jidong Ge, and Vincent Ng. 2019. Predicting licenses for changed source code. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 686–697.
- [34] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *ArXiv abs/1907.11692* (2019).
- [35] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [36] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems* 26 (2013).
- [37] Demetris Paschalides and Georgia M Kapitsaki. 2016. Validate your SPDX files for open source license violations. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*. 1047–1051.
- [38] Ryan Paul. 2009. Cisco settles FSF GPL lawsuit, appoints compliance officer. <http://arstechnica.com/information-technology/2009/05/cisco-settles-fsf-gpl-lawsuit-appoints-compliance-officer>.
- [39] ProgrammerSought. 2021. *The first case of GPL agreement in China is settled. How should the relevant open source software be controlled?*
- [40] Jaideep Reddy. 2015. The Consequences of Violating Open Source Licenses. <https://btlj.org/2015/11/consequences-violating-open-source-licenses/>.
- [41] spacy. 2022. *spacy*. <https://spacy.io/usage/training>.
- [42] Tuunanen Timo, Koskinen Jussi, and Kärkkäinen Tommi. 2009. Automated software license analysis. *Automated Software Engineering* 16 (2009), 455–490.
- [43] Steven Vaughan. 2015. VMware sued for failure to comply with Linux license. <http://www.zdnet.com/article/>.
- [44] David A. Wheeler. 2007. *The free-libre / open source software (FLOSS) license slide*. Retrieved January 26, 2017 from <http://www.dwheeler.com/essays/floss-license-slide.pdf>
- [45] Sihan Xu, Ya Gao, Lingling Fan, Zheli Liu, Yang Liu, and Hua Ji. 2021. LiDetector: License Incompatibility Detection for Open Source Software. *ACM Transactions on Software Engineering and Methodology* (2021).
- [46] Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu, Qingyun Wang, Heng Ji, and Meng Jiang. 2022. A survey of knowledge-enhanced text generation. *ACM Computing Surveys (CSUR)* (2022).