

2017源创会年终盛典

与电子标准院共建开源标准

12月23日 北京万豪酒店

腾讯信鸽海量移动推送服务构建

数据平台部 甘恒通

1 推送系统建设

终端

基础设施

后台

云化治理

2 增值服务

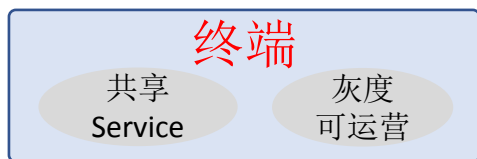
精准推送流程

数据积累

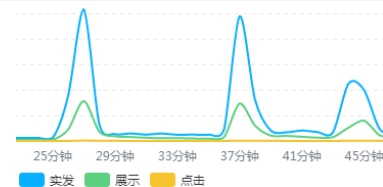
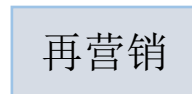
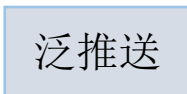
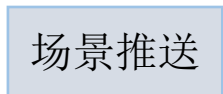
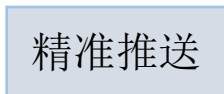
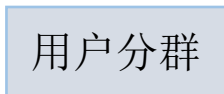
支撑平台

可视化操作

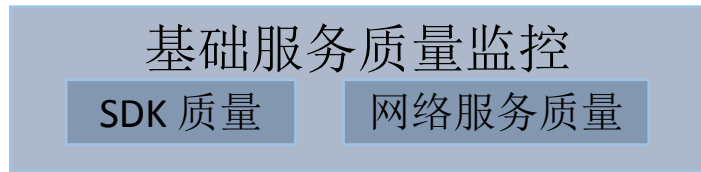
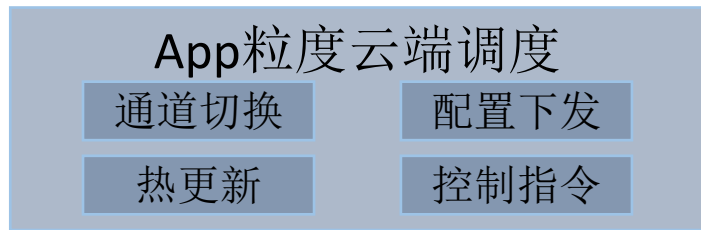
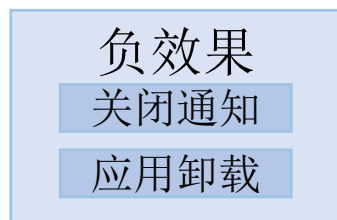
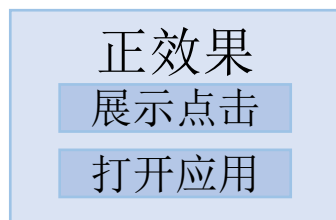
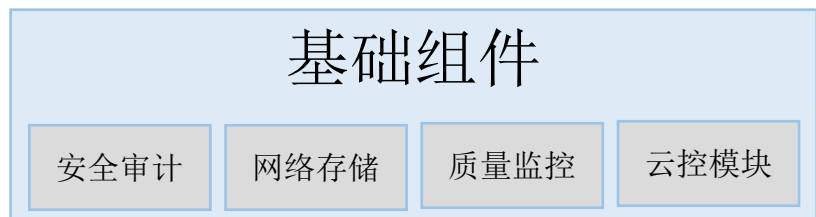
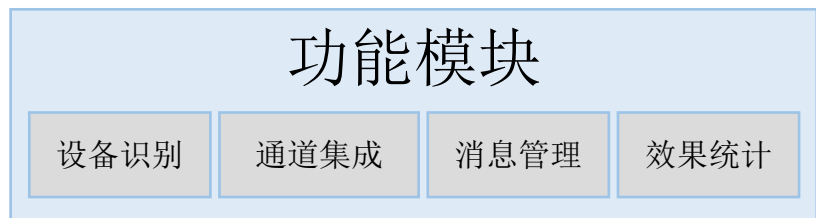
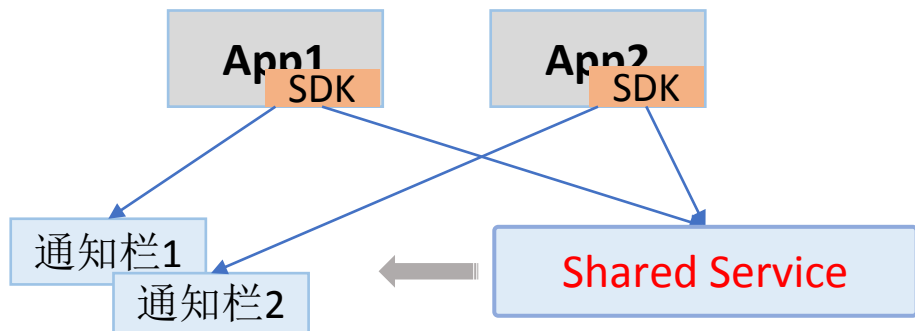
移动推送服务概览



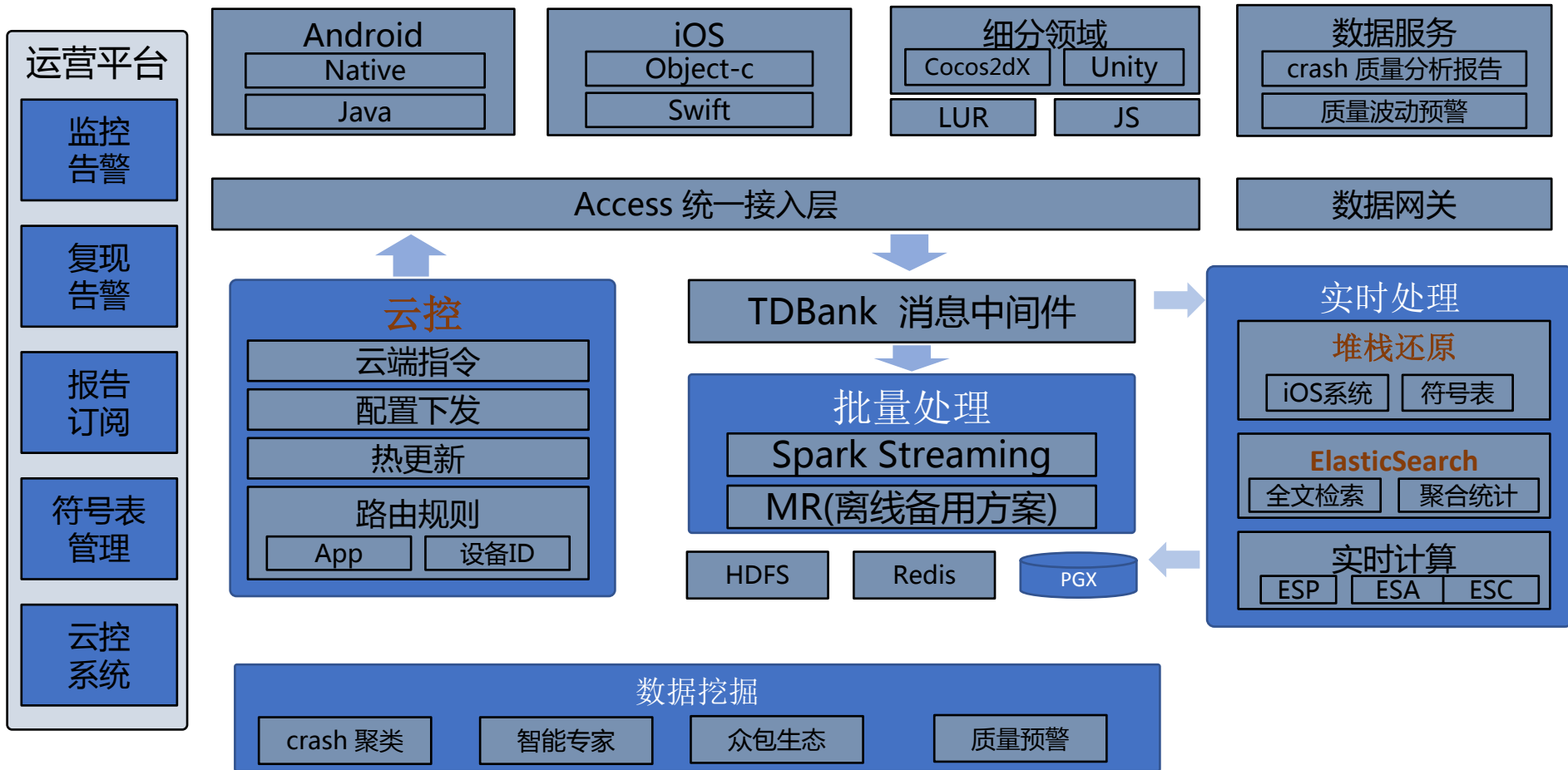
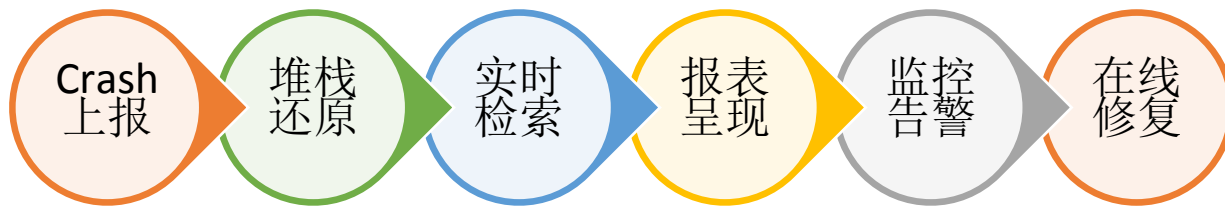
实时效果跟踪、多维运营分析



终端共享通道、服务可运营



终端质量监控



网络基础设施

中国大陆、香港用户

TGW 专线
HAProxy 加速



海外用户

POP 点加速



海外用户

Akamai IPA 加速



北美用户

直连访问



TGW
专线

POP
加速

Akamai
IPA

内网 IDC
深圳
香港
加拿大

腾讯云



AWS

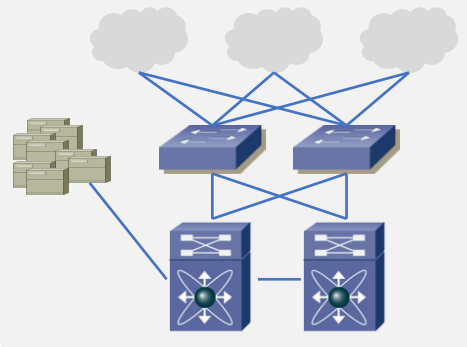


Akamai

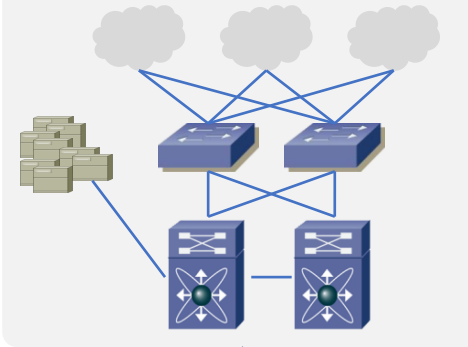


私有IDC

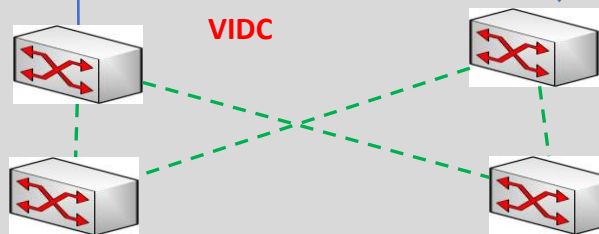
深圳 IDC 园区



上海 IDC 园区

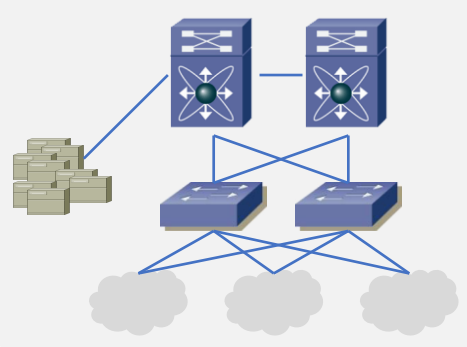


香港 IDC

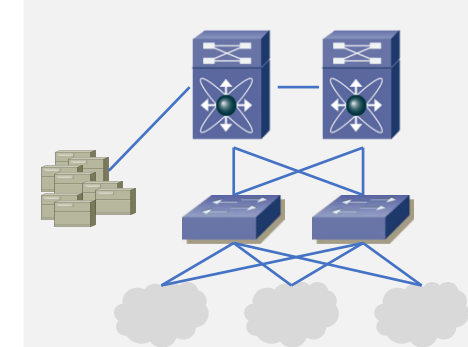


加拿大 IDC

西安 IDC 园区

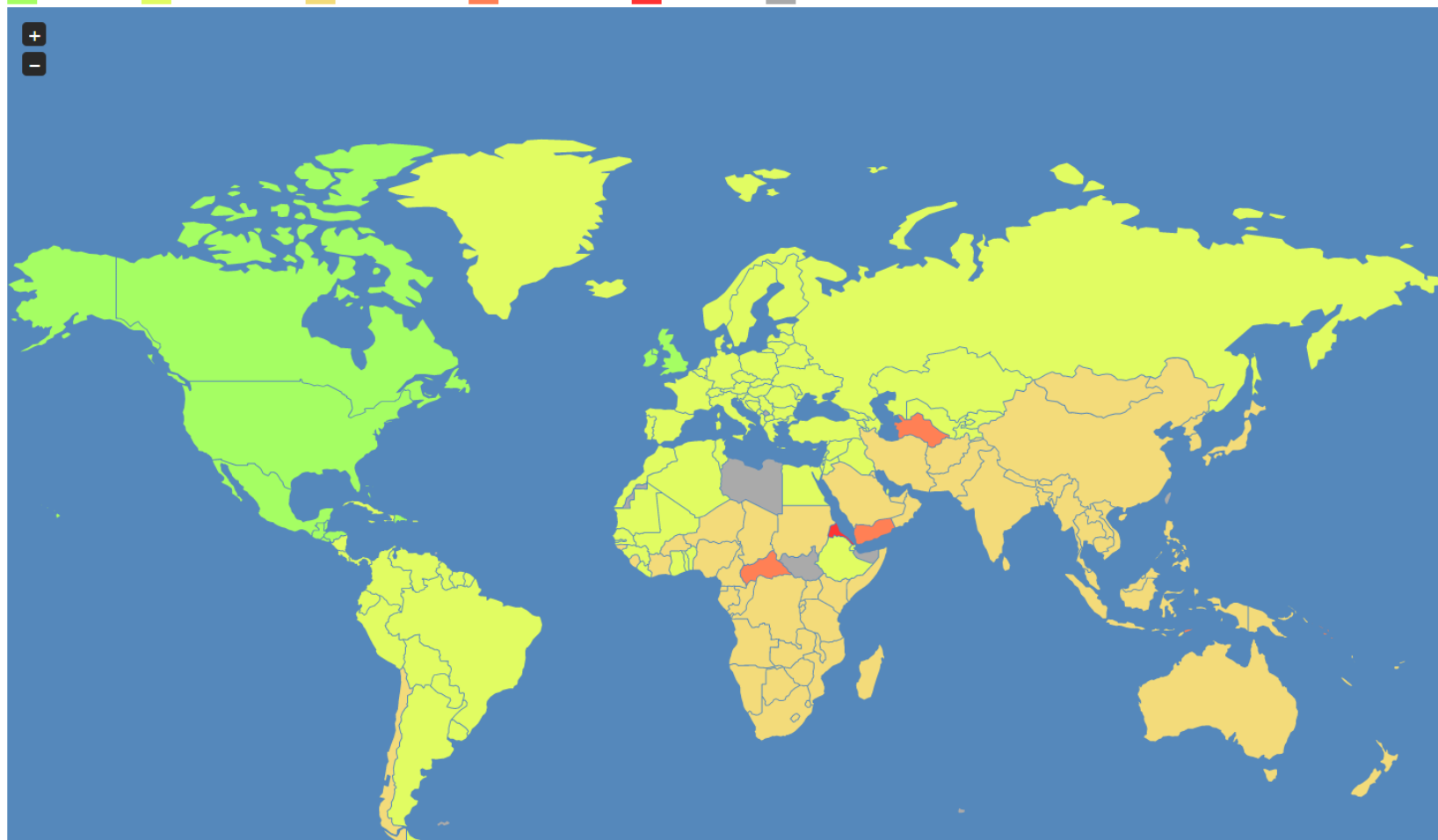


天津 IDC 园区



全球网络服务质量监控 — 海外用户视角

0~99毫秒 100~199毫秒 200~349毫秒 350~499毫秒 500+毫秒 无数据



国内网络服务质量监控 – 国内用户视角

机房所在运营商: 全部 | 机房: 深圳 | ←支持关键字模糊匹配
 用户所在运营商: 中国电信 | 日期范围: 2017-10-03 ~ 2017-10-03 | 探测分类: | 数据类型: 平均值 | 查询 | 导出原始数据

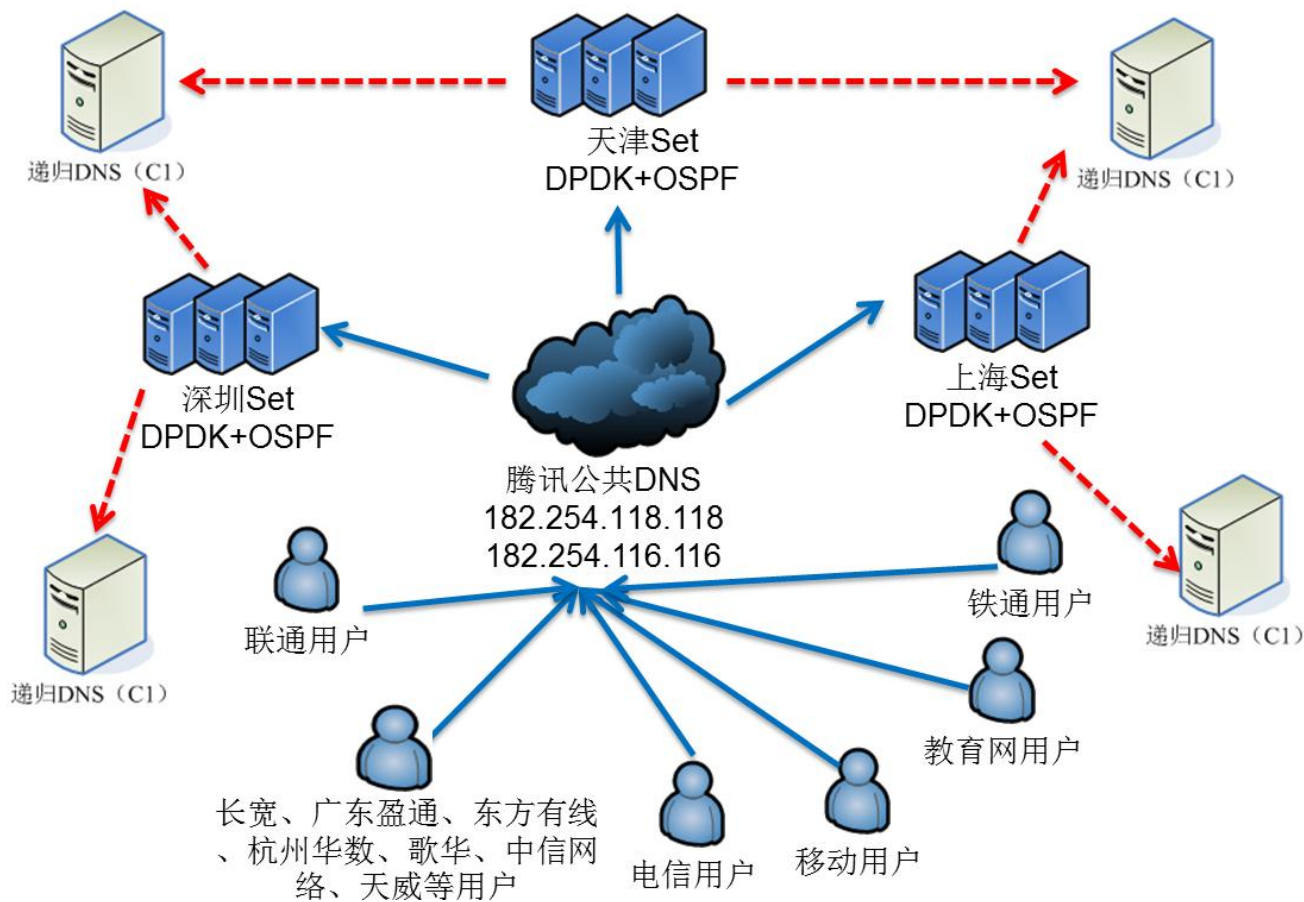
0~99毫秒 | 100~199毫秒 | 200~349毫秒 | 350~499毫秒 | 500+毫秒 | 无数据



| 省份 | 平均延时 | 平均丢包 | 省份 | 平均延时 | 平均丢包 | 省份 | 平均延时 | 平均丢包 |
|----|------|------|----|------|------|----|------|------|
| | | | | | | | | |

自建DNS服务

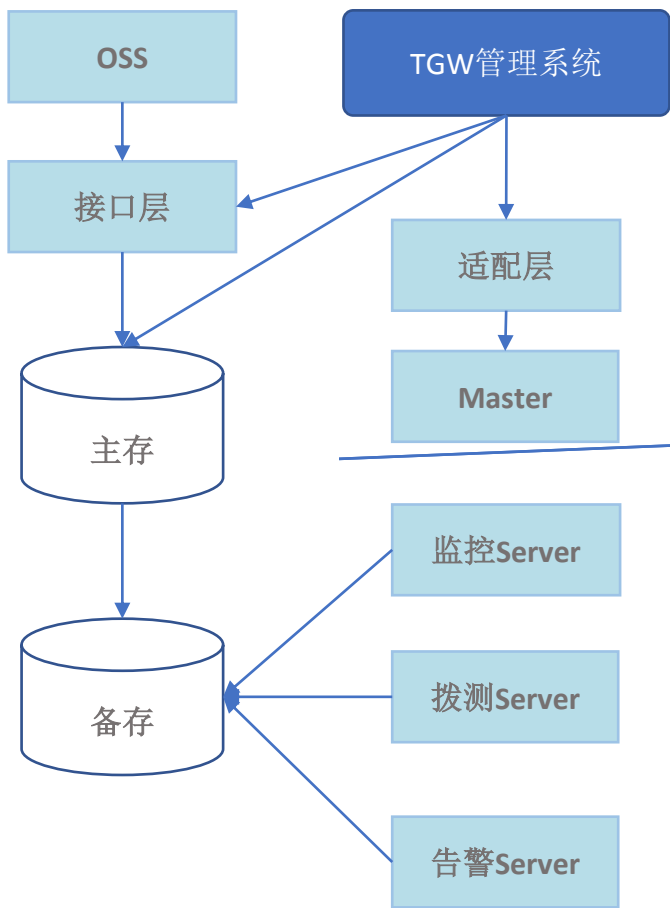
一个IP、多地容灾，OSPF集群IDC独立部署，流量秒级无感知切换。



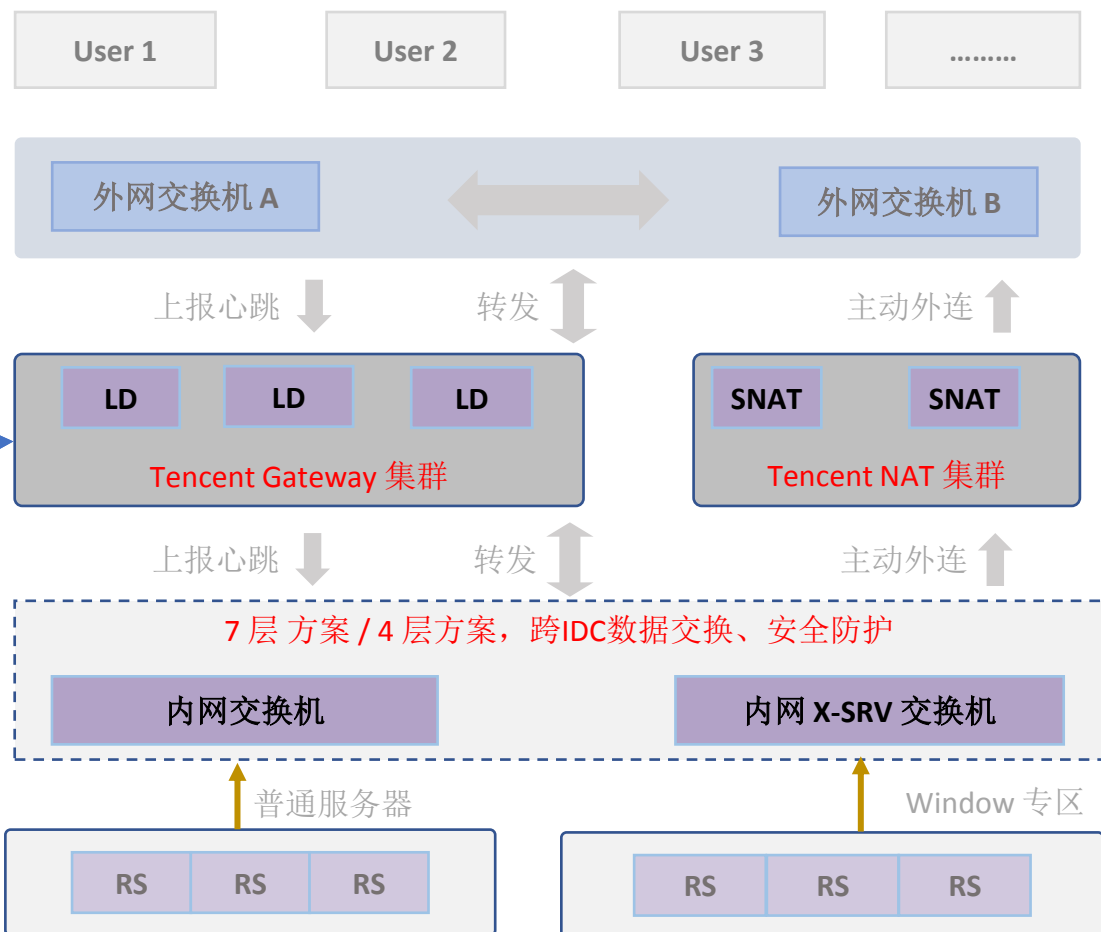
统一接入网关

统一接入、多网合一；负载均衡、平滑扩容；网络隔离防护、IP 收敛。

运营支撑系统

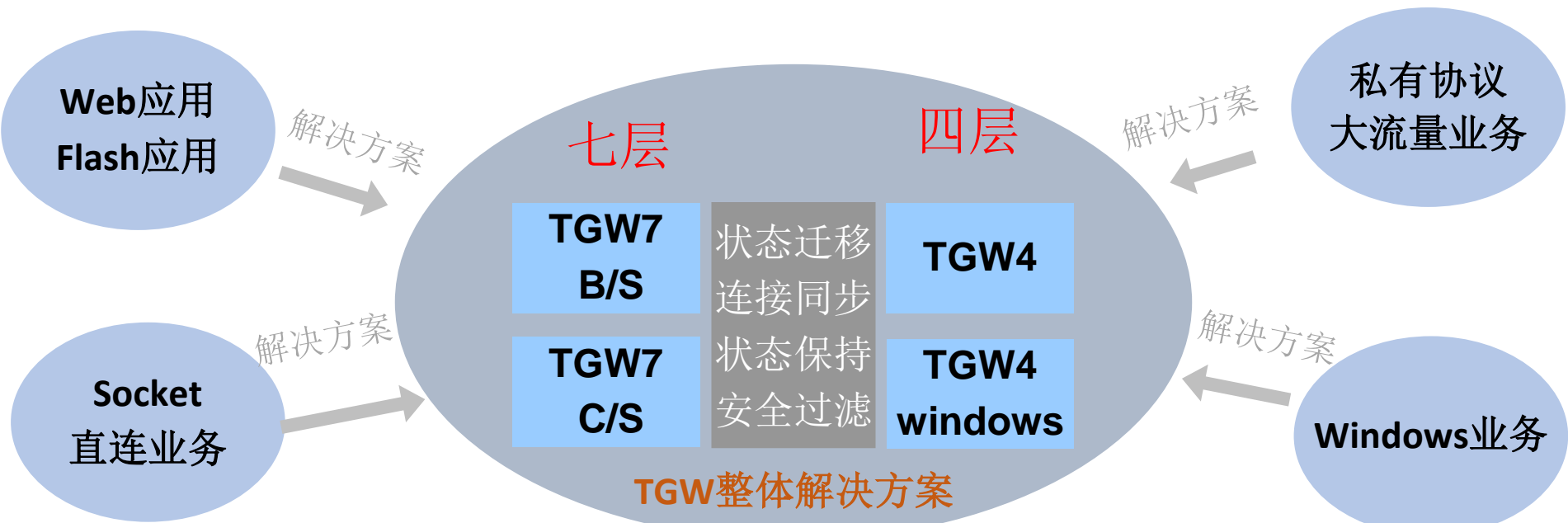


生产系统



四层、七层接入方案

四种方案：七层B/S方案，七层C/S方案，四层方案，四层windows方案。



| TGW SET模型 | LD数 | 最大容量 | 最大包量 | 接入VIP数 | 最多接入域名个数 | 最多接入RS个数 |
|-----------|-----|------|------|--------|----------|----------|
| TGW7-1G | 4 | 2G | 300w | 4 | 800 | 800 |
| TGW7-10G | 4 | 10G | 300w | 20 | 4000 | 4000 |
| TGW4-10G | 4 | 16G | 500w | N/A | N/A | N/A |
| TGW4-NAT | 4 | 1G | 150w | N/A | N/A | N/A |

单集合性能指标

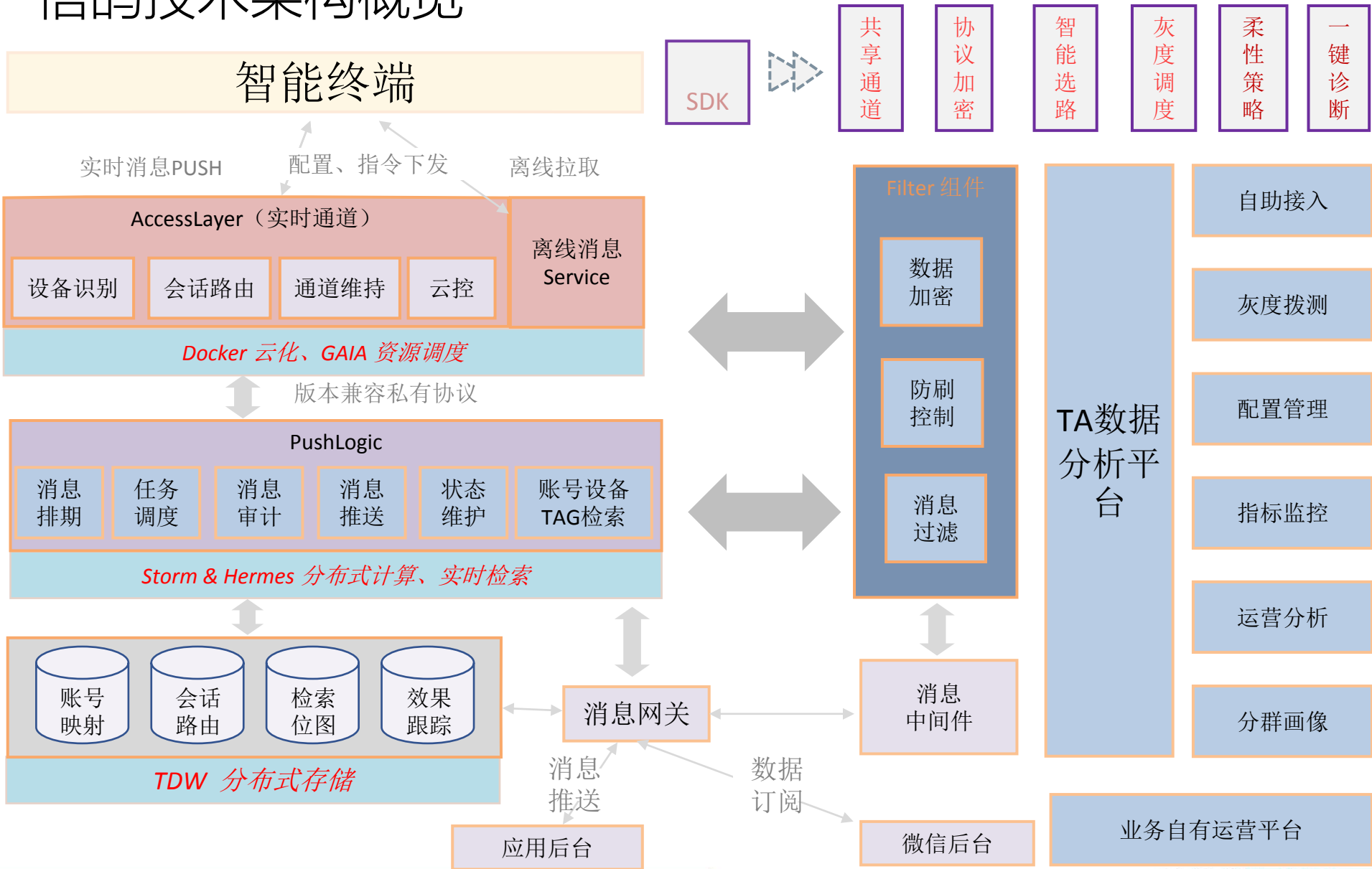
延时: < X ms

并发连接数: XXX万

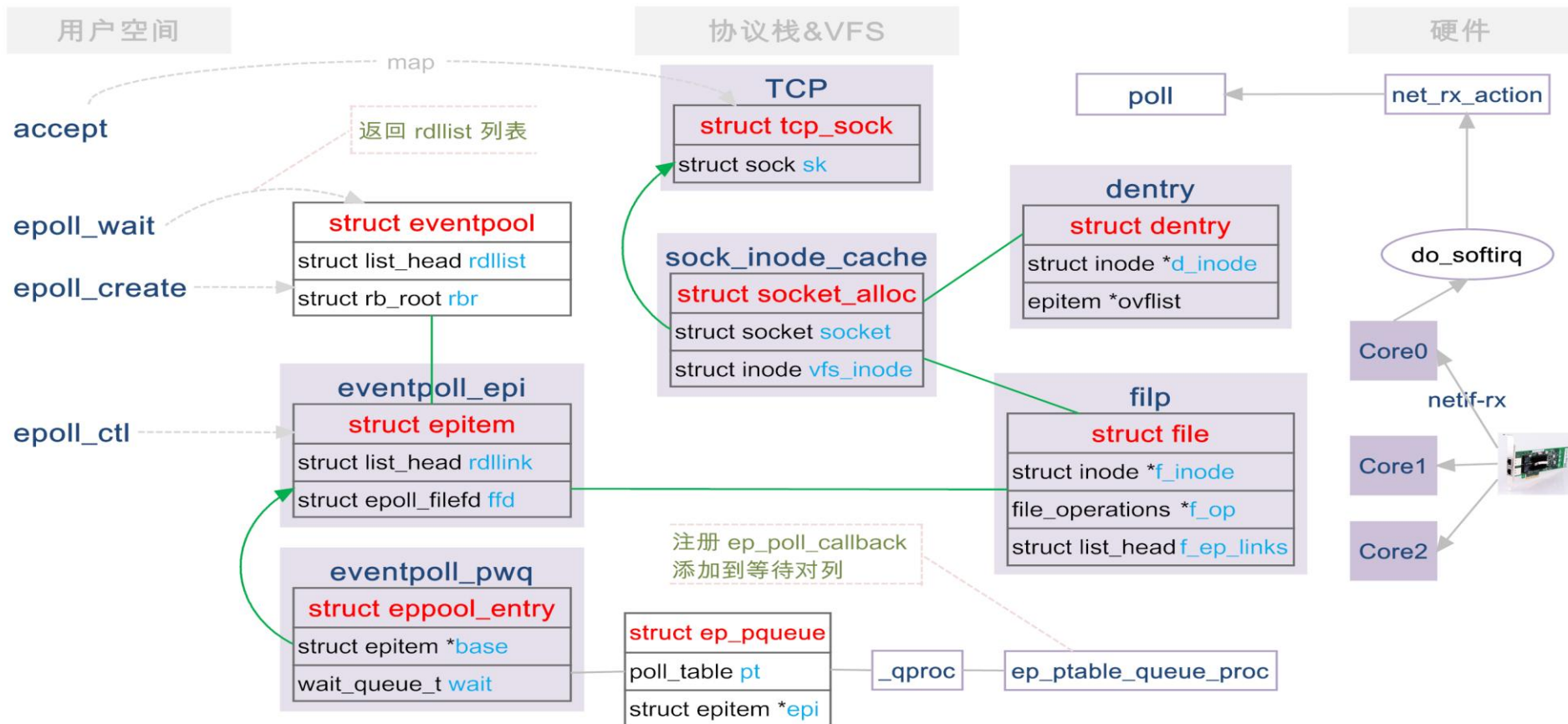
包量: XXX万 PPS

带宽: XX Gps

信鸽技术架构概览



关键的6个 kernel memory cache



单节点接入百万级长连接

C1 机型 8G 内存，支持 200W+ 长连接

```
Active / Total Objects (% used) : 15078299 / 15185017 (99.3%)
Active / Total Slabs (% used)   : 1437021 / 1437093 (100.0%)
Active / Total Caches (% used)  : 90 / 215 (41.9%)
Active / Total Size (% used)    : 6638525.50K / 6651639.44K (99.8%)
Minimum / Average / Maximum Object : 0.02K / 0.44K / 4096.00K
```

| OBJS | ACTIVE | USE | OBJ | SIZE | SLABS | OBJ/SLAB | CACHE | SIZE | NAME |
|---------|---------|-----|-------|--------|-------|----------|------------------|------|------|
| 3492460 | 3491904 | 99% | 0.19K | 174623 | 20 | 698492K | dentry | | |
| 2001840 | 2001624 | 99% | 0.19K | 100092 | 20 | 400368K | filp | | |
| 2000273 | 1999995 | 99% | 0.07K | 37741 | 53 | 150964K | eventpoll_pwq | | |
| 2000100 | 1999995 | 99% | 0.12K | 66670 | 30 | 266680K | eventpoll_epi | | |
| 2000095 | 2000084 | 99% | 0.69K | 400019 | 5 | 1600076K | sock_inode_cache | | |
| 2000020 | 1999996 | 99% | 1.44K | 400004 | 5 | 3200032K | TCP | | |
| 1490376 | 1490039 | 99% | 0.63K | 248396 | 6 | 993584K | proc_inode_cache | | |

- ◆ TCP (1.6K) : 关联了socket、sock、file
- ◆ sock_inode_cache (0.8K) : socket_alloc, 包含 socket 和 inode 节点
- ◆ eventpoll_epi (0.13K) : epitem,
- ◆ eventpoll_pwq (0.07K) : epoll_entry 负责关联fd等待事件和epitem
- ◆ filp (0.2K) : 文件指针
- ◆ dentry (0.2K) : dentry

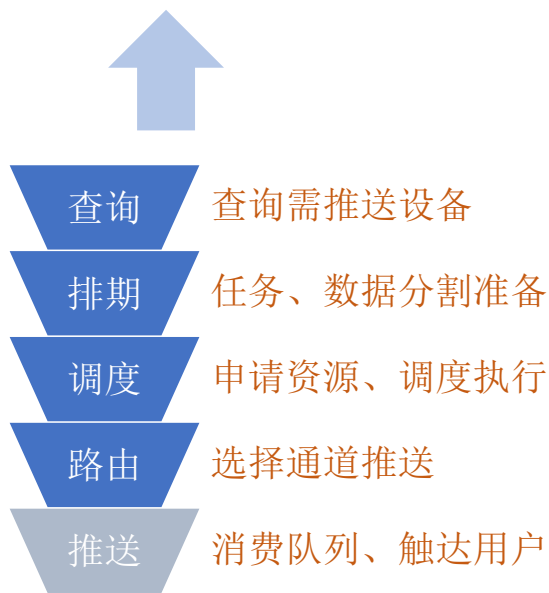
创建一个sock, 到Epoll监听, 内核消耗 3K, 其中结构体占 2.7K, 对齐或 slab 浪费 0.3K

并行检索、推送能力

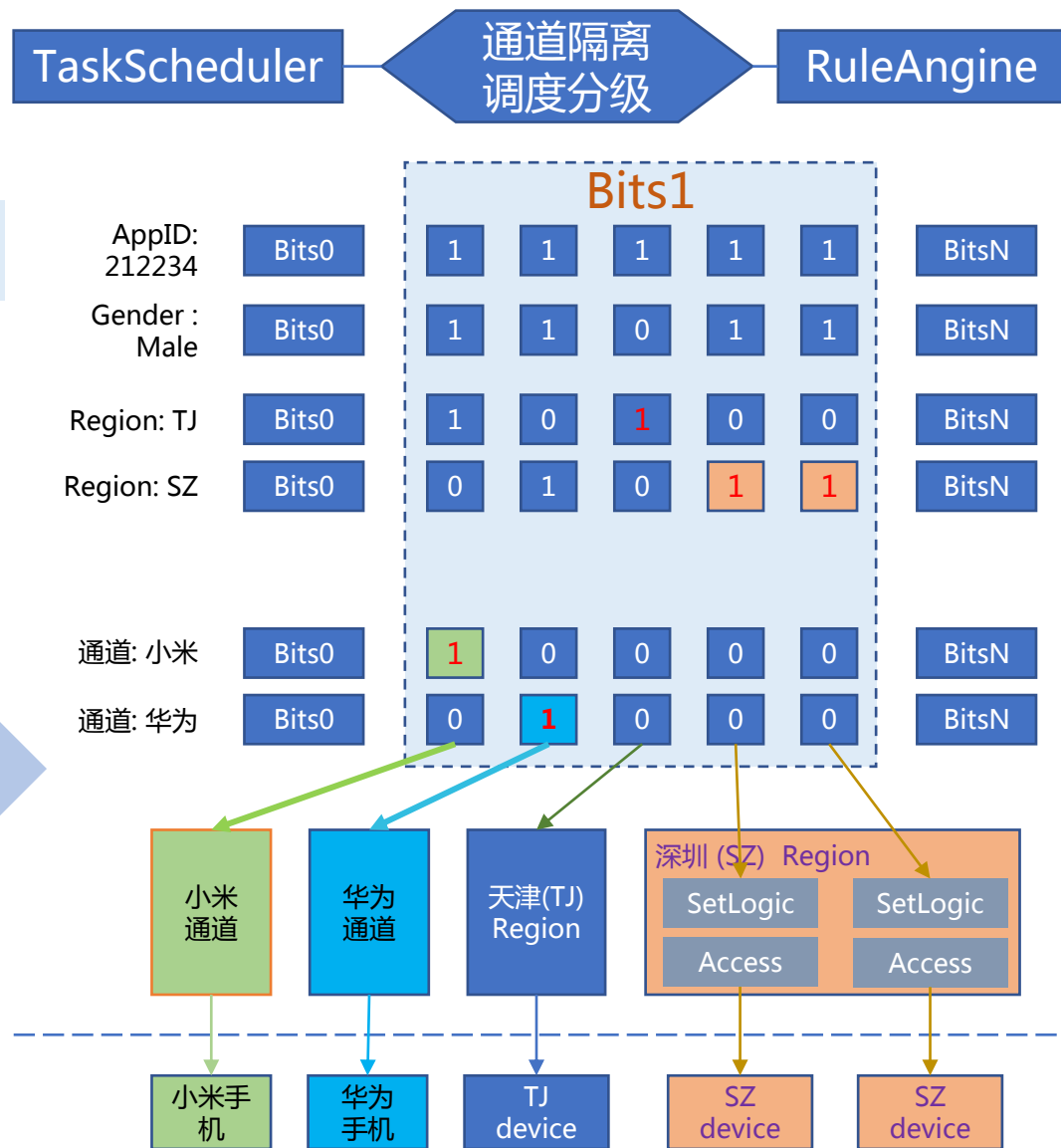
逻辑表达转换为析取范式

$$(A \cup B) \cap C \Rightarrow (A \cap C) \cup (B \cap C)$$

工程问题映射成数学问题

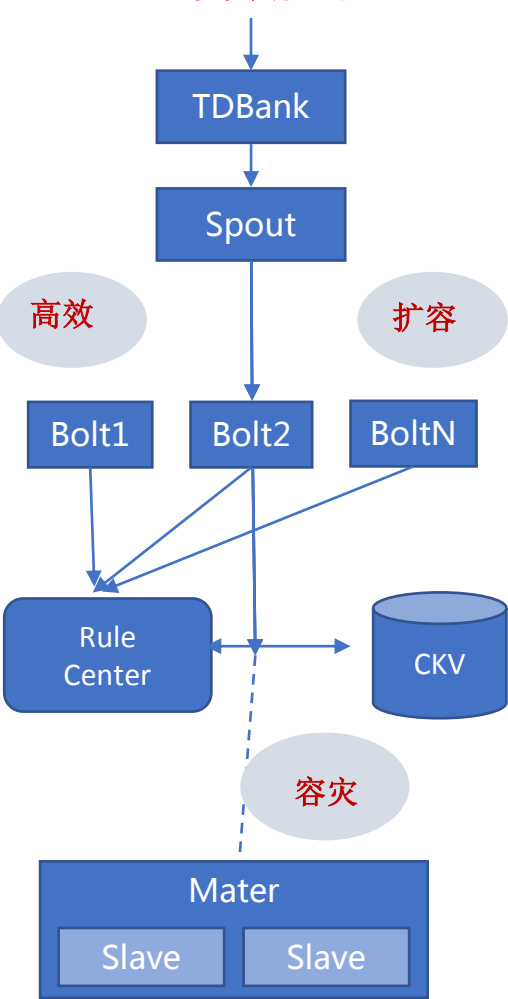


数据传输、计算、查询量大。
逐层下漏，层之间容易形成瓶颈。

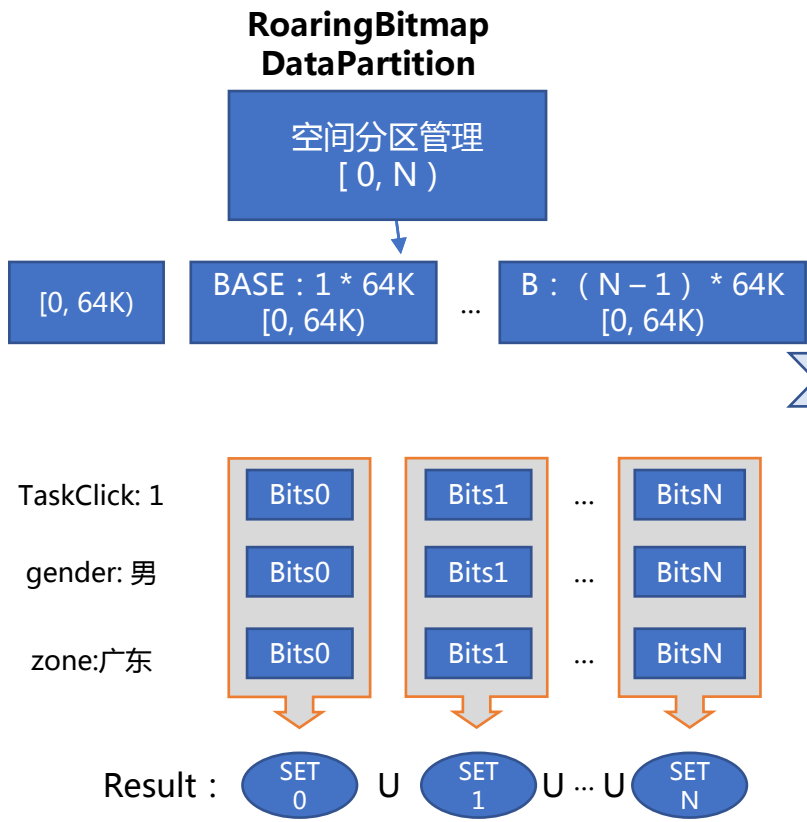


推送效果实时统计

文本数值化



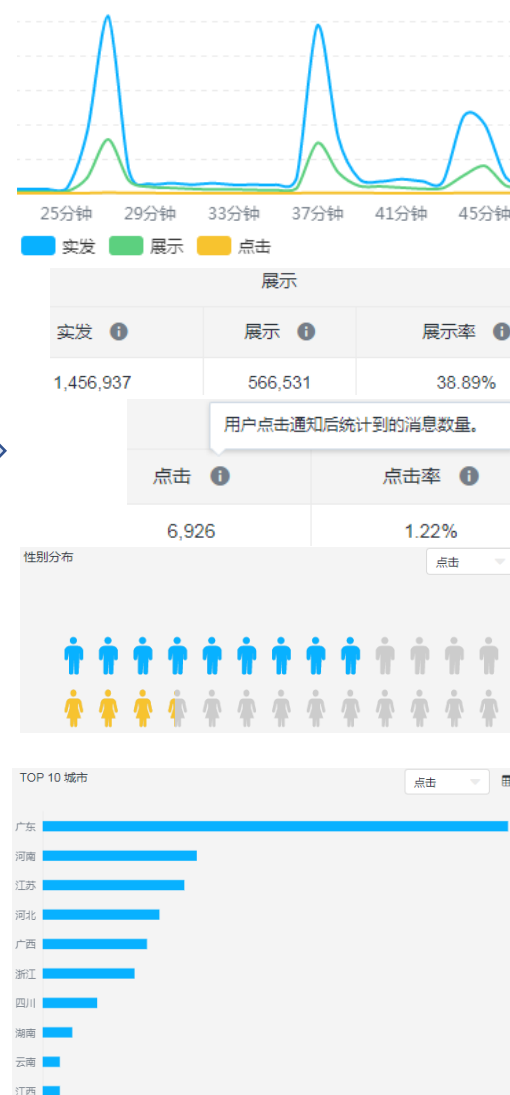
高效数据结构



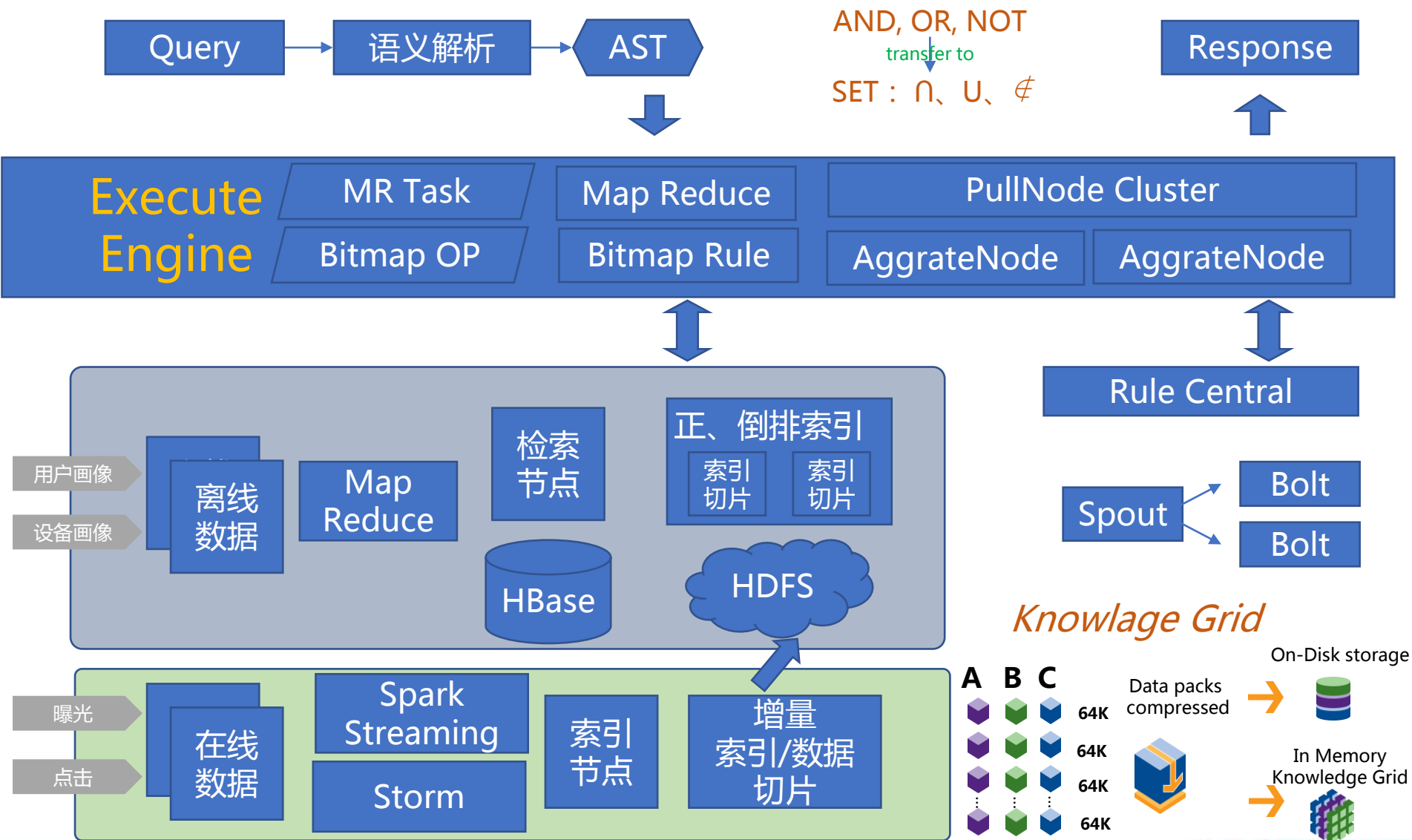
广东、男性用户点击

拿设备标识举例：40 字节，320 bit 到 1 bit 的压缩

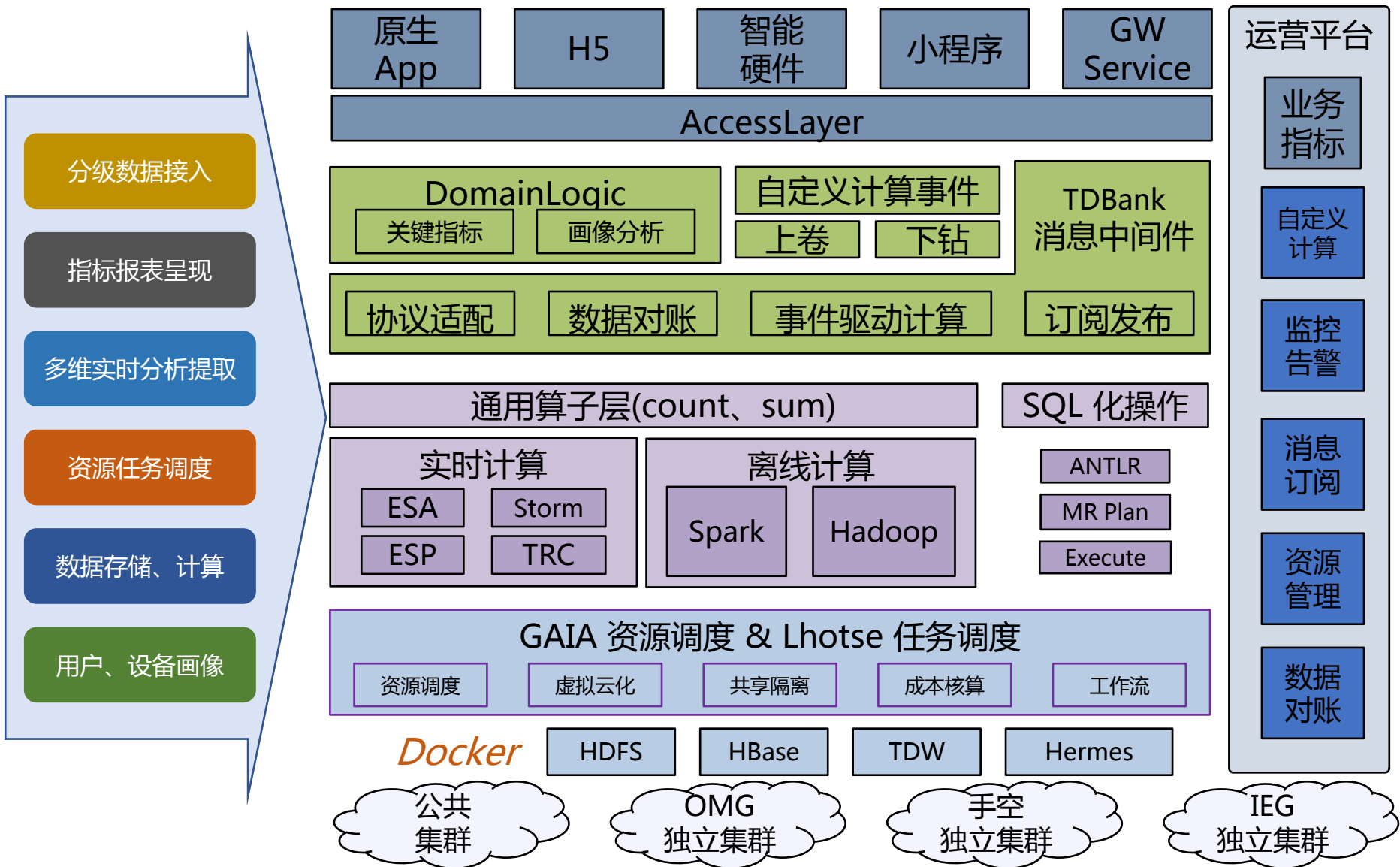
实时效果统计



海量数据，实时多维分析



数据分析平台



DevOps



配置中心

L5

服务发现负载均衡



DevOps



GitLab



Jenkins

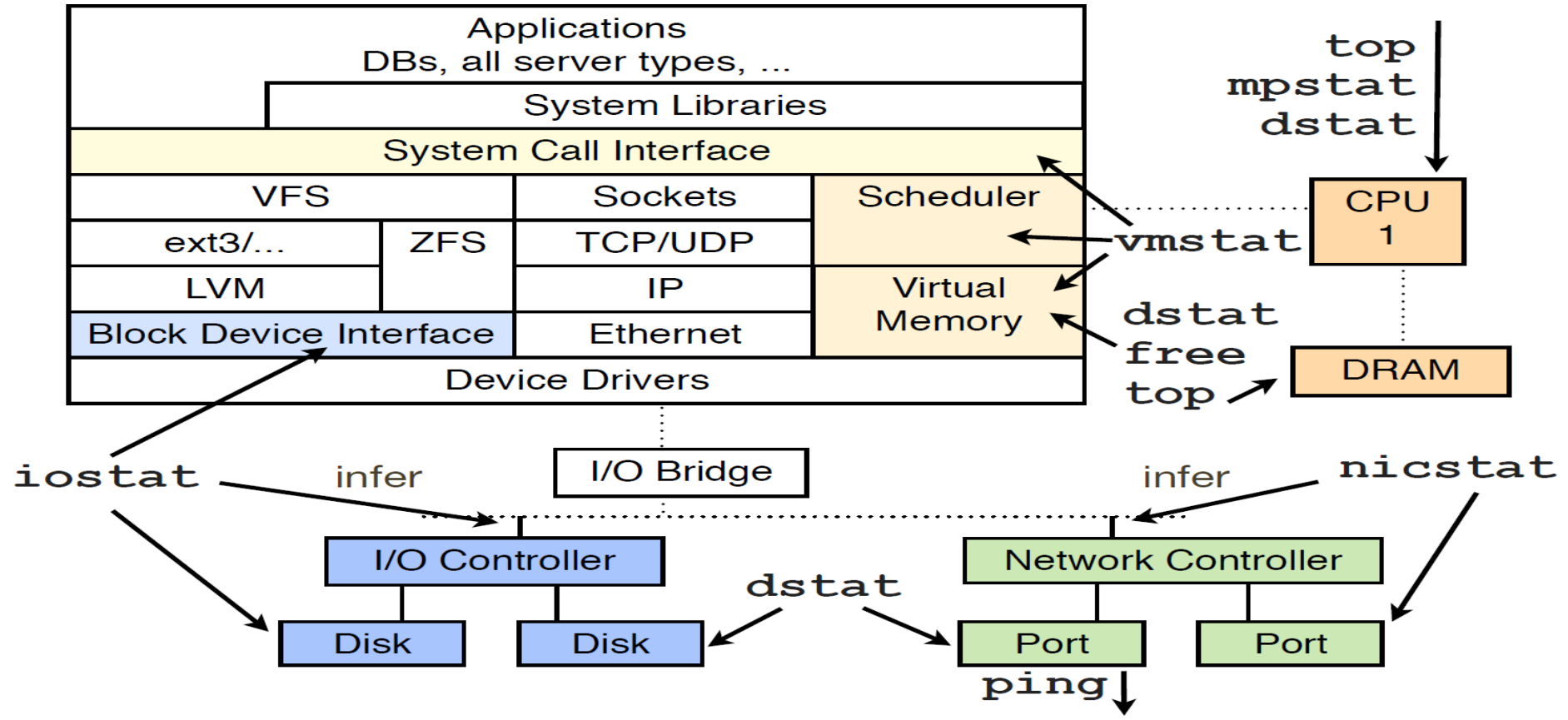


docker

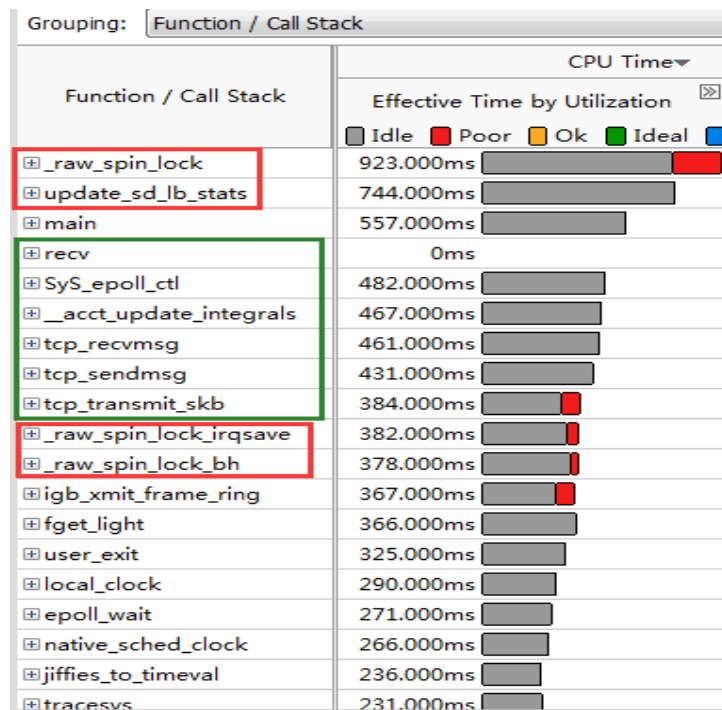
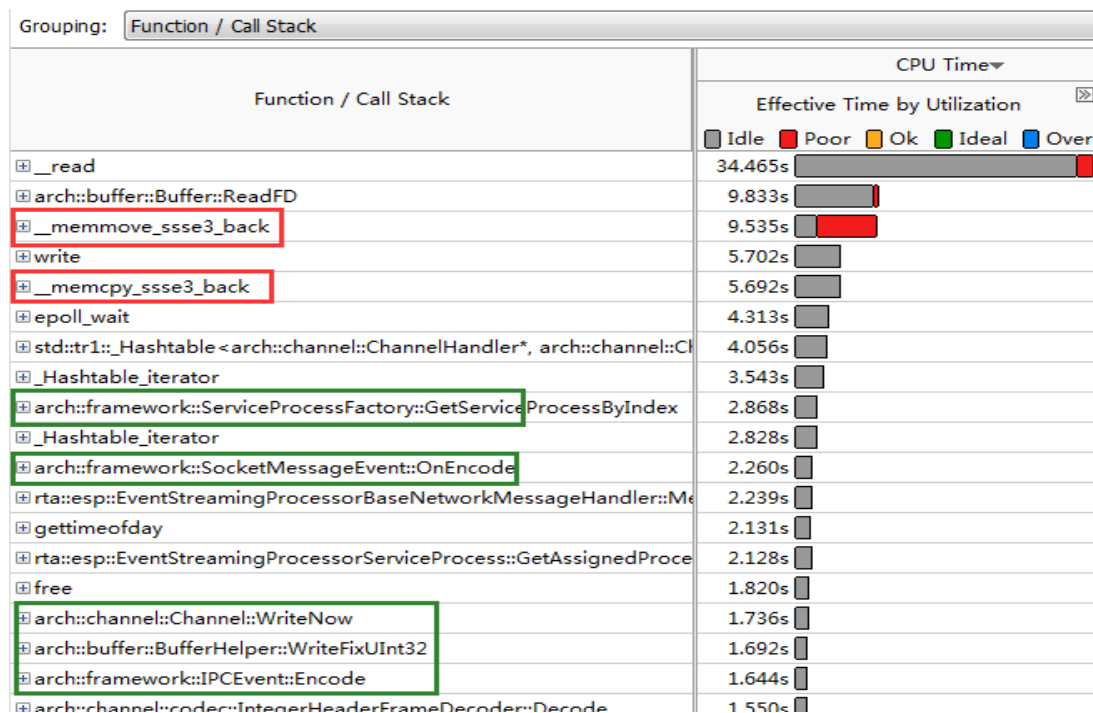


kubernetes





摘自：Systems Performance Enterprise and the Cloud



较好的网络框架，耗时应该主要集中在网络的处理上，而非框架本身。
基于 **Linux协议栈** 网络服务，大部分性能消耗在 **spin_lock & 网络协议栈** 上。

单机性能挖掘 — 评测

amplxe-cl -collect advanced-hotspots -result-dir rs01 ./access_proxy

| Function / Call Stack | CPU Time | | | | CPI Rat. | CPU Fre... Ratio | Module | Function (Full) | |
|---------------------------|-------------------------------|------|----|-------|----------|------------------|--------|-----------------|------------------------|
| | Effective Time by Utilization | | | | | | | | |
| | Idle | Poor | Ok | Ideal | Over | | | | |
| 1 _raw_spin_lock | 722.000ms | | | | | 2.6.. | 0.853 | vmlinux | _raw_spin_lock |
| 2 Sys_epoll_ctl | 379.000ms | | | | | 3.7.. | 0.592 | vmlinux | Sys_epoll_ctl |
| 3 update_sd_lb_stats | 235.000ms | | | | | 1.3.. | 0.682 | vmlinux | update_sd_lb_stats |
| 4 __acct_update_integrals | 221.000ms | | | | | 0.7.. | 0.566 | vmlinux | __acct_update_integ... |
| 5 _raw_spin_lock_bh | 203.000ms | | | | | 1.9.. | 0.582 | vmlinux | _raw_spin_lock_bh |
| 6 user_exit | 191.000ms | | | | | 1.0.. | 0.514 | vmlinux | user_exit |
| 7 kmem_cache_alloc | 170.000ms | | | | | 1.8.. | 0.448 | vmlinux | kmem_cache_alloc |
| 8 main | 170.000ms | | | | | 3.2.. | 0.607 | access_pr... | main |
| 9 _raw_spin_lock_irqsave | 167.000ms | | | | | 1.4.. | 0.743 | vmlinux | _raw_spin_lock_irqs... |
| 10 epoll_wait | 163.000ms | | | | | 43.. | 0.528 | lib... | epoll_wait |
| 11 jiffies_to_timeval | 155.000ms | | | | | 0.8.. | 0.452 | vm... | jiffies_to_timeval |
| 12 native_sched_clock | 127.000ms | | | | | 1.4.. | 0.583 | vm... | native_sched_clock |
| 13 local_clock | 122.000ms | | | | | 0.5.. | 0.591 | vm... | local_clock |
| 14 ip_queue_xmit | 121.000ms | | | | | 3.0.. | 1.100 | vm... | ip_queue_xmit |
| 15 tcp_transmit_skb | 117.000ms | | | | | 1.0.. | 1.138 | vm... | tcp_transmit_skb |
| 16 MD::BinUtil::crc32 | 116.000ms | | | | | 0.9.. | 0.699 | ac... | MD::BinUtil::crc32 |
| 17 igb_xmit_frame_ring | 111.000ms | | | | | 0.6.. | 1.001 | igb... | igb_xmit_frame_ring |
| 18 sched_clock_cpu | 106.000ms | | | | | 0.6.. | 0.482 | vm... | sched_clock_cpu |
| 19 tcp_recvmg | 103.000ms | | | | | 1.6.. | 0.573 | vm... | tcp_recvmg |

perf top -p pid

```

Samples: 377K of event 'cycles', Event count (approx.): 797075664
3.10% [kerne] [k] update_sd_lb_stats
2.46% [kerne] [k] _raw_spin_lock
2.45% access_proxy [.] main
2.20% [kerne] [k] fget_light
2.03% [kerne] [k] sys_epoll_ctl
2.01% [kerne] [k] tcp_sendmsg
1.84% [kerne] [k] tcp_recvmg
1.70% [kerne] [k] _raw_spin_lock_bh
1.55% [kerne] [k] tcp_transmit_skb
1.41% [igb] [k] igb_xmit_frame_ring
1.20% libc-2.17.so [.] __memcpy_ssse3_back
1.01% [kerne] [k] _raw_spin_lock_irqsave
0.93% [kerne] [k] ep_send_events_proc
0.93% [ip_tables] [k] ipt_do_table
0.92% [kerne] [k] tcp_write_xmit
0.90% access_proxy [.] handleMessage(InternetMessag
0.89% [kerne] [k] find_next_bit
0.89% [kerne] [k] system_call
0.87% [kerne] [k] __schedule
0.86% [kerne] [k] update_blocked_averages
0.79% [kerne] [k] tcp_poll
0.76% [kerne] [k] tcp_cleanup_rbuf
    
```

优化目标:

- ✓ 优化函数执行效率
- ✓ 消除内核协议栈消耗
- ✓ 挖掘硬件潜能


```
-----  
Interfaces:  
name: eth0, ifindex: 0, hwaddr: 54:89:98:72:5B:1E, ipaddr: 10.125.225.3, netmask: 255.255.  
Number of NIC queues: 4  
-----  
Loading routing configurations from : config/route.conf  
fopen: No such file or directory  
Skip loading static routing table  
Routes:  
Destination: 10.125.225.0/32, Mask: 255.255.255.192, Masked: 10.125.225.0, Route: ifdx-0  
-----  
Loading ARP table from : config/arp.conf  
fopen: No such file or directory  
Skip loading static ARP table  
ARP Table:  
(blank)  
EAL: Master lcore 0 is ready (tid=ee0d18c0;cpuset=[0])  
EAL: lcore 1 is ready (tid=ebbe7700;cpuset=[1])  
EAL: lcore 2 is ready (tid=eb3e6700;cpuset=[2])  
EAL: lcore 3 is ready (tid=eabe5700;cpuset=[3])  
EAL: PCI device 0000:01:00.0 on NUMA socket -1  
EAL: probe driver: 8086:10c9 rte_igb_pmd  
EAL: PCI memory mapped at 0x7fc2ec800000  
EAL: PCI memory mapped at 0x7fc2ec820000  
EAL: PCI memory mapped at 0x7fc2ec840000  
PMD: eth_igb_dev_init(): port_id 0 vendorID=0x8086 deviceID=0x10c9  
EAL: PCI device 0000:01:00.1 on NUMA socket -1  
EAL: probe driver: 8086:10c9 rte_igb_pmd  
EAL: Not managed by a supported kernel driver, skipped  
Interactive-mode selected  
Configuring Port 0 (socket 0)  
PMD: eth_igb_tx_queue_setup(): To improve 1G driver performance, consider setting the TX W  
PMD: eth_igb_tx_queue_setup(): sw_ring=0x7fc1ec5a7e40 hw_ring=0x7fc1ec5a9e80 dma_addr=0x22  
PMD: eth_igb_rx_queue_setup(): sw_ring=0x7fc1ec5978c0 hw_ring=0x7fc1ec597d00 dma_addr=0x22  
PMD: eth_igb_start(): <<  
Port 0: 54:89:98:72:5B:1E  
Checking link statuses...  
Port 0 Link Up - speed 1000 Mbps - full-duplex  
Done
```


协议栈改造

A : Kernel Socket 函数调用

```
int iSock = socket(PF_INET, SOCK_STREAM, 0);
```

B : mTCP 的函数函数调用

```
int iSock = mtcp_socket(g_pMctx, PF_INET, SOCK_STREAM, 0);
```

C : 业务无感知的 mTCP 协议栈改造

```
typedef int (mTcpHookSock*)(int socket_family, int socket_type, int protocol);  
int realSocketFunc(int socket_family, int socket_type, int protocol) {  
    return mtcp_socket(g_pMctx, socket_family, socket_type, protocol);  
}  
  
typedef struct {  
    mTcpHookSock    hookSockFunc;  
    mTcpHookAccept  hookAcceptFunc;  
} stMTCPHooFuncTable;  
  
g_syncHookTable.hookSockFunc = dlsym(RTLD_NEXT, "socket");
```

拦截系统调用
指向 mtcp_socket WRAP 函数

Table
Lookup

| | | |
|---------------------|-----------|-------------|
| local_clock | 122.000ms | 308,660,000 |
| ip_queue_xmit | 121.000ms | 111,320,000 |
| tcp_transmit_skb | 117.000ms | 331,430,000 |
| MD::BinUtil::crc32 | 116.000ms | 212,520,000 |
| igb_xmit_frame_ring | 111.000ms | 427,570,000 |
| sched_clock_cpu | 106.000ms | 199,870,000 |
| tcp_recvmmsg | 103.000ms | 88,550,000 |

SSE 4.2
_mm_crc32_u32

| | | |
|------------------------|----------|------------|
| _errno_location | 0ms | 5,060,000 |
| update_cfs_shares | 20.000ms | 15,180,000 |
| retransmits_timed_out | 20.000ms | 32,890,000 |
| MD::BinUtil::crc32csse | 20.000ms | 43,010,000 |
| ipv4_dst_check | 20.000ms | 43,010,000 |
| ip_output | 20.000ms | 50,600,000 |
| inode_init_once | 20.000ms | 22,770,000 |

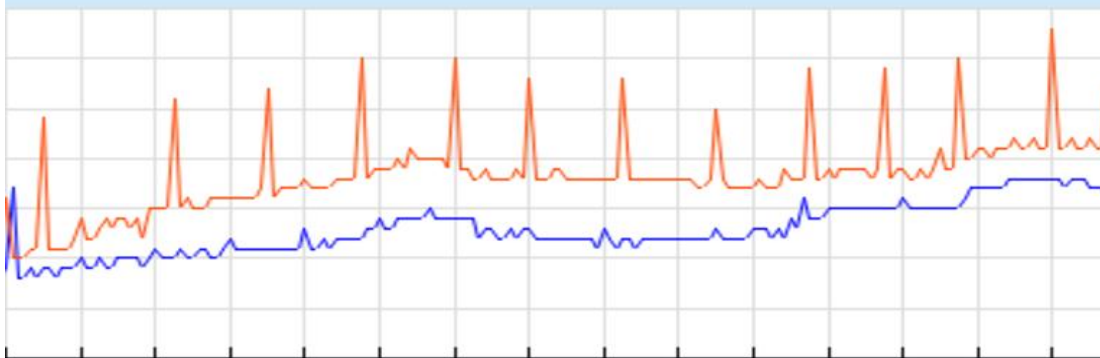
SSE 4.2
Parallel crc32q

执行 35W 次计算结果 (SSE Parallel crc32q 快接近 6 倍)

__memcpy_sse3_back & __memcpy_avx_unaligned、AES_cfb8_encrypt & evp_encrypt (M10 机型)

10. 225. 169. 34-CPU使用率

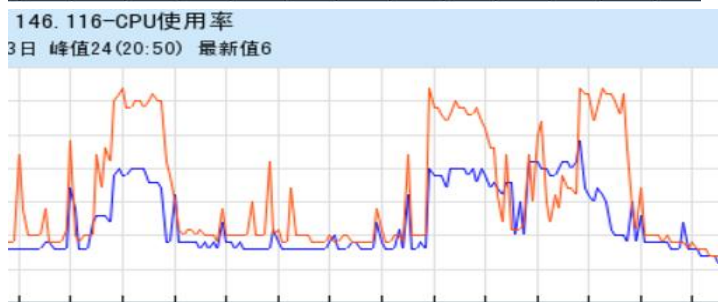
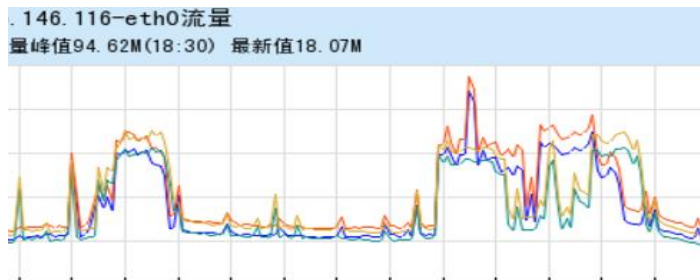
2017年09月01日 峰值18 (20:25) 最新值11



```
typedef uint32_t (CRC32FUNC)(const char* pBuf, int iLen);  
  
typedef struct {  
    CRC32FUNC    crc32;  
} stOptimizeLayer;  
  
uint32_t eax, ebx, ecx, edx;  
  
eax = ebx = ecx = edx = 0;  
__get_cpuid(1, &eax, &ebx, &ecx, &edx);  
if ( ecx & bit_SSE4_2 ) {  
    opts.crc32 = sse42_crc32;  
} else {  
    opts.crc32 = default_crc32;  
}  
  
-----  
  
uint32_t crc32(const char* pBuf, int iLen)  
    __attribute__((ifunc("resolve_crc32")));
```

硬件虚拟化提升效能

docker & SR-IOV, 高峰 CPU 降低约 38.3%



信鸽、MTA Docker 云化部署

1. 混合部署、资源隔离
2. AutoScaling, 自动加入 TGW、L5
3. Floating IP
4. 网络层启用 SR-IOV 特性

xg-access-sz-on-newc1 运行中

实例数及状态数: 40 / 运行中: 40
应用类型: 通用服务器App
访问地址: http://10.175.95.12:8081/abnormal/xg-access-sz-on-newc1/

| 运行状态 | 实例名称 | Pod IP | 策略 | Host IP | 策略 | 部署策略 |
|------|--------------------------|----------------|----|---------------|----|-------------------------------------|
| 运行中 | xg-access-sz-on-newc1-30 | 10.49.94.43 | | 10.49.94.38 | | docker.oa.com:8080/xg/client-access |
| 运行中 | xg-access-sz-on-newc1-32 | 10.120.82.241 | | 10.120.82.207 | | docker.oa.com:8080/xg/client-access |
| 运行中 | xg-access-sz-on-newc1-22 | 10.120.96.176 | | 10.120.96.147 | | docker.oa.com:8080/xg/client-access |
| 运行中 | xg-access-sz-on-newc1-35 | 10.121.124.113 | | 10.121.124.68 | | docker.oa.com:8080/xg/client-access |

xg-access-sz-on-newc1 运行中

实例数及状态数: 40 / 运行中: 40
应用类型: 通用服务器App
访问地址: http://10.175.95.12:8081/abnormal/xg-access-sz-on-newc1/

成功 创建 2017-07-04 17:20:10 patxu 成功创建了应用xg-access-sz-on-newc1

成功 扩容 2017-07-04 17:34:49 dreamxguo 成功扩容(1 → 15)了应用xg-access-sz-on-newc1

成功 扩容 2017-07-07 17:16:02 dreamxguo 成功扩容(15 → 25)了应用xg-access-sz-on-newc1

成功 扩容 2017-07-10 09:05:27 dreamxguo 成功扩容(25 → 27)了应用xg-access-sz-on-newc1

成功 扩容 2017-07-12 09:24:42 dreamxguo 成功扩容(27 → 28)了应用xg-access-sz-on-newc1

成功 扩容 2017-07-12 10:47:43 dreamxguo 成功扩容(28 → 40)了应用xg-access-sz-on-newc1

磁盘 10GB

xg-access-sz-on-newc1-30

运行状态: 运行中

所属: 业务 xg 集群 abnormal 应用 xg-access-sz-on-newc1
部署策略: docker.oa.com:8080/xg/client-access
实例配置: 7500m个 10GB 7680MB 0bit/s 0个
pod IP: 10.49.94.43

```
[root@xg-access-sz-on-newc1-30 /]# ifconfig
eth0: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 1500
    inet 10.49.94.43 netmask 255.255.255.192 broadcast 0.0.0.0
    ether 02:42:0a:31:5e:2b txqueuelen 0 (Ethernet)
    RX packets 39455727 bytes 5592839306 (5.2 GiB)
    RX errors 0 dropped 7708 overruns 0 frame 0
    TX packets 35963968 bytes 4429021642 (4.1 GiB)
    TX errors 0 dropped 17511 overruns 0 carrier 0 collisions 0

lo: flags=73<UP,LOOPBACK,RUNNING> mtu 65536
    inet 127.0.0.1 netmask 255.0.0.0
    loop txqueuelen 0 (Local Loopback)
    RX packets 7840268 bytes 742571057 (708.1 MiB)
    RX errors 0 dropped 0 overruns 0 frame 0
    TX packets 7840268 bytes 742571057 (708.1 MiB)
    TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
```

数据路径: /data

网络类型: Floating IP (浮动IP)

信鸽精准推送服务

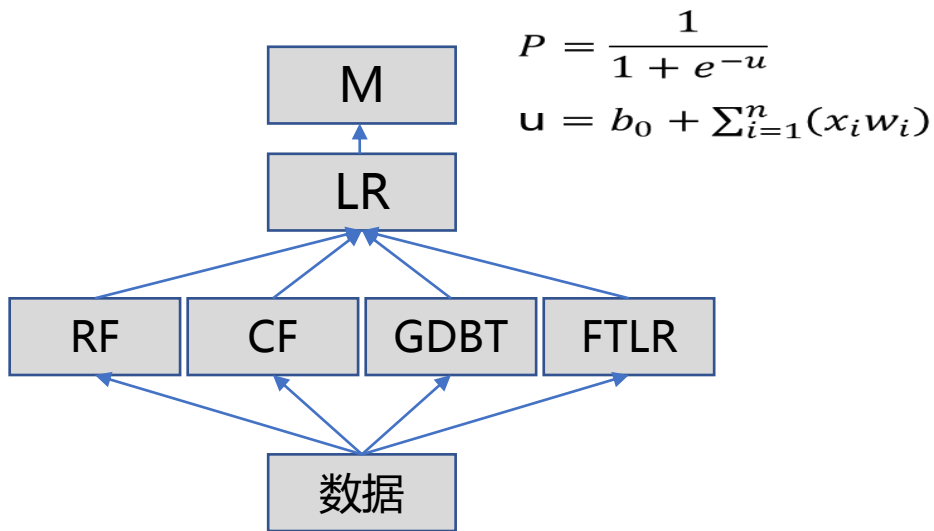
标签推送

A/B 推送

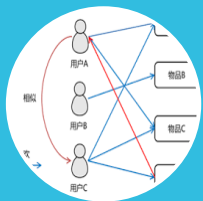
用户分群推送

智能精准推送

Hybrid Model



规则引擎



协同过滤



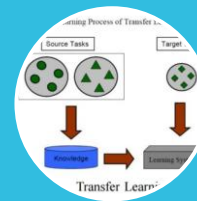
预估模型



深度学习

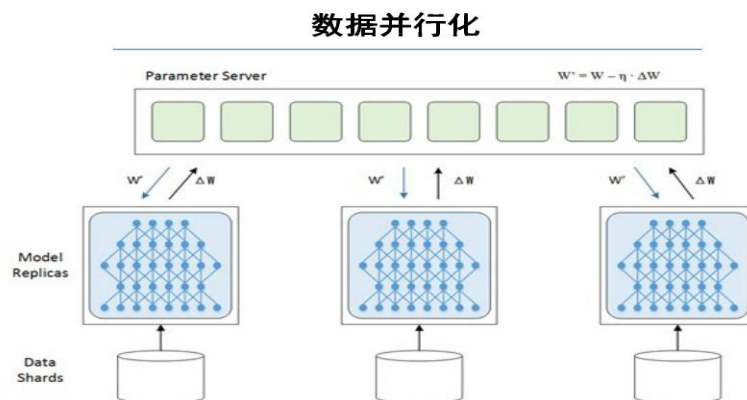
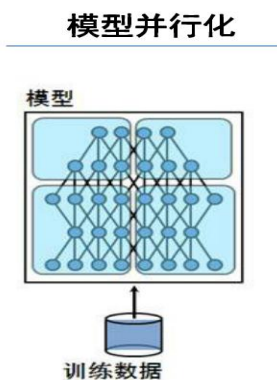
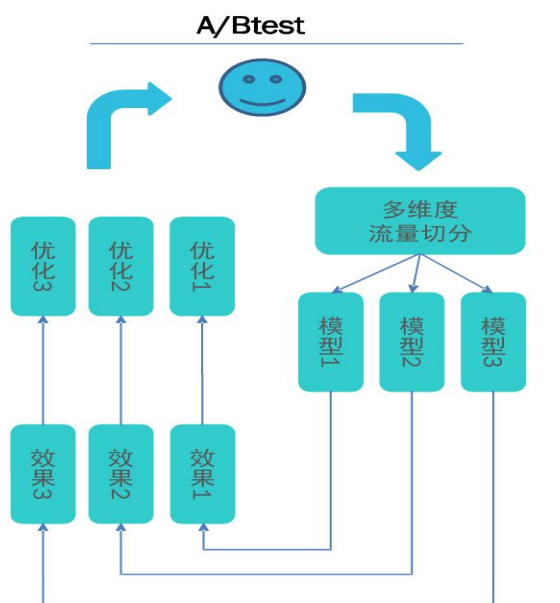


在线学习



迁移学习

精准推送服务基本流程



丰富的产品线和数据积累



月活跃用户超 8 亿
在线人际关系链超 1000 亿

账号信息

年龄, 性别, 学历...

UGC

QZone日志, 图像, 评论...



注册用户数超 10 亿
月活跃超 7 亿



活跃用户数超 6 亿
智能终端活跃用户数超 5 亿



移动支付超过 5 亿笔



其他: 新闻, 视频, 音乐, 游戏, 手机管家

社交属性

QQ, Qzone, 好友关系链...

行为偏好

新闻, 视频, 公众号...

金融

理财, 消费, 信用...

商业兴趣

广告, 商品, 品牌...

游戏爱好

游戏时长, 类型, 付费...

连接

分析

标准化

账号



DNN



结构化



PGM



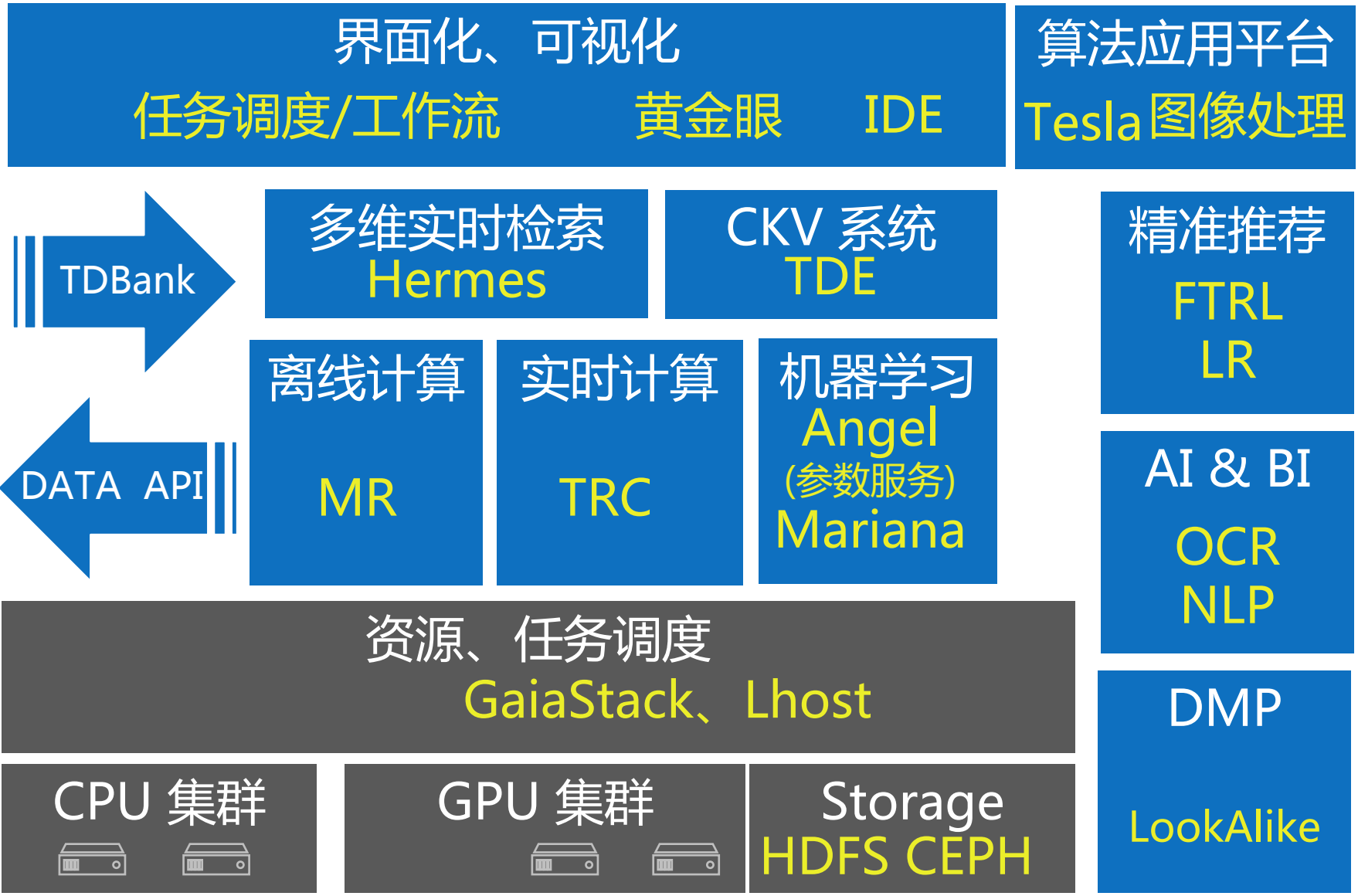
LPA



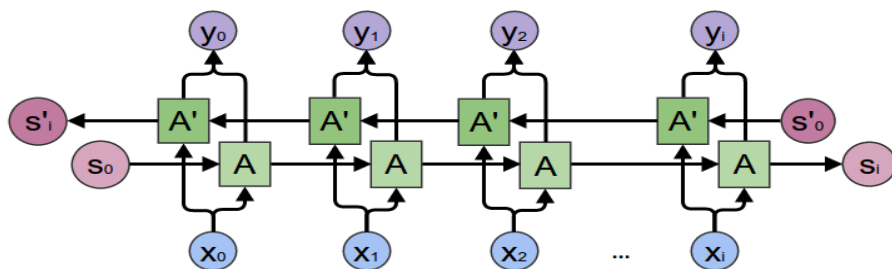
非结构



业务在线系统



机器学习可视化界面操作



for $t = 1$ to T do
 Forward pass for the forward hidden layer, storing activations at each timestep
 for $t = T$ to 1 do
 Forward pass for the backward hidden layer, storing activations at each timestep
 for all t , in any order do
 Forward pass for the output layer, using the stored activations from both hidden layers

for all t , in any order do
 Backward pass for the output layer, storing δ terms at each timestep
 for $t = T$ to 1 do
 BPTT backward pass for the forward hidden layer, using the stored δ terms from the output layer
 for $t = 1$ to T do
 BPTT backward pass for the backward hidden layer, using the stored δ terms from the output layer

$$a_i^t = \sum_{i=1}^I w_{ii} x_i^t + \sum_{h=1}^H w_{hi} b_h^{t-1} + \sum_{c=1}^C w_{ci} s_c^{t-1}$$

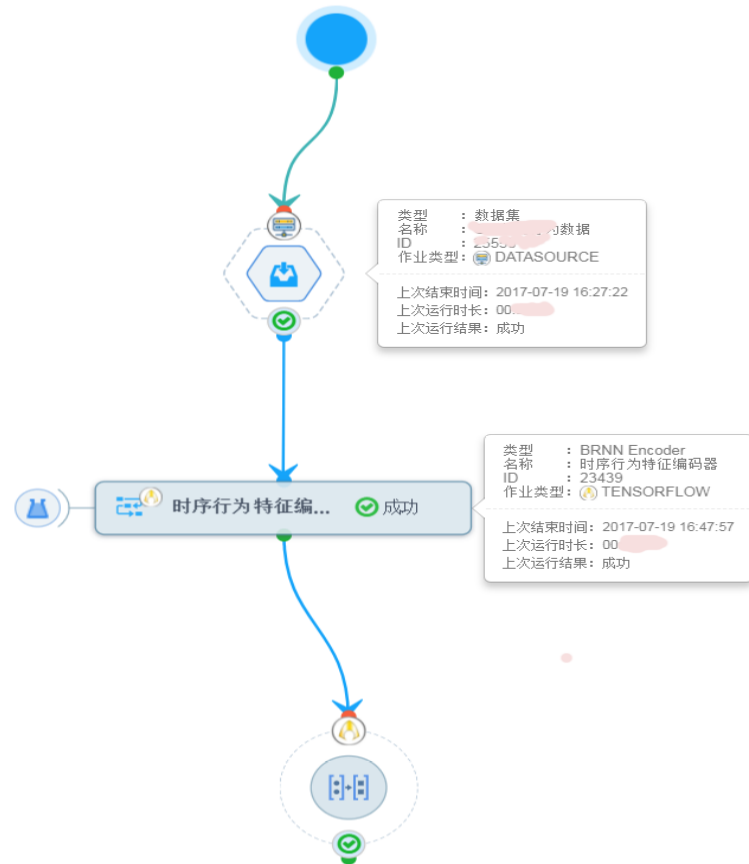
$$a_h^t = \sum_{i=1}^I w_{ih} x_i^t + \sum_{h'=1}^H w_{h'h} b_{h'}^{t-1}$$

$$b_h^t = \theta(a_h^t)$$

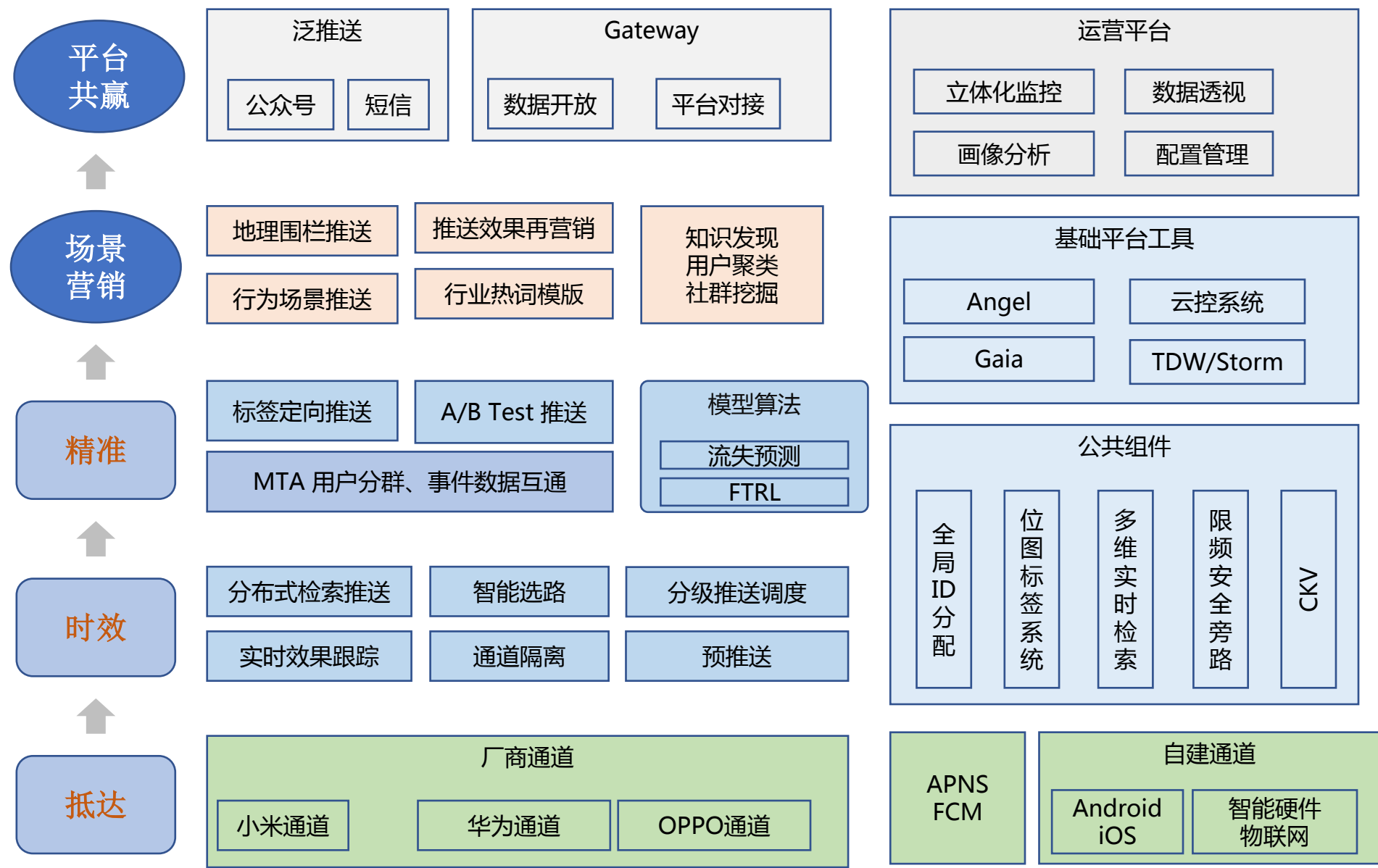
$$a_\phi^t = \sum_{i=1}^I w_{i\phi} x_i^t + \sum_{h=1}^H w_{h\phi} b_h^{t-1} + \sum_{c=1}^C w_{c\phi} s_c^{t-1}$$

$$a_\omega^t = \sum_{i=1}^I w_{i\omega} x_i^t + \sum_{h=1}^H w_{h\omega} b_h^{t-1} + \sum_{c=1}^C w_{c\omega} s_c^t$$

点击行为预测：BRNN (TESLA ML)

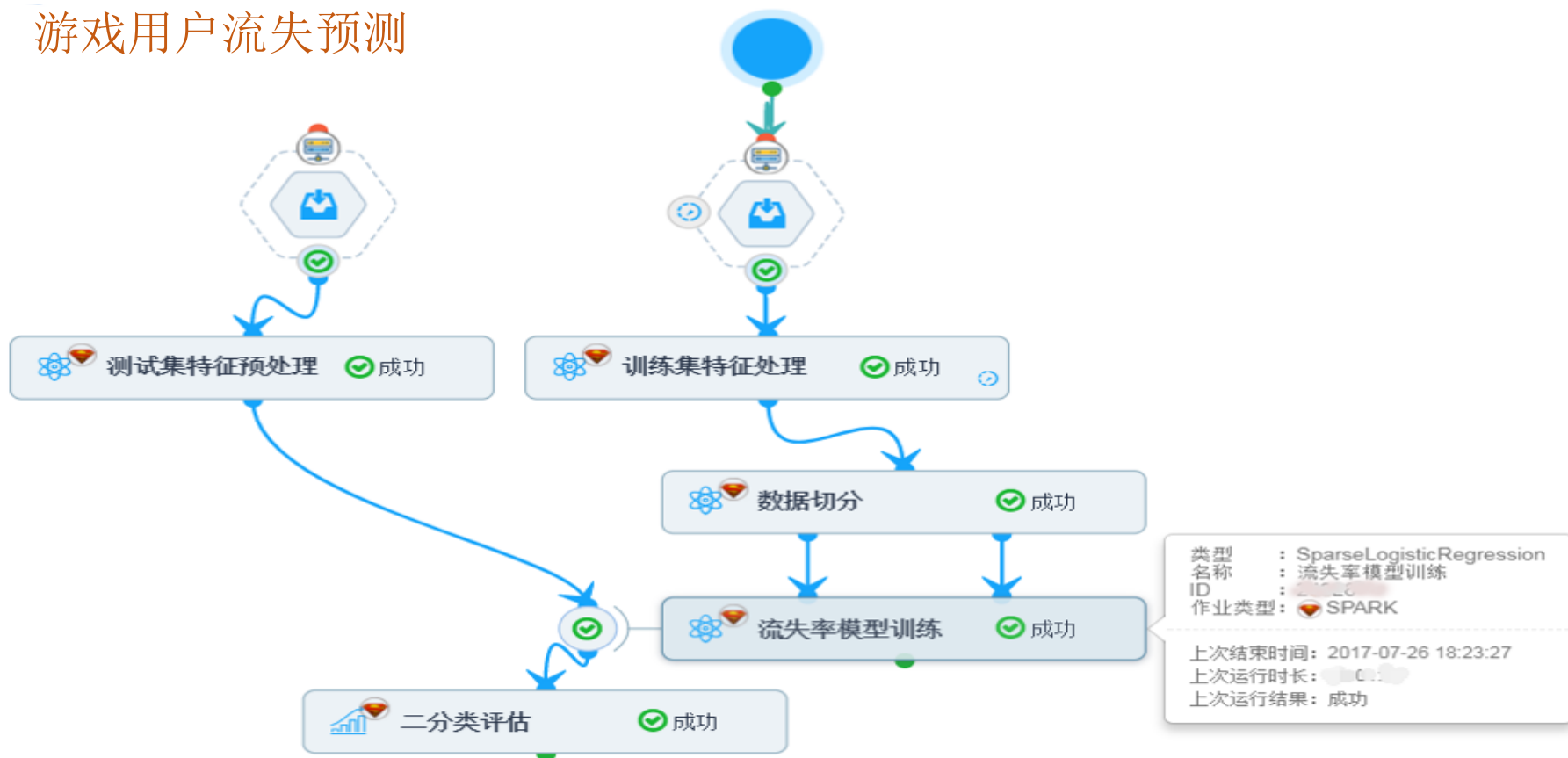


信鸽完整技术解决方案



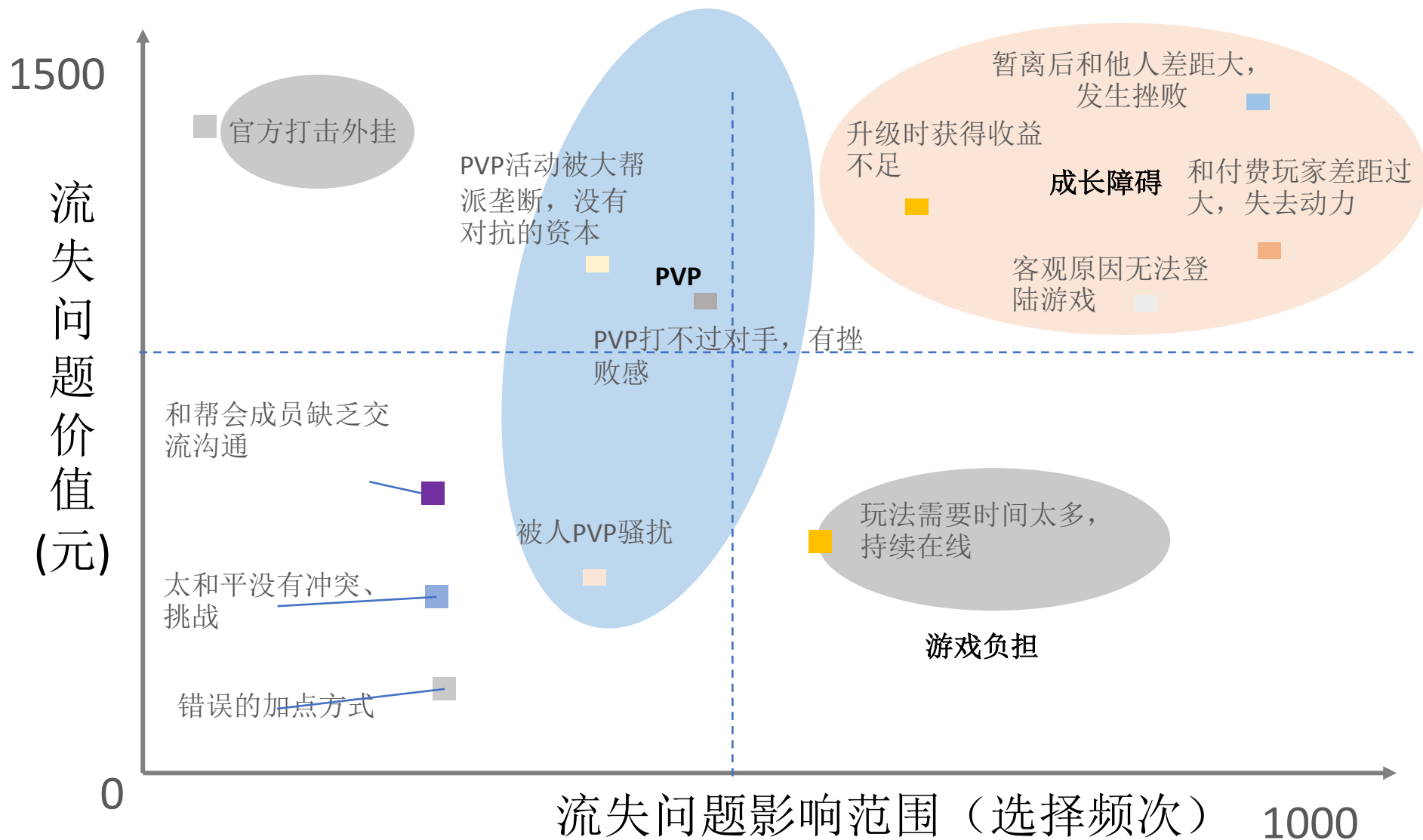
一个案例：提升游戏留存

游戏用户流失预测



一个案例：提升游戏留存

用户流失归因分析



一个案例：提升游戏应用留存

信鸽
移动推送服务

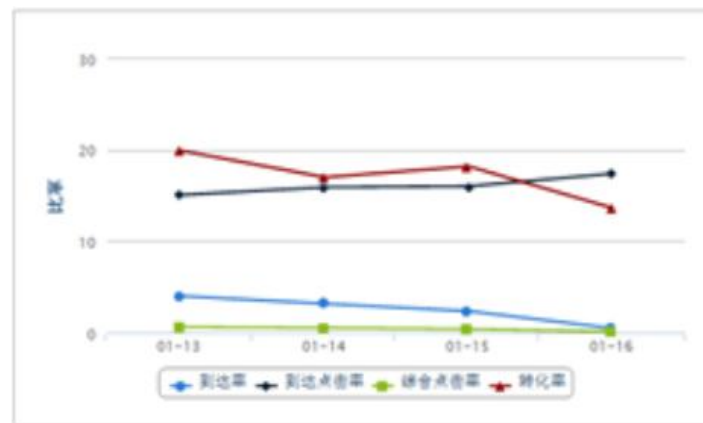
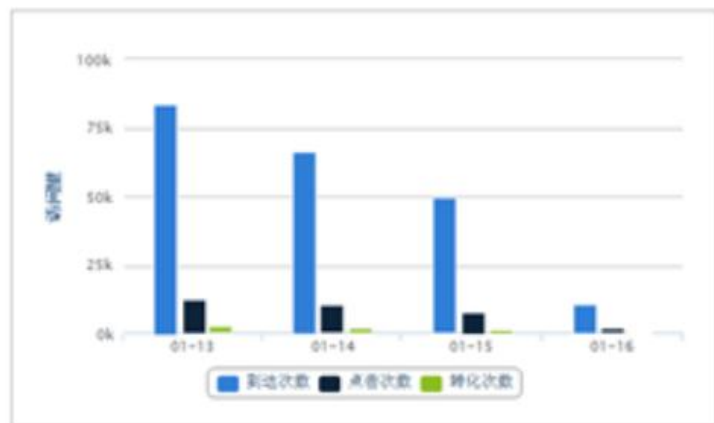


消息精准推送

我有残局，你敢来挑战吗？



生日独享大礼包



KP (开放生态、平台合作)

微信开放平台

腾讯云

腾讯开放平台

应用宝

APICloud

易起秀

效率

Customer Segments

移动开发者

媒体

移动运营者

大数据从业者

Nielsen

成本

增值

KR

用户、设备画像

协同营销

大数据解决方案

流量

资本



易用

Key Activities

Crash 监控

可视化埋点

用户分群

渠道跟踪

精准推送

反作弊

Value Proposition

SaaS

AI

BI



参考资料:

1. 《Systems Performance Enterprise and the Cloud.2013》

Thanks