

SmartX 超融合 技术原理与特性解析 合集（二）

– 管理与运维

深入解读磁盘亚健康检测、存储性能管理、升级、扩容、迁移等关键技术与特性。

关于 SmartX

北京志凌海纳科技有限公司 (SmartX) 成立于 2013 年, 是专业的现代化 IT 基础设施产品与方案提供商。基于自主研发的分布式块存储, SmartX 提供超融合、企业云基础设施、分布式存储、云原生存储等产品和服务, 助力客户构建精简、敏捷、可靠、安全的 IT 基础设施, 实现降本增效, 服务业务创新, 助力数字化转型。SmartX 已服务交通银行、泰康保险集团、国泰君安证券、中信建投证券、海尔、京东方、中山一院、韩国 SBS 电视台、Cafe24 等多个金融、制造、医疗行业领导者, 并先后获评 Gartner 亚太区客户之选、IDC 创新者。

了解更多有关 SmartX 的信息, 请访问官方网站。

www.smartx.com

联系销售了解产品与服务, 请在工作日 9:00 – 18:00 给我们来电。

400-116-5559

发送邮件向我们咨询产品或市场的更多信息。

info@smartx.com

获取最新技术资讯与行业客户实践, 扫码关注微信公众号。



SmartX HCI 是中国独立超融合软件市场份额排名第一的超融合基础设施产品组合。它以弹性、精简的架构一站式提供虚拟化、分布式存储、软件定义网络与安全、容器管理与服务、数据保护与容灾等云基础架构核心组件，关键业务支撑能力经行业头部客户大规模验证，让您以更低的初始投资，逐步、平稳实现云化、国产化替代和容器化转型目标。

SmartX HCI 具备领先的全栈能力、核心自主研发、承载关键业务、方案开放解耦的核心优势。这些优势如何通过具体的技术与特性实现的？能为用户带来哪些好处？与 VMware、Nutanix 等国际厂商对比如何？

为解答以上问题，本文档挑选了 SmartX 超融合所涉及的部分技术原理与特性解析，基于博客内容整理而成，希望能对读者深入了解该产品有所帮助。《管理与运维》部分包含：磁盘亚健康检测、存储性能管理、扩容、升级、迁移等。

目录

亚健康检测 一文了解 SmartX 超融合硬盘健康检测机制与运维实践	2
存储性能管理 如何利用 SmartX 存储性能测试工具 OWL 优化性能管理?	14
扩容 一文了解 SmartX 超融合如何扩容	20
升级 实现 IT 基础架构软硬件升级简单又不停机	26
升级 新建集群 VS 滚动升级: 如何选择服务器硬件平滑升级方案?	29
迁移 一文了解 SMTX 迁移工具原理与实践	35
迁移 从物理机/云平台迁移至超融合? SMTX CloudMove 帮你实现	41

亚健康检测 | 一文了解 SmartX 超融合硬盘健康检测机制与运维实践

[点击链接阅读原文：一文了解 SmartX 超融合硬盘健康检测机制与运维实践](#)

要点总结

硬盘作为消耗品，会随着持续使用出现性能下降及故障，基于此情况，SMTX OS 4.0.10 版本正式增加硬盘健康检查功能，可通过开源工具 S.M.A.R.T. 检测技术、SmartX 自研的 disk-health 硬盘健康检查工具和 SmartX 超融合集群的数据巡检功能三种方式实现。

不健康盘和亚健康盘异常场景均会触发告警，并且在告警信息中体现不同的异常状态分类，均可通过故障确认、更换进行处理。

disk-health 存在局限性，如无法判断是否因为 HBA 卡导致故障，或因网络异常导致节点存储链路异常，从而增加对应的 I/O error 等情况。如遇硬盘告警提示，建议使用前面提供的命令，通过后台进行相关调查、辅助确认。

某客户数据中心集群在一周内收到 17 块数据盘出现亚健康或不健康的告警，涉及 3 套集群。SmartX 工程师协助客户排查问题，并提供了基于节点维度的硬盘批量更换方案，目前所有受影响服务器的硬盘已全部完成替换。

磁盘故障是存储系统内最常见的问题之一，需要在存储系统设计和日常运维管理中重点关注。

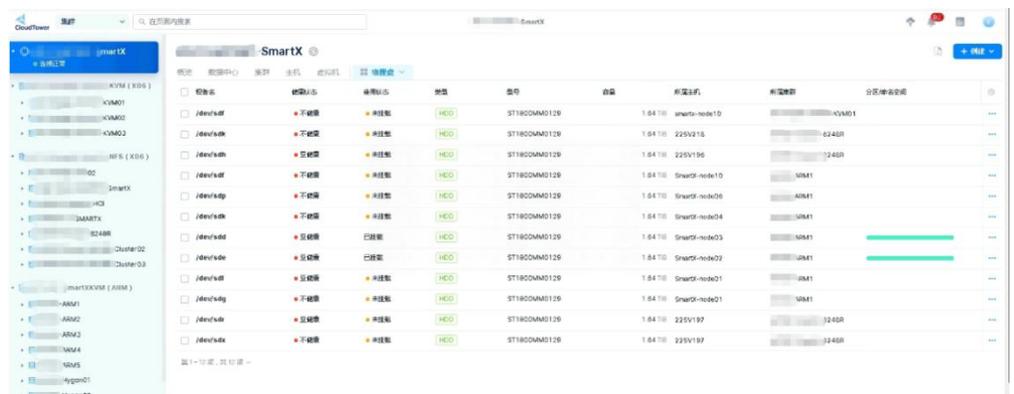
SmartX 超融合集群如何监测磁盘健康状况？

在实际使用时会为运维人员带来哪些帮助？

SmartX 售后技术经理张瑞松在视频中分享了 SmartX 集群硬盘健康检测机制和客户运维实践，希望能为有需求的读者提供参考。

**以下为根据视频整理的文字内容。*

我们这次分享的主题是“集群磁盘亚健康检测机制以及运维的最佳实践”。在集群管理界面上，有时会遇到集群硬盘不健康或者亚健康的告警。为什么会出现这种情况？



我们知道，所有的硬盘都是消耗品，使用寿命是有限的。随着硬盘的持续使用，会出现各种各样的因素导致硬盘性能下降以及故障，因此需要我们被动地进行更换。根据已收集到的数据，一块正常硬盘在持续使用的情况下，出现故障的高峰期一般是在出厂后的第 1.5 年到第 3 年之间。

基于这种情况，SmartX 超融合产品增加了硬盘健康检查的功能。该功能自 SMTX OS 4.0.10 版本正式增加，欲了解功能特性，请阅读：[怎样在“坏盘时刻”保持优雅](#)。我们也会在下面进行详细解读。

SmartX 超融合集群硬盘健康检查功能

技术目标

- 系统可以主动并快速地探测、隔离异常硬盘并开始恢复数据，减少异常硬盘对用户系统的影响。
- 当检测到硬盘异常之后，发出报警并展示硬盘的状态，同时系统进行相应的处理：
 - 减少售后的现场等待时间、简化售后的处理操作。

- 通过报警，明确展示硬盘的当前状态，提示用户在资源许可的情况下进行硬件替换操作。

实现方式

监测磁盘健康状态可通过三种方式实现。一是开源工具 S.M.A.R.T. 检测技术，二是 SmartX 自研的 disk-health 硬盘健康检查工具，三是 SmartX 超融合集群的数据巡检功能。

S.M.A.R.T.

S.M.A.R.T. 是一种使用比较广泛的磁盘分析检测技术，早在 90 年代末就基本普及到每一块硬盘（包括 IDE、SCSI），允许磁盘在运行的时候将自身的若干参数记录下来，包括型号、容量、温度、密度、扇区、寻道时间、传输、误码率等。硬盘开始运行后，部分参数会发生变化，若某一参数超过报警阈值，则说明硬盘接近损坏，此时硬盘虽然在工作，但已经变得不可靠了，随时可能出现故障。

部分参数值如下：

```

-- SMART Attributes --
ID# ATTRIBUTE_NAME          VALUE     THRESH     RAW             CHECK_FIELD   CHECK_THRES  CHECK_RES
  5   Reallocated_Sector_Ct   099       001        12              raw           10          False
194   Temperature_Celsius     100       000        30              raw           45          True
197   Current_Pending_Sector  100       000         0              raw           0           True
198   Offline_Uncorrectable   100       000         0              raw           0           True
253   Media_Wearout_Indicator  100       000       36772991       value         20          True
  
```

其中有几个重要参数：

参数	功能
Raw_Read_Error_Rate	标识磁盘健康与否的关键指标属性
Reallocated_Sector_Ct	有多少数据块被重新 remapping
Current_Pending_Sector	当前有多少数据块处于不可用状态

另外，S.M.A.R.T. 会监控 SSD 剩余使用寿命，默认每 6 小时检查一次。

disk-health

SmartX 自研工具 disk-health 会根据硬盘的异常情况将其划分为不健康盘、亚健康盘、S.M.A.R.T. 自检不通过盘和寿命不足盘。

1.不健康盘

SmartX 分布式块存储 ZBS 内部根据 I/O 的实际返回情况自动探测，在发现如下 I/O 异常的场景会将硬盘标记为不健康盘：

- 坏盘：I/O 返回错误累积超过 100 次（SSD、HDD）或 checksum 错误累积超过 100 次（HDD）。
- 慢盘 I：发生过 I/O 错误且超时 30s（SSD、HDD）。

2.亚健康盘

亚健康盘指没有达到不健康盘程度，但依旧不健康的硬盘。

- **满足以下条件即会被判定为慢盘 II**：通过监控硬盘最近 90s（6 * 15s）内的相关数据进行判断，两种类型的硬盘需在连续 6 个 15s 内分别同时满足以下条件：
 - HDD
 - 平均延迟 > 3s。

- IOPS < 50。
 - 读写速度 <= 50MiB/s。
 - SSD
 - 平均延迟 > 500ms。
 - IOPS < 5000。
 - 读写速度 <= 150MiB/s。
- 满足以下条件也会被标记为亚健康盘 (0 < IO error < 100)：对硬盘的读写出现错误 (小于 100 次)、checksum 出现错误 (小于 100 次)。

3.S.M.A.R.T. 自检不通过盘

smartctl 工具的相关参数显示硬盘自检结果不通过。

4.寿命不足盘

因硬盘寿命不足触发提示性告警，建议用户更换硬盘。

数据巡检

1.功能介绍

SmartX 超融合集群中的数据巡检功能可以主动探测硬盘静默损坏导致的数据不一致问题，并触发数据恢复，保护数据安全。

2.实现机制

首先，数据巡检会周期性检测副本的 Generation 信息。Generation 是块存储的数据版本号，初始状态为 0，每次数据块 (extent) 被写入数据都会触发 Generation 数值加 1，检测周期是 30 秒。由于采用两副本的机制，每隔 30 秒会检测一下这两份数据的 Generation 信息是否一致，如果不一致，就认为 Generation 信息低的数据块可能是发生了预料外的状况，在这种情况下就会触发数据恢复，保证数据的一致性。

另外，数据巡检还会周期性检测副本的 checksum 信息。对集群写入的数据会同步生成 checksum 信息，用于在读取数据时进行校验。检测以磁盘为单位进行扫描，两个硬盘之间的扫描间隔为 5 分钟，扫描速率为 5MB/s、IOPS = 20，这样能够最大限度降低巡检对集群性能的影响，且同一时间只会对一个节点的一个硬盘进行检测，不会占用集群和节点上的更多的 I/O。

以上 Generation 校验和 checksum 校验互相独立，互不影响，以此来保证底层数据的完整性。

3.适用版本

数据巡检功能适用于 SMTX OS 4.0.11 及以上版本、SMTX OS 5.0.1 及以上版本、SMTX ZBS 5.0.0 及以上版本。

异常硬盘处理流程

接下来，跟大家介绍一下应该怎样处理出现问题的硬盘。

处理逻辑

分类		磁盘类型	系统处理方法	人工处理方法
不健康盘	坏盘	SSD HDD	自动隔离	系统隔离完成后进行拔盘 如系统隔离无法完成，则需要人工提前介入
	慢盘 I	SSD HDD		
亚健康盘	慢盘 II	HDD(lsm2)		界面上手动卸载，卸载完成后进行拔盘
		SSD HDD(lsm1)	-	
	0 < IO error < 100	SSD HDD	-	
S.M.A.R.T. 不通过盘		HDD	-	
寿命不足		SSD HDD	-	

当系统将一块磁盘定位为不健康盘后，会自动对其进行隔离，不会再往上写入数据，且会将盘上的数据全部迁移走，等完成隔离操作后，这块盘会被自动从系统里踢掉。

对于被判定为亚健康、S.M.A.R.T. 检测不通过以及寿命不足的硬盘，这些盘只是对比正常硬盘而言性能没有那么好，并不是完全不可用。所以针对这种类型的硬盘，系统会发出相应的告警提示，然后由客户来评判是将这块盘立刻卸载，还是先保留一段时间。

上述不同硬盘异常场景均会触发告警，并且在告警信息中体现不同的异常状态分类。

下面分别介绍不健康盘和亚健康盘出现故障时的处理流程。

不健康盘处理流程

故障确认

通过 Web 管理界面（CloudTower）登录后，可以查看到不健康硬盘被自动卸载完成，处于未挂载状态。



点击“查看故障信息”，显示如下：



我们可以通过硬盘健康检测工具查看出现问题的原因：

- 在故障告警产生的对应节点上执行以下命令，查看 sdc 的状态。
zbs-node show_disk_status sdc
zbs-node show_disk_status /dev/sdc
- 使用 sdc 的序列号，在集群中任意一个节点上输入以下命令，查看 sdc 的状态。
zbs-node show_disk_status -s PHYS826003JB480BGN

```
[root@dogfood-ldc-elf-95 19:14:45 ~]#zbs-node show_disk_status -s PHYS826003JB480BGN
== Base Information ==
is healthy                : False
device name               : /dev/sdc
bus type                  : ata
model                     : SSDSC2KB48067R
firmware                  : SCV1DL58
disk serial               : PHYS826003JB480BGN
last belong to           : 10.255.0.96
== Fault Detection ==
chunk errflag detected    : False
chunk warnflag detected   : False
chunk io error detected   : False
chunk checksum error detected : False
iostat latency detected   : False
smart error detected      : True
== Extra Fault Detection ==
chunk num_io_errors       : -
chunk num_checksum_errors : -
io latency (ms)           : -
smartctl hang process     : -
S.M.A.R.T. assessment error : -
== S.M.A.R.T. Attributes ==
ID#  ATTRIBUTE_NAME      VALUE  THRESH  RAW      CHECK_FIELD  CHECK_THRESH  CHECK_RES
  5  Reallocated_Sector_Ct  099    001     12       raw          10             False
194  Temperature_Celsius    100    000     30       raw          45             True
197  Current_Pending_Sector  100    000     0        raw          0              True
198  Offline_Uncorrectable  100    000     0        raw          0              True
233  Media_Wearout_Indicator  100    000     3672991  value        20             True
[root@dogfood-ldc-elf-95 19:15:01 ~]#
```

其中 S.M.A.R.T. 检测不通过，该项指标显示为 True，可以看到下方指标已经超过了默认的设置阈值 10，达到了 12，导致检测结果不通过。

zbs-node show_disk_status -s PHYS826003JB480BGN 输出字段含义如下：

指标子项	结果值	含义
chunk errflag detected	False	chunk 将磁盘标记为坏盘。
chunk warnflag detected	False	chunk 被 disk-healthd 告知磁盘读写延迟高。
chunk io error detected	False	chunk 探测出磁盘上发生了 I/O 错误。
chunk checksum error detected	False	chunk 检查到磁盘数据不一致。
iostat latency detected	False	disk-healthd 检查到磁盘读写延迟高。
smart error detected	True	disk-healthd 检查到磁盘的 S.M.A.R.T. 属性超出阈值。

也可以通过查看 Message 日志的形式查看问题。

```
[root@SmartX-6248R-Node09 09:03:57 smartx]#grep sdi /var/log/messages
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 BRCM Debug mfi stat 0x2d, data len requested/completed 0x4400/0x0
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 FAILED Result: hostbyte=DID_OK driverbyte=DRIVER_SENSE cmd_age=0s
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 Sense Key : Medium Error [current]
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 Add. Sense: Unrecovered read error
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 CDB: Read(10) 28 00 4d 29 72 0c 00 00 22 00
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: blk_update_request: critical medium error, dev sdi, sector 1294561004
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 BRCM Debug mfi stat 0x2d, data len requested/completed 0x2200/0x0
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 FAILED Result: hostbyte=DID_OK driverbyte=DRIVER_SENSE cmd_age=1s
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 Sense Key : Medium Error [current]
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 Add. Sense: Unrecovered read error
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 CDB: Read(10) 28 00 64 40 f8 35 00 00 11 00
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: blk_update_request: critical medium error, dev sdi, sector 1081979445
[root@SmartX-6248R-Node09 09:04:30 smartx]#
```

更换流程

step 1 确认服务器 SN，定位物理服务器位置。

step 2 Web 界面（CloudTower）点亮磁盘，标记磁盘物理位置。



step 3 线下执行拔出故障盘，更换新硬盘。

step 4 Web UI 识别，点击挂载磁盘。

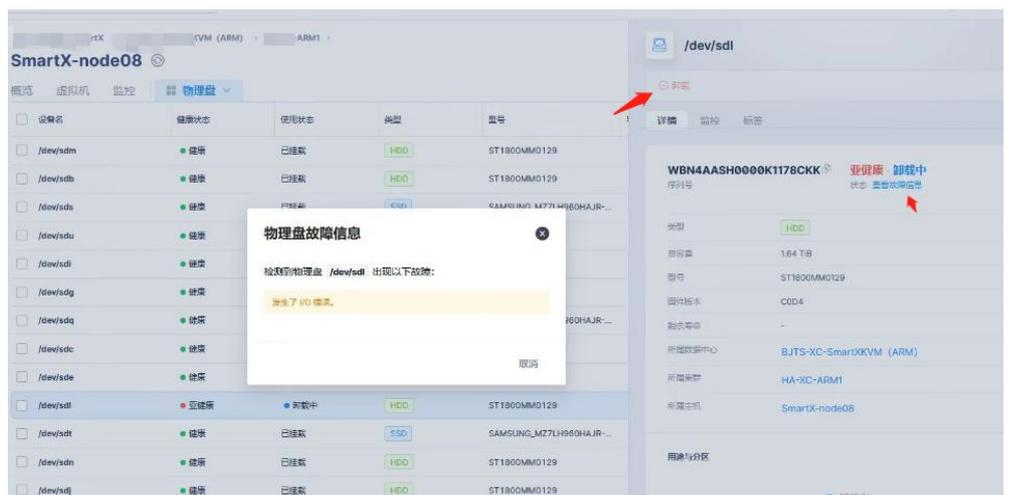
step 5 挂载用途：

a. HDD 选择挂载为“数据盘”，挂载时间通常约 11s 左右。

b. SSD 可选择挂载为“含元数据分区的缓存盘”和“缓存盘”，按照默认挂载即可。“含元数据分区的缓存盘”需要重建软 RAID 1，挂载时间约 10min 左右；“缓存盘”挂载时间通常约 11s 左右。

亚健康盘处理流程

通过 Web 管理界面（CloudTower）登录后，硬盘正常处于“挂载”状态；具体的故障原因，可点击“查看故障信息”显示如下：



亚健康盘只是不符合预期性能的硬盘，因此并不会直接去对它进行卸载操作。

故障确认

流程和方法与不健康盘一致。

更换流程

step 1 确认服务器 SN。

step 2 磁盘卸载前，确认集群无数据恢复。

```
[root@node-01 00:20:28 smartx]$zbs-meta pextent find need_recover
No PExtents found.
[root@node-01 00:20:44 smartx]$
```

step 3 Web 界面（CloudTower）点击“卸载”。



step 4 硬盘卸载完成后，状态变更为“未挂载”（卸载时间和硬盘上数据多少有关）。如果计划换盘，建议提前点击卸载硬盘，等到卸载完成以后，再到机房现场进行更换。

step 5 Web 界面（CloudTower）点击“闪灯”，标记物理磁盘位置。

step 6 线下更换故障盘，插入新硬盘。

step 7 Web 界面（CloudTower），点击“挂载”磁盘。

step 8 挂载用途：

a. HDD 选择挂载为“数据盘”，挂载时间通常约 11s 左右。

b. SSD 可选择挂载为“含元数据分区的缓存盘”和“缓存盘”，按照默认挂载即可。”含元数据分区的缓存盘“需要重建软 RAID 1，挂载时间约 10 min 左右；”缓存盘“，挂载时间通常约 11s 左右。

disk-health 局限性

首先，disk-health 只是基于系统内 I/O 去检测硬盘是否健康，无法判断是否因为 HBA 卡导致了故障。如果 HBA 卡出现故障，可能导致 HBA 卡及 HBA 卡下面的所有硬盘在接触 I/O 时都会被标记为亚健康或

不健康，出现误判现象。

因此如果遇到一个节点上多块硬盘出现问题，建议先联系一下 SmartX 工程师，由工程师协助定位。或者直接联系硬件报修，检查一下 RAID 卡是否出现了问题。

另外，由于我们的集群采用两副本机制，写数据时一份数据写到本地，另一份数据写到另一个节点上。如果因为网络异常导致节点存储链路异常，可能也会增加对应的 I/O error。这种情况可能也需要进行单独排查。

某些软件因素也会导致硬盘无法正常读写，这种情况虽然不属于硬盘故障，但会增加 I/O error 的计数。

总之，如果遇到了类似的硬盘告警提示，建议使用前面提供的命令，通过后台进行相关调查、辅助确认。

硬盘故障运维实践

背景

某客户数据中心集群在一周内收到 17 块数据盘出现亚健康或不健康的告警，涉及 3 套集群。截至集群问题正式处理完成，累计故障硬盘 74 块，涉及 7 套集群。

在故障发生一周内，由于硬盘持续出现亚健康和不健康告警，并非短时间内大量硬盘同时故障，因此亚健康和故障硬盘的告警会自动触发数据迁移操作，避免业务数据的损失。同时由于硬盘故障所在集群有足够的磁盘空间冗余，可以安全快速地实现故障磁盘数据的迁移和平衡，因此没有发生节点级别或集群级别的故障，也没有造成数据损失，对生产无任何影响。

ID	名称	健康状态	使用状态	类型	型号	容量	所属主机	所属组群
/data/dm1		不健康	故障中	HDD		1.64 T	sdm11	KVM01
/data/dm2		不健康	故障中	HDD		1.64 T	sdm9	KVM01
/data/dm3		不健康	故障中	HDD		1.64 T	sdm7	KVM01
/data/dm4		不健康	非狂转	HDD		1.64 T	ss	CSAPR
/data/dm5		不健康	非狂转	HDD		1.64 T	7P	ADAM
/data/dm6		亚健康	已挂载	SDD	SSD	594.25 G	ss	CSAPR
/data/dm7		亚健康	非狂转	HDD		2.16 T	ss	CSAPR
/data/dm8		亚健康	已挂载	HDD		1.64 T	sdm06	ADAM
/data/dm9		亚健康	已挂载	HDD		1.64 T	sdm9	ADAM
/data/dm10		亚健康	非狂转	HDD		1.64 T	ss	ADAM
/data/dm11		亚健康	已挂载	HDD		1.64 T	sdm14	ADAM
/data/dm12		亚健康	非狂转	HDD		1.64 T	ss	ADAM

故障确认

step 1 使用亚健康硬盘检测工具，发现某硬盘出现 I/O error。

```
[root@SmartX-6248R-Node09 09:04:50 smartx]$zbs-node show_disk_status sdi
== Base Information ==
is healthy : False
device name : /dev/sdi
bus type : scsi
model :
firmware :
disk serial :
last belong to : 1.1.1.225
== Fault Detection ==
chunk errflag detected : False
chunk warnflag detected : False
chunk io error detected : True
chunk checksum error detected : False
iostat latency detected : False
smart error detected : False
== Extra Fault Detection ==
chunk num_io_errors : 2
chunk num_checksum_errors : -
io latency (ms) : -
smartctl hang process : -
S.M.A.R.T. assessment error : -
== S.M.A.R.T. Attributes ==
[root@SmartX-6248R-Node09 09:05:12 smartx]$
```

step 2 看 Message 日志，确认 I/O error 出现原因：该硬盘在短时间内连续两次出现 medium error 的情况。

```
[root@SmartX-6248R-Node09 09:03:57 smartx]$grep sdi /var/log/messages
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 BRCM Debug mfi stat 0x2d, data len requested/completed 0x4406/0x0
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 FAILED Result: hostbyte=DID_OK driverbyte=DRIVER_SENSE cmd_age=0s
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 Sense Key : Medium Error [current]
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 Add. Sense: Unrecovered read error
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#0 CDB: Read(10) 28 00 4d 29 72 0c 00 00 22 00
Dec 23 05:41:39 SmartX-6248R-Node09 kernel: blk update_request: critical medium error, dev sdi, sector 1294561804
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 BRCM Debug mfi stat 0x2d, data len requested/completed 0x2280/0x0
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 FAILED Result: hostbyte=DID_OK driverbyte=DRIVER_SENSE cmd_age=1s
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 Sense Key : Medium Error [current]
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 Add. Sense: Unrecovered read error
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: sd 1:2:7:0: [sdi] tag#14 CDB: Read(10) 28 00 64 40 f8 35 00 00 11 00
Dec 23 05:57:50 SmartX-6248R-Node09 kernel: blk update_request: critical medium error, dev sdi, sector 1681979445
[root@SmartX-6248R-Node09 09:04:30 smartx]$
```

step 3 smartctl 工具查看对应的 error 日志，发现有对应的读 error 的情况，也有不正常写入延迟。

```
Error counter log:
      Errors Corrected by           Total   Correction      Gigabytes      Total
      ECC      fast | delayed   rereads/   errors   algorithm  processed   uncorrected
                        rewrites  corrected  invocations [10^9 bytes] errors
read:  3745947516      70      0  3745947586      76      15422.285      4
write:      0      0      0      0      0      12633.650      0
verify: 43955076      41      0  43955117      41      166.627      0

Non-medium error count:      0

[GLTSD (Global Logging Target Save Disable) set. Enable Save with '-S on']
SMART Self-test log
Num Test          Status          segment  LifeTime  LBA_first_err [SK ASC ASQ]
  #  Description  (hours)
# 1 Background short Completed      -      20084      - [- - -]
# 2 Background short Completed      -      20060      - [- - -]
# 3 Background short Completed      -      20036      - [- - -]
# 4 Background short Completed      -      20012      - [- - -]
# 5 Background short Completed      -      19988      - [- - -]
# 6 Background short Completed      -      19964      - [- - -]
# 7 Background short Completed      -      19940      - [- - -]
# 8 Background short Completed      -      19916      - [- - -]
# 9 Background short Completed      -      19892      - [- - -]
#10 Background short Completed      -      19868      - [- - -]
#11 Background short Completed      -      19844      - [- - -]
#12 Background short Completed      -      19819      - [- - -]
#13 Background short Completed      -      19796      - [- - -]
#14 Background short Completed      -      19772      - [- - -]
#15 Background short Completed      -      19748      - [- - -]
#16 Background short Completed      -      19724      - [- - -]
#17 Background short Completed      -      19700      - [- - -]
#18 Background short Completed      -      19676      - [- - -]
#19 Background short Completed      -      19652      - [- - -]
#20 Background short Completed      -      19628      - [- - -]

Long (extended) Self Test duration: 9459 seconds [157.7 minutes]
[root@SmartX-6248R-Node09 09:05:55 smartx]$
```

调研结果

我们将客户 17 块盘出现的问题进行收集整理，发现基本上大多数硬盘型号、固件相同，硬盘通电时间都在 1.9 年以上，故障原因以 medium error 为主。

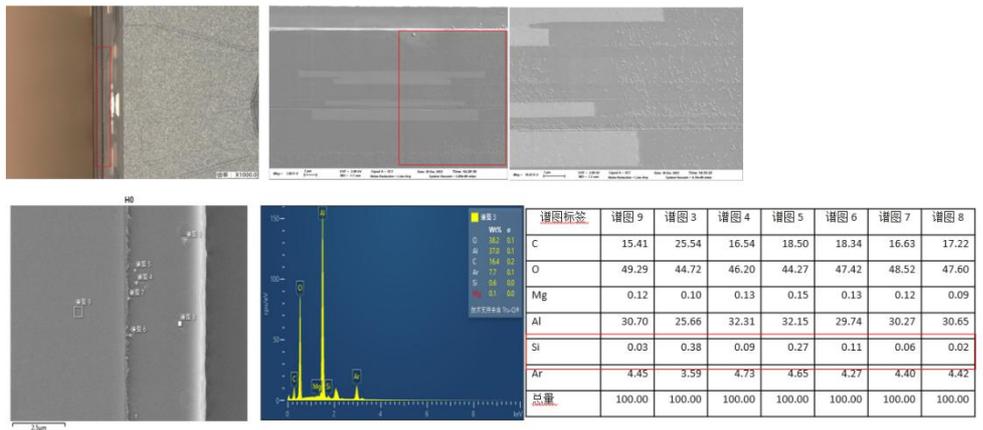
故障集群	主机	IP	盘符	硬盘类型	健康状态	故障告警时间	故障原因	计数	报错信息	通电时间 (时)	通电时间 (天)	通电时间 (年.365)	集群部署时间
M1	SmartX-node04	10.48.76.32	sdk	HDD	不健康	2022/11/27 14:29	io error	127	Medium error	15104.53	629.3554167	1.724281416	Aug-20
M1	SmartX-node01	10.48.76.41	sdg	HDD	不健康	2022/11/30 21:05	io error	130	Medium error	17119.95	713.33125	1.954332192	Aug-20
M1	SmartX-node01	10.48.76.41	sdh	HDD	不健康	2022/12/2 07:42	io error	2	Medium error	17120.18	713.3408333	1.954338447	Aug-20
VM01	SmartX-node10	10.101.225.184	sdj	HDD	不健康	2022/12/2 11:04	io error	58		25388.25	1057.84375	2.888202055	Nov-18
			sdh	HDD	不健康	2022/12/2 12:39	io latency	37695ms		19640.88	818.37	2.242109589	Sep-20
M1	SmartX-node10	10.48.76.38	sdj	HDD	不健康	2022/12/2 17:47	io error	11	aborted Command	no	no	no	Aug-20
M1	SmartX-node06	10.48.76.34	sdp	HDD	不健康	2022/12/3 16:26	io error	4	aborted Command	no	no	no	Aug-20
			sdq	HDD	不健康	2022/12/3 19:37	io latency	18652		19648.37	818.8280833	2.242964612	Sep-20
I1	SmartX-node03	10.48.76.43	sdj	HDD	不健康	2022/12/4 05:29	io error	5	Medium error	17118.33	713.26375	1.95414726	Aug-20
			sdx	HDD	不健康	2022/12/4 10:36	io error	1		19648.2	818.675	2.242945205	Sep-20
I1	SmartX-node03	10.48.76.43	sdj	HDD	不健康	2022/12/4 12:30	io error	10		19646.87	818.6195833	2.242793379	Sep-20
			sdh	HDD	不健康	2022/12/4 16:02	io error	7	Medium error	17119.86	713.3275	1.954321918	Aug-20
M1	SmartX-node02	10.48.76.42	sdh	HDD	不健康	2022/12/5 00:34	io error	2	Medium error	17130.78	713.7825	1.955568493	Aug-20
			sdj	HDD	不健康	2022/12/5 01:09	io latency	13216		818.8283333	2.243365297	Sep-20	
VM01	SmartX-node10	10.101.225.184	sdq	HDD	不健康	2022/12/5 13:20	io error	14		25388.25	1057.84375	2.888202055	Nov-18
			sdv	HDD	不健康	2022/12/5 14:03	io error	12		818.675	2.242945205	Sep-20	
RM1	SmartX-node01	10.48.76.41	sdh	HDD	不健康	2022/12/5 15:03	io error	1	Medium error				Aug-20
			sdj	HDD	不健康	2022/12/7 17:24	io error	1		19720.38	821.6825	2.251184931	Sep-20
VM1	SmartX-node05	10.48.76.33	sdh	HDD	不健康	2022/12/5 21:34	io error	39		19721.22	821.7175	2.251208022	Sep-20
			sdj	HDD	不健康	2022/12/5 20:11	io error	2	Medium error	14512.33	604.68	1.656668675	Aug-20
RM1	SmartX-node01	10.48.76.41	sdk	HDD	不健康	2022/12/8 00:54	io error	66	IO timeout	no	no	no	Aug-20
			sdj	HDD	不健康	2022/12/8 17:31	io error	114	aborted Command	14510.1	604.5875	1.656404109	Aug-20
I1	SmartX-node06	10.48.76.34	sdh	HDD	不健康	2022/12/8 23:16	io latency	22477.554		19755.2	823.13333	2.251598177	Sep-20
			sdj	HDD	不健康	2022/12/9 04:08	io error	1		no	no	no	Sep-20
I1	SmartX-node03	10.48.76.43	sdj	HDD	不健康	2022/12/10 09:54	io error	1		19776.63	824.02925	2.257696164	Sep-20
			sdh	HDD	不健康	2022/12/10 09:54	io error	9		19780.42	824.1841667	2.258038813	Sep-20
ARM1	SmartX-node02	10.48.76.42	sdj	HDD	不健康	2022/12/9 11:17	io error	10	Medium error	17252.53	718.8554167	1.989468894	Aug-20
ARM1	SmartX-node04	10.48.76.32	sdh	HDD	不健康	2022/12/10 09:31	io latency	8	aborted Command	no	no	no	Aug-20
			sdj	HDD	不健康	2022/12/10 09:31	io error	13		20665.23	861.05125	2.35984452	Nov-18
			sdv	HDD	不健康	2022/12/13 21:56	io error	1		32381.4	1349.225	3.696568646	Nov-18
VM01	SmartX-node09	10.101.225.183	sdm	HDD	不健康	2022/12/13 14:54	io error	1		25622.38	1067.599167	2.524592224	Nov-18

我们把调查结果分享给客户，建议协调对应的硬件厂商和专门的硬盘厂商进行分析。

在硬盘分析时，发现其中一块硬盘磁头的脏污染特别严重，SEM的成分分析显示磁头含有 Si 这种物质。硬盘的研发后线深入分析后反馈，碟片工艺在生产时会包含 Si 物质残留，属于正常情况。但是 Si 在低负载的情况下，会吸收空气里的水分而膨胀，导致硬盘上的盘面出现水汽沉积，形成糊状脏污，业务压力增加更容易造成头碟接触，污染磁头，影响磁头飞行高度，增加读写错误率。现场调研后发现，该服务器上的硬盘并没有插满，其中大概一到两个硬盘槽位上使用的是硬盘的占位符，后者对比实际的硬盘来说，会增加空气的流通面积，再叠加硬盘低负载的使用情况，更容易增加水汽沉积，一直到使用大约 1.9 年后，糊状污染真正开始影响硬盘的读写。

硬盘FA分析

开盘分析：H0磁头严重脏污污染，脏污呈糊状，脏污SEM成分分析，脏污为含Si的物质（C/O/Al/Ar分析为磁头本身元素）。



处理结果

最终客户跟硬件厂商协调，将这些故障节点上的所有硬盘进行了批量替换。为了保证客户节点上硬盘批量替换的效率，SmartX 配合客户提供了基于节点维度的硬盘批量更换方案，目前所有受影响服务器的硬盘已全部完成替换。

集群名称	主机名称	节点删除时间	硬盘计划更换时间	节点操作人员	硬盘更换人员	更换进度	硬盘实际更换时间	
	97	2022-12-25	2022-12-26	SmartX	硬件工程师	已完成	2022-12-26	
	96	2022-12-28	2022-12-27	SmartX	硬件工程师	已完成	2022-12-27	
	98	2022-12-27	2022-12-28	SmartX	硬件工程师	已完成	2022-12-29	
	97	2022-12-28	2022-12-29	SmartX	硬件工程师	已完成	2022-12-30	
	98	2022-12-29	2022-12-30	SmartX	硬件工程师	已完成	2022-12-31	
	99	2022-12-30	2023-01-03	SmartX	硬件工程师	已完成	2023-01-01	
	00	2023-01-03	2023-01-04	SmartX	硬件工程师	已完成	2023-01-03	
	01	2023-01-04	2023-01-05	SmartX	硬件工程师	已完成	2023-01-04	
	02	2023-01-05	2023-01-06	SmartX	硬件工程师	已完成	2023-01-05	
	03	2023-01-06	2023-01-07	SmartX	硬件工程师	已完成	2023-01-06	
	04	2023-01-07	2023-01-08	SmartX	硬件工程师	已完成	2023-01-07	
	05	2023-01-08	2023-01-09	SmartX	硬件工程师	已完成	2022-12-28	
	06	2023-01-07	2023-01-08	SmartX	硬件工程师	已完成	2023-01-08	
	07	2023-01-08	2023-01-09	SmartX	硬件工程师	已完成	2023-01-09	
	08	2023-01-09	2023-01-10	SmartX	硬件工程师	已完成	2023-01-10	
	09	2023-01-10	2023-01-11	SmartX	硬件工程师	已完成	2023-01-10	
	01	2023-01-11	2023-01-12	SmartX	硬件工程师	已完成	2023-01-11	
	03	2023-01-12	2023-01-13	SmartX	硬件工程师	已完成	2023-01-11	
	07	2023-01-13	2023-01-14	SmartX	硬件工程师	已完成	2023-01-12	
	09	2023-01-14	2023-01-15	SmartX	硬件工程师	已完成	2023-01-12	
	01	07	2022-12-27	2022-12-28	SmartX	硬件工程师	已完成	2022-12-28
	01	09	2022-12-28	2022-12-29	SmartX	硬件工程师	已完成	2022-12-29
	01	09	2022-12-29	2022-12-30	SmartX	硬件工程师	已完成	2023-01-03
	01	10	2022-12-30	2023-01-03	SmartX	硬件工程师	已完成	2022-12-30
	01	11	2023-01-03	2023-01-04	SmartX	硬件工程师	已完成	2022-12-31
	01	12	2023-01-04	2023-01-05	SmartX	硬件工程师	已完成	2023-01-01

硬盘运维最佳实践建议

- 禁止在有数据恢复的情况下直接对硬盘进行操作。
- 禁止在开、关机的情况下，直接拔插硬盘。
- 同一时间内集群只允许仅对一块硬盘进行卸载，硬盘挂载可同时挂载多块硬盘。
- 两副本仅允许一个节点上硬盘故障，三副本允许同时两个节点上硬盘故障。
- 增加集群告警功能，告警方式支持 snmp、smtp、API 告警功能、CloudTower 管理界面告警等。
- 据不完全统计，服务器运行 3 年及 3 年以上，故障率将会明显增加，其中硬盘故障率增加相对更为明显，我们建议：
 - 超融合集群版本升级到 4.0.10 及之后版本，当前最新版本为 5.0.5 版本。
 - 核心及重要业务尽量避免运行在使用 3 年以上的服务器，SmartX 可提供以下技术支持：
 - i. 服务器老旧硬件替换方案。
 - ii. 相同架构虚拟机迁移方案。
 - iii. VMware 到 SmartX ELF 虚拟化平台 V2V 迁移方案等。

存储性能管理 | 如何利用 SmartX 存储性能测试工具 OWL 优化性能管理?

[点击链接阅读原文：如何利用 SmartX 存储性能测试工具 OWL 优化性能管理?](#)

要点总结

为了帮助用户更好地管理集群存储性能，SmartX 自主研发了自动化存储性能测试工具 OWL。

OWL 可以从三个方面为用户提供帮助：适应不同硬件配置，提供每套存储各自的存储性能“基线”；参考存储性能基线，分类上线业务虚拟机；结合告警功能，主动告警性能瓶颈风险。

利用 OWL 测试结果可以计算告警阈值，从而优化存储性能管理，避免性能瓶颈。其实践案例包括：运维工程师及时收到性能告警，规避业务影响；某国有银行利用 OWL 定制 I/O 模型测试集群性能；某国有银行利用 OWL 评估集群性能是否满足 99th Percentile 要求。

运维人员在日常管理集群时，有时难免会产生这样的困惑：

- 新业务准备上线，在具备多套存储的情况下，应如何选择承载业务的存储环境？
- 业务虚拟机刚上线时运行速度很快，而运行一段时间后，为什么软硬件没有直观的问题但运行还是会变慢？
- 业务反馈虚拟机性能时好时坏，这是怎么回事？
- 新上线的存储与原来的配置不一样，怎样判断两者的性能差别，以及他们分别适合运行什么样的业务？

这些场景都涉及到存储的性能监控，同时也考验运维人员利用监测数据合理安排业务放置、对性能进行调优的能力。

为了帮助用户更好地管理集群存储性能，SmartX 自主研发了自动化存储性能测试工具 OWL。本文中，我们将为大家介绍 OWL 的功能特性和使用方法，并通过实际应用展现如何利用 OWL 测试结果优化性能管理，避免性能瓶颈。

OWL 工具介绍

OWL 是 SmartX 自研的自动化存储性能测试 Web 平台，以 fio 作为性能获取工具，进行集群的性能压力测试。由于 fio 可以调整为多队列、多带宽、多 I/O 模型的测试情景，能够模拟大多数的业务 I/O（如 MySQL 的性能测试调优就经常使用 fio），因此成为了支持 OWL 的最佳选择。另外，OWL 并不绑定 SmartX 超融合集群，[用户也可在其他环境中使用 OWL 进行性能测试](#)。

OWL 可以通过以下三个方面为用户提供帮助：

适应不同硬件配置，提供每套存储各自的存储性能“基线”

为了满足 IT 基础架构信创转型需求，用户可能会购入此前未接触过的国产配件。多种配件组合下，工程师需要了解这些新配置的存储可以达到多少性能、支持哪些应用和数据库。传统的验证方式是直接使用新架构试跑一台业务虚拟机，而使用 OWL 则可以通过模拟类似的 I/O 模型，验证集群的性能情况，从而测试出这套集群存储的性能基线。

参考存储性能基线，分类上线业务虚拟机

用户可基于 OWL 提供的存储性能基线，[为需要上线的业务虚拟机选择合适的存储集群](#)。比如 IOPS 较大的数据库业务，用户可使用全闪集群，而对于 IOPS 比较轻量、数据交互较少的业务，用户可使用性价比较高的混闪集群。

另外，OWL 除了能让用户预先了解每台主机能承载的最大 I/O，还可以通过搭建模拟环境，[帮助用户在业务上线前了解其可能需要的 I/O 大小，合理分配虚拟机放置](#)，避免将多台带宽占用较大的虚拟机放在一台主机上，导致业务正式上线后出现带宽“被业务追着跑”的情况。

结合告警功能，主动告警性能瓶颈风险

用户在使用 OWL 获取性能测试基线后，[可在每套集群的告警规则上设置对应存储性能的读写带宽阈值](#)。当虚拟机的带宽达到了主带宽的 70% 和 80% 时，运维工程师会分别收到告警提示，从而及时观察虚拟机及其他主机的带宽占用情况。这样用户可以在新业务上线之前，将这台虚拟机迁移到相对空闲的主机或集群上。

OWL 使用方法及测试流程

测试前准备

由于 OWL 工具以虚拟机形式运作，用户需要进行 ovf 导入，为 OWL 配置 IP 地址，且保证 OWL 与 test VM ssh 通讯。Test VM 配置要求如下：

- Linux 2c 4G 40G+50G
- 配置 IP 地址，且与 OWL 工具 ssh 通讯
- 安装 FIO 软件

测试流程

- 登录 OWL Web 界面。
- 创建测试模型。
- 添加测试对象。
- 创建测试任务。
- 启动测试任务。
- OWL 结合告警功能，主动告警性能瓶颈风险。

详细测试过程，可观看 [Demo 演示](#)。

利用测试结果优化存储性能管理

常用测试模型

以下为演示中我们常用的 I/O 测试模型。

nPnV	rw	bs	numjobs	iodepth	rwmixread	runtime
nP1V	randread	4k	1	128	-	300
	randwrite	4k	1	128	-	300
	randrw	4k	1	128	70	300
	read	256k	1	128	-	300
	write	256k	1	128	-	300
nPnV	randread	4k	1	128	-	300
	randwrite	4k	1	128	-	300
	randrw	4k	1	128	70	300
	read	256k	1	128	-	300
	write	256k	1	128	-	300

告警阈值计算与设置方式

通过上面的测试拿到性能基线后，用户可以计算出对应的写带宽阈值和读带宽阈值，并在集群里添加告警规则。我们以下图为例介绍阈值的计算方式。

两副本

nPnV	I/O 模型	IOPS	带宽 (MBPS)	延迟 (avg/ms)
8P1V	4k 随机写	41.19k	160.89	3.104
	4k 随机读	122.06k	476.80	1.046
	4k 随机 70% 读 30% 写	65.65k/28.12k	256.43/109.83	1.225/1.681
	256k 顺序写	6.63k	1656.86	19.286
	256k 顺序读	8.65k	2163.88	23.168
8P8V	4k 随机写	264.69k	1033.95	3.862
	4k 随机读	952.25k	3719.74	1.072
	4k 随机 70% 读 30% 写	442.66k/189.67k	1729.16/740.91	1.331/2.28
	256k 顺序写	29.11k	7278.32	35.145
	256k 顺序读	66.46k	16615.33	15.392

以上两组数据，分别是对 8 节点集群中 1 台主机运行 1 台虚拟机，和 8 台主机分别运行 1 虚拟机，进行测试。

我们主要关注带宽。以写带宽为例，在 8P8V 256K 顺序写场景中，写带宽为 7278。我们将 7278 除以 8，得到每一个节点带宽的平均值，然后再将 MBPS 换算成 BPS，该值的 70% 就是我们需要设定为**注意级别**的告警阈值。

写带宽严重告警阈值我们会看 8P1V 256K 场景下的数值。这里写带宽是 1656.86 MBPS，经过单位换算后，这个数值的 80% 将直接作为**严重级别**的告警阈值。由此，我们得到两个写带宽阈值数值，如下图所示。

编辑报警规则

自定义全局规则，或为指定集群定义特例规则

全局规则

启用报警

主机 { .labels.hostname } 的写带宽超过 { .threshold }。

<input checked="" type="checkbox"/>  严重警告	阈值: <input type="text" value="1389874905 Bps"/>
<input checked="" type="checkbox"/>  注意	阈值: <input type="text" value="667788771 Bps"/>
<input type="checkbox"/>  信息	

[重置为默认值](#)

特例规则

[+ 添加特例规则](#)

读带宽告警阈值计算方式与写带宽相同，上述例子中读带宽阈值设置如下图所示。

编辑报警规则

自定义全局规则，或为指定集群定义特例规则

全局规则

启用报警

主机 { .labels.hostname } 的读带宽超过 { .threshold }。

<input checked="" type="checkbox"/>  严重警告	阈值: <input type="text" value="1815194108 Bps"/>
<input checked="" type="checkbox"/>  注意	阈值: <input type="text" value="1524432896 Bps"/>
<input type="checkbox"/>  信息	

[重置为默认值](#)

特例规则

[+ 添加特例规则](#)

用户实践

案例一：运维工程师及时收到性能告警，规避业务影响

某用户使用 OWL 工具进行带宽压测后发现，集群中的某一节点带宽超过了 1.7 GB/s，已超过严重警告级别的阈值。SmartX 后台自动发送告警，提醒运维工程师存储性能已接近极限，从而避免对业务带来直接影响。



案例二：某国有银行利用 OWL 定制 I/O 模型测试集群性能

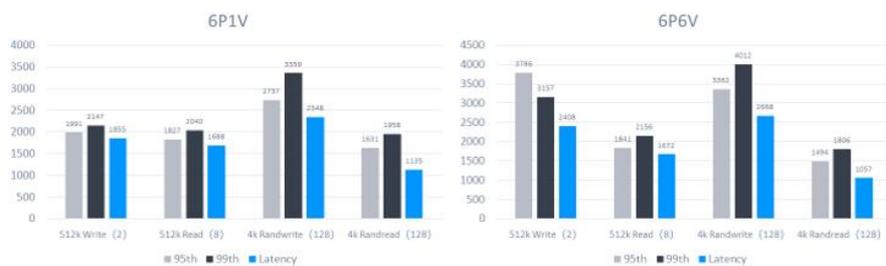
某国有银行为了满足监管需求，利用 OWL 按照定制 I/O 模型（48K，randrw=1:9）连续 12 小时测试集群性能。测试结果显示（如下图），该集群平均 IOPS 标准差可达 54338，延时在 1 毫秒左右。



案例三：某国有银行利用 OWL 评估集群性能是否满足 99th Percentile 要求

某国有银行关注到 99th Percentile 要求，利用 OWL 测试对应的块大小下的存储性能，直观了解该场景下集群的性能情况。测试结果如图所示。

性能测试 - 延时 (3副本)



扩容 | 一文了解 SmartX 超融合如何扩容

[点击查看阅读原文：不止弹性，更加灵活。一文了解 SmartX 超融合如何扩容](#)

要点总结

SmartX 超融合提供三种资源扩容方式：在线添加节点内的存储设备、在集群内在线添加节点扩展计算和存储资源、通过多集群管理平台 CloudTower 实现更大范围的资源池建设和统一管理，赋予用户更大自由度选择服务器和相关设备组件进行扩容。

SmartX 通过自研的存储管理和调度技术，可以弥补节点上的硬盘型号和容量的差异化，在硬盘资源层面上实现“异构”。

SmartX 超融合在自动化智能平衡各个节点、各个硬盘的存储资源方面可以提供有效的方法，因此可以支持将不同型号和配置的服务器组成“异构”集群。

多个超融合集群可以通过 SmartX 管理平台 CloudTower 进行统一管理，实现可持续横向扩展，来提供更大的资源池，而不受单一集群节点数限制。

与 VMware SAN 存储相比，SmartX 超融合在单节点内部增加数据盘和缓存盘的灵活度更高，在同集群不同节点间对数据盘和缓存盘容量不一致的宽容度更大，更好实现多架构集群统一管理。

SmartX 分布式存储技术的实现机制。它基于自主研发的分布式文件系统 LSM，将服务器上本地硬盘（SSD + HDD）资源进行池化，实现弹性灵活扩容，并在实际生产环境中得到部署和应用。

内容导读

传统架构的扩容往往难度高且风险大，主要受制于集中式存储。超融合（HCI）将计算虚拟化和分布式存储进行一体化融合部署，与传统的“虚拟化+集中式存储”模式相比，不仅精简了设备层级和数量，简化了配置管理步骤，良好的弹性带来的按需投资能力也成为用户选择超融合的重要因素之一。

相比于业内大部分超融合产品，SmartX 超融合为用户提供了更加弹性、灵活的扩容选择。作为“软硬件解耦”技术路线的坚定实践者，SmartX 的超融合软件适用于主流品牌的服务器，且支持多种主流硬件的兼容。特别是对于由通用机械盘或 SSD 构成的存储资源池，SmartX 的存储管理和 I/O 加速技术表现得更为灵活，允许用户以更大自由度选择服务器和相关设备组件进行扩容。SmartX 超融合提供了以下资源扩容方式：

- **在线添加节点内的存储设备。**
 - 存储扩容的颗粒度细化，即便只增加 1 块 HDD 或 SSD，也能被顺利融入整体资源池。
 - 采用冷、热数据自动分层模式下，缓存层和数据层可分别扩容，支持在线添加、替换存储设备。
 - 对于采用全闪盘的集群，支持“不分层”模式，提高高性能存储设备的利用率，所有 SSD 均可在线添加和替换。
 - 扩容后的数据存储在不影响业务的情况下自动均衡，无需人工干预。
- **在集群内在线添加节点扩展计算和存储资源。**
 - 同一集群内，支持由不同品牌、不同型号和不同存储配置的服务器组成资源池。
 - 扩容后的数据存储在不影响业务的情况下自动均衡，无需人工干预。
- **通过多集群管理平台 CloudTower 实现更大范围的资源池建设和统一管理。**

以下，将详细介绍对集群中的存储资源进行扩容时的要求和注意事项。

SmartX 超融合存储资源配置和扩容方式

SmartX 超融合集群中，每台服务器对本地盘（HDD 和/或 SSD）配置的最低要求如下：

部件	每节点存储配置要求与建议
系统盘 和 缓存盘	<p>2 * 480 GB 或更大容量的数据中心级 SSD</p> <p>注：</p> <ul style="list-style-type: none"> • 通常，SMTX OS 被安装在 2 块 SSD 上并通过软件 RAID1 进行保护，系统只占用每块 SSD 上 125 GB 左右的存储空间，剩余空间则可用作缓存（分层模式）或容量（不分层模式） • 为保证性能，缓存和数据的容量配比需保持在 10% 以上
数据盘	<p>不分层模式</p> <p>2 * 480 GB 及以上的数据中心级 SSD（全闪）</p> <p>分层模式</p> <p>2 * 1 TB 及以上 HDD（混合）或 NVMe/SATA/SAS SSD（全闪）。</p> <p>全闪配置时，管理员可依据 SSD 型号和参数将其设置为缓存盘或数据盘。</p>
启动盘	<p>1 * SSD/HDD，容量小于 2 TB。</p> <p>推荐使用 2 * 240 GB SSD 并使用独立的控制器构建 RAID1。</p>

集群中应包含至少 3 台如上表配置的服务器，这是最小规模的超融合集群。如果在 SmartX 超融合集群中需要更多存储和缓存容量，可以有以下几种扩容方式：

增加节点上的数据盘和缓存盘

在集群的任一节点上新增一块硬盘作为数据盘，该硬盘的容量都将被加入集群的存储资源池。

- 允许每个节点上存在不同容量的硬盘、允许集群上存在硬盘总容量不同的节点，SmartX 分布式存储技术可以智能调节各个节点、每个数据盘上的存储量，以达到节点和硬盘上的存储量平衡。
- 出于存储性能考虑，推荐各个节点上所有数据盘使用同样性能的产品（比如：同为 7200 RPM 的机械盘）。
- 推荐各个节点上所有用作缓存的 SSD 也使用具有同等性能和耐久度的产品（比如：同为 IOPS=50,000 且 DWPD=3）。

在完全基于新购硬件搭建的集群上，很容易做到所有服务器节点上使用型号和容量完全一致的硬盘。但实际情况是，很多用户需要在已有集群上进行扩容，或利用旧服务器重新搭建集群。这就不太容易在所有节点上实现硬盘的完全一致性。SmartX 通过自研的存储管理和调度技术，可以弥补节点上的硬盘型号和容量的差异化，在硬盘资源层面上实现“异构”。

添加节点上的数据盘时，需要注意以下三点：

- SmartX 超融合集群中，每节点所有数据盘之和不能超过 80 TB。
- 缓存盘与数据盘的容量比例不应低于 10%——如果增加了数据层硬盘的总量，则有可能需要相应地

添加/替换缓存盘。

- 单节点最多支持的缓存层总容量为 16 TB。

在每节点上的缓存盘和数据盘都有冗余保护的情况下，可通过管理界面的配合操作，进行逐盘在线替换，不会导致存储的数据丢失或服务中断。

在集群内增加节点数量

在已有的超融合集群内添加更多服务器节点，则可以同时增加集群内部的 CPU、内存和存储资源。或者，如果单服务器节点内部的硬盘已经无法继续扩容，也可以通过在集群中增加节点的方式来进行扩容。

如果同一集群内的服务器型号及组件能够做到完全一致，当然会具有更好的性能和可维护性。但很多用户在对已有集群扩容时已经无法购买到原有的服务器或组件型号，不得不考虑在集群中混用各种服务器的可能性。如前所述，SmartX 超融合在自动化智能平衡各个节点、各个硬盘的存储资源方面可以提供有效的方法，因此可以支持将不同型号和配置的服务器组成“异构”集群。

对新加入集群的服务器，如果不能与集群中现有节点保持完全一致的配置，至少应符合以下要求：

- 必须与原有服务器采用同样的 CPU 架构，但不强制要求使用同样品牌的服务器。
- 所有相关组件必须符合 SmartX 硬件兼容列表的要求。
- 节点上本地存储设备（HDD 和/或 SSD）组成结构（“全闪”或“混合”、“分层”或“不分层”）应与现有集群内的服务器保持一致，但不要求使用同样的型号和容量的存储设备。

注：SmartX 支持的 CPU 品牌为：Intel、AMD、鲲鹏、海光。SmartX 支持的主流服务器品牌为：戴尔、联想、超微、惠普、华为、神州数码、浪潮、新华三、中科可控、超聚变。

用户可以不断向集群内添加新的服务器硬件来扩充集群规模。新服务器硬件往往具备更高的性能和容量密度。

多集群统一管理和虚拟机迁移

SmartX 超融合软件 SMTX OS 单集群最大支持 255 个节点，最大存储裸容量 6PiB。但在单集群中配置很多节点和存储资源，势必会增加管理复杂度；而且用户往往希望基于业务的类型划分不同资源池或希望控制集群规模，降低单集群内多节点同时出现故障的风险，那么扩展为多集群是很好的选择。多个超融合集群可以通过 SmartX 管理平台 CloudTower 进行统一管理，实现可持续横向扩展，来提供更大的资源池，而不受单一集群节点数限制。

CloudTower 可以在一个集中的管理体系内，通过分集群管理，缩小每个集群进行维护操作时的影响范围，实现集群服务水平的提升。

虚拟机可以在 CloudTower 2.0 统一管理的多个集群之间进行迁移。在不同场景下，可以提供热迁移、分段迁移、冷迁移三种模式，详见《SmartX 发布管理平台 CloudTower 2.0 版本》。

不同超融合厂商的扩容方式对比

不同超融合厂商的技术体系和具体实现方式不同，在进行超融合集群扩容时，面临的选型和限制也不相同。下表对比了 SmartX 与 VMware 在超融合集群扩容方面的异同：

单节点内部

存储资源 (HDD/SSD) 组织结构

SMTX OS 5.0 分布式块存储组件

vSAN 7.0

支持全闪存或 SSD + HDD 混合模式。

全闪存配置下支持“分层”或“不分层”模式。

分层模式下：

- 缓存盘上的 OS 和元数据通过软件 RAID1 进行保护
- 缓存盘上的缓存数据和存储日志采用跨节点副本保护
- 系统盘/缓存盘可二合一

<< 左滑查看对比

支持全闪存或 SSD + HDD 混合“硬盘组”。

每节点 1~5 个 硬盘组：

- 盘组内必须采取“分层”模式
- 每盘组 1 个缓存盘，必须为 SSD，仅对本盘组内的 I/O 操作提供缓存服务
- 盘组内的缓存盘无冗余，因此单节点内推荐配置盘组数量 ≥ 2
- 每盘组必须配置 1~7 个容量盘
- 系统盘不能用作缓存盘

单节点内部

是否允许灵活增加数据盘

SMTX OS 5.0 分布式块存储组件

vSAN 7.0

是。

- 推荐所有节点上的数据盘性能一致
- 允许在每个节点内增加不同容量、不同数量的数据盘，节点内总容量不超过 80 TB
- 节点内缓存盘/数据盘容量比：10~20%，缓存盘不足时可灵活添加
- 可以在不同容量的数据盘间自动均衡分配数据存储量

<< 左滑查看对比

不推荐。

- 允许单独增加某一盘组内的容量盘，每盘组不超过 7 块
- 盘组内的缓存盘/容量盘之比：10~20%，增加容量盘可能导致缓存盘不足，且不能灵活增加缓存盘
- 推荐所有硬盘组内的容量盘数量、型号一致，新增硬盘数量与节点数量 x 硬盘组数量成正比（比如：3 节点，每节点 2 硬盘组，至少需增加 6 硬盘）

注：以预设的硬盘利用率 (%) 阈值作为自动均衡条件，不一致的数据盘容量会导致频繁告警和性能下降。

单节点内部

是否允许灵活增加缓存盘

SMTX OS 5.0 分布式块存储组件

vSAN 7.0

是。

- 支持缓存盘在线替换和添加
- 逐一替换缓存盘过程中，数据不丢失、业务不中断
- 节点内总量不超过 16 TB

<< 左滑查看对比

否。

- 缓存盘数量必与硬盘组数量一致
- 不支持缓存盘在线替换：
 - 需将原缓存盘所属盘组整体离线，再更换为更大容量 SSD
 - 或将盘组进行拆分，添加新的缓存盘后重建盘组
- 以上过程需要预先清空该盘组内数据，需较长的数据准备时间；或完成后从副本恢复，期间数据冗余度下降，有风险

同集群不同节点间

是否允许数据盘容量不一致

SMTX OS 5.0 分布式块存储组件

vSAN 7.0

允许。

可在不同节点间智能分配数据存储空间。

允许。

注：以预设的硬盘利用率 (%) 阈值作为自动均衡条件，不一致的数据盘会导致频繁告警和性能下降。

<< 左滑查看对比

同集群不同节点间

是否允许缓存盘容量不一致

SMTX OS 5.0 分布式块存储组件

vSAN 7.0

允许。

缓存盘与本节点上的数据盘比例适当即可。

允许。

推荐在所有节点的盘组上使用具有同样性能指标和容量的缓存盘。

<< 左滑查看对比

多集群

多架构集群统一管理

SMTX OS 5.0 分布式块存储组件

vSAN 7.0

管理平台 CloudTower，可以同时管理多个超融合集群，每个集群内部所有服务器节点应使用一致的 CPU 架构，但不同集群可以采用不同 CPU 架构 (Intel/AMD/鲲鹏/海光)。

管理中心 vCenter，可以同时管理多个超融合集群，每个集群内部所有服务器节点应使用一致的 CPU 架构 (Intel/AMD)。

<< 左滑查看对比

对比小结：

- vSAN 集群中的服务器节点如果不满足盘组一致性、节点一致性要求，会导致集群整体 I/O 性能下降严重，因此 vSAN 集群设计都是以硬件一致性为前提；虽然集群可以短时间内运行在不一致的硬件配置上，但无法以此支持生产级的性能和可靠性。
- 基于 SmartX 超融合构建的集群则为用户提供了多样化存储扩容的选项，并且得到了实际生产环境的验证。

灵活扩容背后的 SmartX 分布式存储技术

SmartX 超融合的弹性灵活扩容，很大程度上来源于 SmartX 分布式存储技术的实现机制。它基于自主研

发的分布式文件系统 LSM，将服务器上本地硬盘（SSD + HDD）资源进行池化。虚拟化存储资源池通过 SmartX 研发的元数据组件进行管理，元数据记录了集群所有节点本地硬盘资源的信息，使得超融合集群中的存储资源调配可以做到更加细粒度、更加精确的控制。这种控制的优势一方面体现在性能方面，另一方面体现在集群的灵活性，如副本分配策略的调整、副本存放位置的选择、数据保存以及 I/O 访问本地化、节点间数据平衡的控制等。

这些特点加强了超融合服务器上的存储资源的灵活性，提高了对不同服务器节点、不同硬盘容量的综合调度能力。集群规模只需 3 个节点起步，IT 运维人员可以在工作时间插入和添加新硬盘，即可完成存储资源的扩展；也可以在不停机的情况下添加服务器节点，同步扩展计算与存储资源，后台将自动地完成资源的池化和平衡，使得资源可以实现“即插即用”。

用户案例

目前，已有越来越多的用户，在实际部署环境中充分利用 SmartX 超融合的灵活扩展能力，不仅实现了资源池的按需投资和扩展，同时通过灵活的选择获得最优的方案配置。

以五矿期货某超融合资源池的硬件扩容与替换为例，从最开始的 4 节点纯软件（基于 SmartX 原生虚拟化 ELF）逐步扩容到 10 节点，先后使用过的服务器类型包括超微四子星、PowerEdge R740xd、PowerEdge R730。在整个过程中，五矿期货在保障业务“0”中断的情况下，在集群扩容的同时完成了对部分服务器的升级替换。五矿期货的 5 个数据中心内的 7 个集群，也通过 CloudTower 实现了跨地域统一资源管理。可阅读《[五矿期货超融合硬件平滑升级与多数据中心管理实战](#)》了解详情。

升级 | 实现 IT 基础架构软硬件升级简单又不停机

[点击链接阅读原文：如何做到 IT 基础架构软硬件升级简单又不停机？](#)

要点总结

企业中 IT 基础架构升级是必然需求，但同时要尽量降低升级停机对企业业务连续性带来的影响。

传统基础架构下的 VMware vSphere 的热迁移功能具有运维投入大、存在业务中断风险、难以弹性投资等缺陷。

SmartX 超融合架构支持软件一键升级功能，并可通过异构扩容和数据迁移实现硬件动态升级，有效降低软硬件升级带来的停机风险，减轻运维人员压力，让企业 IT 技术轻松迭代，助力企业业务持续升级。

SmartX 超融合核心软件 SMTX OS 的软件升级功能具有自动化升级、无中断升级、兼容性保障、升级期间数据恢复最小化的特点。

SmartX 超融合硬件升级过程具有弹性扩展、无中断升级、数据自动均衡的特点。

在 IT 基础架构日常运维中，升级是最头疼的任务之一。这里的升级既包括硬件的固件升级，也包括软件版本升级，还有补丁的升级。这类工作通常伴随着一些潜在停机或者故障的风险，甚至升级操作本身就要求停机执行，这给企业的关键业务带来了不少的麻烦。因此，运维管理员对于升级操作可以说是慎之又慎，能免则免。

但现实中却存在一些难以避免的升级需求，例如：

- 当前使用的软件版本发现明显的漏洞时，企业需要按照监管要求自行整改升级
- 企业使用的硬件设备达到退役年龄，性能、稳定性明显下降
- 企业使用的基础架构在应对特殊场景时性能不佳

因此，企业需要在进行基础架构必要升级的同时，尽量降低升级停机对企业业务连续性带来的影响。在传统虚拟化架构下，一种可行的不会造成业务中断的升级策略是利用 VMware vSphere 的热迁移功能，将虚拟机在开机状态下从原有存储位置迁移至新的存储位置，在这个过程中升级软件或直接完成硬件升级。但这一策略在具体执行时依旧存在以下问题：

- 运维投入大

传统虚拟化架构下，虚拟机的迁移需要一台一台手动完成，每次操作又包含 5-6 个步骤，对于一些有着两三百台虚拟机的大型企业来说，IT 人员需要消耗相当多的时间精力。同时，对于集中式存储架构，升级操作对于运维人员的技术能力要求较高。由于此类升级需要在命令行里面操作，管理员需要足够了解存储的命令行是如何使用的。况且，即使能够做到在不停机的情况下完成基础架构软硬件升级，多数企业——尤其是金融行业——依旧会准备停机升级的应急方案，以保障业务不会中断。这就要求运维人员花费大量时间做升级计划、等待评审会通过方案，使得每一次升级都变成运维人员的“攻坚战”。

同时，对于企业来说，基于 VMware 热迁移升级基础架构的方案会带来额外的资源投入。由于迁移过程中需要用到更多的交换机端口，企业原有的交换机可能无法支持整个迁移工作，需要进行额外采购。而这些设备一般只作临时使用，升级结束后使用机会较少，易造成 IT 资源浪费。

- 业务中断可能性

通过 VMware 热迁移升级基础架构依旧存在一定的业务中断可能性。由于迁移过程涉及较多手动操作，出现人为失误的可能性也大大增加，并最终导致整个升级的失败。

- 难以弹性投资

在进行硬件升级时，企业常常一次性更新整套新设备，对于资源紧张的企业来说，无法做到按需投资、弹性升级。

那么，如何才能在不停机的前提下简单、高效、灵活地实现 IT 基础架构软硬件平滑升级？这项很多运维人员认为不可能完成的任务，SmartX 已经在诸多客户生产环境中实现。与传统虚拟化架构不同，SmartX 超融合架构支持软件一键升级功能，并可通过异构扩容和数据迁移实现硬件动态升级，有效降低硬件升级带来的停机风险，减轻运维人员压力，让企业 IT 技术轻松迭代，助力企业业务持续升级。

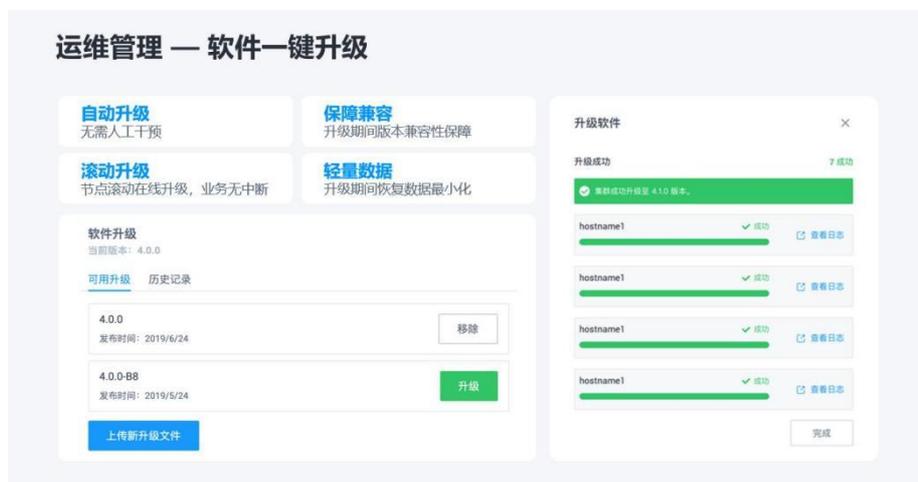
软件一键升级

案例 1

应监管通告要求，某期货公司需要进行 IT 基础架构软件升级。在传统虚拟化架构下，为了不影响业务运行，期货公司需要在深夜或业务外的时间停机并手动完成升级。而 SmartX 超融合软件升级能做到业务“0”中断，支持该期货公司在下午 3 点期货交易结束后的半个小时内开始升级。整个升级过程仅花费 2.5 小时，平均一个节点升级仅需 20 分钟，顺利在下午 6 点下班前完成升级。同时，由于软件升级不需要停机，运维团队仅需内部通过升级方案即可开始升级，免除了复杂的停机审批流程。

这一案例中，客户使用了 SmartX 超融合核心软件 SMTX OS，利用一键升级功能在不停机的情况下完成了集群升级。这一操作的实现有赖于软件升级功能的以下特点：

- **自动化升级**：整个升级过程可在线进行，并预先进行环境检查。能够自动对软件逐一进行升级、重启等操作，减少人工操作带来的差错。
- **无中断升级**：采用滚动升级方式，通过升级控制组件对节点升级进行控制，保证滚动升级正确性且业务无中断。
- **兼容性保障**：SMTX OS 各个版本保证了软件的向后兼容，在升级过程中允许节点间版本不一致，并保证不会对集群已有业务产生影响。
- **升级期间数据恢复最小化**：在保证数据安全性的同时降低数据恢复量，避免集群出现大量数据恢复而造成升级时间过长。



[SmartX 超融合支持软件一键升级](#) (点击阅读案例)

硬件动态升级

案例 2

五矿期货有限公司（以下简称“五矿期货”）是国内注册资本最大的期货公司之一。随着业务的快速发展，五矿期货基于 SmartX 超融合软件先后三次扩容、利旧并升级硬件设备。2018 年，五矿期货利用 SmartX 超融合软件在超微四子星上部署 4 个节点，构建原始集群；2019 年第 1 次扩容，基于 PowerEdge R740xd 部署 2 个节点，实现了不同服务器之间的异构扩容；2020 年基于老旧服务器硬件 PowerEdge R730 部署 4 个节点完成第 2 次扩容；2021 年，通过继续扩容 Dell R740xd，然后利用 SmartX 超融合数据迁移的机制，逐一替换超微四子星。在整个过程中，五矿期货在保障业务“0”中断的情况下，完成硬

件升级替换。



[五矿期货硬件平滑升级流程](#) (点击阅读案例)

这一案例中，企业利用 SmartX 超融合支持集群异构和数据迁移，对节点进行在线扩容并在线替换老旧服务器，实现了基础架构硬件随企业业务发展持续动态升级。这一硬件升级过程包含如下特点：

- **弹性扩展**：3 节点起步，可基于部件或者节点进行扩容，并可整合不同品牌服务器进行异构扩容，整个扩展过程不停机、“0”中断。
- **无中断升级**：利用数据迁移的机制，虚拟机及其副本可快速迁移至其他节点，在全部迁移完成后下线老旧硬件，并在集群中接入新硬件，完成硬件平滑升级。整个过程不停机，且仅在节点迁移及老旧硬件下线时涉及少量手动操作，大幅缩短升级时间，减轻运维压力。
- **数据自动均衡**：新增节点或迁移虚拟机后，动态平衡集群内数据分布，快速恢复分布均衡。

升级 | 新建集群 VS 滚动升级：如何选择服务器硬件平滑升级方案？

点击链接阅读原文：[新建集群 VS 滚动升级：如何选择服务器硬件平滑升级方案？](#)

要点总结

为保证硬件升级时，关键业务可以正常开展、性能和稳定性不受影响，SmartX 提供“新建集群”和“滚动升级”两种方案，帮助企业平稳实现基于超融合架构的服务器硬件替换与升级。

上述两种方案可通过考虑服务器数量、业务连续性、网络资源、集群调整等因素选择适用场景。

在企业 IT 基础架构运维中，经常会遇到以下问题，从而需要对服务器硬件进行更换或升级：

- **服务器达到维护期限**：通常在金融行业中，生产环境的服务器维护期限在 **5 年左右**，超过这一期限，服务器需进行下架。
- **服务器维护成本上升**：服务器使用时间较长，硬件故障或老化会导致性能和稳定性下降，从而增加了企业在人力、物力等方面的运维成本。
- **服务器难以满足业务需求**：随着业务的发展和需求的变化，早期购置的服务器配置无法满足当前的业务需求，升级服务器硬件便需提上日程。

问题是，在硬件升级的同时，运维人员应如何保障关键业务正常开展、性能和稳定性不受到升级影响？针对这一需求，SmartX 为运维人员提供了“新建集群”和“滚动升级”两种方案，帮助企业平稳实现基于超融合架构的服务器硬件替换与升级。下面我们将对两种方案进行详细对比，并通过 2 例实践案例，为用户提供方案选择和落地参考。

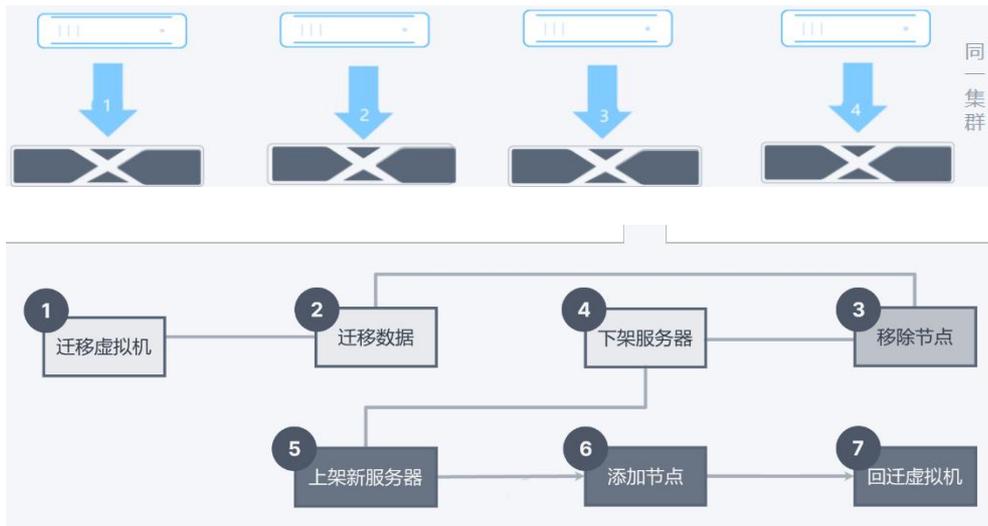
超融合服务器平滑升级方案

方案 1：新建集群



利用新服务器组建一个新集群，将原集群的虚拟机通过跨集群迁移的方式迁移至新集群，从而完成服务器的平滑升级。

方案 2：滚动升级



通过在原有集群中依次对服务器进行替换的方式，实现服务器平滑升级。滚动升级步骤如下：

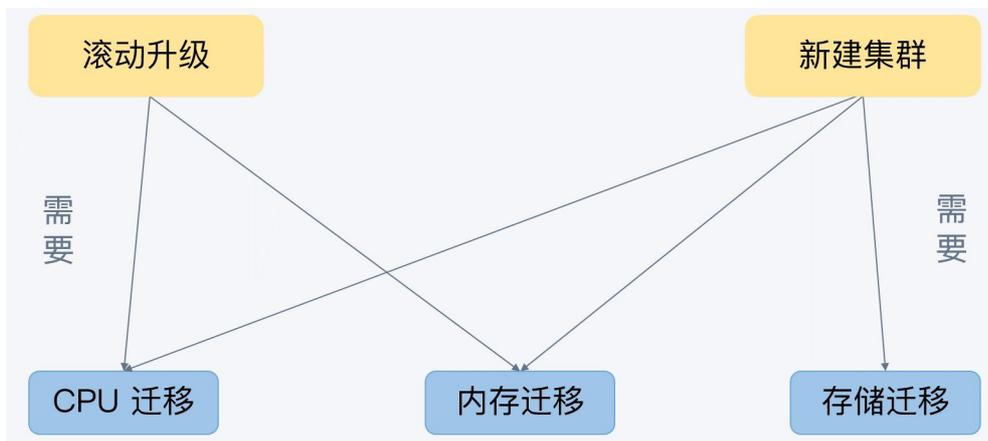
- 迁移虚拟机：将原服务器节点上的虚拟机迁移至集群中其他服务器节点。
- 迁移数据：将原服务器节点上的存储数据迁移至集群中其他服务器节点。
- 移除节点：将原服务器节点从集群中移除。
- 下架服务器：将原服务器节点关机下架。
- 上架新服务器：将新服务器节点加电、连线和上架。
- 添加节点：新服务器节点加入至原集群中。
- 回迁虚拟机：将虚拟机回迁至新服务器节点上。

欲深入了解服务器硬件滚动升级特性与用户实践，请阅读：[如何做到 IT 基础架构软硬件升级简单又不停机？](#)

平滑升级方案对比

以上提到的两种方案皆可实现服务器硬件平滑升级。而两者分别适合什么样的升级环境？企业应如何选择合适的升级方案？我们可以从以下维度进行对比和评估。

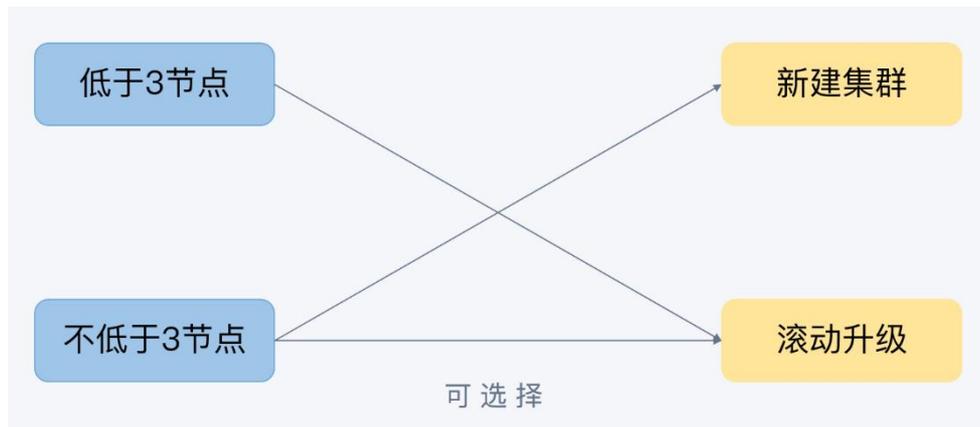
业务连续性



在进行服务器硬件平滑升级时，需保障升级期间集群中的虚拟机业务不受影响。

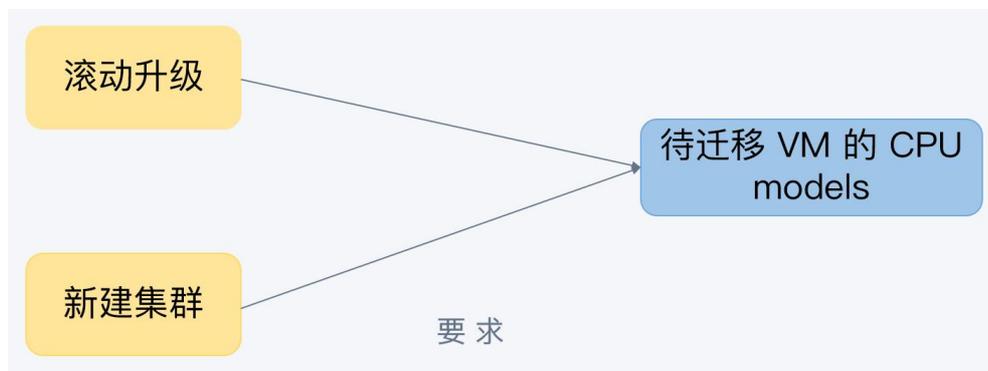
这两种升级方案都涉及了虚拟机迁移操作。在滚动升级方案中，虚拟机迁移仅涉及计算资源迁移；在新建集群的方案中，虚拟机迁移包含了计算资源迁移和存储资源迁移。虽然这两种方案都可做到不影响虚拟机业务，但因新建集群涉及了存储迁移操作，当集群中存在对业务连续性和 I/O 低延迟要求较高的业务时，滚动升级方案会优于新建集群的方式。

服务器数量



滚动升级方案对新服务器数量并无限制，而新建集群方案中，需确保新服务器数量不低于 3 台。因此，当计划对集群中低于 3 台服务器进行升级时，仅能选择滚动升级方案。

虚拟机 CPU 兼容性

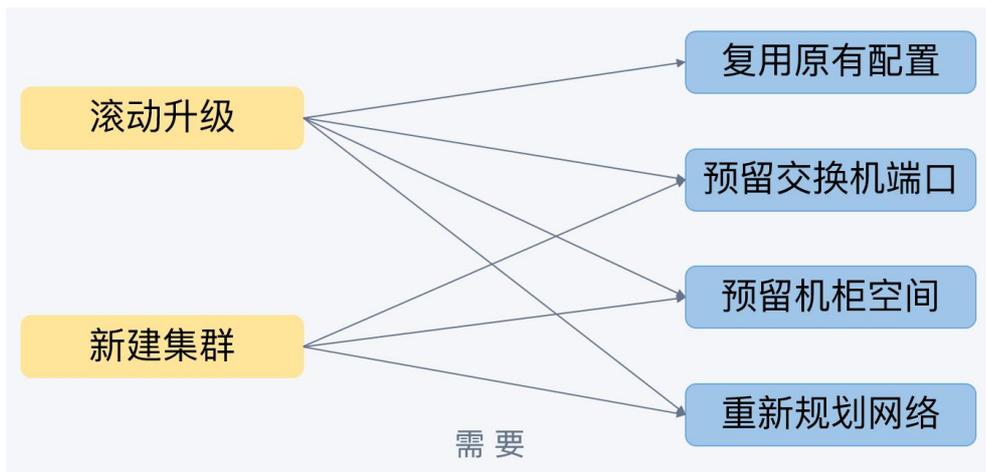


无论是新建集群方案还是滚动升级方案，都需要确保虚拟机可以顺利完成迁移操作。SmartX 集群部署完成后默认会开启虚拟机 CPU 兼容性功能，根据当前宿主机的 CPU 类型和特性，为虚拟机选择一个最接近的 CPU 模型，同时可以让集群中的虚拟机都继承此 CPU 特性。这一功能可以让虚拟机在不同代数 (Generation) 的 CPU 中进行平滑迁移。此外，虚拟机也可自定义选择 CPU 兼容性，比如物理透传或者其他 CPU 的兼容性。

因此，为了确保虚拟机可以顺利完成迁移操作，目标主机或者集群的 CPU model 中必须包含待迁移虚拟机的 CPU model 指令集，并且虚拟机迁移到新集群或者目标主机后，此虚拟机依旧继承迁移前的 CPU model。

如果目标主机或者集群不满足平滑迁移条件，则需要将虚拟机进行关机后再进行迁移。

网络资源

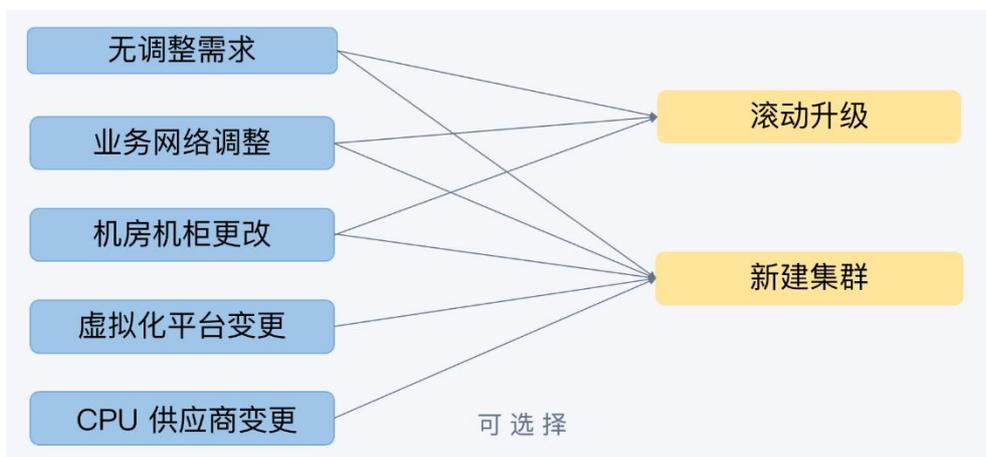


在网络资源方面，滚动升级方案可复用原有配置，而新建集群方案需进行重新配置。这个维度主要考虑，当前集群是否具备新建集群的条件。新建集群需同时满足以下 3 个条件：

- 机房机柜预留了可放置新服务器的空间。
- 交换机预留了管理、存储以及业务网络的端口。
- 新集群有足够的地址为管理、存储以及业务等 IP 地址进行规划。

如果满足，则可以选择新建集群和滚动升级这 2 种方案；如不满足，则选择滚动升级的方式。

集群调整



在进行服务器硬件升级前，用户可能计划对以下方面进行调整，如：

- 业务网络调整：计划将集群中的业务网络和管理网络进行物理层面的隔离。
- 机房机柜更改：计划将服务器放置到 IDC 进行统一管理。
- 虚拟化平台变更：计划将基于 VMware 虚拟化的 SmartX 超融合集群，变更为基于 SmartX 原生虚拟化 ELF 的集群。
- CPU 供应商变更：计划将部分业务迁移至信创集群。

如本次集群调整涉及虚拟化平台和 CPU 供应商的变更，因同一个集群中不能同时存在 2 种虚拟化和 2 种 CPU 供应商，所以需要通过新建集群的方式进行服务器硬件升级。如不涉及这两个方面的变更，那么新建

集群和滚动升级方式皆可供选择。

适用场景

以上提到的两种服务器平滑升级方案并不存在对立的关系，相反，它们在适用场景上存在较多的重合部分。根据以上分析，我们对这两种升级方案在适用场景上的区别进行了以下总结：

考虑维度	滚动升级	新建集群
服务器数量	无节点数要求	不低于 3 节点
业务连续性	业务连续性要求较高且要求 I/O 延迟低	无业务连续性和 I/O 低延迟要求
网络资源	复用现有配置	重新规划
集群调整	无调整需求	存在调整需求

用户案例：方案选择与落地实践

案例一：采用滚动升级方案实现服务器平滑升级

升级背景

- 10 节点 SmartX 超融合（基于原生虚拟化 ELF）集群，单节点存储使用容量为 15TB - 20TB。
- 1 周内需要完成其中 4 台服务器升级。
- 集群存在业务连续性要求较高且要求 I/O 低延迟的业务，升级期间需尽量保障虚拟机业务不受影响。
- 机房无多余机柜空间以及交换机端口，IP 地址段无多余 IP 地址可供分配。

方案选择与实践

用户当前环境**无多余网络资源**，同时由于仅升级集群中的**部分硬件服务器**，应选择滚动升级的方式。采用此方案，一方面可以使新服务器复用原有的服务器网络配置，无需更改网络资源；另一方面，升级部分硬件服务器无需将 1 个集群拆分为 2 个集群，这样可避免增加客户的集群维护工作量。

最终，用户采用滚动升级的方式，在一周时间内，顺利地完成了硬件服务器平滑升级的操作。

案例二：采用新建集群方案实现服务器平滑升级

升级背景

- 8 节点 SmartX 超融合集群，单节点存储使用容量为 12TB - 15TB。
- 3 周内需要完成 8 台服务器升级。
- 8 节点集群被规划为测试集群，机房和集群网络需要被重新调整。
- 在升级期间需尽量保障虚拟机不到影响。

方案选择与实践

用户有**集群调整**的需求，应选择新建集群的方式来进行服务器平滑升级。在这个方案中，新建集群的网络调整以及位置重新放置等操作，对原有集群几乎不产生任何影响，仅需要将原有集群的虚拟机进行跨集群迁移，即可完成全部虚拟机的迁移动作。

最终用户采用此方案，同样在一周时间内，顺利地将 8 个节点的硬件服务器进行了平滑升级。

此外，五矿期货有限公司也利用 SmartX 超融合对异构集群的支持特性，从 4 节点纯软件（基于 SmartX 原生虚拟化 ELF）逐步扩容到 10 节点，同时完成了从超微四子星到 Dell PowerEdge R730xd 的服务器升级替换。更多案例细节，请阅读：[五矿期货超融合硬件平滑升级与多数据中心管理实战](#)。

迁移 | 一文了解 SMTX 迁移工具原理与实践

[点击链接阅读原文：VMware 虚拟机向国产虚拟化平台迁移？一文了解 SMTX 迁移工具原理与实践](#)

要点总结

SmartX 研发跨平台虚拟机迁移工具 SMTX 迁移工具，支持将运行在 VMware 虚拟化平台的虚拟机迁移到基于 SmartX 原生虚拟化 ELF 的超融合集群，顺应国产虚拟化和信创转型趋势。

SMTX 迁移工具支持灵活部署，覆盖虚拟机类型多，迁移过程允许虚拟机保持在线，最大程度降低对业务交付的影响。

SMTX 迁移工具通过自动创建快照的方式完成数据传输，让企业在尽量减少业务中断的情况下进行虚拟机迁移。

SMTX 迁移工具必须连通源端和目标端集群的管理网络，根据不同的网络环境开通相应的 TCP 端口。

随着近些年国产虚拟化和信创转型逐步提上日程，不少客户正在积极寻求 VMware vSphere 虚拟化产品的迁移和替换方案。

作为业内领先的超融合基础设施产品与解决方案提供商，SmartX 为用户提供了跨平台虚拟机迁移工具——SMTX 迁移工具。该工具支持将运行在主流虚拟化平台的虚拟机迁移至基于 SmartX 原生虚拟化 ELF 的超融合集群，帮助用户简单、高效地实现虚拟化平台国产化替代。

为了便于用户更好地理解 SMTX 迁移工具的功能特性，本文将以从 VMware vSphere 虚拟化平台迁移至 ELF 平台为例，浅析迁移原理并展示实践过程。

SMTX 迁移工具

作为一款跨平台虚拟机迁移工具，SMTX 迁移工具具有以下优势：

- SMTX 迁移工具支持**灵活部署**，可选择部署在源端或者目标端虚拟化平台。
- 待迁移的虚拟机**无需安装任何代理插件**，支持的虚拟机类型覆盖 Windows、Linux 等主流的操作系统。
- 使用 SMTX 迁移工具进行迁移时，虚拟机可以**保持在线**，且支持断点续传。
- 帮助客户在有限的业务变更窗口内迁移现有的工作负载，**加速业务交付速度**。

当前发布的 SMTX 迁移工具 1.2.0 版本支持以下 VMware vSphere 和 SMTX OS (ELF) 版本：

平台类型	组件/版本
SMTX OS	3.5.x
	4.0.x
	5.0.x
VMware vSphere	ESXi 组件版本：5.0 / 5.1 / 5.5 / 6.0 / 6.5 / 6.7 / 7.0
	vCenter 组件版本：6.0 / 6.5 / 6.7 / 7.0

迁移技术原理

使用 SMTX 迁移工具迁移虚机的整体流程如下图所示。该工具通过自动创建快照的方式完成数据传输，可让企业尽量减少业务中断的情况下进行虚拟机迁移。在迁移过程中，每一个虚拟机迁移任务都对应一个 task 任务。task 从创建到结束，会经过任务创建、全量迁移、增量迁移、驱动注入、安装 vmtools、配置网络等多个阶段。

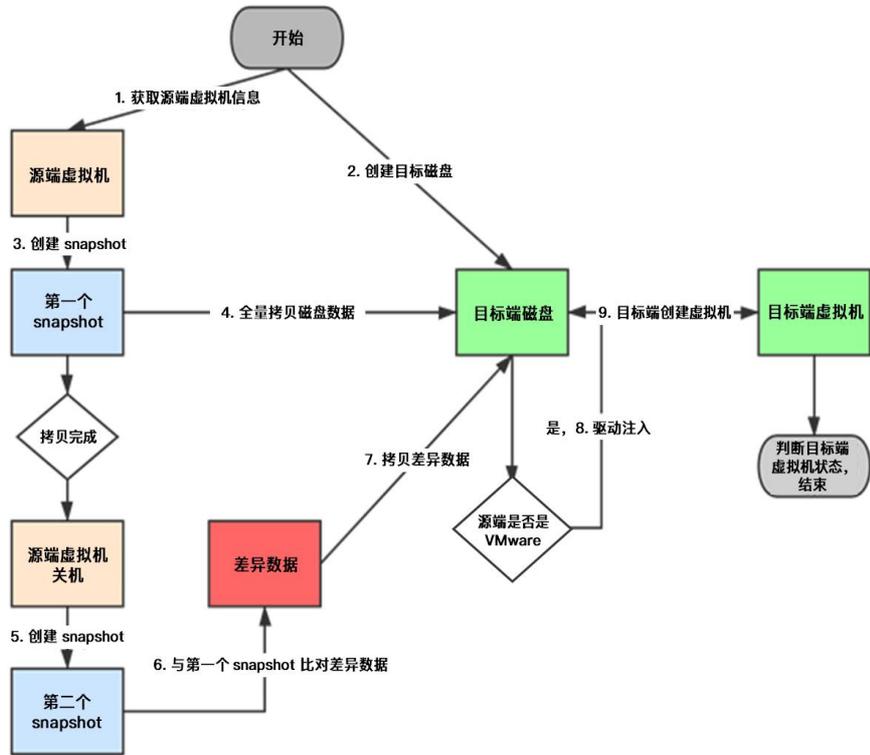


图 1: 迁移整体流程

迁移开始时，首先获取源端待迁移虚机的信息，并在目标端创建目标虚拟磁盘。然后自动创建第一个 snapshot，调用 VMware 的 API 去获取这个快照中有效的磁盘数据区域，执行全量数据迁移至目标虚拟机磁盘。

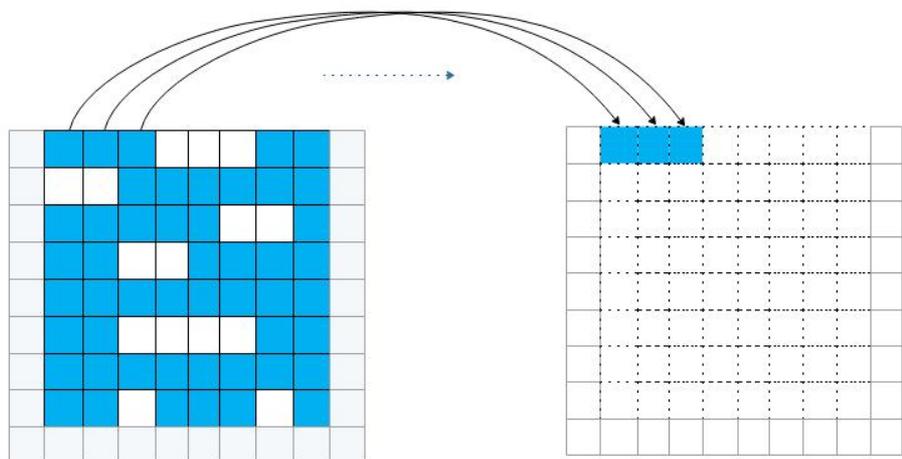


图 2: 全量迁移

全量迁移是首次全量迁移虚拟机快照的磁盘数据的阶段。完成全量迁移后，提示关闭源端虚拟机，创建第二个快照。增量迁移是继全量迁移之后，迁移两次快照的磁盘之间的差异数据的阶段，也叫 Cutover 阶段。当一个 task 的全量迁移阶段结束，此时会判断当前虚拟机的状态，来决定是否现在开启 Cutover 阶段。如果当前虚拟机已经关闭，那么 Cutover 阶段就会立即启动，否则，需要等待用户手动关闭虚拟机之后，再主动发起 http 请求，执行 Cutover 相关逻辑。

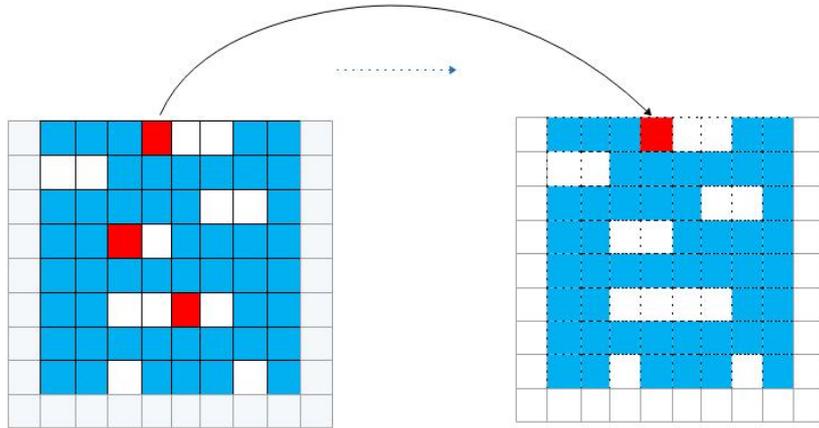


图 3: 增量迁移

关于迁移过程中的数据传输，v2v 工具每次读取 256k 大小的数据，每次读取的数据不会被 v2v 工具缓存，而是立刻被处理。在全量迁移阶段，源端至 v2v 工具需要传输整个有效数据区域，而 v2v 到目标磁盘端只需传输有效数据区域中的非 0 数据块，提升了数据迁移的效率。

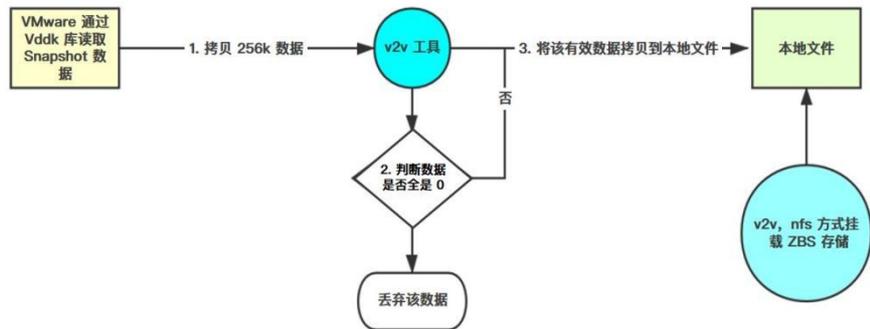


图 4: 数据传输

完成增量数据传输后，判断是否需要注入驱动。VMware 平台上的虚拟机采用的是专属的驱动来支持 Guest OS，ELF 平台的虚拟机采用主流的 Virtio 驱动。迁移工具会自动完成 Virtio 驱动注入。数据迁移完成后，前往目标端站点对虚拟机进行必要的配置和检查，确认迁移后的虚拟机运行正常后，迁移完成。

迁移实践

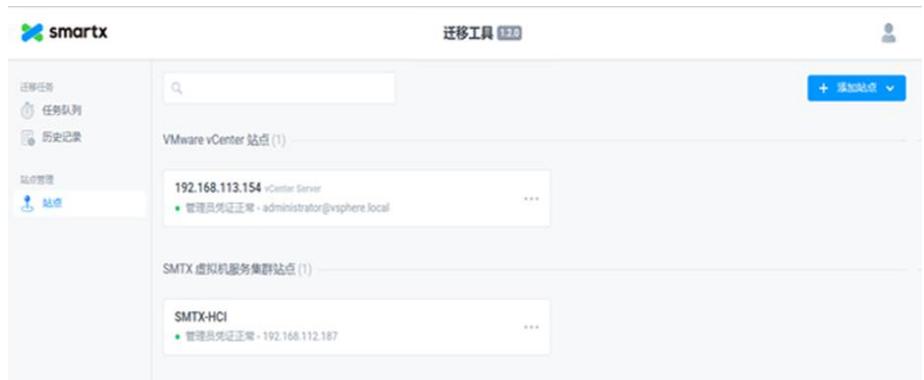
网络环境与要求

SMTX 迁移工具必须连通源端和目标端集群的管理网络。若要加速数据迁移，可以配置 SMTX 迁移工具与源端或目标端的 SMTX OS (ELF) 集群的存储网络连通，以通过存储网络传输数据。若 SMTX 迁移工具与源端或者目标端集群之间存在防火墙，则需要先确保防火墙已开通相应的 TCP 端口（见下表）。

网络环境	需开通的 TCP 端口	描述
SMTX 迁移工具与 SMTX OS (ELF) 集群之间存在防火墙	80 和 443	用于 SMTX 迁移工具与 SMTX 虚拟机服务通信
	10201 - 10206	用于 SMTX 迁移工具与块存储服务之间传输迁移数据
使用 VMware ESXi 平台的集群启用了防火墙	vCenter Server: 80 和 443	用于 SMTX 迁移工具与 vCenter Server 通信
	ESXi: 902	用于 SMTX 迁移工具与 ESXi 之间传输迁移数据

迁移操作

打开 v2v 迁移工具界面，添加源端和目标端站点：选择一个 vCenter Server 站点作为源端，再选择一个 SMTX 虚拟机服务集群，作为目标端。



选择源端站点上待迁移的虚拟机。可以选择按集群或是按主机来检索，已选择的虚拟机会在右侧列出。



迁移工具会自动计算出目标端需要预留的计算和存储资源。

创建迁移任务

选择需要迁移的虚拟机。

The screenshot shows the 'Create Migration Task' interface at Step 2, 'Select VMs to migrate'. On the left, a progress bar indicates the current step. The main area features a search bar and a list of VMs. The selected VM 'Centos7-Test' is highlighted with a blue checkmark. Below the list, 'Data Space' and 'Memory' requirements are shown. On the right, a 'Selected' summary box lists the source and target details.

已选择

源端站点
192.168.113.154
vCenter 数据中心
Datacenter
目标端站点
SMTX-HCI
虚拟机 (1, 共 50 GiB 数据)
主机 192.168.113.153
Centos7-Test

数据空间
预计需要: 200 GiB
2副本 3副本
可用: 23.34 TiB

内存
预计需要: 8 GiB
可分配: 518.05 GiB

< 上一步 取消 下一步

为迁移后的虚拟机网卡指定关联的目标端网络，确保网络的连通性。

创建迁移任务

选择目标端虚拟网络。源端虚拟网络将被映射至目标端虚拟网络。

The screenshot shows the 'Create Migration Task' interface at Step 3, 'Select target virtual network'. The progress bar is updated. The main area has a dropdown menu for 'VM Network' set to 'default'. The 'Selected' summary box on the right remains the same.

已选择

源端站点
192.168.113.154
vCenter 数据中心
Datacenter
目标端站点
SMTX-HCI
虚拟机 (1, 共 50 GiB 数据)
主机 192.168.113.153
Centos7-Test

源端虚拟网络 (共 1 个) 映射的虚拟网络

源端主机 192.168.113.153

VM Network default

< 上一步 取消 下一步

迁移任务创建后会进入任务队列。在等待期间，迁移工具会为虚拟机创建全量快照；等待完毕就会开始数据迁移。这个阶段虚拟机全量快照会被传输至目标端站点。

The screenshot shows the SmartX migration tool interface. The '迁移工具' (Migration Tool) section displays a task 'Centos7-Test' in the '任务队列' (Task Queue) with the status '等待开始' (Waiting to start). The task is associated with the migration from source 192.168.113.154 to target SMTX-HCI.

smartx 迁移工具

任务队列

从 192.168.113.154 迁移至 SMTX-HCI

Centos7-Test 等待开始



当全量快照传输完成后，可以看到关闭源端虚拟机的提示。当确认关闭源端虚拟机后，迁移工具会将上次创建快照后发生的数据变化，以增量快照的形式传输至目标端站点，来完成数据迁移；数据迁移完成后，可前往目标端站点对虚拟机进行必要的配置和检查。



小结

在国产化趋势下，SMTX 迁移工具通过简单高效的迁移方式，帮助客户快速推进国产虚拟化替代 VMware vSphere 的进程。我们将在后续推出更多 SmartX 产品功能展示，敬请期待！

迁移 | 从物理机/云平台迁移至超融合? SMTX CloudMove 帮你实现

[点击链接阅读原文：从物理机/云平台迁移至超融合? SMTX CloudMove 帮你实现](#)

要点总结

为了方便更多用户将不同环境的业务迁移至虚拟化平台，SmartX 发布了 SMTX CloudMove 迁移工具，支持将各类云平台上的虚拟机 (V2V) 或物理机 (P2V) 迁移至基于 ELF 的 SMTX OS 集群，其具备广泛的平台兼容性、迁移不停机、迁移自动化等优势。

SMTX CloudMove 迁移工具适用于以下三大场景：物理机迁移 (P2V)、超融合/虚拟化平台迁移 (V2V)、公有云“下云”迁移。

在《VMware 虚拟机向国产虚拟化平台迁移? 一文了解 SMTX 迁移工具原理与实践》文章中，我们为大家介绍了 SMTX 迁移工具 (支持 V2V 迁移)，以及如何将 VMware vSphere 虚拟化平台上的虚拟机迁移至 SmartX 原生虚拟化 ELF 平台。

为了方便更多用户将不同环境的业务迁移至虚拟化平台，[SmartX 于近期发布了 SMTX CloudMove 迁移工具，支持将各类云平台上的虚拟机 \(V2V\) 或物理机 \(P2V\) 迁移至基于 ELF 的 SMTX OS 集群](#)。下面，我们将为读者详细介绍 SMTX CloudMove 迁移工具及迁移实践。

SMTX CloudMove 迁移工具

产品特性

SMTX CloudMove 迁移工具具备以下优势：

- **广泛的平台兼容性**：相较于 SMTX 迁移工具 (无代理)，SMTX CloudMove (有代理) 不仅支持从虚拟化平台进行迁移，也支持从物理机和公有云平台进行迁移，且不受云平台 Hypervisor 品牌限制，仅对操作系统有兼容性要求，具有更广泛的适用场景和更强的兼容性。
- **迁移不停机**：SMTX CloudMove 可以通过添加源主机信息，对源主机进行连续数据保护 (CDP)，从而实现物理机/虚拟机整机级别的在线迁移。
- **迁移自动化**：迁移时，SMTX CloudMove 控制中心会自动在目标集群创建目标虚拟机，实现源主机到目标虚拟机的自动迁移，人工干预少，自动化效率高。



三大适用场景

物理机迁移 (P2V)

SMTX CloudMove 支持将业务从物理机迁移到 SMTX OS 超融合集群，兼容市面常见的 Windows/Linux 操作系统版本，用户可以通过 SMTX CloudMove 方便地完成 P2V 迁移。

超融合/虚拟化平台迁移 (V2V)

SMTX CloudMove 支持将多种超融合/虚拟化平台 (Nutanix、VMware、华为、新华三、深信服等) 上的虚拟机迁移到 SMTX OS 超融合集群。迁移无需经历导出/导入虚拟机的过程，可降低迁移对存储空间的要求。

公有云“下云”迁移

SMTX CloudMove 支持将公有云 (阿里云、腾讯云、华为云、金山云、天翼云等常见公有云) 云主机下迁移至 SMTX OS 超融合集群。

技术特点

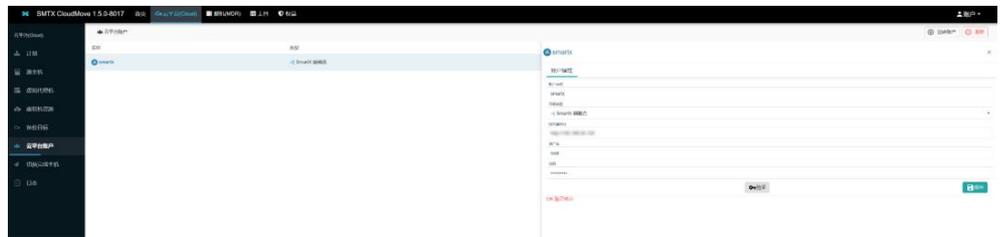
SMTX CloudMove 迁移工具在进行迁移作业时，控制中心会对源主机的配置（CPU、内存、磁盘、网络、操作系统等）进行识别，自动在目标集群创建相同配置的目标虚拟机。用户也可根据自己的需求对硬件配置进行调整，以满足用户定制迁移的需求。SMTX CloudMove 是利用连续数据保护（CDP）技术，实时捕捉源主机上的数据变化，对源主机进行持续数据保护。用户可根据实际业务情况选择合适的时间进行业务切换，切换过程中，控制中心对目标虚拟机自动注入驱动、设置启动选项，配置网络等操作，从而实现物理机/虚拟机整机级别的在线迁移。

迁移实践

下面我们以从阿里云迁移为例，为读者展示如何通过 SMTX CloudMove 将阿里云云主机迁移至 SMTX OS 集群。

step 1 安装控制中心。准备一台物理机或者虚拟机（Ubuntu Server LTS 20.04）安装 SMTX CloudMove 软件，要求能与源主机（阿里云虚拟主机）和 SMTX OS 集群网络连接。

step 2 注册 SMTX OS 集群。登录 SMTX CloudMove 控制中心，添加目标 SMTX OS 集群关联的 CloudTower 信息。



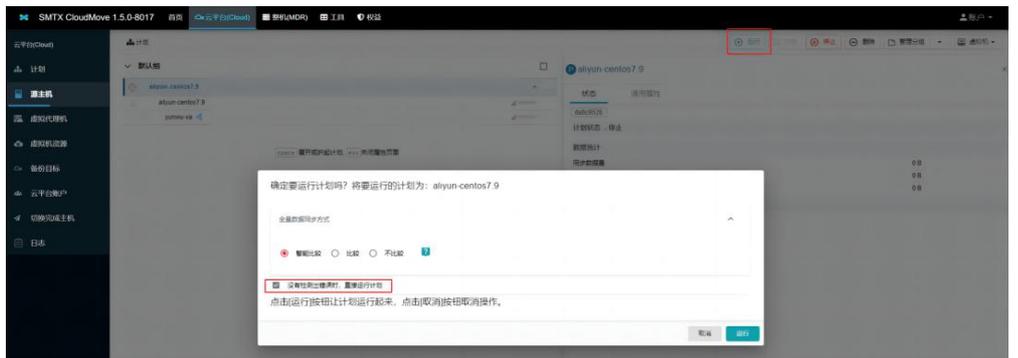
step 3 添加源主机。在 SMTX CloudMove 控制中心添加阿里云云主机的信息，阿里云云主机需启用公网 IP，使之可通过 Internet 被访问。



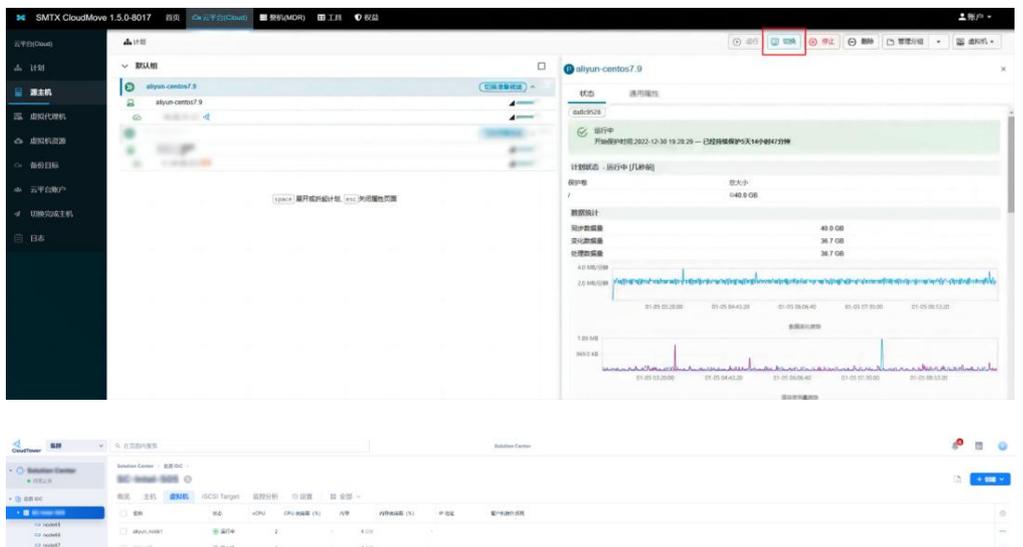
step 4 创建迁移计划。创建阿里云云主机迁移到 SMTX OS 目标集群的迁移计划。



step 5 运行迁移计划。执行迁移计划后，SMTX CloudMove 将自动在 SMTX OS 目标集群创建目标虚拟机。



step 6 切换。当云主机的数据同步完成后，系统将执行切换步骤，完成目标虚拟机驱动自动注入、启动设置，网络配置等操作。最后，开启目标虚拟机电源，关闭源云主机，完成业务系统切换。



step 7 检查业务。检查业务系统是否正常提供服务，如无异常，则迁移成功。

更多资源

下载文档

[SmartX 超融合基础设施及 SMTX Halo 一体机产品介绍](#)

[SmartX 分布式块存储 ZBS 自主研发之旅](#)

[SmartX 行业客户案例集](#)

[行业用户超融合转型实战合集](#)

观看视频

[360 秒了解 SmartX 超融合基础设施](#)

[3 分钟技术解读 —— SMTX OS 副本分配策略](#)

[3 分钟技术解读 —— SMTX OS VM HA 高可用](#)

[SMTX 迁移工具流程讲解与操作实践](#)

阅读博客

[SmartX HCI 5.1 发布：是超融合，更是虚拟化与容器生产级统一架构](#)

更多精彩内容请关注 [SmartX 官网资源中心](#)。

Copyright © 2023 北京志凌海纳科技有限公司 (SmartX) ; 保留所有权利。

本档和本文包含的信息受国际公约下的版权和知识产权的管辖。版权所有。未经 SmartX 事先书面许可, 不得以任何方式, 包含但不限于电子、机械或光学方式对本档的任何部分进行复制, 存储在检索系统中或以任何形式传播。所有非 SmartX 公司名称、产品名称和服务名称仅用于识别目的, 可能是其各自所有者的注册商标、商标或服务标记。所有信息都未获得该所有方的参与、授权或背书。

SmartX 会定期发布产品的新版本。因此, 对于当前使用的某些版本, 本档中介绍的一些功能可能不受支持。有关产品功能的最新信息, 请参阅相关产品的发行说明。如果您的 SmartX 产品未提供本档所述的功能, 请联系 SmartX 以获取硬件升级或软件更新。

您的建议有助于我们提升档内容的准确性或组织结构。将您对本档的意见发送到 info@smartx.com 来帮助我们持续改进本档。