

目录

- 为什么自主开发?
- 我们的系统
- 架构设计
- 特点
- 未来计划

目前开源监控系统存在的问题

- 只有“监”的功能，没有“控”的功能
- 当出现故障时，往往会产生“报警风暴”
- 监控数据都是孤立的，无法产生更多的价值
- 缺少根源分析
- 效率低(性能、配置的工作量)
- ○ ○ ○ ○ ○ ○

- 在50万-100万台设备规模的环境下，上面这些问题尤其突出。我们需要在一个平台中能够支持这么多设备的监和控。

目标

- **眼:**
 - 使得该监控系统能够很好地成为我们的“眼”，帮我们7*24小时监控我们整个系统
- **脑:**
 - 使得该系统成为我们的“脑”，能够根据历史的经验和知识库给我们诊断的建议和历史处理该问题的经验，方便我们快速的诊断和解决问题。
 - 为我们提供问题根源分析。
- **手:**
 - 使得该系统成为我们的“手”，方便我们快速地进行操作
 - 将运维的方式由黑屏变白屏，成为我们的运维操作平台。

目录

- 为什么自主开发?
- 我们的系统
- 架构设计
- 特点
- 未来计划

- ▼ 监控分组
 - ▶ AT-HAII
 - ▶ AY-APP
 - ▶ AY10
 - ▶ DA
 - ▶ 其他
 - ▶ 监控
 - ▶ AT-Search
 - ▶ AT-Search-Perf

设备汇总 设备状态 设备性能

AT04 TOTAL[301] ● UP[301] ● DOWN[0]

状态统计: ● OK[2701] ● WARNING[0] ● CRITICAL[7] ● ERROR[0] ● PENDING[2] ● MISS[0]

主机名	OK	WARNING	CRITICAL	ERROR
r02j11027.yh.aliyun.com	9	0	0	0
r02j11028.yh.aliyun.com	9	0	0	0
r02j11029.yh.aliyun.com	8	0	1	0
r02j11030.yh.aliyun.com	9	0	0	0
r02j11031.yh.aliyun.com	9	0	0	0
r02j11032.yh.aliyun.com	9	0	0	0
r02j11033.yh.aliyun.com	9	0	0	0

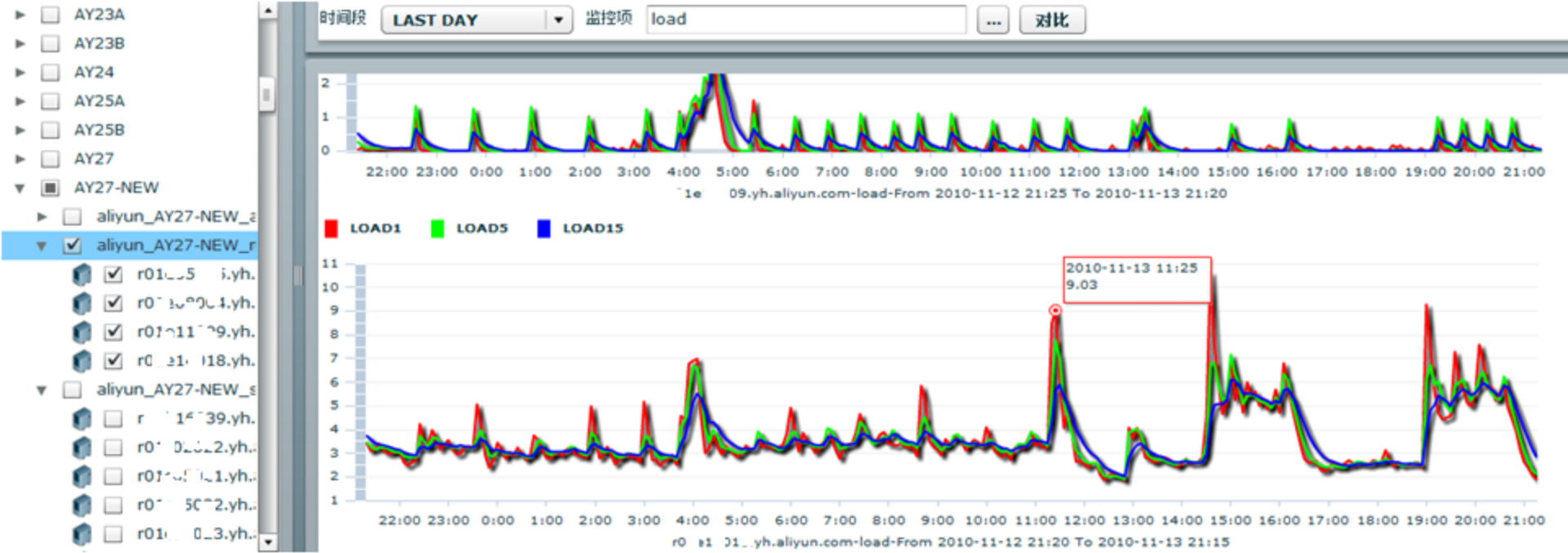
- ▼ 监控分组
 - ▼ AT-HAII
 - ▶ aliyun_AT-HAII-server
 - ▶ AT-ORDER
 - ▶ AT-Search
 - ▶ AT-SEARCH-DATAPLAT
 - ▶ AT-Search-Perf
 - ▶ AT-SQLUSA
 - ▶ AT-staging
 - ▶ AT-yunti2
 - ▶ AT01
 - ▶ AT03
 - ▶ AT04
 - ▶ AY03
 - ▶ AY03-new
 - ▶ AY04-NFW
 - ▶ AY05A
 - ▶ AY05B

设备汇总 设备状态 设备性能

主机名: IP地址: 端口: 报警人:

状态: 报警: 电话: 节点排序:

监控项	状态	状态信息	持续时间	预警次数	定时	报警	电话	最后检测时间
▼ r02j11027.yh.aliyun	UP	设备组: aliyun_AT-HAII-server; IP: 10.						
check_zealot	OK	OK - zealot is running	10d 21h 21m 27s	0/1				10-11-13 00:10:19
check_ethstatus	OK	OK - All eth card is working properly	5d 7h 8m 42s	0/2				10-11-13 21:21:39
check_logwatch	OK	OK-NOT NEED UPGRADE 0.89 - 0.89	23d 2h 16m 32s	0/1				10-10-21 19:17:57
check_ssh	OK	OK - 10.249.69.27 port 22 is ok	11d 8h 47m 5s	0/2				10-11-13 21:32:53
check_server_alive	OK	OK - PING: rtt<1ms	26d 0h 11m 33s	0/2				10-11-13 21:34:17
check_ntp	OK	OK-NTP: offset:0.001559	26d 3h 7m 8s	0/2				10-11-13 21:32:53
check_upgrade		姓名	旺旺	手机	分机	邮箱		12m 38s 0/2
check_hardware								8m 43s 0/2
check_filesystem								21m 6s 0/2
▼ r02j11024.yh.aliyun								m 52s 0/1
check_zealot								10-11-13 17:34:37

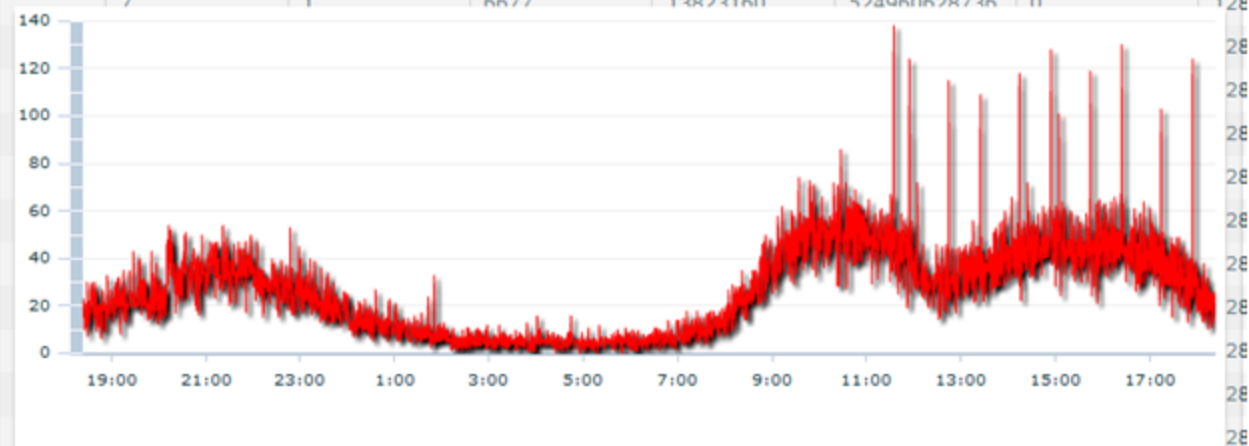


Dashboard 集群状态 集群性能

- AT-Search
- AT-Search-Perf
- AT-staging
- AT-yunti2
- AT01
- AT03
- AT04
- AY03
- AY04-NEW
- AY05A
- AY05B
- AY08
- AY08-NEW
- AY08A
- AY08B
- AY08C
- AY08D

Hosts Up: 41 Hosts Down: 0 CPUs Total: 592 Avg Load(1,5,15): 0.33027,0.302432,0.294595 Host: 搜索 刷新 每页记

Process St...	Pl_g Master	Pan uChun...	K Master	K Server	K fileMeta	SC Master	!_ Worke...	_C Worke...
Host	PartCount	ReadCount	WriteCount	touchaoFileC	CellCount	KVDataSize	KVStartTime	KV...
r031011118.yh.aliy	48	21	1	4974	11132375	388896391168	0	128
r031011115.yh.ali	48	7	1	6677	13823160	524960628736	0	128
r031011104.yh.ali	48							28
r031011103.yh.ali	47							28
r031011111.yh.ali	48							28
r031011102.yh.ali	48							28
r031011103.yh.ali	48							28
r031011107.yh.ali	48							28
r031011103.yh.ali	48							28
r031011100.yh.ali	47							28
r031011100.yh.ali	47							28
r031011101.yh.ali	48							28



ROOM : 1		ROOM : 2		ROOM : 3		ROOM : 4		ROOM : 5	
ROW : a	ROW : b	ROW : c	ROW : d	ROW : e	ROW : f				
0.. msw-1a08-1...	0.. msw-1b08-1...	0.. msw-1c08-1...	0.. msw-1d08-1...	0.. msw-1e08-1...	0.. msw-1f08-11...				
u.. asw-1a08-s5...	u.. asw-1b08-s5...	u.. asw-1c08-s5...	u.. asw-1d08-s5...	u.. asw-1e08-s5...	u.. asw-1f08-s5...				
5.. r01-00021.y...	5.. r01-00021.y...	1.. r01-00017.y...							
6.. r01-00022.y...	6.. r01-00022.y...	2.. r01-00018.y...							
7.. r01-00023.y...	7.. r01-00023.y...	3.. r01-00019.y...							
8.. r01-00024.y...	8.. r01-00024.y...	4.. r01-00020.y...							
9.. r01-00025.y...	9.. r01-00025.y...	5.. r01-00021.y...							
1.. r01-00026.y...	1.. r01-00026.y...	6.. r01-00022.y...							
1.. r01-00027.y...	1.. r01-00027.y...	7.. r01-00023.y...							
1.. r01-00028.y...	1.. r01-00028.y...	8.. r01-00024.y...							
1.. r01-00029.y...	1.. r01-00029.y...	9.. r01-00025.y...							
1.. r01-00030.y...	1.. r01-00030.y...	1.. r01-00026.y...							
1.. r01a-00031.y...	1.. r01b-00031.y...	1.. r01-00027.y...							
1.. r01a-00032.y...	1.. r01b-00032.y...	1.. r01c-00028.y...							
		1.. r01-00028.y...							
		1.. r01-00029.y...							
		1.. r01-00030.y...							

设备状态 设备性能

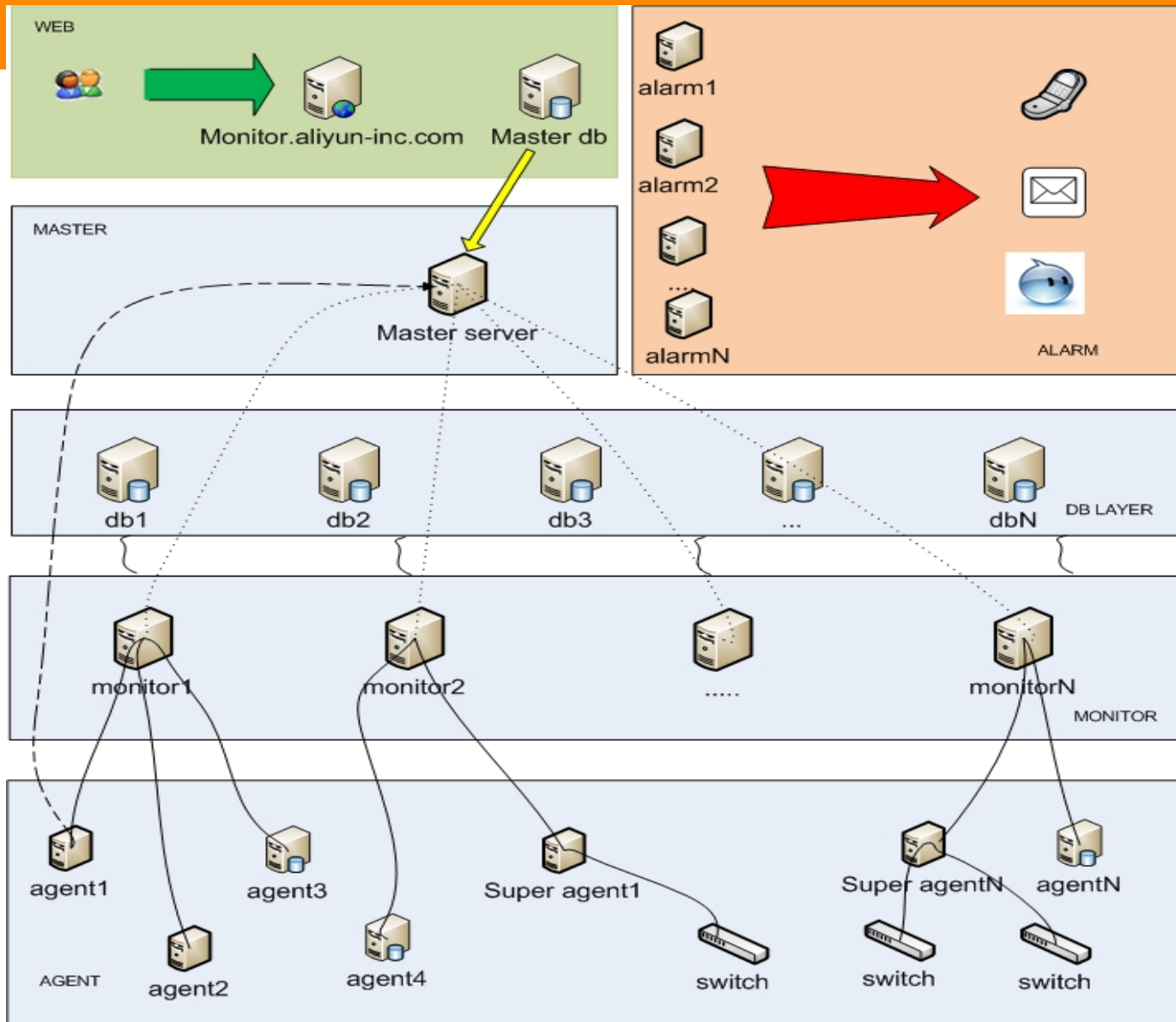
监控项 状态 **ALL**

监控项	状态	状态信息	持续时间	预警次数	报警	电话
check_zealot	OK	OK - zealot is running	24d 18h 55m 25s	0/1		
check_ethstatus	OK	OK - All eth card is v	3d 23h 35m 26s	0/2		
check_upgrade	OK	OK - agent_upgrade	3d 23h 35m 47s	0/2		
check_logwatch	Critical	Critical- mailmonitor	3h 26m 24s	356464/1		
check_server_a	OK	OK - PING: rtt<1ms	8d 21h 38m 13s	0/2		
check_hardware	OK	OK - Health Check p	3d 23h 41m 46s	0/2		
check_ntp	OK	OK-NTP: offset 0.00	24d 19h 33m 55s	0/2		
check_ssh	OK	OK - 10.249.14.28 p	24d 19h 33m 55s	0/2		
check_filesystem	OK	OK - all filesystem fi	3d 23h 48m 12s	0/2		

目录

- 为什么自主开发?
- 我们的系统
- 架构设计
- 特点
- 未来计划

鹰眼架构



整体的架构设计

- master 负责monitor, agent 注册并监控Master DB表的变化并将这个变化传播给monitor。
- monitor 分为多个进程，利用共享内存交互，负责接收agent传递上来的数据，根据状态进行报警，并把异常事件、当前状态、性能数据入库。
- agent 主要负责调度本地的插件，并把插件执行的结果传递给monitor。
- db layer 存放所有的元数据和监控数据，并保证agent 切换到其他monitor 之后数据还是连续的。Agent与DB是多对一的关系，并保持不变。
- alarm 主要是从master 的异常表抢占式的取出异常数据，并报警

监控设计的关键点

- 去单点，可扩展，自我监控
- Master一组，热备，自动切换。
- Monitor 单台，出现了问题，可以快速迁移agent到另外的monitor。
- Master 能够发现monitor 机器上任意进程的问题，通过心跳来实现。发现5次心跳错误就报警。
- Super agent单台，出现了问题，可以快速迁移agent到另外的super agent。
- monitor 会发现super agent 的问题，报警。
- DB采用自主研发的HA Aurora做热备。如果lister进程发现db有问题，数据先记录到文件。
- Monitor 发现 master的故障，报警，继续自己与agent之间的工作，同时启线程连master。受影响的只有新连接的agent

目录

- 为什么自主开发?
- 我们的系统
- 架构设计
- 特点
- 未来计划

监控范围广

- 监控面广泛，从数据中心最低层到最高层应用
 - 机房环境，如电力，空调,能耗，湿度
 - 网络设备
 - 服务器
 - 重要基础设施，DNS，VIP，DB等等
 - 应用监控
 - 集群监控，如集群的各种性能指标，集群的监控状态

高性能及良好的扩展性

- 监控效率非常高，每台监控机可以监控20万监控点，采用主动+被动
- 集成了轻量级的icmp, ntp, snmp, ipmi协议，避免fork进程
- 采用epoll+半异步/半同步的线程池
- 分布式架构，支持100万台设备的规模

监控自动部署

- 定义监控模板，监控模板包括监控项，报警策略，报警人
- 与CMDB打通，所有被监控设备以及设备的组织和关系均来自于CMDB。
- 将设备组与监控模板关联。n:n的关系。所有在该设备组内的机器都具有监控模板里面的监控项。用户只需要操作CMDB即可。
- 允许增加个性化的监控项(不在监控模板中)

根源分析

- 以service监控入手，将与某个service相关联的各种设备和监控项层层建立服务依赖（目前人工建立这种依赖关系，这种服务依赖不是nagios中报警依赖）。
- 将这种依赖关系以拓朴的形式展现在一张图上。
- 当某个服务出现问题时，通过这种依赖，快速的定位到相关联的监控项或者设备组。

控

- 重启物理机
- 增加安全清洗策略
- 重启VM
- 重启云计算OS
- 执行各种事先好的命令

提高报警的准确性

- 合并报警，抑制报警风暴
 - 单位时间内某个报警接收人的报警超过6条时，程序会采取报警合并，将设备组、监控项、状态相同的报警合并成一条报给用户
 - 合并后的旺旺、短信将告知有多少条相似的报警被合并，而邮件中将列出这些报警的详细信息。
- 关联报警
 - 定义一些规则，实现关联报警，比如当load超过某个阈值时，所有监控项的timeout不要报警，但异常信息还是要报。

采用数据库替换RRD文件

- RRD文件具有两个缺点：
 - 一张性能图一个小文件，文件数多
 - 小IO，对IO性能要求较高
- 我们采用JAVA实现RRD算法，用MySQL数据库来代替RRD的存储，大大提升性能，数据得以集中管理，也很容易多张RRD合并成一张

文件形式接口

- 一个开放的接口。它的工作原理是：
鹰眼Agent进程会实时检测一个固定文件的行增量变化，并把增量消息发送给前端monitor。
- 这个日志文件每一行内容的格式如下：
 - Critical-xxxx
 - Warning-xxxx
 - Ok-xxxx
- 如果行头是“Warning -”，“Critical -”的消息，则agent会将这些消息报给logwatch指定的报警组。

报警处理case库和历史报警

- 可以查看每一个监控项的历史报警和历史操作
- 可以查看每一个监控项的历史报警的原因、处理方法，形成case库，帮助处理人员快速定位问题和解决问题

鹰眼Agent升级和插件升级

- 鹰眼Agent可以自动升级。
- 插件也可以自动升级。当插件变更，或者增加新的插件时，通过插件升级机制可以将插件推到所有的机器上面。

监控数据分析

- 由于元数据和异常数据的集中、状态数据和性能数据库化存储，使得监控数据分析非常方便
- 报警分析、异常数据按集群分析、每个监控项的可用率的自动计算、集群平均性能指标等等
-

目录

- 为什么自主开发?
- 我们的系统
- 架构设计
- 特点
- 未来计划

未来计划

- 应用接口级别监控
- 网络设备监控的自动部署
- 取代crontab
- 配置管理的Agent，取代cfengine
-

Q/A?

