

基于图神经网络的切片级漏洞检测及解释方法^{*}



胡雨涛^{1,2,3}, 王溯远^{1,2,3}, 吴月明^{1,2,3}, 邹德清^{1,2,3}, 李文科^{2,3}, 金海^{1,4}

¹(大数据技术与系统国家地方联合工程研究中心 服务计算技术与系统教育部重点实验室, 湖北 武汉 430074)

²(分布式系统安全湖北省重点实验室, 湖北 武汉 430074)

³(华中科技大学 网络空间安全学院, 湖北 武汉 430074)

⁴(华中科技大学 计算机科学与技术学院, 湖北 武汉 430074)

通讯作者: 吴月明, E-mail: wuyueming21@gmail.com

摘要: 随着软件的复杂程度越来越高,对漏洞检测的研究需求也日益增大.软件漏洞的迅速发现和修补可以将漏洞带来的损失降到最低.基于深度学习的漏洞检测方法作为目前新兴的检测手段,可以从漏洞代码中自动学习其隐含的漏洞模式,节省了大量人力投入.但基于深度学习的漏洞检测方法尚未完善,其中函数级别的检测方法存在检测粒度较粗且检测准确率较低的问题,切片级别的检测方法虽然能够有效减少样本噪声,但是仍存在以下两个方面的问题.一方面现有方法大多采用人工漏洞数据集进行实验,因此其在真实环境中的漏洞检测能力仍然存疑;另一方面,相关工作仅致力于检测出切片样本是否存在漏洞,而缺乏对检测结果可解释性的考虑.针对上述问题,本文提出了基于图神经网络的切片级漏洞检测及解释方法.该方法首先对 C/C++源代码进行规范化并提取切片,以减少样本冗余信息干扰;其次采用图神经网络模型进行切片嵌入得到其向量表征,以保留源代码的结构信息和漏洞特征;然后将切片的向量表征输入漏洞检测模型进行训练和预测;最后将训练完成的漏洞检测模型和待解释的漏洞切片输入漏洞解释器得到具体的漏洞代码行.实验结果显示,在漏洞检测方面,该方法对于真实漏洞数据的检测 F1 分数达到 75.1%,相较于对比方法提升了 41.2%-110.4%.在漏洞解释方面,该方法在限定前 10%的关键节点时准确率可以达到 73.6%,相较于两种对比解释器分别提升 8.9%和 24.9%,且时间开销分别缩短 42.5%和 15.4%.最后,该方法正确检测并解释了 4 个开源软件中 59 个真实漏洞,证明了其在现实世界漏洞发掘方面的实用性.

关键词: 漏洞检测;深度学习;图神经网络;人工智能可解释性

中图法分类号: TP311

中文引用格式:胡雨涛,王溯远,吴月明,邹德清,李文科,金海.基于图神经网络的切片级漏洞检测及解释方法.软件学报.
<http://www.jos.org.cn/1000-9825/6849.htm>

英文引用格式: Hu YT, Wang SY, Wu YM, Zou DQ, Li WK, Jin H. A slice-level vulnerability detection and interpretation method based on graph neural network. Ruan Jian Xue Bao/Journal of Software (in Chinese). <http://www.jos.org.cn/1000-9825/6849.htm>

A Slice-level vulnerability detection and interpretation method based on graph neural network

HU Yu-Tao^{1,2,3}, WANG Su-Yuan^{1,2,3}, WU Yue-Ming^{1,2,3}, ZOU De-Qing^{1,2,3}, LI Wen-Ke^{2,3}, JIN Hai^{1,4}

¹(National Engineering Research Center for Big Data Technology and System Services Computing Technology and System Lab, Wuhan, 430074, China)

²(Hubei Key Laboratory of Distributed System Security, Wuhan, 430074, China)

³(School of Cyber Science and Engineering, Huazhong University of Science and Technology, Wuhan 430074, China)

⁴(School of Computer Science and Engineering, Huazhong University of Science and Technology, Wuhan 430074, China)

* 基金项目: 源代码智能漏洞检测模型的抗攻击研究(62172168); 面向云环境的智能防护关键技术研究(2021BAA032)

收稿时间: 2022-09-05; 修改时间: 2022-10-10; 采用时间: 2022-12-14; jos 在线出版时间: 2023-01-13

Abstract: As software becomes more complex, the need for research on vulnerability detection is increasing. The rapid discovery and patching of software vulnerabilities is able to minimize the damage caused by vulnerabilities. As an emerging detection method, deep learning-based vulnerability detection methods can learn from the vulnerability code and automatically generate its implied vulnerability pattern, saving a lot of human effort. However, deep learning-based vulnerability detection methods are not yet perfect; function-level detection methods have a coarse detection granularity with low detection accuracy; slice-level detection methods can effectively reduce sample noise, but there are still the following two aspects of the problem: On the one hand, most of the existing methods use artificial vulnerability datasets for experiments, and the ability to detect vulnerabilities in real environments is still in doubt; on the other hand, the work is only dedicated to detecting the existence of vulnerabilities in the slice samples and the lack of interpretability of the detection results. To address above issues, this paper proposes a slice-level vulnerability detection and interpretation method based on the graph neural network. The method first normalizes the C/C++ source code and extracts slices to reduce the interference of redundant information in the samples; secondly, a graph neural network model is used to embed the slices to obtain their vector representations to preserve the structural information and vulnerability features of the source code; then the vector representations of slices are fed into the vulnerability detection model for training and prediction; finally, the trained vulnerability detection model and the vulnerability slices to be explained are fed into the vulnerability interpreter to obtain the specific lines of vulnerability code. The experimental results show that in terms of vulnerability detection, the method achieves an F1 score of 75.1% for real-world vulnerability, which is 41.2%-110.4% higher than the comparative methods. In terms of vulnerability interpretation, the method can reach 73.6% accuracy when limiting the top 10% of critical nodes, which is 8.9% and 24.9% higher than the other two interpreters, and the time overhead is reduced by 42.5% and 15.4%, respectively. Finally, this method correctly detects and explains 59 real vulnerabilities in the four open-source software, proving its practicality in real-world vulnerability discovery.

Key words: Vulnerability Detection; Deep Learning; Graph Neural Network; Explainable AI

网络空间安全已成为国家安全的关键组成部分和重要保障.近年来,网络黑客勒索、僵尸网络攻击、用户信息泄露等各种网络空间安全事件频发,给用户的人身财产以及国家安全造成了严重威胁.漏洞是各种网络空间安全事件发起的根源,软件系统作为网络空间的一个重要组成部分,其本身的漏洞给网络空间带来了严重的安全威胁.针对源代码进行漏洞检测是发现软件漏洞的主要途径之一^[1],且具有天然优势:一方面,从源代码中可以获得丰富的语义信息,便于表示软件的漏洞特征;另一方面,目前开源软件得到了广泛大量的应用,其中的漏洞也随之广泛传播,因此针对源代码检测具有明显优势和意义.随着机器学习特别是深度学习的兴起,最近的研究开始利用深度学习模型进行源代码漏洞检测^{[2][3][4]},其因为不依赖专家经验而减少了人工干预的主观性,因此能够取得较低的误报率和漏报率.

现阶段基于深度学习的漏洞检测方法按照样本粒度可以分为函数级别和切片级别.其中,函数级别面向整个函数进行漏洞检测,其优势在于函数能够涵盖较为全面的漏洞信息,但同时会引入大量漏洞无关的噪声语句,导致检测的准确率欠佳.切片级别的漏洞检测优化了样本粒度,即去除了样本中漏洞无关的噪声语句,进而使得模型能够学习到更为精准的漏洞特征.然而,已有的切片级漏洞检测工作主要依赖人工构造的漏洞数据集如 Software Assurance Reference Dataset^[5](SARD)以证明其方法的有效性,因此其在真实环境中的漏洞检测能力仍然存疑.此外,相关工作仅致力于检测出切片样本是否存在漏洞,而缺乏对检测结果可解释性的考虑.由于深度学习模型内部涉及到大量复杂的特征计算,其所学习到的漏洞特征难以被安全分析人员理解,导致了检测结果无法被信服.而模型的可解释性不但是决定用户是否能够信任检测模型的关键因素,同时能够帮助用户进一步分析漏洞成因进而快速修补漏洞.

针对上述问题,本文提出了一种基于图神经网络的切片级漏洞检测及解释方法.该方法首先针对 C/C++源代码进行规范化并提取切片,以减少样本冗余信息干扰;其次采用图神经网络模型进行切片嵌入得到其向量表征,以保留源代码的结构信息和漏洞特征;然后将切片的向量表征输入漏洞检测模型进行训练和预测;最后将训练完成的漏洞检测模型和待解释的漏洞代码输入漏洞解释器得到具体的漏洞代码行.

为了验证本方法对于真实漏洞检测和解释的有效性,本文采用真实漏洞数据集进行全部实验.在漏洞检测方面,本文选取了3个基于规则的漏洞检测工具(Checkmark^[6],FlawFinder^[7]和 RATS^[8])和4个基于深度学习的漏洞检测方法(TokednCNN^[9],StatementLSTM^[10],SySeVR^[11]和 Devign^[12])作为对比工具.实验结果表明,本方法 F1 分数可以达到 75.1%,相较于其他漏洞检测方法提升了 41.2%-110.4%.在漏洞解释方面,本文选取了两种先进的

图神经网络解释器(GNNExplainer^[13]和 PGExplainer^[14])作为对比解释器.实验结果表明,本解释方法在限定前10%的关键节点时准确率可以达到73.6%,相较于对比方法分别提升8.9%和24.9%,且在时间开销方面分别缩短42.5%和15.4%.最后,为了检验本方法在现实世界漏洞发掘方面的实用性,本文选取了4个真实的开源软件进行测试,并成功检测和解释了开源软件中59个真实漏洞,实验结果证明本方法能够有效检测真实软件漏洞并提供可靠的检测结果解释.

本文的主要贡献如下:

- (1) 提出了一种基于图神经网络的切片级漏洞检测方法.通过提取代码切片有效降低了代码冗余,减少了无关语句的干扰.同时图神经网络可以更好地保留代码结构和漏洞语法语义信息,使得其针对复杂的真实漏洞数据达到较好的检测效果.
- (2) 利用并改进了一种图神经网络解释器,使通用的图神经网络解释器更适用于漏洞解释任务.本文通过改进 GNNExplainer 使得其能够快速且准确地输出漏洞代码行,增强了基于深度学习漏洞检测方法的可信度,且有助于研究人员分析漏洞成因并修复漏洞.
- (3) 本方法扫描了4个开源软件的3312657个切片,并通过人工分析和与已知漏洞模式匹配的方式验证检测和解释结果,实验结果证明其正确检测并解释了59个真实漏洞,证实了本方法在现实世界漏洞发掘方面的实用性.

本文第1节介绍基于深度学习漏洞检测及解释的相关工作.第2节介绍本文所需的基础知识,包括静态程序分析和图神经网络的相关概念.第3节介绍本文基于图神经网络的切片级漏洞检测及解释方法的具体实现.第4节通过对比实验验证了本方法在漏洞检测、漏洞解释和实际应用三个方面的有效性.最后总结全文并介绍未来工作.

1 相关工作

本节将对本文中涉及的相关工作进行介绍,主要包含基于深度学习的漏洞检测和模型可解释性两个方面.

1.1 基于深度学习漏洞检测相关工作

近年来,由于深度学习模型强大的建模能力和特征学习能力而被应用于漏洞检测领域.基于深度学习的漏洞检测方法可以按照样本粒度和模型类型两个角度进行分类.

基于深度学习的漏洞检测方法根据待处理样本粒度可以分为函数级别和切片级别.在函数级别,Lin^[15]等人针对软件跨项目的特点,首次采用深度学习模型在函数级别进行检测漏洞.Feng^[16]等人设计了一种基于树的漏洞检测器,其首先用静态分析提取函数的抽象语法树,然后应用前序遍历搜索将树转换为序列,最后将该序列输入双向门控循环单元(Biphasic Gated Recurrent Unit,简称 BGRU)模型进行训练和测试.Wang^[17]等人提出了FUNDED漏洞检测模型,其通过提取代码的抽象语法树和程序控制依赖图并输入图神经网络完成漏洞检测.Wu^[18]等人将待测函数转换为图片后送入卷积神经网络进行漏洞检测.为降低样本噪声对检测结果的影响,研究人员优化样本粒度并提出了更细粒度的切片级漏洞检测方法.例如,Li^[19]等人提出的VulDeePecker首先针对待测程序中的敏感API提取程序切片,然后利用双向长短期记忆网络(Bidirectional Long Short-Term Memory,简称 BLSTM)训练漏洞检测器.Zou^[20]等人首次引入了代码注意力的概念,在程序切片构建时考虑程序控制依赖关系,以包含更加全面的漏洞语义和语法信息,并完成多类型的漏洞检测任务.Li^[11]等人提出的SySeVR对漏洞切片方法进一步完善,即在敏感API的基础上补充了数组使用、指针使用和整数使用作为漏洞切片的基准点,该方法使得漏洞切片数量大幅提升,更利于深度学习模型的训练.

基于深度学习的漏洞检测方法根据模型类型可以分为基于文本模型和基于图模型两种方式.基于文本的漏洞检测方法主要采用卷积神经网络模型和循环神经网络模型,其思想是将代码看作普通文本进行向量表征,并将向量输入深度学习模型进行检测.例如,Russell^[9]等人提出的TokenCNN首先提取代码令牌序列,然后使用卷积神经网络进行表征学习并将其输入到全连接层中进行分类,以得到漏洞检测结果.Duan^[21]等人提出的VulSniper通过提取代码属性图并转换为特征张量送入BLSTM模型完成漏洞检测.然而,代码不同于普通文本,

其具有更为丰富的结构和语义信息,因此将代码看作为文本处理的方式会在一定程度上损失代码信息.而代码表征图能够充分表征代码结构和语义^[22],因此研究人员提出了将代码表征图输入图神经网络模型进行漏洞检测.例如,Zhou^[12]等人采用通用的图神经网络模型进行漏洞检测,其模型包含的图卷积模块可以有效地从函数的图表征中学习函数级别的漏洞特征.Cheng^[23]等人提出了 DeepWukong,其利用图神经网络(Graph Neural Network,简称 GNN)实现了面向代码程序依赖图的子图的 C/C++函数过程间漏洞检测,但是过程间的漏洞检测难以应用于真实场景.原因在于真实软件的函数间调用关系的复杂性而难以控制调用链的长度,从而使得图过于复杂而难以有效检测.

函数级漏洞检测的优势在于函数能够覆盖相对完整的漏洞特征,但同时会引入较多与漏洞无关的噪声语句.尤其在真实世界的漏洞检测场景下,函数包含的代码行数过多,样本中存在过多的噪声会对深度学习模型造成干扰,使其无法学习到精确的漏洞特征而影响检测结果的准确性.因此,为了保证本系统对于真实漏洞检测任务的有效性,本文选择构建切片粒度的漏洞检测模型.此外,根据 Chakraborty^[24]等人的调研结果,图比文本能够更好地表达代码的结构和语义信息,因此在漏洞检测领域,基于图模型的检测方式效果明显优于基于文本模型的检测方式.因此本文采用代码切片(即程序依赖图子图)训练图神经网络模型,以实现真实漏洞的准确检测.

1.2 基于深度学习漏洞解释相关工作

目前图神经网络的可解释性工作可以分为基于扰动的解释方法、基于特征的解释方法、基于模型的解释方法和基于分解的解释方法四种类型.其中,基于模型的解释方法是面向模型而设计的解释方法,其他类型的解释方法均面向单个实例提供检测结果的解释.在漏洞解释任务中,研究人员更倾向于了解每一个漏洞的漏洞原理而非理解深度学习模型本身的运作机制.因此,面向实例的解释方法更适用于漏洞检测的解释任务.

基于模型的解释方法侧重于对模型整体进行解释,即对深度学习模型本身提供解释方案,以对模型的工作方式提供高层次的见解和一般性的理解,从而方便研究人员有针对性地对模型进行优化.例如,Yuan^[25]等人提出训练一个图生成器,其生成的图模式可以对深度图模型进行解释,即根据已训练图模型的反馈,使用策略梯度对图生成器进行训练.基于分解的解释方法通过将原始模型预测分解为若干项来衡量输入特征的重要性,这些项被视为相应输入特征的重要性分数.例如,Schnake^[26]等人提出的 GNN-LRP 解释器研究了图神经网络中不同相关步长(记录了层与层之间消息传递的路径)的重要性,通过将分数分配给不同的步长,并从其对应节点上获得分数进行解释,但是其计算复杂度过高,导致实际使用受到限制.基于特征的解释方法依靠梯度或隐藏图特征表示不同输入特征的重要性,其中梯度或特征值越高表示其重要性越高,被广泛运用于图像和文本任务.例如,Selvaraju^[27]等人提出的 Grad-Cam 解释器是一种基于图像分类器的解释模型,该方法利用梯度表示不同输入特征的重要程度.该方法可以拓展应用到图模型以衡量不同节点的重要性,其关键思想是结合隐藏特征图和梯度来指示节点重要性.但是基于特征的解释方法在面向图神经网络解释时对图神经网络的结构有特殊要求,而不具有普适性.相反,基于扰动的解释方法由于不受限于模型的网络结构而被广泛用于解释图神经网络模型.例如,Ying^[13]等人提出的 GNNExplainer 解释器可以从图的结构和节点角度来对任意图神经网络生成解释,探寻与预测结果相关性最大的子图结构以实现对预测的解释.Luo^[14]等人提出的 PGExplainer 解释器主要思想和 GNNExplainer 相似,但其在对图神经网络模型解释过程前需要进行解释器的参数训练,以实现同时对多个实例从模型的全局视角进行解释,并输出对模型预测结果决定关键作用的子图结构.Yuan^[28]等人提出的 SubgraphX 模型旨在解释图神经网络模型在预测过程中重要子图的作用,具体来说,其将预训练的待解释模型和图实例数据,结合蒙特卡罗树搜索法输出图实例中的重要子图.

目前基于深度学习的漏洞解释工作尚处于起步探索阶段,相关工作较少.Li^[29]等人提出的 IVDetect 工具利用图神经网络解释器 GNNExplainer 实现函数级漏洞检测结果解释.不同于已有的 IVDetect 的解释工作,本文将使用并改进 GNNExplainer 对切片级的漏洞检测结果进行解释,一方面切片含有较少的噪声信息,能够获得更好的解释效果.另一方面被解释对象为切片(程序依赖图子图),因其包含的节点数相较完整的程序依赖图较少,以便提高解释效率.

2 基础知识

本节将对本文涉及理论的基础知识进行介绍,主要包括静态程序分析和图神经网络的相关概念。

2.1 程序静态分析相关概念

程序静态分析方法主要基于代码抽象语法树、控制流图、程序依赖图等代码表征图,相关概念介绍如下:

- 抽象语法树(Abstract Syntax Tree,简称 AST):AST 是一种树形结构的源代码抽象表示,其从根节点开始将代码依次分解为代码块、语句、声明、表达式等.AST 主要用于描述代码的语法结构,是构成其他代码表征图的基础.
- 控制流图(Control Flow Graph,简称 CFG):CFG 是一种有向图,其节点表示函数语句,边则表示相邻语句间运行的先后关系,描述了程序在执行过程中可能经过的所有运行路径.
- 数据依赖:设 A 和 B 是两个代码语句,若 A 语句中计算得到的变量在 B 语句中使用,则称 B 语句数据依赖于 A 语句.
- 控制依赖:设 A 和 B 是两个代码语句,若 B 语句能否执行取决于 A 语句执行的结果,则称 B 语句控制依赖于 A 语句.
- 程序依赖图(Program Dependence Graph,简称 PDG):PDG 是一种带有标记的有向多重图,其基于 AST 的节点构建,节点表示代码语句,边表示相邻节点存在数据依赖或控制依赖关系.
- 代码属性图(Code Property Graph,简称 CPG):CPG 将 AST、CFG 和 PDG 整合至一种数据结构,其包含源代码中的所有语法特征和语义特征.
- 代码切片:代码切片由彼此之间具有数据依赖或控制依赖关系的语句构成^[19].

2.2 图神经网络

传统的卷积神经网络和循环神经网络在提取欧氏空间数据(如自然语言)的特征方面取得了巨大的成功,但其处理非欧式空间的图数据(如社交网络)时表现不佳.而图神经网络在处理非欧式空间的图数据的表达能力上相比传统深度学习模型具有明显优势,因此图神经网络在推荐系统、社交网络分析和交通预测等领域被广泛应用.图神经网络模型的工作原理为:首先构建相邻节点之间需要传递的消息;其次收集邻居节点相关的消息;最后更新节点表示.其中信息传递工作主要是通过取出每个节点及其邻居节点的信息,聚合函数聚合以上所有信息,最后更新函数对节点信息进行更新,从而捕获节点间相互依赖的关系,更好地学习图结构数据的相应特征.

图神经网络模型根据具体解决的任务级别可以分为如下三种:

- 图级别任务:图级别任务的目标是预测整个图的属性.例如,本文漏洞检测即为一个图级别任务,其目的是判断输入的代码表征图是否存在漏洞.
- 节点级任务:节点级任务的目标是预测图中每个节点的类型.例如,社交关系图中判断某个节点的所属阵营.
- 边级任务:边级任务的目标是预测每条边的属性.例如,在目标检测的语义分割任务中不仅需要识别每个目标的类型,还需要预测各个目标之间的关系.

目前图神经网络根据模型类型主要分为以下三种^[30]:

- 图卷积网络(Graph Convolution Networks,简称 GCN):GCN 将卷积运算推广到图数据,其核心思想在于学习一个函数映射,通过该映射可以聚合图中节点和其邻居节点特征从而生成该节点的新表示.卷积运算又可以具体分为基于谱的方法和基于空间的方法.其中,基于谱的方法在训练时需要将全图加载入内存,因此在处理复杂图时性能较差;基于空间的图卷积方法中心节点、感受域和聚合函数不确定,导致相互制约依赖,实际使用效果不佳.
- 图注意力网络(Graph Attention Networks,简称 GAT):GAT 在 GNN 上引入注意力机制,在聚合过程使用注意力关注对任务有影响的邻居节点信息,进而为节点分配权重,使得 GNN 重点关注与任务相关的节点和边以提升模型效果.但是其增加了计算邻居之间的注意权重的时间开销和内存消耗,同时注意

力机制的实际作用和网络初始化相关。

- 图循环网络(Graph Recurrent Network,简称 GRN):GRN 将图转换为序列,使用长短期记忆网络(Long Short-Term Memory,简称 LSTM)或门控循环单元(Gated Recurrent Unit,简称 GRU)等循环神经网络作为架构训练.相比于其他 GNN 模型,GRN 在信息传递过程使用门控机制可以提高在整个图上的长期信息传播能力以及模型表征能力.本文使用的门控图神经网络模型(Gated Graph Neural Network,简称 GGNN)即为 GRN 类型下的一种网络结构。

3 基于图神经网络的切片级漏洞检测及解释方法

为了解决现有漏洞检测模型在真实漏洞数据集上效果差且无法给出检测结果解释的问题,本文提出并实现了一种基于图神经网络的切片级漏洞检测及解释方法.如图 1 所示为本方法的整体框架图,本方法由数据预处理模块、漏洞检测模块和漏洞解释模块三个模块组成.系统的输入是待测软件的源代码,输出是目标程序切片中是否存在漏洞以及具体的漏洞语句。

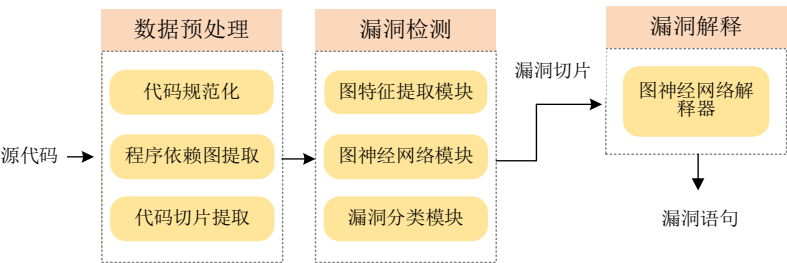


图 1 基于图神经网络的切片级漏洞检测及解释方法整体框架

3.1 数据预处理

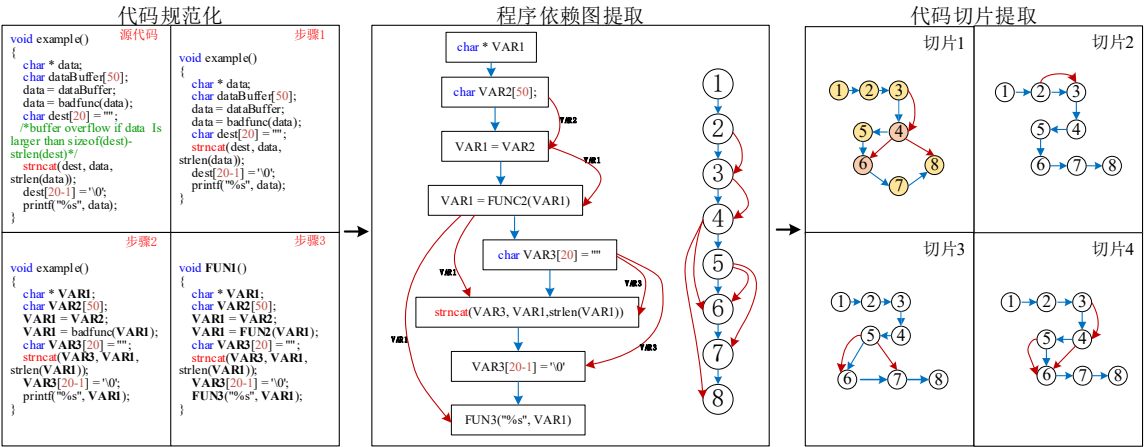


图 2 数据预处理样例

3.1.1 代码规范化

程序源代码含有丰富的语义信息,便于表示漏洞特征.因此,本系统面向源代码进行漏洞检测.但同时,源代码中代码注释、复杂的变量和函数命名等与代码语义无关的信息会干扰深度学习模型的训练和预测.为了降低代码中无关信息的干扰,本系统首先对源代码进行代码规范化处理.考虑到真实代码较为复杂难以展示,因此图 2 示例源代码来自于人工数据集 SARD.如图 2 代码规范化部分所示,代码规范化共包含具体的三个步骤:

- 步骤 1:去除代码中的注释信息(例如/**/);
- 步骤 2:将用户自定义的变量名映射为统一的变量名(例如 VAR1);

- 步骤 3:将用户自定义函数名映射为统一函数名(例如 FUNC1).

3.1.2 程序依赖图提取

代码表征图是一种有效的代码表征方式,可以直观地体现代码的语义、语法和结构信息.其中,程序依赖图是一种由数据依赖关系和控制依赖关系连接抽象语法树节点的数据结构.其不同于抽象语法树、控制流图等仅包含语法或控制流等单一代码信息的图表征方式,且不同于代码属性图包含过多未精简的代码信息而影响模型的检测效率.程序依赖图是一种高效且较为全面的代码图表征方式,且能够有效表征多种类型漏洞^[24].因此,本文选择程序依赖图作为源代码的表征方式.

本文利用开源的静态代码解析工具 Joern^[31]对源代码进行程序依赖图提取.选择该工具的原因在于其无需预编译即可完成代码解析和程序依赖图生成工作,这意味着其适用于任意粒度的待测代码,极大程度降低了用户的使用成本.如图 2 程序依赖图提取部分所示,其中黑色框表示代码行对应的节点,蓝色的边表示控制依赖关系,红色的边表示数据依赖关系,红色边旁的变量表示数据依赖的变量名称.为了更简洁地展示后续工作,本文用编号替换程序依赖图节点中的代码,如图 2 程序依赖图提取的右半部分所示.

3.1.3 代码切片提取

一方面函数级样本中包含大量漏洞无关的噪声语句会干扰模型学习漏洞特征,进而影响检测效果;另一方面真实软件的代码语句较为复杂,因此其程序依赖图中包含节点和边的数量过于庞大,导致训练和解释过程需要大量时间和内存开销.为避免上述问题,本文的漏洞检测和解释过程均面向切片级别,以去除漏洞无关信息从而提高检测和解释效果并降低开销.具体而言,本文采用 Li^[19]等提出的漏洞切片概念作为切片生成的指导思想.其工作总结出软件漏洞主要由指针、数组、表达式运算、敏感 API 函数所在位置引入,并称之为漏洞关注点.如常见的数组越界、整数溢出、空指针、API 函数错误使用等类型漏洞均由以上四类漏洞关注点导致.在程序代码中,与漏洞关注点存在数据依赖或者控制依赖关系的语句集合构成一个可能存在漏洞的程序切片.反之,其他语句被视为会干扰模型训练的漏洞无关语句.不同于 Li^[19]等人构造的文本切片,本文提取程序依赖图的子图作为切片,即选取上述四类漏洞关注点作为程序切片的基准点后,保留与其存在数据及控制依赖关系的节点和边以生成程序依赖图子图.具体来说,生成代码切片可以分为三个步骤.

(1) 漏洞关注点选取

通过遍历程序依赖图节点选取符合四类漏洞关注点的代码元素,并记录该节点为切片基准点.具体来说,通过在标识符声明节点中匹配“[”字符来确定数组元素;通过在标识符声明节点中匹配“*”字符来确定指针元素;通过正则表达式规则来匹配表达式运算节点;通过 Li^[19]等人提供的敏感 API 列表进行敏感 API 节点匹配.如图 2 程序依赖图提取部分所示,1 号节点存在指针元素 *VAR1*,2 号和 5 号节点分别存在数组元素 *VAR2* 和 *VAR3*,6 号存在敏感 API 元素 *strncat*.因此,以上节点被选定为程序切片的基准点.

(2) 程序切片生成

从切片基准点出发分别执行前向及后向切片以生成程序切片.具体而言,在程序依赖图上以漏洞关注点所在节点为起点,分别追溯前向、后向的控制依赖边和数据依赖边,并记录涉及到的节点和边,直到不再出现新增节点和边为止.根据上述步骤得到的程序依赖图子图即为一个程序切片.因为其只包含与漏洞关注点具有依赖关系的节点和边,因此在保留源代码结构信息的同时,排除了图中漏洞无关的信息.图 2 代码切片提取部分展示了从该源代码的程序依赖图中提取的部分切片,切片 1-4 分别为以 *VAR1*、*VAR2*、*VAR3* 和 *strncat* 为切片基准点获得的切片.

(3) 程序切片标注

本文依赖漏洞的补丁信息对切片进行标注,即包含该漏洞补丁中删减行的切片被标注为有漏洞切片,反之被认为是无漏洞切片.根据其提供的漏洞信息,漏洞补丁的删减行为源代码的第 11 行,即对应程序依赖图中的 6 号节点.根据标注规则,由于图 2 的切片 1-4 均包含 6 号节点,因此均被标注为有漏洞切片.

3.2 漏洞检测

本文采用图神经网络模型进行漏洞检测.图神经网络模型自动学习输入图的漏洞特征并对其特征向量进

行分类以检测代码是否含有漏洞,其主要分为图特征提取模块、图神经网络模块和漏洞分类模块,如图 3 所示.

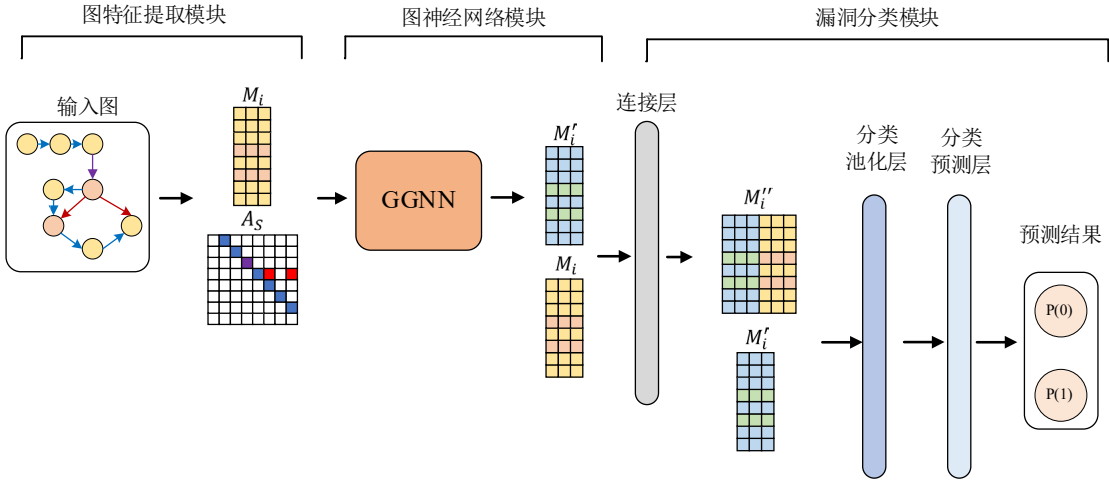


图 3 漏洞检测模型

3.2.1 图特征提取模块

由于本文抽取的代码切片是抽象图格式的文件(程序依赖图子图),无法直接输入图神经网络漏洞检测模型,因此需要提取图的关键特征以获得图特征向量.程序依赖图子图包含两个维度的特征,即节点内的代码特征(下文简称为节点特征)和图结构特征.

针对节点特征,本文对节点内代码进行嵌入表征.具体地,将每个节点中的代码视为一个句子,并分解句子为令牌列表后,嵌入得到一个固定长度的向量.本文使用 word2vec 模型^[32]实现节点嵌入,其采用一种分布式表示的思想,将令牌映射为整数再转换为一个固定长度的向量.该方法在训练嵌入速度上相对于传统嵌入工具速度更快且具有泛用性,因此在文本挖掘领域中被广泛使用.具体来说,本文将所有切片中的令牌列表训练生成一个预训练的 word2vec 模型,再使用该预训练模型对全部节点进行向量嵌入.即将预处理得到的切片输入预训练的 word2vec 模型,输出其形状为 $m \times n$ 的特征矩阵 M_i ,其中 m 表示切片的节点数, n 表示嵌入向量的维度,考虑到实际检测效果和性能,本文设置 n 为 100 维.如图 3 的输入图所示,图中共有 8 个节点,因此其节点特征向量 M_i 的维度为 8×100 .

针对图结构特征,本文对图中的边关系进行嵌入表征.每一条边可以被视为一个三元组(起始节点,终止节点,边类型).其中,起始节点和终止节点均能够从程序依赖图中直接获取,边类型分为数据依赖边和控制依赖边.如图 3 的输入图所示,其具有 10 条边,其中包括 3 条数据依赖边和 7 条控制依赖边(红色边表示数据依赖边,蓝色边表示控制依赖边,紫色边表示既有数据依赖边又有控制依赖边),输出矩阵 A_s 即表示图结构特征矩阵.

3.2.2 图神经网络模块

由于本文将源代码转化为具有数据依赖关系和控制依赖关系的图结构数据,而图神经网络有助于进一步聚合传递更新以便更好地捕捉图的结构和语义信息,因此本文采用图神经网络对图样本进行进一步特征嵌入.目前,深度学习领域提出了多种基于聚合邻域信息的图神经网络模型^{[33][34][35]}.其中 GGNN 增强了网络的长期记忆能力,比传统的图神经网络模型更加深入^[35],更适用于处理既有语义也有图结构的数据,因此本文选用 GGNN 为最终的图神经网络模型.

GGNN 的原理在于聚合节点和其邻居节点的信息,将聚合后的节点和当前节点一起送入 GRU 以得到下一时刻的当前节点,按此过程重复,经过若干个时间步迭代会生成所有节点的最终节点特征.如图 3 图神经网络模块所示,输入图特征 $g_i(M_i, A_s)$ 后,GGNN 通过嵌入每个节点及其邻域将其转换为一个长 $m \times n'$ 的切片特征矩阵 M'_i ,其中 n' 为设置的最终切片特征大小,本文设置为 200 维,则图中的特征矩阵 M'_i 维度为 8×200 .

具体来说,对于图中的每一个节点 v_u ,初始化节点向量 $h_u^{(1)} = [m_u^T, 0]^T$,即复制 v_u 节点的特征向量并额外填充 0. 设 T 为邻域聚合的总时间步数,为了得到整个图传播的信息,对于每个时间步长 $t \leq T$,所有节点根据所依赖的边进行通信,即:

$$a_u^{(t)} = A_u^T (W_u [h_1^{(t-1)}, \dots, h_m^{(t-1)}] + b) \quad (1)$$

其中 W_u 表示可训练的参数, b 为偏差, A_u^T 表示节点 v_u 对应于 A_S 的邻接矩阵, $a_u^{(t)}$ 为当前节点和相邻节点间通过边的相互作用的结果.最后通过聚合函数 AGG 聚合节点 v_u 的信息并与与上一个时间步合并,从而获取该节点的新状态,即:

$$h_u^{(t+1)} = GRU(h_u^{(t)}, AGG(\{a_u^{(t)}\})) \quad (2)$$

3.2.3 漏洞分类模块

漏洞分类模块的核心思想在于选择与漏洞特征相关的特征集合,以完成图级别的漏洞样本分类任务.先前的工作^[36]提出在图卷积层后使用一个分类池化层(SortPooling)从而实现图卷积层输出特征的排序,即可输入传统的神经网络对其进行训练,以提取切片向量嵌入中的有用特征.因此在本文中,节点特征通过 GGNN 层学习,进而使用一维卷积和全连接层学习与图分类任务相关的特征,以便更有效地进行分类.具体来说,本文定义分类池化层 $\tau(M)$ 表述如下式,其中 $Conv$ 表示卷积层, BN 表示BatchNorm层, $Relu$ 表示激活函数, $MaxPool$ 表示最大池化层, M 表示一个特征矩阵.

$$\tau(M) = MaxPool(Relu(BN(Conv(M)))) \quad (3)$$

如图 3 漏洞分类模块所示,本文将切片的节点特征矩阵 M_i 和相应切片特征矩阵 M_i' 连接成一个新的矩阵 M_i'' ,分别对 M_i' 和 M_i'' 执行 τ 分类池化操作得到输出 Y_1 和 Y_2 .接着将 Y_1 和 Y_2 分别送入输出维度为 2 的全连接层中,最后对两个输出值相乘求出平均后进行Sigmoid分类从而得到预测,即图 3 中的分类预测层,如下式所示:

$$P = Sigmoid(Avg(Liner(Y_1) \cdot Liner(Y_2))) \quad (4)$$

其中 Avg 表示平均操作, $Liner$ 表示全连接层, P 是输出的二分类结果,由两个维度组成.其中第一个维度表示结果无漏洞的概率,第二个维度表示结果有漏洞的概率.最后模型采取两者中较大的概率作为漏洞分类的最终结果输出.表 1 给出本文在训练模型中使用到的具体参数信息.本文模型使用交叉熵损失函数 CrossEntropyLoss 纠正错误的分类,同时使用学习率为 0.0001 且权重衰退为 0.001 的 Adam^[37]优化算法,以训练 3.2.2 节式(1)中图神经网络模块的参数 W_u 和 b .获得训练好的模型后,再使用其判断新的代码切片是否存在漏洞.

表 1 模型训练相关参数设置

参数	设置值
损失函数	Cross Entropy Loss
优化算法	Adam
学习率	0.0001
权重衰退	0.001
批量大小	8
训练轮数	500

3.3 漏洞解释

本文在漏洞解释部分使用并改进了 GNNExplainer 图神经网络解释器^[13],从而对漏洞检测模型输出的检测结果提供细粒度的解释,即具体的漏洞代码行,以帮助研究人员完成漏洞成因分析及进一步的漏洞修复.

GNNExplainer 是一种针对 GNN 模型设计的基于扰动的解释方法,其不依赖具体的图神经网络模型,具有很强的泛用性.其根本思想是完成一个最大化互信息的优化任务,对输入的实例图生成相应的边掩码进行更新训练,最终根据预测结果的变化输出图中的关键的信息作为解释结果.在掩码更新训练过程中,每次训练可以得到一个包含关键信息的新图,将新图输入到已训练的 GNN 模型中,根据检测结果评估边掩码的重要程度并继续

重复训练以降低损失,并以找到一个损失最小的关键信息为训练的终点.最后,解释器通过边掩码得到重要边,并根据重要边抽取重要子图作为模型对于该实例的解释.

具体来说,GNNEExplainer 目标是从原图 $G_W(M_i, A_S)$ 中学习一个边掩码 E_M 来获取关键子图 G_S 作为解释结果,它的出发点是最大化子图 G_S 和原图 G_W 互信息 MI ,即保证子图尽可能涵盖原图的重要信息,如下式所示:

$$\text{Max}(H(Y|G=G_S)) = -E_{Y|G_S} \log P(Y|G=G_S) \quad (5)$$

其中 Y 表示漏洞检测模型预测的结果,显然,熵 $H(Y)$ 对于完成预训练的模型而言是一个常数,故上式可等价化为最小化条件熵 $H(Y|G=G_S)$,根据条件熵定义,有下式:

$$\text{Min}(H(Y|G=G_S)) = -E_{Y|G_S} \log P(Y|G=G_S) \quad (6)$$

即表示在解释结果为 G_S 时,预测结果 Y 尚存在多大程度的不确定性,然而上式无法直接优化,GNNEExplainer 在特殊情况下将 $G_S \sim g$ 视为为一个随机的图变量,因此上式变为:

$$\text{Min}_g(E_{G_S \sim g} H(Y|G=G_S)) \quad (7)$$

GNNEExplainer 进一步利用 Jensen 不等式以及凸性假设,将上式进一步转换为:

$$\text{Min}_g(H(Y|G=E_g[G_S])) \quad (8)$$

其中 $E_g[G_S]$ 可以被边掩码所代替,因此解释器最后训练的目标实际为边掩码.训练完成后输出的程序依赖图关键子图 G_S 作为解释结果.从上述描述可知,GNNEExplainer 通过边掩码进行重要边的排序进而生成重要子图.然而,在漏洞检测任务中,研究人员更倾向于知道是哪些代码行(即图节点)存在漏洞,而非程序依赖图中的哪些数据控制依赖边存在漏洞,因此,事实上并不需要额外生成重要子图的操作.为了解决上述问题,本文为 GNNEExplainer 设计并增加了一个节点排序算法使其能够按节点重要性排序并输出图中的关键节点,使得解释器更适用于漏洞检测任务,可以明显提高解释准确率并降低原方法的时间开销(减少了解释器生成子图和无关掩码的计算过程).具体方法是:首先将边掩码按重要程度排序;其次为边掩码对应的边节点累加权重,其权重为边掩码的大小;最后对每个节点按权重大小进行排序,即可得到按重要程度排序的节点序列.具体来说,对于图中的节点 i 其权重大小用 $N_M^{(i)}$ 表示(初始所有节点权重均为 0).用 E_{I_i} 表示含有节点 i 的边号,则 $E_M^{(E_{I_i})}$ 表示边号为 E_{I_i} 的边掩码大小,那么节点 i 的权重更新公式可以表示为:

$$N_M^{(i)} = N_M^{(i)} + E_M^{(E_{I_i})} \quad (9)$$

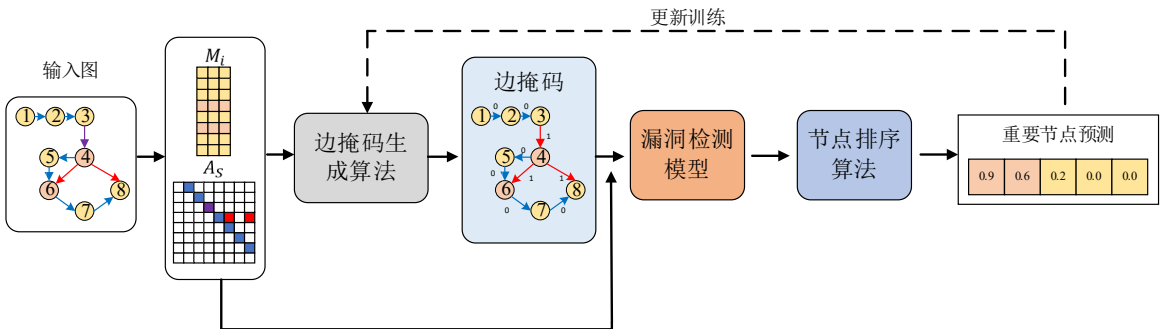


图4 漏洞解释模型

本文的漏洞解释模型如图4所示,解释器的输入为漏洞检测器检测出的漏洞切片,解释器通过扰动输入切片的边掩码以生成新图,并观察漏洞检测模型对新图的检测结果,从而获得该边掩码的重要分数,最后将边掩码的重要分数映射成为节点的重要分数并按照分数排序输出.其中,分数越高的节点则被认为更有可能是一个漏

洞代码行,即对漏洞检测结果的解释.

4 实验分析

本节主要对具体实验环节进行阐述.首先对实验准备进行说明,包括实验数据集、实验具体实施方法以及实验评估指标.其次引入三个研究问题以验证本方法的实验效果,最后对实验结果进行分析.

4.1 实验准备

(1) 实验数据集

本文在广泛使用的漏洞数据集 Big-Vul^[38]的子集上进行全部实验.一方面,Big-Vul 数据集是一个高质量的真实漏洞数据集,其涵盖了从 2002 年到 2019 年间 348 个真实开源项目的全部 CVE(Common Vulnerabilities and Exposures^[39])条目,共有 11834 个漏洞函数和 253096 个无漏洞函数;另一方面,它涵盖了修补前后的完整代码以及具体的漏洞修补位置,以便于进行切片标注.因此,本文选择 Big-Vul 作为实验数据集.

需要说明的是,Big-Vul 数据集爬取了含有具体漏洞补丁的源文件中的所有函数,因此补丁前的函数被标注为有漏洞函数,补丁后的函数被标注为无漏洞函数,对于没有任何修改且仍然保留的函数则标注为无漏洞.在实际预处理过程中,切片标注依赖于漏洞补丁中的修改行信息,导致部分数据由于无法标注而需要过滤.因此,本文最终使用了 Big-Vul 数据集中的 4030 个有漏洞函数和 4716 个无漏洞函数,并从中获得了 16260 个有漏洞的切片和 16301 个无漏洞的切片.

(2) 具体实施

本文实验设备内存大小 128GB,配有 16 核的英特尔 Xeon 处理器和英伟达 Quadro GTX 5000 系列 16GB 显卡.系统各个阶段分别采用 Joern^[31],word2vec^[32]和 Pytorch^[40]等工具实现.在漏洞检测部分,本文根据切片样本的数量将数据集按照 8:1:1 随机划分为训练集、验证集和测试集.在漏洞解释部分,由于只有漏洞切片需要被提供漏洞解释,因此本文随机选择了测试集中预测正确的 1399 个有漏洞切片输入解释器中进行解释.

(3) 评估指标

为评估漏洞检测模型的有效性,本文采用常见的评价模型准确程度的指标对漏洞检测模型进行评估^{[11][17]}.

- 真正例(True Positive, TP):正确预测为有漏洞的样本数量.
- 真负例(True Negative, TN):正确预测为无漏洞的样本数量.
- 假正例(False Positive, FP):错误预测为有漏洞的样本数量.
- 假负例(False Negative, FN):错误预测为无漏洞的样本数量.
- 准确率(Accuracy)=(TP+TN)/(TP+TN+FP+FN)
- 召回率(Recall)=TP/(TP+FN)
- 精确率(Precision)= TP/(TP+FP)
- $F1 = 2 * Precision * Recall / (Precision + Recall)$

为评估解释模型的有效性,本文借鉴 IVDetect^[29]的准确率评价指标对漏洞解释模型进行评估.具体而言,当输出的节点中包含漏洞补丁修改行,则认为解释正确;反之则认为解释错误.解释结果的准确率为解释正确的样本数量在全部样本数量的比值.

4.2 实验结果与分析

本文所设计的实验用于回答以下三个研究问题:

- RQ1:本方法在源代码漏洞检测方面的表现如何?
- RQ2:本方法对于漏洞检测结果的解释效果如何?
- RQ3:本方法是否可以检测并解释真实软件漏洞?

RQ1: 本方法在源代码漏洞检测方面的表现如何?

为回答该研究问题,本文将本方法同几种先进的漏洞检测工具进行比较.对比方法包括 3 个基于规则的漏

洞检测工具(Checkmark^[6],FlawFinder^[7]和 RATS^[8])和 4 个基于深度学习的漏洞检测方法(TokednCNN^[9],StatementLSTM^[10],SySeVR^[11],Devign^[12]).实验结果如表 2 所示.

基于规则的方法是将专家预定义的漏洞模式与源码匹配的漏洞检测方法.如表 2 精确率、召回率和 F1 分数所示,商用静态漏洞检测器(Checkmark)和静态分析系统(FlawFinde 和 RATS)的检测效果均不理想.例如,Checkmark 的召回率只有 31.9%,说明其对于真实世界的漏洞仅能成功检测 31.9%.这些工具依赖人类专家制定的规则和模式,通过词法分析对目标程序的代码进行模式匹配以检测漏洞.但是,专家定义的漏洞模式难以覆盖复杂的真实环境,导致其对于真实漏洞的检测效果不佳.不同于基于规则的方法,本方法使用神经网络模型自动化挖掘并学习训练样本中的漏洞模式,在节省人工开销的同时达到了更低的漏报率和误报率.

基于令牌的方法是将整个源代码片段视为一段自然语言文本,并按照文本的分词方法将代码文本拆解为令牌序列后进行模型训练.其中,TokenCNN 是此类方法中的代表性工作.其首先通过词法分析将源代码转换为令牌序列,然后将其嵌入为向量,最后将向量表征输入卷积神经网络模型进行漏洞预测.因为 TokenCNN 将源代码视为纯文本,且直接将源代码按照文本处理的方式进行模型的训练和预测,缺乏对源代码语义和结构信息的考虑而损失了代码中的大量信息,因此导致检测效果不佳.

表 2 与其他漏洞检测工具对比

类别	工具或方法	精确率(%)	召回率(%)	F1 分数(%)
基于规则方法	Checkmark	40.4	31.9	35.7
基于规则方法	FlawFinder	39.6	33.4	36.2
基于规则方法	RATS	41.6	39.1	40.3
基于令牌方法	TokenCNN	43.2	54.3	48.1
基于语句方法	StatementLSTM	48.9	57.5	52.9
基于函数级方法	Devign	51.8	54.7	53.2
基于切片级方法	SySeVR	49.7	53.8	51.7
基于切片级方法	本方法	70.5	80.3	75.1

基于语句的方法类似于基于令牌的方法,其直接将每行代码进行嵌入代替了基于令牌方法先分词后嵌入的解决方案.例如,StatementLSTM 将每行代码视为一个自然语言句子,将其嵌入固定长度的向量后输入 LSTM 模型训练漏洞检测器.其相较于基于令牌的方式避免了因分词带来的语义损失,所以检测效果相对更好.但是,基于语句的方法同样是在文本角度处理代码,仍没有在根源上解决代码语法语义损失的问题.

基于函数级的方法是将整个函数当作处理对象进行模型训练.例如,Devign 首先使用复杂程序分析提取一个完整函数的代码属性图,接着使用通用的图神经网络检测漏洞.通过完整函数的代码属性图表征使得训练样本能够包含关于代码的全面语义和句法信息.然而函数的代码属性图在一方面包含大量漏洞无关的节点和边干扰模型训练;另一方面代码属性图包含抽象语法树、程序控制流图、程序依赖图等多种数据结构,复杂的图结构使得模型难以学习到精确的漏洞特征,导致模型效果欠佳.

基于切片的方法是指先将待测程序通过程序切片提取出程序关键语义,进而对关键语义相关的语句或节点集合进行训练和测试.其中,SySeVR 是最具代表性的切片级方法之一.其对目标程序执行切片以生成代码切片,然后将其嵌入为相应的向量表征.与其他级别的方法不同,切片级方法相当于对样本进行了预先提纯,即删除样本中的噪声语句.然而 SySeVR 将代码切片视为一段自然语言文本,因此代码行之间的依赖关系不会在其切片中捕获.相反,本方法生成的切片是一个程序依赖图子图,其保留了漏洞相关的语句间完整的语义和代码结构,因此达到了更佳检测结果.

综上,本方法采用源代码切片(程序依赖图子图)训练图神经网络漏洞检测器,实验结果表明其针对真实的漏洞检测可以取得 75.1%的 F1 分数.相较于其他的漏洞检测工具,F1 分数提高 41.2%-110.4%,具有明显优势.

RQ2: 本方法对于漏洞检测结果的解释效果如何?

为回答该研究问题,本文将本方法的解释效果同 2 种先进的图神经网络通用解释器(GNNExplainer^[13]和 PGExplainer^[14])进行比较.表 3 展示了本文改进 GNNExplainer 方法同原 GNNExplainer 和 PGExplainer 的解释

效果对比.其中原 GE 表示 GNNExplainer,原 PE 表示 PGExplainer,本方法则指本文改进的 GNNExplainer.表 3 的第一行分别表示解释器输出待解释图中 2%-20%的节点作为解释结果的准确率.

根据表 3 的实验结果可以看出本文改进后的 GNNExplainer 方法相比于原 GNNExplainer 和 PGExplainer 在实际漏洞解释中准确率分别提升 3.8%-18.5%和 13.4%-29.0%,其中前 10%节点准确率高达 73.6%,前 20%节点准确率高达 91.1%.其原因在于原 GNNExplainer 是以关键边的角度构建关键子图后作为解释结果输出,然而其作为一个通用的图神经网络并非为漏洞样本设计,因此在根据重要边构造出最小联通子图的过程会对漏洞信息造成一定损失.而本方法直接将边掩码的重要程度分数映射给节点,减少了中间步骤,因此较原先的 GNNExplainer 解释效果更佳,且更贴近研究人员分析漏洞成因时倾向于获得哪些代码行是漏洞代码的研究需求.不同于 GNNExplainer 对于单个样本直接进行解释,PGExplainer 在对图神经网络模型解释过程前需要大量样本进行解释器的参数训练,以实现同时对多个实例从模型的全局视角进行解释,且其解释的准确性依赖于解释器训练效果.而本文实验数据均为真实漏洞数据,数据具有高度复杂性,PGExplainer 难以应对大量复杂实例的解释器训练,因此导致解释效果欠佳.相比之下 GNNExplainer 对单个实例进行单独的训练和解释,使得单个实例损失较小,解释准确率更高.

表 3 与其他图神经网络解释器准确率对比

方法	Acc _{2%}	Acc _{4%}	Acc _{6%}	Acc _{8%}	Acc _{10%}	Acc _{12%}	Acc _{14%}	Acc _{16%}	Acc _{18%}	Acc _{20%}
原 GE	25.4%	40.0%	50.3%	58.8%	67.6%	73.8%	77.5%	81.6%	85.1%	87.1%
原 PE	26.2%	41.3%	48.9%	54.5%	58.9%	62.8%	65.6%	67.7%	69.4%	70.6%
本方法	29.7%	47.4%	57.5%	66.0%	73.6%	78.4%	82.9%	85.6%	88.3%	91.1%

根据实验结果可知,随着节点数目的增加,解释语句覆盖漏洞行的准确率越高.但同时,解释结果输出的节点数量越大意味着研究人员需要分析更多的潜在漏洞行,因此带来较大的人工负担.根据表 3 的结果所示,当解释输出节点数大于 10%时,解释的准确率开始增长缓慢,同时考虑到为避免带来过多的无关解释语句,本文最终以前 10%节点作为最后解释器的节点输出数量.

表 4 与其他图神经网络解释器时间开销对比

解释方法	训练时间	解释时间	总时间	平均时间
原 GE	—	136764.1s	136764.1s	97.8s
原 PE	92401.2s	468.7s	92869.9s	66.4s
本方法	—	78605.5s	78605.5s	56.2s

本文对本方法和对比解释器的时间开销进行了实验统计,如表 4 所示.实验随机选取了 1399 个待解释样本进行解释,训练轮数默认均取 100 轮.对比原 GNNExplainer,本方法减少了解释器生成子图和无关掩码的计算过程,因此相比于原方法时间开销明显降低.由于 PGExplainer 需要事先训练解释器,训练完成的解释器再对单个实例依次进行解释,因此其时间开销主要用于解释器训练.而本方法针对每个实例同时进行训练和解释,因此总时间开销明显低于 PGExplainer,解释效率更高.

综上,本文通过改进 GNNExplainer 能够精确且快速地定位漏洞语句,一方面有助于提高漏洞检测模型的可信性,另一方面便于研究人员分析漏洞成因并对漏洞语句有针对性修补.相比于已有的图神经网络通用的解释器,本方法在准确度和时间开销方面均具有明显优势.

RQ3: 本方法是否可以检测并解释真实软件漏洞?

为回答该研究问题,本文对常见的开源软件进行漏洞检测和解释以检验本方法在现实世界漏洞发掘方面的实用性.具体来说,本实验将 4 个目标程序的开源 C/C++代码作为检测目标,待测软件包含 libav、xen、openssl、和 httpd.表 5 给出了目标软件的细节描述,其中共计目标函数 286039 个,并从中获取切片 3312657 个.实验具体过程是首先将预处理后的切片输入训练好的漏洞检测模型进行漏洞检测,接着漏洞切片被输入解释器以获得重要程度前 10%的关键节点对应的漏洞行.由于本方法作为静态检测手段无法直接给出漏洞在程序运行中的

触发信息,即面向NVD申请新的漏洞号(CVE-ID)时需要提交的漏洞触发的程序运行依据,因此实验通过人工分析以验证本方法对于真实软件的检测和解释结果是否正确.具体地,人工分析的过程采用了两种验证方法,首先,通过与 National Vulnerability Database^[41](NVD)已有漏洞的漏洞模式对比,从而确定本方法是否检测出真实漏洞,并核实漏洞解释结果的准确性.其次,通过人工分析漏洞成因的方式进行验证.

表 5 本文选择的开源软件细节

开源产品名称及版本	文件数量	函数数量	代码行数	切片数量
libav-0.8.21/9.21/11.12/12.3	4983	37682	2043116	489866
openssl-1.0.0/1.1.1/3.0.0	3195	24512	1187467	311302
xen-4.12.0/4.14.0/4.15.1	14829	207447	8131195	2281917
httpd-2.0.35/2.2.14/2.4.51	1338	16398	611578	229572
总计	24345	286039	11973356	3312657

本小节对一个 Xen 中关于 USB 仿真程序的漏洞实例进行详细描述,如图 5 所示给出了一个切片 API#2 从切片生成到漏洞解释的完整过程,考虑到实际的切片图占用篇幅较大,此处使用节点对应的代码行展示切片内容.根据分析,该函数中存在一个整数溢出漏洞.该函数允许本地来宾用户通过远程 NDIS 控制消息包造成拒绝服务攻击.具体来说,代码第 7 行声明了两个整数 *bufoffs* 和 *buflen*,*bufoffs* 和 *buflen* 在 11 行和 12 行接收用户输入数值后被赋值,分支条件为 *bufoffs* + *buflen* > *length*(第 13 行).当 *bufoffs* 和 *buflen* 的数值被攻击者故意构造为极大值时会导致整数溢出,从而使得该分支条件不被满足,即无法正常返回.在函数进一步调用第 16 至 18 行函数导致 QEMU 进程崩溃,进而主机内存信息泄露.因此,与 *bufoffs* 和 *buflen* 存在数据和控制依赖关系的代码行,包括 7、11、12、13、16 以及 17,被认为是漏洞相关行.如图 5 中所示,仅 API#2 切片的解释结果就包含正确的漏洞行(17、13、16 行).事实上,原函数的 2-3 个切片的解释结果能够覆盖全部的漏洞引入位置和触发位置,但由于篇幅限制,该漏洞的全部切片和解释结果将公布在本文的代码仓库中^[42].

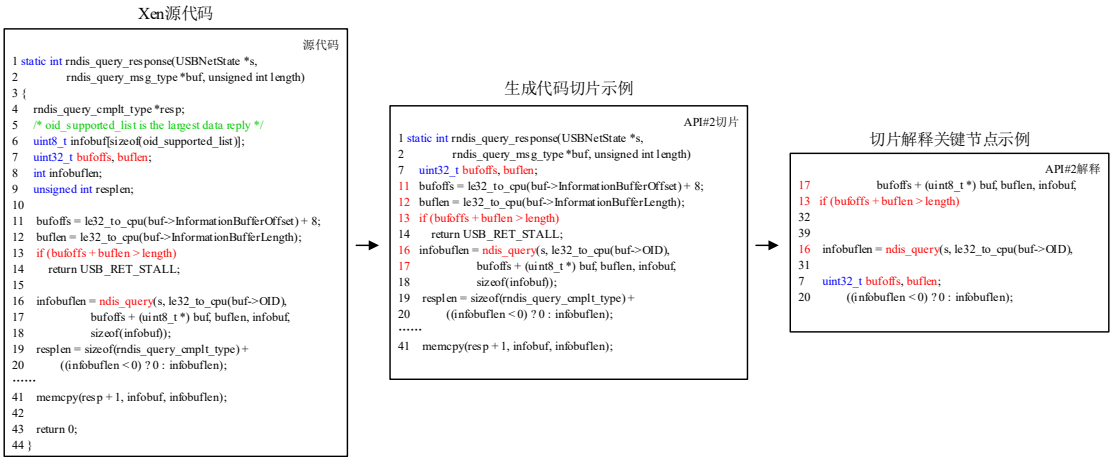


图 5 Xen 真实漏洞检测解释实例分析

根据与已知漏洞模式匹配的验证方式,本方法从待测样本中至少检测到了 54 个新漏洞,对应于 NVD 中 36 种已知的漏洞模式.表 6 的第 1 列和第 2 列显示了本方法检测到的漏洞的详细信息,包括所在软件名称及版本和漏洞所在文件路径,第 3-5 列表示与该漏洞存在相似漏洞模式的已知 NVD 漏洞信息,包含匹配漏洞模式的 CVE 号、报告漏洞的软件名称、漏洞发布时间和漏洞修补情况.其中,21 个漏洞已经在该软件的后续版本中被默默修补或者删除,33 个漏洞仍未进行修补.对于未进行修补的漏洞,本文为了保护软件安全对其详细信息进行了模糊处理.我们已经对其开发厂商提交了漏洞报告和漏洞补丁建议.根据人工分析漏洞成因的验证方式,本方法从待测样本中至少检测到了 5 个新漏洞.表 7 为新漏洞模糊处理后的漏洞信息,我们已将该漏洞提交给

NVD 审核并通知厂商及时修复,待 NVD 审核完成以及漏洞成功修复后,我们将在本文的代码仓库中公布详细的漏洞信息.

表 6 通过与已知漏洞匹配模式验证的漏洞列表

目标软件名称及版本	漏洞文件所在路径	匹配的漏洞模式	报告漏洞软件	发布时间	最新版本是否修复
libav-0.8.21	libavcodec/dsputil.c	CVE-2013-7010	FFmpeg	20131208	修复
	libavcodec/**.c	CVE-2013-****	FFmpeg	20131208	未修复
	libavcodec/error_resilience.c	CVE-2011-3941	FFmpeg	20111001	修复
	libavcodec/**.c	CVE-2015-3395	FFmpeg	20150421	未修复
libav-9.21	libavcodec/**.c	CVE-2013-****	FFmpeg	20131208	未修复
	libavcodec/**.c	CVE-2015-****	FFmpeg	20150421	未修复
	libavcodec/aac_parser.c	CVE-2016-7393	FFmpeg	20160909	修复
	libavcodec/wmalosslessdec.c	CVE-2012-2795	FFmpeg	20120519	修复
	libavcodec/**.c	CVE-2013-****	FFmpeg	20130107	未修复
	libavcodec/dfa.c	CVE-2017-9992	FFmpeg	20170628	修复
	libavcodec/**.c	CVE-2013-****	FFmpeg	20130107	未修复
	libavformat/mov.c	CVE-2016-3062	FFmpeg	20160309	修复
libav-11.12	libavcodec/wmalosslessdec.c	CVE-2012-2795	FFmpeg	20120519	修复
	libavcodec/**.c	CVE-2013-****	FFmpeg	20130107	未修复
	libavformat/**.c	CVE-2017-****	FFmpeg	20170726	未修复
	libavformat/**.c	CVE-2017-****	FFmpeg	20170831	未修复
	libavcodec/h264.c	CVE-2015-3417	FFmpeg	20150424	不存在
	libavcodec/**.c	CVE-2013-****	FFmpeg	20131208	未修复
libav-12.3	libavformat/**.c	CVE-2017-****	FFmpeg	20170726	未修复
	libavcodec/wmalosslessdec.c	CVE-2012-2795	FFmpeg	20120519	修复
	libavcodec/**.c	CVE-2013-****	FFmpeg	20130107	未修复
	libavformat/**.c	CVE-2016-****	FFmpeg	20170201	未修复
	libavformat/**.c	CVE-2018-****	FFmpeg	20180723	未修复
openssl-1.0.0	crypto/asn1/d2i_pr.c	CVE-2015-3195	OpenSSL	20150410	修复
	crypto/asn1/a_d2i_fp.c	CVE-2016-2109	OpenSSL	20160129	修复
	crypto/bn/bn_exp.c	CVE-2016-0702	OpenSSL	20151216	修复
xen-4.12.0	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160302	未修复
	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160223	未修复
	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160112	未修复
	qemu-xen-traditional/**.c	CVE-2016-****	QEMU	20160306	未修复
	qemu-xen/hw/usb/**.c	CVE-2018-****	QEMU	20180911	未修复
	qemu-xen/hw/9pfs/9p.c	CVE-2018-19364	QEMU	20181119	不存在
	qemu-xen/hw/scsi/**.c	CVE-2018-****	QEMU	20181030	未修复
	qemu-xen/hw/scsi/**.c	CVE-2016-****	QEMU	20160523	未修复
	qemu-xen/hw/ppc/**.c	CVE-2018-****	QEMU	20181105	不存在
	qemu-xen/slrp/**.c	CVE-2019-****	QEMU	20190124	未修复
xen-4.14.0	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160302	未修复
	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160223	未修复
	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160112	未修复
	qemu-xen-traditional/**.c	CVE-2016-****	QEMU	20160306	未修复
	qemu-xen/hw/scsi/**.c	CVE-2016-****	QEMU	20160523	未修复
	qemu-xen-traditional/hw/i8254.c	CVE-2015-3214	QEMU	20150410	未修复
xen-4.15.1	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160302	未修复
	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160223	未修复
	qemu-xen-traditional/hw/**.c	CVE-2016-****	QEMU	20160112	未修复
	qemu-xen-traditional/**.c	CVE-2016-****	QEMU	20160306	未修复
	qemu-xen/hw/scsi/**.c	CVE-2016-****	QEMU	20160523	未修复
httpd-2.0.35	lib/xmltok_impl.c	CVE-2009-3720	Expat	20091016	修复
	lib/xmltok_impl.c, xmltok.c	CVE-2016-0718	Expat	20151216	修复
	server/scoreboard.c	CVE-2012-0031	Apache	20111207	修复
httpd-2.2.14	http/byterange_filter.c	CVE-2011-3192	Apache	20110819	修复
	filters/mod_include.c	CVE-2009-1195	Apache	20090331	修复
	arch/win32/mod_isapi.c	CVE-2010-0425	Apache	20100127	修复

综上所述,本方法成功检测并解释了 4 个开源产品中 59 个新漏洞.其中,与现有漏洞模式匹配的新漏洞共计 54 个,人工分析漏洞成因确认新漏洞 5 个.实验结果证明本方法能够有效检测真实软件漏洞并提供可靠的检测结果解释.

表 7 通过漏洞成因分析验证的的漏洞列表

目标软件名称及版本	漏洞文件所在路径	漏洞函数名
libav-0.8.21/9.21/11.12/12.3	libavcodec/**.c	ff_vorbis_**
xen-4.12.0/4.14.0/4.15.1	qemu-xen-traditional/hw/**.c	vmsvga_**
xen-4.14.0	qemu-xen/hw/9pfs/**.c	proxy_**
xen-4.15.1	xen/common/libfdt/**.c	fdt_**
openssl-3.0.0-beta2	crypto/conf/**.c	ossl_**

5 总结与展望

针对已有漏洞检测模型对于真实漏洞的检测效果不佳且无法给出具体漏洞语句的问题,本文提出了基于图神经网络的切片级漏洞检测解释方法.本方法首先针对漏洞关注点提取代码切片,并采用图神经网络和改进的图神经网络解释器对代码切片实现了一个高效准确的漏洞检测和解释方法.本方法在真实漏洞数据集上的测试结果达到 75.1%F1 分数,相较于现有的漏洞检测模型提高了 41.2%-110.4%,实验结果证明了本方法可以有效检测真实漏洞.在对漏洞检测结果的解释方面,本文改进的图神经网络解释器在解释准确率上比其他 2 个通用的图神经网络解释器分别提升 8.9%和 24.9%,时间开销分别缩短 42.5%和 15.4%.实验证明本方法对于真实漏洞能够准确并高效地输出漏洞代码行,即能够为漏洞检测模型提供可信的检测结果解释.最后,本方法扫描 4 个开源软件的千万行代码,成功检测并解释了与 36 种与已知漏洞模式相匹配的真实漏洞,证明了本方法在实际应用方面的有效性.

目前工作仍有具有一定局限性.在漏洞检测方面,并非全部漏洞都是由文中提到的四类漏洞关注点引入,导致部分样本由于不存在漏洞关注点而无法成功提取切片,进而无法进行漏洞检测.因此本工作未来会补充并完善漏洞关注点种类或采用函数级和切片级漏洞检测相结合的方法以使更多样本能够成功检测.其次,本文在漏洞解释部分改进了现有的通用图神经网络解释器 GNNExplainer 作为漏洞解释器,但其对于漏洞样本的特征考虑仍然不够详尽.本工作未来会继续深入分析漏洞样本特征,然后根据漏洞特征设计一种面向漏洞样本的扰动方法并构造解释器.最后,本文对已有真实产品进行漏洞检测解释时,通过人工分析将检测出的漏洞与已知漏洞模式一一比对以确认检测结果的准确性,该过程由 4 位有相关经验的从业者耗时一个月完成.对于本方法检测到的未匹配到已知模式的其他漏洞,则被认为是可能存在的零日漏洞.对于此类漏洞可以采用定向模糊测试技术生成 PoC 以验证其真实存在,并提交 PoC 和漏洞触发信息以申请漏洞号,但是现阶段的模糊测试方法对于该数量级漏洞进行验证将消耗巨大的计算资源.因此,未来我们拟设计一种轻量级的定向模糊测试器以支持本文检测的零日漏洞的后续 PoC 生成等漏洞验证工作.在未来工作中,我们将继续提升漏洞检测模型的检测效果和解释效果,更好地辅助专家进行漏洞挖掘和分析工作.

References:

[1] Li Z, Zou D, Wang Z, Jin H. Survey on static software vulnerability detection for source code. Chinese Journal of Network and Information Security, 2019, 5(1): 1-14.

[2] Chen Z, Zou D, Li Z, Jin H. Intelligent vulnerability detection system based on abstract syntax tree. Journal of Cyber Security, 2020, 5(4): 1-13

[3] Wang J, Kuang H, Li R, Su Y. Software Source Code Vulnerability Detection Based on CNN-GAP Interpretability Model. Journal of Electronics & Information Technology, 2022, 44(7): 2568-2575.

- [4] Duan X, Wu J, Luo T, Yang M, Wu Y. Vulnerability mining method based on code property graph and attention BiLSTM. *Ruan Jian Xue Bao/Journal of Software*, 2020,31(11):3404–3420 (in Chinese). <http://www.jos.org.cn/1000-9825/6061.htm>.
- [5] Software Assurance Reference Dataset, <https://samate.nist.gov/SRD/index.php>.
- [6] “Checkmarx,” <https://www.checkmarx.com/>, 2022.
- [7] “Flawfinder,” <http://www.dwheeler.com/flipfinder/r>, 2022.
- [8] “Rough audit tool for security,” <https://code.google.com/archive/p/rough-auditing-tool-for-security/>, 2022.
- [9] Russell R, Kim L, Hamilton L, Lazovich T, et al. Automated vulnerability detection in source code using deep representation learning. In: *Proceedings of the 17th IEEE international conference on machine learning and applications (ICMLA’18)*, 2018, 757–762.
- [10] Lin G, Xiao W, Zhang J, et al. Deep learning-based vulnerable function detection: A benchmark. In: *Proceedings of the 21st International Conference on Information and Communications Security (ICICS’19)*, 2019, 219–232.
- [11] Li Z, Zou D, Xu S, et al. SySeVR: A framework for using deep learning to detect software vulnerabilities. *IEEE Trans. on Dependable and Secure Computing*, 2021, vol. abs/1807.06756.
- [12] Zhou Y, Liu S, Siow J, et al. Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks. In: *Proceedings of 33rd Annual Conference on Neural Information Processing Systems (NIPS’19)*, 2019, 10197–10207.
- [13] Ying Z, Bourgeois D, You J, et al. Gnnexplainer: Generating explanations for graph neural networks. In: *Proceedings of 33rd Annual Conference on Neural Information Processing Systems (NIPS’19)*, 2019, 32.
- [14] Luo D, Cheng W, Xu D, et al. Parameterized explainer for graph neural network. In: *Proceedings of 34th Annual Conference on Neural Information Processing Systems (NIPS’20)*, 2020, 33: 19620–19631.
- [15] Lin G, Zhang J, Luo W, et al. POSTER: vulnerability discovery with function representation learning from unlabeled projects. In: *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (CCS’17)*, 2017, 2539–2541.
- [16] Feng H, Fu X, Sun H, et al. Efficient Vulnerability Detection based on abstract syntax tree and Deep Learning. In: *Proceedings of IEEE INFOCOM 2020-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS’20)*, 2020: 722–727.
- [17] Wang H, Ye G, Tang Z, et al. Combining graph-based learning with automated data collection for code vulnerability detection. *IEEE Trans. on Information Forensics and Security*, 2020, 16: 1943–1958.
- [18] Wu Y, Zou D, Dou S, et al. VulCNN: An Image-inspired Scalable Vulnerability Detection System. In: *Proceedings of the 44th International Conference on Software Engineering (ICSE’22)*, 2022, 1–12.
- [19] Li Z, Zou D, Xu S, Ou X, et al. Vuldeepecker: A deep learning-based system for vulnerability detection. In: *Proceedings of the 25th Network and Distributed System Security Symposium (NDSS’18)*, 2018, 1–15.
- [20] Zou D, Wang S, Xu S, et al. μ VulDeePecker: A deep learning-based system for multiclass vulnerability detection. *IEEE Trans. on Dependable and Secure Computing*, 2021, 18(5), 2224–2236.
- [21] Duan X, Wu J, Ji S, et al. VulSniper: Focus Your Attention to Shoot Fine-Grained Vulnerabilities. In: *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI’19)*, 2019, 4665–4671.
- [22] Yamaguchi F, Golde N, Arp D, et al. Modeling and discovering vulnerabilities with code property graphs. In: *Proceedings of the 35th IEEE Symposium on Security and Privacy (S&P’14)*, 2014, 590–604.
- [23] Cheng X, Wang H, Hua J, et al. Deepwukong: Statically detecting software vulnerabilities using deep graph neural network. *ACM Trans. on Software Engineering and Methodology*, 2021, 30(3): 1–33.
- [24] Chakraborty S, Krishna R, Ding Y, et al. Deep learning based vulnerability detection: Are we there yet. *IEEE Trans. on Software Engineering*, 2021, 10.1109/TSE.2021.3087402.
- [25] Yuan H, Tang J, Hu X, et al. Xgmn: Towards model-level explanations of graph neural networks. In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (SIGKDD’20)*, 2020, 430–438.
- [26] Schnake T, Eberle O, Lederer J, et al. Higher-order explanations of graph neural networks via relevant walks, 2020, 10.48550/arXiv.2006.03589.
- [27] Selvaraju R R, Cogswell M, Das A, et al. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV’17)*, 2017, 618–626.

- [28] Yuan H, Yu H, Wang J, et al. On explainability of graph neural networks via subgraph explorations. In: Proceedings of the 38th International Conference on Machine Learning (ICML'21), 2021, 12241-12252.
- [29] Li Y, Wang S, Nguyen T N. Vulnerability detection with fine-grained interpretations. In: Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE'21), 2021, 292-303.
- [30] Wu B, Liang X, Zhang S, Xu R. Advances and Applications in Graph Neural Network. Chinese Journal of Computers, 2022, 45(1):34.
- [31] "Open-Source Code Querying Engine for C/C++." <https://joern.io/>, 2020.
- [32] Wolf L, Hanani Y, Bar K, et al. Joint word2vec networks for bilingual semantic representations. International Journal of Computational Linguistics and Applications, 2014, 5(1), 27-42.
- [33] Schlichtkrull M, Kipf T N, Bloem P, et al. Modeling relational data with graph convolutional networks. In: Proceedings of 15th European Semantic Web Conference (ESWC'18), 2018, 593-607.
- [34] Veličković P, Cucurull G, Casanova A, et al. Graph attention networks. arXiv preprint arXiv:1710.10903, 2017.
- [35] Li Y, Tarlow D, Brockschmidt M, et al. Gated graph sequence neural networks. arXiv preprint arXiv:1511.05493, 2015.
- [36] Zhang M, Cui Z, Neumann M, et al. An end-to-end deep learning architecture for graph classification. In: Proceedings of 32nd AAAI Conference on Artificial Intelligence (AAAI'18), 2018, 4438-4445.
- [37] D. P. Kingma and J. Ba, Adam: A method for stochastic optimization, In: Proceedings of the 3rd International Conference on Learning Representations (ICLR'15), 2015. Available: <http://arxiv.org/abs/1412.6980>
- [38] Fan J, Li Y, Wang S, et al. A C/C++ code vulnerability dataset with code changes and CVE summaries. In: Proceedings of the 17th International Conference on Mining Software Repositories (MSR'20), 2020, 508-512.
- [39] Common Vulnerabilities & Exposures, <https://cve.mitre.org/>
- [40] Tensors and Dynamic neural networks in Python with strong GPU acceleration (PyTorch), <https://pytorch.org.>, 2022.
- [41] National Vulnerability Database, <https://nvd.nist.gov/>.
- [42] "Vuldetexp", <https://github.com/Stwsyburg/Vuldetexp>

附中文参考文献:

- [1] 李珍,邹德清,王泽丽,金海.面向源代码的软件漏洞静态检测综述.网络与信息安全学报,2019, 5(1):1-14.
- [2] 陈肇炫,邹德清,李珍,金海.基于抽象语法树的智能化漏洞检测系统.信息安全学报,2020,5(4):1-13.
- [3] 王剑,匡洪宇,李瑞林,苏云飞.基于 CNN-GAP 可解释性模型的软件源码漏洞检测方法.电子与信息学报,2022,44(7): 2568-2575.
- [4] 段旭,吴敬征,罗天悦,杨牧天,武延军.基于代码属性图及注意力双向 LSTM 的漏洞挖掘方法.软件学报,2020,31(11):3404-3420.
- [30] 吴博,梁循,张树森,徐睿.图神经网络前沿进展与应用.计算机学报, 2022, 45(1):34.