

BÁO CÁO TỔNG KẾT ĐỒ ÁN MÔN HỌC

Môn học: **Lập trình an toàn và khai thác lỗ hổng phần mềm**

Tên chủ đề: **SySeVR: Framework for Using Deep Learning to Detect Software Vulnerabilities**

Mã nhóm: G05 Mã đề tài: CK16

Lớp: **NT521.N11.ANTT**

1. THÔNG TIN THÀNH VIÊN NHÓM:

(Sinh viên liệt kê tất cả các thành viên trong nhóm)

STT	Họ và tên	MSSV	Email
1	Võ Anh Kiệt	20520605	20520605@gm.uit.edu.vn
2	Nguyễn Bùi Kim Ngân	20520648	20520648@gm.uit.edu.vn
3	Nguyễn Bình Thục Trâm	20520815	20520815@gm.uit.edu.vn

2. TÓM TẮT NỘI DUNG THỰC HIỆN:¹

A. Chủ đề nghiên cứu trong lĩnh vực An toàn phần mềm:

- ☒ Phát hiện lỗ hổng bảo mật phần mềm
- ☐ Khai thác lỗ hổng bảo mật phần mềm
- ☐ Sửa lỗi bảo mật phần mềm tự động
- ☐ Lập trình an toàn
- ☐ Khác:

B. Tên bài báo tham khảo chính:

Li, Z., Zou, D., Xu, S., Jin, H., Zhu, Y., & Chen, Z. (2021). Sysevr: A framework for using deep learning to detect software vulnerabilities. IEEE Transactions on Dependable and Secure Computing.

C. Dịch tên Tiếng Việt cho bài báo:

SySeVR – Framework cho việc ứng dụng Deep Learning để việc phát hiện lỗ hổng phần mềm

¹ Ghi nội dung tương ứng theo mô tả

D. Tóm tắt nội dung chính:

<mô tả nội dung tóm tắt của bài báo/chủ đề trong vòng 350 từ>

Phát hiện lỗi hỏng phần mềm là vấn đề rất quan trọng và vẫn chưa được giải quyết triệt để bởi nhiều lỗi hỏng luôn được báo cáo hằng ngày. Do đó mở ra hướng đi ứng dụng machine learning vào phát hiện lỗi hỏng.

Deep learning trở nên hấp dẫn cho mục tiêu trên bởi nó có thể giúp giảm bớt các chức năng yêu cầu thực hiện thủ công. Mặc dù đã đạt được nhiều thành tựu, nhưng deep learning chưa được tìm hiểu một cách hệ thống trong lĩnh vực phát hiện lỗi hỏng. Do vậy, nhóm tác giả đề xuất framework có hệ thống đầu tiên dùng deep learning để phát hiện lỗi hỏng.

SySeVR: Syntax-based, Semantics-based, and Vector Representations (biểu diễn vector dựa trên cú pháp và ngữ nghĩa)

E. Tóm tắt các kỹ thuật chính được mô tả sử dụng trong bài báo:

Về tổng quan model :

Bước 1: Thu thập Input

- + Chương trình cho learning (training)
- + Mục tiêu (testing)

Bước 2: Xử lý: qua các giai đoạn sau

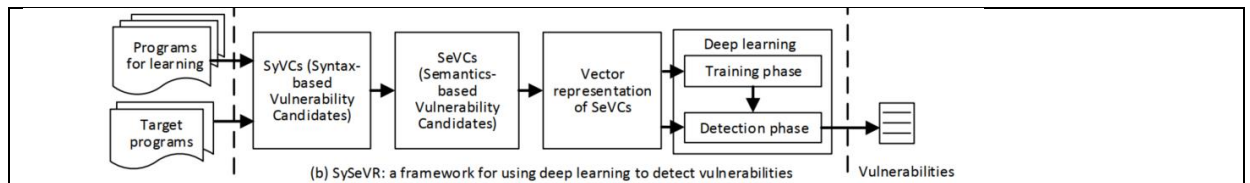
- + Xử lý dựa trên Cú pháp SyVCs
- + Xử lý dựa trên Ngữ nghĩa SeVCs từ SyVCs
- + Biểu diễn Vector từ SeVCs

+ Deep Learning:

- Nếu đi qua pha training thì dữ liệu sẽ hỗ trợ cho phần training model
- Nếu đi qua pha phát hiện lỗi thì sẽ tích hợp model đã được training để thực hiện hỗ trợ trong kiểm tra lỗi và đưa ra báo cáo

Bước 3: Xuất Output:

- + Đưa ra báo cáo về các lỗi đã được tìm thấy



Model gồm 3 thuật toán ứng với cho 3 giai đoạn chính trong bước xử lý

Về chi tiết thuật toán

- Thuật toán 1: Trích xuất đặc điểm cú pháp lỗi hổng từ chương trình

+ Các đặc điểm cú pháp có vai trò xác định các đoạn mã là ứng viên ban đầu của lỗ hổng

Algorithm 1 Extracting SyVCs from a program

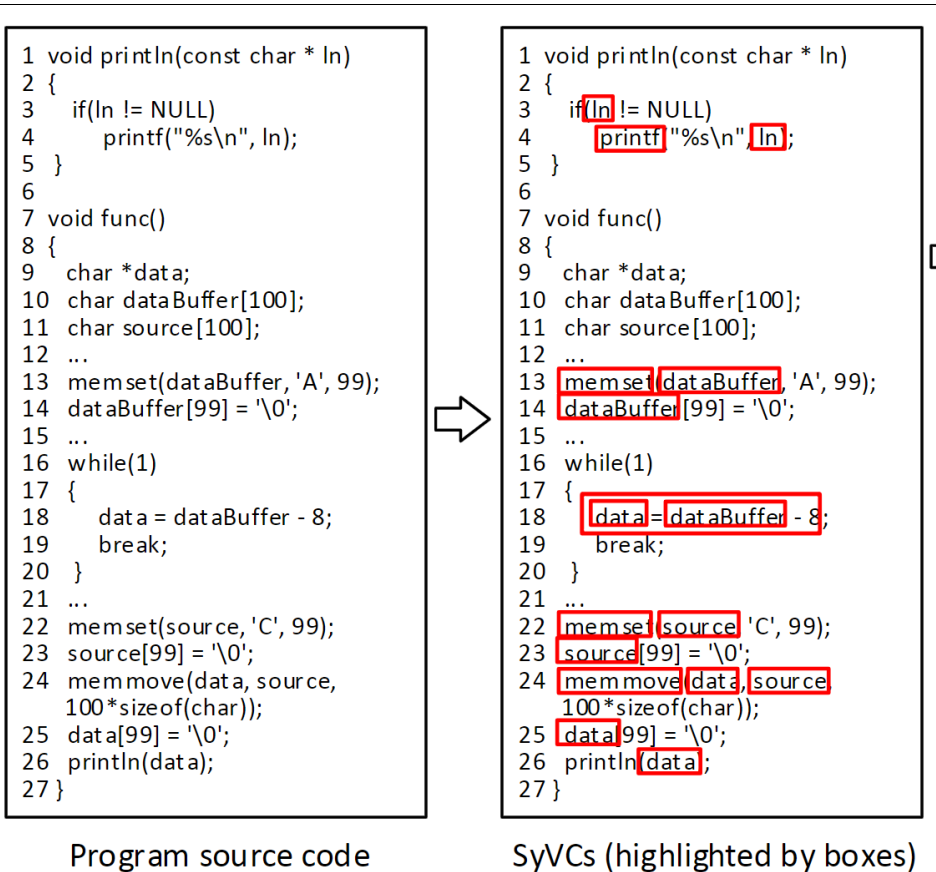
Input: A program $P = \{f_1, \dots, f_n\}$; a set $H = \{h_k\}_{1 \leq k \leq \beta}$ of vulnerability syntax characteristics

Output: A set Y of SyVCs

```

1:  $Y \leftarrow \emptyset$ ;
2: for each function  $f_i \in P$  do
3:   Generate an abstract syntax tree  $T_i$  for  $f_i$ ;
4:   for each code element  $e_{i,j,z}$  in  $T_i$  do
5:     for each  $h_k \in H$  do
6:       if  $e_{i,j,z}$  matches  $h_k$  then
7:          $Y \leftarrow Y \cup \{e_{i,j,z}\}$ ;
8:       end if
9:     end for
10:  end for
11: end for
12: return  $Y$ ; {the set of SyVCs}
    
```

Ví dụ minh họa:



- Thuật toán 2: Chuyển đổi từ SyVCs sang SeVCs

+ Bước này nhằm cung cấp các câu lệnh có liên quan về mặt ngữ nghĩa với SyVCs từ trên

+ Để thực hiện, tác giả sử dụng kỹ thuật cắt lát chương trình. Trong đó, cần dùng Program Dependency Graph (PDG), với yêu cầu dùng data dependency và control dependency (được xác định trên Control Flow Graph CFG)

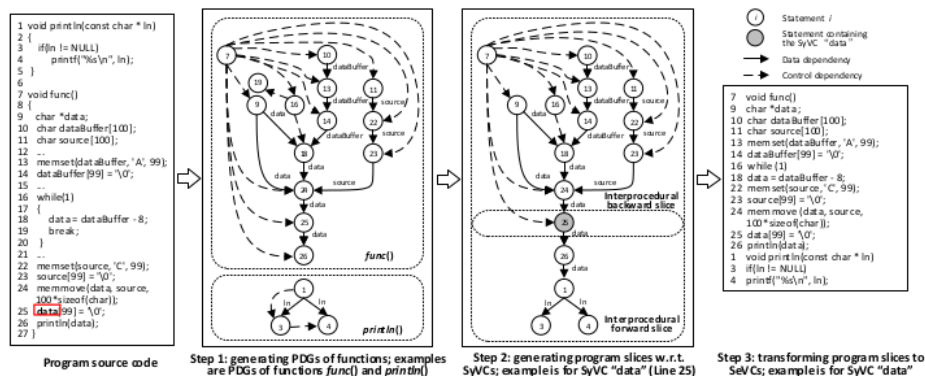
Algorithm 2 Transforming SyVCs to SeVCs

Input: A program $P = \{f_1, \dots, f_n\}$;
a set Y of SyVCs generated by Algorithm 1

Output: The set of SeVCs

- 1: $C \leftarrow \emptyset$;
- 2: **for each** $f_i \in P$ **do**
- 3: Generate a PDG $G'_i = (V_i, E'_i)$ for f_i ;
- 4: **end for**
- 5: **for each** $e_{i,j,z} \in Y$ **in** G'_i **do**
- 6: Generate forward slice $fs_{i,j,z}$ & backward slice $bs_{i,j,z}$ of $e_{i,j,z}$;
- 7: Generate interprocedural forward slice $fs'_{i,j,z}$ by interconnecting $fs_{i,j,z}$ and the forward slices from the functions called by f_i ;
- 8: Generate interprocedural backward slice $bs'_{i,j,z}$ by interconnecting $bs_{i,j,z}$ and the backward slices from both the functions called by f_i and the functions calling f_i ;
- 9: Generate program slice $ps_{i,j,z}$ by connecting $fs'_{i,j,z}$ and $bs'_{i,j,z}$ at $e_{i,j,z}$;
- 10: **for each** statement $s_{i,j} \in f_i$ appearing in $ps_{i,j,z}$ as a node **do**
- 11: $\delta_{i,j,z} \leftarrow \delta_{i,j,z} \cup \{s_{i,j}\}$, according to the order of the appearance of $s_{i,j}$ in f_i ;
- 12: **end for**
- 13: **for two** statements $s_{i,j} \in f_i$ and $s_{a_p,b_q} \in f_{a_p}$ ($i \neq a_p$) appearing in $ps_{i,j,z}$ as nodes **do**
- 14: **if** f_i calls f_{a_p} **then**
- 15: $\delta_{i,j,z} \leftarrow \delta_{i,j,z} \cup \{s_{i,j}, s_{a_p,b_q}\}$, where $s_{i,j} < s_{a_p,b_q}$;
- 16: **else**
- 17: $\delta_{i,j,z} \leftarrow \delta_{i,j,z} \cup \{s_{i,j}, s_{a_p,b_q}\}$, where $s_{i,j} > s_{a_p,b_q}$;
- 18: **end if**
- 19: **end for**
- 20: $C \leftarrow C \cup \{\delta_{i,j,z}\}$;
- 21: **end for**
- 22: **return** C ; {the set of SeVCs}

+ Chi tiết quá trình chuyển đổi SyVCs -> SeVC của SyVC 'data'. Kết quả là SeVC tương ứng với SyVC 'data', cụ thể là là tập hợp thứ tự các câu lệnh có liên quan về mặt ngữ nghĩa với SyVC 'data'.



- Thuật toán 3: Mã hóa SeVCs thành vector

Algorithm 3 Encoding SeVCs into vectors

Input: A set Y of SyVCs generated by Algorithm 1;
a set C of SeVCs corresponding to Y and generated by Algorithm 2;
a threshold θ

Output: The set of vectors corresponding to SeVCs

```

1:  $R \leftarrow \emptyset$ ;
2: for each  $\delta_{i,j,z} \in C$  (corresponding to  $e_{i,j,z} \in Y$ ) do
3:   Remove non-ASCII characters in  $\delta_{i,j,z}$ ;
4:   Map variable names in  $\delta_{i,j,z}$  to symbolic names;
5:   Map function names in  $\delta_{i,j,z}$  to symbolic names;
6: end for
7: for each  $\delta_{i,j,z} \in C$  (corresponding to  $e_{i,j,z} \in Y$ ) do
8:    $R_{i,j,z} \leftarrow \emptyset$ ;
9:   Divide  $\delta_{i,j,z}$  into a set of symbols  $S$ ;
10:  for each  $\alpha \in S$  in order do
11:    Transform  $\alpha$  to a fixed-length vector  $v(\alpha)$ ;
12:     $R_{i,j,z} \leftarrow R_{i,j,z} || v(\alpha)$ , where  $||$  means concatenation;
13:  end for
14:  if  $R_{i,j,z}$  is shorter than  $\theta$  then
15:    Zeroes are padded to the end of  $R_{i,j,z}$ ;
16:  else if the sub-vector (of  $\delta_{i,j,z}$ ) up to the position of the SyVC  $e_{i,j,z}$  is shorter than  $\theta/2$  then
17:    Delete the rightmost portion of  $R_{i,j,z}$  to make the resulting vector of length  $\theta$ ;
18:  else if the sub-vector (of  $\delta_{i,j,z}$ ) next to the the position of the SyVC  $e_{i,j,z}$  is shorter than  $\theta/2$  then
19:    Delete the leftmost portion of  $R_{i,j,z}$  to make the resulting vector of length  $\theta$ ;
20:  else
21:    Keep the sub-vector (in  $\delta_{i,j,z}$ ) immediately left to the position of the SyVC of length  $\lfloor (\theta - 1)/2 \rfloor$ , the sub-vector corresponding to the SyVC, and the sub-vector immediately right to the position of the SyVC of length  $\lceil (\theta - 1)/2 \rceil$  {the resulting vector has length  $\theta$ };
22:  end if
23:   $R \leftarrow R \cup R_{i,j,z}$ ;
24: end for
25: return  $R$ ; {the set of vectors corresponding to SeVCs}

```

- Cuối cùng ta dán nhãn cho vector là có lỗ hổng (1) hoặc không có lỗ hổng (0) để tiến hành learning cho một deep neural network, từ đó nó có thể phát hiện xem SeVCs được cho có chứa lỗ hổng không

F. Môi trường thực nghiệm của bài báo:

- Cấu hình máy tính: NVIDIA GeForce GTX 1080 GPU và Intel Xeon E5-1620 CP 3.50Hz
 - Các công cụ hỗ trợ sẵn có:
- + Checkmarx: giai đoạn extracting SyVCs để phân tích syntax lỗ hổng
- Ngôn ngữ lập trình để hiện thực phương pháp: C/C++
 - Đối tượng nghiên cứu (chương trình phần mềm dùng để kiểm tra tính khả thi của phương pháp/tập dữ liệu – nếu có):
- Vulnerability dataset gồm 2 nguồn từ:
- + NVD: 19 sản phẩm open-source phổ biến, từ đó thu được 1,591 chương trình C/C++, với 874 là có lỗ hổng
 - + SARD: 14,000 chương trình C/C++, với 13,906 là có lỗ hổng
- 4 phần mềm: Libav, Seamonkey, Thunderbird và Xen



- Tiêu chí đánh giá tính hiệu quả của phương pháp: <vd: số lượng lỗ hổng được phát hiện, thời gian chạy, coverage....>
- + False position rate (FPR): tỉ lệ xác định sai mẫu ko có lỗ hổng thành có lỗ hổng (tỉ lệ cảnh báo sai)
- + False negative rate (FNR): tỉ lệ xác định sai mẫu có lỗ hổng thành không có lỗ hổng (tỉ lệ không cảnh báo sai)
- + Accuracy (A): tỉ lệ chính xác (số mẫu đoán đúng/toàn bộ)

G. Kết quả thực nghiệm của bài báo:

Khả năng:

- SySeVR cho phép nhiều loại neural networks phát hiện các loại lỗ hổng bảo mật
- Việc cung cấp thông tin về ngữ nghĩa giúp tăng độ hiệu quả của SySeVR, ví dụ giảm FNR
- Có thể ứng dụng trong các phần mềm thực tế

Ưu:

- Kiểm nghiệm trên 4 phần mềm đã phát hiện được 15 lỗ hổng chưa được báo cáo trong NVD, trong đó 7 cái đã được báo cáo với nhà phát triển, 8 lỗ hổng còn lại đã được lãng lẽ vá lại trong các phiên bản mới
- Thử nghiệm trên nhiều loại neural network model
- Đạt kết quả tốt hơn mô hình trước đó của tác giả là VulDeePecker và các công cụ phát hiện lỗ hổng hiện có.

Nhược:

- Không cho biết thời gian chạy các giai đoạn, chi phí
- Không so sánh với các phương pháp có cùng đề tài ứng dụng deep learning (VulDeePecker là của tác giả) nên thiếu khách quan
- Coverage của đặc điểm cú pháp được xử lí (4 loại) là 93,6% từ dữ liệu của SARD, và dữ liệu SARD không đại diện được cho toàn bộ phần mềm thực tế
- Thí nghiệm chỉ sử dụng đơn model để phát hiện lỗ hổng

H. Công việc/tính năng/kỹ thuật mà nhóm thực hiện lập trình và triển khai cho demo:

Các công việc của nhóm:

Đọc hiểu được ý tưởng, thông tin truyền đạt của bài báo.

Setup môi trường cài đặt.

Chạy file chương trình và đánh giá.

Thực hiện implement và phát triển một số hướng còn hạn chế của bài báo.

Tổng kết và đánh giá

Công việc tương lai:

Cải thiện demo, tìm ra hướng mới cho bài báo

I. Các khó khăn, thách thức hiện tại khi thực hiện:

- Nhóm gặp khó khăn trong đọc hiểu thuật toán
- Thực hiện chạy file từ github của nhóm tác giả có nhiều vấn đề ở cả 2 môi trường window và linux

```

kiet@kiet-Aspire-E5-576: ~/syserv/docker/docker_build
[sudo] password for kiet:
Sending build context to Docker daemon 488.7MB
Step 1/3 : FROM tensorflow/tensorflow:1.6.0-devel-gpu-py3
--> 803f0adee7f
Step 2/3 : COPY ./home/ /home/
--> Using cache
--> 62984a376987
Step 3/3 : RUN mkdir /usr/java && cp -r /home/syserv/softdir/jdk1.8.0_161 /usr/java && mkdir /usr/ant && cp -r /home/syserv/softdir/apache-ant-1.9.14 /usr/ant && rm -rf /etc/apt/sources.list && cp -r /home/syserv/softdir/sources.list /etc/apt/ && rm -rf /etc/apt/sources.list.d && apt-get clean && apt-get update && rm -rf /etc/profile && cp -r /home/syserv/softdir/profile /etc && cd /home/syserv/softdir && chmod +x env.sh && /env.sh && apt-get install -y python-setuptools && apt-get install -y python-dev && apt-get install -y python-pip && cd /home/syserv/softdir/py2neo-py2neo-2.0 && python3 setup.py install && cd /home/syserv/softdir/python-javne-3.1.1 && python3 setup.py install && apt-get install -y graphviz && apt-get install -y libgraphviz-dev && apt-get install -y pkg-config && apt-get install -y python-graph && apt-get install -y python-virtualenv && pip install alrd && pip install gensim-3.4 && pip install pyyaml && rm -rf /home/syserv/softdir
--> Hunting in apt/sources
Err1: http://mirrors.aliyun.com/ubuntu xenial InRelease
Temporary failure resolving 'mirrors.aliyun.com'
Err2: http://mirrors.aliyun.com/ubuntu xenial-updates InRelease
Temporary failure resolving 'mirrors.aliyun.com'
Err3: http://mirrors.aliyun.com/ubuntu xenial-security InRelease
Temporary failure resolving 'mirrors.aliyun.com'
Reading package lists...
W: failed to fetch http://mirrors.aliyun.com/ubuntu/dists/xenial/InRelease. Temporary failure resolving 'mirrors.aliyun.com'
W: failed to fetch http://mirrors.aliyun.com/ubuntu/dists/xenial-updates/InRelease. Temporary failure resolving 'mirrors.aliyun.com'
W: failed to fetch http://mirrors.aliyun.com/ubuntu/dists/xenial-security/InRelease. Temporary failure resolving 'mirrors.aliyun.com'
W: Some index files failed to download. They have been ignored, or old ones used instead.
/home/sh: line 4: /usr/ant/apache-ant-1.9.14/bin/ant: Permission denied
Reading package lists...
Building dependency tree...
Reading state information...
E: Unable to locate package python-setuptools
The command '/bin/sh -c mkdir /usr/java && cp -r /home/syserv/softdir/jdk1.8.0_161 /usr/java && mkdir /usr/ant && cp -r /home/syserv/softdir/apache-ant-1.9.14 /usr/ant && rm -rf /etc/apt/sources.list && cp -r /home/syserv/softdir/sources.list /etc/apt/ && apt-get clean && apt-get update && rm -rf /etc/profile && cp -r /home/syserv/softdir/profile /etc && cd /home/syserv/softdir && chmod +x env.sh && /env.sh && apt-get install -y python-setuptools && apt-get install -y python-dev && apt-get install -y python-pip && cd /home/syserv/softdir/py2neo-py2neo-2.0 && python3 setup.py install && cd /home/syserv/softdir/python-javne-3.1.1 && python3 setup.py install && apt-get install -y graphviz && apt-get install -y libgraphviz-dev && apt-get install -y pkg-config && apt-get install -y python-graph && apt-get install -y python-virtualenv && pip install alrd && pip install gensim-3.4 && pip install pyyaml && rm -rf /home/syserv/softdir' returned a non-zero code: 100
kiet@kiet-Aspire-E5-576: ~/syserv/docker/docker_build
  
```

- Phần code yêu cầu phần cứng hỗ trợ Nvidia

```

1 # syserv Docker image instructions
2
3 *Reminder: the device needs nvidia graphics card, and docker and nvidia-docker2 have been installed in Linux system
4
5 ## 1) Build Image
6
7 The docker_build folder is the working folder where the image is created.
8
9
10
11 bash
12 docker build -t syserv:vi.0 .
13
14
15 "syserv:vi.0" is the name of the created image.
16
17 ## 2) Run container
18
19 execute command:
20
21 bash
22 docker run -it --gpus all --name=syserv -v /home/docker/mapping/implementation:/home/syserv/implementation -v /home/docker/mapping/data:/home/syserv/data syserv:vi.0 /bin/bash
23
24
25 "--name=syserv", syserv is the container name.
26
27 "syserv:vi.0" is the image name obtained in the previous step.
28
29 After entering the container, the folders of /share and /usr/local software required by syserv are under the path of /home/syserv.
30 Other required dependencies have been installed and configured.
31
  
```

3. TỰ ĐÁNH GIÁ MỨC ĐỘ HOÀN THÀNH SO VỚI KẾ HOẠCH THỰC HIỆN:

100%

4. NHẬT KÝ PHÂN CÔNG NHIỆM VỤ:

STT	Công việc	Phân công nhiệm vụ
1	Nghiên cứu và đọc hiểu thuật toán	Nguyễn Bình Thực Trâm Nguyễn Bùi Kim Ngân Võ Anh Kiệt
2	Hiểu luồng dữ liệu, các thành phần và nhiệm vụ của chúng	Nguyễn Bùi Kim Ngân
3	Nghiên cứu các thí nghiệm được trình bày nhằm giải quyết các câu hỏi nghiên cứu được đề ra	Nguyễn Bình Thực Trâm Nguyễn Bùi Kim Ngân Võ Anh Kiệt
4	Triển khai source code, demo	Võ Anh Kiệt Nguyễn Bình Thực Trâm

BÁO CÁO TỔNG KẾT CHI TIẾT

Phần bên dưới của báo cáo này là tài liệu báo cáo tổng kết - chi tiết của nhóm thực hiện cho đề tài này.

Qui định: Mô tả các bước thực hiện/ Phương pháp thực hiện/Nội dung tìm hiểu (Ảnh chụp màn hình, số liệu thống kê trong bảng biểu, có giải thích)

A. Phương pháp thực hiện

Kiến trúc, thành phần của hệ thống trong bài báo:

Về tổng quan model :

Bước 1: Thu thập Input

+ Chương trình cho learning (training)

+ Mục tiêu (testing)

Bước 2: Xử lý: qua các giai đoạn sau

+ Xử lý dựa trên Cú pháp SyVCs

+ Xử lý dựa trên Ngữ nghĩa SeVCs từ SyVCs

+ Biểu diễn Vector từ SeVCs

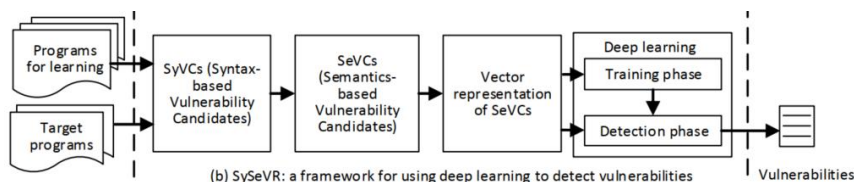
+ Deep Learning:

- Nếu đi qua pha training thì dữ liệu sẽ hỗ trợ cho phần training model

- Nếu đi qua pha phát hiện lỗi thì sẽ tích hợp model đã được training để thực hiện hỗ trợ trong kiểm tra lỗi và đưa ra báo cáo

Bước 3: Xuất Output:

+ Đưa ra báo cáo về các lỗi đã được tìm thấy



Model gồm 3 thuật toán ứng với cho 3 giai đoạn chính trong bước xử lý

Về chi tiết thuật toán

- Thuật toán 1: Trích xuất đặc điểm cú pháp lỗi hỏng từ chương trình
- + Các đặc điểm cú pháp có vai trò xác định các đoạn mã là ứng viên ban đầu của lỗi hỏng

Algorithm 1 Extracting SyVCs from a program

Input: A program $P = \{f_1, \dots, f_n\}$; a set $H = \{h_k\}_{1 \leq k \leq \beta}$ of vulnerability syntax characteristics

Output: A set Y of SyVCs

```

1:  $Y \leftarrow \emptyset$ ;
2: for each function  $f_i \in P$  do
3:   Generate an abstract syntax tree  $T_i$  for  $f_i$ ;
4:   for each code element  $e_{i,j,z}$  in  $T_i$  do
5:     for each  $h_k \in H$  do
6:       if  $e_{i,j,z}$  matches  $h_k$  then
7:          $Y \leftarrow Y \cup \{e_{i,j,z}\}$ ;
8:       end if
9:     end for
10:  end for
11: end for
12: return  $Y$ ; {the set of SyVCs}

```

Ví dụ minh họa:

```

1 void println(const char * ln)
2 {
3   if(ln != NULL)
4     printf("%s\n", ln);
5 }
6
7 void func()
8 {
9   char *data;
10  char dataBuffer[100];
11  char source[100];
12  ...
13  memset(dataBuffer, 'A', 99);
14  dataBuffer[99] = '\0';
15  ...
16  while(1)
17  {
18    data = dataBuffer - 8;
19    break;
20  }
21  ...
22  memset(source, 'C', 99);
23  source[99] = '\0';
24  memmove(data, source,
25    100*sizeof(char));
26  data[99] = '\0';
27  println(data);
28 }

```

Program source code



```

1 void println(const char * ln)
2 {
3   if(ln != NULL)
4     printf("%s\n", ln);
5 }
6
7 void func()
8 {
9   char *data;
10  char dataBuffer[100];
11  char source[100];
12  ...
13  memset(dataBuffer, 'A', 99);
14  dataBuffer[99] = '\0';
15  ...
16  while(1)
17  {
18    data = dataBuffer - 8;
19    break;
20  }
21  ...
22  memset(source, 'C', 99);
23  source[99] = '\0';
24  memmove(data, source,
25    100*sizeof(char));
26  data[99] = '\0';
27  println(data);
28 }

```

SyVCs (highlighted by boxes)

- Thuật toán 2: Chuyển đổi từ SyVCs sang SeVCs

+ Bước này nhằm cung cấp các câu lệnh có liên quan về mặt ngữ nghĩa với SyVCs từ trên

+ Để thực hiện, tác giả sử dụng kỹ thuật cắt lát chương trình. Trong đó, cần dùng Program Dependency Graph (PDG), với yêu cầu dùng data dependency và control dependency (được xác định trên Control Flow Graph CFG)

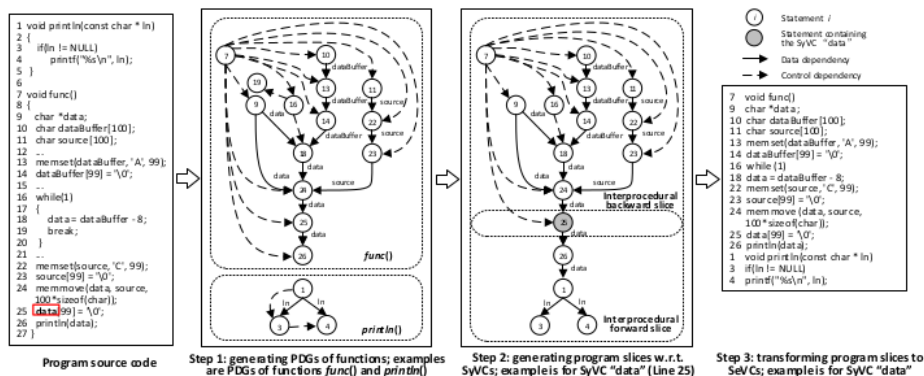
Algorithm 2 Transforming SyVCs to SeVCs

Input: A program $P = \{f_1, \dots, f_n\}$;
a set Y of SyVCs generated by Algorithm 1

Output: The set of SeVCs

- 1: $C \leftarrow \emptyset$;
- 2: **for each** $f_i \in P$ **do**
- 3: Generate a PDG $G'_i = (V_i, E'_i)$ for f_i ;
- 4: **end for**
- 5: **for each** $e_{i,j,z} \in Y$ **in** G'_i **do**
- 6: Generate forward slice $fs_{i,j,z}$ & backward slice $bs_{i,j,z}$ of $e_{i,j,z}$;
- 7: Generate interprocedural forward slice $fs'_{i,j,z}$ by interconnecting $fs_{i,j,z}$ and the forward slices from the functions called by f_i ;
- 8: Generate interprocedural backward slice $bs'_{i,j,z}$ by interconnecting $bs_{i,j,z}$ and the backward slices from both the functions called by f_i and the functions calling f_i ;
- 9: Generate program slice $ps_{i,j,z}$ by connecting $fs'_{i,j,z}$ and $bs'_{i,j,z}$ at $e_{i,j,z}$;
- 10: **for each statement** $s_{i,j} \in f_i$ **appearing in** $ps_{i,j,z}$ **as a node do**
- 11: $\delta_{i,j,z} \leftarrow \delta_{i,j,z} \cup \{s_{i,j}\}$, according to the order of the appearance of $s_{i,j}$ in f_i ;
- 12: **end for**
- 13: **for two statements** $s_{i,j} \in f_i$ **and** $s_{a_p,b_q} \in f_{a_p}$ ($i \neq a_p$) **appearing in** $ps_{i,j,z}$ **as nodes do**
- 14: **if** f_i **calls** f_{a_p} **then**
- 15: $\delta_{i,j,z} \leftarrow \delta_{i,j,z} \cup \{s_{i,j}, s_{a_p,b_q}\}$, where $s_{i,j} < s_{a_p,b_q}$;
- 16: **else**
- 17: $\delta_{i,j,z} \leftarrow \delta_{i,j,z} \cup \{s_{i,j}, s_{a_p,b_q}\}$, where $s_{i,j} > s_{a_p,b_q}$;
- 18: **end if**
- 19: **end for**
- 20: $C \leftarrow C \cup \{\delta_{i,j,z}\}$;
- 21: **end for**
- 22: **return** C ; {the set of SeVCs}

+ Chi tiết quá trình chuyển đổi SyVCs -> SeVC của SyVC 'data'. Kết quả là SeVC tương ứng với SyVC 'data', cụ thể là là tập hợp thứ tự các câu lệnh có liên quan về mặt ngữ nghĩa với SyVC 'data'.



- Thuật toán 3: Mã hóa SeVCs thành vector

Algorithm 3 Encoding SeVCs into vectors

Input: A set Y of SyVCs generated by Algorithm 1;
a set C of SeVCs corresponding to Y and generated by Algorithm 2;
a threshold θ

Output: The set of vectors corresponding to SeVCs

```

1:  $R \leftarrow \emptyset$ ;
2: for each  $\delta_{i,j,z} \in C$  (corresponding to  $e_{i,j,z} \in Y$ ) do
3:   Remove non-ASCII characters in  $\delta_{i,j,z}$ ;
4:   Map variable names in  $\delta_{i,j,z}$  to symbolic names;
5:   Map function names in  $\delta_{i,j,z}$  to symbolic names;
6: end for
7: for each  $\delta_{i,j,z} \in C$  (corresponding to  $e_{i,j,z} \in Y$ ) do
8:    $R_{i,j,z} \leftarrow \emptyset$ ;
9:   Divide  $\delta_{i,j,z}$  into a set of symbols  $S$ ;
10:  for each  $\alpha \in S$  in order do
11:    Transform  $\alpha$  to a fixed-length vector  $v(\alpha)$ ;
12:     $R_{i,j,z} \leftarrow R_{i,j,z} || v(\alpha)$ , where  $||$  means concatenation;
13:  end for
14:  if  $R_{i,j,z}$  is shorter than  $\theta$  then
15:    Zeroes are padded to the end of  $R_{i,j,z}$ ;
16:  else if the sub-vector (of  $\delta_{i,j,z}$ ) up to the position of the SyVC  $e_{i,j,z}$  is shorter than  $\theta/2$  then
17:    Delete the rightmost portion of  $R_{i,j,z}$  to make the resulting vector of length  $\theta$ ;
18:  else if the sub-vector (of  $\delta_{i,j,z}$ ) next to the the position of the SyVC  $e_{i,j,z}$  is shorter than  $\theta/2$  then
19:    Delete the leftmost portion of  $R_{i,j,z}$  to make the resulting vector of length  $\theta$ ;
20:  else
21:    Keep the sub-vector (in  $\delta_{i,j,z}$ ) immediately left to the position of the SyVC of length  $\lfloor (\theta - 1)/2 \rfloor$ , the sub-vector corresponding to the SyVC, and the sub-vector immediately right to the position of the SyVC of length  $\lceil (\theta - 1)/2 \rceil$  {the resulting vector has length  $\theta$ };
22:  end if
23:   $R \leftarrow R \cup R_{i,j,z}$ ;
24: end for
25: return  $R$ ; {the set of vectors corresponding to SeVCs}

```

Cuối cùng ta dán nhãn cho vector là có lỗ hổng (1) hoặc không có lỗ hổng (0) để tiến hành learning cho một deep neural network, từ đó nó có thể phát hiện xem SeVCs được cho có chứa lỗ hổng không

Kiến trúc, thành phần nhóm đã thực hiện:

- Thuật toán 1: Trích xuất đặc điểm cú pháp lỗ hổng từ chương trình

Bước 1: Đầu tiên sử dụng joern để phân tích mã nguồn: đầu vào là tệp mã nguồn và đầu ra là tệp có tên joernIndex.

Bước 2: Tiếp theo sử dụng code get_cfg_relation.py để lấy đồ thị hàm CFG bằng công cụ joern. Đầu vào là kết quả của bước 1 và đầu ra được lưu trữ với các thư mục trong cfg_db

Bước 3: Thực hiện để lấy đồ thị hàm PDG. Đầu vào là các tệp trong cfg_db và đầu ra được lưu trữ với các thư mục trong pdg_db

Bước 4: Thực hiện để lấy đồ thị gọi hàm. Đầu vào là các tệp trong pdg_db và đầu ra được lưu trữ với các thư mục trong dict_call2cfgNodeID_funcID

Bước 5: Thực hiện để nhận bốn loại SyVC. Đầu vào là các tệp trong dict_call2cfgNodeID_funcID và đầu ra là bốn loại SyVC.

Bước 6: Thực hiện để trích xuất các lát. Đầu vào là các tệp được tạo bởi point_get.py và đầu ra là các tệp lát

Bước 7: Thực hiện để lấy số dòng của các dòng dễ bị tấn công trong bộ dữ liệu nvd. Đầu vào là tệp mã nguồn hoặc tệp func và tệp khác

Bước 8: Thực hiện để trích xuất số dòng của các dòng dễ bị tấn công từ SARD_testcaseinfo.xml. "000" là tệp mã nguồn. Đầu ra là SARD_testcaseinfo.txt, sau đó được đổi tên thành contain_all.txt

Bước 9: Tạo nhãn NVD và SARD

Bước 10: Thực hiện để add label vào slide

B. Chi tiết cài đặt, hiện thực

Cài đặt, lập trình trên máy tính, cấu hình máy tính sử dụng, chuẩn bị dữ liệu:

Môi trường:

Hệ điều hành: Ubuntu 18.04

Môi trường cho thuật toán 1:

Joern 0.3.1:

+ JDK 1.7

+ Neo4J 2.1.8 Community Edition

+ Gremlin for Neo4J 2.X

+ Apache Ant build tool

Python 2.7

Môi trường cho thuật toán 2:

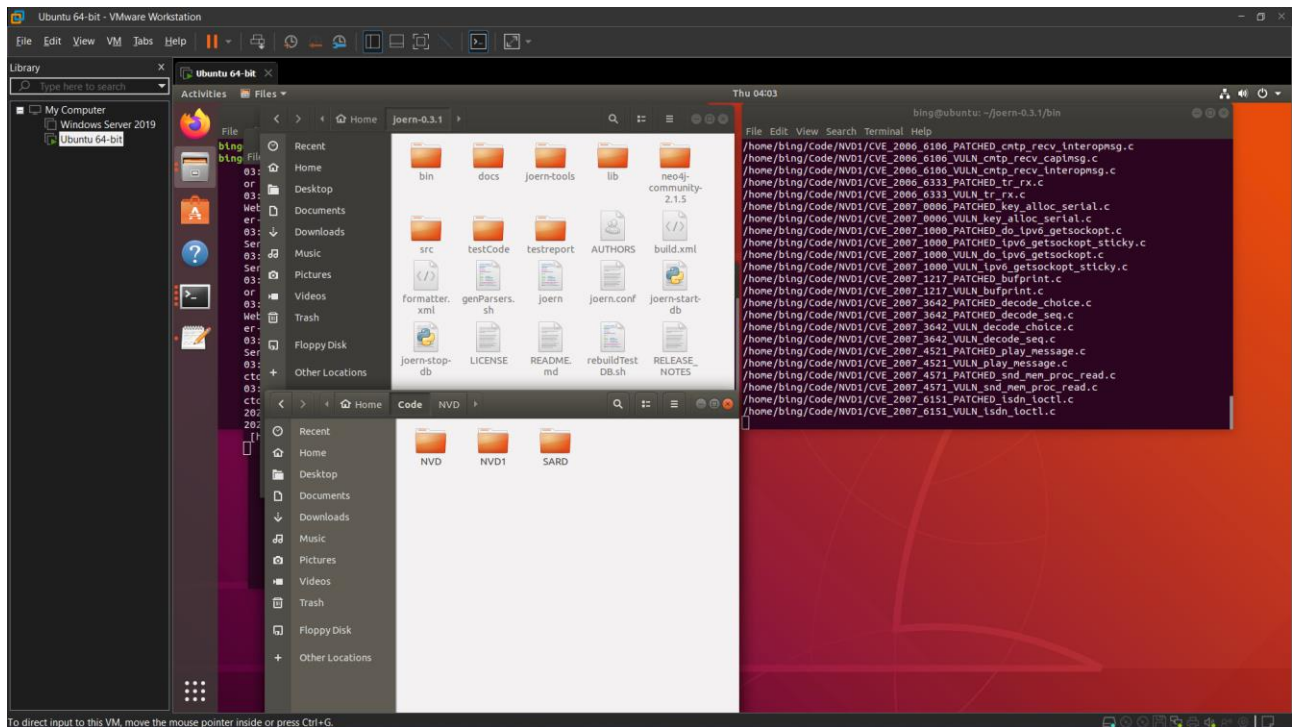
Python 3.6

Tensorflow 1.6

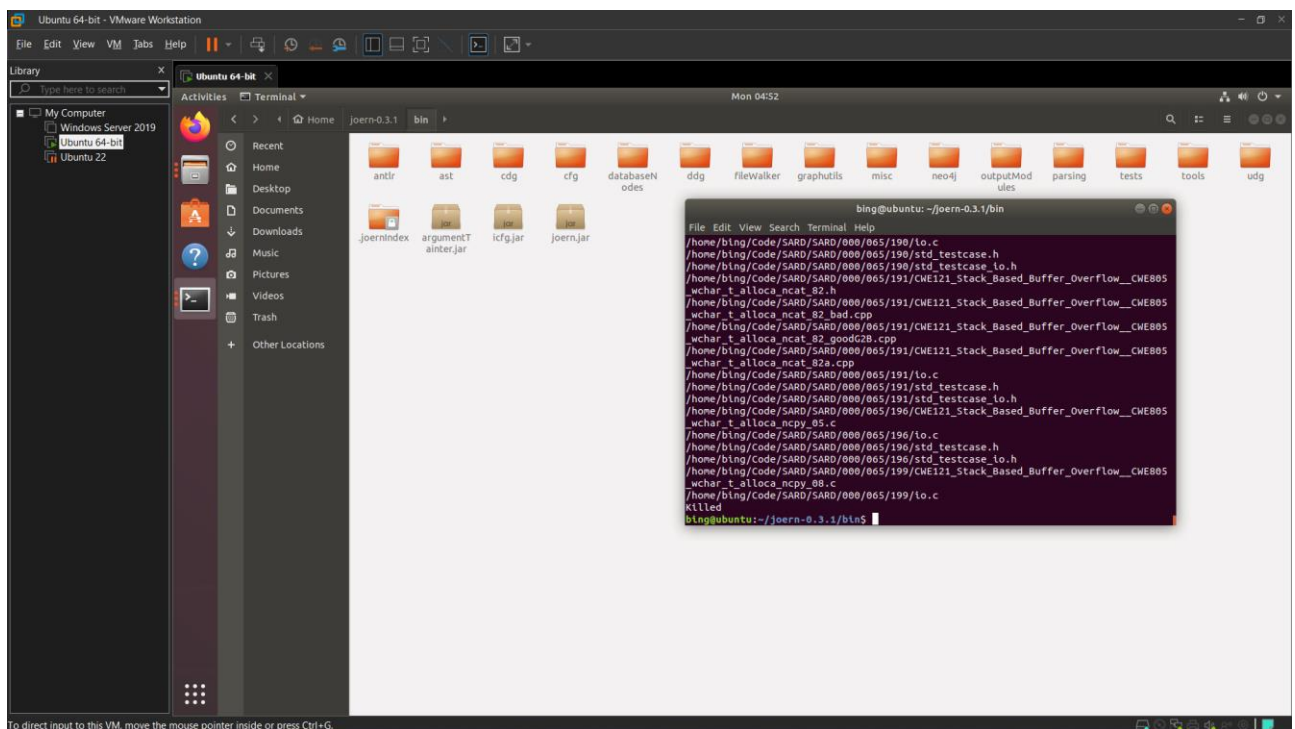
Gensim 3.4

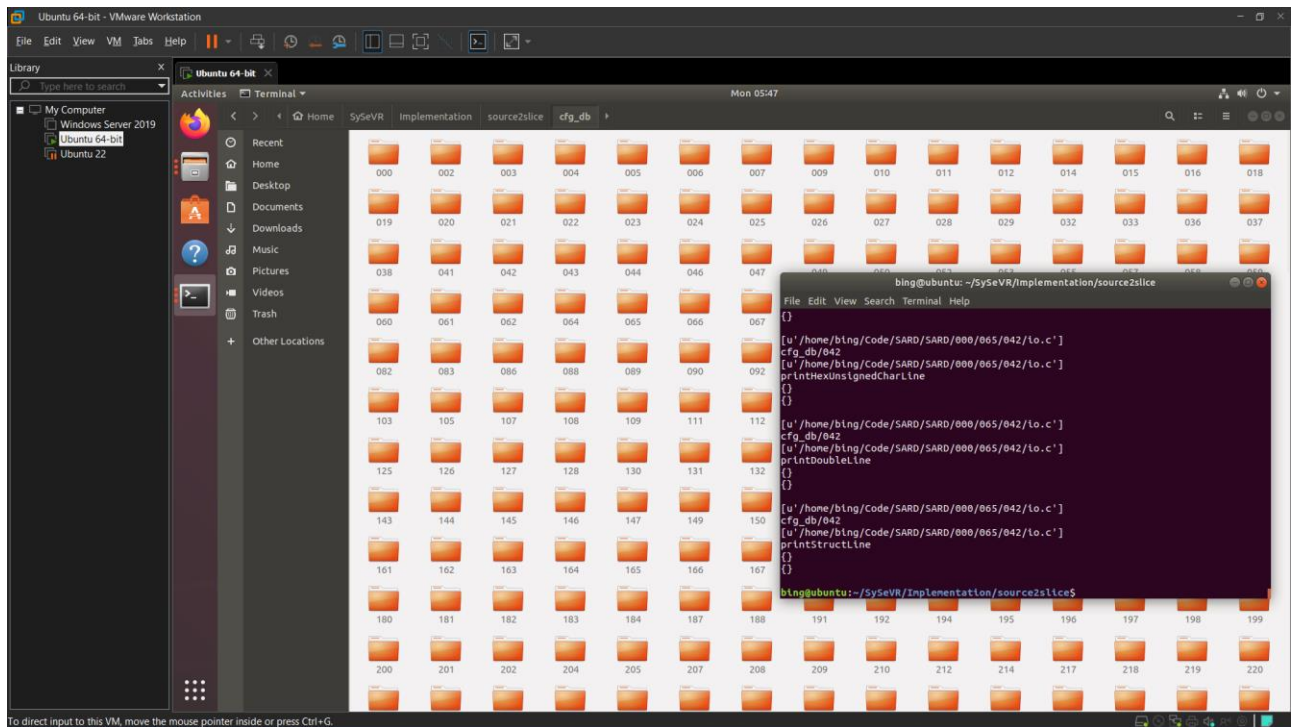
C. Kết quả thực nghiệm

Step 1 bước 1



Step 1 bước 2





To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

Step 1 bước 3



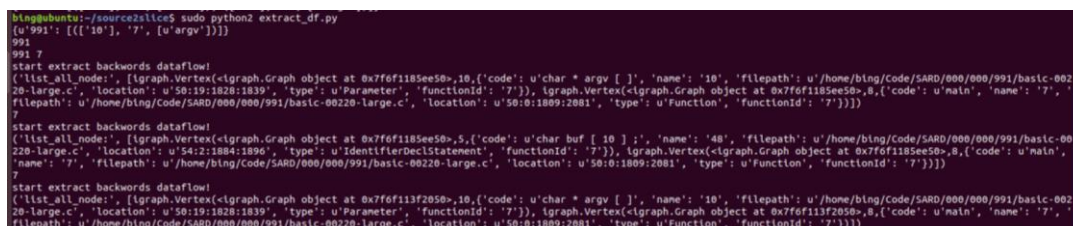
Step 1 bước 4



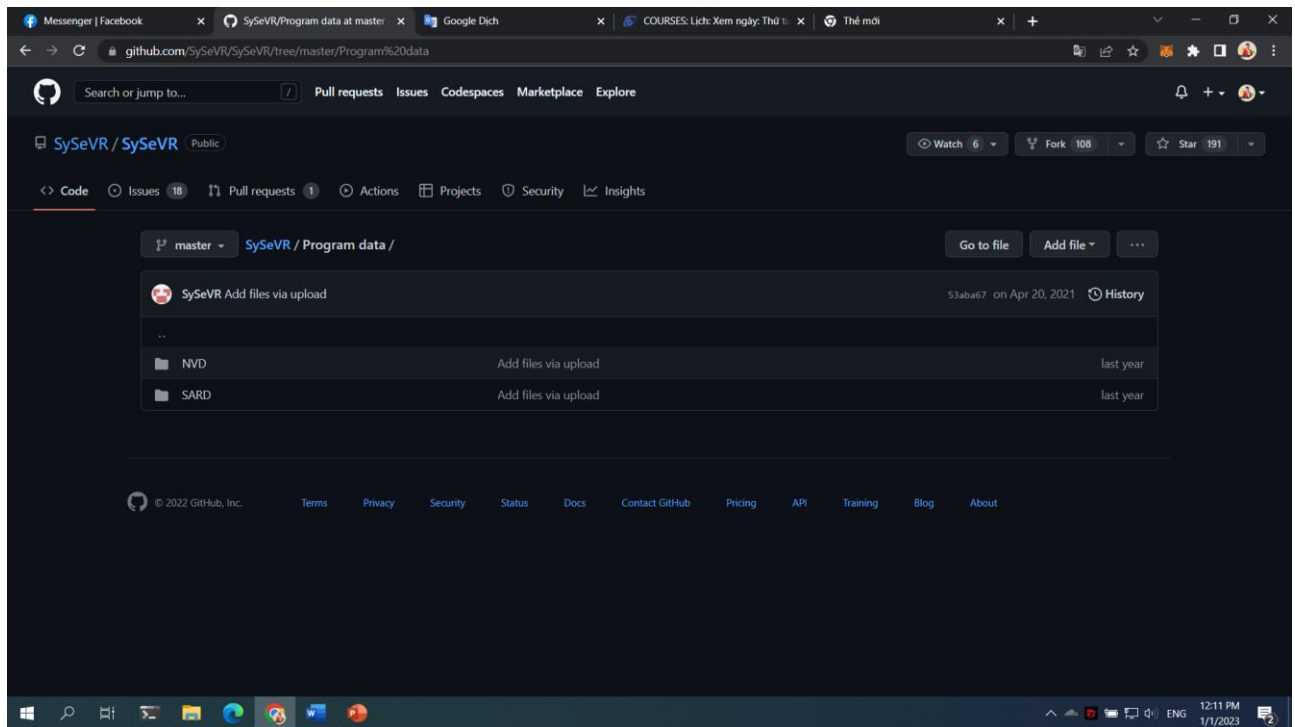
Step 1 bước 5



Step 1 bước 6 7



Step 1 bước 8 sử dụng lại kết quả của bài



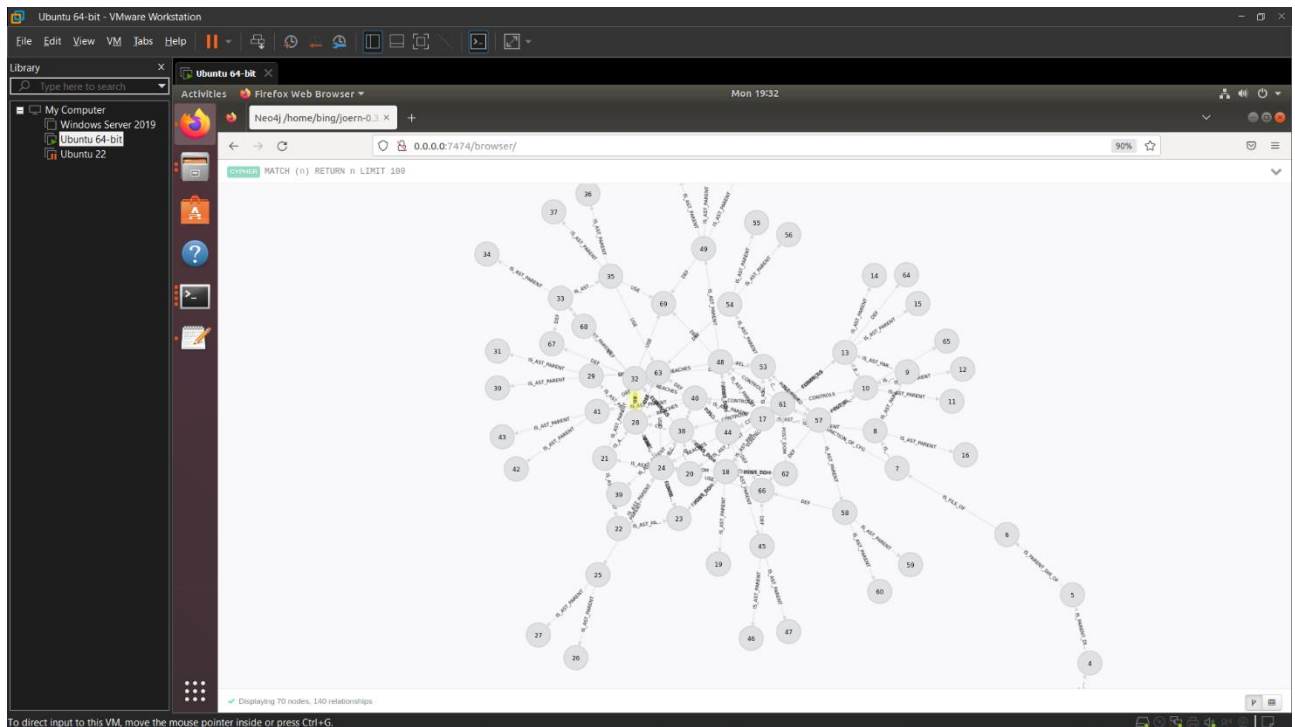
Step 1 bước 9

```
bing@ubuntu:~/source2slice$ sudo python2 make_label.py
1 /home/bing/Code/SARD/000/000/991/basic-00220-large.c [u'buf'] 54
2 /home/bing/Code/SARD/000/000/991/basic-00220-large.c [u'argv'] 50
1 /home/bing/Code/SARD/000/000/991/basic-00220-large.c [u'argv'] 50
bing@ubuntu:~/source2slice$ cp arraysuse_slices.txt integeroverflow_slices.txt pointeruse_slices.txt api_slices.txt ./slices
bing@ubuntu:~/source2slice$ cp api_slices_label.pkl api_slices_vulnline.pkl array_slice_label.pkl expr_slice_label.pkl pointer_slice_label.pkl ./label_source
bing@ubuntu:~/source2slice$ cd label_source
bing@ubuntu:~/source2slice/label_source$ mv expr_slice_label.pkl integeroverflow_slices.pkl
bing@ubuntu:~/source2slice/label_source$ mv api_slices_label.pkl api_slices.pkl
bing@ubuntu:~/source2slice/label_source$ mv array_slice_label.pkl arraysuse_slices.pkl
bing@ubuntu:~/source2slice/label_source$ mv pointer_slice_label.pkl pointeruse_slice.pkl
bing@ubuntu:~/source2slice/label_source$ cd ..
```

Step 1 bước 10

```
bing@ubuntu:~/source2slice$ sudo python2 data_preprocess.py
api_slices.txt
[... ]
integeroverflow_slices.txt
[... ]
arraysuse_slices.txt
[... ]
1 /home/bing/Code/SARD/000/000/991/basic-00220-large.c [u'buf'] 54\nint main(int argc, char *argv[]) 50\nchar buf [ 10 ]; 54\nbuf [ loop_counter ] = 'A'; 62\nloop_counter ++
... ; 63\nif ( loop_counter > test_value ) 64\n, "\n2 /home/bing/Code/SARD/000/000/991/basic-00220-large.c [u'argv'] 50\nint main(int argc, char *argv[]) 50\n", '\n']
Traceback (most recent call last):
  File "data_preprocess.py", line 27, in <module>
    raise ValueError("house pointer inside or press Ctrl+C")
```

Check database



D. Hướng phát triển

Đề tài có nhiều tiềm năng phát triển trong tương lai khi việc thực hiện tự động hóa quá trình thực hiện tìm kiếm lỗ hổng phần mềm, có tính ứng dụng thực tiễn cao trong thực tế. Có thể phát triển đồ án để tiếp tục thực hiện đồ án chuyên ngành hoặc khóa luận tốt nghiệp

Sinh viên báo cáo các nội dung mà nhóm đã thực hiện, có thể là 1 phần hoặc toàn bộ nội dung của bài báo. Nếu nội dung thực hiện có khác biệt với bài báo (như cấu hình, tập dữ liệu, kết quả,...), sinh viên cần chỉ rõ thêm khác biệt đó và nguyên nhân.

Sinh viên đọc kỹ yêu cầu trình bày bên dưới trang này

YÊU CẦU CHUNG

- Sinh viên tìm hiểu và thực hiện bài tập theo yêu cầu, hướng dẫn.
- Nộp báo cáo kết quả chi tiết những việc (**Report**) bạn đã thực hiện, quan sát thấy và kèm ảnh chụp màn hình kết quả (nếu có); giải thích cho quan sát (nếu có).
- Sinh viên báo cáo kết quả thực hiện và nộp bài.

Báo cáo:

- File **.PDF**. Tập trung vào nội dung, không mô tả lý thuyết.
- Đặt tên theo định dạng: [Mã lớp]-Project_Final_NhomX_Madetai. (trong đó X và Madetai là mã số thứ tự nhóm và Mã đề tài trong danh sách đăng ký nhóm đồ án).
Ví dụ: [NT521.N11.ANTT]-Project_Final_Nhom03_CK01.
- Nếu báo cáo có nhiều file, nén tất cả file vào file .ZIP với cùng tên file báo cáo.
- Nộp file báo cáo trên theo thời gian đã thống nhất tại courses.uit.edu.vn.

Đánh giá:

- Hoàn thành tốt yêu cầu được giao.
- Có nội dung mở rộng, ứng dụng.

Bài sao chép, trễ, ... sẽ được xử lý tùy mức độ vi phạm.

HẾT