

阿里云 ZStack for Alibaba Cloud

超融合产品 使用教程

产品版本 : V3.0.0

文档版本 : 20180917

法律声明

阿里云提醒您在阅读或使用本文档之前仔细阅读、充分理解本法律声明各条款的内容。如果您阅读或使用本文档，您的阅读或使用行为将被视为对本声明全部内容的认可。

1. 您应当通过阿里云网站或阿里云提供的其他授权通道下载、获取本文档，且仅能用于自身的合法合规的业务活动。本文档的内容视为阿里云的保密信息，您应当严格遵守保密义务；未经阿里云事先书面同意，您不得向任何第三方披露本手册内容或提供给任何第三方使用。
2. 未经阿里云事先书面许可，任何单位、公司或个人不得擅自摘抄、翻译、复制本文档内容的部分或全部，不得以任何方式或途径进行传播和宣传。
3. 由于产品版本升级、调整或其他原因，本文档内容有可能变更。阿里云保留在没有任何通知或者提示下对本文档的内容进行修改的权利，并在阿里云授权通道中不时发布更新后的用户文档。您应当实时关注用户文档的版本变更并通过阿里云授权渠道下载、获取最新版的用户文档。
4. 本文档仅作为用户使用阿里云产品及服务的参考性指引，阿里云以产品及服务的“现状”、“有缺陷”和“当前功能”的状态提供本文档。阿里云在现有技术的基础上尽最大努力提供相应的介绍及操作指引，但阿里云在此明确声明对本文档内容的准确性、完整性、适用性、可靠性等不作任何明示或暗示的保证。任何单位、公司或个人因为下载、使用或信赖本文档而发生任何差错或经济损失的，阿里云不承担任何法律责任。在任何情况下，阿里云均不对任何间接性、后果性、惩戒性、偶然性、特殊性或刑罚性的损害，包括用户使用或信赖本文档而遭受的利润损失，承担责任（即使阿里云已被告知该等损失的可能性）。
5. 阿里云网站上所有内容，包括但不限于著作、产品、图片、档案、资讯、资料、网站架构、网站画面的安排、网页设计，均由阿里云和/或其关联公司依法拥有其知识产权，包括但不限于商标权、专利权、著作权、商业秘密等。非经阿里云和/或其关联公司书面同意，任何人不得擅自使用、修改、复制、公开传播、改变、散布、发行或公开发表阿里云网站、产品程序或内容。此外，未经阿里云事先书面同意，任何人不得为了任何营销、广告、促销或其他目的使用、公布或复制阿里云的名称（包括但不限于单独为或以组合形式包含“阿里云”、Aliyun”、“万网”等阿里云和/或其关联公司品牌，上述品牌的附属标志及图案或任何类似公司名称、商号、商标、产品或服务名称、域名、图案标示、标志、标识或通过特定描述使第三方能够识别阿里云和/或其关联公司）。
6. 如若发现本文档存在任何错误，请与阿里云取得直接联系。

通用约定

表 -1: 格式约定

格式	说明	样例
	该类警示信息将导致系统重大变更甚至故障，或者导致人身伤害等结果。	 禁止： 重置操作将丢失用户配置数据。
	该类警示信息可能导致系统重大变更甚至故障，或者导致人身伤害等结果。	 警告： 重启操作将导致业务中断，恢复业务所需时间约10分钟。
	用于警示信息、补充说明等，是用户必须了解的内容。	 说明： 导出的数据中包含敏感信息，请妥善保管。
	用于补充说明、最佳实践、窍门等，不是用户必须了解的内容。	 说明： 您也可以通过按 Ctrl + A 选中全部文件。
>	多级菜单递进。	设置 > 网络 > 设置网络类型
粗体	表示按键、菜单、页面名称等UI元素。	单击 确定 。
courier字体	命令。	执行 cd /d C:/windows 命令，进入Windows系统文件夹。
斜体	表示参数、变量。	bae log list --instanceid Instance_ID
[]或者[a b]	表示可选项，至多选择一个。	ipconfig [-all -t]
{ }或者{a b}	表示必选项，至多选择一个。	swich {stand slave}

目录

法律声明	1
通用约定	1
1 安装与部署	1
1.1 概述.....	1
1.1.1 超融合产品.....	1
1.1.2 高可用技术.....	2
1.1.3 网络拓扑规划.....	3
1.2 安装部署.....	6
1.2.1 安装操作系统.....	6
1.2.2 存储网络.....	10
1.2.3 管理网络.....	11
1.2.4 云主机数据网络.....	12
1.2.5 Ceph存储集群.....	12
1.2.6 高可用套件.....	13
1.2.6.1 功能介绍.....	13
1.2.6.2 部署过程.....	13
1.2.7 集群升级.....	20
1.2.7.1 内嵌服务升级.....	20
1.2.7.2 高可用升级.....	20
1.2.8 管理节点迁移.....	21
1.2.9 配置更新.....	21
1.3 其他操作.....	23
1.3.1 卸载操作.....	23
1.3.2 日志输出.....	23
2 高可用测试与恢复	24
2.1 计划运维.....	24
2.1.1 单节点需要维护.....	24
2.1.1.1 单节点关闭.....	24
2.1.1.1.1 该超融合节点运行非管理节点主机.....	24
2.1.1.1.2 该超融合节点运行管理节点主机.....	25
2.1.1.2 单节点启动.....	26
2.1.2 三节点需要维护.....	27
2.1.2.1 三节点关闭.....	27
2.1.2.2 三节点启动.....	28
2.2 异常处理.....	30
2.2.1 单节点异常处理.....	30
2.2.1.1 单节点异常掉电.....	30

2.2.1.1.1 该超融合节点运行非管理节点主机.....	31
2.2.1.1.2 该超融合节点运行管理节点主机.....	32
2.2.1.2 单节点网络异常.....	33
2.2.1.2.1 该超融合节点运行非管理节点主机.....	33
2.2.1.2.2 该超融合节点运行管理节点主机.....	34
2.2.2 两节点异常处理.....	36
2.3 节点修复.....	37
2.3.1 单节点故障修复.....	37
2.3.2 两节点故障无法修复.....	38
3 命令行使用手册.....	39
3.1 简介.....	39
3.2 -h 帮助内容.....	39
3.3 -v 版本信息.....	40
3.4 check-config 配置检查.....	40
3.5 install -p -c 安装命令.....	40
3.6 uninstall 卸载命令.....	42
3.7 config-sample 样本配置生成.....	42
3.8 export-config 当前配置生成.....	43
3.9 migrate 迁移命令.....	43
3.10 import-config 配置升级.....	44
3.11 status 状态信息.....	44
3.12 stop 集群关闭.....	46
3.13 start 集群启动.....	47
3.14 reset-MNVM-password 管理节点主机重置root密码.....	48
专有云术语表.....	49
混合云术语表.....	52

1 安装与部署

1.1 概述

1.1.1 超融合产品

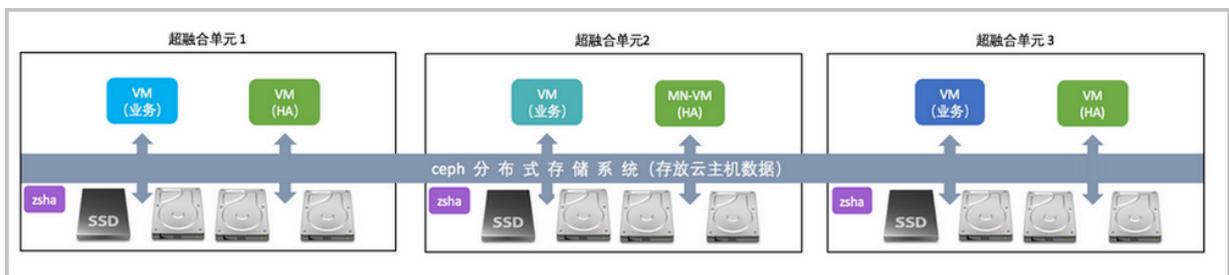
超融合基础架构（Hyper-converged infrastructure，缩写HCI），是一种集成了存储系统及虚拟化系统的信息基础架构框架。

在这样的框架下，每个**超融合单元**将成为存储系统和虚拟化系统的一部分。通过高速以太网网络，将若干个超融合单元**互联**，构建超融合基础架构平台。

ZStack for Alibaba Cloud超融合产品基于超融合基础架构而设计，采用独有的优化技术，聚合通用服务器的CPU、内存、硬盘和网络等资源，实现稳定可靠产品级的交付。

以下是图 1-1: 超融合-三节点经典模型：

图 1-1: 超融合-三节点经典模型



ZStack for Alibaba Cloud超融合产品有以下特点：

- 通过Ceph分布式存储，由各个超融合单元的固态硬盘（SSD）和机械硬盘（HDD）提供统一的存储系统
- 通过虚拟化技术提供云主机服务，并将其数据存放在Ceph分布式存储
- **管理节点主机**以云主机形式运行，管理三个超融合单元的计算、存储和网络资源分配与调度

ZStack for Alibaba Cloud超融合产品，最重要的技术亮点是**高可用服务**，其中：

- Ceph分布式存储提供存储级别高可用能力
- zsha守护服务提供管理节点主机高可用
- 管理节点主机提供业务云主机的高可用

三重高可用技术的保护能有效支撑云主机正常运行，并对云主机数据提供冗余保障（数据两或三副本）。

1.1.2 高可用技术

ZStack for Alibaba Cloud超融合产品提供三重高可用技术，以下分别介绍其功能和作用。

分布式存储高可用

ZStack for Alibaba Cloud超融合产品，是基于**Ceph分布式存储的RBD块存储技术**，实现云主机的云盘存放与访问（包括根云盘和数据云盘）。

针对不同的业务场景，推荐采用合适的数据副本级别：

- **两副本**，一份原始数据以两份存放：偏向成本考虑，适合普通类型业务系统场景。
- **三副本**，一份原始数据以三份存放：偏向数据安全考虑，适合关键类型业务系统场景。

Ceph支持**故障域**（Fail Domain）的保护机制：

- 在相同的故障域内，存储介质和硬件设备可任意离线或失去。
- 在不相同的故障域内，存储介质和硬件设备故障的**位置**和**数量**需依赖故障域的定义。

在三节点经典模型下，部署Ceph时推荐采用**Rack-Node层级**，并设置基于Node的副本冗余机制。

关于存储系统的设定，可参考Ceph官方资料或供应商的产品资料。

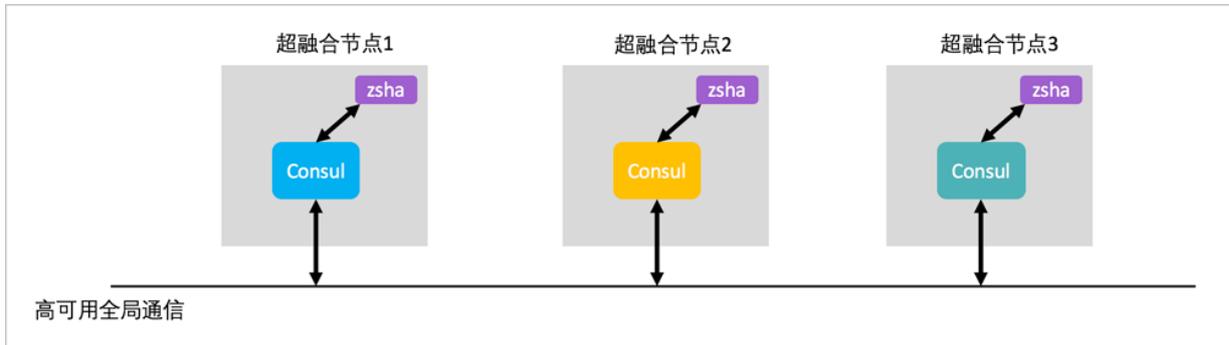
管理节点高可用

ZStack for Alibaba Cloud超融合产品，通过**zsha守护服务**提供强一致选举机制，保障管理节点主机的运行状态。

在三节点经典模型下，每个节点均运行**zsha和consul守护服务**，通过两两间建立**通信服务**（TCP 8300）维护强一致共享数据结构。

其运行通信原理如图 1-2: [zsha与consul服务消息通信原理](#)所示：

图 1-2: zsha与consul服务消息通信原理



业务云主机高可用

ZStack for Alibaba Cloud超融合产品，通过管理节点主机，对其他超融合节点进行**心跳检测**，在网络或节点失效情况下，将业务云主机**自动切换**到其他健康的超融合单元运行，能快速恢复业务系统，提高业务系统的可用性。

1.1.3 网络拓扑规划

ZStack for Alibaba Cloud超融合产品中，超融合单元的**互连**尤其重要。在科学和合适的网络环境下，组建超融合架构，才能获得安全、灵活和高性能的云计算架构。

ZStack for Alibaba Cloud超融合产品，支持以下网络流量模型：

- 云主机数据网络
- 管理网络
- 存储网络

云主机数据网络

云主机数据网络，即云主机向外提供应用服务的网络，或云主机之间相互沟通的网络。

该网络也承载云路由三层流量。根据业务负载类型，建议采用**双链路1GbE或10GbE以太网**。



说明：

- 若计划部署的云主机采用单一扁平网络，建议此网络的交换机端口设定Access模式，以下行文以此为主体描述；
- 若计划部署的云主机考虑多个扁平或云路由网络，建议此网络的交换机端口设定Trunk模式；以下行文并未提及该场景操作过程，如需了解，请联系官方技术支持团队。

管理网络

管理网络主要承载**管理节点主机与物理主机**的消息通信，包括任务下发和云主机迁移。

考虑到管理网络的业务负载需求，建议采用**双链路1GbE或10GbE以太网**。

存储网络

存储网络主要承载**Ceph存储流量和zsha守护服务通信**的网络流量。该网络是保障存储集群、管理节点和业务云主机高可用的关键网络。

考虑到存储网络的业务负载需求，建议采用**双链路10GbE以太网**。

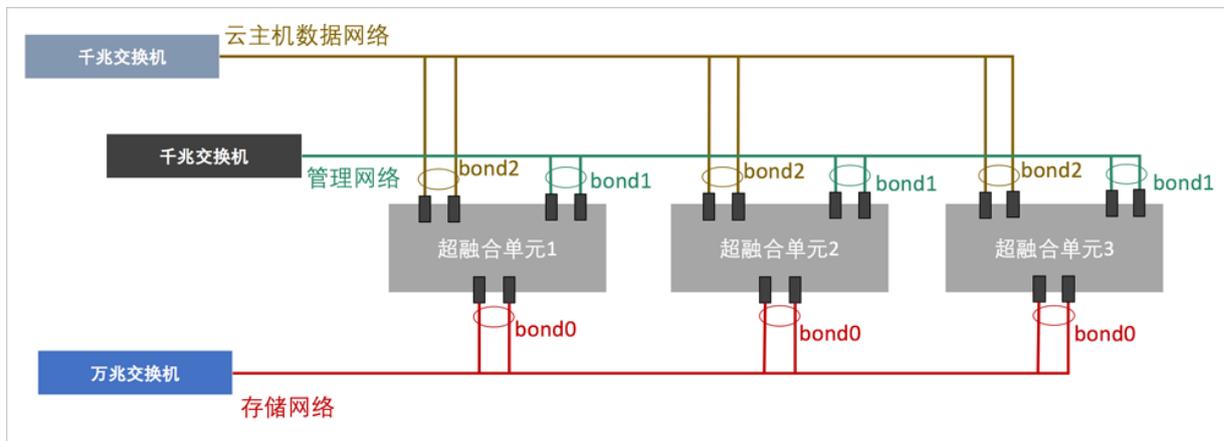


说明：

- 建议此网络的交换机端口设定Access模式（Trunk模式亦支持，但本文未描述）
- 查看交换机是否支持巨型帧（Jumbo Frame），若支持则建议开启，全链路MTU设定值为9000

综上所述，在三节点经典网络中，[图 1-3: 超融合-三节点网络拓扑](#)如下：

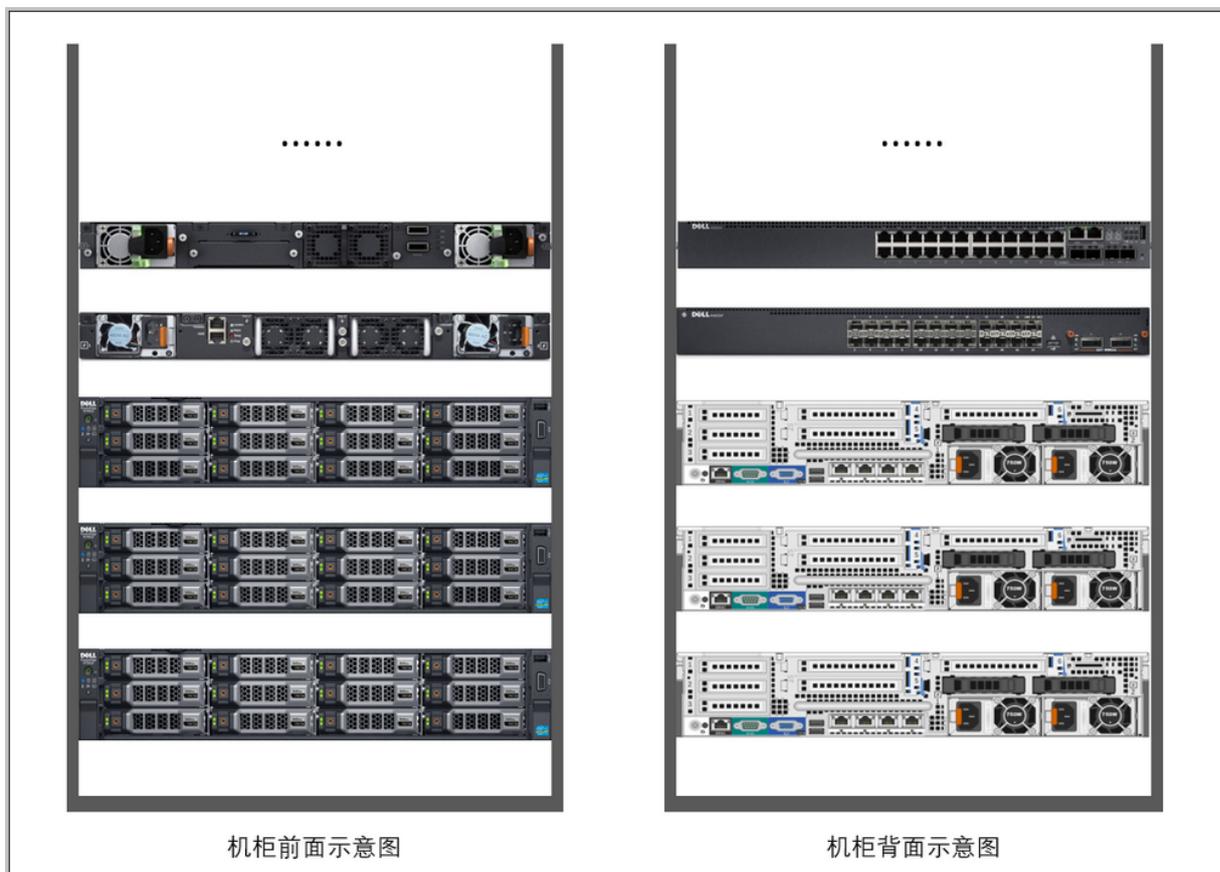
图 1-3: 超融合-三节点网络拓扑



管理员根据上述的网络架构图，对网络设备和服务器进行上架并连线。

以下是[图 1-4: 超融合服务器和网络设备上架后示意图](#)：

图 1-4: 超融合服务器和网络设备上架后示意图



此外，每个超融合节点的配置建议如下。管理员可根据业务性能需求，合理调配CPU、内存和硬盘的容量配比，以达到合适的平衡状态。

表 1-1: 每个超融合节点的配置

	配件	型号	数量	总数
超融合节点	CPU	Intel(R) Xeon(R) CPU E5-2630 v4 @ 2.20GHz	2	3个
	内存	DDR4 16GB	8	
	主板	双路服务器标准主板	1	
	阵列卡	阵列卡支持SAS/SATA RAID 0/1 /10 支持直通模式	1	
	固态硬盘	Intel SSD DC S3610 480GB	2	
	机械硬盘1	SAS HDD 300GB 3.5", 15k rpm	2	

机械硬盘2	NL SAS HDD 2TB 3.5" , 7.2k rpm	6
千兆网口	以太网1GbE , RJ45	4
万兆网口	以太网10GbE , SFP+	2
光电模块	-	
光纤HBA卡	-	
远程管理	DELL iDRAC企业版	1
电源	标准电源1100W	2

1.2 安装部署

本章节描述ZStack for Alibaba Cloud超融合产品的安装部署过程。请管理员准备以下必要的软件包，以便安装部署过程顺利执行：

- ZStack for Alibaba Cloud
- 管理节点镜像
- 高可用服务套件
- 内嵌服务套件

以上安装包最新版本请见[这里](#)

1.2.1 安装操作系统

操作步骤

1. 准备

管理员对上架的网络设备和服务器加载电源，手动启动服务器进入BIOS，检查以下内容：

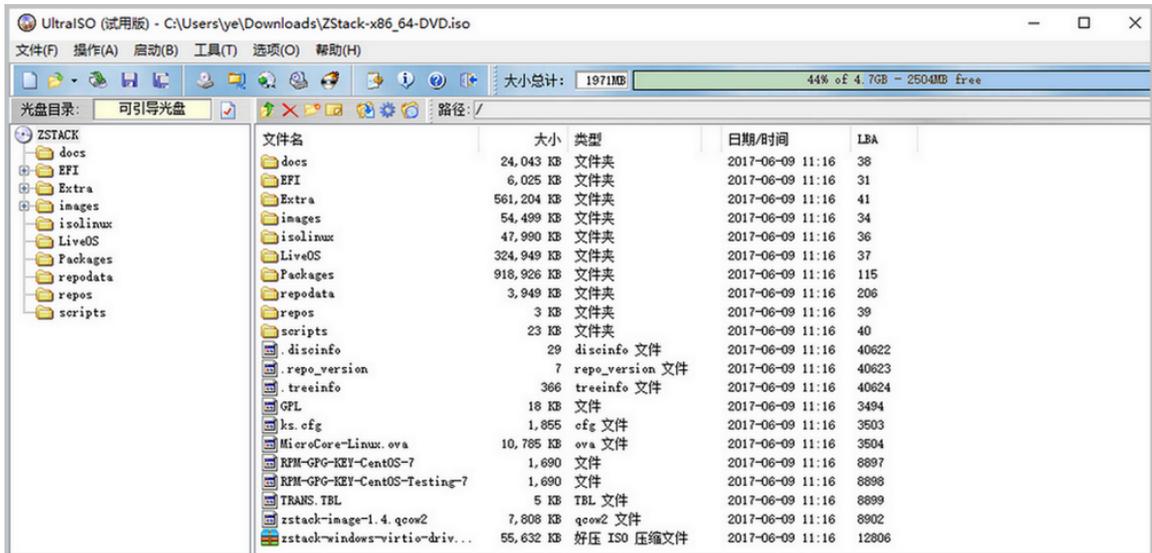
- 激活所有CPU核心和启用超线程功能，设定系统性能为最高性能状态；
- 打开硬件虚拟化VT功能，支持硬件虚拟化技术加速优化功能；
- 进入阵列卡设定，对两块系统硬盘配置RAID1（Mirror），其余硬盘设定直通模式。

2. 在UltraISO打开ZStack for Alibaba Cloud DVD镜像

- ZStack for Alibaba Cloud操作系统ISO镜像可通过DVD-RW设备刻录成安装光盘，也可通过UltraISO工具将把ISO文件刻录到U盘。
- 打开UltraISO，点击**文件**按钮，选择打开已下载好的ISO文件。

如图 1-5: 在UltraISO打开DVD镜像所示：

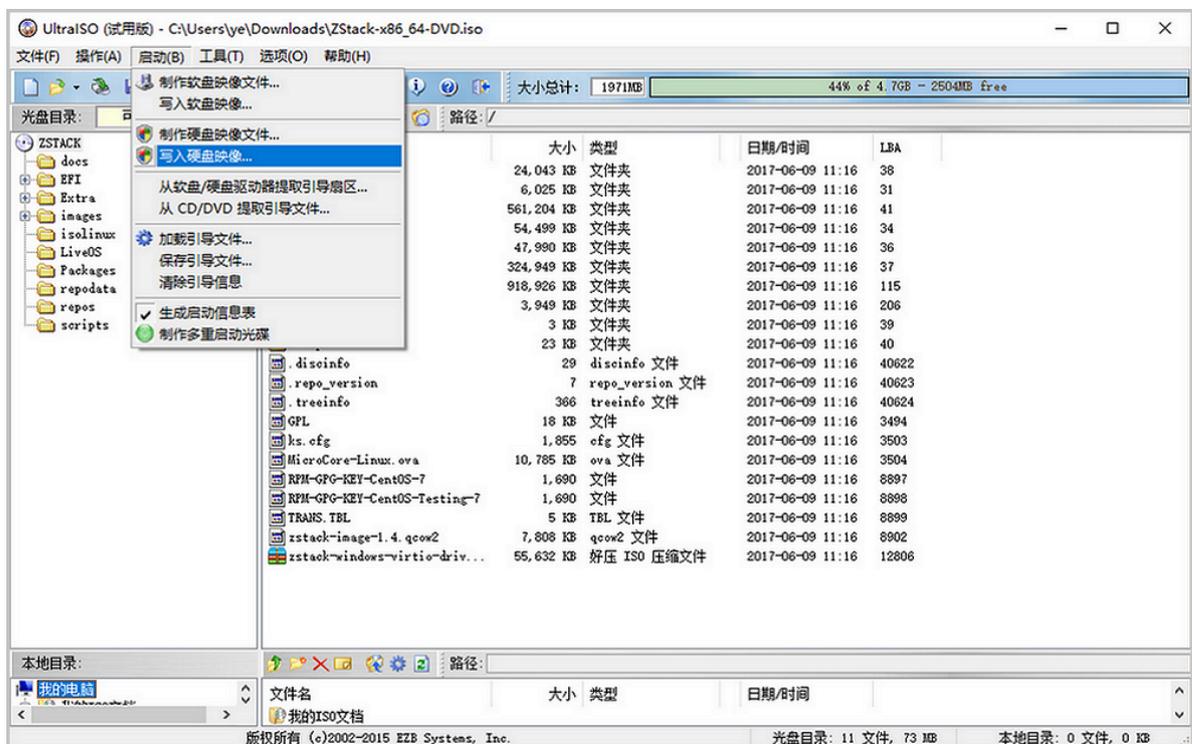
图 1-5: 在UltraISO打开DVD镜像



3. 写入硬盘镜像

在UltraISO点击启动按钮，选择写入硬盘镜像，如图 1-6: 在UltraISO写入DVD镜像所示：

图 1-6: 在UltraISO写入DVD镜像



4. 在UltraISO确认写入ZStack for Alibaba CloudDVD镜像

- 如果系统只插了一个U盘，则默认以此U盘进行刻录和写入，在刻录前，**注意备份U盘之前的内容**。
- 其他选项，按照默认设置，无须额外配置，点击**写入**。

如图 1-7: 在UltraISO确认写入ISO镜像所示：

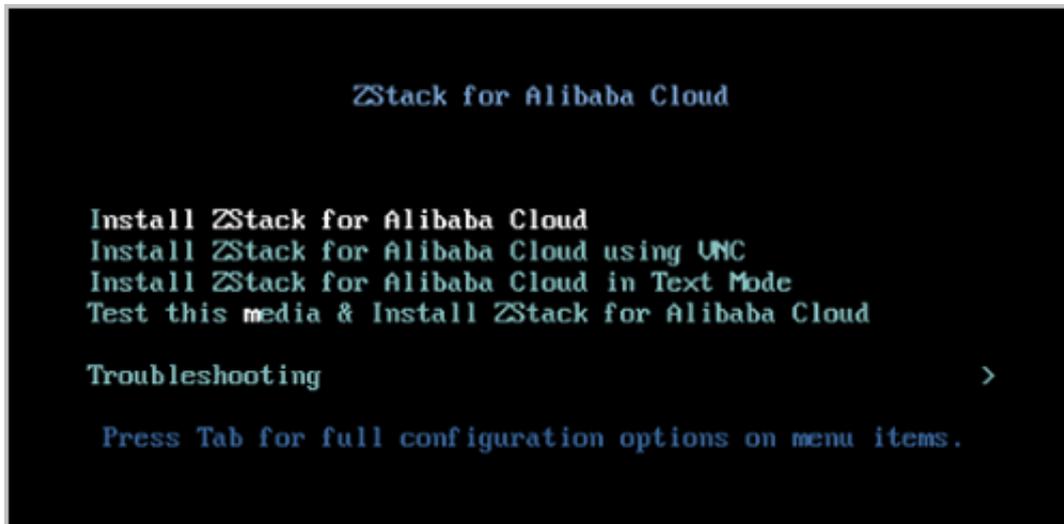
图 1-7: 在UltraISO确认写入ISO镜像



5. 进入安装导航

ISO镜像已经刻录到U盘。此时U盘可用来作为启动盘，支持Legacy模式和UEFI模式引导。管理员通过安装介质，引导节点启动，并进入安装导航，如图 1-8: U盘引导界面所示：

图 1-8: U盘引导界面



6. 安装操作系统

默认选择**Install ZStack for Alibaba Cloud**开始安装操作系统。

在进入安装界面后，已经预先配置默认选项：

- **时区**：亚洲东八区
- **语言**：English(United States)
- **键盘**：English(US)

一般情况下管理员无需更改配置。管理员需自行执行硬盘的分区，推荐分区如下（UEFI 模式）：

- `/boot/efi`：创建分区500MB
- `/boot`：创建分区1GB
- `swap`（交换分区）：创建分区32GB
- `/`（根分区）：配置剩下容量

分区配置完后，选择**Software Selection**进入服务器安装角色候选，选择**ZStack for Alibaba Cloud Compute Node**计算节点模式，确定后回到主界面。

点击**Begin Installation**进行安装。安装过程将会自动进行，管理员需要设定root账户密码。

安装结束后，重新引导服务器并拔掉U盘。如安装成功，则服务器重启后进入操作系统登录提示符，使用root和设置的密码登录到操作系统。

**说明：**

管理员可根据自身需要更改密码。

1.2.2 存储网络

管理员对三超融合节点安装ZStack for Alibaba Cloud操作系统后，可对**存储网络**进行如下设定：

节点	网卡 1	网卡 2	聚合接口	网桥	地址
超融合节点1	p5p1	p5p2	bond0	br_bond0	192.168.93.3/ 24
超融合节点2	p5p1	p5p2	bond0	br_bond0	192.168.93.4/ 24
超融合节点3	p5p1	p5p2	bond0	br_bond0	192.168.93.5/ 24

以上是示例数据，管理员请根据具体部署环境的网卡和网络地址进行配置。以**超融合节点 1**为例，其他节点配置相似。执行配置命令：

```
# 创建聚合网卡bond0
[root@localhost ~]# zs-bond-lacp -c bond0
# 将网卡p5p1与p5p2均添加到bond0
[root@localhost ~]# zs-nic-to-bond -a bond0 p5p1
[root@localhost ~]# zs-nic-to-bond -a bond0 p5p2
# 配置上述链路聚合后,请管理员在对应的交换机网口配置LACP聚合
# 创建网桥br_bond0,指定网络IP、掩码和网关
[root@localhost ~]# zs-network-setting -b bond0 192.168.93.3 255.255.255.0 192.168.93.1
# 查看聚合端口bond0是否创建成功
[root@localhost ~]# zs-show-network
```

```
...
-----
| Bond Name | SLAVE(s)      | BONDING_OPTS                               |
-----
| bond0     | p5p1          | miimon=100 mode=4 xmit_hash_policy=layer2+3 |
|           | p5p2          |                                             |
-----
```

**说明：**

- p5p1和p5p2加载到bond0后，对应交换机的端口需要配置LACP聚合，否则网络通信将异常；如果交换机不支持LACP聚合，请联系网络设备厂商更换设备；
- 通过bond0创建网桥后，网桥命名为**br_bond0**，将提供Ceph存储集群通信和zsha守护服务通信；
- 关于网桥的IP地址、子网掩码和网关参数，用户需按照实际情况填写；

- 存储网络配置完成后，可通过ping命令进行检测；若配置正确，则三节点的存储网络对应的IP地址可两两互ping。

存储网络配置完成后，随之可配置管理网络。

1.2.3 管理网络

管理员对三节点安装ZStack for Alibaba Cloud操作系统后，可对**管理网络**进行如下设定：

节点	网卡 1	网卡 2	聚合接口	网桥	地址
超融合节点1	eth0	eth1	bond1	br_bond1	172.20.198.3/ 24
超融合节点2	eth0	eth1	bond1	br_bond1	172.20.198.4/ 24
超融合节点3	eth0	eth1	bond1	br_bond1	172.20.198.5/ 24

以上是示例数据，管理员请根据具体部署环境的网卡和网络地址进行配置。以**超融合节点 1**为例，其他节点配置相似。执行配置命令：

```
# 创建聚合网卡bond1
[root@localhost ~]# zs-bond-lacp -c bond1
# 将网卡eth0与eth1均添加到bond1
[root@localhost ~]# zs-nic-to-bond -a bond1 eth0
[root@localhost ~]# zs-nic-to-bond -a bond1 eth1
# 配置上述链路聚合后,请管理员在对应的交换机网口配置LACP聚合
# 创建网桥br_bond1,指定网络IP、掩码和网关
[root@localhost ~]# zs-network-setting -b bond1 172.20.198.3 255.255.255.0 172.20.198.1
# 查看聚合端口bond1是否创建成功
[root@localhost ~]# zs-show-network
...
```

```
-----
| Bond Name | SLAVE(s)      | BONDING_OPTS                               |
-----
| bond1     | eth0          | miimon=100 mode=4 xmit_hash_policy=layer2+3 |
|           | eth1          |                                             |
-----
```



说明：

- eth0和eth1加载到bond1后，对应交换机的端口需要配置LACP聚合，否则网络通信将异常；如果交换机不支持 LACP 聚合，请联系网络设备厂商更换设备；
- 通过bond1创建网桥后，网桥命名为br_bond1，将提供管理节点主机与物理主机的消息通信；
- 关于网桥的IP地址、子网掩码和网关参数，用户需按照实际情况填写；

- 管理网络配置完成后，可通过ping命令进行检测；若配置正确，则三节点的管理网络对应的IP地址可两两互ping。

存储网络和管理网络配置完成后，随之可配置云主机数据网络。

1.2.4 云主机数据网络

管理员对三节点安装ZStack for Alibaba Cloud操作系统后，可对**云主机数据网络**进行如下设定：

节点	网卡 1	网卡 2	聚合接口	网桥	地址
超融合节点1	em1	em2	bond2	-	-
超融合节点2	em1	em2	bond2	-	-
超融合节点3	em1	em2	bond2	-	-

以上是示例数据，管理员请根据具体部署环境的网卡和网络地址进行配置。以超融合节点1为例，其他节点配置相似。执行配置命令：

```
# 创建聚合网卡bond2
[root@localhost ~]# zs-bond-lacp -c bond2
# 将网卡em1与em2均添加到bond2
[root@localhost ~]# zs-nic-to-bond -a bond2 em1
[root@localhost ~]# zs-nic-to-bond -a bond2 em2
# 配置上述链路聚合后,请管理员在对应的交换机网口配置LACP聚合
# 云主机数据网络,无需创建网桥
# 查看聚合端口bond2是否创建成功
[root@localhost ~]# zs-show-network
...
-----
| Bond Name | SLAVE(s)      | BONDING_OPTS                               |
-----
| bond2    | em1          | miimon=100 mode=4 xmit_hash_policy=layer2+3 |
|          | em2          |                                             |
-----
```



说明：

em1和em2加载到bond2后，对应交换机的端口需要配置LACP聚合，否则网络通信将异常；如果交换机不支持LACP聚合，请联系网络设备厂商更换设备。

云主机数据网络配置完成后，随之可安装Ceph存储集群。

1.2.5 Ceph存储集群

ZStack for Alibaba Cloud超融合产品，可运行在Ceph社区版之上，具体参照[ZStack官网文档](#)《开源版Ceph分布式存储 部署教程》。

此外，若用户关注数据安全和IO性能，请联系官方技术支持团队，获得Ceph企业版产品相关信息。

1.2.6 高可用套件

1.2.6.1 功能介绍

ZStack for Alibaba Cloud超融合产品提供了专属的高可用套件，保障管理服务长期处于在线运行，与业务云主机高可用机制相结合，最大程度降低物理设备异常情况下的云主机失效带来的业务中断风险。

高可用套件，是基于Consul开发的强一致性消息通信组件，部署在超融合节点中，部署数量满足 $2N+1$ 条件。在三节点经典模型中，能容忍1个超融合节点的离线或失效，保证管理服务正常运行。

在部署高可用套件之前，确保每个超融合节点均安装libvirt和qemu-kvm-ev软件包：

```
# 每个超融合节点均安装libvirt和qemu-kvm-ev软件包
[root@localhost ~]# yum -y --disablerepo=* --enablerepo=zstack-local,qemu-kvm-ev \
install libvirt qemu-kvm-ev qemu-img-ev
```

1.2.6.2 部署过程

操作步骤

1. 导入镜像

管理员已经获得ZStack for Alibaba Cloud管理节点镜像(QCOW2 格式)，可将其导入Ceph存储集群。执行以下命令：

```
# 通过scp工具将QCOW2镜像传输到超融合节点1
[root@localhost ~]# ls
ZStack for Alibaba Cloud-Management-Server-2.x.qcow2

# 创建存储池,例如创建一个名为zstack的存储池,默认传递pg_num参数为128
# 若使用Ceph企业版,则无需创建存储池,直接使用已经创建存储池空间
[root@localhost ~]# ceph osd pool create zstack 128

# 显示当前Ceph存储池中的内容
[root@localhost ~]# ceph osd pool ls
zstack

# 将QCOW2镜像导入到名为zstack的存储池指定路径,显示导入进度条
[root@localhost ~]# qemu-img convert -p -f qcow2 -O raw \
ZStack for Alibaba Cloud-Management-Server-2.x.qcow2 rbd:zstack/mnvm.img

# 导入完成后,显示导入后镜像的信息
[root@localhost ~]# rbd info zstack/mnvm.img
rbd image 'mnvm.img':
size 500 GB in 128001 objects
order 22 (4096 kB objects)
block_name_prefix: rbd_data.10352ae8944a
format: 2
```

```
features: layering, striping
flags:
stripe unit: 4096 kB
stripe count: 1
```



说明：

- ZStack for Alibaba Cloud-Management-Server-2.x.qcow2为镜像文件名，根据版本情况变更
- 示例中，Ceph存储池命名为zstack，可根据具体需要制定命名
- 镜像导入过程中，会显示时间进度，整个过程需若干分钟，请耐心等待
- 若在创建存储池时有误，需要删除存储池，请执行：

```
# 若创建存储池有误，需要删除存储池,可执行以下命令
[root@localhost ~]# ceph osd pool delete zstack zstack --yes-i-really-really-mean-it
```

- 若在导入镜像时，发生网络中断，或者其他异常情况，需删除镜像，请执行：

```
# 若需要删除镜像,可通过以下命令
[root@localhost ~]# rbd rm zstack/mnvm.img
```

2. 编写配置

管理员完成管理服务镜像导入后，接下来编辑高可用套件的初始化配置文件。

执行以下命令，生成并编写配置文件：

```
# 生成的样本配置文件临时保存在/tmp/sample.CephConf.config.json,用户可自行拷贝到其它目录
[root@localhost ~]# chmod +x ZStack for Alibaba Cloud-HA-Suite-2.x.bin
[root@localhost ~]# ./ZStack for Alibaba Cloud-HA-Suite-2.x.bin config-sample CephConf
# 编辑样本配置文件
[root@localhost ~]# cp /tmp/sample.CephConf.config.json /tmp/config.json
[root@localhost ~]# vim /tmp/config.json
...
{
  "Node": [
    "192.168.255.205 ceph-1",
    "192.168.255.221 ceph-2",
    "192.168.255.183 ceph-3"
  ], // 设置三台超融合节点(管理节点主机所在物理机)的IP地址以及主机名称
  "MemorySizeInGB": 8, //管理节点主机的内存配置,最小值为8GB
  "CPU": 4, //管理节点主机的CPU数目,最小值为4核
  "Type": "CephConf", //配置文件的类型
  "MonAddrs": [
    "192.168.200.100:6789",
    "192.168.200.125:6789",
    "192.168.200.160:6789"
  ], // 设置Ceph监控节点IP地址
  "ChronyServers": [
    "192.168.255.205"
  ], // 设置chrony时间服务器IP地址
  "PoolName": "zstack", // 管理节点主机所在的存储池名
  "ImageFolder": "/storage/", // 管理节点镜像存放的本地目录
```

```

"DNS": [
  "172.20.0.1", // 管理节点主机的私网DNS信息
  "223.5.5.5" // 管理节点主机的公网DNS信息
],
"Network": [
  {
    "Bridge": "br_zsn0", //管理节点主机的网卡所挂接的网桥名称
    "MacAddress": "02:98:54:4b:a5:c0", //管理节点主机的网卡MAC地址,单次随机生成
    "Ipaddr": "172.20.198.24@ceph-1, 172.20.200.24@ceph-2, 172.20.210.24@ceph-3",
    //管理节点主机的IP地址
    "Netmask": "255.255.255.0@ceph-1, 255.255.255.0@ceph-2,
    255.255.255.0@ceph-3", //管理节点主机的掩码
    "Gateway": "172.20.198.1@ceph-1, 172.20.200.1@ceph-2, 172.20.210.1@ceph-3", //管
    理节点主机的网关
    //管理节点主机的IP地址以及对应掩码和网关支持per-node格式,可实现管理节点主机跨网
    段启动
    "IsMgmt": true, //管理网络设置, 本示例设为true
    "IsDefRoute": false //默认路由配置, 本示例设为false
  },
  {
    "Bridge": "br_zsn1", //管理节点主机的网卡所挂接的网桥名称
    "MacAddress": "f6:75:d8:5d:73:73", //管理节点主机的网卡MAC地址,单次随机生成
    "Ipaddr": "203.114.54.10", //管理节点主机的IP地址
    "Netmask": "255.255.255.0", //管理节点主机的掩码
    "Gateway": "203.114.54.1", //管理节点主机的网关
    "IsMgmt": false, //管理网络设置, 本示例设为false
    "IsDefRoute": true //默认路由配置, 本示例设为true
  }
],
"ConsoleProxyOverriddenIP": "" //ZStack管理服务的控制台服务地址,支持per-node格式
"AllowResetRootPassword": false //默认为false;设为true时,管理节点主机才允许重置root密
码
}
...

```

管理员需要按照具体部署场景，修改上述参数。



说明：

- CPU和MemorySizeInGB设置值推荐如下：

场景	配置	备注
中小规模	CPU核4个/内存8GB	管理小于100个物理主机与1000个云主机
特大规模	CPU核8个/内存16GB	<ul style="list-style-type: none"> • 管理大于100个物理主机与1000个云主机内 • 管理小于1000个物理主机与10000个云主机

场景	配置	备注
超大规模	CPU核12个/内存32GB	管理大于1000个物理主机与10000个云主机

- 关于管理节点主机的网卡MAC地址，需确保在二层网络里是唯一的。
- 关于MonAddrs的设置：

管理员可单独设置Ceph监控节点IP地址，实现Ceph监控节点与超融合节点分离部署，示例如下：

```
"MonAddrs": [
  "192.168.200.100:6789",
  "192.168.200.125:6789",
  "192.168.200.160:6789"
]
```

- 关于ChronyServers的设置：
 - ChronyServers可使用超融合节点作为时间源，示例如下：

```
"ChronyServers": [
  "192.168.255.205"
],
```

- 管理员也可自行设置chrony时间源，示例如下：

```
"ChronyServers": [
  "time1-7.aliyun.com "
],
```

- ChronyServers中指定的时间服务器不可对自身做时间同步。



说明：

补充说明，对于未部署管理节点高可用服务的场景，默认情况下管理节点作为NTP时间源；若希望指定时间服务器，方法如下：

- 在管理节点的**zstack.properties**文件中添加：

```
chrony.serverIp.0 = xx.xx.xx.xx
# xx.xx.xx.xx为时间服务器IP地址,可为管理节点IP地址,也可自行指定其它chrony时间源
```

- 重连物理主机生效。
- chrony.serverIp.0中指定的时间服务器不可对自身做时间同步。
- 管理节点主机跨网段启动

在**IsMgmt**为**true**情况下，管理节点的IP地址、以及对应的掩码、网关均设置为**per-node**格式，可实现管理节点主机跨网段启动，示例如下：

```
"Network": [  
  {  
    ...  
    "Ipaddr": "172.20.198.24@ceph-1, 172.20.200.24@ceph-2, 172.20.210.24@ceph-3",  
    "Netmask": "255.255.255.0@ceph-1, 255.255.255.0@ceph-2, 255.255.255.0@ceph-3",  
    "Gateway": "172.20.198.1@ceph-1, 172.20.200.1@ceph-2, 172.20.210.1@ceph-3",  
    "IsMgmt": true,  
    ...  
  },  
]
```

ConsoleProxyOverriddenIP也支持**per-node**格式，当管理节点主机跨网段启动后，可打开相应管理服务控制台。

- 管理节点主机允许重置root密码为系统默认值（即**password**）

在**AllowResetRootPassword**为**true**情况下，管理节点主机允许重置root密码为系统默认值（即**password**），详情请参考[配置更新](#)章节。

3. HA初始化

在上一节中，管理员已完成对config.json配置文件的编写，接下来可执行高可用套件的初始化。执行以下命令：

```
# 为三个超融合节点安装高可用套件  
[root@localhost ~]# ./ZStack for Alibaba Cloud-HA-Suite-2.x.bin install \  
-p password -c /tmp/config.json
```



说明：

- 安装命令中，**-p**传递三个超融合节点账号root密码。建议首次初始化时，对三个超融合节点设置一致密码，安装完成后，基于安全考虑可按需修改或关闭密码登陆。
- **/tmp/config.json**为填写完成的配置文件路径，需与上述步骤中真实保存的文件绝对路径保持一致。

高可用套件初始化完成后，可执行以下命令查看集群状态：

```
[root@localhost ~]# zsha status  
# 查看 HA 集群状态  
>>>MN-VM Running On  
192.168.255.205 : running  
  
>>>Last MN-VM Start Record
```

```
>>>ZStack HA Services Running On
192.168.255.205 : ceph-1
192.168.255.221 : ceph-2
192.168.255.183 : ceph-3

>>>Consul Members
Node           Address           Status Type  Build Protocol DC
192.168.255.205 192.168.255.205:8301 alive  server 0.7.2 2    dc1
192.168.255.221 192.168.255.221:8301 alive  server 0.7.2 2    dc1
192.168.255.183 192.168.255.183:8301 alive  server 0.7.2 2    dc1

>>>Management Node IP
172.20.198.3

>>>Public Network IP
203.114.54.10

>>>Ceph Status
cluster 7eafa8a3-7318-4d93-86fe-4186e1509f15
health HEALTH_WARN
  too many PGs per OSD (592 > max 300)
  noout,noscrub,nodeep-scrub flag(s) set
monmap e1: 3 mons at {ceph-1=192.168.255.205:6789/0,ceph-2=192.168.255.221:6789/0,ceph-3=192.168.255.183:6789/0}
election epoch 6, quorum 0,1,2 ceph-3,ceph-1,ceph-2
osdmap e33: 3 osds: 3 up, 3 in
  flags noout,noscrub,nodeep-scrub
pgmap v25752: 592 pgs, 6 pools, 8328 MB data, 2152 objects
  24660 MB used, 860 GB / 884 GB avail
  592 active+clean
client io 10524 B/s rd, 3134 B/s wr, 26 op/s

>>>Checking Host I/O
done.
```

至此，ZStack for Alibaba Cloud超融合产品高可用云管理平台安装完成。

当所有节点正常加入并开始运行高可用服务程序后，将选择其中一个超融合节点运行管理节点主机，高可用保护机制开始正常运转，管理节点也开始正常提供服务。

4. 云平台初始化

管理员可通过浏览器访问管理服务界面（<http://manage-server-ip:5000>），并完成云平台初始化操作。如[登陆界面](#)所示：

图 1-9: 登录界面



初始化结束后，管理员需在**基础设置**中，设定**物理主机保留内存**。

计算公式是Ceph存储保留空间 + 管理节点主机保留空间 + 虚拟化保留内存 = 物理主机保留内存



说明：

- 混合加速场景（SSD+HDD）：

在实际生产环境中，每个OSD的缓存加速内存大小通常在50GB~250GB，相应地，每个OSD的存储保留空间在3.5GB~5.5GB，推荐每个OSD的存储保留空间为6GB

- 全闪存场景（全SSD）：

推荐每个SSD的存储保留空间为4GB

以三节点超融合系统为例，每个服务器有3个OSD，每个OSD的缓存加速内存大小接近250GB，因此，推荐每个服务器的Ceph存储保留空间6GB*3=18GB，虚拟化保留内存1GB，管理节点主机保留空间根据上一章节的配置文件设定内存值而定。故计算结果：

Ceph存储内存 18GB + 管理节点主机内存 8GB + 虚拟化内存 1GB = 物理主机保留内存 27GB
所以设定值为 27G

如图 1-10: 物理主机保留内存所示：

图 1-10: 物理主机保留内存

名称	类别	简介	值	操作
云主机高可用全局开关	高可用	默认为true, 用于设置云主机高可用功能...	true	✎
CPU超分率	物理机	默认为10, 主要用于设置可分配的虚拟C...	10	✎
会话超时时间	会话	默认为7200, 当前会话登录超过该会话...	7200	✎
物理机保留内存	KVM	默认为1G, 用于设置所有KVM物理主机...	27G	✎
云主机缓存模式	KVM	默认为none, 云主机缓存模式设置, 可...	none	✎
云主机CPU模式	KVM	默认为none, 选择云主机的CPU类型是...	none	✎

1.2.7 集群升级

1.2.7.1 内嵌服务升级

管理员获得新的内嵌服务套件后, 可用于升级当前的超融合环境的管理节点主机。

请管理员将新的内嵌服务套件上传到管理节点主机的目录`/root/`。上传完成后, 登录管理节点主机, 执行命令:

```
# 登陆管理节点主机,执行内嵌服务升级操作
[root@managementnode ~]# chmod +x ZStack for Alibaba Cloud-HA-Guest-Agent-2.x.bin
[root@managementnode ~]# ./ZStack for Alibaba Cloud-HA-Guest-Agent-2.x.bin -i
# 升级完成后,检查版本信息
[root@managementnode ~]# zin -v
App Version: 1.5.6.0
Build Date: Fri Jul 28 10:37:49 GMT-8 2017
Branch Name: (detached
Commit ID: c7dc57a012f5bae01eafaea7455bf5784007c172
```

1.2.7.2 高可用升级

管理员获得新的高可用服务套件后, 可用于升级当前的超融合环境。请管理员登录其中一个超融合节点执行命令:

```
# 导出当前的高可用集群配置信息
[root@localhost ~]# zsha export-config
output config json file to /tmp/current.config.json

# 若新版本升级时, 需要增加/删除相关配置信息, 请依照新版本的编写配置说明进行修改, 再执行zsha export-config进行导出

# 关闭集群内管理节点主机以及所有zsha服务
[root@localhost ~]# zsha stop

# 使用新的高可用套件升级管理三个节点
[root@localhost ~]# chmod +x ZStack for Alibaba Cloud-HA-Suite-2.x.bin
[root@localhost ~]# ./ZStack for Alibaba Cloud-HA-Suite-2.x.bin install \
```

```
-p password -c /tmp/current.config.json
```

1.2.8 管理节点迁移

ZStack for Alibaba Cloud超融合集群在正常运行环境下，若管理员需要临时关闭一个超融合节点进行维护，需在管理界面设定物理主机为维护模式。

然后，根据该超融合节点是否运行管理节点主机，进行对应操作。

- 若该超融合节点并非运行管理节点主机，则可以执行关机维护操作。
- 若该超融合节点运行管理节点主机，则需要执行手动热迁移操作。

```
# 登陆到其中一个超融合节点，执行对管理节点主机迁移至目标超融合节点
[root@localhost ~]# zsha migrate target ip
# 用户也可直接填写目标超融合节点的主机名称执行迁移操作
[root@localhost ~]# zsha migrate target hostname
```



说明：

- 其中`target ip`为超融合集群中一个正常超融合节点的 IP 地址。该地址必须是超融合集群内的 IP 地址，否则无法通过检查，将会报错并且不能执行迁移操作。
- 同样的，其中`target hostname`必须为超融合集群中一个正常超融合节点的主机名称。
- 如果目标超融合节点指定为正在运行管理节点主机的 HA 节点，将会报错，报错格式如下：

```
错误:Requested operation is not valid: domain 'ZStack Management Node VM' is already active
```

- 迁移过程需消耗一定时间，迁移过程将会显示进度状态。如果网络为万兆网络且状况良好，一般可在数秒内完成。

1.2.9 配置更新

管理节点高可用机制将全局配置信息存储在 Consul KV 数据库中，如果需要修改诸如管理节点主机配置、管理节点主机网络配置时，执行以下命令即可：

```
# 导出当前的高可用集群配置信息
[root@localhost ~]# zsha export-config
# 修改导出文件的配置
# 更新三个超融合节点的配置信息
[root@localhost ~]# zsha import-config /tmp/current.config.json
# 配置更新后,请登录管理节点主机执行关机操作,高可用服务将会触发启动
[root@managementnode ~]# halt -p
```



说明：

- 修改后的/tmp/current.config.json格式，应与原配置文件格式相同。
- 如需修改超融合节点的IP地址，请按照高可用套件的安装过程执行。
- **MemorySizeinGB**和**CPU**管理节点主机性能配置可以按需修改，但内存不可低于8G，CPU不可低于4核心。
- **Bridge**为三个超融合节点提供给管理节点主机使用的网桥名称。只有当该网桥名称发生变化时，才应修改此参数。
- **NetworkInterface**为管理节点主机内部的网络配置信息，可以按需修改。
- 由于管理节点镜像封装了qemu-guest-agent服务，管理节点主机支持修改root密码。若管理员忘记当前root密码，可通过执行reset-MNVM-password命令重置root密码为系统默认值（即password），但前提是**AllowResetRootPassword**需设置为true。

- 假定管理员使用原root密码登录管理节点主机，希望修改root密码：

```
[root@managementnode ~]# passwd root
Changing password for user root.
New password:
Retype new password:
passwd: all authentication tokens updated successfully.
```

- 假定管理员忘记当前root密码，希望重置root密码为系统默认值（即password）：

```
# 登录到其中一个超融合节点，导出当前的高可用集群配置信息
[root@localhost ~]# zsha export-config
output config json file to /tmp/current.config.json

# 修改导出文件的配置,将AllowResetRootPassword设置为true
[root@localhost ~]# vim /tmp/current.config.json
...
{
  ...
  "AllowResetRootPassword": true //默认为false;设为true时,管理节点主机才允许重置root密码
}
...

# 更新三个超融合节点的配置信息,按照高可用套件的安装过程执行
[root@localhost ~]# zsha install -p password -c /tmp/current.config.json

# 执行以下命令,重置root密码为系统默认值（即password）
[root@localhost ~]# zsha reset-MNVM-password
Reset MNVM root password on nodex IP //nodex IP为MNVM所在节点的IP
```

```
Password set successfully for root in 173de36b-4a5d-458e-a5e3-9a99584551ef
```

1.3 其他操作

1.3.1 卸载操作

管理员若需卸载ZStack for Alibaba Cloud超融合高可用套件，可执行以下命令：

```
# 登陆到其中一个节点,执行关闭高可用套件服务并停止管理节点主机
[root@localhost ~]# zsha stop
# 确认管理节点主机已关闭,避免导致数据损坏或丢失
[root@localhost ~]# zsha status
# 登陆到三个节点后,执行卸载高可用相关服务
[root@localhost ~]# zsha uninstall
```

执行卸载后，管理节点主机数据仍然保留在Ceph分布式存储里。管理员可保留数据，亦可删除该镜像。

若需删除镜像数据和存储池，可参考导入镜像章节关于删除数据的内容。

1.3.2 日志输出

ZStack for Alibaba Cloud超融合管理套件部署后，其运行状态日志信息存放在目录`/var/log/z/`。

高可用套件守护服务在运行状态时，会产生详细的日志内容，并进行按天和容量进行切割。管理员需关注管理节点主机的系统日志分区（`/var/log/`）空间状况，保证预留一定量的存储空间。根据运行估算，每年产生日志10+GB。

2 高可用测试与恢复

2.1 计划运维

2.1.1 单节点需要维护

2.1.1.1 单节点关闭

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。在正常运行环境下，若管理人员需要临时关闭某个超融合节点进行维护。

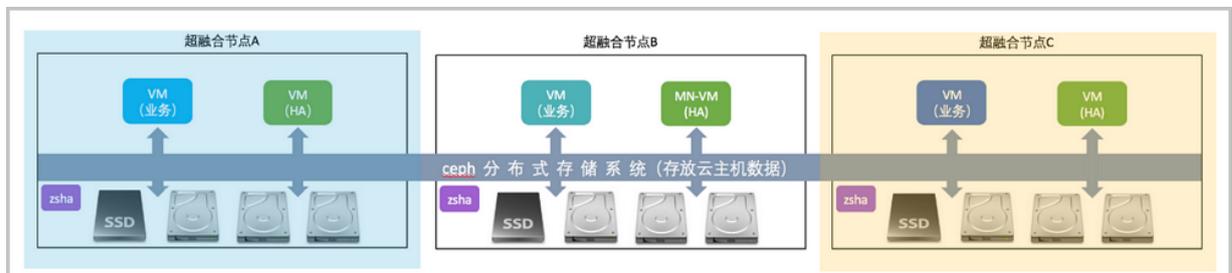
- 该超融合节点运行非管理节点主机。
- 该超融合节点运行管理节点主机。

2.1.1.1.1 该超融合节点运行非管理节点主机

背景信息

需要关闭节点A或C进行维护，如图 2-1: 超融合节点示例图所示：

图 2-1: 超融合节点示例图



以关闭超融合节点A为例。

操作步骤

1. 进入物理主机主机A的维护模式。

通过账户admin登录到ZStack for Alibaba Cloud云管理平台，将节点A的物理主机进入维护模式。

如图 2-2: 维护物理主机A所示：

图 2-2: 维护物理主机A

名称	物理 IP	删除	集群	启用状态	就绪状态	创建日期
物理机C	172.20.14.251		Cluster-1	启用	已连接	2017-07-28 16:01:42
物理机B	172.20.14.215		Cluster-1	启用	已连接	2017-07-28 16:01:15
<input checked="" type="checkbox"/> 物理机A	172.20.13.239		Cluster-1	启用	已连接	2017-07-28 15:58:52

2. 设置Ceph进入维护模式。

```
#Ceph社区进入维护模式。
[root@Ceph-1 ~]# for i in noout nobackfill norecover;do ceph osd set $i;done
```

3. 对接点A进行shutdown关机操作。

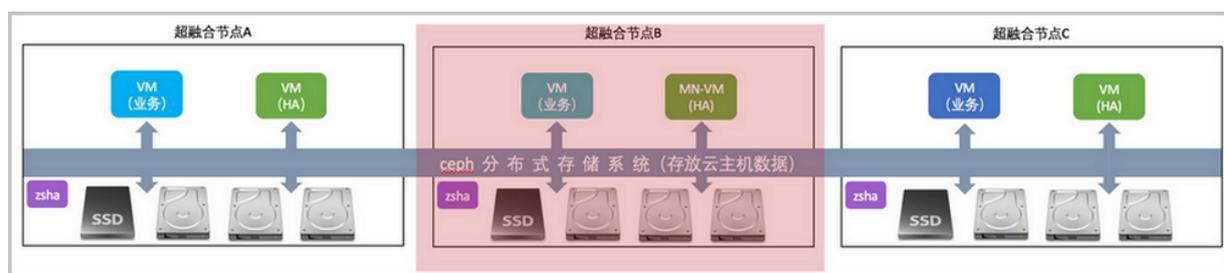
4. 对节点A下电后进行维护。

2.1.1.1.2 该超融合节点运行管理节点主机

背景信息

需要关闭超融合节点B进行维护，如图 2-3: 超融合节点示例图所示：

图 2-3: 超融合节点示例图



操作步骤

1. 管理员进行手动操作，将管理节点迁移至其他健康节点A或C。

将管理节点主机迁移至目标融合节点。

```
# 登陆管理节点主机,执行对管理节点主机迁移至目标节点
[root@managementnode ~]# zsha migrate target ip

# 用户也可直接填写目标节点的主机名称执行迁移操作
[root@managementnode ~]# zsha migrate target hostname
```

2. 进入物理主机B的维护模式。

通过账户admin登录到ZStack for Alibaba Cloud云管理平台，将节点B物理主机进入维护模式。

如图 2-4: 维护物理主机B所示：

图 2-4: 维护物理主机B



名称	物理 IP	删除	集群	启用状态	就绪状态	创建日期
物理机C	172.20.14.251		Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:42
<input checked="" type="checkbox"/> 物理机B	172.20.14.215		Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:15
物理机A	172.20.13.239		Cluster-1	• 启用	◦ 已连接	2017-07-28 15:58:52

3. 设置Ceph进入维护模式。

```
#Ceph社区进入维护模式。
[root@Ceph-1 ~]# for i in noout nobackfill norecover;do ceph osd set $i;done
```

4. 对节点B进行shutdown关机操作。

5. 对节点B下电后进行维护。

2.1.1.2 单节点启动

背景信息

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。若某个超融合节点临时关闭维护完成后，需要启动该节点。

操作步骤

1. 未启动的该超融合节点通电后，通过手动或IPMI启动服务器。

2. 等待超融合节点启动，成功引导操作系统。

3. 设置Ceph退出维护模式。

```
#Ceph社区版退出维护模式。
[root@Ceph-1 ~]# for i in noout nobackfill norecover;do ceph osd unset $i;done
```

4. 设置物理节点为启用状态。

通过admin账户登录到ZStack for Alibaba Cloud云管理平台，设置该节点为启用状态，如图 2-5: 启用物理主机所示：

图 2-5: 启用物理主机

<input type="checkbox"/>	名称	物理机IP	集群	启用状态	就绪状态	创建日期
<input type="checkbox"/>	物理机C	172.20.14.251	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:42
<input type="checkbox"/>	物理机B	172.20.14.215	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:15
<input checked="" type="checkbox"/>	物理机A	172.20.13.239	Cluster-1	• 维护模式	◦ 已连接	2017-07-28 15:58:52

2.1.2 三节点需要维护

2.1.2.1 三节点关闭

背景信息

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。若管理员需要临时关闭整个集群进行维护。

操作步骤

1. 关闭云主机高可用。

在ZStack for Alibaba Cloud专有云主菜单，点击**设置 > 全局设置**，将**云主机高可用全局开关**的值设置为**false**，停止高可用心跳检测，如图 2-6: 关闭云主机高可用所示：

图 2-6: 关闭云主机高可用

名称	类别	简介	值	操作
云主机高可用全局开关	高可用	默认为true, 用于设置云主机高可用功...	false	编辑
CPU超分率	物理机	默认为10, 主要用于设置可分配的虚拟...	10	编辑
会话超时时间	会话	默认为7200, 当前会话登录超过该会话...	7200	编辑
物理机保留内存	KVM	默认为1G, 用于设置所有KVM物理主机...	1G	编辑
云主机缓存模式	KVM	默认为none, 云主机缓存模式设置, 可...	none	编辑
云主机CPU模式	KVM	默认为none, 选择云主机的CPU类型是...	none	编辑
在线迁移	本地存储	默认为false, 本地存储在在线迁移的全局...	false	编辑
内存超分率	系统	默认为1.0. 如果物理内存为4G, 设...	1.0	编辑
主存储超分率	系统	默认为1.0. 如果主存储可用空间为2...	1.0	编辑

2. 将三台物理主机设置为维护模式。

将三台物理主机均设置为维护模式，批量停止云主机，如图 2-7: 维护物理主机集群所示：

图 2-7: 维护物理主机集群



名称	物理 IP	删除	集群	启用状态	就绪状态	创建日期
物理机C	172.20.14.251		Cluster-1	启用	已连接	2017-07-28 16:01:42
物理机B	172.20.14.215		Cluster-1	启用	已连接	2017-07-28 16:01:15
物理机A	172.20.13.239		Cluster-1	启用	已连接	2017-07-28 15:58:52

3. 通过ssh工具登陆到其中一个超融合节点，执行zsha stop命令：

登录到其中一个节点，执行关闭高可用套件服务并停止管理节点主机。

4. 设置Ceph集群进入维护模式。

```
#Ceph社区版进入维护模式：
[root@Ceph-1 ~]# for i in noout nobackfill norecover;do ceph osd set $i;done
```

5. 对三个超融合节点进行shutdown关机操作。

6. 三个超融合节点下电后，维护整个集群。

2.1.2.2 三节点启动

背景信息

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。若整个集群临时关闭维护完成后，需要启动该集群。

操作步骤

1. 三个超融合节点通电后，通过手动或IPMI启动服务器。
2. 操作系统启动引导完成，即可通过ssh工具登陆到三个超融合节点。
3. 集群退出维护模式。

检查Ceph存储集群，确认无误后，存储集群退出维护模式：

```
[root@Ceph-1 ~]# for i in noout nobackfill norecover;do ceph osd unset $i;done
```

4. 执行zsha start命令启动高可用服务套件。

登录到其中一个节点，执行启动高可用套件服务并启动ZStack for Alibaba Cloud管理节点主机。

5. 检查运行状态。

通过`zsha status`检查状态是否健康 并关注管理节点主机是否正常运行。

```
[root@localhost ~]# zsha status
>>>MN-VM Running On
172.20.14.215 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
172.20.13.239 : ceph-1
172.20.14.215 : ceph-2
172.20.14.251 : ceph-3
```

6. 设置超融合节点为启用状态。

通过账户admin登录到ZStack for Alibaba Cloud云管理平台，依次设置三个超融合节点为启用状态，如图 2-8: 启用超融合物理主机所示：

图 2-8: 启用超融合物理主机



<input checked="" type="checkbox"/>	名称	物理机IP	集群	启用状态	就绪状态	创建日期
<input checked="" type="checkbox"/>	物理机C	172.20.14.251	Cluster-1	• 维护模式	◦ 已连接	2017-07-28 16:01:42
<input checked="" type="checkbox"/>	物理机B	172.20.14.215	Cluster-1	• 维护模式	◦ 已连接	2017-07-28 16:01:15
<input checked="" type="checkbox"/>	物理机A	172.20.13.239	Cluster-1	• 维护模式	◦ 已连接	2017-07-28 15:58:52

7. 开启云主机高可用。

通过账户admin登录到ZStack for Alibaba Cloud云管理平台，进入设置 > 全局设置，将云主机高可用全局开关的值设置为true，停止高可用心跳检测，如图 2-9: 开启云主机高可用所示：

图 2-9: 开启云主机高可用

全局设置				
基本设置				
高级设置				
名称	类别	简介	值	操作
云主机高可用全局开关	高可用	默认为true, 用于设置云主机高可用功...	true	
CPU超分率	物理机	默认为10, 主要用于设置可分配的虚拟...	10	
会话超时时间	会话	默认为7200, 当前会话登录超过该会话...	7200	
物理机保留内存	KVM	默认为1G, 用于设置所有KVM物理主机...	1G	
云主机缓存模式	KVM	默认为none, 云主机缓存模式设置, 可...	none	
云主机CPU模式	KVM	默认为none, 选择云主机的CPU类型是...	none	
在线迁移	本地存储	默认为false, 本地存储在线迁移的全局...	false	
内存超分率	系统	默认为1.0. 如果物理内存为4G, 设...	1.0	
主存储超分率	系统	默认为1.0. 如果主存储可用空间为2...	1.0	

**说明：**

- 高可用级别为**NeverStop**的云主机将会自动恢复。
- 高可用级别为**None**的云主机需管理员手动启动，或通知云主机所有者执行启动。

2.2 异常处理

如果节点出现异常掉电/断网，请参考本章[异常处理](#)进行故障恢复。如果节点出现损坏（系统崩溃、磁盘损坏、数据丢失），请参考下章[节点修复](#)进行故障修复，如果依然无法解决问题，请联系官方技术支持团队。。

2.2.1 单节点异常处理

2.2.1.1 单节点异常掉电

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。在正常运行环境下，若其中某个超融合节点突然断电。

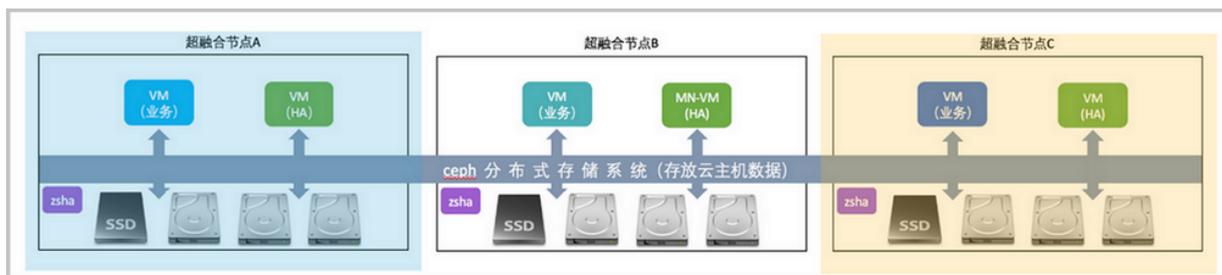
- 该超融合节点运行非管理节点主机。
- 该超融合节点运行管理节点主机。

2.2.1.1.1 该超融合节点运行非管理节点主机

背景信息

超融合节点A或C异常掉电需恢复，如图 2-10: 超融合节点示例图所示：

图 2-10: 超融合节点示例图



操作步骤

1. 该异常超融合节点A上电后，通过手动或IPMI方式启动服务器。
2. 自动引导加载操作系统，确认是否启动成功。
3. 查看集群状态。

检查Ceph存储集群状况，确认存储集群状态健康。

4. 检查异常节点高可用状态。

执行 `zsha status` 命令，检查该异常节点的高可用服务是否正常。

```
[root@localhost ~]# zsha status
>>>MN-VM Running On
172.20.14.215 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
172.20.13.239 : ceph-1
172.20.14.215 : ceph-2
172.20.14.251 : ceph-3
```

5. 检查物理主机连接状态。

通过账号admin登陆ZStack for Alibaba Cloud云管理平台，检查该异常物理主机连接状态是否正常，如图 2-11: 物理主机连接状态界面所示：

图 2-11: 物理主机连接状态界面

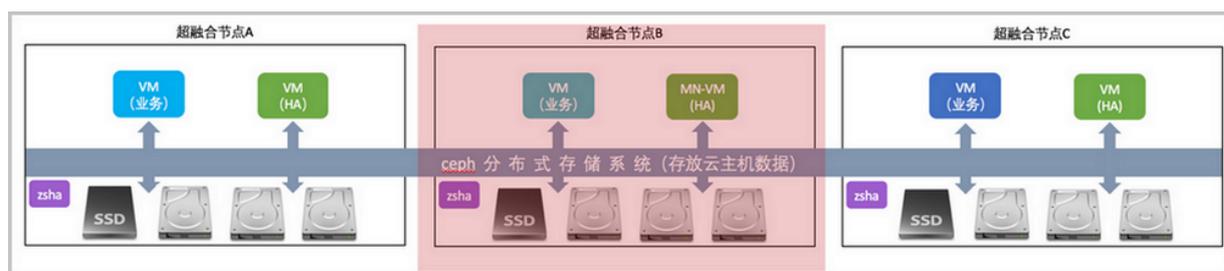
<input type="checkbox"/>	名称	物理机IP	集群	启用状态	就绪状态	创建日期
<input type="checkbox"/>	物理机C	172.20.14.251	Cluster-1	● 启用	○ 已连接	2017-07-28 16:01:42
<input type="checkbox"/>	物理机B	172.20.14.215	Cluster-1	● 启用	○ 已连接	2017-07-28 16:01:15
<input type="checkbox"/>	物理机A	172.20.13.239	Cluster-1	● 启用	○ 已连接	2017-07-28 15:58:52

2.2.1.1.2 该超融合节点运行管理节点主机

背景信息

超融合节点B异常掉电需恢复，如图 2-12: 超融合节点示例图所示：

图 2-12: 超融合节点示例图



操作步骤

1. 检查管理节点主机高可用是否已触发。

通过ssh工具登陆到正在运行的超融合节点A或C，执行**zsha status**命令检查管理节点主机是否正常运行，若正常运行，代表高可用套件已触发高可用切换。

```
[root@localhost ~]# zsha status
>>>MN-VM RUNNING ON
172.20.14.251 : running

>>>LAST START VM RECORD
```

2. 该异常超融合节点B上电后，通过手动或IPMI方式启动服务器。
3. 自动引导加载操作系统，确认是否启动成功。
4. 检查Ceph存储集群状况，确认存储集群状态健康。
5. 检查异常节点高可用状态。

执行`zsha status`命令，检查该异常节点的高可用服务是否正常。

```
[root@localhost ~]# zsha status
>>>MN-VM Running On
172.20.14.251 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
172.20.13.239 : ceph-1
172.20.14.215 : ceph-2
172.20.14.251 : ceph-3
```

2.2.1.2 单节点网络异常

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。在正常运行环境下，若其中某个超融合节点网络异常。

- 该超融合节点运行非管理节点主机。
- 该超融合节点运行管理节点主机。



说明：

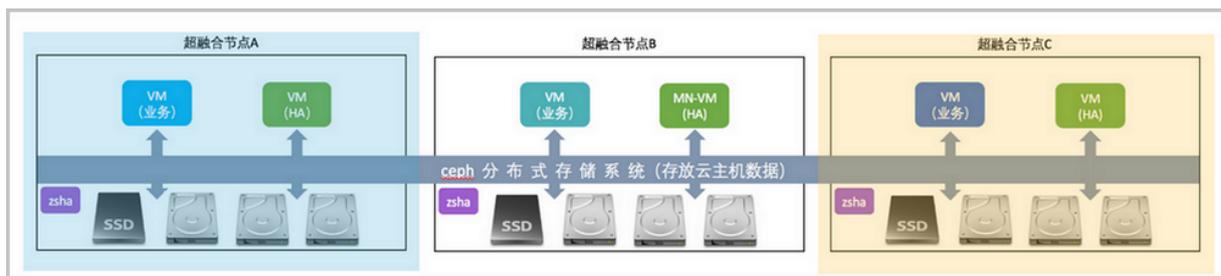
- 本手册中“网络故障”，特指管理与存储网络的故障，云主机数据网络故障，不会触发高可用切换。
- 本手册中“网络故障”，特指两个链路均故障的情况，单链路故障，不会影响超融合节点之间的通信。

2.2.1.2.1 该超融合节点运行非管理节点主机

背景信息

超融合节点A或C网络异常需恢复，如图 2-13: 超融合节点示例图所示：

图 2-13: 超融合节点示例图



以超融合节点A网络异常为例。

操作步骤

1. 检查交换机运行状态，确认网络异常节点A的端口指示灯常亮或闪烁。
2. 通过ssh登录其他正在运行的超融合节点B或C上，检查Ceph存储集群运行状态。
3. 检查网络异常节点A高可用状态。

执行 `zsha status`，检查网络异常节点A的高可用服务是否正常：

```
[root@localhost ~]# zsha status
>>>MN-VM Running On
172.20.14.215 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
172.20.13.239 : ceph-1
172.20.14.215 : ceph-2
172.20.14.251 : ceph-3
```

4. 查看物理主机连接状态。

通过账户admin登录ZStack for Alibaba Cloud云管理平台检查网络异常物理主机连接状态，如果连接正常，则网络异常节点A已恢复，如图 2-14: 物理主机连接状态界面所示：

图 2-14: 物理主机连接状态界面

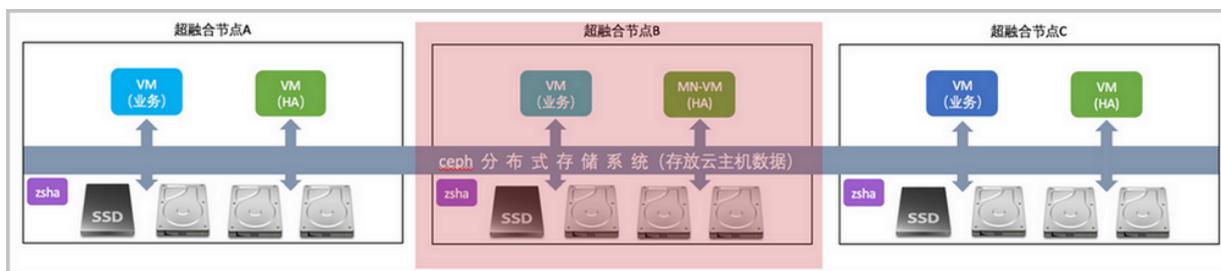
<input type="checkbox"/>	名称	物理机IP	集群	启用状态	就绪状态	创建日期
<input type="checkbox"/>	物理机C	172.20.14.251	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:42
<input type="checkbox"/>	物理机B	172.20.14.215	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:15
<input type="checkbox"/>	物理机A	172.20.13.239	Cluster-1	• 启用	◦ 已连接	2017-07-28 15:58:52

2.2.1.2.2 该超融合节点运行管理节点主机

背景信息

超融合节点B网路异常需恢复，如图 2-15: 超融合节点示例图所示：

图 2-15: 超融合节点示例图



操作步骤

1. 检查管理节点主机高可用是否已触发。

通过ssh工具登陆到正在运行的超融合节点A或C上，执行**zsha status**检查管理节点主机是否正常运行，若正常运行，代表高可用套件已触发高可用切换。

```
[root@localhost ~]# zsha status
>>>MN-VM RUNNING ON
172.20.14.251 : running

>>>LAST START VM RECORD
```

2. 检查交换机运行状态，确认网络异常节点B的端口指示灯常亮或闪烁。
3. 检查网络异常节点高可用状态。

执行**zsha status**检查网络异常节点的高可用服务是否正常。

```
[root@localhost ~]# zsha status
>>>MN-VM Running On
172.20.14.251 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
172.20.13.239 : ceph-1
172.20.14.215 : ceph-2
172.20.14.251 : ceph-3
```

4. 检查物理主机连接状态。

通过账户admin登录ZStack for Alibaba Cloud云管理平台检查网络异常物理主机连接状态，如果连接正常，则网络异常节点B已恢复，如图 2-16: 物理主机连接状态界面所示：

图 2-16: 物理主机连接状态界面



<input type="checkbox"/>	名称	物理机IP	集群	启用状态	就绪状态	创建日期
<input type="checkbox"/>	物理机C	172.20.14.251	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:42
<input type="checkbox"/>	物理机B	172.20.14.215	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:15
<input type="checkbox"/>	物理机A	172.20.13.239	Cluster-1	• 启用	◦ 已连接	2017-07-28 15:58:52

2.2.2 两节点异常处理

背景信息

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。在正常运行环境下，若其中两个超融合节点均断电或网络异常。

操作步骤

1. 两异常融合节点上电后通过手动或IPMI方式启动。
2. 两异常融合节点加载操作系统，确认是否启动成功。
3. 检查集群状态。

检查Ceph存储集群状况，确认存储集群状态健康。

4. 两异常节点高可用状态。

通过`zsha status`检查两异常节点的高可用服务是否正常，检查管理节点主机是否正常运行。

```
[root@localhost ~]# zsha status
>>>MN-VM Running On
172.20.14.215 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
172.20.13.239 : ceph-1
172.20.14.215 : ceph-2
172.20.14.251 : ceph-3
```

5. 查看物理主机连接状态。

通过账号admin登陆ZStack for Alibaba Cloud云管理平台，检查两异常物理主机连接状态，如果连接正常，则两网络异常节点已恢复，如图 2-17: 物理主机连接状态界面所示：

图 2-17: 物理主机连接状态界面



<input type="checkbox"/>	名称	物理机IP	集群	启用状态	就绪状态	创建日期
<input type="checkbox"/>	物理机C	172.20.14.251	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:42
<input type="checkbox"/>	物理机B	172.20.14.215	Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:15
<input type="checkbox"/>	物理机A	172.20.13.239	Cluster-1	• 启用	◦ 已连接	2017-07-28 15:58:52

2.3 节点修复

2.3.1 单节点故障修复

背景信息

ZStack for Alibaba Cloud超融合集群由三个超融合节点A、B、C构建。若其中某个超融合节点损坏后需要执行修复。

操作步骤

1. 调配备用服务器，使得硬件规格与故障节点相近。
2. 参考本文档[安装与部署](#)章节，安装基础操作系统。安装完成后，配置root的密码和网络信息与故障节点一致。
3. 修复Ceph存储集群。

参考Ceph相关手册，修复Ceph存储集群，包括MON和OSD服务。

4. 参考本文档[安装与部署](#)章节，安装高可用套件。

在健康的超融合节点，导出当前的高可用集群配置信息。

```
[root@localhost ~]# zsha export-config  
output config json file to /tmp/current.config.json
```



说明：

若新版本升级时，需要增加/删除相关配置信息，请依照新版本发布说明。

5. 查看管理节点主机状态。

通过**zsha status**检查高可用服务套件运行状态，以及管理节点主机状态。

```
[root@localhost ~]# zsha status  
>>>MN-VM Running On  
172.20.14.215 : running
```

```
>>>Last MN-VM Start Record
```

```
>>>ZStack HA Services Running On  
172.20.13.239 : ceph-1  
172.20.14.215 : ceph-2  
172.20.14.251 : ceph-3
```

6. 对置换节点进行重连操作。

通过账户admin登录ZStack for Alibaba Cloud云管理平台，对置换节点执行重连操作，如图 2-18: 重连超融合物理主机界面所示：

图 2-18: 重连超融合物理主机界面

<input type="checkbox"/>	名称	物理	删除	集群	启用状态	就绪状态	创建日期
<input type="checkbox"/>	物理机C	172.20.14.251		Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:42
<input type="checkbox"/>	物理机B	172.20.14.215		Cluster-1	• 启用	◦ 已连接	2017-07-28 16:01:15
<input checked="" type="checkbox"/>	物理机A	172.20.13.239		Cluster-1	• 启用	◦ 连接中	2017-07-28 15:58:52

2.3.2 两节点故障无法修复

由于ZStack for Alibaba Cloud企业版超融合产品三节点模型，只支持容忍一个节点发生故障，对于两节点损坏情况无法修复。若需满足两个节点可修复的场景，则需构建至少5个节点的超融合集群。

3 命令行使用手册

3.1 简介

zsha是ZStack for Alibaba Cloud针对管理节点高可用场景设计的命令，帮助用户快速完成该场景下的多种操作。

zsha下有多条子命令，本手册将基于超融合产品对**zsha**每条子命令的作用和使用方法进行说明。

3.2 -h 帮助内容

描述

显示帮助，可查看**zsha**全部子命令。

```
[root@localhost ~]# zsha -h

zsha is a proprietary, high-availability suite that protects ZStack management services from
running online for long periods of time.

usage:
  zsha [-h] [-v]

All sub-commands
install [-p] ROOT-PASSWORD [-c] CONFIGURATION-FILE-PATH
    Install Mode, Need a root password for all HA nodes, Need a configuration file path
uninstall --yes-i-really-really-mean-it
    uninstall from all nodes
check-config CONFIGURATION-FILE-PATH
    check configuration file
config-sample [CephConf|FileConf]
    output sample config json file to /tmp/sample.CephConf.config.json
    or /tmp/sample.FileConf.config.json
export-config    output current config json file to /tmp/current.config.json
import-config CONFIGURATION-FILE-PATH
    update configuration to HA cluster
migrate xxx.xxx.xxx.xxx
    migrate to target HA host, ip or hostname should be provided
status          show status
status-conf     show status with current configuration
stop           close MN-VM and all zs-ha services in cluster, consul will continue to operate
start          start all zs-ha services in cluster, if consul service is stopped, it will also be started

optional arguments:
-h            show this help message
```

```
-v show version, print execution details
```

3.3 -v 版本信息

描述

查看版本信息，包括版本号、编译日期、Branch Name和Commit ID。

```
[root@localhost ~]# zsha -v
App Version: 1.5.6.0
Build Date: Fri Jul 28 10:37:49 GMT-8 2017
Branch Name: (detached)
Commit ID: c7dc57a012f5bae01eafaea7455bf5784007c172
```

3.4 check-config 配置检查

描述

检查配置文件格式和内容，该检查在安装和更新时也会被执行。

其中**ConsoleProxyOverriddenIP**用户可按需填写，也可不填，仅检查是否为有效的IPv4格式。其他各个字段必须填写，相应检查格式和内容。

管理节点主机的CPU数目和内存必须满足最低要求，可按需增加。

使用方法

参数	介绍	示例
CONFIGURATION-FILE-PATH	用户需填写完整的配置文件路径	<code>zsha check-config /tmp/current.config.json</code>

```
[root@localhost ~]# zsha check-config /tmp/current.config.json
config file check pass
```

3.5 install -p -c 安装命令

描述

安装命令，对HA集群内所有节点执行安装操作，要求HA集群内所有节点使用同一个password，用户需填写该password，以及配置文件路径。

使用方法

参数	介绍	示例
-p ROOT-PASSWORD	要求HA集群内所有节点使用同一个password，用户需填写该password	<code>./ZStack for Alibaba Cloud -HA-Suite-2.x.bin install -p password</code>
-c CONFIGURATION	用户需填写完整的配置文件路径	<code>./ZStack for Alibaba Cloud -HA-Suite-2.x.bin install -c /tmp/config.json</code>

```
[root@localhost ~]# ./ZStack-Enterprise-HA-Suite-2.x.bin install -p password -c /tmp/config.json
config file check pass
All nodes connect success.
check environment for 172.20.197.242
RBD File zstack/mnvm.img exists
Settings for br_zsn0:
Link detected: yes

check environment for 172.20.198.176
RBD File zstack/mnvm.img exists
Settings for br_zsn0:
Link detected: yes

check environment for 172.20.197.202
RBD File zstack/mnvm.img exists
Settings for br_zsn0:
Link detected: yes

Stop Service consul&z for: 172.20.197.242
Stop Service consul&z for: 172.20.198.176
Stop Service consul&z for: 172.20.197.202
Send /tmp/zstack-ha-installer/zsha to 172.20.197.242
Send /tmp/zstack-ha-installer/config.json to 172.20.197.242
Send /tmp/zstack-ha-installer/zsha to 172.20.198.176
Send /tmp/zstack-ha-installer/config.json to 172.20.198.176
Send /tmp/zstack-ha-installer/zsha to 172.20.197.202
Send /tmp/zstack-ha-installer/config.json to 172.20.197.202
Generate ceph auth finished.
Install for: 172.20.197.242
config file check pass
RBD File zstack/mnvm.img exists
Settings for br_zsn0:
Link detected: yes

Install bash auto-completion
Set iptables finished.
Install to Current Node Finished.

Install for: 172.20.198.176
config file check pass
RBD File zstack/mnvm.img exists
Settings for br_zsn0:
Link detected: yes

Install bash auto-completion
Set iptables finished.
```

```
Install to Current Node Finished.

Install for: 172.20.197.202
config file check pass
RBD File zstack/mnvm.img exists
Settings for br_zsn0:
Link detected: yes

Install bash auto-completion
Set iptables finished.
Install to Current Node Finished.

Upsert configuration to consul kv finished.

Install to All Nodes Finished.
Now, You can monitor zsha's status by using command "zsha status".
More commands are shown in "zsha -h".
```

3.6 uninstall 卸载命令

描述

在HA集群中，对当前节点执行**zsha**卸载操作。建议用户自行备份当前节点内数据库再执行此操作。

```
[root@localhost ~]# zsha uninstall
WARNING: this will *DESTROY* ZStack Management Node VM.
If you are *ABSOLUTELY CERTAIN* that is what you want, followed by --yes-i-really-really-mean-it.
[root@ceph-1 ~]# zsha uninstall --yes-i-really-really-mean-it
start clean
start shutdown cluster
Do for 172.20.197.242
Do for 172.20.198.176
Do for 172.20.197.202
shutdown cluster finished
Do for 172.20.197.242
Do for 172.20.198.176
Do for 172.20.197.202
clean finished
```

3.7 config-sample 样本配置生成

描述

生成样本配置文件并存放在临时路径/tmp/sample.CephConf.config.json或/tmp/sample.FileConf.config.json中。

使用方法

参数	介绍	示例
[CephConf FileConf]	用户根据实际情况选择配置文件类型	zsha config-sample CephConf 或zsha config-sample FileConf

```
[root@localhost ~]# zsha config-sample CephConf
output config json file to /tmp/sample.CephConf.config.json
```

```
[root@localhost ~]# zsha config-sample FileConf
output config json file to /tmp/sample.FileConf.config.json
```

3.8 export-config 当前配置生成

描述

生成当前配置文件并存放在临时路径/tmp/current.config.json。

```
[root@localhost ~]# zsha export-config
output config json file to /tmp/current.config.json
```

3.9 migrate 迁移命令

描述

将管理节点主机从当前运行主机迁移到目标HA主机。

使用方法

参数	介绍	示例
xxx.xxx.xxx.xxx	目标HA主机为HA集群中一个正常的HA节点，xxx.xxx.xxx.xxx为该目标主机的IP地址	zsha migrate xxx.xxx.xxx.xxx
hostname	目标HA主机为HA集群中一个正常的HA节点，hostname为该目标主机的名称	zsha migrate hostname

```
[root@localhost ~]# zsha migrate 10.0.119.130
Migrate from 10.0.134.172 to 10.0.119.130
Migration: [100 %]
```

```
[root@localhost ~]# zsha migrate ceph-3
Migrate from 10.0.119.130 to 10.0.178.146
```

```
Migration: [100 %]
```

3.10 import-config 配置升级

描述

对整个HA集群更新配置文件，需填写完整的配置文件路径。

使用方法

参数	介绍	示例
CONFIGURATION-FILE-PATH	用户需填写完整的配置文件路径	<code>zsha import-config /tmp/current.config.json</code>

```
[root@localhost ~]# zsha import-config /tmp/current.config.json
```

3.11 status 状态信息

描述

显示当前整个HA集群的状态，包括管理服务主机的当前运行状态和上次启动记录、集群中所有HA节点状态、Ceph状态、以及管理服务主机的IP地址。

```
[root@localhost ~]# zsha status
```

使用方法

参数	介绍	示例
-conf	显示整个HA集群状态的同时，显示当前配置情况	<code>zsha status-conf</code>

```
[root@localhost ~]# zsha status-conf
>>>MN-VM Running On
192.168.255.205 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
192.168.255.205 : ceph-1
192.168.255.221 : ceph-2
192.168.255.183 : ceph-3

>>>Consul Members
Node      Address           Status Type  Build Protocol DC
192.168.255.205 192.168.255.205:8301 alive  server 0.7.2 2    dc1
192.168.255.221 192.168.255.221:8301 alive  server 0.7.2 2    dc1
192.168.255.183 192.168.255.183:8301 alive  server 0.7.2 2    dc1
```

```
>>>Management Node IP
172.20.198.3

>>>Public Network IP
203.114.54.10

>>>Configuration
{
  "Node": [
    "192.168.255.205 ceph-1",
    "192.168.255.221 ceph-2",
    "192.168.255.183 ceph-3"
  ],
  "MemorySizeInGB": 8,
  "CPU": 4,
  "Type": "CephConf",
  "MonAddrs": [
    "192.168.255.205:6789",
    "192.168.255.221:6789",
    "192.168.255.183:6789"
  ],
  "ChronyServers": [
    "192.168.255.205",
    "192.168.255.221",
    "192.168.255.183"
  ],
  "PoolName": "zstack",
  "ImageFolder": "/storage/",
  "DNS": [
    "172.20.198.1",
    "223.5.5.5"
  ],
  "Network": [
    {
      "Bridge": "br_zsn0",
      "MacAddress": "02:98:54:4b:a5:c0",
      "Ipaddr": "172.20.198.3",
      "Netmask": "255.255.255.0",
      "Gateway": "172.20.198.1",
      "IsMgmt": true,
      "IsDefRoute": false
    },
    {
      "Bridge": "br_zsn1",
      "MacAddress": "f6:75:d8:5d:73:73",
      "Ipaddr": "203.114.54.10",
      "Netmask": "255.255.255.0",
      "Gateway": "203.114.54.1",
      "IsMgmt": false,
      "IsDefRoute": true
    }
  ],
  "ConsoleProxyOverriddenIP": ""
}

>>>Ceph Status
cluster 7eafa8a3-7318-4d93-86fe-4186e1509f15
health HEALTH_WARN
too many PGs per OSD (592 > max 300)
noout,noscrub,nodeep-scrub flag(s) set
```

```

monmap e1: 3 mons at {ceph-1=192.168.255.205:6789/0,ceph-2=192.168.255.221:6789/0,
ceph-3=192.168.255.183:6789/0}
election epoch 6, quorum 0,1,2 ceph-3,ceph-1,ceph-2
osdmap e33: 3 osds: 3 up, 3 in
flags noout,noscrub,nodeep-scrub
pgmap v25752: 592 pgs, 6 pools, 8328 MB data, 2152 objects
24660 MB used, 860 GB / 884 GB avail
592 active+clean
client io 10524 B/s rd, 3134 B/s wr, 26 op/s

>>>Checking Host I/O
done.

```

3.12 stop 集群关闭

描述

关闭集群内管理节点主机以及所有zsha服务，consul服务依然继续运行。

```

[root@localhost ~]# zsha stop
start shutdown cluster
Do for 192.168.255.205
Do for 192.168.255.221
Do for 192.168.255.183
shutdown cluster finished
[root@localhost ~]# zsha status
>>>MN-VM Running On

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On

>>>Consul Members
Node      Address           Status Type   Build Protocol DC
192.168.255.205 192.168.255.205:8301 alive  server 0.7.2 2    dc1
192.168.255.221 192.168.255.221:8301 alive  server 0.7.2 2    dc1
192.168.255.183 192.168.255.183:8301 alive  server 0.7.2 2    dc1

>>>Management Node IP
172.20.198.3

>>>Public Network IP
203.114.54.10

>>>Ceph Status
cluster 7eafa8a3-7318-4d93-86fe-4186e1509f15
health HEALTH_WARN
too many PGs per OSD (592 > max 300)
noout,noscrub,nodeep-scrub flag(s) set
monmap e1: 3 mons at {ceph-1=192.168.255.205:6789/0,ceph-2=192.168.255.221:6789/0,
ceph-3=192.168.255.183:6789/0}
election epoch 6, quorum 0,1,2 ceph-3,ceph-1,ceph-2
osdmap e33: 3 osds: 3 up, 3 in
flags noout,noscrub,nodeep-scrub
pgmap v25752: 592 pgs, 6 pools, 8328 MB data, 2152 objects
24660 MB used, 860 GB / 884 GB avail

```

```
592 active+clean
client io 10524 B/s rd, 3134 B/s wr, 26 op/s

>>>Checking Host I/O
done.
```

3.13 start 集群启动

描述

启动集群内所有zsha服务，对于处于停止状态的consul服务，也将被启动。

```
[root@localhost ~]# zsha start
Do for 192.168.255.205
Do for 192.168.255.221
Do for 192.168.255.183
[root@localhost ~]# zsha status

>>>MN-VM Running On
192.168.255.205 : running

>>>Last MN-VM Start Record

>>>ZStack HA Services Running On
192.168.255.205 : ceph-1
192.168.255.221 : ceph-2
192.168.255.183 : ceph-3

>>>Consul Members
Node      Address           Status Type  Build Protocol DC
192.168.255.205 192.168.255.205:8301 alive  server 0.7.2 2    dc1
192.168.255.221 192.168.255.221:8301 alive  server 0.7.2 2    dc1
192.168.255.183 192.168.255.183:8301 alive  server 0.7.2 2    dc1

>>>Management Node IP
172.20.198.3

>>>Public Network IP
203.114.54.10

>>>Ceph Status
cluster 7eafa8a3-7318-4d93-86fe-4186e1509f15
health HEALTH_WARN
  too many PGs per OSD (592 > max 300)
  noout,noscrub,nodeep-scrub flag(s) set
  monmap e1: 3 mons at {ceph-1=192.168.255.205:6789/0,ceph-2=192.168.255.221:6789/0,ceph-3=192.168.255.183:6789/0}
  election epoch 6, quorum 0,1,2 ceph-3,ceph-1,ceph-2
  osdmap e33: 3 osds: 3 up, 3 in
  flags noout,noscrub,nodeep-scrub
  pgmap v25752: 592 pgs, 6 pools, 8328 MB data, 2152 objects
  24660 MB used, 860 GB / 884 GB avail
  592 active+clean
  client io 10524 B/s rd, 3134 B/s wr, 26 op/s

>>>Checking Host I/O
```

```
done.
```

3.14 reset-MNVM-password 管理节点主机重置root密码

描述

管理节点主机重置root密码为系统默认值（即password），但前提是配置文件中的**AllowResetRootPassword**需设置为true。

该命令并未显示在[-h 帮助内容](#)中，仅供特殊情况使用。

```
[root@localhost ~]# zsha reset-MNVM-password
Reset MNVM root password on nodex IP //nodex IP为MNVM所在节点的IP
Password set successfully for root in 173de36b-4a5d-458e-a5e3-9a99584551ef
```

专有云术语表

区域 (Zone)

ZStack中最大的一个资源定义，包括集群、二层网络、主存储等资源。

集群 (Cluster)

一个集群是类似物理主机 (Host) 组成的逻辑组。在同一个集群中的物理主机必须安装相同的操作系统 (虚拟机管理程序, Hypervisor)，拥有相同的二层网络连接，可以访问相同的主存储。在实际的数据中心，一个集群通常对应一个机架 (Rack)。

管理节点 (Management Node)

安装系统的物理主机，提供UI管理、云平台部署功能。

计算节点 (Compute Node)

也称之为物理主机 (或物理机)，为云主机实例提供计算、网络、存储等资源的物理主机。

主存储 (Primary Storage)

用于存储云主机磁盘文件的存储服务器。支持本地存储、NFS、Ceph、Shared Mount Point、Shared Block等类型。

镜像服务器 (Backup Storage)

也称之为备份存储服务器，主要用于保存镜像模板文件。建议单独部署镜像服务器。

镜像仓库 (Image Store)

镜像服务器的一种类型，可以为正在运行的云主机快速创建镜像，高效管理云主机镜像的版本变迁以及发布，实现快速上传、下载镜像，镜像快照，以及导出镜像的操作。

云主机 (VM Instance)

运行在物理机上的虚拟机实例，具有独立的IP地址，可以访问公共网络，运行应用服务。

镜像 (Image)

云主机或云盘使用的镜像模板文件，镜像模板包括系统云盘镜像和数据云盘镜像。

云盘 (Volume)

云主机的数据盘，给云主机提供额外的存储空间，共享云盘可挂载到一个或多个云主机共同使用。

计算规格 (Instance Offering)

启动云主机涉及到的CPU数量、内存、网络设置等规格定义。

云盘规格 (Disk Offering)

创建云盘容量大小的规格定义。

二层网络 (L2 Network)

二层网络对应于一个二层广播域，进行二层相关的隔离。一般用物理网络的设备名称标识。

三层网络 (L3 Network)

云主机使用的网络配置，包括IP地址范围、网关、DNS等。

公有网络 (Public Network)

由因特网信息中心分配的公有IP地址或者可以连接到外部互联网的IP地址。

私有网络 (Private Network)

云主机连接和使用的内部网络。

L2NoVlanNetwork

物理主机的网络连接不采用Vlan设置。

L2VlanNetwork

物理主机节点的网络连接采用Vlan设置，Vlan需要在交换机端提前进行设置。

VXLAN网络池 (VXLAN Network Pool)

VXLAN网络中的 Underlay 网络，一个 VXLAN 网络池可以创建多个 VXLAN Overlay 网络 (即 VXLAN 网络)，这些 Overlay 网络运行在同一组 Underlay 网络设施上。

VXLAN网络 (VXLAN)

使用 VXLAN 协议封装的二层网络，单个 VXLAN 网络需从属于一个大的 VXLAN 网络池，不同 VXLAN 网络间相互二层隔离。

云路由 (vRouter)

云路由通过定制的Linux云主机来实现的多种网络服务。

安全组 (Security Group)

针对云主机进行第三层网络的防火墙控制，对IP地址、网络包类型或网络包流向等可以设置不同的安全规则。

弹性IP (EIP)

公有网络接入到私有网络的IP地址。

快照 (Snapshot)

某一个时间点上某一个磁盘的数据备份。包括自动快照和手动快照两种类型。

混合云术语表

访问密钥 (AccessKey)

用于调用阿里云API或大河云联API的唯一凭证，AccessKey包括AccessKeyID（用于标识用户）和AccessKeySecret（用于验证用户密钥）。

数据中心 (Data Center)

包含阿里云的地域和可用区等地域资源，用于匹配阿里云资源的地域属性。

地域 (Region)

物理的数据中心，划分地区的基本单位，ZStack混合云的地域对应了阿里云端的地域。

可用区 (Identity Zone)

在同一地域内，电力和网络互相独立的物理区域，ZStack混合云的可用区对应了阿里云端的可用区 (Zone)。

存储空间 (Bucket)

用于存储对象 (Object) 的容器，ZStack使用对象存储 (OSS) 里的Bucket来上传镜像文件。

ECS云主机 (Elastic Compute Service)

阿里云端创建的ECS实例，可在ZStack混合云界面进行ECS云主机生命周期的管理。

专有网络VPC (Virtual Private Cloud)

用户基于阿里云构建的一个隔离的网络环境，不同的专有网络之间逻辑上彻底隔离。

虚拟交换机 (VSwitch)

组成专有网络VPC的基础网络设备，可以连接不同的云产品实例。ZStack混合云的虚拟交换机对应了阿里云VPC下的虚拟交换机。

虚拟路由器 (VRouter)

专有网络VPC的枢纽，可以连接专有网络的各个虚拟交换机，同时也是连接专有网络与其它网络的网关设备。ZStack支持查看VPC下的虚拟路由器。

路由表 (Route Table)

虚拟路由器上管理路由条目的列表。

路由条目 (Route Entry)

路由表中的每一项是一条路由条目。路由条目定义了通向指定目标网段的网络流量的下一跳地址。

路由条目包括系统路由和自定义路由两种类型。ZStack支持自定义类型的路由条目。

安全组 (Security Group)

针对云主机进行第三层网络的防火墙控制。ZStack混合云的安全组对应了阿里云端ECS云主机三层隔离的防火墙约束。

镜像 (Image)

云主机使用的镜像模板文件，一般包括操作系统和预装的软件。ZStack支持上传本地镜像到阿里云，以及使用阿里云端镜像。

弹性公网IP (EIP)

阿里云端公有网络池中的IP地址，绑定弹性公网IP的ECS实例可以直接使用该IP进行公网通信。

VPN连接 (VPN Connection)

通过建立点对点的IPsec VPN通道，实现企业本地数据中心的私有网络与阿里云端VPN网络进行通信。

VPN网关 (VPN Gateway)

一款基于Internet，通过加密通道将本地数据中心和阿里云专有网络VPC安全可靠连接起来的服务。用户在阿里云VPC创建的IPsec VPN网关，与本地数据中心的用户网关配合使用。

VPN用户网关 (Customer Gateway)

本地数据中心的VPN服务网关。可通过ZStack混合云创建VPN用户网关，并将VPN用户网关与VPN网关连接起来。

高速通道 (Express Connect)

通过物理专线（即租用运营商的专线：电缆或光纤），连通本地数据中心到阿里云专线接入点，与阿里云VPC环境打通，实现云上云下不同网络间高速，稳定，安全的私网通信。

边界路由器 (VBR)

用户申请的物理专线接入交换机的产品映射。用户在物理专线上可以创建边界路由器，边界路由器负责专线上的数据在阿里云上进行转发。通过边界路由器，用户数据可以直达阿里云VPC网络。

路由器接口 (Router Interface)

一种虚拟的网络设备，可以挂载在路由器并与其他路由器接口进行高速通道互联，实现不同网络间的内网互通。