

Detecting Insecure Implementation of Vulnerable Code Snippets found on Stackoverflow

Mazharul Islam

ABSTRACT

Online programming discussion platforms such as Stack Overflow have a rich source of ready to use code snippets for software developers. It is the de-facto place where developers go to find solutions from the coding snippets given as answers to posted problems by the online developer community. However, previous research work has shown that developers have a tendency to directly copy-paste insecure code snippets from Stack Overflow into their production level code. As a result without any countermeasures Stack Overflow is becoming one of major sources of vulnerability of production level code. To address this problem, in this project, we tackle the problem of analyzing code snippets found on Stack overflow. This is challenging since code snippets from Stack overflow are often erroneous, incomplete, and lack dependencies which makes it harder to analyze them by state-of-the-art static tools. Our goal is to build a static analysis tool that can identify common insecure patterns found on a dataset containing a collection of 1.6K code snippets from Stackoverflow.

1 INTRODUCTION

Stack Overflow (SO) is regarded as one of the most popular online helping platforms for software developers [1]. SO seeks to help by creating an eco-system of online developer community. In this eco-system one developer can ask for solutions to the problems one is facing, while other developers can interact by posting code snippets, advices, as solutions to those asked problems. As such, SO has become a rich source of ready-to-use code snippets for software developers. This richness has caused SO to enter the agile software development cycle allowing fast prototyping, and an efficient workflow. Particularly novice developers treasure the quick-direct help from the community providing easy ready-to-use code snippets. Interestingly copying-pasting code snippets into production software is generally practiced not only by the inexperienced developers, but a common practice shared by large parts of the developer community. Furthermore, sometimes experienced developers can potentially promote best-practices by distributing high quality code snippets, and eventually improving code quality on a large basis.

Unfortunately this is not same for secure coding practices. Fischer et al. [5] quantitatively measured that an Android developer seeking help on SO, can find accepted popular code snippets suggesting insecure X.509 certificate validation, misusing Android's cryptographic API as shown in Figure 1. Moreover a developer struggling with handling Java Spring security framework's configuration errors, can find code snippets suggesting turning off CSRF security protection entirely to avoid errors. Insecure code snippets found on Stack Overflow itself is not a serious problem. However these insecure code snippets have naturally entered the development cycle by developers who are copy-pasting them into the production level code – causing vulnerabilities in softwares.

Here is some relevant code:

```
106 // Create a trust manager that does not validate certificate chains
TrustManager[] trustAllCerts = new TrustManager[]{
    new X509TrustManager() {
        public java.security.cert.X509Certificate[] getAcceptedIssuers() {
            return null;
        }
    }
};

public void checkClientTrusted(
    java.security.cert.X509Certificate[] certs, String authType) {
}

public void checkServerTrusted(
    java.security.cert.X509Certificate[] certs, String authType) {
}
```

Figure 1

Moreover to add salt to the problem, these insecure code snippets are being accepted, promoted even by users with high reputation [7] – whereas condescending comments are being directed at security conscious users [8]. This gives the developers a high level of false confidence while copy-pasting insecure code-snippets.

To solve this problem, we need a way to flag the insecure patterns present on code snippets. This will help developers to be more cautious before copy pasting insecure code snippets. This type of flagging can also draw more attention to the developers before upvoting or advising for an insecure code snippet. Eventually this can promote writing secure code snippets in Stack Overflow eco-system. However there is no existing tool used by Stack Overflow to analyze and flag a code snippet when it contains potential insecure patterns.

Towards solving this problem, in this paper, we aim to develop a static analysis tool which can identify which part of the code snippet is insecure and suggest secure alternative suggestions to developers. However analyzing code snippets for insecure patterns presents some unique challenges. Since code snippets are erroneous, incomplete, we have apply repair techniques to be able to convert them to an intermediate representation (IR) before running any static analysis. We also need to identify which insecure patterns are common, and how to identify them. Therefore, we first study the literature to compile a list of 8 most insecure patterns related to code snippets in Java language. We then apply simple parsing repairs to the code snippets before converting them to an intermediate representation (IR) language called Jimple. Finally we use a combination of keyword and backward flow analysis based detection method to identify the insecure patterns in code snippets.

In summary, we make the following the contributions in the paper:

- We survey the literature, and present 8 most insecure patterns appearing recursively in Java language code snippets posted on SO (section 2).
- We develop code repair techniques from simple parsing fixes, and apply them on code snippets from SO to convert them to Jimple intermediate representation (IR) (section 3.1).

- We run keyword and backward flow analysis based detection method to identify the presence of 8 insecure patterns (section 3.3).
- We manually verify our code repair and detection accuracy on a collection of 1.6K code snippets from SO (section 4).

The rest of the paper is organized as the following: section 2 describes 8 most common insecure patterns found on code snippets, and explains why they are insecure, section 3 describes our methodology and section 4 presents experimental evaluation. We conclude the paper by presenting some compelling discussion (section 6), limitations (section 5), and future work (section 8).

2 THREAT MODEL

Rule No	Description	Vulnerability
1	AES default encryption mode ECB	Side channel attack
2	Insecure cryptographic hash	Collision attack
3	Abuse of X509TrustManager Verifier Interface	SSL/TLS MitM attack
4	Weak key length	Brute force attack
5	Static/constant/predictable keys/IV	
7	Presence of AllHostNameVerifier	SSL/TLS MiTM
8	Turning of CSRF protection	CSRF attack

Table 1

We will now summarize the insecure patterns our method aims to detect in the following paragraphs and in Table 1. For each insecure patterns, we will also describe the security risks it presents, and its secure usage from the literature. This will give us a sense what we want our method to detect i.e, the presence of insecure patterns or the absence of secure usage.

2.1 AES default encryption mode ECB

AES is one widely adopted and used encryption standards in the developer community. Therefore it is no surprise that a large number of code snippets uses AES for encryption [FIXME: Add some % number]. In Java an instance of AES can be created using `javax.crypto.Cipher`. However `javax.crypto.Cipher` class uses Electronic Codebook (ECB) as the default mode of operation when "AES" is passed as transformation parameter to `getInstance` method [FIXME: See the code in the Appendix A]. While ECB-encrypted ciphertext allows random access to each block, it can also leak information via side channel attacks [FIXME: Cite source]. However developers being unaware of this default behavior insecure behaviour of AES, share code snippets without any considerations that uses insecure ECB mode for encryption. Instead developers should be using Block Chaining (CBC) or Galois/Counter Mode (GCM) which are not vulnerable to side channel attacks as shown in appendix.

2.2 Abuse of X509TrustManager Verifier Interface

X509TrustManager Verifier interface is popular among developers to instantiate TrustManager class. Ideally a secure implementation of X509TrustManager should i) throw exception after validating a certificate in `checkServerTrusted` method, ii) provide a valid list of certificates in `getAcceptedIssuers` method, and iii) throw

exceptions for self signed certificates. However while writing code snippets developers tend to leave empty methods to implement the X509TrustManager interface. As a result the X509TrustManager Interface accepts any certificate including the ones which are not signed by a trusted certificate authority. This enables a provision for Man-in-the-middle (MitM) attacks.

2.3 Insecure cryptographic hash

A cryptographic hash function produces fixed-length unique alphanumeric string called message digest for any arbitrary message. This unique message digest can be used latter for verifying crypto properties of the message e.g., message integrity, digital signature, and authentication. However if two different messages produces the same message digest i.e., a collusion happens, then attacker can compromise these crypto properties. A cryptographic is broken if attacker has systemic practical way to produces collusion for different message. The list of popular but broken hash functions includes SHA1, MD4, MD5, and MD2. These hash functions produce collisions that cause cryptographic vulnerabilities, and hence should be avoided. However in code snippets developers have been using these popular broken hashes as shown in listing ??.

2.4 Absence of performing hostname verification

Ideally to perform a hostname verification, developer has to implement the `javax.net.ssl.HostnameVerifier` by using `java.net.ssl.SSLSession` parameter inside the `verify` method. However in many cases this `verify` method is always set to return true as shown in listing ??.

The reason being while writing code snippets for brevity this dummy return true will not throw any exceptions. However this type of workaround can cause security threats such as URL spoofing attacks. URL spoofing makes it simpler for numerous cyber-attacks (e.g., identity theft, phishing).

2.5 Weak key length

The strength of asymmetric encryption (e.g., RSA, ECC) depends on using sufficiently large key length. Since 2015, NIST recommends a minimum of 2048-bit keys for RSA, [14] an update to the widely-accepted recommendation of a 1024-bit minimum since at least 2002. This ensures that the key space is large enough to prevent any practical brute force attack. However while writing code snippets developers have been using key length of less than 2048 disregarding this recommendation.

2.6 Static/constant/predictable keys/IV

Predictable keys/ Initialization Vectors (IV) are a major source insecurity in the code snippets. Raw keys and raw IVs created from empty byte arrays are easily guessable by attackers. Additionally some code snippets derive keys directly from simple and insecure passphrases as shown in listing ??.

Static constant keys are susceptible to leaks. As oftentimes attackers can decompile the application and get the static hardcoded keys. To avoid this kind of attacks, developers should avoid using static constant keys. `javax.crypto.spec.SecretKeySpec` and `javax.crypto.spec.PBEKeySpec` are two popular ways to generate secret keys used for encryption. Both of these API takes a

byte array to generate the secret keys. However if the byte array is constant or hardcoded inside the code, the adversary can easily read the cryptographic key and may obtain sensitive information. This is the same case for storing keys in a keystore using `java.security.KeyStore` API. The secret keys by which is key stored is locked for safely storing the keys should take a byte array is not static.

2.7 Presence of AllHostNameVerifier

`org.apache.http.conn.ssl.SSLConnectionSocketFactory` provides a static field `acceptAllCertificates`. This is same as using empty methods as discussed in ???. This time developers can just use `ALLOW_ALL_HOSTNAME_VERIFIER` static field to do this. As this is a very easy way to avoid errors, in code snippets developers insensibly uses them frequently without considering the insecurity associated with using it.

2.8 Turning of CSRF protection

Cross site request forgery (CSRF) is a serious attack that tricks the a web browser by abusing the browser cookie authentication mechanism to execute privilege unwanted actions. To protect against such attacks ideally CSRF-Token should be included in all POST, PUT, DELETE requests. However the from code snippets related to Java Spring security framework, we have found that the developers tend to turn of the CSRF protection forcefully to avoid getting errors.

3 METHODOLOGY

In this section, we will discuss the pipeline we follow to detect insecure patterns in code snippets as shown in Figure 2. We will first discuss about repairing the code snippets (section 3.1), and then converting the repaired to code snippets to an Intermediate representation (IR) to run analysis (section 3.2). Lastly we will finish by describing the techniques we have applied on the converted IR to detect the insecure patterns (section 3.3), which we have previously described in section 2.

3.1 Code Repair

While writing code snippets as answers to posted questions, developers tend to be concise and short. The reason being long code snippets has lower chance of being accepted and upvoted by others in online platforms such as Stack Overflow. [FIXME: Give a statistic on the avg. length of the code snippets of the dataset] Within a few lines of code, developers try to convey the intent hinting at a working solution by assuming everything other are in place to for successful compilation. However this very mindset of developers can leave syntatic errors, missing classes in the code snippets. As a result, converting these code snippets to IR for analysis becomes difficult.

For identifying insecure patterns for which only keyword searching is sufficient (e.g., Rule 7, 8 as shown in Table 2) this is not a problem since we don't need to convert them to any IR. However for identifying insecure pattern (e.g., Rule 1-6 as shown in Table 3) which requires running analysis this poses a problem. Therefore to identify them, we need to add some repairs to the code snippets. For the purpose of this paper, we have applied the following repairs to the code snippets.

3.1.1 Syntatic repair. To remove the syntatic error present on the code snippets we do the following syntatic repairing.

- We remove illegal characters (e.g., ">", "<", "&", """, "'", etc). Many of these illegal characters appeared as the dataset was crawled from Stackoverflow website's raw HTML and HTML sanitizes some characters which are used in the code snippets. Also some code snippets have comments without any comment sign, and dots to imply some code would be here which not relevant to question posted.
- Some code snippets do not have match brackets, and extra quotes for strings.
- Some code snippets have `@Override` notation implying it is implementing an interface. However the partial program analysis tool we will discuss to convert code snippets to IR, can not handle `@Override` notation.

3.1.2 Missing package, class, and method names: Partial Program Analysis (PPA) tool which we have used to convert the code snippets to IR, can not consume a lines of code missing classname, package name, method names. Therefore we applied the following repairs.

- If the code snippet is missing any class name, we wrap the code snippet inside a public class name. If the code snippet already has a public class, we rename the file according to that public class.
- If the code snippet does not have any package name, we place the public class in a package, and add the package name to the code snippet. If the code snippet has package name, we create proper directory structure according to the package name and place the code snippet there before converting them to IR.
- We also add dummy implementation of missing methods as developers tend to have methods names in code snippets but does not give any implementation within the code snippets.
- Finally, we load the some popular crypto classes in Java to the runtime of PPA tool which are imported, implemented by code snippets frequently. This helped us to avoid missing class name, unknown interface error thrown by PPA in many cases.

3.2 Converting code snippets to IR

After repairing the code snippets, we tried to convert them to IR grammar named Jimple. Jimple [15] is a 3-address intermediate representation that has been designed to simplify analysis. Jimple was inspired from SIMPLE an AST to represent C statements. To convert the code snippets we used a tool named partial program analysis (PPA). Dagenais et al. developed PPA [2] with goal of analyzing only subset of a program source code which matches with our use case of analysing code snippets. PPA can infer types where types are not present that subset of the code. In case of failure PPA will place special type "MAGICCLASS", "MAGICCLASS", and "MAGICMETHOD". This is necessary since without types it is not possible to build the abstract syntax tree, and eventually convert the code snippet to Jimple for a strongly typed language such as Java. As PPA can overcome this problem by inferring types of the objects used in the subset of the program source code, it can convert the subset program source code. We leverage PPA after making

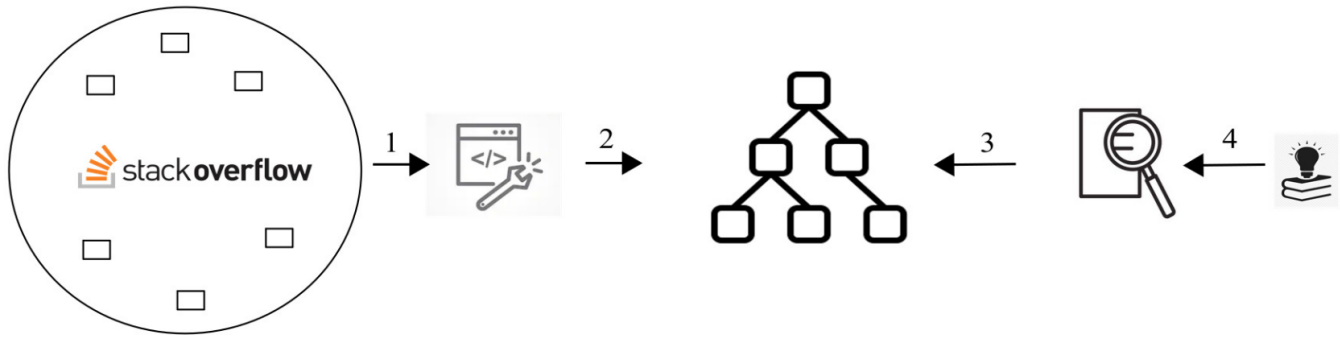


Figure 2

the code repairs presented in previous subsection 3.1. Otherwise a large number of code snippets was throwing errors as PPA can not handle erroneous code snippets.

The idea is to feed the Jimple representation of the code snippet to Soot – a state-of-the-art program analysis tool [13]. Soot API can consume a Jimple representation, and perform data flow analysis which is as discussed in the next subsection 3.3, required for detecting insecure patterns.

3.3 Identifying insecure patterns

In this subsection we will discuss the two techniques we have used to identify the insecure patterns discussed in section 2. Specifically we have used two techniques found in the existing literature. One of them is keyword based detection, and the former one is using back flow sensitive analysis.

Rule No	keywords
7	SSLSocketFactory.ALLOW_ALL_HOSTNAME_VERIFIER
8	*.csrf.disable()

Table 2

Rule No	Slicing criteria
1	KeyGenerator.getInstance(*)
2	MessageDigest.getInstance(*);
3	public void checkClientTrusted public void checkServerTrusted public X509Certificate[] getAcceptedIssuers()
4	keyPairGenerator.initialize(keySize);
5	public boolean verify
6	new SecretKeySpec(keyBytes, "AES") *.load(*.openStream(), new String(keyBytes).toCharArray()); new PBEKeySpec(new String(keyBytes).toCharArray(),*);

Table 3

3.3.1 Key word base analysis. In keyword based analysis, we want to write a regex which can capture the common way developers write the insecure patterns, and then searching in the code snippets

for matching the written code snippets. This method was used by Rahman et al. [11] to detect insecure practices present in Python code snippets on Stackoverflow. Although being a simple technique, it worked surprising well for them i.e, does not introduce any false positives. However, in our case, as we will show in the next, that capturing all the insecure pattern by writing regex can introduce false-positives even simple code snippets.

Therefore, according to our manual observation, we can detect only two insecure patterns using keyword searching. This follows from the reasoning that Rule 7 and 8 can be written by any developer, in the exact pattern as shown in Table 2. For detecting the other 6 insecure patterns, we have to restore to backword flow program analysis as described next.

3.3.2 Backword flow base analysis. We use the backword flow analysis introduced by Rahaman et al. [10] in their CryptoGurd Project. They introduced specialize def-use analysis [16] based on program slicing techniques [3] to detect 16 common cryptographic API misuses in Apache, and Android projects. Def-use dataflow analysis builds a dependency relation based on the definition and use statements. Given a slicing criteria, which is a statement, or a parameters of an API, backword flow analysis computes the set of program statements that affects the slicing criteria in terms of data flow. The key design choice, hence, here is to specify special function invocation places as the slicing criteria. The slicing criteria used for paper are highlighted in Table 3.

Now we will detail why simple keyword based analysis is not sufficient for detecting rules 1-6 as they can introduce FP even for simple rules. To demonstrate this, consider the example code snippet shown in listing 1 corresponding to insecure pattern rule 2 – detecting insecure broken cryptographic hashes MD5, MD4, MD2, SHA1. We can use keyword searching based technique on the name of broken hashes, and successfully detect that the insecure pattern that code snippet in listing 1 have used broken hash. However it will introduce False positive for same insecure pattern on the code snippets shown in listing 1. As there are multiple ways developers can use these broken hashes, unlike rule 7-8, we set `MessageDigest.getInstance(*)`; as the slicing criteria. Then we start backword def-use analysis to see if any of the program sets affects the parameters of `MessageDigest.getInstance(*)`; and has a value of equal to name of the broken hash.


```

2  MessageDigest md = MessageDigest.getInstance("MD5");
3  md.update(str.getBytes());
4  ....

```

Listing 1: A code snippet where keyword based detection work well

```

1  ...
2  int flag = 2;
3  MessageDigest md = MessageDigest.getInstance("MD5");
4  if(choice > 1){
5      md = MessageDigest.getInstance("SHA-256");
6  }
7  md.update(str.getBytes());
8  ....

```

Listing 2: A code snippet where keyword based detection introduces FP

«Mazharul 3.0: A formal proof is given in the Appendix.»

We modified the code base of as CryptoGuard¹ to achieve our analysis on code snippets. This is for two reasons. Firstly, CryptoGuard already has the basic skeleton for use-def analysis, and we just have to change the slicing criteria. Secondly, CryptoGurd uses Soot as the underlying program analysis engine. Hence we can provide our generated Jimple IR using the PPA tool to CryptoGuard, and CryptoGurd's program analysis engine Soot can do backward analysis based on the slicing criteria defined by us.

Figure 3 shows one such example. Taking SecretKeySpec as the slicing criteria, we identify the set of statements which affects the first parameter of SecretKeySpec – which is MyDifficultPassw here. Eventually we will backtrack to the definition program set and can reason that it is randomly generated - rather a hard coded secret. Since this makes the SecretKeySpec class object sks predictable, we have detected the presence of insecure pattern 6.

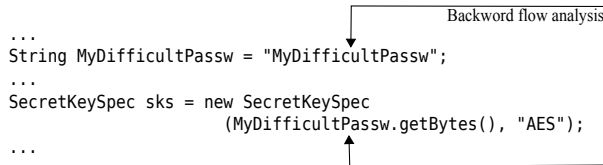


Figure 3

4 RESULTS

In this section, we will first describe our data-collection process. Next we will show the code repair, and detection accuracy techniques on a small subset of this collected dataset.

4.1 Data collection

We used dataset from two previous sources [5, 7]. Both of these dataset contain code snippets posted on Stackoverflow. Fisher et al. crawled 1,161 code snippets posted on Stackoverflow related to Android Security [5]. They considered a code snippet related to Android security if the code snippets makes API calls to one of the security services such as Java cryptography, Java secure Communications, public key infrastructure X.509 certificates, and Java

¹<https://github.com/CryptoGuardOSS/cryptoguard>

authentication - authorization services. The popular crypto libraries used by Andriod developers such as Bouncy Castle, SpongyCastle, Apache TLS/SSL, keyczar, jasypt, and GNU Crypto were also included.

Meng et al. extracted 503 code snippets from 22,195 Stackoverflow posts by filtering the posts based on votes, duplications, and absence of code snipeets [7]. In total our study is baded on the dataset by combining these two. Our dataset contains 1,664 code snippets. The timeline of these code snippets are from 2008-2017.

«Mazharul 4.0: add some more info and some statistics»

As to get the ground truth of the presence of insecure pattern, we have to manully analysis them, it becomes very time consuming. Therefore to make the analysis less consuming time , we only consider randomly sampled 800 code snippets – about half of the 1.6K code snippets available. We then assign the code snippets into one of the 8 insecure patterns.

4.2 Code repair success rate.

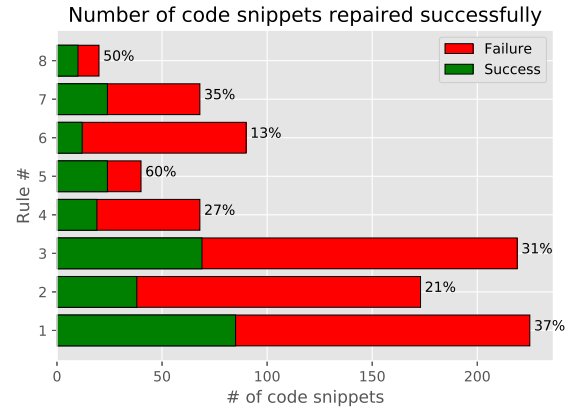


Figure 4: Percentage of code snippets successfully repaired.

We called a code snippet successfully repaired if we can convert it to a Jimple IR. Figure 4 illustrates the percentage of code snippets for each categorize, we have been able to parsed (i.e., convert to Jimple IR), using the code repair techniques discussed in subsection 3.1.

4.3 Insecure pattern detection accuracy.

After converting each successfully repaired code snippets to Jimple IR, we detected rule 7, and 8 using keyword searching as mentioned previously. However for rule 1-6, need to apply backward flow analysis. To do this we give the Jimple IR to our tool. Our tool is built on top of CryptoGurd. CryptoGurd uses Soot as its program analysis engine. Using Soot's JimpleAST API², we can enable backward flow analysis given the slicing criteria from Table 3. The idea is track the special method invokations parameter from the slicing criteria. In this way we can inspect the set of program statements which are affected by this slicing criteria and figure out the presence of insecure patterns.

²<https://www.sable.mcgill.ca/soot/doc/soot/jimple/parser/JimpleAST.html>

Table 4 shows the total number of code snippets for of the 8 rules (they sum up to more than 800 as some code snippets has two or more rules). Parsed column refers to the number of code snippets, we have been able to convert to Jimple IR after applying the code repair discussed in section 3.1. We also shows the TP, FN, FN along side the precision, and recall. As it can seen we have not been able to run backward flow analysis for rule 3. The reasons are explain the in the appendix.

Rule no	Total	Parsed	TP	FP	FN	Precision	Recall
1	225	85					
2	173	38					
3	219	69	-	-	-	-	-
4	68	19					
5	40	24					
6	90	12					
7	68	24					
8	20	10					

Table 4: Table

5 LIMITATIONS

Limitations from conservative datasets. We only considered a small subset of code snippets in Java available on Stackoverflow. As a result the 8 insecure patterns we have considered are tailored to this small subset of code snippets. A more rich dataset would have allowed to capture more naunce insecure patterns. Moreover, we only analzed code snippets with Java/Android tags. If we have analyzed code snippets from other popular languages [9], we would have been able to discover more insecure patterns.

Limitations from program repair techniques. The program repair techniques in section 3.1, we have applied on erroneous code snippets, before converting them to Jimple IR, are simple semi-automated simple parsing repairs. It would have been better to see if state-of-the-art the program repair techniques can be adjusted to repair them. Especially automated program repair tools which can apply quick single edit fixes (e.g. Google's OSS-Fuzz and Microsoft's Springfield project) are quite enticing to apply for fixing erroneous code snippets as discussed in this great survey paper [6].

Limitations due to PPA/Jimple grammer. The partial program analysis tool (PPA) [2] we have used theto run analysis depends on Jimple. Jimple, a 3-address IR, was designed to handle, and simplify several difficulties while performing optimizations on the stack-based Java bytecode directly [15]. However, the simple grammer of Jimple on which we are running backward flow analysis can not handle anonymous class. Anonymous classes are the most common way developers declare the X509TrustManager interface (insecure pattern # 3). Because of Jimple inability to express anonymous classes, we can not detect the insecure pattern #3.

Limitations from disregarding comments. Insecurity in production level code can also come from insecure advices given in text

forms on Stackoverflow. As we are only considering code snippets we are missing out on them. Ideally we would like to use NLP techniques to detect the insecure advices. Also some code snippets do contains insecure patterns but the developers has warned against blindly copy pasting this insecure code snippets as comments in the code snippets. As we are disregarding comments, it is quite debateable that considering the code snippets would make sense at any one before copy pasting this code snippet would read the comment and already know about the insecurity of the code snippet.

6 DISCUSSIONS ON ML BASED DETECTION TECHNIQUES

Identification of insecure patterns, vulnerabilities, and code-smells using ML based techniques have been gaining tractions recently. This is mainly for two reasons i) leveraging the wealth of open source code available ii) adaptive neural network models which can capture/represent the complex patterns of source code. While giving a overview of ongoing research in this area is outside the scope of the paper, we can mention some work existing research work closely related to ours that uses ML techniques. For example, Zhou et al. proposed [17] vulnerability detection model "Devign" for code snippets in C language. Their key idea to detect vulnerability is by capturing the abstract syntax tree structure of the source code via Gated-Graph Neural Network (GGNN). The study in [12] proposed automated vulnerability detection tools using deep feature representation learning.

One problem associated with using machine learning models is that either they are function level detection model [12], do have limited ability to reason why source code is vulnerable [17], or suffers from subjectivity [4].

Static analysis based approaches such as ours, can overcome for two reasons i) it can reason about the source code, and ii) as insecure patterns are repetitive. As a result we can build a static analysis tool by observing the common insecure trends to reason about the source. We agree this can be an oversteared claim, and leave it as a future work this paper.

7 RELATED WORK

Our work is motivated by the following related research work. Subramanian, et al. [14] used Eclipse Java Development Tools (JDT) ³ to find structural models of code snippets in Stack Overflow. Consequently, by analyzing *solved* Stack Overflow questions having *Android* tag they present a common list of Android API types and methods – something which normal lexical parsers are unable to detect. Fischer et al. [5] quantitatively evaluated the observation that a large number of insecure code snippets are being directly copy-pasted, repeatedly reused. They showed that a simple stochastic gradient descent based classifier can confirm that among 1.3 million Google Play Android applications, 15.4% contains security-related code snippets from from Stack Overflow – out of which 97.9% contain at least one insecure code snippet. Meng et al. [7] did an empirical study on the on StackOverflow posts, aiming to understand developers' concerns on Java secure coding. This study highlights a number of popular-accepted insecure suggestions on

³<http://www.eclipse.org/jdt>

StackOverflow including suggestions to disabling the default protection against Cross-Site Request Forgery (CSRF) attacks, breaking SSL/TLS security through bypassing certificate validation, and using insecure cryptographic hash functions. These harmful insecure suggestions can easily misguide developer – the extend of which is still unknown today. Interestingly, Rahman et al. [11] did a study similar to Meng et al. [7], but for code snippets for Python language. They observed that 9.8% of the 7,444 accepted answers to include at least one insecure code block. Most importantly they also find user reputation not translate to the presence of insecure code blocks, implying that both high and low-reputed users are likely to introduce insecure code blocks.

8 FUTURE WORK ON SYNTHESIZING SECURE CODE

9 CONCLUSION

In this paper, we explore the problem of identifying insecure patterns on source code snippets given in response to Stack Overflow (SO) questions. We motivate the problem by emphasizing that insecurity in code snippets present in SO, can potentially trickled down to production level source code – introducing vulnerabilities. We believe that by flagging the common 8 insecure patterns compiled from existing literature can encourage developers to accept, upvote, share secure code snippets. We propose a static analysis tool based on keyword searching and backward flow analysis to do so. Our experimental result on 1.6K code snippets written in Java language demonstrates the promising outcome of our proposed solution.

REFERENCES

- [1] S. Baltes, C. Treude, and S. Diehl. 2019. SOTorrent: Studying the Origin, Evolution, and Usage of Stack Overflow Code Snippets. In *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*. 191–194. <https://doi.org/10.1109/MSR.2019.00038>
- [2] Barthélémy Dagenais and Laurie Hendren. 2008. Enabling static analysis for partial java programs. In *Proceedings of the 23rd ACM SIGPLAN conference on Object-oriented programming systems languages and applications*. 313–328.
- [3] A. De Lucia. 2001. Program slicing: methods and applications. In *Proceedings First IEEE International Workshop on Source Code Analysis and Manipulation*. 142–149. <https://doi.org/10.1109/SCAM.2001.972675>
- [4] D. Di Nucci, F. Palomba, D. A. Tamburri, A. Serebrenik, and A. De Lucia. 2018. Detecting code smells using machine learning techniques: Are we there yet?. In *2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. 612–621. <https://doi.org/10.1109/SANER.2018.8330266>
- [5] Felix Fischer, Konstantin Böttinger, Huang Xiao, Christian Stransky, Yasemin Acar, Michael Backes, and Sascha Fahl. 2017. Stack overflow considered harmful? the impact of copy&paste on android application security. In *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 121–136.
- [6] Claire Le Goues, Michael Pradel, and Abhik Roychoudhury. 2019. Automated Program Repair. *Commun. ACM* 62, 12 (Nov. 2019), 56–65. <https://doi.org/10.1145/3318162>
- [7] Na Meng, Stefan Nagy, Danfeng Yao, Wenjie Zhuang, and Gustavo Arango Argoty. 2018. Secure coding practices in java: Challenges and vulnerabilities. In *Proceedings of the 40th International Conference on Software Engineering*. 372–383.
- [8] Stack Overflow. 2015. java - When I try to convert a string with certificate, Exception is raised. <https://stackoverflow.com/questions/10594000/when-i-try-to-convert-a-string-with-certificate-exception-is-raised>.
- [9] Stack Overflow Developer Survey 2020. 2020. - Most Loved, Dreaded, and Wanted Languages. <https://insights.stackoverflow.com/survey/2020#technology-most-loved-dreaded-and-wanted-languages-loved>.
- [10] Sazzadur Rahaman, Ya Xiao, Sharmin Afrose, Fahad Shaon, Ke Tian, Miles Frantz, Murat Kantarcioglu, and Danfeng (Daphne) Yao. 2019. CryptoGuard: High Precision Detection of Cryptographic Vulnerabilities in Massive-Sized Java Projects. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (London, United Kingdom) (CCS '19). Association for Computing Machinery, New York, NY, USA, 2455–2472. <https://doi.org/10.1145/3319535.3345659>

- [11] A. Rahman, E. Farhana, and N. Imtiaz. 2019. Snakes in Paradise?: Insecure Python-Related Coding Practices in Stack Overflow. In *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*. 200–204. <https://doi.org/10.1109/MSR.2019.00040>
- [12] R. Russell, L. Kim, L. Hamilton, T. Lazovich, J. Harer, O. Ozdemir, P. Ellingwood, and M. McConley. 2018. Automated Vulnerability Detection in Source Code Using Deep Representation Learning. In *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*. 757–762. <https://doi.org/10.1109/ICMLA.2018.00120>
- [13] Soot. 1999. - A framework for analyzing and transforming Java and Android applications. <https://soot-oss.github.io/soot/>.
- [14] Siddharth Subramanian and Reid Holmes. 2013. Making sense of online code snippets. In *2013 10th Working Conference on Mining Software Repositories (MSR)*. IEEE, 85–88.
- [15] Raja Vallee-Rai and Laurie J Hendren. 1998. Jimple: Simplifying Java bytecode for analyses and transformations. (1998).
- [16] H. Y. Yang, E. Tempero, and H. Melton. 2008. An Empirical Study into Use of Dependency Injection in Java. In *19th Australian Conference on Software Engineering (aswec 2008)*. 239–247. <https://doi.org/10.1109/ASWEC.2008.4483212>
- [17] Yaqin Zhou, Shangqing Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective Vulnerability Identification by Learning Comprehensive Program Semantics via Graph Neural Networks. In *Advances in Neural Information Processing Systems (NeurIPS)* 32. Curran Associates, Inc., 10197–10207.

10 APPENDIX

```

1 private byte[] encrypt(byte[] raw, byte[] clear) {
2     ...
3     CIPHER_MODE = "AES"
4     Cipher cipher = Cipher.getInstance(CIPHER_MODE);
5     ....
6     return encrypted
7 }
8

```

Listing 3: A real code snippet taken from Stackoverflow. I want to build a tool which after analyzing the code snippet will highlight the part of the code that is insecure and suggest an alternative secure implementation as showed in the figure.

```

1 private byte[] encrypt(byte[] raw, byte[] clear) {
2     ...
3     Cipher cipher = Cipher.getInstance("AES");
4     ....
5     return encrypted
6 }
7

```

Listing 4: A real code snippet taken from Stackoverflow. I want to build a tool which after analyzing the code snippet will highlight the part of the code that is insecure and suggest an alternative secure implementation as showed in the figure.

```

1 static void decrypt() {
2     ...
3     String MyDifficultPassw = "MyDifficultPassw";
4     ...
5     SecretKeySpec sks = new SecretKeySpec(
6         MyDifficultPassw.getBytes(), "AES");
7     ...
8

```

```
7 }
```

Listing 5: A real code snippet taken from Stackoverflow. I want to build a tool which after analyzing the code snippet will highlight the part of the code that is insecure and suggest an alternative secure implementation as showed in the figure.

```
1 private byte[] encrypt(byte[] raw, byte[] clear) {  
2     ...
```

```
3     MessageDigest md = MessageDigest.getInstance("MD5")  
4     ;  
5     ....  
6     return encrypted  
7 }
```

Listing 6: A real code snippet taken from Stackoverflow. I want to build a tool which after analyzing the code snippet will highlight the part of the code that is insecure and suggest an alternative secure implementation as showed in the figure.