



破解阿喀琉斯之踵

--利用系统开发破解架构中的硬伤



运动啦！ 七、八年就来一次~！

我们是谁？

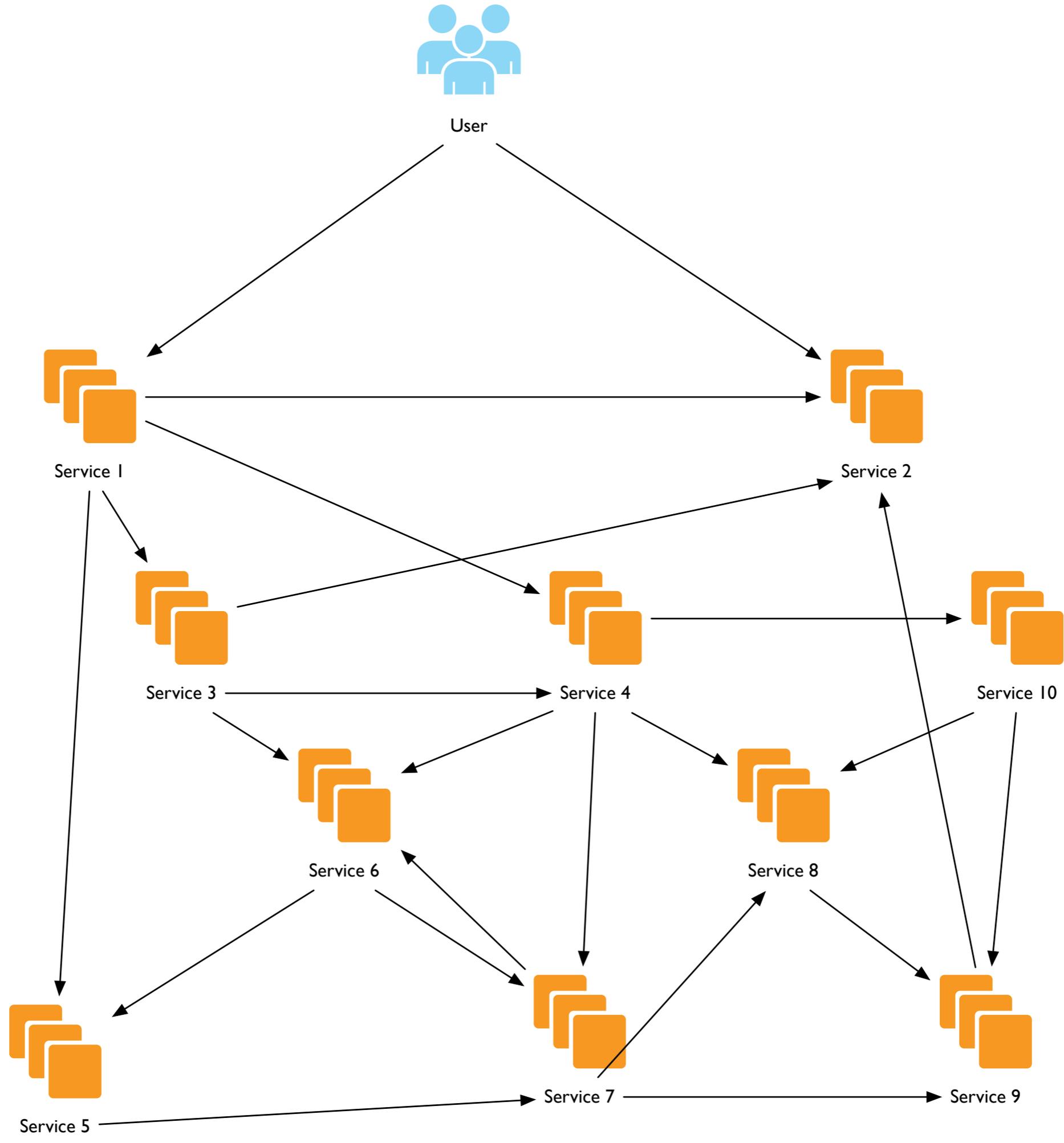
- DevOps @Sina
- Lamp项目托管平台
- 已托管项目500+
- 日均请求10B+

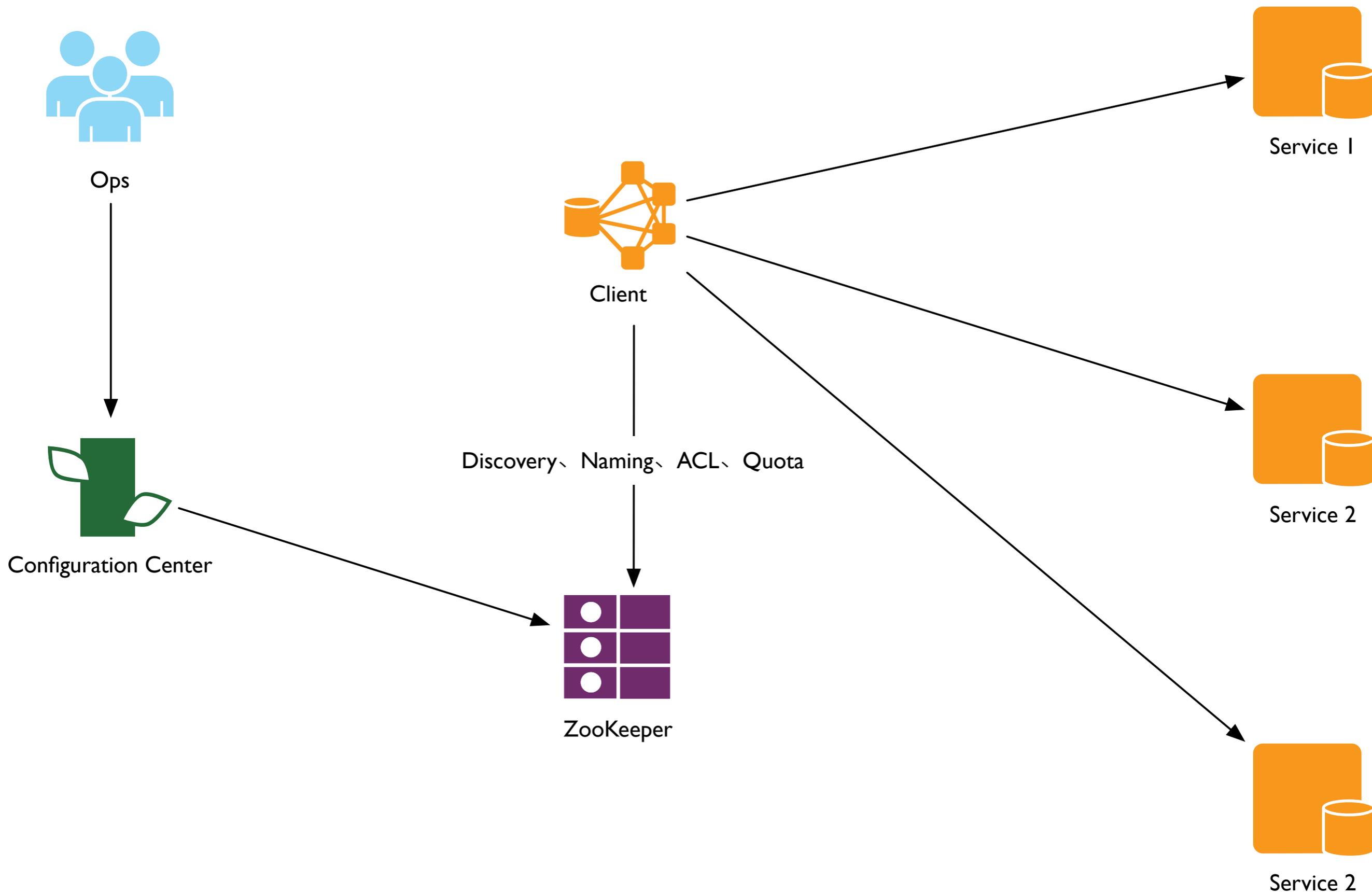
今天要破解哪些硬伤

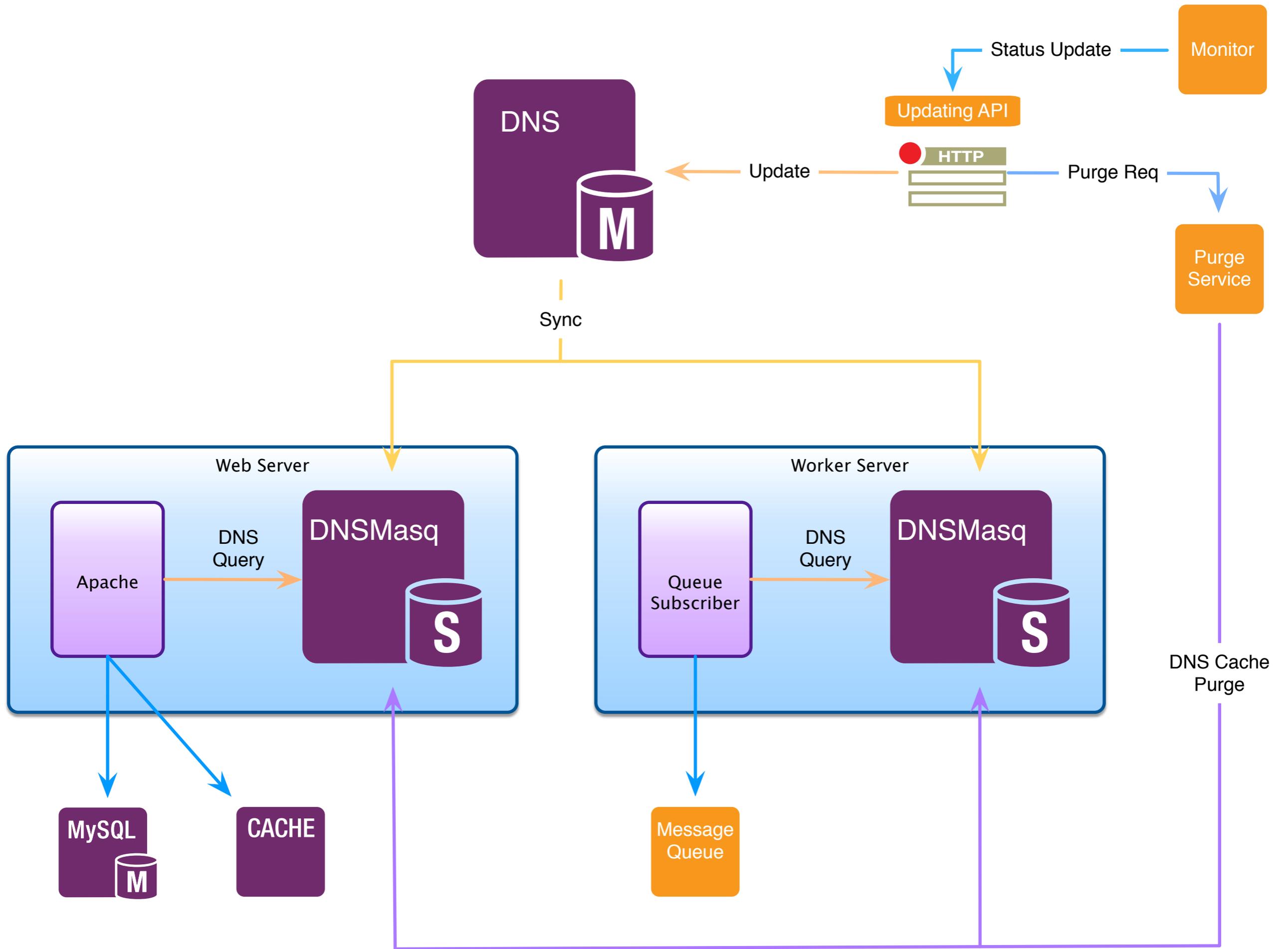
- SOA陷阱
- Memcached “太”好用了
- RDBMS难以扩展
- 困扰我们的“阿喀琉斯之踵”

我们都爱SOA

- 解耦
- 分而治之
- ...太多名词来描述优点和必要性
- 当目标太多的时候，是不是容易变成为了做而做？







通用命名服务--*Macedonia*

- C、Python、Node.js、PowerDNS、DNSMasq、MongoDB
- DNS协议，更新时通知所有客户端，缓存一致性强
- 多样化API
- 无监控框架，需结合API设计开发
- A记录无法返回端口号



陈尔冬 LV3
1902 关注 2811 粉丝 3955 微博

- 首页
- 提到我的
- 评论
- 私信
- 收藏
- 相册
- 微音乐
- 微活动
- 微群
- 微公益
- 位置
- 更多»

- 小小乐园
- 达人麻将
- 咖啡恋人
- 神仙道

- iPad客户端
- 新浪视野
- 爱婚恋
- 微漫画
- 虾米音乐

Memcached... Wow!
Memcached... OMG!

蘇打綠福嘎: 節目很用心~團隊很棒~主持人也很好, 讓整個節目很精彩。但是出來新聞怎會這樣, QQ

@吳青峰: 關於今天的新聞, 我的聲明。如果有耐心也有空閒的朋友再看吧。我永遠愛著真心愛著也溫柔對待這個世界的你們。



7月4日22:51 来自新浪微博

转发(10657) | 评论(5027)

2分钟前 来自新浪微博

转发(1) | 收藏 | 评论(3)

shineyear: 这个好 // @笑多了会怀孕: 有人在背后议论你, 是因为你走在了他们的前面

@治愈系英文签名: People talk behind your back because you're ahead of them.

Cheer up! 有人在背后议论你, 是因为你走在了他们的前面。

写心情

勋章赢iPad2!

刘强东 + 加关注
22个间接关注人

炎亞綸 + 加关注
蔡康永等也关注TA

[活动] 赢总冠军T恤

我关注的人中: 新浪NBA、RikaV王麗嘉OctoBeez关注了Love_NBA小秘书



重度应用 Memcached 带来的挑战

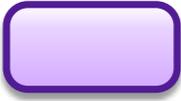
- 我的首页最多要取 $50+1+4+2$ 个用户的信息
- 我的首页最少要取 $50 \times 2+3$ 个计数器
- 为了访问速度大量使用并依赖 Memcached
- 如果有一台 Memcached 服务器宕机的话...

Legend

Server



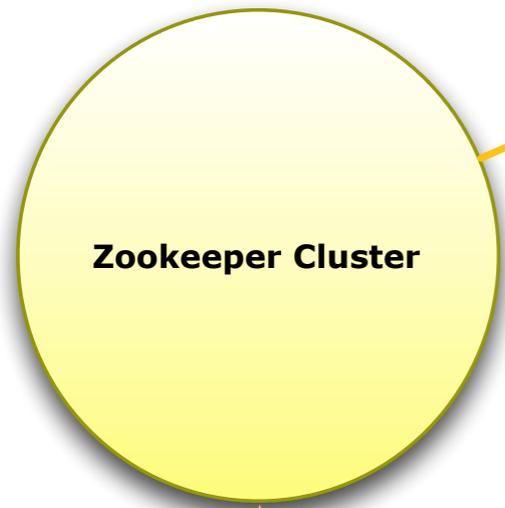
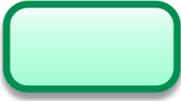
Process



Cluster

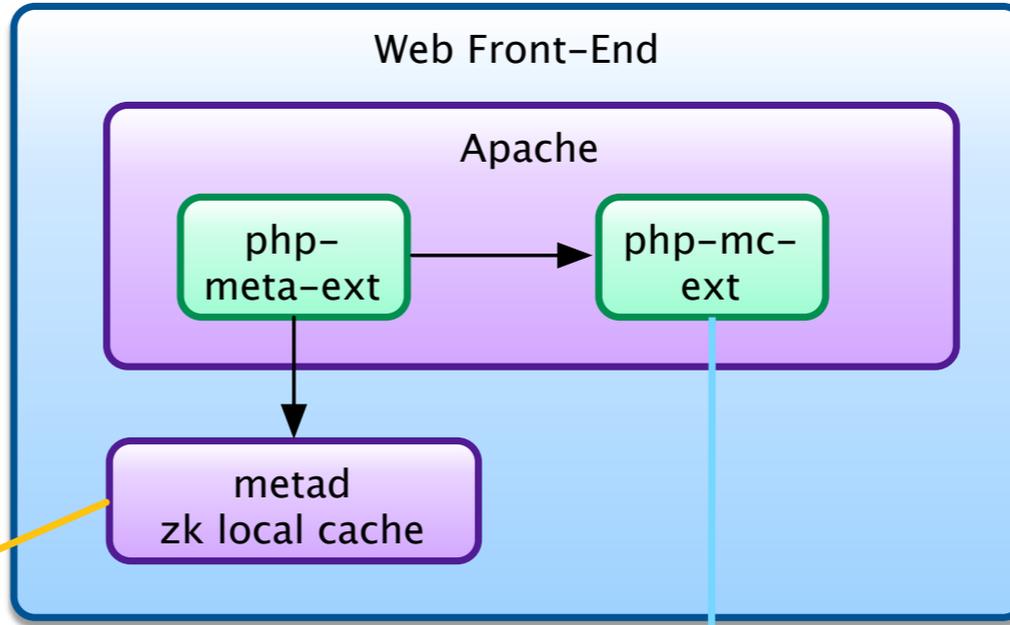


Extension



Zookeeper Cluster

Updating Notification



Web Front-End

Apache

php-
meta-ext

php-mc-
ext

metad
zk local cache

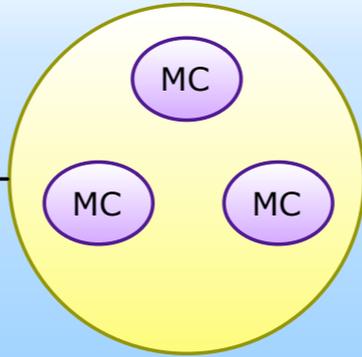
Memcached Command

Node Register
&
Heartbeat

Cache Node



Watcher



MC

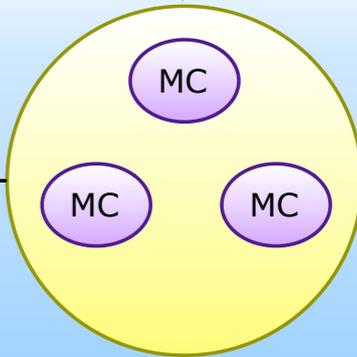
MC

MC

Cache Node



Watcher



MC

MC

MC

分布式缓存服务--*Babylon*

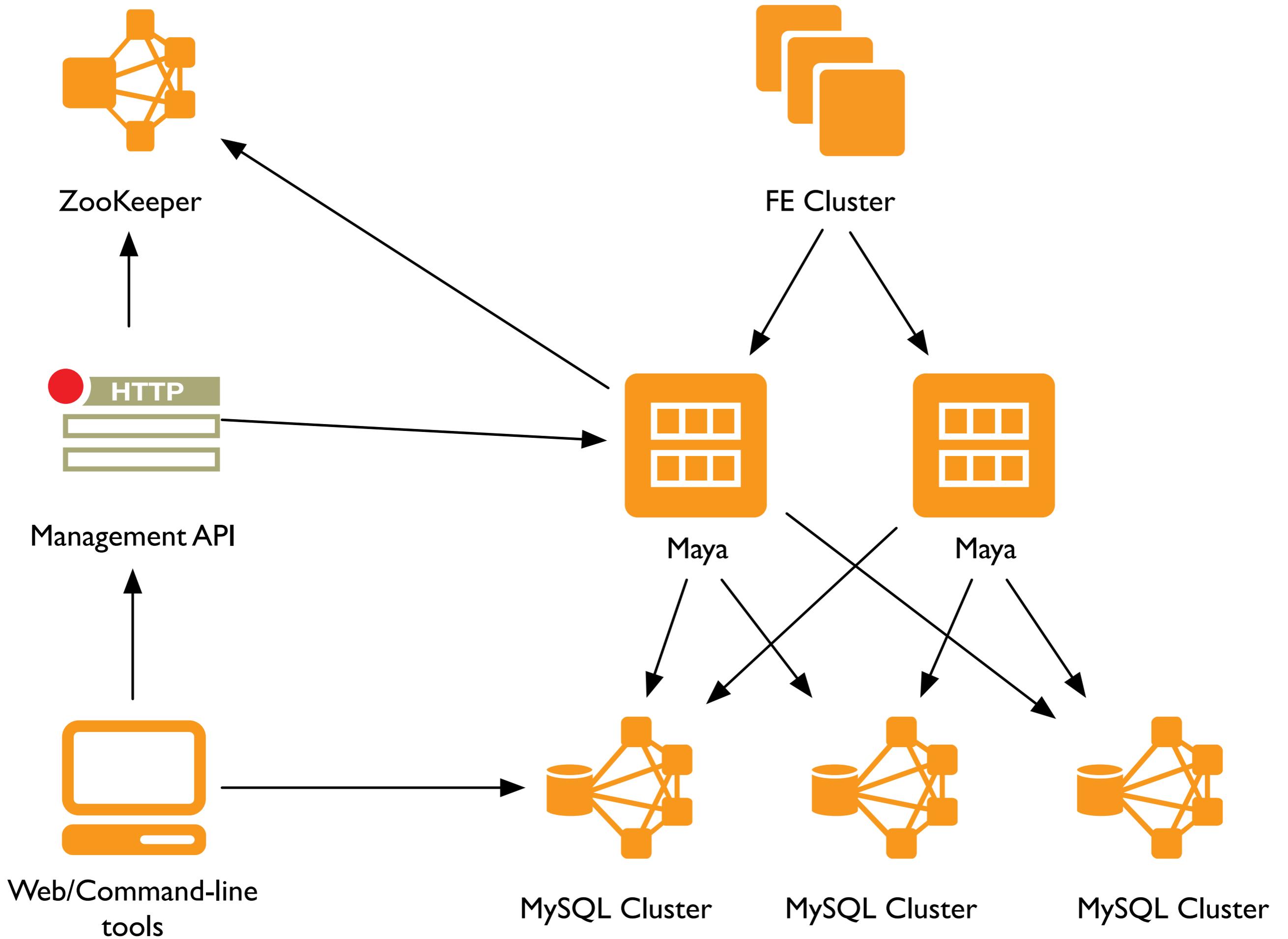
- Thrift、C、C++、Python、ZooKeeper
- 分片信息保存在ZooKeeper中，并在Metad缓存
- 支持引擎级主-主复制
- 连接数LRU（Patch from SAE）
- Memcached-Like SDK

分布式缓存服务--*Babylon*

- 未来特性：
 - 跨IDC复制（开发中）
 - 数据迁移和无缝扩容（也许）
 - FlashCached（也许）
 - 更好的管理性

DBA的烦恼

- MySQL很棒，但是...
- Range VS Hash
- 开发改造成本
- 不合“规范”的SQL
- Master HA



数据库中间层--*Maya*

- Node.JS、C、Flex & Bison、Python、
- SQL路由
- 流量统计、配额、限制
- 数据库预拆分（不是自动拆分）
- 主库半自动切换、从库自动切换
- SQL审计、自动DDL

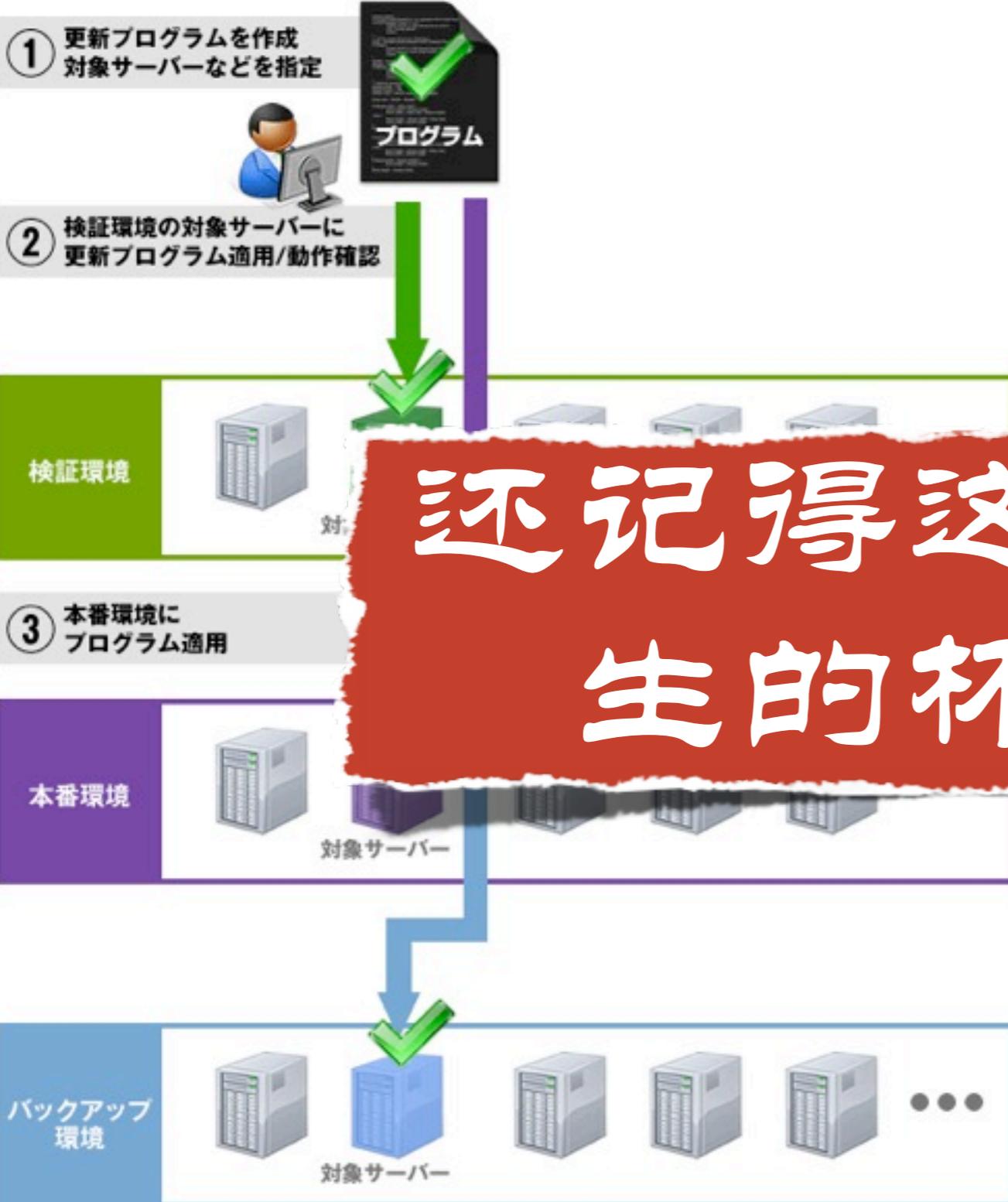


我们的 “阿喀琉斯之踵”

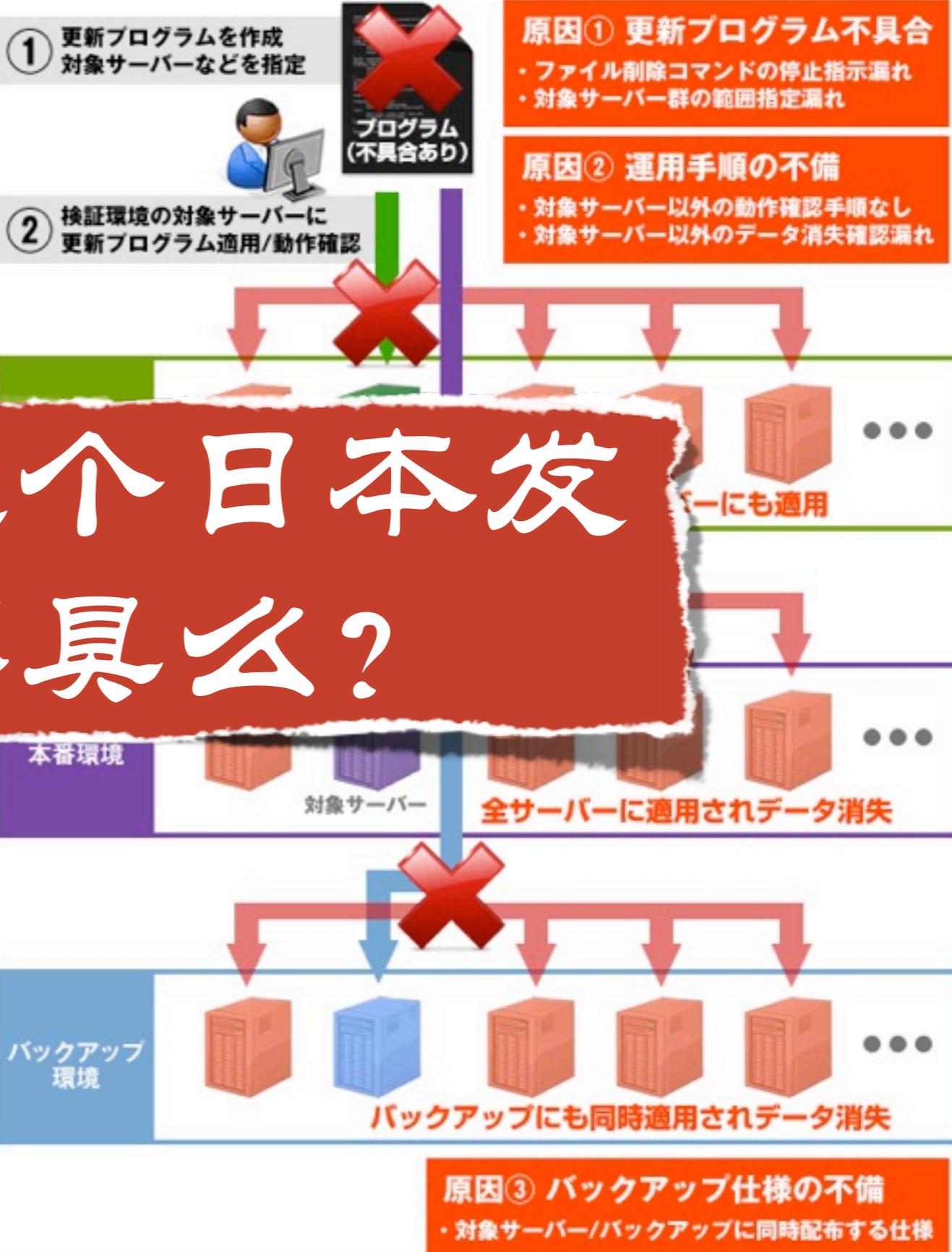
- 每年新增PHP项目100+
- 3000+台前端服务器 100亿+ Hits/天
- 45000+行虚拟主机配置 30G代码
(GZip后)
- 3000+台数据库服务器 200亿+ 数据库请求/天
- 6000+ 个MySQL实例 总存储量2P+



通常の状態



今回の事故の原因



还记得这个日本发生的杯具么？

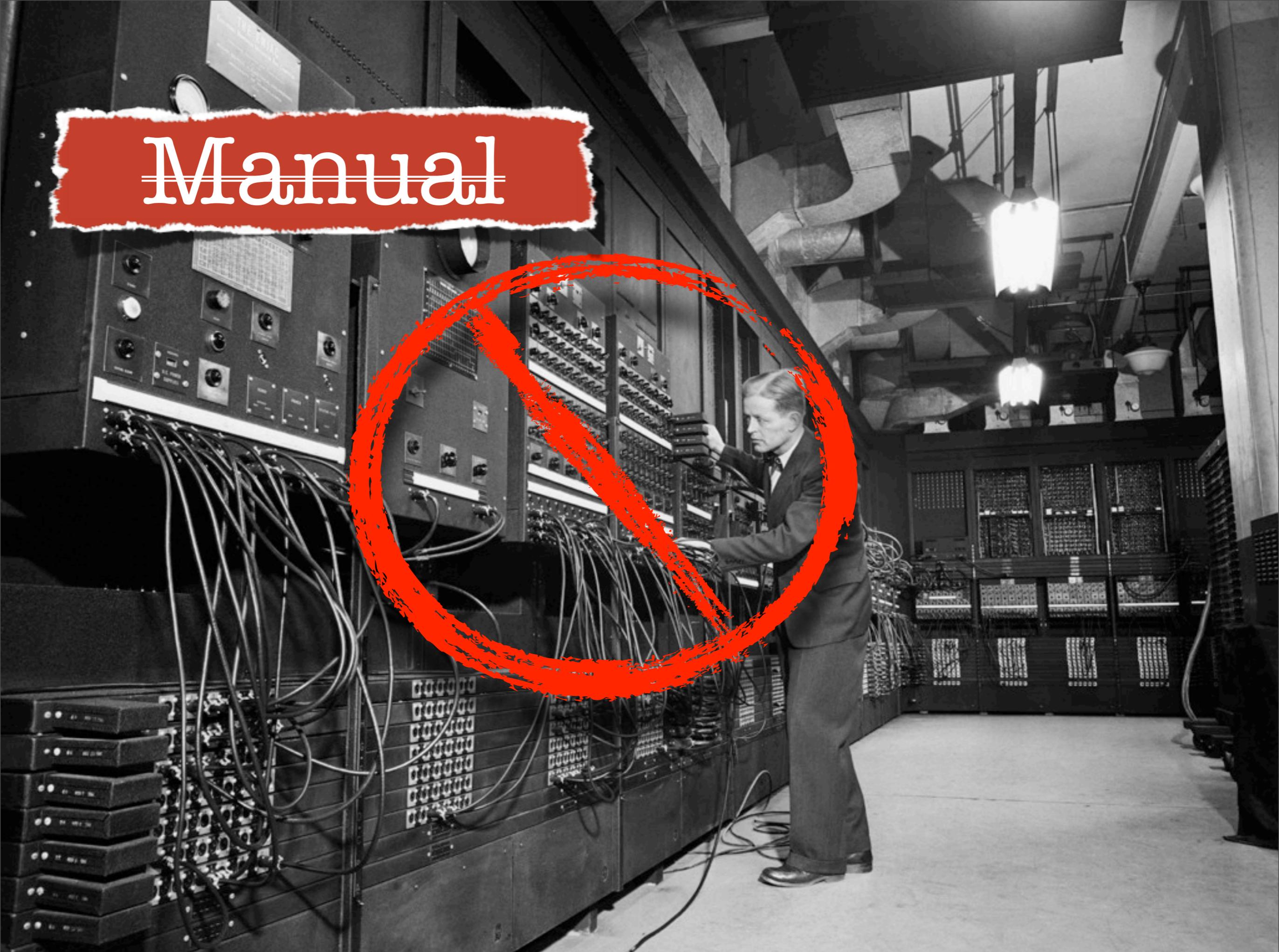
OPERATOR!

Make
Disasters

Western Electric
MICROPHONIC
SOUND SYSTEM

PRODUCED BY
AUDIO PRODUCTIONS, INC.

Manual



自动化运维萌芽

- 配置与部署越统一越容易自动化
- 将服务器按照功能划分角色
- 为每一个角色指定相应的配置文件
- 角色和配置文件的分类尽量简洁
- 差异配置利用程序来自动化处理

httpd-vhost.conf_tmpl

```
...  
SetEnv SINASRV_ZONE_IDC @@XD@@  
SetEnv SINASRV_ZONE_ISP @@CNC@@  
SetEnv SINASRV_ZONE_ID @@010201@@  
...
```

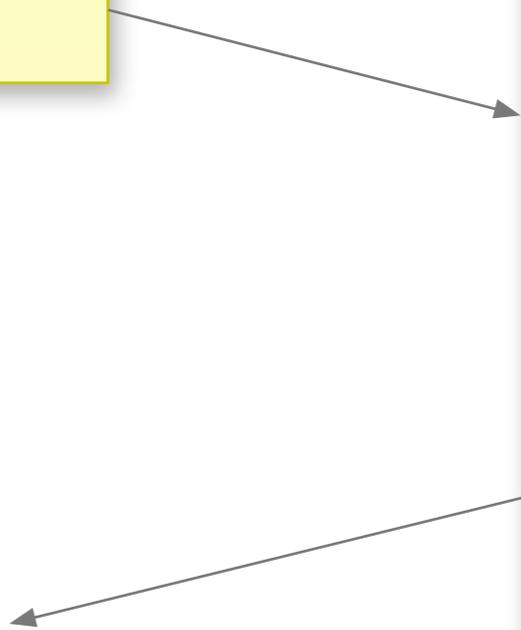
vhost_config.pl

```
...  
%idc_rule = (  
    zjm => {  
        SINASRV_ZONE_IDC => "ZJM",  
        SINASRV_ZONE_ISP => "CNC",  
        SINASRV_ZONE_ID => "010101",  
    }  
)  
...  
elsif ($line =~ m!\s*SetEnv\s+SINASRV_ZONE_IDC.+!) {  
    push @vhost_configfile, "SetEnv SINASRV_ZONE_IDC ". $idc_rule{$idc}{SINASRV_ZONE_IDC};  
    next;  
}  
elsif ($line =~ m!\s*SetEnv\s+SINASRV_ZONE_ISP.+!) {  
    push @vhost_configfile, "SetEnv SINASRV_ZONE_ISP ". $idc_rule{$idc}{SINASRV_ZONE_ISP};  
    next;  
}  
elsif ($line =~ m!\s*SetEnv\s+SINASRV_ZONE_ID.+!) {  
    push @vhost_configfile, "SetEnv SINASRV_ZONE_ID ". $idc_rule{$idc}{SINASRV_ZONE_ID};  
    next;  
}  
)  
...
```



FE Node

/etc/httpd-vhost.conf



对于管理性低的 开源软件

- 那就创建一个配置文件吧
- 当然要尽量创建一个全局统一的
- 索性与监控程序结合起来怎么样?
- INI VS YAML
- 就比如Memcached!

mc.conf

```
[Mblog_Userinfo]
ip=10.55.22.100:10000 10.55.22.101:10000 10.55.22.102:10000 10.55.22.103:10000
memsize=10.55.22.100:1G 10.55.22.101:1G 10.55.22.102:1G 10.55.22.103:1G
conn_limit=10.55.22.100:1W 10.55.22.101:1W 10.55.22.102:1W 10.55.22.103:1W

[Mblog_Counter]
ip=10.55.22.100:10001 10.55.22.101:10001 10.55.22.102:10001 10.55.22.103:10001
memsize=10.55.22.100:1G 10.55.22.101:1G 10.55.22.102:1G 10.55.22.103:1G
conn_limit=10.55.22.100:1W 10.55.22.101:1W 10.55.22.102:1W 10.55.22.103:1W
```



MC Node

```
/etc/init.d/memcached
/etc/cron.d/check_mc
```



Monitor Node

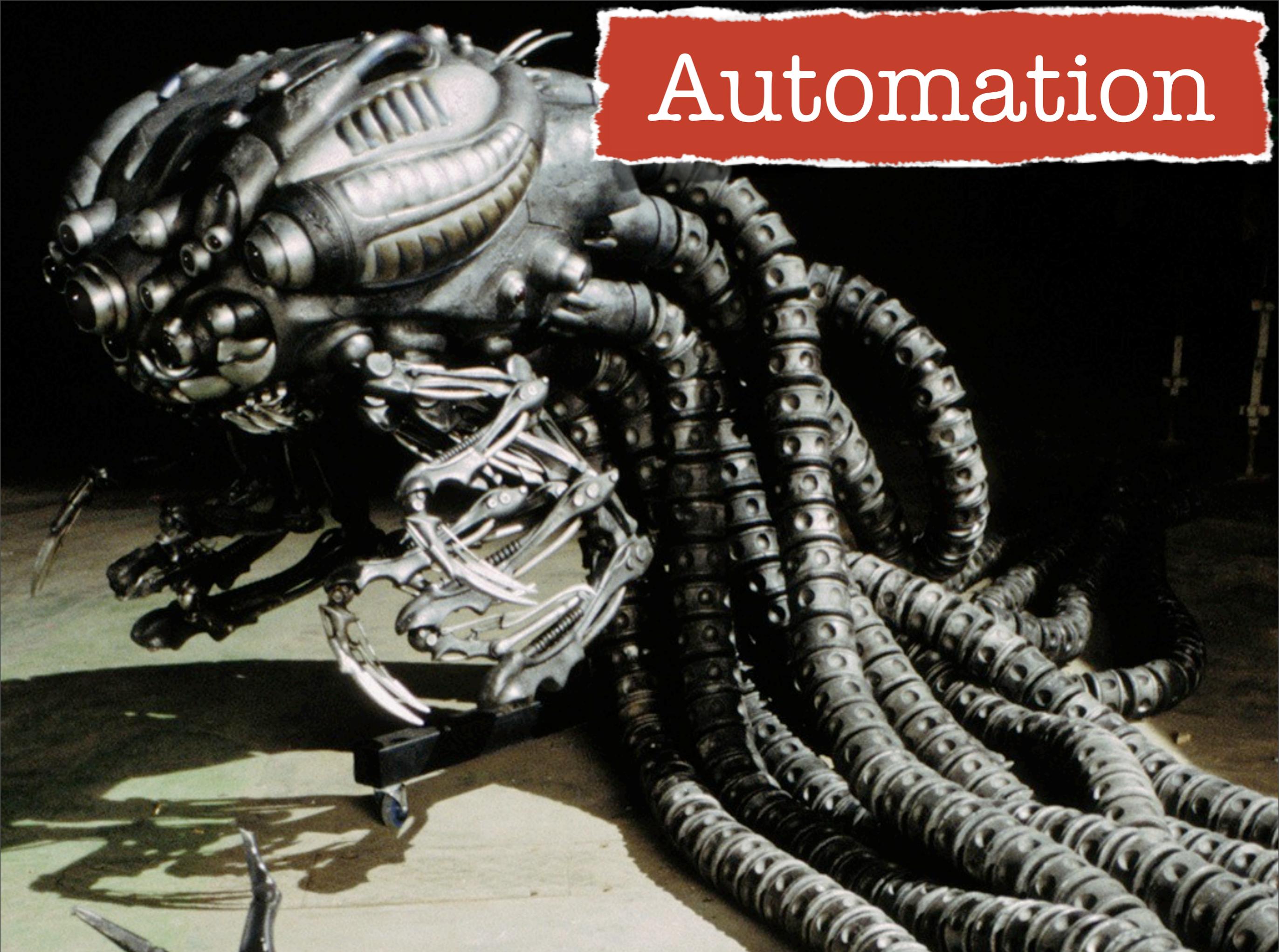
```
/etc/monitor.d/check_mc.py
```



Reporting Node

```
/etc/reporting.d/get_mc_stats.py
```

Automation



一年后



线上变更管理

- ▣ 工作流程管理
- ▣ workflow 日志
- ▣ workflow 执行日志
- ▣ 配置文件管理
- ▣ 包部署管理
- ▣ WebShell
- ▣ 用户临时授权

运维信息管理

- ▣ 产品单元管理
- ▣ 环境设置
- ▣ 模块设置
- ▣ IDC 设置
- ▣ 负载均衡

监控报警管理

- ▣ 性能监控
- ▣ 报警管理
- ▣ 基础监控模板设置

资源资产管理

- ▣ 节点管理
- ▣ DB 资产管理

数据库应用

- ▣ 端口管理
- ▣ 域名管理
- ▣ 备份管理

实验特性专区

- ▣ Cron 管理

系统管理专区

- ▣ 管理日志
- ▣ 账户设置
- ▣ 修改资料
- ▣ 注销登录

隐藏左栏

» 编辑配置文件 - wptest-httpd-vhost.conf

名称: (该名称是为了方便查看, 与实际部署无关) wptest-httpd-vhost.conf

目标路径: (绝对路径, 包括文件名, 指在目标机器上的存放路径) /etc/dAppCluster/httpd-vhost.conf_dpool2

文件权限: (默认"0644") 文件所有者: (默认"root:root")

关联模块: [添加一行](#) [添加DB组模块](#) [添加DPOOL组模块](#)

模块0: 类型:

模块1: 类型:

模块2: 类型:

配置文件下发前要执行的检查命令(可选, 用于执行配置检查, 如果该命令未能正确执行(即\$?不为0), 则拒绝替换, 配置文件地址请用\$(conf_path)代替):

配置文件下发后要执行的命令(可选, 建议有的操作: ps查看进程, kill -HUP重启服务, tail查看访问/错误日志):

配置文件内容

```

MaxClients @@800@@

<IfModule vhost_limit_module>
    MaxVhostClients @@200@@
</IfModule>

SetEnv SINASRV_GLOBAL_MEMCACHED_SERVERS "@@10.44.6.44:7601 10.44.6.45:7601 10.44.6.46:7601 10.44.6.47:7601 10.44.6.48:7601 10.44.6.49:7601 10.44.6.50:7601"
SetEnv SINASRV_MEMCACHED_SERVERS "@@10.44.6.232:7601 10.44.6.234:7601 10.44.6.238:7601 10.44.6.240:7601@@@"
SetEnv SINASRV_MEMCACHED_HOST 127.0.0.1
SetEnv SINASRV_MEMCACHED_PORT 7600
SetEnv SINASRV_DATA_DIST_SERVER @@10.44.6.43@@
SetEnv SINASRV_DATA_DIST_PORT 8080

Header add DPOOL_HEADER "Hello World!"

SetEnv SINASRV_ZONE_IDC @@XD@@
SetEnv SINASRV_ZONE_ISP @@CNC@@
SetEnv SINASRV_ZONE_ID @@010201@@
SetEnv SINASRV_ROLE web
SetEnv SINASRV_OUTIP outip
  
```

和最新线上版本的配置差异

```

--- /var/tmp/httpd-vhost.conf_dpool2.old.tmp      2012-05-24 18:00:02.000000000 +0800
+++ /var/tmp/httpd-vhost.conf_dpool2.new.tmp      2012-05-24 18:00:02.000000000 +0800
@@ -205,6 +205,7 @@
     SetEnv SINASRV_NDATA_CACHE_URL "http://cache.@@mars@@.weibo.com/nd/photoweibo/"

     SetEnv SINASRV_MEMCACHED_ALBUMS_SERVERS "10.75.31.39:7122 10.75.31.38:7122 10.75.31.37:7122 10.75.31.36:7122 10.77.7.85:7122 10.77.7.86:7122"
+    SetEnv SINASRV_MEMCACHED_NEW_ALBUMS_SERVERS "10.75.31.40:7122 10.75.31.39:7122 10.75.31.38:7122 10.75.31.37:7122 10.75.31.36:7122 10.77.7.85:7122 10.77.7.86:7122"
  
```

线上变更管理

- ▣ workflows管理
- ▣ workflows日志
- ▣ workflows执行日志
- ▣ 配置文件管理
- ▣ 包部署管理
- ▣ WebShell
- ▣ 用户临时授权

运维信息管理

- ▣ 产品单元管理
- ▣ 环境设置
- ▣ 模块设置
- ▣ IDC设置
- ▣ 负载均衡

监控报警管理

- ▣ 性能监控
- ▣ 报警管理
- ▣ 基础监控模板设置

资源资产管理

- ▣ 节点管理

数据库应用

- ▣ 端口管理
- ▣ 域名管理
- ▣ 备份管理

实验特性专区

- ▣ Cron管理

系统管理专区

- ▣ 管理日志
- ▣ 账户设置
- ▣ 修改资料
- ▣ 注销登录

隐藏左栏

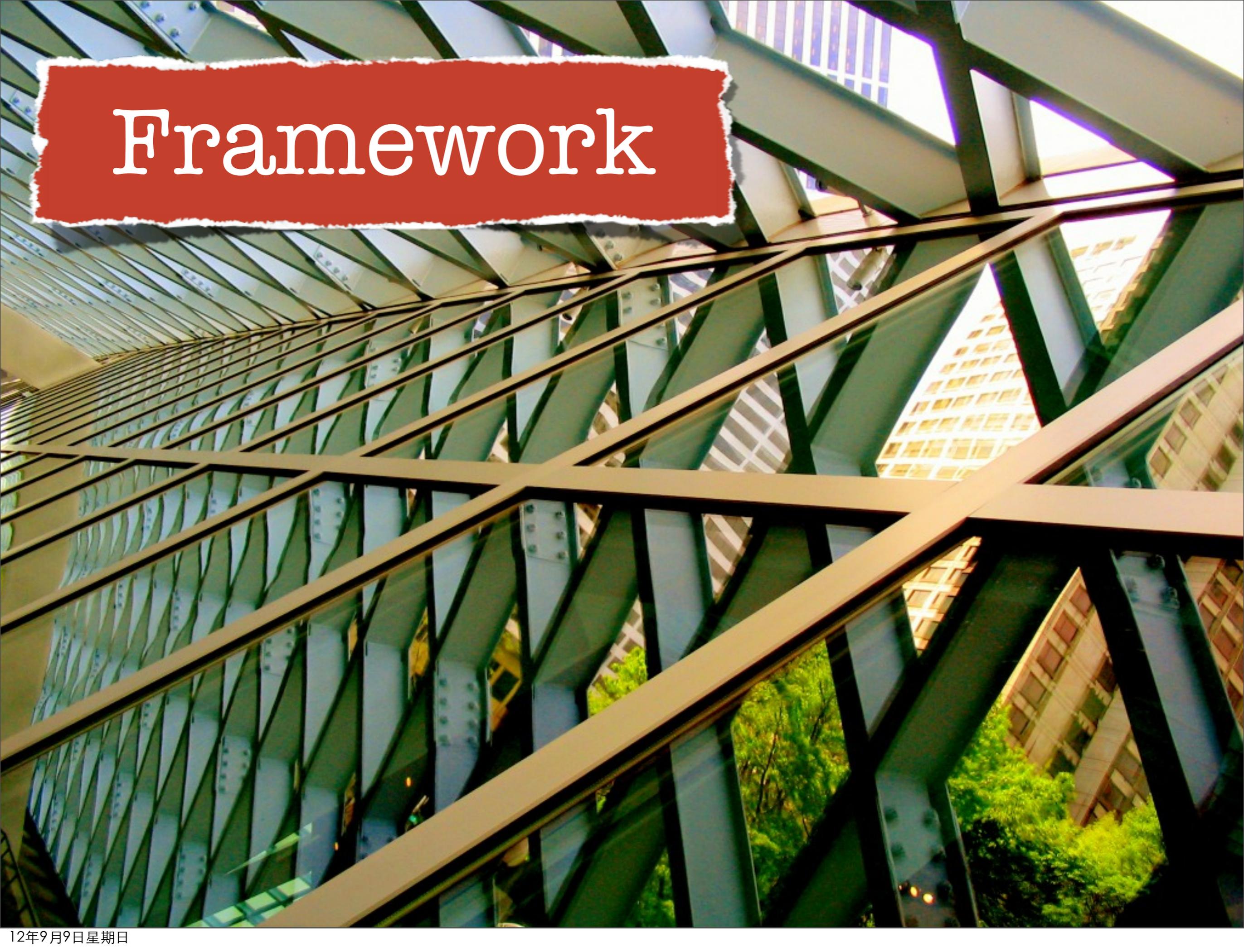
» 节点管理(默认只显示30条数据, 请善用"查找"功能)

环境: 所有环境 模块: 数据库从库 IDC: 北京土城 关键字: IP 查找 共166条结果									
状态及操作	内网IP	外网IP	环境	初始化状态	模块				
					模块操作	模块名	端口	产品单元	
<input type="checkbox"/> ✓ 操作 <div style="border: 1px solid gray; padding: 2px; margin-top: 5px;"> [编辑] [下线] [批注] [删除] [收集硬件信息] [下发配置] [撤销初始化状态] [监控设置] [发起报修] </div>	10.73.11.119	10.72.11.119	线上	<input checked="" type="checkbox"/>	[初始化]	Fusion_io(fusion_io)		-----	
					[初始化]	硬件RAID(hwr aid)		-----	
					查看端口 待监测 扩容 [初始化]	MytriggerQ MySQL 从库(mysql-innodb)	3618	数据库平台 / 微博客 / 基础服务 / myfeed	
					查看端口 待监测 扩容 [初始化]	数据库从库(mysql-innodb)	3618	数据库平台 / 微博客 / 基础服务 / myfeed	
					查看端口 待监测 扩容 [撤销初始化状态]	mytriggerMySQL从库(t)	3618	数据库平台 / 微博客 / 基础服务 / myfeed	
<input type="checkbox"/> ✓ 操作	10.73.25.86	10.72.25.86	线上	<input checked="" type="checkbox"/>	查看端口 待监测 扩容 [初始化]	主模块: 数据库从库(mysql-innodb)	4624	数据库平台 / 微博客 / 基础服务 / 关注	
					[初始化]	SSD群集(ssd)		-----	
					[初始化]	硬件RAID(hwr aid)		-----	
					查看端口 待监测	主模块: 数据库从库(mysql-innodb)	462	数据库平台 / 微博客 /	

足够自动化了吗？

- 自动化了很多个点
- 但是运维是一个体系
- 但当我们尝试把点联起来的时候
- 我们发现不是点不够多
- 就是重复发明了轮子

Framework



我们发现问题在出发点

- 自动化运维不是“另一个”系统
- 它是对你架构可运维性的更高要求
- 需要将可运维性作为软件和架构设计的主要考量因素之一
- 运维工具也是软件和架构的一部分

it's

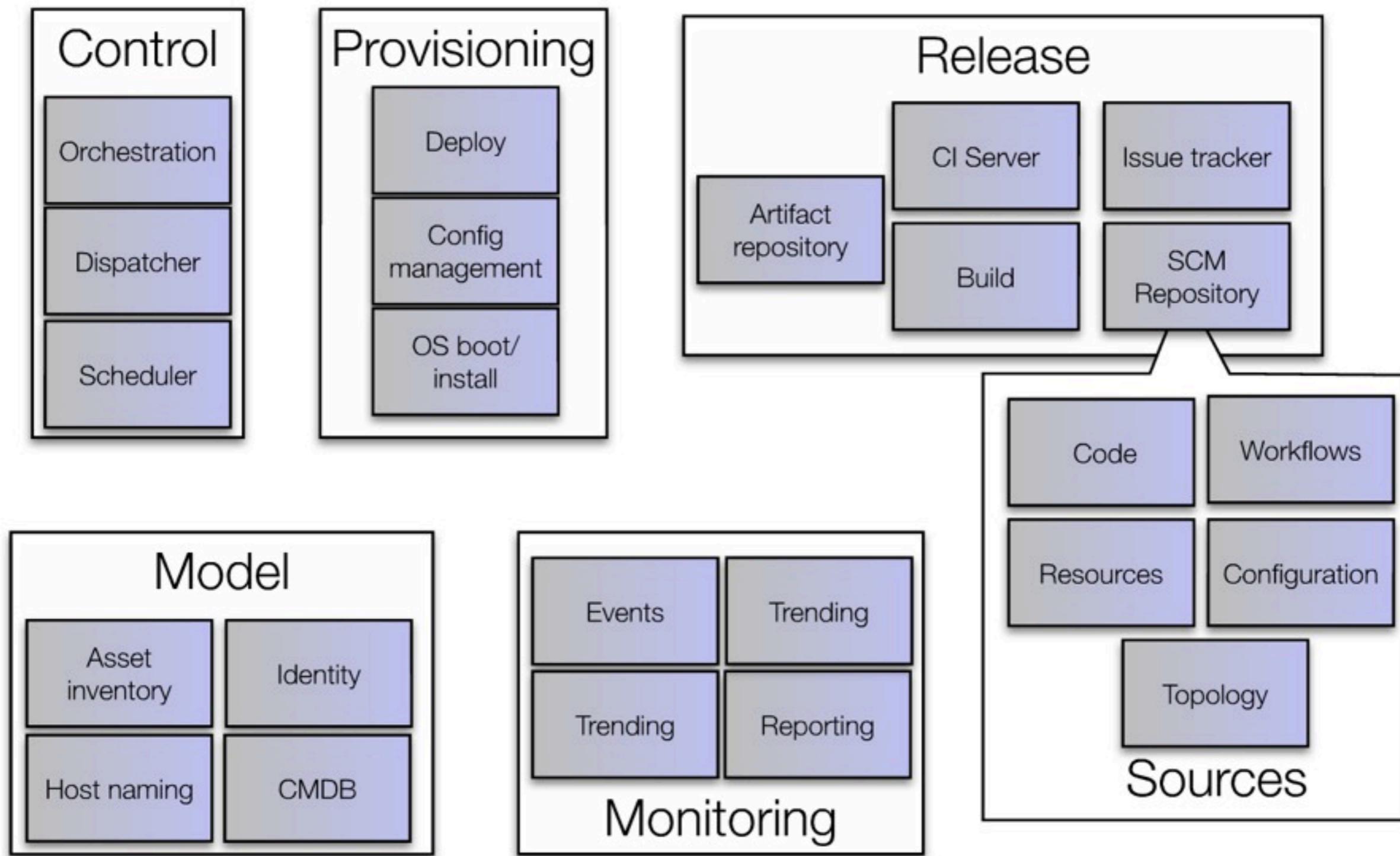
Revolution **O**f **M**artyred **E**lites





ROME

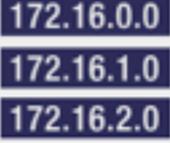
Generalized architecture




Operation as a Service


Workflow


Administration


172.16.0.0
172.16.1.0
172.16.2.0
Naming


Command-line tools


Orchestration


Load Balancing


Reporting


Asset Inventory


Auditing


MySQL
Database


Dashboards


HTTP

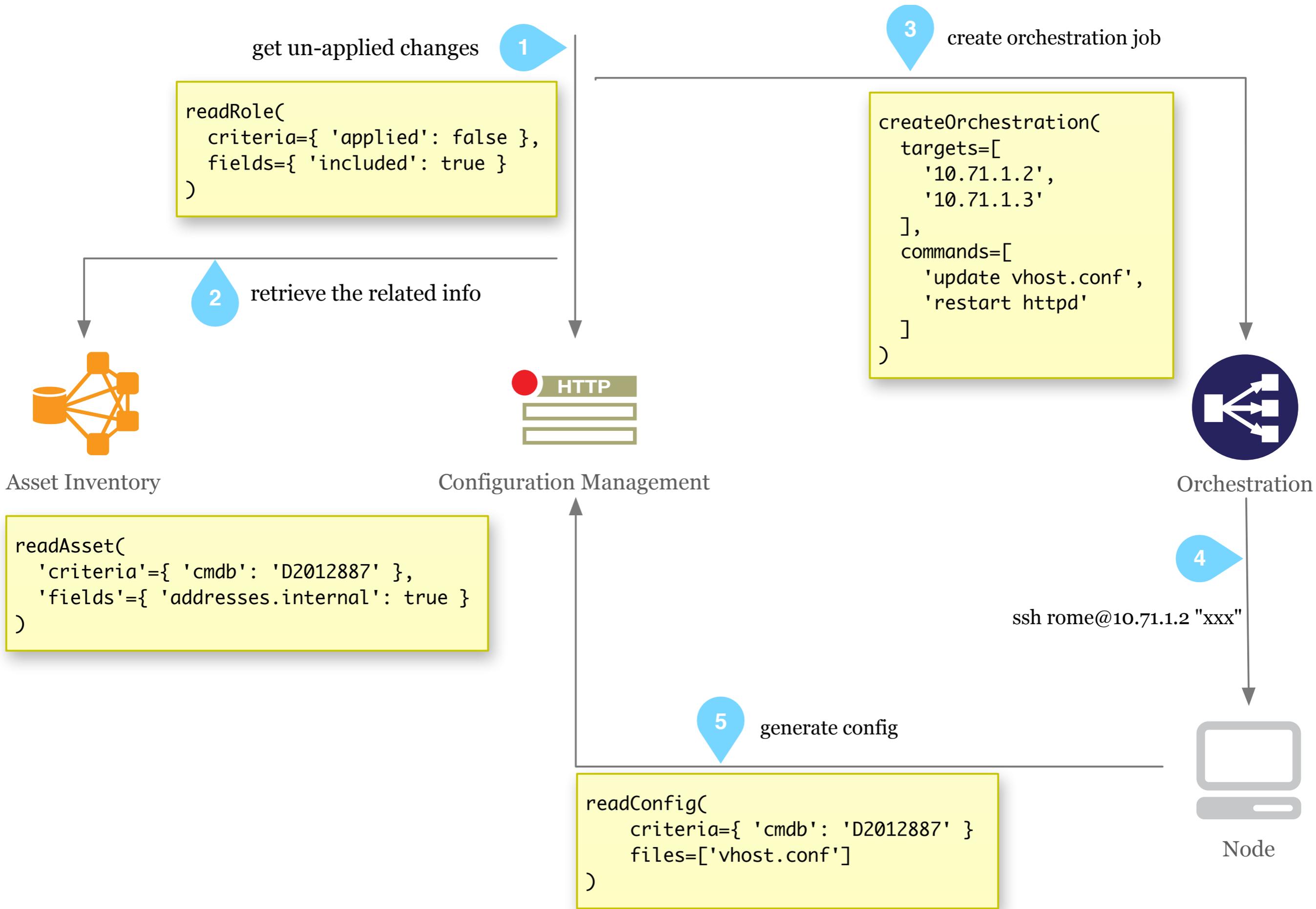
Configuration Management

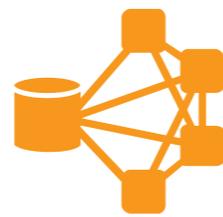

CACHE
Caching


Logging

More ...


Mobile App





Asset Inventory

2 get the related info

```
readAsset(
  'criteria'={ 'cmdb': 'D20120601x' },
  'fields'={ 'addresses.internal': true }
)
```

1 some update

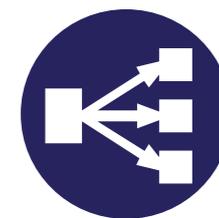
```
createRole(
  name='D20120601x',
  includes=['mysql', 'xd'],
  ports=[3310, 3324]
)
```



Configuration Management

3 create orchestration job

```
createOrchestration(
  target='10.71.1.2',
  commands=[
    'update all',
    'initialize'
  ]
)
```



Orchestration

4

ssh rome@10.71.1.2 "xxx"



Node

5 generate config

```
readConfig(
  criteria={ 'cmdb': 'D20120601x' }
)
```

我们用了什么技术？

- HTTP-based RESTful API with JSON payload
- A variety of MongoDB-like CRUD
- Node.JS
- Connect & Express
- MongoDB

我们如何开发？

- Ops->DevOps
- Scrum
- BDD
- CI
- 工具大部分来自于Atlassian(比如Jira)

谢谢

Q&A

