

中大型网站的运维体系

--诸超

zhuchao@gmail.com

Agenda

- PPTV网站运维体系
- PPTV运维碰到过的问题
- 理想网站运维体系/我们的努力方向
 - 标准化: ops infrastructure & Dev Env
 - 技术组件化 & 服务化, 人员专业化
 - 自动化运维: provision, release, mon;
 - 数据化运维: cmdb, log analysis,
 - 监控: technical/(server/app), end user point, business
 - 安全: network/server/app
 - Review架构设计, 推动架构优化(简化)
- About DevOps, 应用运维
- 我的另外一些感想

中大型网站

- 200-5000 服务器
- 5-50个运维员工
- 组织结构：
 - 系统运维(IDC/Network/Server)
 - 应用运维(Web/CDN Operations)
 - 数据库
 - 运维/平台开发
 - 监控值班团队
 - 安全
 - 流程

互联网运维需要

- ✓ 稳定
- ✓ 低成本
- ✓ 快速响应

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

How

— 标准化，规范化：

- Kickstart/Puppet/LDAP/Zabbix/DNS

— 服务化

- 搭建运维基础架构，应用基础架构

— 监控

- 系统监控
- 应用监控
- 业务监控

— 自动化

- 装机
- 发布/Release
- 监控

— 数据化运维

- CMDB
- 日志分析
- 容量规划

— 网站安全

- 网络，系统，应用

— 推动架构优化/简化

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

标准化规范化

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

统一规范化-Dev

- 统一公司层面运维支持技术开发平台
 - 控制研发技术使用
 - 不是想用什么技术就可以用什么技术的，
 - 要经过arch review讨论
 - 简化技术平台，有利于网站稳定和技术共享
 - 适当允许引入新技术
 - 维护一个稳定版，一个推进版
- 搞深搞好技术，不要啥都搞不懂
- 不随便引入新技术

统一规范化--Ops

- 操作系统：
 - OS Kickstart全部从Cobbler安装,取消光盘安装
 - 硬件规格尽量统一，几个template(VM)
 - OS内核，package 统一
 - **Puppet 控制系统和CDN应用**
 - 服务器命名
 - 统一LB配置
- 统一了，才不会出现诡异配置问题

统一应用部署模式

- 各类应用基础架构
 - 日志格式，路径
 - 空间管理，日志删除
 - 参数配置文件
 - Java, PHP, MySQL, Redis 安装模板
- guideline/Principle
 - IDC consideration
 - App VM min 2;
 - LB 公用：不许新建LB服务
 - IP 分配：默认一律内网
 - DB read高可用，写不管
 - Cache Consideration
 - Logging Consideration
- 统一了，才能做自动化部署，分析，排障

运维基础服务化

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

运维基础服务化

■ 为什么要服务化

- 避免重复开发/部署，加快项目部署
- 更好积累经验，用精技术,不用每个业务有自己的Solution
- 避免单点，复用服务
- 简化故障排查，简化系统架构

基础运维基础组件

- Kickstart/Cobbler
- LDAP
- DNS (in/out)
- Zabbix
- Puppet
- CMDB

Puppet

- 大规模运维必须的利器！
 - 统一
 - 标准
 - 快速
- 控制
 - OS初始化
 - App环境初始化
 - LDAP , SSH , 等安全加固
 - 全网/部分推送某个配置变化
 - 应用直接部署
- 注意控制节奏！
- CDN应用

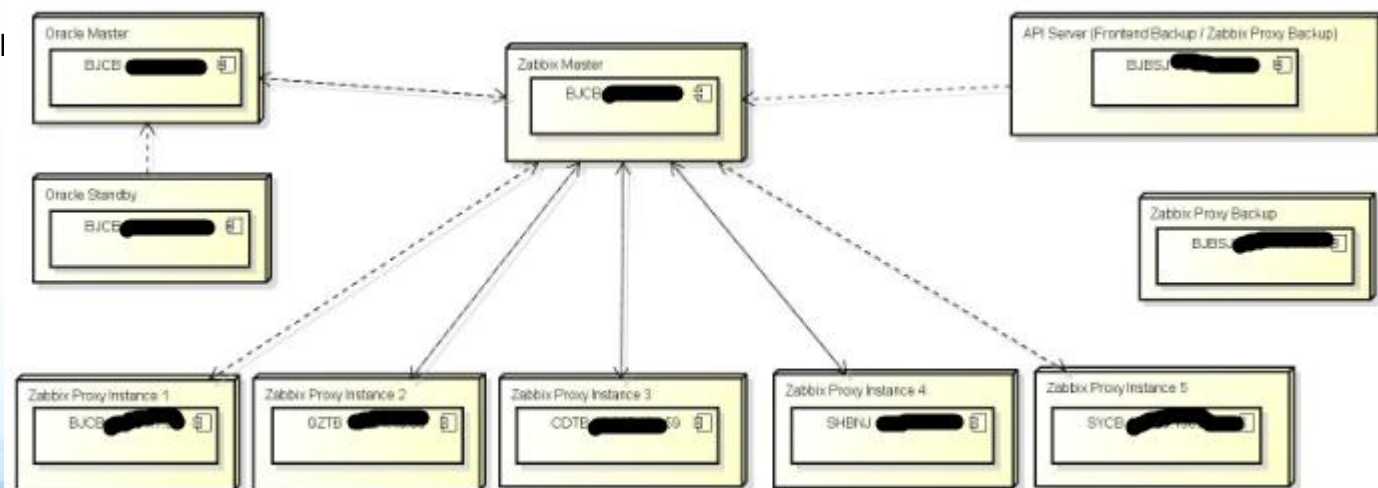
Zabbix

■ Why Zabbix

- 基于数据库，可以分析数据
- 有Agent，可以跑任务

■ Template

- Linux base template
 - Cpu/mem/net/disk/space/
- Def App template
 - Squid/Php/I
- Self-defined



CMDB

■ 资产：

- 多少机房，带宽，ip段，服务器，
- 各自型号分配，
- IDC分布，
- 过保分布，
- 基础运维服务透明化

■ 应用：

- 多少应用，各自是什么业务负责人（研发，运维），各自用了多少机器
- 有什么监控点(URLMON来读取)，
- 告警邮件直接抄送研发
- 上线自动添加监控，下线自动取消监控，维护自动暂停监控

各类服务

- App日志接收平台
- 各类Hadoop日志收集平台:Rsync
- 系统/安全日志: rsyslog
- Hive 日志分析
- 实时日志分析平台：fluentd+mongo+zabbix
- 监控平台：Zabbix+URLMon
- 自动化Release 平台
- 单点登录平台:CAS-SSO
- Web源站统一：统一管理优化
- 负载均衡服务化：LB Pair
- Cache Tier: shared cache service(MC,Redis)
- Database as a service (MySQL)

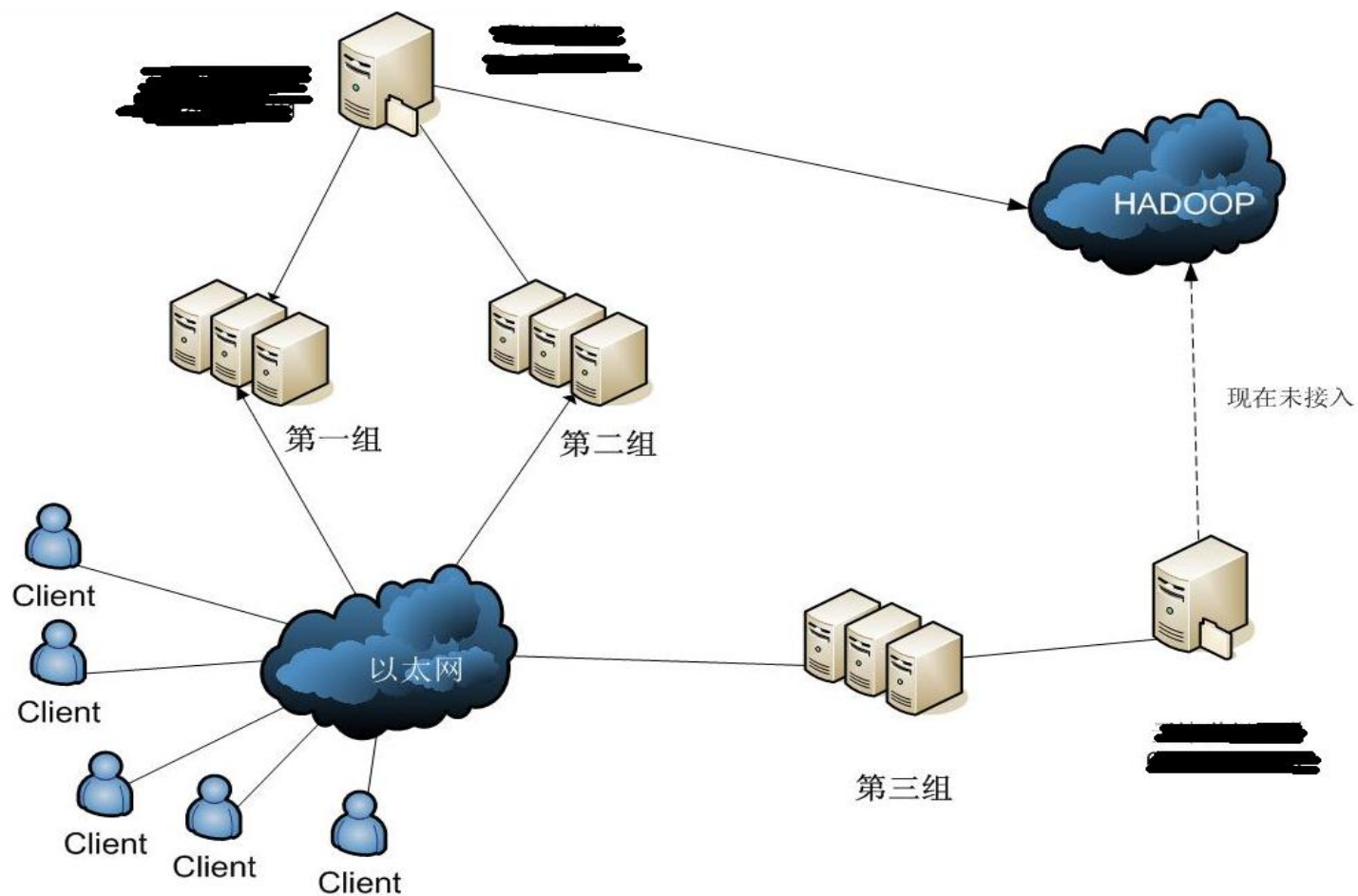
SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

App/客户端/Web 日志收集



日志传输平台

Why

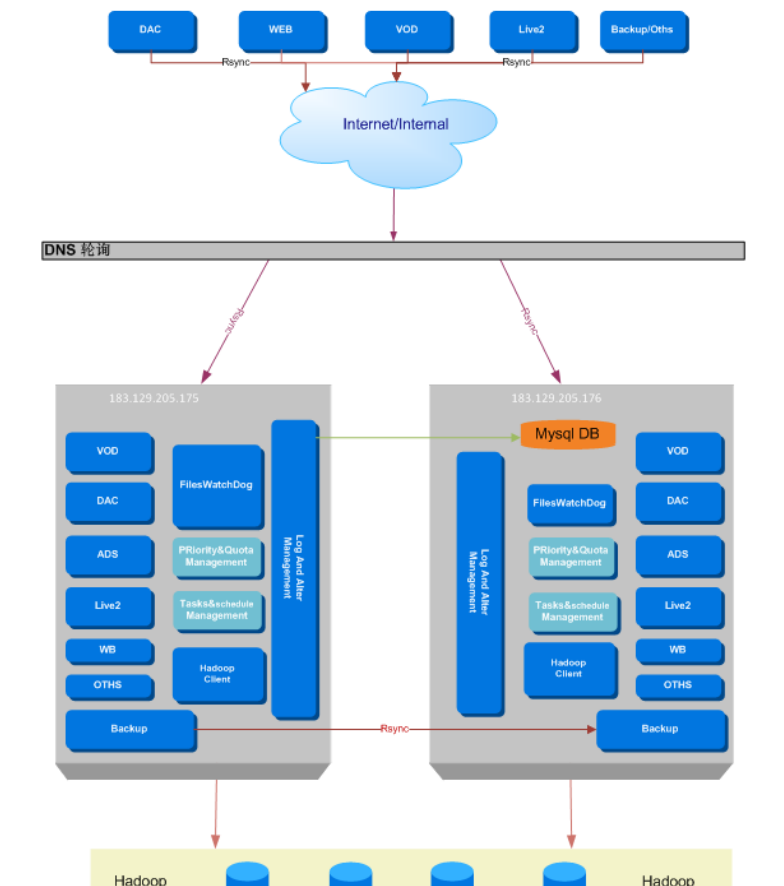
- 有10+ 不同日志进Hadoop需求
- 有5+ 不同进Hadoop的方案
- 断断续续轮流出问题

How :

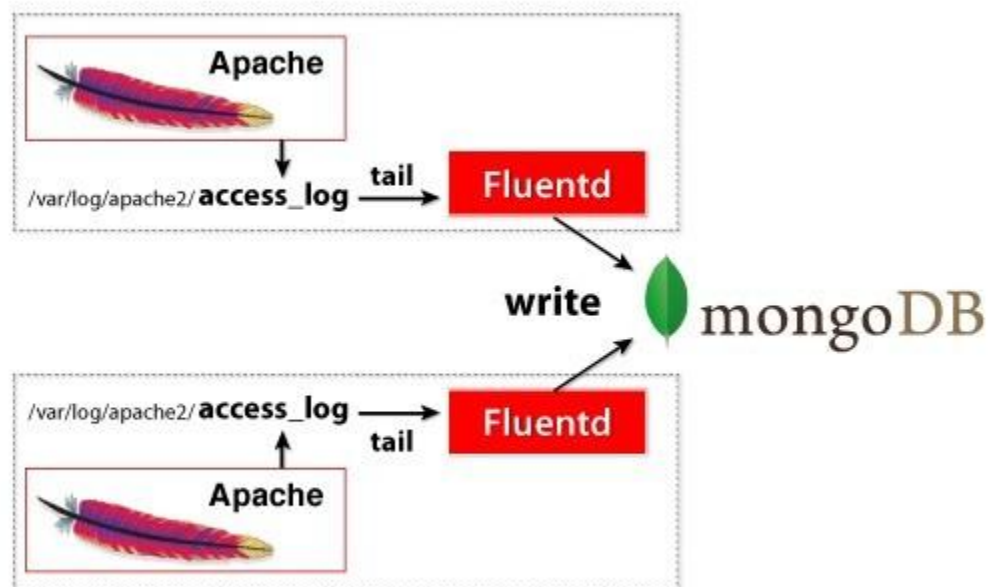
- 一个统一方案
- 吸取各个方案lesson
- 尽量简单稳定
- 非实时
- 新应用要用：配置客户端即可

2 基于Rsync Client+Server日志收集和管理方案

系统逻辑图



实时日志分析平台



离线日志分析

- Hadoop!
- Hive

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

安全日志分析(WIP)

- 系统: Rsyslog + **PHP**
 - Bash 日志
 - /var/log/message
 - 登陆日志
 - Sudo 日志
- Web安全 : fluentd+Mongo
 - Fluentd + mongo
 - LB post日志

监控

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

监控

- 定义KPI，定义Threshold
- 基础监控
 - 服务器
 - 网络：多机房，CDN，
 - App metrics：收集为主
- 应用监控/用户体验监控
 - 用户端URL监控
 - 程序接口监控
 - LB Traffic/code/Exceptions
- 业务数据监控

告警

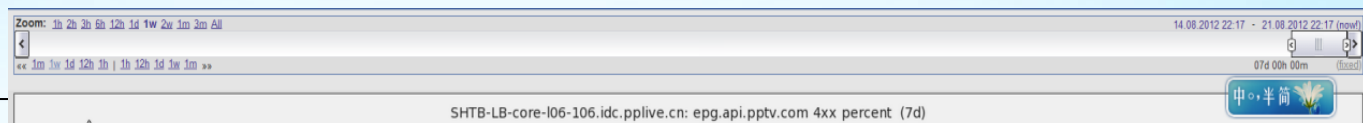
- 监控容易告警难

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算



LATEST DATA

ITEMS

Description

Agent (1 Items)

CPU (6 Items)

Disk Health (4 Items)

Disk_Performance (3 Items)

Filesystem (16 Items)

General (2 Items)

Integrity (1 Items)

Memory (2 Items)

Network (6 Items)

Nginx (1 Items)

OS (3 Items)

Performance (9 Items)

Processes (2 Items)

Security (4 Items)

Swap (2 Items)

SyslogMon (1 Items)

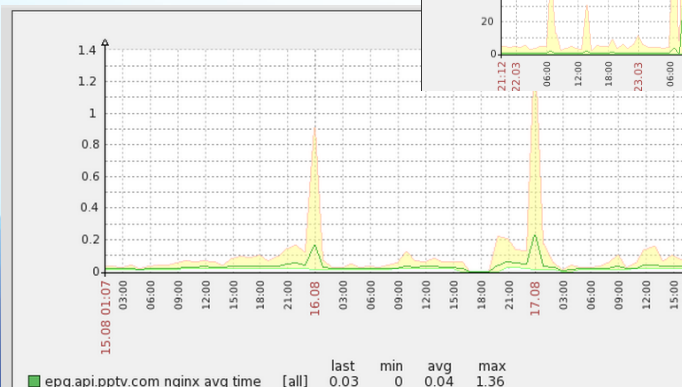
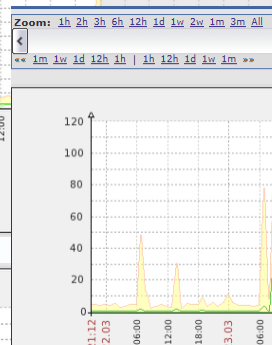
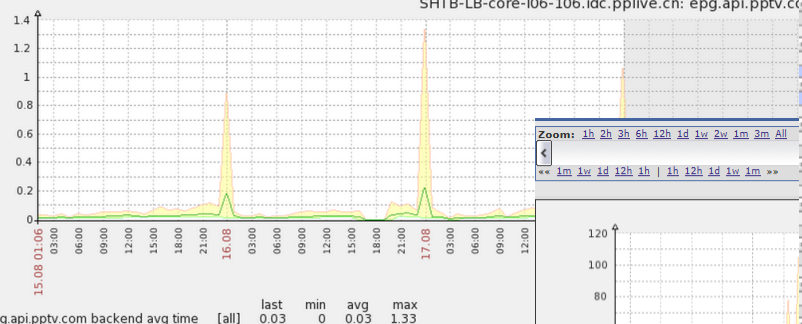
network_ss (1 Items)

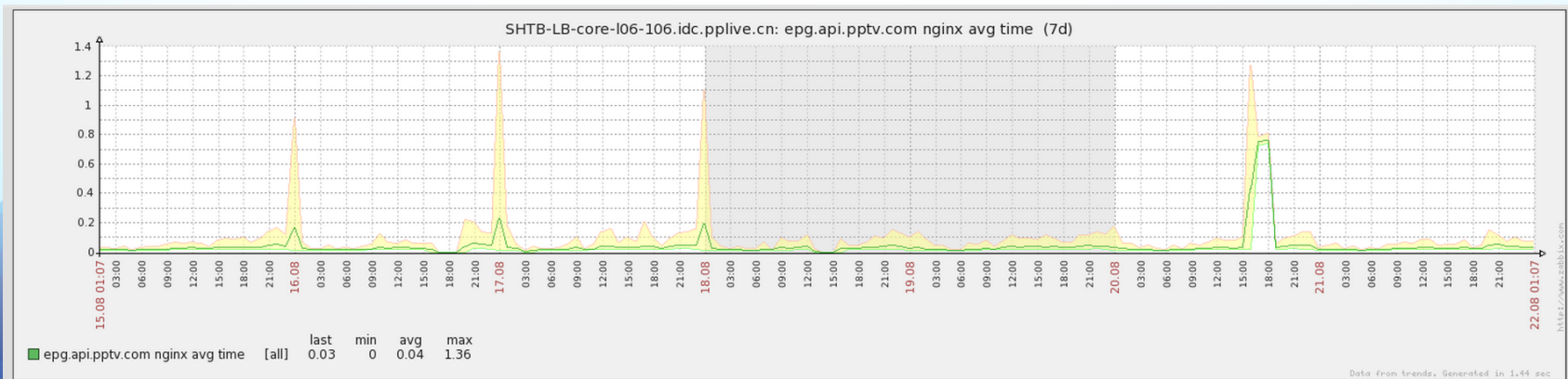
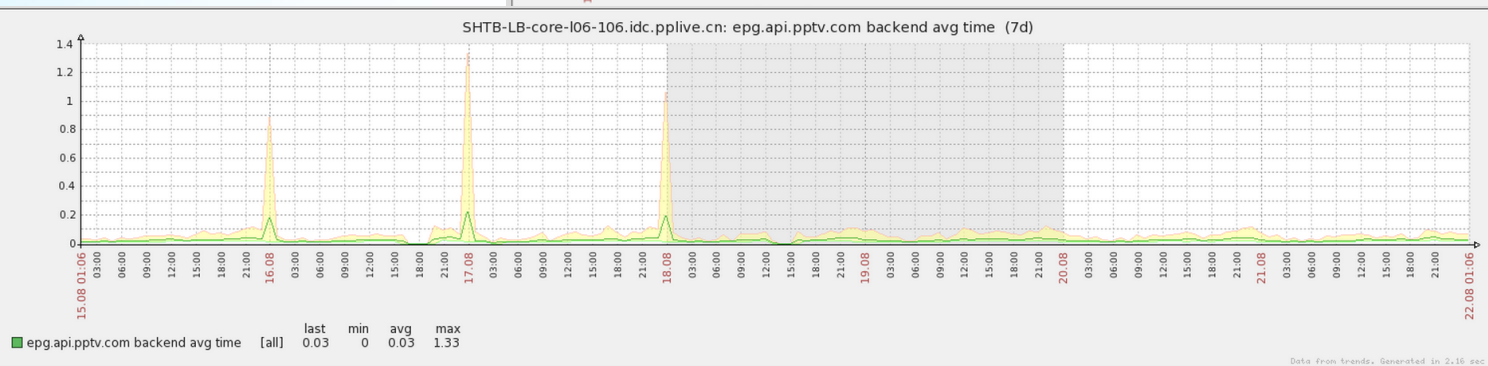
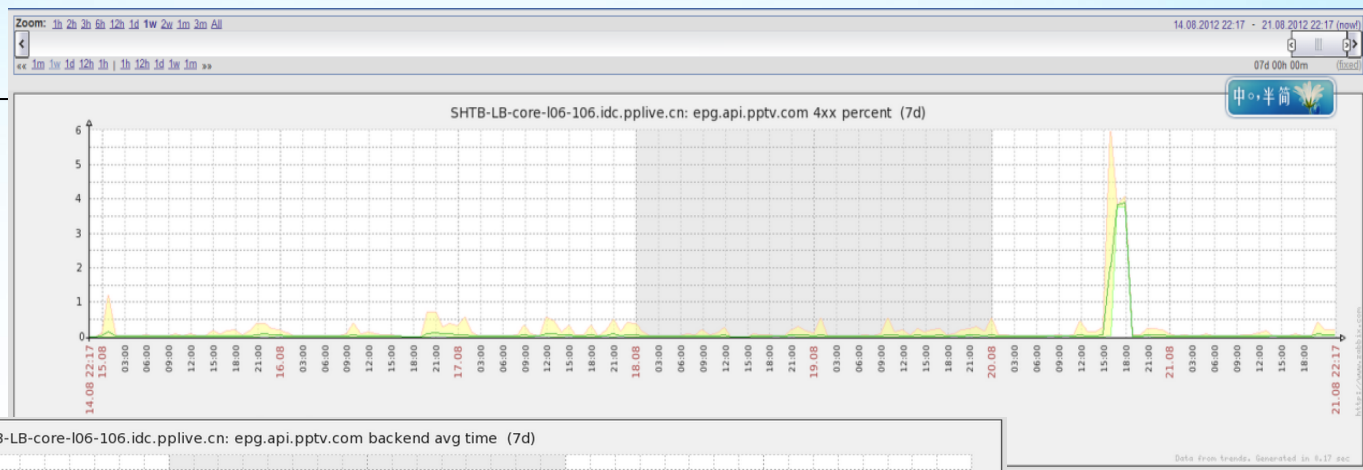
puppet (2 Items)

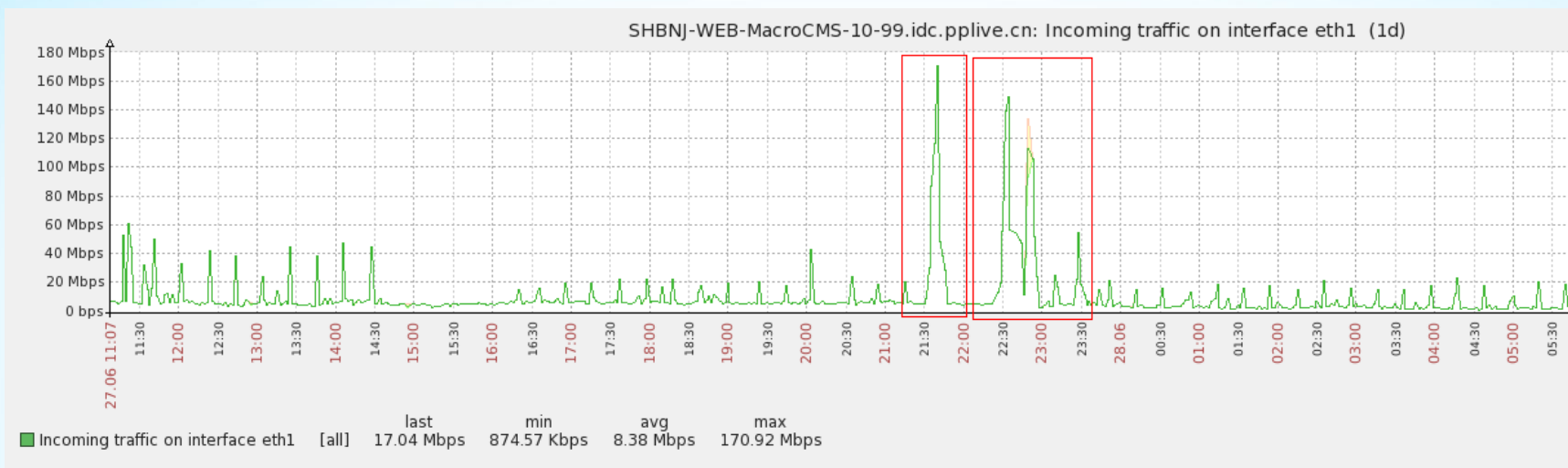
safe_check (1 Items)

systeminfo (15 Items)

- other - (2 Items)







答复: [Notify][P2]OutOfMemoryError find in /home/logs/resin/stderr.log on BJCB-WEB_searchtips-resin-101-96.idc.pplive.cn

答复: [Notify][P1][Site Alarm]Web:client-searchtips.pptv.com 4xx percent > 5%

应用日志监控

- Linux系统/Nginx/PHP日志（关键字规律）
- Java OOM
- LB 4xx/5xx错误统计
- CDN 错误统计

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

业务数据监控

- 实时业务数据的监控是最重要的！
 - 服务器down 了没关系
 - App down 了没关系
 - 在赚钱就行！
- 介于运营统计和运维监控之间
 - 模糊地带

自动化运维

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

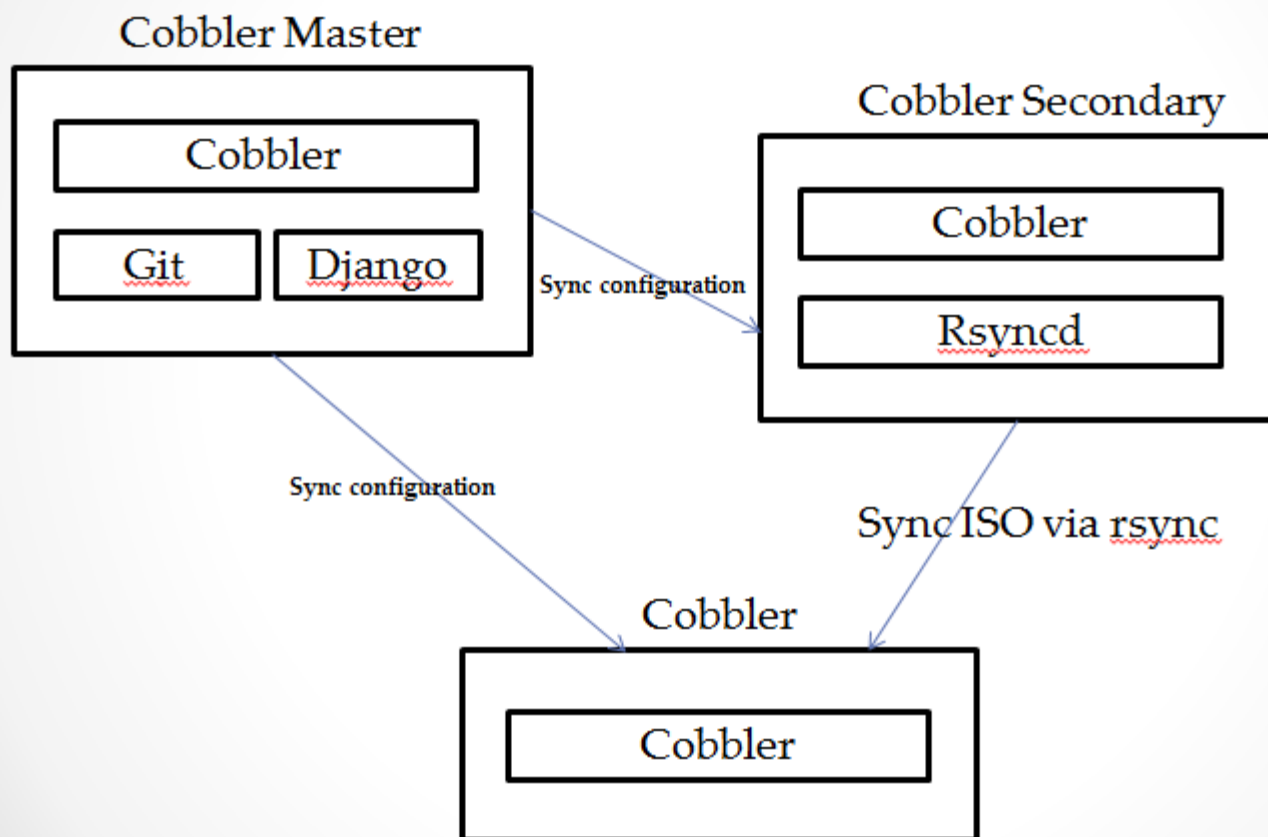
架构设计 · 自动化运维 · 云计算

运维自动化

- 统一标准的基础上才能实现自动化
- 自动化才能保证标准化
- 自动化：
 - 应用标准，系统标准，减少诡异问题
 - 大大节约人力:减少对工程师个人经验的依赖
 - 大大提高响应速度: **4小时**减少到**5分钟**
- 全网装机自动化(CDN/核心网)
- Puppet 保证:应用配置与环境统一 (90%以上)
- Web应用部署90%自动化
- CDN应用100%自动化(Puppet+ControlTier)
- 上包自动化(p2p, 多终端,Client)
- 监控自动化
 - 上线自动添加基础监控
 - 下线，维护自动暂停/删除监控
- 排错自动化 -简单排错

CDN机房自动化装机

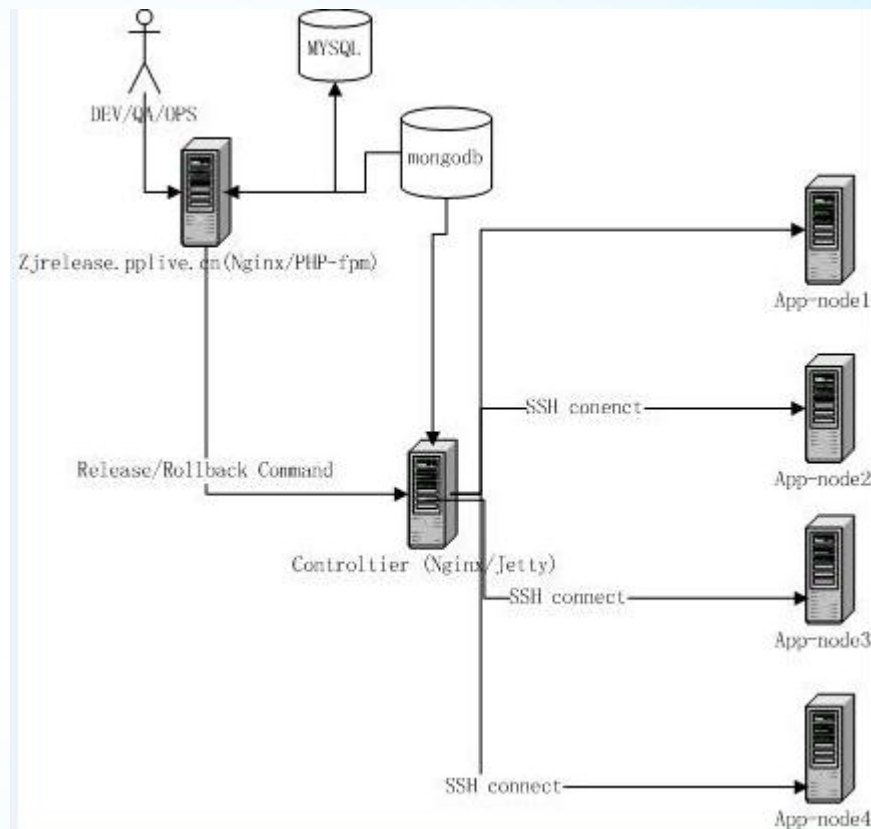
■ ControlTier + Git + Cobbler



一键发布

■ PHP + ControlTier

- 逐个发布
- 与cmdb集成
- 与Zabbix/URLmon集成
- CAS认证
- 借助LB日志分析
- 一键回滚



自动化排错

■ Zabbix Rule Engine

- 适应删除日志
- 自动重启:php/Resin
- 告警邮件自动分析
- 自动Java thread dump
- ...

http code by each backend server

```
1 "10.204.101.111:8080" "-"
3 "10.204.101.112:8080" "-"
328 "10.204.101.112:8080" "404"
342 "10.204.101.111:8080" "404"
```

http code by URL:

```
8 Code: 404 URL: /%b2%bd%b2%bd%be%aa%d0%c4(%b5%da32%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%b2%bd%b2%bd%be%aa%d0%c4(%b5%da33%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%b2%bd%b2%bd%be%aa%d0%c4(%b5%da34%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%b2%bd%b2%bd%be%aa%d0%c4(%b5%da35%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%c9%f1%d6%ae%cd%ed%b2%cd(%b5%da17%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%c9%f1%d6%ae%cd%ed%b2%cd(%b5%da18%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%c9%f1%d6%ae%cd%ed%b2%cd(%b5%da30%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%c9%f1%d6%ae%cd%ed%b2%cd(%b5%da31%bc%af).mp4 Refer: "-"
8 Code: 404 URL: /%c9%f1%d6%ae%cd%ed%b2%cd(%b5%da32%bc%af).mp4 Refer: "-"
12 Code: 404 URL: /%5Bmobile%5D%BC%E0%D3%FC%B7%E7%D4%C6.mp4dt?type=aphone Refer: "-"
```

http code by Server Name:

```
1 Code: 499 Server Name: jump.gld.net
3 Code: 499 Server Name: jump.synacast.com
30 Code: 404 Server Name: jump.gld.net
640 Code: 404 Server Name: jump.synacast.com
```

http code by Client IP:

```
2 Code: 404 Client IP: 180.153.106.10
2 Code: 404 Client IP: 183.9.165.44
```

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

数据化运维

-----监控和日志数据的分析

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

数据化运维

- 监控 → 告警
 - 告警阈值的制定
- 日志分析
 - 服务域名可用性
 - CDN命中率，性能，流量，
 - Top-down 分析
- 容量分析和优化
 - 提前发现问题
- 结合CMDB
 - 机器究竟都在干嘛？数据！

关于网站架构

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

按照数据流提炼网站通用基础架构

- 系统稳定性取决于应用架构
 - Dev 了解 IDC Strategy
 - Ops了解 App 数据流
- Ops总结几个数据流的Pattern，归纳几个通用解决方案
 - 简化
 - 服务话
 - 消除单点
 - Scalable service

架构优化

- 推动应用架构简化
 - More tiers, more tears
- 通过事故推动研发优化/简化架构
 - 不然没有人鸟你
- 每一层有高可用方案，每一层可以扩展
 - 尤其是数据库层面

搞好核心数据库平台

■ 为什么要整合数据库

- 数据库用外网IP！
- 数据库用虚拟机！
- 130+ MySQL
- 太多管理overhead, 浪费机器, 太多问题

■ 目标：

- 内网，完整监报告警
- 解决关键业务的高读可用问题(HA自动failover)
- 所有MySQL都有实时Replication机器

■ MYSQL的运维自动化和细粒度监控：

- mycon(所有端口的10s指标)
- mysql_monitor(高负载自动截取数据),
- Slow_Query_Report(每日慢查询报表)
- mora(自动化运维和分析)

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

缓存方案

■ Principle

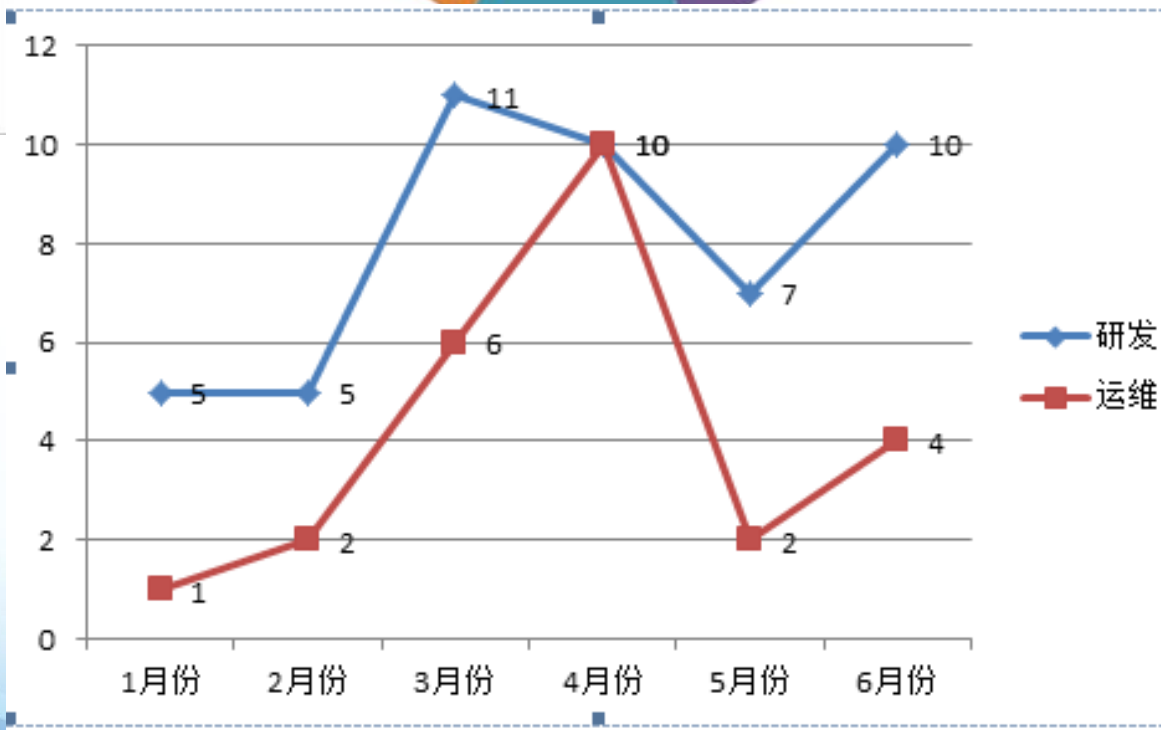
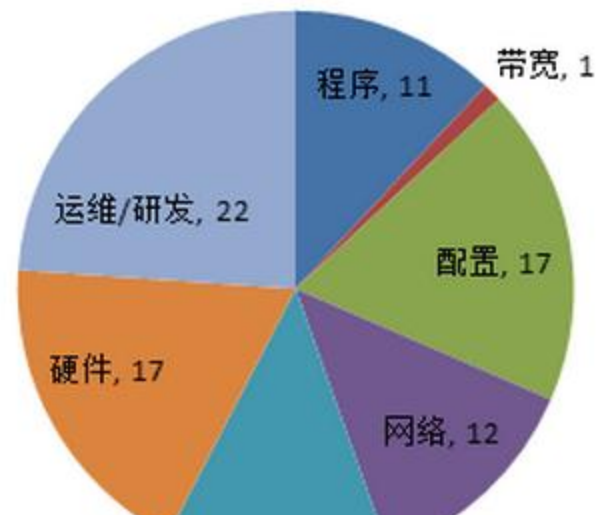
- 只支持MC/Redis
- 一致性hash，推荐3个部署
- 不命中，要回源
- 保持足够回源能力

■ 确保

- 足够带宽
- 交换机互联
- 监控，容量规划
- 不用太新feature

流程

- 变更管理流程
- Incident Report
- 成本Report



VAS业务问题汇总

About DevOps

- 知道IDC布局
- 知道IDC 策略
- 知道基本硬件平台
- 可以看到服务器状态，性能指标，看到日志
- Should Dev has production access?
- Platform for Dev/Ops
- Dev should understand the environment their code runs on
- Ops should understand the Dev code/logic

我的另外一些想法

- Don' t re-invent the wheel
- Embracing the open-source
- Don' t be on your own
- Adapt the mainstream
- Don' t be too aggressive
- Simple is beauty

SACC

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算

2012

Monitoring

Windows平台的集中配置管理和监控

实时业务数据监控

网络质量/安全监控

运维日志监控/分析平台

BI日志收集平台

运维资源平台/CMDB

Automation

资源使用分析和容量优化

应用Release自动化

核心网交换机高可用和千兆互联解
决

内部DNS

VM 管理平台

排错分析自动化

Application

应用标准化

Layered Service

Redis/Mongo/Hadoop

应用安全

攻击报警系统

2012中国系统架构师大会

SYSTEM ARCHITECT CONFERENCE CHINA 2012

架构设计 · 自动化运维 · 云计算