

OpenShift Container Platform 4.6

スケーラビリティーおよびパフォーマンス

実稼働環境における OpenShift Container Platform クラスターのスケーリングおよび パフォーマンスチューニング

Last Updated: 2022-11-30

OpenShift Container Platform 4.6 スケーラビリティーおよびパフォーマンス

実稼働環境における OpenShift Container Platform クラスターのスケーリングおよびパフォーマンスチューニング

Enter your first name here. Enter your surname here.

Enter your organisation's name here. Enter your organisational division here.

Enter your email address here.

法律上の通知

Copyright © 2022 | You need to change the HOLDER entity in the en-US/Scalability_and_performance.ent file |.

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

本書では、クラスターをスケーリングし、OpenShift Container Platform 環境のパフォーマンスを 最適化する方法について説明します。

目次

| 第1章 大規模なクラスターのインストールに推奨されるプラクティス | . 6 |
|--|---|
| 第2章 ホストについての推奨されるプラクティス 2.1. ノードホストについての推奨プラクティス 2.2. KUBELET パラメーターを編集するための KUBELETCONFIG CRD の作成 2.3. コントロールプレーンノードのサイジング 2.3.1. Amazon Web Services (AWS) マスターインスタンスのフレーバーサイズを増やす 2.4. ETCD についての推奨されるプラクティス 2.5. ETCD データのデフラグ 2.6. OPENSHIFT CONTAINER PLATFORM インフラストラクチャーコンポーネント 2.7. モニタリングソリューションの移動 2.8. デフォルトレジストリーの移行 2.9. ルーターの移動 2.10. インフラストラクチャーノードのサイジング 2.11. 関連情報 | 7 7 8 10 11 12 13 16 17 18 19 21 |
| 第3章 クラスタースケーリングに関する推奨プラクティス 3.1. クラスターのスケーリングに関する推奨プラクティス 3.2. マシンセットの変更 3.3.1. ベアメタル上の MachineHealthCheck 3.3.2. マシンヘルスチェックのデプロイ時の制限 3.4. サンプル MACHINEHEALTHCHECK リソース 3.4.1. マシンヘルスチェックによる修復の一時停止 (short-circuiting) 3.4.1.1. 絶対値を使用した maxUnhealthy の設定 3.4.1.2. パーセンテージを使用した maxUnhealthy の設定 3.5. MACHINEHEALTHCHECK リソースの作成 | 23 23 24 25 25 26 27 28 28 28 |
| 第4章 NODE TUNING OPERATOR の使用 4.1. NODE TUNING OPERATOR について 4.2. NODE TUNING OPERATOR 仕様サンプルへのアクセス 4.3. クラスターに設定されるデフォルトのプロファイル 4.4. TUNED プロファイルが適用されていることの確認 4.5. カスタムチューニング仕様 4.6. カスタムチューニングの例 4.7. サポートされている TUNED デーモンプラグイン 第5章 クラスターローダーの使用 | 30 30 31 32 34 38 38 |
| 5.1. クラスターローダーのインストール 5.2. クラスターローダーの実行 5.3. クラスターローダーの設定 5.3.1. クラスターローダー設定ファイルの例 5.3.2. 設定フィールド 5.4. 既知の問題 | 40 40 40 41 42 45 |
| 第6章 CPU マネージャーの使用 6.1. CPU マネージャーの設定 | 46 46 |
| 第7章 TOPOLOGY MANAGER の使用 7.1. TOPOLOGY MANAGER ポリシー 7.2. TOPOLOGY MANAGER のセットアップ 7.3. POD の TOPOLOGY MANAGER ポリシーとの対話 | 51 51 52 52 |

| 第8章 CLUSTER MONITORING OPERATOR のスケーリング 8.1. PROMETHEUS データベースのストレージ要件 8.2. クラスターモニターリングの設定 | 54 54 55 |
|---|--|
| 第9章 オブジェクトの最大値に合わせた環境計画 9.1. メジャーリリースについての OPENSHIFT CONTAINER PLATFORM のテスト済みクラスターの最大値 9.2. クラスターの最大値がテスト済みの OPENSHIFT CONTAINER PLATFORM 環境および設定 9.3. テスト済みのクラスターの最大値に基づく環境計画 9.4. アプリケーション要件に合わせて環境計画を立てる方法 | 57 57 59 59 60 |
| 第10章 ストレージの最適化 10.1. 利用可能な永続ストレージオプション 10.2. 設定可能な推奨のストレージ技術 10.2.1. 特定アプリケーションのストレージの推奨事項 10.2.1.1. レジストリー 10.2.1.2. スケーリングされたレジストリー 10.2.1.3. メトリクス 10.2.1.4. ロギング 10.2.1.5. アプリケーション 10.2.2. 特定のアプリケーションおよびストレージの他の推奨事項 10.3. データストレージ管理 | 64 64 65 66 66 66 67 67 67 |
| 第11章 ルーティングの最適化 11.1. ベースライン INGRESS コントローラー (ルーター) のパフォーマンス 11.2. INGRESS コントローラー (ルーター) のパフォーマンスの最適化 | 69 69 70 |
| 第12章 ネットワークの最適化 12.1. ネットワークでの MTU の最適化 12.2. 大規模なクラスターのインストールに推奨されるプラクティス 12.3. IPSEC の影響 | 71 71 72 72 |
| 第13章 ベアメタルホストの管理 13.1. ベアメタルホストおよびノードについて 13.2. ベアメタルホストのメンテナーンス 13.2.1. Web コンソールを使用したベアメタルホストのクラスターへの追加 13.2.2. Web コンソールの YAML を使用したベアメタルホストのクラスターへの追加 13.2.3. 利用可能なベアメタルホストの数へのマシンの自動スケーリング | 73 73 73 73 74 75 |
| 第14章 HUGE PAGE の機能およびそれらがアプリケーションによって消費される仕組み 14.1. HUGE PAGE の機能 14.2. HUGE PAGE がアプリケーションによって消費される仕組み 14.3. HUGE PAGE の設定 14.3.1. ブート時 | 77 77 77 78 78 |
| 第15章低レイテンシーのノード向けの PERFORMANCE ADDON OPERATOR 15.1. 低レイテンシー 15.2. PERFORMANCE ADDON OPERATOR のインストール 15.2.1. CLI を使用した Operator のインストール 15.2.2. Web コンソールを使用した Performance Addon Operator のインストール 15.3. PERFORMANCE ADDON OPERATOR のアップグレード 15.3.1. Performance Addon Operator のアップグレードについて 15.3.1.1. Performance Addon Operator のクラスターへの影響 15.3.1.2. Performance Addon Operator の次のマイナーバージョンへのアップグレード 15.3.2. アップグレードステータスの監視 15.4. リアルタイムおよび低レイテンシーワークロードのプロビジョニング | 81 81 81 83 84 84 84 85 86 |

| 15.4.1. リアルタイムの既知の制限 | 86 |
|---|-----|
| 15.4.2. リアルタイム機能のあるワーカーのプロビジョニング | 87 |
| 15.4.3. リアルタイムカーネルのインストールの確認 | 88 |
| 15.4.4. リアルタイムで機能するワークロードの作成 | 88 |
| 15.4.5. QoS クラスの Guaranteed を指定した Pod の作成 | 89 |
| 15.4.6. オプション: DPDK 用の CPU 負荷分散の無効化 | 90 |
| 15.4.7. 適切なノードセレクターの割り当て | 91 |
| 15.4.8. リアルタイム機能を備えたワーカーへのワークロードのスケジューリング | 91 |
| 15.5. HUGE PAGE の設定 | 91 |
| 15.6. 複数の HUGE PAGE サイズの割り当て | 92 |
| 15.7. INFRA およびアプリケーションコンテナーの CPU の制限 | 93 |
| 15.8. パフォーマンスプロファイルによる低レイテンシーを実現するためのノードのチューニング | 95 |
| 15.9. プラットフォーム検証のエンドツーエンドテストの実行 | 96 |
| 15.9.1. 前提条件 | 96 |
| 15.9.2. テストの実行 | 97 |
| 15.9.3. イメージパラメーター | 97 |
| 15.9.3.1. ginkgo パラメーター | 98 |
| 15.9.3.2. 利用可能な機能 | 98 |
| 15.9.4. ドライラン | 98 |
| 15.9.5. 非接続モード | 98 |
| 15.9.5.1. クラスターからアクセスできるカスタムレジストリーへのイメージのミラーリング | 99 |
| 15.9.5.2. テストに対してカスタムレジストリーからのそれらのイメージを使用するように指示します。 | 99 |
| 15.9.5.3. クラスター内部レジストリーへのミラーリング | 99 |
| 15.9.5.4. イメージの異なるセットのミラーリング | 100 |
| 15.9.6. 検出モード | 101 |
| 15.9.6.1. 必要な環境設定の前提条件 | 101 |
| 15.9.6.2. テスト中に使用されるノードの制限 | 102 |
| 15.9.6.3. 単一のパフォーマンスプロファイルの使用 | 102 |
| 15.9.6.4. パフォーマンスプロファイルのクリーンアップの無効化 | 103 |
| 15.9.7. トラブルシューティング | 103 |
| 15.9.8. テストレポート | 103 |
| 15.9.8.1. JUnit テスト出力 | 103 |
| 15.9.8.2. テスト失敗レポート | 104 |
| 15.9.8.3. podman に関する注記 | 104 |
| 15.9.8.4. OpenShift Container Platform 4.4 での実行 | 104 |
| 15.9.8.5. 単一のパフォーマンスプロファイルの使用 | 104 |
| 15.9.9. クラスターへの影響 | 105 |
| 15.9.9.1. SCTP | 105 |
| 15.9.9.2. SR-IOV | 105 |
| 15.9.9.3. PTP | 105 |
| 15.9.9.4. パフォーマンス | 105 |
| 15.9.9.5. DPDK | 106 |
| 15.9.9.6. クリーンアップ | 106 |
| 15.10. 低レイテンシー CNF チューニングステータスのデバッグ | 106 |
| 15.10.1. マシン設定プール | 107 |
| 15.11. RED HAT サポート向けの低レイテンシーのチューニングデバッグデータの収集 | 108 |
| 15.11.1. must-gather ツールについて | 108 |
| 15.11.2. 低レイテンシーチューニングデータの収集について | 109 |
| 15.11.3. 特定の機能に関するデータ収集 | 109 |
| 第16章 INTEL FPGA PAC N3000 および INTEL VRAN DEDICATED ACCELERATOR ACC100 でのデータプンのパフォーマンスの最適化 | |

16.1. OPENSHIFT CONTAINER PLATFORM 向け INTEL ハードウェアアクセラレーターカードについて

111

| Intel FPGA PAC N3000 | 111 |
|--|-----|
| vRAN Dedicated Accelerator ACC100 | 111 |
| 16.2. INTEL FPGA PAC N3000 向け OPENNESS OPERATOR のインストール | 111 |
| 16.2.1. CLI を使用した Operator のインストール | 112 |
| 16.2.2. Web コンソールを使用した Intel FPGA PAC N3000 Operator 向け OpenNESS Operator のイン: | ストー |
| ル | 113 |
| 16.3. INTEL FPGA PAC N3000 向けの OPENNESS OPERATOR のプログラミング | 114 |
| 16.3.1. vRAN ビットストリームを持つ N3000 のプログラミング | 114 |
| 16.4. ワイヤレス FEC ACCELERATOR 向けの OPENNESS SR-IOV OPERATOR のインストール | 121 |
| 16.4.1. CLI の使用によるワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator のインストール | 121 |
| 16.4.2. Web コンソールを使用したワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator のイン. | ストー |
| ル | 123 |
| 16.4.3. Intel FPGA PAC N3000 向け SR-IOV-FEC Operator の設定 | 124 |
| 16.4.4. Intel vRAN Dedicated Accelerator ACC100 向け SR-IOV-FEC Operator の設定 | 131 |
| 16.4.5. アプリケーション Pod アクセスおよび OpenNESS での FPGA 使用量の 確認 | 139 |
| 16.5. 関連情報 | 143 |

第1章 大規模なクラスターのインストールに推奨されるプラクティス

大規模なクラスターをインストールする場合や、クラスターを大きなノード数にスケーリングする際 に、以下のプラクティスを適用します。

1.1. 大規模なクラスターのインストールに推奨されるプラクティス

大規模なクラスターをインストールする場合や、クラスターを大規模なノード数に拡張する場合、クラスターをインストールする前に、install-config.yaml ファイルに適宜クラスターネットワーク cidr を設定します。

networking:

clusterNetwork:
- cidr: 10.128.0.0/14
hostPrefix: 23

machineCIDR: 10.0.0.0/16 networkType: OpenShiftSDN

serviceNetwork: - 172.30.0.0/16

クラスターのサイズが 500 を超える場合、デフォルトのクラスターネットワーク cidr 10.128.0.0/14 を使用することはできません。500 ノードを超えるノード数にするには、10.128.0.0/12 または 10.128.0.0/10 に設定する必要があります。

第2章 ホストについての推奨されるプラクティス

このトピックでは、OpenShift Container Platform のホストについての推奨プラクティスについて説明します。



重要

これらのガイドラインは、Open Virtual Network (OVN) ではなく、ソフトウェア定義ネットワーク (SDN) を使用する OpenShift Container Platform に該当します。

2.1. ノードホストについての推奨プラクティス

OpenShift Container Platform ノードの設定ファイルには、重要なオプションが含まれています。たとえば、**podsPerCore** および **maxPods** の 2 つのパラメーターはノードにスケジュールできる Pod の最大数を制御します。

両方のオプションが使用されている場合、2つの値の低い方の値により、ノード上の Pod 数が制限されます。これらの値を超えると、以下の状態が生じる可能性があります。

- CPU 使用率の増大。
- Pod のスケジューリングの速度が遅くなる。
- (ノードのメモリー量によって)メモリー不足のシナリオが生じる可能性。
- IP アドレスのプールを消費する。
- リソースのオーバーコミット、およびこれによるアプリケーションのパフォーマンスの低下。



重要

Kubernetes では、単一コンテナーを保持する Pod は実際には 2 つのコンテナーを使用します。2 つ目のコンテナーは実際のコンテナーの起動前にネットワークを設定するために使用されます。そのため、10 の Pod を使用するシステムでは、実際には 20 のコンテナーが実行されていることになります。

podsPerCore は、ノードのプロセッサーコア数に基づいてノードが実行できる Pod 数を設定します。たとえば、4 プロセッサーコアを搭載したノードで **podsPerCore** が **10** に設定される場合、このノードで許可される Pod の最大数は **40** になります。

kubeletConfig: podsPerCore: 10

podsPerCore を 0 に設定すると、この制限が無効になります。デフォルトは 0 です。**podsPerCore** は **maxPods** を超えることができません。

maxPods は、ノードのプロパティーにかかわらず、ノードが実行できる Pod 数を固定値に設定します。

kubeletConfig: maxPods: 250

2.2. KUBELET パラメーターを編集するための KUBELETCONFIG CRD の作成

kubelet 設定は、現時点で Ignition 設定としてシリアル化されているため、直接編集することができます。ただし、新規の **kubelet-config-controller** も Machine Config Controller (MCC) に追加されます。これにより、**KubeletConfig** カスタムリソース (CR) を作成して kubelet パラメーターを編集することができます。



注記

kubeletConfig オブジェクトのフィールドはアップストリーム Kubernetes から kubelet に直接渡されるため、kubelet はそれらの値を直接検証します。**kubeletConfig** オブジェクトに無効な値により、クラスターノードが利用できなくなります。有効な値は、Kubernetes ドキュメント を参照してください。

手順

1. これは、選択可能なマシン設定オブジェクトを表示します。

\$ oc get machineconfig

デフォルトで、2 つの kubelet 関連の設定である **01-master-kubelet** および **01-worker-kubelet** を選択できます。

2. ノードあたりの最大 Pod の現在の値を確認するには、以下を実行します。

oc describe node <node-ip> | grep Allocatable -A6

value: pods: <value> を検索します。

以下は例になります。

oc describe node ip-172-31-128-158.us-east-2.compute.internal | grep Allocatable -A6

出力例

Allocatable:

attachable-volumes-aws-ebs: 25

cpu: 3500m hugepages-1Gi: 0 hugepages-2Mi: 0

memory: 15341844Ki

pods: 250

3. ワーカーノードでノードあたりの最大の Pod を設定するには、kubelet 設定を含むカスタムリソースファイルを作成します。たとえば、**change-maxPods-cr.yaml** を使用します。

apiVersion: machineconfiguration.openshift.io/v1

kind: KubeletConfig

metadata:

name: set-max-pods

spec:

machineConfigPoolSelector:

matchLabels:

custom-kubelet: large-pods

kubeletConfig: maxPods: 500

kubelet が API サーバーと通信する速度は、1秒あたりのクエリー (QPS) およびバースト値により異なります。デフォルト値の **50** (**kubeAPIQPS** の場合) および **100** (**kubeAPIBurst** の場合) は、各ノードで制限された Pod が実行されている場合には十分な値です。ノード上に CPU およびメモリーリソースが十分にある場合には、kubelet QPS およびバーストレートを更新することが推奨されます。

apiVersion: machineconfiguration.openshift.io/v1

kind: KubeletConfig

metadata:

name: set-max-pods

spec:

machineConfigPoolSelector:

matchLabels:

custom-kubelet: large-pods

kubeletConfig:

maxPods: <pod_count> kubeAPIBurst: <burst_rate> kubeAPIQPS: <QPS>

a. ラベルを使用してワーカーのマシン設定プールを更新します。

\$ oc label machineconfigpool worker custom-kubelet=large-pods

b. KubeletConfig オブジェクトを作成します。

\$ oc create -f change-maxPods-cr.yaml

c. KubeletConfig オブジェクトが作成されていることを確認します。

\$ oc get kubeletconfig

これにより set-max-pods が返されるはずです。

クラスター内のワーカーノードの数によっては、ワーカーノードが1つずつ再起動されるのを待機します。3つのワーカーノードを持つクラスターの場合は、10分から15分程度かかる可能性があります。

4. ワーカーノードを変更する maxPods の有無を確認します。

\$ oc describe node

a. 以下を実行して変更を確認します。

\$ oc get kubeletconfigs set-max-pods -o yaml

これは True と type:Success のステータスを表示します。

手順

デフォルトでは、kubelet 関連の設定を利用可能なワーカーノードに適用する場合に1つのマシンのみを利用不可の状態にすることが許可されます。大規模なクラスターの場合、設定の変更が反映されるまでに長い時間がかかる可能性があります。プロセスのスピードを上げるためにマシン数の調整をいつでも実行することができます。

1. worker マシン設定プールを編集します。

\$ oc edit machineconfigpool worker

2. maxUnavailable を必要な値に設定します。

spec:

maxUnavailable: <node_count>



重要

値を設定する際に、クラスターで実行されているアプリケーションに影響を与えずに利用不可にできるワーカーノードの数を検討してください。

2.3. コントロールプレーンノードのサイジング

コントロールプレーンノードリソースの要件は、クラスター内のノード数によって異なります。コントロールプレーンノードのサイズについての以下の推奨内容は、テストに重点を置いた場合のコントロールプレーンの密度の結果に基づいています。コントロールプレーンのテストでは、ノード数に応じて各namespaceでクラスター全体に展開される以下のオブジェクトを作成します。

- 12 イメージストリーム
- 3 ビルド設定
- 6ビルド
- それぞれに2つのシークレットをマウントする2Podレプリカのある1デプロイメント
- 2つのシークレットをマウントする 1Pod レプリカのある 2 デプロイメント
- 先のデプロイメントを参照する3つのサービス
- 先のデプロイメントを参照する3つのルート
- 10 のシークレット (それらの内の2つは先ののデプロイメントでマウントされる)
- 10 の設定マップ (それらの内の2つは先のデプロイメントでマウントされる)

| ワーカーノードの数 | クラスターの負荷 (namespace) | CPU コア数 | メモリー (GB) |
|-----------|-------------------------|---------|-----------|
| 25 | 500 | 4 | 16 |
| 100 | 1000 | 8 | 32 |
| 250 | 4000 | 16 | 96 |

3つのコントロールプレーンノード (またはマスターノード) がある大規模で高密度のクラスターでは、いずれかのノードが停止、起動、または障害が発生すると、CPU とメモリーの使用量が急上昇します。障害は、コストを節約するためにシャットダウンした後にクラスターが再起動する意図的なケースに加えて、電源、ネットワーク、または基礎となるインフラストラクチャーの予期しない問題が発生することが原因である可能性があります。残りの2つのコントロールプレーンノードは、高可用性を維持するために負荷を処理する必要があります。これにより、リソースの使用量が増えます。これは、マスターが遮断 (cordon)、ドレイン (解放) され、オペレーティングシステムおよびコントロールプレーン Operator の更新を適用するために順次再起動されるため、アップグレード時に想定される動作になります。障害が繰り返し発生しないようにするには、コントロールプレーンノードでの全体的な CPU およびメモリーリソース使用状況を、利用可能な容量の最大 60% に維持し、使用量の急増に対応できるようにします。リソース不足による潜在的なダウンタイムを回避するために、コントロールプレーンノードの CPU およびメモリーを適宜増やします。



重要

ノードのサイジングは、クラスター内のノードおよびオブジェクトの数によって異なります。また、オブジェクトがそのクラスター上でアクティブに作成されるかどうかによっても異なります。オブジェクトの作成時に、コントロールプレーンは、オブジェクトが running フェーズにある場合と比較し、リソースの使用状況においてよりアクティブな状態になります。



重要

インストール手法にインストーラーでプロビジョニングされるインフラストラクチャーを使用した場合には、実行中の OpenShift Container Platform 4.6 クラスターでコントロールプレーンノードのサイズを変更できません。その代わりに、ノードの合計数を見積もり、コントロールプレーンノードの推奨サイズをインストール時に使用する必要があります。



重要

この推奨内容は、OpenShiftSDN がネットワークプラグインとして設定されている OpenShift Container Platform クラスターでキャプチャーされるデータポイントに基づく ものです。



注記

OpenShift Container Platform 4.6 では、デフォルトで CPU コア (500 ミリコア) の半分がシステムによって予約されます (OpenShift Container Platform 3.11 以前のバージョンと比較)。サイズはこれを考慮に入れて決定されます。

2.3.1. Amazon Web Services (AWS) マスターインスタンスのフレーバーサイズを増やす

クラスター内の AWS マスターノードに過負荷がかかり、マスターノードがより多くのリソースを必要とする場合は、マスターインスタンスのフレーバーサイズを増加させることができます。



注記

AWS マスターインスタンスのフレーバーサイズを増やす前に、etcd をバックアップすることをお勧めします。

前提条件

AWS に IPI (installer-provisioned infrastructure) または UPI (user-provisioned infrastructure) クラスターがある。

手順

- 1. AWS コンソールを開き、マスターインスタンスを取得します。
- 2. 1つのマスターインスタンスを停止します。
- 3. 停止したインスタンスを選択し、Actions → Instance Settings → Change instance type をクリックします。
- 4. インスタンスをより大きなタイプに変更し、タイプが前の選択と同じベースであることを確認して、変更を適用します。たとえば、m5.xlarge を m5.2xlarge または m5.4xlarge に変更できます。
- 5. インスタンスをバックアップし、次のマスターインスタンスに対して手順を繰り返します。

関連情報

● etcd のバックアップ

2.4. ETCD についての推奨されるプラクティス

大規模で密度の高いクラスターの場合に、キースペースが過剰に拡大し、スペースのクォータを超過すると、etcd は低下するパフォーマンスの影響を受ける可能性があります。etcd を定期的に維持および最適化して、データストアのスペースを解放します。Prometheus で etcd メトリックをモニターし、必要に応じてデフラグします。そうしないと、etcd はクラスター全体のアラームを発生させ、クラスターをメンテナーンスモードにして、キーの読み取りと削除のみを受け入れる可能性があります。

これらの主要な指標をモニターします。

- etcd server quota backend bytes、これは現在のクォータ制限です
- etcd_mvcc_db_total_size_in_use_in_bytes、これはヒストリーコンパクション後の実際のデータベース使用状況を示します。
- etcd_debugging_mvcc_db_total_size_in_bytes、これはデフラグを待機している空き領域を含む、データベースのサイズを示します。

etcd の最適化の詳細については、etcd データの最適化を参照してください。

etcd はデータをディスクに書き込み、プロポーザルをディスクに保持するため、そのパフォーマンスはディスクのパフォーマンスに依存します。遅いディスクと他のプロセスからのディスクアクティビティーは、長い fsync 待ち時間を引き起こす可能性があります。これらの待ち時間により、etcd はハートビートを見逃し、新しいプロポーザルを時間どおりにディスクにコミットせず、最終的にリクエストのタイムアウトと一時的なリーダーの喪失を経験する可能性があります。低遅延と高スループットのSSD または NVMe ディスクでバックアップされたマシンで etcd を実行します。シングルレベルセル(SLC) ソリッドステートドライブ (SSD) を検討してください。これは、メモリーセルごとに1ビットを提供し、耐久性と信頼性が高く、書き込みの多いワークロードに最適です。

デプロイされた OpenShift Container Platform クラスターでモニターする主要なメトリクスの一部は、etcd ディスクの write ahead log 期間の p99 と etcd リーダーの変更数です。Prometheus を使用してこれらのメトリクスを追跡します。

- etcd_disk_wal_fsync_duration_seconds_bucket メトリックは、etcd ディスクの fsync 期間を報告します。
- etcd server leader changes seen totalメトリックは、リーダーの変更を報告します。
- 遅いディスクを除外し、ディスクが適度に速いことを確認するには、etcd_disk_wal_fsync_duration_seconds_bucket の 99 パーセンタイルが 10 ミリ秒未満であることを確認します。

OpenShift Container Platform クラスターを作成する前または後に etcd のハードウェアを検証するには、fio と呼ばれる I/O ベンチマークツールを使用できます。

前提条件

- Podman や Docker などのコンテナーランタイムは、テストしているマシンにインストールされます。
- データは /var/lib/etcd パスに書き込まれます。

手順

- fio を実行し、結果を分析します。
 - o Podman を使用する場合は、次のコマンドを実行します。

\$ sudo podman run --volume /var/lib/etcd:/var/lib/etcd:Z quay.io/openshift-scale/etcd-perf

o Docker を使用する場合は、次のコマンドを実行します。

\$ sudo docker run --volume /var/lib/etcd:/var/lib/etcd:Z quay.io/openshift-scale/etcd-perf

この出力では、実行からキャプチャーされた fsync メトリクスの 99 パーセンタイルの比較でディスクが 10 ms 未満かどうかを確認して、ディスクの速度が etcd をホストするのに十分であるかどうかを報告します。

etcd はすべてのメンバー間で要求を複製するため、そのパフォーマンスはネットワーク入出力 (I/O) のレイテンシーによって大きく変わります。ネットワークのレイテンシーが高くなると、etcd のハートビートの時間は選択のタイムアウトよりも長くなり、その結果、クラスターに中断をもたらすリーダーの選択が発生します。デプロイされた OpenShift Container Platform クラスターでのモニターの主要なメトリクスは、各 etcd クラスターメンバーの etcd ネットワークピアレイテンシーの 99 番目のパーセンタイルになります。Prometheus を使用してメトリクスを追跡します。

histogram_quantile(0.99, rate(etcd_network_peer_round_trip_time_seconds_bucket[2m])) メトリックは、etcd がメンバー間でクライアントリクエストの複製を完了するまでのラウンドトリップ時間をレポートします。50 ミリ秒未満であることを確認してください。

2.5. ETCD データのデフラグ

etcd 履歴の圧縮および他のイベントによりディスクの断片化が生じた後にディスク領域を回収するために、手動によるデフラグを定期的に実行する必要があります。

履歴の圧縮は5分ごとに自動的に行われ、これによりバックエンドデータベースにギャップが生じます。この断片化された領域は etcd が使用できますが、ホストファイルシステムでは利用できません。ホストファイルシステムでこの領域を使用できるようにするには、etcd をデフラグする必要があります。

etcd はデータをディスクに書き込むため、そのパフォーマンスはディスクのパフォーマンスに大きく依存します。etcd のデフラグは、毎月、月に 2 回、またはクラスターでの必要に応じて行うことを検討してください。etcd_db_total_size_in_bytes メトリクスをモニターして、デフラグが必要であるかどうかを判別することもできます。



警告

etcd のデフラグはプロセスを阻止するアクションです。etcd メンバーはデフラグが完了するまで応答しません。このため、各 Pod のデフラグアクションごとに少なくとも 1 分間待機し、クラスターが回復できるようにします。

以下の手順に従って、各 etcd メンバーで etcd データをデフラグします。

前提条件

• cluster-admin ロールを持つユーザーとしてクラスターにアクセスできる。

手順

- 1. リーダーを最後にデフラグする必要があるため、どの etcd メンバーがリーダーであるかを判別します。
 - a. etcd Pod の一覧を取得します。

\$ oc get pods -n openshift-etcd -o wide | grep -v quorum-guard | grep etcd

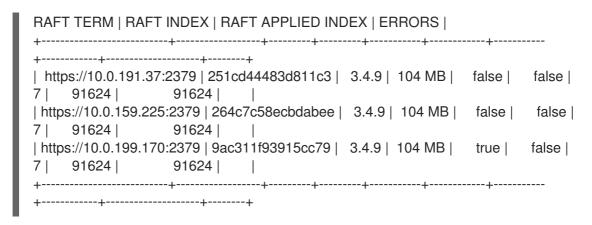
出力例

etcd-ip-10-0-159-225.example.redhat.com 3/3 Running 0 10.0.159.225 ip-10-0-159-225.example.redhat.com <none> <none> etcd-ip-10-0-191-37.example.redhat.com 3/3 Running 0 173m 10.0.191.37 ip-10-0-191-37.example.redhat.com <none> <none> etcd-ip-10-0-199-170.example.redhat.com 3/3 Running 0 176m 10.0.199.170 ip-10-0-199-170.example.redhat.com <none> <none>

b. Pod を選択し、以下のコマンドを実行して、どの etcd メンバーがリーダーであるかを判別します。

\$ oc rsh -n openshift-etcd etcd-ip-10-0-159-225.example.redhat.com etcdctl endpoint status --cluster -w table

出力例



この出力の IS LEADER 列に基づいて、https://10.0.199.170:2379 エンドポイントがリーダーになります。このエンドポイントを直前の手順の出力に一致させると、リーダーの Pod 名は etcd-ip-10-0-199-170.example.redhat.com になります。

- 2. etcd メンバーのデフラグ。
 - a. 実行中の etcd コンテナーに接続し、リーダーではない Pod の名前を渡します。

\$ oc rsh -n openshift-etcd etcd-ip-10-0-159-225.example.redhat.com

b. ETCDCTL ENDPOINTS環境変数の設定を解除します。

sh-4.4# unset ETCDCTL_ENDPOINTS

c. etcd メンバーのデフラグを実行します。

sh-4.4# etcdctl --command-timeout=30s --endpoints=https://localhost:2379 defrag

出力例

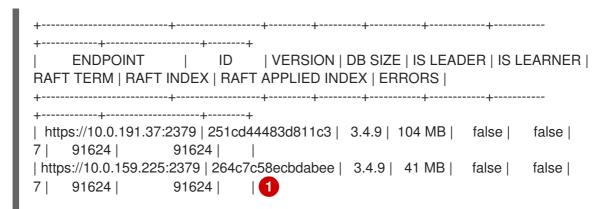
Finished defragmenting etcd member[https://localhost:2379]

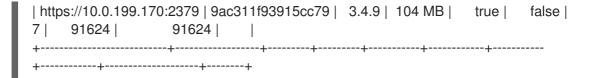
タイムアウトエラーが発生した場合は、コマンドが正常に実行されるまで --command-timeout の値を増やします。

d. データベースサイズが縮小されていることを確認します。

sh-4.4# etcdctl endpoint status -w table --cluster

出力例





この例では、この etcd メンバーのデータベースサイズは、開始時のサイズの 104~MB ではなく 41~MB です。

- e. これらの手順を繰り返して他の etcd メンバーのそれぞれに接続し、デフラグします。常に 最後にリーダーをデフラグします。 etcd Pod が回復するように、デフラグアクションごとに1分以上待機します。etcd Pod が
- 3. 領域のクォータの超過により **NOSPACE** アラームがトリガーされる場合、それらをクリアします。
 - a. NOSPACE アラームがあるかどうかを確認します。

回復するまで、etcd メンバーは応答しません。

sh-4.4# etcdctl alarm list

出力例

memberID:12345678912345678912 alarm:NOSPACE

b. アラームをクリアします。

sh-4.4# etcdctl alarm disarm

2.6. OPENSHIFT CONTAINER PLATFORM インフラストラクチャーコンポーネント

以下のインフラストラクチャーワークロードでは、OpenShift Container Platform ワーカーのサブスクリプションは不要です。

- マスターで実行される Kubernetes および OpenShift Container Platform コントロールプレーン サービス
- デフォルトルーター
- 統合コンテナーイメージレジストリー
- HAProxy ベースの Ingress Controller
- ユーザー定義プロジェクトのモニターリング用のコンポーネントを含む、クラスターメトリクスの収集またはモニターリングサービス
- クラスター集計ロギング
- サービスブローカー
- Red Hat Quay
- Red Hat OpenShift Container Storage

- Red Hat Advanced Cluster Manager
- Kubernetes 用 Red Hat Advanced Cluster Security
- Red Hat OpenShift GitOps
- Red Hat OpenShift Pipelines

他のコンテナー、Pod またはコンポーネントを実行するノードは、サブスクリプションが適用される必要のあるワーカーノードです。

関連情報

● インフラストラクチャーノードおよびインフラストラクチャーノードで実行できるコンポーネントの詳細は、OpenShift sizing and subscription guide for enterprise Kubernetes の"Red Hat OpenShift control plane and infrastructure nodes"セクションを参照してください。

2.7. モニタリングソリューションの移動

デフォルトでは、Prometheus、Grafana、および AlertManager が含まれる Prometheus Cluster Monitoring スタックはクラスターモニターリングをデプロイするためにデプロイされます。これは Cluster Monitoring Operator によって管理されます。このコンポーネント異なるマシンに移行するには、カスタム設定マップを作成し、これを適用します。

手順

1. 以下の ConfigMap 定義を cluster-monitoring-configmap.yaml ファイルとして保存します。

```
apiVersion: v1
kind: ConfigMap
metadata:
 name: cluster-monitoring-config
 namespace: openshift-monitoring
 config.yaml: |+
  alertmanagerMain:
   nodeSelector:
    node-role.kubernetes.io/infra: ""
  prometheusK8s:
   nodeSelector:
    node-role.kubernetes.io/infra: ""
  prometheusOperator:
   nodeSelector:
    node-role.kubernetes.io/infra: ""
  grafana:
   nodeSelector:
     node-role.kubernetes.io/infra: ""
  k8sPrometheusAdapter:
   nodeSelector:
     node-role.kubernetes.io/infra: ""
  kubeStateMetrics:
   nodeSelector:
     node-role.kubernetes.io/infra: ""
  telemeterClient:
   nodeSelector:
```

node-role.kubernetes.io/infra: ""
openshiftStateMetrics:
nodeSelector:
node-role.kubernetes.io/infra: ""
thanosQuerier:
nodeSelector:
node-role.kubernetes.io/infra: ""

この設定マップを実行すると、モニタリングスタックのコンポーネントがインフラストラクチャーノードに再デプロイされます。

2. 新規の設定マップを適用します。

\$ oc create -f cluster-monitoring-configmap.yaml

3. モニタリング Pod が新規マシンに移行することを確認します。

\$ watch 'oc get pod -n openshift-monitoring -o wide'

4. コンポーネントが **infra** ノードに移動していない場合は、このコンポーネントを持つ Pod を削除します。

\$ oc delete pod -n openshift-monitoring <pod>

削除された Pod からのコンポーネントが infra ノードに再作成されます。

2.8. デフォルトレジストリーの移行

レジストリー Operator を、その Pod を複数の異なるノードにデプロイするように設定します。

前提条件

● 追加のマシンセットを OpenShift Container Platform クラスターに設定します。

手順

1. config/instance オブジェクトを表示します。

\$ oc get configs.imageregistry.operator.openshift.io/cluster -o yaml

出力例

apiVersion: imageregistry.operator.openshift.io/v1

kind: Config metadata:

creationTimestamp: 2019-02-05T13:52:05Z

finalizers:

- imageregistry.operator.openshift.io/finalizer

generation: 1 name: cluster

resourceVersion: "56174"

selfLink: /apis/imageregistry.operator.openshift.io/v1/configs/cluster

uid: 36fd3724-294d-11e9-a524-12ffeee2931b

```
spec:
httpSecret: d9a012ccd117b1e6616ceccb2c3bb66a5fed1b5e481623
logging: 2
managementState: Managed
proxy: {}
replicas: 1
requests:
    read: {}
    write: {}
    storage:
        s3:
        bucket: image-registry-us-east-1-c92e88cad85b48ec8b312344dff03c82-392c
    region: us-east-1
status:
...
```

2. config/instance オブジェクトを編集します。

\$ oc edit configs.imageregistry.operator.openshift.io/cluster

3. オブジェクトの spec セクションを以下の YAML のように変更します。

```
spec:
affinity:
podAntiAffinity:
preferredDuringSchedulingIgnoredDuringExecution:
- podAffinityTerm:
namespaces:
- openshift-image-registry
topologyKey: kubernetes.io/hostname
weight: 100
logLevel: Normal
managementState: Managed
nodeSelector:
node-role.kubernetes.io/infra: ""
```

- 4. レジストリー Pod がインフラストラクチャーノードに移動していることを確認します。
 - a. 以下のコマンドを実行して、レジストリー Pod が置かれているノードを特定します。

\$ oc get pods -o wide -n openshift-image-registry

b. ノードに指定したラベルがあることを確認します。

\$ oc describe node <node_name>

コマンド出力を確認し、**node-role.kubernetes.io/infra** が **LABELS** 一覧にあることを確認します。

2.9. ルーターの移動

ルーター Pod を異なるマシンセットにデプロイできます。デフォルトで、この Pod はワーカーノードにデプロイされます。

前提条件

• 追加のマシンセットを OpenShift Container Platform クラスターに設定します。

手順

1. ルーター Operator の IngressController カスタムリソースを表示します。

\$ oc get ingresscontroller default -n openshift-ingress-operator -o yaml

コマンド出力は以下のテキストのようになります。

```
apiVersion: operator.openshift.io/v1
kind: IngressController
metadata:
 creationTimestamp: 2019-04-18T12:35:39Z
 finalizers:
 - ingresscontroller.operator.openshift.io/finalizer-ingresscontroller
 generation: 1
 name: default
 namespace: openshift-ingress-operator
 resourceVersion: "11341"
 selfLink: /apis/operator.openshift.io/v1/namespaces/openshift-ingress-
operator/ingresscontrollers/default
 uid: 79509e05-61d6-11e9-bc55-02ce4781844a
spec: {}
status:
 availableReplicas: 2
 conditions:
 - lastTransitionTime: 2019-04-18T12:36:15Z
  status: "True"
  type: Available
 domain: apps.<cluster>.example.com
 endpointPublishingStrategy:
  type: LoadBalancerService
 selector: ingresscontroller.operator.openshift.io/deployment-ingresscontroller=default
```

2. ingresscontroller リソースを編集し、 nodeSelector を infra ラベルを使用するように変更します。

\$ oc edit ingresscontroller default -n openshift-ingress-operator

以下に示すように、infra ラベルを参照する nodeSelector スタンザを spec セクションに追加します。

```
spec:
nodePlacement:
nodeSelector:
matchLabels:
node-role.kubernetes.io/infra: ""
```

- 3. ルーター Pod が infra ノードで実行されていることを確認します。
 - a. ルーター Pod の一覧を表示し、実行中の Pod のノード名をメモします。

\$ oc get pod -n openshift-ingress -o wide

出力例

NAME READY STATUS RESTARTS AGE IΡ **NODE** NOMINATED NODE READINESS GATES router-default-86798b4b5d-bdlvd 1/1 10.130.2.4 ip-10-28s Running 0 0-217-226.ec2.internal <none> <none> router-default-955d875f4-255g8 0/1 Terminating 0 19h 10.129.2.4 ip-10-0-148-172.ec2.internal <none> <none>

この例では、実行中の Pod は ip-10-0-217-226.ec2.internal ノードにあります。

b. 実行中の Pod のノードのステータスを表示します。

\$ oc get node <node_name> 1

111Pod の一覧より取得した **<node_name>** を指定します。

出力例

NAME STATUS ROLES AGE VERSION ip-10-0-217-226.ec2.internal Ready infra,worker 17h v1.19.0

ロールの一覧に infra が含まれているため、Pod は正しいノードで実行されます。

2.10. インフラストラクチャーノードのサイジング

これらの要素により、Prometheusのメトリクスまたは時系列の数が増加する可能性があり、インフラストラクチャーノードのリソース要件はクラスターのクラスターの使用年数、ノード、およびオブジェクトによって異なります。以下のインフラストラクチャーノードのサイズの推奨内容は、クラスターの最大値およびコントロールプレーンの密度に重点を置いたテストの結果に基づいています。



重要

以下のサイジングの推奨内容は、クラスターのインストール時にインストールされるインフラストラクチャーコンポーネント (Prometheus、ルーターおよびレジストリー) についてのみ適用されます。ロギングは Day 2 オペレーションであり、そのサイジングの推奨内容は最後の 2 つの列に記載されます。

| ワーカーノードの 数 | CPU コア数 | メモリー (GB) | CPU コア (Logging) | メモリー (GB) (Logging) |
|---------------|---------|-----------|---------------------|------------------------|
| 25 | 4 | 16 | 4 | 64 |
| 100 | 8 | 32 | 8 | 128 |
| 250 | 16 | 128 | 16 | 128 |

| ワーカーノードの 数 | CPU コア数 | メモリー (GB) | CPU コア (Logging) | メモリー (GB) (Logging) |
|---------------|---------|-----------|---------------------|------------------------|
| 500 | 32 | 128 | 32 | 192 |



重要

これらのサイジングの推奨内容は、クラスター全体に多数のオブジェクトを作成するスケーリングのテストに基づいています。これらのテストでは、一部のクラスターの最大値に達しいます。OpenShift Container Platform 4.6 クラスターでノード数が 250 および 500 の場合、これらの最大値は、10000 namespace および 61000 Pod、10000 デプロイメント、181000 シークレット、400 設定マップなどになります。Prometheus はメモリー集約型のアプリケーションであり、リソースの使用率はノード数、オブジェクト数、Prometheus メトリクスの収集間隔、メトリクスまたは時系列、クラスターの使用年数などのさまざまな要素によって異なります。ディスクサイズは保持期間によっても変わります。これらの要素を考慮し、これらに応じてサイズを設定する必要があります。



注記

OpenShift Container Platform 4.6 では、デフォルトで CPU コア (500 ミリコア) の半分がシステムによって予約されます (OpenShift Container Platform 3.11 以前のバージョンと比較)。これは、上記のサイジングの推奨内容に影響します。

2.11. 関連情報

● OpenShift Container Platform クラスターの最大値

第3章 クラスタースケーリングに関する推奨プラクティス



重要

本セクションのガイダンスは、クラウドプロバイダーの統合によるインストールにのみ関連します。

これらのガイドラインは、Open Virtual Network (OVN) ではなく、ソフトウェア定義ネットワーク (SDN) を使用する OpenShift Container Platform に該当します。

以下のベストプラクティスを適用して、OpenShift Container Platform クラスター内のワーカーマシンの数をスケーリングします。ワーカーのマシンセットで定義されるレプリカ数を増やしたり、減らしたりしてワーカーマシンをスケーリングします。

3.1. クラスターのスケーリングに関する推奨プラクティス

クラスターをノード数のより高い値にスケールアップする場合:

- 高可用性を確保するために、ノードを利用可能なすべてのゾーンに分散します。
- 1度に25未満のマシンごとに50マシンまでスケールアップします。
- 定期的なプロバイダーの容量関連の制約を軽減するために、同様のサイズの別のインスタンスタイプを使用して、利用可能なゾーンごとに新規のマシンセットを作成することを検討してください。たとえば、AWS で、m5.large および m5d.large を使用します。



注記

クラウドプロバイダーは API サービスのクォータを実装する可能性があります。そのため、クラスターは段階的にスケーリングします。

マシンセットのレプリカが1度に高い値に設定される場合に、コントローラーはマシンを作成できなくなる可能性があります。OpenShift Container Platform が上部にデプロイされているクラウドプラットフォームが処理できる要求の数はプロセスに影響を与えます。コントローラーは、該当するステータスのマシンの作成、確認、および更新を試行する間に、追加のクエリーを開始します。OpenShift Container Platform がデプロイされるクラウドプラットフォームには API 要求の制限があり、過剰なクエリーが生じると、クラウドプラットフォームの制限によりマシンの作成が失敗する場合があります。

大規模なノード数にスケーリングする際にマシンヘルスチェックを有効にします。障害が発生する場合、ヘルスチェックは状態を監視し、正常でないマシンを自動的に修復します。



注記

大規模で高密度のクラスターをノード数を減らしてスケールダウンする場合には、長い時間がかかる可能性があります。このプロセスで、終了するノードで実行されているオブジェクトのドレイン (解放) またはエビクトが並行して実行されるためです。また、エビクトするオブジェクトが多過ぎる場合に、クライアントはリクエストのスロットリングを開始する可能性があります。デフォルトのクライアント QPS およびバーストレートは、現時点で 5 と 10 にそれぞれ設定されています。これらは OpenShift Container Platform で変更することはできません。

3.2. マシンセットの変更

マシンセットを変更するには、MachineSet YAML を編集します。次に、各マシンを削除するか、またはマシンセットを0レプリカにスケールダウンしてマシンセットに関連付けられたすべてのマシンを削除します。レプリカは必要な数にスケーリングします。マシンセットへの変更は既存のマシンに影響を与えません。

他の変更を加えずに、マシンセットをスケーリングする必要がある場合、マシンを削除する必要はありません。



注記

デフォルトで、OpenShift Container Platform ルーター Pod はワーカーにデプロイされます。ルーターは Web コンソールなどの一部のクラスターリソースにアクセスすることが必要であるため、 ルーター Pod をまず再配置しない限り、ワーカーのマシンセットを $\mathbf{0}$ にスケーリングできません。

前提条件

- OpenShift Container Platform クラスターおよび oc コマンドラインをインストールすること。
- cluster-admin パーミッションを持つユーザーとして、oc にログインする。

手順

- 1. マシンセットを編集します。
 - \$ oc edit machineset <machineset> -n openshift-machine-api
- 2. マシンセットを 0 にスケールダウンします。
 - \$ oc scale --replicas=0 machineset <machineset> -n openshift-machine-api

または、以下を実行します。

\$ oc edit machineset <machineset> -n openshift-machine-api

マシンが削除されるまで待機します。

3. マシンセットを随時スケールアップします。

\$ oc scale --replicas=2 machineset <machineset> -n openshift-machine-api

または、以下を実行します。

\$ oc edit machineset <machineset> -n openshift-machine-api

マシンが起動するまで待ちます。新規マシンにはマシンセットに加えられた変更が含まれます。

3.3. マシンのヘルスチェック

MachineHealthCheck リソースを使用して、クラスター内のマシンが正常ではないとみなされる条件を定義できます。条件に一致するマシンは自動的に修復されます。

マシンの正常性を監視するには、監視する一連のマシンのラベルや、NotReady ステータスの期間を 15分にしたり、 node-problem-detector に永続的な条件を表示したりするなど、チェックする条件を含む MachineHealthCheck カスタムリソース (CR) を作成します。

MachineHealthCheck CR を観察するコントローラーは定義した条件の有無をチェックします。マシンがヘルスチェックに失敗した場合、このマシンは自動的に検出され、新規マシンが代わりに作成されます。マシンが削除されると、machine deleted イベントが表示されます。



注記

マスターロールを持つマシンの場合、マシンのヘルスチェックは正常でないノードの数を報告しますが、マシンは削除されません。以下は例になります。

出力例

\$ oc get machinehealthcheck example -n openshift-machine-api

NAME MAXUNHEALTHY EXPECTEDMACHINES CURRENTHEALTHY example 40% 3 1

マシンの削除による破壊的な影響を制限するために、コントローラーは1度に1つのノードのみをドレイン (解放) し、これを削除します。マシンのターゲットプールで許可される maxUnhealthy しきい値を上回る数の正常でないマシンがある場合、コントローラーはマシンの削除を停止し、手動で介入する必要があります。

チェックを停止するには、カスタムリソースを削除します。

3.3.1. ベアメタル上の MachineHealthCheck

ベアメタルクラスターでのマシンの削除により、ベアメタルホストの再プロビジョニングがトリガーされます。通常、ベアメタルの再プロビジョニングは長いプロセスで、クラスターにコンピュートリソースがなくなり、アプリケーションが中断される可能性があります。デフォルトの修復プロセスをマシンの削除からホストの電源サイクルに切り換えるには、MachineHealthCheck リソースにmachine.openshift.io/remediation-strategy: external-baremetal アノテーションを付けます。

アノテーションの設定後に、BMC 認証情報を使用して正常でないマシンの電源が入れ直されます。

3.3.2. マシンヘルスチェックのデプロイ時の制限

マシンヘルスチェックをデプロイする前に考慮すべき制限事項があります。

- マシンセットが所有するマシンのみがマシンヘルスチェックによって修復されます。
- コントロールプレーンマシンは現在サポートされておらず、それらが正常でない場合にも修正 されません。
- マシンのノードがクラスターから削除される場合、マシンヘルスチェックはマシンが正常では ないとみなし、すぐにこれを修復します。
- **nodeStartupTimeout** の後にマシンの対応するノードがクラスターに加わらない場合、マシン は修復されます。
- Machine リソースフェーズが Failed の場合、マシンはすぐに修復されます。

3.4. サンプル MACHINEHEALTHCHECK リソース

MachineHealthCheck リソースは以下の YAML ファイルのようになります。

ベアメタルの MachineHealthCheck

apiVersion: machine.openshift.io/v1beta1

kind: MachineHealthCheck

metadata:

name: example 1

namespace: openshift-machine-api

annotations:

machine.openshift.io/remediation-strategy: external-baremetal 2

spec:

selector:

matchLabels:

machine.openshift.io/cluster-api-machine-role: <role> 3

machine.openshift.io/cluster-api-machine-type: <role> 4

machine.openshift.io/cluster-api-machineset: <cluster_name>-<label>-<zone> 5

unhealthy Conditions:

- type: "Ready"

timeout: "300s" 6

status: "False"

- type: "Ready"

timeout: "300s" 7

status: "Unknown"

maxUnhealthy: "40%" 8

nodeStartupTimeout: "10m" 9

- デプロイするマシンヘルスチェックの名前を指定します。
- ベアメタルクラスターの場合、電源サイクルの修復を有効にするために machine.openshift.io/remediation-strategy: external-baremetal アノテーションを annotations セクションに含める必要があります。この修復ストラテジーにより、正常でないホストはクラス ターから削除される代わりに、再起動されます。
- 341チェックする必要のあるマシンプールのラベルを指定します。
- 5 追跡するマシンセットを <cluster_name>-<label>-<zone> 形式で指定します。たとえば、prodnode-us-east-1a とします。
- 67ノードの状態のタイムアウト期間を指定します。タイムアウト期間の条件が満たされると、マシンは修正されます。タイムアウトの時間が長くなると、正常でないマシンのワークロードのダウンタイムが長くなる可能性があります。
- 8 ターゲットプールで同時に修復できるマシンの数を指定します。これはパーセンテージまたは整数として設定できます。正常でないマシンの数が maxUnhealthy で設定された制限を超える場合、修復は実行されません。
- マシンが正常でないと判別される前に、ノードがクラスターに参加するまでマシンヘルスチェック が待機する必要のあるタイムアウト期間を指定します。



注記

matchLabels はあくまでもサンプルであるため、特定のニーズに応じてマシングループをマッピングする必要があります。

他のすべてのインストールタイプの MachineHealthCheck

apiVersion: machine.openshift.io/v1beta1

kind: MachineHealthCheck

metadata:

name: example 1

namespace: openshift-machine-api

spec:

selector:

matchLabels:

machine.openshift.io/cluster-api-machine-role: <role> 2

machine.openshift.io/cluster-api-machine-type: <role> 3

machine.openshift.io/cluster-api-machineset: <cluster_name>-<label>-<zone> 4 unhealthyConditions:

- type: "Ready"

timeout: "300s" 5

status: "False"

- type: "Ready"

timeout: "300s" 6

status: "Unknown"

maxUnhealthy: "40%" 7

nodeStartupTimeout: "10m" (8)

- デプロイするマシンヘルスチェックの名前を指定します。
- 23チェックする必要のあるマシンプールのラベルを指定します。
- 4 追跡するマシンセットを <cluster_name>-<label>-<zone> 形式で指定します。たとえば、prodnode-us-east-1a とします。
- 56/1ードの状態のタイムアウト期間を指定します。タイムアウト期間の条件が満たされると、マシンは修正されます。タイムアウトの時間が長くなると、正常でないマシンのワークロードのダウンタイムが長くなる可能性があります。
- ターゲットプールで同時に修復できるマシンの数を指定します。これはパーセンテージまたは整数 として設定できます。正常でないマシンの数が maxUnhealthy で設定された制限を超える場合、 修復は実行されません。
- 8 マシンが正常でないと判別される前に、ノードがクラスターに参加するまでマシンヘルスチェックが待機する必要のあるタイムアウト期間を指定します。



注記

matchLabels はあくまでもサンプルであるため、特定のニーズに応じてマシングループをマッピングする必要があります。

3.4.1. マシンヘルスチェックによる修復の一時停止 (short-circuiting)

一時停止 (short-circuiting) が実行されることにより、マシンのヘルスチェックはクラスターが正常な場合にのみマシンを修復するようになります。一時停止 (short-circuiting) は、**MachineHealthCheck** リソースの **maxUnhealthy** フィールドで設定されます。

ユーザーがマシンの修復前に maxUnhealthy フィールドの値を定義する場合、 MachineHealthCheck は maxUnhealthy の値を、正常でないと判別するターゲットプール内のマシン数と比較します。正常でないマシンの数が maxUnhealthy の制限を超える場合、修復は実行されません。



重要

maxUnhealthy が設定されていない場合、値は 100% にデフォルト設定され、マシンは クラスターの状態に関係なく修復されます。

適切な maxUnhealthy 値は、デプロイするクラスターの規模や、MachineHealthCheck が対応するマシンの数によって異なります。たとえば、maxUnhealthy 値を使用して複数のアベイラビリティーゾーン間で複数のマシンセットに対応でき、ゾーン全体が失われると、maxUnhealthy の設定によりクラスター内で追加の修復を防ぐことができます。

maxUnhealthy フィールドは整数またはパーセンテージのいずれかに設定できます。maxUnhealthy の値によって、修復の実装が異なります。

3.4.1.1. 絶対値を使用した maxUnhealthy の設定

maxUnhealthy が 2 に設定される場合:

- 2つ以下のノードが正常でない場合に、修復が実行されます。
- 3つ以上のノードが正常でない場合は、修復は実行されません。

これらの値は、マシンヘルスチェックによってチェックされるマシン数と別個の値です。

3.4.1.2. パーセンテージを使用した maxUnhealthy の設定

maxUnhealthy が 40% に設定され、25 のマシンがチェックされる場合:

- 10 以下のノードが正常でない場合に、修復が実行されます。
- 11以上のノードが正常でない場合は、修復は実行されません。

maxUnhealthy が 40% に設定され、6 マシンがチェックされる場合:

- 2つ以下のノードが正常でない場合に、修復が実行されます。
- 3つ以上のノードが正常でない場合は、修復は実行されません。



注記

チェックされる maxUnhealthy マシンの割合が整数ではない場合、マシンの許可される数は切り捨てられます。

3.5. MACHINEHEALTHCHECK リソースの作成

クラスターに、すべての MachineSets の MachineHealthCheck リソースを作成できます。コントロールプレーンマシンをターゲットとする MachineHealthCheck リソースを作成することはできません。

前提条件

• oc コマンドラインインターフェイスをインストールします。

手順

- 1. マシンヘルスチェックの定義を含む healthcheck.yml ファイルを作成します。
- 2. healthcheck.yml ファイルをクラスターに適用します。

\$ oc apply -f healthcheck.yml

第4章 NODE TUNING OPERATOR の使用

Node Tuning Operator について説明し、この Operator を使用し、Tuned デーモンのオーケストレーションを実行してノードレベルのチューニングを管理する方法について説明します。

4.1. NODE TUNING OPERATOR について

Node Tuning Operator は、Tuned デーモンのオーケストレーションによるノードレベルのチューニングの管理に役立ちます。ほとんどの高パフォーマンスアプリケーションでは、一定レベルのカーネルのチューニングが必要です。Node Tuning Operator は、ノードレベルの sysctl の統一された管理インターフェイスをユーザーに提供し、ユーザーが指定するカスタムチューニングを追加できるよう柔軟性を提供します。

Operator は、コンテナー化された OpenShift Container Platform の Tuned デーモンを Kubernetes デーモンセットとして管理します。これにより、カスタムチューニング仕様が、デーモンが認識する形式でクラスターで実行されるすべてのコンテナー化された Tuned デーモンに渡されます。デーモンは、ノードごとに1つずつ、クラスターのすべてのノードで実行されます。

コンテナー化された Tuned デーモンによって適用されるノードレベルの設定は、プロファイルの変更をトリガーするイベントで、または終了シグナルの受信および処理によってコンテナー化された Tuned デーモンが正常に終了する際にロールバックされます。

Node Tuning Operator は、バージョン 4.1 以降における標準的な OpenShift Container Platform インストールの一部となっています。

4.2. NODE TUNING OPERATOR 仕様サンプルへのアクセス

このプロセスを使用して Node Tuning Operator 仕様サンプルにアクセスします。

手順

1. 以下を実行します。

\$ oc get Tuned/default -o yaml -n openshift-cluster-node-tuning-operator

デフォルトの CR は、OpenShift Container Platform プラットフォームの標準的なノードレベルのチューニングを提供することを目的としており、Operator 管理の状態を設定するためにのみ変更できます。デフォルト CR へのその他のカスタム変更は、Operator によって上書きされます。カスタムチューニングの場合は、独自のチューニングされた CR を作成します。新規に作成された CR は、ノード/Pod ラベルおよびプロファイルの優先順位に基づいて OpenShift Container Platform ノードに適用されるデフォルトの CR およびカスタムチューニングと組み合わされます。



警告

特定の状況で Pod ラベルのサポートは必要なチューニングを自動的に配信する便利な方法ですが、この方法は推奨されず、とくに大規模なクラスターにおいて注意が必要です。デフォルトの調整された CR は Pod ラベル一致のない状態で提供されます。カスタムプロファイルが Pod ラベル一致のある状態で作成される場合、この機能はその時点で有効になります。Pod ラベル機能は、Node Tuning Operatorの今後のバージョンで非推奨になる場合があります。

4.3. クラスターに設定されるデフォルトのプロファイル

以下は、クラスターに設定されるデフォルトのプロファイルです。

```
apiVersion: tuned.openshift.io/v1
kind: Tuned
metadata:
 name: default
 namespace: openshift-cluster-node-tuning-operator
spec:
 profile:
 - name: "openshift"
  data: |
   [main]
   summary=Optimize systems running OpenShift (parent profile)
   include=${f:virt_check:virtual-guest:throughput-performance}
   [selinux]
   avc_cache_threshold=8192
   nf conntrack hashsize=131072
   [sysctl]
   net.ipv4.ip forward=1
   kernel.pid max=>4194304
   net.netfilter.nf conntrack max=1048576
   net.ipv4.conf.all.arp_announce=2
   net.ipv4.neigh.default.gc_thresh1=8192
   net.ipv4.neigh.default.gc_thresh2=32768
   net.ipv4.neigh.default.gc thresh3=65536
   net.ipv6.neigh.default.gc_thresh1=8192
   net.ipv6.neigh.default.gc_thresh2=32768
   net.ipv6.neigh.default.gc_thresh3=65536
   vm.max map count=262144
   [sysfs]
   /sys/module/nvme_core/parameters/io_timeout=4294967295
   /sys/module/nvme_core/parameters/max_retries=10
```

name: "openshift-control-plane"

data: |

```
[main]
  summary=Optimize systems running OpenShift control plane
  include=openshift
  [sysctl]
  # ktune sysctl settings, maximizing i/o throughput
  # Minimal preemption granularity for CPU-bound tasks:
  # (default: 1 msec# (1 + ilog(ncpus)), units: nanoseconds)
  kernel.sched min granularity ns=10000000
  # The total time the scheduler will consider a migrated process
  # "cache hot" and thus less likely to be re-migrated
  # (system default is 500000, i.e. 0.5 ms)
  kernel.sched_migration_cost_ns=5000000
  # SCHED_OTHER wake-up granularity.
  # Preemption granularity when tasks wake up. Lower the value to
  # improve wake-up latency and throughput for latency critical tasks.
  kernel.sched wakeup granularity ns=4000000
- name: "openshift-node"
 data: |
  [main]
  summary=Optimize systems running OpenShift nodes
  include=openshift
  [sysctl]
  net.ipv4.tcp_fastopen=3
  fs.inotify.max user watches=65536
  fs.inotify.max_user_instances=8192
recommend:
- profile: "openshift-control-plane"
 priority: 30
 match:
 - label: "node-role.kubernetes.io/master"
 - label: "node-role.kubernetes.io/infra"
- profile: "openshift-node"
 priority: 40
```

4.4. TUNED プロファイルが適用されていることの確認

この手順を使用して、すべてのノードに適用される Tuned プロファイルを確認します。

手順

1. 各ノードで実行されている Tuned Pod を確認します。

\$ oc get pods -n openshift-cluster-node-tuning-operator -o wide

出力例

NAME

READY STATUS RESTARTS AGE IP

NODE

| NOMINATED NODE READINE cluster-node-tuning-operator-599 ip-10-0-145-113.eu-west-3.comp | 9489d4f7 | -k4hw4 | | Running <none< th=""><th></th><th>d2h</th><th>10.129.0.76</th></none<> | | d2h | 10.129.0.76 |
|--|--|--|---|--|------------|-------|---------------------------|
| tuned-2jkzp | | | 1 | _ | | 113 | ip-10-0-145- |
| 113.eu-west-3.compute.internal | | 0 | one> | odon | 10.0.145.1 | 110 1 | ip-10-0-1 4 5- |
| tuned-g9mkx | 1/1 | Running | 1 | 6d3h | 10.0.147 | '.108 | ip-10-0- |
| 147-108.eu-west-3.compute.inte | rnal <no< td=""><td>one></td><td><nc< td=""><td>ne></td><td></td><td></td><td></td></nc<></td></no<> | one> | <nc< td=""><td>ne></td><td></td><td></td><td></td></nc<> | ne> | | | |
| tuned-kbxsh | 1/1 | Running | 1 | 6d3h | 10.0.132. | 143 | ip-10-0-132- |
| 143.eu-west-3.compute.internal | <none></none> | <n< td=""><td>one></td><td></td><td></td><td></td><td></td></n<> | one> | | | | |
| tuned-kn9x6 | 1/1 | Running | 1 | 6d3h | 10.0.163. | .177 | ip-10-0-163- |
| 177.eu-west-3.compute.internal | <none></none> | <n< td=""><td>one></td><td></td><td></td><td></td><td></td></n<> | one> | | | | |
| tuned-vvxwx | 1/1 | Running | 1 | 6d3h | 10.0.131. | .87 | ip-10-0-131- |
| 87.eu-west-3.compute.internal | <none></none> | <no< td=""><td>ne></td><td></td><td></td><td></td><td></td></no<> | ne> | | | | |
| tuned-zqrwq | 1/1 | Running | 1 | 6d3h | 10.0.161. | .51 | ip-10-0-161- |
| 51.eu-west-3.compute.internal | <none></none> | <nc< td=""><td>ne></td><td></td><td></td><td></td><td></td></nc<> | ne> | | | | |

2. 各 Pod から適用されるプロファイルを抽出し、それらを直前の一覧に対して一致させます。

 $\$ for p in `oc get pods -n openshift-cluster-node-tuning-operator -l openshift-app=tuned -o=jsonpath='{range .items[*]}{.metadata.name} {end}'`; do printf "\n*** \$p ***\n" ; oc logs pod/\$p -n openshift-cluster-node-tuning-operator | grep applied; done

出力例

| *** tuned-2jkzp *** 2020-07-10 13:53:35,368 INFO 'openshift-control-plane' applied | tuned.daemon.daemon: static tuning from profile |
|---|---|
| *** tuned-g9mkx *** 2020-07-10 14:07:17,089 INFO 'openshift-node' applied 2020-07-10 15:56:29,005 INFO 'openshift-node-es' applied 2020-07-10 16:00:19,006 INFO 'openshift-node' applied 2020-07-10 16:00:48,989 INFO 'openshift-node-es' applied | tuned.daemon.daemon: static tuning from profile |
| *** tuned-kbxsh *** 2020-07-10 13:53:30,565 INFO 'openshift-node' applied 2020-07-10 15:56:30,199 INFO 'openshift-node-es' applied | tuned.daemon.daemon: static tuning from profile tuned.daemon.daemon: static tuning from profile |
| *** tuned-kn9x6 *** 2020-07-10 14:10:57,123 INFO 'openshift-node' applied 2020-07-10 15:56:28,757 INFO 'openshift-node-es' applied | tuned.daemon.daemon: static tuning from profile tuned.daemon.daemon: static tuning from profile |
| *** tuned-vvxwx *** 2020-07-10 14:11:44,932 INFO 'openshift-control-plane' applied | tuned.daemon.daemon: static tuning from profile |

*** tuned-zqrwq ***
2020-07-10 14:07:40,246 INFO tuned.daemon.daemon: static tuning from profile 'openshift-control-plane' applied

4.5. カスタムチューニング仕様

Operator のカスタムリソース (CR) には 2 つの重要なセクションがあります。1 つ目のセクションの **profile:** は Tuned プロファイルおよびそれらの名前の一覧です。2 つ目の **recommend:** は、プロファイル選択ロジックを定義します。

複数のカスタムチューニング仕様は、Operator の namespace に複数の CR として共存できます。新規 CR の存在または古い CR の削除は Operator によって検出されます。既存のカスタムチューニング仕様はすべてマージされ、コンテナー化された Tuned デーモンの適切なオブジェクトは更新されます。

管理状態

Operator 管理の状態は、デフォルトの Tuned CR を調整して設定されます。デフォルトで、Operator は Managed 状態であり、**spec.managementState** フィールドはデフォルトの Tuned CR に表示されません。Operator Management 状態の有効な値は以下のとおりです。

- Managed: Operator は設定リソースが更新されるとそのオペランドを更新します。
- Unmanaged: Operator は設定リソースへの変更を無視します。
- Removed: Operator は Operator がプロビジョニングしたオペランドおよびリソースを削除します。

プロファイルデータ

profile: セクションは、Tuned プロファイルおよびそれらの名前を一覧表示します。

```
profile:
- name: tuned_profile_1
data: |
# Tuned profile specification
[main]
summary=Description of tuned_profile_1 profile

[sysctl]
net.ipv4.ip_forward=1
# ... other sysctl's or other Tuned daemon plugins supported by the containerized Tuned

# ...
- name: tuned_profile_n
data: |
# Tuned profile specification
[main]
summary=Description of tuned_profile_n profile

# tuned_profile_n profile_settings
```

推奨プロファイル

profile: 選択ロジックは、CR の **recommend:** セクションによって定義されます。**recommend:** セクションは、選択基準に基づくプロファイルの推奨項目の一覧です。

recommend: <recommend-item-1> # ... <recommend-item-n>

一覧の個別項目:

- machineConfigLabels: 1

 <mcLabels> 2

 match: 3

 <match> 4

 priority: <pri>priority> 5
 - profile: <tuned_profile_name> 6
- 🚹 オプション:
- 🔈 キー/値の MachineConfig ラベルのディクショナリー。キーは一意である必要があります。
- 3 省略する場合は、優先度の高いプロファイルが最初に一致するか、または machineConfigLabels が設定されていない限り、プロファイルの一致が想定されます。
- 4 オプションの一覧。
- 5 プロファイルの順序付けの優先度。数値が小さいほど優先度が高くなります (**0** が最も高い優先度になります)。

<match>は、以下のように再帰的に定義されるオプションの一覧です。

- label: <label_name> 1
 value: <label_value> 2
 type: <label_type> 3
 <match> 4
- ✓ ノードまたは Pod のラベル名。
- 2 オプションのノードまたは Pod のラベルの値。省略されている場合も、**<label_name>** があるだけで一致条件を満たします。
- 3 オプションのオブジェクトタイプ (node または pod)。省略されている場合は、node が想定されます。
- 🕢 オプションの <match> 一覧。

<match> が省略されない場合、ネストされたすべての **<match>** セクションが **true** に評価される必要 もあります。そうでない場合には **false** が想定され、それぞれの **<match>** セクションのあるプロファイルは適用されず、推奨されません。そのため、ネスト化 (子の **<match>** セクション) は論理 AND 演算子として機能します。これとは逆に、**<match>** 一覧のいずれかの項目が一致する場合、**<match>** の一覧全体が **true** に評価されます。そのため、一覧は論理 OR 演算子として機能します。

machineConfigLabels が定義されている場合、マシン設定プールベースのマッチングが指定のrecommend: 一覧の項目に対してオンになります。<mcLabels> はマシン設定のラベルを指定します。マシン設定は、プロファイル <tuned_profile_name> についてカーネル起動パラメーターなどのホスト設定を適用するために自動的に作成されます。この場合、マシン設定セレクターが <mcLabels> に一致するすべてのマシン設定プールを検索し、プロファイル <tuned_profile_name> を確認されるマシン設定プールが割り当てられるすべてのノードに設定する必要があります。マスターロールとワーカーのロールの両方を持つノードをターゲットにするには、マスターロールを使用する必要があります。

一覧項目の match および machineConfigLabels は論理 OR 演算子によって接続されます。 match 項目は、最初にショートサーキット方式で評価されます。そのため、true と評価される場合、machineConfigLabels 項目は考慮されません。



重要

マシン設定プールベースのマッチングを使用する場合、同じハードウェア設定を持つノードを同じマシン設定プールにグループ化することが推奨されます。この方法に従わない場合は、チューニングされたオペランドが同じマシン設定プールを共有する2つ以上のノードの競合するカーネルパラメーターを計算する可能性があります。

例: ノード/Pod ラベルベースのマッチング

- match:
- label: tuned.openshift.io/elasticsearch match:
- label: node-role.kubernetes.io/masterlabel: node-role.kubernetes.io/infra

type: pod priority: 10

profile: openshift-control-plane-es

- match:
- label: node-role.kubernetes.io/masterlabel: node-role.kubernetes.io/infra

priority: 20

profile: openshift-control-plane

- priority: 30

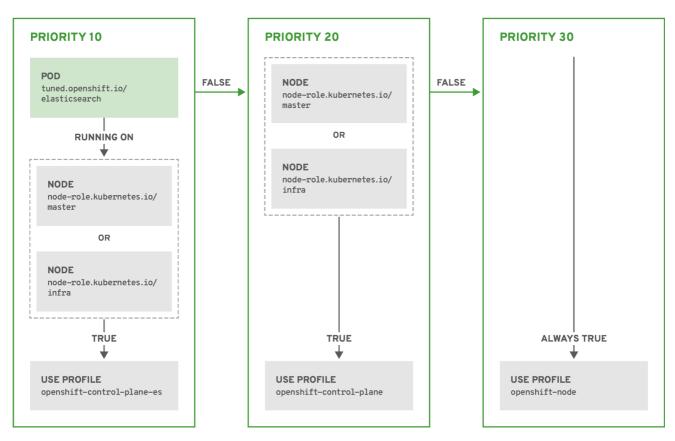
profile: openshift-node

上記のコンテナー化された Tuned デーモンの CR は、プロファイルの優先順位に基づいてその recommend.conf ファイルに変換されます。最も高い優先順位 (10) を持つプロファイルは openshift-control-plane-es であるため、これが最初に考慮されます。指定されたノードで実行されるコンテナー 化された Tuned デーモンは、同じノードに tuned.openshift.io/elasticsearch ラベルが設定された Pod が実行されているかどうかを確認します。これがない場合、 <match> セクション全体が false として 評価されます。このラベルを持つこのような Pod がある場合、 <match> セクションが true に評価されるようにするには、ノードラベルは node-role.kubernetes.io/master または node-role.kubernetes.io/infra である必要もあります。

優先順位が **10** のプロファイルのラベルが一致した場合、**openshift-control-plane-es** プロファイルが適用され、その他のプロファイルは考慮されません。ノード/Pod ラベルの組み合わせが一致しない場合、2 番目に高い優先順位プロファイル (**openshift-control-plane**) が考慮されます。このプロファイルは、コンテナー化されたチューニング済み Pod が **node-role.kubernetes.io/master** または **node-role.kubernetes.io/infra** ラベルを持つノードで実行される場合に適用されます。

最後に、プロファイル openshift-node には最低の優先順位である 30 が設定されます。これには <match> セクションがないため、常に一致します。これは、より高い優先順位の他のプロファイルが

指定されたノードで一致しない場合に **openshift-node** プロファイルを設定するために、最低の優先順位のノードが適用される汎用的な (catch-all) プロファイルとして機能します。



OPENSHIFT_10_0319

例: マシン設定プールベースのマッチング

apiVersion: tuned.openshift.io/v1

kind: Tuned metadata:

name: openshift-node-custom

namespace: openshift-cluster-node-tuning-operator

spec: profile:

- data: | [main]

summary=Custom OpenShift node profile with an additional kernel parameter

include=openshift-node

[bootloader]

cmdline openshift node custom=+skew tick=1

name: openshift-node-custom

recommend:

- machineConfigLabels:

machineconfiguration.openshift.io/role: "worker-custom"

priority: 20

profile: openshift-node-custom

ノードの再起動を最小限にするには、ターゲットノードにマシン設定プールのノードセレクターが一致するラベルを使用してラベルを付け、上記の Tuned CR を作成してから、最後にカスタムのマシン設定プール自体を作成します。

4.6. カスタムチューニングの例

以下の CR は、ラベル **tuned.openshift.io/ingress-node-label** を任意の値に設定した状態で OpenShift Container Platform ノードのカスタムノードレベルのチューニングを適用します。管理者として、以下のコマンドを使用してカスタムの Tune CR を作成します。

カスタムチューニングの例

```
$ oc create -f- << EOF
apiVersion: tuned.openshift.io/v1
kind: Tuned
metadata:
 name: ingress
 namespace: openshift-cluster-node-tuning-operator
spec:
 profile:
 - data: |
   [main]
   summary=A custom OpenShift ingress profile
   include=openshift-control-plane
   [sysctl]
   net.ipv4.ip local port range="1024 65535"
   net.ipv4.tcp tw reuse=1
  name: openshift-ingress
 recommend:
 - match:
  - label: tuned.openshift.io/ingress-node-label
  priority: 10
  profile: openshift-ingress
EOF_
```



重要

カスタムプロファイル作成者は、デフォルトの Tuned CR に含まれるデフォルトの調整されたデーモンプロファイルを組み込むことが強く推奨されます。上記の例では、デフォルトの openshift-control-plane プロファイルを使用してこれを実行します。

4.7. サポートされている TUNED デーモンプラグイン

[main] セクションを除き、以下の Tuned プラグインは、Tuned CR の profile: セクションで定義されたカスタムプロファイルを使用する場合にサポートされます。

- audio
- cpu
- disk
- eeepc_she
- modules
- mounts
- net

- scheduler
- scsi_host
- selinux
- sysctl
- sysfs
- usb
- video
- vm

これらのプラグインの一部によって提供される動的チューニング機能の中に、サポートされていない機能があります。以下の Tuned プラグインは現時点でサポートされていません。

- bootloader
- script
- systemd

詳細は、利用可能な Tuned プラグイン および Tuned の使用 を参照してください。

第5章 クラスターローダーの使用

クラスターローダーとは、クラスターに対してさまざまなオブジェクトを多数デプロイするツールであり、ユーザー定義のクラスターオブジェクトを作成します。クラスターローダーをビルド、設定、実行して、さまざまなクラスターの状態にある OpenShift Container Platform デプロイメントのパフォーマンスメトリクスを測定します。

5.1. クラスターローダーのインストール

手順

1. コンテナーイメージをプルするには、以下を実行します。

\$ podman pull quay.io/openshift/origin-tests:4.6

5.2. クラスターローダーの実行

前提条件

● リポジトリーは認証を要求するプロンプトを出します。レジストリーの認証情報を使用する と、一般的に利用できないイメージにアクセスできます。インストールからの既存の認証情報 を使用します。

手順

1. 組み込まれているテスト設定を使用してクラスターローダーを実行し、5つのテンプレートビルドをデプロイして、デプロイメントが完了するまで待ちます。

\$ podman run -v \${LOCAL_KUBECONFIG}:/root/.kube/config:z -i \ quay.io/openshift/origin-tests:4.6 /bin/bash -c 'export KUBECONFIG=/root/.kube/config && \ openshift-tests run-test "[sig-scalability][Feature:Performance] Load cluster \ should populate the cluster [Slow][Serial] [Suite:openshift]"'

または、VIPERCONFIG の環境変数を設定して、ユーザー定義の設定でクラスターローダーを 実行します。

 $\label{local_KUBECONFIG} $$ podman run -v $\{LOCAL_KUBECONFIG\}:/root/.kube/config:z \ \\ \\$

- -v \${LOCAL_CONFIG_FILE_PATH}:/root/configs/:z \
- -i quay.io/openshift/origin-tests:4.6 \

openshift-tests run-test "[sig-scalability][Feature:Performance] Load cluster \

should populate the cluster [Slow][Serial] [Suite:openshift]"

この例では、 \${LOCAL_KUBECONFIG} はローカルファイルシステムの kubeconfig のパスを参照します。さらに、 \${LOCAL_CONFIG_FILE_PATH} というディレクトリーがあり、これは test.yaml という設定ファイルが含まれるコンテナーにマウントされます。また、test.yaml が外部テンプレートファイルや podspec ファイルを参照する場合、これらもコンテナーにマウントされる必要があります。

5.3. クラスターローダーの設定

このツールは、複数のテンプレートや Pod を含む namespace (プロジェクト) を複数作成します。

5.3.1. クラスターローダー設定ファイルの例

クラスターローダーの設定ファイルは基本的な YAML ファイルです。

```
provider: local 1
ClusterLoader:
 cleanup: true
 projects:
  - num: 1
   basename: clusterloader-cakephp-mysql
   tuning: default
   ifexists: reuse
   templates:
    - num: 1
      file: cakephp-mysql.json
  - num: 1
   basename: clusterloader-dancer-mysql
   tuning: default
   ifexists: reuse
   templates:
    - num: 1
      file: dancer-mysql.json
  - num: 1
   basename: clusterloader-django-postgresql
   tuning: default
   ifexists: reuse
   templates:
    - num: 1
      file: django-postgresql.json
  - num: 1
   basename: clusterloader-nodejs-mongodb
   tuning: default
   ifexists: reuse
   templates:
    - num: 1
      file: quickstarts/nodejs-mongodb.json
  - num: 1
   basename: clusterloader-rails-postgresql
   tuning: default
   templates:
    - num: 1
      file: rails-postgresql.json
 tuningsets: 2
  - name: default
   pods:
    stepping: 3
      stepsize: 5
      pause: 0 s
     rate_limit: 4
      delay: 0 ms
```

- 1 エンドツーエンドテストのオプション設定。local に設定して、過剰に長いログメッセージを回避します。
- 2 このチューニングセットでは、速度制限やステップ設定、複数の Pod バッチ作成、セット間での 一時停止などが可能になります。クラスターローダーは、以前のステップが完了したことをモニ ターリングしてから、続行します。
- 👩 ステップ設定では、オブジェクトが N 個作成されてから、M 秒間一時停止します。
- 速度制限は、次のオブジェクトを作成するまで M ミリ秒間待機します。

この例では、外部テンプレートファイルや Pod 仕様ファイルへの参照もコンテナーにマウントされていることを前提とします。



重要

Microsoft Azure でクラスターローダーを実行している場合、AZURE_AUTH_LOCATION 変数を、インストーラーディレクトリーにある terraform.azure.auto.tfvars.json の出力が含まれるファイルに設定する必要があります。

5.3.2. 設定フィールド

表5.1クラスターローダーの最上位のフィールド

| フィールド | 説明 |
|------------|--|
| cleanup | true または false に設定します。設定ごとに1つの 定義を設定します。true に設定すると、cleanup は、テストの最後にクラスターローダーが作成した namespace (プロジェクト) すべてを削除します。 |
| projects | 1つまたは多数の定義が指定されたサブオブジェクト。 projects の下に、作成する各 namespace が定義され、 projects には必須のサブヘッダーが複数指定されます。 |
| tuningsets | 設定ごとに1つの定義が指定されたサブオブジェクト。 tuningsets では、チューニングセットを定義して、プロジェクトやオブジェクト作成に対して設定可能なタイミングを追加することができます(Pod、テンプレートなど)。 |
| sync | 設定ごとに1つの定義が指定されたオプションのサブ オブジェクト。オブジェクト作成時に同期できるか どうかについて追加します。 |

表5.2 projects の下にあるフィールド

| フィールド | 説明 |
|------------|--|
| num | 整数。作成するプロジェクト数の1つの定義。 |
| basename | 文字列。プロジェクトのベース名の定義。競合が発生しないように、同一の namespace の数が Basename に追加されます。 |
| tuning | 文字列。オブジェクトに適用するチューニングセットの1つの定義。 これは対象の namespace にデプロイします。 |
| ifexists | reuse または delete のいずれかが含まれる文字列。 ツールが実行時に作成するプロジェクトまたは namespace の名前と同じプロジェクトまたは namespace を見つける場合のツールの機能を定義し ます。 |
| configmaps | キーと値のペア一覧。キーは設定マップの名前で、 値はこの設定マップの作成元のファイルへのパスで す。 |
| secrets | キーと値のペア一覧。キーはシークレットの名前 で、値はこのシークレットの作成元のファイルへの パスです。 |
| pods | デプロイする Pod の1つまたは多数の定義を持つサ ブオブジェクト |
| templates | デプロイするテンプレートの1つまたは多数の定義を 持つサブオブジェクト |

表5.3 pods および templates のフィールド

| フィールド | 説明 |
|----------|---|
| num | 整数。デプロイする Pod またはテンプレート数。 |
| image | 文字列。プルが可能なリポジトリーに対する Docker イメージの URL |
| basename | 文字列。作成するテンプレート (または Pod) のベース名の1つの定義。 |
| file | 文字列。ローカルファイルへのパス。 作成する Pod 仕様またはテンプレートのいずれかです。 |

| フィールド | 説明 |
|------------|--|
| parameters | キーと値のペア。 parameters の下で、Pod または テンプレートでオーバーライドする値の一覧を指定 できます。 |

表5.4 tuningsets の下にあるフィールド

| フィールド | 説明 |
|-----------|---|
| name | 文字列。チューニングセットの名前。 プロジェクト のチューニングを定義する時に指定した名前と一致 します。 |
| pods | Pod に適用される tuningsets を特定するサブオブ ジェクト |
| templates | テンプレートに適用される tuningsets を特定する サブオブジェクト |

表5.5 tuningsets pods または tuningsets templates の下にあるフィールド

| フィールド | 説明 |
|------------|---|
| stepping | サブオブジェクト。ステップ作成パターンでオブ ジェクトを作成する場合に使用するステップ設定。 |
| rate_limit | サブオブジェクト。オブジェクト作成速度を制限す るための速度制限チューニングセットの設定。 |

表5.6 tuningsets pods または tuningsets templates、stepping の下にあるフィールド

| フィールド | 説明 |
|----------|--|
| stepsize | 整数。オブジェクト作成を一時停止するまでに作成 するオブジェクト数。 |
| pause | 整数。 stepsize で定義したオブジェクト数を作成後 に一時停止する秒数。 |
| timeout | 整数。オブジェクト作成に成功しなかった場合に失 敗するまで待機する秒数。 |
| delay | 整数。次の作成要求まで待機する時間 (ミリ秒)。 |

表5.7 sync の下にあるフィールド

| フィールド | 説明 |
|-----------|---|
| server | enabled および port フィールドを持つサブオブジェクト。ブール値 enabled を指定すると、Pod を同期するために HTTP サーバーを起動するかどうか定義します。port の整数はリッスンする HTTP サーバーポートを定義します (デフォルトでは 9090)。 |
| running | ブール値。 selectors に一致するラベルが指定された Pod が Running の状態になるまで待機します。 |
| succeeded | ブール値。 selectors に一致するラベルが指定された Pod が Completed の状態になるまで待機します。 |
| selectors | Running または Completed の状態の Pod に一致 するセレクター一覧 |
| timeout | 文字列。 Running または Completed の状態の Pod を待機してから同期をタイムアウトするまでの 時間。 0 以外の値は、単位 [ns us ms s m h] を使用してください。 |

5.4. 既知の問題

- クラスターローダーは設定なしで呼び出される場合に失敗します。(BZ#1761925)
- IDENTIFIER パラメーターがユーザーテンプレートで定義されていない場合には、テンプレートの作成は error: unknown parameter name "IDENTIFIER" エラーを出して失敗します。テンプレートをデプロイする場合は、このエラーが発生しないように、以下のパラメーターをテンプレートに追加してください。

```
{
  "name": "IDENTIFIER",
  "description": "Number to append to the name of resources",
  "value": "1"
}
```

Pod をデプロイする場合は、このパラメーターを追加する必要はありません。

第6章 CPU マネージャーの使用

CPU マネージャーは、CPU グループを管理して、ワークロードを特定の CPU に制限します。

CPU マネージャーは、以下のような属性が含まれるワークロードに有用です。

- できるだけ長い CPU 時間が必要な場合
- プロセッサーのキャッシュミスの影響を受ける場合
- レイテンシーが低いネットワークアプリケーションの場合
- 他のプロセスと連携し、単一のプロセッサーキャッシュを共有することに利点がある場合

6.1. CPU マネージャーの設定

手順

1. オプション: ノードにラベルを指定します。

oc label node perf-node.example.com cpumanager=true

2. CPU マネージャーを有効にする必要のあるノードの **MachineConfigPool** を編集します。この 例では、すべてのワーカーで CPU マネージャーが有効にされています。

oc edit machineconfigpool worker

3. ラベルをワーカーのマシン設定プールに追加します。

metadata: creationTimestamp: 2020-xx-xxx generation: 3 labels:

custom-kubelet: cpumanager-enabled

4. KubeletConfig、cpumanager-kubeletconfig.yaml、カスタムリソース (CR) を作成します。 直前の手順で作成したラベルを参照し、適切なノードを新規の kubelet 設定で更新しま す。machineConfigPoolSelector セクションを参照してください。

apiVersion: machineconfiguration.openshift.io/v1

kind: KubeletConfig

metadata:

name: cpumanager-enabled

machineConfigPoolSelector:

matchLabels:

custom-kubelet: cpumanager-enabled

kubeletConfig:

cpuManagerPolicy: static 1

cpuManagerReconcilePeriod: 5s 2

ポリシーを指定します。

- **none**このポリシーは、既存のデフォルト CPU アフィニティースキームを明示的に有効にし、スケジューラーが自動的に実行するもの以外のアフィニティーを提供しません。
- **static**このポリシーにより、特定のリソース特性を持つ Pod の CPU アフィニティー を増やし、これらの Pod のノードにおける排他性を付与することができます。
- オプション:CPU マネージャーの調整頻度を指定します。デフォルトは 5s です。
- 5. 動的な kubelet 設定を作成します。

oc create -f cpumanager-kubeletconfig.yaml

これにより、CPU マネージャー機能が kubelet 設定に追加され、必要な場合には Machine Config Operator (MCO) がノードを再起動します。CPU マネージャーを有効にするために再起動する必要はありません。

6. マージされた kubelet 設定を確認します。

oc get machineconfig 99-worker-XXXXXX-XXXXX-XXXXX-kubelet -o json | grep ownerReference -A7

出力例

```
"ownerReferences": [
{
          "apiVersion": "machineconfiguration.openshift.io/v1",
          "kind": "KubeletConfig",
          "name": "cpumanager-enabled",
          "uid": "7ed5616d-6b72-11e9-aae1-021e1ce18878"
        }
}
```

7. ワーカーで更新された kubelet.conf を確認します。

oc debug node/perf-node.example.com sh-4.2# cat /host/etc/kubernetes/kubelet.conf | grep cpuManager

出力例

cpuManagerPolicy: static 1 cpuManagerReconcilePeriod: 5s 2

- 1つこれらの設定は、KubeletConfig CR を作成する際に定義されたものです。
- 8. コア1つまたは複数を要求する Pod を作成します。制限および要求の CPU の値は整数にする 必要があります。これは、対象の Pod 専用のコア数です。

cat cpumanager-pod.yaml

出力例

```
apiVersion: v1
kind: Pod
metadata:
 generateName: cpumanager-
 containers:
 - name: cpumanager
  image: gcr.io/google_containers/pause-amd64:3.0
  resources:
   requests:
    cpu: 1
    memory: "1G"
   limits:
    cpu: 1
    memory: "1G"
 nodeSelector:
  cpumanager: "true"
```

9. Pod を作成します。

oc create -f cpumanager-pod.yaml

10. Pod がラベル指定されたノードにスケジュールされていることを確認します。

oc describe pod cpumanager

出力例

```
Name:
             cpumanager-6cqz7
                default
Namespace:
Priority:
            0
PriorityClassName: <none>
Node: perf-node.example.com/xxx.xx.xx.xxx
Limits:
   cpu:
        1
   memory: 1G
  Requests:
   cpu:
   memory:
             1G
QoS Class:
             Guaranteed
Node-Selectors: cpumanager=true
```

11. **cgroups** が正しく設定されていることを確認します。**pause** プロセスのプロセス ID (PID) を取得します。

QoS (quality of service) 階層 **Guaranteed** の Pod は、**kubepods.slice** に配置されます。他の QoS の Pod は、**kubepods** の子である **cgroups** に配置されます。

```
# cd /sys/fs/cgroup/cpuset/kubepods.slice/kubepods-
pod69c01f8e_6b74_11e9_ac0f_0a2b62178a22.slice/crio-
b5437308f1ad1a7db0574c542bdf08563b865c0345c86e9585f8c0b0a655612c.scope
# for i in `ls cpuset.cpus tasks`; do echo -n "$i "; cat $i ; done
```

出力例

cpuset.cpus 1 tasks 32706

12. 対象のタスクで許可される CPU 一覧を確認します。

grep ^Cpus_allowed_list /proc/32706/status

出力例

Cpus_allowed_list: 1

13. システム上の別の Pod (この場合は **burstable** QoS 階層にある Pod) が、**Guaranteed** Pod に割り当てられたコアで実行できないことを確認します。

cat /sys/fs/cgroup/cpuset/kubepods.slice/kubepods-besteffort.slice/kubepods-besteffort-podc494a073_6b77_11e9_98c0_06bba5c387ea.slice/crio-c56982f57b75a2420947f0afc6cafe7534c5734efc34157525fa9abbf99e3849.scope/cpuset.cpus

0

oc describe node perf-node.example.com

出力例

```
Capacity:
attachable-volumes-aws-ebs: 39
cpu:
ephemeral-storage:
                    124768236Ki
hugepages-1Gi:
                     0
hugepages-2Mi:
                    0
                 8162900Ki
memory:
pods:
                 250
Allocatable:
attachable-volumes-aws-ebs: 39
cpu:
                1500m
ephemeral-storage:
                    124768236Ki
hugepages-1Gi:
                     0
hugepages-2Mi:
                     0
                7548500Ki
memory:
pods:
                 250
 default
                       cpumanager-6cqz7 1 (66%)
                                                      1 (66%) 1G (12%)
```

Allocated resources:
(Total limits may be over 100 percent, i.e., overcommitted.)
Resource Requests Limits
-----cpu 1440m (96%) 1 (66%)

この仮想マシンには、2 つの CPU コアがあります。**system-reserved** 設定は 500 ミリコアを予約し、**Node Allocatable** の量になるようにノードの全容量からコアの半分を引きます。ここで **Allocatable CPU** は 1500 ミリコアであることを確認できます。これは、それぞれがコアを 1 つ受け入れるので、CPU マネージャー Pod の 1 つを実行できることを意味します。1 つのコア全体は 1000 ミリコアに相当します。2 つ目の Pod をスケジュールしようとする場合、システムは Pod を受け入れますが、これがスケジュールされることはありません。

NAME READY STATUS RESTARTS AGE cpumanager-6cqz7 1/1 Running 0 33m cpumanager-7qc2t 0/1 Pending 0 11s

第7章 TOPOLOGY MANAGER の使用

Topology Manager は、CPU マネージャー、デバイスマネージャー、およびその他の Hint Provider からヒントを収集し、同じ Non-Uniform Memory Access (NUMA) ノード上のすべての QoS (Quality of Service) クラスについて CPU、SR-IOV VF、その他デバイスリソースなどの Pod リソースを調整します。

Topology Manager は、収集したヒントのトポロジー情報を使用し、設定される Topology Manager ポリシーおよび要求される Pod リソースに基づいて、pod がノードから許可されるか、または拒否されるかどうかを判別します。

Topology Manager は、ハードウェアアクセラレーターを使用して低遅延 (latency-critical) の実行と高スループットの並列計算をサポートするワークロードの場合に役立ちます。



注記

Topology Manager を使用するには、**static** ポリシーで CPU マネージャーを使用する必要があります。CPU マネージャーの詳細は、CPU マネージャーの使用 を参照してください。

7.1. TOPOLOGY MANAGER ポリシー

Topology Manager は、CPU マネージャーやデバイスマネージャーなどの Hint Provider からトポロ ジーのヒントを収集し、収集したヒントを使用して **Pod** リソースを調整することで、すべての QoS (Quality of Service) クラスの **Pod** リソースを調整します。



注記

CPU リソースを **Pod** 仕様の他の要求されたリソースと調整するには、CPU マネージャーを **static** CPU マネージャーポリシーで有効にする必要があります。

Topology Manager は、**cpumanager-enabled** カスタムリソース (CR) で割り当てる 4 つの割り当てポリシーをサポートします。

none ポリシー

これはデフォルトのポリシーで、トポロジーの配置は実行しません。

best-effort ポリシー

best-effortトポロジー管理ポリシーを持つ Pod のそれぞれのコンテナーの場合、kubelet は 各 Hint Provider を呼び出してそれらのリソースの可用性を検出します。この情報を使用して、Topology Manager は、そのコンテナーの推奨される NUMA ノードのアフィニティーを保存します。アフィニティーが優先されない場合、Topology Manager はこれを保管し、ノードに対して Pod を許可します。

restricted ポリシー

restricted トポロジー管理ポリシーを持つ Pod のそれぞれのコンテナーの場合、kubelet は 各 Hint Provider を呼び出してそれらのリソースの可用性を検出します。この情報を使用して、Topology Manager は、そのコンテナーの推奨される NUMA ノードのアフィニティーを保存します。アフィニティーが優先されない場合、Topology Manager はこの Pod をノードから拒否します。これにより、Pod が Pod の受付の失敗により **Terminated** 状態になります。

single-numa-node ポリシー

single-numa-node トポロジー管理ポリシーがある Pod のそれぞれのコンテナーの場合、kubelet は 各 Hint Provider を呼び出してそれらのリソースの可用性を検出します。この情報を使用して、 Topology Manager は単一の NUMA ノードのアフィニティーが可能かどうかを判別します。可能で

ある場合、Pod はノードに許可されます。単一の NUMA ノードアフィニティーが使用できない場合には、Topology Manager は Pod をノードから拒否します。これにより、Pod は Pod の受付失敗と共に Terminated (終了) 状態になります。

7.2. TOPOLOGY MANAGER のセットアップ

Topology Manager を使用するには、 **cpumanager-enabled** カスタムリソース (CR) で割り当てポリシーを設定する必要があります。 CPU マネージャーをセットアップしている場合は、このファイルが存在している可能性があります。ファイルが存在しない場合は、作成できます。

前提条件

● CPU マネージャーのポリシーを **static** に設定します。スケーラビリティーおよびパフォーマンスセクションの CPU マネージャーの使用を参照してください。

手順

Topololgy Manager をアクティブにするには、以下を実行します。

1. **cpumanager-enabled** カスタムリソース (CR) で Topology Manager 割り当てポリシーを設定します。

\$ oc edit KubeletConfig cpumanager-enabled

apiVersion: machineconfiguration.openshift.io/v1

kind: KubeletConfig

metadata:

name: cpumanager-enabled

spec:

machineConfigPoolSelector:

matchLabels:

custom-kubelet: cpumanager-enabled

kubeletConfig:

cpuManagerPolicy: static 1 cpuManagerReconcilePeriod: 5s

topologyManagerPolicy: single-numa-node 2

- 🚹 このパラメーターは static である必要があります。
- 選択した Topology Manager 割り当てポリシーを指定します。このポリシーは single-numa-node になります。使用できる値は、default、best-effort、restricted、single-numa-node です。

関連情報

● CPU マネージャーの詳細は、CPU マネージャーの使用 を参照してください。

7.3. POD の TOPOLOGY MANAGER ポリシーとの対話

以下のサンプル **Pod** 仕様は、Pod の Topology Manger との対話について説明しています。

以下の Pod は、リソース要求や制限が指定されていないために BestEffort QoS クラスで実行されます。

spec:

containers:
- name: nginx
image: nginx

以下の Pod は、要求が制限よりも小さいために Burstable QoS クラスで実行されます。

spec: containers: - name: nginx image: nginx resources: limits: memory: "200Mi" requests: memory: "100Mi"

選択したポリシーが **none** 以外の場合は、Topology Manager はこれらの **Pod** 仕様のいずれかも考慮しません。

以下の最後のサンプル Pod は、要求が制限と等しいために Guaranteed QoS クラスで実行されます。

spec: containers: - name: nginx image: nginx resources: limits: memory: "200Mi" cpu: "2" example.com/device: "1" requests: memory: "200Mi" cpu: "2" example.com/device: "1"

Topology Manager はこの Pod を考慮します。Topology Manager は、利用可能な CPU のトポロジー を返す CPU マネージャーの静的ポリシーを確認します。また Topology Manager はデバイスマネー ジャーを確認し、example.com/device の利用可能なデバイスのトポロジーを検出します。

Topology Manager はこの情報を使用して、このコンテナーに最適なトポロジーを保管します。この Pod の場合、CPU マネージャーおよびデバイスマネージャーは、リソース割り当ての段階でこの保存 された情報を使用します。

第8章 CLUSTER MONITORING OPERATOR のスケーリング

OpenShift Container Platform は、Cluster Monitoring Operator が収集し、Prometheus ベースのモニターリングスタックに保存するメトリクスを公開します。管理者は、Grafana という1つのダッシュボードインターフェイスでシステムリソース、コンテナーおよびコンポーネントのメトリクスを表示できます。

8.1. PROMETHEUS データベースのストレージ要件

Red Hat では、異なるスケールサイズに応じて各種のテストが実行されました。



注記

以下の Prometheus ストレージ要件は規定されていません。ワークロードのアクティビティーおよびリソースの使用に応じて、クラスターで観察されるリソースの消費量が大きくなる可能性があります。

表8.1 クラスター内のノード/Pod の数に基づく Prometheus データベースのストレージ要件

| ノード数 | Pod 数 | 1日あたりの Prometheus ス トレージの増 加量 | 15 日ごとの Prometheus ス トレージの増 加量 | RAM 領域 (ス ケールサイズ に基づく) | ネットワーク (tsdb チャンク に基づく) |
|------|-------|--|---|------------------------------|-------------------------------|
| 50 | 1800 | 6.3 GB | 94 GB | 6 GB | 16 MB |
| 100 | 3600 | 13 GB | 195 GB | 10 GB | 26 MB |
| 150 | 5400 | 19 GB | 283 GB | 12 GB | 36 MB |
| 200 | 7200 | 25 GB | 375 GB | 14 GB | 46 MB |

ストレージ要件が計算値を超過しないようにするために、オーバーヘッドとして予期されたサイズのおよそ 20% が追加されています。

上記の計算は、デフォルトの OpenShift Container Platform Cluster Monitoring Operator についての計算です。



注記

CPU の使用率による影響は大きくありません。比率については、およそ 50 ノードおよび 1800 Pod ごとに 1 コア (/40) になります。

OpenShift Container Platform についての推奨事項

- 3つ以上のインフラストラクチャー (infra) ノードを使用します。
- NVMe (non-volatile memory express) ドライブを搭載した 3 つ以上の **openshift-container-storage** ノードを使用します。

8.2. クラスターモニターリングの設定

クラスターモニターリングスタック内の Prometheus コンポーネントのストレージ容量を増やすことができます。

手順

Prometheus のストレージ容量を拡張するには、以下を実行します。

1. YAML 設定ファイル cluster-monitoring-config.yml を作成します。以下に例を示します。

```
apiVersion: v1
kind: ConfigMap
data:
 config.yaml: |
  prometheusK8s:
   retention: {{PROMETHEUS RETENTION PERIOD}}
   nodeSelector:
    node-role.kubernetes.io/infra: ""
   volumeClaimTemplate:
    spec:
     storageClassName: {{STORAGE CLASS}} 2
     resources:
      requests:
       storage: {{PROMETHEUS_STORAGE_SIZE}} 3
  alertmanagerMain:
   nodeSelector:
    node-role.kubernetes.io/infra: ""
   volumeClaimTemplate:
     storageClassName: {{STORAGE CLASS}} 4
     resources:
      requests:
       storage: {{ALERTMANAGER STORAGE SIZE}} 5
metadata:
 name: cluster-monitoring-config
 namespace: openshift-monitoring
```

- 標準の値は PROMETHEUS_RETENTION_PERIOD=15d になります。時間は、接尾辞 s、m、h、d のいずれかを使用する単位で測定されます。
- 24 クラスターのストレージクラス。
- 標準の値は PROMETHEUS_STORAGE_SIZE=2000Gi です。ストレージの値には、接尾 辞 E、P、T、G、M、K のいずれかを使用した単純な整数または固定小数点整数を使用できます。 また、2 のべき乗の値 (Ei、Pi、Ti、Gi、Mi、Ki) を使用することもできます。
- 標準の値は **ALERTMANAGER_STORAGE_SIZE=20Gi** です。ストレージの値には、接尾 辞 E、P、T、G、M、K のいずれかを使用した単純な整数または固定小数点整数を使用できます。 また、2 のべき乗の値 (Ei、Pi、Ti、Gi、Mi、Ki) を使用することもできます。
- 2. 保存期間、ストレージクラス、およびストレージサイズの値を追加します。
- 3. ファイルを保存します。

4. 以下を実行して変更を適用します。

\$ oc create -f cluster-monitoring-config.yaml

第9章 オブジェクトの最大値に合わせた環境計画

OpenShift Container Platform クラスターの計画時に以下のテスト済みのオブジェクトの最大値を考慮します。

これらのガイドラインは、最大規模のクラスターに基づいています。小規模なクラスターの場合、最大値はこれより低くなります。指定のしきい値に影響を与える要因には、etcd バージョンやストレージデータ形式などの多数の要因があります。



重要

これらのガイドラインは、Open Virtual Network (OVN) ではなく、ソフトウェア定義ネットワーク (SDN) を使用する OpenShift Container Platform に該当します。

ほとんど場合、これらの制限値を超えると、パフォーマンスが全体的に低下します。ただし、これによって必ずしもクラスターに障害が発生する訳ではありません。

9.1. メジャーリリースについての **OPENSHIFT CONTAINER PLATFORM** のテスト済みクラスターの最大値

OpenShift Container Platform 3.x のテスト済みクラウドプラットフォーム: Red Hat OpenStack (RHOSP)、Amazon Web Services および Microsoft AzureOpenShift Container Platform 4.x のテスト済み Cloud Platform : Amazon Web Services、Microsoft Azure および Google Cloud Platform

| 最大値のタイプ | 3.x テスト済みの最大値 | 4.x テスト済みの最大値 |
|--|--|--|
| ノード数 | 2,000 | 2,000 |
| Pod 数 ^[1] | 150,000 | 150,000 |
| ノードあたりの Pod 数 | 250 | 500 [2] |
| コアあたりの Pod 数 | デフォルト値はありません。 | デフォルト値はありません。 |
| namespace 数 ^[3] | 10,000 | 10,000 |
| ビルド数 | 10,000(デフォルト Pod RAM 512 Mi)- Pipeline ストラテジー | 10,000(デフォルト Pod RAM 512 Mi)- Source-to-Image (S2I) ビル ドストラテジー |
| namespace ごとの Pod 数 ^[4] | 25,000 | 25,000 |
| Ingress Controller ごとのルートと バックエンドの数 | ルーターあたり 2,000 | ルーターあたり 2,000 |
| シークレットの数 | 80,000 | 80,000 |
| config map の数 | 90,000 | 90,000 |

| 最大値のタイプ | 3.x テスト済みの最大値 | 4.x テスト済みの最大値 |
|--|---------------|--------------------|
| サービス数 ^[5] | 10,000 | 10,000 |
| namespace ごとのサービス数 | 5,000 | 5,000 |
| サービスごとのバックエンド数 | 5,000 | 5,000 |
| namespace ごとのデプロイメント 数 ^[4] | 2,000 | 2,000 |
| ビルド設定の数 | 12,000 | 12,000 |
| シークレットの数 | 40,000 | 40,000 |
| カスタムリソース定義 (CRD) の数 | デフォルト値はありません。 | 512 ^[6] |

- 1. ここで表示される Pod 数はテスト用の Pod 数です。実際の Pod 数は、アプリケーションのメモリー、CPU、ストレージ要件により異なります。
- 2. これは、ワーカーノードごとに 500 の Pod を持つ 100 ワーカーノードを含むクラスターでテストされています。デフォルトの maxPods は 250 です。500 maxPods に到達するには、クラスターはカスタム kubelet 設定を使用し、maxPods が 500 に設定された状態で作成される必要があります。500 ユーザー Pod が必要な場合は、ノード上に 10-15 のシステム Pod がすでに実行されているため、hostPrefix が 22 である必要があります。永続ボリューム要求(PVC)が割り当てられている Pod の最大数は、PVC の割り当て元のストレージバックエンドによって異なります。このテストでは、OpenShift Container Storage v4 (OCS v4) のみが本書で説明されているノードごとの Pod 数に対応することができました。
- 3. 有効なプロジェクトが多数ある場合、キースペースが過剰に拡大し、スペースのクォータを超過すると、etcd はパフォーマンスの低下による影響を受ける可能性があります。etcd ストレージを解放するために、デフラグを含む etcd の定期的なメンテナーンスを行うことを強くお勧めします。
- 4. システムには、状態の変更に対する対応として特定の namespace にある全オブジェクトに対して反復する多数のコントロールループがあります。単一の namespace に特定タイプのオブジェクトの数が多くなると、ループのコストが上昇し、特定の状態変更を処理する速度が低下します。この制限については、アプリケーションの各種要件を満たすのに十分な CPU、メモリー、およびディスクがシステムにあることが前提となっています。
- 5. 各サービスポートと各サービスのバックエンドには、iptables の対応するエントリーがあります。特定のサービスのバックエンド数は、エンドポイントのオブジェクトサイズに影響があり、その結果、システム全体に送信されるデータサイズにも影響を与えます。
- 6. OpenShift Container Platform には、OpenShift Container Platform によってインストールされたもの、OpenShift Container Platform と統合された製品、およびユーザー作成の CRD を含め、合計 512 のカスタムリソース定義 (CRD) の制限があります。512 を超える CRD が作成されている場合は、oc コマンドリクエストのスロットリングが適用される可能性があります。



注記

Red Hat は、OpenShift Container Platform クラスターのサイズ設定に関する直接的なガイダンスを提供していません。これは、クラスターが OpenShift Container Platform のサポート範囲内にあるかどうかを判断するには、クラスターのスケールを制限するすべての多次元な要因を慎重に検討する必要があるためです。

9.2. クラスターの最大値がテスト済みの OPENSHIFT CONTAINER PLATFORM 環境および設定

AWS クラウドプラットフォーム:

| ノード | フレー バー | vCPU | RAM(GiB) | ディスク タイプ | ディスク サイズ (GiB)/IO S | カウント | リージョ ン |
|---------------------------|-----------------|------|----------|-------------|------------------------------|---------------------------------|-----------|
| マス ター/etcd [1] | r5.4xlarge | 16 | 128 | io1 | 220 / 3000 | 3 | us-west-2 |
| インフラ [2] | m5.12xlarg e | 48 | 192 | gp2 | 100 | 3 | us-west-2 |
| ワーク ロード ^[3] | m5.4xlarg e | 16 | 64 | gp2 | 500 [4] | 1 | us-west-2 |
| ワーカー | m5.2xlarg e | 8 | 32 | gp2 | 100 | 3/25/250 /500 ^[5] | us-west-2 |

- 1. 3000 IOPS を持つ io1 ディスクは、etcd が I/O 集約型であり、かつレイテンシーの影響を受けやすいため、マスター/etcd ノードに使用されます。
- 2. インフラストラクチャーノードは、モニターリング、Ingress およびレジストリーコンポーネントをホストするために使用され、これにより、それらが大規模に実行する場合に必要とするリソースを十分に確保することができます。
- 3. ワークロードノードは、パフォーマンスとスケーラビリティーのワークロードジェネレーター を実行するための専用ノードです。
- 4. パフォーマンスおよびスケーラビリティーのテストの実行中に収集される大容量のデータを保存するのに十分な領域を確保できるように、大きなディスクサイズが使用されます。
- 5. クラスターは反復的にスケーリングされ、パフォーマンスおよびスケーラビリティーテストは 指定されたノード数で実行されます。

9.3. テスト済みのクラスターの最大値に基づく環境計画



重要

ノード上で物理リソースを過剰にサブスクライブすると、Kubernetes スケジューラーが Pod の配置時に行うリソースの保証に影響が及びます。メモリースワップを防ぐために 実行できる処置について確認してください。

一部のテスト済みの最大値については、単一の namespace/ユーザーが作成するオブジェクトでのみ変更されます。これらの制限はクラスター上で数多くのオブジェクトが実行されている場合には異なります。

本書に記載されている数は、Red Hat のテスト方法、セットアップ、設定、および チューニングに基づいています。これらの数は、独自のセットアップおよび環境に応じ て異なります。

環境の計画時に、ノードに配置できる Pod 数を判別します。

required pods per cluster / pods per node = total number of nodes needed

ノードあたりの現在の Pod の最大数は 250 です。ただし、ノードに適合する Pod 数はアプリケーション自体によって異なります。アプリケーション要件に合わせて環境計画を立てる方法で説明されているように、アプリケーションのメモリー、CPU およびストレージの要件を検討してください。

シナリオ例

クラスターごとに 2200 の Pod のあるクラスターのスコープを設定する場合、ノードごとに最大 500 の Pod があることを前提として、最低でも 5 つのノードが必要になります。

2200 / 500 = 4.4

ノード数を 20 に増やす場合は、Pod 配分がノードごとに 110 の Pod に変わります。

2200 / 20 = 110

ここで、

required pods per cluster / total number of nodes = expected pods per node

9.4. アプリケーション要件に合わせて環境計画を立てる方法

アプリケーション環境の例を考えてみましょう。

| Pod タイプ | Pod 数 | 最大メモリー | CPUコア数 | 永続ストレージ |
|------------|-------|--------|--------|---------|
| apache | 100 | 500 MB | 0.5 | 1GB |
| node.js | 200 | 1GB | 1 | 1GB |
| postgresql | 100 | 1GB | 2 | 10 GB |
| JBoss EAP | 100 | 1GB | 1 | 1GB |

推定要件: CPU コア 550 個、メモリー 450GB およびストレージ 1.4TB

ノードのインスタンスサイズは、希望に応じて増減を調整できます。ノードのリソースはオーバーコミットされることが多く、デプロイメントシナリオでは、小さいノードで数を増やしたり、大きいノードで数を減らしたりして、同じリソース量を提供することもできます。このデプロイメントシナリオでは、小さいノードで数を増やしたり、大きいノードで数を減らしたりして、同じリソース量を提供することもできます。運用上の敏捷性やインスタンスあたりのコストなどの要因を考慮する必要があります。

| ノードのタイプ | 数量 | CPU | RAM (GB) |
|---------------|-----|-----|----------|
| ノード (オプション 1) | 100 | 4 | 16 |
| ノード (オプション 2) | 50 | 8 | 32 |
| ノード (オプション 3) | 25 | 16 | 64 |

アプリケーションによってはオーバーコミットの環境に適しているものもあれば、そうでないものもあります。たとえば、Java アプリケーションや Huge Page を使用するアプリケーションの多くは、オーバーコミットに対応できません。対象のメモリーは、他のアプリケーションに使用できません。上記の例では、環境は一般的な比率として約30%オーバーコミットされています。

アプリケーション Pod は環境変数または DNS のいずれかを使用してサービスにアクセスできます。環境変数を使用する場合、それぞれのアクティブなサービスについて、変数が Pod がノードで実行される際に kubelet によって挿入されます。クラスター対応の DNS サーバーは、Kubernetes API で新規サービスの有無を監視し、それぞれに DNS レコードのセットを作成します。 DNS がクラスター全体で有効にされている場合、すべての Pod は DNS 名でサービスを自動的に解決できるはずです。 DNS を使用したサービス検出は、5000 サービスを超える使用できる場合があります。サービス検出に環境変数を使用する場合、引数の一覧は namespace で 5000 サービスを超える場合の許可される長さを超えると、Pod およびデプロイメントは失敗します。デプロイメントのサービス仕様ファイルのサービスリンクを無効にして、以下を解消します。

```
apiVersion: v1
kind: Template
metadata:
 name: deployment-config-template
 creationTimestamp:
 annotations:
  description: This template will create a deploymentConfig with 1 replica, 4 env vars and a service.
  tags: "
objects:
- apiVersion: v1
 kind: DeploymentConfig
 metadata:
  name: deploymentconfig${IDENTIFIER}
 spec:
  template:
   metadata:
    labels:
      name: replicationcontroller${IDENTIFIER}
```

enableServiceLinks: false

```
containers:
    - name: pause${IDENTIFIER}
      image: "${IMAGE}"
      ports:
      - containerPort: 8080
       protocol: TCP
      env:
      - name: ENVVAR1_${IDENTIFIER}
       value: "${ENV VALUE}"
      - name: ENVVAR2 ${IDENTIFIER}
       value: "${ENV_VALUE}"
      - name: ENVVAR3_${IDENTIFIER}
       value: "${ENV_VALUE}"
      - name: ENVVAR4_${IDENTIFIER}
       value: "${ENV_VALUE}"
      resources: {}
      imagePullPolicy: IfNotPresent
      capabilities: {}
      securityContext:
       capabilities: {}
       privileged: false
    restartPolicy: Always
    serviceAccount: "
  replicas: 1
  selector:
   name: replicationcontroller${IDENTIFIER}
  triggers:
  - type: ConfigChange
  strategy:
   type: Rolling
- apiVersion: v1
 kind: Service
 metadata:
  name: service${IDENTIFIER}
 spec:
  selector:
   name: replicationcontroller${IDENTIFIER}
  - name: serviceport${IDENTIFIER}
   protocol: TCP
   port: 80
   targetPort: 8080
  portallP: "
  type: ClusterIP
  sessionAffinity: None
 status:
  loadBalancer: {}
parameters:
- name: IDENTIFIER
 description: Number to append to the name of resources
 value: '1'
 required: true
- name: IMAGE
 description: Image to use for deploymentConfig
 value: gcr.io/google-containers/pause-amd64:3.0
 required: false
```

- name: ENV_VALUE

description: Value to use for environment variables

generate: expression from: "[A-Za-z0-9]{255}"

required: false

labels:

template: deployment-config-template

namespace で実行できるアプリケーション Pod の数は、環境変数がサービス検出に使用される場合にサービスの数およびサービス名の長さによって異なります。システムの ARG_MAX は、新規プロセスの引数の最大の長さを定義し、デフォルトで 2097152 KiB に設定されます。Kubelet は、以下を含むnamespace で実行するようにスケジュールされる各 Pod に環境変数を挿入します。

- <SERVICE_NAME>_SERVICE_HOST=<IP>
- <SERVICE_NAME>_SERVICE_PORT=<PORT>
- <SERVICE_NAME>_PORT=tcp://<IP>:<PORT>
- <SERVICE_NAME>_PORT_<PORT>_TCP=tcp://<IP>:<PORT>
- <SERVICE_NAME>_PORT_<PORT>_TCP_PROTO=tcp
- <SERVICE_NAME>_PORT_<PORT>_TCP_PORT=<PORT>
- <SERVICE NAME> PORT <PORT> TCP ADDR=<ADDR>

引数の長さが許可される値を超え、サービス名の文字数がこれに影響する場合、namespace の Pod は起動に失敗し始めます。たとえば、5000 サービスを含む namespace では、サービス名の制限は 33 文字であり、これにより namespace で 5000 Pod を実行できます。

第10章 ストレージの最適化

ストレージを最適化すると、すべてのリソースでストレージの使用を最小限に抑えることができます。 管理者は、ストレージを最適化することで、既存のストレージリソースが効率的に機能できるようにす ることができます。

10.1. 利用可能な永続ストレージオプション

永続ストレージオプションについて理解し、OpenShift Container Platform 環境を最適化できるようにします。

表10.1利用可能なストレージオプション

| ストレー ジタイプ | 説明 | 例 |
|--------------|--|--|
| ブロック | ブロックデバイスとしてオペレーティングシステムに公開されます。 ストレージを完全に制御し、ファイルシステムを通過してファイルの低いレベルで操作する必要のあるアプリケーションに適しています。 ストレージエリアネットワーク (SAN) とも呼ばれます。 共有できません。一度に1つのクライアントだけがこのタイプのエンドポイントをマウントできるという意味です。 | AWS EBS および VMware vSphere は、OpenShift Container Platform で 永続ボリューム (PV) の動的なプロビ ジョニングをサポートします。 |
| ファイル | マウントされるファイルシステムのエクスポートとして、OS に公開されます。 ネットワークアタッチストレージ (NAS) とも呼ばれます。 同時実行、レイテンシー、ファイルロックのメカニズムその他の各種機能は、プロトコルおよび実装、ベンダー、スケールによって大きく異なります。 | RHEL NFS、NetApp NFS ^[1] 、および Vendor NFS |
| オブジェクト | REST API エンドポイント経由でアクセスできます。 OpenShift Container Platform レジストリーで使用するために設定できます。 アプリケーションは、ドライバーをアプリケーションやコンテナーに組み込む必要があります。 | AWS S3 |

1. NetApp NFS は Trident プラグインを使用する場合に動的 PV のプロビジョニングをサポートします。



重要

現時点で、CNS は OpenShift Container Platform 4.6 ではサポートされていません。

10.2. 設定可能な推奨のストレージ技術

以下の表では、特定の OpenShift Container Platform クラスターアプリケーション向けに設定可能な推奨のストレージ技術についてまとめています。

表10.2 設定可能な推奨ストレージ技術

| ストレー ジタイプ | ROX1 | RWX2 | レジスト リー | スケーリ ングされ たレジス トリー | メトリク ス3 | ロギング | アプリ |
|--------------|-----------------|------|------------|-----------------------------|-------------------|-------------------|-------------------|
| ブロック | はい ⁴ | いいえ | 設定可能 | 設定不可 | 推奨 | 推奨 | 推奨 |
| ファイル | はい4 | 0 | 設定可能 | 設定可能 | 設定可能 ⁵ | 設定可能 ⁶ | 推奨 |
| オブジェ クト | 0 | 0 | 推奨 | 推奨 | 設定不可 | 設定不可 | 設定不可 ⁷ |

¹ReadOnlyMany

⁴ これは、物理ディスク、VM 物理ディスク、VMDK、NFS 経由のループバック、AWS EBS、および Azure Disk には該当しません。

⁵メトリクスの場合、**ReadWriteMany** (RWX) アクセスモードのファイルストレージを信頼できる方法で使用することはできません。ファイルストレージを使用する場合、メトリクスと共に使用されるように設定される永続ボリューム要求 (PVC) で RWX アクセスモードを設定しないでください。

⁶ ロギングの場合、共有ストレージを使用することはアンチパターンとなります。 elasticsearch ごとに 1つのボリュームが必要です。

⁷ オブジェクトストレージは、OpenShift Container Platform の PV/PVC で消費されません。アプリは、オブジェクトストレージの REST API と統合する必要があります。



注記

スケーリングされたレジストリーとは、2 つ以上の Pod レプリカが稼働する OpenShift Container Platform レジストリーのことです。

² ReadWriteMany

³ Prometheus はメトリクスに使用される基礎となるテクノロジーです。

10.2.1. 特定アプリケーションのストレージの推奨事項



重要

テストにより、NFS サーバーを Red Hat Enterprise Linux (RHEL) でコアサービスのストレージバックエンドとして使用することに関する問題が検出されています。これには、OpenShift Container レジストリーおよび Quay、メトリクスストレージのPrometheus、およびロギングストレージの Elasticsearch が含まれます。そのため、コアサービスで使用される PV をサポートするために RHEL NFS を使用することは推奨されていません。

他の NFS の実装ではこれらの問題が検出されない可能性があります。OpenShift Container Platform コアコンポーネントに対して実施された可能性のあるテストに関する詳細情報は、個別の NFS 実装ベンダーにお問い合わせください。

10.2.1.1. レジストリー

スケーリングなし/高可用性 (HA) ではない OpenShift Container Platform レジストリークラスターのデプロイメント:

- ストレージ技術は、RWX アクセスモードをサポートする必要はありません。
- ストレージ技術は、リードアフターライト (Read-After-Write) の一貫性を確保する必要があります。
- 推奨されるストレージ技術はオブジェクトストレージであり、次はブロックストレージです。
- ファイルストレージは、実稼働環境のワークロードを処理する OpenShift Container Platform レジストリークラスターのデプロイメントには推奨されません。

10.2.1.2. スケーリングされたレジストリー

スケーリングされた/高可用性 (HA) の OpenShift Container Platform レジストリーのクラスターデプロイメント:

- ストレージ技術は、RWX アクセスモードをサポートする必要があります。
- ストレージ技術は、リードアフターライト (Read-After-Write) の一貫性を確保する必要があります。
- Amazon Simple Storage Service (Amazon S3)、Google Cloud Storage (GCS)、Microsoft Azure Blob Storage、および OpenStack Swift はサポートされます。
- オブジェクトストレージは S3 または Swift に準拠する必要があります。
- vSphere やベアメタルインストールなどのクラウド以外のプラットフォームの場合、設定可能な技術はファイルストレージのみです。
- ブロックストレージは設定できません。

10.2.1.3. メトリクス

OpenShift Container Platform がホストするメトリクスのクラスターデプロイメント:

- 推奨されるストレージ技術はブロックストレージです。
- オブジェクトストレージは設定できません。



重要

実稼働ワークロードがあるホスト型のメトリクスクラスターデプロイメントにファイルストレージを使用することは推奨されません。

10.2.1.4. ロギング

OpenShift Container Platform がホストするロギングのクラスターデプロイメント:

- 推奨されるストレージ技術はブロックストレージです。
- オブジェクトストレージは設定できません。

10.2.1.5. アプリケーション

以下の例で説明されているように、アプリケーションのユースケースはアプリケーションごとに異なります。

- 動的な PV プロビジョニングをサポートするストレージ技術は、マウント時のレイテンシーが低く、ノードに関連付けられておらず、正常なクラスターをサポートします。
- アプリケーション開発者はアプリケーションのストレージ要件や、それがどのように提供されているストレージと共に機能するかを理解し、アプリケーションのスケーリング時やストレージレイヤーと対話する際に問題が発生しないようにしておく必要があります。

10.2.2. 特定のアプリケーションおよびストレージの他の推奨事項



重要

etcd などの Write 集中型ワークロードで RAID 設定を使用することはお勧めしません。 RAID 設定で etcd を実行している場合、ワークロードでパフォーマンスの問題が発生するリスクがある可能性があります。

- Red Hat OpenStack Platform (RHOSP) Cinder: RHOSP Cinder は ROX アクセスモードのユースケースで適切に機能する傾向があります。
- データベース: データベース (RDBMS、NoSQL DB など) は、専用のブロックストレージで最適に機能することが予想されます。
- etcd データベースには、大規模なクラスターを有効にするのに十分なストレージと十分なパフォーマンス容量が必要です。十分なストレージと高性能環境を確立するための監視およびベンチマークツールに関する情報は、**推奨される etcd プラクティス** に記載されています。

10.3. データストレージ管理

以下の表は、OpenShift Container Platform コンポーネントがデータを書き込むメインディレクトリーの概要を示しています。

表10.3 OpenShift Container Platform データを保存するメインディレクトリー

| ディレクトリー | 注記 | サイジング | 予想される拡張 |
|---------------------|--|--|---|
| /var/lib/etcd | データベースを保存する 際に etcd ストレージに 使用されます。 | 20 GB 未満。 データベースは、最大 8 GB まで拡張できます。 | 環境と共に徐々に拡張します。メタデータのみを格納します。 メモリーに 8 GB が追加されるたびに 20-25 GB を追加します。 |
| /var/lib/containers | これは CRI-O ランタイムのマウントポイントです。アクティブなコンテナーランタイム (Pod を含む) およびローカルイメージのストレージに使用されるストレージです。レジストリーストレージには使用されません。 | 16 GB メモリーの場合、 1 ノードにつき 50 GB。 このサイジングは、クラ スターの最小要件の決定 には使用しないでください。 メモリーに 8 GB が追加 されるたびに 20-25 GB を追加します。 | 拡張は実行中のコンテナーの容量によって制限されます。 |
| /var/lib/kubelet | Pod の一時ボリュームストレージです。これには、ランタイムにコンテナーにマウントされるすいです。環境変数、kubeシークレット、および永続ボリュームが含まれていデータボリュームが含まれます。 | 変動あり。 | ストレージを必要とする Podが永続ボリュームを 使用している場合は最小 になります。一時スト レージを使用する場合は すぐに拡張する可能性が あります。 |
| /var/log | すべてのコンポーネント のログファイルです。 | 10 から 30 GB。 | ログファイルはすぐに拡 張する可能性がありま す。サイズは拡張する ディスク別に管理する か、ログローテーション を使用して管理できま す。 |

第11章 ルーティングの最適化

OpenShift Container Platform HAProxy ルーターは、パフォーマンスを最適化するためにスケーリングします。

11.1. ベースライン INGRESS コントローラー (ルーター) のパフォーマンス

OpenShift Container Platform Ingress コントローラーまたはルーターは、宛先が OpenShift Container Platform サービスのすべての外部トラフィックに対する Ingress ポイントです。

1秒に処理される HTTP 要求について、単一の HAProxy ルーターを評価する場合に、パフォーマンスは 多くの要因により左右されます。特に以下が含まれます。

- HTTP keep-alive/close モード
- ルートタイプ
- TLS セッション再開のクライアントサポート
- ターゲットルートごとの同時接続数
- ターゲットルート数
- バックエンドサーバーのページサイズ
- 基礎となるインフラストラクチャー (ネットワーク/SDN ソリューション、CPU など)

特定の環境でのパフォーマンスは異なりますが、Red Hat ラボはサイズが 4 vCPU/16 GB RAM のパブリッククラウドインスタンスでテストしています。1 kB 静的ページを提供するバックエンドで終端する 100 ルートを処理する単一の HAProxy ルーターは、 1 秒あたりに以下の数のトランザクションを処理できます。

HTTP keep-alive モードのシナリオの場合:

| 暗号化 | LoadBalancerService | HostNetwork |
|-------------|---------------------|-------------|
| なし | 21515 | 29622 |
| edge | 16743 | 22913 |
| passthrough | 36786 | 53295 |
| re-encrypt | 21583 | 25198 |

HTTP close (keep-alive なし) のシナリオの場合:

| 暗号化 | LoadBalancerService | HostNetwork |
|------|---------------------|-------------|
| なし | 5719 | 8273 |
| edge | 2729 | 4069 |

| 暗号化 | LoadBalancerService | HostNetwork |
|-------------|---------------------|-------------|
| passthrough | 4121 | 5344 |
| re-encrypt | 2320 | 2941 |

ROUTER_THREADS=4 が設定されたデフォルトの Ingress コントローラー設定が使用され、2 つの異なるエンドポイントの公開ストラテジー (LoadBalancerService/HostNetwork) がテストされています。 TLS セッション再開は暗号化ルートについて使用されています。 HTTP keep-alive の場合は、単一の HAProxy ルーターがページサイズが 8kB でも、1 Gbit の NIC を飽和させることができます。

最新のプロセッサーが搭載されたベアメタルで実行する場合は、上記のパブリッククラウドインスタンスのパフォーマンスの約2倍のパフォーマンスになることを予想できます。このオーバーヘッドは、パブリッククラウドにある仮想化層により発生し、プライベートクラウドベースの仮想化にも多くの場合、該当します。以下の表は、ルーターの背後で使用するアプリケーション数についてのガイドです。

| アプリケーション数 | アプリケーションタイプ |
|-----------|-------------------------------|
| 5-10 | 静的なファイル/Web サーバーまたはキャッシュプロキシー |
| 100-1000 | 動的なコンテンツを生成するアプリケーション |

通常、HAProxy は、使用されるて技術に応じて 5 から 1000 のアプリケーションのルーターをサポートします。Ingress コントローラーのパフォーマンスは、言語や静的コンテンツと動的コンテンツの違いを含め、その背後にあるアプリケーションの機能およびパフォーマンスによって制限される可能性があります。

Ingress またはルーターのシャード化は、アプリケーションに対してより多くのルートを提供するために使用され、ルーティング層の水平スケーリングに役立ちます。

Ingress のシャード化についての詳細は、ルートラベルを使用した Ingress コントローラーのシャード化 の設定 および namespace ラベルを使用した Ingress コントローラーのシャード化の設定 を参照してください。

11.2. INGRESS コントローラー (ルーター) のパフォーマンスの最適化

OpenShift Container Platform では、環境変数

(ROUTER_THREADS、ROUTER_DEFAULT_TUNNEL_TIMEOUT、ROUTER_DEFAULT_CLIENT_TIMEOUT、ROUTER_DEFAULT_SERVER_TIMEOUT、および RELOAD_INTERVAL) を設定して Ingress コントローラーのデプロイメントを変更することをサポートしていません。

Ingress コントローラーのデプロイメントは変更できますが、Ingress Operator が有効にされている場合、設定は上書きされます。

第12章 ネットワークの最適化

OpenShift SDN は OpenvSwitch、VXLAN (Virtual extensible LAN) トンネル、OpenFlow ルール、iptables を使用します。このネットワークは、ジャンボフレーム、ネットワークインターフェイスコントローラー (NIC) オフロード、マルチキュー、および ethtool 設定を使用して調整できます。

OVN-Kubernetes は、トンネルプロトコルとして VXLAN ではなく Geneve (Generic Network Virtualization Encapsulation) を使用します。

VXLAN は、4096 から 1600 万以上にネットワーク数が増え、物理ネットワーク全体で階層 2 の接続が追加されるなど、VLAN での利点が提供されます。これにより、異なるシステム上で実行されている場合でも、サービスの背後にある Pod すべてが相互に通信できるようになります。

VXLAN は、User Datagram Protocol (UDP) パケットにトンネル化されたトラフィックをすべてカプセル化しますが、CPU 使用率が上昇してしまいます。これらの外部および内部パケットは、移動中にデータが破損しないようにするために通常のチェックサムルールの対象になります。これらの外部および内部パケットはどちらも、移動中にデータが破損しないように通常のチェックサムルールの対象になります。CPU のパフォーマンスによっては、この追加の処理オーバーヘッドによってスループットが減り、従来の非オーバーレイネットワークと比較してレイテンシーが高くなります。

クラウド、仮想マシン、ベアメタルの CPU パフォーマンスでは、1 Gbps をはるかに超えるネットワークスループットを処理できます。10 または 40 Gbps などの高い帯域幅のリンクを使用する場合には、パフォーマンスが低減する場合があります。これは、VXLAN ベースの環境では既知の問題で、コンテナーや OpenShift Container Platform 固有の問題ではありません。VXLAN トンネルに依存するネットワークも、VXLAN 実装により同様のパフォーマンスになります。

1Gbps 以上にするには、以下を実行してください。

- Border Gateway Protocol (BGP) など、異なるルーティング技術を実装するネットワークプラグインを評価する。
- VXLAN オフロード対応のネットワークアダプターを使用します。VXLAN オフロードは、システムの CPU から、パケットのチェックサム計算と関連の CPU オーバーヘッドを、ネットワークアダプター上の専用のハードウェアに移動します。これにより、CPU サイクルを Pod やアプリケーションで使用できるように開放し、ネットワークインフラストラクチャーの帯域幅すべてをユーザーは活用できるようになります。

VXLAN オフロードはレイテンシーを短縮しません。ただし、CPU の使用率はレイテンシーテストでも 削減されます。

12.1. ネットワークでの MTU の最適化

重要な Maximum Transmission Unit (MTU) が 2 つあります。1 つはネットワークインターフェイスコントローラー (NIC) MTU で、もう 1 つはクラスターネットワーク MTU です。

NIC MTU は OpenShift Container Platform のインストール時にのみ設定されます。MTU は、お使いのネットワークの NIC でサポートされる最大の値以下でなければなりません。スループットを最適化する場合は、可能な限り大きい値を選択します。レイテンシーを最低限に抑えるために最適化するには、より小さい値を選択します。

SDN オーバーレイの MTU は、最低でも NIC MTU より 50 バイト少なくなければなりません。これは、SDN オーバーレイのヘッダーに相当します。そのため、通常のイーサネットネットワークでは、この値を **1450** に設定します。ジャンボフレームのイーサネットネットワークの場合は、これを **8950** に設定します。

OVN および Geneve については、MTU は最低でも NIC MTU より 100 バイト少なくなければなりません。



注記

50 バイトのオーバーレイヘッダーは OpenShift SDN に関連します。他の SDN ソリューションの場合はこの値を若干変動させる必要があります。

12.2. 大規模なクラスターのインストールに推奨されるプラクティス

大規模なクラスターをインストールする場合や、クラスターを大規模なノード数に拡張する場合、クラスターをインストールする前に、install-config.yaml ファイルに適宜クラスターネットワーク cidr を設定します。

networking:

clusterNetwork:

 cidr: 10.128.0.0/14 hostPrefix: 23

machineCIDR: 10.0.0.0/16 networkType: OpenShiftSDN

serviceNetwork: - 172.30.0.0/16

クラスターのサイズが 500 を超える場合、デフォルトのクラスターネットワーク cidr 10.128.0.0/14 を使用することはできません。500 ノードを超えるノード数にするには、10.128.0.0/12 または 10.128.0.0/10 に設定する必要があります。

12.3. IPSEC の影響

ノードホストの暗号化、復号化に CPU 機能が使用されるので、使用する IP セキュリティーシステムにかかわらず、ノードのスループットおよび CPU 使用率の両方でのパフォーマンスに影響があります。

IPSec は、NIC に到達する前に IPペイロードレベルでトラフィックを暗号化して、NIC オフロードに使用されてしまう可能性のあるフィールドを保護します。つまり、IPSec が有効な場合には、NIC アクセラレーション機能を使用できない場合があり、スループットの減少、CPU 使用率の上昇につながります。

関連情報

- 高度なネットワーク設定パラメーターの変更
- OVN-Kubernetes デフォルト CNI ネットワークプロバイダーの設定パラメーター
- OpenShift SDN デフォルト CNI ネットワークプロバイダーの設定パラメーター

第13章 ベアメタルホストの管理

OpenShift Container Platform をベアメタルクラスターにインストールする場合、クラスターに存在するベアメタルホストの **machine** および **machineset** カスタムリソース (CR) を使用して、ベアメタルノードをプロビジョニングし、管理できます。

13.1. ベアメタルホストおよびノードについて

Red Hat Enterprise Linux CoreOS (RHCOS) ベアメタルホストをクラスター内のノードとしてプロビジョニングするには、まずベアメタルホストハードウェアに対応する **MachineSet** カスタムリソース (CR) オブジェクトを作成します。ベアメタルホストマシンセットは、お使いの設定に固有のインフラストラクチャーコンポーネントを記述します。特定の Kubernetes ラベルをこれらのマシンセットに適用してから、インフラストラクチャーコンポーネントを更新して、それらのマシンでのみ実行されるようにします。

Machine CR は、**metal3.io/autoscale-to-hosts** アノテーションを含む関連する **MachineSet** をスケールアップする際に自動的に作成されます。OpenShift Container Platform は **Machine** CR を使用して、**MachineSet** CR で指定されるホストに対応するベアメタルノードをプロビジョニングします。

13.2. ベアメタルホストのメンテナーンス

OpenShift Container Platform Web コンソールからクラスター内のベアメタルホストの詳細を維持することができます。Compute \rightarrow Bare Metal Hosts に移動し、Actions ドロップダウンメニューからタスクを選択します。ここでは、BMC の詳細、ホストの起動 MAC アドレス、電源管理の有効化などの項目を管理できます。また、ホストのネットワークインターフェイスおよびドライブの詳細を確認することもできます。

ベアメタルホストをメンテナーンスモードに移行できます。ホストをメンテナーンスモードに移行すると、スケジューラーはすべての管理ワークロードを対応するベアメタルノードから移動します。新しいワークロードは、メンテナーンスモードの間はスケジュールされません。

Web コンソールでベアメタルホストのプロビジョニングを解除することができます。ホストのプロビジョニング解除により以下のアクションが実行されます。

- 1. ベアメタルホスト CR に cluster.k8s.io/delete-machine: true のアノテーションを付けます。
- 2. 関連するマシンセットをスケールダウンします。



注記

デーモンセットおよび管理対象外の静的 Pod を別のノードに最初に移動することなく、ホストの電源をオフにすると、サービスの中断やデータの損失が生じる場合があります。

関連情報

コンピュートマシンのベアメタルへの追加

13.2.1. Web コンソールを使用したベアメタルホストのクラスターへの追加

Web コンソールのクラスターにベアメタルホストを追加できます。

前提条件

- RHCOS クラスターのベアメタルへのインストール
- cluster-admin 権限を持つユーザーとしてログインしている。

手順

- 1. Web コンソールで、Compute → Bare Metal Hosts に移動します。
- 2. Add Host → New with Dialog を選択します。
- 3. 新規ベアメタルホストの一意の名前を指定します。
- 4. Boot MAC address を設定します。
- 5. Baseboard Management Console (BMC) Addressを設定します。
- 6. 必要に応じて、ホストの電源管理を有効にします。これにより、OpenShift Container Platform はホストの電源状態を制御できます。
- 7. ホストのベースボード管理コントローラー (BMC) のユーザー認証情報を入力します。
- 8. 作成後にホストの電源をオンにすることを選択し、Create を選択します。
- 9. 利用可能なベアメタルホストの数に一致するようにレプリカ数をスケールアップします。Compute → MachineSets に移動し、Actions ドロップダウンメニューから Edit Machine count を選択してクラスター内のマシンレプリカ数を増やします。



注記

oc scale コマンドおよび適切なべアメタルマシンセットを使用して、ベアメタルノードの数を管理することもできます。

13.2.2. Web コンソールの YAML を使用したベアメタルホストのクラスターへの追加

ベアメタルホストを記述する YAML ファイルを使用して、Web コンソールのクラスターにベアメタルホストを追加できます。

前提条件

- クラスターで使用するために RHCOS コンピュートマシンをベアメタルインフラストラクチャーにインストールします。
- cluster-admin 権限を持つユーザーとしてログインしている。
- ベアメタルホストの Secret CR を作成します。

手順

- 1. Web コンソールで、Compute → Bare Metal Hosts に移動します。
- 2. Add Host → New from YAML を選択します。
- 3. 以下の YAML をコピーして貼り付け、ホストの詳細で関連フィールドを変更します。

apiVersion: metal3.io/v1alpha1

kind: BareMetalHost

metadata:

name: <bare_metal_host_name>

spec:

online: true bmc:

address:
bmc address>

credentialsName: <secret credentials name> 1

disableCertificateVerification: True

bootMACAddress: < host boot mac address>

hardwareProfile: unknown

- credentialsName は有効な Secret CR を参照する必要があります。baremetal-operator は、credentialsName で参照される有効な Secret なしに、ベアメタルホストを管理できません。シークレットの詳細および作成方法については、シークレットについて を参照してください。
- 4. Create を選択して YAML を保存し、新規ベアメタルホストを作成します。
- 5. 利用可能なベアメタルホストの数に一致するようにレプリカ数をスケールアップします。Compute → MachineSets に移動し、Actions ドロップダウンメニューから Edit Machine count を選択してクラスター内のマシン数を増やします。



注記

oc scale コマンドおよび適切なべアメタルマシンセットを使用して、ベアメタルノードの数を管理することもできます。

13.2.3. 利用可能なベアメタルホストの数へのマシンの自動スケーリング

利用可能な BareMetalHost オブジェクトの数に一致する Machine オブジェクトの数を自動的に作成するには、metal3.io/autoscale-to-hosts アノテーションを MachineSet オブジェクトに追加します。

前提条件

- クラスターで使用する RHCOS ベアメタルコンピュートマシンをインストールし、対応する **BareMetalHost** オブジェクトを作成します。
- OpenShift Container Platform CLI (oc) をインストールします。
- cluster-admin 権限を持つユーザーとしてログインしている。

手順

1. **metal3.io/autoscale-to-hosts** アノテーションを追加して、自動スケーリング用に設定するマシンセットにアノテーションを付けます。**<machineset>** を、マシンセット名に置き換えます。

\$ oc annotate machineset <machineset> -n openshift-machine-api 'metal3.io/autoscale-to-hosts=<any_value>'

新しいスケーリングされたマシンが起動するまで待ちます。



注記

BareMetalHost オブジェクトを使用してクラスター内にマシンを作成し、その後ラベルまたはセレクターが BareMetalHost で変更される場合、BareMetalHost オブジェクトは Machine オブジェクトが作成された MachineSet に対して引き続きカウントされます。

関連情報

- クラスターの拡張
- ベアメタル上の MachineHealthCheck

第14章 HUGE PAGE の機能およびそれらがアプリケーションに よって消費される仕組み

14.1. HUGE PAGE の機能

メモリーは Page と呼ばれるブロックで管理されます。多くのシステムでは、1ページは 4Ki です。メモリー 1Mi は 256 ページに、メモリー 1Gi は 256,000 ページに相当します。CPU には、内蔵のメモリー管理ユニットがあり、ハードウェアでこのようなページリストを管理します。トランスレーションルックアサイドバッファー (TLB: Translation Lookaside Buffer) は、仮想から物理へのページマッピングの小規模なハードウェアキャッシュのことです。ハードウェアの指示で渡された仮想アドレスが TLBにあれば、マッピングをすばやく決定できます。そうでない場合には、TLBミスが発生し、システムは速度が遅く、ソフトウェアベースのアドレス変換にフォールバックされ、パフォーマンスの問題が発生します。TLBのサイズは固定されているので、TLBミスの発生率を減らすには Page サイズを大きくする必要があります。

Huge Page とは、4Ki より大きいメモリーページのことです。 $x86_64$ アーキテクチャーでは、2Mi と 1Gi の 2 つが一般的な Huge Page サイズです。別のアーキテクチャーではサイズは異なります。Huge Page を使用するには、アプリケーションが認識できるようにコードを書き込む必要があります。 Transparent Huge Pages (THP) は、アプリケーションによる認識なしに、Huge Page の管理を自動化しようとしますが、制約があります。特に、ページサイズは 2Mi に制限されます。THP では、THP のデフラグが原因で、メモリー使用率が高くなり、断片化が起こり、パフォーマンスの低下につながり、メモリーページがロックされてしまう可能性があります。このような理由から、アプリケーションは THP ではなく、事前割り当て済みの Huge Page を使用するように設計 (また推奨) される場合があります。

OpenShift Container Platform では、Pod のアプリケーションが事前に割り当てられた Huge Page を割り当て、消費することができます。

14.2. HUGE PAGE がアプリケーションによって消費される仕組み

ノードは、Huge Page の容量をレポートできるように Huge Page を事前に割り当てる必要があります。ノードは、単一サイズの Huge Page のみを事前に割り当てることができます。

Huge Page は、リソース名の **hugepages-<size>** を使用してコンテナーレベルのリソース要件で消費可能です。この場合、サイズは特定のノードでサポートされる整数値を使用した最もコンパクトなバイナリー表記です。たとえば、ノードが 2048KiB ページサイズをサポートする場合、これはスケジュール可能なリソース **hugepages-2Mi** を公開します。CPU やメモリーとは異なり、Huge Page はオーバーコミットをサポートしません。

apiVersion: v1 kind: Pod metadata:

generateName: hugepages-volume-

spec:

containers:

 securityContext: privileged: true image: rhel7:latest command:

- sleep

- inf

name: example volumeMounts:

- mountPath: /dev/hugepages

name: hugepage
resources:
limits:
hugepages-2Mi: 100Mi 1
memory: "1Gi"
cpu: "1"
volumes:
- name: hugepage
emptyDir:
medium: HugePages

1 hugepages のメモリー量は、実際に割り当てる量に指定します。この値は、ページサイズで乗算した hugepages のメモリー量に指定しないでください。たとえば、Huge Page サイズが 2MB と仮定し、アプリケーションに Huge Page でバックアップする RAM 100 MB を使用する場合には、Huge Page は 50 に指定します。OpenShift Container Platform により、計算処理が実行されます。上記の例にあるように、100MB を直接指定できます。

指定されたサイズの Huge Page の割り当て

プラットフォームによっては、複数の Huge Page サイズをサポートするものもあります。特定のサイズの Huge Page を割り当てるには、Huge Page の起動コマンドパラメーターの前に、Huge Page サイズの選択パラメーター hugepagesz=<size> を指定してください。<size> の値は、バイトで指定する必要があります。その際、オプションでスケール接尾辞 [kKmMgG] を指定できます。デフォルトの Huge Page サイズは、default hugepagesz=<size> の起動パラメーターで定義できます。

Huge page の要件

- Huge Page 要求は制限と同じでなければなりません。制限が指定されているにもかかわらず、 要求が指定されていない場合には、これがデフォルトになります。
- Huge Page は、Pod のスコープで分割されます。コンテナーの分割は、今後のバージョンで予定されています。
- Huge Page がサポートする **EmptyDir** ボリュームは、Pod 要求よりも多くの Huge Page メモリーを消費することはできません。
- **shmget()** で **SHM_HUGETLB** を使用して Huge Page を消費するアプリケーション は、proc/sys/vm/hugetlb_shm_group に一致する補助グループで実行する必要があります。

関連情報

● Transparent Huge Page の設定

14.3. HUGE PAGE の設定

ノードは、OpenShift Container Platform クラスターで使用される Huge Page を事前に割り当てる必要があります。Huge Page を予約する方法は、ブート時とランタイム時に実行する 2 つの方法があります。ブート時の予約は、メモリーが大幅に断片化されていないために成功する可能性が高くなります。Node Tuning Operator は、現時点で特定のノードでの Huge Page のブート時の割り当てをサポートします。

14.3.1. ブート時

手順

ノードの再起動を最小限にするには、以下の手順の順序に従う必要があります。

1. ラベルを使用して同じ Huge Page 設定を必要とするすべてのノードにラベルを付けます。

\$ oc label node <node_using_hugepages> node-role.kubernetes.io/worker-hp=

2. 以下の内容でファイルを作成し、これに hugepages-tuned-boottime.yaml という名前を付けます。

```
apiVersion: tuned.openshift.io/v1
kind: Tuned
metadata:
 name: hugepages 1
 namespace: openshift-cluster-node-tuning-operator
spec:
 profile: 2
 - data: |
   [main]
   summary=Boot time configuration for hugepages
   include=openshift-node
   [bootloader]
   cmdline_openshift_node_hugepages=hugepagesz=2M hugepages=50
  name: openshift-node-hugepages
 recommend:
 - machineConfigLabels: 4
   machineconfiguration.openshift.io/role: "worker-hp"
  priority: 30
  profile: openshift-node-hugepages
```

- 🚹 チューニングされたリソースの name を hugepages に設定します。
- 2 Huge Page を割り当てる **profile** セクションを設定します。
- 一部のプラットフォームではさまざまなサイズの Huge Page をサポートするため、パラメーターの順序に注意してください。
- マシン設定プールベースのマッチングを有効にします。
- 3. チューニングされた hugepages プロファイルの作成

\$ oc create -f hugepages-tuned-boottime.yaml

4. 以下の内容でファイルを作成し、これに hugepages-mcp.yaml という名前を付けます。

apiVersion: machineconfiguration.openshift.io/v1
kind: MachineConfigPool
metadata:
name: worker-hp
labels:
worker-hp: ""
spec:
machineConfigSelector:
matchExpressions:

- {key: machineconfiguration.openshift.io/role, operator: In, values: [worker,worker-hp]} nodeSelector:

matchLabels:

node-role.kubernetes.io/worker-hp: ""

5. マシン設定プールを作成します。

\$ oc create -f hugepages-mcp.yaml

断片化されていないメモリーが十分にある場合、worker-hp マシン設定プールのすべてのノードには 50 2Mi の Huge Page が割り当てられているはずです。

\$ oc get node <node_using_hugepages> -o jsonpath="{.status.allocatable.hugepages-2Mi}" 100Mi



警告

この機能は、現在 Red Hat Enterprise Linux CoreOS (RHCOS) 8.x ワーカーノードでのみサポートされています。Red Hat Enterprise Linux (RHEL) 7.x ワーカーノードでは、チューニングされた **[bootloader]** プラグインは現時点でサポートされていません。

第15章 低レイテンシーのノード向けの PERFORMANCE ADDON OPERATOR

15.1. 低レイテンシー

Telco / 5G の領域でのエッジコンピューティングの台頭は、レイテンシーと輻輳を軽減し、アプリケーションのパフォーマンスを向上させる上で重要なロールを果たします。

簡単に言うと、レイテンシーは、データ (パケット) が送信側から受信側に移動し、受信側の処理後に送信側に戻るスピードを決定します。レイテンシーによる遅延を最小限に抑えた状態でネットワークアーキテクチャーを維持することが 5 G のネットワークパフォーマンス要件を満たすのに鍵となります。 4G テクノロジーと比較し、平均レイテンシーが 50ms の 5 G では、レイテンシーの数値を 1ms 以下にするようにターゲットが設定されます。このレイテンシーの減少により、ワイヤレスのスループットが10 倍向上します。

Telco 領域にデプロイされるアプリケーションの多くは、ゼロパケットロスに耐えられる低レイテンシーを必要とします。パケットロスをゼロに調整すると、ネットワークのパフォーマンス低下させる固有の問題を軽減することができます。詳細は、Tuning for Zero Packet Loss in Red Hat OpenStack Platform (RHOSP) を参照してください。

エッジコンピューティングの取り組みは、レイテンシーの削減にも役立ちます。コンピュート能力が文字通りクラウドのエッジ上にあり、ユーザーの近く置かれること考えてください。これにより、ユーザーと離れた場所にあるデータセンター間の距離が大幅に削減されるため、アプリケーションの応答時間とパフォーマンスのレイテンシーが短縮されます。

管理者は、すべてのデプロイメントを可能な限り低い管理コストで実行できるように、多数のエッジサイトおよびローカルサービスを一元管理できるようにする必要があります。また、リアルタイムの低レイテンシーおよび高パフォーマンスを実現するために、クラスターの特定のノードをデプロイし、設定するための簡単な方法も必要になります。低レイテンシーノードは、Cloud-native Network Functions (CNF) や Data Plane Development Kit (DPDK) などのアプリケーションに役立ちます。

現時点で、OpenShift Container Platform はリアルタイムの実行および低レイテンシーを実現するため に OpenShift Container Platform クラスターでソフトウェアを調整するメカニズムを提供します (約 20 マイクロ秒未満の応答時間)。これには、カーネルおよび OpenShift Container Platform の設定値の チューニング、カーネルのインストール、およびマシンの再設定が含まれます。ただし、この方法では 4 つの異なる Operator を設定し、手動で実行する場合に複雑であり、間違いが生じる可能性がある多くの設定を行う必要があります。

OpenShift Container Platform は、OpenShift アプリケーションの低レイテンシーパフォーマンスを実現するために自動チューニングを実装する Performance Addon Operator を提供します。クラスター管理者は、このパフォーマンスプロファイル設定を使用することにより、より信頼性の高い方法でこれらの変更をより容易に実行することができます。管理者は、カーネルを kernel-rt、ハウスキーピング用に予約される CPU、ワークロードの実行に使用される CPU に更新するかどうかを指定できます。

15.2. PERFORMANCE ADDON OPERATOR のインストール

Performance Addon Operator は、一連のノードで高度なノードのパフォーマンスチューニングを有効にする機能を提供します。クラスター管理者は、OpenShift Container Platform CLI または Web コンソールを使用して Performance Addon Operator をインストールできます。

15.2.1. CLI を使用した **Operator** のインストール

クラスター管理者は、CLI を使用して Operator をインストールできます。

前提条件

- ▼アメタルハードウェアにインストールされたクラスター。
- OpenShift CLI (oc) をインストールしている。
- cluster-admin 権限を持つユーザーとしてログインしている。

手順

- 1. 以下のアクションを実行して、Performance Addon Operator の namespace を作成します。
 - a. **openshift-performance-addon-operator** namespace を定義する以下の Namespace カス タムリソース (CR) を作成し、YAML を **pao-namespace.yaml** ファイルに保存します。

apiVersion: v1 kind: Namespace metadata:

name: openshift-performance-addon-operator

b. 以下のコマンドを実行して namespace を作成します。

\$ oc create -f pao-namespace.yaml

- 2. 以下のオブジェクトを作成して、直前の手順で作成した namespace に Performance Addon Operator をインストールします。
 - a. 以下の **OperatorGroup** CR を作成し、YAML を **pao-operatorgroup.yaml** ファイルに保存します。

apiVersion: operators.coreos.com/v1

kind: OperatorGroup

metadata:

name: openshift-performance-addon-operator namespace: openshift-performance-addon-operator spec:

targetNamespaces:

- openshift-performance-addon-operator
- b. 以下のコマンドを実行して Operator Group CR を作成します。

\$ oc create -f pao-operatorgroup.yaml

c. 以下のコマンドを実行して、次の手順に必要な channel の値を取得します。

\$ oc get packagemanifest performance-addon-operator -n openshift-marketplace -o jsonpath='{.status.defaultChannel}'

出力例

4.6

d. 以下の Subscription CR を作成し、YAML を **pao-sub.yaml** ファイルに保存します。

Subscription の例

apiVersion: operators.coreos.com/v1alpha1

kind: Subscription

metadata:

name: openshift-performance-addon-operator-subscription

namespace: openshift-performance-addon-operator

spec:

channel: "<channel>" 1

name: performance-addon-operator

source: redhat-operators 2

sourceNamespace: openshift-marketplace

- **status.defaultChannel** パラメーターの直前の手順で取得した値を指定します。
- redhat-operators 値を指定する必要があります。
- e. 以下のコマンドを実行して Subscription オブジェクトを作成します。

\$ oc create -f pao-sub.yaml

f. openshift-performance-addon-operator プロジェクトに切り替えます。

\$ oc project openshift-performance-addon-operator

15.2.2. Web コンソールを使用した Performance Addon Operator のインストール

クラスター管理者は、Web コンソールを使用して Performance Addon Operator をインストールできます。



注記

先のセクションで説明されているように **Namespace** CR および **OperatorGroup** CR を作成する必要があります。

手順

- 1. OpenShift Container Platform Web コンソールを使用して Performance Addon Operator をインストールします。
 - a. OpenShift Container Platform Web コンソールで、**Operators** → **OperatorHub** をクリックします。
 - b. 利用可能な Operator の一覧から **Performance Addon Operator** を選択してから **Install** を クリックします。
 - c. Install Operator ページの A specific namespace on the clusterの下で、 openshift-performance-addon-operator を選択します。次に、Install をクリックします。
- 2. オプション: performance-addon-operator が正常にインストールされていることを確認します。
 - a. Operators → Installed Operators ページに切り替えます。

b. Performance Addon Operator が openshift-operators プロジェクトに Succeeded の Status でリストされていることを確認します。



注記

インストール時に、 Operator は **Failed** ステータスを表示する可能性があります。インストールが成功し、**Succeeded** メッセージが表示された場合は、**Failed** メッセージを無視できます。

Operator がインストール済みとして表示されない場合は、さらにトラブルシューティングを行うことができます。

- Operators → Installed Operators ページに移動し、Operator Subscriptions および Install Plans タブで Status にエラーがあるかどうかを検査します。
- Workloads → Pods ページに移動し、openshift-operators プロジェクトで Pod のログ を確認します。

15.3. PERFORMANCE ADDON OPERATOR のアップグレード

次のマイナーバージョンの Performance Addon Operator に手動でアップグレードし、Web コンソールを使用して更新のステータスをモニターできます。

15.3.1. Performance Addon Operator のアップグレードについて

- OpenShift Web コンソールを使用して Operator サブスクリプションのチャネルを変更することで、Performance Addon Operator の次のマイナーバージョンにアップグレードできます。
- Performance Addon Operator のインストール時に z-stream の自動更新を有効にできます。
- 更新は、OpenShift Container Platform のインストール時にデプロイされる Marketplace Operator 経由で提供されます。 Marketplace Operator は外部 Operator をクラスターで利用可能にします。
- 更新の完了までにかかる時間は、ネットワーク接続によって異なります。ほとんどの自動更新は15分以内に完了します。

15.3.1.1. Performance Addon Operator のクラスターへの影響

- 低レイテンシーのチューニング Huge Page は影響を受けません。
- Operator を更新しても、予期しない再起動は発生しません。

15.3.1.2. Performance Addon Operator の次のマイナーバージョンへのアップグレード

OpenShift Container Platform Web コンソールを使用して Operator サブスクリプションのチャネルを変更することで、Performance Addon Operator を次のマイナーバージョンに手動でアップグレードできます。

前提条件

• cluster-admin ロールを持つユーザーとしてのクラスターへのアクセスがあること。

手順

- 1. OpenShift Web コンソールにアクセスし、 Operators → Installed Operators に移動します。
- 2. Performance Addon Operator をクリックし、Operator Details ページを開きます。
- 3. Subscription タブをクリックし、Subscription Overview ページを開きます。
- 4. Channel ペインで、バージョン番号の右側にある鉛筆アイコンをクリックし、Change Subscription Update Channel ウィンドウを開きます。
- 5. 次のマイナーバージョンを選択します。たとえば、Performance Addon Operator 4.6 にアップグレードする場合は、**4.6** を選択します。
- 6. **Save** をクリックします。
- 7. Operators → Installed Operators に移動してアップグレードのステータスを確認します。以下 の oc コマンドを実行してステータスを確認することもできます。

\$ oc get csv -n openshift-performance-addon-operator

15.3.2. アップグレードステータスの監視

Performance Addon Operator アップグレードステータスをモニターする最適な方法として、**ClusterServiceVersion** (CSV) **PHASE** を監視できます。Web コンソールを使用するか、または**oc get csv** コマンドを実行して CSV の状態をモニターすることもできます。



注記

PHASE および状態の値は利用可能な情報に基づく近似値になります。

前提条件

- cluster-admin ロールを持つユーザーとしてクラスターにアクセスできる。
- OpenShift CLI (oc) をインストールしている。

手順

1. 以下のコマンドを実行します。

\$ oc get csv

2. 出力を確認し、PHASE フィールドをチェックします。以下に例を示します。

VERSION REPLACES PHASE
4.6.0 performance-addon-operator.v4.5.0 Installing
4.5.0 Replacing

3. get csv を再度実行して出力を確認します。

oc get csv

出力例

NAME

DISPLAY

VERSION REPLACES

PHASE

performance-addon-operator.v4.5.0 Performance Addon Operator 4.6.0 performance-addon-operator.v4.5.0 Succeeded

15.4. リアルタイムおよび低レイテンシーワークロードのプロビジョニング

多くの企業や組織は、非常に高性能なコンピューティングを必要としており、とくに金融業界や通信業界では、低い、予測可能なレイテンシーが必要になる場合があります。このような固有の要件を持つ業界では、OpenShift Container Platform は Performance Addon Operator を提供して、OpenShift Container Platform アプリケーションの低レイテンシーのパフォーマンスと一貫性のある応答時間を実現するための自動チューニングを実装します。

クラスター管理者は、このパフォーマンスプロファイル設定を使用することにより、より信頼性の高い方法でこれらの変更をより容易に実行することができます。管理者は、カーネルを kernel-rt (リアルタイム)、ハウスキーピング用に予約される CPU、ワークロードの実行に使用される CPU に更新するかどうかを指定できます。



警告

保証された CPU を必要とするアプリケーションと組み合わせて実行プローブを使用すると、レイテンシースパイクが発生する可能性があります。代わりに、適切に設定されたネットワークプローブのセットなど、他のプローブを使用することをお勧めします。

15.4.1. リアルタイムの既知の制限



注記

RT カーネルはワーカーノードでのみサポートされます。

リアルタイムモードを完全に使用するには、コンテナーを昇格した権限で実行する必要があります。権限の付与についての情報は、Set capabilities for a Container を参照してください。

OpenShift Container Platform は許可される機能を制限するため、**SecurityContext** を作成する必要がある場合もあります。



注記

この手順は、Red Hat Enterprise Linux CoreOS (RHCOS) システムを使用したベアメタルのインストールで完全にサポートされます。

パフォーマンスの期待値を設定する必要があるということは、リアルタイムカーネルがあらゆる問題の解決策ではないということを意味します。リアルタイムカーネルは、一貫性のある、低レイテンシーの、決定論に基づく予測可能な応答時間を提供します。リアルタイムカーネルに関連して、追加のカーネルオーバーヘッドがあります。これは、主に個別にスケジュールされたスレッドでハードウェア割り

込みを処理することによって生じます。一部のワークロードのオーバーヘッドが増加すると、スループット全体が低下します。ワークロードによって異なりますが、パフォーマンスの低下の程度は 0% から 30% の範囲になります。ただし、このコストは決定論をベースとしています。

15.4.2. リアルタイム機能のあるワーカーのプロビジョニング

- 1. Performance Addon Operator をクラスターにインストールします。
- 2. オプション: ノードを OpenShift Container Platform クラスターに追加します。BIOS パラメーターの設定 について参照してください。
- 3. ocコマンドを使用して、リアルタイム機能を必要とするワーカーノードにラベルworker-rtを追加します。
- 4. リアルタイムノード用の新しいマシン設定プールを作成します。

```
apiVersion: machineconfiguration.openshift.io/v1
kind: MachineConfigPool
metadata:
 name: worker-rt
 labels:
  machineconfiguration.openshift.io/role: worker-rt
 machineConfigSelector:
  matchExpressions:
   - {
       key: machineconfiguration.openshift.io/role.
       operator: In.
      values: [worker, worker-rt],
    }
 paused: false
 nodeSelector:
  matchLabels:
   node-role.kubernetes.io/worker-rt: ""
```

マシン設定プール worker-rt は、worker-rt というラベルを持つノードのグループに対して作成されることに注意してください。

5. ノードロールラベルを使用して、ノードを適切なマシン設定プールに追加します。



注記

リアルタイムワークロードで設定するノードを決定する必要があります。クラスター内のすべてのノード、またはノードのサブセットを設定できます。 Performance Addon Operator は、すべてのノードが専用のマシン設定プールの一部であることを想定します。すべてのノードを使用する場合は、Performance Addon Operator がワーカーノードのロールラベルを指すようにする必要があり

ます。サブセットを使用する場合、ノードを新規のマシン設定プールにグループ

化する必要があります。

- 6. ハウスキーピングコアの適切なセットと realTimeKernel: enabled: true を設定して PerformanceProfile を作成します。
- 7. PerformanceProfile で machineConfigPoolSelector を設定する必要があります:

apiVersion: performance.openshift.io/v2

kind: PerformanceProfile

metadata:

name: example-performanceprofile

spec:

...

realTimeKernel:

enabled: true nodeSelector:

node-role.kubernetes.io/worker-rt: ""

machineConfigPoolSelector:

machineconfiguration.openshift.io/role: worker-rt

8. 一致するマシン設定プールがラベルを持つことを確認します。

\$ oc describe mcp/worker-rt

出力例

Name:

worker-rt

Namespace:

_abels: machineconfiguration.openshift.io/role=worker-rt

- 9. OpenShift Container Platform はノードの設定を開始しますが、これにより複数の再起動が伴う可能性があります。ノードが起動し、安定するのを待機します。特定のハードウェアの場合に、これには長い時間がかかる可能性がありますが、ノードごとに 20 分の時間がかかることが予想されます。
- 10. すべてが予想通りに機能していることを確認します。

15.4.3. リアルタイムカーネルのインストールの確認

以下のコマンドを使用して、リアルタイムカーネルがインストールされていることを確認します。

\$ oc get node -o wide

文字列 **4.18.0-211.rt5.23.el8.x86_64** が含まれる、ロール worker-rt を持つワーカーに留意してください。

NAME STATUS ROLES AGE VERSION INTERNAL-IP

EXTERNAL-IP OS-IMAGE KERNEL-VERSION

CONTAINER-RUNTIME

rt-worker-0.example.com Ready worker,worker-rt 5d17h v1.22.1

128.66.135.107 <none> Red Hat Enterprise Linux CoreOS 46.82.202008252340-0 (Ootpa)

4.18.0-211.rt5.23.el8.x86 64 cri-o://1.19.0-90.rhaos4.6.git4a0ac05.el8-rc.1

[...]

15.4.4. リアルタイムで機能するワークロードの作成

リアルタイム機能を使用するワークロードを準備するには、以下の手順を使用します。

手順

- 1. QoS クラスの Guaranteed を指定して Pod を作成します。
- 2. オプション: DPDK の CPU 負荷分散を無効にします。
- 3. 適切なノードセレクターを割り当てます。

アプリケーションを作成する場合には、アプリケーションのチューニングとデプロイメント に記載されている一般的な推奨事項に従ってください。

15.4.5. QoS クラスの Guaranteed を指定した Pod の作成

QoS クラスの **Guaranteed** が指定されている Pod を作成する際には、以下を考慮してください。

- Pod のすべてのコンテナーにはメモリー制限およびメモリー要求があり、それらは同じである 必要があります。
- Pod のすべてのコンテナーには CPU の制限と CPU 要求が必要であり、それらは同じである必要があります。

以下の例は、1つのコンテナーを持つ Pod の設定ファイルを示しています。コンテナーにはメモリー制限とメモリー要求があり、どちらも 200 MiB に相当します。コンテナーには CPU 制限と CPU 要求があり、どちらも 1 CPU に相当します。

apiVersion: v1
kind: Pod
metadata:
name: qos-demo
namespace: qos-example
spec:
containers:
- name: qos-demo-ctr
image: <image-pull-spec>
resources:
limits:
memory: "200Mi"
cpu: "1"

memory: "200Mi" cpu: "1"

requests:

1. Pod を作成します。

\$ oc apply -f gos-pod.yaml --namespace=gos-example

2. Pod についての詳細情報を表示します。

\$ oc get pod qos-demo --namespace=qos-example --output=yaml

出力例

spec: containers: ... status:

qosClass: Guaranteed



注記

コンテナーが独自のメモリー制限を指定するものの、メモリー要求を指定しない場合、OpenShift Container Platform は制限に一致するメモリー要求を自動的に割り当てます。同様に、コンテナーが独自の CPU 制限を指定するものの、CPU要求を指定しない場合、OpenShift Container Platform は制限に一致する CPU要求を自動的に割り当てます。

15.4.6. オプション: DPDK 用の CPU 負荷分散の無効化

CPU 負荷分散を無効または有効にする機能は CRI-O レベルで実装されます。CRI-O のコードは、以下の要件を満たす場合にのみ CPU の負荷分散を無効または有効にします。

● Pod は **performance-<profile-name>** ランタイムクラスを使用する必要があります。以下に示すように、パフォーマンスプロファイルのステータスを確認して、適切な名前を取得できます。

```
apiVersion: performance.openshift.io/v1 kind: PerformanceProfile ... status: ... runtimeClass: performance-manual
```

● Pod には cpu-load-balancing.crio.io: true アノテーションが必要です。

Performance Addon Operator は、該当するノードで高パフォーマンスのランタイムハンドラー設定スニペットの作成や、クラスターで高パフォーマンスのランタイムクラスの作成を行います。これには、CPU 負荷分散の設定機能を有効にすることを除くと、デフォルトのランタイムハンドラーと同じ内容が含まれます。

Pod の CPU 負荷分散を無効にするには、 Pod 仕様に以下のフィールドが含まれる必要があります。

```
apiVersion: v1
kind: Pod
metadata:
...
annotations:
...
cpu-load-balancing.crio.io: "true"
...
spec:
...
runtimeClassName: performance-<profile_name>
...
```



注記

CPU マネージャーの静的ポリシーが有効にされている場合に、CPU 全体を使用する Guaranteed QoS を持つ Pod について CPU 負荷分散を無効にします。これ以外の場合 に CPU 負荷分散を無効にすると、クラスター内の他のコンテナーのパフォーマンスに影響する可能性があります。

15.4.7. 適切なノードセレクターの割り当て

Pod をノードに割り当てる方法として、以下に示すようにパフォーマンスプロファイルが使用するものと同じノードセレクターを使用することが推奨されます。

apiVersion: v1
kind: Pod
metadata:
name: example
spec:
...
nodeSelector:
node-role.kubernetes.io/worker-rt: ""

ノードセレクターの詳細は、Placing pods on specific nodes using node selectors を参照してください。

15.4.8. リアルタイム機能を備えたワーカーへのワークロードのスケジューリング

Performance Addon Operator によって低レイテンシーを確保するために設定されたマシン設定プールに割り当てられるノードに一致するラベルセレクターを使用します。詳細は、Assigning pods to nodesを参照してください。

15.5. HUGE PAGE の設定

ノードは、OpenShift Container Platform クラスターで使用される Huge Page を事前に割り当てる必要があります。Performance Addon Operator を使用し、特定のノードで Huge Page を割り当てます。

OpenShift Container Platform は、Huge Page を作成し、割り当てる方法を提供します。Performance Addon Operator は、パーマンスプロファイルを使用してこれを実行するための簡単な方法を提供します。

たとえば、パフォーマンスプロファイルの hugepages pages セクションで、size、count、およびオプションで node の複数のブロックを指定できます。

hugepages:

defaultHugepagesSize: "1G"

pages:
- size: "1G"
count: 4
node: 0 1

node は、Huge Page が割り当てられる NUMA ノードです。 **node** を省略すると、ページはすべて の NUMA ノード間で均等に分散されます。



注記

更新が完了したことを示す関連するマシン設定プールのステータスを待機します。

これらは、Huge Page を割り当てるのに必要な唯一の設定手順です。

検証

• 設定を確認するには、ノード上の /proc/meminfo ファイルを参照します。

\$ oc debug node/ip-10-0-141-105.ec2.internal

grep -i huge /proc/meminfo

出力例

AnonHugePages: ##### ##
ShmemHugePages: 0 kB
HugePages_Total: 2
HugePages_Free: 2
HugePages_Rsvd: 0
HugePages_Surp: 0
Hugepagesize: #####
Hugetlb: ######

● 新規サイズを報告するには、oc describe を使用します。

\$ oc describe node worker-0.ocp4poc.example.com | grep -i huge

出力例

hugepages-1g=true

hugepages-###: ### hugepages-###: ###

15.6. 複数の **HUGE PAGE** サイズの割り当て

同じコンテナーで異なるサイズの Huge Page を要求できます。これにより、Huge Page サイズのニーズの異なる複数のコンテナーで設定されるより複雑な Pod を定義できます。

たとえば、サイズ **1G** と **2M** を定義でき、Performance Addon Operator は以下に示すようにノード上に両方のサイズを設定できます。

spec:

hugepages:

defaultHugepagesSize: 1G

pages:

- count: 1024 node: 0 size: 2M - count: 4 node: 1 size: 1G

15.7. INFRA およびアプリケーションコンテナーの CPU の制限

一般的なハウスキーピングおよびワークロードタスクは、レイテンシーの影響を受けやすいプロセスに影響を与える可能性のある方法で CPU を使用します。デフォルトでは、コンテナーランタイムはすべてのオンライン CPU を使用して、すべてのコンテナーを一緒に実行します。これが原因で、コンテキストスイッチおよびレイテンシーが急増する可能性があります。 CPU をパーティション化することで、ノイズの多いプロセスとレイテンシーの影響を受けやすいプロセスを分離し、干渉を防ぐことができます。以下の表は、Performance Addon Operator を使用してノードを調整した後、CPU でプロセスがどのように実行されるかを示しています。

表15.1 プロセスの CPU 割り当て

| プロセスタイプ | Details |
|------------------------------|---|
| Burstable および BestEffort Pod | 低レイテンシーのワークロードが実行されている場合を除き、任意の CPU で実行されます。 |
| インフラストラクチャー Pod | 低レイテンシーのワークロードが実行されている場合を除き、任意の CPU で実行されます。 |
| 割り込み | 予約済み CPU にリダイレクトします (OpenShift Container Platform 4.6 以降ではオプション) |
| カーネルプロセス | 予約済み CPU へのピン |
| レイテンシーの影響を受けやすいワークロード Pod | 分離されたプールからの排他的 CPU の特定のセット へのピン |
| OS プロセス/systemd サービス | 予約済み CPU へのピン |

すべての QoS プロセスタイプ (**Burstable、BestEffort、**または **Guaranteed**) の Pod に割り当て可能 なノード上のコアの容量は、分離されたプールの容量と同じです。予約済みプールの容量は、クラス ターおよびオペレーティングシステムのハウスキーピング業務で使用するためにノードの合計コア容量 から削除されます。

例 1

ノードは 100 コアの容量を備えています。クラスター管理者は、パフォーマンスプロファイルを使用して、50 コアを分離プールに割り当て、50 コアを予約プールに割り当てます。クラスター管理者は、25 コアを QoS **Guaranteed** Pod に割り当て、25 コアを **BestEffort** または **Burstable** Pod に割り当てます。これは、分離されたプールの容量と一致します。

例 2

ノードは 100 コアの容量を備えています。クラスター管理者は、パフォーマンスプロファイルを使用して、50 コアを分離プールに割り当て、50 コアを予約プールに割り当てます。クラスター管理者は、50 個のコアを QoS **Guaranteed** Pod に割り当て、1 個のコアを **BestEffort** または **Burstable** Pod に割り当てます。これは、分離されたプールの容量を 1 コア超えています。CPU 容量が不十分なため、Podのスケジューリングが失敗します。

使用する正確なパーティショニングパターンは、ハードウェア、ワークロードの特性、予想されるシステム負荷などの多くの要因によって異なります。いくつかのサンプルユースケースは次のとおりです。

- レイテンシーの影響を受けやすいワークロードがネットワークインターフェイスカード (NIC) などの特定のハードウェアを使用する場合は、分離されたプール内の CPU が、このハードウェ アにできるだけ近いことを確認してください。少なくとも、ワークロードを同じ Non-Uniform Memory Access (NUMA) ノードに配置する必要があります。
- 予約済みプールは、すべての割り込みを処理するために使用されます。システムネットワークに依存する場合は、すべての着信パケット割り込みを処理するために、十分なサイズの予約プールを割り当てます。4.6 以降のバージョンでは、ワークロードはオプションで機密としてラベル付けできます。予約済みパーティションと分離パーティションにどの特定の CPU を使用するかを決定するには、詳細な分析と測定が必要です。デバイスやメモリーの NUMA アフィニティーなどの要因が作用しています。選択は、ワークロードアーキテクチャーと特定のユースケースにも依存します。



重要

予約済みの CPU プールと分離された CPU プールは重複してはならず、これらは共に、ワーカーノードの利用可能なすべてのコアに広がる必要があります。

ハウスキーピングタスクとワークロードが相互に干渉しないようにするには、パフォーマンスプロファイルの **spec** セクションで CPU の 2 つのグループを指定します。

- **isolated** アプリケーションコンテナーワークロードの CPU を指定します。これらの CPU で 実行しているワークロードは、最小のレイテンシーとゼロの中断を経験し、たとえば、高いゼロパケット損失帯域幅に到達できます。
- reserved クラスターおよびオペレーティングシステムのハウスキーピング業務用の CPU を 指定します。reserved グループのスレッドは、ビジーであることが多いです。reserved グ ループでレイテンシーの影響を受けやすいアプリケーションを実行しないでください。遅延の 影響を受けやすいアプリケーションは、isolated グループで実行しています。
 - 1. 環境のハードウェアとトポロジーに適したパフォーマンスプロファイルを作成します。
 - 2. infra およびアプリケーションコンテナー用に予約して分離する CPU で、 **reserved** および **isolated** パラメーターを追加します。

apiVersion: performance.openshift.io/v2

kind: PerformanceProfile

metadata:

name: infra-cpus

spec:

cpu:

reserved: "0-4,9" 1

isolated: "5-8" (2)

nodeSelector: 3

node-role.kubernetes.io/worker: ""

- 1 クラスターおよびオペレーティングシステムのハウスキーピングタスクを実行する infra コンテナーの CPU を指定します。
- 🥱 アプリケーションコンテナーがワークロードを実行する CPU を指定します。
- 3 ノードセレクターを指定してパフォーマンスプロファイルを特定のノードに適用します。

15.8. パフォーマンスプロファイルによる低レイテンシーを実現するための ノードのチューニング

パフォーマンスプロファイルを使用すると、特定のマシン設定プールに属するノードのレイテンシーの 調整を制御できます。設定を指定すると、**PerformanceProfile** オブジェクトは実際のノードレベルの チューニングを実行する複数のオブジェクトにコンパイルされます。

- ノードを操作する MachineConfig ファイル。
- Topology Manager、CPU マネージャー、および OpenShift Container Platform ノードを設定する **KubeletConfig** ファイル。
- Node Tuning Operator を設定する Tuned プロファイル。

手順

- 1. クラスターを準備します。
- 2. マシン設定プールを作成します。
- 3. Performance Addon Operator をインストールします。
- 4. ハードウェアとトポロジーに適したパフォーマンスプロファイルを作成します。パフォーマンスプロファイルでは、カーネルを kernel-rt、hugepages の割り当て、オペレーティングシステムのハウスキーピング用に予約される CPU、およびワークロードの実行に使用される CPU に更新するかどうかを指定できます。

これは、典型的なパフォーマンスプロファイルです。

apiVersion: performance.openshift.io/v1 kind: PerformanceProfile metadata: name: performance spec: cpu: isolated: "5-15" reserved: "0-4" hugepages: defaultHugepagesSize: "1G" pages: -size: "1G" count: 16 node: 0 realTimeKernel: enabled: true 1 numa: (2) topologyPolicy: "best-effort" nodeSelector: node-role.kubernetes.io/worker-cnf: ""

有効な値は true または false です。true 値を設定すると、ノード上にリアルタイムカーネルがインストールされます。

Topology Manager ポリシーを設定するには、このフィールドを使用します。有効な値は **none** (デフォルト)、 **best-effort、restricted**、および **single-numa-node** です。詳細は、Topology

15.9. プラットフォーム検証のエンドツーエンドテストの実行

Cloud-native Network Functions (CNF) テストイメージは、CNF ペイロードの実行に必要な機能を検証するコンテナー化されたテストスイートです。このイメージを使用して、CNF ワークロードの実行に必要なすべてのコンポーネントがインストールされている CNF 対応の OpenShift クラスターを検証できます。

イメージで実行されるテストは、3つの異なるフェーズに分かれます。

- 単純なクラスター検証
- セットアップ
- エンドツーエンドテスト

検証フェーズでは、テストに必要なすべての機能がクラスターに正しくデプロイされていることを確認 します。

検証には以下が含まれます。

- テストするマシンに属するマシン設定プールのターゲット設定
- ノードでの SCTP の有効化
- Performance Addon Operator がインストールされていること
- SR-IOV Operator がインストールされていること
- PTP Operator がインストールされていること
- OVN kubernetes の SDN としての使用

テストは、実行されるたびに環境設定を実行する必要があります。これには、SR-IOV ノードポリシー、パフォーマンスプロファイル、または PTP プロファイルの作成などの項目が関係します。テストですでに設定されているクラスターを設定できるようにすると、クラスターの機能が影響を受ける可能性があります。また、SR-IOV ノードポリシーなどの設定項目への変更により、設定の変更が処理されるまで環境が一時的に利用できなくなる可能性があります。

15.9.1. 前提条件

- テストエントリーポイントは /usr/bin/test-run.sh です。設定されたテストセットと実際の適合 テストスイートの両方を実行します。最小要件として、これをボリューム経由でマウントされ る kubeconfig ファイルおよび関連する **\$KUBECONFIG** 環境変数と共に指定します。
- このテストでは、特定の機能が Operator、クラスターで有効にされるフラグ、またはマシン設定の形式でクラスターで利用可能であることを前提としています。
- テストによっては、変更を追加するための既存のマシン設定プールが必要です。これは、テストを実行する前にクラスター上に作成する必要があります。 デフォルトのワーカープールは worker-cnf であり、以下のマニフェストを使用して作成できます。

```
apiVersion: machineconfiguration.openshift.io/v1
kind: MachineConfigPool
metadata:
 name: worker-cnf
 labels:
  machineconfiguration.openshift.io/role: worker-cnf
 machineConfigSelector:
  matchExpressions:
      key: machineconfiguration.openshift.io/role,
      operator: In,
      values: [worker-cnf, worker],
    }
 paused: false
 nodeSelector:
  matchLabels:
   node-role.kubernetes.io/worker-cnf: ""
```

ROLE_WORKER_CNF変数を使用して、ワーカープール名を上書きできます。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig -e ROLE_WORKER_CNF=custom-worker-pool registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh



注記

現時点で、プールに属するノードですべてのテストが選択的に実行される訳では ありません。

15.9.2. テストの実行

ファイルが現在のフォルダーにある場合、テストスイートを実行するコマンドは以下のようになります。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh

これにより、kubeconfig ファイルを実行中のコンテナー内から使用できます。

15.9.3. イメージパラメーター

要件に応じて、テストは異なるイメージを使用できます。以下の環境変数を使用して変更できるテストで、以下の2つのイメージが使用されます。

- CNF TESTS IMAGE
- DPDK TESTS IMAGE

たとえば、カスタムレジストリーを使用して CNF_TESTS_IMAGE を変更するには、以下のコマンドを実行します。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig -e CNF_TESTS_IMAGE="custom-cnf-tests-image:latests" registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh

15.9.3.1. ginkgo パラメーター

テストスイートは、ginkgo BDD フレームワーク上に構築されます。これは、テストをフィルターするか、または省略するためのパラメーターを受け入れることを意味します。

-ginkgo.focus パラメーターを使用してテストセットをフィルターできます。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh -ginkgo.focus="performance|sctp"



注記

特定のテストでは、SR-IOV と SCTP の両方が必要になります。focus パラメーターの選択的な性質を考慮すると、このテストは sriov Matcher のみを配置してトリガーできます。テストが SR-IOV がインストールされているクラスターに対して実行されるものの、SCTP がはない場合、-ginkgo.skip=SCTP パラメーターを追加すると、テストはSCTP テストを省略します。

15.9.3.2. 利用可能な機能

フィルターに使用できる機能のセットは以下になります。

- performance
- sriov
- ptp
- sctp
- dpdk

15.9.4. ドライラン

このコマンドを使用してドライランモードで実行します。これは、テストスイートの内容を確認し、イメージが実行するすべてのテストの出力を提供します。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh -ginkgo.dryRun -ginkgo.v

15.9.5. 非接続モード

CNFは、非接続クラスター、つまり外部レジストリーに到達できないクラスターでテストを実行してイメージのサポートをテストします。これは2つの手順で行います。

- 1. ミラーリングの実行。
- 2. テストに対してカスタムレジストリーからのイメージを使用するように指示します。

15.9.5.1. クラスターからアクセスできるカスタムレジストリーへのイメージのミラーリング

mirror 実行可能ファイルはイメージに含まれ、テストをローカルレジストリーに対して実行するために必要なテストをミラーリングするために oc で必要な入力を提供します。

クラスターおよびインターネット経由で registry.redhat.io にアクセスできる中間マシンから、以下のコマンドを実行します。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/mirror -registry my.local.registry:5000/ | oc image mirror -f -

次に、イメージの取得に使用されるレジストリーを上書きする方法について、以下のセクションの手順 に従います。

15.9.5.2. テストに対してカスタムレジストリーからのそれらのイメージを使用するように指示します。

これは、IMAGE REGISTRY環境変数を設定して実行されます。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig -e IMAGE_REGISTRY="my.local.registry:5000/" -e CNF_TESTS_IMAGE="custom-cnf-tests-image:latests" registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh

15.9.5.3. クラスター内部レジストリーへのミラーリング

OpenShift Container Platform は、クラスター上の標準ワークロードとして実行される組み込まれたコンテナーイメージレジストリーを提供します。

手順

1. レジストリーをルートを使用して公開し、レジストリーへの外部アクセスを取得します。

\$ oc patch configs.imageregistry.operator.openshift.io/cluster --patch '{"spec": {"defaultRoute":true}}' --type=merge

2. レジストリーエンドポイントを取得します。

REGISTRY=\$(oc get route default-route -n openshift-image-registry --template='{{ .spec.host }}')

3. イメージを公開する namespace を作成します。

\$ oc create ns cnftests

4. イメージストリームを、テストに使用されるすべての namespace で利用可能にします。これは、テスト namespace が **cnftests** イメージストリームからイメージを取得できるようにするために必要です。

\$ oc policy add-role-to-user system:image-puller system:serviceaccount:sctptest:default -- namespace=cnftests

\$ oc policy add-role-to-user system:image-puller system:serviceaccount:cnf-features-testing:default --namespace=cnftests

\$ oc policy add-role-to-user system:image-puller system:serviceaccount:performance-addon-operators-testing:default --namespace=cnftests

\$ oc policy add-role-to-user system:image-puller system:serviceaccount:dpdk-testing:default --namespace=cnftests

\$ oc policy add-role-to-user system:image-puller system:serviceaccount:sriov-conformance-testing:default --namespace=cnftests

5. docker シークレット名と認証トークンを取得します。

SECRET=\$(oc -n cnftests get secret | grep builder-docker | awk {'print \$1'}
TOKEN=\$(oc -n cnftests get secret \$SECRET -o jsonpath="{.data['\.dockercfg']}" | base64 -decode | jq '.["image-registry.openshift-image-registry.svc:5000"].auth')

6. 以下のような dockerauth.json を作成します。

echo "{\"auths\": { \"\$REGISTRY\": { \"auth\": \$TOKEN } }}" > dockerauth.json

7. ミラーリングを実行します。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/mirror -registry \$REGISTRY/cnftests oc image mirror --insecure=true -a=\$(pwd)/dockerauth.json -f -

8. テストを実行します。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig -e IMAGE_REGISTRY=image-registry.openshift-image-registry.svc:5000/cnftests cnf-tests-local:latest/usr/bin/test-run.sh

15.9.5.4. イメージの異なるセットのミラーリング

手順

1. **mirror** コマンドは、デフォルトで u/s イメージのミラーリングを試行します。これは、以下の 形式のファイルをイメージに渡すことで上書きできます。

2. これを mirror コマンドに渡します。たとえば、images.json としてローカルに保存できます。 以下のコマンドでは、ローカルパスはコンテナー内の /kubeconfig にマウントされ、これを mirror コマンドに渡すことができます。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/mirror --registry "my.local.registry:5000/" --images "/kubeconfig/images.json" | oc image mirror -f -

15.9.6. 検出モード

検出モードでは、設定を変更せずにクラスターの機能を検証できます。既存の環境設定はテストに使用されます。テストは必要な設定項目の検索を試行し、それらの項目を使用してテストを実行します。特定のテストの実行に必要なリソースが見つからない場合、テストは省略され、ユーザーに適切なメッセージが表示されます。テストが完了すると、事前に設定された設定項目のクリーンアップは行われず、テスト環境は別のテストの実行にすぐに使用できます。

一部の設定項目は、テストで引き続き作成されます。これらの項目は、テストを実行するために必要な特定の項目です (例: SR-IOV ネットワーク)。これらの設定項目はカスタム namespace に作成され、テストの実行後にクリーンアップされます。

さらに、これによりテストの実行時間が削減されます。設定項目はすでに存在しているので、環境設定 および安定化に追加の時間は取る必要はありません。

検出モードを有効にするには、以下のように **DISCOVERY_MODE** 環境変数を設定してテストに対して 指示する必要があります。

\$ docker run -v \$(pwd)/:/kubeconfig:Z -e KUBECONFIG=/kubeconfig/kubeconfig -e DISCOVERY MODE=true registry.redhat.io/openshift-kni/cnf-tests /usr/bin/test-run.sh

15.9.6.1. 必要な環境設定の前提条件

SR-IOV テスト

ほとんどの SR-IOV テストには以下のリソースが必要です。

- SriovNetworkNodePolicy
- **SriovNetworkNodePolicy** で指定されるリソースの1つ以上が割り当て可能な状態であること。リソース数が5以上の場合に十分であると見なされます。

テストによっては、追加の要件があります。

- 利用可能なポリシーリソースがあるノード上の未使用のデバイス。リンク状態が **DOWN** であり、ブリッジスレーブではない。
- MTU 値が 9000 の SriovNetworkNodePolicy。

DPDK テスト

DPDK 関連のテストには、以下が必要です。

- パフォーマンスプロファイル
- SR-IOV ポリシー。

● SRIOV ポリシーで利用可能なリソースが含まれ、**PerformanceProfile** ノードセレクターで利用可能なノード。

PTP テスト

- スレーブ PtpConfig (ptp4lOpts="-s" ,phc2sysOpts="-a -r")。
- スレーブ PtpConfig に一致するラベルの付いたノード。

SCTP テスト

- SriovNetworkNodePolicy
- SriovNetworkNodePolicy および SCTP を有効にする MachineConfig の両方に一致するノード。

Performance Operator のテスト

各種のテストにはそれぞれ異なる要件があります。以下はその一部になります。

- パフォーマンスプロファイル
- profile.Spec.CPU.Isolated = 1 を持つパフォーマンスプロファイル。
- profile.Spec.RealTimeKernel.Enabled == true を持つパーマンスプロファイル。
- Huge Page が使用されていないノード。

15.9.6.2. テスト中に使用されるノードの制限

テストが実行されるノードは、**NODES_SELECTOR** 環境変数を指定して制限できます。テストによって作成されるリソースは、指定されるノードに制限されます。

\$ docker run -v \$(pwd)/:/kubeconfig:Z -e KUBECONFIG=/kubeconfig/kubeconfig -e NODES_SELECTOR=node-role.kubernetes.io/worker-cnf registry.redhat.io/openshift-kni/cnf-tests/usr/bin/test-run.sh

15.9.6.3. 単一のパフォーマンスプロファイルの使用

DPDKテストで必要なリソースは、パフォーマンスのテストスイートで必要なリソースのレベルよりも高くなります。実行をより高速にするために、テストで使用されるパフォーマンスプロファイルは、DPDKテストスイートも提供するプロファイルを使用して上書きできます。

これを実行するには、コンテナー内にマウントできる以下のようなプロファイルを使用し、パフォーマンステストに対してこれをデプロイするように指示できます。

apiVersion: performance.openshift.io/v1

kind: PerformanceProfile

metadata:

name: performance

spec: cpu:

isolated: "4-15" reserved: "0-3" hugepages:

defaultHugepagesSize: "1G"
pages:
- size: "1G"
count: 16
node: 0
realTimeKernel:
enabled: true
nodeSelector:
node-role.kubernetes.io/worker-cnf: ""

使用されるパフォーマンスプロファイルを上書きするには、マニフェストをコンテナー内にマウントし、**PERFORMANCE_PROFILE_MANIFEST_OVERRIDE** パラメーターを設定してテストに対して指示する必要があります。

\$ docker run -v \$(pwd)/:/kubeconfig:Z -e KUBECONFIG=/kubeconfig/kubeconfig -e PERFORMANCE_PROFILE_MANIFEST_OVERRIDE=/kubeconfig/manifest.yaml registry.redhat.io/openshift-kni/cnf-tests /usr/bin/test-run.sh

15.9.6.4. パフォーマンスプロファイルのクリーンアップの無効化

検出モードで実行されない場合、スイートは作成されるすべてのアーティファクトおよび設定をクリーンアップします。これには、パフォーマンスプロファイルが含まれます。

パフォーマンスプロファイルを削除する際に、マシン設定プールが変更され、ノードが再起動されます。新規の繰り返しの後に、新規プロファイルが作成されます。これにより、長いテストサイクルが実行間に生じます。

このプロセスを迅速化するには、CLEAN_PERFORMANCE_PROFILE="false" を設定し、パフォーマンスプロファイルをクリーンアップしないようにテストに指示します。これにより、次の反復でこれを作成し、これが適用されるのを待機する必要がなくなります。

\$ docker run -v \$(pwd)/:/kubeconfig:Z -e KUBECONFIG=/kubeconfig/kubeconfig -e CLEAN_PERFORMANCE_PROFILE="false" registry.redhat.io/openshift-kni/cnf-tests /usr/bin/test-run.sh

15.9.7. トラブルシューティング

クラスターはコンテナー内からアクセスできる必要があります。これを確認するには、以下を実行します。

\$ docker run -v \$(pwd)/:/kubeconfig -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift-kni/cnf-tests oc get nodes

これが機能しない場合は、原因が DNS、MTU サイズ、またはファイアウォールの問題に関連している可能性があります。

15.9.8. テストレポート

CNF エンドツーエンドテストにより、JUnit テスト出力とテスト失敗レポートという 2 つの出力が生成されます。

15.9.8.1. JUnit テスト出力

レポートがダンプされるパスと共に **--junit** パラメーターを渡すと、JUnit 準拠の XML が作成されます。

\$ docker run -v \$(pwd)/:/kubeconfig -v \$(pwd)/junitdest:/path/to/junit -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh --junit /path/to/junit

15.9.8.2. テスト失敗レポート

クラスターの状態とトラブルシューティング用のリソースに関する情報が含まれるレポートは、レポートがダンプされるパスと共に --report パラメーターを渡すことで生成できます。

\$ docker run -v \$(pwd)/:/kubeconfig -v \$(pwd)/reportdest:/path/to/report -e KUBECONFIG=/kubeconfig/kubeconfig registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh --report /path/to/report

15.9.8.3. podman に関する注記

非 root として、また権限なしで podman を実行する場合、permission denied というエラーを出してパスのマウントに失敗する可能性があります。これを機能させるには、:Z をボリューム作成に追加します。たとえば、-v \$(pwd)/:/kubeconfig:Z を使用して podman が SELinux のラベルを適切に書き換えることができます。

15.9.8.4. OpenShift Container Platform 4.4 での実行

以下を除き、CNF のエンドツーエンドのテストは OpenShift Container Platform 4.4 と互換性があります。

[test_id:28466][crit:high][vendor:cnf-qe@redhat.com][level:acceptance] Should contain configuration injected through openshift-node-performance profile [test_id:28467][crit:high][vendor:cnf-qe@redhat.com][level:acceptance] Should contain configuration injected through the openshift-node-performance profile

これらのテストを省略するには、-ginkgo.skip "28466|28467" パラメーターを追加します。

15.9.8.5. 単一のパフォーマンスプロファイルの使用

DPDK テストには、パフォーマンステストスイートに必要なリソースよりも多くのリソースが必要です。実行をより迅速にするには、DPDK テストスイートを提供するプロファイルを使用して、テストが使用するパフォーマンスプロファイルを上書きすることができます。

これを実行するには、コンテナー内にマウントできる以下のようなプロファイルを使用し、パフォーマンステストに対してこれをデプロイするように指示できます。

apiVersion: performance.openshift.io/v1

kind: PerformanceProfile

metadata:

name: performance

spec: cpu:

isolated: "5-15" reserved: "0-4" hugepages:

defaultHugepagesSize: "1G"

pages: -size: "1G" count: 16 node: 0

realTimeKernel: enabled: true

numa:

topologyPolicy: "best-effort"

nodeSelector:

node-role.kubernetes.io/worker-cnf: ""

パフォーマンスプロファイルを上書きするには、マニフェストをコンテナー内にマウントし、**PERFORMANCE_PROFILE_MANIFEST_OVERRIDE** を設定してテストに対して指示する必要があります。

\$ docker run -v \$(pwd)/:/kubeconfig:Z -e KUBECONFIG=/kubeconfig/kubeconfig -e PERFORMANCE_PROFILE_MANIFEST_OVERRIDE=/kubeconfig/manifest.yaml registry.redhat.io/openshift4/cnf-tests-rhel8:v4.6 /usr/bin/test-run.sh

15.9.9. クラスターへの影響

機能によっては、テストスイートを実行すると、クラスターに異なる影響が及ぶ可能性があります。通常、SCTP テストのみではクラスター設定を変更しません。他のすべての機能は設定にさまざまな影響を与えます。

15.9.9.1. SCTP

SCTP テストは、接続性をチェックするために異なるノードで異なる Pod を実行するのみです。クラスターへの影響は、2 つのノードでの単純な Pod の実行に関連します。

15.9.9.2. SR-IOV

SR-IOV テストには、SR-IOV ネットワーク設定での変更が必要になります。この場合、テストは異なるタイプの設定を作成し、破棄します。

これは、既存の SR-IOV ネットワーク設定がクラスターにインストールされている場合に影響を与える可能性があります。これらの設定の優先順位によっては競合が生じる可能性があるためです。

同時に、テストの結果は既存の設定による影響を受ける可能性があります。

15.9.9.3. PTP

PTP テストは、PTP 設定をクラスターのノードセットに適用します。SR-IOV と同様に、これはすでに 有効な既存の PTP 設定と競合する可能性があります。これにより予測不可能な結果が出される可能性 があります。

15.9.9.4. パフォーマンス

パフォーマンステストにより、パフォーマンスプロファイルがクラスターに適用されます。これにより、ノード設定の変更、CPU の予約、メモリー Huge Page の割り当て、およびカーネルパッケージのリアルタイムの設定が行われます。**performance** という名前の既存のプロファイルがクラスターですでに利用可能な場合、テストはこれをデプロイしません。

15.9.9.5. DPDK

DPDK はパフォーマンスおよび SR-IOV 機能の両方に依存するため、テストスイートはパフォーマンス プロファイルと SR-IOV ネットワークの両方を設定します。そのため、これによる影響は SR-IOV テストおよびパフォーマンステストで説明されているものと同じになります。

15.9.9.6. クリーンアップ

テストスイートの実行後に、未解決のリソースすべてがクリーンアップされます。

15.10. 低レイテンシー CNF チューニングステータスのデバッグ

PerformanceProfile カスタムリソース (CR) には、チューニングのステータスを報告し、レイテンシーのパフォーマンスの低下の問題をデバッグするためのステータスフィールドが含まれます。これらのフィールドは、Operator の調整機能の状態を記述する状態について報告します。

パフォーマンスプロファイルに割り当てられるマシン設定プールのステータスが degraded 状態になる と典型的な問題が発生する可能性があり、これにより **PerformanceProfile** のステータスが低下しま す。この場合、マシン設定プールは失敗メッセージを発行します。

Performance Addon Operator には **performanceProfile.spec.status.Conditions** ステータスフィールドが含まれます。

Status:

Conditions:

Last Heartbeat Time: 2020-06-02T10:01:24Z Last Transition Time: 2020-06-02T10:01:24Z

Status: True Type: Available

Last Heartbeat Time: 2020-06-02T10:01:24Z Last Transition Time: 2020-06-02T10:01:24Z

Status: True

Type: Upgradeable

Last Heartbeat Time: 2020-06-02T10:01:24Z Last Transition Time: 2020-06-02T10:01:24Z

Status: False
Type: Progressing

Last Heartbeat Time: 2020-06-02T10:01:24Z Last Transition Time: 2020-06-02T10:01:24Z

Status: False Type: Degraded

Status フィールドには、 パフォーマンスプロファイルのステータスを示す Type 値を指定する Conditions が含まれます。

Available

すべてのマシン設定および Tuned プロファイルが正常に作成され、クラスターコンポーネントで利用可能になり、それら (NTO、MCO、Kubelet) を処理します。

Upgradeable

Operator によって維持されるリソースは、アップグレードを実行する際に安全な状態にあるかどうかを示します。

Progressing

パフォーマンスプロファイルからのデプロイメントプロセスが開始されたことを示します。

Degraded

以下の場合にエラーを示します。

- パーマンスプロファイルの検証に失敗しました。
- すべての関連するコンポーネントの作成が完了しませんでした。

これらのタイプには、それぞれ以下のフィールドが含まれます。

Status

特定のタイプの状態 (true または false)。

Timestamp

トランザクションのタイムスタンプ。

Reason string

マシンの読み取り可能な理由。

Message string

状態とエラーの詳細を説明する人が判読できる理由(ある場合)。

15.10.1. マシン設定プール

パフォーマンスプロファイルとその作成される製品は、関連付けられたマシン設定プール (MCP) に従ってノードに適用されます。MCP は、カーネル引数、kube 設定、Huge Page の割り当て、およびrt-kernel のデプロイメントを含むパフォーマンスアドオンが作成するマシン設定の適用についての進捗に関する貴重な情報を保持します。パフォーマンスアドオンコントローラーは MCP の変更を監視し、それに応じてパフォーマンスプロファイルのステータスを更新します。

MCP がパフォーマンスプロファイルのステータスに返す状態は、MCP が **Degraded** の場合のみとなり、この場合、**performaceProfile.status.condition.Degraded** = true になります。

例

以下の例は、これに作成された関連付けられたマシン設定プール (worker-cnf) を持つパフォーマンスプロファイルのサンプルです。

1. 関連付けられたマシン設定プールの状態は degraded (低下) になります。

oc get mcp

出力例

NAME CONFIG UPDATED UPDATING DEGRADED MACHINECOUNT READYMACHINECOUNT UPDATEDMACHINECOUNT DEGRADEDMACHINECOUNT AGE master rendered-master-2ee57a93fa6c9181b546ca46e1571d2d True False False 3 2d21h 0 rendered-worker-d6b2bdc07d9f5a59a6b68950acf25e5f worker True False 2d21h worker-cnf rendered-worker-cnf-6c838641b8a08fff08dbd8b02fb63f7c False True True 1 1 2d20h

2. MCP の describe セクションには理由が示されます。

oc describe mcp worker-cnf

出力例

Message: Node node-worker-cnf is reporting: "prepping update: machineconfig.machineconfiguration.openshift.io \"rendered-worker-cnf-40b9996919c08e335f3ff230ce1d170\" not

found"

Reason: 1 nodes are reporting degraded status on sync

3. degraded (低下) の状態は、**degraded = true** とマークされたパフォーマンスプロファイルの **status** フィールドにも表示されるはずです。

oc describe performanceprofiles performance

出力例

Message: Machine config pool worker-cnf Degraded Reason: 1 nodes are reporting degraded status on sync.

Machine config pool worker-cnf Degraded Message: Node yquinn-q8s5v-w-b-z5lqn.c.openshift-gce-devel.internal is

reporting: "prepping update: machineconfig.machineconfiguration.openshift.io \"rendered-worker-cnf-40b9996919c08e335f3ff230ce1d170\" not found". Reason:

MCPDegraded Status: True Type: Degraded

15.11. RED HAT サポート向けの低レイテンシーのチューニングデバッグデータの収集

サポートケースを作成する際、ご使用のクラスターについてのデバッグ情報を Red Hat サポートに提供していただくと Red Hat のサポートに役立ちます。

must-gather ツールを使用すると、ノードのチューニング、NUMA トポロジー、および低レイテンシーの設定に関する問題のデバッグに必要な OpenShift Container Platform クラスターについての診断情報を収集できます。

迅速なサポートを得るには、OpenShift Container Platform と低レイテンシーチューニングの両方の診断情報を提供してください。

15.11.1. must-gather ツールについて

oc adm must-gather CLI コマンドは、以下のような問題のデバッグに必要となる可能性のあるクラスターからの情報を収集します。

- リソース定義
- 監査ログ
- サービスログ

--image 引数を指定してコマンドを実行する際にイメージを指定できます。イメージを指定する際、 ツールはその機能または製品に関連するデータを収集します。oc adm must-gather を実行すると、新

しい Pod がクラスターに作成されます。データは Pod で収集され、must-gather.local で始まる新規 ディレクトリーに保存されます。このディレクトリーは、現行の作業ディレクトリーに作成されます。

15.11.2. 低レイテンシーチューニングデータの収集について

oc adm must-gather CLI コマンドを使用してクラスターについての情報を収集できます。これには、 以下を始めとする低レイテンシーチューニングに関連する機能およびオブジェクトが含まれます。

- Performance Addon Operator namespace および子オブジェクト
- MachineConfigPool および関連付けられた MachineConfig オブジェクト
- Node Tuning Operator および関連付けられた Tuned オブジェクト
- Linux カーネルコマンドラインオプション
- CPU および NUMA トポロジー
- 基本的な PCI デバイス情報と NUMA 局所性

must-gatherを使用して Performance Addon Operator のデバッグ情報を収集するには、Performance Addon Operator の**must-gather**イメージを指定する必要があります。

--image=registry.redhat.io/openshift4/performance-addon-operator-must-gather-rhel8:v4.6.

15.11.3. 特定の機能に関するデータ収集

oc adm must-gather CLI コマンドを --image または --image-stream 引数と共に使用して、特定に機 能についてのデバッグ情報を収集できます。must-gather ツールは複数のイメージをサポートするた め、単一のコマンドを実行して複数の機能についてのデータを収集できます。



注記

特定の機能データに加えてデフォルトの must-gather データを収集するには、--imagestream=openshift/must-gather 引数を追加します。

前提条件

- cluster-admin ロールを持つユーザーとしてクラスターにアクセスできる。
- OpenShift Container Platform CLI (oc) がインストールされている。

手順

- 1. must-gather データを保存するディレクトリーに移動します。
- 2. oc adm must-gather コマンドを1つまたは複数の --image または --image-stream 引数と共に 実行します。たとえば、以下のコマンドは、デフォルトのクラスターデータと Performance Addon Operator に固有の情報の両方を収集します。

\$ oc adm must-gather \

--image-stream=openshift/must-gather \ 1



--image=registry.redhat.io/openshift4/performance-addon-operator-must-gather-rhel8:v4.6



- ______デフォルトの OpenShift Container Platform **must-gather** イメージ。
- **仮**レイテンシーチューニングの診断用の must-gather イメージ。
- 3. 作業ディレクトリーに作成された must-gather ディレクトリーから圧縮ファイルを作成します。たとえば、Linux オペレーティングシステムを使用するコンピューターで以下のコマンドを実行します。
 - \$ tar cvaf must-gather.tar.gz must-gather.local.5421342344627712289/
 - must-gather-local.5421342344627712289/ を実際のディレクトリー名に置き換えます。
- 4. 圧縮ファイルを Red Hat カスタマーポータル で作成したサポートケースに添付します。

関連情報

- MachineConfig および KubeletConfig についての詳細は、ノードの管理 を参照してください。
- Node Tuning Operator についての詳細は、Node Tuning Operator の使用を参照してください。
- PerformanceProfile についての詳細は、Huge Page の設定 を参照してください。
- コンテナーからの Huge Page の消費に関する詳細は、Huge Page がアプリケーションによって消費される仕組み を参照してください。

第16章 INTEL FPGA PAC N3000 および INTEL VRAN DEDICATED ACCELERATOR ACC100 でのデータプレーンのパフォーマンスの最適化

16.1. OPENSHIFT CONTAINER PLATFORM 向け INTEL ハードウェアアクセラレーターカードについて

Intel 高速 4G/LTE および 5G Virtualized Radio Access Networks (vRAN) ワークロードからのハードウェアアクセラレーターカード。これにより、市販の既製のプラットフォームの全体的な計算能力が向上します。

Intel FPGA PAC N3000

Intel FPGA PAC N3000 はリファレンス FPGA であり、5G または 4G/LTE RAN レイヤー 1(L1) 基地局ネットワーク機能を高速化するワークロードの例として 4G/LTE または 5G フォワードエラー補正 (FEC) を使用します。vRAN ワークロードをサポートするために、4G/LTE または 5G ビットストリームを持つ Intel FPGA PAC N3000 カードをフラッシュします。

Intel FPGA PAC N3000 は、マルチワークロードネットワーキングアプリケーションアクセラレーション向け全二重、100 Gbps インシステムの再プログラム可能なアクセラレーションカードです。

Intel FPGA PAC N3000 が 4G/LTE または 5G ビットストリームでプログラム化されると、vRAN ワークロードの FEC を加速するのに使用する SR-IOV (Single Root I/O Virtualization) の仮想関数 (VF) デバイスを公開します。クラウドネイティブなデプロイメントのためにこの機能を活用するには、複数の VF を作成するために、デバイスの物理関数 (PF) を **pf-pci-stub** ドライバーにバインドする必要があります。VF の作成後、VF は、vRAN ワークロードを実行する特定の Pod に割り当てるために DPDK ユーザー空間ドライバー (vfio) にバインドする必要があります。

OpenShift Container Platform での Intel FPGA PAC N3000 サポートは 2 つの Operator に依存します。

- Intel FPGA PAC N3000 (プログラミング) 向け OpenNESS Operator
- ワイヤレス FEC Accelerator 向け OpenNESS Operator

vRAN Dedicated Accelerator ACC100

Intel の eASIC テクノロジーに基づく vRAN Dedicated Accelerator ACC100 は、4G/LTE および 5G テクノロジーに対する前方誤り訂正のコンピューター集約型プロセスをオフロードし、処理能力を解放するように設計されています。Intel eASIC デバイスは設定された ASIC で、FPGA と標準のアプリケーション固有の統合サーキット (ASIC) 間の中間テクノロジーです。

OpenShift Container Platform での Intel vRAN Dedicated Accelerator ACC100 サポートは1つの Operator を使用します。

ワイヤレス FEC Accelerator 向け OpenNESS Operator

16.2. INTEL FPGA PAC N3000 向け OPENNESS OPERATOR のインストール

Intel FPGA PAC N3000 向け OpenNESS Operator は、OpenShift Container Platform クラスター内の Intel FPGA PAC N3000 カードによって公開されるリソースまたはデバイスをオーケストレーションし、管理します。

vRAN ユースケースの場合、Intel FPGA PAC N3000 向け OpenNESS Operator はワイヤレス FEC Accelerator の OpenNESS Operator で使用されます。

クラスター管理者は、OpenShift Container Platform CLI または Web コンソールを使用して Intel FPGA PAC N3000 の OpenNESS Operator をインストールできます。

16.2.1. CLI を使用した **Operator** のインストール

クラスター管理者は、CLI を使用して Operator をインストールできます。

前提条件

- ▼アメタルハードウェアにインストールされたクラスター。
- OpenShift CLI (oc) をインストールしている。
- cluster-admin 権限を持つユーザーとしてログインしている。

手順

- 1. 以下のアクションを実行して、N3000 Operator の namespace を作成します。
 - a. 以下の例のように n3000-namespace.yaml ファイルを作成して、vran-acceleration-operators namespace を定義します。

apiVersion: v1 kind: Namespace metadata:

name: vran-acceleration-operators

labels:

openshift.io/cluster-monitoring: "true"

b. 以下のコマンドを実行して namespace を作成します。

\$ oc create -f n3000-namespace.yaml

- 2. 直前の手順で作成した namespace に N3000 Operator をインストールします。
 - a. 以下の **OperatorGroup** CR を作成し、YAML を **n3000-operatorgroup.yaml** ファイルに保存します。

apiVersion: operators.coreos.com/v1

kind: OperatorGroup

metadata:

name: n3000-operators

namespace: vran-acceleration-operators

spec:

targetNamespaces:

- vran-acceleration-operators
- b. 以下のコマンドを実行して Operator Group CR を作成します。

\$ oc create -f n3000-operatorgroup.yaml

c. 以下のコマンドを実行して、次の手順に必要な channel の値を取得します。

\$ oc get packagemanifest n3000 -n openshift-marketplace -o jsonpath='{.status.defaultChannel}'

出力例

stable

d. 以下の Subscription CR を作成し、YAML を n3000-sub.yaml ファイルに保存します。

apiVersion: operators.coreos.com/v1alpha1

kind: Subscription

metadata:

name: n3000-subscription

namespace: vran-acceleration-operators

spec:

channel: "<channel>"

1

name: n3000

source: certified-operators 2

sourceNamespace: openshift-marketplace

- **1 .status.defaultChannel** パラメーターの直前の手順で取得した値からチャネルの値を 指定します。
- 👱 certified-operators 値を指定する必要があります。
- e. 以下のコマンドを実行して Subscription CR を作成します。

\$ oc create -f n3000-sub.yaml

検証

● Operator がインストールされていることを確認します。

\$ oc get csv

出力例

NAME DISPLAY VERSION REPLACES PHASE n3000.v1.1.0 OpenNESS Operator for Intel® FPGA PAC N3000 1.1.0 Succeeded

Operator が正常にインストールされました。

16.2.2. Web コンソールを使用した Intel FPGA PAC N3000 Operator 向け OpenNESS Operator のインストール

クラスター管理者は、Web コンソールを使用して Intel FPGA PAC N3000 の OpenNESS Operator をインストールできます。



注記

先のセクションで説明されているように **Namespace** および **OperatorGroup** CR を作成する必要があります。

手順

- 1. OpenShift Container Platform Web コンソールを使用して、Intel FPGA PAC N3000 向け OpenNESS Operator をインストールします。
 - a. OpenShift Container Platform Web コンソールで、**Operators** → **OperatorHub** をクリックします。
 - b. 利用可能な Operator の一覧から Intel FPGA PAC N3000 向け OpenNESS Operatorを選択し、Install をクリックします。
 - c. Install Operator ページで、All namespaces on the cluster ϵ 選択します。次に、Install ϵ クリックします。
- 2. オプション: N3000 Operator が正常にインストールされていることを確認します。
 - a. Operators → Installed Operators ページに切り替えます。
 - b. Intel FPGA PAC N3000 向け OpenNESS Operatorが Status が InstallSucceeded の vran-acceleration-operators プロジェクトに一覧表示されていることを確認します。



注記

インストール時に、Operator は **Failed** ステータスを表示する可能性があります。インストールが後に **InstallSucceeded** メッセージを出して正常に実行される場合は、**Failed** メッセージを無視できます。

コンソールが Operator がインストールされていることを示していない場合は、以下のトラブルシューティング手順を実行します。

- Operators → Installed Operators ページに移動し、Operator Subscriptions および Install Plans タブで Status にエラーがあるかどうかを検査します。
- Workloads → Pods ページに移動し、vran-acceleration-operators プロジェクトで Pod のログを確認します。

16.3. INTEL FPGA PAC N3000 向けの OPENNESS OPERATOR のプログラミング

Intel FPGA PAC N3000 が vRAN 5G ビットストリームでプログラムされると、ハードウェアは vRAN 5G ビットストリームで Intel FPGA PAC N3000 を公開します。このビットストリームは、vRAN ワークロードの FEC を加速するために使用される SR-IOV (Single Root I/O Virtualization) 仮想機能 (VF) デバイスを公開します。

クラスター管理者は、OpenShift Container Platform CLI または Web コンソールを使用して Intel FPGA PAC N3000 の OpenNESS Operator をインストールできます。

16.3.1. vRAN ビットストリームを持つ N3000 のプログラミング

クラスター管理者は、vRAN 5G ビットストリームを持つ Intel FPGA PAC N3000 をプログラムできます。このビットストリームは、vRAN ワークロードの前方誤り訂正 (FEC) を加速するために使用される SR-IOV (Single Root I/O Virtualization) 仮想機能 (VF) デバイスを公開します。

前方誤り訂正 (FEC) のロールは、メッセージ内の特定のビットが失われたり、文字化けしている可能性がある転送エラーの修正です。伝送メディアのノイズ、干渉、または信号強度の低下により、メッセージが失われたり文字化けしたりする可能性があります。FEC を使用しないと、文字化けしたメッセージは、ネットワーク負荷に加え、スループットとレイテンシーの両方に影響を与える必要があります。

前提条件

- Intel FPGA PAC N3000 カード
- RT カーネル設定のある Performance Addon Operator
- Intel FPGA PAC N3000 向け OpenNESS Operator でインストールされる1つまたは複数の ノード
- cluster-admin 権限を持つユーザーとしてログインします。



注記

すべてのコマンドは **vran-acceleration-operators** namespace で実行されます。

手順

1. vran-acceleration-operators プロジェクトに切り替えます。

\$ oc project vran-acceleration-operators

2. Pod が実行されていることを確認します。

\$ oc get pods

出力例

| NAME | READY S | STATUS | RESTA | ARTS AG | GE |
|------------------------------|---------------|---------|---------|---------|-----|
| fpga-driver-daemonset-8xz4c | 1/1 | Runnin | g 0 | 15d | |
| fpgainfo-exporter-vhvdq | 1/1 | Running | 1 | 15d | |
| N3000-controller-manager-b68 | 3475c76-gcc6v | 2/2 | Running | 1 | 15d |
| N3000-daemonset-5k55l | 1/1 | Running | , 1 | 15d | |
| N3000-discovery-blmjl | 1/1 | Running | 1 | 15d | |
| N3000-discovery-lblh7 | 1/1 | Running | 1 | 15d | |

以下のセクションでは、インストールされた Pod に関する情報を提供します。

- **fpga-driver-daemonset** は必要な Open Programmable Accelerator Engine (OPAE) ドライバーを提供し、読み込みます。
- fpgainfo-exporter は Prometheus に N3000 Telemetry データを提供します。
- N3000-controller-manager は N3000Node CR をクラスターに適用し、すべてのオペランドコンテナーを管理します。

- N3000-daemonset が主要なワーカーアプリケーションです。これは各ノードの CR の変更を監視し、変更に基づいて動作します。このデーモンに実装されたロジックは、カードの FPGA ユーザーイメージおよび NIC ファームウェアの更新を行います。また、ノードをドレイン (解放) し、更新で必要になる場合に提出します。
- **N3000-discovery** は、デバイスが存在する場合にワーカーノードがインストールされ、ラベルのワーカーノードがある場合に N3000 Accelerator デバイスを検出します。
- 3. Intel FPGA PAC N3000 カードを含むすべてのノードを取得します。

\$ oc get n3000node

出力例

NAME FLASH node1 NotRequested

4. 各ノードのカードに関する情報を取得します。

\$ oc get n3000node node1 -o yaml

出力例

status:

conditions:

- lastTransitionTime: "2020-12-15T17:09:26Z"

message: Inventory up to date

observedGeneration: 1 reason: NotRequested

status: "False" type: Flashed

fortville:

- N3000PCI: 0000:1b:00.0

NICs:

- MAC: 64:4c:36:11:1b:a8

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1a:00.0

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card

N3000 for Networking

- MAC: 64:4c:36:11:1b:a9

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1a:00.1

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card

N3000 for Networking

- MAC: 64:4c:36:11:1b:ac

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1c:00.0

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card

N3000 for Networking

- MAC: 64:4c:36:11:1b:ad

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1c:00.1

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card

N3000 for Networking

fpga:

- PCIAddr: 0000:1b:00.0 1

bitstreamld: "0x23000410010310" 2

bitstreamVersion: 0.2.3 deviceId: "0x0b30"

- ↑ PCIAddr フィールドは、カードの PCI アドレスを表示します。
- **2** bitstreamId フィールドは、現在フラッシュに保存されるビットストリームを示します。
- 5. 現在の bitstreamld、PCIAddr、名前、および deviceld を 0x パディングなしで保存します。

\$ oc get n3000node -o json

- 6. Intel FPGA PAC N3000 カードのユーザービットストリームを更新します。
 - a. 以下の例のように **n3000-cluster.yaml** という名前のファイルを作成し、N3000 クラスターリソースをプログラムに定義します。

apiVersion: fpga.intel.com/v1

kind: N3000Cluster

metadata:

name: n3000 1

namespace: vran-acceleration-operators

spec:

nodes:

- nodeName: "node1" 2

fpga:

- userImageURL: "http://10.10.10.122:8000/pkg/20ww27.5-2x2x25G-5GLDPC-

v1.6.1-3.0.0_unsigned.bin" (3)

PCIAddr: "0000:1b:00.0" 4

checksum: "0b0a87b974d35ea16023ceb57f7d5d9c" 5

- 🚹 名前を指定します。名前は n3000 である必要があります。
- プログラムするノードを指定します。
- 3 ユーザービットストリームの URL を指定します。このビットストリームファイルは、 HTTP または HTTPS サーバーでアクセスできる必要があります。
- 🕜 プログラムするカードの PCI アドレスを指定します。
- **userImageURL** フィールドに指定されるビットストリームの MD5 チェックサムを指定します。

N3000 デーモンは、Open Programmable Acceleration Engine (OPAE) ツールを使用して FPGA ユーザービットストリームを更新し、PCI デバイスをリセットします。FPGA ユーザービットストリームの更新では、カードごとに最大 40 分かかる場合があります。複数の ノードでカードをプログラミングする場合、プログラミングは一度に 1 つのノードで設定されます。

b. 更新を適用して、ビットストリームでカードのプログラミングを開始します。

\$ oc apply -f n3000-cluster.yaml

N3000 デーモンは、適切な 5G FEC ユーザービットストリーム (この例では **20ww27.5-2x2x25G-5GLDPC-v1.6.1-3.0.0_unsigned.bin** などがプロビジョニングされた後、および CR が作成された後に) ビットストリームのプログラミングを開始します。

c. ステータスを確認します。

oc get n3000node

出力例

NAME FLASH node1 InProgress

7. ログを確認します。

a. N3000 デーモンの Pod 名を判別します。

\$ oc get pod -o wide | grep n3000-daemonset | grep node1

出力例

n3000-daemonset-5k55l 1/1 Running 0 15d

b. ログを表示します。

\$ oc logs n3000-daemonset-5k55l

出力例

{"level":"info", "ts":1608054338.8866854, "logger": "daemon.drainhelper.cordonAndDrain()", "msg": "node drained"}

 $\label{logger} $$ {\tt "level":"info","ts":1608054338.9003832,"logger":"daemon.fpgaManager.ProgramFPGAs","msg":"Start program","PCIAddr":"0000:1b:00.0"}$

{"level":"info","ts":1608054338.9004142,"logger":"daemon.fpgaManager.ProgramFPGA"," msg":"Starting","pci":"0000:1b:00.0"}

 $\label{logger} $$ {\tt "level":"info","ts":1608056309.9367146,"logger":"daemon.fpgaManager.ProgramFPGA"," msg":"Program FPGA completed, start new power cycle N3000 ...","pci":"0000:1b:00.0"} {\tt "level":"info","ts":1608056333.3528838,"logger":"daemon.drainhelper.Run()","msg":"work er function - end","performUncordon":true}$

ログファイルは、以下のイベントのフローを示します。

- ビットストリームがダウンロードされ、検証されています。
- ノードはドレイン (解放) され、現時点でワークロードを実行できません。
- フラッシュが開始します。

- o ビットストリームはカードにフラッシュされます。
- o ビットストリームが適用されます。
- フラッシュが完了すると、1つまたは複数のノードの PCI デバイス (1つまたは複数) が 再読み込みされます。ワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator が、新しいフラッシュデバイス (1つまたは複数) を検索できるようになりました。

検証

1. FPGA のユーザービットストリームの更新が完了した後にステータスを確認します。

oc get n3000node

出力例

NAME FLASH node1 Succeeded

2. カードのビットストリーム ID が変更されたことを確認します。

oc get n3000node node1 -o yaml

出力例

status:

conditions:

- lastTransitionTime: "2020-12-15T18:18:53Z"

message: Flashed successfully 1

observedGeneration: 2 reason: Succeeded status: "True" type: Flashed

fortville:

- N3000PCI: 0000:1b:00.0

NICs:

- MAC: 64:4c:36:11:1b:a8

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1a:00.0

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card

N3000 for Networking

- MAC: 64:4c:36:11:1b:a9

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1a:00.1

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card

N3000 for Networking

- MAC: 64:4c:36:11:1b:ac

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1c:00.0

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card

N3000 for Networking

- MAC: 64:4c:36:11:1b:ad

NVMVersion: 7.00 0x800052b0 0.0.0

PCIAddr: 0000:1c:00.1

name: Ethernet Controller XXV710 Intel(R) FPGA Programmable Acceleration Card N3000 for Networking

fpga:

- PCIAddr: 0000:1b:00.0 2

bitstreamld: "0x2315842A010601" 3

bitstreamVersion: 0.2.3 deviceId: "0x0b30" 4

- 🚹 message フィールドは、デバイスが正常にフラッシュされたことを示します。
- 🥠 PCIAddr フィールドは、カードの PCI アドレスを表示します。
- **3 bitstreamId** フィールドは、更新されたビットストリーム ID を示します。
- 4 **deviceID** フィールドは、システムに公開されるカード内のビットストリームのデバイス ID を示します。
- 3. ノード上の FEC PCI デバイスを確認します。
 - a. ノードの設定が正しく適用されていることを確認します。

\$ oc debug node/node1

予想される出力

Starting pod/<node-name>-debug ...
To use host binaries, run `chroot /host`

Pod IP: <ip-address>

If you don't see a command prompt, try pressing enter.

sh-4.4#

b. ノードのファイルシステムを使用できることを確認します。

sh-4.4# chroot /host

予想される出力

sh-4.4#

c. システムのアクセラレーターに関連付けられた PCI デバイスを一覧表示します。

\$ Ispci | grep accelerators

予想される出力

1b:00.0 Processing accelerators: Intel Corporation Device 0b30

1d:00.0 Processing accelerators: Intel Corporation Device 0d8f (rev 01)

FPGA に属するデバイスは出力で報告されます。デバイス ID **0b30** はカードのプログラム に使用される RSU インターフェイスです。**0d8f** は、新規プログラムの 5G デバイスの物理 機能です。

16.4. ワイヤレス FEC ACCELERATOR 向けの OPENNESS SR-IOV OPERATOR のインストール

ワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator のロールは、OpenShift Container Platform クラスター内のさまざまな Intel vRAN FEC アクセラレーションハードウェアにより公開されるデバイスをオーケストレーションし、管理することです。

最もコンピューター集約型な 4G/LTE および 5G ワークロードの 1 つは、RAN レイヤー 1(L1) 前方誤り訂正 (FEC) になります。FEC は、信頼性が低い通信チャネルまたはノイズー通信チャネルでデータ転送エラーを解決します。FEC テクノロジーは、再送信を必要とせずに 4G/LTE または 5G データ内のエラーの一部を検出して修正します。

FEC デバイスは、vRAN ユースケースとして Intel FPGA PAC N3000 および Intel vRAN Dedicated Accelerator ACC100 により提供されます。



注記

Intel FPGA PAC N3000 FPGA には、4G/LTE または 5G ビットストリームを持つフラッシュが必要です。

ワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator は、FEC デバイスの仮想機能 (VF) を作成し、それらを適切なドライバーにバインドし、4G/LTE または 5G デプロイメントの機能について VF キューを設定します。

クラスター管理者は、OpenShift Container Platform CLI または Web コンソールを使用して、ワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator をインストールできます。

16.4.1. CLI の使用によるワイヤレス **FEC Accelerator** の **OpenNESS SR-IOV Operator** のインストール

クラスター管理者は、CLI を使用してワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator をインストールできます。

前提条件

- ベアメタルハードウェアにインストールされたクラスター。
- OpenShift CLI (oc) をインストールしている。
- cluster-admin 権限を持つユーザーとしてログインしている。

手順

- 1. 以下のアクションを実行して、ワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator の namespace を作成します。
 - a. 以下の例のように **sriov-namespace.yaml** という名前のファイルを作成し、**vran-acceleration-operators** namespace を定義します。

apiVersion: v1

kind: Namespace

metadata:

name: vran-acceleration-operators

labels:

openshift.io/cluster-monitoring: "true"

b. 以下のコマンドを実行して namespace を作成します。

\$ oc create -f sriov-namespace.yaml

- 2. 以下のオブジェクトを作成して、ワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator を直前の手順で作成した namespace にインストールします。
 - a. 以下の **OperatorGroup** CR を作成し、YAML を **sriov-operatorgroup.yaml** ファイルに保存します。

apiVersion: operators.coreos.com/v1

kind: OperatorGroup

metadata:

name: vran-operators

namespace: vran-acceleration-operators

spec:

targetNamespaces:

vran-acceleration-operators

b. 以下のコマンドを実行して Operator Group CR を作成します。

\$ oc create -f sriov-operatorgroup.yaml

c. 以下のコマンドを実行して、次の手順に必要な channel の値を取得します。

\$ oc get packagemanifest sriov-fec -n openshift-marketplace -o jsonpath='{.status.defaultChannel}'

出力例

stable

d. 以下の Subscription CR を作成し、YAML を **sriov-sub.yaml** ファイルに保存します。

apiVersion: operators.coreos.com/v1alpha1

kind: Subscription

metadata:

name: sriov-fec-subscription

namespace: vran-acceleration-operators

spec:

channel: "<channel>"

name: sriov-fec

source: certified-operators 2

sourceNamespace: openshift-marketplace

1 .status.defaultChannel パラメーターの直前の手順で取得した値からチャネルの値を 指定します。

- 2
- certified-operators 値を指定する必要があります。
- e. 以下のコマンドを実行して Subscription CR を作成します。

\$ oc create -f sriov-sub.yaml

検証

• Operator がインストールされていることを確認します。

\$ oc get csv -n vran-acceleration-operators -o custom-columns=Name:.metadata.name,Phase:.status.phase

出力例

Name sriov-fec.v1.1.0

Phase Succeeded

16.4.2. Web コンソールを使用したワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator のインストール

クラスター管理者は、Web コンソールを使用してワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator をインストールできます。



注記

先のセクションで説明されているように **Namespace** および **OperatorGroup** CR を作成する必要があります。

手順

- 1. OpenShift Container Platform Web コンソールを使用して、ワイヤレス FEC Accelerator の OpenNESS SR-IOV Operator をインストールします。
 - a. OpenShift Container Platform Web コンソールで、**Operators** → **OperatorHub** をクリックします。
 - b. 利用可能な Operator の一覧から **ワイヤレス FEC Accelerator 向け OpenNESS SR-IOV Operator** を選択し、**Install** をクリックします。
 - c. Install Operator ページで、All namespaces on the clusterを選択します。次に、Install を クリックします。
- 2. オプション: SRIOV-FEC Operator が正常にインストールされていることを確認します。
 - a. Operators → Installed Operators ページに切り替えます。
 - b. **ワイヤレス FEC Accelerator 向け OpenNESS SR-IOV Operator**が **Status** が **InstallSucceeded** の **vran-acceleration-operators** プロジェクトに一覧表示されていることを確認します。



注記

インストール時に、 Operator は **Failed** ステータスを表示する可能性があります。インストールが後に **InstallSucceeded** メッセージを出して正常に実行される場合は、**Failed** メッセージを無視できます。

コンソールが Operator がインストールされていることを示していない場合は、以下のトラブルシューティング手順を実行します。

- Operators → Installed Operators ページに移動し、Operator Subscriptions および Install Plans タブで Status にエラーがあるかどうかを検査します。
- Workloads → Pods ページに移動し、vran-acceleration-operators プロジェクトで Pod のログを確認します。

16.4.3. Intel FPGA PAC N3000 向け SR-IOV-FEC Operator の設定

このセクションでは、Intel FPGA PAC N3000 向け SR-IOV-FEC Operator をプログラムする方法を説明します。SR-IOV-FEC Operator は、vRAN L1 アプリケーションの FEC プロセスを加速するために使用される前方誤り訂正 (FEC) デバイスの管理を処理します。

SR-IOV-FEC Operator を設定するには、以下のことを行う必要があります。

- FEC デバイスに必要な仮想機能 (VF) の作成
- VF を適切なドライバーにバインド
- 4G または 5G デプロイメントで希望する機能向け VF キューの設定

前方誤り訂正 (FEC) のロールは、メッセージ内の特定のビットが失われたり、文字化けしている可能性がある転送エラーの修正です。伝送メディアのノイズ、干渉、または信号強度の低下により、メッセージが失われたり文字化けしたりする可能性があります。FEC を使用しないと、文字化けしたメッセージは、ネットワーク負荷に加え、スループットとレイテンシーに影響を与える必要があります。

前提条件

- Intel FPGA PAC N3000 カード
- Intel FPGA PAC N3000 (プログラミング) 向け OpenNESS Operator でインストールされる1 つまたは複数のノード
- ワイヤレス FEC Accelerator 向け OpenNESS Operator でインストールされる 1 つまたは複数のノード
- Performance Addon Operator で設定された RT カーネル

手順

- 1. vran-acceleration-operators プロジェクトに切り替えます。
 - \$ oc project vran-acceleration-operators
- 2. SR-IOV-FEC Operator がインストールされていることを確認します。

\$ oc get csv -o custom-columns=Name:.metadata.name,Phase:.status.phase

出力例

| Name | Phase |
|------------------|-----------|
| sriov-fec.v1.1.0 | Succeeded |
| n3000.v1.1.0 | Succeeded |

3. N3000 および sriov-fec Pod が実行していることを確認します。

\$ oc get pods

出力例

| NAME | READY | STATUS | RESTAI | RTS A | GE |
|----------------------------------|-------------|---------|---------|-------|-----|
| fpga-driver-daemonset-8xz4c | 1/1 | Runni | ng 0 | 15d | |
| fpgainfo-exporter-vhvdq | 1/1 | Running | 1 | 15d | |
| N3000-controller-manager-b6847 | 75c76-gcc6v | 2/2 | Running | 1 | 15d |
| N3000-daemonset-5k55l | 1/1 | Runnin | g 1 | 15d | |
| N3000-discovery-blmjl | 1/1 | Running | 1 | 15d | |
| N3000-discovery-lblh7 | 1/1 | Running | 1 | 15d | |
| sriov-device-plugin-j5jlv | 1/1 | Running | 1 1 | 15d | |
| sriov-fec-controller-manager-85b | 6b8f4d4-gd2 | qg 1/1 | Running | 1 | 15d |
| sriov-fec-daemonset-kqqs6 | 1/1 | Runnin | g 1 | 15d | |

以下のセクションでは、インストールされた Pod に関する情報を提供します。

- **fpga-driver-daemonset** は必要な Open Programmable Accelerator Engine (OPAE) ドライバーを提供し、読み込みます。
- fpgainfo-exporter は Prometheus に N3000 Telemetry データを提供します。
- **N3000-controller-manager** は N3000Node CR をクラスターに適用し、すべてのオペランドコンテナーを管理します。
- N3000-daemonset が主要なワーカーアプリケーションです。
- **N3000-discovery** は、デバイスが存在する場合にワーカーノードがインストールされ、ラベルのワーカーノードがある場合に N3000 Accelerator デバイスを検出します。
- sriov-device-plugin は、FEC 仮想機能をノードの下にあるリソースとして公開します。
- **sriov-fec-controller-manager** は CR をノードに適用し、オペランドコンテナーを維持します。
- sriov-fec-daemonset は、次のことを行います。
 - 各ノードで SRIOV NIC の検出
 - ステップ6で定義されたカスタムリソース (CR) のステータスの同期
 - CRの spec を入力として実行し、検出された NIC の設定
- 4. サポート対象の vRAN FEC アクセラレーターデバイスのいずれかを含むすべてのノードを取得します。

\$ oc get sriovfecnodeconfig

出力例

NAME CONFIGURED node1 Succeeded

5. 設定する SR-IOV FEC アクセラレーターデバイスの Physical Function (PF) を検索します。

\$ oc get sriovfecnodeconfig node1 -o yaml

出力例

status:

conditions:

- lastTransitionTime: "2021-03-19T17:19:37Z"

message: Configured successfully

observedGeneration: 1

reason: ConfigurationSucceeded

status: "True" type: Configured

inventory:

sriovAccelerators:

- deviceID: 0d5c

driver: ""

maxVirtualFunctions: 16 pciAddress: 0000.1d.00.0 1

vendorID: "8086" virtualFunctions: [] 2

- 🚹 このフィールドは、カードの PCI Address を示します。
- 2 このフィールドは、仮想関数が空であることを示します。
- 6. 希望の設定で FEC デバイスを設定します。
 - a. 以下のカスタムリソース (CR) を作成し、YAML を **sriovfec_n3000_cr.yaml** ファイルに保存します。

apiVersion: sriovfec.intel.com/v1 kind: SriovFecClusterConfig metadata:

name: config

namespace: vran-acceleration-operators

spec: nodes:

nodeName: node1 1
 physicalFunctions:

- pciAddress: 0000:1d:00.0 2

pfDriver: pci-pf-stub vfDriver: vfio-pci vfAmount: 2 3 bbDevConfig: n3000:

Network Type: either "FPGA_5GNR" or "FPGA_LTE"

```
networkType: "FPGA_5GNR"
pfMode: false
flrTimeout: 610
downlink:
 bandwidth: 3
 loadBalance: 128
 queues: 4
  vf0: 16
  vf1: 16
  vf2: 0
  vf3: 0
  vf4: 0
  vf5: 0
  vf6: 0
  vf7: 0
uplink:
 bandwidth: 3
 loadBalance: 128
 queues: 5
  vf0: 16
  vf1: 16
  vf2: 0
  vf3: 0
  vf4: 0
  vf5: 0
  vf6: 0
  vf7: 0
```

- 1 ノード名を指定します。
- 2 SR-IOV-FEC Operator がインストールされているカードの PCI Address を指定します。
- 仮想関数の数を指定します。2つの仮想関数を作成します。
- vf0 で、16 バス (ダウンリンクおよびアップリンク) を使用してキューを1つ作成します。
- **5 vf1** で、16 バス (ダウンリンクおよびアップリンク) を使用してキューを1つ作成します。



注記

vRAN Acceleration 向け Intel PAC N3000 では、ユーザーが最大 8 個の VF デバイスを作成できます。各 FEC PF デバイスでは、設定可能な合計 64 個のキュー (アップリンク用に 32 個と、ダウンリンク用に 32 個) が提供されます。キューは通常、VF 全体に均等に分散されます。

b. CR を適用します。

\$ oc apply -f sriovfec_n3000_cr.yaml

CRの適用後、SR-IOV FEC デーモンは FEC デバイスの設定を開始します。

検証

1. ステータスを確認します。

\$ oc get sriovfecclusterconfig config -o yaml

出力例

status:

conditions:

- lastTransitionTime: "2020-12-15T17:19:37Z"

message: Configured successfully

observedGeneration: 1

reason: ConfigurationSucceeded

status: "True" type: Configured

inventory:

sriovAccelerators:
- deviceID: 0d8f
driver: pci-pf-stub

maxVirtualFunctions: 8 pciAddress: 0000:1d:00.0

vendorID: "8086" virtualFunctions: - deviceID: 0d90 driver: vfio-pci

pciAddress: 0000:1d:00.1

 deviceID: 0d90 driver: vfio-pci

pciAddress: 0000:1d:00.2

- 2. ログを確認します。
 - a. SR-IOV デーモン Pod の名前を確認します。

\$ oc get pod | grep sriov-fec-daemonset

出力例

sriov-fec-daemonset-kqqs6 1/1 Running 0 19h

b. ログを表示します。

\$ oc logs sriov-fec-daemonset-kqqs6

出力例

2020-12-16T12:46:47.720Z INFO daemon.NodeConfigurator.applyConfig configuring PF {"requestedConfig": {"pciAddress":"0000:1d:00.0","pfDriver":"pci-pf-stub","vfDriver":"vfio-pci","vfAmount":2,"bbDevConfig":{"n3000":{
"networkType":"FPGA_5GNR","pfMode":false,"flrTimeout":610,"downlink":
{"bandwidth":3,"loadBalance":128,"queues":{"vf0":16,"vf1":16}},"uplink":
{"bandwidth":3,"loadBalance":128,"queues":{"vf0":16,"vf1":16}}}}}}
2020-12-16T12:46:47.720Z INFO daemon.NodeConfigurator.loadModule

```
{"cmd": "/usr/sbin/chroot /host/ modprobe pci-pf-stub"}
executing command
2020-12-16T12:46:47.724Z
                              INFO daemon.NodeConfigurator.loadModule
commands output {"output": ""}
2020-12-16T12:46:47.724Z
                              INFO daemon.NodeConfigurator.loadModule
executing command {"cmd": "/usr/sbin/chroot /host/ modprobe vfio-pci"}
                              INFO daemon.NodeConfigurator.loadModule
2020-12-16T12:46:47.727Z
commands output {"output": ""}
                              INFO daemon.NodeConfigurator device's
2020-12-16T12:46:47.727Z
driver_override path {"path": "/sys/bus/pci/devices/0000:1d:00.0/driver_override"}
2020-12-16T12:46:47.727Z
                              INFO daemon.NodeConfigurator driver bind path
{"path": "/sys/bus/pci/drivers/pci-pf-stub/bind"}
2020-12-16T12:46:47.998Z
                             INFO daemon.NodeConfigurator device's
driver_override path {"path": "/sys/bus/pci/devices/0000:1d:00.1/driver_override"}
                              INFO daemon.NodeConfigurator driver bind path
2020-12-16T12:46:47.998Z
{"path": "/sys/bus/pci/drivers/vfio-pci/bind"}
2020-12-16T12:46:47.998Z
                              INFO daemon.NodeConfigurator device's
driver_override path {"path": "/sys/bus/pci/devices/0000:1d:00.2/driver_override"}
2020-12-16T12:46:47.998Z
                              INFO daemon.NodeConfigurator driver bind path
{"path": "/sys/bus/pci/drivers/vfio-pci/bind"}
                              INFO daemon.NodeConfigurator.applyConfig
2020-12-16T12:46:47.999Z
                      {"cmd": "/sriov workdir/pf bb config FPGA 5GNR -c
executing command
/sriov artifacts/0000:1d:00.0.ini -p 0000:1d:00.0"}
2020-12-16T12:46:48.017Z
                             INFO daemon.NodeConfigurator.applyConfig
commands output {"output": "ERROR: Section (FLR) or name (flr time out) is not valid.
FEC FPGA RTL v3.0
UL.DL Weights = 3.3
UL.DL Load Balance = 1
28.128
Queue-PF/VF Mapping Table = READY
Ring Descriptor Size = 256 bytes
    | PF | VF0 | VF1 | VF2 | VF3 | VF4 | VF5 | VF6 | VF7 |
UL-Q'00 |
UL-Q'01 |
           | X |
UL-Q'02 |
           | X |
          | X |
UL-Q'03 |
UL-Q'04 | X |
UL-Q'05 | X |
UL-Q'06 |
         | X
UL-Q'07 |
           | X
          | X |
UL-Q'08 |
UL-Q'09 | X |
UL-Q'10 | X |
UL-Q'11 |
             Χ
UL-Q'12 |
           | X
UL-Q'13 |
           | X |
UL-Q'14 |
             Χ
UL-Q'15 | X |
UL-Q'16 |
              | X
UL-Q'17 | |
                Χ
UL-Q'18 | |
                Χ |
UL-Q'19 |
                Χ |
UL-Q'20 | |
                Χ |
UL-Q'21 |
                Χ |
```

| UL-Q'22 | ļ | X | | ļ | ļ | | [| [|
|------------------------|------|------|------|---|---|------|------|------|
| UL-Q'23 UL-Q'24 | | | | | | | | |
| UL-Q'25 | į | | įį | į | į | į | į | į |
| UL-Q'26 UL-Q'27 | | | | | | | | |
| UL-Q'28 | İ | Χ | İ | į | į | į | į | į |
| UL-Q'29 UL-Q'30 | | | | | | | | |
| UL-Q'31 | j | | | į | İ | į | İ | į |
| DL-Q'32 DL-Q'33 | X | | | | | | | |
| DL-Q'34 | X | : | | | | | | |
| DL-Q'35 DL-Q'36 | X | | | | | | | |
| DL-Q'37 DL-Q'38 | | | | | | | | |
| DL-Q'38 DL-Q'39 | X | | | | | | | |
| DL-Q'40 DL-Q'41 | X | | | | | | | |
| DL-Q'42 | X | | | | | | | |
| DL-Q'43 DL-Q'44 | X | | | | | | | |
| DL-Q'45 | X | į | | į | | į | | |
| DL-Q'46 DL-Q'47 | | | | | | | | |
| DL-Q'48 | | Χ | | į | į | į | į | İ |
| DL-Q'49 DL-Q'50 | | | | | | | | |
| DL-Q'51 DL-Q'52 | | | | | | | | |
| DL-Q 52 DL-Q'53 | | Χ | | | | | | |
| DL-Q'54 DL-Q'55 | | | | | | | | |
| DL-Q'56 | | Χ | | | | | | |
| DL-Q'57 DL-Q'58 | | X | | | | | | |
| DL-Q'59 | į | Χ | | | | į | | |
| DL-Q'60 DL-Q'61 | | X | | | | | | |
| DL-Q'62 | į | X | | į | į | į | į | į |
| DL-Q'63 | - 1 | Χ | | - | | | | |

Mode of operation = VF-mode

FPGA_5GNR PF [0000:1d:00.0] configuration complete!"}

2020-12-16T12:46:48.017Z INFO daemon.NodeConfigurator.enableMasterBus executing command {"cmd": "/usr/sbin/chroot /host/ setpci -v -s 0000:1d:00.0 COMMAND"}

2020-12-16T12:46:48.037Z INFO daemon.NodeConfigurator.enableMasterBus executing command {"cmd": "/usr/sbin/chroot /host/ setpci -v -s 0000:1d:00.0 COMMAND=0106"}

2020-12-16T12:46:48.054Z INFO daemon.NodeConfigurator.enableMasterBus commands output {"output": "0000:1d:00.0 @04 0106\n"}

2020-12-16T12:46:48.054Z INFO daemon.NodeConfigurator.enableMasterBus MasterBus set {"pci": "0000:1d:00.0", "output": "0000:1d:00.0 @04 0106\n"} 2020-12-16T12:46:48.160Z INFO daemon.drainhelper.Run() worker function end {"performUncordon": true}

3. カードの FEC 設定を確認します。

\$ oc get sriovfecnodeconfig node1 -o yaml

出力例

status:

conditions:

- lastTransitionTime: "2020-12-15T17:19:37Z"

message: Configured successfully

observedGeneration: 1

reason: ConfigurationSucceeded

status: "True" type: Configured

inventory:

sriovAccelerators:
- deviceID: 0d8f 1

driver: pci-pf-stub maxVirtualFunctions: 8 pciAddress: 0000:1d:00.0 vendorID: "8086"

virtualFunctions:
- deviceID: 0d90 2

pciAddress: 0000:1d:00.1

deviceID: 0d90
 driver: vfio-pci

driver: vfio-pci

pciAddress: 0000:1d:00.2

- **1 0d8f** は、FEC デバイスの **deviceID** 物理機能です。
- od90 は、FEC デバイスの deviceID 仮想機能です。

16.4.4. Intel vRAN Dedicated Accelerator ACC100 向け SR-IOV-FEC Operator の設定

Intel vRAN Dedicated Accelerator ACC100 のプログラミングは、vRAN ワークロードで FEC を加速するのに使用される SR-IOV (Single Root I/O Virtualization) 仮想関数 (VF) デバイスを公開します。Intel vRAN Dedicated Accelerator ACC100 は 4G および 5G vRAN (Virtualized Radio Access Networks) ワークロードを加速します。これにより、市販の既製のプラットフォームの全体的な計算能力が向上します。このデバイスは Mount Bryce としても知られています。

SR-IOV-FEC Operator は、vRAN L1 アプリケーションの FEC プロセスを加速するために使用される前方誤り訂正 (FEC) デバイスの管理を処理します。

SR-IOV-FEC Operator を設定するには、以下のことを行う必要があります。

● FEC デバイスの VF (Virtual Function) の作成

- VF を適切なドライバーにバインド
- 4G または5G デプロイメントで希望する機能向け VF キューの設定

前方誤り訂正 (FEC) のロールは、メッセージ内の特定のビットが失われたり、文字化けしている可能性がある転送エラーの修正です。伝送メディアのノイズ、干渉、または信号強度の低下により、メッセージが失われたり文字化けしたりする可能性があります。FEC を使用しないと、文字化けしたメッセージは、ネットワーク負荷に加え、スループットとレイテンシーに影響を与える必要があります。

前提条件

- Intel FPGA ACC100 5G/4G カード
- ワイヤレス FEC Accelerator 向け OpenNESS Operator でインストールされる 1 つまたは複数のノード
- ノードの BIOS でグローバル SR-IOV および VT-d 設定を有効にします。
- Performance Addon Operator で設定された RT カーネル
- cluster-admin 権限を持つユーザーとしてログインします。

手順

1. vran-acceleration-operators プロジェクトに切り替えます。

\$ oc project vran-acceleration-operators

2. SR-IOV-FEC Operator がインストールされていることを確認します。

\$ oc get csv -o custom-columns=Name:.metadata.name,Phase:.status.phase

出力例

Name Phase sriov-fec.v1.1.0 Succeeded

3. **sriov-fec** Pod が実行していることを確認します。

\$ oc get pods

出力例

NAME READY STATUS RESTARTS AGE sriov-device-plugin-j5jlv 1/1 Running 1 15d sriov-fec-controller-manager-85b6b8f4d4-gd2qg 1/1 Running 15d 1 Running sriov-fec-daemonset-kggs6 1/1 1 15d

- sriov-device-plugin は、FEC 仮想機能をノードの下にあるリソースとして公開します。
- **sriov-fec-controller-manager** は CR をノードに適用し、オペランドコンテナーを維持します。
- sriov-fec-daemonset は、次のことを行います。

- o 各ノードで SRIOV NIC の検出
- o ステップ6で定義されたカスタムリソース (CR) のステータスの同期
- CRの spec を入力として実行し、検出された NIC の設定
- 4. サポート対象の vRAN FEC アクセラレーターデバイスのいずれかを含むすべてのノードを取得します。

\$ oc get sriovfecnodeconfig

出力例

NAME CONFIGURED node1 Succeeded

5. 設定する SR-IOV FEC アクセラレーターデバイスの Physical Function (PF) を検索します。

\$ oc get sriovfecnodeconfig node1 -o yaml

出力例

status:

conditions:

- lastTransitionTime: "2021-03-19T17:19:37Z"

message: Configured successfully

observedGeneration: 1

reason: ConfigurationSucceeded

status: "True" type: Configured

inventory:

sriovAccelerators: - deviceID: 0d5c

driver: ""

maxVirtualFunctions: 16 pciAddress: 0000:af:00.0 1

vendorID: "8086" virtualFunctions: [] 2

- ↑ このフィールドは、カードの PCI アドレスを表示します。
- 2 このフィールドは、仮想関数が空であることを示します。
- 6. FEC デバイスの仮想機能とキューグループの数を設定します。
 - a. 以下のカスタムリソース (CR) を作成し、YAML を **sriovfec_acc100cr.yaml** ファイルに保存します。



注記

この例では、5G の ACC100 8/8 キューグループ、Uplink 向け 4 キューグループ、および Downlink の別の 4 つのキューグループを設定します。

-

```
apiVersion: sriovfec.intel.com/v1
kind: SriovFecClusterConfig
metadata:
 name: config 1
spec:
 nodes:
 - nodeName: node1 2
  physicalFunctions:
    - pciAddress: 0000:af:00.0 3
     pfDriver: "pci-pf-stub"
     vfDriver: "vfio-pci"
     vfAmount: 16 4
     bbDevConfig:
      acc100:
       # Programming mode: 0 = VF Programming, 1 = PF Programming
       pfMode: false
       numVfBundles: 16
       maxQueueSize: 1024
       uplink4G:
        numQueueGroups: 0
        numAqsPerGroups: 16
        aqDepthLog2: 4
       downlink4G:
       numQueueGroups: 0
        numAqsPerGroups: 16
        aqDepthLog2: 4
       uplink5G:
       numQueueGroups: 4
        numAqsPerGroups: 16
       aqDepthLog2: 4
       downlink5G:
       numQueueGroups: 4
       numAqsPerGroups: 16
        aqDepthLog2: 4
```

- 🚹 CR オブジェクトの名前を指定します。指定できる唯一の名前は **config** です。
- 2 ノード名を指定します。
- 3 SR-IOV-FEC Operator がインストールされているカードの PCI アドレスを指定します。
- 4 作成する仮想関数の数を指定します。Intel vRAN Dedicated Accelerator ACC100 の場合は、16 個の VF をすべて作成します。



注記

カードは、グループごとに最大 16 個のキューを持つ最大 8 個のキューを提供するように設定されます。キューグループは、5G と 4G と Uplink と Downlink に割り当てられたグループ間で分割できます。Intel vRAN Dedicated Accelerator ACC100 を設定することができます。

- 4G または 5G のみ
- 4G と 5G を同時に

設定した各 VF は、すべてのキューにアクセスできます。各キューグループの優先度は、それぞれ個別の優先レベルを持ちます。特定のキューグループの要求はアプリケーションレベル (FEC デバイスを利用する vRAN アプリケーション) から行われます。

b. CR を適用します。

\$ oc apply -f sriovfec_acc100cr.yaml

CR の適用後、SR-IOV FEC デーモンは FEC デバイスの設定を開始します。

検証

1. ステータスを確認します。

\$ oc get sriovfecclusterconfig config -o yaml

出力例

status:

conditions:

- lastTransitionTime: "2021-03-19T11:46:22Z"

message: Configured successfully

observedGeneration: 1 reason: Succeeded status: "True" type: Configured

inventory:

sriovAccelerators:
- deviceID: 0d5c
driver: pci-pf-stub
maxVirtualFunctions: 16
pciAddress: 0000:af:00.0

vendorID: "8086" virtualFunctions: - deviceID: 0d5d driver: vfio-pci

pciAddress: 0000:b0:00.0

deviceID: 0d5d
 driver: vfio-pci

pciAddress: 0000:b0:00.1

 deviceID: 0d5d driver: vfio-pci

pciAddress: 0000:b0:00.2

deviceID: 0d5d
 driver: vfio-pci

pciAddress: 0000:b0:00.3

deviceID: 0d5d
 driver: vfio-pci

pciAddress: 0000:b0:00.4

2. ログを確認します。

a. SR-IOV デーモンの Pod 名を判別します。

\$ oc get po -o wide | grep sriov-fec-daemonset | grep node1

出力例

sriov-fec-daemonset-kggs6

1/1 Running 0

19h

b. ログを表示します。

\$ oc logs sriov-fec-daemonset-kqqs6

出力例

```
{"level":"Level(-
2)","ts":1616794345.4786215,"logger":"daemon.drainhelper.cordonAndDrain()","msg":"no
de drained"}
{"level":"Level(-
4)","ts":1616794345.4786265,"logger":"daemon.drainhelper.Run()","msg":"worker
function - start"}
{"level":"Level(-
4)","ts":1616794345.5762916,"logger":"daemon.NodeConfigurator.applyConfig","msg":"cur
rent node status", "inventory": {"sriovAccelerat
ors":
[{"vendorID":"8086","deviceID":"0b32","pciAddress":"0000:20:00.0","driver":"","maxVirtualF
unctions":1,"virtualFunctions":[]},{"vendorID":"8086"
,"deviceID":"0d5c","pciAddress":"0000:af:00.0","driver":"","maxVirtualFunctions":16,"virtualF
unctions":[]}]}}
{"level":"Level(-
4)","ts":1616794345.5763638,"logger":"daemon.NodeConfigurator.applyConfig","msg":"co
nfiguring PF", "requestedConfig": {"pciAddress":"
0000:af:00.0", "pfDriver": "pci-pf-stub", "vfDriver": "vfio-pci", "vfAmount": 2, "bbDevConfig":
{"acc100":{"pfMode":false,"numVfBundles":16,"maxQueueSize":1
024,"uplink4G":
{"numQueueGroups":4,"numAqsPerGroups":16,"aqDepthLog2":4},"downlink4G":
{"numQueueGroups":4,"numAqsPerGroups":16,"aqDepthLog2":4},"uplink5G":
{"numQueueGroups":0,"numAqsPerGroups":16,"aqDepthLog2":4},"downlink5G":
{"numQueueGroups":0,"numAqsPerGroups":16,"aqDepthLog2":4}}}}
{"level":"Level(-
4)","ts":1616794345.5774765,"logger":"daemon.NodeConfigurator.loadModule","msg":"ex
ecuting command", "cmd": "/usr/sbin/chroot /host/ modprobe pci-pf-stub"}
{"level":"Level(-
4)","ts":1616794345.5842702,"logger":"daemon.NodeConfigurator.loadModule","msg":"co
mmands output","output":""}
{"level":"Level(-
```

```
4)","ts":1616794345.5843055,"logger":"daemon.NodeConfigurator.loadModule","msg":"ex
ecuting command","cmd":"/usr/sbin/chroot /host/ modprobe vfio-pci"}
{"level":"Level(-
4)","ts":1616794345.6090655,"logger":"daemon.NodeConfigurator.loadModule","msg":"co
mmands output","output":""}
{"level":"Level(-
2)","ts":1616794345.6091156,"logger":"daemon.NodeConfigurator","msg":"device's
driver override path", "path": "/sys/bus/pci/devices/0000:af:00.0/driver override"}
2)","ts":1616794345.6091807,"logger":"daemon.NodeConfigurator","msg":"driver bind
path","path":"/sys/bus/pci/drivers/pci-pf-stub/bind"}
{"level":"Level(-
2)","ts":1616794345.7488534,"logger":"daemon.NodeConfigurator","msg":"device's
driver_override path","path":"/sys/bus/pci/devices/0000:b0:00.0/driver override"}
{"level":"Level(-
2)","ts":1616794345.748938,"logger":"daemon.NodeConfigurator","msg":"driver bind
path", "path": "/sys/bus/pci/drivers/vfio-pci/bind"}
{"level":"Level(-
2)","ts":1616794345.7492096,"logger":"daemon.NodeConfigurator","msg":"device's
driver override path", "path": "/sys/bus/pci/devices/0000:b0:00.1/driver override"}
{"level":"Level(-
2)","ts":1616794345.7492566,"logger":"daemon.NodeConfigurator","msg":"driver bind
path","path":"/sys/bus/pci/drivers/vfio-pci/bind"}
{"level":"Level(-
4)","ts":1616794345.74968,"logger":"daemon.NodeConfigurator.applyConfig","msg":"exec
uting command","cmd":"/sriov_workdir/pf_bb_config ACC100 -c
/sriov_artifacts/0000:af:00.0.ini -p 0000:af:00.0"}
{"level":"Level(-
4)","ts":1616794346.5203931,"logger":"daemon.NodeConfigurator.applyConfig","msg":"co
mmands output", "output": "Queue Groups: 0 5GUL, 0 5GDL, 4 4GUL, 4 4GDL\nNumber
of 5GUL engines 8\nConfiguration in VF mode\nPF ACC100 configuration
complete\nACC100 PF [0000:af:00.0] configuration complete!\n\n"}
{"level":"Level(-
4)","ts":1616794346.520459,"logger":"daemon.NodeConfigurator.enableMasterBus","msg"
:"executing command","cmd":"/usr/sbin/chroot /host/ setpci -v -s 0000:af:00.0
COMMAND"}
{"level":"Level(-
4)","ts":1616794346.5458736,"logger":"daemon.NodeConfigurator.enableMasterBus","ms
g":"commands output","output":"0000:af:00.0 @04 = 0142\n"
{"level":"Level(-
4)","ts":1616794346.5459251,"logger":"daemon.NodeConfigurator.enableMasterBus","ms
g":"executing command","cmd":"/usr/sbin/chroot /host/ setpci -v -s 0000:af:00.0
COMMAND=0146"}
{"level":"Level(-
4)","ts":1616794346.5795262,"logger":"daemon.NodeConfigurator.enableMasterBus","ms
g":"commands output","output":"0000:af:00.0 @04 0146\n"}
{"level":"Level(-
2)","ts":1616794346.5795407,"logger":"daemon.NodeConfigurator.enableMasterBus","ms
g":"MasterBus set","pci":"0000:af:00.0","output":"0000:af:00.0 @04 0146\n"}
{"level":"Level(-
4)","ts":1616794346.6867144,"logger":"daemon.drainhelper.Run()","msg":"worker
function - end", "performUncordon": true}
{"level":"Level(-
4)","ts":1616794346.6867719,"logger":"daemon.drainhelper.Run()","msg":"uncordoning
node"}
{"level":"Level(-
```

```
4)","ts":1616794346.6896322,"logger":"daemon.drainhelper.uncordon()","msg":"starting
uncordon attempts"}
{"level":"Level(-
2)","ts":1616794346.69735,"logger":"daemon.drainhelper.uncordon()","msg":"node
uncordoned"}
{"level":"Level(-
4)","ts":1616794346.6973662,"logger":"daemon.drainhelper.Run()","msg":"cancelling the
context to finish the leadership"}
{"level":"Level(-
4)","ts":1616794346.7029872,"logger":"daemon.drainhelper.Run()","msg":"stopped
leading"}
{"level":"Level(-
4)","ts":1616794346.7030034,"logger":"daemon.drainhelper","msg":"releasing the lock
(bug mitigation)"}
{"level":"Level(-
4)","ts":1616794346.8040674,"logger":"daemon.updateInventory","msg":"obtained
inventory", "inv": {"sriovAccelerators":
[{"vendorID":"8086","deviceID":"0b32","pciAddress":"0000:20:00.0","driver":"","maxVirtualF
unctions":1,"virtualFunctions":[]},
{"vendorID":"8086","deviceID":"0d5c","pciAddress":"0000:af:00.0","driver":"pci-pf-
stub","maxVirtualFunctions":16,"virtualFunctions":
[{"pciAddress":"0000:b0:00.0","driver":"vfio-pci","deviceID":"0d5d"},
{"pciAddress":"0000:b0:00.1","driver":"vfio-pci","deviceID":"0d5d"}]}}}
{"level":"Level(-4)", "ts":1616794346.9058325, "logger": "daemon", "msg": "Update ignored,
generation unchanged"}
{"level":"Level(-
2)","ts":1616794346.9065044,"logger":"daemon.Reconcile","msg":"Reconciled","namespac
e":"vran-acceleration-operators","name":"pg-itengdvs02r.altera.com"}
```

3. カードの FEC 設定を確認します。

- deviceID: 0d5d

\$ oc get sriovfecnodeconfig node1 -o yaml

出力例

status: conditions: - lastTransitionTime: "2021-03-19T11:46:22Z" message: Configured successfully observedGeneration: 1 reason: Succeeded status: "True" type: Configured inventory: sriovAccelerators: - deviceID: 0d5c 1 driver: pci-pf-stub maxVirtualFunctions: 16 pciAddress: 0000:af:00.0 vendorID: "8086" virtualFunctions: - deviceID: 0d5d 2 driver: vfio-pci pciAddress: 0000:b0:00.0

driver: vfio-pci

pciAddress: 0000:b0:00.1

deviceID: 0d5d
 driver: vfio-pci

pciAddress: 0000:b0:00.2

deviceID: 0d5d
 driver: vfio-pci

pciAddress: 0000:b0:00.3

deviceID: 0d5d
 driver: vfio-pci

pciAddress: 0000:b0:00.4

- **Od5c** の値は、FEC デバイスの **deviceID** の物理機能です。
- Od5d の値は、FEC デバイスの deviceID 仮想機能です。

16.4.5. アプリケーション Pod アクセスおよび OpenNESS での FPGA 使用量の 確認

OpenNESS は、あらゆるタイプのネットワーク上のアプリケーションおよびネットワーク機能をオンボーディングおよび管理するために使用できるエッジコンピューティングソフトウェアツールキットです。

SR-IOV バインディング、デバイスプラグイン、ワイヤレス Base Band Device (bbdev) 設定、root 以外の Pod 内での SR-IOV (FEC) VF 機能など、すべての OpenNESS 機能がすべて連携していることを確認するには、イメージをビルドし、デバイスの単純な検証アプリケーションを実行できます。

詳細は、openess.org にアクセスします。

前提条件

- オプション: Intel FPGA PAC N3000 カード
- n3000-operator でインストールされる1つまたは複数のノード
- SR-IOV-FEC Operator でインストールされる1つまたは複数のノード
- Performance Addon Operator で設定されたリアルタイムカーネルおよび Huge Page
- cluster-admin 権限を持つユーザーとしてログインします。

手順

- 1. 以下のアクションを実行して、テスト用の namespace を作成します。
 - a. 以下の例のように **test-bbdev-namespace.yaml** という名前のファイルを作成し、**test-bbdev** namespace を定義します。

apiVersion: v1 kind: Namespace metadata:

name: test-bbdev

labels:

openshift.io/run-level: "1"

b. 以下のコマンドを実行して namespace を作成します。

\$ oc create -f test-bbdev-namespace.yaml

2. 以下の Pod 仕様を作成してから、YAML を pod-test.yaml ファイルに保存します。

```
apiVersion: v1
kind: Pod
metadata:
 name: pod-bbdev-sample-app
 namespace: test-bbdev 1
spec:
 containers:
 - securityContext:
   privileged: false
   capabilities:
    add:
     - IPC LOCK
     - SYS NICE
  name: bbdev-sample-app
  image: bbdev-sample-app:1.0 2
  command: [ "sudo", "/bin/bash", "-c", "--" ]
  runAsUser: 0 3
  resources:
   requests:
    hugepages-1Gi: 4Gi 4
    memory: 1Gi
    cpu: "4" 5
    intel.com/intel_fec_5g: '1' 6
     #intel.com/intel fec acc100: '1'
     #intel.com/intel fec Ite: '1'
   limits:
     memory: 4Gi
    cpu: "4"
     hugepages-1Gi: 4Gi
    intel.com/intel_fec_5g: '1'
     #intel.com/intel_fec_acc100: '1'
     #intel.com/intel fec Ite: '1
```

- 手順1で作成した namespace を指定します。
- 2 コンパイルされた DPDK を含むテストイメージを定義します。
- 3 コンテナーが root ユーザーとして内部で実行するようにします。
- 4 hugepage サイズ **hugepages-1Gi** および Pod に割り当てられる hugepage の量を指定します。hugepages および分離された CPU は、Performance Addon Operator を使用して設定する必要があります。
- **G** CPU の数を指定します。
- 6 N3000 5G FEC 設定のテストは、intel.com/intel_fec_5g でサポートされています。



注記

ACC100 設定をテストするには、#記号を削除して、intel.com/intel_fec_acc100 をコメント解除します。N3000 4G/LTE 設定をテストするには、#記号を削除して intel.com/intel_fec_Ite をコメント解除します。一度にアクティブにできるリソースは1つだけです。

3. Pod を作成します。

\$ oc apply -f pod-test.yaml

4. Pod が作成されていることを確認します。

\$ oc get pods -n test-bbdev

出力例

NAME READY STATUS RESTARTS AGE pod-bbdev-sample-app 1/1 Running 0 80s

5. リモートシェルを使用して pod-bbdev-sample-app にログインします。

\$ oc rsh pod-bbdev-sample-app

出力例

sh-4.4#

6. 環境変数の一覧を出力します。

sh-4.4# env

出力例

N3000_CONTROLLER_MANAGER_METRICS_SERVICE_PORT_8443_TCP_ADDR=172.3 0.133.131

SRIOV_FEC_CONTROLLER_MANAGER_METRICS_SERVICE_PORT_8443_TCP_PROT O=tcp

DPDK VERSION=20.11

PCIDEVICE_INTEL_COM_INTEL_FEC_5G=0.0.0.0:1d.00.0 1

~/usr/bin/env

HOSTNAME=fec-pod

- 1 これは、仮想関数の PCI アドレスです。**pod-test.yaml** ファイルで要求したリソースに応じて、以下の 3 つの PCI アドレスのいずれかを使用することができます。
 - PCIDEVICE_INTEL_COM_INTEL_FEC_ACC100
 - PCIDEVICE_INTEL_COM_INTEL_FEC_5G
 - PCIDEVICE_INTEL_COM_INTEL_FEC_LTE

7. test-bbdev ディレクトリーに移動します。

sh-4.4# cd test/test-bbdev/



注記

ディレクトリーは Pod にあり、ローカルコンピューター上にはない。

8. Pod に割り当てられている CPU を確認します。

sh-4.4# export CPU=\$(cat /sys/fs/cgroup/cpuset/cpuset.cpus) sh-4.4# echo \${CPU}

これにより、fec.pod に割り当てられた CPU が出力されます。

出力例

24,25,64,65

9. test-bbdev アプリケーションを実行してデバイスをテストします。

sh-4.4# ./test-bbdev.py -e="-I \${CPU} -a \${PCIDEVICE_INTEL_COM_INTEL_FEC_5G}" -c validation \ -n 64 -b 32 -I 1 -v ./test_vectors/*"

出力例

Executing: ../../build/app/dpdk-test-bbdev -l 24-25,64-65 0000:1d.00.0 -- -n 64 -l 1 -c validation -v ./test_vectors/bbdev_null.data -b 32

EAL: Detected 80 lcore(s)

EAL: Detected 2 NUMA nodes

Option -w, --pci-whitelist is deprecated, use -a, --allow option instead

EAL: Multi-process socket /var/run/dpdk/rte/mp socket

EAL: Selected IOVA mode 'VA'

EAL: Probing VFIO support...

EAL: VFIO support initialized

EAL: using IOMMU type 1 (Type 1)

EAL: Probe PCI driver: intel_fpga_5ngr_fec_vf (8086:d90) device: 0000:1d.00.0 (socket 1)

EAL: No legacy callbacks, legacy socket not created

._____

Starting Test Suite: BBdev Validation Tests

Test vector file = ldpc_dec_v7813.data

Device 0 queue 16 setup failed

Allocated all queues (id=16) at prio0 on dev0

Device 0 queue 32 setup failed

Allocated all queues (id=32) at prio1 on dev0

Device 0 queue 48 setup failed

Allocated all queues (id=48) at prio2 on dev0

Device 0 queue 64 setup failed

Allocated all queues (id=64) at prio3 on dev0

Device 0 queue 64 setup failed

```
All queues on dev 0 allocated: 64
== test: validation
dev:0000:b0:00.0, burst size: 1, num ops: 1, op type: RTE_BBDEV_OP_LDPC_DEC
Operation latency:
   avg: 23092 cycles, 10.0838 us
   min: 23092 cycles, 10.0838 us
   max: 23092 cycles, 10.0838 us
TestCase [0]: validation to passed
+ Test Suite Summary : BBdev Validation Tests
+ Tests Total:
+ Tests Skipped: 0
+ Tests Passed: 1 1
+ Tests Failed: 0
+ Tests Lasted: 177.67 ms
```

一部のテストはスキップできますが、ベクターテストに合格していることを確認してください。

16.5. 関連情報

- OpenNESS Operator for Intel® FPGA PAC N3000 (Programming)
- OpenNESS Operator for Wireless FEC Accelerators