

Traffic Characterization of Large-scale Access Network and Implications

Li-hua MIAO, Wei DING

College of Computer Science and Engineering, Southeast University Nanjing, Jiangsu, China
{lh-miao, wding}@njnet.edu.cn

Abstract—this paper focus on analyzing traffic characterization of a large-scale access network based on active IP address count and link utilization. We find out that the behavior of internal and external hosts is asymmetric, the utilization of internal IP addresses per day is about 18.64% and the maximum link utilization is about 53%. Based on these, we proposed a scheme for detecting traffic anomalies which is proved to be powerful than another work.

Keywords—access network; traffic characterization; active IP address; link utilization; traffic anomalies detection

大规模接入网边界流量特性分析及其应用

缪丽华¹, 丁伟²

东南大学计算机科学与工程学院, 江苏南京, 中国, 211189
{lh-miao, wding}@njnet.edu.cn

【摘要】本文从活跃 IP 地址数目和链路带宽两个角度分析大规模接入网边界的流量特性, 得出接入网内外部主机活跃行为不对称, 一天内网内 IP 地址使用率仅达 18.64%, 链路带宽最大利用率为 53%等结论. 通过对流量特性的分析, 本文提出一个异常流量报警方案, 并与其他研究成果进行了比较.

【关键词】接入网; 流量特性分析; 活跃 IP 地址; 链路带宽; 异常流量发现

1 引言

广义地讲, 网络流量特性分析包括采集统计数据, 统计数据建模和流量特征化等多个方面. 网络流量特性分析对异常流量报警, 用户行为分析, 网络管理等多个研究领域起着重要的作用. 目前, 很多组织都在从事网络流量特性分析的相关研究工作, 如 CAIDA, NLANR 等.

为了特定的目的, 研究人员往往从不同角度对流量特性进行研究. K.C. Claffy[1]选择“延迟”, “丢包”, “使用率”等多个测度对 NSFNET 骨干网流量进行分析, 得出多个基于流的流量特性结论. 李永才[2]等选择“包到达时间间隔”, “包长”等测度对 NLANR 的一组 OC192 链路流量进行分析, 得出“高速链路流量仍具有很强的突发性”等多个结论. Phillipa Gill 等[3]选择使用率等测度对 YouTube 视频流量特性进行了分析, 并将其与传统 web 和媒体流的特性进行了比较. Ilmari Juva 等[4]对 Funet 网的流量特性进行了研究, 以期更好地规划网络流量. 杨艳等[5]通过对 TCP 双向流特性的理论分析, 提出了一个用于接入网边界异常流量报警的指标.

本文从活跃 IP 地址数目及链路带宽两个角度, 选择多个测度对大规模接入网边界流量的特性进行分析. 本研究的侧重点是 (i) 接入网内部与外部主机活跃的规律; (ii) 接入网资源的使用情况与管理; (iii) 接入网边界异常流量的监测.

本文使用 CERNET 江苏省网边界的流量进行研究, 主要发现有: (i)该接入网内部/外部主机的活跃行为存在不对称性; (ii)该接入网边界出/入两方向的 IP 流超时条件的设置应不同; (iii)该接入网边界流量仍符合周期性特性; (iv)该接入网内/外 IP 地址的使用符合 ON/OFF 模型; (v)网内部分 IP 地址的 ON/OFF 模型存在同步现象; (vi)网内 IP 地址的使用存在局部性; (vii)该接入网边界链路带宽在一天中流量最高峰时达到 53%的使用率; (viii)一天中网内 IP 地址平均使用率仅为 18.64%. 这些结论有助于更好地管理该接入网. 此外, 本文提出了一个易于实现的可用于接入网边界异常流量报警的方案, 该方案较文献[5]中的方案功能性更强.

本文第二章主要介绍实验数据和实验方案; 第三章介绍基于活跃 IP 地址数, 链路带宽两类测度的实验结果及结论; 第四章讨论了网内 IP 地址使用率并提出异常流量报警方案; 第五章对全文进行了总结并指出未来的工作方向.

2 测量方法

2.1 数据采集

作者所在的 CERNET 华东北地区网络中心从 2005 年开始对 CERNET 江苏省网与 CERNET 国家骨干网之间的 10G 信道进行监测. 网络中心研发的 WATCHER 采集系统通过分光技术, 按照 1/4 流抽样, 在选定时段对信道上的报

文进行采集并按时戳升序存储为 *IP Trace* 文件[6, 7]. 采集时, WATCHER 系统将该报文的方向(接入网内到网外, 或相反)记录在报文头部. 根据报文的方向标记, 可判别报文的源宿 IP 地址分别属于接入网内或网外.

本文的实验数据为 5 个持续时间为 24 小时的 *IP Trace* 文件, 具体信息见表 1. 将持续时间为 (t_0, t_1) 的 *IP Trace* 文件记为 *IP Tracei(t_0, t_1)*.

Table 1 IP Trace Information

名称	采集日期	持续时间	数据量 (GB)
IP Trace1(0,24)	2009-12-17(周四)	00:00~24:00	856
IP Trace2(0,24)	2010-03-18(周四)	00:00~24:00	913
IP Trace3(0,24)	2010-06-15(端午节假期)	00:00~24:00	844
IP Trace4(0,24)	2010-09-11(周六)	00:00~24:00	831
IP Trace5(0,24)	2011-01-16(周日)	00:00~24:00	873

2.2 活跃 IP 地址集

定义 1: 对 *IP Tracei(t_0, t_1)* 及 $t \in (t_0, t_1)$, 在 (t_0, t) 内所有报文中源 IP 地址集合称为 *IP Tracei* 在 t 时刻的活跃 IP 地址集合, 记为 $AIP(Tracei, t)$. 根据地理位置可分为网内活跃 IP 地址集合和网外活跃 IP 地址集合, 分别记为 $IAIP(Tracei, t)$ 和 $EAIP(Tracei, t)$.

定义 2: 对 *IP Tracei(t_0, t_1)* 及 $t \in (t_0, t_1)$. 设不活跃超时时间为 t' , 在 (t_0, t) 内所有报文中未超时的源 IP 地址集合称为 *IP Tracei* 在 t 时刻未超时活跃 IP 地址集合, 记为 $NTASIP(Tracei, t, t')$; 称已超时的源 IP 地址列表为 $TASIP(Tracei, t, t')$. 注意, $TASIP(Tracei, t, t')$ 中存在重复的 IP. 同理, 可分为 $INTASIP(Tracei, t, t')$ 和 $ENTASIP(Tracei, t, t')$, $ITASIP(Tracei, t, t')$ 和 $ETASIP(Tracei, t, t')$.

2.3 测度

2.3.1 基于活跃 IP 地址数的测度

定义 3: 称 $IAIP(Tracei, t)/EAIP(Tracei, t)$ 的大小分别为 t 时刻 *IP Tracei* 中“类型 1 的累计网内/网外活跃 IP 数”, 记为 $T_{in}(i, t) / T_{ex}(i, t)$.

定义 4: 称 $INTASIP(Tracei, t, t')$ 与 $ENTASIP(Tracei, t, t')$ 的大小总和为 t 时刻 *IP Tracei* 中“类型 2 的累计网内活跃 IP 数”, 记为 $T_{in}(i, t, t')$. 称 $ITASIP(Tracei, t, t')$ 与 $ETASIP(Tracei, t, t')$ 的大小总和为 t 时刻 *IP Tracei* 中“类型 2 的累计网外活跃 IP 数”, 记为 $T_{ex}(i, t, t')$.

定义 5: 称 $INTASIP(Tracei, t, t')/ENTASIP(Tracei, t, t')$ 的大小分别为 t 时刻 *IP Tracei* 中“网内/网外实际活跃 IP 数”, 记为 $A_{in}(i, t, t') / A_{ex}(i, t, t')$.

2.3.2 基于链路带宽的测度

定义 6: 对 *IP Tracei(t_0, t_1)*, 将 (t_0, t_1) 划分为 n 个长度相等的时间片, 在第 t 个时间片内接入网接收到的报文总数, 称为第 t 个时间片内 *IP Tracei* 中“入报文数”, 记为 $P_{in}(i, t, n)$; 接入网接收到的字节总数, 称为“入字节数” $B_{in}(i, t, n)$; 接入网发出去的报文总数, 称为“出报文数” $P_{out}(i, t, n)$; 接入网发出去的字节总数, 称为“出字节数” $B_{out}(i, t, n)$.

定义 7: 对 *IP Tracei(t_0, t_1)*, 将 (t_0, t_1) 划分为 n 个长度相等的时间片, 称 $(B_{in}(i, t, n) + B_{out}(i, t, n)) / (P_{in}(i, t, n) + P_{out}(i, t, n))$ 为第 t 个时间片内 *IP Tracei* 中“总字节数/总报文数”, 记为 $\lambda(i, t, n)$.

2.4 实验方案

为了有效地输出测度计算结果, 本实验中将 24 小时划分为 1440 个长 1 分钟的时间片.

2.4.1 基于活跃 IP 数测度的计算方案

● $T_{in}(i, t), T_{ex}(i, t)$ 计算方案:

输入: *IP Tracei(0,24)* $i=1, 2, 3, 4, 5$

输出: $T_{in}(i, t), T_{ex}(i, t)$ 长 1440 的统计值序列.

初始化: 建立网内/网外活跃 IP 地址表 *Internal_table/External_table*; 将 $T_{in}(i, t), T_{ex}(i, t)$ 清零.

while(还有报文未读) {

 取一个未读报文;

 if(该报文的源 IP 地址为初次到达)

 根据报文方向标识将 $T_{in}(i, t) / T_{ex}(i, t)$ 加 1;

 if(当前时间片结束)

 输出 $T_{in}(i, t), T_{ex}(i, t)$; }

● $T_{in}(i, t, t'), T_{ex}(i, t, t'), A_{in}(i, t, t'), A_{ex}(i, t, t')$ 计算方案:

输入: *IP Tracei(0,24)* $i=1, 2, 3, 4, 5$

输出: $T_{in}(i, t, t'), T_{ex}(i, t, t'), A_{in}(i, t, t'), A_{ex}(i, t, t')$ 长 1440 的统计值序列.

初始化: 建立网内/网外实际活跃 IP 地址表 *Internal_table/External_table*; 将 $T_{in}(i, t, t'), T_{ex}(i, t, t'), A_{in}(i, t, t'), A_{ex}(i, t, t')$ 清零; 设不活跃超时时间为 t' 为 a 秒.

while(还有报文未读) {

 取一个未读报文;

 if(该报文的源 IP 地址初次到达)

 根据报文方向标识将 $T_{in}(i, t, t') / T_{ex}(i, t, t'), A_{in}(i, t, t') / A_{ex}(i, t, t')$ 加 1;

 if(该报文的源 IP 地址已存在且已超时)

 根据报文方向标识将 $A_{in}(i, t, t') / A_{ex}(i, t, t')$ 减 1;

 if(当前时间片结束)

 输出 $T_{in}(i, t, t'), T_{ex}(i, t, t'), A_{in}(i, t, t'), A_{ex}(i, t, t')$;

$$A_{ex}(i, t, t'); \}$$

2.4.2 基于链路带宽的测度计算方案

输入: IP Trace_i(0,24) i=1, 2, 3, 4, 5

输出: $P_{in}(i, t, n)$, $B_{in}(i, t, n)$, $P_{out}(i, t, n)$, $B_{out}(i, t, n)$, $\lambda(i, t, n)$ 长 1440 的统计值序列.

初始化: 将 $P_{in}(i, t, n)$, $B_{in}(i, t, n)$, $P_{out}(i, t, n)$, $B_{out}(i, t, n)$, $\lambda(i, t, n)$ 清零, n=1440

while(还有报文未读) {

 取一个未读报文, 其字节数为 pkt_length;

 根据报文方向标识将 $P_{in}(i, t, n)$ / $P_{out}(i, t, n)$ 加 1, 将 pkt_length 累加到 $B_{in}(i, t, n)$ / $B_{out}(i, t, n)$;

 if(当前时间片结束)

 输出 $P_{in}(i, t, n)$, $B_{in}(i, t, n)$, $P_{out}(i, t, n)$, $B_{out}(i, t, n)$, $\lambda(i, t, n)$ 并将其清零; }

3 流量统计特性

3.1 基于 $T_{in}(i, t)$, $T_{ex}(i, t)$, $T_{in}(i, t, t')$, $T_{ex}(i, t, t')$ 的特性

本小节中, $T_{in}(i, t, t')$ 与 $T_{ex}(i, t, t')$ 的计算使用 $t'=64$ 秒的不活跃超时条件.

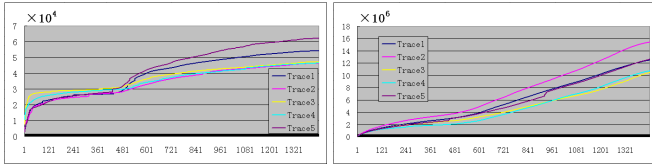


Fig. 1 (a) $T_{in}(i, t)$

Fig. 1 (b) $T_{ex}(i, t)$

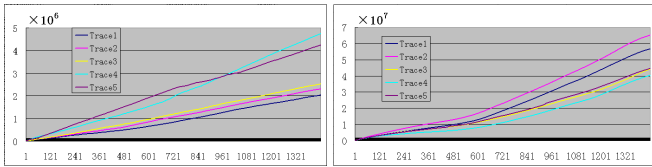


Fig. 1 (c) $T_{in}(i, t, t')$

Fig. 1 (d) $T_{ex}(i, t, t')$

由图 1(a)可知, 对 5 天实验数据, $T_{in}(i, t)$ 表现出相似的特性, 均在 1 点左右达到 0 点到 7:30 时间段的平稳值, 但到达稳定的速率和稳定集合的大小随实验数据的不同而有所不同. 这表明在 1:00~7:30 这段时间内, 新增加的网内活跃 IP 地址很少, 符合接入网内用户的上网习惯. 2010-06-15 为端午节假期, 到达稳定的速率及稳定值均是最大值. 这表明在假期里, 部分网民喜欢熬夜上网. 2010-09-11 到达稳定的速率及稳定值位居第二, 证明了上述结论. 虽然 2011-01-16 为周日, 但由于接近寒假, 网内不少学生用户已经回家, 导致其水平跟工作日接近. 在 7:30 以后, 新增的网内活跃 IP 数呈递增趋势. 2011-01-16 这天网内 IP 的使用率为最高. 由图 1(b)可知, $T_{ex}(i, t)$ 表现出了类似趋势.

由图 1(a)和(c)可知, 在 $t=24$ 时, 5 天实验数据的平均

$T_{in}(i, t, t')$ 约为 $T_{in}(i, t)$ 的 60 倍. 即在不活跃超时为 64 秒时, 一个网内活跃 IP 地址平均被重复计算了 60 遍左右. 这表明(i)网内 IP 地址的活跃行为仍符合 ON/OFF 模型, 且(ii)网内 IP 地址的使用存在局部性, 即部分 IP 地址被不断重复利用. 这符合网络配置的实际情况, 因为接入网内存在很多学校, 每个学校对外使用的 IP 地址较为固定, 且数目较少. 由图 1(b)和(d)可知, 在 $t=24$ 时, 5 天实验数据的平均 $T_{ex}(i, t, t')$ 平均约为 $T_{ex}(i, t)$ 的 4 倍. 即在超时条件为 64 秒时, 一个网外活跃 IP 平均约被重复计算了 4 遍. 这表明(i)网外 IP 地址的活跃行为也符合 ON/OFF 模型但 ON/OFF 周期较长, (ii)网外 IP 地址的使用也存在一定的局部性, 部分 IP 地址被重复利用. 上述结论表明网内/网外主机活跃行为不对称性, 但均符合 ON/OFF 模型. 此外, 网内 IP 地址的使用存在局部性.

3.2 基于 $A_{in}(i, t, t')$, $A_{ex}(i, t, t')$ 的特性

本小节中, $A_{in}(i, t, t')$ 与 $A_{ex}(i, t, t')$ 的计算分别使用 $t'=64$ 秒, 16 秒, 1 秒的不活跃超时条件.

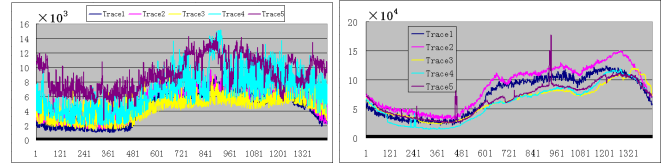


Fig. 2(a) $A_{in}(i, t, t')$, $t'=64$ sec

Fig. 2(b) $A_{ex}(i, t, t')$, $t'=64$ sec

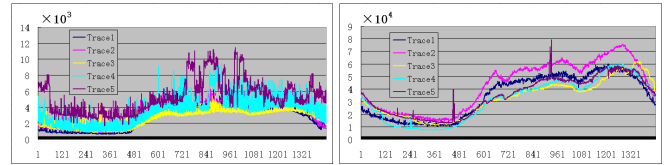


Fig. 2(c) $A_{in}(i, t, t')$, $t'=16$ sec

Fig. 2(d) $A_{ex}(i, t, t')$, $t'=16$ sec

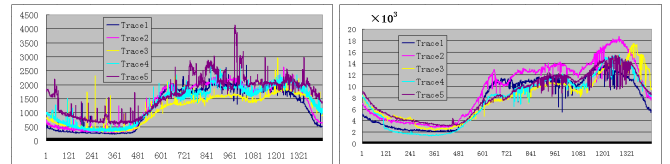


Fig. 2(e) $A_{in}(i, t, t')$, $t'=1$ sec

Fig. 2(f) $A_{ex}(i, t, t')$, $t'=1$ sec

由图 2 可知, 5 天实验数据的 $A_{in}(i, t, t')$ 和 $A_{ex}(i, t, t')$ 表现出了相似的特性. 网内/网外实际活跃的 IP 数目均具有季节性(周期性). 由图 2(a), (c)和(e)可知, 当 t' 变短时, 网内实际活跃 IP 数时间序列中突变点变少. 突变点的存在, 可能是由网内存在大量同时活跃或同时超时的 IP 地址引起的(即部分网内 IP 地址的 ON/OFF 模型同步). t' 越长, 同时活跃或同时超时的 IP 地址数越多. 由图 2(b), (d)和(f)可知, t' 对 $A_{ex}(i, t, t')$ 趋势的影响较小. $t'=1$ 秒时, $A_{ex}(i, t, t')$ 序列的突变点增多, 原因可能是超时时间过短. 此外, $A_{ex}(i, t, t')$ 平均约为 $A_{in}(i, t, t')$ 的 10 倍左右, 说明网内主机与网外主

机间存在并发连接,这与网内 IP 地址的重用有一定关系。

由以上结论可知,网内/网外主机活跃行为存在不对称性;网内 IP 地址发起的流超时条件设置与网外 IP 地址发卡的流超时设置应不同。

3.3 基于链路带宽类测度的流量特性

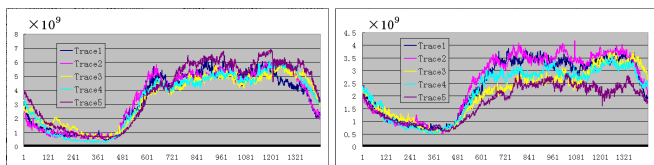


Fig. 3(a) $B_{in}(i, t, n)$

Fig. 3(b) $B_{out}(i, t, n)$

由图 3 所示可知,本接入网边界的流量仍符合具有季节性。在一天流量的高峰期,监测链路上的进出总字节数约为 5.3Gbps(注意:本文 IP Trace 为 1/4 流抽样获得)。监测链路的总容量为 10Gbps,因此链路带宽利用率最高达 53%左右。 $P_{in}(i, t, n)$, $P_{out}(i, t, n)$ 的趋势与 $B_{in}(i, t, n)$, $B_{out}(i, t, n)$ 的类似,省略。 $\lambda(i, t, n)$ 的计算结果将在 4.2 节中讨论。

4 讨论

4.1 网内 IP 地址使用率

2011 年 2 月 3 日(美国当地时间),互联网编号分配机构(IANA)在迈阿密举行新闻发布会,正式宣布最后五个 IPv4 地址大组分配完毕。各地区地址分配机构在分配完所持地址后,将不再有可分配的 IPv4 地址[8]。发布会预测,亚太地区将成为 IPv4 地址最早耗尽的地区。CERNET 江苏省网内的 IP 地址总数为 1098368 个。5 天实验数据中活跃源 IP 的平均最大值为 51180,则一天内网内 IP 的平均最大利用率约为 18.64%(1/4 流抽样)。因此,解决 IP 地址紧缺的一个可能方案是将目前未使用的 IP 地址回收并重新分配。

4.2 异常流量报警方案

定义 8: 对 IP Trace $i(t_0, t_1)$, 将 (t_0, t_1) 划分为 n 个长度相等的时间片,称 $B_{in}(i, t, n) / P_{in}(i, t, n)$ 为第 t 个时间片内 IP Trace i 中“入字节数/入报文数”,记为 $\lambda_{in}(i, t_0, n)$ 。称 $B_{out}(i, t, n) / P_{out}(i, t, n)$ 为“出字节数/出报文数”,记为 $\lambda_{out}(i, t_0, n)$ 。计算方案类似于 $\lambda(i, t, n)$ 。

为了满足实时性,异常报警的时间粒度不宜过短。因此,将 2.4.2 节中计算方案的时间片修改为 5 分钟,即 $n=288$ 。本文的异常流量报警方案为(i)将 t 时间片内 $\lambda(i, t, n)$ 的值与最近的正常历史时间片的 $\lambda(i, t_0, n)$ 比较,如果变化幅度超过阈值 α ,则判定 t 时间片流量异常;(ii)使用 $\lambda_{in}(i, t_0, n)$, $\lambda_{out}(i, t_0, n)$ 判断异常流量的方向(原理与(i)类似)。

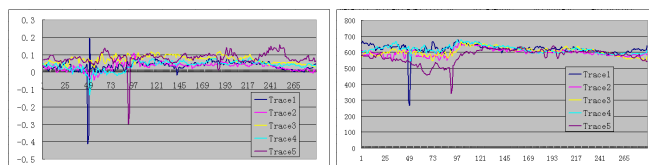


Fig. 4(a) metric in reference[5]

Fig.4(b) $\lambda(i, t, n)$

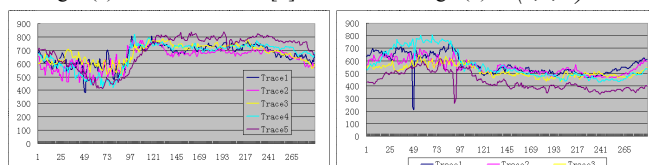


Fig.4(c) $\lambda_{in}(i, t_0, n)$

Fig.4(d) $\lambda_{out}(i, t_0, n)$

图 4(a)为文献[5]中测度的计算结果。由图 4(a)和(b)可知, $\lambda(i, t_0, n)$ 能够判断出 Trace1 的 48, 49, 50 三个时间片, Trace5 的 91, 92 两个时间片异常。由图 4(c)和(d)可知, Trace1 中 48, 49 两时间片出方向流量异常; 49, 50 两时间片入方向流量异常。Trace5 的 91, 92 两时间片出方向流量异常。因此,本文的方案不仅能够正确地识别出异常时间片,且能够判断异常流量的方向,比文献[5]的功能性更强。

5 结论和未来工作

本文从活跃 IP 地址数和链路带宽两个角度,共选择 13 个测度对接入网边界流量特性进行研究,验证了接入网边缘流量仍具有周期性,并得出一天中网内 IP 地址使用率仅为 18.64%,链路带宽利用率在流量高峰期达到 53%,网内外主机活跃行为不对称,网内 IP 地址使用存在局域性,IP 地址的使用符合 ON/OFF 模型等多个结论。此外,本文提出了一个简单易行的异常流量报警方案,既可准确报警异常还能检测异常流量的方向,比文献[5]的功能性更强。未来工作将围绕以下几个方面展开:① 研究网内外活跃 IP 超时条件的设置;② 研究网内 IP 地址 ON/OFF 模型同步的原因;③异常流量报警方案中报警阈值的设定。

References (参考文献)

- [1] K. C. Claffy. Internet Traffic Characterization. PhD thesis, University of California, San Diego, 1994.
- [2] Li Yongcai, Yan Jun, Liu Dapeng, Huang Jianhua. Analysis of Internet High Speed Link Traffic. Journal of Information Engineering University, 2008, 9(4).
- [3] Phillipa Gill, Martin Arlitt, Zongpeng Li, Anirban Mahanti. Youtube traffic characterization: a view from the edge. IMC'07.
- [4] I. Juva, R. Susitaival, M. Peuhkuri, and S. Aalto. Traffic characterization for traffic engineering purpose: Analysis of Funet data. NGI 2005.
- [5] Yang Yan, Ding Wei, Cheng Guang, Gong Jian. A Metric Model for Access Network Management. Computer Science, 2008, 35(5).
- [6] The IPTAS project. <http://iptas.edu.cn>
- [7] Li-hua MIAO, Wei DING, Hai-ting ZHU, Qing XIA. Cost-effective IP Trace Publishing Using Data Sketch. NCIS'11.
- [8] <http://www.iana.org/>