



Red Hat OpenStack Platform 17.0

インスタンス作成のための Compute サービスの 設定

インスタンスを作成するための Red Hat OpenStack Platform Compute (nova) サービスの設定および管理に関するガイド

Red Hat OpenStack Platform 17.0 インスタンス作成のための Compute サービスの設定

インスタンスを作成するための Red Hat OpenStack Platform Compute (nova) サービスの設定および管理に関するガイド

OpenStack Team
rhos-docs@redhat.com

法律上の通知

Copyright © 2023 Red Hat, Inc.

The text of and illustrations in this document are licensed by Red Hat under a Creative Commons Attribution–Share Alike 3.0 Unported license ("CC-BY-SA"). An explanation of CC-BY-SA is available at

<http://creativecommons.org/licenses/by-sa/3.0/>

. In accordance with CC-BY-SA, if you distribute this document or an adaptation of it, you must provide the URL for the original version.

Red Hat, as the licensor of this document, waives the right to enforce, and agrees not to assert, Section 4d of CC-BY-SA to the fullest extent permitted by applicable law.

Red Hat, Red Hat Enterprise Linux, the Shadowman logo, the Red Hat logo, JBoss, OpenShift, Fedora, the Infinity logo, and RHCE are trademarks of Red Hat, Inc., registered in the United States and other countries.

Linux[®] is the registered trademark of Linus Torvalds in the United States and other countries.

Java[®] is a registered trademark of Oracle and/or its affiliates.

XFS[®] is a trademark of Silicon Graphics International Corp. or its subsidiaries in the United States and/or other countries.

MySQL[®] is a registered trademark of MySQL AB in the United States, the European Union and other countries.

Node.js[®] is an official trademark of Joyent. Red Hat is not formally related to or endorsed by the official Joyent Node.js open source or commercial project.

The OpenStack[®] Word Mark and OpenStack logo are either registered trademarks/service marks or trademarks/service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation, or the OpenStack community.

All other trademarks are the property of their respective owners.

概要

本ガイドでは、クラウド管理者が OpenStack クライアント CLI を使用して Red Hat OpenStack Platform Compute (nova) サービスを設定および管理するための概念および手順について説明します。

目次

多様性を受け入れるオープンソースの強化	4
RED HAT ドキュメントへのフィードバック (英語のみ)	5
第1章 COMPUTE サービス (NOVA) の機能	6
第2章 COMPUTE サービス (NOVA) の設定	8
2.1. オーバーコミットのためのメモリー設定	9
2.2. コンピュートノード上で確保するホストメモリーの計算	10
2.3. スワップサイズの計算	10
第3章 パフォーマンス改善のためのコンピュートノードの設定	12
3.1. コンピュートノードでの CPU ピニングの設定	12
3.2. エミュレータースレッドの設定	19
3.3. コンピュートノードでのヒュージページの設定	20
3.4. インスタンスにファイルベースのメモリーを使用するコンピュートノードの設定	25
第4章 COMPUTE サービスのストレージの設定	28
4.1. イメージキャッシュの設定オプション	28
4.2. インスタンスの一時ストレージ属性の設定オプション	30
4.3. 共有インスタンスストレージの設定	33
4.4. RED HAT CEPH RADOS BLOCK DEVICE (RBD) からの直接イメージダウンロードの設定	35
4.5. 関連情報	36
第5章 PCI パススルーの設定	37
5.1. PCI パススルー用コンピュートノードの指定	37
5.2. PCI パススルー用コンピュートノードの設定	40
5.3. PCI パススルーデバイス種別フィールド	43
5.4. NOVAPCIPASSTHROUGH 設定のガイドライン	43
第6章 VDPA ポートを使用するインスタンスを有効にするための VDPA コンピュートノードの設定	46
6.1. VDPA 用コンピュートノードの指定	46
6.2. VDPA コンピュートノードの設定	49
第7章 ホストアグリゲートの作成と管理	52
7.1. ホストアグリゲートでのスケジューリングの有効化	52
7.2. ホストアグリゲートの作成	53
7.3. アベイラビリティゾーンの作成	54
7.4. ホストアグリゲートの削除	55
7.5. プロジェクト分離ホストアグリゲートの作成	56
第8章 インスタンスのスケジューリングと配置の設定	59
8.1. PLACEMENT サービスを使用した事前絞り込み	59
8.2. COMPUTE スケジューラーサービス用フィルターおよび重みの設定	65
8.3. COMPUTE スケジューラーのフィルター	66
8.4. COMPUTE スケジューラーの重み	72
第9章 インスタンス起動用のフレーバーの作成	81
9.1. フレーバーの作成	81
9.2. フレーバーの引数	83
9.3. フレーバーのメタデータ	84
第10章 インスタンスへのメタデータの追加	99
10.1. インスタンスメタデータの種別	99

10.2. 全インスタンスへのコンフィグドライブの追加	99
10.3. インスタンスへの動的メタデータの追加	101
第11章 インスタンスの CPU 機能フラグの設定	103
11.1. 前提条件	103
11.2. インスタンスの CPU 機能フラグの設定	103
第12章 KERNELARGS を定義するための手動でのノード再起動の設定	106
12.1. KERNELARGS を定義するための手動でのノード再起動の設定	106
第13章 インスタンス用の仮想 GPU の設定	108
13.1. サポートされる設定および制限	108
13.2. コンピュートノードでの仮想 GPU の設定	109
13.3. カスタム GPU インスタンスイメージの作成	113
13.4. インスタンス用の仮想 GPU フレーバーの作成	114
13.5. 仮想 GPU インスタンスの起動	115
13.6. GPU デバイスの PCI パススルーの有効化	116
第14章 インスタンスの管理	120
14.1. インスタンスの VNC コンソールへの接続のセキュリティー保護	120
14.2. データベースのクリーニング	121
14.3. コンピュートノード間の仮想マシンインスタンスの移行	125

多様性を受け入れるオープンソースの強化

Red Hat では、コード、ドキュメント、Web プロパティにおける配慮に欠ける用語の置き換えに取り組んでいます。まずは、マスター (master)、スレーブ (slave)、ブラックリスト (blacklist)、ホワイトリスト (whitelist) の 4 つの用語の置き換えから始めます。この取り組みは膨大な作業を要するため、今後の複数のリリースで段階的に用語の置き換えを実施して参ります。詳細は、[Red Hat CTO である Chris Wright のメッセージ](#) を参照してください。

RED HAT ドキュメントへのフィードバック (英語のみ)

Red Hat ドキュメントに対するご意見をお聞かせください。改善点などがあれば、ぜひお知らせください。

ドキュメント直接フィードバック (DDF) 機能の使用 (英語版のみ)

特定の文章、段落、またはコードブロックに対して直接コメントを送付するには、DDF の **Add Feedback** 機能を使用してください。なお、この機能は英語版のドキュメントでのみご利用いただけます。

1. **Multi-page HTML** 形式でドキュメントを表示します。
2. ドキュメントの右上隅に **Feedback** ボタンが表示されていることを確認してください。
3. コメントするテキスト部分をハイライト表示します。
4. **Add Feedback** をクリックします。
5. **Add Feedback** フィールドにコメントを入力します。
6. (オプション) ご自身のメールアドレスを追加してください。ドキュメントチームが問題の詳細を確認するために連絡を取ることがあります。
7. **Submit** をクリックします。

第1章 COMPUTE サービス (NOVA) の機能

Compute (nova) サービスを使用して、Red Hat OpenStack Platform (RHOSP) 環境で仮想マシンインスタンスおよびベアメタルサーバーを作成、プロビジョニング、および管理します。Compute サービスは、ベースとなるホストプラットフォームの詳細を公開するのではなく、Compute サービスを実行しているベースとなるハードウェアを抽象化します。たとえば、ホスト上で動作中の CPU の種別およびトポロジを公開するのではなく、Compute サービスは多数の仮想 CPU (vCPU) を公開し、これらの仮想 CPU のオーバーコミットに対応します。

Compute サービスは、KVM ハイパーバイザーを使用して Compute サービスのワークロードを実行します。libvirt ドライバーは QEMU と協調して KVM との相互のやり取りをすべて処理し、仮想マシンインスタンスの作成を可能にします。インスタンスを作成およびプロビジョニングするために、Compute サービスは以下の RHOSP サービスと協調します。

- 認証のための Identity (keystone) サービス
- リソースインベントリをトラッキングおよび選択するための Placement サービス
- ディスクおよびインスタンスイメージのための Image サービス (glance)
- インスタンスがブート時に接続する仮想ネットワークまたは物理ネットワークをプロビジョニングするための Networking (neutron) サービス

Compute サービスは、**nova-*** という名前のデーモンプロセスおよびサービスで設定されます。コアの Compute サービスを以下に示します。

Compute サービス (nova-compute)

このサービスは、KVM または QEMU ハイパーバイザー API の libvirt を使用してインスタンスを作成、管理、および削除し、インスタンスの状態をデータベースを更新します。

Compute コンダクター (nova-conductor)

このサービスは、Compute サービスとデータベースの協調を仲介します。これにより、コンピュータードをデータベースへの直接アクセスから隔離します。このサービスを **nova-compute** サービスが実行されているノードにデプロイしないでください。

Compute スケジューラー (nova-scheduler)

このサービスはキューからインスタンスの要求を受け取り、インスタンスをホストするコンピュータードを決定します。

Compute API (nova-api)

このサービスは、ユーザーに外部 REST API を提供します。

API データベース

このデータベースはインスタンスの場所の情報をトラッキングします。また、ビルドされているがスケジュールされていないインスタンスの一時的な場所を提供します。マルチセルのデプロイメントでは、このデータベースには各セルのデータベース接続を指定するセルのマッピングも含まれます。

セルデータベース

このデータベースには、インスタンスに関するほとんどの情報が含まれます。API データベース、コンダクター、および Compute サービスによって使用されます。

メッセージキュー

このメッセージングサービスは、セル内の各サービス間の通信およびグローバルなサービスとの通信のために、すべてのサービスによって使用されます。

Compute メタデータ

このサービスは、インスタンス固有のデータを保管します。インスタンス

は、<http://169.254.169.254> または IPv6 を通じてリンクローカルアドレス 80::a9fe:a9fe から、メタデータサービスにアクセスします。Networking (neutron) サービスは、要求をメタデータ API サーバーに転送するロールを持ちます。**NeutronMetadataProxySharedSecret** パラメーターを使用して、Networking サービスと Compute サービス両方の設定でシークレットキーワードを設定する必要があります。これにより、これらのサービスが通信を行うことができます。Compute メタデータサービスは、Compute API の一部としてグローバルに実行することや、それぞれのセルで実行することができます。

複数のコンピュートノードをデプロイすることができます。インスタンスを操作するハイパーバイザーは、それぞれのコンピュートノードで実行されます。それぞれのコンピュートノードには、少なくとも2つのネットワークインターフェイスが必要です。コンピュートノードでは、インスタンスを仮想ネットワークに接続し、セキュリティーグループを介してインスタンスにファイアウォールサービスを提供する Networking サービスエージェントも実行されます。

デフォルトでは、director はすべてのコンピュートノードを単一のセルとしてオーバークラウドをインストールします。このセルには、仮想マシンインスタンスを制御および管理するすべての Compute サービスおよびデータベース、ならびにすべてのインスタンスおよびインスタンスのメタデータが含まれます。大規模なデプロイメントでは、複数のセルでオーバークラウドをデプロイし、多数のコンピュートノードに対応することができます。新しいオーバークラウドをインストールする際に、またはその後の任意の時に、セルを環境に追加することができます。

第2章 COMPUTE サービス (NOVA) の設定

クラウド管理者は、環境ファイルを使用して Compute (nova) サービスをカスタマイズします。Puppet は、この設定を生成して `/var/lib/config-data/puppet-generated/<nova_container>/etc/nova/nova.conf` ファイルに保存します。Compute サービスの設定をカスタマイズするには、以下の設定方法を順番どおりに使用します。

1. **heat パラメーター - オーバークラウドパラメーター** ガイドの [Compute \(nova\) パラメーター](#) セクションに詳細が記載されています。以下の例では、heat パラメーターを使用して、デフォルトのスケジューラーフィルターを設定し、Compute サービスの NFS バックエンドを設定します。

```
parameter_defaults:
  NovaNfsEnabled: true
  NovaNfsOptions: "context=system_u:object_r:nfs_t:s0"
  NovaNfsShare: "192.0.2.254:/export/nova"
  NovaNfsVersion: "4.2"
  NovaSchedulerEnabledFilters:
    - AggregateInstanceExtraSpecsFilter
    - ComputeFilter
    - ComputeCapabilitiesFilter
    - ImagePropertiesFilter
```

2. **Puppet パラメーター:** `/etc/puppet/modules/nova/manifests/*` で定義されます。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::compute::force_raw_images: True
```



注記

この方法は、等価な heat パラメーターが存在しない場合にのみ使用してください。

3. **手動での hieradata の上書き:** heat パラメーターまたは Puppet パラメーターが存在しない場合の、パラメーターカスタマイズ用。たとえば、以下の設定により Compute ロールの **[DEFAULT]** セクションに **timeout_nbd** が定義されます。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::config::nova_config:
      DEFAULT/timeout_nbd:
        value: '20'
```



警告

heat パラメーターが存在する場合は、Puppet パラメーターではなく heat パラメーターを使用します。Puppet パラメーターは存在するが heat パラメーターは存在しない場合は、手動で上書きする方法ではなく Puppet パラメーターを使用します。等価な heat パラメーターまたは Puppet パラメーターがない場合に限り、手動で上書きする方法を使用してください。

ヒント

[変更するパラメーター](#) の特定のガイダンスに従って、特定の設定をカスタマイズするために、heat または Puppet パラメーターを使用できるかどうかを判断します。

オーバークラウドサービスの設定方法に関する詳細は、[director のインストールおよび使用方法 ガイド](#) の [Heat パラメーター](#) を参照してください。

2.1. オーバーコミットのためのメモリー設定

メモリーのオーバーコミットを使用する場合 (**NovaRAMAllocationRatio** ≥ 1.0)、割り当て率をサポートするのに十分なスワップ領域を確保してオーバークラウドをデプロイする必要があります。



注記

NovaRAMAllocationRatio パラメーターが **1 より小さい値** に設定されている場合は、RHEL でのスワップサイズの推奨事項に従ってください。詳しくは、RHEL の [ストレージデバイスの管理](#) の [システムの推奨スワップ領域](#) を参照してください。

前提条件

- ノードに必要なスワップサイズが計算されている。詳細は、[スワップサイズの計算](#) を参照してください。

手順

1. `/usr/share/openstack-tripleo-heat-templates/environments/enable-swap.yaml` ファイルを環境ファイルのディレクトリーにコピーします。

```
$ cp /usr/share/openstack-tripleo-heat-templates/environments/enable-swap.yaml
/home/stack/templates/enable-swap.yaml
```

2. 以下のパラメーターを **enable-swap.yaml** ファイルに追加して、スワップサイズを設定します。

```
parameter_defaults:
  swap_size_megabytes: <swap size in MB>
  swap_path: <full path to location of swap, default: /swap>
```

3. その他の環境ファイルと共に **enable_swap.yaml** 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
  -e [your environment files] \
  -e /home/stack/templates/enable-swap.yaml
```

2.2. コンピュートノード上で確保するホストメモリーの計算

ホストのプロセス用に確保する RAM 容量の合計を判断するには、以下の各項目に十分なメモリーを割り当てる必要があります。

- ホスト上で実行されるリソース (たとえば、OSD は 3 GB のメモリーを消費します)
- ホストインスタンスに必要なエミュレーターのオーバーヘッド
- 各インスタンスのハイパーバイザー

メモリーへの追加要求を計算したら、以下の式を使用して各ノードのホストプロセス用に確保するメモリーの容量を決定します。

$$\text{NovaReservedHostMemory} = \text{total_RAM} - ((\text{vm_no} * (\text{avg_instance_size} + \text{overhead})) + (\text{resource1} * \text{resource_ram}) + (\text{resourcen} * \text{resource_ram}))$$

- **vm_no** は、インスタンスの数に置き換えてください。
- **avg_instance_size** は、各インスタンスが使用できるメモリーの平均容量に置き換えてください。
- **overhead** は、各インスタンスに必要なハイパーバイザーのオーバーヘッドに置き換えてください。
- **resource1** および **<resourcen>** までのすべてのリソースを、ノード上のリソース種別の数に置き換えてください。
- **resource_ram** は、この種別の各リソースに必要な RAM の容量に置き換えてください。

2.3. スワップサイズの計算

割り当てるスワップサイズは、メモリーのオーバーコミットを処理するのに十分な容量でなければなりません。以下の式を使用して、ノードに必要なスワップサイズを計算することができます。

- $\text{overcommit_ratio} = \text{NovaRAMAllocationRatio} - 1$
- 最小スワップサイズ (MB) = $(\text{total_RAM} * \text{overcommit_ratio}) + \text{RHEL_min_swap}$
- 推奨 (最大) スワップサイズ (MB) = $\text{total_RAM} * (\text{overcommit_ratio} + \text{percentage_of_RAM_to_use_for_swap})$

percentage_of_RAM_to_use_for_swap 変数は、QEMU のオーバーヘッドおよびオペレーティングシステムまたはホストのサービスが消費するその他のリソースに対応するバッファです。

たとえば、RAM 容量の合計が 64 GB で **NovaRAMAllocationRatio** が 1 に設定されている場合、利用可能な RAM の 25% をスワップに使用するには、

- 推奨 (最大) スワップサイズ = $64000 \text{ MB} * (0 + 0.25) = 16000 \text{ MB}$

NovaReservedHostMemory の値を計算する方法は、[コンピュートノード上で確保するホストメモリーの計算](#) を参照してください。

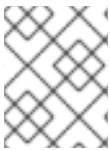
RHEL_min_swap の値を決定する方法は、RHEL のストレージデバイスの管理の [システムの推奨スワップ領域](#) を参照してください。

第3章 パフォーマンス改善のためのコンピュートノードの設定

クラウドユーザーは、インスタンスのスケジューリングおよび配置を設定して、最大のパフォーマンスを得ることができます。そのためには、NFV や高性能コンピューティング (HPC) などの特化されたワークロードを対象にするカスタムフレーバーを作成します。

以下の機能を使用して、最大のパフォーマンスを得るためにインスタンスを調整します。

- **CPU ピニング:** 仮想 CPU を物理 CPU に固定します。
- **エミュレータスレッド:** インスタンスに関連付けられたエミュレータスレッドを物理 CPU に固定します。
- **ヒュージページ:** 通常のメモリー (4 KB ページ) とヒュージページ (2 MB または 1 GB ページ) の両方について、インスタンスのメモリー割り当てポリシーを調整します。



注記

これらの機能のいずれかを設定すると、インスタンス上に NUMA トポロジが存在しない場合、暗黙的な NUMA トポロジが作成されます。

3.1. コンピュートノードでの CPU ピニングの設定

コンピュートノードで CPU ピニングを有効化することで、各インスタンスの CPU プロセスを専用のホスト CPU で実行するように設定することができます。インスタンスが CPU ピニングを使用する場合には、各インスタンスの仮想 CPU プロセスには、他のインスタンスの仮想 CPU プロセスが使用できない独自のホストの物理 CPU が割り当てられます。CPU ピニングが設定されたコンピュートノード上で動作するインスタンスには、NUMA トポロジがあります。インスタンスの NUMA トポロジの各 NUMA ノードは、ホストコンピュートノード上の NUMA ノードにマッピングされます。

専用の (ピンングされた) CPU を持つインスタンスと共有 (フローティング) の CPU を持つインスタンスを同じコンピュートノード上にスケジューリングするように、Compute のスケジューラーを設定することができます。NUMA トポロジを持つコンピュートノード上で CPU ピニングを設定するには、以下の手順を実施する必要があります。

1. CPU ピニング用のコンピュートノードを指定する。
2. ピニングされたインスタンス仮想 CPU プロセス、フローティングのインスタンス仮想 CPU プロセス、およびホストのプロセス用にホストコアを確保するようにコンピュートノードを設定する。
3. オーバークラウドをデプロイする。
4. CPU ピニングを要求するインスタンスを起動するためのフレーバーを作成する。
5. 共有 (あるいはフローティング) の CPU を使用するインスタンスを起動するためのフレーバーを作成する。

3.1.1. 前提条件

- コンピュートノードの NUMA トポロジを把握している。

3.1.2. CPU ピニング用コンピュートノードの指定

ピンングされた CPU を持つインスタンスのコンピュートノードを指定するには、新しいロールファイルを作成して CPU ピンングロールを設定し、CPU ピンング用のコンピュートノードにタグを付けるために使用する CPU ピンングリソースクラスを使用してベアメタルノードを設定する必要があります。



注記

以下の手順は、まだプロビジョニングされていない新しいオーバークラウドノードに適用されます。すでにプロビジョニングされている既存のオーバークラウドノードにリソースクラスを割り当てるには、スケールダウン手順を使用してノードのプロビジョニングを解除してから、スケールアップ手順を使用して新しいリソースクラスの割り当てでノードを再プロビジョニングする必要があります。詳細は、[オーバークラウドノードのスケールリング](#)を参照してください。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. **stackrc** ファイルを取得します。

```
[stack@director ~]$ source ~/stackrc
```

3. **Controller**、**Compute**、**ComputeCPUPinning** ロール、およびオーバークラウドに必要なその他のロールを含む、**roles_data_cpu_pinning.yaml** という名前の新しいロールデータファイルを生成します。

```
(undercloud)$ openstack overcloud roles \
generate -o /home/stack/templates/roles_data_cpu_pinning.yaml \
Compute:ComputeCPUPinning Compute Controller
```

4. **roles_data_cpu_pinning.yaml** を開き、以下のパラメーターおよびセクションを編集または追加します。

セクション/パラメーター	現在の値	新しい値
ロールのコメント	Role: Compute	Role: ComputeCPUPinning
ロール名	Compute	name: ComputeCPUPinning
description	Basic Compute Node role	CPU Pinning Compute Node role
HostnameFormatDefault	%stackname%- novacompute-%index%	%stackname%- novacomputepinning- %index%
deprecated_nic_config_name	compute.yaml	compute-cpu- pinning.yaml

5. オーバークラウド用の CPU ピニングコンピュートノードをノード定義のテンプレート **node.json** または **node.yaml** に追加して、そのノードを登録します。詳細は、**director** のインストールと使用方法 の [オーバークラウドノードの登録](#) を参照してください。
6. ノードのハードウェアを検査します。

```
(undercloud)$ openstack overcloud node introspect \
--all-manageable --provide
```

詳細については、**Director のインストールと使用** ガイドの [ベアメタルノードハードウェアのインベントリーの作成](#) を参照してください。

7. CPU ピニング用に指定する各ベアメタルノードに、カスタムの CPU ピニングリソースクラスをタグ付けします。

```
(undercloud)$ openstack baremetal node set \
--resource-class baremetal.CPU-PINNING <node>
```

<node> をベアメタルノードの ID に置き換えてください。

8. ノード定義ファイル **overcloud-baremetal-deploy.yaml** に **ComputeCPUPinning** ロールを追加し、予測ノード配置、リソースクラス、ネットワークポロジー、またはノードに割り当てるその他の属性を定義します。

```
- name: Controller
  count: 3
- name: Compute
  count: 3
- name: ComputeCPUPinning
  count: 1
  defaults:
    resource_class: baremetal.CPU-PINNING
    network_config:
      template: /home/stack/templates/nic-config/myRoleTopology.j2 1
```

- 1** 既存のネットワークポロジーを再利用するか、ロール用の新しいカスタムネットワークインターフェイステンプレートを作成できます。詳細は、**Director のインストールと使用** ガイドの [カスタムネットワークインターフェイステンプレート](#) を参照してください。**network_config** プロパティを使用してネットワーク定義を定義しない場合、デフォルトのネットワーク定義が使用されます。

ノード定義ファイルでノード属性を設定するために使用できるプロパティについて詳しくは、[ベアメタルノードのプロビジョニング属性](#) を参照してください。ノード定義ファイルの例は、[ノード定義ファイルの例](#) を参照してください。

9. プロビジョニングコマンドを実行して、ロールの新しいノードをプロビジョニングします。

```
(undercloud)$ openstack overcloud node provision \
--stack <stack> \
[--network-config \]
--output /home/stack/templates/overcloud-baremetal-deployed.yaml \
/home/stack/templates/overcloud-baremetal-deploy.yaml
```

- **<stack>** を、ベアメタルノードがプロビジョニングされるスタックの名前に置き換えます。指定しない場合、デフォルトは **overcloud** です。
 - **--network-config** オプションの引数を含めて、**cli-overcloud-node-network-config.yaml** Ansible Playbook にネットワーク定義を提供します。**network_config** プロパティを使用してネットワーク定義を定義しない場合、デフォルトのネットワーク定義が使用されます。
10. 別のターミナルでプロビジョニングの進捗をモニターリングします。プロビジョニングが成功すると、ノードの状態が **available** から **active** に変わります。

```
(undercloud)$ watch openstack baremetal node list
```

11. **--network-config** オプションを指定してプロビジョニングコマンドを実行しなかった場合は、**network-environment.yaml** ファイルで **<Role>NetworkConfigTemplate** パラメーターを設定して、NIC テンプレートファイルを指すようにします。

```
parameter_defaults:
  ComputeNetworkConfigTemplate: /home/stack/templates/nic-configs/compute.j2
  ComputeCPUPinningNetworkConfigTemplate: /home/stack/templates/nic-
  configs/<cpu_pinning_net_top>.j2
  ControllerNetworkConfigTemplate: /home/stack/templates/nic-configs/controller.j2
```

<cpu_pinning_net_top> を **ComputeCPUPinning** ロールのネットワークトポロジが含まれるファイルの名前に置き換えます。たとえば、デフォルトのネットワークトポロジを使用する場合は **compute.yaml** です。

3.1.3. CPU ピニング用コンピュートノードの設定

ノードの NUMA トポロジに基づいて、コンピュートノードでの CPU ピニングを設定します。効率を高めるために、全 NUMA ノードにわたって、CPU コアの一部をホストのプロセス用に確保します。残りの CPU コアをインスタンスの管理に割り当てます。

以下の手順では、以下の NUMA トポロジ (8 つの CPU コアを 2 つの NUMA ノードに分散) を使用して、CPU ピニングの設定方法を説明します。

表3.1 NUMA トポロジの例

NUMA ノード 0		NUMA ノード 1	
コア 0	コア 1	コア 2	コア 3
コア 4	コア 5	コア 6	コア 7

以下の手順では、コア 0 および 4 をホストのプロセス用に、コア 1、3、5、および 7 を CPU ピニングが必要なインスタンス用に、そしてコア 2 および 6 を CPU ピニングが不要なフローティングインスタンス用に、それぞれ確保します。

手順

1. ピニングされたインスタンス、フローティングのインスタンス、およびホストプロセス用にコアを確保するようにコンピュートノードを設定する環境ファイルを作成します (例: **cpu_pinning.yaml**)。

2. NUMA 対応コンピュートノードに NUMA トポロジーが設定されたインスタンスをスケジュールするには、Compute 環境ファイルの **NovaSchedulerEnabledFilters** パラメーターに **NUMATopologyFilter** がなければ、このフィルターを追加します。

```
parameter_defaults:
  NovaSchedulerEnabledFilters:
    - AvailabilityZoneFilter
    - ComputeFilter
    - ComputeCapabilitiesFilter
    - ImagePropertiesFilter
    - ServerGroupAntiAffinityFilter
    - ServerGroupAffinityFilter
    - PciPassthroughFilter
    - NUMATopologyFilter
```

NUMATopologyFilter の詳細は、[Compute scheduler filters](#) を参照してください。

3. 専用のインスタンス用に物理 CPU コアを確保するには、以下の設定を **cpu_pinning.yaml** に追加します。

```
parameter_defaults:
  ComputeCPUPinningParameters:
    NovaComputeCpuDedicatedSet: 1,3,5,7
```

4. 共有のインスタンス用に物理 CPU コアを確保するには、以下の設定を **cpu_pinning.yaml** に追加します。

```
parameter_defaults:
  ComputeCPUPinningParameters:
    ...
    NovaComputeCpuSharedSet: 2,6
```

5. ホストのプロセス用に確保する RAM 容量を指定するには、以下の設定を **cpu_pinning.yaml** に追加します。

```
parameter_defaults:
  ComputeCPUPinningParameters:
    ...
    NovaReservedHostMemory: <ram>
```

<ram> を、確保するメモリー容量 (MB 単位) に置き換えます。

6. インスタンス用に確保した CPU コアでホストプロセスが実行されないようにするには、**IsolCpusList** パラメーターに、インスタンス用に確保した CPU コアを設定します。

```
parameter_defaults:
  ComputeCPUPinningParameters:
    ...
    IsolCpusList: 1-3,5-7
```

コンマ区切りの CPU インデックスのリストまたは範囲を使用して、**IsolCpusList** パラメーターの値を指定します。

7. 新しいファイルを他の環境ファイルと一緒にスタックに追加し、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-r /home/stack/templates/roles_data_cpu_pinning.yaml \
-e /home/stack/templates/network-environment.yaml \
-e /home/stack/templates/cpu_pinning.yaml \
-e /home/stack/templates/overcloud-baremetal-deployed.yaml \
-e /home/stack/templates/node-info.yaml
```

3.1.4. インスタンス用の専用 CPU フレーバーの作成

クラウドユーザーが専用の CPU を持つインスタンスを作成できるようにするには、インスタンス起動用の専用 CPU ポリシーが設定されたフレーバーを作成します。

前提条件

- ホストで同時マルチスレッド (SMT) が有効である。
- コンピュートノードが CPU ピニングを許可するように設定されている。詳しくは、[コンピュートノードでの CPU ピニングの設定](#) を参照してください。

手順

1. source コマンドで **overcloudrc** ファイルを読み込みます。

```
(undercloud)$ source ~/overcloudrc
```

2. CPU ピニングを要求するインスタンス用のフレーバーを作成します。

```
(overcloud)$ openstack flavor create --ram <size_mb> \
--disk <size_gb> --vcpus <no_reserved_vcpus> pinned_cpus
```

3. ピニングされた CPU を要求するには、フレーバーの **hw:cpu_policy** 属性を **dedicated** に設定します。

```
(overcloud)$ openstack flavor set \
--property hw:cpu_policy=dedicated pinned_cpus
```

4. それぞれの仮想 CPU をスレッドシブリングに配置するには、フレーバーの **hw:cpu_thread_policy** 属性を **require** に設定します。

```
(overcloud)$ openstack flavor set \
--property hw:cpu_thread_policy=require pinned_cpus
```



注記

- ホストに SMT アーキテクチャがない場合や、スレッドシブリングが利用可能な CPU コアが十分でない場合には、スケジューリングが失敗します。これを回避するには、**hw:cpu_thread_policy** を **require** ではなく **prefer** に設定します。**prefer** ポリシーは、スレッドシブリングが利用可能な場合に使用されるデフォルトのポリシーです。
- **hw:cpu_thread_policy=isolate** を使用する場合は、SMT を無効にするか、SMT をサポートしないプラットフォームを使用する必要があります。

検証

1. フレーバーにより専用の CPU を持つインスタンスが作成されることを確認するには、新しいフレーバーを使用してインスタンスを起動します。

```
(overcloud)$ openstack server create --flavor pinned_cpus \
--image <image> pinned_cpu_instance
```

2. 新規インスタンスが正しく配置されていることを確認するには、以下のコマンドを入力し、その出力で **OS-EXT-SRV-ATTR:hypervisor_hostname** の箇所を確認します。

```
(overcloud)$ openstack server show pinned_cpu_instance
```

3.1.5. インスタンス用の共有 CPU フレーバーの作成

クラウドユーザーが共有の (あるいはフローティング) CPU を使用するインスタンスを作成できるようにするには、インスタンス起動用の共有 CPU ポリシーが設定されたフレーバーを作成します。

前提条件

- コンピュートノードが共有 CPU 用に物理 CPU コアを確保するように設定されている。詳しくは、[コンピュートノードでの CPU ピニングの設定](#) を参照してください。

手順

1. source コマンドで **overcloudrc** ファイルを読み込みます。

```
(undercloud)$ source ~/overcloudrc
```

2. CPU ピニングを要求しないインスタンス用のフレーバーを作成します。

```
(overcloud)$ openstack flavor create --ram <size_mb> \
--disk <size_gb> --vcpus <no_reserved_vcpus> floating_cpus
```

3. フローティング CPU を要求するには、フレーバーの **hw:cpu_policy** 属性を **shared** に設定します。

```
(overcloud)$ openstack flavor set \
--property hw:cpu_policy=shared floating_cpus
```

検証

1. フレーバーにより共有 CPU を使用するインスタンスが作成されることを確認するには、新しいフレーバーを使用してインスタンスを起動します。

```
(overcloud)$ openstack server create --flavor floating_cpus \
--image <image> floating_cpu_instance
```

2. 新規インスタンスが正しく配置されていることを確認するには、以下のコマンドを入力し、その出力で **OS-EXT-SRV-ATTR:hypervisor_hostname** の箇所を確認します。

```
(overcloud)$ openstack server show floating_cpu_instance
```

3.1.6. 同時マルチスレッド (SMT) 対応のコンピュートノードでの CPU ピニングの設定

コンピュートノードが同時マルチスレッド (SMT) をサポートする場合、スレッドシブリングを専用または共有セットのいずれかにグルーピングします。スレッドシブリングは共通のハードウェアを共有するため、あるスレッドシブリング上で動作しているプロセスが、他のスレッドシブリングのパフォーマンスに影響を与える可能性があります。

たとえば、ホストは、SMT 対応のデュアルコア CPU に 4 つの論理 CPU コア (0、1、2、および 3) を認識します。この 4 つの CPU に対して、スレッドシブリングのペアが 2 つあります。

- スレッドシブリング 1: 論理 CPU コア 0 および 2
- スレッドシブリング 2: 論理 CPU コア 1 および 3

このシナリオでは、論理 CPU コア 0 および 1 を専用として、2 および 3 を共有として割り当てないでください。そうではなく、0 および 2 を専用として、1 および 3 を共有として割り当てます。

`/sys/devices/system/cpu/cpuN/topology/thread_siblings_list` のファイル。N は論理 CPU 番号で、スレッドペアが含まれます。以下のコマンドを使用して、スレッドシブリングである論理 CPU コアを特定できます。

```
# grep -H . /sys/devices/system/cpu/cpu*/topology/thread_siblings_list | sort -n -t ':' -k 2 -u
```

以下の出力は、論理 CPU コア 0 と論理 CPU コア 2 が同じコア上のスレッドであることを示しています。

```
/sys/devices/system/cpu/cpu0/topology/thread_siblings_list:0,2
/sys/devices/system/cpu/cpu2/topology/thread_siblings_list:1,3
```

3.1.7. 関連情報

- Network Functions Virtualization Planning and Configuration ガイドの [NUMA ノードトポロジーの検出](#)。
- Network Functions Virtualization Product ガイドの [CPU と NUMA ノード](#)。

3.2. エミュレータースレッドの設定

コンピュートノードには、エミュレータースレッドと呼ばれる各インスタンスのハイパーバイザーとリンクするオーバーヘッドタスクがあります。デフォルトでは、エミュレータースレッドはインスタンスと同じ CPU で実行され、インスタンスのパフォーマンスに影響を及ぼします。

エミュレータースレッドポリシーを設定して、インスタンスが使用する CPU とは別の CPU でエミュレータースレッドを実行することができます。



注記

パケットロスを避けるために、NFV デプロイメントでは絶対に仮想 CPU のプリエンブションを行わないでください。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。

2. Compute 環境ファイルを開きます。
3. CPU ピニングを必要とするインスタンス用に物理 CPU コアを確保するには、Compute 環境ファイルで **NovaComputeCpuDedicatedSet** パラメーターを設定します。たとえば、以下の設定では、32 コア CPU を持つコンピュートノードに専用の CPU を設定します。

```
parameter_defaults:
...
NovaComputeCpuDedicatedSet: 2-15,18-31
...
```

詳しくは、[コンピュートノードでの CPU ピニングの設定](#) を参照してください。

4. エミュレータスレッド用に物理 CPU コアを確保するには、Compute 環境ファイルで **NovaComputeCpuSharedSet** パラメーターを設定します。たとえば、以下の設定では、32 コア CPU を持つコンピュートノードに共有の CPU を設定します。

```
parameter_defaults:
...
NovaComputeCpuSharedSet: 0,1,16,17
...
```



注記

Compute スケジューラーは、共有 (またはフローティング) CPU 上で動作するインスタンス用にも共有セット内の CPU を使用します。詳しくは、[コンピュートノードでの CPU ピニングの設定](#) を参照してください。

5. **NovaSchedulerEnabledFilters** パラメーターにまだ **NUMATopologyFilter** がなければ、この Compute スケジュールフィルターを追加します。
6. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

7. インスタンスのエミュレータスレッドを **NovaComputeCpuSharedSet** を使用して設定した共有 CPU から選択した専用の CPU 上で実行するフレーバーを設定します。

```
(overcloud)$ openstack flavor set --property hw:cpu_policy=dedicated \
--property hw:emulator_threads_policy=share \
dedicated_emulator_threads
```

hw:emulator_threads_policy の設定オプションについての詳しい情報は、[Flavor metadata の Emulator threads policy](#) を参照してください。

3.3. コンピュートノードでのヒュージページの設定

クラウド管理者は、インスタンスがヒュージページを要求できるようにコンピュートノードを設定することができます。

手順

1. Compute 環境ファイルを開きます。
2. インスタンスではないプロセス用に各 NUMA ノードで確保するヒュージページのメモリー量を設定します。

```
parameter_defaults:
  ComputeParameters:
    NovaReservedHugePages: ["node:0,size:1GB,count:1","node:1,size:1GB,count:1"]
```

- 各ノードの **size** の値を、割り当てられたヒュージページのサイズに置き換えます。以下の有効な値のいずれかに設定します。
 - 2048 (2 MB 用)
 - 1GB
- 各ノードの **count** の値を、NUMA ノードごとに OVS が使用するヒュージページの数に置き換えます。たとえば、Open vSwitch が 4096 のソケットメモリーを使用する場合、この属性を 2 に設定します。

3. コンピュートノードで Huge Page を設定します。

```
parameter_defaults:
  ComputeParameters:
    ...
    KernelArgs: "default_hugepagesz=1GB hugepagesz=1G hugepages=32"
```



注記

複数の Huge Page サイズを設定する場合、初回の起動時に Huge Page フォルダーもマウントする必要があります。詳細は、[Mounting multiple huge page folders during first boot](#) を参照してください。

4. (オプション) インスタンスが 1GB のヒュージページを割り当ててるのを許可するには、CPU 機能フラグ **NovaLibvirtCPUModelExtraFlags** を設定して **pdpe1gb** を指定します。

```
parameter_defaults:
  ComputeParameters:
    NovaLibvirtCPUMode: 'custom'
    NovaLibvirtCPUModels: 'Haswell-noTSX'
    NovaLibvirtCPUModelExtraFlags: 'vmx, pdpe1gb'
```



注記

- インスタンスが 2 MB のヒュージページしか要求しない場合、CPU 機能フラグを設定する必要はありません。
- ホストが 1GB ヒュージページの割り当てをサポートする場合に限り、インスタンスに 1GB のヒュージページを割り当てることができます。
- **NovaLibvirtCPUMode** が **host-model** または **custom** に設定されている場合にのみ、**NovaLibvirtCPUModelExtraFlags** を **pdpe1gb** に設定する必要があります。
- ホストが **pdpe1gb** をサポートし、**NovaLibvirtCPUMode** に **host-passthrough** が使用される場合、**NovaLibvirtCPUModelExtraFlags** に **pdpe1gb** を設定する必要はありません。**pdpe1gb** フラグは Opteron_G4 および Opteron_G5 CPU モデルにのみ含まれ、QEMU がサポートする Intel CPU モデルには含まれません。
- Microarchitectural Data Sampling (MDS) 等の CPU ハードウェアの問題を軽減するには、他の CPU フラグを設定しなければならない場合があります。詳しくは、[RHOS Mitigation for MDS \("Microarchitectural Data Sampling"\) Security Flaws](#) を参照してください。

5. Meltdown に対する保護の適用後にパフォーマンスが失われないようにするには、CPU 機能フラグ **NovaLibvirtCPUModelExtraFlags** を設定して **+pcid** を指定します。

```
parameter_defaults:
  ComputeParameters:
    NovaLibvirtCPUMode: 'custom'
    NovaLibvirtCPUModels: 'Haswell-noTSX'
    NovaLibvirtCPUModelExtraFlags: 'vmx, pdpe1gb, +pcid'
```

ヒント

詳しくは、[Reducing the performance impact of Meltdown CVE fixes for OpenStack guests with "PCID" CPU feature flag](#) を参照してください。

6. **NovaSchedulerEnabledFilters** パラメーターにまだ **NUMATopologyFilter** がなければ、このフィルターを追加します。
7. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
  -e [your environment files] \
  -e /home/stack/templates/<compute_environment_file>.yaml
```

3.3.1. インスタンス用のヒュージページフレイバーの作成

クラウドユーザーがヒュージページを使用するインスタンスを作成できるようにするには、インスタンス起動用の **hw:mem_page_size** 追加スペックキーが設定されたフレイバーを作成します。

前提条件

- コンピュートノードがヒュージページに対応する設定である。詳しくは、[Configuring huge pages on Compute nodes](#) を参照してください。

手順

1. ヒュージページを要求するインスタンス用のフレーバーを作成します。

```
$ openstack flavor create --ram <size_mb> --disk <size_gb> \
--vcpus <no_reserved_vcpus> huge_pages
```

2. ヒュージページを要求するには、フレーバーの **hw:mem_page_size** 属性を必要なサイズに設定します。

```
$ openstack flavor set huge_pages --property hw:mem_page_size=1GB
```

hw:mem_page_size を以下の有効な値のいずれかに設定します。

- **large**: ホストでサポートされる最大のページサイズを選択します。x86_64 システムでは 2 MB または 1 GB です。
 - **small**: (デフォルト) ホストでサポートされる最小のページサイズを選択します。X86_64 システムでは、4 kB (通常のページ) です。
 - **any**: libvirt ドライバーで決定される、利用可能な最大のヒュージページサイズを選択します。
 - **<pagesize>**: (文字列) ワークロードに具体的な要件がある場合、ページサイズを明示的に設定します。ページサイズには整数値を使用し、kB またはその他の標準単位で指定します。(例: 4kB、2MB、2048、1GB)。
3. フレーバーによりヒュージページを使用するインスタンスが作成されることを確認するには、新しいフレーバーを使用してインスタンスを起動します。

```
$ openstack server create --flavor huge_pages \
--image <image> huge_pages_instance
```

Compute スケジューラーは、インスタンスのメモリーをサポートするのに十分なサイズの空きヒュージページを持つホストを特定します。スケジューラーが十分なページを持つホストおよび NUMA ノードを検出できない場合、リクエストは失敗して **NoValidHost** エラーが報告されます。

3.3.2. 初回起動時に複数の Huge Page フォルダをマウント

Compute サービス (nova) を設定して、初回起動プロセスの一環として複数のページサイズを処理することができます。初回起動プロセスでは、初めてノードを起動する際に heat テンプレート設定をすべてのノードに追加します。これ以降は (たとえば、オーバークラウドスタックの更新時)、これらのテンプレートを追加してもこれらのスクリプトは実行されません。

手順

1. Huge Page フォルダのマウントを作成するためにスクリプトを実行する初回起動テンプレートファイル **hugepages.yaml** を作成します。**OS::TripleO::MultipartMime** リソースタイプを使用して、この設定スクリプトを送信することができます。

```
heat_template_version: <version>
```

```

description: >
  Huge pages configuration

resources:
  userdata:
    type: OS::Heat::MultipartMime
    properties:
      parts:
        - config: {get_resource: hugepages_config}

  hugepages_config:
    type: OS::Heat::SoftwareConfig
    properties:
      config: |
        #!/bin/bash
        hostname | grep -qiE 'co?mp' || exit 0
        systemctl mask dev-hugepages.mount || true
        for pagesize in 2M 1G;do
          if ! [ -d "/dev/hugepages${pagesize}" ]; then
            mkdir -p "/dev/hugepages${pagesize}"
            cat << EOF > /etc/systemd/system/dev-hugepages${pagesize}.mount
[Unit]
Description=${pagesize} Huge Pages File System
Documentation=https://www.kernel.org/doc/Documentation/vm/hugetlbpage.txt
Documentation=https://www.freedesktop.org/wiki/Software/systemd/APIFileSystems
DefaultDependencies=no
Before=sysinit.target
ConditionPathExists=/sys/kernel/mm/hugepages
ConditionCapability=CAP_SYS_ADMIN
ConditionVirtualization=!private-users

[Mount]
What=hugetlbfs
Where=/dev/hugepages${pagesize}
Type=hugetlbfs
Options=pagesize=${pagesize}

[Install]
WantedBy = sysinit.target
EOF
          fi
        done
        systemctl daemon-reload
        for pagesize in 2M 1G;do
          systemctl enable --now dev-hugepages${pagesize}.mount
        done

  outputs:
    OS::stack_id:
      value: {get_resource: userdata}

```

このテンプレートの **config** スクリプトは、以下のタスクを実行します。

- a. **'co?mp'** に一致するホスト名を指定することで、ホストをフィルターリンクして、huge page フォルダのマウントを作成します。必要に応じて、特定のコンピュートの filter grep パターンを更新できます。
 - b. デフォルトの **dev-hugepages.mount systemd** ユニットファイルをマスクし、ページサイズを使用して新規マウントを作成できるようにします。
 - c. フォルダが最初に作成されていることを確認します。
 - d. **pagesize** ごとに **systemd** マウントユニットを作成します。
 - e. 最初のループ後に **systemd daemon-reload** を実行し、新たに作成されたユニットファイルを追加します。
 - f. 2M および 1G のページサイズのマウントを有効にします。このループを更新して、必要に応じて追加のページサイズを含めることができます。
2. オプション: **/dev** フォルダは、自動的に **nova_compute** および **nova_libvirt** のコンテナにバインドマウントされます。Huge Page マウントに異なる宛先を使用した場合は、そのマウントを **nova_compute** および **nova_libvirt** コンテナに渡す必要があります。

```
parameter_defaults
NovaComputeOptVolumes:
- /opt/dev:/opt/dev
NovaLibvirtOptVolumes:
- /opt/dev:/opt/dev
```

3. **~/templates/firstboot.yaml** 環境ファイルの **OS::TripleO::NodeUserData** リソースタイプとして heat テンプレートを登録します。

```
resource_registry:
OS::TripleO::NodeUserData: ./hugepages.yaml
```



重要

NodeUserData リソースを登録することができるのは、1つのリソースにつき1つの heat テンプレートだけです。別の heat テンプレートに登録すると、使用する heat テンプレートがそのテンプレートに変わります。

4. 初回起動環境ファイルを他の環境ファイルと一緒にスタックに追加し、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/firstboot.yaml \
...
```

3.4. インスタンスにファイルベースのメモリーを使用するコンピュートノードの設定

ファイルベースのメモリーを使用して、コンピュートノードのメモリー容量を拡張することができます。この場合、libvirt メモリーバックングディレクトリ内のファイルをインスタンスのメモリーとして割り当てます。インスタンスのメモリーとして使用できるホストディスクの容量、およびインスタンスメモリーファイルのディスク上の場所を設定することができます。

Compute サービスは、ファイルベースのメモリーに設定された容量をシステムメモリーの合計容量として Placement サービスに報告します。これにより、コンピュートノードは通常システムメモリー内で対応できるよりも多くのインスタンスをホストすることができます。

インスタンスにファイルベースのメモリーを使用するには、コンピュートノードでファイルベースのメモリーを有効にする必要があります。

制限

- ファイルベースのメモリーが有効なコンピュートノードとファイルベースのメモリーが有効ではないコンピュートノード間で、インスタンスのライブマイグレーションを行うことはできません。
- ファイルベースのメモリーとヒュージページとの間に互換性はありません。ファイルベースのメモリーが有効なコンピュートノード上で、ヒュージページを使用するインスタンスを起動することはできません。ホストアグリゲートを使用して、ヒュージページを使用するインスタンスがファイルベースのメモリーが有効なコンピュートノードに配置されないようにします。
- ファイルベースのメモリーとメモリーのオーバーコミットとの間に互換性はありません。
- **NovaReservedHostMemory** を使用してホストのプロセス用にメモリーを確保することはできません。ファイルベースのメモリーが使用されている場合、確保されるメモリーはファイルベースのメモリー用に用意されないディスク領域を表します。ファイルベースのメモリーは、キャッシュメモリーとして使用される RAM と共に、システムメモリーの合計として Placement サービスに報告されます。

前提条件

- ノードおよびノードが追加されたすべてのホストアグリゲートで、**NovaRAMAllocationRatio** を 1.0 に設定する必要があります。
- **NovaReservedHostMemory** を 0 に設定する必要があります。

手順

1. Compute 環境ファイルを開きます。
2. Compute 環境ファイルに以下のパラメーターを追加して、インスタンスの RAM 用に利用可能なホストディスク容量を MiB 単位で設定します。

```
parameter_defaults:
  NovaLibvirtFileBackedMemory: 102400
```

3. オプション: メモリーバックギングファイルを保存するディレクトリーを設定するには、Compute 環境ファイルに **QemuMemoryBackingDir** パラメーターを設定します。設定しなければ、メモリーバックギングディレクトリーはデフォルトの **/var/lib/libvirt/qemu/ram/** に設定されます。



注記

デフォルトのディレクトリー (**/var/lib/libvirt/qemu/ram/**) またはそれより上の階層のディレクトリーに、バックギングストアを配置する必要があります。

バックギングストアのホストディスクを変更することもできます。詳細は、[Changing the memory backing directory host disk](#) を参照してください。

4. 更新内容を Compute 環境ファイルに保存します。
5. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
  -e [your environment files] \
  -e /home/stack/templates/<compute_environment_file>.yaml
```

3.4.1. メモリーバックングディレクトリーのホストディスクの変更

デフォルトのプライマリーディスクから別のディスクに、メモリーバックングディレクトリーを変更することができます。

手順

1. 代替のバックングデバイスにファイルシステムを作成します。たとえば、**/dev/sdb** に **ext4** ファイルシステムを作成するには、以下のコマンドを入力します。

```
# mkfs.ext4 /dev/sdb
```

2. バックングデバイスをマウントします。たとえば、**/dev/sdb** をデフォルトの libvirt メモリーバックングディレクトリーにマウントするには、以下のコマンドを入力します。

```
# mount /dev/sdb /var/lib/libvirt/qemu/ram
```



注記

マウントポイントは、**QemuMemoryBackingDir** パラメーターの値と一致する必要があります。

第4章 COMPUTE サービスのストレージの設定

Compute サービスが Image (glance) サービスからコピーしてコンピュートノード上でローカルにキャッシュしたベースイメージから、インスタンスを作成します。インスタンスのバックエンドであるインスタンスディスクも、ベースイメージに基づいています。

Compute サービスを設定して、一時インスタンスのディスクデータをホストのコンピュートノード上にローカルに保存するか、または NFS 共有または Ceph クラスタでリモートで保存することができます。あるいは、Compute サービスを設定して、Block Storage (Cinder) サービスが提供する永続ストレージにインスタンスディスクデータを保存することもできます。


環境のイメージキャッシュを設定し、インスタンスディスクのパフォーマンスおよびセキュリティーを設定することができます。Image サービス (glance) がバックエンドとして Red Hat Ceph RADOS Block Device (RBD) を使用する場合に、Image サービス API を使用せずに RBD イメージリポジトリから直接イメージをダウンロードするように Compute サービスを設定することもできます。

4.1. イメージキャッシュの設定オプション

以下の表で詳細を説明するパラメーターを使用して、Compute サービスがどのようにコンピュートノードでのイメージのキャッシュを実装して管理するかを設定します。

表4.1 Compute (nova) サービスのイメージキャッシュパラメーター

設定方法	パラメーター	説明
Puppet	nova::compute::image_cache::manager_interval	<p>コンピュートノードでのベースイメージのキャッシュを管理するイメージキャッシュマネージャーの実行間で待機する時間を秒単位で指定します。 nova::compute::image_cache::remove_unused_base_images が True に設定されている場合、Compute サービスはこの期間を使用して使用されていないキャッシュイメージの自動削除を実行します。</p> <p>デフォルトのメトリック間隔である 60 秒で実行するには、0 に設定します (推奨されません)。イメージキャッシュマネージャーを無効にするには、-1 に設定します。</p> <p>デフォルト: 2400</p>

設定方法	パラメーター	説明
Puppet	nova::compute::image_cache::precache_concurrency	<p>イメージを並行して事前キャッシュできるコンピュートノードの数の最大値を指定します。</p> <div>  <div> <p>注記</p> <ul style="list-style-type: none"> このパラメーターに大きな数値を設定すると、事前キャッシュの処理が遅くなり、Image サービスで DDoS が発生する場合があります。 このパラメーターに小さな数値を設定すると、Image サービスの負荷は軽減されますが、事前キャッシュがより連続的な操作で実行されるため、完了までの実行時間が長くなる場合があります。 </div> </div> <p>デフォルト: 1</p>
Puppet	nova::compute::image_cache::remove_unused_base_images	<p>manager_interval を使用して設定された間隔で使用されていないベースイメージをキャッシュから自動的に削除するには、True に設定します。NovalImageCacheTTL を使用して指定された期間アクセスされなかった場合、イメージは使用されていないと定義されます。</p> <p>デフォルト: True</p>
Puppet	nova::compute::image_cache::remove_unused_resized_minimum_age_seconds	<p>未使用のサイズ変更されたベースイメージをキャッシュから削除する最短の期間を指定します (秒単位)。未使用のサイズ変更されたベースイメージは、この期間削除されません。無効にする場合は undef に設定します。</p> <p>デフォルト: 3600</p>

設定方法	パラメーター	説明
Puppet	nova::compute::image_cache::subdirectory_name	<p>キャッシュされたイメージを保存するフォルダーの名前を、\$instances_path との相対パスで指定します。</p> <p>デフォルト: _base</p>
heat	NovalImageCacheTTL	<p>コンピュートノード上のどのインスタンスもイメージを使用しなくなった場合に、Compute サービスがイメージをキャッシュし続ける期間を秒単位で指定します。Compute サービスは、コンピュートノードにキャッシュされたイメージのうち、ここで設定したライフタイムを過ぎたイメージを、再度必要になるまでキャッシュディレクトリーから削除します。</p> <p>デフォルト: 86400 (24 時間)</p>

4.2. インスタンスの一時ストレージ属性の設定オプション

以下の表で詳細を説明するパラメーターを使用して、インスタンスが使用する一時ストレージのパフォーマンスおよびセキュリティーを設定します。



注記


Red Hat OpenStack Platform (RHOSP) は、インスタンスディスクの LVM イメージ種別をサポートしません。したがって、インスタンスの削除時に一時ディスクを消去する **[libvirt]/volume_clear** 設定オプションは、インスタンスディスクイメージ種別が LVM の場合にのみ適用されるため、サポートされません。

表4.2 Compute (nova) サービスのインスタンス一時ストレージパラメーター

設定方法	パラメーター	説明
Puppet	nova::compute::default_ephemeral_format	<p>新規の一時ボリュームに使用されるデフォルトの形式を指定します。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● ext2 ● ext3 ● ext4 <p>ext4 形式では、ext3 に比べ、サイズの大きい新規ディスクを初期化する時間が大幅に短縮されます。</p> <p>デフォルト: ext4</p>

設定方法	パラメーター	説明
Puppet	nova::compute::force_raw_images	<p>raw 以外の形式でキャッシュされたベースイメージを raw 形式に変換するには、True に設定します。raw イメージ形式は、qcow2 等の他のイメージ形式よりも多くの領域を使用します。raw 以外のイメージ形式は、圧縮により多くの CPU パワーを使用します。False に設定すると、CPU のボトルネックを防ぐために、Compute サービスは圧縮時にベースイメージからすべての圧縮を削除します。システムの I/O 処理が遅い場合や空き容量が少ない場合は、入力の帯域幅を削減するために False に設定します。</p> <p>デフォルト: True</p>
Puppet	nova::compute::use_cow_images	<p>インスタンスのディスクに qcow2 形式の CoW (Copy on Write) イメージを使用するには、True に設定します。CoW を使用する場合には、バックングストアとホストのキャッシュによっては、各インスタンスが独自のコピー上で稼働することで、並行処理が改善される場合があります。</p> <p>raw 形式を使用するには、False に設定します。raw 形式は、ディスクイメージの共通の部分により多くの領域を使用します。</p> <p>デフォルト: True</p>
Puppet	nova::compute::libvirt::preallocate_images	<p>インスタンスディスクに対する事前割り当てモードを指定します。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● none: インスタンスの起動時にはストレージがプロビジョニングされません。 ● space: インスタンスのディスクイメージで fallocate(1) を実行することで、Compute サービスはインスタンスの起動時にストレージを完全に割り当てます。これにより CPU のオーバーヘッドおよびファイルの断片化が軽減され、I/O パフォーマンスが向上し、要求されたディスク領域を確保するのに役立ちます。 <p>デフォルト: none</p>

設定方法	パラメーター	説明
hieradata override	DEFAULT/resize_fs_using_block_device	<p>ブロックデバイスを使用してイメージにアクセスすることで、ベースイメージのサイズを直接変更できるようにするには、True に設定します。これは、cloud-init のバージョンがより古いイメージ (それ自体ではサイズの変更ができない) でのみ必要です。</p> <p>このパラメーターによりイメージの直接マウントが可能になるため (セキュリティ上の理由により無効にされる場合がある)、このパラメーターはデフォルトでは有効化されていません。</p> <p>デフォルト: False</p>

設定方法	パラメーター	説明
hieradata override	[libvirt]/images_type	<p>インスタンスのディスクに使用するイメージ種別を指定します。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● raw ● qcow2 ● flat ● rbd ● default <div>  <div> <p>注記</p> <p>RHOSP は、インスタンスディスクの LVM イメージ種別をサポートしません。</p> </div> </div> <p>default 以外の有効な値に設定すると、イメージ種別は use_cow_images の設定よりも優先されます。default が指定されると、use_cow_images の設定によりイメージ種別が決定されます。</p> <ul style="list-style-type: none"> ● use_cow_images が True(デフォルト) に設定されると、イメージ種別は qcow2 になります。 ● use_cow_images を False に設定すると、イメージ種別は Flat になります。 <p>デフォルト値は NovaEnableRbdBackend の設定により決定されます。</p> <ul style="list-style-type: none"> ● NovaEnableRbdBackend: False デフォルト: default ● NovaEnableRbdBackend: True デフォルト: rbd

4.3. 共有インスタンスストレージの設定

デフォルトでは、インスタンスの起動時に、インスタンスのディスクはインスタンスディレクトリー **/var/lib/nova/instances** にファイルとして保存されます。Compute サービスの NFS ストレージバックエンドを設定して、これらのインスタンスファイルを共有 NFS ストレージに保存することができます。

前提条件

- NFSv4 以降を使用している。Red Hat OpenStack Platform (RHOSP) は、以前のバージョンの NFS をサポートしません。詳細は、Red Hat ナレッジベースのソリューション [RHOS NFSv4-Only Support Notes](#) を参照してください。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. **stackrc** ファイルを取得します。

```
[stack@director ~]$ source ~/stackrc
```

3. 共有インスタンスストレージを設定するための環境ファイルを作成します (例: **nfs_instance_disk_backend.yaml**)。
4. インスタンスファイル用に NFS バックエンドを設定するには、以下の設定を **nfs_instance_disk_backend.yaml** に追加します。

```
parameter_defaults:
...
NovaNfsEnabled: True
NovaNfsShare: <nfs_share>
```

<nfs_share> をインスタンスのファイルストレージ用にマウントする NFS 共有ディレクトリーに置き換えます (例: **'192.168.122.1:/export/nova'** または **'192.168.24.1:/var/nfs'**)。IPv6 を使用している場合は、二重と単一引用符の両方を使用してください (例: **""[fdd0::1]:/export/nova""**)。

5. オプション: NFS バックエンドストレージが有効な場合、NFS ストレージ用のデフォルトのマウント SELinux コンテキストは **'context=system_u:object_r:nfs_t:s0'** です。以下のパラメーターを追加して、NFS インスタンスファイルストレージのマウントポイントのマウントオプションを変更します。

```
parameter_defaults:
...
NovaNfsOptions: 'context=system_u:object_r:nfs_t:s0,<additional_nfs_mount_options>'
```

<additional_nfs_mount_options> を、NFS インスタンスファイルストレージに使用するマウントオプションのコンマ区切りリストに置き換えます。利用可能なマウントオプションの詳細は、**mount** の man ページを参照してください。

```
$ man 8 mount.
```

6. 更新を環境ファイルに保存します。
7. その他の環境ファイルと共に新しい環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/nfs_instance_disk_backend.yaml
```

4.4. RED HAT CEPH RADOS BLOCK DEVICE (RBD) からの直接イメージダウンロードの設定

Image サービス (glance) がバックエンドとして Red Hat Ceph RADOS Block Device (RBD) を使用し、Compute サービスがローカルのファイルベースの一時ストレージを使用する場合、Image サービス API を使用せずに RBD イメージリポジトリから直接イメージをダウンロードするように Compute サービスを設定することができます。これにより、インスタンスのブート時にコンピュートノードイメージキャッシュにイメージをダウンロードする時間が短縮されます。これにより、インスタンスの起動時間が短縮されます。

前提条件

- Image サービスのバックエンドが、Red Hat Ceph RADOS Block Device (RBD) である。
- Compute サービスが、イメージキャッシュおよびインスタンスのディスクにローカルのファイルベースの一時ストアを使用している。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. Compute 環境ファイルを開きます。
3. RBD バックエンドから直接イメージをダウンロードするには、以下の設定を Compute 環境ファイルに追加します。

```
parameter_defaults:
  ComputeParameters:
    NovaGlanceEnableRbdDownload: True
    NovaEnableRbdBackend: False
  ...
```

4. (オプション) Image サービスが複数の Red Hat Ceph Storage バックエンドを使用するように設定されている場合には、Compute 環境ファイルに以下の設定を追加して、イメージをダウンロードする RBD バックエンドを特定します。

```
parameter_defaults:
  ComputeParameters:
    NovaGlanceEnableRbdDownload: True
    NovaEnableRbdBackend: False
    NovaGlanceRbdDownloadMultistoreID: <rbd_backend_id>
  ...
```

<rbd_backend_id> を **GlanceMultistoreConfig** 設定のバックエンドを指定するために使用される ID(例: **rbd2_store**) に置き換えます。

5. 以下の設定を Compute 環境ファイルに追加して Image サービス RBD バックエンドを指定し、Compute サービスが Image サービス RBD バックエンドへの接続を待機する最大期間 (秒単位) を指定します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::config::nova_config:
      glance/rbd_user:
        value: 'glance'
```

```
glance/rbd_pool:
  value: 'images'
glance/rbd_ceph_conf:
  value: '/etc/ceph/ceph.conf'
glance/rbd_connect_timeout:
  value: '5'
```

6. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

7. Compute サービスが RBD から直接イメージをダウンロードすることを確認するには、インスタンスを作成してインスタンスのデバッグログで Attempting to export RBD image: のエントリーを確認します。

4.5. 関連情報

- [Compute サービス \(nova\) の設定](#)

第5章 PCI パススルーの設定

PCI パススルーを使用して、グラフィックカードまたはネットワークデバイス等の物理 PCI デバイスをインスタンスにアタッチすることができます。デバイスに PCI パススルーを使用する場合、インスタンスはタスクを実行するためにデバイスへの排他的アクセスを確保し、ホストはデバイスを利用することができません。

重要

ルーティング対応プロバイダーネットワークでの PCI パススルーの使用

Compute サービスは、複数のプロバイダーネットワークにまたがる単一のネットワークをサポートしません。ネットワークに複数の物理ネットワークが含まれる場合、Compute サービスは最初の物理ネットワークだけを使用します。したがって、ルーティング対応プロバイダーネットワークを使用する場合は、すべてのコンピュートノードで同じ **physical_network** 名を使用する必要があります。

VLAN またはフラットネットワークのルーティング対応プロバイダーネットワークを使用する場合は、すべてのセグメントで同じ **physical_network** 名を使用する必要があります。その後、ネットワークに複数のセグメントを作成し、そのセグメントを適切なサブネットにマッピングします。

クラウドユーザーが PCI デバイスがアタッチされたインスタンスを作成できるようにするには、以下の手順を実施する必要があります。

1. PCI パススルー用のコンピュートノードを指定する。
2. 必要な PCI デバイスを持つ PCI パススルー用のコンピュートノードを設定する。
3. オーバークラウドをデプロイする。
4. PCI デバイスがアタッチされたインスタンスを起動するためのフレーバーを作成する。

前提条件

- 必要な PCI デバイスを持つコンピュートノード

5.1. PCI パススルー用コンピュートノードの指定

物理 PCI デバイスが接続されたインスタンスのコンピュートノードを指定するには、新しいロールファイルを作成して PCI パススルーロールを設定し、PCI パススルーのコンピュートノードにタグを付けるために使用する PCI パススルーリソースクラスを使用してベアメタルノードを設定する必要があります。

注記

以下の手順は、まだプロビジョニングされていない新しいオーバークラウドノードに適用されます。すでにプロビジョニングされている既存のオーバークラウドノードにリソースクラスを割り当てるには、スケールダウン手順を使用してノードのプロビジョニングを解除してから、スケールアップ手順を使用して新しいリソースクラスの割り当てでノードを再プロビジョニングする必要があります。詳細は、[オーバークラウドノードのスケールリング](#)を参照してください。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. **stackrc** ファイルを取得します。

```
[stack@director ~]$ source ~/stackrc
```

3. **Controller**、**Compute**、および **ComputePCI** ロールを含む、**roles_data_pci_passthrough.yaml** という名前の新しいロールデータファイルを、オーバークラウドに必要なその他のロールとともに生成します。

```
(undercloud)$ openstack overcloud roles \
generate -o /home/stack/templates/roles_data_pci_passthrough.yaml \
Compute:ComputePCI Compute Controller
```

4. **roles_data_pci_passthrough.yaml** を開き、以下のパラメーターおよびセクションを編集または追加します。

セクション/パラメーター	現在の値	新しい値
ロールのコメント	Role: Compute	Role: ComputePCI
ロール名	Compute	name: ComputePCI
description	Basic Compute Node role	PCI パススルー用コンピュートノードロール
HostnameFormatDefault	%stackname%-novacompute-%index%	%stackname%-novacomputepci-%index%
deprecated_nic_config_name	compute.yaml	compute-pci-passthrough.yaml

5. オーバークラウドの PCI パススルー用コンピュートノードをノード定義のテンプレート **node.json** または **node.yaml** に追加して、そのノードを登録します。詳細は、**director** のインストールと使用方法の [オーバークラウドノードの登録](#) を参照してください。
6. ノードのハードウェアを検査します。

```
(undercloud)$ openstack overcloud node introspect \
--all-manageable --provide
```

詳細については、**Director のインストールと使用** ガイドの [ベアメタルノードハードウェアのインベントリーの作成](#) を参照してください。

7. PCI パススルー用に指定する各ベアメタルノードに、カスタムの PCI パススルーリソースクラスをタグ付けします。

```
(undercloud)$ openstack baremetal node set \
--resource-class baremetal.PCI-PASSTHROUGH <node>
```

<node> をベアメタルノードの ID に置き換えてください。

8. ノード定義ファイル **overcloud-baremetal-deploy.yaml** に **ComputePCI** ロールを追加し、予測ノード配置、リソースクラス、ネットワークトポロジー、またはノードに割り当てたいその他の属性を定義します。

```
- name: Controller
  count: 3
- name: Compute
  count: 3
- name: ComputePCI
  count: 1
  defaults:
    resource_class: baremetal.PCI-PASSTHROUGH
    network_config:
      template: /home/stack/templates/nic-config/myRoleTopology.j2 ❶
```

- ❶ 既存のネットワークトポロジーを再利用するか、ロール用の新しいカスタムネットワークインターフェーステンプレートを作成できます。詳細は、**Director のインストールと使用ガイド**の [カスタムネットワークインターフェーステンプレート](#) を参照してください。**network_config** プロパティを使用してネットワーク定義を定義しない場合、デフォルトのネットワーク定義が使用されます。

ノード定義ファイルでノード属性を設定するために使用できるプロパティについて詳しくは、[ベアメタルノードのプロビジョニング属性](#) を参照してください。ノード定義ファイルの例は、[ノード定義ファイルの例](#) を参照してください。

9. プロビジョニングコマンドを実行して、ロールの新しいノードをプロビジョニングします。

```
(undercloud)$ openstack overcloud node provision \
--stack <stack> \
[--network-config \]
--output /home/stack/templates/overcloud-baremetal-deployed.yaml \
/home/stack/templates/overcloud-baremetal-deploy.yaml
```

- **<stack>** を、ベアメタルノードがプロビジョニングされるスタックの名前に置き換えます。指定しない場合、デフォルトは **overcloud** です。
 - **--network-config** オプションの引数を含めて、**cli-overcloud-node-network-config.yaml** Ansible Playbook にネットワーク定義を提供します。**network_config** プロパティを使用してネットワーク定義を定義しない場合、デフォルトのネットワーク定義が使用されます。
10. 別のターミナルでプロビジョニングの進捗をモニターリングします。プロビジョニングが成功すると、ノードの状態が **available** から **active** に変わります。

```
(undercloud)$ watch openstack baremetal node list
```

11. **--network-config** オプションを指定してプロビジョニングコマンドを実行しなかった場合は、**network-environment.yaml** ファイルで **<Role>NetworkConfigTemplate** パラメーターを設定して、NIC テンプレートファイルを指すようにします。

```
parameter_defaults:
  ComputeNetworkConfigTemplate: /home/stack/templates/nic-configs/compute.j2
  ComputePCINetworkConfigTemplate: /home/stack/templates/nic-
```

```
configs/<pci_passthrough_net_top>.j2
ControllerNetworkConfigTemplate: /home/stack/templates/nic-configs/controller.j2
```

<pci_passthrough_net_top> を **ComputePCI** ロールのネットワークトポロジを含むファイルの名前に置き換えます。たとえば、デフォルトのネットワークトポロジを使用するには、**compute.yaml** のようにします。

5.2. PCI パススルー用コンピュートノードの設定

クラウドユーザーが PCI デバイスがアタッチされたインスタンスを作成できるようにするには、PCI デバイスを持つコンピュートノードとコントローラーノードの両方を設定する必要があります。

手順

1. PCI パススルー用にオーバークラウド上のコントローラーノードを設定するには、環境ファイル (例: **pci_passthru_controller.yaml**) を作成します。
2. **pci_passthrough_controller.yaml** の **NovaSchedulerEnabledFilters** パラメーターに **PciPassthroughFilter** を追加します。

```
parameter_defaults:
  NovaSchedulerEnabledFilters:
    - AvailabilityZoneFilter
    - ComputeFilter
    - ComputeCapabilitiesFilter
    - ImagePropertiesFilter
    - ServerGroupAntiAffinityFilter
    - ServerGroupAffinityFilter
    - PciPassthroughFilter
    - NUMATopologyFilter
```

3. コントローラーノード上のデバイスの PCI エイリアスを指定するには、以下の設定を **pci_passthrough_controller.yaml** に追加します。

```
parameter_defaults:
  ...
  ControllerExtraConfig:
    nova::pci::aliases:
      - name: "a1"
        product_id: "1572"
        vendor_id: "8086"
        device_type: "type-PF"
```

device_type フィールドの設定に関する詳細は、[PCI パススルーデバイス種別フィールド](#) を参照してください。



注記

nova-api サービスが **Controller** 以外のロールで実行されている場合は、**ControllerExtraConfig** を **<Role>ExtraConfig** の形式でユーザーロールに置き換えます。

4. (オプション): PCI パススルーデバイスにデフォルトの NUMA アフィニティーポリシーを設定するには、ステップ 3 の **numa_policy** 設定に **nova::pci::aliases:** を追加します。

```
parameter_defaults:
...
ControllerExtraConfig:
  nova::pci::aliases:
    - name: "a1"
      product_id: "1572"
      vendor_id: "8086"
      device_type: "type-PF"
      numa_policy: "preferred"
```

5. PCI パススルー用にオーバークラウド上のコンピュートノードを設定するには、環境ファイル (例: **pci_passthrough_compute.yaml**) を作成します。
6. コンピュートノード上のデバイスの利用可能な PCI を指定するには、**vendor_id** および **product_id** オプションを使用して、インスタンスへのパススルーに使用できる PCI デバイスのプールに、一致するすべての PCI デバイスを追加します。たとえば、すべての Intel® Ethernet Controller X710 デバイスをインスタンスへのパススルーに使用できる PCI デバイスのプールに追加するには、以下の設定を **pci_passthrough_compute.yaml** に追加します。

```
parameter_defaults:
...
ComputePCIParameters:
  NovaPCIPassthrough:
    - vendor_id: "8086"
      product_id: "1572"
```

NovaPCIPassthrough の設定方法の詳細は、[NovaPCIPassthrough 設定のガイドライン](#) を参照してください。

7. インスタンスの移行およびサイズ変更の操作を行うために、コンピュートノードの PCI エイリアスのコピーを作成する必要があります。PCI パススルー用コンピュートノード上のデバイスの PCI エイリアスを指定するには、以下の設定を **pci_passthrough_compute.yaml** に追加します。

```
parameter_defaults:
...
ComputePCIExtraConfig:
  nova::pci::aliases:
    - name: "a1"
      product_id: "1572"
      vendor_id: "8086"
      device_type: "type-PF"
```



注記

コンピュートノードのエイリアスは、コントローラーノードのエイリアスと同じでなければなりません。したがって、**pci_passthrough_controller.yaml** の **nova::pci::aliases** に **numa_affinity** を追加した場合は、**pci_passthrough_compute.yaml** の **nova::pci::aliases** にも追加する必要があります。

8. PCI パススルーをサポートするためにコンピュートノードのサーバー BIOS で IOMMU を有効にするには、**pci_passthrough_compute.yaml** に **KernelArgs** パラメーターを追加します。たとえば、Intel IOMMU を有効にするには、以下の **KernelArgs** 設定を使用します。

```
parameter_defaults:
...
ComputePCIParameters:
  KernelArgs: "intel_iommu=on iommu=pt"
```

AMD IOMMU を有効にするには、**KernelArgs** を "**amd_iommu=on iommu=pt**" に設定します。



注記

KernelArgs パラメーターをロールの設定に初めて追加すると、オーバークラウドノードが自動的に再起動されます。必要に応じて、ノードの自動再起動を無効にし、代わりに各オーバークラウドのデプロイ後にノードの再起動を手動で実行できます。詳細は、[KernelArgs を定義するための手動でのノード再起動の設定](#)を参照してください。

9. その他の環境ファイルと共にこれらのカスタム環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-r /home/stack/templates/roles_data_pci_passthrough.yaml \
-e /home/stack/templates/network-environment.yaml \
-e /home/stack/templates/pci_passthrough_controller.yaml \
-e /home/stack/templates/pci_passthrough_compute.yaml \
-e /home/stack/templates/overcloud-baremetal-deployed.yaml \
-e /home/stack/templates/node-info.yaml
```

10. クラウドユーザーが PCI デバイスを要求するのに使用できるフレーバーを作成および設定します。以下の例では、ステップ 7 で定義したエイリアスを使用して、それぞれベンダー ID および製品 ID が **8086** および **1572** の 2 つのデバイスを要求します。

```
(overcloud)$ openstack flavor set \
--property "pci_passthrough:alias"="a1:2" device_passthrough
```

11. (オプション) フレーバーまたはイメージに NUMA アフィニティポリシーの属性キーを追加して、PCI パススルーデバイスのデフォルト NUMA アフィニティポリシーをオーバーライドすることができます。

- フレーバーを使用してデフォルトの NUMA アフィニティポリシーをオーバーライドするには、**hw:pci_numa_affinity_policy** 属性キーを追加します。

```
(overcloud)$ openstack flavor set \
--property "hw:pci_numa_affinity_policy"="required" \
device_passthrough
```

hw:pci_numa_affinity_policy の有効な値についての詳しい情報は、[フレーバーのメタデータ](#)を参照してください。

- イメージを使用してデフォルトの NUMA アフィニティポリシーをオーバーライドするには、**hw_pci_numa_affinity_policy** 属性キーを追加します。

```
(overcloud)$ openstack image set \
  --property hw_pci_numa_affinity_policy=required \
  device_passthrough_image
```



注記

イメージとフレーバーの両方で NUMA アフィニティポリシーを設定する場合には、属性の値が一致している必要があります。フレーバーの設定は、イメージおよびデフォルトの設定よりも優先されます。したがって、イメージの NUMA アフィニティポリシーの設定は、フレーバーで属性が設定されていない場合に限り効果を持ちます。

検証

1. PCI パススルーデバイスを設定してインスタンスを作成します。

```
$ openstack server create --flavor device_passthrough \
  --image <image> --wait test-pci
```

2. クラウドユーザーとしてインスタンスにログインします。詳しくは、[インスタンスへの接続](#) を参照してください。
3. インスタンスが PCI デバイスにアクセスできることを確認するには、インスタンスから以下のコマンドを入力します。

```
$ lspci -nn | grep <device_name>
```

5.3. PCI パススルーデバイス種別フィールド

Compute サービスでは、デバイスが報告する機能に応じて、PCI デバイスは 3 つの種別のいずれかに分類されます。**device_type** フィールドに設定することのできる有効な値を、以下の一覧に示します。

type-PF

デバイスは、SR-IOV をサポートする親またはルートデバイスです。SR-IOV を完全にサポートするデバイスをパススルーするには、このデバイス種別を指定します。

type-VF

デバイスは、SR-IOV をサポートするデバイスの子デバイスです。

type-PCI

デバイスは SR-IOV をサポートしません。**device_type** フィールドが設定されていない場合は、これがデフォルトのデバイス種別です。



注記

コンピュートノードとコントローラーノードの **device_type** を、同じ値に設定する必要があります。

5.4. NOVAPCIPASSTHROUGH 設定のガイドライン

- NIC のデバイス名は変更される可能性があるため、PCI パススルーを設定する場合は **devname** パラメーターを使用しないでください。代わりに、**vendor_id** と **product_id** の方が安定しているため使用するか、NIC の **address** を使用してください。

- 特定の Physical Function (PF) をパススルーするには、PCI アドレスが各デバイスに固有であるので、**address** パラメーターを使用できます。または、**product_id** パラメーターを使用して PF をパススルーすることもできますが、同じ種別の PF が複数ある場合には、PF の **address** も指定する必要があります。
- すべての Virtual Function (VF) をパススルーするには、PCI パススルーに使用する VF の **product_id** および **vendor_id** のみを指定します。NIC の分割に SRIOV を使用し、VF 上で OVS を実行している場合は、VF の **address** も指定する必要があります。
- PF の VF のパススルーだけを設定し、PF そのもののパススルーは設定しない場合は、**address** パラメーターを使用して PF の PCI アドレスを指定し、**product_id** を使用して VF の 製品 ID を指定することができます。

address パラメーターの設定

address パラメーターは、デバイスの PCI アドレスを指定します。String または **dict** マッピングのいずれかを使用して、**address** パラメーターの値を設定できます。

文字列形式

文字列を使用してアドレスを指定する場合は、以下の例のようにワイルドカード (*) を含めることができます。

```
NovaPCIPassthrough:
-
  address: "*:0a:00.*"
  physical_network: physnet1
```

ディクショナリー形式

ディクショナリー形式を使用してアドレスを指定する場合は、以下の例のように正規表現構文を含めることができます。

```
NovaPCIPassthrough:
-
  address:
    domain: ".*"
    bus: "02"
    slot: "01"
    function: "[0-2]"
  physical_network: net1
```



注記

Compute サービスは、**アドレス フィールド**の設定を次の最大値に制限します。

- domain - 0xFFFF
- bus - 0xFF
- slot - 0x1F
- function - 0x7

Compute サービスは、16 ビットアドレスドメインを持つ PCI デバイスをサポートします。Compute サービスは、アドレスドメインが 32 ビットの PCI デバイスを無視します。

第6章 VDPA ポートを使用するインスタンスを有効にするための VDPA コンピュートノードの設定



重要

この機能は、本リリースでは [テクノロジープレビュー](#) として提供しているため、Red Hat では全面的にはサポートしていません。これは、テスト用途にのみご利用いただく機能で、実稼働環境にデプロイすべきではありません。テクノロジープレビュー機能についての詳しい情報は、[対象範囲の詳細](#) を参照してください。

VIRTIO データパスアクセラレーション (VDPA) は、VIRTIO を介したワイヤースピードのデータ転送を提供します。VDPA デバイスは、SR-IOV 仮想機能 (VF) に対する VIRTIO 抽象化を提供します。これにより、インスタンスでベンダー固有のドライバーなしで VF を使用できます。



注記

NIC を VDPA インターフェイスとして使用する場合は、VDPA インターフェイス専用にする必要があります。NIC の物理機能 (PF) を **switchdev** モードで設定し、ハードウェアオフロード OVS を使用して PF を管理する必要があるため、NIC を他の接続に使用することはできません。

クラウドユーザーが VDPA ポートを使用するインスタンスを作成できるようにするには、次のタスクを完了します。

1. オプション: VDPA のコンピュートノードを指定します。
2. 必要な VDPA ドライバーを持つ VDPA のコンピュートノードを設定します。
3. オーバークラウドをデプロイします。

ヒント

VDPA ハードウェアが制限されている場合は、ホストアグリゲートを設定して VDPA コンピュートノードでのスケジューリングを最適化することもできます。VDPA コンピュートノードで VDPA を要求するインスタンスのみをスケジューリングするには、VDPA ハードウェアを備えたコンピュートノードのホストアグリゲートを作成し、VDPA インスタンスのみをホストアグリゲートに配置するようにコンピュートスケジューラーを設定します。詳細は、[ホストアグリゲートの分離による絞り込み](#) および [ホストアグリゲートの作成と管理](#) を参照してください。

前提条件

- コンピュートノードには、必要な VDPA デバイスとドライバーがあります。
- Compute ノードには Mellanox NIC があります。
- オーバークラウドは OVS ハードウェアオフロード用に設定されています。詳細は、[OVS ハードウェアオフロードの設定](#) を参照してください。
- オーバークラウドは ML2/OVN を使用するよう設定されています。

6.1. VDPA 用コンピュートノードの指定

VIRTIO データパスアクセラレーション (VDPA) インターフェイスを要求するインスタンスのコンピュータノードを指定するには、新しいロールファイルを作成して VDPA ロールを設定し、VDPA リソースクラスを使用してベアメタルノードを設定し、VDPA のコンピュータノードにタグを付けます。



注記

以下の手順は、まだプロビジョニングされていない新しいオーバークラウドノードに適用されます。すでにプロビジョニングされている既存のオーバークラウドノードにリソースクラスを割り当てるには、オーバークラウドをスケールダウンしてノードのプロビジョニングを解除してから、オーバークラウドをスケールアップして、新しいリソースクラスの割り当てでノードを再プロビジョニングします。詳細は、[オーバークラウドノードのスケールリング](#)を参照してください。

手順

1. アンダークラウドホストに **stack** ユーザーとしてログインします。
2. **stackrc** アンダークラウド認証情報ファイルを入手します。

```
[stack@director ~]$ source ~/stackrc
```

3. **Controller**、**Compute**、および **ComputeVdpa** ロールを含む、**roles_data_vdpa.yaml** という名前の新しいロールデータファイルを生成します。

```
(undercloud)$ openstack overcloud roles \
generate -o /home/stack/templates/roles_data_vdpa.yaml \
ComputeVdpa Compute Controller
```

4. VDPA ロールの **roles_data_vdpa.yaml** ファイルを更新します。

```
#####
####
# Role: ComputeVdpa                                     #
#####
####
- name: ComputeVdpa
  description: |
    VDPA Compute Node role
  CountDefault: 1
  # Create external Neutron bridge
  tags:
    - compute
    - external_bridge
  networks:
    InternalApi:
      subnet: internal_api_subnet
    Tenant:
      subnet: tenant_subnet
    Storage:
      subnet: storage_subnet
  HostnameFormatDefault: '%stackname%-computevdpa-%index%'
  deprecated_nic_config_name: compute-vdpa.yaml
```

5. オーバークラウド用の VDPA コンピュートノードをノード定義のテンプレート **node.json** または **node.yaml** に追加して、そのノードを登録します。詳細は、**director のインストールと使用方法** の [オーバークラウドノードの登録](#) を参照してください。

6. ノードのハードウェアを検査します。

```
(undercloud)$ openstack overcloud node introspect \
--all-manageable --provide
```

詳細については、**Director のインストールと使用** ガイドの [ベアメタルノードハードウェアのインベントリーの作成](#) を参照してください。

7. カスタム VDPA リソースクラスを使用して、VDPA 用に指定する各ベアメタルノードにタグを付けます。

```
(undercloud)$ openstack baremetal node set \
--resource-class baremetal.VDPA <node>
```

<node> は、ベアメタルノードの名前または UUID に置き換えます。

8. ノード定義ファイル **overcloud-baremetal-deploy.yaml** に **ComputeVdpa** ロールを追加し、予測ノード配置、リソースクラス、ネットワークポロジ、またはノードに割り当てるその他の属性を定義します。

```
- name: Controller
  count: 3
- name: Compute
  count: 3
- name: ComputeVdpa
  count: 1
  defaults:
    resource_class: baremetal.VDPA
    network_config:
      template: /home/stack/templates/nic-config/<role_topology_file>
```

- **<role_topology_file>** を ComputeVdpa ロールに使用するトポロジーファイルの名前 (**myRoleTopology.j2** など) に置き換えます。既存のネットワークポロジを再利用するか、ロール用の新しいカスタムネットワークインターフェイステンプレートを作成できます。詳細は、**director のインストールと使用** ガイドの [カスタムネットワークインターフェイステンプレート](#) を参照してください。デフォルトのネットワーク定義設定を使用するには、ロール定義に **network_config** を含めないでください。

ノード定義ファイルでノード属性を設定するために使用できるプロパティについて詳しくは、[ベアメタルノードのプロビジョニング属性](#) を参照してください。ノード定義ファイルの例は、[ノード定義ファイルの例](#) を参照してください。

9. ロールの新しいノードをプロビジョニングします。

```
(undercloud)$ openstack overcloud node provision \
[--stack <stack>] \
[--network-config \]
--output <deployment_file> \
/home/stack/templates/overcloud-baremetal-deploy.yaml
```

- オプション: **<stack>** をベアメタルノードがプロビジョニングされるスタックの名前に置き換えます。デフォルトは **overcloud** です。
 - オプション: **--network-config** オプションの引数を含めて、Ansible Playbook **cli-overcloud-node-network-config.yaml** にネットワーク定義を提供します。**network_config** プロパティを使用してノード定義ファイルにネットワーク定義を定義していない場合は、デフォルトのネットワーク定義が使用されます。
 - **<deployment_file>** は、デプロイメントコマンドに含めるために生成する heat 環境ファイルの名前に置き換えます (例: **/home/stack/templates/overcloud-baremetal-deployed.yaml**)。
10. 別のターミナルでプロビジョニングの進捗をモニタリングします。プロビジョニングが成功すると、ノードの状態が **available** から **active** に変わります。

```
(undercloud)$ watch openstack baremetal node list
```

11. **--network-config** オプションを指定せずにプロビジョニングコマンドを実行した場合は、**network-environment.yaml** ファイルで **<Role>NetworkConfigTemplate** パラメーターを設定して、NIC テンプレートファイルを指すようにします。

```
parameter_defaults:
  ComputeNetworkConfigTemplate: /home/stack/templates/nic-configs/compute.j2
  ComputeVdpaNetworkConfigTemplate: /home/stack/templates/nic-configs/<vdpa_net_top>.j2
  ControllerNetworkConfigTemplate: /home/stack/templates/nic-configs/controller.j2
```

<vdpa_net_top> を **ComputeVdpa** ロールのネットワークポロジーマスターが含まれるファイルの名前に置き換えます。たとえば、デフォルトのネットワークポロジーマスターを使用する場合は **compute.yaml** です。

6.2. VDPA コンピュートノードの設定

クラウドユーザーが VIRTIO データパスアクセラレーション (VDPA) ポートを使用するインスタンスを作成できるようにするには、VDPA デバイスを持つコンピュートノードを設定します。

手順

1. VDPA コンピュートノードを設定するための新しいコンピューティング環境ファイル (**vdpa_compute.yaml** など) を作成します。
2. **PciPassthroughFilter** と **NUMATopologyFilter** を **vdpa_compute.yaml** の **NovaSchedulerDefaultFilters** パラメーターに追加します。

```
parameter_defaults:
  NovaSchedulerDefaultFilters:
    ['AvailabilityZoneFilter','ComputeFilter','ComputeCapabilitiesFilter','ImagePropertiesFilter','ServerGroupAntiAffinityFilter','ServerGroupAffinityFilter','PciPassthroughFilter','NUMATopologyFilter']
```

3. **NovaPCIPassthrough** パラメーターを **vdpa_compute.yaml** に追加して、コンピュートノード上の VDPA デバイスで使用可能な PCI を指定します。たとえば、NVIDIA® ConnectX®-6 Dx デバイスを、インスタンスへのパススルーに使用できる PCI デバイスのプールに追加するには、次の設定を **vdpa_compute.yaml** に追加します。

```
parameter_defaults:
```

```
...
ComputeVdpaParameters:
  NovaPCIPassthrough:
    - vendor_id: "15b3"
      product_id: "101d"
      address: "06:00.0"
      physical_network: "tenant"
    - vendor_id: "15b3"
      product_id: "101d"
      address: "06:00.1"
      physical_network: "tenant"
```

NovaPCIPassthrough の設定方法の詳細は、[NovaPCIPassthrough 設定のガイドライン](#) を参照してください。

4. **KernelArgs** パラメーターを **vdpa_compute.yaml** に追加して、各コンピュートノード BIOS で入出力メモリー管理ユニット (IOMMU) を有効にします。たとえば、Intel Corporation IOMMU を有効にするには、次の **KernelArgs** 設定を使用します。

```
parameter_defaults:
  ...
  ComputeVdpaParameters:
    ...
    KernelArgs: "intel_iommu=on iommu=pt"
```

AMD IOMMU を有効にするには、**KernelArgs** を "**amd_iommu=on iommu=pt**" に設定します。



注記

KernelArgs パラメーターをロールの設定に初めて追加すると、オーバークラウドのデプロイメント中にオーバークラウドノードが自動的に再起動します。必要に応じて、ノードの自動再起動を無効にし、代わりに各オーバークラウドのデプロイ後にノードの再起動を手動で実行できます。詳細は、[KernelArgs を定義するための手動でのノード再起動の設定](#) を参照してください。

5. ネットワーク環境ファイルを開き、次の設定を追加して物理ネットワークを定義します。

```
parameter_defaults:
  ...
  NeutronBridgeMappings:
    - <bridge_map_1>
    - <bridge_map_n>
  NeutronTunnelTypes: '<tunnel_types>'
  NeutronNetworkType: '<network_types>'
  NeutronNetworkVLANRanges:
    - <network_vlan_range_1>
    - <network_vlan_range_n>
```

- **<bridge_map_1>**、および **<bridge_map_n>** までのすべてのブリッジマッピングを、VDPA ブリッジに使用する論理から物理へのブリッジマッピングに置き換えます。たとえば、**tenant:br-tenant** です。
- **<tunnel_types>** を、プロジェクトネットワークのトンネルタイプのコンマ区切りリストに置き換えます。たとえば、**geneve** です。

- **<network_types>** を Networking サービス (neutron) のプロジェクトネットワークタイプのコンマ区切りリストに置き換えます。利用可能なネットワークがすべてなくなるまで、最初に指定した種別が使用されます。その後、次の種別が使用されます。たとえば、**geneve,vlan** です。
 - **<network_vlan_range_1>**、および **<network_vlan_range_n>** までのすべての物理ネットワークと VLAN 範囲を、サポートする ML2 および OVN VLAN マッピング範囲に置き換えます。たとえば、**datacenter:1:1000**、**tenant:100:299**。
6. ネットワークインターフェイステンプレートを開き、次の設定を追加して、VDPA 対応のネットワークインターフェイスを OVN ブリッジのメンバーとして指定します。

```
- type: ovn_bridge
  name: br-tenant
  members:
    - type: sriov_pf
      name: enp6s0f0
      numvfs: 8
      use_dhcp: false
      vdpa: true
      link_mode: switchdev
    - type: sriov_pf
      name: enp6s0f1
      numvfs: 8
      use_dhcp: false
      vdpa: true
      link_mode: switchdev
```

7. その他の環境ファイルと共にこれらのカスタム環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
  -e [your environment files] \
  -r /home/stack/templates/roles_data_vdpa.yaml \
  -e /home/stack/templates/network-environment.yaml \
  -e /home/stack/templates/vdpa_compute.yaml \
  -e /home/stack/templates/overcloud-baremetal-deployed.yaml \
  -e /home/stack/templates/node-info.yaml
```

検証

1. VDPA デバイスでインスタンスを作成します。詳細は、[インスタンスの作成と管理](#) ガイドの [VDPA インターフェイスを使用したインスタンスの作成](#) を参照してください。
2. クラウドユーザーとしてインスタンスにログインします。詳細については、[インスタンスの作成と管理](#) ガイドの [インスタンスへの接続](#) を参照してください。
3. インスタンスから VDPA デバイスにアクセスできることを確認します。

```
$ openstack port show vdpa-port
```

第7章 ホストアグリゲートの作成と管理

クラウド管理者は、パフォーマンスおよび管理目的で、コンピュートのデプロイメントを論理グループに分割することができます。Red Hat OpenStack Platform (RHOSP) では、論理グループへの分割に以下のメカニズムを使用することができます。

ホストアグリゲート

ホストアグリゲートとは、ハードウェアやパフォーマンス特性などの属性に基づいてコンピュートノードを論理的なユニットにグループ化したものです。コンピュートノードを1つまたは複数のホストアグリゲートに割り当てることができます。

ホストアグリゲートでメタデータを設定してフレーバーおよびイメージをホストアグリゲートにマッピングし、続いてフレーバーの追加スペックまたはイメージのメタデータ属性をホストアグリゲートのメタデータにマッチングすることができます。必要なフィルターが有効な場合、Compute スケジューラーはこのメタデータを使用してインスタンスのスケジューリングを行うことができます。ホストアグリゲートで指定するメタデータは、ホストの使用をフレーバーまたはイメージで指定するメタデータと同じメタデータのインスタンスに限定します。

ホストアグリゲートのメタデータで **xxx_weight_multiplier** 設定オプションを定義することで、それぞれのホストアグリゲートに重みの乗数を設定することができます。

ホストアグリゲートを使用して、負荷分散の処理、物理的な分離または冗長性の適用、共通の属性を持つサーバーのグループ化、ハードウェアのクラス分け等を行うことができます。

ホストアグリゲートを作成する際に、ゾーン名を指定することができます。この名前は、クラウドユーザーが選択することのできるアベイラビリティゾーンとして提示されます。

アベイラビリティゾーン

アベイラビリティゾーンは、ホストアグリゲートのクラウドユーザー側のビューです。クラウドユーザーは、アベイラビリティゾーンに属するコンピュートノードやアベイラビリティゾーンのメタデータを把握することはできません。クラウドユーザーは、アベイラビリティゾーンの名前を見ることしかできません。

それぞれのコンピュートノードは、1つのアベイラビリティゾーンにしか割り当てることができません。デフォルトのアベイラビリティゾーンを設定することができます。クラウドユーザーがゾーンを指定しない場合、インスタンスはこのアベイラビリティゾーンにスケジューリングされます。特定の機能を持つアベイラビリティゾーンを使用するように、クラウドユーザーに指示することができます。

7.1. ホストアグリゲートでのスケジューリングの有効化

特定の属性を持つホストアグリゲートにインスタンスをスケジュールするには、Compute スケジューラーの設定を更新し、ホストアグリゲートのメタデータに基づく絞り込みを有効にします。

手順

1. Compute 環境ファイルを開きます。
2. **NovaSchedulerEnabledFilters** パラメーターにまだ以下の値がなければ、値を追加します。
 - **AggregateInstanceExtraSpecsFilter**: フレーバーの追加スペックに一致するホストアグリゲートメタデータでコンピュートノードを絞り込むには、この値を追加します。



注記

このフィルターが想定どおりに機能するには、**extra_specs** キーに **aggregate_instance_extra_specs**: 名前空間の接頭辞を指定して、フレーバーの追加スペックのスコープを定義する必要があります。

- **AggregateImagePropertiesIsolation**: イメージメタデータ属性に一致するホストアグリゲートメタデータでコンピュートノードを絞り込むには、この値を追加します。



注記

イメージメタデータ属性を使用してホストアグリゲートのメタデータを絞り込むには、ホストアグリゲートメタデータキーが有効なイメージメタデータ属性と一致する必要があります。有効なイメージメタデータ属性に関する情報は、[イメージ設定パラメーター](#) を参照してください。

- **AvailabilityZoneFilter**: インスタンスの起動時にアベイラビリティゾーンで絞り込むには、この値を追加します。



注記

Compute スケジューラーサービスのフィルター **AvailabilityZoneFilter** を使用する代わりに、Placement サービスを使用してアベイラビリティゾーンの要求を処理することができます。詳細は、[Filtering by availability zone using the Placement service](#) を参照してください。

3. 更新内容を Compute 環境ファイルに保存します。
4. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

7.2. ホストアグリゲートの作成

クラウド管理者は、ホストアグリゲートを必要なだけ作成することができます。

手順

1. ホストアグリゲートを作成するには、以下のコマンドを入力します。

```
(overcloud)# openstack aggregate create <aggregate_name>
```

<aggregate_name> をホストアグリゲートに割り当てる名前に置き換えてください。

2. ホストアグリゲートにメタデータを追加します。

```
(overcloud)# openstack aggregate set \
--property <key=value> \
--property <key=value> \
<aggregate_name>
```

- **<key=value>** はメタデータのキー/値のペアに置き換えてください。**AggregateInstanceExtraSpecsFilter** フィルターを使用している場合、キーは任意の文字列 (例: **ssd=true**) にすることができます。**AggregateImagePropertiesIsolation** フィルターを使用している場合は、キーは有効なイメージメタデータ属性と一致する必要があります。有効なイメージメタデータプロパティの詳細は、[イメージ設定パラメーター](#) を参照してください。
- **<aggregate_name>** をホストアグリゲートの名前に置き換えてください。

3. コンピュートノードをホストアグリゲートに追加します。

```
(overcloud)# openstack aggregate add host \
  <aggregate_name> \
  <host_name>
```

- **<aggregate_name>** は、コンピュートノードを追加するホストアグリゲートの名前に置き換えます。
- **<host_name>** は、ホストアグリゲートに追加するコンピュートノードの名前に置き換えてください。

4. ホストアグリゲート用のフレーバーまたはイメージを作成します。

- フレーバーの作成

```
(overcloud)$ openstack flavor create \
  --ram <size_mb> \
  --disk <size_gb> \
  --vcpus <no_reserved_vcpus> \
  host-aggr-flavor
```

- イメージの作成

```
(overcloud)$ openstack image create host-aggr-image
```

5. フレーバーまたはイメージに、ホストアグリゲートのキー/値のペアに一致するキー/値のペアを1つまたは複数設定します。

- フレーバーにキー/値のペアを設定するには、スコープ **aggregate_instance_extra_specs** を使用します。

```
(overcloud)# openstack flavor set \
  --property aggregate_instance_extra_specs:ssd=true \
  host-aggr-flavor
```

- イメージにキー/値のペアを設定するには、有効なイメージメタデータ属性をキーとして使用します。

```
(overcloud)# openstack image set \
  --property os_type=linux \
  host-aggr-image
```

7.3. アベイラビリティーゾーンの作成

クラウド管理者は、クラウドユーザーがインスタンスを作成する際に選択できるアベイラビリティーゾーンを作成することができます。

手順

1. アベイラビリティーゾーンを作成するには、新しいアベイラビリティーゾーンホストアグリゲートを作成するか、既存のホストアグリゲートをアベイラビリティーゾーンにすることができます。

- a. 新しいアベイラビリティーゾーンホストアグリゲートを作成するには、以下のコマンドを入力します。

```
(overcloud)# openstack aggregate create \
--zone <availability_zone> \
<aggregate_name>
```

- **<availability_zone>** をアベイラビリティーゾーンに割り当てる名前に置き換えてください。
- **<aggregate_name>** をホストアグリゲートに割り当てる名前に置き換えてください。

- b. 既存のホストアグリゲートをアベイラビリティーゾーンにするには、以下のコマンドを入力します。

```
(overcloud)# openstack aggregate set --zone <availability_zone> \
<aggregate_name>
```

- **<availability_zone>** をアベイラビリティーゾーンに割り当てる名前に置き換えてください。
- **<aggregate_name>** をホストアグリゲートの名前に置き換えてください。

2. オプション: アベイラビリティーゾーンにメタデータを追加します。

```
(overcloud)# openstack aggregate set --property <key=value> \
<aggregate_name>
```

- **<key=value>** をメタデータのキー/値のペアに置き換えてください。キー/値の属性は、必要なだけ追加することができます。
- **<aggregate_name>** をアベイラビリティーゾーンホストアグリゲートの名前に置き換えてください。

3. コンピュートノードをアベイラビリティーゾーンホストアグリゲートに追加します。

```
(overcloud)# openstack aggregate add host <aggregate_name> \
<host_name>
```

- **<aggregate_name>** は、コンピュートノードを追加するアベイラビリティーゾーンホストアグリゲートの名前に置き換えてください。
- **<host_name>** は、アベイラビリティーゾーンに追加するコンピュートノードの名前に置き換えてください。

7.4. ホストアグリゲートの削除

ホストアグリゲートを削除するには、まずホストアグリゲートからすべてのコンピュートノードを削除します。

手順

1. ホストアグリゲートに割り当てられたすべてのコンピュートノードの一覧を表示するには、以下のコマンドを入力します。

```
(overcloud)# openstack aggregate show <aggregate_name>
```

2. ホストアグリゲートから割り当てられたすべてのコンピュートノードを削除するには、それぞれのコンピュートノードで以下のコマンドを入力します。

```
(overcloud)# openstack aggregate remove host <aggregate_name> \
<host_name>
```

- **<aggregate_name>** は、コンピュートノードを削除するホストアグリゲートの名前に置き換えてください。
 - **<host_name>** は、ホストアグリゲートから削除するコンピュートノードの名前に置き換えてください。
3. ホストアグリゲートからすべてのコンピュートノードを削除したら、以下のコマンドを入力してホストアグリゲートを削除します。

```
(overcloud)# openstack aggregate delete <aggregate_name>
```

7.5. プロジェクト分離ホストアグリゲートの作成

特定のプロジェクトでのみ利用可能なホストアグリゲートを作成することができます。ホストアグリゲートに割り当てたプロジェクトだけが、ホストアグリゲートでインスタンスを起動することができます。



注記

プロジェクト分離では、Placement サービスを使用して各プロジェクトのホストアグリゲートを絞り込みます。このプロセスは、**AggregateMultiTenancyIsolation** フィルターの機能に優先します。したがって、**AggregateMultiTenancyIsolation** フィルターを使用する必要はありません。

手順

1. Compute 環境ファイルを開きます。
2. プロジェクト分離ホストアグリゲートでプロジェクトインスタンスをスケジュールするには、Compute 環境ファイルの **NovaSchedulerLimitTenantsToPlacementAggregate** パラメーターを **True** に設定します。
3. オプション: ホストアグリゲートに割り当てたプロジェクトだけがクラウド上でインスタンスを作成できるようにするには、**NovaSchedulerPlacementAggregateRequiredForTenants** パラメーターを **True** に設定します。



注記

NovaSchedulerPlacementAggregateRequiredForTenants のデフォルト値は **False** です。このパラメーターが **False** の場合、ホストアグリゲートに割り当てられていないプロジェクトは、任意のホストアグリゲートでインスタンスを作成することができます。

4. 更新内容を Compute 環境ファイルに保存します。
5. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml \
```

6. ホストアグリゲートを作成します。
7. プロジェクト ID の一覧を取得します。

```
(overcloud)# openstack project list
```

8. **filter_tenant_id<suffix>** メタデータキーを使用して、プロジェクトをホストアグリゲートに割り当てます。

```
(overcloud)# openstack aggregate set \
--property filter_tenant_id<ID0>=<project_id0> \
--property filter_tenant_id<ID1>=<project_id1> \
...
--property filter_tenant_id<IDn>=<project_idn> \
<aggregate_name>
```

- **<ID0>**、**<ID1>**、および **<IDn>** までのすべての ID を、作成する各プロジェクトフィルターの一意の値に置き換えてください。
- **<project_id0>**、**<project_id1>**、および **<project_idn>** までのすべてのプロジェクト ID を、ホストアグリゲートに割り当てる各プロジェクトの ID に置き換えてください。
- **<aggregate_name>** をプロジェクト分離ホストアグリゲートの名前に置き換えてください。

たとえば、プロジェクト **78f1**、**9d3t**、および **aa29** をホストアグリゲート **project-isolated-aggregate** に割り当てるには、以下の構文を使用します。

```
(overcloud)# openstack aggregate set \
--property filter_tenant_id=78f1 \
--property filter_tenant_id1=9d3t \
--property filter_tenant_id2=aa29 \
project-isolated-aggregate
```

ヒント

filter_tenant_id メタデータキーの接尾辞を省略することで、単一の特定プロジェクトでのみ利用可能なホストアグリゲートを作成することができます。

```
(overcloud)# openstack aggregate set \  
--property filter_tenant_id=78f1 \  
single-project-isolated-aggregate
```

関連情報

- ホストアグリゲートの作成に関する詳細は、[Creating and managing host aggregates](#) を参照してください。

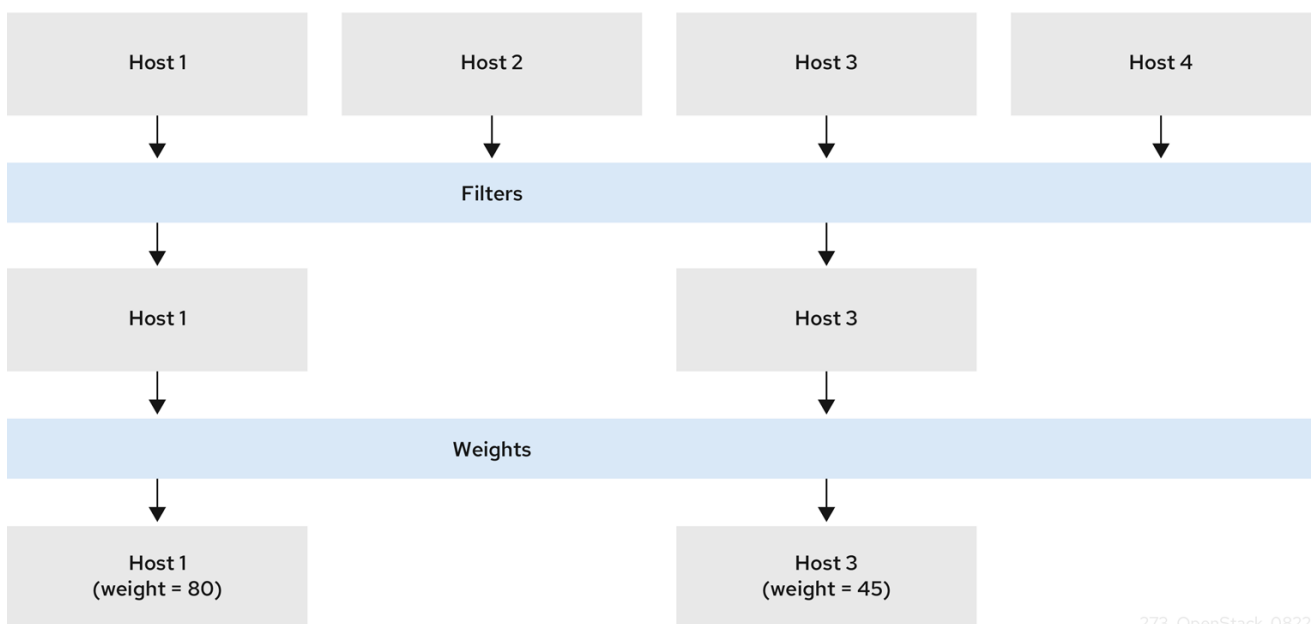
第8章 インスタンスのスケジューリングと配置の設定

Compute スケジューラーサービスは、インスタンスの配置先となるコンピュートノードまたはホストアグリゲートを決定します。Compute (nova) サービスがインスタンスの起動または移動に関するリクエストを受け取ると、リクエスト、フレーバー、およびイメージで提供される仕様を使用して適切なホストを決定します。たとえば、フレーバーでは、ストレージディスクの種別や Intel CPU 拡張命令セットなど、インスタンスがホストに要求する特性を指定することができます。

Compute スケジューラーサービスは、以下の順序で以下のコンポーネントの設定を使用して、インスタンスを起動または移動するコンピュートノードを決定します。

1. **Placement サービスのプレフィルター:** Compute スケジューラーサービスは Placement サービスを使用して、特定の属性に基づいて候補のコンピュートノードのセットを絞り込みます。たとえば、Placement サービスは無効な状態のコンピュートノードを自動的に除外します。
2. **フィルター:** Compute スケジューラーサービスは、これを使用してインスタンスを起動するコンピュートノードの初期セットを決定します。
3. **重み:** Compute スケジューラーサービスは、重み付けシステムを使用して絞り込まれたコンピュートノードの優先順位付けを行います。最も高い重みが最も優先されます。

下図では、絞り込み後、Host1 および 3 が条件を満たしています。Host1 の重みが最も高いため、スケジューリングで最も優先されます。



273_OpenStack_0822

8.1. PLACEMENT サービスを使用した事前絞り込み

Compute サービス (nova) は、Placement サービスと協調してインスタンスを作成および管理します。Placement サービスは、コンピュートノード、共有ストレージプール、または IP 割り当てプールなど、リソースプロバイダーのインベントリおよび使用状況、ならびに利用可能な仮想 CPU 数などのリソースの量的情報を追跡します。リソースの選択および消費を管理する必要があるサービスは、すべて Placement サービスを使用することができます。

Placement サービスは、リソースプロバイダーのストレージディスク特性の種別など、リソースの機能的情報とリソースプロバイダー間のマッピングも追跡します。

Placement サービスは、Placement サービスリソースプロバイダーインベントリーおよび特性に基づいて、候補のコンピュートノードセットにプレフィルタを適用します。以下の尺度に基づいてプレフィルタを作成することができます。

- サポートされるイメージ種別
- 特性
- プロジェクトまたはテナント
- アベイラビリティーゾーン

8.1.1. 要求されたイメージ種別のサポートによる絞り込み

インスタンスの起動に使用するイメージのディスク形式をサポートしないコンピュートノードを除外することができます。これは、環境の一時バックエンドに QCOW2 イメージをサポートしない Red Hat Ceph Storage が使用される場合に有用です。この機能を有効にすると、スケジューラーは QCOW2 イメージを使用するインスタンスの起動要求を Red Hat Ceph Storage ベースのコンピュートノードに送信しないようになります。

手順

1. Compute 環境ファイルを開きます。
2. インスタンスの起動に使用するイメージのディスク形式をサポートしないコンピュートノードを除外するには、Compute 環境ファイルの **NovaSchedulerQueryImageType** パラメーターを **True** に設定します。
3. 更新内容を Compute 環境ファイルに保存します。
4. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

8.1.2. リソースプロバイダー特性による絞り込み

各リソースプロバイダーには特性のセットがあります。特性は、ストレージディスクの種別や Intel CPU 拡張命令セットなど、リソースプロバイダーの機能的な要素です。

コンピュートノードは、その機能を特性として Placement サービスに報告します。インスタンスは、要求する特性またはリソースプロバイダーにあってはいけない特性を指定することができます。Compute スケジューラーは、これらの特性を使用して、インスタンスをホストするのに適したコンピュートノードまたはホストアグリゲートを特定することができます。

クラウドユーザーが特定の特性を持つホストにインスタンスを作成できるようにするには、特定の特性を要求または禁止するフレーバーを定義して、その特性を要求または禁止するイメージを作成することができます。

利用可能な特性の一覧は、[os-traits ライブラリー](#) を参照してください。必要に応じて、カスタムの特性を作成することもできます。

8.1.2.1. リソースプロバイダー特性を要求または禁止するイメージの作成

クラウドユーザーが特定の特性を持つホストでインスタンスを起動するのに使用することのできるインスタンスイメージを作成することができます。

手順

1. 新規イメージを作成します。

```
(overcloud)$ openstack image create ... trait-image
```

2. ホストまたはホストアグリゲートに必要な特性を識別します。既存の特性を選択するか、新たな特性を作成することができます。

- 既存の特性を使用するには、既存特性の一覧を表示して特性名を取得します。

```
(overcloud)$ openstack --os-placement-api-version 1.6 trait list
```

- 新規特性を作成するには、以下のコマンドを入力します。

```
(overcloud)$ openstack --os-placement-api-version 1.6 trait \
create CUSTOM_TRAIT_NAME
```

カスタムの特性は接頭辞 **CUSTOM_** で始まり、A から Z までの文字、0 から 9 までの数字、およびアンダースコア_だけを使用する必要があります。

3. 各ホストの既存のリソースプロバイダー特性を収集します。

```
(overcloud)$ existing_traits=$(openstack --os-placement-api-version 1.6 resource provider
trait list -f value <host_uuid> | sed 's/^/--trait /')
```

4. 既存のリソースプロバイダー特性に、ホストまたはホストアグリゲートに必要な特性があることを確認します。

```
(overcloud)$ echo $existing_traits
```

5. 必要な特性がまだリソースプロバイダーに追加されていない場合は、既存の特性と必要な特性を各ホストのリソースプロバイダーに追加してください。

```
(overcloud)$ openstack --os-placement-api-version 1.6 \
resource provider trait set $existing_traits \
--trait <TRAIT_NAME> \
<host_uuid>
```

<TRAIT_NAME> を、リソースプロバイダーに追加する特性の名前に置き換えます。必要に応じて、**--trait** オプションを複数回使用して、さらに特性を追加することができます。



注記

このコマンドは、リソースプロバイダーの特性をすべて置き換えます。したがって、ホスト上の既存のリソースプロバイダー特性の一覧を取得して、削除されないように再度設定する必要があります。

6. 要求された特性を持つホストまたはホストアグリゲートにインスタンスをスケジュールするには、イメージの追加スペックに特性を追加します。たとえば、AVX-512 をサポートするホスト

またはホストアグリゲートにインスタンスをスケジュールするには、イメージの追加スペックに以下の特性を追加します。

```
(overcloud)$ openstack image set \
  --property trait:HW_CPU_X86_AVX512BW=required \
  trait-image
```

7. 禁止された特性を持つホストまたはホストアグリゲートを除外するには、イメージの追加スペックに特性を追加します。たとえば、ボリュームの複数接続をサポートするホストまたはホストアグリゲートにインスタンスがスケジュールされるのを防ぐには、イメージの追加スペックに以下の特性を追加します。

```
(overcloud)$ openstack image set \
  --property trait:COMPUTE_VOLUME_MULTI_ATTACH=forbidden \
  trait-image
```

8.1.2.2. リソースプロバイダー特性を要求または禁止するフレーバーの作成

クラウドユーザーが特定の特性を持つホストでインスタンスを起動するのに使用することのできるフレーバーを作成することができます。

手順

1. フレーバーを作成します。

```
(overcloud)$ openstack flavor create --vcpus 1 --ram 512 \
  --disk 2 trait-flavor
```

2. ホストまたはホストアグリゲートに必要な特性を識別します。既存の特性を選択するか、新たな特性を作成することができます。

- 既存の特性を使用するには、既存特性の一覧を表示して特性名を取得します。

```
(overcloud)$ openstack --os-placement-api-version 1.6 trait list
```

- 新規特性を作成するには、以下のコマンドを入力します。

```
(overcloud)$ openstack --os-placement-api-version 1.6 trait \
  create CUSTOM_TRAIT_NAME
```

カスタムの特性は接頭辞 **CUSTOM_** で始まり、A から Z までの文字、0 から 9 までの数字、およびアンダースコア_だけを使用する必要があります。

3. 各ホストの既存のリソースプロバイダー特性を収集します。

```
(overcloud)$ existing_traits=$(openstack --os-placement-api-version 1.6 resource provider
  trait list -f value <host_uuid> | sed 's/^/--trait /')
```

4. 既存のリソースプロバイダー特性に、ホストまたはホストアグリゲートに必要な特性があることを確認します。

```
(overcloud)$ echo $existing_traits
```

5. 必要な特性がまだリソースプロバイダーに追加されていない場合は、既存の特性と必要な特性を各ホストのリソースプロバイダーに追加してください。

```
(overcloud)$ openstack --os-placement-api-version 1.6 \
resource provider trait set $existing_traits \
--trait <TRAIT_NAME> \
<host_uuid>
```

<TRAIT_NAME> を、リソースプロバイダーに追加する特性の名前に置き換えます。必要に応じて、**--trait** オプションを複数回使用して、さらに特性を追加することができます。



注記

このコマンドは、リソースプロバイダーの特性をすべて置き換えます。したがって、ホスト上の既存のリソースプロバイダー特性の一覧を取得して、削除されないように再度設定する必要があります。

6. 要求された特性を持つホストまたはホストアグリゲートにインスタンスをスケジュールするには、フレーバーの追加スペックに特性を追加します。たとえば、AVX-512 をサポートするホストまたはホストアグリゲートにインスタンスをスケジュールするには、フレーバーの追加スペックに以下の特性を追加します。

```
(overcloud)$ openstack flavor set \
--property trait:HW_CPU_X86_AVX512BW=required \
trait-flavor
```

7. 禁止された特性を持つホストまたはホストアグリゲートを除外するには、フレーバーの追加スペックに特性を追加します。たとえば、ボリュームの複数接続をサポートするホストまたはホストアグリゲートにインスタンスがスケジュールされるのを防ぐには、フレーバーの追加スペックに以下の特性を追加します。

```
(overcloud)$ openstack flavor set \
--property trait:COMPUTE_VOLUME_MULTI_ATTACH=forbidden \
trait-flavor
```

8.1.3. ホストアグリゲートの分離による絞り込み

ホストアグリゲートへのスケジューリングを、フレーバーおよびイメージの特性がホストアグリゲートのメタデータと一致するインスタンスだけに制限することができます。フレーバーとイメージのメタデータの組み合わせでは、そのホストアグリゲートに属するコンピュータードへのスケジューリングを有効にするホストアグリゲート特性をすべて要求する必要があります。

手順

1. Compute 環境ファイルを開きます。
2. ホストアグリゲートを分離してフレーバーおよびイメージの特性がアグリゲートのメタデータと一致するインスタンスだけをホストするには、Compute 環境ファイルの **NovaSchedulerEnableIsolatedAggregateFiltering** パラメーターを **True** に設定します。
3. 更新内容を Compute 環境ファイルに保存します。
4. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

5. ホストアグリゲートを分離する対象の特性を特定します。既存の特性を選択するか、新たな特性を作成することができます。

- 既存の特性を使用するには、既存特性の一覧を表示して特性名を取得します。

```
(overcloud)$ openstack --os-placement-api-version 1.6 trait list
```

- 新規特性を作成するには、以下のコマンドを入力します。

```
(overcloud)$ openstack --os-placement-api-version 1.6 trait \
create CUSTOM_TRAIT_NAME
```

カスタムの特性は接頭辞 **CUSTOM_** で始まり、A から Z までの文字、0 から 9 までの数字、およびアンダースコア_だけを使用する必要があります。

6. 各コンピュートノードの既存のリソースプロバイダー特性を収集します。

```
(overcloud)$ existing_traits=$(openstack --os-placement-api-version 1.6 resource provider
trait list -f value <host_uuid> | sed 's/^/--trait /')
```

7. 既存のリソースプロバイダー特性で、ホストアグリゲートを分離する特性を確認します。

```
(overcloud)$ echo $existing_traits
```

8. 必要な特性がまだリソースプロバイダーに追加されていない場合は、既存の特性と必要な特性をホストアグリゲートの各コンピュートノードのリソースプロバイダーに追加してください。

```
(overcloud)$ openstack --os-placement-api-version 1.6 \
resource provider trait set $existing_traits \
--trait <TRAIT_NAME> \
<host_uuid>
```

<**TRAIT_NAME**> を、リソースプロバイダーに追加する特性の名前に置き換えます。必要に応じて、**--trait** オプションを複数回使用して、さらに特性を追加することができます。



注記

このコマンドは、リソースプロバイダーの特性をすべて置き換えます。したがって、ホスト上の既存のリソースプロバイダー特性の一覧を取得して、削除されないように再度設定する必要があります。

9. ホストアグリゲートに属する各コンピュートノードで、ステップ 6 - 8 を繰り返します。
10. 特性のメタデータ属性をホストアグリゲートに追加します。

```
(overcloud)$ openstack --os-compute-api-version 2.53 aggregate set \
--property trait:<TRAIT_NAME>=required <aggregate_name>
```

11. フレーバーまたはイメージに特性を追加します。

```
(overcloud)$ openstack flavor set \
--property trait:<TRAIT_NAME>=required <flavor>
(overcloud)$ openstack image set \
--property trait:<TRAIT_NAME>=required <image>
```

8.1.4. Placement サービスを使用したアベイラビリティゾーンによる絞り込み

Placement サービスを使用して、アベイラビリティゾーンの要求を適用することができます。Placement サービスを使用してアベイラビリティゾーンで絞り込むには、アベイラビリティゾーンホストアグリゲートのメンバーシップおよび UUID と一致する配置アグリゲートが存在する必要があります。

手順

1. Compute 環境ファイルを開きます。
2. Placement サービスを使用してアベイラビリティゾーンで絞り込むには、Compute 環境ファイルの **NovaSchedulerQueryPlacementForAvailabilityZone** パラメーターを **True** に設定します。
3. **NovaSchedulerEnabledFilters** パラメーターから **AvailabilityZoneFilter** フィルターを削除します。
4. 更新内容を Compute 環境ファイルに保存します。
5. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

関連情報

- アベイラビリティゾーンとして使用するホストアグリゲートの作成に関する詳細は、[Creating an availability zone](#) を参照してください。

8.2. COMPUTE スケジューラーサービス用フィルターおよび重みの設定

インスタンスを起動するコンピュータノードの初期セットを決定するには、Compute スケジューラーサービス用にフィルターおよび重みを設定する必要があります。

手順

1. Compute 環境ファイルを開きます。
2. スケジューラーが使用するフィルターを **NovaSchedulerEnabledFilters** パラメーターに追加します。以下に例を示します。

```
parameter_defaults:
  NovaSchedulerEnabledFilters:
    - AggregateInstanceExtraSpecsFilter
```

- ComputeFilter
- ComputeCapabilitiesFilter
- ImagePropertiesFilter

3. 各コンピュートノードの重みを計算するのに使用する属性を指定します。以下に例を示します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::config::nova_config:
      filter_scheduler/weight_classes:
        value: nova.scheduler.weights.all_weighers
```

利用可能な属性についての詳しい情報は、[Compute スケジューラーの重み](#)を参照してください。

4. オプション: 各重み付け関数に適用する重みの乗数を設定します。たとえば、コンピュートノードの利用可能な RAM が他のデフォルトの重み付け関数よりも高い重みを持つように指定し、Compute スケジューラーが、利用可能な RAM が少ないコンピュートノードよりも、利用可能な RAM が多いコンピュートノードを優先するには、以下の設定を使用します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::config::nova_config:
      filter_scheduler/weight_classes:
        value: nova.scheduler.weights.all_weighers
      filter_scheduler/ram_weight_multiplier:
        value: 2.0
```

ヒント

また、重みの乗数を負の値に設定することもできます。上記の例では、利用可能な RAM が多いコンピュートノードよりも利用可能な RAM が少ないノードを優先するには、**ram_weight_multiplier** を **-2.0** に設定します。

5. 更新内容を Compute 環境ファイルに保存します。
6. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

関連情報

- 利用可能な Compute スケジューラーサービスのフィルター一覧は、[Compute scheduler filters](#)を参照してください。
- 利用可能な重みの設定オプション一覧は、[Compute scheduler weights](#)を参照してください。

8.3. COMPUTE スケジューラーのフィルター

インスタンスをホストするのに適切なコンピュートノードを選択する際に Compute スケジューラーが適用しなければならないフィルターを指定するには、Compute 環境ファイルの **NovaSchedulerEnabledFilters** パラメーターを設定します。デフォルト設定では、以下のフィルターが適用されます。

- **AvailabilityZoneFilter**: コンピュートノードは要求されたアベイラビリティゾーンに属していない。
- **ComputeFilter**: コンピュートノードは要求に対応することができる。
- **ComputeCapabilitiesFilter**: コンピュートノードはフレーバーの追加スペックを満足する。
- **ImagePropertiesFilter**: コンピュートノードは要求されたイメージ属性を満足する。
- **ServerGroupAntiAffinityFilter**: コンピュートノードは、まだ指定されたグループに属するインスタンスをホストしていない。
- **ServerGroupAffinityFilter**: コンピュートノードは、すでに指定されたグループに属するインスタンスをホストしている。

フィルターを追加および削除することができます。利用可能なすべてのフィルターの詳細を以下の表に示します。

表8.1 Compute スケジューラーのフィルター

フィルター	説明
AggregateImagePropertiesIsolation	このフィルターを使用して、インスタンスのイメージメタデータとホストアグリゲートのメタデータを照合します。いずれかのホストアグリゲートのメタデータがイメージのメタデータと一致する場合は、そのホストアグリゲートに属するコンピュートノードはそのイメージからインスタンスを起動する候補となります。スケジューラーは、有効なイメージメタデータ属性のみを認識します。有効なイメージメタデータプロパティの詳細は、 イメージ設定パラメーター を参照してください。
AggregateInstanceExtraSpecsFilter	このフィルターを使用して、インスタンスのフレーバー追加スペックで定義された名前空間属性とホストアグリゲートのメタデータを照合します。 フレーバー extra_specs キーのスコープは、 aggregate_instance_extra_specs: namespace の前に付けて定義する必要があります。 いずれかのホストアグリゲートのメタデータがフレーバー追加スペックのメタデータと一致する場合は、そのホストアグリゲートに属するコンピュートノードはそのイメージからインスタンスを起動する候補となります。
AggregateIopsFilter	このフィルターを使用して、ホストアグリゲートごとの filter_scheduler/max_io_ops_per_host 値を条件に I/O 操作でホストを絞り込みます。ホストアグリゲートごとの値が確認されない場合は、値はグローバル設定にフォールバックします。ホストが複数のアグリゲートに属し、複数の値が確認された場合、スケジューラーは最小の値を使用します。

フィルター	説明
AggregateMultiTenancyIsolation	<p>このフィルターを使用して、プロジェクト分離ホストアグリゲートに属するコンピュートノードの可用性を、指定したプロジェクトのセットに制限します。filter_tenant_id メタデータキーを使用して指定したプロジェクトだけが、ホストアグリゲートに属するコンピュートノードでインスタンスを起動することができます。詳しくは、Creating a project-isolated host aggregate を参照してください。</p> <div>  <div> <p>注記</p> <p>プロジェクトが他のホストにインスタンスを配置することは可能です。これを制限するには、NovaSchedulerPlacementAggregateRequiredForTenants パラメーターを使用します。</p> </div> </div>
AggregateNumInstancesFilter	<p>このフィルターを使用して、アグリゲートに属する各コンピュートノードでホスト可能なインスタンスの数を制限します。filter_scheduler/max_instances_per_host パラメーターを使用して、ホストアグリゲートごとのインスタンスの最大数を設定することができます。ホストアグリゲートごとの値が確認されない場合は、値はグローバル設定にフォールバックします。コンピュートノードが複数のホストアグリゲートに属する場合、スケジューラーは最小の max_instances_per_host 値を使用します。</p>
AggregateTypeAffinityFilter	<p>フレーバーメタデータキーが設定されていない場合や、フレーバーアグリゲートメタデータの値に要求するフレーバーの名前が含まれる場合には、このフィルターを使用してホスト渡します。フレーバーメタデータエントリーの値は、単一のフレーバー名またはフレーバー名のコンマ区切りリストのいずれかを含む文字列です (例: m1.nano または m1.nano,m1.small)。</p>
AllHostsFilter	<p>このフィルターを使用して、利用可能なすべてのコンピュートノードをインスタンスのスケジューリング対象として考慮します。</p> <div>  <div> <p>注記</p> <p>このフィルターを使用しても、他のフィルターは無効になりません。</p> </div> </div>
AvailabilityZoneFilter	<p>このフィルターを使用して、インスタンスが指定するアベイラビリティーゾーンに属するコンピュートノードでインスタンスを起動します。</p>

フィルター	説明
ComputeCapabilitiesFilter	<p>このフィルターを使用して、インスタンスのフレーバー追加スペックで定義した名前空間属性とコンピュートノードのキャパビリティを照合します。フレーバーの追加スペックに capabilities: 名前空間の接頭辞を指定する必要があります。</p> <p>ComputeCapabilitiesFilter フィルターを使用するよりも効率的な方法は、フレーバーで CPU 特性を使用することです。これは、Placement サービスに報告されます。特性により、CPU 機能に一貫性のある名前が付けられます。詳細は、Filtering by using resource provider traits を参照してください。</p>
ComputeFilter	<p>このフィルターを使用して、稼働中で有効なすべてのコンピュートノードを渡します。このフィルターは常に設定されている必要があります。</p>
DifferentHostFilter	<p>このフィルターを使用して、特定のインスタンスセットとは異なるコンピュートノードへのインスタンスのスケジューリングを有効にします。インスタンスの起動時にこれらのインスタンスを指定するには、--hint 引数を使用して different_host およびインスタンスの UUID をキー/値のペアとして指定します。</p> <pre>\$ openstack server create --image cedef40a-ed67-4d10-800e-17455edce175 \ --flavor 1 --hint different_host=a0cf03a5-d921-4877-bb5c-86d26cf818e1 \ --hint different_host=8c19174f-4220-44f0-824a-cd1eeef10287 server-1</pre>
ImagePropertiesFilter	<p>このフィルターを使用して、インスタンスイメージで定義された以下の属性に基づいてコンピュートノードを絞り込みます。</p> <ul style="list-style-type: none"> ● hw_architecture: ホストのアーキテクチャーを定義します (例: x86、ARM、および Power)。 ● img_hv_type: ハイパーバイザーの種別を定義します (例: KVM、QEMU、Xen、および LXC)。 ● img_hv_requested_version: Compute サービスが報告するハイパーバイザーのバージョンを定義します。 ● hw_vm_mode: ハイパーバイザーの種別を定義します (例: hvm、xen、uml、または exe)。 <p>インスタンスに含まれる指定のイメージ属性をサポートできるコンピュートノードが、スケジューラーに渡されます。イメージプロパティの詳細は、イメージ設定パラメーター を参照してください。</p>

フィルター	説明
IsolatedHostsFilter	<p>このフィルターを使用して、分離したコンピュートノード上で分離したイメージだけを使用してインスタンスをスケジュールします。また、filter_scheduler/restrict_isolated_hosts_to_isolated_imagesを設定して、分離したコンピュートノード上でのインスタンスのビルドに分離していないイメージが使用されるのを防ぐこともできます。</p> <p>分離されたイメージとホストのセットを指定するには、filter_scheduler/isolated_hosts および filter_scheduler/isolated_images 設定オプションを使用します。以下に例を示します。</p> <pre>parameter_defaults: ComputeExtraConfig: nova::config::nova_config: filter_scheduler/isolated_hosts: value: server1, server2 filter_scheduler/isolated_images: value: 342b492c-128f-4a42-8d3a-c5088cf27d13, ebd267a6-ca86-4d6c-9a0e-bd132d6b7d09</pre>
IoOpsFilter	<p>このフィルターを使用して、(ホストで実行可能な高 I/O 負荷インスタンスの最大数を指定する) filter_scheduler/max_io_ops_per_host の設定値を超える同時 I/O 操作があるホストを除外します。</p>
MetricsFilter	<p>このフィルターを使用して、metrics/weight_setting を使用して設定されたメトリックを報告するコンピュートノードにスケジューリングを制限します。</p> <p>このフィルターを使用するには、Compute 環境ファイルに以下の設定を追加します。</p> <pre>parameter_defaults: ComputeExtraConfig: nova::config::nova_config: DEFAULT/compute_monitors: value: 'cpu.virt_driver'</pre> <p>デフォルトでは、Compute スケジューラーサービスは 60 秒ごとにメトリックを更新します。メトリックが最新の状態になるようにするには、update_resources_interval 設定オプションを使用してメトリックデータの更新頻度を上げることができます。たとえば、以下の設定を使用すると、2 秒ごとにメトリックデータが更新されます。</p> <pre>parameter_defaults: ComputeExtraConfig: nova::config::nova_config: DEFAULT/update_resources_interval: value: '2'</pre>

フィルター	説明
NUMATopologyFilter	<p>このフィルターを使用して、NUMA 対応コンピュータノードに NUMA トポロジーが設定されたインスタンスをスケジューリングします。フレーバー extra_specs およびイメージ属性を使用して、インスタンスの NUMA トポロジーを指定します。このフィルターは、各ホストの NUMA セルのオーバーサブスクリプション限度を考慮して、インスタンスの NUMA トポロジーをコンピュータノードのトポロジーに一致させるを試みます。</p>
NumInstancesFilter	<p>このフィルターを使用して、max_instances_per_host オプションで指定した以上のインスタンスを実行中のコンピュータノードを除外します。</p>
PciPassthroughFilter	<p>このフィルターを使用して、フレーバー extra_specs を使用してインスタンスが要求するデバイスを持つコンピュータノードにインスタンスをスケジューリングします。</p> <p>(通常高価で使用が制限される) PCI デバイスを要求するインスタンス用に、そのデバイスを持つノードを確保する場合に、このフィルターを使用します。</p>
SameHostFilter	<p>このフィルターを使用して、特定のインスタンスセットと同じコンピュータノードへのインスタンスのスケジューリングを有効にします。インスタンスの起動時にこれらのインスタンスを指定するには、--hint 引数を使用して same_host およびインスタンスの UUID をキー/値のペアとして指定します。</p> <pre>\$ openstack server create --image cedef40a-ed67-4d10-800e-17455edce175 \ --flavor 1 --hint same_host=a0cf03a5-d921-4877-bb5c-86d26cf818e1 \ --hint same_host=8c19174f-4220-44f0-824a-cd1eeef10287 server-1</pre>
ServerGroupAffinityFilter	<p>このフィルターを使用して、同じコンピュータノード上でアフィニティーサーバーグループに属するインスタンスをスケジューリングします。サーバーグループを作成するには、以下のコマンドを入力します。</p> <pre>\$ openstack server group create --policy affinity <group_name></pre> <p>このグループに属するインスタンスを起動するには、--hint 引数を使用して group およびグループの UUID をキー/値のペアとして指定します。</p> <pre>\$ openstack server create --image <image> \ --flavor <flavor> \ --hint group=<group_uuid> <instance_name></pre>

フィルター	説明
ServerGroupAntiAffinity Filter	<p>このフィルターを使用して、異なるコンピューターノード上で非アフィニティーサーバーグループに属するインスタンスをスケジュールします。サーバーグループを作成するには、以下のコマンドを入力します。</p> <pre>\$ openstack server group create --policy anti-affinity <group_name></pre> <p>このグループに属するインスタンスを起動するには、--hint 引数を使用して group およびグループの UUID をキー/値のペアとして指定します。</p> <pre>\$ openstack server create --image <image> \ --flavor <flavor> \ --hint group=<group_uuid> <instance_name></pre>
SimpleCIDRAffinityFilter	<p>このフィルターを使用して、特定の IP サブネット範囲を持つコンピューターノードにインスタンスをスケジュールします。必要な範囲を指定するには、--hint 引数を使用してインスタンスの起動時にキー build_near_host_ip および cidr を渡します。</p> <pre>\$ openstack server create --image <image> \ --flavor <flavor> \ --hint build_near_host_ip=<ip_address> \ --hint cidr=<subnet_mask> <instance_name></pre>

8.4. COMPUTE スケジューラーの重み

それぞれのコンピューターノードは重みを持ち、スケジューラーはこれを使用してインスタンスのスケジューリングの優先度を決定します。フィルターを適用した後、Compute スケジューラーは残った候補のコンピューターノードから最大の重みを持つコンピューターノードを選択します。

Compute スケジューラーは、以下のタスクを実行することにより、各コンピューターノードの重みを決定します。

1. スケジューラーは、各重みを 0.0 から 1.0 までの値に正規化します。
2. スケジューラーは、正規化された重みを重み付け関数の乗数で乗算します。

Compute スケジューラーは、候補のコンピューターノード全体でリソースの可用性の低い値および高い値を使用して、各リソース種別の重みの正規化を計算します。

- 最も低いリソース可用性 (minval) を持つノードには、0 が割り当てられます。
- 最も高いリソース可用性 (maxval) を持つノードには 1 が割り当てられます。
- minval と maxval 内の範囲のリソース可用性を持つノードには、以下の式を使用して正規化された重みが割り当てられます。

$$(\text{node_resource_availability} - \text{minval}) / (\text{maxval} - \text{minval})$$

すべてのコンピューターノードが同じリソース可用性を持つ場合、それらはすべて 0 に正規化されます。

たとえば、スケジューラーは、利用可能な仮想 CPU の数がそれぞれ異なる 10 個のコンピュータノードに関して、利用可能な仮想 CPU の正規化された重みを以下のように計算します。

コンピュータノード	1	2	3	4	5	6	7	8	9	10
仮想 CPU の数	5	5	10	10	15	20	20	15	10	5
正規化された重み	0	0	0.33	0.33	0.67	1	1	0.67	0.33	0

Compute スケジューラーは、以下の式を使用してコンピュータノードの重みを計算します。

$$(w1_multiplier * norm(w1)) + (w2_multiplier * norm(w2)) + ...$$

重みに使用することのできる設定オプションの詳細を以下の表に示します。



注記

以下の表で説明するオプションと同じ名前のアグリゲートメタデータのキーを使用して、ホストアグリゲートに重みを設定することができます。ホストアグリゲートに設定すると、ホストアグリゲートの値が優先されます。

表8.2 Compute スケジューラーの重み


設定オプション	型	説明
---------	---	----


設定オプション	型	説明
filter_scheduler/weight_classes	文字列	<p>このパラメーターを使用して、各コンピュートノードの重みを計算するのに以下の属性のどれを使用するかを設定します。</p> <ul style="list-style-type: none"> ● nova.scheduler.weights.ram.RAMWeigher: コンピュートノードで利用可能な RAM を重み付けします。 ● nova.scheduler.weights.cpu.CPUWeigher: コンピュートノードで利用可能な CPU の数を重み付けします。 ● nova.scheduler.weights.disk.DiskWeigher: コンピュートノードで利用可能なディスクを重み付けします。 ● nova.scheduler.weights.metrics.MetricsWeigher: コンピュートノードのメトリクスを重み付けします。 ● nova.scheduler.weights.affinity.ServerGroupSoftAffinityWeigher: 指定したインスタンスグループの他のノードとコンピュートノードの近接性を重み付けします。 ● nova.scheduler.weights.affinity.ServerGroupSoftAntiAffinityWeigher: 指定したインスタンスグループの他のノードとコンピュートノードの近接性を重み付けします。 ● nova.scheduler.weights.compute.BuildFailureWeigher: 直近ブート試行の失敗回数でコンピュートノードを重み付けします。 ● nova.scheduler.weights.io_ops.IoOpsWeigher: ワークロードでコンピュートノードを重み付けします。 ● nova.scheduler.weights.pci.PCIWeigher: PCI の可用性でコンピュートノードを重み付けします。 ● nova.scheduler.weights.cross_cell.CrossCellWeigher: 置かれているセルに基づいてコンピュートノードを重み付けします。インスタンスを移動する際に、移行元セル内にあるコンピュートノードを優先します。 ● nova.scheduler.weights.all_weighters: (デフォルト) 上記の重み付け関数をすべて使用します。

設定オプション	型	説明
filter_scheduler/ram_weight_multiplier	浮動小数点	<p>このパラメーターを使用して、利用可能な RAM 容量に基づいてホストを重み付けするのに使用する重みの乗数を指定します。</p> <p>利用可能な RAM 容量がより大きいホストを優先するには、正の値に設定します。この場合、インスタンスは多くのホストに分散されます。</p> <p>利用可能な RAM 容量がより小さいホストを優先するには、負の値に設定します。この場合、可能な限り多くのインスタンスをホストに分担 (スタック) させた後に、使用率が低いホストをスケジューリングします。</p> <p>正または負の絶対値で、他の重み付け関数と比べて {b}RAM の重み付け関数をどれだけ優先するかを指定します。</p> <p>デフォルトは 1.0 で、スケジューラーはインスタンスをすべてのホストに均等に分散します。</p>
filter_scheduler/disk_weight_multiplier	浮動小数点	<p>このパラメーターを使用して、利用可能なディスク容量に基づいてホストを重み付けするのに使用する重みの乗数を指定します。</p> <p>利用可能なディスク容量がより大きいホストを優先するには、正の値に設定します。この場合、インスタンスは多くのホストに分散されます。</p> <p>利用可能なディスク容量がより小さいホストを優先するには、負の値に設定します。この場合、可能な限り多くのインスタンスをホストに分担 (スタック) させた後に、使用率が低いホストをスケジューリングします。</p> <p>正または負の絶対値で、他の重み付け関数と比べてディスクの重み付け関数をどれだけ優先するかを指定します。</p> <p>デフォルトは 1.0 で、スケジューラーはインスタンスをすべてのホストに均等に分散します。</p>

設定オプション	型	説明
filter_scheduler/cpu_weight_multiplier	浮動小数点	<p>このパラメーターを使用して、利用可能な仮想 CPU の数に基づいてホストを重み付けするのに使用する重みの乗数を指定します。</p> <p>利用可能な仮想 CPU の数がより多いホストを優先するには、正の値に設定します。この場合、インスタンスは多くのホストに分散されます。</p> <p>利用可能な仮想 CPU の数がより少ないホストを優先するには、負の値に設定します。この場合、可能な限り多くのインスタンスをホストに分担 (スタック) させた後に、使用率が低いホストをスケジューリングします。</p> <p>正または負の絶対値で、他の重み付け関数と比べて仮想 CPU の重み付け関数をどれだけ優先するかを指定します。</p> <p>デフォルトは 1.0 で、スケジューラーはインスタンスをすべてのホストに均等に分散します。</p>
filter_scheduler/io_ops_weight_multiplier	浮動小数点	<p>このパラメーターを使用して、負荷に基づいてホストを重み付けするのに使用する重みの乗数を指定します。</p> <p>負荷がより軽いホストを優先するには、負の値に設定します。この場合、負荷はより多くのホストに分散されます。</p> <p>負荷がより重いホストを優先するには、正の値に設定します。この場合、インスタンスはすでにビジー状態にあるホストにスケジューリングされます。</p> <p>正または負の絶対値で、他の重み付け関数と比べて I/O 操作の重み付け関数をどれだけ優先するかを指定します。</p> <p>デフォルトは -1.0 で、スケジューラーは負荷をより多くのホストに分散します。</p>
filter_scheduler/build_failure_weight_multiplier	浮動小数点	<p>このパラメーターを使用して、直近のビルド失敗回数に基づいてホストを重み付けするのに使用する重みの乗数を指定します。</p> <p>ホストから報告される直近のビルド失敗回数をより重要視するには、正の値に設定します。直近ビルドに失敗したホストは、選択されにくくなります。</p> <p>直近の失敗回数でコンピュートホストを重み付けするのを無効にするには、0 に設定します。</p> <p>デフォルト: 1000000.0</p>

設定オプション	型	説明
filter_scheduler/cross_cell_move_weight_multiplier	浮動小数点	<p>このパラメーターを使用して、セルを越えて移動する際にホストを重み付けするのに使用する重みの乗数を指定します。このオプションは、インスタンスを移動する際に、同じ移動元セル内にあるホストに加える重みを決定します。インスタンスを移行する場合、デフォルトではスケジューラーは同じ移行元セル内にあるホストを優先します。</p> <p>現在インスタンスを実行中のセル内にあるホストを優先するには、正の値に設定します。現在インスタンスを実行中のセルとは別のセルにあるホストを優先するには、負の値に設定します。</p> <p>デフォルト: 1000000.0</p>
filter_scheduler/pci_weight_multiplier	正の浮動小数点	<p>このパラメーターを使用して、ホスト上の PCI デバイスの数とインスタンスが要求する PCI デバイスの数に基づいてホストを重み付けするのに使用する重みの乗数を指定します。インスタンスが PCI デバイスを要求する場合、より多くの PCI デバイスを持つコンピュータードにより高い重みが割り当てられます。</p> <p>たとえば、ホストが 3 台利用可能で、PCI デバイスが 1 つのホストが 1 台、複数の PCI デバイスがあるホストが 1 台、PCI デバイスがないホストが 1 台の場合には、Compute スケジューラーはインスタンスの需要に基づいてこれらのホストの優先順位付けを行います。スケジューラーは、インスタンスが PCI デバイスを 1 つ要求している場合には 1 番目のホストを、複数の PCI デバイスを要求している場合には 2 番目のホストを、PCI デバイスを要求していない場合には 3 番目のホストを、それぞれ優先するはずでず。</p> <p>このオプションを設定して、PCI を要求しないインスタンスが PCI デバイスを持つホストのリソースを占有するのを防ぎます。</p> <p>デフォルト: 1.0</p>

設定オプション	型	説明
filter_scheduler/host_subset_size	Integer	<p>このパラメーターを使用して、ホストを選択するサブセット (絞り込まれたホストのサブセット) のサイズを指定します。このオプションを1以上に設定する必要があります。値を1に指定した場合には、重み付け関数によって最初に返されるホストが選択されます。スケジューラーは1未満の値を無視し、代わりに1を使用します。</p> <p>類似の要求を処理する複数のスケジューラープロセスが同じホストを選択して競合状態が生じるのを防ぐには、1より大きい値に設定します。要求に最も適したN台のホストからホストを無作為に選択することで、競合の可能性が低減されます。ただし、この値を高く設定すると、選択されるホストが特定の要求に対して最適ではない可能性が高くなります。</p> <p>デフォルト:1</p>
filter_scheduler/soft_affinity_weight_multiplier	正の浮動小数点	<p>このパラメーターを使用して、グループのソフトアフィニティーのホストを重み付けするのに使用する重みの乗数を指定します。</p> <div>  <div> <p>注記</p> <p>このポリシーでグループを作成する場合は、マイクロバージョンを指定する必要があります。</p> <pre>\$ openstack --os-compute-api-version 2.15 server group create --policy soft-affinity <group_name></pre> </div> </div> <p>デフォルト:1.0</p>

設定オプション	型	説明
filter_scheduler/soft_anti_affinity_weight_multiplier	正の浮動小数点	<p>このパラメーターを使用して、グループのソフト非アフィニティーのホストを重み付けするのに使用する重みの乗数を指定します。</p> <div>  <div> <p>注記</p> <p>このポリシーでグループを作成する場合は、マイクロバージョンを指定する必要があります。</p> <pre>\$ openstack --os-compute-api-version 2.15 server group create --policy soft-affinity <group_name></pre> </div> </div> <p>デフォルト: 1.0</p>
metrics/weight_multiplier	浮動小数点	<p>このパラメーターを使用して、メトリックの重み付けに使用する重みの乗数を指定します。デフォルトでは weight_multiplier=1.0 に設定されており、使用可能な全ホストの間でインスタンスを分散します。</p> <p>重み全体でメトリックの影響を増大させるには、1.0 を超える数値に設定します。</p> <p>重み全体でメトリックの影響を減少させるには、0.0 と 1.0 の間の数値に設定します。</p> <p>メトリックの値を無視して weight_of_unavailable オプションの値を返すには、0.0 に設定します。</p> <p>低いメトリックのホストを優先してインスタンスをホストにスタックするには、負の数値に設定します。</p> <p>デフォルト: 1.0</p>

設定オプション	型	説明
metrics/weight_setting	metric=ratio ペアの コンマ区切りリスト	<p>このパラメーターを使用して、重み付けに使用するメトリック、および各メトリックの重みを計算するのに使用する比率を指定します。有効なメトリック名は以下のとおりです。</p> <ul style="list-style-type: none"> ● cpu.frequency: CPU の周波数 ● cpu.user.time: CPU のユーザーモード時間 ● cpu.kernel.time: CPU のカーネル時間 ● cpu.idle.time: CPU のアイドル時間 ● cpu.iowait.time: CPU の I/O 待機時間 ● cpu.user.percent: CPU のユーザーモード率 ● cpu.kernel.percent: CPU のカーネル率 ● cpu.idle.percent: CPU のアイドル率 ● cpu.iowait.percent: CPU の I/O 待機率 ● cpu.percent: 汎用 CPU の使用率 <p>例: weight_setting=cpu.user.time=1.0</p>
metrics/required	Boolean	<p>このパラメーターを使用して、設定した metrics/weight_setting メトリックが使用できない場合の処理方法を指定します。</p> <ul style="list-style-type: none"> ● True: メトリックは必須です。メトリックが使用できない場合には、例外が発生します。この例外を回避するには、NovaSchedulerEnabledFilters の MetricsFilter フィルターを使用します。 ● False: 使用できないメトリックは、重み付け処理において負の係数として扱われます。weight_of_unavailable 設定オプションを使用して、戻り値を設定します。
metrics/weight_of_unavailable	浮動小数点	<p>このパラメーターを使用して、metrics/weight_setting メトリックが使用できず、かつ metrics/required=False の場合に用いる重みを指定します。</p> <p>デフォルト: -10000.0</p>

第9章 インスタンス起動用のフレーバーの作成

インスタンスのフレーバーは、インスタンス用の仮想ハードウェアプロファイルを指定するリソースのテンプレートです。クラウドユーザーは、インスタンスを起動する際にフレーバーを指定する必要があります。

フレーバーにより、Compute サービスがインスタンスに割り当てる必要のある以下のリソースの量を指定することができます。

- 仮想 CPU の数
- RAM (MB 単位)
- ルートディスク (GB 単位)
- セカンダリー一時ストレージおよびスワップディスクを含む仮想ストレージ

フレーバーを全プロジェクトに公開したり、特定のプロジェクトまたはドメインを対象にプライベートに設定したりすることで、フレーバーを使用できるユーザーを指定することができます。

フレーバーでは、メタデータ (追加スペックとも呼ばれる) を使用して、インスタンス用ハードウェアのサポートおよびクォータを指定することができます。フレーバーのメタデータは、インスタンスの配置、リソースの使用上限、およびパフォーマンスに影響を及ぼします。利用可能なメタデータ属性の完全なリストは、[Flavor metadata](#) を参照してください。

フレーバーメタデータキーを使用することで、ホストアグリゲートで設定した **extra_specs** メタデータを照合して、インスタンスをホストするのに適したホストアグリゲートを探すこともできます。ホストアグリゲートでインスタンスをスケジュールするには、**extra_specs** キーに **aggregate_instance_extra_specs:** 名前空間の接頭辞を指定して、フレーバーのメタデータのスコープを定義する必要があります。詳細は、[Creating and managing host aggregates](#) を参照してください。

Red Hat OpenStack Platform (RHOSP) のデプロイメントには、クラウドユーザーが使用可能な以下のデフォルトパブリックフレーバーのセットが含まれています。

表9.1 デフォルトのフレーバー

名前	仮想 CPU の数	メモリー	ルートディスクのサイズ
m1.nano	1	128 MB	1 GB
m1.micro	1	192 MB	1 GB



注記

フレーバー属性を使用して設定した動作は、イメージを使用して設定した動作よりも優先されます。クラウドユーザーがインスタンスを起動する際、指定したフレーバーの属性がイメージの属性よりも優先されます。

9.1. フレーバーの作成

特定の機能や動作のために特化したフレーバーを作成および管理することができます。以下に例を示します。

- 基になるハードウェアの要件に応じて、デフォルトのメモリーと容量を変更する

- インスタンスに特定の I/O レートを強制するためのメタデータ、またはホストアグリゲートと一致させるためのメタデータを追加する

手順

1. インスタンスが利用可能な基本的なリソースを指定するフレーバーを作成します。

```
(overcloud)$ openstack flavor create --ram <size_mb> \
--disk <size_gb> --vcpus <no_vcpus> \
[--private --project <project_id>] <flavor_name>
```

- **<size_mb>** を、このフレーバーで作成するインスタンスに割り当てる RAM の容量に置き換えます。
- **<size_gb>** を、このフレーバーで作成するインスタンスに割り当てるルートディスクのサイズに置き換えます。
- **<no_vcpus>** を、このフレーバーで作成するインスタンスに確保する仮想 CPU の数に置き換えます。
- (オプション) **--private** および **--project** オプションを指定して、特定のプロジェクトまたはユーザーグループだけがフレーバーにアクセスできるようにします。**<project_id>** を、このフレーバーを使用してインスタンスを作成できるプロジェクトの ID に置き換えます。アクセシビリティを指定しない場合は、フレーバーはデフォルトで public に設定されます。つまり、すべてのプロジェクトで使用可能になります。



注記

パブリックフレーバーの作成後は、そのフレーバーをプライベートにすることはできません。

- **<flavor_name>** を、一意のフレーバー名に置き換えます。
フレーバーの引数についての詳細は、[Flavor arguments](#) を参照してください。
2. (オプション) フレーバーのメタデータを指定するには、キー/値のペアを使用して必要な属性を設定します。

```
(overcloud)$ openstack flavor set \
--property <key=value> --property <key=value> ... <flavor_name>
```

- **<key>** を、このフレーバーで作成するインスタンスに割り当てる属性のメタデータキーに置き換えます。利用可能なメタデータキーの一覧は、[Flavor metadata](#) を参照してください。
- **<value>** を、このフレーバーで作成するインスタンスに割り当てるメタデータキーの値に置き換えます。
- **<flavor_name>** を、フレーバーの名前に置き換えます。
たとえば、以下のフレーバーを使用して起動されるインスタンスには 2 つの CPU ソケットがあり、それぞれに 2 つの CPU が割り当てられます。

```
(overcloud)$ openstack flavor set \
--property hw:cpu_sockets=2 \
--property hw:cpu_cores=2 processor_topology_flavor
```

9.2. フレーバーの引数

openstack flavor create コマンドには、新しいフレーバーの名前を指定する1つの位置引数 **<flavor_name>** を使用することができます。

以下の表には、新規フレーバーを作成する際に必要に応じて指定することのできる、オプションの引数の詳細をまとめています。

表9.2 オプションのフレーバー引数

オプションの引数	説明
--id	一意のフレーバー ID。デフォルト値は auto で、UUID4 値を生成します。この引数を使用して、整数値を手動で指定することや、UUID4 値を生成することができます。
--ram	(必須) インスタンスが利用可能なメモリーの容量 (MB 単位) デフォルト: 256 MB
--disk	<p>(必須) ルート (/) パーティションに使用するディスク容量 (GB 単位)。ルートディスクは、ベースイメージがコピーされる一時ディスクです。インスタンスが永続ボリュームからブートする場合、ルートディスクは使用されません。</p> <div>  <div> <p>注記</p> <p>--disk が 0 に設定されたフレーバーでインスタンスを作成する場合、インスタンスはボリュームからブートする必要があります。</p> </div> </div> <p>デフォルト: 0 GB</p>
--ephemeral	<p>一時ディスクに使用するディスク容量 (GB 単位) デフォルト値は 0 GB で、セカンダリー一時ディスクは作成されません。一時ディスクは、インスタンスのライフサイクルにリンクしたマシンのローカルディスクストレージを提供します。一時ディスクは、いずれのスナップショットにも含まれません。インスタンスが削除されると、このディスクは破棄されすべてのデータが失われます。</p> <p>デフォルト: 0 GB</p>
--swap	<p>スワップディスクのサイズ (MB 単位) Compute サービスのバックエンドストレージがローカルストレージでない場合は、フレーバーに swap を指定しないでください。</p> <p>デフォルト: 0 GB</p>
--vcpus	<p>(必須) インスタンス用の仮想 CPU の数</p> <p>デフォルト: 1</p>

オプションの引数	説明
--public	フレーバーは、すべてのプロジェクトで利用可能です。デフォルトでは、フレーバーはパブリックですべてのプロジェクトで利用することができます。
--private	フレーバーは、 --project オプションで指定したプロジェクトでのみ利用可能です。プライベートのフレーバーを作成し、フレーバーにプロジェクトを追加しない場合、クラウド管理者だけがそのフレーバーを利用することができます。
--property	以下の形式のキー/値のペアで指定されるメタデータまたは追加スペック --property <key=value> 複数の属性を設定するには、このオプションを繰り返します。
--project	プライベートのフレーバーを使用することができるプロジェクトを指定します。この引数は、 --private オプションと共に使用する必要があります。プロジェクトを指定しないと、フレーバーは admin ユーザーだけに表示されます。 複数のプロジェクトにアクセスを許可するには、このオプションを繰り返します。
--project-domain	プライベートのフレーバーを使用することができるプロジェクトドメインを指定します。この引数は、 --private オプションと共に使用する必要があります。 複数のプロジェクトドメインにアクセスを許可するには、このオプションを繰り返します。
--description	フレーバーの説明。長さは 65535 文字に制限されます。使用できるのは印刷可能文字だけです。

9.3. フレーバーのメタデータ

フレーバーを作成する際に、**--property** オプションを使用してフレーバーのメタデータを指定します。フレーバーのメタデータは **追加スペック** と呼ばれます。フレーバーのメタデータで指定するインスタンス用ハードウェアのサポートおよびクォータは、インスタンスの配置、インスタンスの制限、およびパフォーマンスに影響を及ぼします。

インスタンスによるリソースの使用

以下の表に示す属性キーを使用して、インスタンスによる CPU、メモリー、およびディスク I/O の使用に制限を設定します。



注記

インスタンスの CPU リソース使用量を制限するための追加の仕様は、libvirt に直接渡されるホスト固有の調整可能なプロパティであり、libvirt は制限をホスト OS に渡します。したがって、サポートされるインスタンスの CPU リソース制限の設定は、基盤となるホスト OS によって異なります。

RHOSP デプロイメントでコンピュートノードのインスタンス CPU リソース使用を設定する方法の詳細については、RHEL 9 ドキュメントの [cgroup について](#) および Libvirt ドキュメントの [CPU チューニング](#) を参照してください。

表9.3 リソースの使用を制限するためのフレーバーメタデータ

キー	説明
quota:cpu_shares	ドメイン内の CPU 時間の配分に使用する重みを指定します。デフォルトは OS が提供するデフォルト値です。Compute スケジューラーは、この属性の設定と同じドメイン内にある他のインスタンスの設定を比較して、相対的な重み付けを行います。たとえば、設定が quota:cpu_shares=2048 のインスタンスには、設定が quota:cpu_shares=1024 のインスタンスの 2 倍の CPU 時間が割り当てられます。
quota:cpu_period	cpu_quota を適用する期間を指定します (マイクロ秒単位)。この cpu_period の期間、各仮想 CPU は cpu_quota を超えるランタイムを使用することはできません。1000 - 10000000 の範囲で値を設定します。無効にするには 0 に設定します。
quota:cpu_quota	<p>各 cpu_period における仮想 CPU の最大許容帯域幅を指定します (マイクロ秒単位)。</p> <ul style="list-style-type: none"> 1000 - 18446744073709551 の範囲で値を設定します。 無効にするには 0 に設定します。 帯域幅に制限を設けない場合は、負の値に設定します。 <p>cpu_quota および cpu_period を使用して、全仮想 CPU が同じ速度で実行されるようにすることができます。たとえば、以下のフレーバーを使用して、物理 CPU の処理能力の最大 50% しか消費できないインスタンスを起動することができます。</p> <pre>\$ openstack flavor set cpu_limits_flavor \ --property quota:cpu_quota=10000 \ --property quota:cpu_period=20000</pre>

インスタンスディスクのチューニング

以下の表に示す属性キーを使用して、インスタンスのディスクのパフォーマンスをチューニングします。



注記

Compute サービスは、以下の QoS 設定を、Compute サービスがプロビジョニングしたストレージ (一時ストレージなど) に適用します。Block Storage (cinder) ボリュームのパフォーマンスを調整するには、ボリュームタイプのサービス品質 (QoS) 仕様を設定して関連付ける必要もあります。詳細は、[ストレージガイドの Block Storage サービス \(cinder\) のサービス品質の仕様](#) を参照してください。

表9.4 ディスクチューニング用のフレーバーメタデータ

キー	説明
quota:disk_read_bytes_sec	インスタンスが利用可能な最大ディスク読み取りを指定します (バイト毎秒単位)。
quota:disk_read_iops_sec	インスタンスが利用可能な最大ディスク読み取りを指定します (IOPS 単位)。
quota:disk_write_bytes_sec	インスタンスが利用可能な最大ディスク書き込みを指定します (バイト毎秒単位)。
quota:disk_write_iops_sec	インスタンスが利用可能な最大ディスク書き込みを指定します (IOPS 単位)。
quota:disk_total_bytes_sec	インスタンスが利用可能な最大 I/O 操作を指定します (バイト毎秒単位)。
quota:disk_total_iops_sec	インスタンスが利用可能な最大 I/O 操作を指定します (IOPS 単位)。

インスタンスのネットワークトラフィックの帯域幅

以下の表に示す属性キーを使用して、VIF I/O オプションの設定により、インスタンスのネットワークトラフィックの帯域幅上限を設定します。



注記

quota:vif_* 属性は非推奨になりました。この属性の代わりに、Networking (neutron) サービスの Quality of Service (QoS) ポリシーを使用する必要があります。QoS ポリシーの詳細については、[ネットワーキングガイドの Quality of Service \(QoS\) ポリシーの設定](#) を参照してください。**quota:vif_*** プロパティは、**NeutronOVSEnvironmentDriver** が **iptables_hybrid** に設定されている ML2/OVS メカニズムドライバーを使用する場合にのみサポートされます。

表9.5 帯域幅を制限するためのフレーバーメタデータ

キー	説明
quota:vif_inbound_average	(非推奨) インスタンスに送付されるトラフィックに要求される平均ビットレートを指定します (kbps 単位)。

キー	説明
quota:vif_inbound_burst	(非推奨) ピークの速度でバースト処理することができる受信トラフィックの最大量を指定します (KB 単位)。
quota:vif_inbound_peak	(非推奨) インスタンスが受信することのできるトラフィックの最大レートを指定します (kbps 単位)。
quota:vif_outbound_average	(非推奨) インスタンスから送信されるトラフィックに要求される平均ビットレートを指定します (kbps 単位)。
quota:vif_outbound_burst	(非推奨) ピークの速度でバースト処理することができる送信トラフィックの最大量を指定します (KB 単位)。
quota:vif_outbound_peak	(非推奨) インスタンスが送信することのできるトラフィックの最大レートを指定します (kbps 単位)。

ハードウェアビデオ RAM

以下の表に示す属性キーを使用して、ビデオデバイスに使用するインスタンス RAM の上限を設定します。

表9.6 ビデオデバイス用のフレーバーメタデータ

キー	説明
hw_video:ram_max_mb	ビデオデバイスに使用する最大 RAM を指定します (MB 単位)。 hw_video_ram イメージ属性で使用します。 hw_video_ram は hw_video:ram_max_mb 以下でなければなりません。

ウォッチドッグの動作

以下の表に示す属性キーを使用して、インスタンスで仮想ハードウェアのウォッチドッグデバイスを有効にします。

表9.7 ウォッチドッグの動作を設定するためのフレーバーメタデータ

キー	説明
----	----

キー	説明
hw:watchdog_action	<p>仮想ハードウェアのウォッチドッグデバイスを有効にするかどうかを指定し、その動作を設定します。インスタンスがハングアップまたはエラーが発生した場合に、ウォッチドッグデバイスは設定されたアクションを実行します。ウォッチドッグは i6300esb デバイスを使用し、PCI Intel 6300ESB をエミュレートします。hw:watchdog_action が指定されていない場合には、ウォッチドッグは無効になります。</p> <p>以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● disabled: (デフォルト) デバイスは接続されません。 ● reset: インスタンスを強制的にリセットします。 ● poweroff: インスタンスを強制的にシャットダウンします。 ● pause: インスタンスを一時停止します。 ● none: ウォッチドッグを有効にしますが、インスタンスがハングアップまたはエラーが発生しても何も実行しません。 <div>  <div> <p>注記</p> <p>特定のイメージの属性を使用して設定するウォッチドッグの動作は、フレーバーを使用して設定する動作よりも優先されます。</p> </div> </div>

乱数ジェネレーター (RNG)

以下の表に示す属性キーを使用して、インスタンスで RNG デバイスを有効にします。

表9.8 RNG 用のフレーバーメタデータ

キー	説明
hw_rng:allowed	<p>イメージ属性によりインスタンスに追加された RNG デバイスを無効にするには、False に設定します。</p> <p>デフォルト: True</p>
hw_rng:rate_bytes	<p>期間中、インスタンスがホストのエントロピーから読み取ることのできる最大バイト数を指定します。</p>
hw_rng:rate_period	<p>読み取り期間を指定します (ミリ秒単位)。</p>

仮想パフォーマンス監視ユニット (vPMU)

以下の表に示す属性キーを使用して、インスタンスの仮想 PMU を有効にします。

表9.9 仮想 PMU 用のフレーバーメタデータ

キー	説明
hw:pmu	<p>インスタンスの仮想 PMU を有効にするには、True に設定します。</p> <p>perf 等のツールは、インスタンスの仮想 PMU を使用して、インスタンスのパフォーマンスをプロファイリングおよびモニターリングするためのより正確な情報を提供します。リアルタイム負荷のケースでは、仮想 PMU のエミュレーションにより、望ましくないレイテンシーが付加される場合があります。テレメトリーの提供が不要な場合は、hw:pmu=False に設定します。</p>

インスタンスの CPU トポロジー

以下の表に示す属性キーを使用して、インスタンス内のプロセッサのトポロジーを定義します。

表9.10 CPU トポロジー用のフレーバーメタデータ

キー	説明
hw:cpu_sockets	<p>インスタンスの推奨ソケット数を指定します。</p> <p>デフォルト: 要求される仮想 CPU の数</p>
hw:cpu_cores	<p>インスタンスの1ソケットあたりの推奨コア数を指定します。</p> <p>デフォルト: 1</p>
hw:cpu_threads	<p>インスタンスの1コアあたりの推奨スレッド数を指定します。</p> <p>デフォルト: 1</p>
hw:cpu_max_sockets	<p>イメージ属性を使用してユーザーがインスタンスに選択できる最大ソケット数を指定します。</p> <p>例: hw:cpu_max_sockets=2</p>
hw:cpu_max_cores	<p>イメージ属性を使用してユーザーがインスタンスに選択できる、1ソケットあたりの最大コア数を指定します。</p>
hw:cpu_max_threads	<p>イメージ属性を使用してユーザーがインスタンスに選択できる、1コアあたりの最大スレッド数を指定します。</p>

シリアルポート

以下の表に示す属性キーを使用して、1インスタンスあたりのシリアルポートの数を設定します。

表9.11 シリアルポート用のフレーバーメタデータ

キー	説明
<code>hw:serial_port_count</code>	1 インスタンスあたりの最大シリアルポート数

CPU ピニングポリシー

デフォルトでは、インスタンスの仮想 CPU (vCPU) は 1 コア 1 スレッドのソケットです。属性を使用して、インスタンスの仮想 CPU をホストの物理 CPU コア (pCPU) に固定するフレーバーを作成することができます。1 つまたは複数のコアがスレッドシブリングを持つ同時マルチスレッド (SMT) アーキテクチャーで、ハードウェア CPU スレッドの動作を設定することもできます。

以下の表に示す属性キーを使用して、インスタンスの CPU ピニングポリシーを定義します。

表9.12 CPU ピニング用のフレーバーメタデータ

キー	説明
<code>hw:cpu_policy</code>	<p>使用する CPU ポリシーを指定します。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● shared: (デフォルト) インスタンスの仮想 CPU は、ホストの物理 CPU 全体で共有されます。 ● dedicated: インスタンスの仮想 CPU をホストの物理 CPU のセットに固定します。これにより、インスタンスが固定される CPU のトポロジに一致するインスタンス CPU のトポロジが作成されます。このオプションの場合、必然的にオーバーコミット比は 1.0 になります。

キー	説明
hw:cpu_thread_policy	<p>設定が hw:cpu_policy=dedicated の場合に使用する CPU スレッドポリシーを指定します。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● prefer: (デフォルト) ホストは SMT アーキテクチャーを持つ場合と持たない場合があります。SMT アーキテクチャーが存在する場合、Compute スケジューラーはスレッドシブリングを優先します。 ● isolate: ホストは SMT アーキテクチャーを持たないか、あるいは SMT 以外のアーキテクチャーをエミュレートする必要があります。このポリシーにより、Compute スケジューラーは HW_CPU_HYPERTHREADING 特性が設定されていないホストを要求して、SMT を持たないホストにインスタンスを配置するようになります。以下の属性を使用して、この特性を明示的に要求することもできます。 <pre>--property trait:HW_CPU_HYPERTHREADING=forbidden</pre> <p>ホストが SMT アーキテクチャーを持たない場合、Compute サービスはそれぞれの仮想 CPU を想定どおりに異なるコアに配置します。ホストが SMT アーキテクチャーを持つ場合は、動作は [workarounds]/disable_fallback_pcpu_query パラメーターの設定により決定されます。</p> <ul style="list-style-type: none"> ○ True: SMT アーキテクチャーを持つホストは使用されず、スケジューリングに失敗します。 ○ False: Compute サービスはそれぞれの仮想 CPU を異なる物理コアに配置します。Compute サービスは、別のインスタンスからの仮想 CPU を同じコア上に配置しません。したがって、使用される各コアのスレッドシブリングは、1つを除きすべて使用できなくなります。 ● require: ホストは SMT アーキテクチャーを持つ必要があります。このポリシーにより、Compute スケジューラーは HW_CPU_HYPERTHREADING 特性が設定されたホストを要求して、SMT を持つホストにインスタンスを配置するようになります。以下の属性を使用して、この特性を明示的に要求することもできます。 <pre>--property trait:HW_CPU_HYPERTHREADING=required</pre> <p>Compute サービスは、それぞれの仮想 CPU をスレッドシブリングに割り当てます。ホストが SMT アーキテクチャーを持たない場合、そのホストは使用されません。ホストは SMT アーキテクチャーを持つが、スレッドシブリングが使用されていないコアが十分でない場合、スケジューリングに失敗します。</p>

インスタンス PCI NUMA アフィニティーポリシー

以下の表に示す属性キーを使用して、PCI パススルーデバイスおよび SR-IOV インターフェイスの NUMA アフィニティーポリシーを指定するフレーバーを作成します。

表9.13 PCI NUMA アフィニティーポリシー用のフレーバーメタデータ

キー	説明
hw:pci_numa_affinity_policy	<p>PCI パススルーデバイスおよび SR-IOV インターフェイスの NUMA アフィニティーポリシーを指定します。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● required: インスタンスの NUMA ノードの少なくとも 1 つが PCI デバイスとのアフィニティーを持つ場合に限り、Compute サービスは PCI デバイスを要求するインスタンスを作成します。このオプションは、最高のパフォーマンスを提供します。 ● preferred: Compute サービスは、NUMA アフィニティーに基づきベストエフォートで PCI デバイスの選択を試みます。これができない場合には、Compute サービスは PCI デバイスとのアフィニティーを持たない NUMA ノード上でインスタンスをスケジュールします。 ● legacy: (デフォルト) 以下のどちらかのケースで、Compute サービスは PCI デバイスを要求するインスタンスを作成します。 <ul style="list-style-type: none"> ○ PCI デバイスが少なくとも 1 つの NUMA ノードとのアフィニティーを持つ。 ○ PCI デバイスが NUMA アフィニティーに関する情報を提供しない。

インスタンスの NUMA トポロジー

属性を使用して、インスタンスの仮想 CPU スレッドのホスト NUMA 配置、ならびにホスト NUMA ノードからのインスタンスの仮想 CPU およびメモリーの割り当てを定義するフレーバーを作成することができます。

メモリーおよび仮想 CPU の割り当てがコンピュータホスト内の NUMA ノードのサイズよりも大きいフレーバーの場合、インスタンスの NUMA トポロジーを定義するとインスタンス OS のパフォーマンスが向上します。

Compute スケジューラーは、これらの属性を使用してインスタンスに適したホストを決定します。たとえば、クラウドユーザーは以下のフレーバーを使用してインスタンスを起動します。

```
$ openstack flavor set numa_top_flavor \
--property hw:numa_nodes=2 \
--property hw:numa_cpus.0=0,1,2,3,4,5 \
--property hw:numa_cpus.1=6,7 \
--property hw:numa_mem.0=3072 \
--property hw:numa_mem.1=1024
```

Compute スケジューラーは、2 つの NUMA ノード (1 つは 3 GB の RAM を持ち 6 つの CPU を実行できるノード、もう 1 つは 1 GB の RAM を持ち 2 つの CPU を実行できるノード) を持つホストを探しま

す。4 GB の RAM を持ち 8 つの CPU を実行できる 1 つの NUMA ノードを持つホストの場合、Compute スケジューラーは有効な一致とは見なしません。



注記

フレーバーで定義された NUMA トポロジは、イメージで定義された NUMA トポロジでオーバーライドされることはありません。イメージの NUMA トポロジがフレーバーの NUMA トポロジと競合する場合、Compute サービスは **ImageNUMATopologyForbidden** エラーを報告します。

注意

この機能を使用して、インスタンスを特定のホスト CPU または NUMA ノードに制限することはできません。包括的なテストおよびパフォーマンス計測が完了した後にのみ、この機能を使用してください。代わりに **hw:pci_numa_affinity_policy** プロパティを使用することができます。

以下の表に示す属性キーを使用して、インスタンスの NUMA トポロジを定義します。

表9.14 NUMA トポロジ用のフレーバーメタデータ

キー	説明
hw:numa_nodes	<p>インスタンスの仮想 CPU スレッドの実行先として制限するホスト NUMA ノードの数を指定します。指定しない場合は、利用可能な任意の数のホスト NUMA ノード上で仮想 CPU スレッドを実行することができます。</p>
hw:numa_cpus.N	<p>インスタンスの NUMA ノード N にマッピングするインスタンスの仮想 CPU のコンマ区切りリスト。このキーを指定しない場合、仮想 CPU は利用可能な NUMA ノード間で均等に配分されます。</p> <p>N は 0 から始まります。*.N 値を使用する場合には注意が必要です。NUMA ノードが少なくとも 2 つある場合に限り使用してください。</p> <p>この属性は hw:numa_nodes を設定している場合にのみ有効で、インスタンスの NUMA ノードに CPU および RAM が対称的に割り当てられていない場合 (一部の NFV 負荷で重要) にのみ必要です。</p>
hw:numa_mem.N	<p>インスタンスの NUMA ノード N にマッピングするインスタンスのメモリー容量 (MB 単位)。このキーを指定しない場合、メモリーは利用可能な NUMA ノード間で均等に配分されます。</p> <p>N は 0 から始まります。*.N 値を使用する場合には注意が必要です。NUMA ノードが少なくとも 2 つある場合に限り使用してください。</p> <p>この属性は hw:numa_nodes を設定している場合にのみ有効で、インスタンスの NUMA ノードに CPU および RAM が対称的に割り当てられていない場合 (一部の NFV 負荷で重要) にのみ必要です。</p>

**警告**

hw:numa_cpus.N で指定する仮想 CPU の総数または **hw:numa_mem.N** で指定するメモリー容量が、それぞれ利用可能な CPU の数またはメモリー容量よりも大きい場合、Compute サービスは例外を発生させます。

CPU リアルタイムポリシー

以下の表に示す属性キーを使用して、インスタンス内のプロセッサのリアルタイムポリシーを定義します。

**注記**

- インスタンスのほとんどの仮想 CPU は、リアルタイムポリシーを設定して実行することができますが、非リアルタイムのゲストプロセスとエミュレーターのオーバーヘッドプロセスの両方に使用するために、少なくとも1つの仮想 CPU を非リアルタイムと識別する必要があります。
- この追加スペックを使用するには、ピンングされた CPU を有効にする必要があります。

表9.15 CPU リアルタイムポリシー用のフレーバーメタデータ

キー	説明
hw:cpu_realtime	<p>リアルタイムポリシーをインスタンスの仮想 CPU に割り当てるフレーバーを作成するには、yes に設定します。</p> <p>デフォルト: no</p>
hw:cpu_realtime_mask	<p>リアルタイムポリシーを割り当てない仮想 CPU を指定します。マスクする値の前にキャレット記号 (^) を追加する必要があります。仮想 CPU 0 および 1 を除くすべての仮想 CPU にリアルタイムポリシーを設定する場合の例を以下に示します。</p> <pre>\$ openstack flavor set <flavor> \ --property hw:cpu_realtime="yes" \ --property hw:cpu_realtime_mask="^0-1"</pre> <div>  <p>注記</p> <p>hw_cpu_realtime_mask 属性がイメージで設定されている場合、フレーバーで設定した hw:cpu_realtime_mask 属性よりも優先されます。</p> </div>

エミュレータースレッドポリシー

物理 CPU をインスタンスに割り当てて、エミュレータースレッドに使用することができます。エミュレータースレッドとは、インスタンスと直接関係しないエミュレータープロセスを指します。リアルタイム負荷には、専用のエミュレータースレッド用物理 CPU が必要です。エミュレータースレッドポリシーを使用するには、以下の属性を設定してピンングされた CPU を有効にする必要があります。

```
--property hw:cpu_policy=dedicated
```

以下の表に示す属性キーを使用して、インスタンスのエミュレータースレッドポリシーを定義します。

表9.16 エミュレータースレッドポリシー用のフレーバーメタデータ

キー	説明
hw:emulator_threads_policy	<p>インスタンスに使用するエミュレータースレッドポリシーを指定します。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● share: エミュレータースレッドは、NovaComputeCpuSharedSet heat パラメーターで定義される物理 CPU 全体で共有されます。NovaComputeCpuSharedSet が設定されていない場合、エミュレータースレッドはインスタンスに関連付けられたピンングされた CPU 全体で共有されます。 ● isolate: エミュレータースレッド用に、インスタンスごとに専用の物理 CPU をさらに確保します。このポリシーは過度にリソースを消費するので、使用には注意が必要です。 ● unset: (デフォルト) エミュレータースレッドポリシーは有効ではなく、エミュレータースレッドはインスタンスに関連付けられたピンングされた CPU 全体で共有されます。

インスタンスのメモリーページサイズ

以下の表に示す属性キーを使用して、明示的なメモリーページサイズでインスタンスを作成します。

表9.17 メモリーページサイズ用のフレーバーメタデータ

キー	説明
----	----

キー	説明
hw:mem_page_size	<p>インスタンスをサポートするのに使用するラージページのサイズを指定します。このオプションを使用すると、hw:numa_nodes で特に指定しない限り、1 NUMA ノードの暗黙的な NUMA トポロジーが作成されます。以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● large: ホストでサポートされる最小のページサイズより大きいページサイズを選択します。x86_64 システムでは 2 MB または 1 GB です。 ● small: ホストでサポートされる最小のページサイズを選択します。X86_64 システムでは、4 kB (通常のページ) です。 ● any: libvirt ドライバーで決定される、利用可能な最大のヒュージページサイズを選択します。 ● <pagesize>: (文字列) ワークロードに具体的な要件がある場合、ページサイズを明示的に設定します。ページサイズには整数値を使用し、kB またはその他の標準単位で指定します。(例: 4KB、2MB、2048、1GB)。 ● unset: (デフォルト) インスタンスをサポートするのにラージページは使用されず、暗黙的な NUMA トポロジーは生成されません。

PCI パススルー

以下の表に示す属性キーを使用して、グラフィックカードやネットワークデバイス等の物理 PCI デバイスをインスタンスにアタッチします。PCI パススルーの使用に関する詳細は、[Configuring PCI passthrough](#) を参照してください。

表9.18 PCI パススルー用のフレーバーメタデータ

キー	説明
pci_passthrough:alias	<p>以下の形式を使用して、インスタンスに割り当てる PCI デバイスを指定します。</p> <pre><alias>:<count></pre> <ul style="list-style-type: none"> ● <alias> を、特定の PCI デバイスクラスに対応するエイリアスに置き換えます。 ● <count> を、インスタンスに割り当てる種別 <alias> の PCI デバイスの数に置き換えます。

ハイパーバイザーの署名

以下の表に示す属性キーを使用して、ハイパーバイザーの署名をインスタンスからは非表示にします。

表9.19 ハイパーバイザーの署名を非表示にするためのフレーバーメタデータ

キー	説明
<code>hide_hypervisor_id</code>	ハイパーバイザーの署名をインスタンスからは非表示にするには、 True に設定します。これにより、すべてのドライバーがインスタンスで読み込みおよび操作を行うことができます。

UEFI セキュアブート

次の表のプロパティキーを使用して、UEFI セキュアブートで保護されたインスタンスを作成します。



注記

UEFI セキュアブートを使用するインスタンスは、UEFI および GUID パーティションテーブル (GPT) 標準をサポートし、EFI システムパーティションを含める必要があります。

表9.20 UEFI セキュアブートのフレーバーメタデータ

キー	説明
<code>os:secure_boot</code>	このフレーバーで起動されたインスタンスのセキュアブートを有効にするには、 required に設定します。デフォルトでは無効にされています。

インスタンスのリソース特性

各リソースプロバイダーには特性のセットがあります。特性は、ストレージディスクの種別や Intel CPU 拡張命令セットなど、リソースプロバイダーの機能的な要素です。インスタンスは、これらの中から要求する特性を指定することができます。

指定することのできる特性は **os-traits** ライブラリーで定義されます。特性の例を以下に示します。

- **COMPUTE_TRUSTED_CERTS**
- **COMPUTE_NET_ATTACH_INTERFACE_WITH_TAG**
- **COMPUTE_IMAGE_TYPE_RAW**
- **HW_CPU_X86_AVX**
- **HW_CPU_X86_AVX512VL**
- **HW_CPU_X86_AVX512CD**

os-traits ライブラリーの使用方法の詳細は、[Usage](#) を参照してください。

以下の表に示す属性キーを使用して、インスタンスのリソース特性を定義します。

表9.21 リソース特性用のフレーバーメタデータ

キー	説明
trait:<trait_name>	<p>コンピュータノードの特性を指定します。特性を、以下の有効な値のいずれかに設定します。</p> <ul style="list-style-type: none"> ● required: インスタンスをホストするために選択したコンピュータノードになければいけない特性 ● forbidden: インスタンスをホストするために選択したコンピュータノードにあってはいけない特性 <p>例:</p> <pre>\$ openstack flavor set --property trait:HW_CPU_X86_AVX512BW=required avx512- flavor</pre>

インスタンスのベアメタルリソースクラス

以下の表に示す属性キーを使用して、インスタンスのベアメタルリソースクラスを要求します。

表9.22 ベアメタルリソースクラス用のフレーバーメタデータ

キー	説明
resources:<resource_class_name>	<p>この属性を使用して、値をオーバーライドする標準のベアメタルリソースクラスを指定するか、インスタンスが要求するカスタムのベアメタルリソースクラスを指定します。</p> <p>オーバーライドすることができる標準のリソースクラスは VCPU、MEMORY_MB、および DISK_GB です。Compute スケジューラーがインスタンスのスケジューリングにベアメタルフレーバー属性を使用するのを防ぐには、標準のリソースクラスの値を 0 に設定します。</p> <p>カスタムリソースクラスの名前は、CUSTOM_ で始まる必要があります。Bare Metal サービスノードのリソースクラスに対応するカスタムリソースクラスの名前を指定するには、リソースクラスを大文字に変換し、すべての句読点をアンダースコアに置き換え、CUSTOM_ の接頭辞を追加します。</p> <p>たとえば、--resource-class baremetal.SMALL のノードにインスタンスをスケジューリングするには、以下のフレーバーを作成します。</p> <pre>\$ openstack flavor set \ --property resources:CUSTOM_BAREMETAL_SMALL=1 \ --property resources:VCPU=0 --property resources:MEMORY_MB=0 \ --property resources:DISK_GB=0 compute-small</pre>

第10章 インスタンスへのメタデータの追加

Compute (nova) サービスは、メタデータを使用してインスタンスの起動時に設定情報を渡します。インスタンスは、コンフィグドライブまたはメタデータサービスを使用してメタデータにアクセスすることができます。

コンフィグドライブ

コンフィグドライブは、インスタンスのブート時にアタッチすることのできる特別なドライブです。コンフィグドライブは読み取り専用ドライブとしてインスタンスに提示されます。インスタンスはこのドライブをマウントしてそこからファイルを読み取り、通常メタデータサービスから利用する情報を取得することができます。

メタデータサービス

Compute サービスは、REST API としてメタデータサービスを提供します。これを使用して、インスタンス固有のデータを取得することができます。インスタンスは、**169.254.169.254** または **fe80::a9fe:a9fe** からこのサービスにアクセスします。

10.1. インスタンスメタデータの種別

クラウドユーザー、クラウド管理者、および Compute サービスは、メタデータをインスタンスに渡すことができます。

クラウドユーザーが提供するデータ

クラウドユーザーは、インスタンスがブート時に実行するシェルスクリプトなど、インスタンスを起動する際に使用する追加データを指定することができます。クラウドユーザーは、インスタンスを作成または更新する際に、ユーザーデータ機能を使用し、キー/値のペアを必要な属性として渡すことで、データをインスタンスに渡すことができます。

クラウド管理者が提供するデータ

RHOSP 管理者は、ベンダーデータ機能を使用してデータをインスタンスに渡します。Compute サービスの提供するベンダーデータモジュール **StaticJSON** および **DynamicJSON** により、管理者はメタデータをインスタンスに渡すことができます。

- **StaticJSON**:(デフォルト) 全インスタンスで共通のメタデータに使用します。
- **DynamicJSON**: 各インスタンスで異なるメタデータに使用します。このモジュールは外部の REST サービスにリクエストを行い、インスタンスに追加するメタデータを決定します。

ベンダーデータの設定は、インスタンスの以下の読み取り専用ファイルのいずれかにあります。

- `/openstack/{version}/vendor_data.json`
- `/openstack/{version}/vendor_data2.json`

Compute サービスが提供するデータ

Compute サービスはメタデータサービスの内部実装を使用して、要求されたインスタンスのホスト名やインスタンスが属するアベイラビリティゾーン等の情報をインスタンスに渡します。この操作はデフォルトで実施され、クラウドユーザーまたは管理者が設定を行う必要はありません。

10.2. 全インスタンスへのコンフィグドライブの追加

管理者は Compute サービスを設定し、常にインスタンス用のコンフィグドライブを作成し、コンフィグドライブにデプロイメント固有のメタデータを設定することができます。たとえば、以下の理由によりコンフィグドライブを使用する場合があります。

- デプロイメントにおいて、インスタンスへの IP アドレスの割り当てに DHCP を使用しない場合に、ネットワーク設定を渡すため。インスタンスのネットワーク設定を行う前に、コンフィグドライブを通じてインスタンスの IP アドレス設定を渡すことができます。インスタンスは、コンフィグドライブをマウントして設定にアクセスすることができます。
- Active Directory ポストブートにインスタンスを登録するのに使用される暗号化トークン等、インスタンスを起動するユーザーがアクセスできないデータをインスタンスに渡すため。
- ローカルにキャッシュされたディスク読み取りを作成し、インスタンスのリクエストの負荷を管理するため。これにより、ファクトのチェックインおよびビルドのために定期的にメタデータサーバーにアクセスするインスタンスの影響が軽減されます。

ISO 9660 または VFAT ファイルシステムをマウントできるインスタンスのオペレーティングシステムは、すべてコンフィグドライブを使用することができます。

手順

1. Compute 環境ファイルを開きます。
2. インスタンスの起動時に常にコンフィグドライブをアタッチするには、以下のパラメーターを **True** に設定します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::compute::force_config_drive: 'true'
```

3. (オプション) コンフィグドライブの形式をデフォルト値の **iso9660** から **vfat** に変更するには、設定に **config_drive_format** パラメーターを追加します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::compute::force_config_drive: 'true'
    nova::compute::config_drive_format: vfat
```

4. 更新内容を Compute 環境ファイルに保存します。
5. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml \
```

検証

1. インスタンスを作成します。

```
(overcloud)$ openstack server create --flavor m1.tiny \
--image cirros test-config-drive-instance
```

2. インスタンスにログインします。
3. コンフィグドライブをマウントします。

- インスタンスの OS が **udev** を使用する場合:

```
# mkdir -p /mnt/config
# mount /dev/disk/by-label/config-2 /mnt/config
```

- インスタンスの OS が **udev** を使用しない場合は、まずコンフィグドライブに対応するブロックデバイスを特定する必要があります。

```
# blkid -t LABEL="config-2" -o device
/dev/vdb
# mkdir -p /mnt/config
# mount /dev/vdb /mnt/config
```

4. マウントされたコンフィグドライブディレクトリー **mnt/config/openstack/{version}/** で、メタデータのファイルを検査します。

10.3. インスタンスへの動的メタデータの追加

デプロイメントを設定してインスタンス固有のメタデータを作成し、そのインスタンスが JSON ファイルを使用してメタデータを利用できるようにすることができます。

ヒント

アンダークラウド上で動的メタデータを使用して、director を Red Hat Identity Management (IdM) サーバーと統合することができます。IdM サーバーは認証局として使用することができます、オーバークラウドで SSL/TLS が有効な場合にオーバークラウドの証明書を管理することができます。詳細は、[セキュリティおよび強化ガイドの Ansible を使用した TLS-e の実装](#) を参照してください。

手順

1. Compute 環境ファイルを開きます。
2. ベンダーデータプロバイダーモジュールに **DynamicJSON** を追加します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::config::nova_config:
    ...
    api/vendordata_providers:
      value: StaticJSON,DynamicJSON
```

3. メタデータを生成するためにアクセスする REST サービスを指定します。必要な数だけ目的の REST サービスを指定することができます。以下に例を示します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::config::nova_config:
    ...
    api/vendordata_providers:
      value: StaticJSON,DynamicJSON
    api/vendordata_dynamic_targets:
      value: target1@http://127.0.0.1:125
    api/vendordata_dynamic_targets:
      value: target2@http://127.0.0.1:126
```

Compute サービスは設定されたターゲットサービスから取得したメタデータが含まれる JSON ファイル **vendordata2.json** を生成し、それをコンフィグドライブディレクトリーに保存します。



注記

ターゲットサービスに同じ名前を複数回使用しないでください。

4. 更新内容を Compute 環境ファイルに保存します。
5. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \  
-e [your environment files] \  
-e /home/stack/templates/<compute_environment_file>.yaml
```

第11章 インスタンスの CPU 機能フラグの設定

ホストコンピュータノードの設定を変更してコンピュータノードをリブートすることなく、インスタンスの CPU 機能フラグを有効または無効にすることができます。インスタンスに適用される CPU 機能フラグの標準的なセットを設定することで、コンピュータノード間でライブマイグレーションの互換性を実現するのに役立ちます。また、特定の CPU モデルにおいてインスタンスのセキュリティーやパフォーマンスに悪影響を与えるフラグを無効にしたり、セキュリティーやパフォーマンスの問題を軽減するフラグを有効にしたりして、インスタンスのパフォーマンスおよびセキュリティーを管理するのにも役立ちます。

11.1. 前提条件

- CPU モデルおよび機能フラグが、ホストコンピュータノードのハードウェアおよびソフトウェアでサポートされる必要があります。
 - ホストが対応しているハードウェアを確認するには、コンピュータノードで以下のコマンドを入力します。

```
$ cat /proc/cpuinfo
```

- ホストが対応している CPU モデルを確認するには、コンピュータノードで以下のコマンドを入力します。

```
$ sudo podman exec -it nova_libvirt virsh cpu-models <arch>
```

<arch> をアーキテクチャーの名前に置き換えます (例: **x86_64**)。

11.2. インスタンスの CPU 機能フラグの設定

Compute サービスを設定し、特定の仮想 CPU モデルのインスタンスに CPU 機能フラグを適用します。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. **stackrc** ファイルを取得します。

```
[stack@director ~]$ source ~/stackrc
```

3. Compute 環境ファイルを開きます。
4. インスタンスの CPU モードを設定します。

```
parameter_defaults:
  ComputeParameters:
    NovaLibvirtCPUMode: <cpu_mode>
```

<cpu_mode> をコンピュータノード上の各インスタンスの CPU モードに置き換えます。以下の有効な値のいずれかに設定します。

- **host-model**:(デフォルト) ホストコンピュートノードの CPU モデルを使用します。この CPU モードを使用して、重要な CPU フラグをインスタンスに自動的に追加し、セキュリティ上の欠陥の軽減策を提供します。
- **custom**: 各インスタンスが使用する特定の CPU モデルを設定するのに使用します。



注記

CPU モードを **host-passthrough** に設定すると、コンピュートノードでホストされるインスタンスにそのコンピュートノードと同じ CPU モデルおよび機能フラグを使用することができます。

5. (オプション) **NovaLibvirtCPUMode** を **custom** に設定した場合は、カスタマイズするインスタンス CPU モデルを設定します。

```
parameter_defaults:
  ComputeParameters:
    NovaLibvirtCPUMode: 'custom'
    NovaLibvirtCPUModels: <cpu_model>
```

<cpu_model> を、ホストがサポートする CPU モデルのコンマ区切りリストに置き換えます。CPU モデルを順に一覧表示します。この際、より一般的で高度ではない CPU モデルは最初に一覧に配置され、より機能が充実した CPU モデルが最後に置かれます (例:

SandyBridge,IvyBridge,Haswell,Broadwell)。モデル名の一覧

は、**/usr/share/libvirt/cpu_map.xml** を参照してください。または、ホストコンピュートノードで以下のコマンドを入力します。

```
$ sudo podman exec -it nova_libvirt virsh cpu-models <arch>
```

<arch> をコンピュートノードのアーキテクチャー名に置き換えてください (例: **x86_64**)。

6. 指定の CPU モデルのインスタンスの CPU 機能フラグを設定します。

```
parameter_defaults:
  ComputeParameters:
    ...
    NovaLibvirtCPUModelExtraFlags: <cpu_feature_flags>
```

<cpu_feature_flags> を、有効または無効にする機能フラグのコンマ区切りリストに置き換えます。フラグを有効にするには各フラグの前に+を付け、無効にするには-を付けます。接頭辞が指定されていない場合、フラグが有効になります。特定の CPU モデルで利用可能な機能フラグの一覧は、**/usr/share/libvirt/cpu_map/*.xml** を参照してください。

以下の例では、**IvyBridge** および **Cascadelake-Server** モデルの CPU 機能フラグ **pcid** および **ssbd** を有効にし、機能フラグ **mtrr** を無効にします。

```
parameter_defaults:
  ComputeParameters:
    NovaLibvirtCPUMode: 'custom'
    NovaLibvirtCPUModels: 'IvyBridge','Cascadelake-Server'
    NovaLibvirtCPUModelExtraFlags: 'pcid,+ssbd,-mtrr'
```

7. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \  
-e [your environment files] \  
-e /home/stack/templates/<compute_environment_file>.yaml
```

第12章 KERNELARGS を定義するための手動でのノード再起動の設定

オーバークラウドのデプロイメントに **KernelArgs** の初回設定が含まれる場合に、オーバークラウドノードは自動的に再起動されます。**KernelArgs** をすでに運用中のデプロイメントに追加する場合には、ノードを再起動することで、既存のワークロードに対して問題となる可能性があります。デプロイメントの更新時にノードの自動再起動を無効にし、代わりに各オーバークラウドのデプロイメント後にノードの再起動を手動で実行できます。



注記

自動再起動を無効にしてから新しいコンピュータノードをデプロイに追加すると、新しいノードは初期プロビジョニング中に再起動されません。**KernelArgs** の設定は再起動後にのみ適用されるため、デプロイメントエラーが発生する可能性があります。

12.1. KERNELARGS を定義するための手動でのノード再起動の設定

KernelArgs を初めて設定するときにノードの自動再起動を無効にし、代わりにノードを手動で再起動できます。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. **stackrc** ファイルを取得します。

```
[stack@director ~]$ source ~/stackrc
```

3. カスタム環境ファイル (**kernelargs_manual_reboot.yaml** など) で **KernelArgsDeferReboot** ロールパラメーターを有効にします。

```
parameter_defaults:
  <Role>Parameters:
    KernelArgsDeferReboot: True
```

4. その他の環境ファイルと共にこれらのカスタム環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/kernelargs_manual_reboot.yaml
```

5. コンピュータノードのリストを取得して、再起動するノードのホスト名を特定します。

```
(undercloud)$ source ~/overcloudrc
(overcloud)$ openstack compute service list
```

6. 再起動するコンピュータノードで Compute サービスを無効にして、Compute スケジューラーが新しいインスタンスをノードに割り当てないようにします。

```
(overcloud)$ openstack compute service set <node> nova-compute --disable
```

<node> を、Compute サービスを無効にするノードのホスト名に置き換えます。

7. 移行するコンピュートノードでホストされているインスタンスのリストを取得します。

```
(overcloud)$ openstack server list --host <node_UUID> --all-projects
```

8. インスタンスを別のコンピュートノードに移行する。インスタンスの移行については、[コンピュートノード間での仮想マシンインスタンスの移行](#) を参照してください。

9. リブートするノードにログインします。

10. ノードをリブートします。

```
[tripleo-admin@overcloud-compute-0 ~]$ sudo reboot
```

11. ノードがブートするまで待ちます。

12. コンピュートノードを再度有効にします。

```
(overcloud)$ openstack compute service set <node_UUID> nova-compute --enable
```

13. コンピュートノードが有効であることを確認します。

```
(overcloud)$ openstack compute service list
```

第13章 インスタンス用の仮想 GPU の設定

重要

この機能は、本リリースでは **テクノロジープレビュー** として提供しているため、Red Hat では全面的にはサポートしていません。これは、テスト用途にのみご利用いただく機能で、実稼働環境にデプロイすべきではありません。テクノロジープレビュー機能についての詳しい情報は、[対象範囲の詳細](#) を参照してください。

インスタンスで GPU ベースのレンダリングをサポートするには、利用可能な物理 GPU デバイスおよびハイパーバイザーの種別に応じて、仮想 GPU (vGPU) リソースを定義し、管理できます。この設定を使用して、レンダリングの負荷をすべての物理 GPU デバイス間でより効果的に分割し、仮想 GPU 対応のインスタンスをスケジューリングする際の制御性を向上させることができます。

Compute (nova) サービスで仮想 GPU を有効にするには、クラウドユーザーが仮想 GPU デバイスの設定された Red Hat Enterprise Linux (RHEL) インスタンスを作成するのに使用できるフレーバーを作成します。これにより、各インスタンスは物理 GPU デバイスに対応する仮想 GPU デバイスで GPU 負荷に対応することができます。

Compute サービスは、各ホストに定義する GPU プロファイルで利用可能な仮想 GPU デバイスの数を追跡します。Compute サービスはフレーバーに基づいてこれらのホストにインスタンスをスケジューリングし、デバイスをアタッチし、使用状況を継続的に監視します。インスタンスが削除されると、Compute サービスは仮想 GPU デバイスを利用可能なプールに戻します。

重要

Red Hat では、サポート例外を要求せずに RHOSP での NVIDIA 仮想 GPU の使用を有効にしています。ただし、Red Hat は、NVIDIA 仮想 GPU ドライバーのテクニカルサポートを提供しません。NVIDIA 仮想 GPU ドライバーは、NVIDIA により提供され、サポートされます。NVIDIA 仮想 GPU ソフトウェアの NVIDIA Enterprise サポートを取得するには、NVIDIA 認定サポートサービスサブスクリプションが必要です。サポートされるコンポーネントで問題を再現できない NVIDIA 仮想 GPU の使用から生じる問題については、以下のサポートポリシーが適用されます。

- サードパーティーコンポーネントが問題に関与していないと Red Hat が考える場合は、通常の [サポート対象範囲](#) および [Red Hat SLA](#) が適用されます。
- サードパーティーコンポーネントが問題に関与していると Red Hat が考える場合は、お客様は Red Hat の [サードパーティーサポートおよび認定ポリシー](#) に従って NVIDIA に問い合わせを依頼されます。詳細は、ナレッジベースの記事 [Obtaining Support from NVIDIA](#) を参照してください。

13.1. サポートされる設定および制限

サポートされる GPU カード

サポートされる NVIDIA GPU カードの一覧については、NVIDIA の Web サイトで [Virtual GPU Software Supported Products](#) を参照してください。

仮想 GPU デバイスを使用する際の制限

- 各コンピュータノードで有効にできる仮想 GPU の種別は1つだけです。
- 各インスタンスが使用できる仮想 GPU のリソースは1つだけです。

- ホスト間の vGPU インスタンスのライブマイグレーションはサポートされていません。
- vGPU インスタンスの退避はサポートされていません。
- 仮想 GPU インスタンスをホストするコンピュータノードをリブートする必要がある場合、仮想 GPU は自動的に再作成されたインスタンスに再割り当てされません。コンピュータノードをリブートする前にインスタンスのコールドマイグレーションを行うか、リブート後に各仮想 GPU を正しいインスタンスに手動で割り当てる必要があります。各仮想 GPU を手動で割り当てるには、リブートする前にコンピュータノードで実行される各仮想 GPU インスタンスのインスタンス XML から **mdev** UUID を取得する必要があります。以下のコマンドを使用して、各インスタンスの **mdev** UUID を検出することができます。

```
# virsh dumpxml <instance_name> | grep mdev
```

<instance_name> を、Compute API への **/servers** リクエストで返される libvirt インスタンス名 (**OS-EXT-SRV-ATTR:instance_name**) に置き換えます。

- libvirt の制限により、仮想 GPU 対応インスタンスでの休止操作はサポートされていません。代わりに、インスタンスのスナップショット作成またはシェルフ処理が可能です。
- デフォルトでは、コンピュータホストの仮想 GPU の種別は API ユーザーに公開されません。アクセス権限を付与するには、ホストをホストアグリゲートに追加します。詳細は、[Creating and managing host aggregates](#) を参照してください。
- NVIDIA アクセラレーターハードウェアを使用する場合は、NVIDIA ライセンス要件に従う必要があります。たとえば、NVIDIA vGPU GRID にはライセンスサーバーが必要です。NVIDIA のライセンス要件の詳細は、NVIDIA の Web サイトで [Virtual GPU License Server Release Notes](#) を参照してください。

13.2. コンピュータノードでの仮想 GPU の設定

クラウドユーザーが仮想 GPU (vGPU) を使用するインスタンスを作成できるようにするには、物理 GPU を持つコンピュータノードを設定する必要があります。

1. vGPU 用のコンピュータノードを指定する。
2. 仮想 GPU 用のコンピュータノードを設定する。
3. オーバークラウドをデプロイします。
4. vGPU を持つインスタンスを起動するための vGPU フレーバーを作成します。

ヒント

GPU ハードウェアが制限されている場合は、ホストアグリゲートを設定して vGPU コンピュータノードでのスケジューリングを最適化することもできます。仮想 GPU を要求するインスタンスのみを仮想 GPU コンピュータノードにスケジューリングするには、仮想 GPU が設定されたコンピュータノードのホストアグリゲートを作成し、Compute スケジューラーが仮想 GPU インスタンスのみをホストアグリゲートに配置するように設定します。詳細は、[Creating and managing host aggregates](#) および [Filtering by isolating host aggregates](#) を参照してください。



注記

NVIDIA GRID vGPU を使用するには、NVIDIA GRID ライセンス要件に従う共に、セルフホストライセンスサーバーの URL が必要です。詳細は、[Virtual GPU License Server Release Notes](#) の Web ページを参照してください。

13.2.1. 前提条件

- NVIDIA の Web サイトから、GPU デバイスに対応する NVIDIA GRID ホストドライバー RPM パッケージをダウンロードしている。必要なドライバーを確認するには、[NVIDIA ドライバーダウンロードポータル](#) を参照してください。ポータルからドライバーをダウンロードするには、NVIDIA カスタマーとして登録されている必要があります。
- NVIDIA GRID ホストドライバーがインストールされているカスタムオーバークラウドイメージをビルドしている。

13.2.2. 仮想 GPU 用コンピュータノードの指定

vGPU ワークロードのコンピュータノードを指定するには、新しいロールファイルを作成して vGPU ロールを設定し、GPU 対応のコンピュータノードのタグ付けに使用する GPU リソースクラスを使用してベアメタルノードを設定する必要があります。



注記

以下の手順は、まだプロビジョニングされていない新しいオーバークラウドノードに適用されます。すでにプロビジョニングされている既存のオーバークラウドノードにリソースクラスを割り当てるには、スケールダウン手順を使用してノードのプロビジョニングを解除してから、スケールアップ手順を使用して新しいリソースクラスの割り当てでノードを再プロビジョニングする必要があります。詳細は、[オーバークラウドノードのスケールリング](#) を参照してください。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. **stackrc** ファイルを取得します。

```
[stack@director ~]$ source ~/stackrc
```

3. **Controller**、**Compute**、**ComputeGpu** ロールを含む、**roles_data_gpu.yaml** という名前の新しいロールデータファイルを、オーバークラウドに必要なその他のロールとともに生成します。

```
(undercloud)$ openstack overcloud roles \
generate -o /home/stack/templates/roles_data_gpu.yaml \
Compute:ComputeGpu Compute Controller
```

4. **roles_data_gpu.yaml** を開き、以下のパラメーターおよびセクションを編集または追加します。

セクション/パラメーター	現在の値	新しい値
ロールのコメント	Role: Compute	Role: ComputeGpu

セクション/パラメーター	現在の値	新しい値
name (ロール名)	Compute	ComputeGpu
description	Basic Compute Node role	GPU Compute Node role
HostnameFormatDefault	-compute-	-computegpu-
deprecated_nic_config_name	compute.yaml	compute-gpu.yaml

5. オーバークラウド用の GPU 対応コンピュートノードをノード定義のテンプレート **node.json** または **node.yaml** に追加して、そのノードを登録します。詳細は、**director のインストールと使用方法** の [オーバークラウドノードの登録](#) を参照してください。
6. ノードのハードウェアを検査します。

```
(undercloud)$ openstack overcloud node introspect --all-manageable \
--provide
```

詳細については、**Director のインストールと使用 ガイド**の [ベアメタルノードハードウェアのインベントリーの作成](#) を参照してください。

7. GPU 負荷用に指定する各ベアメタルノードに、カスタムの GPU リソースクラスをタグ付けします。

```
(undercloud)$ openstack baremetal node set \
--resource-class baremetal.GPU <node>
```

<node> をベアメタルノードの ID に置き換えてください。

8. ノード定義ファイル **overcloud-baremetal-deploy.yaml** に **ComputeGpu** ロールを追加し、予測ノード配置、リソースクラス、ネットワークポロジ、またはノードに割り当てるその他の属性を定義します。

```
- name: Controller
  count: 3
- name: Compute
  count: 3
- name: ComputeGpu
  count: 1
defaults:
  resource_class: baremetal.GPU
  network_config:
    template: /home/stack/templates/nic-config/myRoleTopology.j2 1
```

- 1** 既存のネットワークポロジを再利用するか、ロール用の新しいカスタムネットワークインターフェイステンプレートを作成できます。詳細は、**Director のインストールと使用ガイド**の [カスタムネットワークインターフェイステンプレート](#) を参照してください。**network_config** プロパティを使用してネットワーク定義を定義しない場合、デフォルトのネットワーク定義が使用されます。

ノード定義ファイルでノード属性を設定するために使用できるプロパティについて詳しくは、[ベアメタルノードのプロビジョニング属性](#)を参照してください。ノード定義ファイルの例は、[ノード定義ファイルの例](#)を参照してください。

9. プロビジョニングコマンドを実行して、ロールの新しいノードをプロビジョニングします。

```
(undercloud)$ openstack overcloud node provision \
--stack <stack> \
[--network-config \]
--output /home/stack/templates/overcloud-baremetal-deployed.yaml \
/home/stack/templates/overcloud-baremetal-deploy.yaml
```

- **<stack>** を、ベアメタルノードがプロビジョニングされるスタックの名前に置き換えます。指定しない場合、デフォルトは **overcloud** です。
 - **--network-config** オプションの引数を含めて、**cli-overcloud-node-network-config.yaml** Ansible Playbook にネットワーク定義を提供します。**network_config** プロパティを使用してネットワーク定義を定義しない場合、デフォルトのネットワーク定義が使用されます。
10. 別のターミナルでプロビジョニングの進捗をモニターリングします。プロビジョニングが成功すると、ノードの状態が **available** から **active** に変わります。

```
(undercloud)$ watch openstack baremetal node list
```

11. **--network-config** オプションを指定してプロビジョニングコマンドを実行しなかった場合は、**network-environment.yaml** ファイルで **<Role>NetworkConfigTemplate** パラメーターを設定して、NIC テンプレートファイルを指すようにします。

```
parameter_defaults:
  ComputeNetworkConfigTemplate: /home/stack/templates/nic-configs/compute.j2
  ComputeGpuNetworkConfigTemplate: /home/stack/templates/nic-configs/<gpu_net_top>.j2
  ControllerNetworkConfigTemplate: /home/stack/templates/nic-configs/controller.j2
```

<gpu_net_top> を **ComputeGpu** ロールのネットワークトポロジが含まれるファイルの名前に置き換えます。たとえば、デフォルトのネットワークトポロジを使用する場合は **compute.yaml** です。

13.2.3. 仮想 GPU 用コンピュートノードの設定およびオーバークラウドのデプロイ

環境内の物理 GPU デバイスに対応する仮想 GPU の種別を取得して割り当て、仮想 GPU 用コンピュートノードを設定するための環境ファイルを準備する必要があります。

手順

1. Red Hat Enterprise Linux と NVIDIA GRID ドライバーを一時コンピュートノードにインストールし、そのノードを起動します。
2. コンピュートノードで、有効にする物理 GPU デバイスに対応する仮想 GPU の種別を確認します。libvirt の場合、仮想 GPU は仲介デバイスまたは **mdev** 種別のデバイスです。サポートされている **mdev** デバイスを検出するには、以下のコマンドを入力します。

```
[root@overcloud-computegpu-0 ~]# ls
/sys/class/mdev_bus/0000:06:00.0/mdev_supported_types/
nvidia-11 nvidia-12 nvidia-13 nvidia-14 nvidia-15 nvidia-16 nvidia-17 nvidia-18 nvidia-19
```

```
nvidia-20 nvidia-21 nvidia-210 nvidia-22
```

```
[root@overcloud-computegpu-0 ~]# cat
/sys/class/mdev_bus/0000\:06\:00.0/mdev_supported_types/nvidia-18/description
num_heads=4, frl_config=60, framebuffer=2048M, max_resolution=4096x2160,
max_instance=4
```

3. **gpu.yaml** ファイルを作成し、GPU デバイスの仮想 GPU 種別を指定します。

```
parameter_defaults:
  ComputeGpuExtraConfig:
    nova::compute::vgpu::enabled_vgpu_types:
      - nvidia-18
```



注記

各物理 GPU がサポートするのは1つの仮想 GPU 種別だけです。この属性で複数の仮想 GPU 種別を指定した場合は、最初の種別だけが使用されます。

4. 更新内容を Compute 環境ファイルに保存します。
5. その他の環境ファイルと共に新しいロールファイルおよび環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-r /home/stack/templates/roles_data_gpu.yaml \
-e /home/stack/templates/network-environment.yaml \
-e /home/stack/templates/gpu.yaml \
-e /home/stack/templates/overcloud-baremetal-deployed.yaml \
-e /home/stack/templates/node-info.yaml
```

13.3. カスタム GPU インスタンスイメージの作成

クラウドユーザーが仮想 GPU (vGPU) を使用するインスタンスを作成できるようにするには、インスタンス起動用のカスタムの仮想 GPU 対応イメージを作成します。NVIDIA GRID ゲストドライバおよびライセンスファイルを使用してカスタムの仮想 GPU 対応インスタンスイメージを作成するには、以下の手順を使用します。

前提条件

- GPU 対応のコンピュータードと共にオーバークラウドを設定およびデプロイしている。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。
2. `source` コマンドで **overcloudrc** 認証情報ファイルを読み込みます。

```
$ source ~/overcloudrc
```

3. 仮想 GPU インスタンスが必要とするハードウェアおよびソフトウェアプロファイルでインスタンスを作成します。

```
(overcloud)$ openstack server create --flavor <flavor> \
--image <image> temp_vgpu_instance
```

- **<flavor>** を、仮想 GPU インスタンスが必要とするハードウェアプロファイルを持つフレーバーの名前または ID に置き換えてください。仮想 GPU フレーバー作成に関する詳細は、[Creating a vGPU flavor for instances](#) を参照してください。
 - **<image>** を、仮想 GPU インスタンスが必要とするソフトウェアプロファイルを持つイメージの名前または ID に置き換えてください。RHEL クラウドイメージのダウンロードに関する情報は、[イメージの管理](#) を参照してください。
4. cloud-user としてインスタンスにログインします。
 5. NVIDIA のガイダンス ([Licensing an NVIDIA vGPU on Linux by Using a Configuration File](#)) に従って、インスタンス上に **gridd.conf** NVIDIA GRID ライセンスファイルを作成します。
 6. インスタンスに GPU ドライバーをインストールします。NVIDIA ドライバーのインストールについての詳細は、[Installing the NVIDIA vGPU Software Graphics Driver on Linux](#) を参照してください。



注記

hw_video_model イメージ属性を使用して GPU ドライバーの種別を定義します。仮想 GPU インスタンスのエミュレートされた GPU を無効にする場合は、**none** を選択します。サポートされているドライバーについての詳しい情報は、[イメージ設定パラメーター](#) を参照してください。

7. インスタンスのイメージスナップショットを作成します。

```
(overcloud)$ openstack server image create \
--name vgpu_image temp_vgpu_instance
```

8. オプション: インスタンスを削除します。

13.4. インスタンス用の仮想 GPU フレーバーの作成

クラウドユーザーが GPU 負荷用のインスタンスを作成できるようにするには、仮想 GPU インスタンスを起動するための GPU フレーバーを作成し、仮想 GPU のリソースをそのフレーバーに割り当てます。

前提条件

- GPU 対応コンピュートノードと共にオーバークラウドを設定およびデプロイしている。

手順

1. NVIDIA GPU フレーバーを作成します。以下に例を示します。

```
(overcloud)$ openstack flavor create --vcpus 6 \
--ram 8192 --disk 100 m1.small-gpu
```

Field	Value
OS-FLV-DISABLED:disabled	False
OS-FLV-EXT-DATA:ephemeral	0

```

| disk          | 100          |
| id            | a27b14dd-c42d-4084-9b6a-225555876f68 |
| name          | m1.small-gpu |
| os-flavor-access:is_public | True        |
| properties    |              |
| ram           | 8192         |
| rxtx_factor   | 1.0          |
| swap          |              |
| vcpus         | 6            |
+-----+-----+

```

- 作成したフレーバーに仮想 GPU のリソースを割り当てます。各インスタンスに割り当てられる仮想 GPU は1つだけです。

```

(overcloud)$ openstack flavor set m1.small-gpu \
--property "resources:VGPU=1"

(overcloud)$ openstack flavor show m1.small-gpu
+-----+-----+
| Field          | Value          |
+-----+-----+
| OS-FLV-DISABLED:disabled | False          |
| OS-FLV-EXT-DATA:ephemeral | 0              |
| access_project_ids      | None           |
| disk                   | 100            |
| id                     | a27b14dd-c42d-4084-9b6a-225555876f68 |
| name                   | m1.small-gpu   |
| os-flavor-access:is_public | True          |
| properties             | resources:VGPU='1' |
| ram                    | 8192           |
| rxtx_factor            | 1.0            |
| swap                   |                |
| vcpus                  | 6              |
+-----+-----+

```

13.5. 仮想 GPU インスタンスの起動

GPU 負荷用の GPU 対応インスタンスを作成することができます。

手順

- GPU フレーバーおよびイメージを使用して、インスタンスを作成します。以下に例を示します。

```

(overcloud)$ openstack server create --flavor m1.small-gpu \
--image vgpu_image --security-group web --nic net-id=internal0 \
--key-name lambda vgpu-instance

```

- cloud-user としてインスタンスにログインします。
- インスタンスが GPU にアクセスできることを確認するには、インスタンスから以下のコマンドを入力します。

```

$ lspci -nn | grep <gpu_name>

```

13.6. GPU デバイスの PCI パススルーの有効化

PCI パススルーを使用して、グラフィックカード等の物理 PCI デバイスをインスタンスに接続することができます。デバイスに PCI パススルーを使用する場合、インスタンスはタスクを実行するためにデバイスへの排他的アクセスを確保し、ホストはデバイスを利用することができません。

前提条件

- **pciutils** パッケージが PCI カードを持つ物理サーバーにインストールされている。
- GPU デバイスのドライバーが、デバイスをパススルーするインスタンスにインストールされている必要があります。したがって、必要な GPU ドライバーがインストールされたカスタムのインスタンスイメージを作成している必要があります。GPU ドライバーがインストールされたカスタムのインスタンスイメージを作成する方法についての詳細は、[Creating a custom GPU instance image](#) を参照してください。

手順

1. 各パススルーデバイス種別のベンダー ID および製品 ID を確認するには、PCI カードを持つ物理サーバーで以下のコマンドを入力します。

```
# lspci -nn | grep -i <gpu_name>
```

たとえば、NVIDIA GPU のベンダーおよび製品 ID を確認するには、以下のコマンドを入力します。

```
# lspci -nn | grep -i nvidia
3b:00.0 3D controller [0302]: NVIDIA Corporation TU104GL [Tesla T4] [10de:1eb8] (rev a1)
d8:00.0 3D controller [0302]: NVIDIA Corporation TU104GL [Tesla T4] [10de:1db4] (rev a1)
```

2. 各 PCI デバイスに Single Root I/O Virtualization (SR-IOV) 機能があるかどうかを確認するには、PCI カードを持つ物理サーバーで以下のコマンドを入力します。

```
# lspci -v -s 3b:00.0
3b:00.0 3D controller: NVIDIA Corporation TU104GL [Tesla T4] (rev a1)
...
Capabilities: [bcc] Single Root I/O Virtualization (SR-IOV)
...
```

3. PCI パススルー用にオーバークラウド上のコントローラーノードを設定するには、環境ファイル (例: **pci_passthru_controller.yaml**) を作成します。
4. **pci_passthru_controller.yaml** の **NovaSchedulerEnabledFilters** パラメーターに **PciPassthroughFilter** を追加します。

```
parameter_defaults:
  NovaSchedulerEnabledFilters:
    - AvailabilityZoneFilter
    - ComputeFilter
    - ComputeCapabilitiesFilter
    - ImagePropertiesFilter
    - ServerGroupAntiAffinityFilter
```

- ServerGroupAffinityFilter
- PciPassthroughFilter
- NUMATopologyFilter

5. コントローラーノード上のデバイスの PCI エイリアスを指定するには、以下の設定を **pci_passthru_controller.yaml** に追加します。

- PCI デバイスに SR-IOV 機能がある場合:

```
ControllerExtraConfig:
nova::pci::aliases:
- name: "t4"
  product_id: "1eb8"
  vendor_id: "10de"
  device_type: "type-PF"
- name: "v100"
  product_id: "1db4"
  vendor_id: "10de"
  device_type: "type-PF"
```

- PCI デバイスに SR-IOV 機能がない場合:

```
ControllerExtraConfig:
nova::pci::aliases:
- name: "t4"
  product_id: "1eb8"
  vendor_id: "10de"
- name: "v100"
  product_id: "1db4"
  vendor_id: "10de"
```

device_type フィールドの設定に関する詳細は、[PCI passthrough device type field](#) を参照してください。



注記

nova-api サービスが Controller 以外のロールで実行されている場合は、**ControllerExtraConfig** を **<Role>ExtraConfig** の形式でユーザーロールに置き換えます。

6. PCI パススルー用にオーバークラウド上のコンピュートノードを設定するには、環境ファイル (例: **pci_passthru_compute.yaml**) を作成します。

7. コンピュートノード上のデバイスで利用可能な PCI を指定するには、以下の設定を **pci_passthru_compute.yaml** に追加します。

```
parameter_defaults:
NovaPCIPassthrough:
- vendor_id: "10de"
  product_id: "1eb8"
```

8. インスタンスの移行およびサイズ変更の操作を行うために、コンピュートノードの PCI エイリアスのコピーを作成する必要があります。コンピュートノード上のデバイスの PCI エイリアスを指定するには、以下の設定を **pci_passthru_compute.yaml** に追加します。

- PCI デバイスに SR-IOV 機能がある場合:

```
ComputeExtraConfig:
nova::pci::aliases:
- name: "t4"
  product_id: "1eb8"
  vendor_id: "10de"
  device_type: "type-PF"
- name: "v100"
  product_id: "1db4"
  vendor_id: "10de"
  device_type: "type-PF"
```

- PCI デバイスに SR-IOV 機能がない場合:

```
ComputeExtraConfig:
nova::pci::aliases:
- name: "t4"
  product_id: "1eb8"
  vendor_id: "10de"
- name: "v100"
  product_id: "1db4"
  vendor_id: "10de"
```

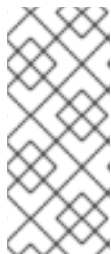


注記

コンピュートノードのエイリアスは、コントローラーノードのエイリアスと同じでなければなりません。

9. PCI パススルーをサポートするためにコンピュートノードのサーバー BIOS で IOMMU を有効にするには、**pci_passthru_compute.yaml** に **KernelArgs** パラメーターを追加します。

```
parameter_defaults:
...
ComputeParameters:
  KernelArgs: "intel_iommu=on iommu=pt"
```



注記

KernelArgs パラメーターをロールの設定に初めて追加すると、オーバークラウドノードが自動的に再起動されます。必要に応じて、ノードの自動再起動を無効にし、代わりに各オーバークラウドのデプロイ後にノードの再起動を手動で実行できます。詳細は、[KernelArgs を定義するための手動でのノード再起動の設定](#)を参照してください。

10. その他の環境ファイルと共にこれらのカスタム環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/pci_passthru_controller.yaml \
-e /home/stack/templates/pci_passthru_compute.yaml
```

11. PCI デバイスを要求するためのフレーバーを設定します。以下の例では、それぞれベンダー ID および製品 ID が **10de** および **13f2** の 2 つのデバイスをリクエストします。

```
# openstack flavor set m1.large \
--property "pci_passthrough:alias"="t4:2"
```

検証

1. PCI パススルーデバイスを設定してインスタンスを作成します。

```
# openstack server create --flavor m1.large \
--image <custom_gpu> --wait test-pci
```

<custom_gpu> を、必要な GPU ドライバーがインストールされたカスタムインスタンスイメージの名前に置き換えます。

2. クラウドユーザーとしてインスタンスにログインします。
3. インスタンスが GPU にアクセスできることを確認するには、インスタンスから以下のコマンドを入力します。

```
$ lspci -nn | grep <gpu_name>
```

4. NVIDIA System Management Interface のステータスを確認するには、インスタンスから以下のコマンドを入力します。

```
$ nvidia-smi
```

出力例:

```
-----
| NVIDIA-SMI 440.33.01    Driver Version: 440.33.01    CUDA Version: 10.2    |
|-----+-----+
| GPU Name      Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
|=====+=====+
| 0 Tesla T4               Off | 00000000:01:00:0 Off |             0      |
| N/A   43C    P0   20W /  70W |  0MiB / 15109MiB |      0%    Default |
|-----+-----+

| Processes:
| GPU      PID   Type   Process name                      Usage      |
|=====+=====+
| No running processes found                |
|-----+-----+
```

第14章 インスタンスの管理

クラウド管理者は、クラウド上で実行されているインスタンスを監視および管理することができます。

14.1. インスタンスの VNC コンソールへの接続のセキュリティ保護

VNC プロキシサービスへの受信クライアントの接続用に適用する許可される TLS 暗号と最小プロトコルバージョンを設定することで、インスタンスの VNC コンソールへの接続をセキュアにすることができます。

手順

1. アンダークラウドに **stack** ユーザーとしてログインします。

2. **stackrc** ファイルを取得します。

```
[stack@director ~]$ source ~/stackrc
```

3. Compute 環境ファイルを開きます。

4. インスタンスへの VNC コンソール接続に使用する最小プロトコルバージョンを設定します。

```
parameter_defaults:
...
NovaVNCProxySSLMinimumVersion: <version>
```

<version> を、許可される最小の SSL/TLS プロトコルバージョンに置き換えます。以下の有効な値のいずれかに設定します。

- **default**: 基礎となるシステム OpenSSL のデフォルトを使用します。
- **tlsv1_1**: 新しいバージョンをサポートしていないクライアントがある場合に使用します。



注記

TLS 1.0 および TLS 1.1 は RHEL 8 で廃止され、RHEL 9 ではサポートされていません。

- **tlsv1_2**: インスタンスへの VNC コンソール接続に使用する SSL/TLS 暗号を設定する場合に使用します。
 - **tlsv1_3**: TLSv1.3 の標準暗号ライブラリーを使用する場合に使用します。 **NovaVNCProxySSLCiphers** パラメーターの設定は無視されます。
5. 使用可能な最小の SSL/TLS プロトコルバージョンを **tlsv1_2** に設定する場合は、インスタンスへの VNC コンソール接続に使用する SSL/TLS 暗号を設定します。

```
parameter_defaults:
NovaVNCProxySSLCiphers: <ciphers>
```

<ciphers> を、許可する暗号スイートのコロン区切りリストに置き換えます。 **openssl** から利用可能な暗号の一覧を取得します。

6. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

14.2. データベースのクリーニング

Compute サービスには管理ツール **nova-manage** が含まれています。このツールを使用して、データベーススキーマの適用、アップグレード中のオンラインデータ移行の実行、データベースの管理およびクリーンアップ等の、デプロイメント、アップグレード、クリーンアップ、およびメンテナンス関連のタスクを実行することができます。

director は、cron を使用してオーバークラウドでの以下のデータベース管理タスクを自動化します。

- 削除された行を実稼働テーブルからシャドウテーブルに移動して、削除されたインスタンスレコードをアーカイブする。
- アーカイブ処理が完了した後に、シャドウテーブルから削除された行をパージする。

14.2.1. データベース管理の設定

cron ジョブは、デフォルト設定を使用してデータベース管理タスクを実行します。デフォルトでは、データベースをアーカイブする cron ジョブは毎日 00:01 に実行され、データベースをパージする cron ジョブは毎日 05:00 に実行されます。共にジッターは 0 秒から 3600 秒の間です。必要に応じて、これらの設定は heat パラメーターを使用して変更することができます。

手順

1. Compute 環境ファイルを開きます。
2. 追加または変更する cron ジョブを制御する heat パラメーターを追加します。たとえば、シャドウテーブルをアーカイブ直後にパージするには、次のパラメーターを True に設定します。

```
parameter_defaults:
...
NovaCronArchiveDeleteRowsPurge: True
```

データベースの cron ジョブを管理する heat パラメーターの完全一覧は、[Configuration options for the Compute service automated database management](#) を参照してください。

3. 更新内容を Compute 環境ファイルに保存します。
4. その他の環境ファイルと共に Compute 環境ファイルをスタックに追加して、オーバークラウドをデプロイします。

```
(undercloud)$ openstack overcloud deploy --templates \
-e [your environment files] \
-e /home/stack/templates/<compute_environment_file>.yaml
```

14.2.2. Compute サービスのデータベース自動管理用設定オプション

以下の heat パラメーターを使用して、データベースを管理する自動 cron ジョブを有効化および変更します。

表14.1 Compute (nova) サービスの cron パラメーター

パラメーター	説明
NovaCronArchiveDeleteAllCells	すべてのセルから削除されたインスタンスレコードをアーカイブするには、このパラメーターを True に設定します。 デフォルト: True
NovaCronArchiveDeleteRowsAge	このパラメーターを使用して、削除されたインスタンスレコードを経過時間 (日数単位) に基づいてアーカイブします。 シャドウテーブル内の本日以前のデータをアーカイブするには、 0 に設定します。 デフォルト: 90
NovaCronArchiveDeleteRowsDestination	このパラメーターを使用して、削除されたインスタンスレコードを記録するファイルを設定します。 デフォルト: /var/log/nova/nova-rowsflush.log
NovaCronArchiveDeleteRowsHour	このパラメーターを使用して、削除されたインスタンスレコードを別のテーブルに移動する cron コマンドを実行する時刻の時間部分を設定します。 デフォルト: 0
NovaCronArchiveDeleteRowsMaxDelay	このパラメーターを使用して、削除されたインスタンスレコードを別のテーブルに移動するまでの最大遅延時間 (秒単位) を設定します。 デフォルト: 3600
NovaCronArchiveDeleteRowsMaxRows	このパラメーターを使用して、別のテーブルに移動することのできる削除インスタンスレコード数の最大値を設定します。 デフォルト: 1000
NovaCronArchiveDeleteRowsMinute	このパラメーターを使用して、削除されたインスタンスレコードを別のテーブルに移動する cron コマンドを実行する時刻の分部分を設定します。 デフォルト: 1

パラメーター	説明
NovaCronArchiveDeleteRowsMonthday	<p>このパラメーターを使用して、削除されたインスタンスレコードを別のテーブルに移動する cron コマンドを実行する日にちを設定します。</p> <p>デフォルト: * (毎日)</p>
NovaCronArchiveDeleteRowsMonth	<p>このパラメーターを使用して、削除されたインスタンスレコードを別のテーブルに移動する cron コマンドを実行する月を設定します。</p> <p>デフォルト: * (毎月)</p>
NovaCronArchiveDeleteRowsPurge	<p>スケジュールされたアーカイブ処理の直後にシャドウテーブルをパージするには、このパラメーターを True に設定します。</p> <p>デフォルト: False</p>
NovaCronArchiveDeleteRowsUntilComplete	<p>すべてのレコードを移動するまで削除されたインスタンスレコードを別のテーブルに移動し続けるには、このパラメーターを True に設定します。</p> <p>デフォルト: True</p>
NovaCronArchiveDeleteRowsUser	<p>このパラメーターを使用して、削除されたインスタンスレコードをアーカイブする crontab の所有権を持ち、crontab が使用するログファイルにアクセスできるユーザーを設定します。</p> <p>デフォルト: nova</p>
NovaCronArchiveDeleteRowsWeekday	<p>このパラメーターを使用して、削除されたインスタンスレコードを別のテーブルに移動する cron コマンドを実行する曜日を設定します。</p> <p>デフォルト: * (毎日)</p>
NovaCronPurgeShadowTablesAge	<p>このパラメーターを使用して、シャドウテーブルを経過時間 (日数単位) に基づいてパージします。</p> <p>本日以前のシャドウテーブルをパージするには、0 に設定します。</p> <p>デフォルト: 14</p>
NovaCronPurgeShadowTablesAllCells	<p>すべてのセルからシャドウテーブルをパージするには、このパラメーターを True に設定します。</p> <p>デフォルト: True</p>

パラメーター	説明
NovaCronPurgeShadowTablesDestination	<p>このパラメーターを使用して、パージされたシャドウテーブルを記録するファイルを設定します。</p> <p>デフォルト: /var/log/nova/nova-rowspurge.log</p>
NovaCronPurgeShadowTablesHour	<p>このパラメーターを使用して、シャドウテーブルをパージする cron コマンドを実行する時刻の時間部分を設定します。</p> <p>デフォルト: 5</p>
NovaCronPurgeShadowTablesMaxDelay	<p>このパラメーターを使用して、シャドウテーブルをパージするまでの最大遅延時間 (秒単位) を設定します。</p> <p>デフォルト: 3600</p>
NovaCronPurgeShadowTablesMinute	<p>このパラメーターを使用して、シャドウテーブルをパージする cron コマンドを実行する時刻の分部分を設定します。</p> <p>デフォルト: 0</p>
NovaCronPurgeShadowTablesMonth	<p>このパラメーターを使用して、シャドウテーブルをパージする cron コマンドを実行する月を設定します。</p> <p>デフォルト: *(毎月)</p>
NovaCronPurgeShadowTablesMonthday	<p>このパラメーターを使用して、シャドウテーブルをパージする cron コマンドを実行する日にちを設定します。</p> <p>デフォルト: *(毎日)</p>
NovaCronPurgeShadowTablesUser	<p>このパラメーターを使用して、シャドウテーブルをパージする crontab の所有権を持ち、crontab が使用するログファイルにアクセスできるユーザーを設定します。</p> <p>デフォルト: nova</p>
NovaCronPurgeShadowTablesVerbose	<p>このパラメーターを使用して、パージされたシャドウテーブルに関するログファイルの詳細ロギングを有効にします。</p> <p>デフォルト: False</p>

パラメーター	説明
NovaCronPurgeShadowTablesWeekday	このパラメーターを使用して、シャドウテーブルをパージする cron コマンドを実行する曜日を設定します。 デフォルト: *(毎日)

14.3. コンピュートノード間の仮想マシンインスタンスの移行

メンテナンスを実行する場合やワークロードのリバランスを行う場合、あるいは障害が発生した/障害が発生しつつあるノードを置き換える場合に、あるコンピュートノードからオーバークラウド内の別のコンピュートノードにインスタンスを移行しなければならない場合があります。

コンピュートノードのメンテナンス

ハードウェアのメンテナンスや修理、カーネルのアップグレードおよびソフトウェアの更新を行うなどの理由により、コンピュートノードを一時的に停止する必要がある場合、コンピュートノード上で実行中のインスタンスを別のコンピュートノードに移行することができます。

障害が発生しつつあるコンピュートノード

コンピュートノードで障害が発生する可能性があり、ノードのサービスまたは置き換えが必要な場合、障害が発生しつつあるコンピュートノードから正常なコンピュートノードにインスタンスを移行することができます。

障害が発生したコンピュートノード

コンピュートノードですでに障害が発生している場合には、インスタンスを退避させることができます。同じ名前、UUID、ネットワークアドレス、およびコンピュートノードに障害が発生する前にインスタンスに割り当てられていたその他すべてのリソースを使用して、元のイメージから別のコンピュートノードにインスタンスを再ビルドすることができます。

ワークロードのリバランス

ワークロードをリバランスするために、1つまたは複数のインスタンスを別のコンピュートノードに移行することができます。たとえば、コンピュートノード上のインスタンスを1つにまとめて電力を節約する、他のネットワークリソースに物理的に近いコンピュートノードにインスタンスを移行してレイテンシーを低減する、インスタンスを全コンピュートノードに分散してホットスポットをなくし復元力を向上させる、等が可能です。

director は、すべてのコンピュートノードがセキュアな移行を提供するように設定します。すべてのコンピュートノードには、移行プロセス中それぞれのホストのユーザーが他のコンピュートノードにアクセスできるように、共有 SSH キーも必要です。director は、**OS::TripleO::Services::NovaCompute** コンポーザブルサービスを使用してこのキーを作成します。このコンポーザブルサービスは、すべての Compute ロールにデフォルトで含まれているメインのサービスの1つです。詳細は、**Director インストールおよび使用方法** ガイドの [コンポーザブルサービスおよびカスタムロール](#) を参照してください。



注記

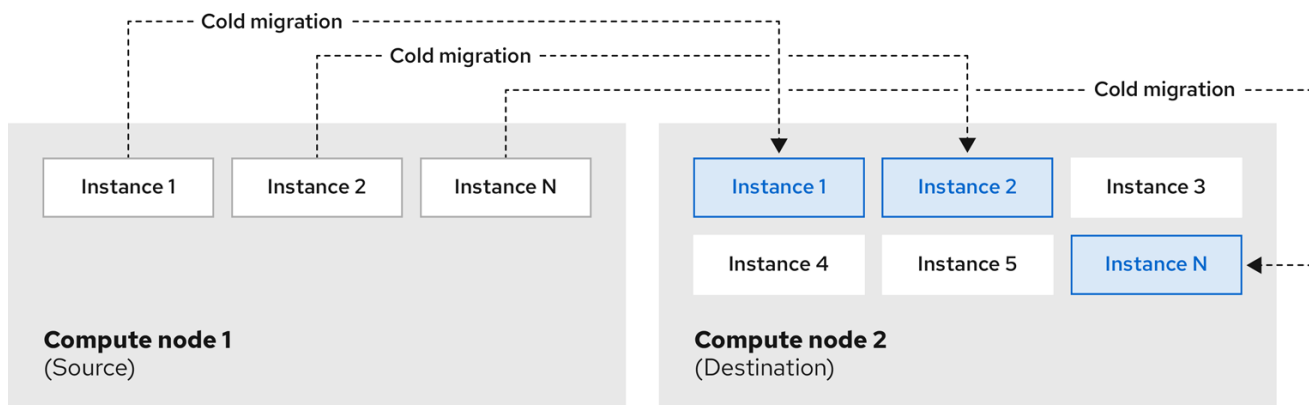
機能しているコンピュートノードがあり、バックアップ目的でインスタンスのコピーを作成する場合、またはインスタンスを別の環境にコピーする場合は、**director** のインストールと使用 ガイドの [オーバークラウドへの仮想マシンのインポート](#) の手順に従います。

14.3.1. 移行の種別

Red Hat OpenStack Platform (RHOSP) では、以下の移行種別がサポートされています。

コールドマイグレーション

コールドマイグレーション (あるいは、非ライブマイグレーション) では、動作中のインスタンスをシャットダウンしてから、移行元コンピュートノードから移行先コンピュートノードに移行します。



273_OpenStack_0822

コールドマイグレーションでは、インスタンスに多少のダウンタイムが発生します。移行したインスタンスは、引き続き同じボリュームおよび IP アドレスにアクセスすることができます。

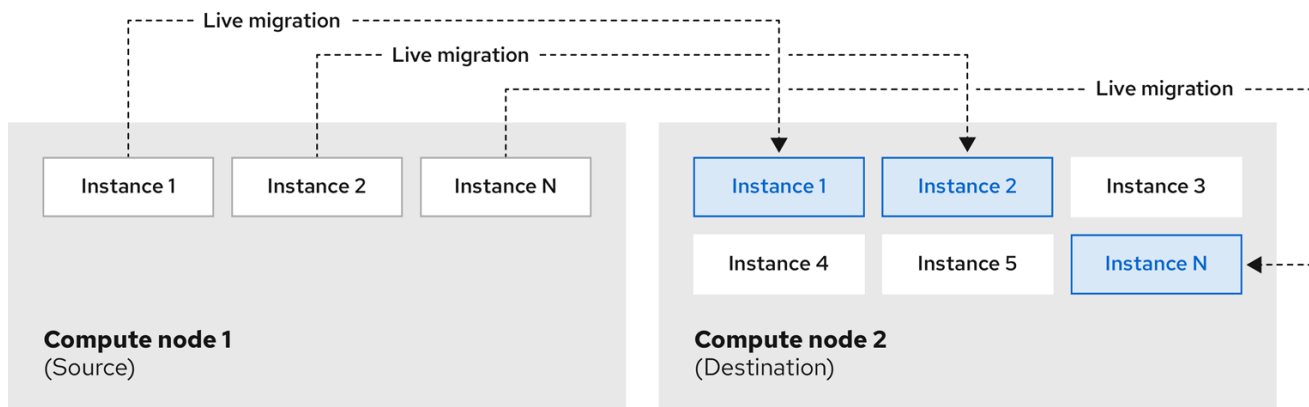


注記

コールドマイグレーションを行うためには、移行元および移行先両方のコンピュートノードが動作状態になければなりません。

ライブマイグレーション

ライブマイグレーションでは、インスタンスをシャットダウンせずに、動作状態を維持しながら移行元コンピュートノードから移行先コンピュートノードに移行します。



273_OpenStack_0822

インスタンスのライブマイグレーションを行う場合、ダウンタイムはほとんど発生しません。ただし、ライブマイグレーションは、移行操作中のパフォーマンスに影響を及ぼします。したがって、移行中のインスタンスは重要なパスから除外する必要があります。



重要

ライブマイグレーションは、移動されるワークロードのパフォーマンスに影響を与えます。Red Hat は、ライブマイグレーション中のパケット損失、ネットワーク遅延、メモリー遅延の増加、またはネットワーク帯域幅、メモリー帯域幅、ストレージ IO、または CPU パフォーマンスの低下をサポートしません。



注記

ライブマイグレーションを行うためには、移行元および移行先両方のコンピュートノードが動作状態になければなりません。

状況によっては、インスタンスのライブマイグレーションを行うことができない場合があります。詳細は、[移行の制約](#)を参照してください。

退避

コンピュートノードですでに障害が発生しているためインスタンスを移行する必要がある場合、インスタンスを退避させることができます。

14.3.2. 移行の制約

移行の制約は通常、ブロックマイグレーション、設定ディスク、またはいずれかのインスタンスがコンピュートノード上の物理ハードウェアにアクセスする場合に生じます。

CPU に関する制約

移行元および移行先コンピュートノードの CPU アーキテクチャーは、同一であることが必須です。たとえば、Red Hat では、**ppc64le** CPU から **x86_64** CPU へのインスタンスの移行をサポートしません。

異なる CPU モデル間の移行はサポートされていません。CPU ホストパススルーを使用するインスタンス等の場合には、移行元および移行先コンピュートノードの CPU は、完全に同一でなければなりません。すべてのケースで、移行先ノードの CPU 機能は、移行元ノードの CPU 機能の上位セットであることが必須です。

メモリーに関する制約

移行先コンピュートノードでは、十分な RAM が利用可能でなければなりません。メモリーのオーバースブスクリプションが、移行失敗の原因となる可能性があります。

ブロックマイグレーションに関する制約

インスタンスの使用するディスクがコンピュートノード上にローカルに格納されている場合、その移行には、共有ストレージ (Red Hat Ceph Storage 等) を使用するボリュームベースのインスタンスよりもはるかに長い時間がかかります。このレイテンシーは、OpenStack Compute (nova) がコンピュートノード間でローカルディスクをブロックごとに移行するために発生します。この処理は、デフォルトではコントロールプレーンネットワークを通じて行われます。これとは対照的に、Red Hat Ceph Storage 等の共有ストレージを使用するボリュームベースのインスタンスでは、ボリュームを移行する必要がありません。それぞれのコンピュートノードが、共有ストレージにアクセスできるためです。



注記

大量の RAM を消費するローカルディスクまたはインスタンスの移行により、コントロールプレーンネットワークに輻輳が生じ、コントロールプレーンネットワークを使用する他のシステム (RabbitMQ 等) のパフォーマンスに悪影響を及ぼす場合があります。

読み取り専用ドライブの移行に関する制約

ドライブの移行は、ドライブに読み取りおよび書き込み両方の機能がある場合に限りサポートされます。たとえば、OpenStack Compute (nova) は CD-ROM ドライブまたは読み取り専用のコンフィグドライブを移行することはできません。ただし、OpenStack Compute (nova) は、**vfat** 等のドライブ形式を持つコンフィグドライブなど、読み取りおよび書き込み両方の機能を持つドライブを移行することができます。

ライブマイグレーションに関する制約

インスタンスのライブマイグレーションでは、さらに制約が生じる場合があります。



重要

ライブマイグレーションは、移動されるワークロードのパフォーマンスに影響を与えます。Red Hat は、ライブマイグレーション中のパケット損失、ネットワーク遅延、メモリー遅延の増加、またはネットワーク帯域幅、メモリー帯域幅、ストレージ IO、または CPU パフォーマンスの低下をサポートしません。

移行中に新しい操作を行うことができない

移行元および移行先ノードのインスタンスのコピー間で状態の整合性を確保するために、RHOSP ではライブマイグレーション中の新規操作を拒否する必要があります。拒否しないと、ライブマイグレーションでメモリーの状態を複製する前にメモリーへの書き込みが行われた場合、ライブマイグレーションに長い時間がかかる状況や、永久に完了しない状況が発生する可能性があります。

NUMA を使用した CPU ピニング

Compute 設定の **NovaSchedulerEnabledFilters** パラメーターには、**AggregateInstanceExtraSpecsFilter** および **NUMATopologyFilter** の値が含まれている必要があります。

マルチセルクラウド

マルチセルクラウドでは、同じセル内の別のホストにインスタンスのライブマイグレーションを行うことができますが、セル全体では移行できません。

フローティングインスタンス

フローティングインスタンスのライブマイグレーションを行う場合、移行先コンピュートノードの **NovaComputeCpuSharedSet** の設定と移行元コンピュートノードの **NovaComputeCpuSharedSet** の設定が異なると、移行先コンピュートノードでは、インスタンスは共有の (ピンングされていない) インスタンス用に設定された CPU には割り当てられません。したがって、フローティングインスタンスのライブマイグレーションを行う必要がある場合は、専用の (ピンングされた) インスタンスおよび共有の (ピンングされていない) インスタンスに関して、すべてのコンピュートノードに同じ CPU マッピングを設定する必要があります。あるいは、共有のインスタンスにホストアグリゲートを使用します。

移行先コンピュートノードの容量

移行先コンピュートノードには、移行するインスタンスをホストするのに十分な空き容量が必要です。

SR-IOV ライブマイグレーション

SR-IOV ベースのネットワークインターフェイスを使用するインスタンスは、ライブマイグレーションが可能です。ダイレクトモードの SR-IOV ネットワークインターフェイスを持つインスタンスのライブマイグレーションでは、ネットワークのダウンタイムが発生します。これは、移行時に、ダイレクトモードのインターフェイスの接続を解除し再び接続する必要があるためです。

ML2/OVN デプロイでのパケット損失

ML2/OVN は、パケット損失のないライブマイグレーションをサポートしていません。これは、OVN が複数のポートバインディングを処理できないため、ポートがいつ移行されるかがわからないためです。

ライブマイグレーション中のパケット損失を最小限に抑えるには、移行が完了したら、宛先ホストでインスタンスをアナウンスするように、ML2/OVN デプロイメントを設定します。

```
parameter_defaults:
  ComputeExtraConfig:
    nova::workarounds::enable_qemu_monitor_announce_self: true
```

ML2/OVS デプロイメントでのライブマイグレーション

ML2/OVS デプロイメントでインスタンスをライブマイグレーションする場合のパケット損失を最小限に抑えるには、ネットワーキングサービス (neutron) ライブマイグレーションイベントを有効にするように、ML2/OVS デプロイメントを設定し、移行が完了したら、宛先ホストでインスタンスをアナウンスします。

```
parameter_defaults:
  NetworkExtraConfig:
    neutron::server::notifications::nova::live_migration_events: true
  ComputeExtraConfig:
    nova::workarounds::enable_qemu_monitor_announce_self: true
```

ライブマイグレーションの妨げとなる制約

以下の機能を使用するインスタンスのライブマイグレーションを行うことはできません。

PCI パススルー

QEMU/KVM ハイパーバイザーでは、コンピュータノード上の PCI デバイスをインスタンスにアタッチすることができます。PCI パススルーを使用すると、インスタンスは PCI デバイスに排他的にアクセスすることができ、これらのデバイスがインスタンスのオペレーティングシステムに物理的に接続されているかのように表示され、動作します。ただし、PCI パススルーには物理デバイスへの直接アクセスが必要なため、QEMU/KVM は PCI パススルーを使用するインスタンスのライブマイグレーションをサポートしません。

ポートリソースの要求

最小帯域幅を確保する QoS ポリシーなど、リソース要求が設定されたポートを使用するインスタンスのライブマイグレーションを行うことはできません。ポートにリソース要求があるかどうかを確認するには、以下のコマンドを使用します。

```
$ openstack port show <port_name/port_id>
```

14.3.3. 移行の準備

1つまたは複数のインスタンスを移行する前に、コンピュータノード名および移行するインスタンスの ID を把握する必要があります。

手順

1. 移行元コンピュータノードのホスト名および移行先コンピュータノードのホスト名を特定します。

```
(undercloud)$ source ~/overcloudrc
(overcloud)$ openstack compute service list
```

2. 移行元コンピュートノード上のインスタンスを一覧表示し、移行するインスタンスの ID を特定します。

```
(overcloud)$ openstack server list --host <source> --all-projects
```

<source> を移行元コンピュートノードの名前または ID に置き換えてください。

3. オプション: ノードのメンテナンส์を行うためにインスタンスを移行元コンピュートノードから移行する場合、ノードを無効にして、メンテナンส์中にスケジューラーがノードに新規インスタンスを割り当てるのを防ぐ必要があります。

```
(overcloud)$ openstack compute service set <source> nova-compute --disable
```

<source> を移行元コンピュートノードのホスト名に置き換えてください。

これで移行を行う準備が整いました。[Cold migrating an instance](#) または [Live migrating an instance](#) で詳しく説明されている必須手順に従います。

14.3.4. インスタンスのコールドマイグレーション

インスタンスのコールドマイグレーションでは、インスタンスを停止して別のコンピュートノードに移動します。コールドマイグレーションは、PCI パススルーを使用するインスタンスの移行など、ライブマイグレーションでは対応することのできない移行シナリオに対応します。移行先コンピュートノードは、スケジューラーにより自動的に選択されます。詳細は、[移行の制約](#) を参照してください。

手順

1. インスタンスのコールドマイグレーションを行うには、以下のコマンドを入力してインスタンスの電源をオフにして移動します。

```
(overcloud)$ openstack server migrate <instance> --wait
```

- **<instance>** を移行するインスタンスの名前または ID に置き換えてください。
 - ローカルに確保されたボリュームを移行する場合には、**--block-migration** フラグを指定します。
2. 移行が完了するまで待ちます。インスタンスの移行が完了するのを待つ間、移行のステータスを確認することができます。詳細は、[Checking migration status](#) を参照してください。
 3. インスタンスのステータスを確認します。

```
(overcloud)$ openstack server list --all-projects
```

ステータスが VERIFY_RESIZE と表示される場合は、移行を確認する、または元に戻す必要があることを示しています。

- 予想どおりに機能している場合は、移行を確認します。

```
(overcloud)$ openstack server resize --confirm <instance>
```

<instance> を移行するインスタンスの名前または ID に置き換えてください。ステータスが ACTIVE と表示される場合は、インスタンスを使用する準備が整っていることを示しています。

- 予想どおりに機能していない場合は、移行を元に戻します。

```
(overcloud)$ openstack server resize --revert <instance>
```

<instance> をインスタンスの名前または ID に置き換えてください。

4. インスタンスを再起動します。

```
(overcloud)$ openstack server start <instance>
```

<instance> をインスタンスの名前または ID に置き換えてください。

5. オプション: メンテナンスのために移行元コンピュートノードを無効にした場合は、新規インスタンスがノードに割り当てられるようにノードを再度有効にする必要があります。

```
(overcloud)$ openstack compute service set <source> nova-compute --enable
```

<source> を移行元コンピュートノードのホスト名に置き換えてください。

14.3.5. インスタンスのライブマイグレーション

ライブマイグレーションでは、ダウンタイムを最小限に抑えて、インスタンスを移行元コンピュートノードから移行先コンピュートノードに移動します。ライブマイグレーションがすべてのインスタンスに適しているとは限りません。詳細は、[移行の制約](#) を参照してください。

手順

1. インスタンスのライブマイグレーションを行うには、インスタンスおよび移行先コンピュートノードを指定します。

```
(overcloud)$ openstack server migrate <instance> --live-migration [--host <dest>] --wait
```

- **<instance>** をインスタンスの名前または ID に置き換えてください。
- **<dest>** を移行先コンピュートノードの名前または ID に置き換えてください。



注記

openstack server migrate コマンドは、共有ストレージを持つインスタンスの移行が対象です。これがデフォルトの設定です。ローカルに確保されたボリュームを移行するには、**--block-migration** フラグを指定します。

```
(overcloud)$ openstack server migrate <instance> --live-migration [--host <dest>] --wait --block-migration
```

2. インスタンスが移行されていることを確認します。

```
(overcloud)$ openstack server show <instance>
```

Field	Value
...	...
status	MIGRATING
...	...

- 移行が完了するまで待ちます。インスタンスの移行が完了するのを待つ間、移行のステータスを確認することができます。詳細は、[Checking migration status](#) を参照してください。
- インスタンスのステータスをチェックして、移行が成功したかどうかを確認します。

```
(overcloud)$ openstack server list --host <dest> --all-projects
```

<dest> を移行先コンピュートノードの名前または ID に置き換えてください。

- オプション: メンテナンスのために移行元コンピュートノードを無効にした場合は、新規インスタンスがノードに割り当てられるようにノードを再度有効にする必要があります。

```
(overcloud)$ openstack compute service set <source> nova-compute --enable
```

<source> を移行元コンピュートノードのホスト名に置き換えてください。

14.3.6. 移行ステータスの確認

移行が完了するまでに、さまざまな移行状態を遷移します。正常な移行では、通常、移行状態は以下のようになります。

- Queued:** Compute サービスはインスタンス移行の要求を受け入れ、移行は保留中です。
- Preparing:** Compute サービスはインスタンス移行の準備中です。
- Running:** Compute サービスはインスタンスを移行中です。
- Post-migrating:** Compute サービスはインスタンスを移行先コンピュートノードにビルドし、移行元コンピュートノードのリソースを解放しています。
- Completed:** Compute サービスはインスタンスの移行を完了し、移行元コンピュートノードのリソース解放を終了しています。

手順

- インスタンスの移行 ID の一覧を取得します。

```
$ nova server-migration-list <instance>
```

```
+-----+-----+-----+ (...)  
| Id | Source Node | Dest Node | (...)  
+-----+-----+-----+ (...)  
| 2 | -          | -          | (...)  
+-----+-----+-----+ (...)
```

<instance> をインスタンスの名前または ID に置き換えてください。

2. 移行のステータスを表示します。

```
$ nova server-migration-show <instance> <migration_id>
```

- **<instance>** をインスタンスの名前または ID に置き換えてください。
- **<migration_id>** を移行の ID に置き換えてください。
nova server-migration-show コマンドを実行すると、以下の例に示すような出力が返されます。

```
+-----+-----+
| Property      | Value                                |
+-----+-----+
| created_at    | 2017-03-08T02:53:06.000000         |
| dest_compute  | controller                          |
| dest_host     | -                                  |
| dest_node     | -                                  |
| disk_processed_bytes | 0                                |
| disk_remaining_bytes | 0                                |
| disk_total_bytes | 0                                |
| id            | 2                                  |
| memory_processed_bytes | 65502513                        |
| memory_remaining_bytes | 786427904                       |
| memory_total_bytes | 1091379200                       |
| server_uuid   | d1df1b5a-70c4-4fed-98b7-423362f2c47c |
| source_compute | compute2                           |
| source_node   | -                                  |
| status        | running                            |
| updated_at    | 2017-03-08T02:53:47.000000         |
+-----+-----+
```

ヒント

OpenStack Compute サービスは、コピーする残りのメモリーのバイト数によって移行の進捗を測定します。時間が経過してもこの数字が減少しない場合、移行を完了することができず、Compute サービスは移行を中止する可能性があります。

インスタンスの移行に長い時間がかかったり、エラーが発生したりする場合があります。詳細は、[Troubleshooting migration](#) を参照してください。

14.3.7. インスタンスの退避

インスタンスを障害の発生したコンピュートノードまたはシャットダウンしたコンピュートノードから同じ環境内の新しいホストに移動する場合、インスタンスを退避させることができます。

退避プロセスにより元のインスタンスは破棄され、元のイメージ、インスタンス名、UUID、ネットワークアドレス、およびインスタンスに割り当てられていたその他すべてのリソースを使用して、別のコンピュートノードに元のインスタンスが再ビルドされます。

インスタンスが共有ストレージを使用する場合、インスタンスのルートディスクは退避プロセス中に再ビルドされません。移行先コンピュートノードが引き続きこのディスクにアクセス可能なためです。インスタンスが共有ストレージを使用しない場合は、インスタンスのルートディスクも移行先コンピュートノードに再ビルドされます。



注記

- コンピュートノードがフェンシングされ、API が報告するコンピュートノードの状態が down または forced-down である場合に限り、退避を行うことができます。コンピュートノードが down または forced-down と報告されない場合、**evacuate** コマンドは失敗します。
- クラウド管理者でなければ、退避を行うことはできません。

14.3.7.1. 単一のインスタンスの退避

インスタンスを一度に1つずつ退避させることができます。

手順

1. 障害の発生したコンピュートノードに管理者としてログインします。
2. コンピュートノードを無効にします。

```
(overcloud)[stack@director ~]$ openstack compute service set \
  <host> <service> --disable
```

- **<host>** をインスタンスの退避元となるコンピュートノードの名前に置き換えてください。
- **<service>** を無効にするサービスの名前 (例: **nova-compute**) に置き換えてください。

3. インスタンスを退避させるには、以下のコマンドを入力します。

```
(overcloud)[stack@director ~]$ nova evacuate [--password <pass>] <instance> [<dest>]
```

- **<pass>** を退避するインスタンスに設定する admin パスワードに置き換えてください。パスワードを指定しなかった場合には、無作為に生成され、退避の完了時に出力されます。
- **<instance>** を退避させるインスタンスの名前または ID に置き換えてください。
- **<dest>** をインスタンスの退避先となるコンピュートノードの名前に置き換えてください。退避先コンピュートノードを指定しなかった場合には、Compute のスケジューラーがノードを選択します。退避先に指定可能なコンピュートノードを確認するには、以下のコマンドを使用します。

```
(overcloud)[stack@director ~]$ openstack hypervisor list
```

14.3.7.2. ホスト上の全インスタンスの退避

指定したコンピュートノード上の全インスタンスを退避させることができます。

手順

1. 障害の発生したコンピュートノードに管理者としてログインします。
2. コンピュートノードを無効にします。

```
(overcloud)[stack@director ~]$ openstack compute service set \
  <host> <service> --disable
```

- **<host>** をインスタンスの退避元となるコンピュートノードの名前に置き換えてください。
- **<service>** を無効にするサービスの名前 (例: **nova-compute**) に置き換えてください。

3. 指定したコンピュートノード上の全インスタンスを退避させます。

```
(overcloud)[stack@director ~]$ nova host-evacuate [--target_host <dest>] <host>
```

- **<dest>** をインスタンスの退避先となるコンピュートノードの名前に置き換えてください。退避先を指定しなかった場合には、Compute のスケジューラーがノードを選択します。退避先に指定可能なコンピュートノードを確認するには、以下のコマンドを使用します。

```
(overcloud)[stack@director ~]$ openstack hypervisor list
```

- **<host>** をインスタンスの退避元となるコンピュートノードの名前に置き換えてください。

14.3.8. 移行に関するトラブルシューティング

インスタンスの移行時に、以下の問題が発生する可能性があります。

- 移行プロセスでエラーが生じる。
- 移行プロセスが終了しない。
- 移行後にインスタンスのパフォーマンスが低下する。

14.3.8.1. 移行中のエラー

以下の問題が発生すると、移行操作が **error** 状態に遷移します。

- 実行しているクラスターに異なるバージョンの Red Hat OpenStack Platform (RHOSP) が存在する。
- 指定したインスタンス ID が見つからない。
- 移行を試みているインスタンスが **error** 状態にある。
- Compute サービスが停止している。
- 競合状態が発生する。
- ライブマイグレーションが **failed** 状態に移行する。

ライブマイグレーションが **failed** 状態に移行すると、通常は **error** 状態になります。**failed** 状態の原因となる可能性のある典型的な問題を以下に示します。

- 移行先コンピュートホストが利用可能な状態にない。
- スケジューラーの例外が発生する。
- コンピュティングリソースが不十分なため、再ビルドプロセスに失敗する。
- サーバグループの確認に失敗する。
- 移行先コンピュートノードへの移行が完了する前に、移行元コンピュートノードのインスタンスが削除される。

14.3.8.2. ライブマイグレーションのスタック

ライブマイグレーションが完了せず、永久に **running** 状態のままになる可能性があります。ライブマイグレーションが永久に完了しない一般的な理由は、Compute サービスがインスタンスの変更を移行先コンピュートノードに複製するより早く、クライアントのリクエストにより移行元コンピュートノード上で実行中のインスタンスに変更が生じることです。

この状況に対処するには、以下のいずれかの方法を使用します。

- ライブマイグレーションを中止する。
- ライブマイグレーションを強制的に完了させる。

ライブマイグレーションの中止

移行プロセスがインスタンスの状態の変化を移行先ノードにコピーするより早くインスタンスの状態が変化する状況で、インスタンスの動作を一時的に中断したくない場合には、ライブマイグレーションを中止することができます。

手順

1. インスタンスの移行の一覧を取得します。

```
$ nova server-migration-list <instance>
```

<instance> をインスタンスの名前または ID に置き換えてください。

2. ライブマイグレーションを中止します。

```
$ nova live-migration-abort <instance> <migration_id>
```

- **<instance>** をインスタンスの名前または ID に置き換えてください。
- **<migration_id>** を移行の ID に置き換えてください。

ライブマイグレーション完了の強制

移行プロセスがインスタンスの状態の変化を移行先ノードにコピーするより早くインスタンスの状態が変化する状況で、インスタンスの動作を一時的に中断して移行を強制的に完了させたい場合には、ライブマイグレーションの手順を強制的に完了させることができます。



重要

ライブマイグレーションを強制的に完了させると、かなりのダウンタイムが発生する可能性があります。

手順

1. インスタンスの移行の一覧を取得します。

```
$ nova server-migration-list <instance>
```

<instance> をインスタンスの名前または ID に置き換えてください。

2. ライブマイグレーションを強制的に完了させます。

```
$ nova live-migration-force-complete <instance> <migration_id>
```

- **<instance>** をインスタンスの名前または ID に置き換えてください。
- **<migration_id>** を移行の ID に置き換えてください。

14.3.8.3. 移行後のインスタンスパフォーマンスの低下

NUMA トポロジーを使用するインスタンスの場合、移行元および移行先コンピュータノードの NUMA トポロジーおよび設定は同一でなければなりません。移行先コンピュータノードの NUMA トポロジーでは、十分なリソースが利用可能でなければなりません。移行元および移行先コンピュータノード間で NUMA 設定が同一でない場合、ライブマイグレーションは成功するがインスタンスのパフォーマンスが低下する可能性があります。たとえば、移行元コンピュータノードは NIC 1 を NUMA ノード 0 にマッピングするが、移行先コンピュータノードは NIC 1 を NUMA ノード 5 にマッピングする場合、移行後にインスタンスはバス内の最初の CPU からのネットワークトラフィックを NUMA ノード 5 の別の CPU にルーティングし、トラフィックを NIC 1 にルーティングする可能性があります。その結果、予想されたとおりに動作はしますが、パフォーマンスが低下します。同様に、移行元コンピュータノードの NUMA ノード 0 では十分な CPU および RAM が利用可能だが、移行先コンピュータノードの NUMA ノード 0 にリソースの一部を使用するインスタンスがすでに存在する場合、インスタンスは正しく動作するがパフォーマンスが低下する可能性があります。詳細は、[移行の制約](#) を参照してください。