

安全容器服务 使用手册

产品版本: v6.2.1

发布日期: 2023-07-01

目录

1 版本说明	1
1.1 版本说明书	1
2 产品介绍	3
2.1 什么是安全容器服务	3
2.2 使用场景	5
2.3 基本概念	6
2.4 产品获取	9
2.5 权限说明	10
2.6 使用限制	14
2.7 与其他服务的关系	16
3 快速入门	17
3.1 操作指引	17
3.2 前置条件准备	19
3.3 创建命名空间	21
3.4 创建工作负载	22
3.5 创建Ingress（可选）	33
3.6 使用Yaml创建资源（可选）	34
4 用户指南	36
4.1 集群管理	36

4.2 节点管理	37
4.3 命名空间	39
4.4 存储管理	40
4.5 业务概览	42
4.6 工作负载	43
4.7 持久卷声明	50
4.8 配置中心	52
4.9 网络管理	55
4.10 自定义资源管理	59
5 最佳实践	61
5.1 配置SR-IOV网络	61
5.1.1 前置条件准备	61
5.1.2 使用Yaml配置SR-IOV网络	70
5.1.3 runc容器和InfiniBand卡IB模式场景	77
5.1.4 runc容器和InfiniBand卡ETH模式场景	80
5.1.5 rune容器和InfiniBand卡IB模式场景	83
5.1.6 rune容器和InfiniBand卡ETH模式场景	86
5.2 创建GPU资源的容器实例	89
6 常见问题	95
6.1 如何理解安全容器网络方案	95
6.2 容器状态为未知错误，如何排查解决	101
6.3 容器状态为失败，如何排查解决	102

7 部署指南	103
7.1 GPU解决方案对接包使用指南	103
8 API参考	113
8.1 API文档模板	113

1 版本说明

1.1 版本说明书

版本信息

产品名称	产品版本	发布日期
安全容器服务	V6.2.1	2023-06-30

更新说明

新增功能

- 支持SR-IOV网卡功能，且支持使用切分后的Virtual Functions作为容器（Pod）的第二块网卡。
- 支持Ethernet和InfiniBand两种模式的SR-IOV网卡。
- 支持安全容器和传统runc容器使用SR-IOV网卡能力。
- 支持由Mellanox公司生产的，支持InfiniBand通信标准的网络接口卡。
- 支持不再通过virtio-fs文件系统，而是直接通过virtio-blk块设备将PVC提供给安全容器Guest VM，这项改进将提升存储性能。
- 支持百度昆仑XPU设备能力。

优化功能

- QEMU组件版本升级。
- 安全容器运行时版本升级。
- 安全容器Guest VM Kernel版本升级。
- Kube-OVN CNI VPC的对应路由器已默认设置为连接外部网络。
- 优化了在部署中使用单节点读写（RWO）数据卷或设置自动弹性伸缩策略时的副本数校验逻辑，提升操作体验。

已修复问题

- 修复了[存储管理]-[存储类]界面中内部存储local-disk的显示问题。
- 修复了有状态副本集创建时，挂载高性能存储卷类型数据卷配额显示错误的问题。
- 修复代码，解决因为设置sandbox_cgroup_only导致的cgroup泄露问题。
- 修复VM镜像，解决SR-IOV网卡驱动和GPU驱动冲突的问题。
- 修复了安全容器挂载GPU后，负载显示GPU未使用的问题。
- 修复了部署创建成功后，Pod状态持续显示为启动中的问题。

依赖说明

- 平台版本至少为v6.1.1。
- SDN网络服务版本至少为v6.2.2。

2 产品介绍

2.1 什么是安全容器服务

安全容器服务基于成熟、轻量的安全容器运行时和SDN网络服务，提供卓越的不可信应用隔离、故障隔离、性能隔离以及多租户应用网络隔离等能力，以使用户轻松高效地在云端运行安全容器化应用。

产品优势

- **安全且故障隔离**

基于安全容器运行时，提供超强的不可信应用隔离、故障隔离等能力。

- **云资源网络互通**

安全容器与计算、存储、网络等资源内网互通，以便容器可以分配到虚拟网卡、公网IP和负载均衡等资源，同时也方便传统云主机应用与容器应用之间网络互通。

- **标准适配**

在网络、日志、监控、存储等方面有着和普通容器一样的用户体验，并具备极速启动和优秀的兼容性、稳定性等特点。

- **网络隔离**

基于SDN网络服务，在安全容器运行时上增加多租户应用网络隔离能力。

- **统一权限管理**

为防止资源误操作，将资源的操作能力和云平台的授权管理服务结合，一体化实现。

- **统一配额管理**

为防止资源滥用，将资源的配额管理能力和云平台的配额管理结合，一体化管理。

主要功能

- **配额管理**

为防止资源滥用，云平台支持设置安全容器相关资源的配额，对各项目的可用资源数量和容量做出限制，配额项包括CPU、内存、存储容量、GPU等。

- **容器负载**

支持对部署容器实例的全生命周期管理，包括启动/停止、重新部署、配置更新、历史版本回滚、终端操作、查看监控与日志、删除等操作。

- **弹性伸缩**

基于HPA（Horizontal Pod Autoscaler）能力，根据容器当前使用的CPU与内存压力自动扩缩容。

- **滚动升级**

当通过控制器部署多副本的工作负载时，支持设置自定义滚动更新策略。

- **负载均衡**

当用户服务是通过控制器部署时，使用负载均衡可将传入流量分配到部署中的各个容器实例，当部署发生变化时，云平台会自动从负载均衡器中添加和删除实例。

- **GPU调度**

提供NVIDIA GPU设备的发现与管理能力，云平台在创建容器时将依据指定GPU使用需求自动调度GPU资源。

- **持久化存储**

提供数据持久化存储满足容器运行过程中需要保存数据的需求，并支持普通容量型以及高性能型两种存储类型（使用高性能存储类型时需要搭配高性能云存储产品）。在创建工作负载时支持添加多个存储卷，以及为每个存储卷指定存储类型和容量，并支持挂载存储卷到容器的指定路径。

2.2 使用场景

- **替换传统虚拟机业务**

传统虚拟机部署应用，虽然安全性较高，但无法享用到镜像和容器带来的技术红利，并且传统虚拟机的损耗开销较大，交付效率分钟级，难以脱离虚拟机镜像，网络自建和交付不统一等难题，安全容器为解决以上痛点，通过精简的虚拟机，启动在秒级内，复用容器管理平台上的多种资源，包括CNI，CSI等通用化网络，存储方案。

- **隔离不可信应用与故障**

由于在同一节点中，普通容器通常都混部着不同的业务和租户应用，这些容器都共享同一内核。所以，当内核或者运行时出现漏洞时，恶意代码将会逃逸到对宿主机产生不可逆影响，甚至会导致系统瘫痪。安全容器服务提供超强的不可信应用隔离、故障隔离、性能隔离以及多租户应用网络隔离等能力，保障用户轻松高效地在云端运行安全容器化应用。

- **业务应用运行独占操作系统内核**

安全容器服务提供内核级的进程隔离机制，天然满足容器独占内核的需求。

2.3 基本概念

集群 (Cluster)

一个集群指容器运行所需要的云资源组合，关联了若干服务器节点、存储、网络等基础资源。

节点 (Node)

安全容器集群中的节点包括Master节点和Worker节点两种类型，每一个节点对应一个云主机。Master节点是安全容器集群的管理者，运行着一些用于保证集群正常工作的组件，如 kube-apiserver、kube-scheduler等。Worker节点是安全容器集群中承担工作负载的节点，承担实际的 Pod 调度以及与管理节点的通信等。一个 Worker节点上运行的组件包括containerd运行时组件、kubelet、Kube-Proxy等。

命名空间 (Namespace)

在同一个集群内可以创建不同的命名空间，不同命名空间中的数据彼此隔离，使它们既可以共享同一个集群的服务，也能够互不干扰，为集群提供资源逻辑隔离作用。

容器组 (Pod)

容器组即Pod，是安全容器服务部署应用或服务的最小的基本单位。一个容器组封装多个容器（也可以只有一个容器）、存储资源、网络资源以及管理控制容器运行方式的策略选项。

工作负载

工作负载是安全容器服务对一组Pod的抽象模型，用于描述业务的运行载体，包括部署 (Deployment)、有状态副本集 (StatefulSet)、守护进程集 (DaemonSet)、任务 (Job)、定时任务 (CronJob)。

- 部署：即Kubernetes中的“Deployment”，部署支持弹性伸缩与滚动升级，适用于容器组完全独立、功能相同的场景，如nginx。
- 有状态副本集：即Kubernetes中的“StatefulSet”，有状态副本集支持容器组有序部署和删除，支持持久化存储，适用于实例间存在互访的场景，如ETCD等。
- 守护进程集：即Kubernetes中的“DaemonSet”，守护进程集确保全部（或者某些）节点都运行一个容器组，支持容器组动态添加到新节点，适用于容器组在每个节点上都需要运行的场景，如fluentd、Prometheus

Node Exporter等。

- 任务：即Kubernetes中的“Job”，任务是一次性运行的短任务，部署完成后即刻执行。
- 定时任务：即Kubernetes中的“CronJob”，定时任务是按照指定时间周期运行的任务。

安全工作负载

安全工作负载拥有独立的操作系统内核以及安全隔离的虚拟化层。通过安全工作负载，不同容器之间的内核、计算和网络资源均相互隔离，保护Pod的资源 and 数据不被其他Pod抢占和窃取。

服务 (Service)

由于每个容器组都有自己的IP地址，并且可能随时被删除重建，如果这个容器组要为其它容器组提供服务，则如何找出并跟踪要连接的IP地址会非常麻烦。安全容器服务针对这个问题给出的方案是服务 (Service)。Service是将运行在一组Pods上的应用程序公开为网络服务的抽象方法。使用安全容器服务，您无需修改应用程序即可使用不熟悉的服务发现机制。安全容器服务为Pods提供自己的IP地址和一组Pod的单个DNS名称，并且可以在它们之间进行负载均衡。

路由 (Ingress)

Ingress是一组将集群内服务暴露给集群外服务的路由规则集合。一个ingress对象能够配置具备为服务提供外部可访问的URL、负载均衡流量、卸载 SSL/TLS，以及提供基于名称的虚拟主机等能力。

持久化存储

持久卷 (PV)

持久卷描述的是持久化存储卷，主要定义的是一个持久化存储在宿主机上的目录，独立于容器组生命周期。具体到本平台，一个持久卷对应一个云硬盘。

持久卷声明 (PVC)

持久卷是存储资源，而持久卷声明 (PVC) 是对持久卷的请求。持久卷声明跟容器组类似：容器组消费节点资源，而持久卷声明消费持久卷资源；容器组能够请求CPU和内存资源，而持久卷声明请求特定大小和访问模式的持久卷。

存储类 (StorageClass)

存储类可以实现动态供应持久卷，即能够按照用户的需要，自动创建其所需的存储。

配置（ConfigMap）

ConfigMap用于将非机密性的数据保存到键值对中。使用时，容器组可以将其用作环境变量、命令行参数或者存储卷中的配置文件。ConfigMap将环境配置信息和容器镜像解耦，便于应用配置的修改。

密钥（Secret）

密钥（Secret）是一种包含认证信息、密钥等敏感信息的资源类型，可以用作工作负载的环境变量、加密配置文件。将数据放在密钥对象中，可以更好地控制它的用途，并降低意外暴露的风险。

标签（Label）

标签是一对 key/value，被关联到对象上，比如节点、容器组。通过标签可以方便地标识及筛选对象。

2.4 产品获取

前提条件

在执行下述产品获取操作步骤前，请确保以下条件均已满足：

- 已成功获取并安装“计算服务”、“块存储”、“SDN网络服务”和“容器镜像服务”云产品。获取并安装云产品的具体操作说明，请参考“产品与服务管理”帮助中的相关内容。
- 如需获取正式版云产品，请提前将已获取的许可文件准备就绪。

操作步骤

1. 获取并安装“安全容器服务”云产品。

在顶部导航栏中，依次选择[产品与服务]-[产品与服务管理]-[云产品]，进入“云产品”页面获取并安装“安全容器服务”云产品。具体的操作说明，请参考“产品与服务管理”帮助中“云产品”的相关内容。

2. 访问安全容器服务。

在顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[任意子菜单]，即可访问该服务的各项功能。

2.5 权限说明

本章节主要用于说明安全容器服务各功能的用户权限范围。其中，√代表该类用户可对云平台内所有项目的操作对象执行此功能，**XX项目**代表该类用户仅支持对XX项目内的操作对象执行此功能，未标注代表该类用户无权限执行此功能。

功能		云管理员	部门管理员/项目管理员	普通用户
集群管理	信息展示	√		
节点管理	信息展示	√		
	开始/停止调度			
	标签管理			
	污点管理			
命名空间	信息展示	√	仅已加入项目	
	创建命名空间			
	删除			
存储管理	信息展示	√	仅已加入项目	
	查看存储类Yaml			
	查看持久卷Yaml			
	删除持久卷			
工作负载	信息展示	√	仅已加入项目	仅已加入项目
	创建部署			
	创建有状态副本集			
	创建守护进程集			
	创建任务			

	功能	云管理员	部门管理员/项目管理员	普通用户
	创建定时任务			
	容器配置			
	手动伸缩			
	版本回滚			
	升级策略			
	伸缩策略			
	调度策略			
	网络设置			
	标签设置			
	编辑Yaml			
	启动/停止			
	重新部署			
	删除			
	运行/停止定时任务			
	查看容器组Yaml			
	容器日志			
	容器终端			
	删除容器组			
持久卷声明	信息展示	√	仅已加入项目	仅已加入项目
	创建持久卷声明			
	编辑Yaml			
	删除			

功能		云管理员	部门管理员/项目管理员	普通用户
配置中心	信息展示	√	仅已加入项目	仅已加入项目
	创建配置			
	更新配置			
	编辑配置Yaml			
	删除配置			
	创建密钥			
	更新密钥			
	编辑密钥Yaml			
	删除密钥			
网络管理	信息展示	√	仅已加入项目	仅已加入项目
	创建服务			
	更新服务			
	编辑服务Yaml			
	删除服务			
	创建Ingress			
	更新Ingress			
	编辑Ingress Yaml			
	删除Ingress			
自定义资源管理	信息展示	√	仅已加入项目	仅已加入项目
	使用Yaml导入自定义资源描述			

	功能	云管理员	部门管理员/项目 管理员	普通用户
	使用Yaml导入自定义资源		仅已加入项目	仅已加入项目
	删除自定义资源		仅已加入项目	仅已加入项目

2.6 使用限制

- 在Arm架构的云平台上，容器不支持使用GPU。
- 目前云平台支持在安全容器中使用英伟达GPU和百度昆仑XPU。

英伟达 (NVIDIA) GPU

英伟达 (NVIDIA) GPU均直接采用预装的450.80.02版本的NVIDIA GPU驱动，并且该驱动不支持卸载。该GPU驱动兼容的CUDA版本和支持的GPU设备如下表所示：

类型	型号
英伟达 (NVIDIA) GPU设备	Tesla V100
	Tesla P100
	Tesla P40
	Tesla P6
	Tesla P4
	Tesla M60
	Tesla M10
	Tesla M6
	Tesla T4
	Quadro RTX 8000
	Quadro RTX 6000
CUDA版本	CUDA 11.2及以下

百度昆仑XPU

百度昆仑XPU驱动支持的GPU设备如下表所示：

类型	型号
百度昆仑XPU设备	R200-8F

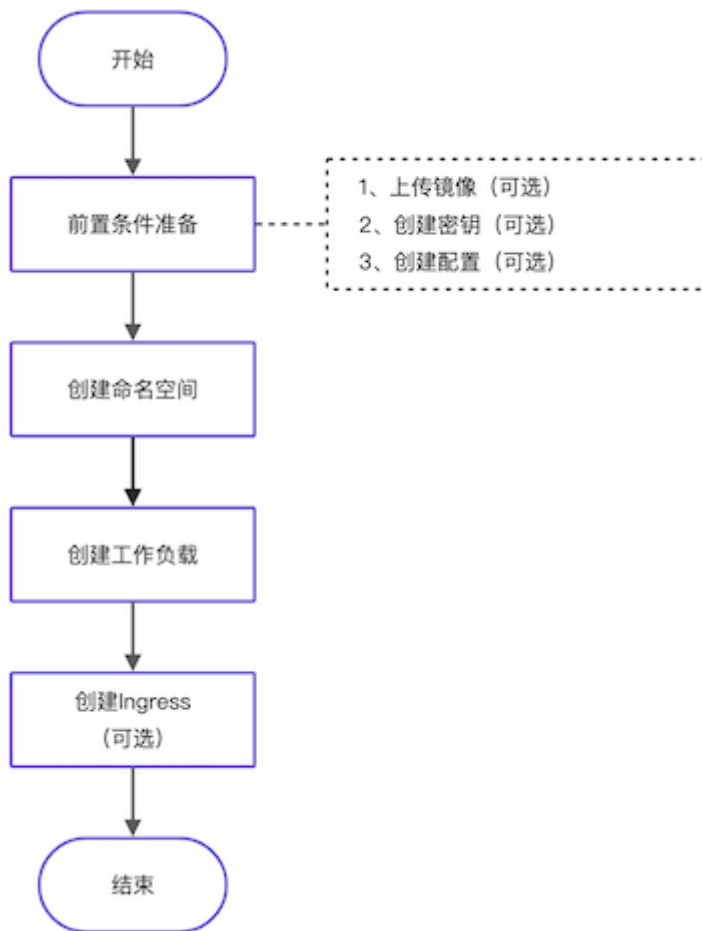
2.7 与其他服务的关系

服务	关系说明
容器镜像服务	创建工作负载时需要为容器指定所使用的容器镜像。
块存储	块存储为容器集群提供持久化存储资源。
SDN网络服务	为安全容器服务提供网络、公网IP、负载均衡等网络资源及相关服务。

3 快速入门

3.1 操作指引

安全容器服务云产品的主线使用流程及具体说明如下：



操作流程		描述
前置条件准备	上传镜像 (可选)	预先上传工作负载的容器创建时所需要的镜像文件。请根据客户实际业务需求酌情创建。如已有可用镜像或使用第三方镜像时，可跳过本步骤。

操作流程		描述
	创建密钥（可选）	预先创建工作负载创建时所需要的密钥。 请根据客户实际业务需求酌情创建。如不使用开启密钥认证的第三方镜像，且数据卷、环境变量都不选择“密钥”类型时，可跳过本步骤。
	创建配置（可选）	预先创建工作负载创建时所需要的配置。 请根据客户实际业务需求酌情创建。如数据卷和环境变量都不选择“配置”类型时，可跳过本步骤。
创建命名空间		通过命名空间实现同一集群内不同资源之间的隔离。
创建工作负载		工作负载是对一组Pod的逻辑抽象，用于承载业务运行。
创建Ingress（可选）		通过Ingress为工作负载的服务提供外部访问时所需的路由规则集合。 请根据客户实际业务需求酌情配置。当工作负载已添加服务，且该服务需要配置对外访问的路由规则时，才需执行此操作。

3.2 前置条件准备

在创建工作负载前，请先完成以下准备工作。

上传镜像（可选）

本操作用于预先上传工作负载的容器创建时所需要的镜像文件，请根据客户实际业务需求酌情创建。如已有可用镜像或使用第三方镜像时，可跳过本步骤。

1. 在云平台的顶部导航栏中，依次选择[产品与服务]-[容器镜像服务]-[镜像管理]，进入“镜像管理”页面。
2. 单击 **上传镜像** 或 **Push镜像** ，弹出对应的对话框。
3. 配置参数后，完成操作。各参数的具体说明，请参考“容器镜像服务”帮助中“镜像管理”的相关内容。

创建密钥（可选）

本操作用于预先创建工作负载创建时所需要的密钥，请根据客户实际业务需求酌情创建。如不使用开启密钥认证的第三方镜像，且数据卷、环境变量都不选择“密钥”类型时，可跳过本步骤。

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[配置中心]，进入“配置中心”页面。
2. 在左侧导航栏选择[业务视图]，选择目标命名空间，选择[配置中心]-[密钥]，进入“密钥”页面。
3. 单击 **创建密钥** ，进入“创建密钥”页面。
4. 配置参数后，单击 **创建** 完成操作。

参数	说明
----	----

参数	说明
密钥类型	<ul style="list-style-type: none"> * Opaque：一般密钥类型。 * TLS：存放7层负载均衡服务所需的证书。 * 镜像访问密钥：存放拉取私有仓库镜像所需的认证信息。
密钥数据	<ul style="list-style-type: none"> * 当密钥类型为Opaque时，单击“添加密钥数据”，输入键、值。 * 当密钥类型为TLS时，上传证书和私钥文件。 * 当密钥类型为镜像访问密钥时，输入镜像仓库地址、用户名、密码和邮箱。

创建配置（可选）

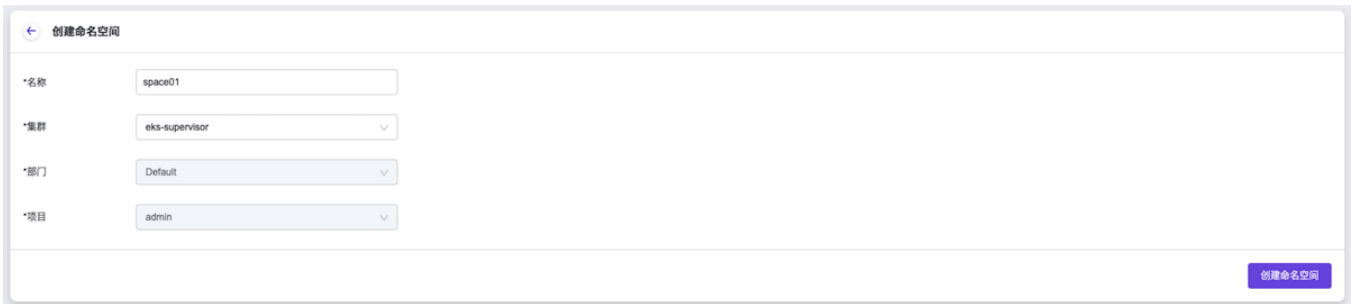
本操作用于预先创建工作负载创建时所需要的配置，请根据客户实际业务需求酌情创建。如数据卷和环境变量都不选择“配置”类型时，可跳过本步骤。

1. 在左侧导航栏选择[业务视图]，选择目标命名空间，选择[配置中心]-[配置]，进入“配置”页面。
2. 单击 **创建配置**，进入“创建配置”页面。
3. 填写参数后，单击 **创建配置**，完成操作。

3.3 创建命名空间

通过命名空间实现同一集群内不同项目资源之间的隔离。

1. 在顶部导航栏选择[产品与服务]-[安全容器服务]-[命名空间]，进入“命名空间”管理页面。
2. 单击 **创建命名空间** ，跳转至“创建命名空间”页面。
3. 配置参数，单击 **创建命名空间** 完成操作。



参数	说明
名称	选择命名空间的名称。
集群	选择命名空间所属集群。
部门/项目	用户所在部门和项目，不支持修改。

3.4 创建工作负载

工作负载是对一组Pod的逻辑抽象，用于承载业务运行。其类型包括部署（Deployment）、有状态副本集（StatefulSet）、守护进程集（DaemonSet）、任务（Job）、定时任务（CronJob），请根据客户实际业务需求酌情创建。创建方式支持界面创建和Yaml创建，本节将介绍界面创建方式，Yaml创建方式请参见 [如何使用Yaml创建资源](#)。

1. 在左侧导航栏选择[业务视图]页签-选择目标命名空间后，依据工作负载类型选择对应子菜单，进入对应页面。
2. 单击 **创建部署/有状态副本集/守护进程集/任务/定时任务** ，进入对应创建页面的“容器配置”页面。
3. 在“容器配置”页面中，配置参数后，单击 **下一步：访问方式** ，进入“访问方式”配置页面。其中，在“容器配置”区域框中，单击 **添加容器** ，可在该容器实例中添加多个容器，但是在容器添加过程中，请先确保已完成当前容器的配置。

← 创建部署

① 容器配置 ② 访问方式 ③ 高级配置

*容器运行时 安全运行时 runc运行时

*安全负载名称

*副本数

容器配置

*容器名称

容器类型 业务容器 初始化容器

镜像来源

*镜像 [选择镜像](#)

*镜像版本

拉取镜像策略 本地不存在时拉取 总是拉取

*资源预留 CPU 内存 Mi

*资源限制 CPU 内存 Mi

GPU 使用GPU
开启后，所有增量安全容器默认开启 GPU 能力，共享 GPU 卡资源配置。

环境变量 [+ 添加环境变量](#)

数据卷 [+ 添加数据卷](#)

健康检查 [展开](#)

安全设置 [展开](#)

命令 [展开](#)

日志采集 [展开](#)

配额

下一步: 访问方式 [Yami](#)

参数	说明
----	----

参数		说明
容器运行时		<p>该工作负载中安全容器的运行时类型。该参数值可选“安全运行时”或“runc运行时”。</p> <p>运行安全运行时的工作负载与运行runc运行时的工作负载相比，其进程隔离机制为内核级隔离，安全容器间的计算资源、网络资源具有更为彻底的隔离性。</p> <p>runc运行时不支持内核隔离，不支持使用GPU。</p>
副本数		<p>仅当负载类型为“部署”或“有状态副本集”时可配置此参数。</p> <p>表示该工作负载包括的容器组个数。每个容器组都由相同的容器部署而成。设置多个容器组主要用于实现高可靠性，当某个实例故障时，工作负载还能正常运行。</p>
容器配置	容器类型	<p>包括业务容器和初始化容器。业务容器即真正运行业务的容器，初始化容器则运行于业务容器启动期间。若容器组中有多个初始化容器，这些容器会按顺序逐个运行，每个初始化容器必须运行成功，下一个才能够运行，当所有初始化容器运行完成时，集群才会正常运行业务容器。由于一个容器组中的存储卷是共享的，所以初始化容器中产生的数据可以被业务容器使用到。由于初始化容器提供了一种机制来阻塞或延迟业务容器的启动，可以应用于有启动顺序要求的容器组之间。</p>
	镜像来源	<p>包括镜像仓库和第三方镜像两种来源。选择镜像仓库则使用本集群对接的镜像仓库，选择第三方镜像则需要输入第三方镜像地址且保证网络可达。</p>
	密钥认证	<p>仅当镜像来源为“第三方镜像”时可配置。</p>
	密钥	<p>仅当镜像来源为“第三方镜像”且密钥认证为“是”时可配置。</p>
	镜像	<p>若镜像来源为“镜像仓库”，则单击 选择镜像，弹出选择镜像对话框。选择目标镜像，单击 确定 完成操作。若镜像来源为“第三方镜像”，则输入格式为ip:port/path/name的镜像地址。</p>
	镜像版本	<p>若镜像来源为“镜像仓库”，则在下拉框中选择目标版本；若镜像来源为“第三方镜像”，则手动输入目标版本。</p>

参数		说明
	拉取镜像策略	包括“本地不存在时拉取”和“总是拉取”两种策略。
	资源预留	保证容器成功调度到节点的最小资源。 当需要勾选“使用GPU”时，建议此参数值的CPU大于等于1，内存大于等于1024MiB。
	资源限制	容器运行中允许使用的最大资源。
	使用GPU	仅当容器运行时为“安全运行时”时可配置此参数。表示容器是否使用GPU资源。 “守护进程集（DaemonSet）”和“定时任务（CronJob）”类型的工作负载不支持使用GPU。
	环境变量（可选）	容器在启动过程中需要的一些配置信息如启动命令、证书等，这类信息需要在容器组故障重启后仍然存在并重新加载到新容器组中，这类信息可以通过环境变量的形式单独存储。当前支持以下类型： * 普通变量：普通变量不需提前创建，直接输入即可。 * 配置：选择已创建好的配置。 * 密钥：选择已创建好的密钥。 * Pod字段：直接选择具体字段即可。 * 容器资源：直接选择具体资源即可。

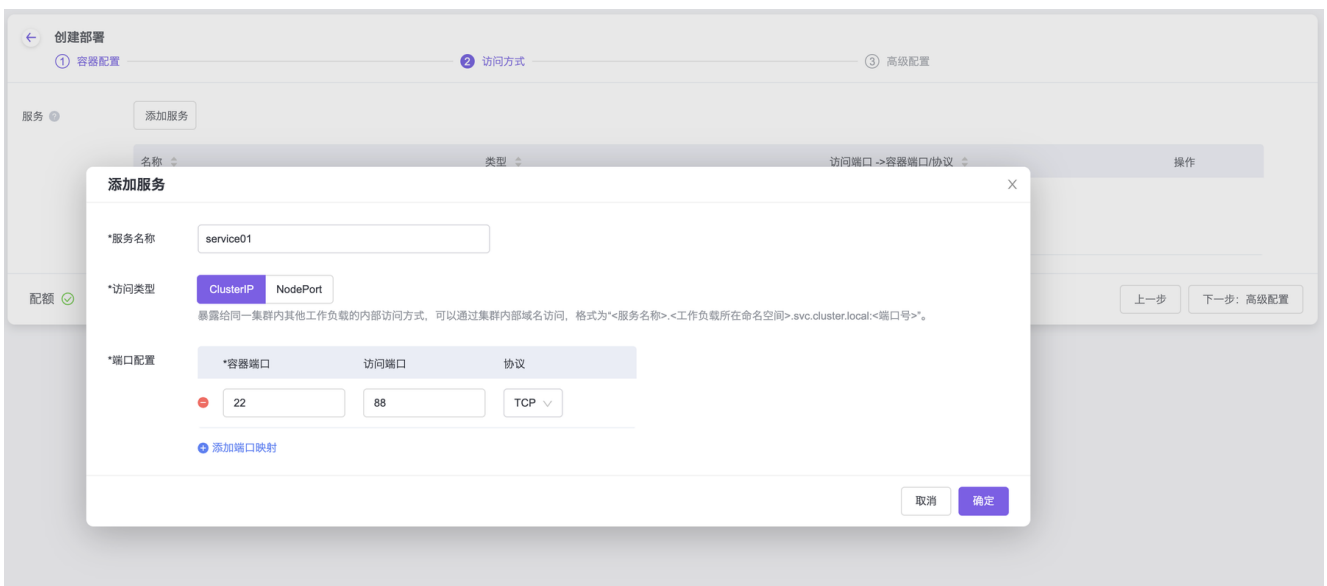
参数		说明
	数据卷（可选）	<p>单击 <code>添加数据卷</code>，弹出“添加数据卷”对话框。配置参数，单击 <code>确定</code> 完成操作。参数说明如下：</p> <p>* 类型：</p> <ul style="list-style-type: none"> - 持久卷声明：仅工作负载类型为部署、任务、定时任务时可选择本类型。给容器挂载持久化存储，数据不会因容器的销毁或节点异常而消失。适用于需持久化存储、高磁盘IO等场景。持久卷声明需要事先创建，相关介绍请参考 创建持久卷声明。 - 存储类：仅工作负载类型为有状态副本集时可选择。不需事先创建持久卷声明，可直接通过指定存储类及所需存储容量创建持久卷，并挂载到指定的容器路径。各参数的具体说明，请参考 创建持久卷声明。 - 临时路径：将容器所在宿主机的临时目录挂载到容器的指定路径。 - 配置：选择已创建好的配置。 - 密钥：选择已创建好的密钥。 <p>* 挂载路径(可选)：所选数据卷挂载至容器的绝对路径。</p>
	健康检查（可选）	<p>健康检查包括存活检查、就绪检查和启动检查功能。存活检查用于检测容器是否正常，如果容器的存活检查失败，集群会对该容器执行重启操作；若容器的存活检查成功则不执行任何操作。就绪检查用于检查用户业务是否就绪，如果容器的就绪检查失败，则不转发流量到当前容器组；若检查成功，则会开放对该容器组的访问。启动检查用于保护慢启动容器有充足时间完成启动，避免死锁状况发生。</p> <p>* 检查方式：</p> <ul style="list-style-type: none"> - HTTP/HTTPS方式：适用于提供HTTP/HTTPS服务的容器，集群周期性地对该容器发起HTTP/HTTPS GET请求，如果HTTP/HTTPS 返回状态码小于400，则证明检查成功、容器健康，否则检查失败。例如，方式选择HTTP，路径为/check，端口为80，则集群周期性向容器发起如下请求：<code>GET http://容器IP:80/check</code>。 - TCP方式：适用于提供TCP通信服务的容器，集群周期性地检测端口是否为打开状态，若端口为打开状态，则检查成

参数		说明
		<p>功、容器健康；若端口为关闭或进程为停止状态，则检查失败。例如：一个提供nginx服务的容器，服务端口为80，则配置TCP检查端口为80，那么集群会周期性检测该容器的80端口打开状态。</p> <ul style="list-style-type: none"> - 容器命令方式：该方式要求用户指定一个容器内的可执行命令，集群会周期性地在容器内执行该命令，若进程退出状态码为 0则检查成功、容器健康，否则检查失败。 <p>* 公共参数：</p> <ul style="list-style-type: none"> - 首次检查延时：容器启动后第一次进行健康检查的延迟时间，这段时间为预留给业务程序正常启动。例如，设置为10，表明容器启动后10秒才开始健康检查。 - 检查间隔：执行健康检查的时间间隔。例如，设置为30，则每间隔30秒执行一次健康检查。 - 超时时间：检查超时后的等待时间。例如，设置为10，表明执行健康检查的超时等待时间为10秒，如果超过这个时间，本次健康检查就被视为失败。 - 健康认定（）次成功：假设本参数设置为N，健康检查失败后，至少连续成功N次会认为容器健康。 - 不健康认定（）次失败：假设本参数设置为X，健康检查失败后，集群将继续尝试X次健康检查，若仍不符合健康条件，则放弃该容器。对于存活检查，放弃意味着重启容器；对于就绪检查，放弃意味着容器组将被标记为未就绪。
	安全设置（可选）	<ul style="list-style-type: none"> * 非root用户运行：要求容器组具有非零runAsUser值，或在镜像中定义了USER环境变量。 * 只读root文件系统：是否必须使用一个只读的root文件系统。 * runAsUser：用户ID。容器中的进程都以该用户ID运行。 * runAsGroup：Group ID。容器中的进程都以该Group ID运行。

参数		说明
	命令（可选）	<ul style="list-style-type: none"> * 启动命令：容器启动时运行的第一条命令，将覆盖镜像中的Entrypoint指令。 * 启动命令参数：覆盖镜像中的CMD执行，如已设置了运行命令，该条指令将被附加到运行命令的参数中。 * 启动后执行命令：该命令在创建容器之后立即执行。 * 停止前执行命令：这个命令在停止容器前执行，是否立即调用此命令取决于 API 的请求或者管理事件。
	日志采集（可选）	<p>采集应用的运行日志，实现平台内应用与组件日志的全量采集与查询。</p> <ul style="list-style-type: none"> * 日志源：选择需采集日志的资源，支持容器标准日志（默认）和应用日志两种。 * 日志文件路径：当“日志源”选择“应用日志”时需配置。日志文件存放路径，需注意，该路径要求已挂载数据卷。

4. 在“访问方式”页面中，配置参数后，单击 **下一步：高级配置**，进入“高级配置”配置页面。

当该工作负载需要提供对外访问的服务时，请单击 **添加服务**，在弹出的对话框中配置参数后，单击 **确定**，完成服务添加。否则，可直接跳过本步骤。



参数		说明
服务（可选）	服务名称	该工作负载中用于提供外部访问的服务的名称。
	访问类型	<ul style="list-style-type: none">* ClusterIP：适用于集群内部访问场景，集群为服务分配一个固定的集群内虚拟IP，集群内其它pod可以通过集群内部域名访问，格式为“<服务名称>.<工作负载所在命名空间>.svc.cluster.local:<端口号>”。集群外无效。* NodePort：适用于集群外部访问场景，集群除了会给服务分配一个内部的虚拟IP，还会在每个节点上为服务分配静态端口号，集群外部可通过集群任一节点IP和静态端口号访问服务。
	端口配置	<ul style="list-style-type: none">* 容器端口：容器镜像中工作负载实际监听的端口。* 访问端口：容器端口映射到节点IP上的端口。当访问方式为“NodePort”时，支持随机生成。* 协议：包括TCP、UDP，根据业务类型选择。

5. 在“高级配置”页面中，配置参数后，单击 **确定** ，完成操作。

← 创建部署
① 容器配置
② 访问方式
③ 高级配置

升级策略 展开 ▾

伸缩策略 开启

最小实例数

最大实例数

CPU使用率阈值 % 开启

内存使用率阈值 % 开启

调度策略

主机调度

Pod亲和性

Pod反亲和性

污点容忍

网络设置 收起 ▲

主机别名 + 添加主机别名

*VPC ↻

*子网

出口IP

标签 + 添加标签

配额 ✔

上一步
确定

参数	说明
<p>升级策略</p> <p style="text-align: center; margin-top: 20px;">工作负载类型为“部署”</p>	<p>* 先启动新Pod，再停止旧Pod/先停止旧Pod，再启动新Pod：可定义每次启动或停止Pod的数量。例如选择先启动新Pod，再停止旧Pod，批量大小设置为1，则每次先启动1个新的Pod，新的Pod成功后停止1个旧Pod，以此类推。</p> <p>* 停止所有Pod，再启动新Pod：先停止所有老版本容器组，再启动新版本容器组，升级过程中业务会中断。</p> <p>* 自定义：“最大超量”表示更新过程中容器组数量可以超过期望副本的数量或百分比。“最多不可用数”表示升级过程中允许的最多不可用容器组数量，如果等于期望副本数量有业务中断风险（最小存活容器组数量=期望副本数量-最多不可用数）。</p>

参数		说明
	工作负载类型为“有状态副本集”或“守护进程集”	* 滚动：滚动升级将逐步用新版本的实例替换旧版本的实例，升级的过程中，业务流量会同时负载均衡分布到新老的实例上，因此业务不会中断。其中，“最多不可用数”表示升级过程中允许的最多不可用容器组数量，如果等于期望副本数量则有业务中断风险（最小存活容器组数量=期望副本数量-最多不可用数）。 * 手动删除时更新：集群不会自动更新工作负载中的容器组，需手动删除容器组以使集群创建新的容器组。
伸缩策略		仅当工作负载类型为“部署”时可配置。当达到设置的条件后自动扩展或收缩容器组数量。 * 最小实例数：期望容器组数量的最小值。 * 最大实例数：期望容器组数量的最大值。 * CPU使用率阈值：所有容器组平均cpu使用率超过阈值自动扩展，n-1（n为容器组总数）个容器组平均内存使用率低于阈值自动收缩。需勾选“开启”后才能输入阈值。 * 内存使用率阈值：所有容器组平均内存使用率超过阈值自动扩展，n-1（n为容器组总数）个容器组平均内存使用率低于阈值自动收缩。需勾选“开启”后才能输入阈值。
调度策略（可选）	主机调度	* 指定主机：可选择集群内任一节点，该工作负载内的容器将被调度到所选节点上。 * 自定义规则：包括必须满足条件和尽量满足条件。必须满足条件是硬性要求，必须满足才能成功调度，支持添加多条规则，多条规则间是“且”的关系，即需要满足所有规则才可以调度；尽量满足条件表示集群会尽量将容器调度到符合规则的主机上，支持添加多条规则，多条规则间是“或”的关系，不满足规则的主机也会进行调度，根据规则的权重值，权重值越高越会被优先调度。

参数		说明
	Pod亲和性/Pod反亲和性	Pod亲和性决定哪些工作负载的Pod部署在同一个拓扑域，可根据业务需求进行工作负载的就近部署，容器间通信就近路由，减少网络消耗。Pod反亲和性决定工作负载的Pod不和哪些工作负载的Pod部署在同一个拓扑域，互相干扰的工作负载反亲和部署，避免干扰，减少宕机影响。拓扑域是由一个或多个节点组成的，这些节点在所指定的属性上具有相同的值，例如拓扑域为kubernetes.io/hostname，则具有相同hostname的节点成为一个拓扑域（即同一节点）。必须满足条件是硬性要求，支持添加多条规则，多条规则间是“且”的关系，即需要满足所有规则才可以调度；尽量满足条件表示集群会尽量将容器调度到符合规则的主机上，多条规则间是“或”的关系，不满足规则的主机也会进行调度，根据规则的权重值，权重值越高越会被优先调度。
	污点容忍	调度工作负载时能够容忍具有指定污点的节点。支持添加多条污点规则，多条规则间是“或”的关系，即满足任一规则即可调度。
网络设置		<ul style="list-style-type: none"> * 主机别名：添加主机别名后即可通过域名访问对应IP地址的主机。 * VPC：该工作负载的私有网络名称。支持选择系统默认生成的ovn-cluster，也支持自定义VPC和各网络的子网，具体操作请参考 如何自定义VPC网络。 * 子网：该工作负载的子网名称。 * 出口IP：该工作负载提供对外访问服务时，所使用的公网IP地址。
标签（可选）		通过标签可以方便地标识及筛选对象。

3.5 创建Ingress（可选）

通过Ingress可以为工作负载的服务提供外部访问时所需的路由规则集合，请根据客户实际业务需求酌情创建。当工作负载已添加服务，且该服务需要配置对外访问的路由规则时，才需执行此操作。

1. 在左侧导航栏选择[业务视图]，选择目标命名空间，选择[网络管理]-[Ingress]，进入“Ingress”页面。
2. 单击 **创建Ingress** ，进入“创建Ingress”页面。
3. 配置参数后，单击 **创建** 完成操作。

参数	说明
Ingress规则	是一种HTTP方式的路由转发机制。例如域名填写为example.com，路径填写为/path，服务选择已创建的名称为“app”的服务，则外部可通过 <code>http://example.com/path</code> 访问名称为“app”的服务。
注解	Ingress经常使用注解（annotations）来配置一些选项，具体取决于Ingress控制器。

3.6 使用Yaml创建资源（可选）

本章节主要描述如何使用Yaml创建安全容器服务资源。工作负载、持久卷声明和配置、密钥、服务、Ingress、自定义资源描述等资源不仅可以通过云平台页面进行创建，还可以通过Yaml文件进行创建。用户可以根据实际的业务场景，灵活选择适合的资源创建方案。

1. 在顶部导航栏选择[产品与服务]-[安全容器服务]任意一个子菜单，进入“安全容器服务”页面。
2. 在左侧导航栏选择[业务视图]页签-选择目标命名空间。
3. 单击页面右下角的“Yaml”图标，进入“导入Yaml”页面。



4. 直接粘贴Yaml文件内容，或单击编辑区域右上角的“导入”图标，选择本地存储的Yaml文件。

说明：

请关注调试结果。该调试主要针对格式校验，若有错误可点击错误信息，跳至目标行进行修改。

为了更清晰地说明，下面以创建名为“demo”的“部署”类型的工作负载为例，介绍如何使用Yaml创建资源。以下是一个示例Yaml内容：

```
apiVersion: apps/v1
kind: Deployment
metadata:
  creationTimestamp: null
  labels:
    app: demo
  name: demo
spec:
  replicas: 1
  selector:
    matchLabels:
      app: demo
  strategy: {}
  template:
    metadata:
      creationTimestamp: null
      labels:
        app: demo
    spec:
      containers:
      - image: nginx:latest
        name: nginx
        resources: {}
status: {}
```

5. 待调试通过后，单击 **导入** ，完成操作。

4 用户指南

4.1 集群管理

本章节主要介绍在“集群管理”页面中，针对集群的运维管理操作。“集群管理”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[集群管理]，进入“集群管理”页面。

查看集群详情

可以查看集群内资源使用情况及当前运行情况和集群事件等。

1. 进入“集群管理”页面。
2. 在集群列表中单击目标集群名称链接，进入集群详情页面。
3. 选择[概览]页签，查看集群概览信息；选择[集群事件]页签，查看集群中对象的错误或警示信息。

4.2 节点管理

本章节主要介绍在“节点管理”页面中，针对节点的运维管理操作。“节点管理”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[节点管理]，进入“节点管理”页面。

说明：

云平台控制平面的节点，不支持开始调度、停止调度、标签管理和污点管理操作。

查看节点详情

1. 进入“节点管理”页面。
2. 单击节点名称链接，进入节点详情页面，查看详细信息。

开始调度/停止调度

开始调度后，新创建的容器组可以调度到节点上，停止调度后不可以。

1. 进入“节点管理”页面。
2. 单击目标节点操作栏的 **开始调度** 或 **停止调度**，弹出“开始调度”或“停止调度”提示框。
3. 单击 **确定** 完成操作。

标签管理

本功能用于增加或删除节点标签，系统标签不支持编辑和删除。

1. 进入“节点管理”页面。
2. 单击目标节点操作栏的 **更多** - **标签管理**，弹出“标签管理”对话框。
3. 添加或移除标签。
4. 单击 **确定** 完成操作。

污点管理

节点设置上污点之后就与容器组之间存在着相斥的关系，可以让节点拒绝容器组的调度，甚至将节点上已经存在的容器组驱逐出去。例如，当已知某个节点资源不足且无法为其扩容时，可以给节点打上污点标记，使容器组不再调度到该节点或者驱逐节点上的容器组，隔离该节点。

1. 进入“节点管理”页面。
2. 单击目标节点操作栏的 **更多** - **污点管理** ，弹出“污点管理”对话框。
3. 添加或删除污点。
4. 单击 **确定** 完成操作。

参数		说明
调度策略	不允许调度 (NoSchedule)	新创建的容器组不会调度到该节点。
	尽量不调度 (PreferNoSchedule)	新创建的容器组尽量不调度到该节点。
	不允许并驱逐已有容器组 (NoExecute)	新创建的容器组不会调度到该节点，已经运行在节点上的容器组也会被驱逐。

4.3 命名空间

本章节主要介绍在“命名空间”页面中，针对命名空间的运维管理操作。“命名空间”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[命名空间]，进入“命名空间”页面。

创建命名空间

1. 进入“命名空间”页面。
2. 单击 **创建命名空间** ，进入“创建命名空间”页面。
3. 配置参数，单击 **创建命名空间** 完成操作。

参数	说明
名称	选择命名空间的名称。
集群	选择命名空间所属集群。
部门/项目	用户所在部门和项目，不支持修改。

删除命名空间

警告：

删除命名空间会同时删除该命名空间下的所有资源。

1. 进入“命名空间”页面。
2. 选择目标命名空间，单击 **删除** ，弹出“删除命名空间”提示框。
3. 单击 **删除** 完成操作。

4.4 存储管理

存储类

存储类可以实现动态供应持久卷，即能够按照用户的需要，自动创建其所需的存储。本章节主要介绍在“存储类”页面中，针对存储的运维管理操作。“存储类”页面的进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[存储管理]，进入“存储管理”页面。
2. 在左侧导航栏选择[存储管理]-[存储类]，进入“存储类”页面。

查看Yaml

1. 进入“存储类”页面。
2. 单击目标存储类操作栏的 **查看Yaml** ，弹出“查看Yaml”提示框。
3. 查看信息后，单击 **关闭** 完成操作。

持久卷

持久卷描述的是持久化存储卷，主要定义的是一个持久化存储在宿主机上的目录，独立于容器组生命周期。具体到本平台，一个持久卷对应一个云硬盘。本章节主要介绍在“持久卷”页面中，针对存储的运维管理操作。“持久卷”页面的进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[存储管理]，进入“存储管理”页面。
2. 在左侧导航栏选择[存储管理]-[持久卷]，进入“持久卷”页面。

查看Yaml

1. 进入“持久卷”页面。
2. 单击目标持久卷操作栏的 **查看Yaml** ，弹出“查看Yaml”提示框。
3. 查看信息后，单击 **关闭** 完成操作。

删除持久卷

说明：

已被绑定至持久卷声明的持久卷无法删除。

1. 进入“持久卷”页面。
2. 单击目标持久卷操作栏的 **删除** ，弹出“删除持久卷”提示框。
3. 单击 **删除** 完成操作。

批量删除持久卷

说明：

已被绑定至持久卷声明的持久卷无法删除。

1. 进入“持久卷”页面。
2. 勾选目标持久卷后，单击 **删除** ，弹出“删除持久卷”提示框。
3. 单击 **删除** 完成操作。

4.5 业务概览

本功能用于查看各命名空间下的业务运行状况，如容器组状态、容器配额使用情况等。

1. 在顶部导航栏单击[产品与服务]-[安全容器服务]进入“安全容器服务”页面。
2. 在左侧导航栏选择[业务视图]页签，进入业务视图页面。
3. 在左侧导航栏选择目标命名空间，切换至目标命名空间视图。
4. 在左侧导航栏选择[概览]，进入概览页面即可查看信息。

4.6 工作负载

本章节主要介绍在相应类型工作负载页面中，针对工作负载的运维管理操作。相应类型工作负载页面的进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[工作负载]，进入“工作负载”页面。
2. 在左侧导航栏选择目标命名空间后，选择对应子菜单，进入对应页面。

创建部署/有状态副本集/守护进程集

部署即kubernetes中的Deployment控制器，一个“部署”可以包含一个或多个容器组副本，这些容器组是无状态的（即完全相同、相互独立、可被替换），系统会自动为Deployment的多个Pod副本分发请求。通过定义期望的副本数、容器属性等，“部署”会保证实际状态与所需状态一致，即使发生意外情况也可以将容器组恢复到期望状态。通过“部署”可以实现上线部署、滚动升级（不停止旧服务的状态下升级）、回滚应用（将应用回滚到之前的版本）、平滑扩缩容功能。

有状态副本集即kubernetes中的StatefulSet控制器，一个“有状态副本集”可以包含一个或多个容器组副本，这些容器组是有状态的（运行过程中会保存数据或状态），支持有序部署和删除，支持持久化存储，适用于容器组间存在主从关系、主备关系、互相访问等关系的场景。

守护进程集即kubernetes中的DaemonSet控制器。守护进程集确保全部（或者某些）节点都运行一个容器组，支持实例动态添加到新节点，适用于实例在每个节点上都需要运行的场景，例如在每个节点上运行日志收集程序、节点监视程序等。

1. 进入相应类型工作负载的页面。
2. 单击 **创建部署/有状态副本集/守护进程集** ，进入对应页面。
3. 配置参数，参数说明请参考 [创建工作负载](#) 。
4. 单击 **确认** 完成操作。

创建任务

任务会创建一个或者多个容器组，并将持续重试容器组的执行，直到指定数量的容器组成功终止。随着容器组成功结束，任务跟踪记录成功完成的容器组个数。当数量达到指定的成功个数阈值时，任务（即 Job）结束。

1. 进入“任务”页面。

- 单击 **创建任务** ，进入“创建任务”的“基础配置”页面。
- 填写基础配置参数。

参数	说明
目标完成次数	当成功完成的容器组达到该值时认为任务完成。
并行实例数	每次创建的容器组数量。
失败重试次数	失败容器组的最大重试次数，超过这个次数不会继续重试。
超时时间	任务运行的超时时间。如果任务运行的时间超过了设定的时间，此任务将自动停止运行所有容器组。
重启策略	容器组内容器的重启策略，包括“不重启”和“失败时重启”。
调度策略	容器组内容器的调度策略。即调度工作负载时，是否能够容忍具有污点的节点。

- 单击 **下一步：容器配置** ，进入“创建任务”的“容器配置”页面。
- 填写容器配置参数，参数说明请参考 [创建工作负载](#)。
- 单击 **创建** 完成操作。

创建定时任务

定时任务即Kubernetes中的CronJob，是基于时间的“任务”，在指定的时间周期运行指定的“任务”。

- 进入“定时任务”页面。
- 单击 **创建定时任务** ，进入“创建定时任务”的“基础配置”页面。
- 填写基础配置参数。

参数	说明
定时规则	指定任务运行周期。

参数	说明
并发策略	* Forbid: 在前一个任务未完成时, 不创建新任务。 * Allow: 当到达新任务创建时间点, 而前一个任务未完成时, 新的任务会取代前一个任务。 * Replace: 定时任务不断创建新的任务, 会抢占集群资源。
目标完成次数	当成功完成的容器组达到该值时认为任务完成。
并行实例数	每次创建的容器组数量。
失败重试次数	失败容器组的最大重试次数, 超过这个次数不会继续重试。
超时时间	任务运行的超时时间。如果任务运行的时间超过了设定的时间, 此任务将自动停止运行所有容器组。
重启策略	容器组内容器的重启策略, 包括“不重启”和“失败时重启”。
调度策略	容器组内容器的调度策略。即调度工作负载时, 是否能够容忍具有污点的节点。

- 单击 **下一步: 容器配置**, 进入“创建定时任务”的“容器配置”页面。
- 填写容器配置参数, 参数说明请参考 [创建工作负载](#)。
- 单击 **创建** 完成操作。

管理工作负载

说明:

各类型工作负载支持的操作不尽相同, 请根据实际页面显示和业务需求酌情配置。

查看工作负载详情

- 进入相应类型工作负载的页面。
- 找到目标工作负载, 单击工作负载名称链接, 进入工作负载详情页。
- 查看工作负载详细信息。

容器配置

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **容器配置**，进入“容器配置”页面。
3. 配置参数，参数说明请参考 [创建工作负载](#)。
4. 单击 **确认** 完成操作。

手动伸缩

说明：

- 处于“已停止”状态的的工作负载不支持手动伸缩。
- 针对部署类型的工作负载，若设置了弹性伸缩策略，则不支持进行手动伸缩。

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **手动伸缩**，弹出“手动伸缩”对话框。
3. 默认展示当前工作负载副本数量，可手动修改。此数量为目标值而非差值。
4. 单击 **确认** 完成操作。

版本回滚

1. 进入“部署”页面。
2. 找到目标工作负载，单击操作栏的 **更多 - 版本回滚**，弹出“版本回滚”对话框。
3. 选择需要回滚到的历史版本。
4. 单击 **确认** 完成操作。

升级策略

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多 - 升级策略**，弹出“升级策略”对话框。
3. 配置参数，参数说明请参考 [创建工作负载](#)。
4. 单击 **确认** 完成操作。

伸缩策略

1. 进入“部署”页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **伸缩策略** ，弹出“伸缩策略”对话框。
3. 配置参数，参数说明请参考 [创建工作负载](#) 。
4. 单击 **确认** 完成操作。

调度策略

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **调度策略** ，弹出“调度策略”对话框。
3. 配置参数，参数说明请参考 [创建工作负载](#) 。
4. 单击 **确认** 完成操作。

网络设置

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **网络设置** ，弹出“网络设置”对话框。
3. 配置参数，参数说明请参考 [创建工作负载](#) 。
4. 单击 **确认** 完成操作。

标签设置

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **标签设置** ，弹出“标签设置”对话框。
3. 增加或移除标签。
4. 单击 **确认** 完成操作。

编辑Yaml

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **编辑Yaml** 或直接单击操作栏的 **编辑Yaml** ，弹出“编辑Yaml”对话框。
3. 修改信息。
4. 单击 **确认** 完成操作。

启动

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **启动** ，弹出对应提示框。
3. 单击 **启动** 完成操作。

停止

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **停止** ，弹出对应提示框。
3. 单击 **停止** 完成操作。

重新部署

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **重新部署** ，弹出对应提示框。
3. 单击 **重新部署** 完成操作。

删除

1. 进入相应类型工作负载的页面。
2. 找到目标工作负载，单击操作栏的 **更多** - **删除** 或直接单击操作栏的 **删除** ，弹出对应提示框。
3. 单击 **删除** 完成操作。

运行/停止定时任务

1. 进入“定时任务”页面。
2. 找到目标工作负载，单击操作栏的 **运行** 或 **停止** ，弹出对应提示框。
3. 单击 **运行** 或 **停止** 完成操作。

管理容器组

查看容器组详情

支持查看容器组基本信息、容器配置、状态、事件、监控、日志和终端。

1. 进入“容器组”页面。
2. 单击容器组名称链接，进入容器组详情页面，查看信息。

查看Yaml

1. 进入“容器组”页面。
2. 单击目标容器组操作栏的 **查看Yaml** ，查看信息。

查看日志

1. 进入“容器组”页面。
2. 单击目标容器组操作栏的 **日志** ，查看信息。

终端

1. 进入“容器组”页面。
2. 单击目标容器组操作栏的 **更多** - **终端** ，进入终端页面。

删除

1. 进入“容器组”页面。
2. 单击目标容器组操作栏的 **更多** - **删除** ，弹出“删除容器组”对话框。
3. 根据需要确认是否勾选“强制删除”。例如，目标容器组因所在节点已经停止或者无法连接API Server等异常情况无法被正常删除，此时可进行强制删除。
4. 单击 **删除** 完成操作。

4.7 持久卷声明

本章节主要介绍在“持久卷声明”页面中，针对持久卷声明的运维管理操作。“持久卷声明”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[持久卷声明]，进入“持久卷声明”页面。
2. 在左侧导航栏选择目标命名空间。

创建持久卷声明

容器可通过持久卷声明请求使用持久化存储。

1. 进入“持久卷声明”页面。
2. 单击 **创建持久卷声明**，弹出“创建持久卷声明”对话框。
3. 配置参数。
4. 单击 **创建** 完成操作。

参数	说明
存储类	即[管理视图]-[存储管理]中管理的存储类，详细介绍请参考 存储管理-存储类 。
大小	所需存储卷的容量。
访问模式	包括三种模式，需根据存储类的能力选择其支持的模式： * 单节点读写（RWO）：卷可以被一个节点以读写方式挂载。 * 多节点读写（RWX）：卷可以被多个节点以读写方式挂载。 * 多节点只读（ROX）：卷可以被多个节点以只读方式挂载。若“部署”类型的工作负载需挂载单节点读写（RWO）模式的卷，其副本数需为1；若“任务”、“定时任务”类型的工作负载需挂载单节点读写（RWO）模式的卷，其并行实例数需为1。

编辑Yaml

1. 进入“持久卷声明”页面。
2. 单击目标持久卷声明操作栏的 **编辑Yaml**，弹出“编辑Yaml”对话框。
3. 修改信息。

4. 单击 **确认** 完成操作。

删除

说明：

已关联容器组的持久卷声明不支持删除。

1. 进入“持久卷声明”页面。
2. 单击目标持久卷声明操作栏的 **删除** ，弹出“删除持久卷声明”提示框。
3. 单击 **删除** 完成操作。

4.8 配置中心

配置

配置用于保存配置数据，可以用作工作负载的环境变量、命令行参数或者存储卷中的配置文件。使用配置实现容器化应用的配置管理，可以使配置与镜像内容分离，保持容器化应用的可移植性。本章节主要介绍在“配置”页面中，针对“配置”的运维管理操作。“配置”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[配置中心]，进入“配置中心”页面。
2. 在左侧导航栏选择目标命名空间，选择[配置中心]-[配置]，进入“配置”页面。

创建配置

1. 进入“配置”页面。
2. 单击 **创建配置** ，进入“创建配置”页面。
3. 填写配置名称和配置项内容。
4. 单击 **创建配置** 完成操作。

编辑Yaml

1. 进入“配置”页面。
2. 单击目标配置操作栏的 **编辑Yaml** ，弹出“编辑Yaml”对话框。
3. 修改信息。
4. 单击 **确认** 完成操作。

更新

1. 进入“配置”页面。
2. 单击目标配置操作栏的 **更新** ，进入“更新”页面。
3. 填写配置名称和配置项内容。
4. 单击 **更新** 完成操作。

删除

1. 进入“配置”页面。
2. 单击目标配置操作栏的 **删除** ，弹出“删除配置”提示框。
3. 单击 **删除** 完成操作。

密钥

密钥（Secret）是一种包含认证信息、密钥等敏感信息的资源类型，可以用作工作负载的环境变量、加密配置文件。将数据放在密钥对象中，可以更好地控制它的用途，并降低意外暴露的风险。本章节主要介绍在“密钥”页面中，针对“密钥”的运维管理操作。“密钥”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[配置中心]，进入“配置中心”页面。
2. 在左侧导航栏选择目标命名空间，选择[配置中心]-[密钥]，进入“密钥”页面。

创建密钥

1. 进入“密钥”页面。
2. 单击 **创建密钥** ，进入“创建密钥”页面。
3. 配置参数。
4. 单击 **创建** 完成操作。

参数	说明
密钥类型	<ul style="list-style-type: none">* Opaque：一般密钥类型。* TLS：存放7层负载均衡服务所需的证书。* 镜像访问密钥：存放拉取私有仓库镜像所需的认证信息。
密钥数据	<ul style="list-style-type: none">* 当密钥类型为Opaque时，单击“添加密钥数据”，输入键、值。* 当密钥类型为TLS时，上传证书和私钥文件。* 当密钥类型为镜像访问密钥时，输入镜像仓库地址、用户名、密码和邮箱。

编辑Yaml

1. 进入“密钥”页面。
2. 单击目标密钥操作栏的 **编辑Yaml** ，弹出“编辑Yaml”对话框。

3. 修改信息。
4. 单击 **确认** 完成操作。

更新

1. 进入“密钥”页面。
2. 单击目标密钥操作栏的 **更新** ，进入“更新”页面。
3. 修改信息。
4. 单击 **更新** 完成操作。

删除

1. 进入“密钥”页面。
2. 单击目标密钥操作栏的 **删除** ，弹出“删除密钥”提示框。
3. 单击 **删除** 完成操作。

4.9 网络管理

服务

服务（Service）是容器服务的基本操作单元，是将请求进行负载分发到后端的各个容器应用上的控制器。对外表现为一个单一访问接口，外部不需要了解后端如何运行，这给扩展或维护后端带来很大的好处。本章节主要介绍在“服务”页面中，针对“服务”的运维管理操作。“服务”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[网络管理]，进入“网络管理”页面。
2. 在左侧导航栏选择目标命名空间，选择[网络管理]-[服务]，进入“服务”页面。

创建服务

1. 进入“服务”页面。
2. 单击 **创建服务** ，进入“创建服务”页面。
3. 配置参数。
4. 单击 **创建** 完成操作。

参数		说明
类型	ClusterIP	适用于集群内部访问场景，集群为服务分配一个固定的集群内虚拟IP，集群内其它pod可以通过集群内部域名访问，格式为“<服务名称>.<工作负载所在命名空间>.svc.cluster.local:<端口号>”。集群外无效。
	NodePort	适用于集群外部访问场景，集群除了会给服务分配一个内部的虚拟IP，还会在每个节点上为服务分配静态端口号，集群外部可通过集群任一节点IP和静态端口号访问服务。
	ExternalName	用于将服务请求指向一个自定义的域名。
容器端口	容器镜像中工作负载实际监听的端口。	
访问端口	容器端口映射到节点IP上的端口。当访问方式为“NodePort”时，支持随机生成。	
协议	包括TCP、UDP，根据业务类型选择。	

参数	说明
关联工作负载	选择服务需关联的工作负载。当服务类型为“ExternalName”无此参数。

编辑Yaml

1. 进入“服务”页面。
2. 单击目标服务操作栏的编辑Yaml，弹出“编辑Yaml”对话框。
3. 修改信息。
4. 单击 **确认** 完成操作。

更新

1. 进入“服务”页面。
2. 单击目标服务操作栏的 **更新**，进入“更新”页面。
3. 修改参数。
4. 单击 **保存** 完成操作。

删除

1. 进入“服务”页面。
2. 单击目标服务操作栏的 **删除**，弹出“删除服务”提示框。
3. 单击 **删除** 完成操作。

Ingresses

Ingress是一组将集群内服务暴露给集群外服务的路由规则集合。一个Ingress对象能够配置具备为服务提供外部可访问的URL、负载均衡流量、卸载SSL/TLS，以及提供基于名称的虚拟主机等能力。本章节主要介绍在“Ingresses”页面中，针对“Ingress”的运维管理操作。“Ingresses”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[网络管理]，进入“网络管理”页面。
2. 在左侧导航栏选择目标命名空间，选择[网络管理]-[Ingresses]，进入“Ingresses”页面。

创建Ingress

1. 进入“Ingresses”页面。
2. 单击 **创建Ingress** ，进入“创建Ingress”页面。
3. 配置参数。
4. 单击 **创建** 完成操作。

参数	说明
Ingress规则	是一种HTTP方式的路由转发机制。例如域名填写为example.com，路径填写为/path，服务选择已创建的名称为“app”的服务，则外部可通过 <code>http://example.com/path</code> 访问名称为“app”的服务。
注解	Ingress经常使用注解（annotations）来配置一些选项，具体取决于Ingress控制器。

编辑Yaml

1. 进入“Ingresses”页面。
2. 单击目标Ingress操作栏的 **编辑Yaml** ，弹出“编辑Yaml”对话框。
3. 修改信息。
4. 单击 **确认** 完成操作。

更新

1. 进入“Ingresses”页面。
2. 单击目标Ingress操作栏的 **更新** ，进入“更新”页面。
3. 修改参数。
4. 单击 **保存** 完成操作。

删除

1. 进入“Ingresses”页面。
2. 单击目标Ingress操作栏的 **删除** ，弹出“删除Ingress”提示框。
3. 单击 **删除** 完成操作。

4.10 自定义资源管理

本章节主要介绍在“自定义资源管理”页面中，针对自定义资源的运维管理操作。“自定义资源管理”页面进入路径如下：

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[安全容器服务]-[自定义资源CRD管理]或[自定义资源CRD]，进入“自定义资源管理”页面。
2. 在左侧导航栏选择目标命名空间。

导入自定义资源描述/自定义资源

1. 进入“自定义资源管理”页面。
2. 单击页面右下角的“Yaml”图标，进入“导入Yaml”页面。
3. 直接粘贴Yaml文件内容，或单击编辑区域右上角的“导入”图标，选择本地存储的Yaml文件。

说明：

- 请关注调试结果。该调试主要针对格式校验，若有错误可点击错误信息，跳至目标行进行修改。
- 针对自定义资源描述，只能云管理员在[管理视图]中导入。

4. 待调试通过后，单击 **导入** ，完成操作。

查看自定义资源描述详情

1. 进入“自定义资源管理”页面。
2. 在自定义资源描述列表中单击目标资源描述的名称链接，进入资源描述详情页面。
3. 选择[Yaml]页签，查看其Yaml信息；选择[自定义资源]页签，查看自定义资源的信息。

删除自定义资源

1. 进入“自定义资源管理”页面。
2. 在自定义资源描述列表中单击目标资源描述的名称链接，进入资源描述详情页面。
3. 选择[自定义资源]页签后，勾选目标资源，单击 **删除** ，弹出“删除自定义资源”提示框。

4. 单击 **删除** 完成操作。

5 最佳实践

5.1 配置SR-IOV网络

5.1.1 前置条件准备

为了使用SR-IOV功能，需要满足以下前提条件：

1. 服务器BIOS中已开启SR-IOV功能。
2. 启用Intel IOMMU功能，并以Pass-Through模式进行配置。
3. 已安装NVIDIA MLNX_OFED驱动。
4. 网卡固件启用SR-IOV功能。
5. 在MLNX_OFED驱动上启用SR-IOV。

这些前提条件的满足将确保系统具备使用SR-IOV功能所需的硬件支持和驱动程序。请确保在使用SR-IOV功能之前，在执行后续的配置操作之前，先完成以下准备工作：

操作步骤

1.服务器BIOS中已开启SR-IOV功能

请确保服务器BIOS中已开启SR-IOV功能。每个服务器都有不同的虚拟化BIOS配置选项。BIOS configuration examples可参考：[BIOS configuration examples](#)

2.服务器BIOS中已开启SR-IOV功能

请确保将"intel_iommu=on"和"iommu=pt"添加到grub文件的配置中。

```
# cat /boot/grub2/grub.cfg

# GRUB Environment Block
saved_entry=0
kernelopts=root=/dev/mapper/os-root ro edd=off kvm.halt_poll_ns=400000
cgroup.memory=nokmem intel_iommu=on iommu=pt pci=realloc
```

```
ixgbe.allow_unsupported_sfp=1 rootdelay=90 nomodeset intel_idle.max_cstate=0
processor.max_cstate=0 crashkernel=300M rd.lvm.lv=os/root biosdevname=0
net.ifnames=1
boot_success=0
```

要了解更多关于iommu grub参数的信息，请参阅：[Understanding the iommu Linux grub File Configuration](#)

3.安装NVIDIA MLNX_OFED驱动

3.1 下载驱动包

NVIDIA驱动官方下载：[Linux InfiniBand Drivers](#) 请根据您所选择的操作系统和版本，以及CPU架构，下载相应的tar包。例如，选择MLNX_OFED_LINUX-5.9-0.5.6.0-rhel8.4-x86_64.tgz进行下载。

MLNX_OFED Download Center

Version (Current)	OS Distribution	OS Distribution Version	Architecture	Download/ Documentation
5.9-0.5.6.0.107-for DGX H100 Systems Only	Ubuntu	RHEL/Rocky 9.1	x86_64	ISO: MLNX_OFED_LINUX-5.9-0.5.6.0-rhel8.4-x86_64.iso
	UOS	RHEL/Rocky 9.0	ppc64le	SHA256: ac7b98491b803530fb1696b52afaa618aa8963562f398a39f472a6
	SLES	RHEL/Rocky 8.7	aarch64	Size: 263M
	RHEL/CentOS/Rocky	RHEL/Rocky 8.6		tgz: MLNX_OFED_LINUX-5.9-0.5.6.0-rhel8.4-x86_64.tgz
5.9-0.5.6.0	Oracle Linux	RHEL/CentOS/Rocky 8.5		SHA256: a7034c7fc978ffb36eee7213d0e8aca5119f62b8830860a35dcd18
5.8-2.0.3.0-LTS	OPENEULER	RHEL/CentOS 8.4		Size: 260M
5.4-3.6.8.1-LTS	KYLIN	RHEL/CentOS 8.3		SOURCES: MLNX_OFED_SRC-5.9-0.5.6.0.tgz
4.9-6.0.6.0-LTS	EulerOS	RHEL/CentOS 8.2		
	Debian	RHEL/CentOS 8.1		
	Community	RHEL/CentOS 8.0		
	Citrix XenServer	RHEL/CentOS 7.0		

3.2 执行安装操作

```
# ./mlnxofedinstall
```

3.3 查看驱动版本信息

```
# ofed_info -s
MLNX_OFED_LINUX-5.9-0.5.6.0:
```

3.4 网卡状态

3.4.1 首次安装网卡时，如果未连接到交换机，网卡的状态将如下所示：

```
[root@nodeb ~]# ibdev2netdev
mlx5_0 port 1 ==> ib0 (Down)
[root@nodeb ~]#
[root@nodeb ~]# ibstat
CA 'mlx5_0'
    CA type: MT4119
    Number of ports: 1
    Firmware version: 16.28.1002
    Hardware version: 0
    Node GUID: 0x0c42a10300b643c6
    System image GUID: 0x0c42a10300b643c6
    Port 1:
        State: Down
        Physical state: Polling
        Rate: 10
        Base lid: 65535
        LMC: 0
        SM lid: 0
        Capability mask: 0x2651e848
        Port GUID: 0x0c42a10300b643c6
        Link layer: InfiniBandSinoinfo
```

（如上图，ib0状态Down（双口卡还有ib1），网卡的SM lid为0，Base lid为65535，Link layer为IB模式）

3.4.2 使用线缆将网卡连接到交换机，大概30秒，端口UP，如下所示：

```
[root@nodeb ~]# ibdev2netdev
mlx5_0 port 1 ==> ib0 (Up)
[root@nodeb ~]#
[root@nodeb ~]# ibstat
CA 'mlx5_0'
    CA type: MT4119
    Number of ports: 1
    Firmware version: 16.28.1002
    Hardware version: 0
    Node GUID: 0x0c42a10300b643c6
    System image GUID: 0x0c42a10300b643c6
    Port 1:
        State: Active
        Physical state: LinkUp
        Rate: 56
        Base lid: 3
        LMC: 0
        SM lid: 2
        Capability mask: 0x2651e848
        Port GUID: 0x0c42a10300b643c6
        Link layer: InfiniBandSinoinfo
```

(如上图，网卡ib0物理状态LinkUp, 逻辑状态Active，子网管理器SM lid为 2，网卡Base lid为 3)

说明：

本案例交换机SB7800已经开启子网管理器功能，交换机端口速率为EDR 100Gb/s，网卡CX5也是 EDR 100Gb/s，由于线缆是FDR 56G，所以上图的Rate: 56，并没有达到100

3.4.3 启动opensmd服务（视交换机型号选装）

如果交换机是SB7890或者是其它不带子网管理器功能的交换机（SB7890，SB7790，SX6025，IS5025），必须在服务器上开启子网管理器，安装好网卡驱动后就可以使用下面命令：

```
# systemctl start opensmd
# systemctl enable opensmd
```

3.4.4 网卡运行环境监测

完成以上步骤，使用命令'mlnx_tune'检查，所有状态都正常，则网卡运行环境正常。

```
# mlnx_tune
```

4.网卡固件启用SR-IOV功能

4.1 下载Mellanox Management Tools (MFT)工具

下载链接：[Mellanox Management Tools \(MFT\)](#)

4.2 使用Mellanox Firmware Tools包在固件中启用和配置SR-IOV

```
# mst start
Starting MST (Mellanox Software Tools) driver set
Loading MST PCI module - Success
Loading MST PCI configuration module - Success
Create devices
```

4.2.1 在需要的PCI插槽上找到Connect-IB设备

```
# mst status
MST modules:
-----
    MST PCI module loaded
    MST PCI configuration module loaded
MST devices:
-----
/dev/mst/mt4115_pciconf0      - PCI configuration cycles access.
...
```

4.2.2 查询设备状态

```
# mlxconfig -d /dev/mst/mt4115_pciconf0 q
```

4.2.3 开启SR-IOV, 设置可以切分的vf数量上限

```
# mlxconfig -d /dev/mst/mt4115_pciconf0 set SRIOV_EN=1 NUM_OF_VFS=8
...
Apply new Configuration? ? (y/n) [n] : y
Applying... Done!
-I- Please reboot machine to load new configurations.
```

4.2.4 设置网卡的工作模式。

Infiniband卡支持两种工作模式：IB模式和Ethernet模式

```
# mlxconfig -d /dev/mst/mt4115_pciconf0 q

Configurations:                                Next Boot
...
LINK_TYPE_P1                                IB(1)    #当前工作模式 IB
...
```

修改网卡的工作模式：

```
Ethernet模式: mlxconfig -d /dev/mst/mt4115_pciconf0 set LINK_TYPE_P1=2
IB模式: mlxconfig -d /dev/mst/mt4115_pciconf0 set LINK_TYPE_P1=1
```

重启机器生效

```
# reboot
```

此时，通过lspci无法看到vf。只有当SR-IOV在MLNX_OFED驱动上启用时，您才能看到它们。

4.2.5 查看节点物理网卡是否支持SR-IOV

```
查看所有的pci设备
# lspci -nnn |grep -i mellanox
```

```
61:00.0 Infiniband controller [0207]: Mellanox Technologies MT27700 Family
[ConnectX-4] [15b3:1013]
```

查看设备详情

```
# lspci -vvs
```

```
61:00.0 Infiniband controller: Mellanox Technologies MT27700 Family
[ConnectX-4]
```

```
...
```

```
Capabilities: [180 v1] Single Root I/O Virtualization (SR-IOV)
IOVCap: Migration-, Interrupt Message Number: 000
IOVctl: Enable- Migration- Interrupt- MSE- ARIHierarchy+
IOVSta: Migration-
Initial VFs: 3, Total VFs: 3, Number of VFs: 0, Function Dependency
```

```
Link: 00
```

```
VF offset: 1, stride: 1, Device ID: 1014
Supported Page Size: 000007ff, System Page Size: 00000001
Region 0: Memory at 000005464e000000 (64-bit, prefetchable)
VF Migration: offset: 00000000, BIR: 0
```

```
Capabilities: [1c0 v1] Secondary PCI Express
LnkCtl3: LnkEquIntrruptEn- PerformEqu-
LaneErrStat: 0
```

```
Kernel driver in use: mlx5_core
Kernel modules: mlx5_core
```

显示如下代码段，说明固件已启动SR-IOV功能

 [/doc/SecureContainerService/6.2.1/zh-cn/images/scs_gs_image_13.png](#)

5.在MLNX_OFED驱动上启用SR-IOV

5.1 定位设备（通常为 'mlx5_0'）

```
# ibstat
CA 'mlx5_0'
  CA type: MT4115
  Number of ports: 1
  Firmware version: 12.28.2006
  Hardware version: 0
  Node GUID: 0xe41d2d03006768fa
  System image GUID: 0xe41d2d03006768fa
  Port 1:
```

```
State: Down
Physical state: Disabled
Rate: 10
Base lid: 65535
LMC: 0
SM lid: 0
Capability mask: 0x2650e848
Port GUID: 0xe41d2d03006768fa
Link layer: InfiniBand
```

5.2 获取此设备上的当前vf数

```
# cat /sys/class/infiniband/mlx5_0/device/mlx5_num_vfs
0
```

如果命令失败，这可能意味着驱动程序没有加载。

5.3 设置所需的vf数量

```
# echo 2 > /sys/class/infiniband/mlx5_0/device/mlx5_num_vfs
```

更改mlx5_num_vfs不是持久的，并且不能在服务器重新启动后继续存在。

5.4 检查PCI总线

```
# lspci -nnn |grep -i mellanox
61:00.0 Infiniband controller [0207]: Mellanox Technologies MT27700 Family
[ConnectX-4] [15b3:1013]
61:00.1 Infiniband controller [0207]: Mellanox Technologies MT27700 Family
[ConnectX-4 Virtual Function] [15b3:1014]
61:00.2 Infiniband controller [0207]: Mellanox Technologies MT27700 Family
[ConnectX-4 Virtual Function] [15b3:1014]
```

5.5 恢复vf数量为0

```
# echo 0 > /sys/class/infiniband/mlx5_0/device/mlx5_num_vfs
```


5.1.2 使用Yaml配置SR-IOV网络

背景描述

SR-IOV (Single Root I/O Virtualization) 是一种允许物理设备如网络接口卡(NIC)在没有软件的情况下将其资源分割成多个虚拟设备的技术。使用SR-IOV可以在性能和资源隔离方面提供显著的优势，因为它减少了网络数据在主机和虚拟机之间传输的开销。安全容器服务支持SR-IOV网络设备插件，用于发现、公告和分配SR-IOV网络虚拟功能 (VF) 资源。本章节通过使用Yaml创建功能，配置SR-IOV网络。

前置条件

为了使用SR-IOV功能，需要满足以下前提条件：

1. 服务器BIOS中已开启SR-IOV功能。
2. 启用Intel IOMMU功能，并以Pass-Through模式进行配置。
3. 已安装NVIDIA MLNX_OFED驱动。
4. 网卡固件启用SR-IOV功能。
5. 在MLNX_OFED驱动上启用SR-IOV。详细操作步骤请参见[前置条件](#)。这些前提条件的满足将确保系统具备使用SR-IOV功能所需的硬件支持和驱动程序。请确保在使用SR-IOV之前，你的系统已按照要求进行相应的设置和安装。

操作步骤

前置操作完成后，需重启环境中sriov-config-daemon的Pod

因为 `sriov-config-daemon` 是 `daemonSet` 部署，则删除 `namespace` 为 `eks-managed`，`label` 为 `app=sriov-network-config-daemon` 的Pod即可。

创建SriovNetworkNodePolicy对象：

SriovNetworkNodePolicy是SR-IOV Network Operator的一部分，用于定义如何配置SR-IOV网络。它是一个Kubernetes自定义资源(CR)，用于指定SR-IOV网络配置策略。可以通过定义SriovNetworkNodePolicy对象来指定节点的SR-IOV网络设备配置：

1.查看节点SR-IOV设备信息 节点sriov-config-daemon上报上来的SR-IOV设备情况会更新到节点对应的

`sriovnetworknodestates.status` , 即查看 `sriovnetworknodestates.status` 内容

```
# kubectl get sriovnetworknodestates.sriovnetwork.openshift.io -n eks-managed -o yaml
```

```
status:
  interfaces:
  - deviceID: "1013"
    driver: mlx5_core
    linkType: IB
    mac: 00:00:07:ff:fe:80:00:00:00:00:00:00:e4:1d:2d:03:00:67:68:fa
    mtu: 2048
    name: ib0
    pciAddress: 0000:61:00.0
    totalvfs: 3
    vendor: 15b3
    syncStatus: Succeeded
```

2.根据节点SR-IOV设备情况, 创建SriovNetworkNodePolicy资源。

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovNetworkNodePolicy
metadata:
  name: 1
  namespace: eks-managed
spec:
  resourceName: 2
  nodeSelector:
    : 3
  priority: 3
  mtu: 4
  numVfs: 5
  nicSelector: 6
  vendor: "" 7
  deviceID: "" 8
  pfNames: ["", "..."] 9
  rootDevices: ["", "..."] 10
```

```
deviceType: vfio-pci
isRdma: false 11
```

参数	说明
name	CR对象的名称。
resourceName	SR-IOV设备插件的资源名称。您可以为一个资源名称创建多个"SriovNetworkNode Policy"对象。
priority (可选)	"0"到"99"之间的整数。较小的数值具有较高的优先权，优先级"10"高于优先级"99"。默认值为"99"。
mtu (可选)	为虚拟功能 (VF) 的最大传输单位 (MTU) 指定一个值。MTU的值必须是在"1"到"9 000"之间。如果您不需要指定MTU，请指定一个"值。
numVfs	为SR-IOV物理网络设备指定要创建的虚拟功能 (VF) 的数量。对于Intel网络接口卡 (NIC)，VF的数量不能超过该设备支持的VF总数。对于Mellanox NIC，VF的数量不能超过"128"。
nicSelector	nicSelector 映射为Operator选择要配置的设备。您不需要为所有参数指定值。建议您以足够的准确度来识别网络设备，以便尽量减小意外选择其他网络设备的可能性。如果指定了rootDevices，则必须同时为vendor、deviceID或pfNames指定一个值。如果同时指定了"pfNames"和"rootDevices"，请确保它们指向同一个设备。
vendor (可选)	SR-IOV网络设备的厂商十六进制代码。允许的值只能是"8086"和"15b3"。
deviceID (可选)	SR-IOV网络设备的设备十六进制代码。允许的值只能是"158b"、"1015"和"1017"。
pfNames (可选)	该设备的一个或多个物理功能 (PF) 名称的数组。
rootDevices	用于该设备的PF的一个或多个PCI总线地址的数组。使用以下格式提供地址: "000 0:02:00.1"。
isRdma (可选)	是否启用远程直接访问 (RDMA) 模式。默认值为"false"。

说明：

写 `SriovNetworkPolicy` 的 `spec` 的定义，要根据 `sriovNetworkNodeState.Status.interface` 去写。若 `SriovNetworkPolicy` 的 `spec` 定义，不在 `sriovNetworkNodeState.Status.interface` 的设备范围内，则生成只有 `DpConfigVersion` 的 `spec`。

创建 SriovIBNetwork 资源：

您可以通过定义对象来配置 InfiniBand (IB) 网络设备 "SriovIBNetwork"。

1. 下方示例Yaml展示了如何通过 SriovIBNetwork 对象，配置 Infiniband network attachment

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovIBNetwork
metadata:
  name: 1
  namespace: eks-managed
spec:
  resourceName: 2
  networkNamespace: 3
  vlan: 4
  spoofChk: "" 5
  ipam: |- 6
    {}
  linkState: 7
  maxTxRate: 8
  minTxRate: 9
  vlanQoS: 10
  trust: "" 11
  capabilities: 12
```

参数	说明
name	对象的名称。SR-IOV Network Operator创建一个名称相同的"NetworkAttachment Definition"对象。

参数	说明
resourceName	定义附加网络的SR-IOV硬件"SriovNetworkNodePolicy"对象中的"spec.resourceName"参数的值。
networkNames pace	SriovIBNetwork对象的目标命名空间。只有目标命名空间中的Pod才能连接到网络设备。
vlan (可选)	额外网络的虚拟LAN (VLAN) ID。它需要是一个从"0"到"4095"范围内的一个整数值。默认值为"0"。
spoofChk (可选)	VF的"spoof"检查模式。允许的值是字符串"on"和"off"。重要: 指定的值必须由引号包括, 否则SR-IOV Network Operator将拒绝对象。
ipam	为IPAM CNI插件指定一个配置对象做为一个Yaml块scalar。该插件管理网络附加定义的IP地址分配。
linkState (可选)	虚拟功能 (VF) 的链接状态。允许的值是"enable"、"disable"和"auto"。
maxTxRate (可选)	VF的最大传输率 (以Mbps为单位)。
minTxRate (可选)	VF的最低传输率 (以Mbps为单位)。这个值必须小于或等于最大传输率。说明: Intel NIC不支持minTxRate参数。如需更多信息, 请参阅 BZ#1772847
vlanQoS (可选)	VF的"IEEE 802.1p"优先级级别。默认值为"0"。
trust (可选)	VF的信任模式。允许的值是字符串"on"和"off"。说明:您必须在引号中包含指定的值, 或者SR-IOV Network Operator拒绝对象。
capabilities (可选)	为这个额外网络配置功能。您可以指定 "{ "ips": true }" 来启用IP地址支持, 或指定 "{ "mac": true }" 来启用MAC地址支持。

2.使用Whereabouts进行动态IP地址分配配置

```
{
  "ipam": {
    "type": "whereabouts",
```

```
"range": "", 1
"exclude": ["", "..."], 2
}
}
```

参数	说明
range	以CIDR标记指定一个IP地址范围。IP地址是通过这个地址范围来分配的。
exclude (可选)	在CIDR标记中指定IP地址和范围列表。包含在排除地址范围中的IP地址。

创建业务容器 (Pod) 资源:

指定业务容器的规范, 包括容器运行时 (如runc、rune)、容器镜像、资源需求 (如CPU和内存)、环境变量、命令等。

```
apiVersion: v1
kind: Pod
metadata:
  name: sriov-rune-pod
  annotations:
    io.katacontainers.config.runtime.enable_sriov: "true"
    k8s.v1.cni.cncf.io/networks: eks-managed/nics
    #v1.multus-cni.io/default-network: eks-managed/kube-ovn
spec:
  runtimeClassName: rune
  containers:
  - name: appcntr1
    image: hub.easystack.io/captain/nginx-ingress-controller:v0.49.3
    imagePullPolicy: IfNotPresent
    command: [ "/bin/bash", "-c", "--" ]
    args: [ "while true; do sleep 300000; done;" ]
    resources:
      requests:
        : '1'
      limits:
        : '1'
```

注意：业务Pod的resources.requests、resources.limits的资源名称必须与sriovNetworkNodePolicy.spec.resourceName、sriovNetwork.spec.resourceName相等。

5.1.3 runc容器和InfiniBand卡IB模式场景

基于runc运行时容器和InfiniBand卡（IB模式）的组合场景，主要原理是将容器运行时环境与高性能的InfiniBand网络卡相结合。通过利用runc容器的轻量级和可移植性优势，与InfiniBand卡（IB模式）相结合，适用于安全隔离的需求较低，对网络性能要求较高的应用场景。本文将通过Yaml配置信息和参数，演示如何定义SR-IOV网络节点的策略。

操作步骤

配置SriovNetworkNodePolicy对象

指定切分 `kubernetes.io/hostname=node-10` 节点上, `rootDevices: 0000:71:00.0` 的PF设备

警告:

创建SR-IOV SriovNetworkNodePolicy对象时，节点应用修改会重启。

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovNetworkNodePolicy
metadata:
  name: node-policy-10
  namespace: eks-managed
spec:
  resourceName: mlxnic
  nodeSelector:
    kubernetes.io/hostname: node-10
  nicSelector:
    vendor: "15b3"
    deviceID: "1017"
    rootDevices:
      - 0000:71:00.0
  deviceType: netdevice
  numVfs: 3
  priority: 50
  isRdma: true
  linkType: IB
```

配置SriovIBNetwork对象

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovIBNetwork
metadata:
  name: ibnics
  namespace: eks-managed
spec:
  ipam: |-
    {
      "type": "whereabouts",
      "range": "192.168.100.0/24",
      "gateway": "192.168.100.1",
      "exclude": [
        "192.168.100.0/26"
      ]
    }
  resourceName: mlxnic
  linkState: auto
```

配置runc运行时环境中的业务Pod

```
apiVersion: v1
kind: Pod
metadata:
  name: sriov-runc-pod-demo
  annotations:
    k8s.v1.cni.cncf.io/networks: eks-managed/ibnics
spec:
  containers:
  - name: app-demo
    image: hub.ecns.io/test/nginx:latest
    imagePullPolicy: Always
    command: [ "/bin/bash", "-c", "--" ]
    args: [ "while true; do sleep 300000; done;" ]
    resources:
      requests:
        ecnf.io/mlxnic: "1"
      limits:
```

```
ecnf.io/mlxnic: "1"  
nodeName: node-10
```

5.1.4 runc容器和InfiniBand卡ETH模式场景

基于runc运行时容器和InfiniBand卡（ETH模式）的组合场景，主要原理是将容器运行时环境与高性能的InfiniBand网络卡相结合。通过利用runc容器的轻量级和可移植性优势，与InfiniBand卡（ETH模式）相结合，适用于安全隔离的需求较低，需要高性能数据传输和低延迟的应用场景。本文将通过Yaml配置信息和参数，演示如何定义SR-IOV网络节点的策略。

操作步骤

配置SriovNetworkNodePolicy对象

指定切分 `kubernetes.io/hostname=node-10` 节点上, `rootDevices: 0000:71:00.0` 的PF设备.

警告:

创建SR-IOV SriovNetworkNodePolicy对象时，节点应用修改会重启。

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovNetworkNodePolicy
metadata:
  name: node-policy-10
  namespace: eks-managed
spec:
  resourceName: mlxnic
  nodeSelector:
    kubernetes.io/hostname: node-10
  nicSelector:
    vendor: "15b3"
    deviceID: "1017"
    rootDevices:
      - 0000:71:00.0
  deviceType: netdevice
  numVfs: 3
  priority: 50
  isRdma: false
  linkType: ETH
```

配置SriovIBNetwork对象

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovNetwork
metadata:
  name: nics
  namespace: eks-managed
spec:
  ipam: |-
    {
      "type": "whereabouts",
      "range": "192.168.100.0/24",
      "gateway": "192.168.100.1",
      "exclude": [
        "192.168.100.0/26"
      ]
    }
  resourceName: mlxnic
```

配置runc运行时环境中的业务Pod

```
apiVersion: v1
kind: Pod
metadata:
  name: sriov-runc-pod-demo
  annotations:
    k8s.v1.cni.cncf.io/networks: eks-managed/nics
spec:
  containers:
    - name: app-demo
      image: hub.ecns.io/test/nginx:latest
      imagePullPolicy: Always
      command: [ "/bin/bash", "-c", "--" ]
      args: [ "while true; do sleep 300000; done;" ]
      resources:
        requests:
          ecnf.io/mlxnic: "1"
        limits:
```

```
ecnf.io/mlxnic: "1"  
nodeName: node-10
```

5.1.5 rune容器和InfiniBand卡IB模式场景

基于rune（安全运行时）容器和InfiniBand卡（IB模式）的组合场景，主要原理是将容器运行时环境与高性能的InfiniBand网络卡相结合。通过利用rune容器的安全性和隔离性优势，与InfiniBand卡（IB模式）相结合，适用于对安全性和隔离性有一定需求的轻量级传输应用场景。本文将通过Yaml配置信息和参数，演示如何定义SR-IOV网络节点的策略。

操作步骤

配置SriovNetworkNodePolicy对象：

指定切分 `kubernetes.io/hostname=node-10` 节点上，`rootDevices: 0000:71:00.0` 的PF设备

警告：

创建SR-IOV SriovNetworkNodePolicy对象时，节点应用修改会重启。

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovNetworkNodePolicy
metadata:
  name: node-policy-5
  namespace: eks-managed
spec:
  resourceName: mlxnic
  nodeSelector:
    kubernetes.io/hostname: node-5
  nicSelector:
    vendor: "15b3"
    deviceID: "1017"
    rootDevices:
      - 0000:71:00.0
  deviceType: vfio-pci
  numVfs: 3
  priority: 50
  isRdma: false
  linkType: IB
```

配置SriovIBNetwork对象：

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovIBNetwork
metadata:
  name: ibnics
  namespace: eks-managed
spec:
  ipam: |-
    {
      "type": "whereabouts",
      "range": "192.168.100.0/24",
      "gateway": "192.168.100.1",
      "exclude": [
        "192.168.100.0/26"
      ]
    }
  resourceName: mlxnic
  linkState: auto
```

配置 rune(安全运行时) 环境中的业务 Pod：

如果需要对容器（Pod）进行资源限制（limit）的设置，您可以在Pod的request字段中设置limit值。为了实现这个需求，您需要为Pod添加以下的annotation配置：

```
io.katacontainers.config.runtime.sandbox_cgroup_only: "false"
```

```
apiVersion: v1
kind: Pod
metadata:
  name: sriov-rune-pod-demo
  annotations:
    k8s.v1.cni.cncf.io/networks: eks-managed/ibnics
    io.katacontainers.config.runtime.enable_sriov: "true"
spec:
  runtimeClassName: rune
  containers:
  - name: app-demo
    image: hub.ecns.io/test/nginx:latest
    imagePullPolicy: Always
```



```
command: [ "/bin/bash", "-c", "--" ]
args: [ "while true; do sleep 300000; done;" ]
resources:
  requests:
    ecnf.io/mlxnic: "1"
  limits:
    ecnf.io/mlxnic: "1"
nodeName: node-10
```

5.1.6 rune容器和InfiniBand卡ETH模式场景

基于rune（安全运行时）容器和InfiniBand卡（ETH模式）的组合场景，主要原理是将容器运行时环境与高性能的InfiniBand网络卡相结合。通过利用rune容器的安全性和隔离性优势，与InfiniBand卡（ETH模式）相结合，适用于对安全性和隔离性有一定需求的，需要快速数据传输和低延迟的应用场景。本文将通过Yaml配置信息和参数，演示如何定义SR-IOV网络节点的策略。

操作步骤

配置SriovNetworkNodePolicy对象：

指定切分 `kubernetes.io/hostname=node-10` 节点上， `rootDevices: 0000:71:00.0` 的PF设备

警告：

创建SR-IOV SriovNetworkNodePolicy对象时，节点应用修改会重启。

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovNetworkNodePolicy
metadata:
  name: node-policy-10
  namespace: eks-managed
spec:
  resourceName: mlxnic
  nodeSelector:
    kubernetes.io/hostname: node-10
  nicSelector:
    vendor: "15b3"
    deviceID: "1017"
    rootDevices:
      - 0000:71:00.0
  deviceType: vfio-pci
  numVfs: 3
  priority: 50
  isRdma: false
  linkType: ETH
```

配置SriovIBNetwork对象:

```
apiVersion: sriovnetwork.openshift.io/v1
kind: SriovNetwork
metadata:
  name: nics
  namespace: eks-managed
spec:
  ipam: |-
    {
      "type": "whereabouts",
      "range": "192.168.100.0/24",
      "gateway": "192.168.100.1",
      "exclude": [
        "192.168.100.0/26"
      ]
    }
  resourceName: mlxnic
```

配置 rune(安全运行时) 环境中的业务 Pod:

如果需要对容器（Pod）进行资源限制（limit）的设置，您可以在Pod的request字段中设置limit值。为了实现这个需求，您需要为Pod添加以下的annotation配置：

```
io.katacontainers.config.runtime.sandbox_cgroup_only: "false"
```

```
apiVersion: v1
kind: Pod
metadata:
  name: sriov-rune-pod-demo
  annotations:
    k8s.v1.cni.cncf.io/networks: eks-managed/nics
    io.katacontainers.config.runtime.enable_sriov: "true"
spec:
  runtimeClassName: rune
  containers:
    - name: app-demo
      image: hub.ecns.io/test/nginx:latest
      imagePullPolicy: Always
      command: [ "/bin/bash", "-c", "--" ]
```

```
args: [ "while true; do sleep 300000; done;" ]
resources:
  requests:
    ecnf.io/mlxnic: "1"
  limits:
    ecnf.io/mlxnic: "1"
nodeName: node-10
```

5.2 创建GPU资源的容器实例

背景描述

GPU在工作负载中的主要作用是利用其出色的并行计算能力。在安全容器服务中，将GPU卡用于运行计算密集型工作负载对于某些特定的业务场景非常有价值，包括高性能计算、数据处理和分析加速、图形渲染和可视化优化等。通过将GPU与容器技术结合，可以实现灵活、高效且可扩展的计算环境，提供更好的用户体验和性能表现。本章节将详细介绍如何在部署工作负载时使用GPU卡能力。

前提条件

- **上传对应的GPU解决方案对接包：**通过上传解决方案对接包，系统能够正确识别和使用GPU资源，以便在相应的工作负载中进行分配和调度。请确保在需要使用GPU资源时，已上传并配置了对应的GPU解决方案对接包。
- **容器运行时：**当容器运行时为"安全运行时 (rune)"时，支持使用GPU资源。（注意："守护进程集 (DaemonSet)"和"定时任务 (CronJob)"类型的工作负载不支持使用GPU。）
- **资源预留：**当需要勾选"使用GPU"时，建议设置CPU的参数值大于等于1，内存的参数值大于等于1024MiB。这样可以确保在使用GPU时，为工作负载分配足够的计算资源和内存资源，以获得良好的性能和稳定性。

操作步骤

目前平台中支持在安全容器中使用英伟达GPU设备和百度昆仑XPU设备。具体支持型号参见[使用限制](#)。下文将分别列出不同品牌GPU的操作步骤和示例，以帮助用户正确配置和使用GPU资源。

英伟达 (NVIDIA) GPU设备

在支持使用GPU能力的工作负载类型创建界面中，勾选"使用GPU"选项，并选择"nvidia.com/gpu"作为GPU资源的类型。这样配置后，系统将为工作负载分配相应的GPU资源，以提供所需的计算能力和性能。创建容器组配置时，配置好容器 (Pod) GPU参数，所有增量容器 (Pod) 默认共享GPU卡资源。

← 创建部署
① 容器配置
② 访问方式
③ 高级配置

*容器运行时 安全运行时 runc运行时

*安全负载名称

*副本数

容器配置 container1 | x 添加安全容器

*容器名称

容器类型 业务容器 初始化容器

镜像来源 镜像仓库 第三方镜像

*镜像 选择镜像

*镜像版本

拉取镜像策略 本地不存在时拉取 总是拉取

*资源预留 CPU 内存 Mi

*资源限制 CPU 内存 Mi

GPU 使用GPU
开启后，所有增量安全容器默认开启 GPU 能力，共享 GPU 卡资源配置。

*型号

*数量 张

环境变量 + 添加环境变量

如果需要对容器（Pod）进行资源限制（limit）的设置，您可以在Pod的request字段中设置limit值。为了实现这个需求，您需要为Pod添加以下的annotation配置：

```
io.katacontainers.config.runtime.sandbox_cgroup_only: "false"
```

Yaml示例：

```
template:
  metadata:
    labels:
      my-app: perf-server0
    annotations:
      # 在limit中配置了memory时，需要添加如下注解
```

```
io.katacontainers.config.runtime.sandbox_cgroup_only: "false"
spec:
  runtimeClassName: rune
  containers:
  - name: perf-server0
    image: docker.io/test/perf:0.0.1
    command:
    - iperf3
    args:
    - -s
    resources:
      limits:
        cpu: "2"
        memory: 2Gi
        # 配置GPU资源数量
        nvidia.com/gpu: "1"
      requests:
        cpu: "2"
        memory: 2Gi
        nvidia.com/gpu: "1"
```

百度昆仑XPU

在支持使用GPU能力的工作负载类型创建界面中，勾选“使用GPU”选项，并选择“baidu.com/XPU”作为GPU资源的类型。这样配置后，系统将为工作负载分配相应的GPU资源，以提供所需的计算能力和性能。创建容器组配置时，配置好容器（Pod）GPU参数后，所有增量容器（Pod）默认共享GPU卡资源。

← 创建部署
① 容器配置
② 访问方式
③ 高级配置

*容器运行时 安全运行时 runc运行时

*安全负载名称

*副本数

容器配置 container1 | x 添加安全容器

*容器名称

容器类型 业务容器 初始化容器

镜像来源 镜像仓库 第三方镜像

*镜像 选择镜像

*镜像版本

拉取镜像策略 本地不存在时拉取 总是拉取

*资源预留 CPU 内存 Mi

*资源限制 CPU 内存 Mi

GPU 使用GPU
开启后，所有增量安全容器默认开启 GPU 能力，共享 GPU 卡资源配置。

*型号

*数量 张

环境变量 + 添加环境变量

如果需要对容器（Pod）进行资源限制（limit）的设置，您可以在容器（Pod）的“request”字段中设置“limit”值。为了实现这个需求，您需要为容器（Pod）添加以下的“annotation”配置：

```
io.katacontainers.config.runtime.sandbox_cgroup_only: "false"
```

Yaml示例：

```
template:
  metadata:
    labels:
      my-app: perf-server0
    annotations:
      # 在limit中配置了memory时，需要添加如下注解
```



```
    io.katacontainers.config.runtime.sandbox_cgroup_only: "false"
spec:
  runtimeClassName: rune
  containers:
  - name: perf-server0
    image: docker.io/test/perf:0.0.1
    command:
    - iperf3
    args:
    - -s
    resources:
      limits:
        cpu: "2"
        memory: 2Gi
        # 配置XPU资源数量
        baidu.com/XPU: "1"
      requests:
        cpu: "2"
        memory: 2Gi
        baidu.com/XPU: "1"
```

此外，百度昆仑XPU支持指定某个容器对象独享某个具体的GPU卡（其他容器不可见），可通过环境变量：

`CXPU_VISIBLE_DEVICES: 1` 设置GPU卡隔离。当环境可用GPU卡未被指定完，配置后且未被隔离的GPU卡为所有容器共享资源。当环境可用GPU卡被指定完，剩余的容器无可用的GPU卡资源。

百度昆仑XPU在各个容器（Pod）中的可见性环境变量：`CXPU_VISIBLE_DEVICES` 可以设置值为 `all`或者是卡序号（1, 2, 3）

Yaml示例：

```
template:
  metadata:
    labels:
      my-app: perf-server0
    annotations:
      # 在limit中配置了memory时，需要添加如下注解
      io.katacontainers.config.runtime.sandbox_cgroup_only: "false"
  spec:
    runtimeClassName: rune
    containers:
```

```
- name: perf-server0
  image: docker.io/test/perf:0.0.1
  command:
  - iperf3
  args:
  - -s
  env:
  - name: CXPU_VISIBLE_DEVICES
    value: ALL
  resources:
    limits:
      cpu: "2"
      memory: 2Gi
      # 配置GPU资源数量
      nvidia.com/gpu: "1"
    requests:
      cpu: "2"
      memory: 2Gi
      nvidia.com/gpu: "1"
```

附录

为了正确使用GPU资源，请确保准确指定GPU卡的型号对应的资源名称：

GPU卡型号	资源名称
英伟达 (NVIDIA) GPU设备	nvidia.com/gpu
百度昆仑XPU设备	baidu.com/XPU

6 常见问题

6.1 如何理解安全容器网络方案

问题描述

安全容器使用的CNI插件是什么？有什么重要特性？用户侧如何使用VPC、Subnet、FIP等各种网络资源？

解决方案

安全容器6.1.1版本CNI是基于易捷行云软SDN网络服务作为后端，以kube-ovn（v1.9.0版本）作为API接口，为不同算力(容器、云原生云主机、裸金属)提供统一网络方案；为实现这个设计目标，对kube-ovn 作为容器CNI 做了一些改动和限制；和社区版本的主要区别体现在：

- 1、当前版本跨节点通信 不支持 underlay 模式，只支持geneve 隧道模式
- 2、VPC 功能增强 1) 容器网络使用的VPC 依赖于软SDN的 router，所以默认VPC的更新配置，及 自定义VPC的创建 需要结合软SDN router页面操作完成，具体步骤见操作手册。
2) 容器网络的VPC 网关出外网配置由软SDN router 自动完成，对应社区kube-ovn文档中 VPC 网关配置相关的操作，这部分无需用户再做配置。同时VPC只支持集中式网关(对应 VPC spec中 gatewaytype 字段)。
3) 支持自定义VPC，且自定义VPC内容器网络也支持 nodeport、探针等功能(当前版本限制：不同vpc间 subnet cidr 不能有 overlap)，同时支持自定义VPC的网关节点的配置。
4) 容器VPC与其他算力的互联需要在 软SDN router页面完成配置。
- 3、新增Floating IP（简称fip）功能 1) fip 作为软SDN管理的统一资源，可以同时被 不同算力(vm、pod、裸金属) 使用； fip 在pods中的具体使用方式请参照用户手册。
2) 产品中使用内部CRD fips.neutron.io 管理记录当前容器可用的fip 资源池，用户可以k8s API通过查看对应的 fip CR查看可用的 fip address 使用，fip的记录与回收会自动处理。
- 4、subnet功能增强 支持subnet 中定义projectIDs，表示该subnet只能被哪些projects使用；如果不配置表示该subnet为share subnet，所有用户可用。

如何启用kube-ovn 组件

通过安全容器服务页面创建的命名空间，默认所有负载都是使用kube-ovn作为CNI。

API方式，需要对命名空间打上标签 `managed.es.io/resource=namespace` ；例如：

```
apiVersion: v1
kind: Namespace
metadata:
  name: easystack
  managed.es.io/resource: namespace
```

如何自定义VPC网络

在创建工作负载时，其网络不仅支持使用系统默认生成的ovn-cluster和子网，还支持使用自定义的VPC网络或子网。用户可根据实际业务场景，酌情选择对应方案。

1. 配置外部网络。

1. （可选）创建外部网络。

本操作用于预创建外部网络，以便在自定义VPC网络时能够为其建立外部连接。如使用已有外部网络时，可跳过本步骤。

1. 在云平台的顶部导航栏中，依次选择[产品与服务]-[网络]-[网络]，进入“网络”页面。

2. 单击 `创建网络` ，进入“创建网络”页面。

3. “网络类型”请选择“外部网络”，并配置其他参数后，单击 `创建网络` ，完成操作。其中，其他参数的具体参数说明，请参考“SDN网络服务”帮助中“网络”的相关内容。

2. 查看外部网络的ID。

在“网络”页面中，单击待操作网络名称，进入其详情页面。在该页面中，查看并记录该网络的ID（即UUID参数的值）。

2. 配置路由器。

1. （可选）创建路由器。

本操作用于预创建路由器，以便在自定义VPC网络时能够为其建立外部连接。如使用已有路由器时，可跳过本步骤。

1. 在云平台的顶部导航栏中，依次选择[产品与服务]-[网络]-[路由器]，即可进入“路由器”页面。
2. 单击 **创建路由器** ，弹出“创建路由器”对话框。
3. 配置参数后，单击 **创建** ，完成操作。

2. （可选）设置路由器网关。

本操作用于预设置路由器网关，以便在自定义VPC网络时能够为其建立外部连接。如路由器已设置网关时，可跳过本步骤。

1. 在“路由器”页面中，勾选待操作路由器后，单击 **更多** - **设置网关** ，弹出“设置路由器网关”对话框。
2. “分配外部IP”选择“手动选择”，“子网”选择上述外部网络子网，并配置其他参数后，单击 **设置** ，完成操作。其中，其他参数的具体参数说明，请参考“SDN网络服务”帮助中“路由器”的相关内容。

3. 查看路由器ID和外部网络IP地址。

在“路由器”页面中，单击上述路由器名称，进入其详情页面。在该页面中，查看并记录该路由器的ID（即UUID参数的值）和外部网络IP地址（即外部IP参数的值）。

3. 导入VPC资源。

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[容器服务]-[安全容器服务]，进入“安全容器服务”页面。
2. 在左侧导航栏选择[管理视图]-[自定义资源管理]，或在左侧导航栏选择[业务视图]，并选择目标命名空间后，选择[自定义资源管理]，进入“自定义资源管理”页面。
3. 单击页面右下角的“Yaml”图标，进入“导入Yaml”页面。
4. 依据实际业务情况，输入Yaml文件内容，或直接单击编辑区域右上角的“导入”图标，导入预先配置的Yaml文件。Yaml文件格式如下（其中，name为自定义输入的VPC名称，externalGatewayIp为外部网络的IP地址，externalNetworkID为外部网络的ID，neutronRouter为路由器的ID）：

```
apiVersion: kubeovn.io/v1
kind: Vpc
metadata:
  name: test-vpc2
```

```
spec:
  externalGatewayIp: 172.110.0.160
  externalNetworkID: c9a831a0-6298-44c6-a664-2f362e60e419
  neutronRouter: c761887f-40bd-4df0-9b43-7c6bb009aab7
```

5. 待调试通过后，单击 **导入** ，完成操作。

6. spec关键字段说明：

externalNetworkID：外部网络ID(必填)

externalNetworkName：外部网络名字(选填)

外部网关IP(选填)；注：该ip为vpc下pod访问外网的默认snat ip，允许pod访问外网有两种方式：vpc yaml中定义该字段，或者网络页面设置路由网关开启snat

neutronRouter：路由ID(必填)

gatewayNode(选填, ovn pod访问节点的出入口节点，默认为网络节点)

自定义子网

> 警告：

>

> * 在同一VPC内，请确保各子网不发生冲突。

> * 在不同VPC之间，如需配置多个用户路由器之间的云内对等连接，请确保各路由器的子网不发生冲突。

1. 在云平台顶部导航栏中，依次选择[产品与服务]-[容器服务]-[安全容器服务]，进入“安全容器服务”页面。
2. 在左侧导航栏选择[管理视图]-[自定义资源管理]，或在左侧导航栏选择[业务视图]，并选择目标命名空间后，选择[自定义资源管理]，进入“自定义资源管理”页面。
3. 单击页面右下角的“Yaml”图标，进入“导入Yaml”页面。
4. 依据实际业务情况，输入子网内容，或直接单击编辑区域右上角的“导入”图标，导入预先配置的Yaml文件。Yaml文件格式如下（其中，name为自定义输入的子网名称，vpc为该子网所属VPC的名称，cidrBlock为该子网的网段，natOutgoing为是否允许访问外部网络）：

```
kind: Subnet
apiVersion: kubeovn.io/v1
metadata:
  name: net192
spec:
  vpc: test-vpc-2
  cidrBlock: 192.168.100.0/24
  natOutgoing: true
```

5. 待调试通过后，单击 **导入** ，完成操作。

6. spec 关键字段说明

vpc: 属于哪个vpc

cidrBlock: subnet cidr

natOutgoing: 是否可以访问外部网络

projectIDs: 配置该subnet只能被哪些projects使用，如果不配置表示 该subnet为share subnet 所用用户可用

7. 如果通过页面部署工作负载时不选择Subnet子网，会使用vpc中默认子网进行部署；如果通过yaml定义，spec样例如下：

```
apiVersion: v1
kind: Pod
metadata:
  annotations:
    ovn.kubernetes.io/logical_switch: another-subnet
  namespace: default
  name: another-subnet-pod
```

Floating IP 使用

fip 在容器产品中使用包括两个场景：pods 指定使用 SNAT作为出口IP；允许kubevirt 云原生云主机 使用 EIP。

使用方式有两种：页面操作，用户可在网络高级配置中可以选择正确的fip资源；或者yaml定义。

两种使用场景yaml spec定义样例如下： 1) SNAT出口IP

```
apiVersion: v1
kind: Pod
metadata:
  annotations:
    ovn.kubernetes.io/snat: 172.35.0.18
  name: snat-busybox
```

2) EIP

```
apiVersion: kubevirt.io/v1
kind: VirtualMachine
metadata:
  name: kubevirt-eip
spec:
  template:
    metadata:
      annotations:
        ovn.kubernetes.io/eip: 172.35.0.20
```


6.2 容器状态为未知错误，如何排查解决

问题描述

在云平台中，查看到容器的状态为“未知错误”。具体使用现象可能表现为：

- 在连接容器终端正常使用过程中，断开连接后无法再次建立连接。
- 在云监控服务中，上报微服务管理停止告警。
- 维护此容器时，维护失败。

问题原因

容器发生未知错误。

解决方案

1. 在顶部导航栏选择[产品与服务]-[容器服务]-[安全容器服务]，进入“安全容器服务”页面。
2. 在左侧导航栏选择[业务视图]页签-选择目标命名空间后，选择“容器组”，进入“容器组”页面。
3. 单击目标容器所在容器组的名称链接，进入详情页面。
4. 在[终端]页签中，选择目标容器后，执行以下命令：

```
ctr -n k8s.io t ls | grep UNKNOWN |awk '{print $1}'|xargs -I % ctr -n k8s.io t rm %
```

5. 确认问题已解决。具体命令如下（无返回数据，即表示问题已解决）：

```
ctr -n k8s.io t ls | grep UNKNOWN
```

6.3 容器状态为失败，如何排查解决

问题描述

在云平台中，查看到容器的状态为“失败”。

问题原因

容器命令执行错误。

解决方案

1. 在顶部导航栏选择[产品与服务]-[容器服务]-[安全容器服务]，进入“安全容器服务”页面。
2. 在左侧导航栏选择[业务视图]页签-选择目标命名空间后，选择“容器组”，进入“容器组”页面。
3. 单击目标容器所在容器组的名称链接，进入详情页面。
4. 在[事件]或[终端]页签中，排查具体执行错误的命令并解决。

7 部署指南

7.1 GPU解决方案对接包使用指南

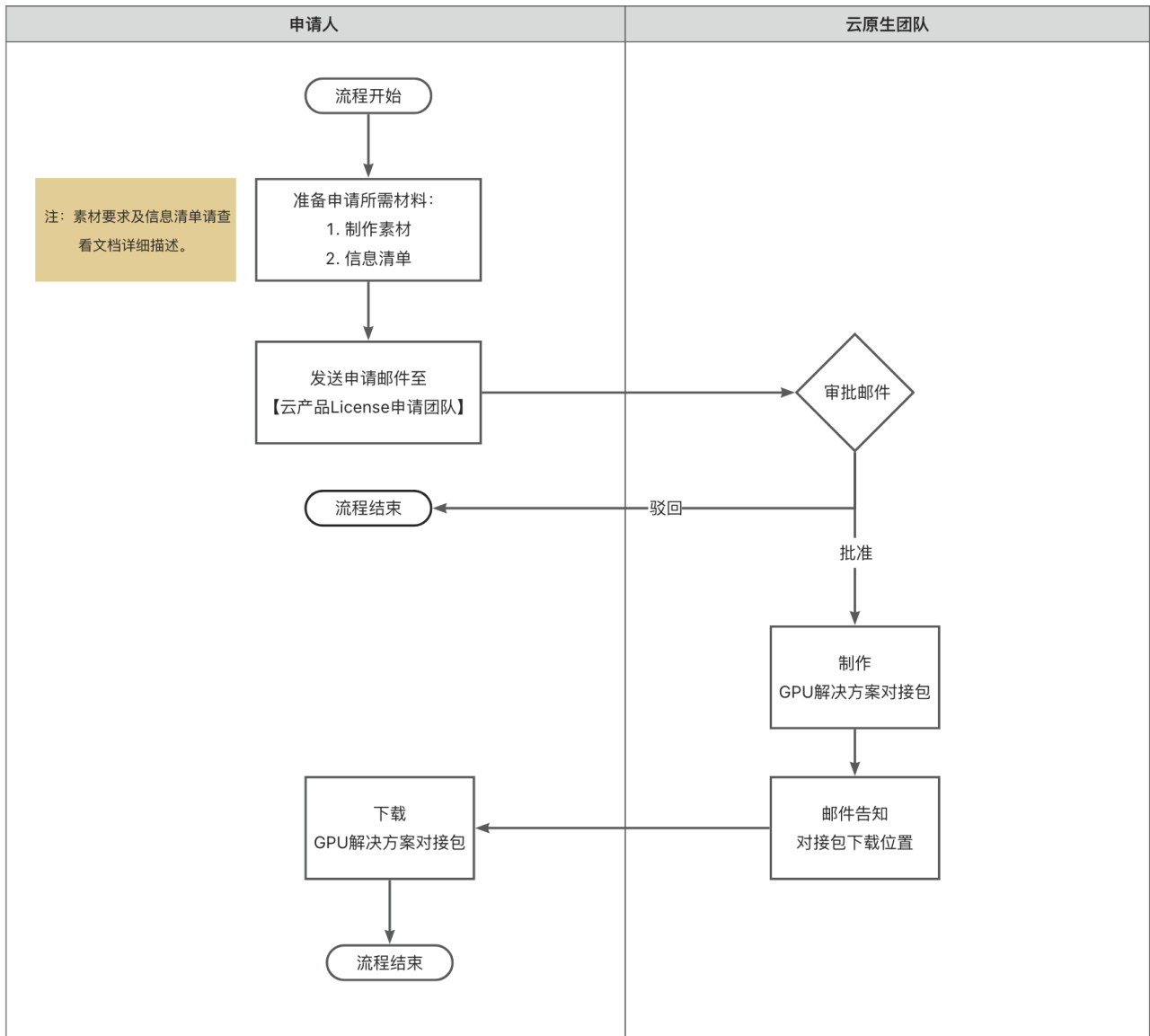
解决方案配置中心提供通过解决方案对接包，在环境部署完成后，经过配置方式实现环境配置、代码、组件变更的功能。

安全容器服务**「GPU解决方案对接包」**主要职责包括以下几个方面：

1. **安装GPU/XPU驱动YAML：**负责安装GPU或XPU驱动程序的YAML文件，确保系统正确识别和管理这些设备。
2. **上传驱动镜像：**将驱动程序镜像上传到相应的容器平台或云服务，供后续部署和使用。
3. **部署安全容器VM镜像：**负责部署安全容器VM镜像，实现与GPU或XPU卡的有效通信和资源管理。
4. **修改GPU配额相关设置：**根据需求，对与GPU配额相关的配置进行修改，以合理分配和管理GPU资源。

申请流程

安全容器服务的GPU解决方案支持通过**「邮件」**进行申请和审批流程(与云产品License申请流程一致)。根据要求，需求方可以通过邮件向云原生团队提交申请，并经过审批后，团队将制作相应的GPU解决方案对接包，并提供下载地址。详细的流程请参考下图：



操作步骤

1. 申请GPU解决方案对接包

提供以下素材和信息后，发邮件至cloudproducts_license@easystack.cn申请：

1.1 需要提供的素材

素材	获取方式
----	------

素材	获取方式
GPU/XPU Driver Yaml	GPU/XPU厂商提供
Driver镜像	GPU/XPU厂商提供

1.2 需要的信息

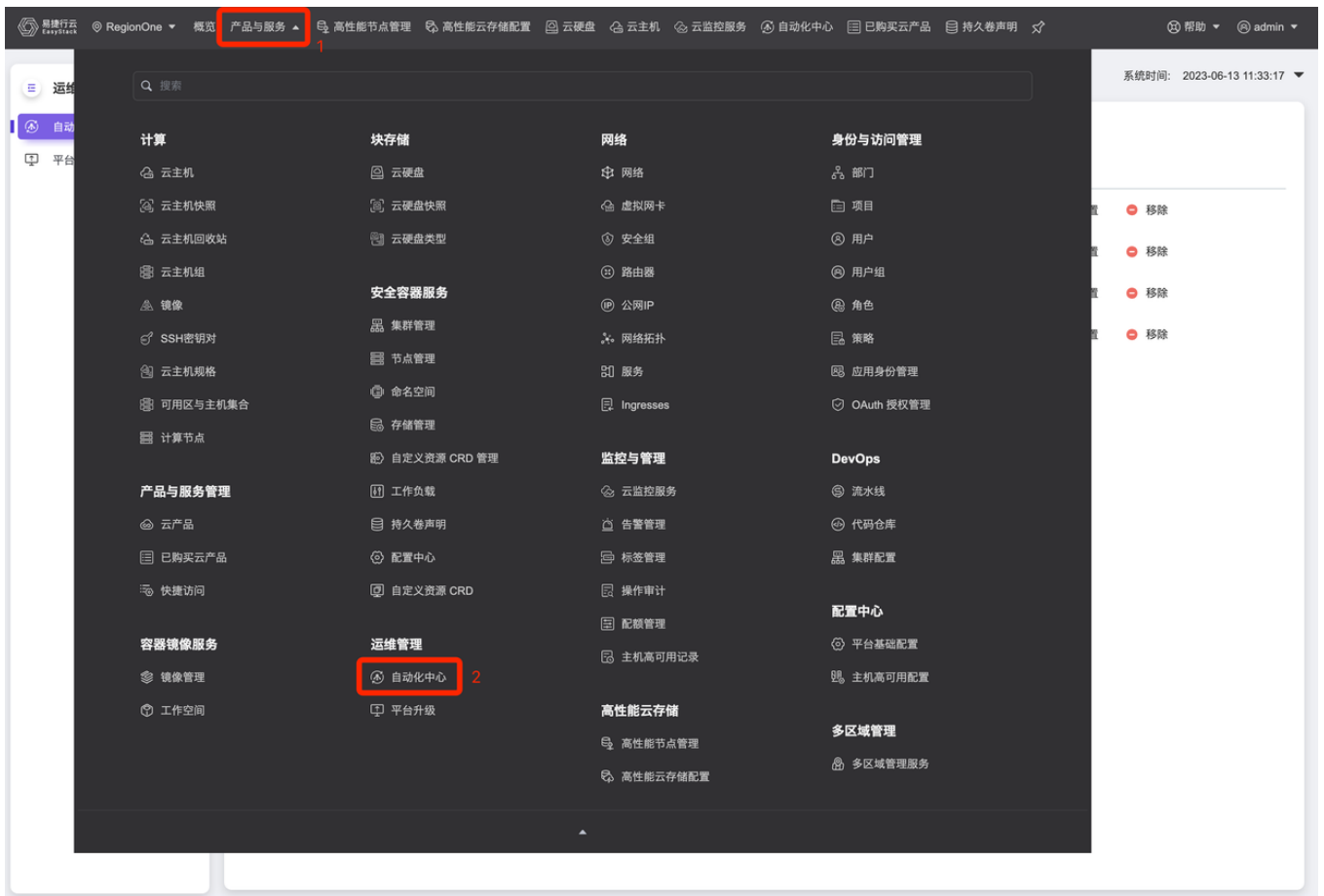
信息	获取方式
对接包部署平台的CPU架构	部署人员提供
GPU/XPU厂商	GPU/XPU厂商提供
GPU/XPU型号	GPU/XPU厂商提供(型号对应的是GPU在Kubernetes Pod Resouce中的字段)

2.获取GPU解决方案对接包

在对接包申请并审批通过后，云原生团队将根据提供的信息进行对接包制作。一旦制作完成，申请人将通过电子邮件收到对接包所在位置的通知。申请人可前往该位置获取对接包以进行后续操作。

3.GPU解决方案对接包使用指南

3.1 访问顶部导航，点击[产品与服务]，点击进入[自动化中心]界面



3.2 在[自动化中心]界面点击[高级配置]进入配置页

控制台 / 云环境 系统时间: 2023-06-13 11:35:10

云环境信息 许可信息 环境配置 存储配置 高级配置

57

当前版本 6.1.1

✔ 部署成功

客户名称	北京物理环境57	授权存储容量(TB)	193.00	3	2200-01-01	正式版
项目名称	QA	授权CPU(核)	192	许可节点数	维保有效期	许可类型
环境序列号	996adb2e-eced...					

[查看详情](#) [更新许可](#)

节点信息

节点状态

在线 **6**

离线 **2**

[节点管理](#) [导出配置](#)

名称	序列号	节点类型	角色	状态	授权状态	运行状态
node-1	210235A2CQ6214F002P3	物理节点	融合节点	● 已部署	● 已授权	✔ 在线
node-2	210235A2CQ6214F002P5	物理节点	融合节点	● 已部署	● 已授权	✔ 在线
node-3	210235A2CQ6214F002P4	物理节点	融合节点	● 已部署	● 已授权	✔ 在线
node-4	6b93d5a4-5bb8-4224-90fa-e6f785ecb89d	虚拟节点	云产品节点	● 已部署	● 已授权	✔ 在线
node-5	a9838939-f3e3-4b7d-a6f8-80c10e440f92	虚拟节点	云产品节点	● 已部署	● 已授权	⊘ 离线
node-6	c10ff929-bef4-43b3-a027-085a792d6b3f	虚拟节点	云产品节点	● 已部署	● 已授权	⊘ 离线

3.3 在[高级配置]子项，点击进入[解决方案配置中心]

控制台 / 云环境 系统时间: 2023-06-13 11:38:30

云环境信息 许可信息 环境配置 存储配置 **高级配置**

57

当前版本 6.1.1

部署成功

存储池

解决方案配置中心

节点信息

节点状态

在线 **6** 离线 **2**

节点管理 导出配置

名称	序列号	节点类型	角色	状态	授权状态	运行状态
node-1	210235A2CQ6214F002P3	物理节点	融合节点	已部署	已授权	在线
node-2	210235A2CQ6214F002P5	物理节点	融合节点	已部署	已授权	在线
node-3	210235A2CQ6214F002P4	物理节点	融合节点	已部署	已授权	在线
node-4	6b93d5a4-5bb8-4224-90fa-e6f785ecb89d	虚拟节点	云产品节点	已部署	已授权	在线
node-5	a9838939-f3e3-4b7d-a6f8-80c10e440f92	虚拟节点	云产品节点	已部署	已授权	离线
node-6	c10ff929-bef4-43b3-a027-085a792d6b3f	虚拟节点	云产品节点	已部署	已授权	离线

1 2

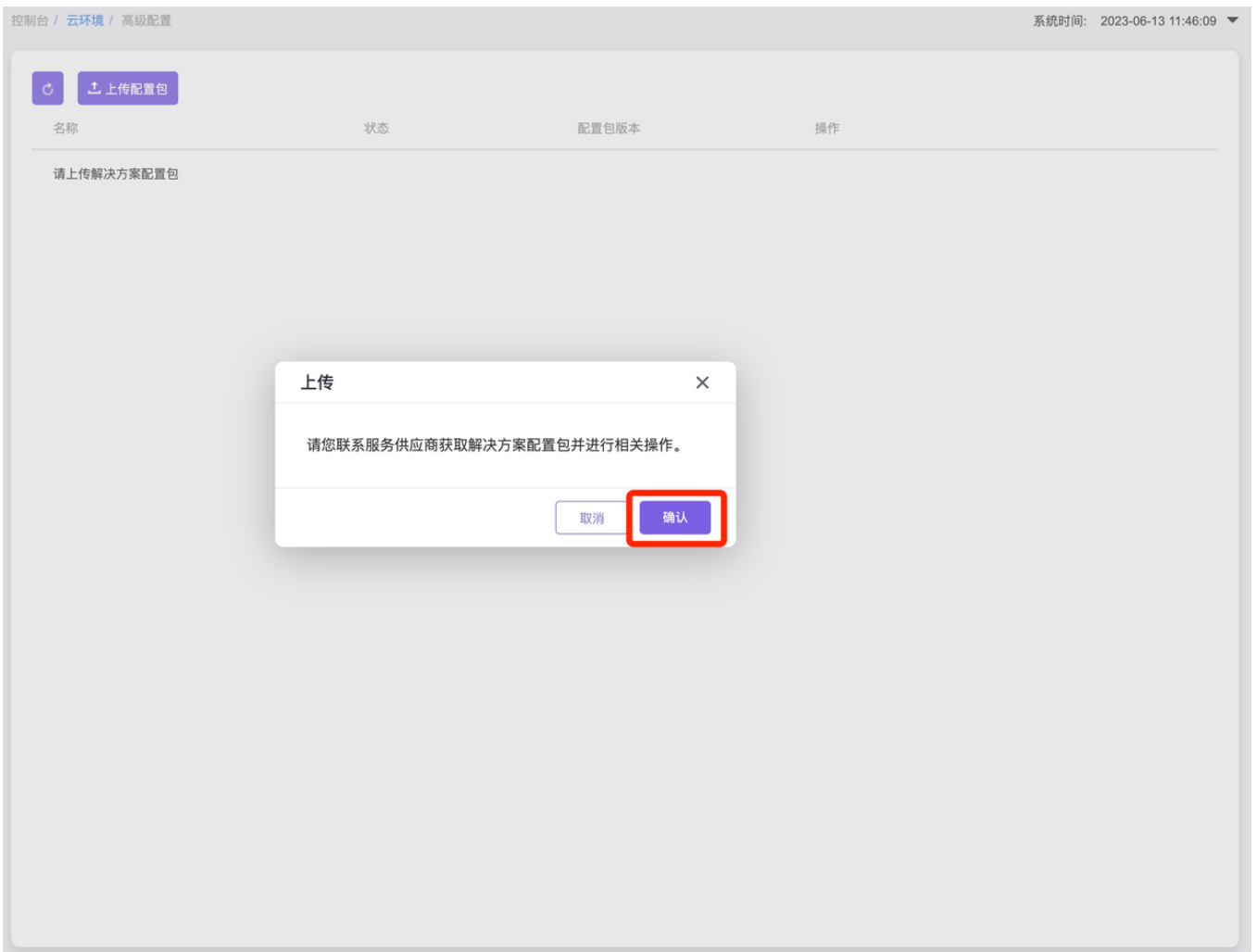
3.4 点击[上传配置包]

控制台 / 云环境 / 高级配置 系统时间: 2023-06-13 11:44:10

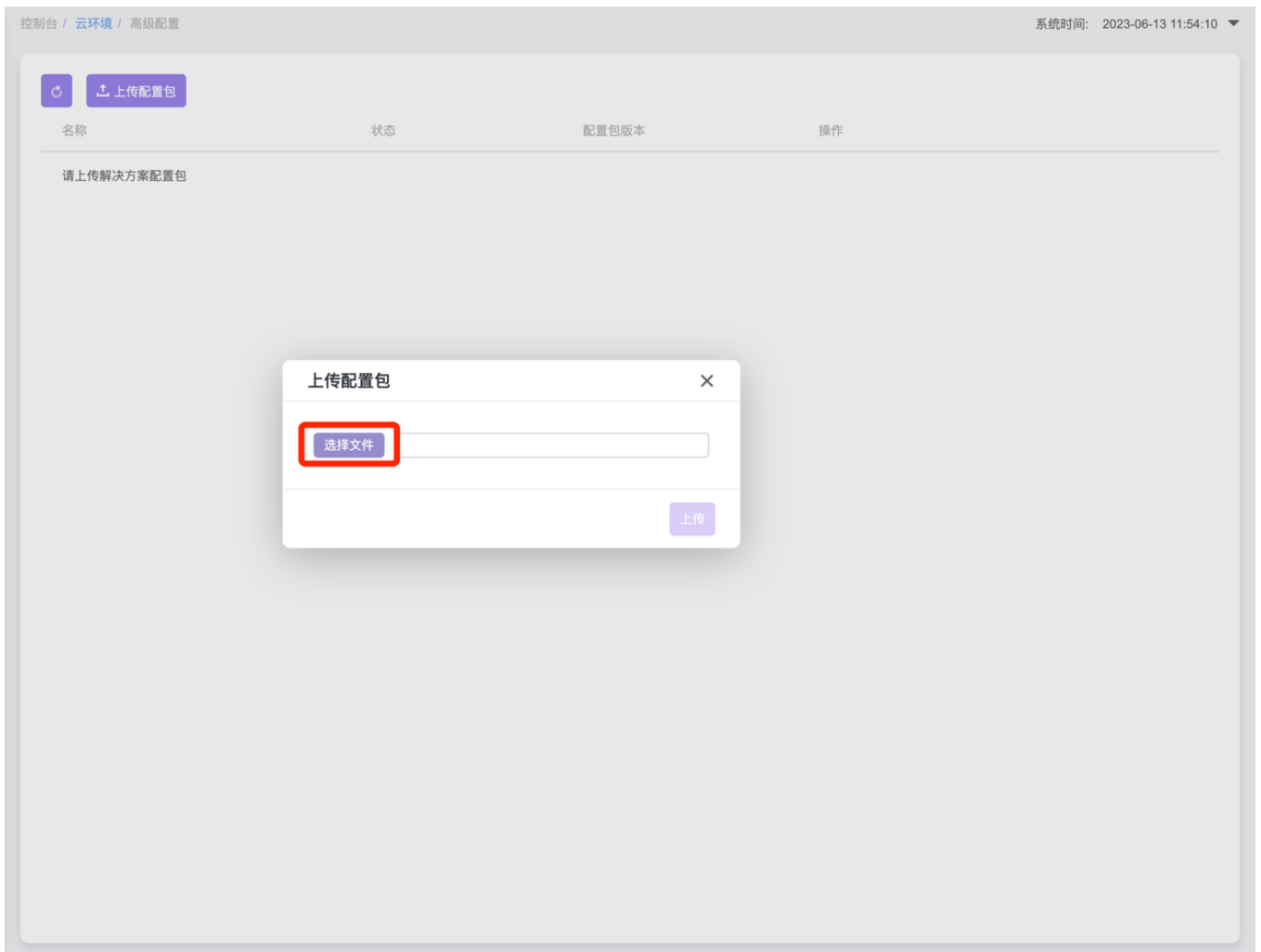
↻↑ 上传配置包

名称	状态	配置包版本	操作
请上传解决方案配置包			

3.5 点击[确认]



3.6 选择对应对接包文件后并点击[上传]



3.7 上传成功后需要等待配置加载完成。当状态更改为"已配置"时，表示GPU解决方案对接包已成功加载并可供使用。

上传配置包

名称	状态	配置包版本	操作
ecr-containers-nvidiaGpu	已配置	6.1.1	更新配置包 重新配置 移除
ecr-containers-baiduXpu	已配置	6.1.1	更新配置包 重新配置 移除

8 API参考

8.1 API文档模板

一级规格（例：云主机）

二级规格（例：启动云主机）

功能介绍

说明该操作实现的功能或效果，例：启动已停止的云主机并将其状态更改为“ACTIVE”。

前提条件（可选）

若执行本操作前存在必要的前提条件，请说明；若无，则删除。

接口约束（可选）

若执行本操作存在限制或注意事项，请说明；若无，则删除。

注意接口约束与前提条件的区别：

- 前提条件强调必须先做了什么才能执行本操作；
- 接口约束强调与本操作相关的注意，例如本操作带来的重要影响，执行本操作时不宜进行的其它操作等。

URI

示例： `POST /v2.1/{project-id}/servers/{server_id}/action`

说明：需使用“行内代码”样式。

参数	是否必选	描述

请求消息

参数	参数类型	是否必选	描述

响应消息

参数	参数类型	描述

请求示例

```
POST https://{endpoint}/v1/{project_id}/cloudservers/delete
```

```
{
  "servers": [
    {
      "id": "616fb98f-46ca-475e-917e-2563e5a8cd19"
    }
  ],
  "delete_publicip": false,
  "delete_volume": false
}
```

正常响应示例

```
{
  "user": {
    "name": "test_user",
    "links": {
```

```
      "self": "http://keystone-  
api.openstack.svc.cluster.local:35357/v3/users/5df4ae79648b4d7e954382da88cc6  
9ef"  
    },  
    "extra": {  
      "user_type": "individual",  
      "user_role": "domain_member"  
    },  
    "enabled": true,  
    "user_type": "individual",  
    "email": null,  
    "user_role": "domain_member",  
    "id": "5df4ae79648b4d7e954382da88cc69ef",  
    "domain_id": "default",  
    "password_expires_at": null  
  }  
}
```

正常响应代码

例：200

错误码

例：400, 401

咨询热线：400-100-3070

北京易捷思达科技发展有限公司：

北京市海淀区西北旺东路10号院东区1号楼1层107-2号

南京易捷思达软件科技有限公司：

江苏省南京市雨花台区软件大道168号润和创智中心4栋109-110

邮箱：

contact@easystack.cn (业务咨询)

partners@easystack.cn(合作伙伴咨询)

marketing@easystack.cn (市场合作)