

1
2
3
4
5
6
7
8
9
10
11
12
13
14

```
CSCS 2022 'Conference' {
```

```
[VulDetective]
```

```
##Towards a Block-Level ML-Based Python  
Vulnerability Detection Tool
```

```
##Machine Learning Based Vulnerability Detection
```

```
University of Szeged :
```

```
#Peter Hegedus
```

```
#Amirreza-Bagheri
```

```
}
```

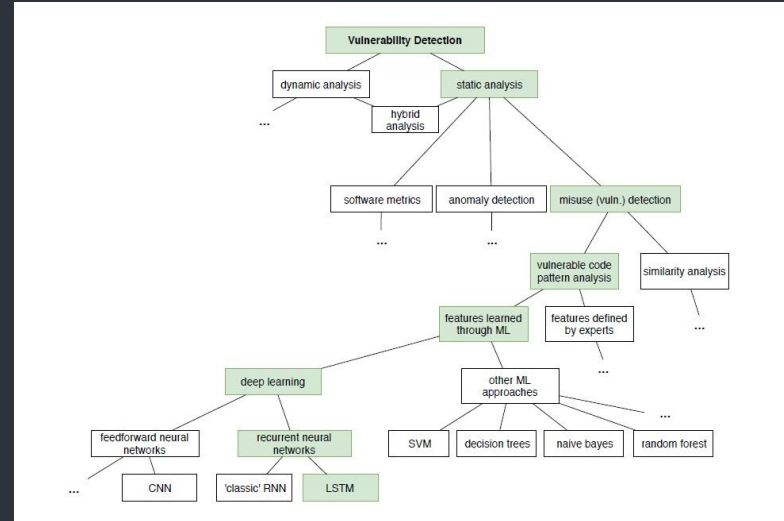
Introduction; {

'Creating safe, dependable, and secure software is no easy undertaking.'

<p Software architects and engineers' oversights and errors can quickly lead to software vulnerabilities, with serious repercussions >

</p>

}



```
1 01 {  
2  
3
```

```
4  
5 [Dataset Collecting]  
6
```

```
7 < We need to use a data set of real  
8 source codes for our training for this  
9 purpose we used Github API with more  
10 than 127 million repositories>  
11
```

```
12 }  
13  
14
```

1 Not Vulnerable{



2
3 < In general, we'll require non-vulnerable
4 data for two purposes: first, we'll need a
5 validated not-vulnerable data set to train our
6 embedding layer with, and second, we'll need
7 data for model training >

8 }
9

10 Vulnerable{



11 < We need it for the susceptible data, and we
12 need to gather as much as possible for the
13 primary model training to cover all distinct
14 patterns of vulnerable source code >

15 }

1
2
3
4
5
6
7
8
9
10
11
12
13
14

02 {

[Filtering Dataset]

< We need to filter our data because
we couldn't use any relevant filters
when we got it from Github >

}

Filtering 'Dataset' {

Step 01 Only files with a character count of fewer than 10,000 were evaluated.

Step 02 Commits that altered or removed a file in more than ten separate places were removed.

Step 03 Remove any projects that act as indicators of security problems or exploits.

Step 04 Read commit messages to aid in the filtration of specific vulnerabilities.

}

1
2
3
4
5
6
7
8
9
10
11
12
13
14

03 {

[Embedding Layer]


< We look into various embedding layers to see which one is ideal for conveying our model requirements >

}

Source code 'representative' {

Embedding layers

W2V  87%

FastText  83%

BERT  91%

< Word2Vec is
state of art
embedding layer
that usually
used in NLP
models >

< FastText is
basically the
same as
word2vec but it
uses different
me >

< BERT is a
pretrained
embedding layer
from microsoft >

}

1
2
3
4
5
6
7
8
9
10
11
12
13
14

04 {

[Labeling]

< We look into various embedding
layers to see which one is ideal for
conveying our model requirements >

}

Labeling {

The data is tagged using information from the commit context :

Misguided attempt:

- * Using only diffs
- * subtle errors in creating the data

Processing the data :

- * vulnerable and non-vulnerable treated the same

}

1
2
3
4
5
6
7
8
9
10
11
12
13
14

05 {

[Training]

< After we labeled our dataset and both embedding layer and model hyperparameters were set we can start training our model >

}

Training process {



Embedding Layer Training

< A substantial training base of clean and working code is necessary to train >



Model Training

< labeled data used for the model training to produce a forecast between 0 and 1>



Different Hyperparameters

< We tested different hyperparameters both for embedding layer and model training>

}

1
2
3
4
5
6
7
8
9
10
11
12
13
14

06 {

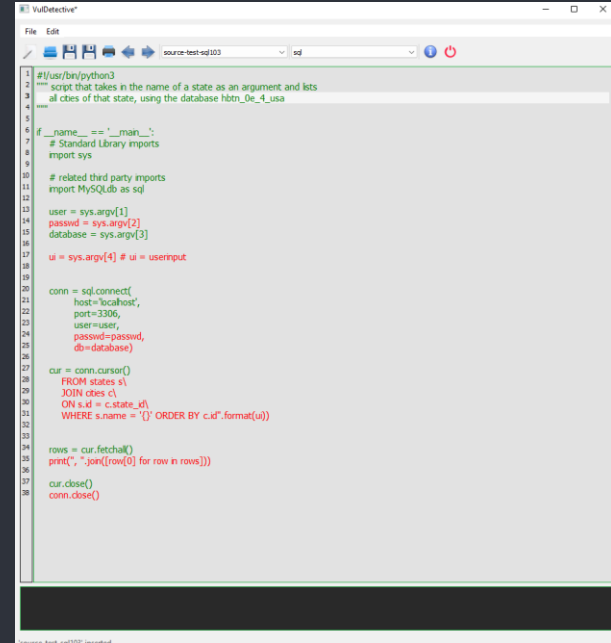
[Application]

< Finally, once the model has been built, we can use it to detect and estimate the percentage of vulnerability for each token, and we may display the results in text edit mode or as an image output >

}

'VulDetective' Application {

< The sections of the
source code that our
model identifies as
vulnerable can be seen >



```
1 #!/usr/bin/python3
2 --- script that takes in the name of a state as an argument and lists
3 all cities of that state, using the database hbtn_0e_4_usa
4 ---
5
6 # __name__ == '__main__':
7 # Standard Library imports
8 import sys
9
10 # related third party imports
11 import MySQLdb as sql
12
13 user = sys.argv[1]
14 passwd = sys.argv[2]
15 database = sys.argv[3]
16
17 ui = sys.argv[4] # ui = userinput
18
19
20 conn = sql.connect(
21     host='localhost',
22     port=3306,
23     user=user,
24     passwd=passwd,
25     db=database)
26
27 cur = conn.cursor()
28 FROM states s
29 JOIN cities c
30 ON s.id = c.state_id
31 WHERE s.name = '{}' ORDER BY c.id'.format(ui)
32
33
34 rows = cur.fetchall()
35 print(", ".join([row[0] for row in rows]))
36
37 cur.close()
38 conn.close()
```

1
2
3
4 ThankYou {
5
6

7 The presented work was carried out
8 within the SETIT Project (2018-1.2.1-
9 NKP-2018-00004). Project no. 2018-
10 1.2.1-NKP-2018-00004 has been
11 implemented with the support provided
12 from the National Research,
13 Development and Innovation Fund of
14 Hungary, financed under the 2018-
1.2.1-NKP funding scheme.

}

