

**SACC** 2014中国系统架构师大会  
SYSTEM ARCHITECT CONFERENCE CHINA 2014

发现架构之美

# 网易私有云网络虚拟化实践

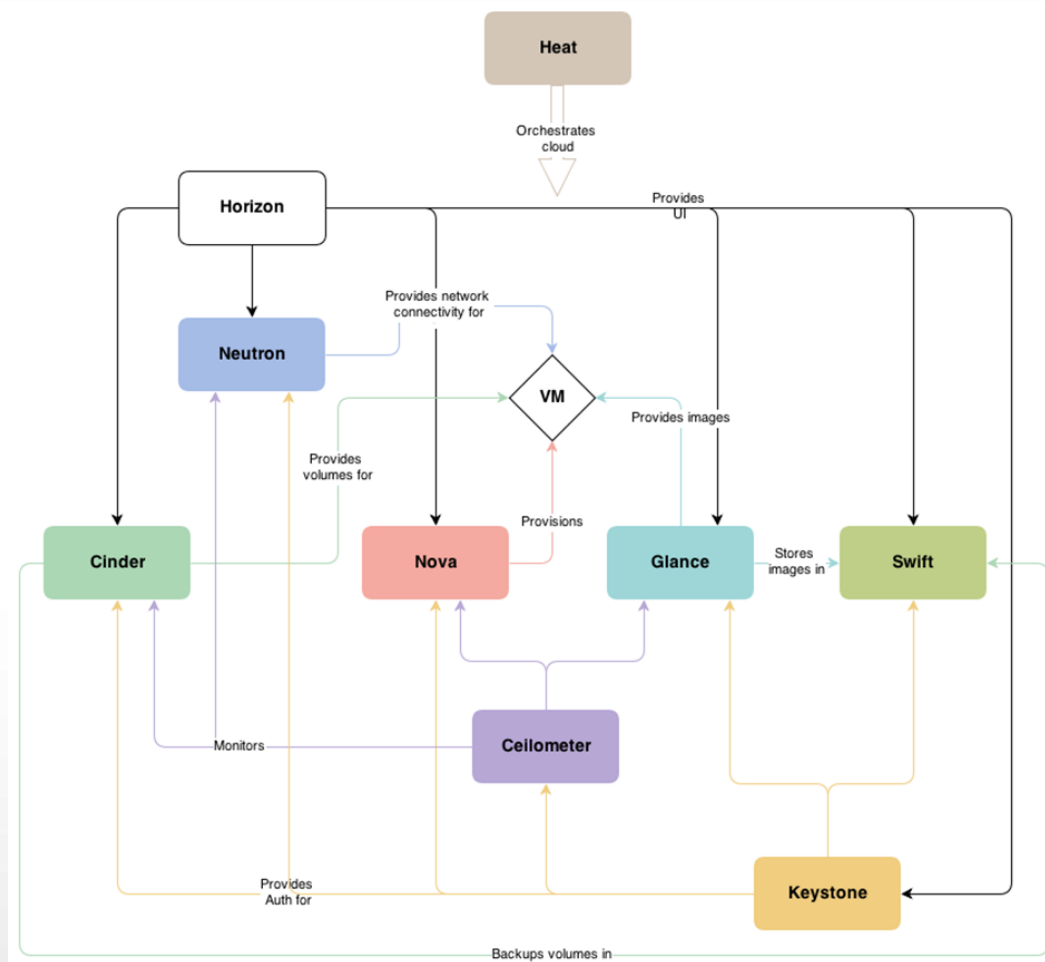
网易杭州研究院  
徐城利

# 提纲

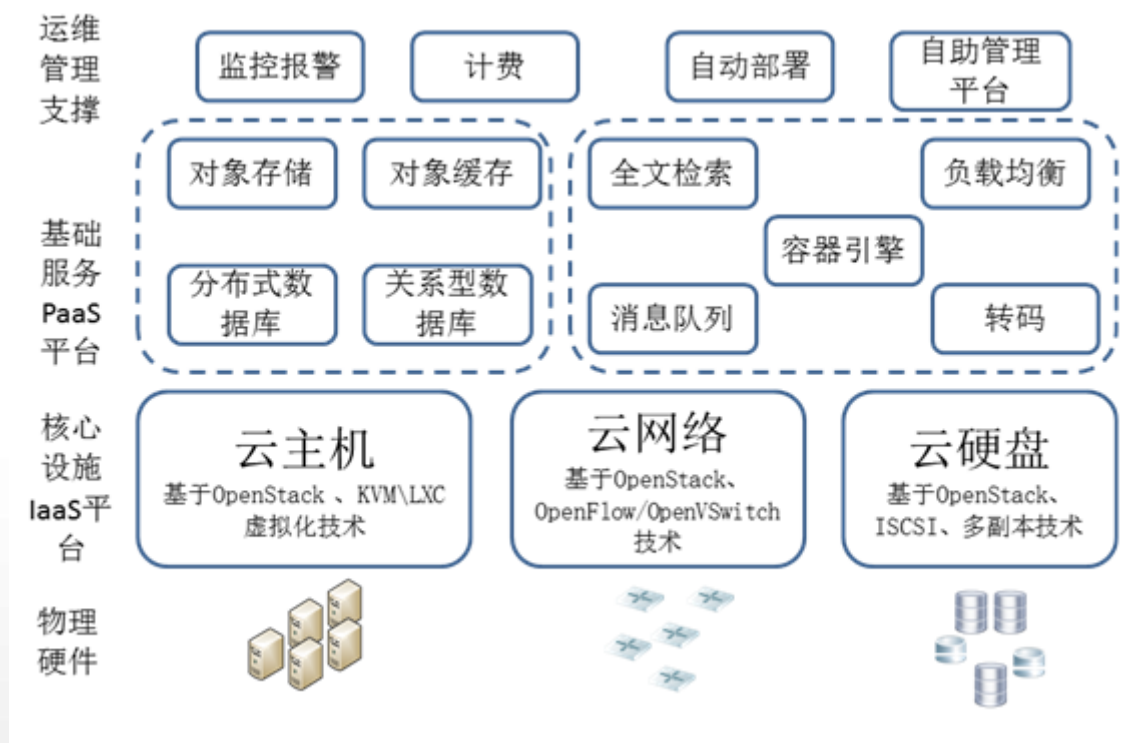
- OpenStack介绍
- 网易私有云简介
- 第一代网络虚拟化服务
- 第二代网络虚拟化服务
- 碰到的问题及一些经验总结

# OpenStack

- 计算 (nova)
- 网络 (neutron)
- 认证 (keystone)
- 镜像 (glance)
- 块存储 (cinder)
- 对象存储 (swift)
- 其他模块



# 网易私有云



# 网易私有云(Cont'd)

- 12年年底正式上线，已稳定运行近2年
- 部署两个Region
- 100+物理节点
- 已广泛服务于网易公司产品（30+）如：门户、易信、网易新闻客户端、云音乐、云课堂/公开课、云阅读等



# 网易私有云(Cont'd)

- 提高硬件资源使用率  
如CPU: ~10% -> ~50%
- 提高物理资源管理及运维自动化
  - 自助服务
  - 运维人员减少1/2
- 提高基础资源弹性
  - 资源池: 快速申请, 按需使用
  - 适应业务波动
  - 满足开发、测试等零碎需求



# IaaS平台研发

- 三个基础服务
  - 云主机：提供可扩展、安全可靠的弹性计算
  - 云网络：提供动态、安全、灵活的网络服务
  - 云硬盘：提供可扩展、安全稳定的块存储服务
- 研发工作
  - 基于OpenStack的nova/neutron/keystone/glance/cinder开发
  - 对OpenStack作充分的功能、性能、稳定性及异常测试
  - 发现并修复OpenStack社区的bug (提交90+、修复50+)
  - 根据公司需求，研发新功能及优化30+
  - 支持整合网易私有云10多个上层服务

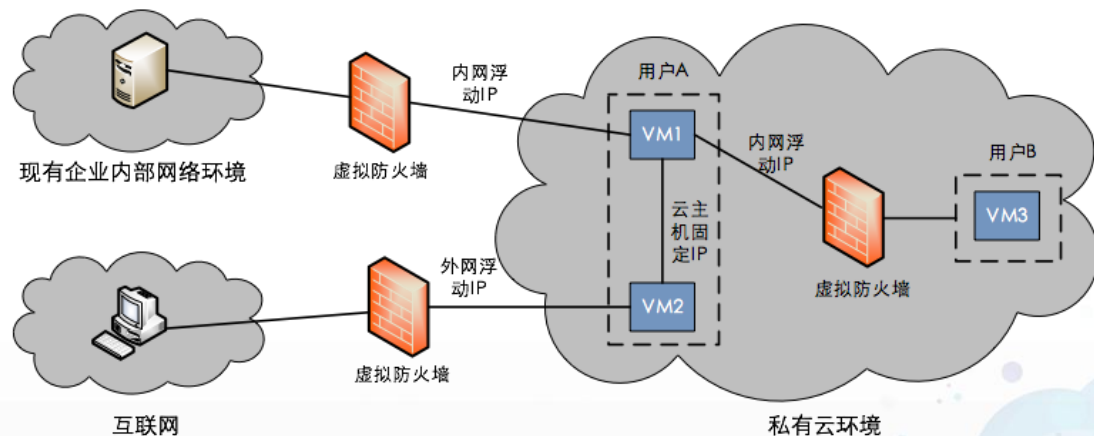
# 云网络演进

- 第一代网络虚拟化服务
  - 基于nova-network
  - 使用Flat DHCP Network Manager
  - 定制开发
- 第二代网络虚拟化服务
  - 基于neutron
  - 使用ML2 + Open vSwitch Agnt
  - 应用SDN (Software-Defined Networking)技术
  - 定制开发



# 第一代云网络

- 基于nova-network，简单可靠
- 固定IP (Fixed IP)、浮动IP (Floating IP)
- 安全组 (Security Groups)
- 租户网络隔离
- QoS



# 第一代云网络(Cont'd)

- 技术选型
  - Flat DHCP Network Manager
  - Multi-host
- 自研功能及优化
  - 内网浮动IP
  - 租户网络优化与隔离
  - 网络QoS

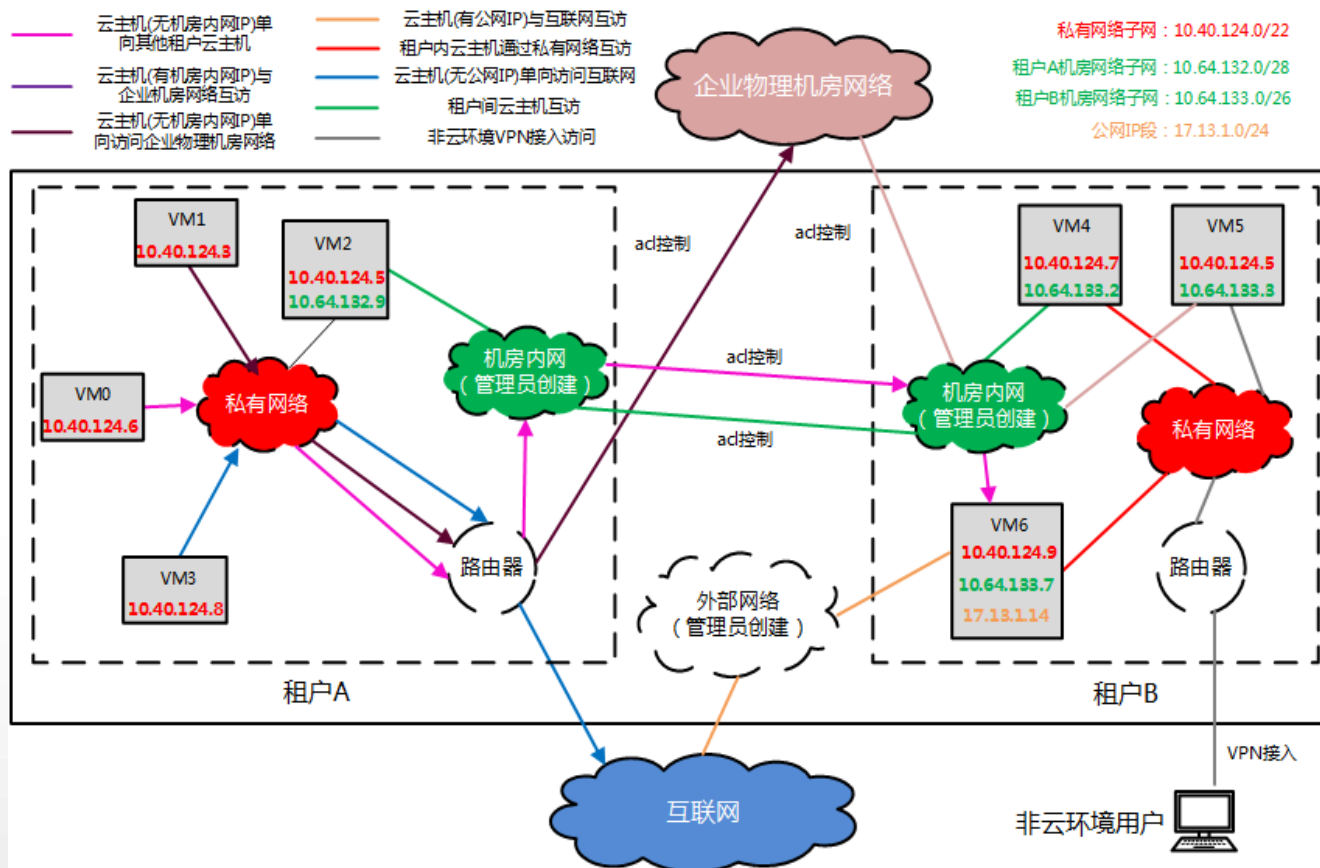
# 不足

- 所有云主机位于同个物理L2
  - 广播域扩大
  - 网络边界
  - IP资源管理
- 云主机和宿主机部分重叠，管理成本高
  - NAT、安全组等数千条iptables规则
  - 路由干扰
- 用户使用习惯上不同
- 技术上扩展性不够

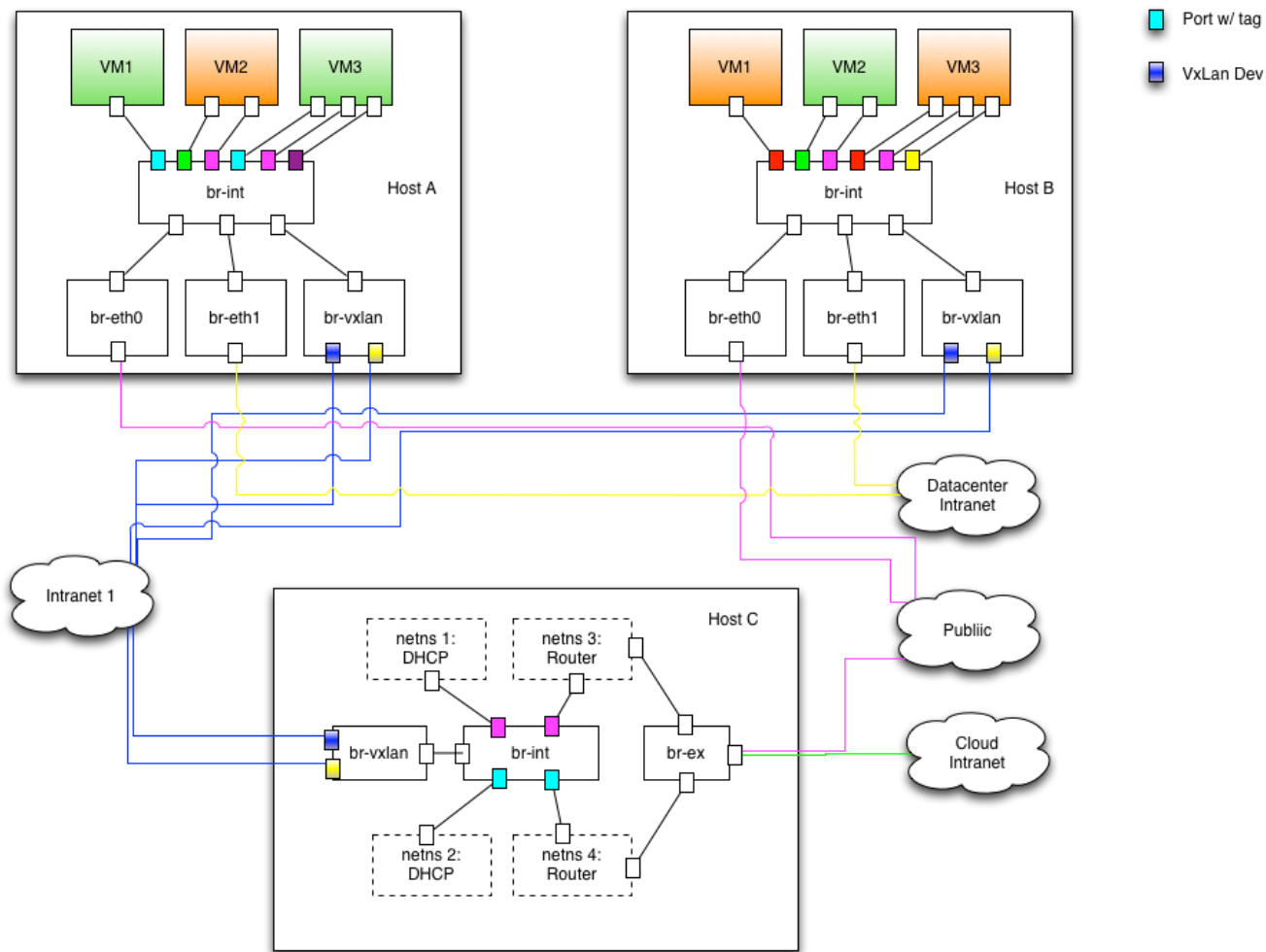
# 第二代云网络

- 基于neutron开发
- 设计上更加合理
  - 资源抽象更彻底
  - 可扩展性高，支持各种网络技术
- 灵活性好，易于实现不同的网络拓扑结构
- 根据业务需求定制
  - 私有网络
  - 机房内网
  - 公网

# 第二代云网络架构



# 网络拓扑示意图

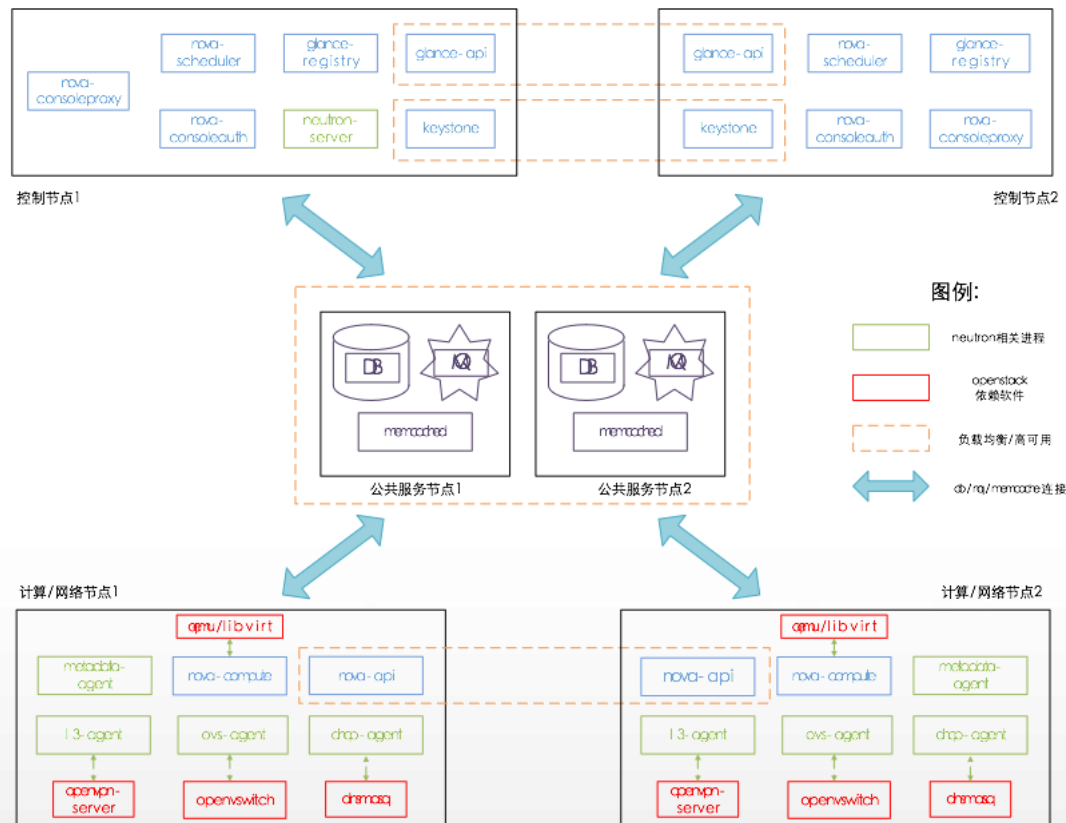




# 第二代云网络(Cont'd)

- 技术选型
  - VXLAN (Virtual Extensible LAN)
  - Flat网络
  - ML2 Plugin + Open vSwitch Agent
  - L2 Population
- 自研功能
  - 新网络类型SVLAN (Shared VLAN)
  - 支持ACL
  - L3使用两次NAT
  - 网络监控
  - 租户OpenVPN接入
  - L3高可用

# 服务部署



# 问题及总结

- Neutron还在不断的改进中，建议使用新版本，并关注最新功能与bug修复
- 基于neutron的网络服务稳定性上仍有不足，需要完整、严格的测试
- L2 Population在实际使用中容易出问题且后果严重，功能上需要增强
- OVS Agent重启导致网络中断
- 3.9以前内核及低版本iproute容易导致namespace锁死
- 性能上相对原来有所下降，有进一步优化的余地
- 网络问题相对不容易定位和排查，可以定制一些自动化工具
- ...

# 后续工作

- 继续提高稳定性
- 功能上的扩展
  - 安全组
  - QoS
  - L3的SNAT支持DVR
  - ...
- 性能优化
  - 软件上的调优
  - 尝试一些硬件方案 (OpenFlow交换机等)
- 加强和社区协作

# 联系方式

- [xuchengli@corp.netease.com](mailto:xuchengli@corp.netease.com)



徐城利

superekcah



扫描二维码, 立即加我为易信好友



superekcah

Hangzhou, Zhejiang



Scan the QR code above to add me on WeChat

# Q&A

# THANKS

SequeMedia  
盛拓传媒

IT168.com  
www.it168.com

ChinaUnix

ITPUB