

面向智能手机的隐私保护系统

参赛队员：闵大为 任勇勋 张梦月 指导老师：付安民

南京理工大学



目录

1. 背景



2. 蠕虫检测思想的应用



3. 机器学习思想的应用



4. 总结



背景



在人们都越来越依赖于手机的强大功能时，有没有人想过，它实际并不安全。它就像是一个时时贴近我们的小偷，悄悄的将我们的个人隐私、手机通讯录、身份信息向外送走。

——《央视曝光安卓手机应用偷偷上传用户隐私》

一些解决方案

❖ 方案一：安装时注意权限

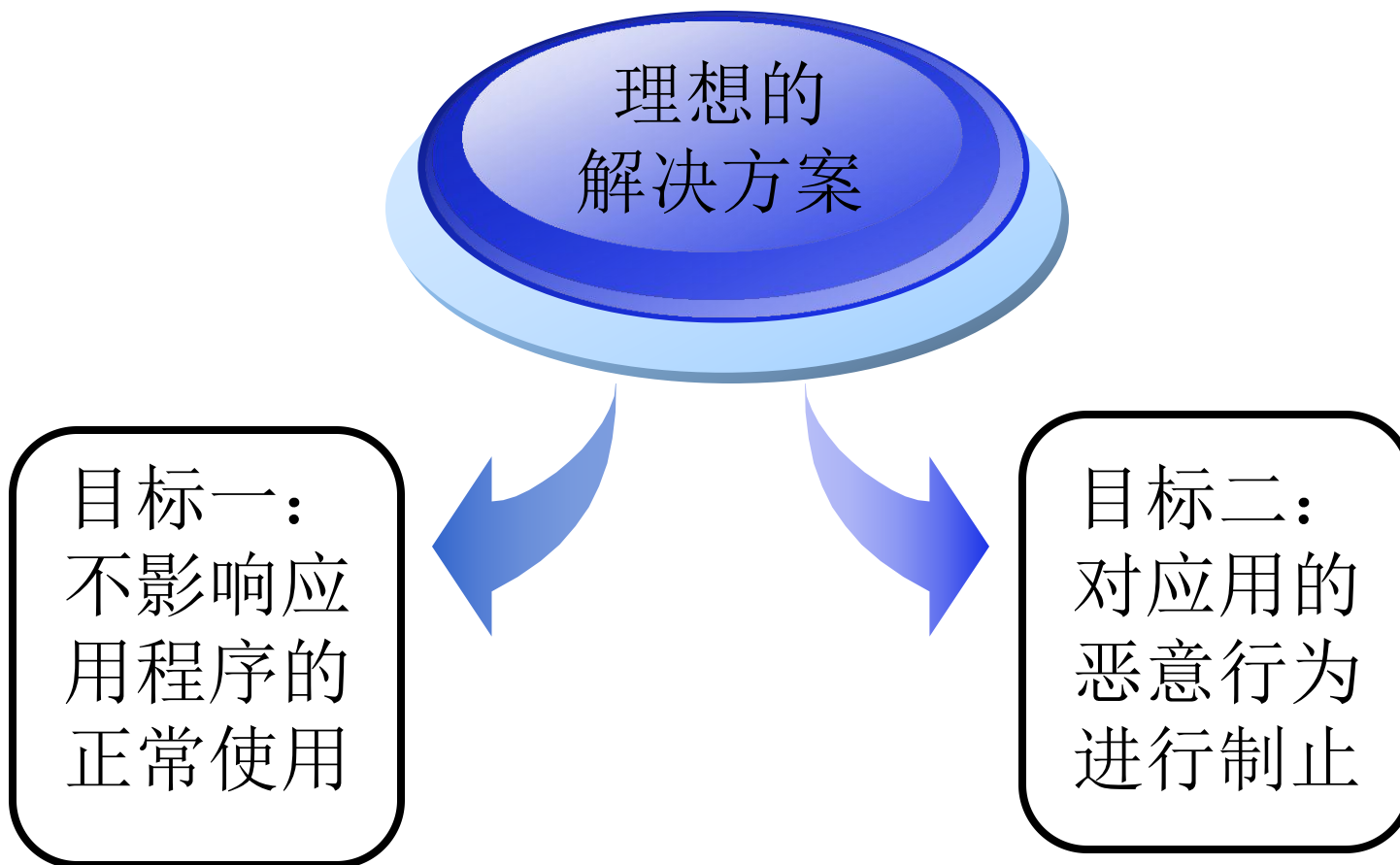
不足：一次性判断，如果用户不小心安装了，则无法提供保护；强迫安装

❖ 方案二：加密信息

不足：使用不便，如要保护通讯录则需要加密通讯录，要保护短信又必须加密短信，使用加密信息时需要用相关软件进行解密

❖ 方案三：禁止应用上网

不足：不让应用访问网络，应用无法实现自己的功能，例如即时通，浏览器等。所以这种方法不是很合理。





我们的解决方案：

将蠕虫检测的思想
移植到
手机隐私保护上



为什么可以将
蠕虫检测的思想
应用到
手机隐私保护上？

蠕虫检测

- ❖ **基本原理：**感染蠕虫后，网络流量会出现异常变化，以此进行检测和防护。
- ❖ **成熟性：**在这个方面已经有大量的研究，该检测思想已有一定的积累，比较稳定可靠。

[基于流流量的蠕虫检测思想](#) [相关论文\(共41篇\)](#) [百度学术](#)

[未知网络蠕虫检测与特征码自动生成的研究](#) 《华中科技大学》

被引频次: 2

[基于异常流里的蠕虫检测系统研究与实现](#) 《华中科技大学》

被引频次: 1

[一种基于进程流里行为的蠕虫检测系统](#) 计算机工程与科学 ISTIC PKU

被引频次: 1

[查看更多相关论文>>](#)

xueshu.baidu.com

手机上隐私泄露的特点

- ❖ 上传数据一般加密，且加密方式多种多样。

（说明基于文本内容的行为分析不太可行）

- ❖ 上传隐私数据时，**会提交大量数据**。而大部分的时间内，手机都是从服务器端取数据，如获取网页文字，图片等。不会连续地提交大量数据。

（说明隐私泄露过程中上传数据量会激增）

相似性

- ❖ 蠕虫检测中：由于感染蠕虫，网络流量会激增
- ❖ 手机隐私保护中：由于要上传隐私数据，相对于平常，提交的数据会激增



基于流量变化的检测模型

我们给出了手机端隐私泄露行为的定义：

在一定时间内，手机端向某个远程服务器提交的数据包总数超过正常值（阈值）时，则认定该手机中有应用在上传用户隐私数据。

模型一 算法

算法:

for each Δt

1 capture and analyze the data

2(1) calculate $\sum_{t=T}^{T+\Delta t} S_{IP_i, t} \quad (i = 1, 2 \cdots n)$

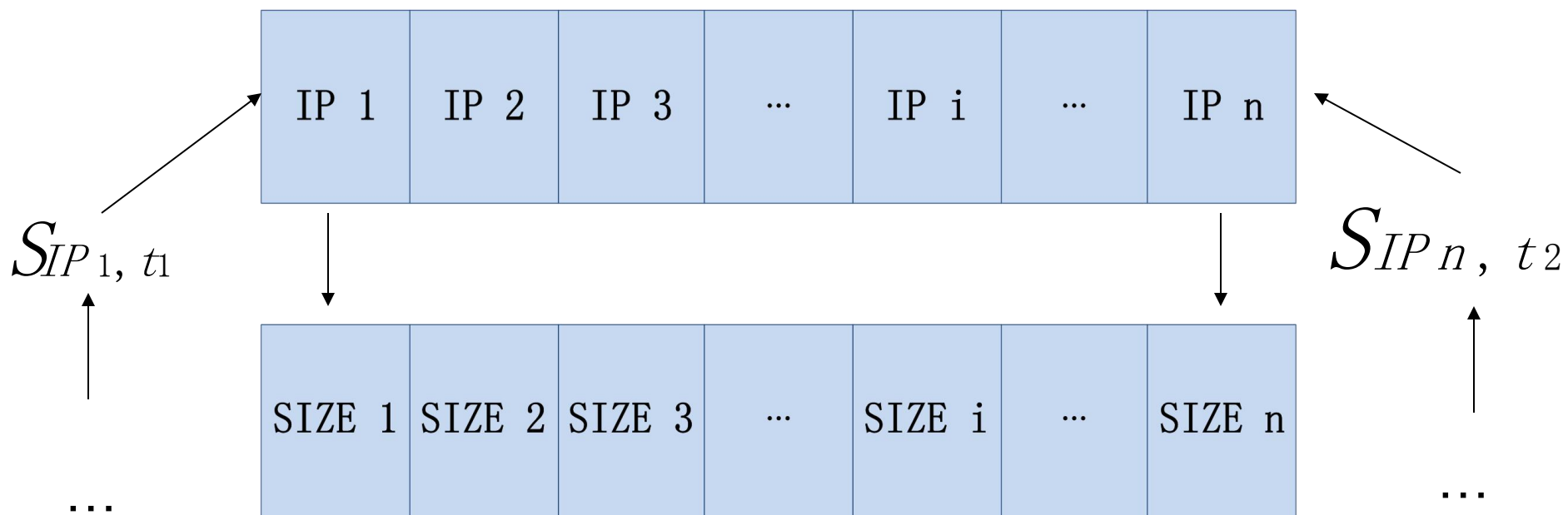
(2) if $\exists j(1 \leq j \leq n), \sum_{t=T}^{T+\Delta t} S_{IP_j, t} \geq T$

then take action

next Δt

IP地址和对应数据包大小变量储存格式

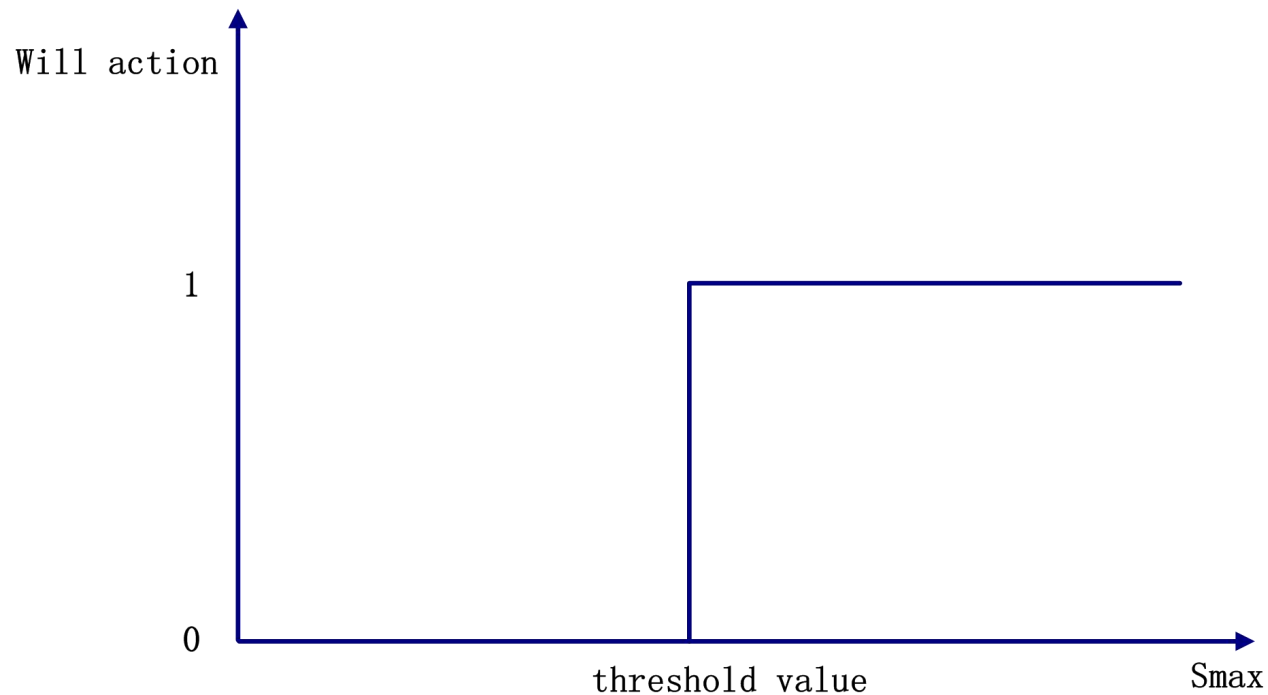
$$\sum_{t=T}^{T+\Delta t} S_{IP_i, t} \quad (i = 1, 2 \cdots n)$$



模型一 閾值

$$\text{if } \exists j(1 \leq j \leq n), \sum_{t=T}^{T+\Delta t} S_{IP_j, t} \geq T$$

then take action



模型一 新的问题



有些应用在使用的过程中会上传一定量的数据包，例如微信。

但是用户希望信任该程序，也就是说如果是微信在上传大量数据，我们不采取行动。

如何考虑用户信任的情况？

难点分析

数据包总和变量是属于哪个应用的？

要信任某个程序，我们必须知道超过阈值的数据包总和变量是属于哪个应用的，只有这样才能选择是否忽略。但是抓到的数据包并不提供与应用对应的信息。

IP 1	IP 2	IP 3	...	IP i	...	IP n
------	------	------	-----	------	-----	------

SIZE 1	SIZE 2	SIZE 3	...	SIZE i	...	SIZE n
--------	--------	--------	-----	--------	-----	--------



我们的解决方案：

将机器学习的思想
应用到
SIZE 变量和应用的对应上

机器学习的思想

❖ 机器学习的一种定义：

学习是系统所作的适应性变化，使系统在下次完成同样或类似的任务时更为有效。

- ## ❖ 特点：
1. 每次行动都会学习
 2. 下次更为有效

启发：学习的过程

- ❖ 借鉴机器学习的两个特点：
 1. 每次行动都会学习
 2. 下次更为有效
- ❖ 我们给出了基于学习的解决方案，该方案不是直接给出变量**SIZE** i 对应的应用，而是不断学习，下一次给出更好的答案。

模型二 问题分析

问题转换： 由于**SIZE**与**IP**是一一对应的， **SIZE** 和应用的对应问题就等价于**IP**与应用的对应问题。

关系分析：

(1) 一个应用可以使用多个**IP**地址：

$$G(App) = \{ Ip_1, Ip_2, \dots, Ip_n \}$$

(2) 一个**IP**地址只给一个应用使用：

$$F(Ip) = App$$

引入机器学习思想前后对比

❖ 假设：集合 A 表示{可疑应用1,可疑应用2...可疑应用 n }

$F(Ip)$ 表示求 Ip 在集合 A 中对应的可疑应用

❖ 引入机器学习思想之前：

由 Ip ， 集合 A ， 求 $F(Ip)$ 条件不足，难以解决。

❖ 引入机器学习思想之后：

由 IP ， 集合 A ， 以及”以前出现过的对应关系”， 求 $F(Ip)$ 。

因为“以前出现过的对应关系”这个条件的内容是在增加的，所以解决问题的可能性也在增加。

模型二 学习过程

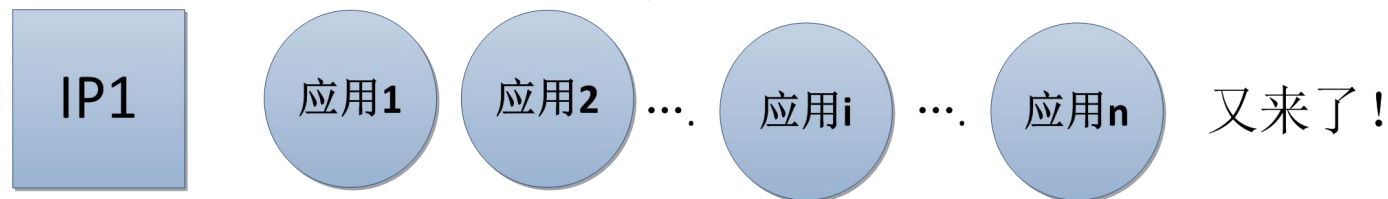
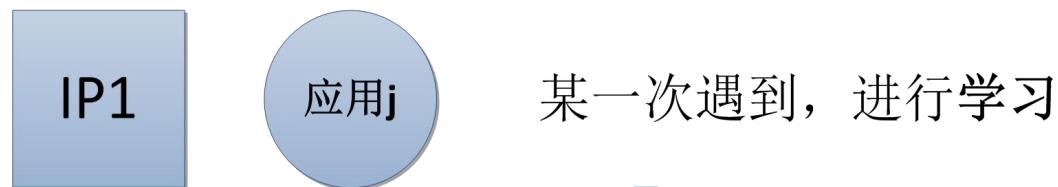
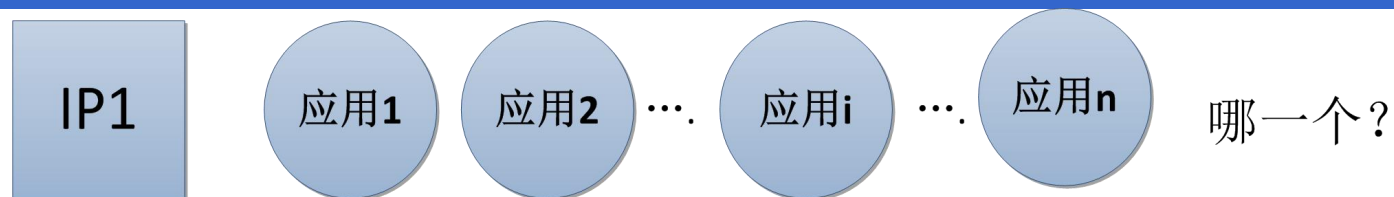
❖ 在用户使用的过程中，会出现：

$F(Ip)$ 所在的可疑应用集合A里只有一个应用，
也就是说出现了一一对应的情况：

$$F(Ip) = App$$

❖ 通过学习，我们可以把这个关系记录下来，下次再遇到这个Ip的时候就可以排除干扰项，找到目标了。

模型二 学习过程



迷惑不了我了, 是



模型二 算法

for each Δt

1 capture and analyze the data

2(1) calculate $\sum_{t=T}^{T+\Delta t} S_{IP_i, t} \quad (i = 1, 2 \cdots n)$

(2) if $\exists j(1 \leq j \leq n), \sum_{t=T}^{T+\Delta t} S_{IP_j, t} \geq T \{$

learn to find App

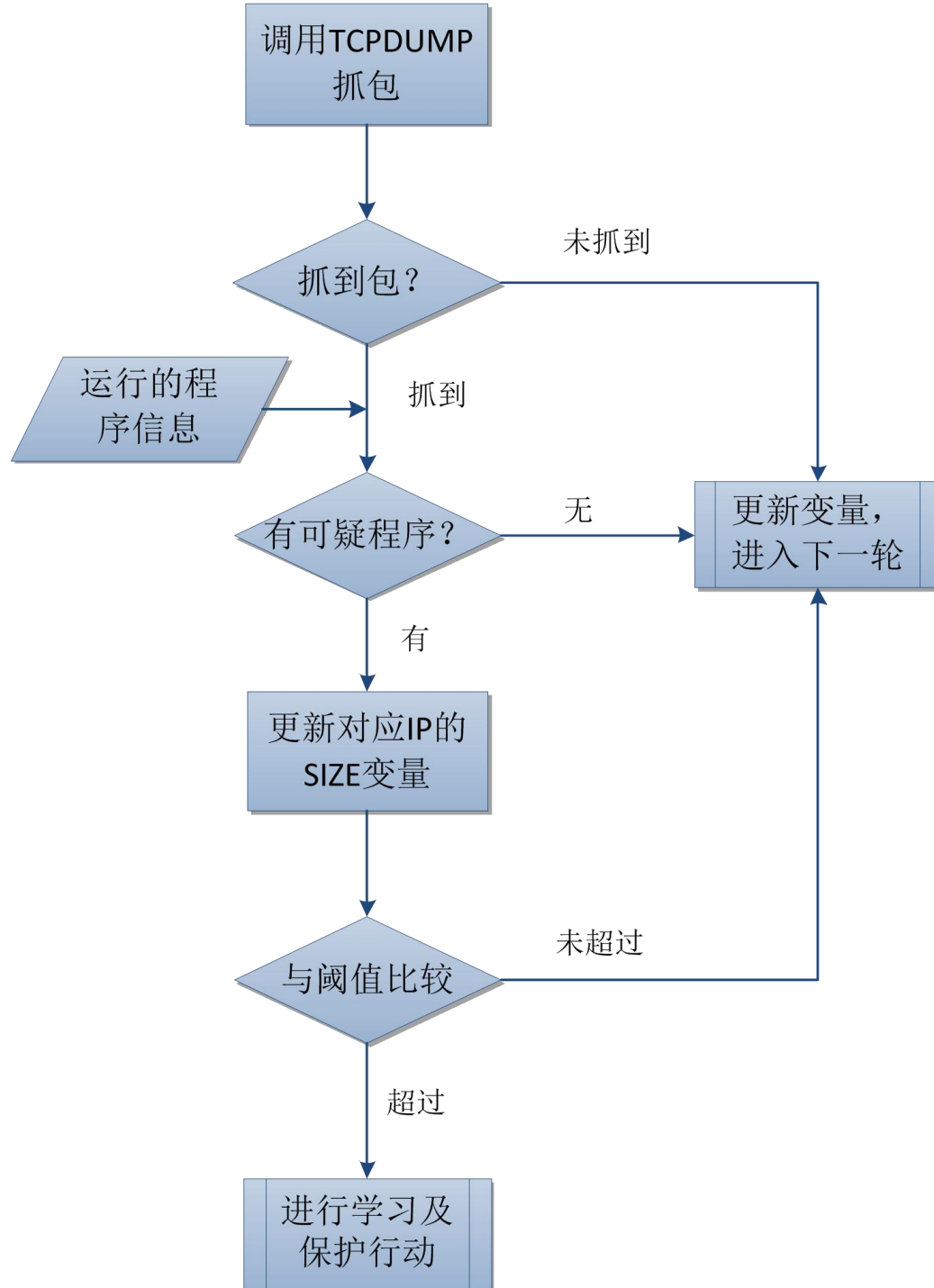
if trustless

then take action

}

next Δt

简化流程图



效果



总结

- ❖ 特色：与市场上的一些静态解决方案不同，该方案能够在不影响用户正常使用的前提下进行动态监控并提供保护。
- ❖ 创新一：将基于蠕虫检测的思想应用在手机隐私保护上。
- ❖ 创新二：将机器学习的思想应用在数据包大小和可疑应用的对应上。

Thank You !

