

# 基于流的僵尸网络检测框架\*

张军<sup>1,2</sup>, 程光<sup>1,2</sup>

<sup>1</sup>(东南大学 计算机科学与工程学院, 江苏 南京 211189)

<sup>2</sup>(东南大学 教育部计算机网络和信息集成重点实验室, 江苏 南京 211189)

**摘要:** 僵尸网络是目前互联网安全最大安全问题之一, 随着僵尸网络的变化, 僵尸网络的检测方法也在不断更新, 从基于数据包检测到现在基于僵尸网络通信行为检测。本文提出一种僵尸网络检测系统框架, 该系统仅使用高质量网络流就可以对僵尸感染主机数据流进行检测。该检测框架主要分为僵尸网络特征获取模块和僵尸网络检测模块。僵尸网络特征获取模块是基于主动方式的僵尸网络通信流量获取系统, 模块主体是利用主动技术的恶意代码捕获部分和僵尸网络通信流量获取部分。僵尸网络检测模块是通过对获取的僵尸网络通信流量进行特征分析提取完成僵尸网络检测。

**关键词:** 僵尸网络检测; 网络流量特征; 无监督聚类; 网络安全;

中图法分类号: TP393 文献标识码: A

## The research on the framework of botnet detection\*

ZHANG Jun<sup>1,2</sup>, CHENG Guang<sup>1,2</sup>

<sup>1</sup>(School of Computer Science and Engineering, Southeast University, Nanjing 211189, China)

<sup>2</sup>(Key Laboratory of Computer Network and Information Integration, Ministry of Education, Southeast University, Nanjing 211189, China)

**Abstract:** Botnet is one of the largest Internet security problem at present. With the change of the Botnet, the method of Botnet detection has been updated, which is changing from DPI detection to DFI detection. This paper presents a botnet detection system framework, that can be used to detect botnet traffic,ection with high quality network flow. The framework of this system is mainly divided into feature creation module and botnet detection module. Feature creation module is based on the active technology that can actively get the traffic of botnet communication. Botnet detection module analysis and extract the characteristics of Botnet communication traffic to detect botnet.

**Key words:** botnet detection; network measurement; traffic classification; network security

---

\*作者简介: 张军(1989—), 男, 江苏南京人, 硕士生, 主要研究领域为网络安全, E-mail: jzhang@njnet.edu.cn; 程光(1973—), 男, 博士, 教授, 博士生导师, 主要研究领域为网络测量、网络安全、网络管理等, E-mail: gcheng@njnet.edu.cn。

# 1 引言

僵尸网络是现代网络中威胁最大的安全问题之一，僵尸程序感染网络中存在安全隐患的主机。僵尸网络控制者集中控制僵尸网络的计算能力和数据资源，并利用这些资源进行非法活动，如信息窃取、垃圾邮件、钓鱼攻击以及 DDoS 攻击等。

文献[1]对僵尸网络进行了基本定义：僵尸网络是一组被控制的主机，这些主机通过 C&C 服务器获取命令，共同完成指定任务。僵尸网络中的主机必须与 C&C 服务器保持联系，确保能够成功获得命令。相同僵尸网络中不同主机在与 C&C 服务器通信时会表现出相似特征。

本文提出一种僵尸网络追踪检测系统框架，该检测框架主要分为两个部分：第一部分为僵尸网络通信数据获取模块，搭建虚拟网络环境，利用主动技术获取的僵尸程序，并在模拟环境中运行，获取通信数据，分析其网络通信流量通信特征。第二部分为僵尸网络检测模块，检测模块有两种工作模式：第一种模式为训练模式，输入标记数据，进行训练过程，产生决策模型；第二种模式是检测模式，利用决策模型检测未知通信数据。

本文安排如下：第二节介绍僵尸网络检测研究的相关现状，第三节介绍僵尸网络分析检测平台的主体框架和功能模块，第四节介绍僵尸网络分析检测平台的分析检测结果，最后是总结。

# 2 相关研究

目前，僵尸网络在不断变化，其检测方法也在不断变化。最早僵尸网络检测是基于数据包的负载检测，近年开始基于流量行为特征对僵尸网络进行检测。

Masud 在文献[3]中提出了一种 P2P 僵尸网络检测方法，这种方法将僵尸网络流量视为流数据。文中的评估方法使用僵尸网络流量和模拟数据，检测正确率优于流数据分类技术。Gu 先后提出了两种僵尸网络检测系统分别命名为 BotHunter<sup>[4]</sup>和 BotMiner<sup>[1]</sup>。BotHunter 利用 Snort 入侵检测系统产生的关联分析报告进行僵尸网络检测。BotHunter 指出僵尸网络中的主机具有相似的行为特征，而这些行为都属于生命周期的一部分：扫描、感染、二进制文件下载、C&C 通信和外向扫描。BotHunter 在僵尸主机执行完整生命周期行为时是最有效的，而且加密的 C&C 通信会影响 BotHunter 的检测效果。BotMiner 是 Gu 在 2008 年提出的僵尸网络检测系统，基于僵尸网络主机行为的群组性进行检测。BotMiner 揭示僵尸网络中主机行为模式的潜在一致性，并对僵尸主机行为进行聚类分析检测僵尸网络。实验表明，BotMiner 对部分流行的变异僵尸程序能达到 99% 的检测率，而误报率仅为 1%。Zeidanloo 和 Rouhani 在 2012

年也提出了一种基于检测网络流量特征的僵尸网络监测系统<sup>[5]</sup>，通过三个阶段处理（过滤、恶意行为检测和流量监听）并根据僵尸主机群组性行为对僵尸主机进行分组。

基于群组行为的聚类方法都是检测正在进行某种恶意行为的僵尸主机，无法在僵尸主机早期的 C&C 通信阶段进行有效检测。另外，群组行为检测方法是以网络中具有多个同种僵尸主机为前提，当监测网络仅有一个受感染主机时，这种检测方法的有效性将受到影响。随着机器学习的不断发展，现在越来越多的僵尸网络检测系统融入了机器学习，利用标记好的僵尸网络流量，进行学习分析，再对未作标记的流量进行分析检测，如文献[6][7][8][9]等所述的检测方法。文献[10][11][12]则是从信息论的角度研究僵尸网络通信特征。

以上文献所用分析方法基本思想都是将僵尸网络通信流量的行为特征作为研究对象，减少僵尸主机与 C&C 服务器之间加密通信的干扰。

# 3 僵尸网络分析检测平台

## 3.1 僵尸网络通信机制

僵尸主机与 C&C 服务器以某种方式进行通信，以便控制者能够控制整个僵尸网络。僵尸网络有两种通信机制，即“pull”机制和“push”机制，具体实现方式如图 1 所示。因为大多主机处于防火墙的保护之中，被动接受命令的“push”机制会受到防火墙阻拦。通常僵尸网络采用“pull”方式进行通信连接，即僵尸主机主动向 C&C 服务器发起通信连接。

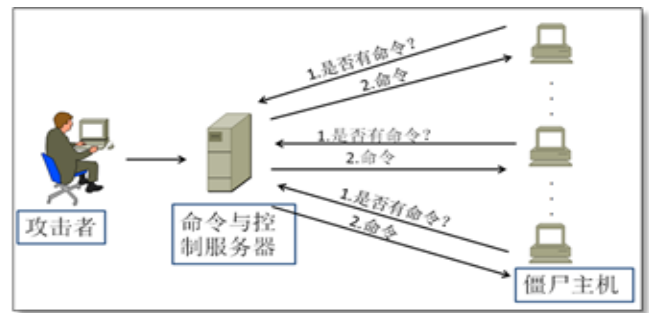


图 1(1) pull 机制

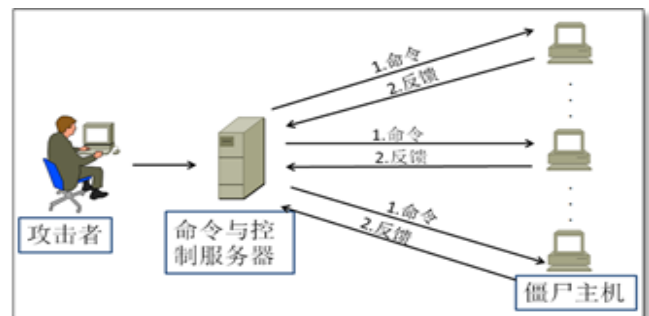


图 1(2) push 机制

早期的僵尸网络检测主要都是基于负载分析，分析 TCP 和 UDP 数据包是否含有恶意结构。负载分析方法具有良好的正确率，但是随着限制因素的增加其检测效果会收到严重影响。另外，新型僵尸程序与 C&C 服务器通常采用加密或其他混淆技术进行通信，负载检测无法有效检测这类数据包。基于僵尸网络通信行为特征的分析可以解决僵尸网络加密通信的问题。同种僵尸网络的主机通常会表现出一致性通信行为和相似的行为特征，而某些测度可以用来表示这些行为特征。通信行为特征分析并不是分析数据包内容，所以加密或混淆通信不会对检测产生影响。但是，基于通信行为特征的僵尸网络检测技术具有一定的滞后性，流根据不同的定义方式滞后程度有所不同。流的结束取决于下一个数据包到达时的状态，以及相关定义，但是这种滞后性并不影响对离线记录的检测。

### 3.2 系统结构与功能概述

僵尸网络分析检测平台主要分为两个部分，第一部分为僵尸网络通信数据获取模块，第二部分为僵尸网络检测模块。

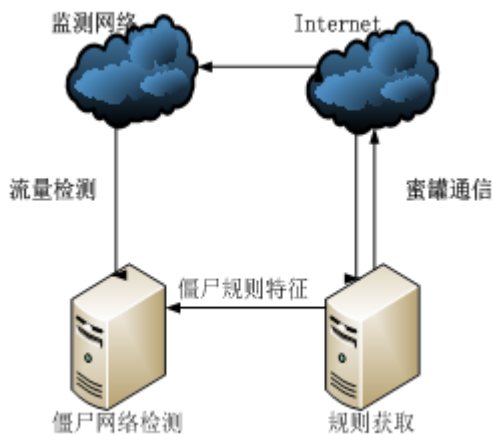


图2 系统结构图

如上图2所示，僵尸网络特征获取模块将捕获的僵尸样本程序在可控制环境中运行，获取僵尸主机与C&C服务器的通信数据进行分析，将分析获得的僵尸网络特征数据传输给僵尸网络检测模块，僵尸网络检测模块利用接受到的特征数据，对网络流量进行检测分析。

本僵尸网络流量分析检测平台主要用C语言进行编写，并在Linux系统上运行。整个分析检测平台主要由可执行功能模块和插件组成，系统流程如图3所示。

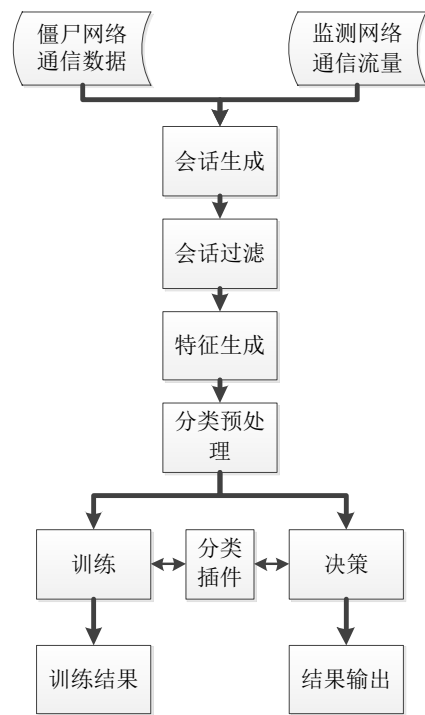


图3 僵尸网络分析检测系统流程图

**会话生成** 本阶段解析数据包并完成组流，该部分根据网络通信数据生成会话数据。解析数据包的开发主要是基于Libpcap开发包。本文中，流采用的是传统的五元组定义，即(源IP,宿IP,源端口,宿端口,协议)，超时时间设置为64s。一个流为单向流，每个流记录统计如下信息：报文数、字节数、流开始时间、流结束时间，如表1中基本特征。本阶段使用哈希链表追踪记录，哈希算法可以使得系统能够适应高速网络的实时通信数据。

**会话过滤** 会话数据(即流记录)生成后进入过滤阶段，这个阶段对会话数据根据配置规则完成以下过滤：首先，过滤外部主机发起的通信流；其次，过滤单向数据流；最后，白名单过滤。其中，白名单是根据Alexa.com公布的中国前一百访问量网站得出。

**特征生成** 处理过滤后的会话数据，生成特征。族(trace)是僵尸网络检测系统BotFinder<sup>[13]</sup>中提出的重要概念，族保存记录两个终端间通信的流序列信息。族表示在一个时间窗口内主机与服务器在某个端口的通信情况，按照(源IP,宿IP,宿端口)三元组对会话数据进行聚合，生成高级特征。注：表1中高级特征都是根据基本特征计算而来。

表1 基本特征与高级特征

特征类型	特征名	包检查方法
基本特征	字节数	统计每个数据包
	报文数	统计每个数据包
	流开始时间	统计流的第一个数据包
	流结束时间	统计流的最后一个数据包
高级特征	流持续时间	每个流的开始与结束时间差
	子流开始时间间隔	族中流开始时间序列

在聚合处理中，使用平衡二叉树按照流开始时间时间顺序保存族中流信息，平衡二叉树中每个节点保存单

个流的基本统计信息。僵尸网络流量分析检测平台中所用特征测度可以根据需要进行修改添加，便于检测方法的改进。

**分类预处理** 输入为僵尸网络通信数据时，本阶段生成格式化数据并标记僵尸种类。输入为监测网络通信数据时，本模块仅生成未标记的格式化数据。

**分类插件** 本插件对数据进行训练与分类，可以根据不同的方法完成上述工作，如聚类方法和机器学习算法等。目前分类插件所用为无监督聚类 Xmeans，利用 Xmeans 算法可以在不确定分类情况进行聚类。也可以使机器学习算法，如决策树等，利用 weka 提供的 java 接口，编写相应的算法插件，实现机器学习进行决策分析。可以实现含有多种聚类或机器学习方法的分类插件，生成综合性分类结果，提交给决策阶段进行分析。

**训练阶段** 输入数据是训练数据时，分析平台调用分类插件，将标记的通信数据输入分类插件进行训练。并根据分类插件的分类结果产生相应的训练结果集。目前本平台采用 Xmeans 无监督聚类算法，根据分类插件聚类的结果，产生相应聚类模型以及聚类中心和聚类质量等参数。

**决策阶段** 当分析检测平台对未知网络通信流量进行分析时，决策阶段使用分类插件，根据训练阶段产生的分类模型对未知流量进行检测分析，对分类结果进行决策评估。

**结果输出** 本阶段输出检测结果，结果包括每个族的处理信息以及相应的分类结果，从而判断是否僵尸网络通信流量，以及所属僵尸种类。

### 3.3 僵尸网络通信数据获取模块

图 3 中所用僵尸网络通信数据由僵尸网络特征获取模块产生，并根据僵尸网络报告，对产生的僵尸网络通信数据进行标记。

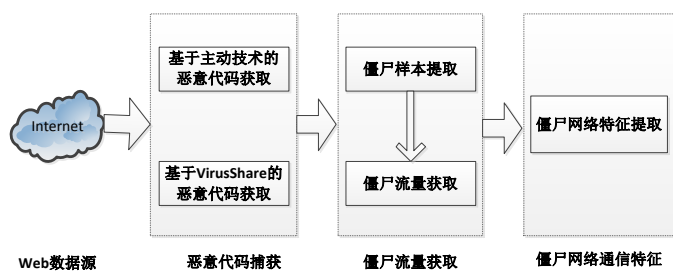


图 4 僵尸网络特征获取模块流程图

图 4 为僵尸网络通信数据的获取流程，首先基于主动方式获取恶意代码，然后从中提取出僵尸样本程序，最后在模拟环境中自动运行僵尸样本程序，获取僵尸主机与 C&C 服务器的通信数据。整个模块主要分为三个部分：恶意代码样本获取、僵尸样本程序分析提取、僵尸网络通信数据获取与特征提取。

**恶意代码获取** 从大量的 Web 数据中获取恶意代码。通过两种方法获取恶意代码，一种是基于主动技术的恶意代码捕获；另一种是通过 VirusShare 网站获取恶

意代码。前者是利用网络爬虫捕获网络中传播的恶意代码，通过网络爬虫能够获取未知恶意代码。而 VirusShare 网站分享了大量的恶意代码文件，可以高效的获取恶意代码。

**僵尸样本提取** 本阶段在大量的恶意代码中筛选出僵尸样本程序，VirusTotal 网站提供恶意代码样本分析服务，通过 VirusTotal 回馈的分析报告，获得恶意代码种类，筛选出僵尸样本程序以及完成僵尸样本分类。

**获取僵尸样本通信数据** 在搭建的虚拟环境中运行筛选出的僵尸样本程序，主动运行僵尸样本程序，获取僵尸网络通信数据。分析通信数据有效性的，获取特征。

## 4 僵尸网络分析检测平台检测结果

本文下面所述僵尸网络通信流量均来自僵尸网络分析检测平台中僵尸网络通信数据获取模块。目前，本平台通过网络爬虫和 VirusShare 网站获取了大量的僵尸程序，下面进行简要分析。

文献[15][16][17]中，对流行和最新出现的僵尸程序进行的统计，可以分为以下常见僵尸程序：sdbot、Agobot、GT-bot、Rbot、Bobax、Rustock、Clickbot、Phatbot、Sinit、Phatbot、SpamThru、Nugache、Peacomm、poebot、gaobot、zbot、spybot、kwbot、shellbot、backdoor.bot、spambot、antbot、dsbot，统计结果如下：

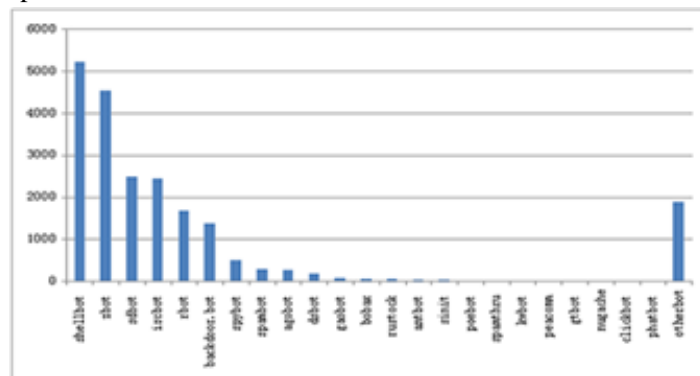


图 5 僵尸程序统计情况

图 6 是 IRC 僵尸控制服务器端口分布图。对 IRC 僵尸网络通信数据分析得到 C&C 服务器的端口分布情况。IRC 僵尸网络中 6667 端口依然是主要的端口，占全部端口的 17%；。而 80 号端口占了所有端口的 14%。IRC 僵尸网络中的非标准端口占了全部的 69%，例如 6969 号端口、445 号端口、8080 号端口。攻击者在 IRC 僵尸网络中大量的使用非标准端口，是为了混淆端口，从而躲避基于端口的检测方法。

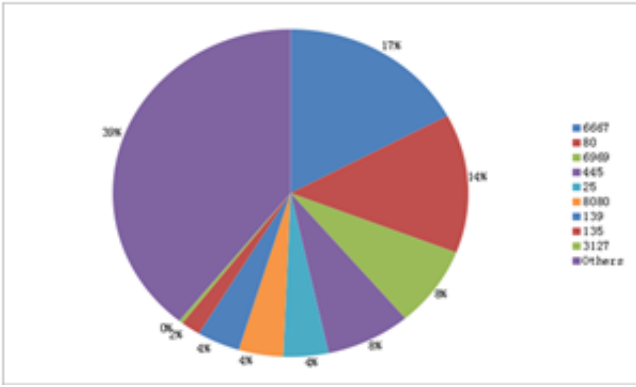


图6 IRC僵尸C&C服务器端口分布

图7是HTTP僵尸C&C服务器的端口分布图，从图中可以发现，80号端口占有所有端口数量的69%。但是，仍然有很多其他类型的端口，3300端口占了8%，25端口占了总数的6%。当僵尸主机通过HTTP通信获得了控制命令后，就会表现出其他行为，例如25号端口是SMTP(简单邮件传输协议)，C&C服务器可能利用HTTP僵尸网络发送垃圾邮件。

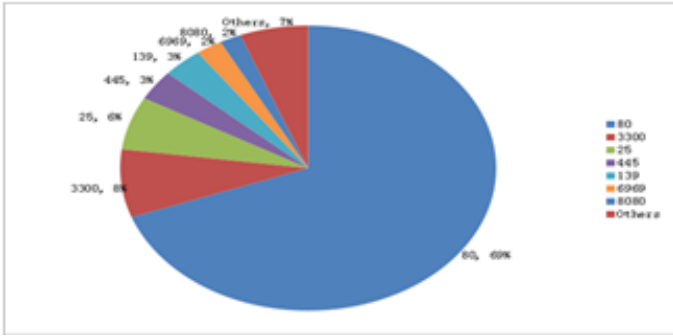


图7 HTTP僵尸C&C服务器端口分布

本文使用僵尸网络检测模块采用僵尸网络检测系统 BotFinder 所述方法，将获取的僵尸网络通信数据分为两部分，一部分作为训练数据，另一部分作为测试数据。使用 Xmeans 无监督聚类算法对上述僵尸网络通信流量进行训练分析，计算出不同方向特征的聚类中心和聚类质量等参数。选取的僵尸程序种类如下表所示：

表2 本文实验所用僵尸程序

僵尸程序类型	僵尸样本数	产生的流序列数
SDbot	19	56
Shadowbot	10	40
BlackEnergy	10	38
Zbot	23	46

通过 Xmeans 聚类的训练结果，构建聚类模型，在未知流量进入分析检测平台后，通过分类插件对各个特征聚类后进行决策。检测结果如图8所示。

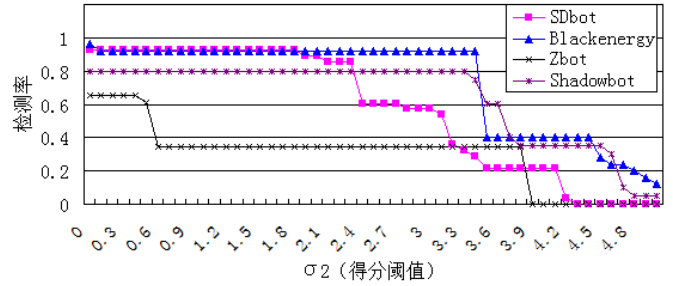


图8 匹配特征阈值为5时，得分阈值的检测情况

图8为检测模块检测结果，当匹配特征阈值为5时，即至少匹配至少5个特征时，才会被判定为僵尸网络通信数据。如图8所示，在匹配特征阈值为5时，对得分低于2.0的检测率均在80%以上(除Zbot)，可以看出使用本系统对Blackenergy和SDbot的检测效果最佳，Shadowbot次之，对Zbot的检测效果并不理想。其中，botFinder对Blackenergy检测率为85%，如图本系统对该种僵尸检测率约为90%。

图9是各类僵尸误报率和漏报率的对比，分别使用柱状图和折线图表示四类僵尸网络的误报率和漏报率。图9左边数值表示误报率，其中误报率Blackenergy最高接近0.016%，Shadowbot最低，而Zbot因其HTTP请求流序列的特殊性，误报率则较低。图9右边数值则表示漏报率情况，漏报率Zbot最高约65%，明显高于其他三种。由上述检测结果可以看出，本系统所用检测方法对Zbot检测效果劣于其他三种僵尸网络。

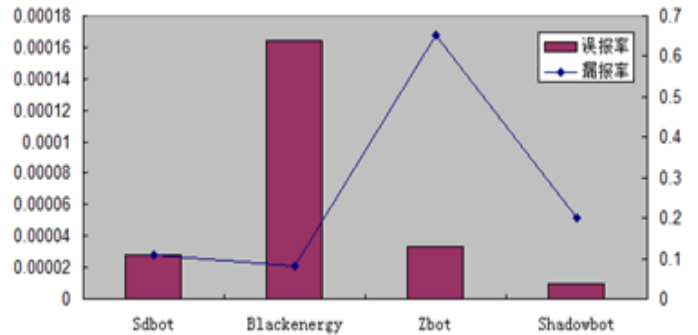


图9 各类bot误报率/漏报率对比

有上述分析，本节使用 BotFinder 中的方法对 SDbot、Shadowbot、BlackEnergy、Zbot 四种僵尸网络进行检测分析，以流序列为研究对象，提取流间隔时间、流持续时间、流字节数、交互频率等特征，采用一维聚类及相应的评分方式，建立僵尸模型并与未知流量进行匹配，从而达到检测僵尸网络的目的，总体上取得了良好的效果。

## 5 结束语

僵尸网络对互联网威胁的日益加剧，迫切需要一种完整的僵尸网络检测框架。现有僵尸网络检测方法都是直接对僵尸网络进行分析研究，忽视了僵尸网络通信样本数据的重要性，僵尸网络研究流程缺少有效的分析数

据来源。

本文认为在僵尸网络分析检测中,僵尸样本程序和僵尸网络通信数据也是僵尸网络研究中不可或缺的一部分,只有保证获取的僵尸网络通信数据的有效性才能更好的进行僵尸网络研究。本文并不是提出具体的检测方法,是在传统僵尸网络检测技术的基础上提出了一种完整的僵尸网络通信数据分析框架与流程,该框架由僵尸网络特征获取模块与僵尸网络检测模块组成。僵尸网络特征获取模块是以主动获取恶意代码技术为基础,构建的僵尸网络通信流量获取系统。僵尸网络检测模块根据特征获取模块提供的特征数据对监测网络进行僵尸网络检测。

在后续研究中需要进一步对比分析不同的分类方法和不同僵尸网络类型检测的有效性。另外,特征获取模块中的僵尸网络数据通信的有效性直接影响系统的检测效果,在以后的工作中需要对获取到的僵尸网络通信数据进行筛选,进一步提高训练数据的质量。

### 参考文献

- [1] Gu G, Perdisci R, Zhang J, et al. BotMiner: Clustering Analysis of Network Traffic for Protocol-and Structure-Independent Botnet Detection[C]//USENIX Security Symposium. 2008: 139-154.
- [2] AsSadhan B, Moura J M F, Lapsley D. Periodic behavior in botnet command and control channels traffic[C]//Global Telecommunications Conference, 2009. GLOBECOM 2009. IEEE. IEEE, 2009: 1-6.
- [3] Masud M M, Gao J, Khan L, et al. Mining concept-drifting data stream to detect peer to peer botnet traffic[J]. Univ. of Texas at Dallas Tech. Report# UTDCS-05-08, 2008.
- [4] Gu G, Porras P A, Yegneswaran V, et al. BotHunter: Detecting Malware Infection Through IDS-Driven Dialog Correlation[C]//USENIX Security. 2007, 7: 1-16.
- [5] Zeidanloo HR, Rouhani S. Botnet detection by monitoring common network behaviors. Lambert Academic Publishing, ISBN 9783848404759; March 2012.
- [6] Saad S, Traore I, Ghorbani A, et al. Detecting P2P botnets through network behavior analysis and machine learning[C]//Privacy, Security and Trust (PST), 2011 Ninth Annual International Conference on. IEEE, 2011: 174-180.
- [7] Bilge L, Balzarotti D, Robertson W, et al. Disclosure: detecting botnet command and control servers through large-scale netflow analysis[C]//Proceedings of the 28th Annual Computer Security Applications Conference. ACM, 2012: 129-138.
- [8] Stevanovic M, Pedersen J M. An efficient flow-based botnet detection using supervised machine learning[C]//Computing, Networking and Communications (ICNC), 2014 International Conference on. IEEE, 2014: 797-801.
- [9] Zhao D, Traore I, Sayed B, et al. Botnet detection based on traffic behavior analysis and flow intervals[J]. Computers & Security, 2013, 39: 2-16.
- [10] Zhigang J, Ying W, Bo W. P2P Botnets detection based on user behavior sociality and traffic entropy function[C]//Consumer Electronics, Communications and Networks (CECNet), 2012 2nd International Conference on. IEEE, 2012: 1953-1955.
- [11] Kang J, Zhang J Y. Application entropy theory to detect new peer-to-peer botnet with multi-chart CUSUM[C]//Electronic Commerce and Security, 2009. ISECS'09. Second International Symposium on. IEEE, 2009, 1: 470-474.
- [12] Husna H, Phithakkitnukoon S, Dantu R. Traffic shaping of spam botnets[C]//Consumer Communications and Networking Conference, 2008. CCNC 2008. 5th IEEE. IEEE, 2008: 786-787.
- [13] Tegeler F, Fu X, Vigna G, et al. Botfinder: Finding bots in network traffic without deep packet inspection[C]//Proceedings of the 8th international conference on Emerging networking experiments and technologies. ACM, 2012: 349-360.
- [14] Donato W D, Dainotti A. Traffic identification engine: an open platform for traffic classification[J]. Network, IEEE, 2014, 28(2): 56-64.
- [15] 诸葛建伟, 韩心慧, 周勇林, 等. 僵尸网络研究[J]. 软件学报, 2008, 19(3): 702-715.
- [16] Grizzard J B, Sharma V, Nunnery C, et al. Peer-to-peer botnets: Overview and case study[C]//Proceedings of the first conference on First Workshop on Hot Topics in Understanding Botnets. 2007: 1-1.
- [17] Li W M, Xie S L, Luo J, et al. A Detection Method for Botnet based on Behavior Features[J]. Advanced Materials Research, 2013, 765: 1512-1517.