裸金属主机监控架构设计文档

Review 记录信息

时间	参与人员	主要议题	后续review计划

问题描述 ≥

目前不支持裸金属主机监控功能,用户只能登录裸金属主机操作系统查看CPU,内存,磁盘,网络等信息,不方便进行裸金属主机的监控和运维。

方案提议♂

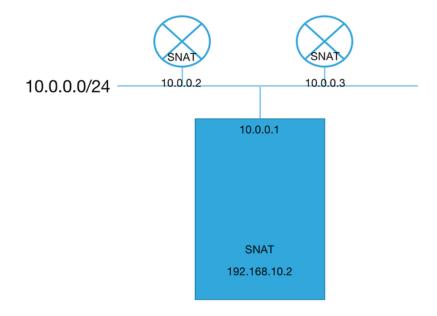
总体方案设计 🔗

裸金属云产品部署时在云产品节点部署pushgateway服务;配置裸机通过业务网络访问pushgateway服务;创建裸机时选择已安装监控agent的镜像,或勾选安装监控agent(创建裸机时通过userdata注入安装监控agent脚本,在裸机运行时自动从pushgateway服务下载和安装监控agent),或在已创建的裸机中安装和运行监控agent;监控agent获取裸金属主机CPU,内存,磁盘,网络等监控数据,并上报给pushgateway服务;云监控服务从pushgateway服务获取并保存监控数据供用户查询。

网络方案设计 🔗

裸金属主机监控agent上报监控数据到pushgateway服务,需要从裸机业务网络访问pushgateway服务的网络。考虑几种网络方案:

1、ovn实现分布式SNAT路由,完成业务网络到管理网络的路由转发,如图所示:

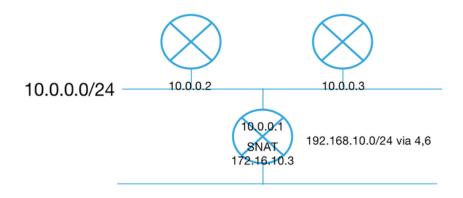


优点:

用户只需将裸机业务子网绑定到用户路由器,操作简单

问题:

- (1) 产品上在用户透明的条件下将业务网络与管理网络打通,存在安全合规风险
- (2) 方案实现需要网络产品开发工作量,发布时间长
- (3) 依赖路由器,只能在geneve网络场景下使用
- 2、通过一系列网络配置完成业务网络到外部网络,再到管理网络的路由转发,如图所示:



172.16.10.4

172.16.10.5

SNAT

192.168.10.2

网络配置步骤如下:

- (1) 用户创建用户路由器,设置网关,并开启SNAT
- (2) 用户将业务子网绑定到用户路由器
- (3) 管理员在用户路由器上添加静态路由,下一跳是网络节点br-pub的ip,目的CIDR是管理网络网段
- (4) 管理员在网络节点上配置br-pub地址到管理网络网段的SNAT转发规则

优点:

实现方案不需要产品开发

问题:

- (1) 网络配置涉及普通用户和管理员,操作复杂
- (2) 用户路由器静态路由配置对普通用户可见,存在安全风险
- (3) 依赖路由器,只能在geneve网络场景下使用
- 3、给pushgateway服务添加ingress,需要做的事情如下:
- (1) 裸机业务子网绑定用户路由器,使裸机可以访问外部网络
- (2) 客户防火墙打通外部网络和控制台网络(小规模项目一般外部网络和控制台网络是一个网段,可以直接通;邮储这种项目是通过客户防火墙打通外部网络和控制台网络)

优点:

模型成熟,用户操作简单

问题:

- (1) 依赖路由器,只能在geneve网络场景下使用
- (2) 上报监控数据是否会给nginx产生大的压力,初步评估,单台裸机一次推送数据大概在200KB左右,按单集群1000台裸机来算,集群一次推送大概在200MB左右,实际网络传输时会有数据压缩,200KB数据传输带宽在10KB左右,从而集群一次推送带宽在10MB左右,压力不大。
- 4、给pushgateway服务加LB svc,需要做的事情如下:
- (1) LB通过外部网络FIP访问
- (2) 裸机业务子网绑定用户路由器,使裸机可以访问外部网络

优点:

用户操作简单

问题:

- (1) 产品服务上未使用过LB svc,且产品上容器的LB service还在开发过程中,只在项目的客户应用上用过
- (2) 需要支持路由器,只能在geneve网络场景下使用
- 5、实现一个类似于nova metadata service的服务,该服务完成内部网段加端口到管理网段加端口代理,裸机向该服务的内部网段和端口推送数据,该服务转发数据到pushgateway服务,需要做的事情如下:

产品上需要实现这样一个代理服务,能将内部网段加端口代理到代理服务管理网ip加端口,然后可以通过管理网络访问pushgateway服务

优点:

- (1) 用户无感知
- (2) vlan和geneve都可以支持

问题:

需要站在平台视角开发一个统一的代理服务,开发工作量较大

- 6、给pushgateway加NodePort svc,需要做的事情如下:
- (1) 裸机业务子网绑定用户路由器,使裸机可以访问外部网络
- (2) 客户防火墙打通外部网络和控制台网络(小规模项目一般外部网络和控制台网络是一个网段,可以直接通;邮储这种项目是通过客户防火墙打通外部网络和控制台网络)

优点:

- 1、模型成熟,用户操作简单
- 2、通过node port端口访问,可进行访问控制,安全性较好
- 3、避免了ingress方式可能对平台nginx产生压力的问题

问题:

依赖路由器,只能在geneve网络场景下使用

综合上述分析,并为解决pushgateway多副本数据一致性问题,选用给pushgateway服务添加ingress的方案。

其它可选方案 ≥

方案	方案描述	代价	代表厂商
----	------	----	------

ipmi/redfish	通过ipmi/redfish接口访问裸 金属主机,查询CPU,内 存,磁盘和网络信息	1、不能准确收集 CPU,内存,磁盘和网 络使用量信息。	无
		2、有服务器型号限制,有些服务器不能查到相应信息。	

竞品类似方案 ≥

华为云



只支持安装agent的操作系统内部监控。

监控指标:CPU使用率,CPU 5分钟平均负载,内存使用率,磁盘剩余存储量,磁盘使用率,磁盘I/O使用率,inode已使用占比,入网带宽,出网带宽,网卡包接收速率,网卡包发送速率,所有状态的TCP连接数总和,处于ESTABLISHED状态的TCP连接数量

ZStack ∂

▲ 定时运维 - 平台运维 - ZStack用户手册 - 产品手册 - ZStack Cloud 4.7.21 - ZStack

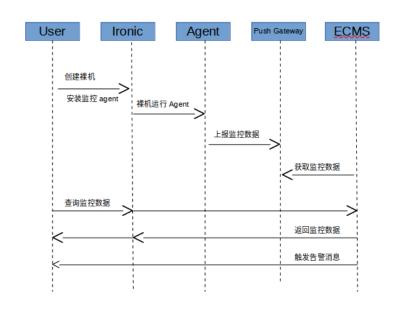
支持提供CPU、内存、磁盘、网络的监控和报警功能。

agent是安装在弹性裸金属实例内的代理,定时与管理节点通信。弹性裸金属实例需预先安装agent,才能获取硬件信息、查看内部监控数据、打开控制台、修改密码、加/卸载云盘、以及加/卸载网络。

监控指标:磁盘读IOPS,磁盘写IOPS,磁盘已用量百分比,磁盘用量,网卡入包速率,网卡出包速率,网卡入错误率,网卡出错误率,CPU平均使用率,内存使用率,磁盘读速度,磁盘写速度,网卡入速度,网卡出速度

重要图表 🔗

流程图:



可升级影响 🔗

可随云产品进行升级,存量裸机需要安装并运行监控agent才能获取到监控和告警信息。

稳定可靠性影响 🖉

pushgateway服务默认3副本部署,可用license最多设置9副本,每个副本承载部分裸机监控数据。由于监控服务会从所有pushgateway服务副本获取监控数据,所以单副本故障只会影响部分裸机监控数据获取。

性能影响 ②

1、监控agent运行,收集和发送监控数据会占用裸金属主机一部分CPU,内存和网络带宽资源。

监控agent的CPU占用限制在5%,内存占用限制在50M,单台裸机一次推送数据占用带宽在10KB左右

2、大规模裸机监控对裸机业务网络性能影响,对pushgateway服务压力影响。

pushgateway服务采用默认3副本部署,由ingress提供负载均衡,每个副本承载部分裸机监控数据,同时提供定时数据清理机制

安装部署影响 🔗

无影响。

API影响 🔗

新增查询裸机监控信息的Django API

修改创建裸机的Django API

修改重建裸机的Django API

前端界面影响 🖉

新增查询裸机监控按钮和监控信息界面

新增裸机监控概览页面

裸机告警信息前端适配

安全性影响 🔗

裸机业务网络通过控制台网络访问pushgateway服务的安全性:

pushgateway服务下载agent和上报监控数据增加用户名密码访问认证

文档影响 🔗

用户手册,API文档,技术白皮书,运维手册内容补充

监控 日志 告警方案 ≥

不涉及

License 方案 &

新增pushgateway服务的license控制

其它影响 🔗

无

实现♂

主要实现成员 🔗

于尚斌,潘广超,王壮年

JIRA 任务 ID 🔗

□ EAS-78591: 裸金属主机监控功能 已完成

测试 🖉

功能测试

查询裸机监控信息功能测试

查询裸机告警信息功能测试

性能测试

- 1、获取不同时间段监控数据的性能测试
- 2、大规模裸机监控的性能测试

压力测试 ⊘

大规模裸机监控对服务的压力测试

场景测试

正常场景:裸机监控agent正常安装和运行,查询监控和告警信息成功,关闭/卸载监控agent,查询监控和告警信息无数据,开启监控agent,

查询监控和告警信息有数据

异常场景:裸机监控agent异常停止运行,监控服务(含pushgateway服务)异常,监控网络异常

特殊场景:裸金属主机关机、重启、重建等操作

longrun测试

裸机监控稳定可靠性测试

组件间依赖 ♂

依赖组件:

ECMS 6.2.1

参考链接 ፟

待讨论确认的内容 ♂

修订记录 ≥

当前版本	主要修改人	主要签署人	修改时间	描述
v0.1	于尚斌		2022.6.16	初稿
v0.2	于尚斌		2023.4.10	修改稿
v0.3	于尚斌		2023.5.4	修改方案描述和流 程图
v1.0	于尚斌		2024.1.3	终稿