

# 阿里智能运维平台 如何助力研发应对双11挑战

运维中台 如柏

2017.12

1

阿里运维历程

2

基础运维平台

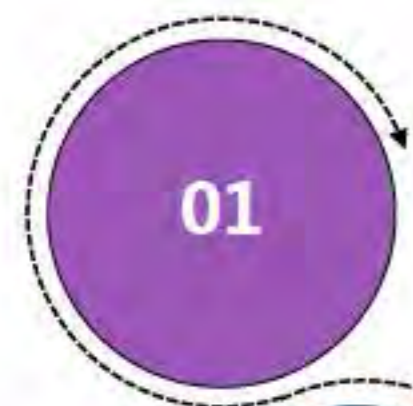
3

应用运维平台

4

AIOPS

# 阿里运维历程



## 命令行工具

2007 ~ 2010

Excel , SSH , PGM



## 系统化工具

2010 ~ 2013

Armory , Nagios , Dragon...



## 自动化平台

2013 ~ 2016

StarAgent , PSP , Alimonitor...



## 智能化平台

2016 ~

StarAgent , 蜻蜓 , Normandy , Sunfire...

## 规模化

一键建站

搬迁

腾挪

单元调整

## 稳定性

多活

故障修复

故障定位

故障注入

全链路压测

## 监控

基础监控

业务监控

链路监控

报警

视图

## 变更

变更信息

应用变更

基础软件变更

网络变更

IDC变更

## 资源

Quota管理

资源规划

资源采购

资源调度

bootstrap



# 基础运维平台



**ELECTRIC GRID**



**BRIDGES**



**RAILROADS**



**AIRPORTS**



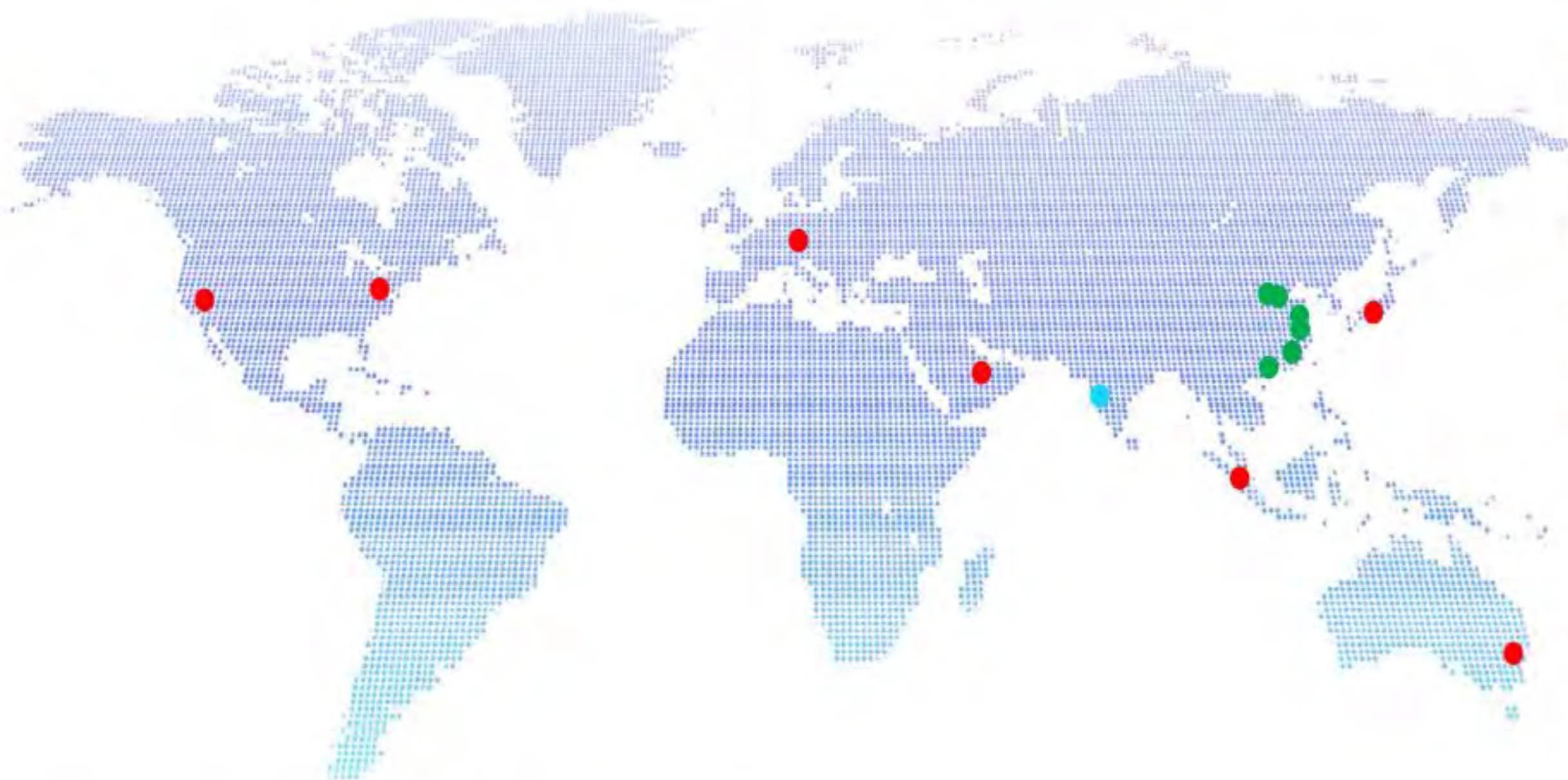
**PIPELINES**



**WATERWAYS**



## 阿里堡垒机



账号 / 权限 / 密码 / 密钥 / Root  
访问控制  
操作日志  
高危拦截  
审计 / 行为分析  
合规 / 认证  
操作录屏 / 回放  
SSH非法入侵审计

- 1.高承载，可达5000人同时在线；对管理对象没有限制
- 2.经历多种认证合规最佳实践考验 ( sox404 , iso27001 , DSS-PCI )

## 运维通道 - StarAgent

### 稳定性 / 性能

去数据库依赖

单元化

### 插件平台

低耦合、高可扩展性

### 安全

高危命令审计 / 拦截

命令映射

全链路加密签名

### 自动化运维 95%

IP段自动关联

自动扩缩容 / 自动负载均衡

Agent / 插件版本自动更新

Agent自检 / 自愈

**50万**  
每分钟API调用

**99.995%**  
系统稳定性



### 其他功能

基础信息采集

Web 终端

分布式定时任务



服务器生命周期

应用生命周期



## 命令通道

同步命令

异步命令

查询

## 数据通道

核心系统指标配置

核心系统数据采集

数据通道

## 插件平台

静态：命令、脚本

动态：常驻进程

## Portal

主机管理

分布式定时任务

主机账号

插件市场

Web 终端

文件分发

## API

cmd

File

Plugin

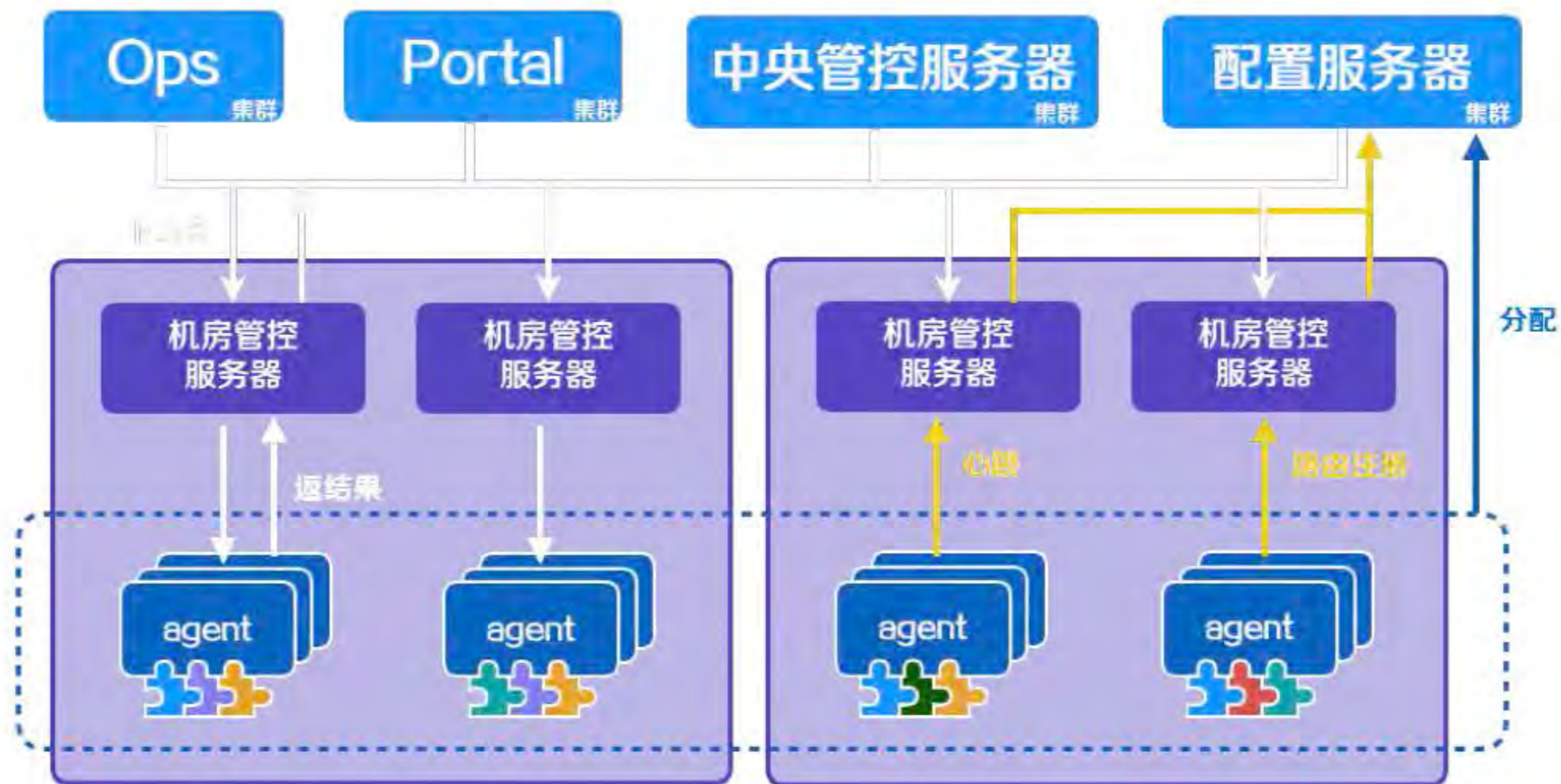
Store

Cron

Action

## SDK

## 移动端



## 完善错误码体系

系统错误

环境错误

用户错误

第三方依赖错误

未知错误

## 环境监测

监控ping报警

ssh报警

磁盘报警

syslog中机器宕机记录

## 去依赖

## 定期演练

## 命令通道

快车、慢车分道

同步任务异步化 - DeferredResult

心跳、注册 批量化写入，只更新需要更新的

### 第一重保护

核心问题：API密钥泄露、API权限过大

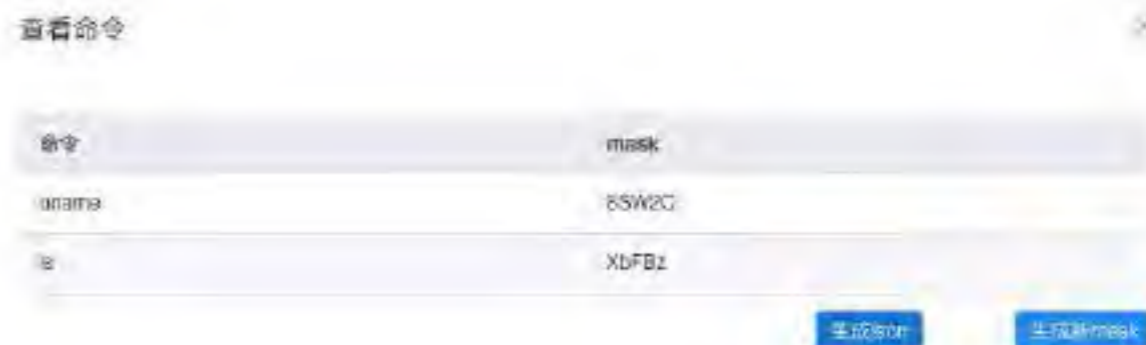
解法：API KEY生命周期管理，提供更新、重置。另外要有禁用策略

解法：引入安全负责人、Bu稳定性负责人审批。其中降权的变更直接免批。

### 第二重保护

核心问题：人工命令过多且命令可被伪造

解法：命令沙盒模块化处理，命令白名单提供主命令映射方式，以mask方式执行，参数保持不变。



### 第三重保护

核心问题：高危命令执行风险过大

解法：针对高危命令进行正则匹配记录拦截，固化每个API帐号可执行高危命令的机器数量、频率，定义高危级别并实时监控预警。

### 第四重保护

核心问题：命令被拦截和篡改

解法：TCP Socket 全链路加密、校验



**协议：启动、停止、配置、强制结束、同步执行**

**服务：守护、额度、部署、升级**

**价值：**

安全：统一管理、审计

稳定性：保护服务器正常运行

成本/效率：提升agent运维效率、  
避免重复建设agent、  
避免重复数据采集

```
{
  "name": "plugin_name",
  "version": "1.0.1-build15",
  "schedule": "daemon",
  "daemon": {
    "exec": "xxx/bin/xxx.py",
    "user": "admin",
    "group": "admin",
    "env": {
      "PATH": "bin:$PATH",
      "CWD": "${plugin.name}"
    },
    "log_dir": "${CWD}/logs"
  },
  "limit": {
    "cpu": "5%",
    "disk": "200m",
    "mem": "50m"
  }
}
```



## 快速复制与输出的能力

### 阿里云专有云

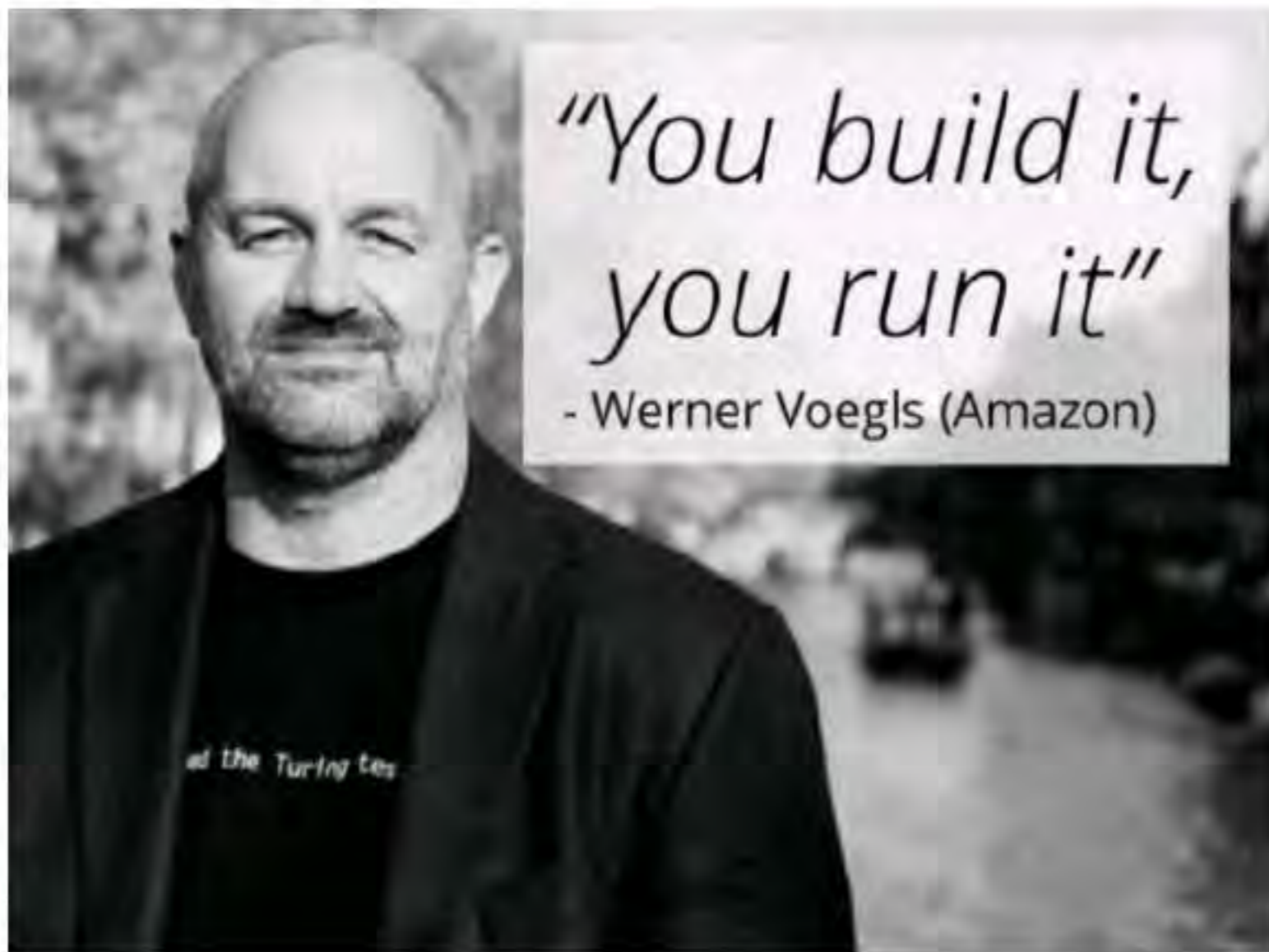
- V1/V2 40个左右的输出项目 (StarAgent 1.0)
- V3 天基 开发中...

### 蚂蚁金融云

- 公有云 30+ 客户；
- 专有云 / 专有域：网商银行、Alipay+、新加坡等几个 (SSH / StarAgent 1.0)

### 中间件EDAS

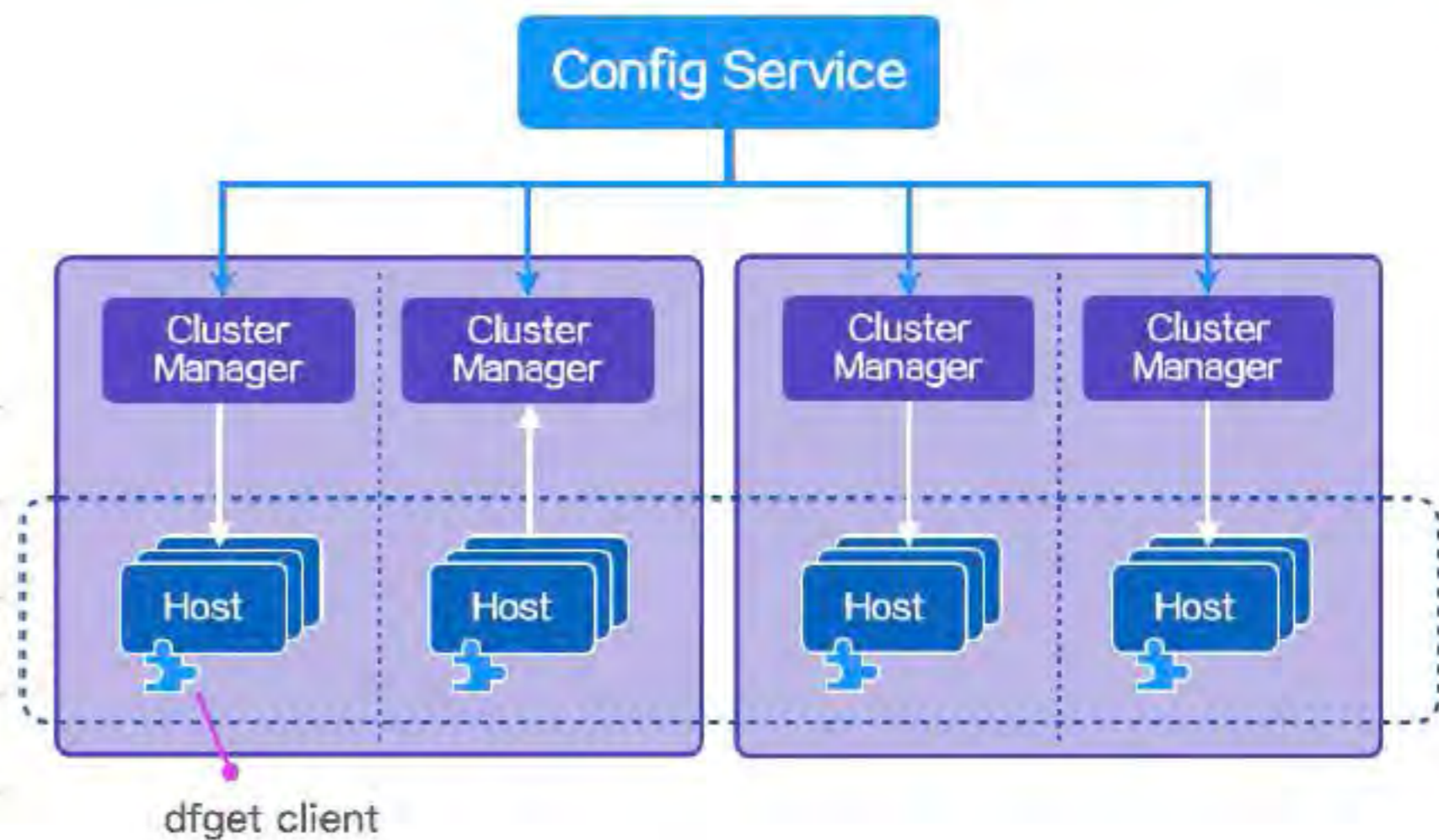
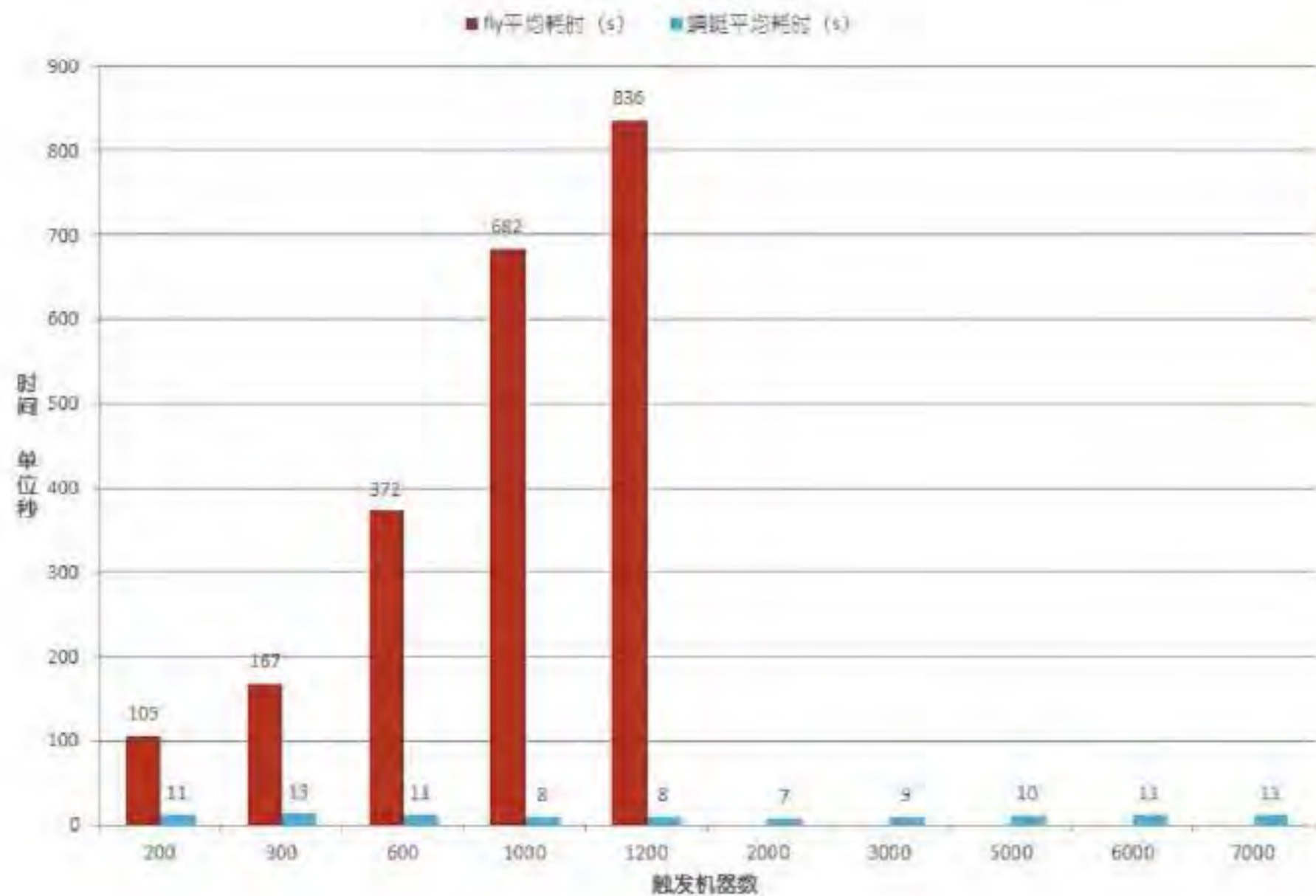
- 几百个输出项目 (StarAgent 1.0)



Werner Vogels of Amazon.com in [a blog comment](#) back in 2006:

"The best way to completely automate operations is to have to developers be responsible for running the software they develop. It is painful at times, but also means considerable creativity gets applied to a very important aspect of the software stack. It also brings developers into direct contact with customers and a very effective feedback loop starts. There is no separate operations department at Amazon: you build it; you run it."

### fly与蜻蜓性能对比



保护数据源 Protect Source

加速分发速度 Boost Delivery Speed

节省跨IDC带宽 Save cross IDC bandwidth

节省跨国带宽 Save cross board bandwidth

## 场景 Scenarios :

安装包

packages for deploy system

配置文件

config files

数据文件

AI, Search Index, business hot data

静态文件

image files, js files, CDN resources etc.

镜像

container images, VM images, Physical images etc.

其他文件

any other files served by http

## 功能特性 Features:

断点续传

Resume from break point

智能网络/磁盘IO控制

Intelligent Network/Disk I/O control

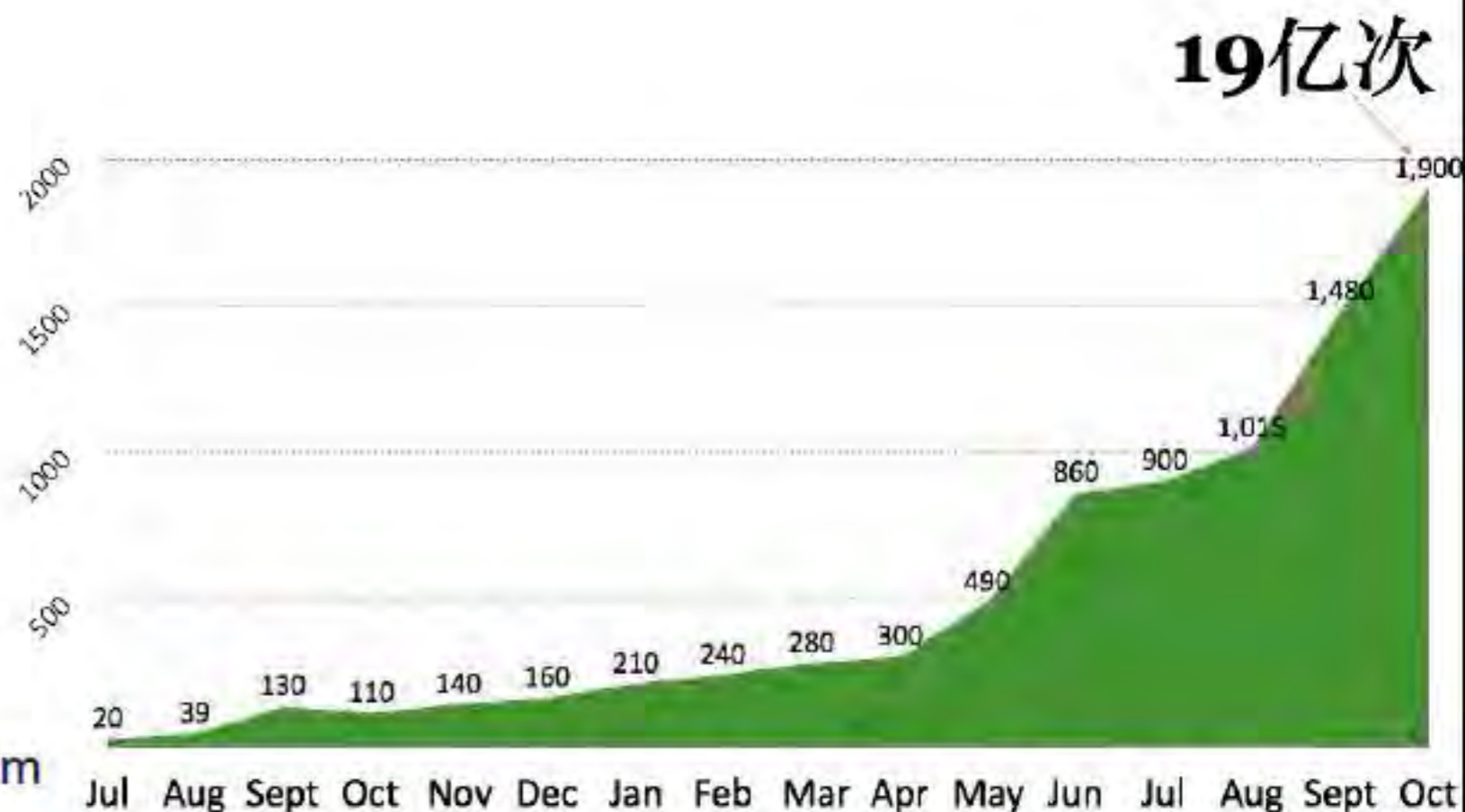
智能动态压缩

Intelligent Dynamic compression **10%** best **86%** for images

镜像预热

Container Image Warmup (Pouch, Docker, Hyper) **67%** average saving **120X** max

2017.thegiach.com



全面：支持pouch、docker、hyper等各种容器镜像

无缝：0侵入，无需daemon改造

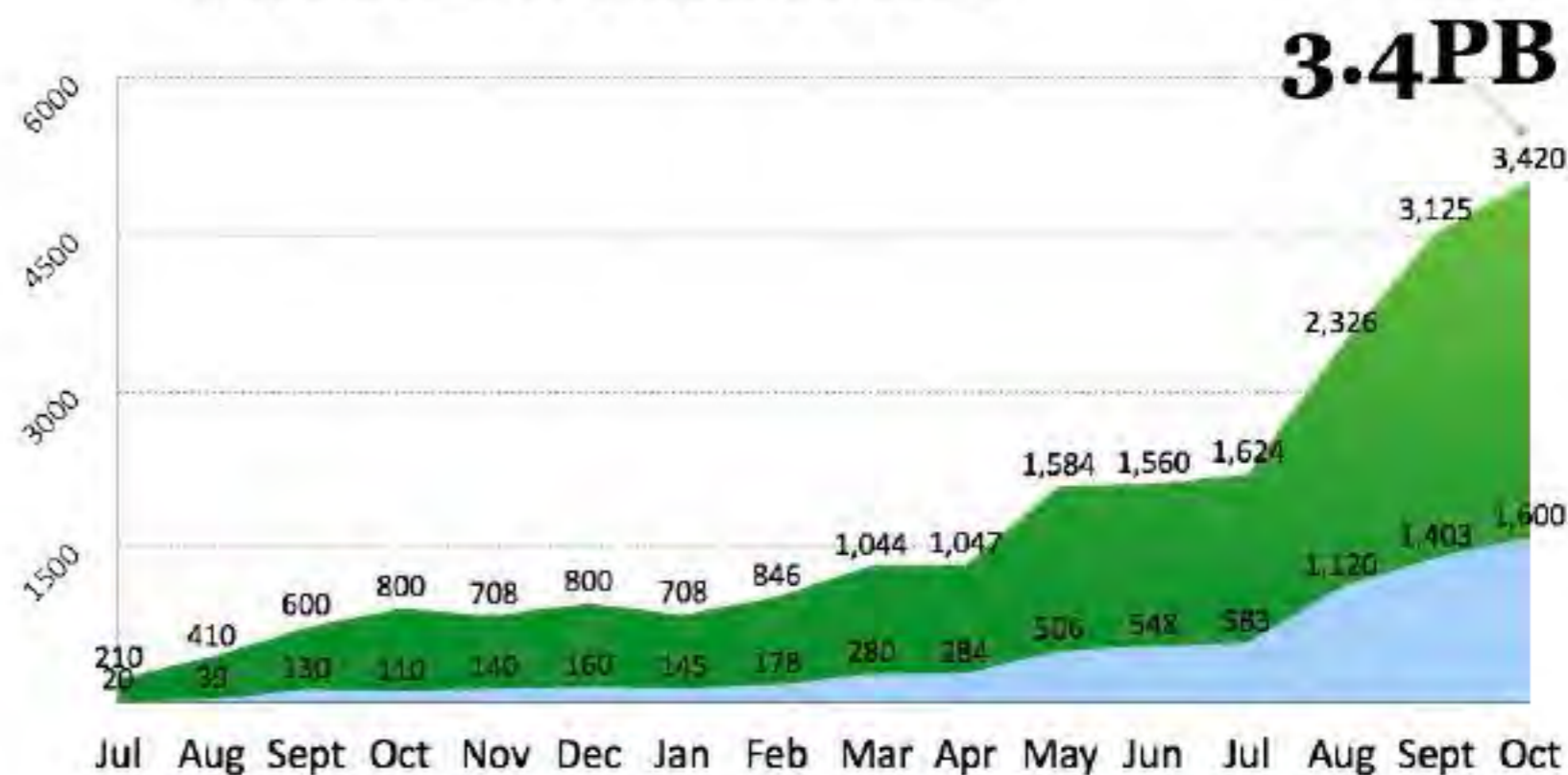
极速：支持镜像预热，push时完成并行分层预热

稳定：支持30GB甚至更大的镜像分发

高效：支持内存文件系统，极致的资源利用和极致的效率提升

极简：完美替代mirror，并在窄带、跨洋、远距离传输等极端条件下具备世界级竞争力

## 总分发量 v.s 镜像分发量



30%

动态压缩提升整体速度

## 开源：



协议 Apache 2.0

<https://github.com/alibaba/dragonfly>

功能：P2P文件分发，容器镜像分发、局部限速、磁盘容量预检

## 企业版：



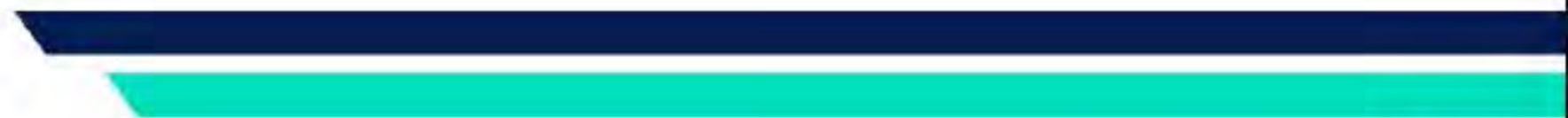
功能：断点续传、全局限速、镜像预热、支持内存文件系统、智能网络流控、智能动态压缩、智能调度策略

渠道：云效

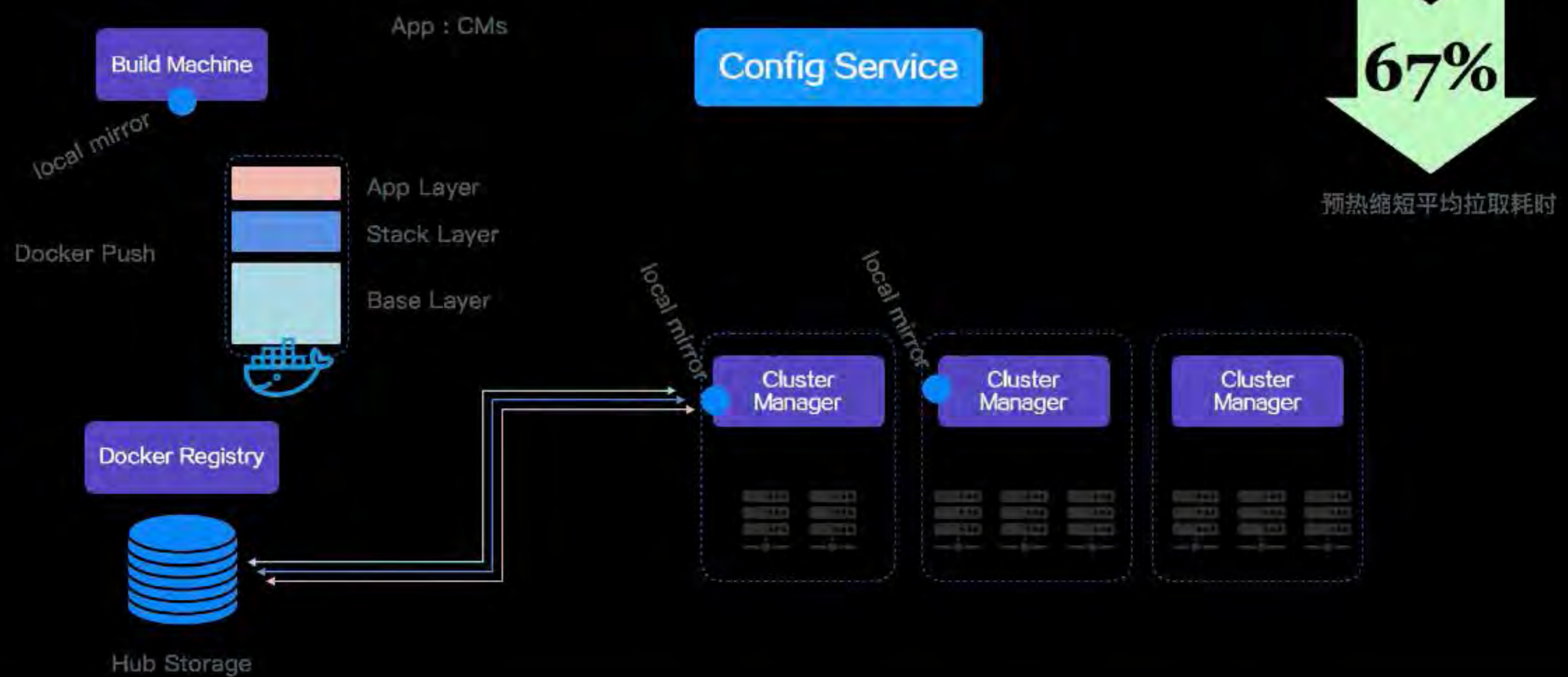
阿里云容器服务（公有云、专有云）

<https://www.aliyun.com/product/container-service>

<https://www.aliyun.com/solution/dedicatedcloud>



## 镜像预热 Container Image Warmup (Docker compatible)



大促期间（11.10~11.11）共执行文件分发任务7600W次，其中镜像分发17.6W次。总流量118TB。

10点23分触发万台服务器下载2G的预热文件，总计耗时111秒。100%下载成功



# 应用运维平台

## 业务关系



一体化DevOps平台

## 环境类型

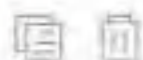
生产环境

预发环境

测试环境

复制环境

正式

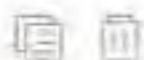


机器数量 10394  
当前版本 7579146  
基线一致性 99%

应用重启 资源变更

发布 查看拓扑

正式

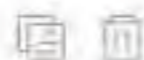


机器数量 10394  
当前版本 7579146  
基线一致性 99%

应用重启 资源变更

发布 查看拓扑

beta

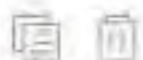


机器数量 10394  
当前版本 7579146  
基线一致性 99%

应用重启 资源变更

发布 查看拓扑

灰度



机器数量 10394  
当前版本 7579146  
基线一致性 99%

应用重启 资源变更

发布 查看拓扑



新建环境  
(建设中)

## 我的告警

最近12时告警数 21 条

## 变更记录

时间	级别	摘要	状态	智能推荐操作
11-04 00:02	Warning	buy2 共有1235台机器触发, cpu利用率>70%, <a href="#">查看详情</a>	进行中	机器重启
11-04 00:02	Critical	buy2 共有1235台机器触发, cpu利用率>70%, <a href="#">查看详情</a>	待处理	清理磁盘
11-04 00:02	Critical	buy2 共有1235台机器触发, cpu利用率>70%, <a href="#">查看详情</a>	成功	Java dump
11-04 00:02	Warning	buy2 共有1235台机器触发, cpu利用率>70%, <a href="#">查看详情</a>	进行中	暂无推荐
11-04 00:02	Warning	buy2 共有1235台机器触发, cpu利用率>70%, <a href="#">查看详情</a>	失败	机器重启

查看更多

-  buy2  
预发环境发布预发环境发布预发...  
2017-11-09 13:05:22 香子
-  buy2  
预发环境发布预发环境发布预发...  
2017-11-09 13:05:22 香子
-  buy2  
预发环境发布预发环境发布预发...  
2017-11-09 13:05:22 香子

## 基础设施即代码

混合云一站式资源申请

资源基线

自动化

弹性伸缩

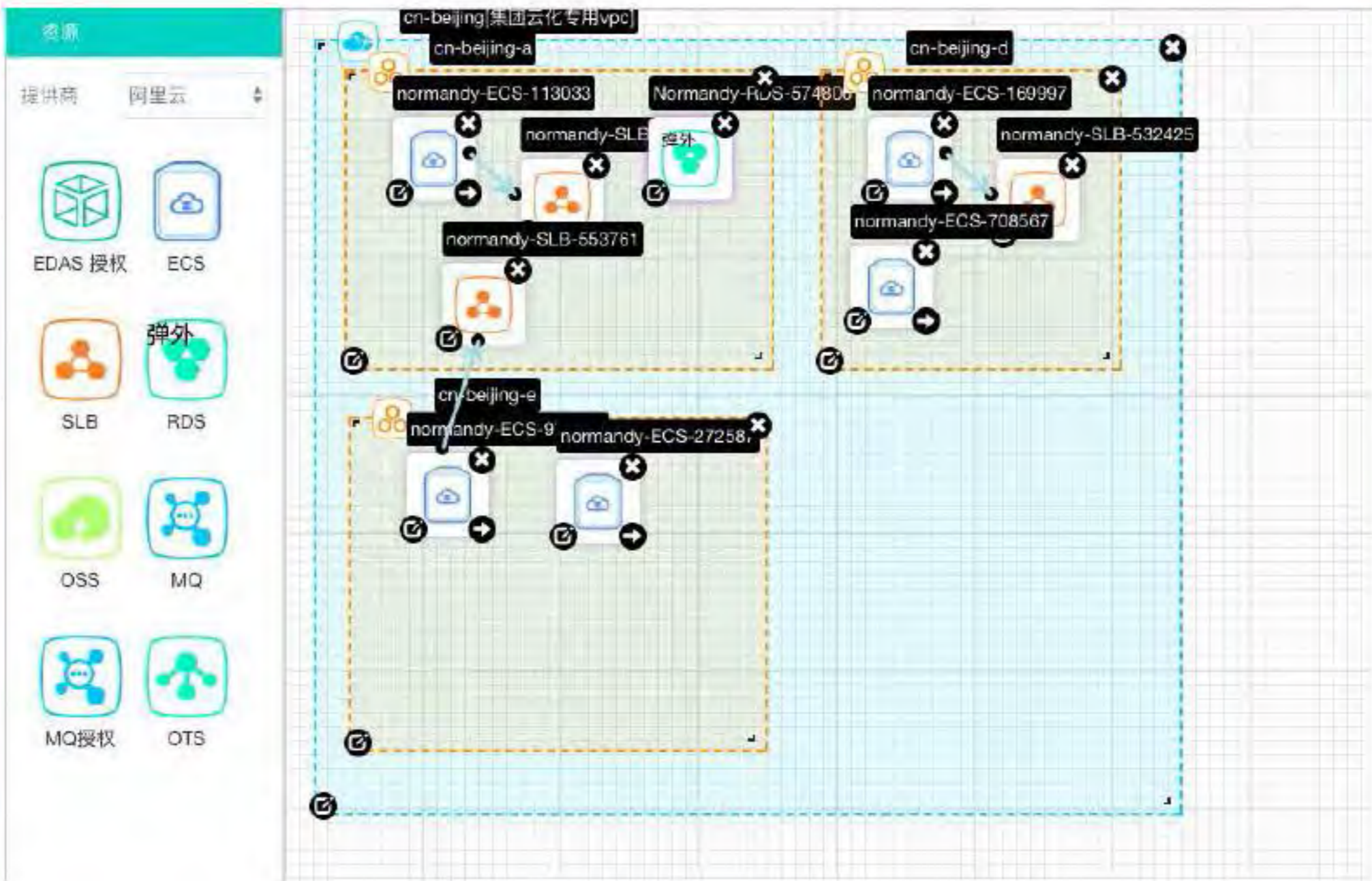


Time spent on  
resource provisioning

```
{
  "provider": {
    "name": "aliyun",
    "version": "1.0",
    "auth": {
      "accessId": "L*****d",
      "accessSecretKey": "d0*****NWer",
      "uid": "197*****858"
    }
  },
  "resource": [
    {
      "name": "hello_ecs",
      "ay_ecs": {
        "InstanceType": "ecs.g5.large",
        "ImageId": "centos5u20170401.20170401.20170401.50130.vhd",
        "count": 2
      }
    },
    {
      "name": "hello_slb",
      "sy_slb": {
        "memberList": {
          "ipPort": 8080,
          "rsPort": 8080,
          "hcType": "HTTPS",
          "hcURL": "s:hello.com",
          "resName": "${hello_ecs.*}"
        }
      }
    }
  ]
}
```

第一次使用, 不知道如何申请资源? 没关系, 请先阅读 [【使用指南】](#)。

点击 **【执行资源变更】** 后, 您将可以预览资源的生产 / 销毁情况。如果出现了预期之外的资源销毁提示, 继续操作将可能导致线上服务受到影响! 为避免出现故障, 此时请暂停操作, 请进入支持群进行咨询。



### 资源

提供商 阿里云

- EDAS 授权
- ECS
- SLB
- 弹外
- RDS
- OSS
- MQ
- MQ授权
- OTS

### RDS 参数文件

参数不是必须的, 您可以忽略; 若感兴趣, 可点击下面的问号了解

选择套餐

点击选择套餐

保存为新套餐

删除套餐

名称 Normandy-RDS-574806

只读实例  是  否

白名单 127.0.0.1

**【白名单】** 允许访问该实例下所有数据库的IP名单, 以逗号隔开, 不可重复, 最多1000个; 支持格式: 10.23.12.24 (IP), 或者 10.23.12.24/24 (CIDR模式, 无类别间路由, /24表示了地址中前缀的长度, 范围[1,32])

Region 北京(cn-beijing)

VPC cn-beijing(集团云化专用)

Zone cn-beijing-a

Vswitch beijing-a-1

应用分组 liverecordhost.cm12

100.67.31.113

## 创建发布单

\*选择应用

环境级别  预发  测试  正式

### 发布信息

\*环境      [更多](#)

\*发布类型

\*分批方式  机房均分  分组均分

\*暂停策略  第一批暂停  每批暂停

\*分批数量

\*计划发布时间

### 版本信息

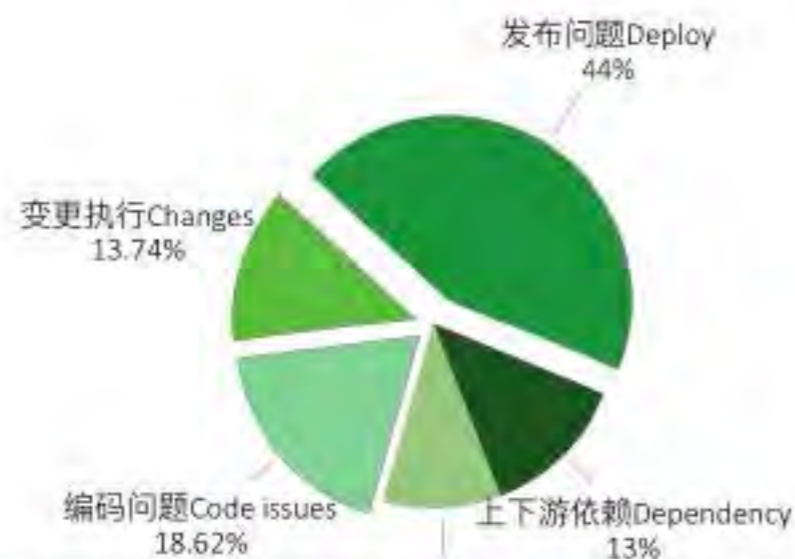
V-1.6.8 薇子 2017-1-23 18:00:00  
 V-1.6.7 薇子 2017-1-23 18:00:00  
 V-1.6.6 薇子 2017-1-23 18:00:00

V-1.6.5 薇子 2017-1-23 18:00:00  
 V-1.6.4 薇子 2017-1-23 18:00:00  
[查看更多版本](#)

### 机器信息



## 无人值守发布变更



发布/变更故障 deploy failure



**总览**

6个异常

应用: ump  
是否异常: 异常  
异常评分: 86.78417  
分析区间: 17-07-04 17:35:58到 17-07-04 17:39:58  
结果评价: 正确 不正确

**异常检测报告**

- ⚠ HSF 异常 查看详情
- ⊗ 指标middleware.hsf.provider.rt偏离基线 75.54385494183828 查看详情
- ⚠ TAIR 异常 查看详情
- ⊗ 指标middleware.tair.write.qps偏离基线 86.79417069929711 查看详情
- ⚠ TDDL 异常 查看详情
- ⊗ 指标middleware.tddl.read.qps偏离基线 76.55289727034106 查看详情
- ⊗ 指标middleware.tddl.read.rt偏离基线 55.94238708736593 查看详情
- ⊗ 指标middleware.tddl.write.qps偏离基线 76.44219531384623 查看详情
- ⊗ 指标middleware.tddl.write.rt偏离基线 55.94057535957258 查看详情

**发布单信息**

应用	ump	发布单ID	5134692
发布批次	1	发布单状态	FINISH
此次开始	17-07-04 17:23:44	此次结束	17-07-04 17:31:05
发布机器列表	<a href="#">查看</a>	对账机器列表	<a href="#">查看</a>

**实时数据对比**

前后对比 | 未发布对比

**middleware.tddl.read.qps** 发布前 发布后

**middleware.tddl.write.rt** 发布前 发布后

System Monitoring	CPU, Memory, IO	User, System, Load etc.
Container	Tomcat, JVM etc.	HTTP failure/error, thread etc.
Middleware	HSF, Notify, TDDL, MetaQ etc.	QPS, RT etc.
Application	Log	Exception count, rate, trend etc.
Business	Login, Transaction etc.	count, rate, trend etc.

## 基本策略

新增的Exception

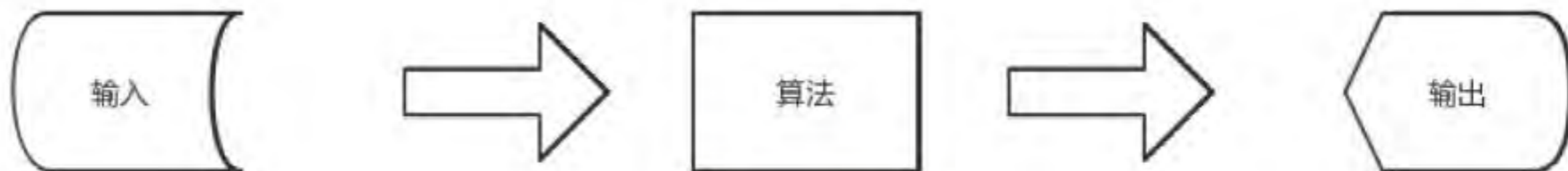
Fatal Exception (NoClassDefine, NPE...)

旧的Exception, 但是Rate飙高

业务指标  
系统指标

} 异常检测

趋势 同比 环比 静态阈值 分段静态阈值 动态阈值





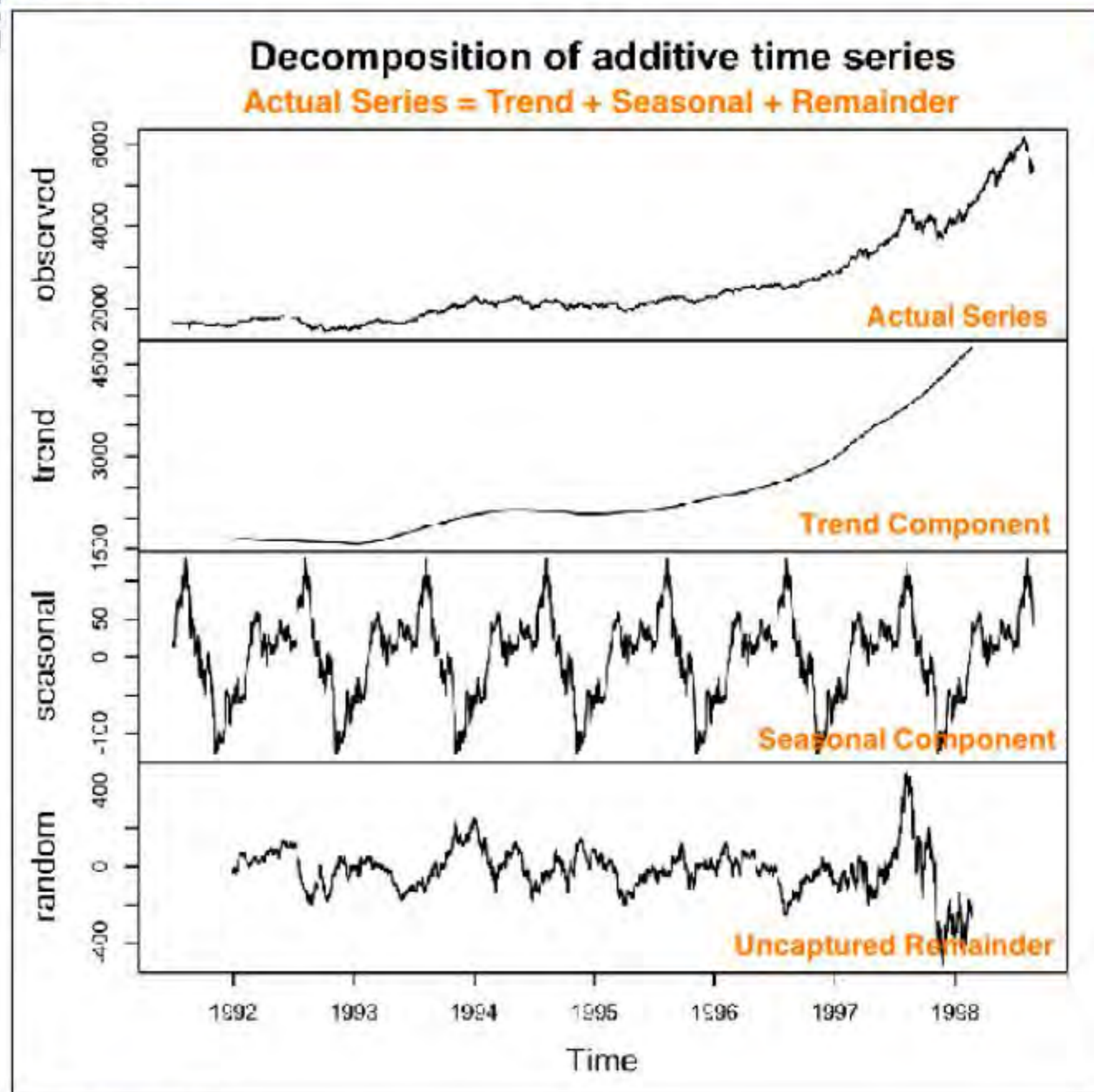
## 时间序列预测

分段历史平均

Arima

Holt-Winters

STL



## 预测：三次指数平滑Holt-Winters

适用于趋势性(Trend)和周期性(Seasonality)的时序指标

模型参数： $\alpha / \beta / \gamma$  截距/斜率/周期平滑系数

参数确定：

◆ 人工配置

◆ 自动训练：排除异常点 -> 最大化拟合度

## 异常判定：

明确上下界：预测值 $\pm \delta$

固定阈值

历史周期点的指数平滑

滑动窗口的偏差标准差

• Prediction:

$$\hat{y}_{t+1} = a_t + b_t + c_{t+1-m}$$

• Baseline ("intercept"):

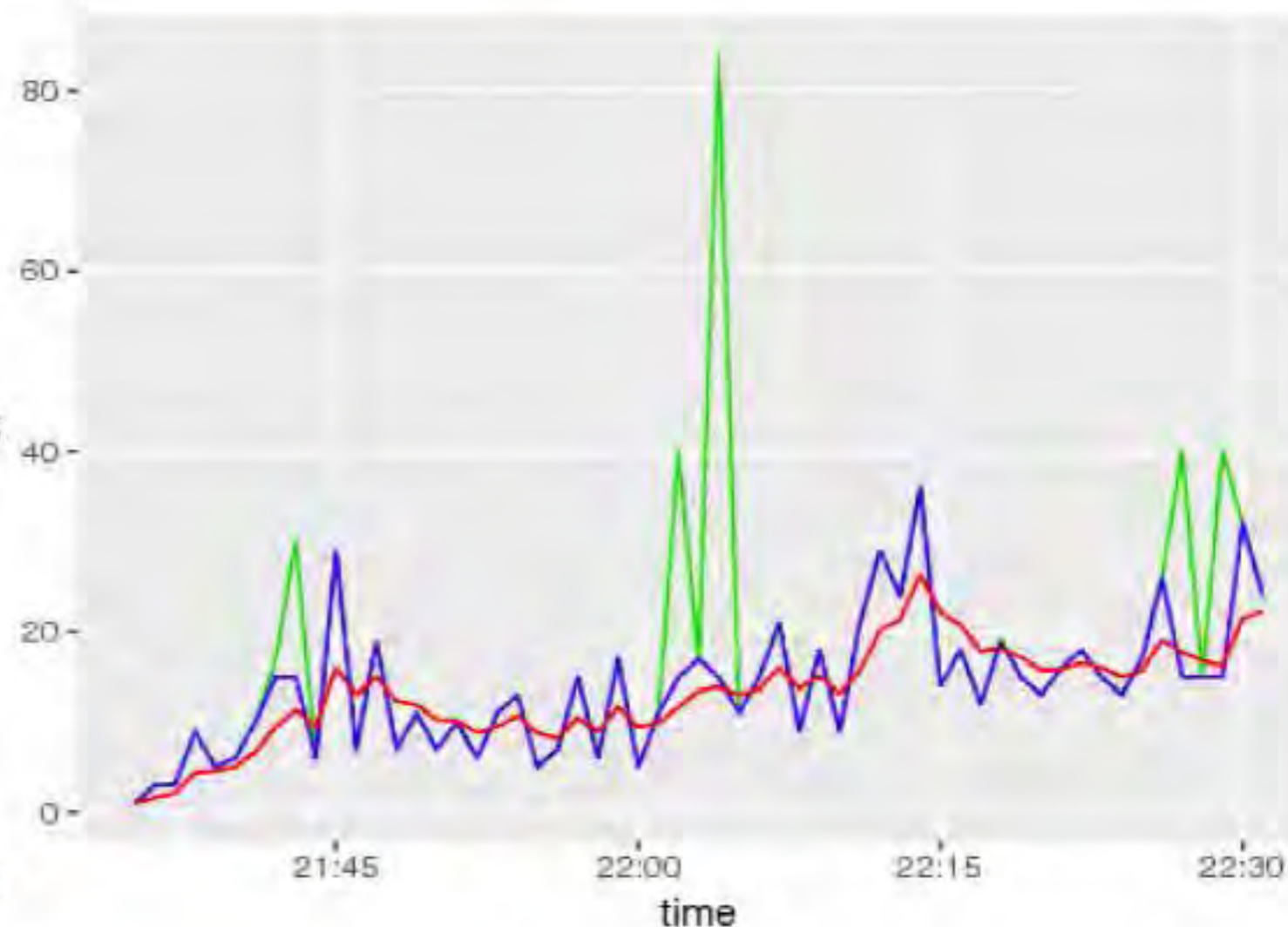
$$a_{t+1} = \alpha(y_{t+1} - c_{t+1-m}) + (1-\alpha)(a_t + b_t)$$

• Linear Trend ("slope"):

$$b_{t+1} = \beta(a_{t+1} - a_t) + (1-\beta)b_t$$

• Seasonal Trend:

$$c_{t+1} = \gamma(y_{t+1} - a_{t+1}) + (1-\gamma)c_{t+1-m}$$



2017.thegiacy

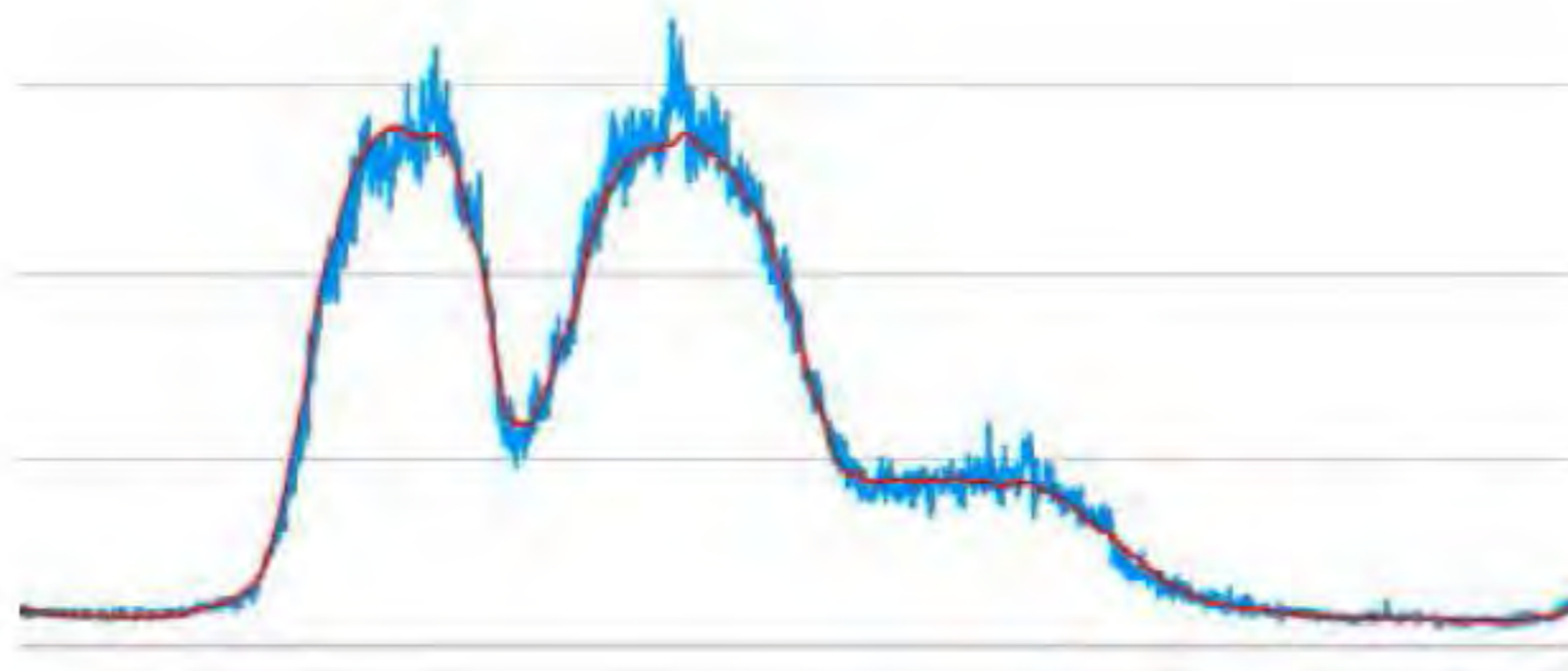
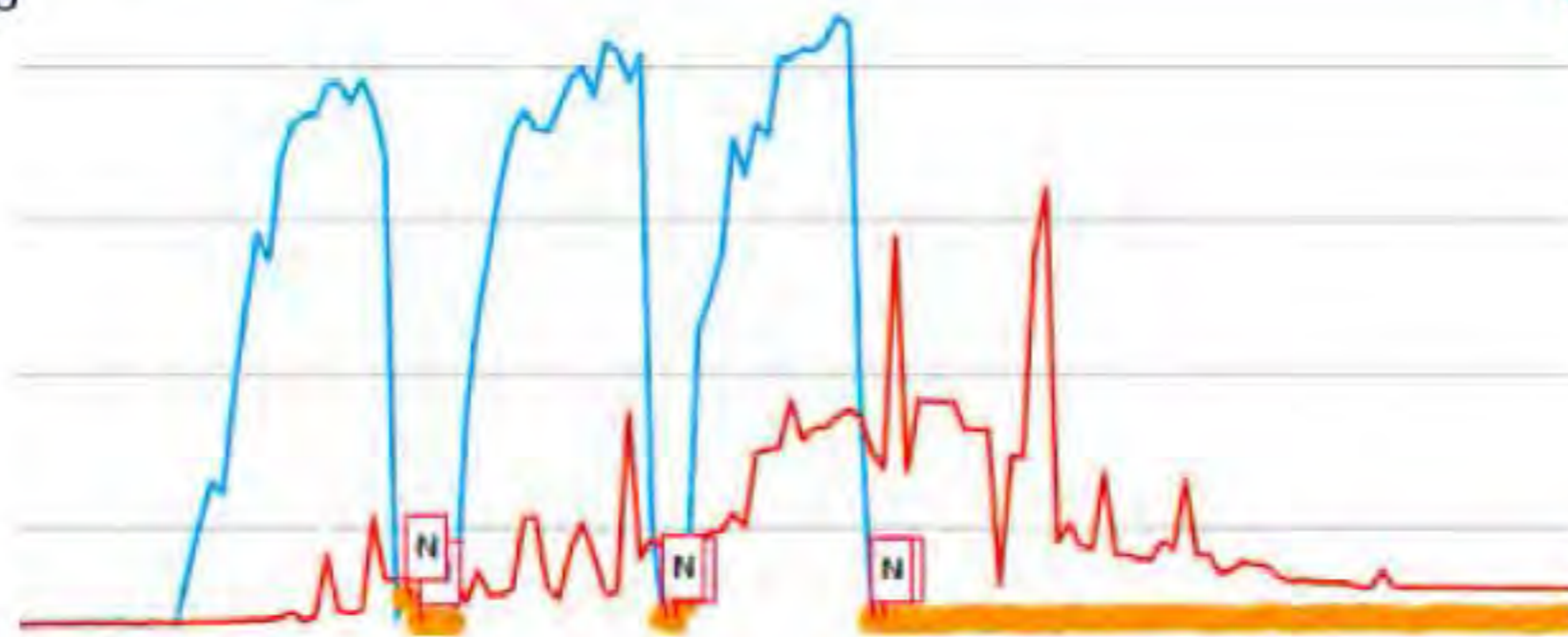
## 关键步骤

数据补全

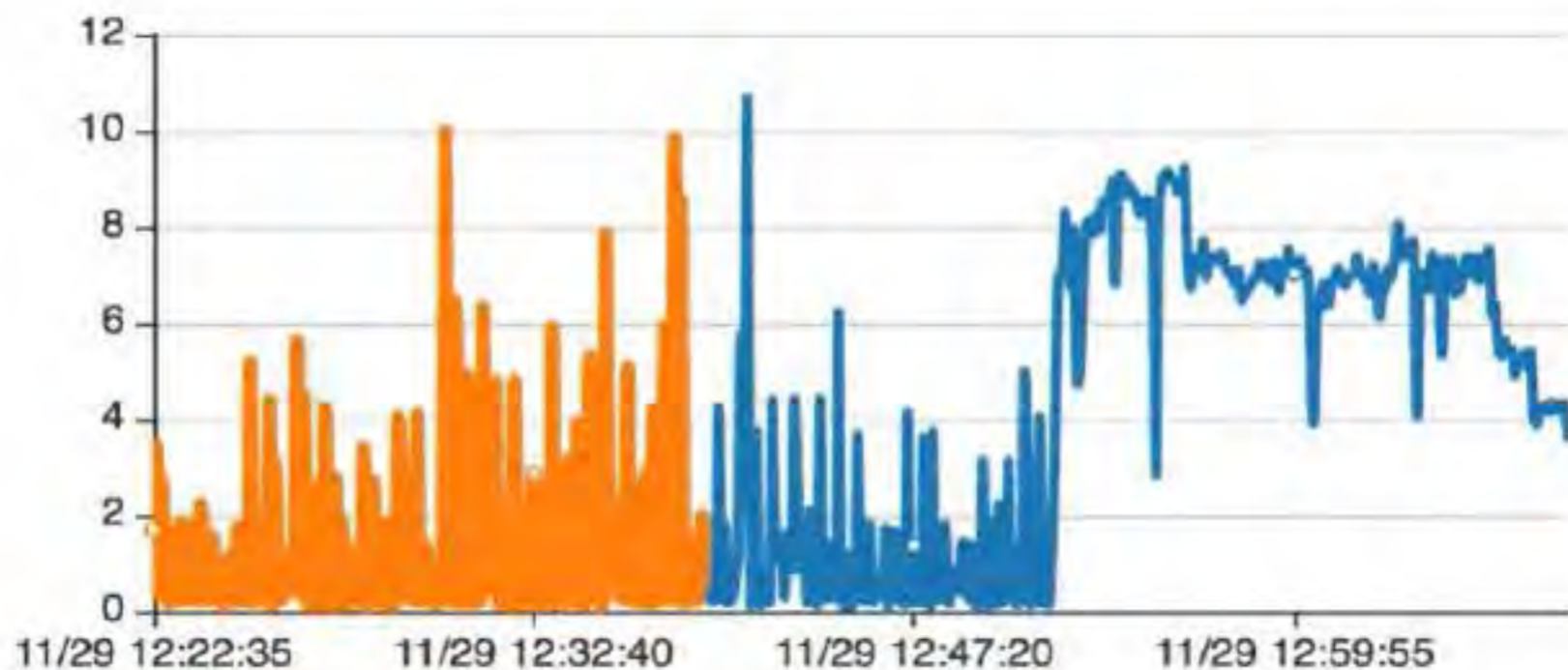
日期类型划分

预测“趋势”

局部趋势反馈



# 系统指标



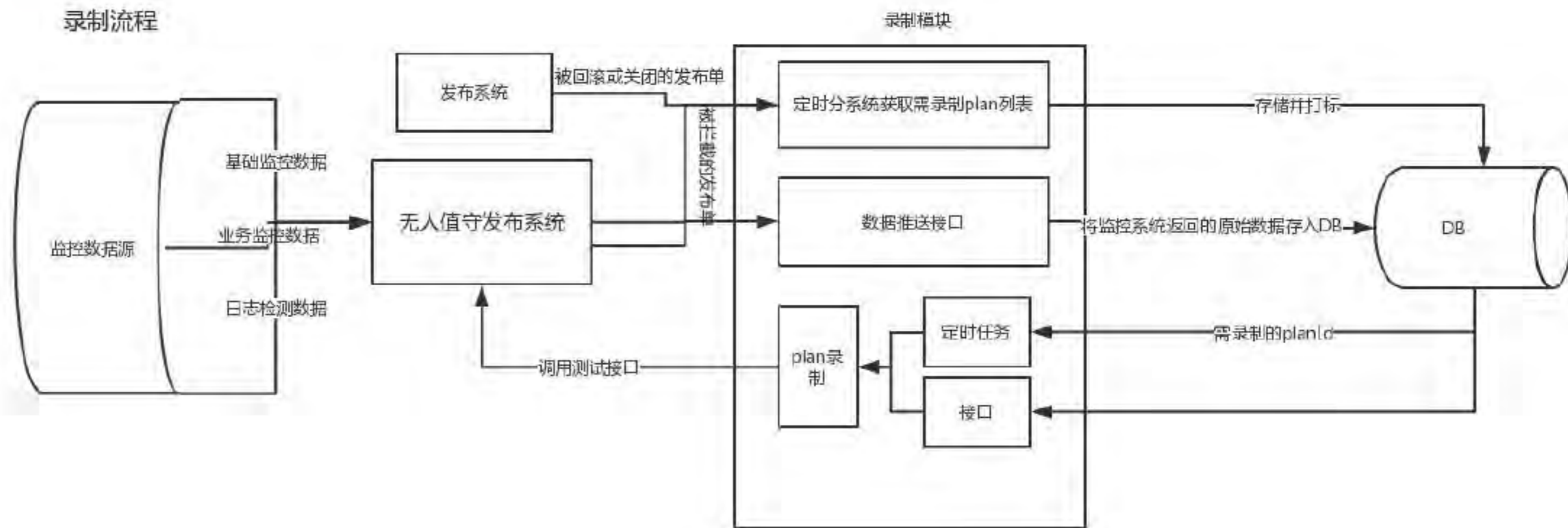
middleware.metaq-clnt.receive.rt

before after

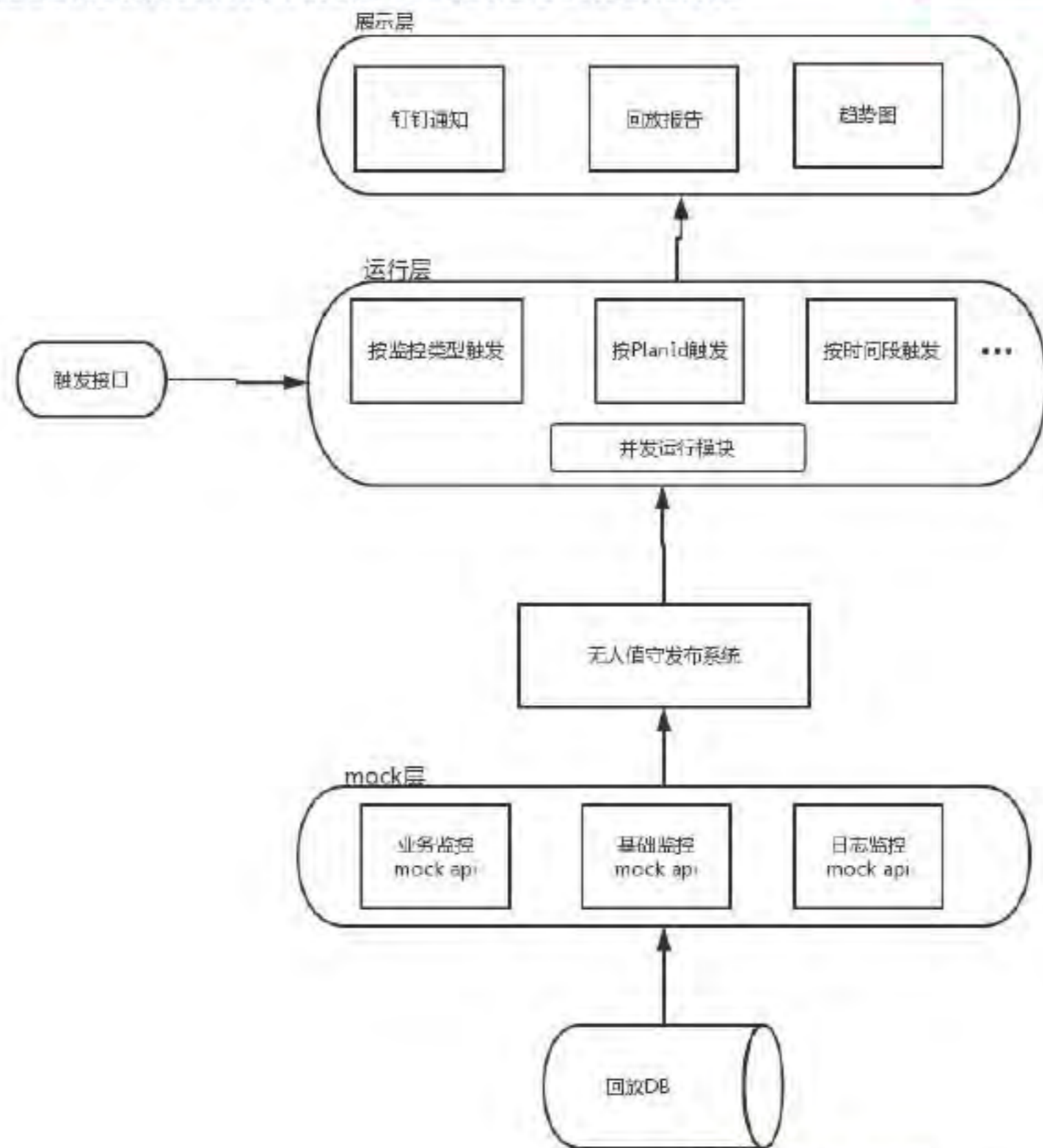


Method	Precision	Recall	F1	False Positives	False Negatives	Total Positive
AUC	0.33	0.90	0.49	18	1	9
DTW	0.57	0.82	0.67	7	2	11
AUC+DTW	0.80	0.80	0.80	2	2	10
Risk-free	0.33	1.00	0.40	24	0	12

## 反馈 - 录制



## 反馈 - 回放



## 问题和挑战

### 应用画像

日志异常/error库

### 接入成本

日志路径很难自动探测到

业务指标缺失，需要配置，甚至要改代码（500个应用 1500个指标）；机器维度的聚合

### 数据质量

数据打标少、准确度低

监控数据准确性，断图

GOC故障分类、变更关联

### 数据数量

应用发布时多数业务量较少甚至没有

### 判定标准

固定阈值