

云网分析与可视化

发掘网络数据的真正价值

向阳@云杉网络
QCon2017@北京



促进软件开发领域知识与创新的传播



关注InfoQ官方信息
及时获取QCon软件开发者
大会演讲视频信息



扫码，获取限时优惠



全球架构师峰会 2017 [深圳站]

2017年7月7-8日 深圳·华侨城洲际酒店

咨询热线：010-89880682



全球软件开发大会 [上海站]

2017年10月19-21日

咨询热线：010-64738142

云网：云数据中心网络

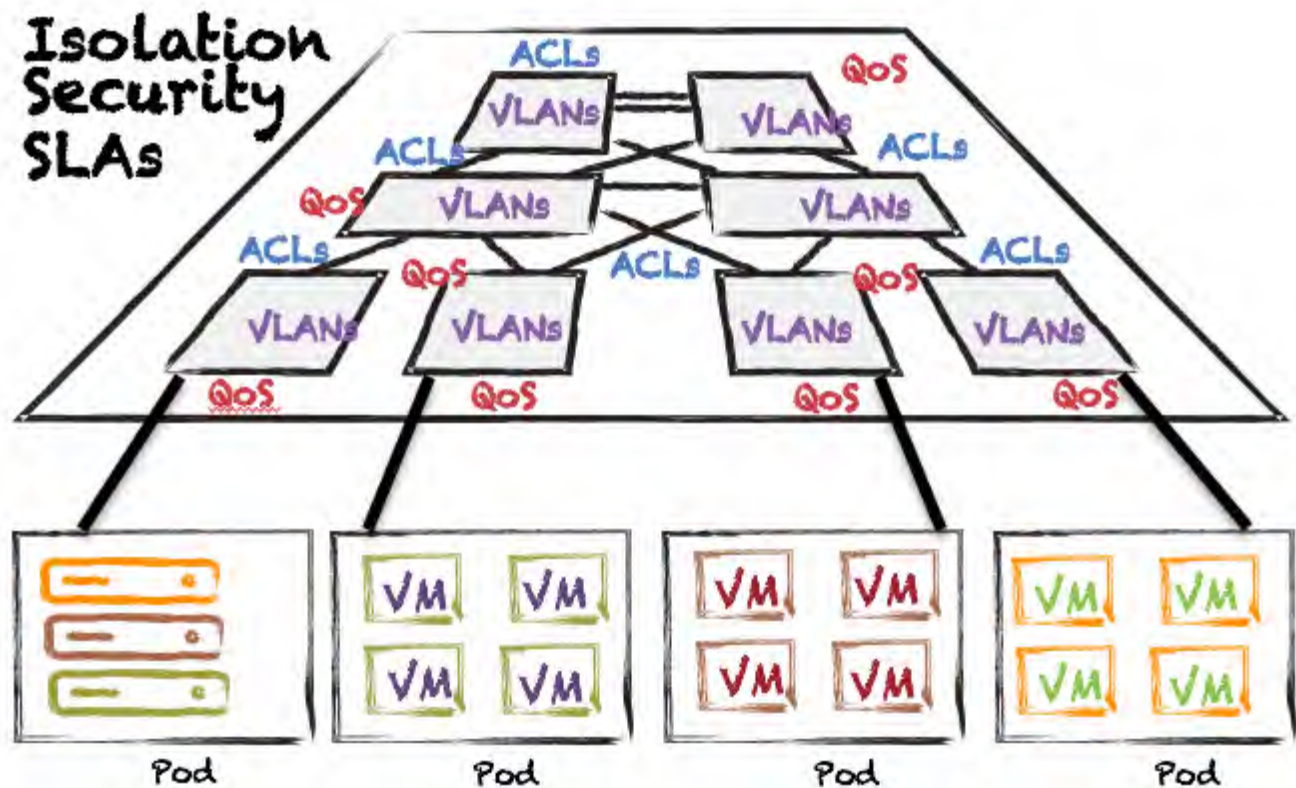
- 云网的变迁
- 云网运维的痛
- 摆脱云网运维的困境
 - 数据采集
 - 数据分析
 - 数据存储



云网的变迁

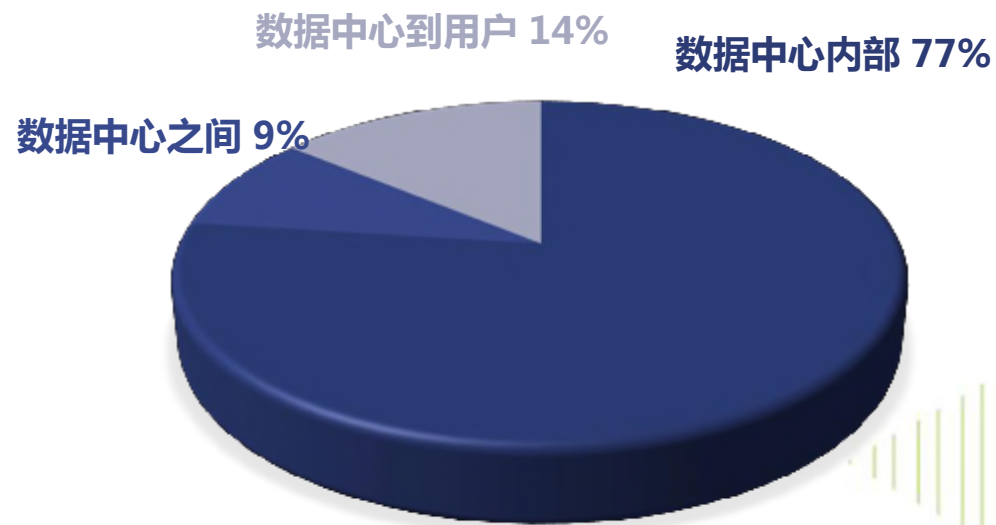
由实到虚

- 从**烟囱式**到**集中化**
 - 单租户独享 → 多租户共享
- 从**Underlay**到**Overlay**
 - VLAN
 - VXLAN
 - EVPN



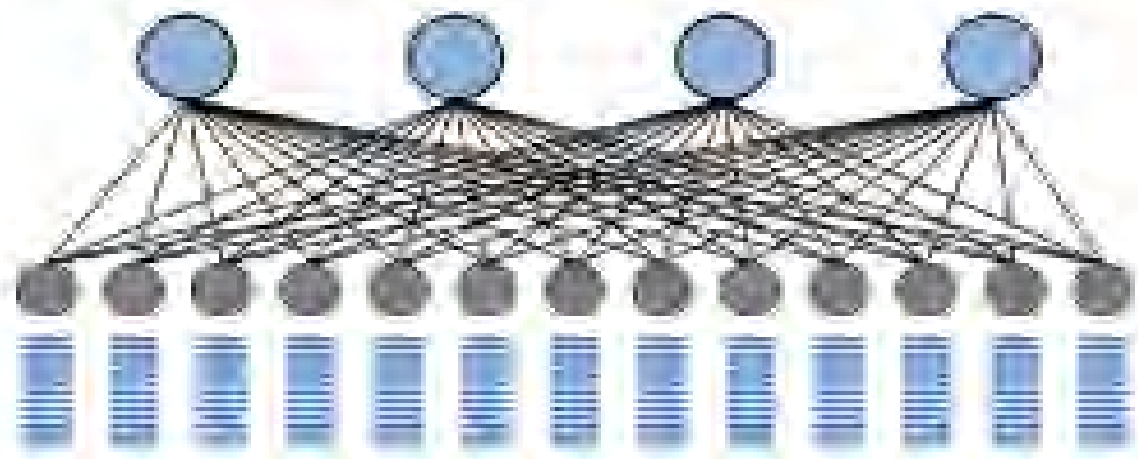
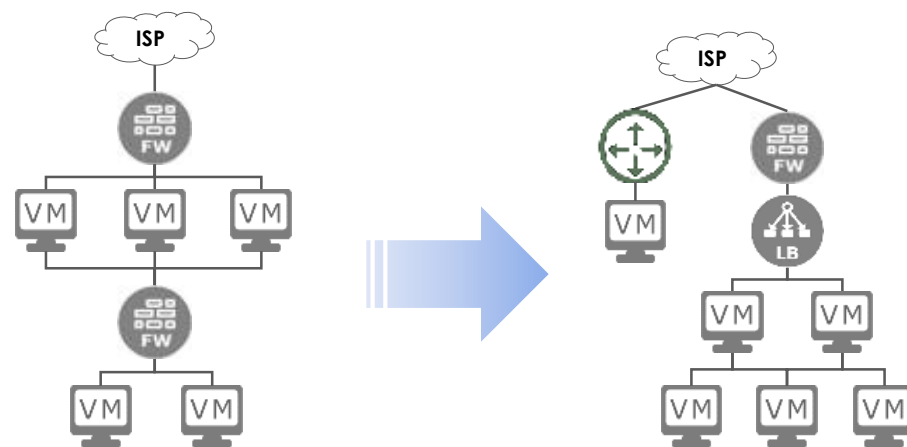
由外到内

- 2020年，云数据中心将处理**92%**的工作负载
- **数据中心东西向流量**将占据总流量的**77%**
 - 多租户共享资源池
 - 分布式系统
 - 数据备份



由静到动

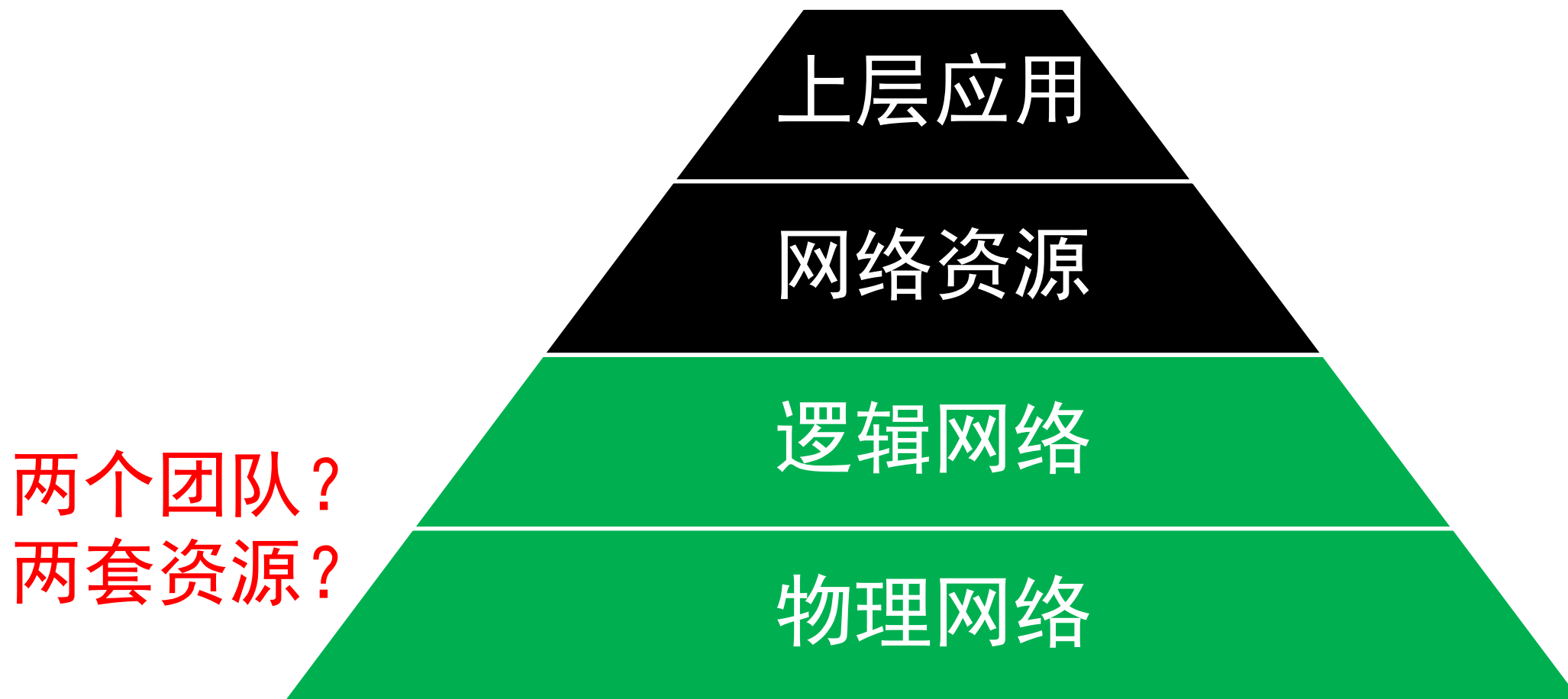
- 业务负载弹性伸缩，网络随之动态配置
- 物理基础网络（Fabric）极少变更
- 虚拟网络需要随业务秒级变更



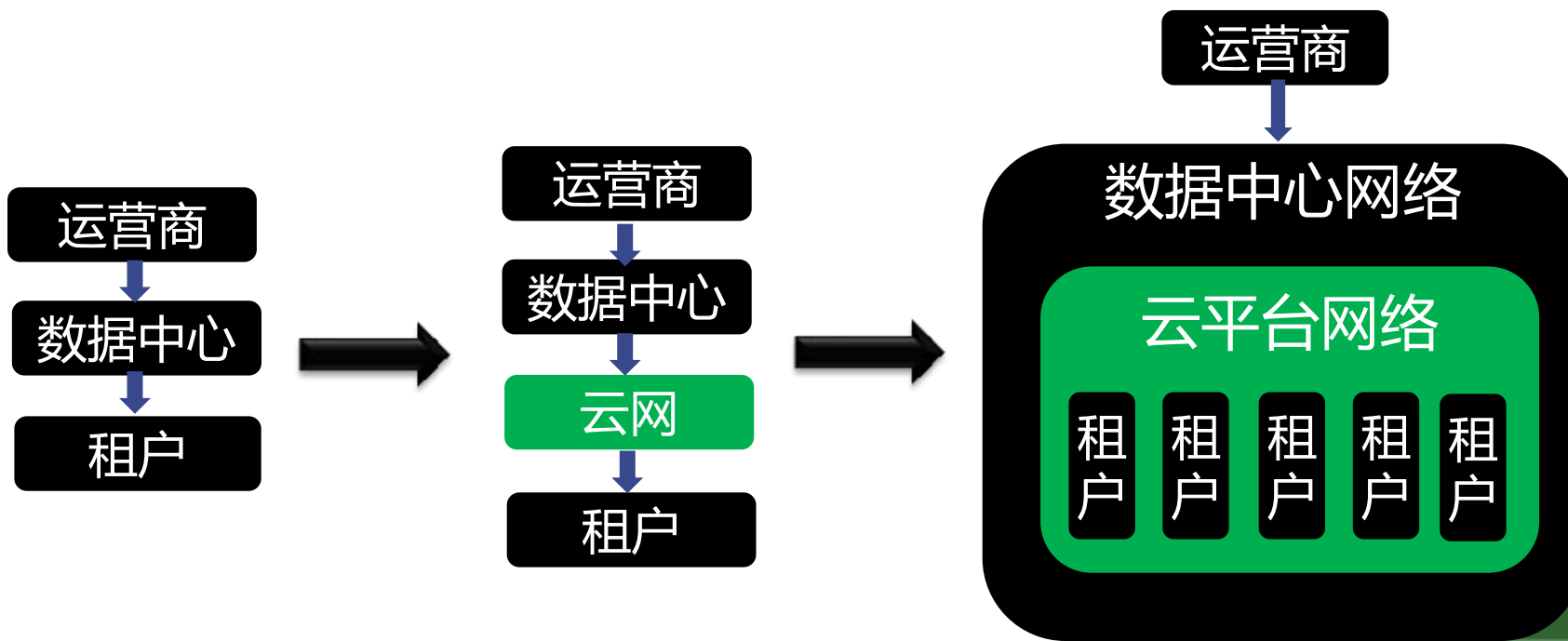


云网运维的痛

物理网络和逻辑网络的边界



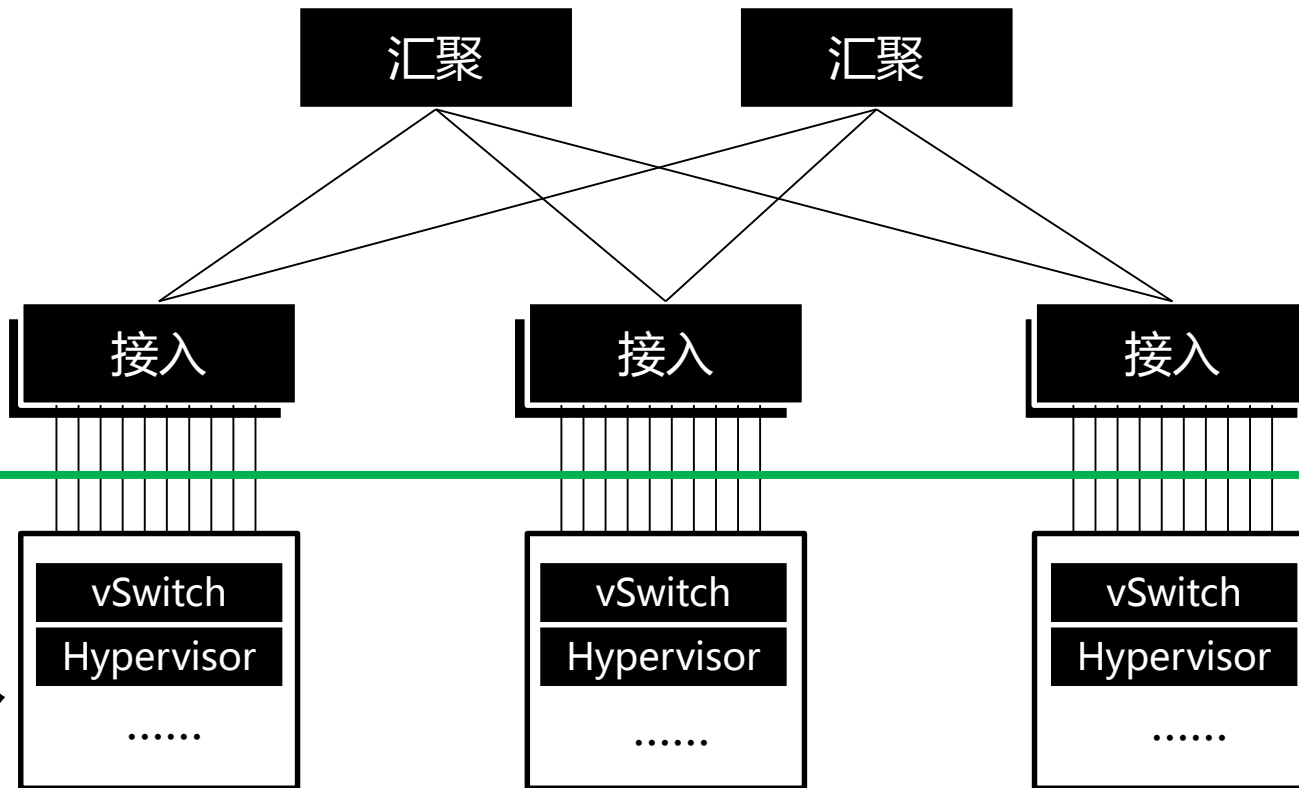
故障总是被用户先发现



服务器和网络团队相关吗

网络团队

IP地址重叠
隧道封装 (UDP 4789)

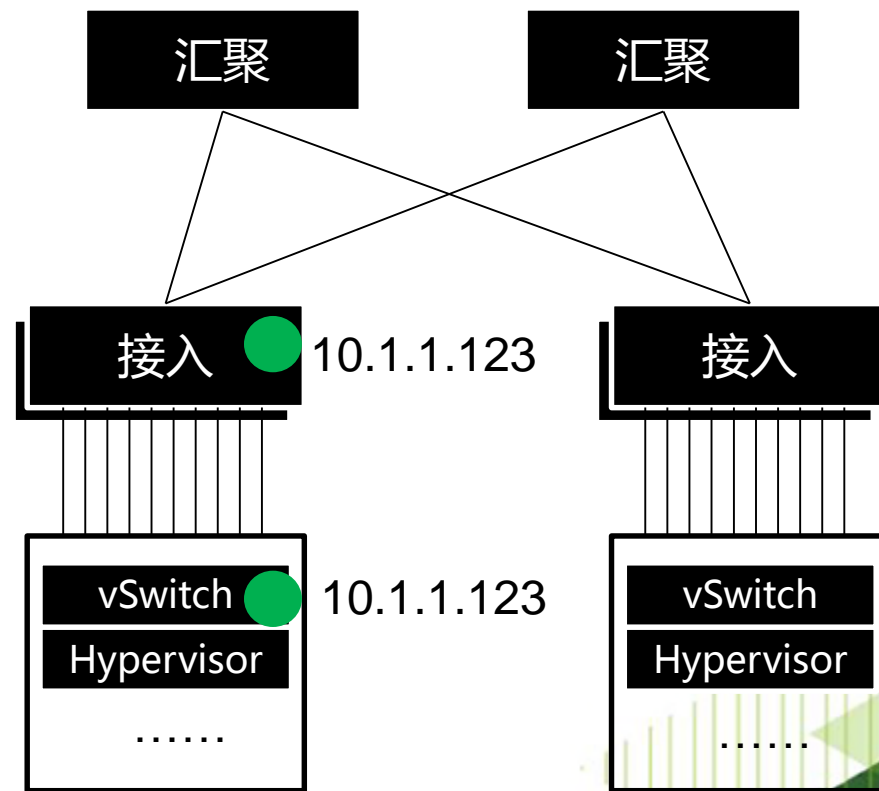


服务器团队

HTTP / DNS / ICMP / ARP ...

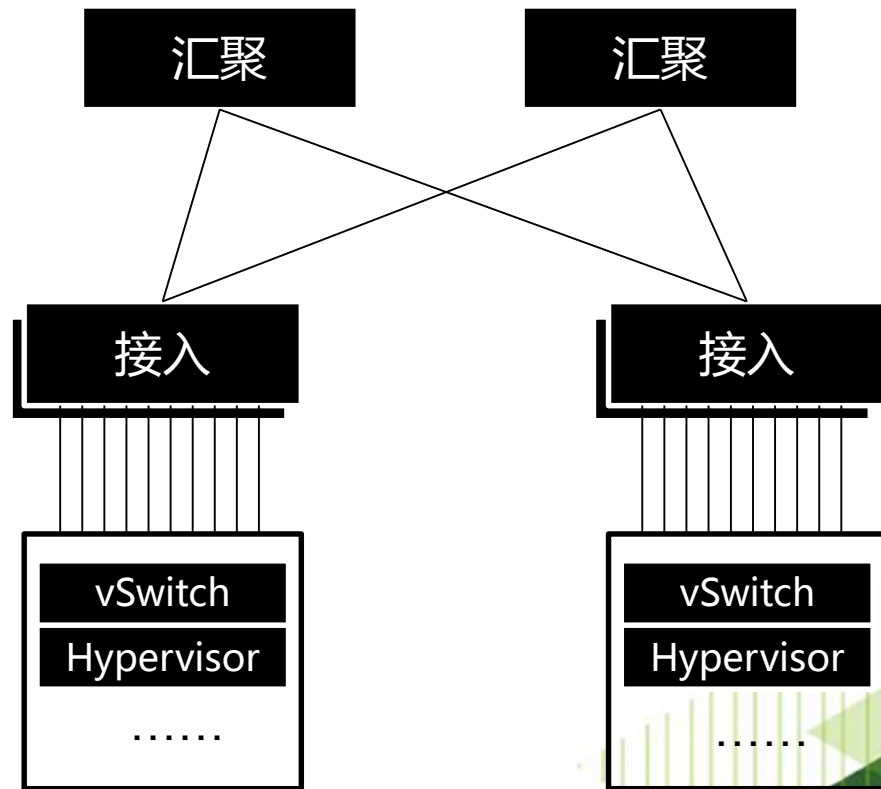
物理网络和逻辑网络的隔离和协作

- 那些Debug以后留下的坑
 - 在交换机上创建了一个SVI
 - 影响了租户的虚拟机
 - 只在一个MLAG Peer 交换机上放行了VLAN
 - 租户网络时通时断



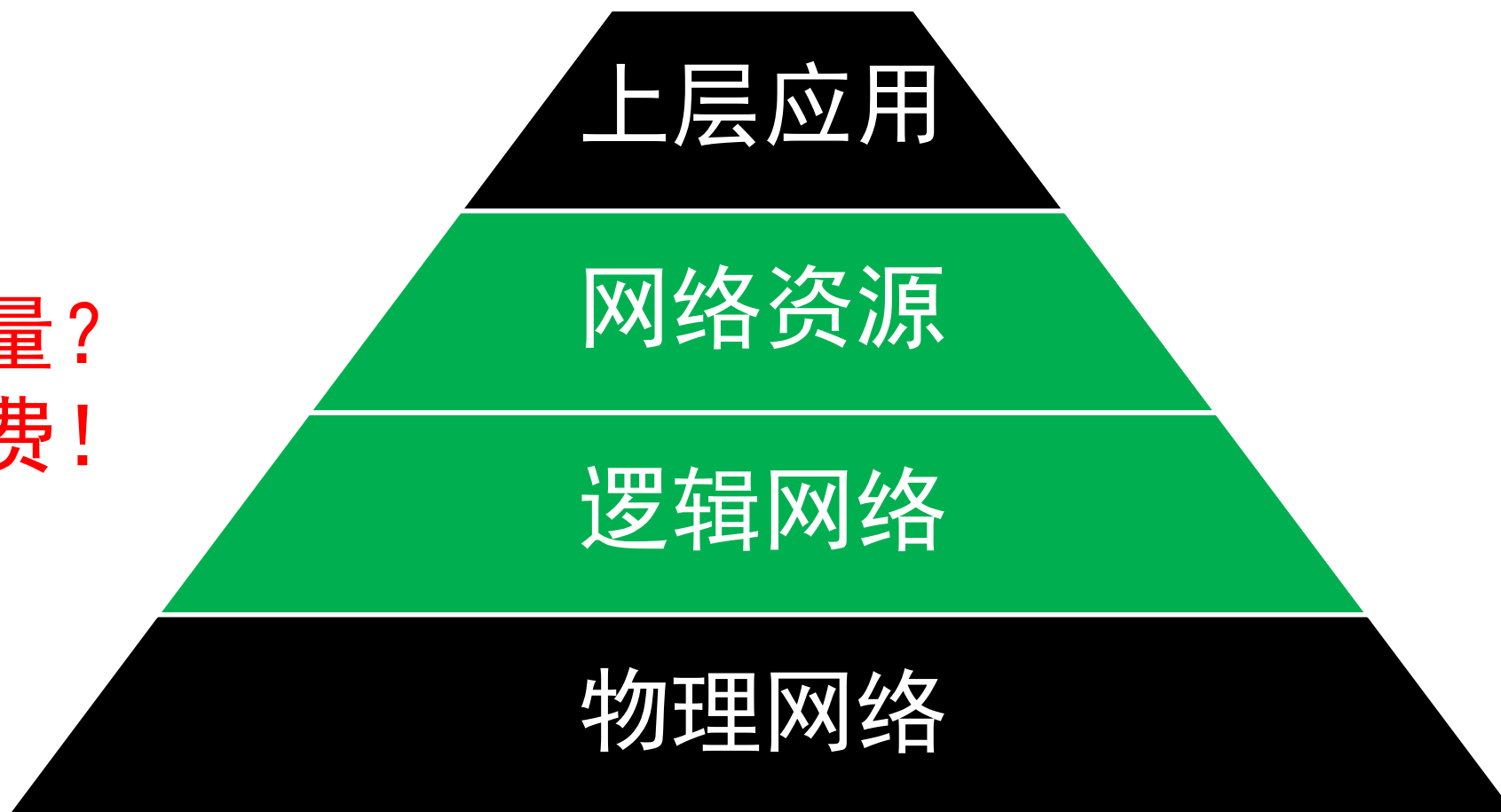
如何设置风暴抑制

- 硬件交换机不怕风暴
 - 硬件交换机擅长转发
 - 设置明确的BUM风暴速率
- 复制是虚拟交换机的灾难
 - 硬件交换机放行的风暴流量能被虚拟交换机放大几十倍



网络资源的监控

没有计量?
就有浪费!



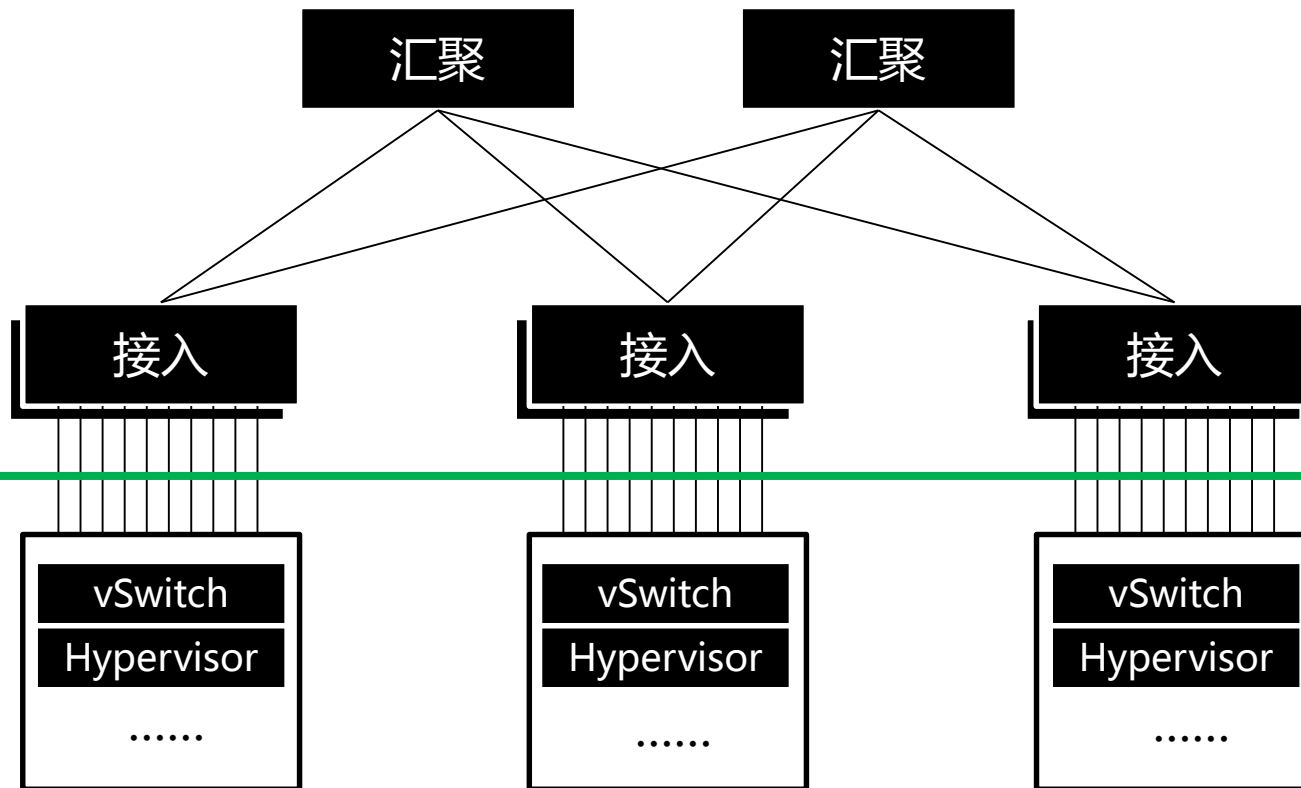
IP地址不够用

- 租户自由申请
 - IP地址自动分配
- 如何分配
 - 网关下沉：地址碎片化
 - 内网地址：体验差
 - /32地址：运维复杂
- 分配 == 使用？
 - 如何发现Dark Host



内网带宽共享

服务质量
如何监控

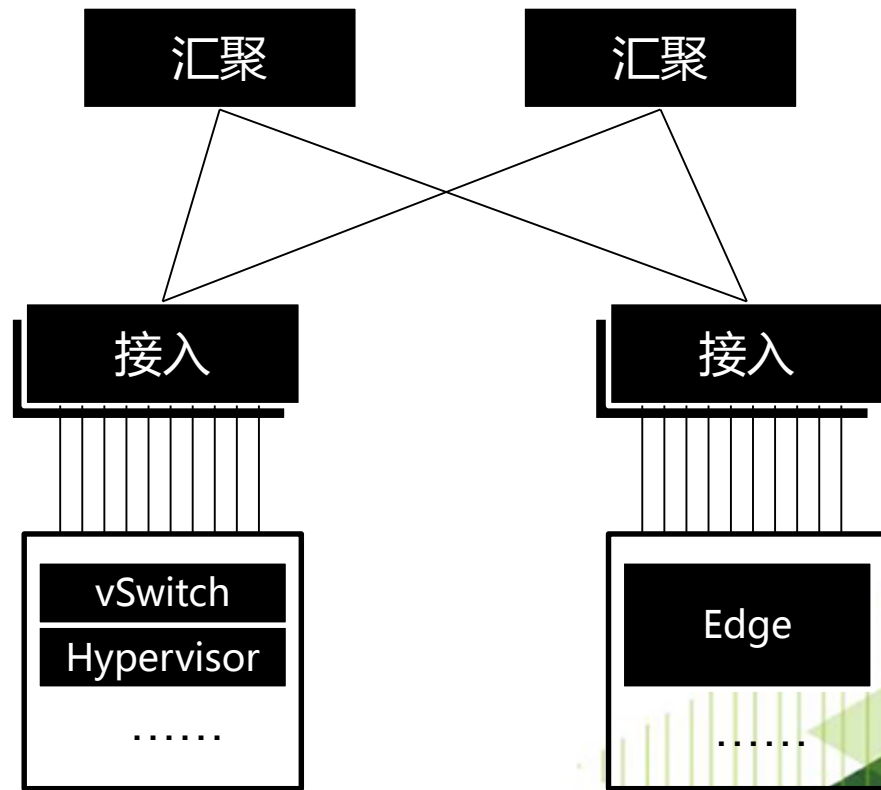


IP地址重叠
隧道封装 (UDP 4789)

HTTP / DNS / ICMP / ARP ...

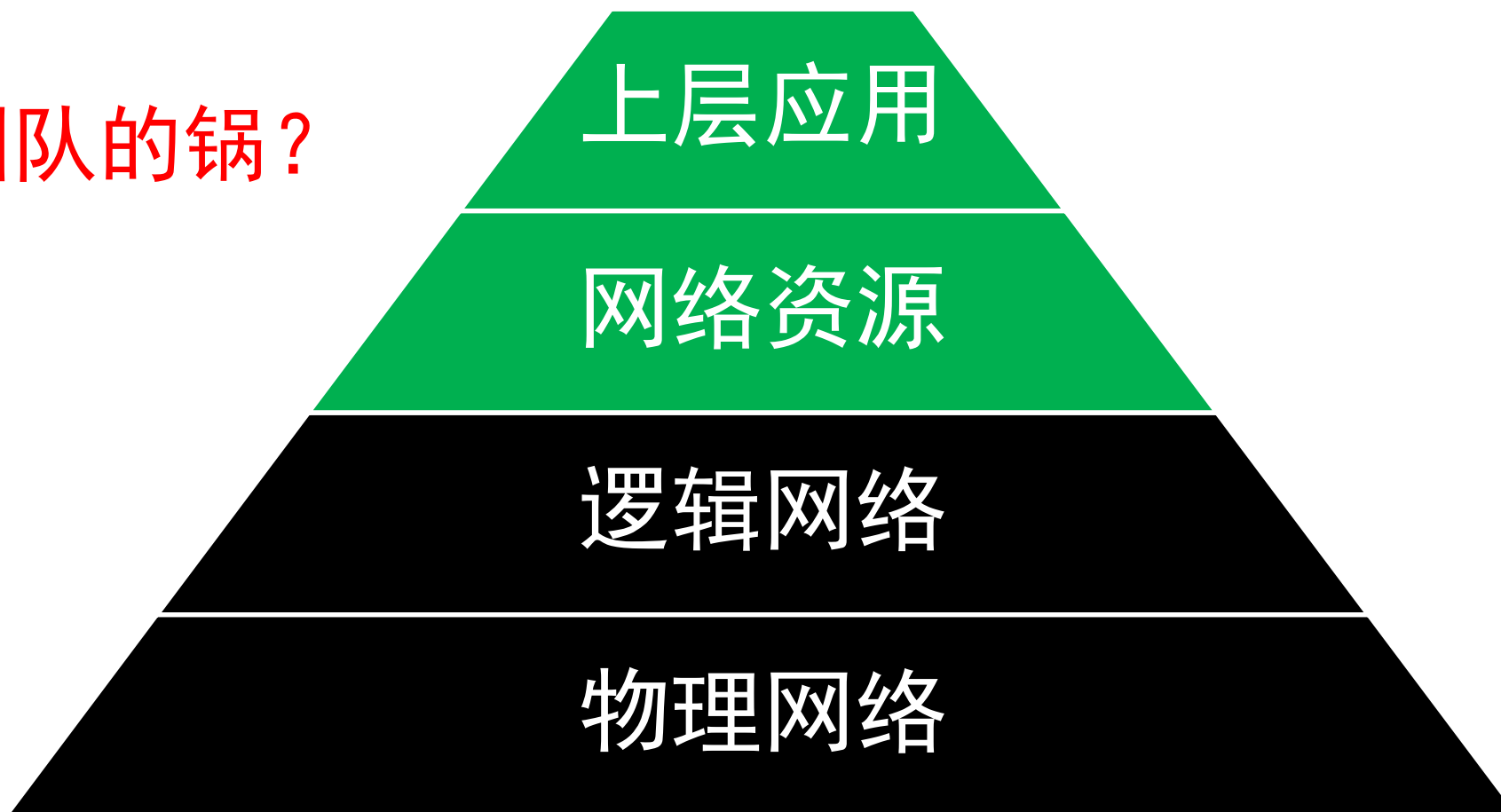
外网带宽计量

- 多维度
 - IP地址
 - L4 Port: 租户
 - VLAN: 不同运营商线路
 - VXLAN VNI: DCI虚拟机专线
- 细粒度
 - 95峰
 - 第三峰
 - ...
- Netflow?
 - 数据糙
 - 负载高



是应用，还是网络？

哪个团队的锅？



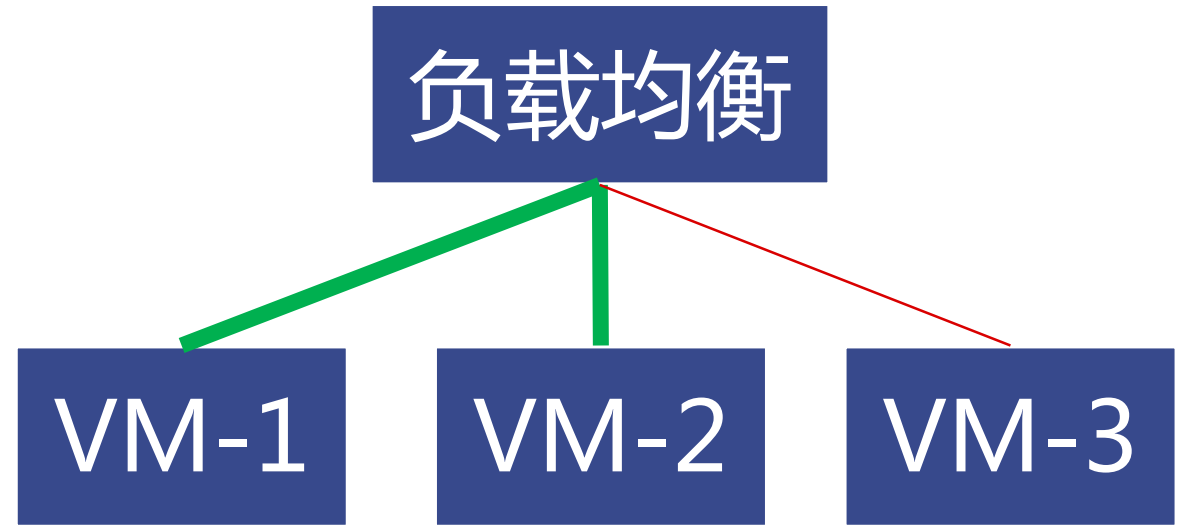
是应用的问题还是网络的问题



Loading...

一个抽风的后端主机

- 哪个阶段出了问题
 - TCP三次握手
 - HTTP Response
 - TCP Reset
- 其他配置是否正常
 - VIP/RIP、MAC地址是否正常
 - DNS配置是否正常



如果历史可以重来

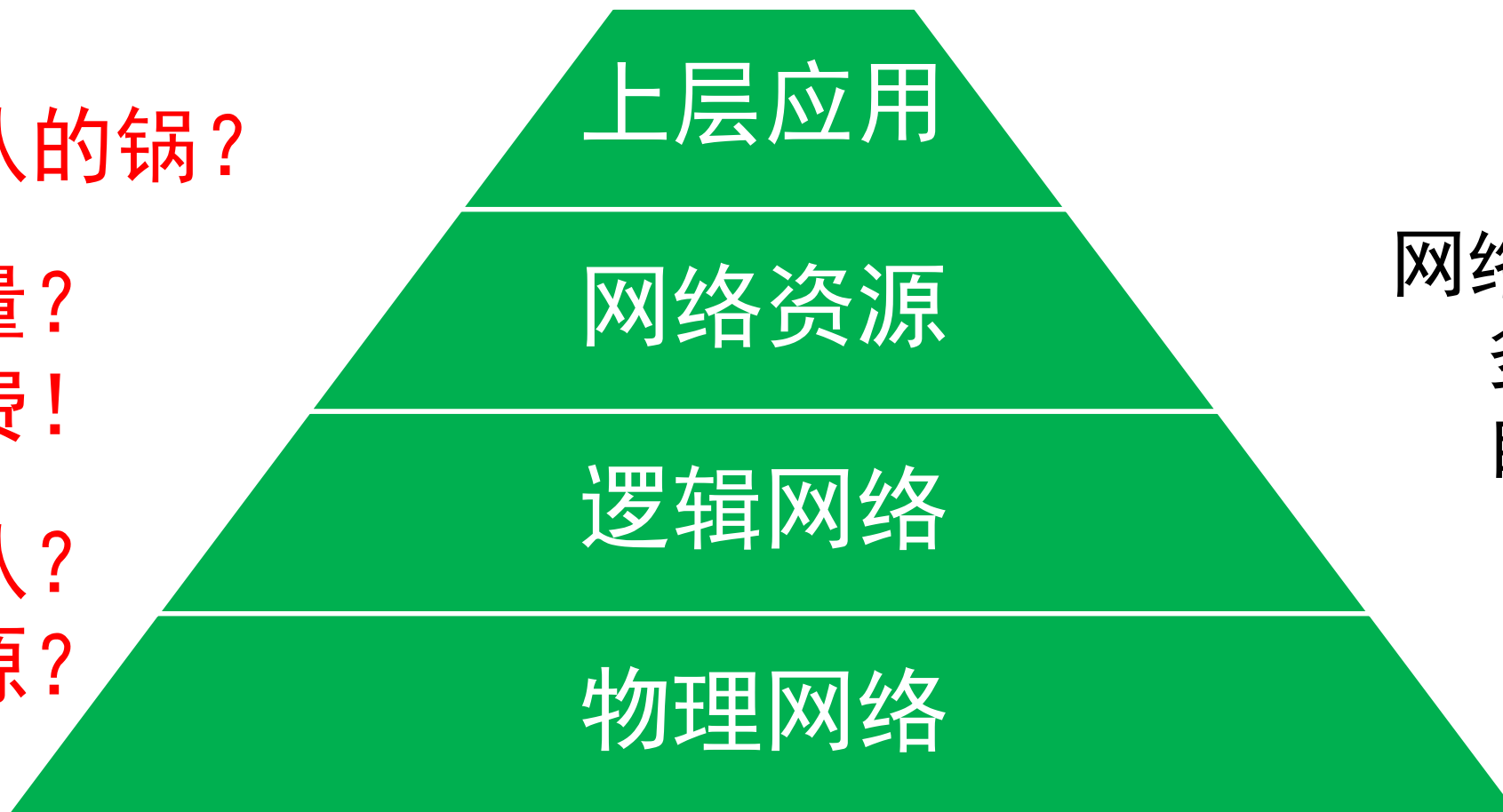
- 排查网络故障的时效性
 - 那些莫名其妙恢复的故障
- 历史证据
 - 昨天？ 上周？ 上个月

云网运维的痛

哪个团队的锅？

没有计量？
就有浪费！

两个团队？
两套资源？



网络Overlay
多租户
自服务

如何解决这些问题

网络运维内外交困



Loading...

老问题灭不完

应用问题 or 网络问题
内部网络问题 or 运营商网络问题



救火经验和本领失效

云平台无报警，但用户报障
更早的网络故障能追溯吗



新问题层出不穷

全是UDP4789，傻傻分不清楚
虚拟网络故障怕不怕，看不见的黑洞

如何解决这些问题

规章制度 员工规章制度

8、保持清洁，良好的工作环境，不得在办公区域吸烟、随地吐痰、大声喧哗；

9、员工违纪处分：迟到、早退、旷工、脱岗等四种，管理程序如下：

迟到：指未按规定到达工作岗位（或作业地点）；迟到30分钟以内，每次扣十元；迟到30分钟以上的扣半天基本工资；迟到一小时的扣全天工资；（人力不可抗拒因素造成的迟到除外）。

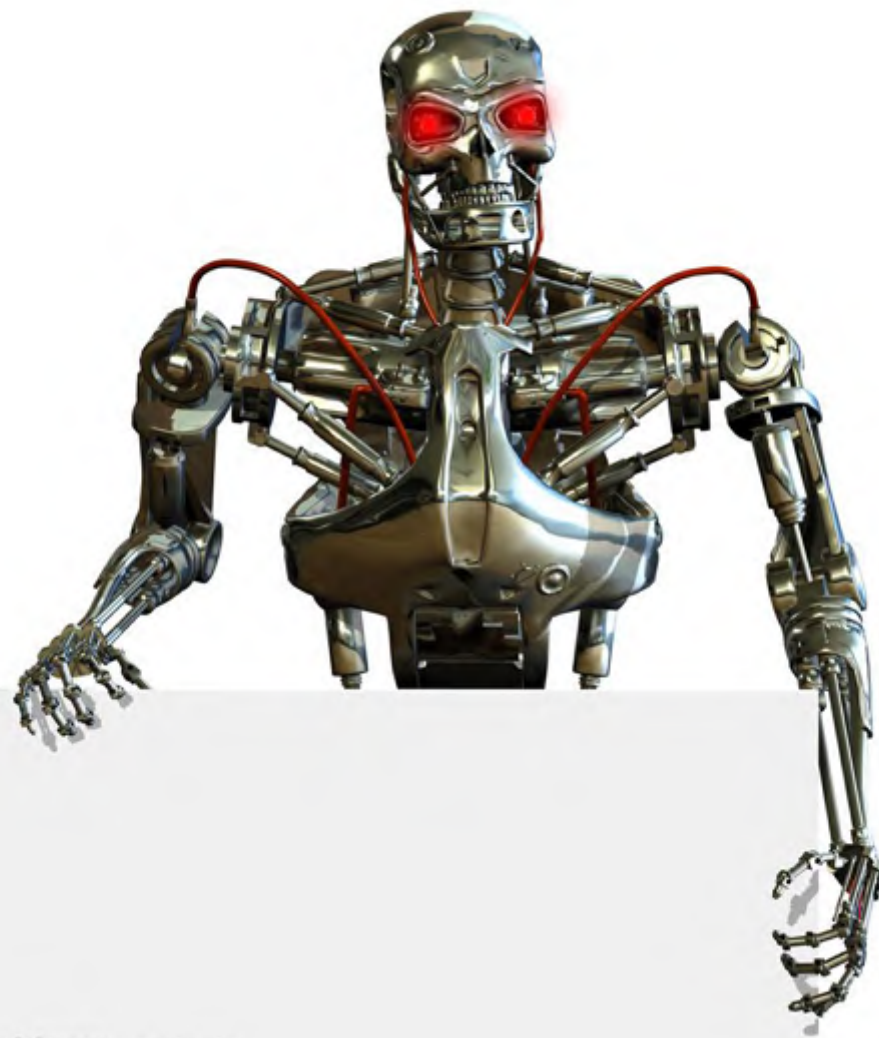
早退：指提前离开工作岗位下班；早退3分钟以内，每次扣十元；30分钟以上按旷工半天处理。

旷工：指未经同意或按规定程序办理请假手续而未正常上班的；旷工半天扣一天工资，旷工一天扣罚2天工资；

脱岗：指员工在上班期间未履行任何手续擅自离开工作岗位的，脱岗一次罚款十元；

10、假别分为：病假、事假、婚假、产假、年假、工伤假。

※质量第一※



络

火

平台
更早

楚
黑洞

内部网



打造云网大数据平台
摆脱云网运维困境

技术上的挑战

- 数据采集

- 虚拟交换机：与服务器/云平台运维团队协作
- 全网全时：全量采集的能力 + 细粒度策略控制

- 数据分析

- 海量流量：高性能协议栈、分布式处理
- 多租户网络：VLAN/VXLAN识别，流量和资源/租户进行映射
- 数据高效存储：聚合、压缩

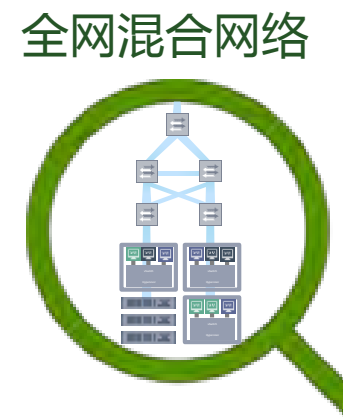
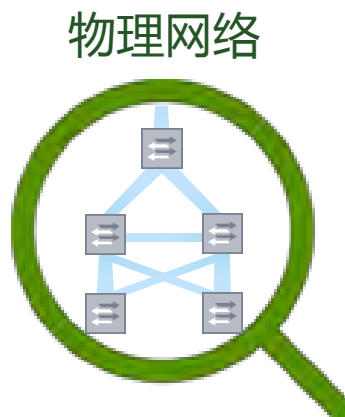
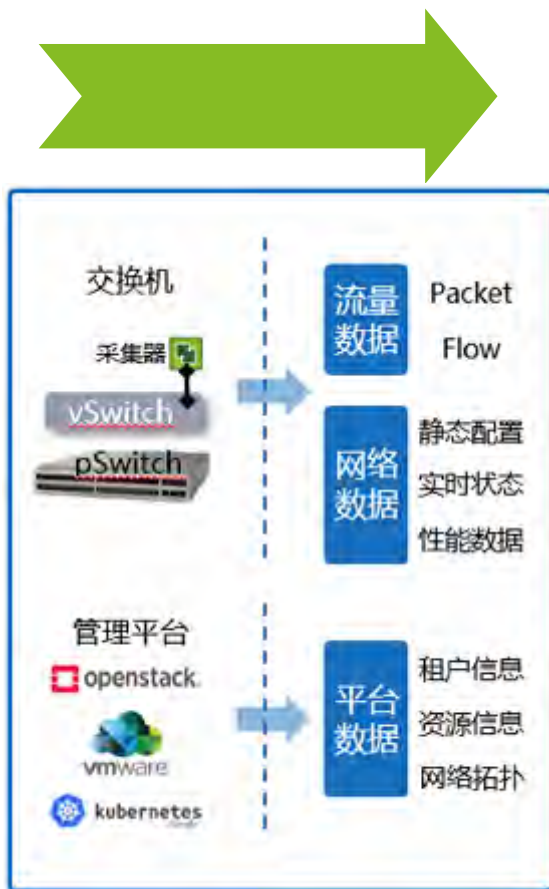
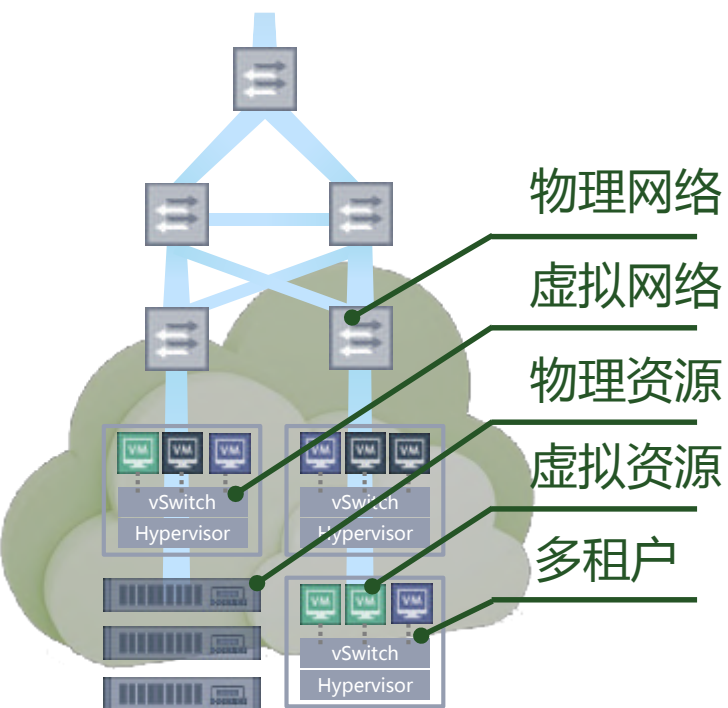
- 数据展现

- 数据关联，多维度展现
- 实时报警：发现毫秒级毛刺、与业务关联



大范围 and 细粒度 数据采集

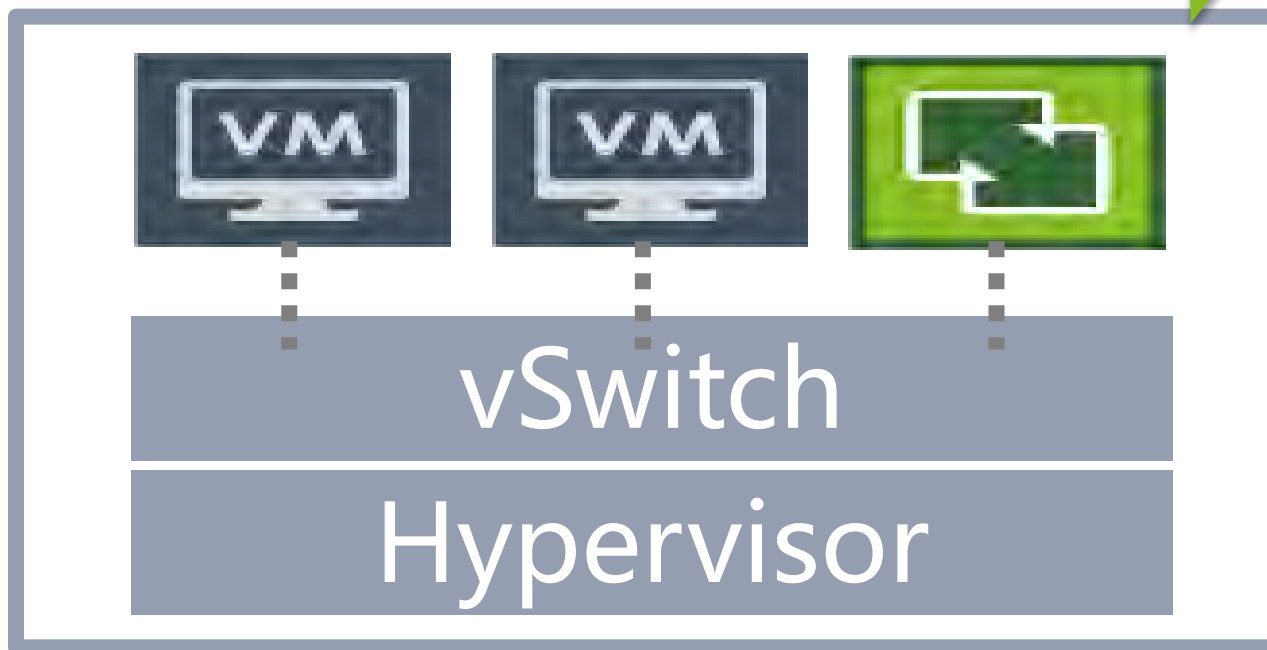
多维度采集



采集虚拟网络

- OpenStack Tap-as-a-Service

采集器



控制投入

- 生产网络和监控网络的投入比例
 - 包处理的瓶颈：Linux内核 1.488 Mpps
 - 数据存储的瓶颈：CPU/内存/硬盘的速度
- 全量采集的能力
 - 采集点全覆盖
- 细粒度策略控制
 - 物理交换机：ACL过滤流量
 - 虚拟交换机：OpenFlow过滤流量

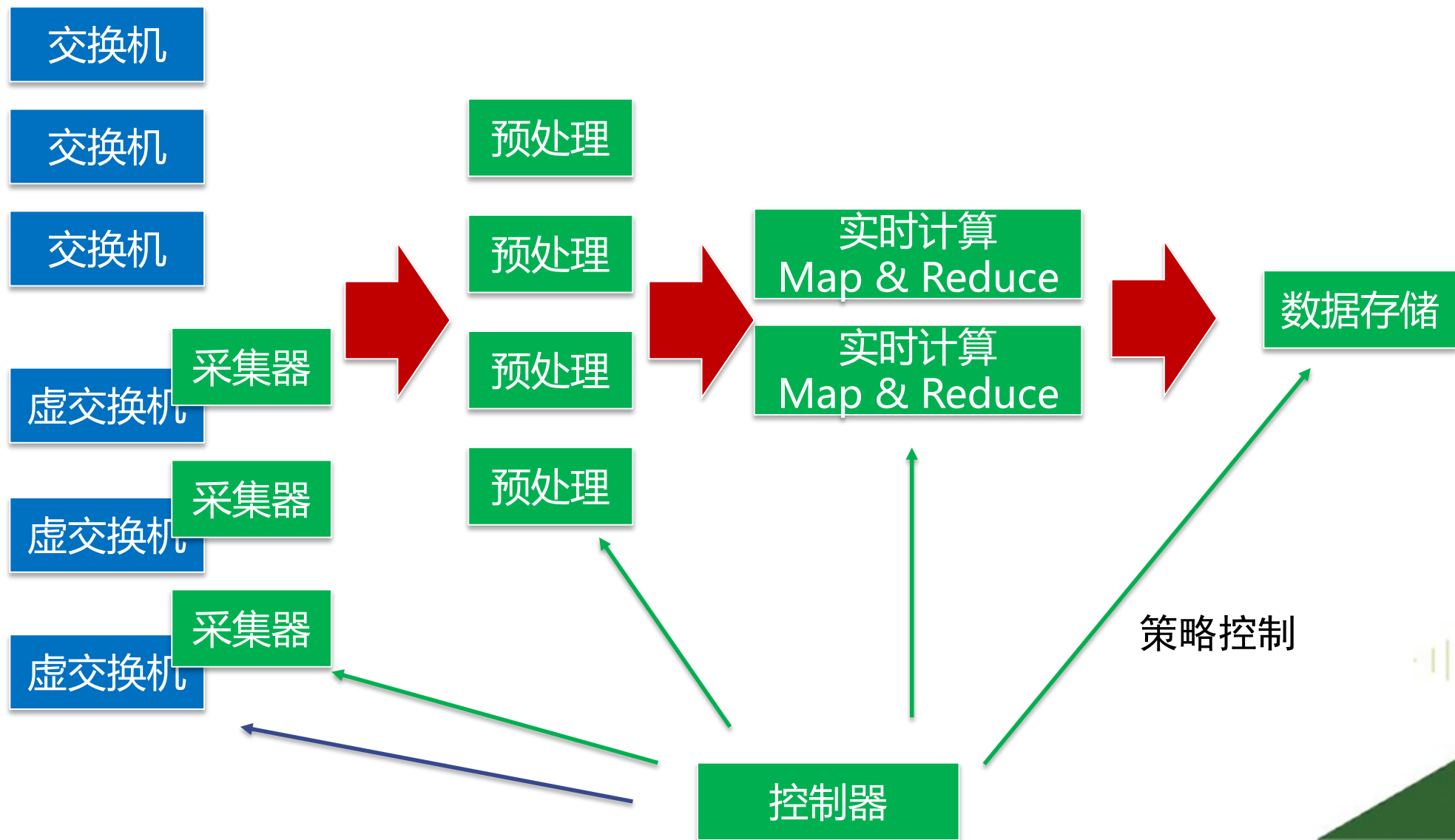


高性能 数据分析

发挥最大的能力

- 不同的流量，需要不同的处理
 - 任意流量：基于Flow元组的分析
 - TCP流量：基于四层包头的分析
 - HTTP流量：基于HTTP头的分析
 - 私有协议流量：基于前N个字节的分析
 - 重要流量：全量记录，用于审核

策略控制



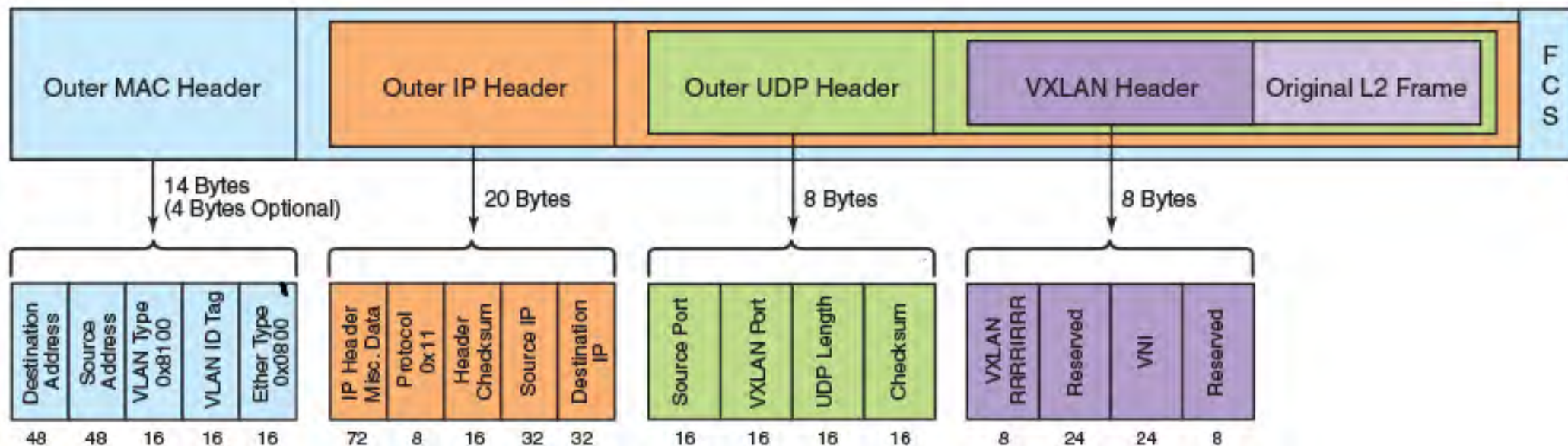
高性能协议栈：网络层

- 聚合Packet为Flow元组
- 保留所有4层包头
- 处理性能
 - Linux内核：1.488 Mpps
 - DPDK 14.88 Mpps **起**
- 存储性能
 - 包头差量压缩，90%压缩率
 - 一个惊人的事实：包头存储空间 ~ Flow元组



识别多租户网络

- VLAN - 租户映射关系
- VXLAN VNI - 租户映射关系：网包重新解析



高性能协议栈：应用层

- 常见应用的元数据提取
 - HTTP头：非侵入式Web Server日志
 - DNS请求日志





多维度
数据展现

网络配置的Git log

- 记录配置历史
- 追溯人为的失误

```
[root@NSP_Controller_Server ~]# mt history.switch mip=172.17.1.102
--- 35b4895e733d1c004eb961c5dd609c35 SEQ-28 2016-12-28 13:44:21
+++ 95dc542ab31049a6baf8ba12a39b0aaa SEQ-29 2016-12-29 15:00:27
@@ -102,10 +102,10 @@
     policer-aggregate ingress_qos_eth-0-2
     !
     interface eth-0-1
-   - description VIF933_IFINDEX1_Nets
     jumboframe enable
     ip arp inspection trust
     switchport access allowed vlan add 1024,1025
+   + shutdown
     lldp enable txrx
     !
     interface eth-0-2

--- 0a103cbd52092c53e5a110b8607d0702 SEQ-27 2016-12-27 23:43:07
+++ 35b4895e733d1c004eb961c5dd609c35 SEQ-28 2016-12-28 13:44:21
@@ -105,7 +105,7 @@
     description VIF933_IFINDEX1_Nets
     jumboframe enable
     ip arp inspection trust
-   - switchport access allowed vlan add 1024,4090
+   + switchport access allowed vlan add 1024,1025
     lldp enable txrx
     !
     interface eth-0-2
```

MAC震荡和ARP欺骗

- 覆盖物理、虚拟网络
 - ARP流量分析
 - 流量中的MAC入口分析

The screenshot displays the DeepFlow network analysis dashboard. The main content area is titled '云网分析 / 流量日志 / ARP欺骗' (Cloud Network Analysis / Traffic Log / ARP Spoofing). It includes a navigation sidebar on the left with options like '云网总览', '云网服务', '云网分析', '物理网络', '项目网络', and '流量分析'. The main panel shows a table of ARP spoofing events with columns for '序号' (Serial Number), '项目' (Project), 'IP地址' (IP Address), 'MAC地址' (MAC Address), '关联设备' (Associated Device), '关联接口' (Associated Interface), '发送网包数量' (Number of Sent Packets), '首次监控时间' (First Monitoring Time), and '最近监控时间' (Last Monitoring Time). Two entries are visible: one for IP 10.33.0.3 with MAC 14:cf:92:88:e3:2e (81 packets) and another for the same IP with MAC 28:2c:b2:ff:60:db (15 packets).

序号	项目	IP地址	MAC地址	关联设备	关联接口	发送网包数量	首次监控时间	最近监控时间
01	demo	10.33.0.3	14:cf:92:88:e3:2e			81	2017/0/2 10:12:16	2017/0/4 20:16:41
02	demo	10.33.0.3	28:2c:b2:ff:60:db			15	2017/0/4 20:16:2	2017/0/4 20:16:4

抽取服务日志

- Telnet/SSH/RDP日志：发现暴力破解
- DHCP/DNS：发现集群配置差异



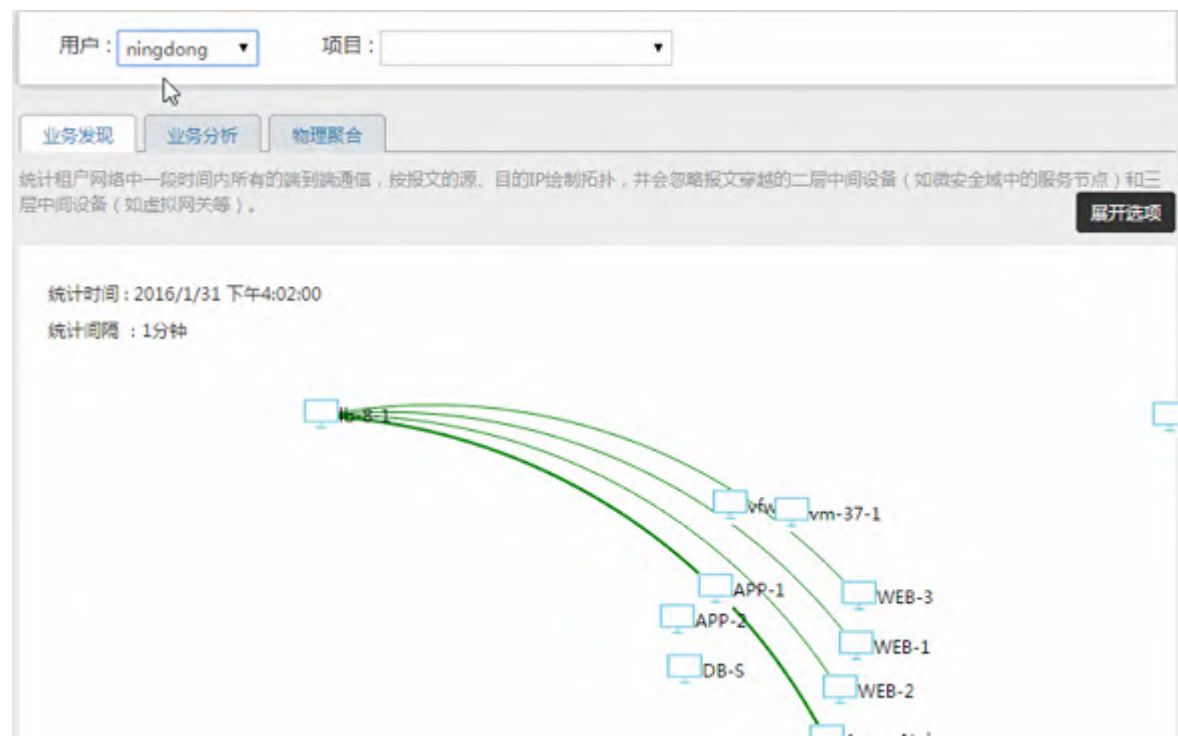
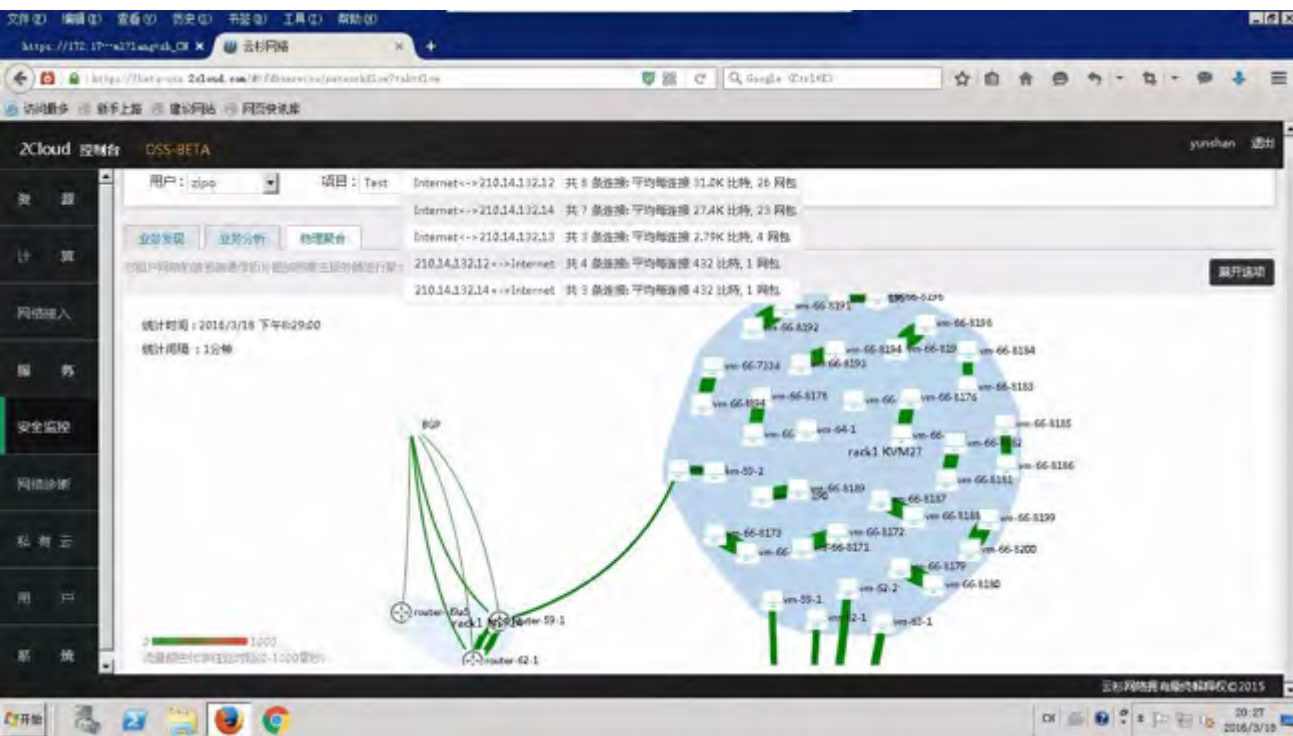
The screenshot shows the DeepFlow interface with the following components:

- Header:** DeepFlow logo and user 'admin'.
- Breadcrumbs:** 云网分析 / 流量日志 / 控制台登录
- Time Range:** 时间粒度: 1分钟, 时间间隔: 2017/01/05 19:40:21 - 2017/01/05 20:40:21
- Navigation:** MAC震荡, ARP欺骗, 控制台登录 (selected), DHCP服务, DNS服务, IP活跃度
- Table:** 控制台登录列表

序号	项目	服务端IP	服务端端口	客户端IP	客户端网包数量	平均会话时长
01		192.168.90.150	22	10.33.0.73	3	8490
02	demo	10.33.0.110	22	10.33.2.200	808212	94815
03	demo	10.33.49.5	22	10.33.2.200	8	205530

是网络的问题还是应用的问题

- 租户流量拓扑



细粒度计量

- 自定义计量的聚合字段
 - IP地址、VLAN、端口号、...
- 基于Apache Storm的实时流式计算



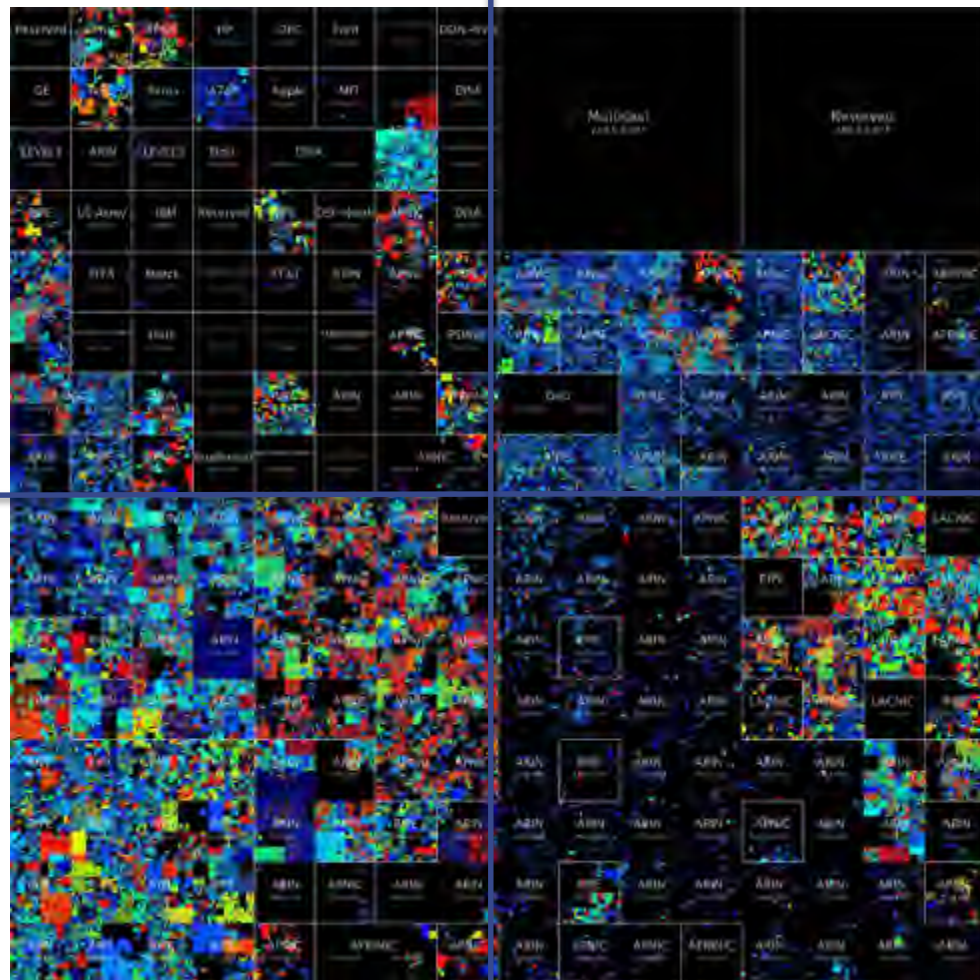
到底哪些IP在用

0.0.0.0/2

64.0.0.0/2

192.0.0.0/2

128.0.0.0/2



数据存储：用包头还原历史会话

The screenshot displays the Wireshark interface with a packet list table and two sequence number graphs. The packet list table shows the following data:

No.	Time	Source	Destination	Protocol	Length	Info
1	0.000000	172.16.187.246	172.16.1.131	TCP	74	50380→22 [SYN] Seq=0 Win=14600 Len=0 MSS=1460 SACK_PERM=1 TSval=3457413568 TSecr=0 WS=128
2	0.000660	172.16.1.131	172.16.187.246	TCP	74	22→50380 [SYN, ACK] Seq=0 Ack=1 Win=14480 Len=0 MSS=1460 SACK_PERM=1 TSval=3029870721 TSecr=3457413568 WS=128
3	0.000686	172.16.187.246	172.16.1.131	TCP	66	50380→22 [ACK] Seq=1 Ack=1 Win=14720 Len=0 TSval=3457413569 TSecr=3029870721
4	0.001098	172.16.187.246	172.16.1.131	SSH	87	Client: Encrypted packet (len=21)
5	0.001418	172.16.1.131	172.16.187.246	TCP	66	22→50380
6	0.000539	172.16.1.131	172.16.187.246	SSH	87	Serv
7	0.008602	172.16.187.246	172.16.1.131	TCP	66	50380
8	0.009463	172.16.1.131	172.16.187.246	SSH	1682	Serv
9	0.009475	172.16.187.246	172.16.1.131	TCP	66	50380
10	0.009771	172.16.187.246	172.16.1.131	SSH	1514	Clie
11	0.009787	172.16.187.246	172.16.1.131	SSH	466	Clie
12	0.009904	172.16.1.131	172.16.187.246	TCP	66	22→50380
13	0.011872	172.16.187.246	172.16.1.131	SSH	146	Clie

Below the packet list, the details pane shows the structure of the selected packet (Frame 1):

- Ethernet II, Src: RealtekU_cc:58:e6 (52:54:00:12:35:00)
- Internet Protocol Version 4, Src: 172.16.187.246, Dst: 172.16.1.131
- Transmission Control Protocol, Src Port: 50380, Dst Port: 22

Two sequence number graphs are overlaid on the interface:

- Sequence Number (Stevens) for 172.16.187.246:50380 → 172.16.1.131:22**: This graph shows the sequence number on the y-axis (0 to 3200) and time on the x-axis (0 to 0.09 seconds). The sequence number starts at approximately 1600 and increases in steps to about 2000 by 0.015 seconds, then jumps to 2400 and continues to rise to 3200 by 0.09 seconds.
- Sequence Number (tcptrace) for 172.16.187.246:50380 → 172.16.1.131:22**: This graph shows the sequence number on the y-axis (0 to 27000) and time on the x-axis (0 to 0.09 seconds). The sequence number starts at 0 and increases in steps to about 18000 by 0.015 seconds, then jumps to 22000 and continues to rise to 27000 by 0.09 seconds.

The graphs include controls for zooming, panning, and saving the image.

我们的技术栈

展现



数据存储



数据分析



运维



数据传输



策略控制



采集/预处理



数据源



期待你的加入



技术创造价值

www.yunshan.net.cn
info@yunshan.net.cn