

第一章 概述

1. **网络安全威胁的主体有哪些？**公共基础设施安全运行、国家军事安全、经济金融安全、文化安全。**业余黑客 黑产组织 网络犯罪团伙或黑客组织 网络恐怖组织 一般/高级/超高能能力国家/地区行为体**
2. **什么是网络空间？网络空间的四要素包括哪些？**网络空间就是所有可对外交换信息的电磁设备作为载体，通过与人与人互动而形成的虚拟空间。包括互联网、通信网、广电网、物联网、社交网络、计算机系统、通信系统、控制系统等。四要素：**“网络角色”依托“信息通信技术系统”以进行“广义信号”的“交互”**。“信息通信技术系统”包括互联网、各种电信网与通信系统、各种传播系统与广电网、各种计算机系统、各类关键工业设施中的嵌入式处理器和控制器等声光电或数字信息处理设备，用“设施”来表征。“广义信号”是指基于声、光、电、磁等各类能够用于表达、存储、加工、传输的电磁信号，以及量子信号、生物信号等能够与电磁信号进行交互的信号形态，这些信号通过在信息通信技术系统中进行加工处理而成为“信息”，用“数据”来表征。“设施”与“数据”反映的是信息通信技术基础设施。“网络角色”是指生产、传输广义信号的主体，用“用户”来表征；“交互”是指用户借助广义信号，以信息通信技术设施为平台，以信息通信技术为手段，达到产生信号、保存数据、修改状态、传输信息、输出结果等表达人类意志的行为，用“操作”来表征。“用户”与“操作”反映的是“与信息通信技术相关的活动”。
3. **《网络安全法》中的主体、客体主要有哪些？试列举各主体的基本责任和义务。**
- 主体：**人：管理者；政府部门；网络运营者；是指网络的所有者、管理者和网络服务提供者；产品和服务提供者(360、华为、阿里等公司)；信息传播者(微博、微信等)；信息提供者(新浪网、凤凰网等)；**客体：**财物；网络数据；是指通过网络收集、存储、传输、处理和产生的各种电子数据。公民个人信息：是指以电子或者其他方式记录的能够单独或者与其他信息结合识别公民个人身份的各种信息。产品（软件、硬件）、服务。网络基础设施。责任义务：**（一）制定内部安全管理制度和操作规程，确定网络安全负责人，落实网络安全保护责任；（二）采取防范计算机病毒和网络攻击、网络侵入等危害网络安全行为的技术措施；（三）采取监测、记录网络运行状态、网络安全事件的技术措施，并按照规定留存相关的网络日志不少于六个月；（四）采取数据分类、重要数据备份和加密等措施；
4. **什么是网络空间？什么是网络空间安全？网络空间安全研究方向有哪些？**
- 网络空间**定义见第二题。**网络空间安全**主要是在信息通信技术的电磁设备、电子信息系统、运行数据、系统应用等系统与应用层面上，围绕信息获取、信息传输、信息处理、信息利用等核心功能，针对网络空间的设施、数据、用户、操作等核心要素来采取安全保护措施，以确保网络空间中信息通信技术系统及其所承载数据的机密性、可鉴别性（包含完整性、真实性、不可抵赖性）、可用性、可控性等元安全属性得到保障，从而保证信息通信系统能够提供安全、可信、可靠、可控的服务。**网络空间安全研究方向：网络空间安全基础、密码学及应用、系统安全、网络安全、应用安全、信息内容安全**

5. **国家网络强国战略的主要内容？我国政府中网络安全部门主要有哪些？**
- 主要内容：**重视互联网、发展互联网、治理互联网、人才培养**主要部门：**国家网信部门、国务院电信主管部门、公安部门和人民政府有关部门（网信办、工信部、公安部门）
6. **什么是网络空间主权？基本原则是什么？**（网络独立权，平等，自卫，管辖）
- 一个国家的网络空间主权建立在本国所管辖的信息通信技术系统之上（领网）；其作用边界为由直接连向他国网络设备的本国网络设备端口集合所构成（疆界）；用于保护虚拟角色对数据的各种操作（政权、用户、数据）。网络空间的构成平台、承载数据及其活动受所属国家的司法与行政管辖（管辖权）；各国可以在国际网络互联中平等参与治理（平等权）；位于本国领土内的信息通信基础设施的运行不能被他国所干预（独立权）；国家拥有保护本国网络空间不被侵犯的权力及其军事能力（自卫权）。网络空间主权应该受到相互尊重（尊重主权）；国家间互不侵犯他国的网络空间（互不侵犯）；互不干涉他国的网络空间管理事务（不干涉他国内政）；各国网络空间主权在国际网络空间治理活动中具有平等地位（主权平等）。
- 基本原则：网络空间主权的基本原则也是来自于国家主权。尊重网络主权，就是要尊重网络独立权，不采取导致主权网络空间无法自主运行的行为；互不侵犯，就是不能对他国的网络空间实施网络攻击；互不干涉网络内政，就是对主权网络空间的管辖权不指手画脚；网络主权平等，就是主权国家之间具有平等共治网络空间的权力，而不是依靠“利益攸关方”的模式导致一些国家失去了参与网络共治的权力，而另一些国家则掌控了全球的网络空间。
7. **什么是信息内容安全？**主要涉及对传播信息的有效审查监管，剔除非授权信息（非法信息、泄密信息、垃圾邮件等），保护授权信息。
8. **信息内容安全的主要威胁有哪些？有哪些典型事件？**【净网 反恐 反腐 知识产权 邪教 反分裂 隐私保护】政治方面：防止来自国内外反动势力的分裂、邪教、暴恐行为，虚假内容传播；健康方面：剔除色情、淫秽和暴力内容等；保密方面：防止国家和企业机密被窃取、泄露和流失；隐私方面：防止个人隐私被篡改、倒卖、滥用和扩散；产权方面：防止知识产权被剽窃、盗用等；破坏方面：防止病毒、垃圾邮件、网络蠕虫等恶意信息耗费网络资源；经济方面：防止赌博、诈骗、骗贷等；净网行动 全民反恐 全民反腐 反邪教 反分裂 知识产权 隐私保护典型事件：全球反恐打黑、信息战、网络诈骗与舆论安全、BT 侵权、影视盗版
9. **信息内容安全技术主要包括哪些？**信息理解、发现与追踪技术，信息识别、检索与筛选技术，信息监测与阻断技术，内容信息版权保护技术。大规模串匹配技术、面向海量文本信息的分类聚类技术、语音/图像识别技术、信息渗透及检测
10. **信息内容安全技术面临的挑战是什么？**（数据量大，计算复杂度高，网络技术新，社会矛盾深）**规模巨大：**非结构化数据的超大规模，比结构化数据增长快 10 倍到 50 倍；产生高速：实时分析而非批量式分析，数据输入、处理与丢弃，立竿见影而非事后见效；形式多样：异构性（文本、图像、视频、机器数据），模式不明显，语法语义不连贯；信息价值：大量的不相关信息，对未来趋势与模式的深度复杂分析（机器学习、人工智能）
11. **网络与信息安全问题的本质根源是什么？**物理安全问题、方案设计的缺陷、系统的安全漏洞、TCP/IP 协议的安全问题、人为的无意失误、人为的恶意攻击、管理上的因素

第二章 网络信息获取

	被动获取	主动获取	<b>网络信息获取方式分为哪两类？试比较各自的主要特点。</b>  被动获取需要获得管理员授权，协助接入网络设备；且要求高实时性，否则丢包。主动获取无需协助接入，只要客户端接入互联网即可获取，且数据处理可离线进行，无需实时性。被动获取需要对应层协议进行解析，一事一议；主动获取需要针对信息发布方式设计爬虫程序，一事一议
是否需要协作	需要授权	无需授权	
获取范围	仅仅包括途径机器的流量(局部网络)	网上的开放信息(互联网)	
数据处理实时性	实时性要求高(在线)	无需实时性要求(离线)	
技术指标	吞吐量、丢包率	爬取效率	
技术难点	理论分析	适应不同的发布方式	

1. **BPF 捕包原理是什么？** Berkeley Packet Filter 是一个高效的数据包捕获机制，工作在操作系统的内核层。BPF 主要由网络转发部分和数据包过滤两部分组成。网络转发部分是从事链路层捕获数据包并把它们转发给数据过滤部分，数据包过滤部分是从接收到的数据包中接收过滤规则决定的网络数据包，其他数据包被丢弃在操作系统的内核中完成，效率很高。使用了数据缓存机制，使捕获数据包缓存存在内核中，达到一定数量再传递给应用程序。实际应用中，使用 libpcap。
2. **被动捕包程序的主要流程是什么？** (1) 把网卡等同于文件进行 I/O，查找网卡：Find\_all\_devices()，打开网卡：open ()；(2) 从网卡中读取数据，监听：loop()，数据回传给用户变量(3) 处理获取的数据，转用户程序执行：Handler()；(4) 释放 I/O 资源
3. **关于各种数据包：IP 首部：版本 4 首部长度 4 服务类型 8 总长 16 标识 16 分段偏移 16 生存期 8 协议 8 头部校验和 16 源 IP 32 目的 IP 32 选项 TCP 首部：源 16 目的 16 序列号 32 确认号 32 数据偏移 4 保留 6 Flags(6) 窗口 16 校验和 16 紧急指针 16 选项(可变) 填充 UDP 首部：源 16 目的 16 数据偏移 16 校验值 16；[0~31，单位都是 bit 比特]**
4. **IP 首部主要包含哪些信息？在数据包分析时如何提取？** 版本，首部长度，区分服务，总长度，标识，标志，片偏移，生存时间，协议，首部校验和，填充字段，源 IP，目标 IP。用 struct ip\_header 类型的结构体指针去做一个指针强制类型转换，然后就可以通过结构体的字段偏移提取出 ip 首部的各个字段。
4. **TCP 首部主要包含哪些信息？在分析时如何提取？** 源端口号、目的端口号、序列号、确认号、数据偏移、URG/ACK/PSH/RST/SYN/FIN 标识、窗口大小、校验和、紧急指针、选项。用 struct tcp\_header 去提取。
5. **大流量网络环境下如何提高捕包系统的性能？**在操作系统的内核中完成，效率很高。使用了数据缓存机制，使捕获数据包缓存存在内核中，达到一定数量再传递给应用程序。零拷贝技术。
6. **如何提高单机捕包能力？零拷贝捕包主要解决哪些核心问题？**减少系统调用和内存操作。数据报从网络设备到用户程序空间传递的过程中，减少数据拷贝次数，减少系统调用，降低 CPU 在这方面的负载。(2)数据先从网卡到内核，再到用户空间，这个拷贝过程仅仅起到“传输”作用，而数据包内容没有任何变化，零拷贝技术的核心：取消这两次拷贝。
7. **网站信息爬取的主要思路是什么？用伪代码描述单机网页信息爬取算法。**
- (1) 选取一部分种子 URL 放入队列。(2) 取出 URL，解析 DNS，得到主机 ip，下载存储入库。此外，将这些 URL 放进已抓取 URL 队列。(3) 分析已抓取 URL 队列中的 URL 得到新 URL，进入下一个循环。  
DS:URL\_QUEUE uqueue; History\_LIST hlist; PROCEDURE: Crawler(seed\_url) { in\_queue (uqueue, seed\_url); while (u=out\_queue(uqueue)) { wpage=http\_get(u); save wpage; for each url in wpage { if url not in hlist then in\_queue (uqueue, url); } }}
8. **多机爬取多个网站，核心解决哪些问题？主要策略有哪几种？**控制爬取范围、爬取深度、负载均衡、定期更新、爬取实时性；如何将爬取任务均匀的分配到各个机器上，如何让各个机器同步工作。按 URL 的散列值均匀分配抓取任务。深度优先；广度优先：在目前为覆盖尽可能多的网页，一般使用广度优先搜索方法。缺点：随着抓取网页的增多，大量的无关网页将被下载并过滤，算法的效率将变低；最佳优先：预测候选 URL 与目标网页的相似度，或与主题的相关性，并选取评价最好的一个或几个 URL 进行抓取。它只访问过网页分析算法预测为“有用”的网页。缺点：在爬虫抓取路径上的很多相关网页可能被忽略，因为最佳优先策略是一种局部最优搜索算法；部分的 PageRank 的策略
9. **用伪代码描述多机单网页信息爬取算法。**
- DS:URL\_QUEUE uqueue; History\_LIST hlist; PROCEDURE: Crawler(seed\_url) { in\_queue (uqueue, seed\_url); while (u=out\_queue(uqueue)) { wpage=http\_get(u); save wpage; for each url in wpage { if url not in hlist then { n=hash(url) mod n; if (n==my\_node\_ID) then in\_queue (uqueue, url); else sendurl(url,n); } }}
10. **简述基于主动获取技术的网站信息监测系统的构成。** 搜索器、索引器、检索器和用户接口
11. **简述 PageRank 算法的主要思路。** 如果网页 B 存在一个指向网页 A 的连接，则表明 B 的所有者认为 A 比较重要，从而把 B 的一部分重要性得分赋予 A。这个重要性得分值为：PR(B)/(L(B)，PR(B)为 B 的 PageRank 值,L(B)为 B 的出链数；A 的 PageRank 值为一系列类似于 B 的页面重要性得分值的累加。PR(A)=(PR(B)/L(B)+PR(C)/L(C)+...)q + (1-q) q 为阻尼系数(一般为 0.85);在任意时刻,用户到达某页面后并继续向后浏览的概率。
12. **网页抓取的搜索策略有哪些？各自适合哪些应用场景？**深度优先在很多情况下会导致爬取的陷入(trapped)问题，目前常见的是**广度优先**和**最佳优先**方法。广度优先搜索策略是指在抓取过程中，在完成当前层次的搜索后，才进行下一层次的搜索。该算法的设计和实现相对简单。在目前为覆盖尽可能多的网页，一般使用广度优先搜索方法。缺点：随着抓取网页的增多，大量的无关网页将被下载并过滤，算法的效率将变低。最佳优先：按照一定的网页分析算法，预测候选 URL 与目标网页的相似度，或与主题的相关性，并选取评价最好的一个或几个 URL 进行抓取。它只访问过网页分析算法预测为“有用”的网页。缺点：在爬虫抓取路径上的很多相关网页可能被忽略，因为最佳优先策略是一种局部最优搜索算法。
13. **试比较主动获取与被动获取技术的异同。** 如何网上违法活动的追踪溯源？收集证据,数据处理,信息提取,数据分析,溯源追踪,联合打击

第三章 串匹配算法

1. **串匹配算法分类和评价指标有哪些？**分类指标：匹配的模式数目(单模式/多模式)、匹配方式(精确/近似)、匹配的具体内容(单词/字符/正则表达式)、实时性(实时文本:动态更新~网络入侵检测/非实时文本:被查找文本是静态的-搜索引擎查找的数据) 评价指标：匹配次数、时间复杂度、是否需要回溯
2. **串匹配算法优化的主要思路？**KMP 算法是充分利用已经比较过的字符信息来提高效率，而 BM 算法则是从利用匹配失败时获得的信息出发提高效率,基于(倒序)自动机/哈希/状态合并(自动机)/分治/预处理/并行的匹配算法，
3. **试比较不同串匹配算法的优缺点。** KMP 算法时间复杂度 o(m+n)，并且不需要回溯；BM 算法时间复杂度 o(mn)，但是实际效果却常常是最快的。二者:每一次比较后，有明确的信息记录下一次比较的位置.AC 复杂度 o(n),时间复杂度与关键字的数目和长度无关,扫描文本时不需要回溯,wm 算法的复杂度 o(B\*N/m):B:块字长度 N:文本长 m:模式最短长度,对最短模式长度敏感。
4. **KMP 算法的基本流程？**能够针对实际数据给出求解步骤。先求 next 数组：next[i]=[-1,i=0;max{k|(1<k<i&& x1...xk-1== xi-k+1...xi-1' )} ,此集合不空;0,其他情况] (优化：若 p[i]=p[next[i]],那么 next[i]=next[next[i]]) 每次的跳跃长度: i-next[i] (yj 和 xi 匹配失败,就和 x\_next[i]匹配) 注意: next 数组最后比模式串多一个
- ```
void preKmp-opt(char *x, int m, int kmpNext[]) {int i, j; i = 0, j = -1; kmpNext[0] = -1; while (i < m) { while (i > -1 && x[i] != x[j]) j = kmpNext[j]; i++; j++; if (x[i] == x[j]) kmpNext[i] = kmpNext[j]; else kmpNext[i] = j; } for (i from 0 to m-1) if (x[i] == x[i] && next[i] == 0) next[i] = -1}
```
- ```
void kmp-opt(char *x, int m, char *y, int n) { int i, j, kmpNext[XSIZ]; preKmp(x, m, kmpNext); i = j = 0; while (j < n) { while (i > -1 && x[i] != y[j]) (i = kmpNext[i]); if (i == -1 then i=0; j++; ++j; if (i >= m) {OUTPUT(i - j); i = kmpNext[i];}}
```
5. **BM 算法的基本流程？** [取 max]（主要数据结构）能够针对实际数据给出求解步骤。算法从正文左端开始扫描,对齐后从模式最右端开始自右向左诸字符比较。在不匹配（或完全匹配）时,用两个预先计算的函数 bad character 和 good suffix 来确定将模式向右移动的距离。计算出字符串集中每个字符的偏移值 bmBC[i],对于未在模式中出现的字符，偏移为 m,否则取 m-i-1,(其中 i 为字符在模式中的位置,即取此字符在模式中出现的最右位置和文本中此字符位置对齐；若字符未在模式中出现,则可将模式向前推 m 个字符位置,但是,在某些情况下这种偏移可能是负的,实际的偏移值取得是两者中值较大者 “坏字符转移的时候转移到的是最右边的和坏字符相同的字符处,而且坏字符数组是针对每一个不同的字符的移位数为下标” eg: GCAGAGAG bmBC[C]: A-1 C-6 G-2 T-8; suffix[i]:10020408;bmGs[i]:7727471 shift(坏字符) = bmBC[T[i]]-(m-1-i) \*suffix[i]数组是为了找最长后缀长度在 i 为末尾分界的情况下)
- bmGs 构造方法: void preBmGs(char \*x, int m, int bmGs[]) {int i, j, suffix[XSIZ]; suffixes(x, m, suffix); for (i = 0; i < m; ++i) bmGs[i] = m; j = 0; for (i = m - 1; i >= 0; --i) if (suffix[i] == i + 1) for (j <= m - 1 - i; ++j) if (bmGs[j] == m) bmGs[j] = m - 1 - i; for (i = 0; i <= m - 2; ++i) bmGs[i] = m - 1 - i; } bmGs 三种情况:模式串中没有子串匹配上好后缀,但找不到一个最大前缀; 模式串中没有子串匹配上好后缀,但找到一个最大前缀; 模式串中有子串匹配上好后缀。 bmBC 构造: bmBC[m] = m - 1 - i; 没出现就是 m 且 i≠m-1
- AC 算法在千万量级的模式集上如何并行加速?:分组多线程并行,并行 tries 树,分布式计算(多计算节点),GPU 加速



